

UNIVERSITY OF SOUTHAMPTON

# A Globally Wireless Locally Wired Hybrid Clock Distribution Network for Many-core Systems

by

Qian Ding

A thesis submitted in partial fulfillment for the  
degree of Doctor of Philosophy

in the  
Faculty of Engineering and Physical Science  
School of Electronics and Computer Science

February 2021



UNIVERSITY OF SOUTHAMPTON

ABSTRACT

FACULTY OF ENGINEERING AND PHYSICAL SCIENCES  
SCHOOL OF ELECTRONICS AND COMPUTER SCIENCE

Doctor of Philosophy

**A Globally Wireless Locally Wired Hybrid Clock Distribution Network for  
Many-core Systems**

by Qian Ding

Modern ICs are now facing critical issues on generating power-efficient and globally interconnected clock networks, as the clock distribution network might contribute to more than 50% of the overall power consumption. Besides, due to the increasing wire delay caused by shrinking interconnect dimensions, synchronous many-core systems are now facing challenges such as to propagate high-frequency clock signals across the chip with limited power budget. Delivering a clock with low uncertainties across active dies with large chip density has also become one of the major tasks using conventional metallic interconnects.

This thesis presents a novel architecture and comprehensive evaluations for a power-efficient and extremely low-delay approach using hybrid wire-wireless clock distribution network (CDN). The proposed hybrid CDN adopts wireless on-chip clock transmitters and receivers for broadcasting the clock signal globally. It then incorporates with conventional metal-based clock tree or mesh for local clock distribution. Comparisons between the proposed approach and two baseline architectures, namely a full fan-out tree and a global tree local mesh (TLM) structure, have been presented. Also, an accurate mathematical model with interconnect RLC parameters for the local clock distribution is employed.

The hybrid CDN has shown its superiority in terms of low clock delay, low clock skew and high energy efficiency compared with conventional solutions, which is evaluated via an industrial standard Arm Mali G77 GPU case study. Experimental results indicate that the proposed clock distribution network can achieve a significant global delay reduction of up to 28.8%. Also, on average, an up to 62.8% and 42.7% reduction in clock skew and power consumption, are identified, respectively, in our proposed test bench. Hence, our proposed approach offers a promising solution to clock distribution in future many-core integrated circuits, especially for high-performance systems.





# Contents

<b>Nomenclature</b>	<b>xv</b>
<b>Acknowledgements</b>	<b>xvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Challenges of the Conventional Interconnect . . . . .	1
1.2 Motivations for the Proposed Hybrid Clock Distribution Scheme . . . . .	6
1.2.1 The Impact of Wire Delay . . . . .	6
1.2.2 Synchronisation in Modern Many-Core Systems . . . . .	8
1.3 Thesis Contributions and Structures . . . . .	11
<b>2 Literature Review</b>	<b>15</b>
2.1 Reviewing of the Emerging Interconnect Techniques . . . . .	15
2.1.1 Wireless Interconnect . . . . .	15
2.1.2 Other Radio-frequency Based Interconnects . . . . .	17
2.1.3 Optical Interconnect . . . . .	19
2.2 Conventional CDN with Broadcast Architecture . . . . .	20
2.2.1 Combination of Binary-tree, H-tree, X-tree and Clock Grid . . . . .	21
2.2.2 Challenges of Existing Metallic Wire-based CDN . . . . .	24
2.3 Optical CDN with Broadcast Architecture . . . . .	26
2.3.1 Guided and Free-space Optical CDN . . . . .	26
2.3.2 Challenges of Existing Optical CDN . . . . .	28
2.4 Wireless CDN with Broadcast Architecture . . . . .	28
2.4.1 Wireless CDN with VCO and Frequency Divider . . . . .	29
2.4.2 Wireless CDN with Digital Modulation Techniques . . . . .	30
2.4.3 Challenges of Existing Wireless CDN . . . . .	32
2.5 Summary . . . . .	33
<b>3 Proposed Models and Algorithms for Local Wired CDN</b>	<b>35</b>
3.1 Delay Models of Conventional CDN . . . . .	35
3.2 Model of the Tree-based CDN . . . . .	41
3.2.1 Balanced Bi-partitioning Algorithm for Clock Tree Synthesis . . . . .	44
3.2.2 K-means Bi-partitioning (KBP) Algorithm for Clock Tree Synthesis . . . . .	47
3.2.3 Merging Tree Generation . . . . .	51
3.2.4 Buffer Insertion with Inverted Slew-estimation . . . . .	53
3.3 Model of the Mesh-based CDN . . . . .	55
3.3.1 Isolated Mesh Model . . . . .	55
3.3.2 Coupled Mesh Model . . . . .	57

3.3.3	Clock Receiver Planning . . . . .	59
3.3.4	Local Mesh Design . . . . .	60
3.4	Model of the Proposed Hybrid Wireless-Wired CDN . . . . .	62
3.5	Summary . . . . .	66
<b>4</b>	<b>Proposed Clock Transmitter and Receiver for Global Wireless CDN</b>	<b>67</b>
4.1	Clock Transmitter Design . . . . .	69
4.1.1	Voltage-Controlled Oscillator Design . . . . .	69
4.1.2	On-Off Keying Modulator Design . . . . .	71
4.2	Clock Receiver Design . . . . .	75
4.2.1	Low Noise Amplifier Design . . . . .	75
4.2.2	On-Off Keying Demodulator Design . . . . .	78
4.3	On-chip Antenna Design . . . . .	80
4.4	Summary . . . . .	83
<b>5</b>	<b>Experimental Setup and Evaluations of the Proposed Hybrid CDN</b>	<b>85</b>
5.1	Results of the Proposed Antenna for Hybrid CDN . . . . .	86
5.2	Results of the Proposed Clock Transmitter . . . . .	92
5.3	Results of the Proposed Clock Receiver . . . . .	95
5.4	Impacts of the Local Parameters on Global CDN Performance . . . . .	100
5.5	Summary . . . . .	108
<b>6</b>	<b>Case Study: the Proposed CDN Verified by Testbench Circuits</b>	<b>111</b>
6.1	Architectures of the Proposed Test Case . . . . .	111
6.2	Test Case Circuit Generations and Experimental Setup . . . . .	114
6.3	Test Case Results and Analysis . . . . .	116
6.4	Summary . . . . .	118
<b>7</b>	<b>Conclusions and Future Works</b>	<b>121</b>
7.1	Conclusions and Contributions of the Thesis . . . . .	121
7.2	Challenges of the Thesis . . . . .	123
7.3	Potential Future Works . . . . .	124
7.3.1	Short Term Future Works . . . . .	124
7.3.2	Long Term Future Works . . . . .	125
	<b>Bibliography</b>	<b>127</b>

# List of Figures

1.1	Simplified interconnect structure. . . . .	3
1.2	Simplified interconnect structures with wiring capacitance. . . . .	4
1.3	The propagation delay vs shrinking technology node for transistor gate delay, local communication delay and global communication delay, respectively, according to [1]. . . . .	5
1.4	Increasing interconnect delay within a specific 4-level H-tree network with (a) interconnect delay under different technology nodes using the same geometry, and (b) Increasing delay gap between 90 nm and 10 nm process. ©2021 <i>IEEE</i> . . . . .	6
1.5	A typical structure of (a) plan view of a conventional mesh-based CDN and (b) $\pi$ -model of the mesh wires. . . . .	7
1.6	The proposed hybrid clock distribution network which combines both global wireless interconnect and local metallic wires. . . . .	10
2.1	Wireless interconnect with WCube as application [2]. Wireless interconnect serves as the global communication channel. . . . .	16
2.2	A typical structure of (a) surface-wave interconnect (SWI) and (b) transmission line interconnect (Tx-line) for global communication. . . . .	18
2.3	An global H-tree constructed using optical interconnect to distribute clock signals. . . . .	19
2.4	Identifying the start and the endpoints of a typical clock tree in a design. . . . .	21
2.5	A typical structure of (a) single and big buffer/driver insertion strategy in a CDN and (b) the distributed and small buffer/driver insertion strategy in a CDN. . . . .	22
2.6	Buffered-tree with local clock grid for DEC Alpha 21064 clock distribution structure. . . . .	23
2.7	A typical structure of global tree and local mesh (TLM) clock distribution scheme. . . . .	23
2.8	A typical structure of interconnect model from PTM with (a) global and (b) local wiring structure [3]. . . . .	25
2.9	CDN architecture for global optical guided interconnect and local metallic H-tree. . . . .	26
2.10	Stack-up schematic of the optical CDN proposed in [4]. . . . .	26
2.11	Free-space optical interconnect with (a) unfocused and (b) focused clock broadcast architecture [5]. . . . .	27
2.12	General architecture of a wireless CDN with integrated wireless clock transmitter and receivers. . . . .	29
2.13	Block diagrams of wireless CDN [6] with (a) clock transmitter and (b) clock receiver. . . . .	30

2.14	Wireless CDN transceiver architecture in [7] with 5.6mm propagation distance. . . . .	31
2.15	Simplified 64-QAM wireless interconnect architecture [8] with (a) transmitter and (b) receiver. . . . .	32
3.1	Conventional metallic interconnect structure and its lumped RC model. .	36
3.2	Distributed lumped circuit structure for improved accuracy of interconnect model. . . . .	37
3.3	Conventional tapered H-tree network as baseline architecture with 16 fan-out nodes. . . . .	39
3.4	A $k$ -level H-tree modelled as a folded binary RC tree with $k^2$ fan-out. . .	40
3.5	Transmission line model with distributed elements $R_u$ , $L_u$ , $G_u$ , and $C_u$ , respectively. . . . .	41
3.6	Equivalent ABCD representation of an interconnect in transmission line model. . . . .	42
3.7	Top view of the Balanced Bi-partitioning (BB) Algorithm with half of the sinks on the boundary as a reference set. The clock sinks within the octagon that are closer to the reference set A are partitioned into subset A, the rest are grouped into subset B. Also, the summation of the capacitance of clock sinks are minimised. . . . .	45
3.8	Example of the merging segment generation, with $MS_a$ , $MS_b$ as a tuple element and $MS_c$ , $MS_d$ as a tuple element, respectively. The generated parent merging segments are the intersection region of the two TRRs with the respective child merging segments as cores and calculated/balanced edges as radius. . . . .	46
3.9	Example of the topology tree generation with a sink set containing 8 clock endpoints. The layer 1 MS are constructed in iteration 1 with tuple 1 (8 sinks) as input, the layer 2 MS are constructed in iteration 2 with tuple 2 (layer 1 MS) as input. The iterations are executed from a bottom-up fashion, until the root merging segment has been generated. . . . .	46
3.10	Illustration of the xy-cut, using (a) a conventional balanced bi-partition in [9], and (b) our proposed xy-cut algorithm for minimum wire usage. The total wire reduction between the proposed method and the conventional method is around 15.1% in this example, thus improving energy efficiency. .	49
3.11	Example of the merging tree generation with two sink sets, each of which contains 8 clock endpoints (not to scale). The two subtrees are merged with root segment $q$ . . . . .	51
3.12	Illustration of the (a) proposed buffer insertion method with inverted slew-estimation and (b) its time-domain behavior when the estimated slew is smaller than the slew lower bound of the driver. Without a clock driver, the child clock end point might violate the slew constraint. . . . .	54
3.13	Partial top views of the proposed isolated model, including transfer function zone (TF zone) and lumped zone. Unit cells are considered independent of each other. ©2021 IEEE . . . . .	56
3.14	Partial top views of the proposed coupled model (not to scale), including transfer function zone (TF zone), $3-\pi$ zone and lumped zone. Unit cells are considered coupled to each other shown in the overlapped region. ©2021 IEEE . . . . .	57

3.15	Typical structure of an RLC $3-\pi$ model [10] which has been used in our proposed model for representing interconnect segments longer than 100 $\mu\text{m}$ inside $3-\pi$ region. . . . .	59
3.16	Top view of (a) distributed planning and (b) concentrated planning in an arbitrary $i \times j$ local mesh network. . . . .	59
3.17	The proposed hybrid clock distribution network which combines both global RF-I and local metallic wires. . . . .	63
3.18	Estimated results using the proposed hybrid model and a conventional H-tree CDN with 96% average delay reduction. . . . .	65
4.1	Block diagrams of the proposed (a) clock transmitter and (b) clock receiver. . . . .	68
4.2	General oscillatory condition of feedback system. . . . .	69
4.3	Circuit schematics of the proposed (a) tuned stage and (b) the proposed balanced VCO with cross-coupled NMOS pair and source follower buffers. . . . .	70
4.4	Typical structure of the (a) single and (b) complementary implementation of a MOS switch. . . . .	72
4.5	1-to-2 Clock buffer with (a) unbalanced propagation delay and (b) balanced propagation delay. . . . .	72
4.6	The schematic of MOS switch-based OOK modulator with differential signaling. . . . .	73
4.7	Schematic of the proposed OOK modulator with differential signaling and leakage compensation cross-coupling NMOS pairs. . . . .	74
4.8	Schematics of the proposed (a) pseudo-differential Inductive degeneration LNA and (b) its high frequency single-stage small signal model. . . . .	75
4.9	Schematic of the proposed gain-boosting load, with a conventional NMOS rectifier to rectify the input radio frequency signals. . . . .	78
4.10	Schematic of the proposed Proposed OOK demodulator with gain-boosting technique and noise-compensation output buffer. . . . .	79
4.11	Top view of the proposed meander monopole antenna (MMA) structure with top copper layer. . . . .	81
4.12	Proposed meandering monopole antenna (MMA) structure with EM simulation model setup (not to scale). . . . .	81
4.13	Top view of the proposed meander dipole antenna (MDA) structure with top copper layer. . . . .	82
4.14	Proposed meandering dipole antenna (MDA) structure with EM simulation model setup (not to scale). ©2021 IEEE . . . . .	82
5.1	Complete antenna-circuits co-simulation flow diagram for global wireless CDN verification. . . . .	87
5.2	Top view of the test structure for the study of EM crosstalk between TRx antenna and the nearby interconnects with different separation distance. The antenna and nearby lines act as the victim of the crosstalk in the experiment, respectively. ©2021 IEEE . . . . .	88
5.3	EM crosstalks versus the frequency of interest in (a) victim interconnects and (b) victim antenna, measured in $S_{21}$ . The differential signaling can reduce the common-mode noise by 35 dB averagely, in the victim lines. . . . .	89
5.4	Antenna gain under matched condition with 3 separation distances ranging from 1.5 mm to 5 mm, which shows it's suitable for the proposed hybrid CDN. ©2021 IEEE . . . . .	90

5.5	EM simulation in terms of S-parameter results for (a) the proposed meander monopole antenna (MMA) antenna and (b) the proposed meander dipole antenna (MDA) antenna, with 3 separation distances ranging from 1.5 mm to 5 mm. . . . .	91
5.6	Directivity of the proposed meander dipole antenna (MDA) with 1.5 mm Tx and Rx separation, which shows it's suitable for our proposed hybrid CDN. . . . .	92
5.7	Output clock signal with an input 2.5 GHz nominal clock for the proposed (a) switch-based modulator (SWM) and (b) leakage-compensation modulator (LCM). . . . .	93
5.8	65 nm leakage compensation modulator (LCM) and switching modulator (SWM) output in frequency domain with the proposed cancelling gate M5 and M6. . . . .	94
5.9	Rx front-end results in terms of the S-parameters of the proposed pseudo-differential low-noise amplifier (lna) and noise figure at 69 GHz. . . . .	95
5.10	Waveform of the (a) input 2.5 GHz nominal clock signal, (b) rectifier output with inverted signal polarization, (c) the full-swing demodulator output signal at 1.2 V and (d) effective signal stack-up. . . . .	96
5.11	Example of 7 GHz recovered clock signal with 20.8% reduced signal swing. Clock amplitude will keep degrading with the increase of input clock frequency, transistors in clock Rx circuit cannot switch on/off completely in a reducing clock cycle, hence generating attenuated recovered clocks. . . .	97
5.12	Comparisons between the measured eye diagrams of the recovered clock signal from conventional wired global tree local mesh (TLM) CDN and our proposed wireless global CDN at 2.5 GHz, with robust eye height/clock amplitude around 1.2 V with 50 mV additive noise at the output of clock buffers. ©2021 <i>IEEE</i> . . . . .	98
5.13	Clock skew with a pseudo-mesh based topology at 2.5 GHz using 65 nm Technology. (a) An overall 16 fan-out architecture is adopted with an average skew reduction of 45.5% under random unbalanced load ranging from 10 fF to 500 fF, and (b) unbalanced load allocation (NBFS) for skew minimisation. . . . .	99
5.14	Propagation delay comparison between H-tree and proposed wireless CDN. ©2021 <i>IET</i> . . . . .	100
5.15	Power comparison between global H-tree and proposed global wireless CDN. ©2021 <i>IET</i> . . . . .	101
5.16	Comparison between the TLM approach and the proposed hybrid architecture of (a) interconnect delay and (b) normalised delay reduction according to size and load variation with a maximum 28.8% decrease under a large model scale. ©2021 <i>IEEE</i> . . . . .	103
5.17	Comparison between the TLM and the proposed hybrid architecture of (a) global clock skew and (b) normalised skew reduction according to size and load variation with a maximum 68.1% decrease under large load unbalance, within one clock domain. ©2021 <i>IEEE</i> . . . . .	104
5.18	Comparisons between the TLM approach and the proposed hybrid architecture of (a) overall power consumption, (b) power reduction according to local clock frequency and load variation with a maximum 32% decrease under large local clock frequency. ©2021 <i>IEEE</i> . . . . .	105

5.19	Different parameters in terms of (a) output skew, (b) clock latency and (c) power consumption in the test case for both distributed and concentrated planning. Concentrated planning trade higher propagation delay with better skew and power performance using coupled model. . . . .	106
5.20	Total skew and normalised cost in the local region with increasing iterations. ©2021 <i>IEEE</i> . . . . .	107
5.21	Output skew in local CDN under different input global skew caused by different communication distance or buffer mismatch. The coupled model shows an average 5.7% error comparing to SPICE simulation. ©2021 <i>IEEE</i>	108
6.1	Models of a typical conventional global-tree local-mesh structure with (a) top view, (b) cross-sectional view, which consumes a large amount of metal. The flipped die could use a different technology node than the active silicon interposer, e.g. the flipped die is in 7 nm process and the interposer is in 65 nm process, which provides flexibility to the integration of different technologies. . . . .	112
6.2	Models of our proposed hybrid clock distribution network, with (a) top view and (b) cross-sectional view, respectively. The flipped die could use a different technology node than the active silicon interposer, e.g. the flipped die is in 7 nm process and the interposer is in 65 nm process, which provides flexibility to the integration of different technologies. . . .	113
6.3	Unevenly distributed test case $T_{GPU}$ extracted from Arm Mali-G77 GPU with (a) an example top view and (b) time-domain responses for $T_{GPU}$ with 1.7 GHz recovered global clock input in an 8-core system. ©2021 <i>IEEE</i>	115
6.4	Power output in local region with incremental iterations in the proposed test case $T_{GPU}$ , with an wireless/wire crossing point near 3.5 mm. . . . .	118





# List of Tables

2.1	Interconnect parameters for global and local wires from 180 nm process to 7 nm process in terms of both resistance and capacitance per unit length.	25
2.2	Qualitative analysis of the CDN using the conventional and emerging interconnect techniques.	34
5.1	Comparisons between different implementation of the wireless transmitter in related researches and our proposed LCM-based clock transmitter.	94
5.2	Comparisons between different implementation of global CDNs using wire and wireless approaches	99
6.1	Overall CDN Performance Comparison for Test Case $T_{GPU}$ using Concentrated Receiver Planning	117
6.2	Overall CDN Performance Comparison for Test Case $T_{GPU}$ using Distributed Receiver Planning	117



# Nomenclature

$A_k$	Amplitude of the $k$ th-order harmonic of a Fourier series representation of the local clock signal.
$\alpha$	Coefficient for horizontal interconnect segment.
$\beta$	Coefficient for vertical interconnect segment.
$cent_a$	Centroid vector of sink set $a$ .
$cent_b$	Centroid vector of sink set $b$ .
$C_{\alpha,\beta,\gamma}$	Cost coefficients for power, area and skew, respectively.
$C_i$	Capacitance per unit length for interconnect segment $i$ .
$C_L$	Lumped capacitive loading of local clock sinks.
$C_t$	Number of all clock endpoints within current TF zone.
$D(i, j)$	50% logic transition delay from point $i$ to point $j$ .
$d_{m,n}$	An arbitrary clock sink in a $m \times n$ local mesh network.
$e_a$	Edge length connecting set $a$ to its parent set.
$e_b$	Edge length connecting set $b$ to its parent set.
$G_i$	Admittance per unit length for interconnect segment $i$ .
$H_i(s)$	Transfer function of the interconnect segment $i$ .
$Horz(i)$	Horizontal segment condition for the $i$ -th layer in a CDN.
$I$	Iteration count to increment mesh cells in Algorithm 3.
$I_{max}$	Max iteration count to avoid deadlock in Algorithm 3.
$l_i$	Physical length of the interconnect segment $i$ .
$l_{ub}$	Overall interconnect length upper bound of a design.
$i_{max}$	Max iteration count to avoid deadlock in Algorithm 1.
$index_a$	Index vector of sink set $a$ .
$index_b$	Index vector of sink set $b$ .
$L_i$	Inductance per unit length for interconnect segment $i$ .
$N_0$	Root node of the given CDN.
$n_b$	Number of neighbor meshes of an arbitrary mesh cell.
$num_x$	Number of single mesh segments in x direction.
$num_y$	Number of single mesh segments in y direction.
$PC_k$	Power cost evaluated by wire length for sink set $k$ .
$P_i$	Power consumption of a single interconnect segment.
$P_{sum}$	Overall power consumption of a local network.

---

$P_{ub}$	Overall power upper bound of a design.
$P$	Probability vector in $xy$ -cut procedures.
$P_x$	Probability interval vector in $xy$ -cut procedures.
$R_d$	Output resistance of an arbitrary local clock buffer.
$Re(Z_{in})$	Real part of the complex input impedance.
$R_i$	Resistance per unit length for interconnect segment $i$ .
$s_{m,n}$	An arbitrary clock source in a $m \times n$ local mesh.
$sk_{ub}$	Overall skew upper bound of a design.
$T$	Period of the clock signal.
$Vert(i)$	Vertical segment condition for the $i$ -th layer in a CDN.
$v$	Velocity of the electromagnetic wave.
$WL_m(a, b)$	Merging length of the sink set $m$ with subset $a$ and $b$ .
$WL_s(m)$	Fusion length of the sink set $m$ .
$Z_{\alpha, \beta, \gamma}$	Effective impedance for $3\text{-}\pi$ interconnect model.
$Z_d$	Effective output impedance of local buffer.
$Z_{in}$	Input impedance seen at an arbitrary local clock buffer.
$Z_{load_i}$	Effective load impedance at the end of interconnect $i$ .
$Z_{o_i}$	Characteristic impedance of interconnect segment $i$ .
$\delta$	Convergence coefficient for topology tree generation.
$\gamma_i$	Propagation constant of interconnect segment $i$ .
$\theta_{m,n}$	Set of all unit mesh regions for a $m \times n$ mesh network.
$\tau$	Logic transition time of the clock signal.

## **Acknowledgements**

I would like to thank my supervisor and friend, Dr Terrence Mak, for his support and invaluable guidance as well as our friendship through my PhD period. I have been extremely fortunate to have such supervisor who cared so much about my work and responded to my questions so promptly.

Also, I would like to offer my thanks to the staffs and students in CPS group, for their kind help and advice, which can always encourage me when needed. Finally, my wife and my parents have always been my strongest supporters. I am greatly thankful for their love, and I would like to dedicate this work to them.



*To my family and friends who always support me during these  
meaningful and exciting years.*





# Chapter 1

## Introduction

Clock signal is the heartbeat of a chip. It is one of the most commonly adopted global signals in a synchronous system, served as a timing reference to control the timing behaviour of all sequential logic on a silicon die. To deliver the clock signal across the entire core area, a unique routing network is used, namely clock distribution network (CDN). Conventional CDNs adopt a tree or mesh-based topology to evenly deliver a global clock signal along the designated wires inside CDN to each of the local logic which requires a clock signal. Previously, the design of a CDN is less of an overall concern during the VLSI design stages, however, with the shrinking of the CMOS fabrication process, it now becomes one of the major challenges for very large scale integrated circuits (VLSI) designers.

With the fast development of technology, modern ICs have already beyond billion transistors scale. As the process technology scaling down, it becomes more and more sophisticated for chip designers to make trade-offs between power and system performance. In the state-of-art technology nodes such as 7nm or 5nm, dimensions of conventional metallic wires are getting smaller and raised a lot of new physical challenges. As the signal propagation delay has already exceeded circuitry gate delay, according to [1] and [11], the performance of interconnects inside an IC becomes the dominant factor which limits the performance of a chip. Hence, it is essential to consider the interconnect delay when designing a combination or sequential logic if there's any potential global signal, such as the clock signal input to this logic.

### 1.1 Challenges of the Conventional Interconnect

Different researchers such as [12] and [13], have suggested that the increasing RC characteristic of the interconnect has become the major constraint for clock distribution in synchronous VLSI systems, as the clock distribution network will consume more than

half of the overall power budget. Furthermore, another challenge is that the different arrival time of the clock signal (clock skew) significantly affects the timing performance of a synchronous system, especially for the logic cells which are correlated on the data path. Fixing timing convergence issues caused by clock skew via inserting clock buffers and delay elements in one region will result in new timing violations in other regions, thus producing design iterations and extra works. Major reasons for the time convergence challenge can be elaborated as follows:

1. The increasing propagation delay of the interconnect wires becomes the major delay contribution of the overall delay elements inside the chip, because of the increasing product of unit length resistance and capacitance. Since wiring density keeps increasing with the shrinking technology, space between interconnects is getting smaller, which will in turn, increases the coupling capacitance and cross-talk noise between wires. Besides, since the resistance per unit length of a conductor is inversely proportional to the cross-sectional area of an interconnect segment, the electron scattering effect [14], [15] of the thinner wires also leads to an increasing overall wiring resistance. Hence, for the first order approximation of the propagation delay of an interconnect segment, the delay constant  $\tau$  [14] essentially becomes bigger, thus limiting overall clock period between two subsequent sequential elements.
2. CMOS process variations are becoming unpredictable and difficult to control, as the physical size of the fabricated chips is approaching to their physical limit. For example, in the state-of-art 7nm process, the minimum interconnect pitch width is around  $0.03 \mu\text{m}$ . A slight fabrication variation near 2 nm will cause a  $\pm 7\%$  variation compared to the standard wire dimension, thus cause unnecessary signal distortion in this specific interconnect segment. In addition, the thickness of the gate oxide of a typical transistor is almost near the width of a few atoms. Even a slight change of one atom in width can cause different switching speed among all transistors, thus producing unpredictable circuit delays.
3. The ever-increasing clock frequency and design complexity, especially for high-performance systems, makes it difficult to get a positive slack during static timing analysis (STA) stage. Thus, designers need to spend more time on iteratively checking the timing constraints and the design itself.

From the above points, to provide an accurate timing reference for system timing convergence, it becomes one of the top challenges when designing a clock distribution network using the conventional interconnects. The basic structure of a conventional wire model is given as follows shown in Figure 1.1. As interconnect unit length resistance is inversely proportional to the cross-section area of a metallic conductor according to Equation 1.1, the shrinking wires essentially lead to a situation with increasing wire RC characteristics.

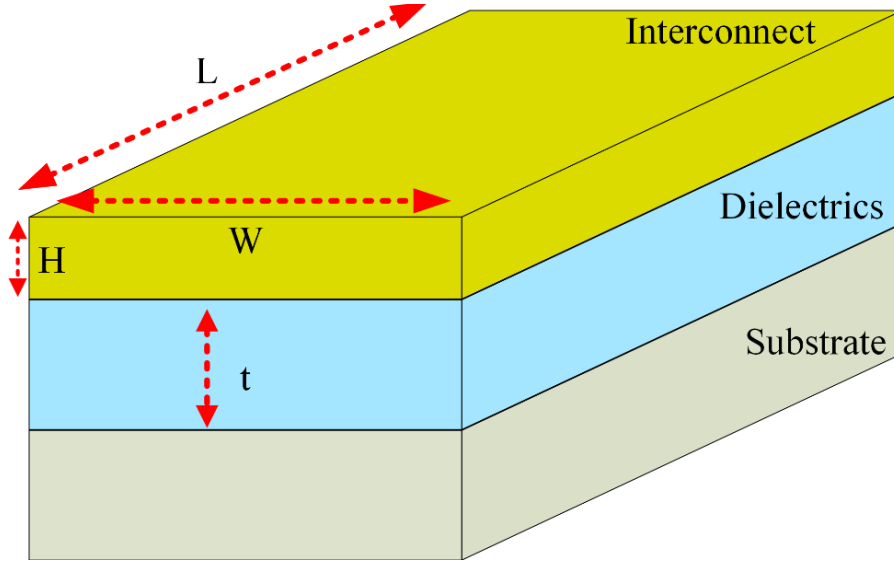


Figure 1.1: Simplified interconnect structure.

Wire delay, especially global interconnect suffers from this phenomenon significantly. Although there are several solutions trying to tackle out this problem, such as adopting new interconnect materials or inserting global repeaters to reduce the total length of the global wires, the overall delay gap between gate and wire delay is still increasing.

For normalised wire resistance according to 1.1, it can be given that:

$$R = \rho \cdot \frac{L}{A} = \rho \cdot \frac{L}{W \cdot H} \quad (1.1)$$

where  $\rho$  is the resistivity of a metallic conductor,  $L$  represent the total length of this wire and  $A$  denotes the cross-section area of the wire with the dimension of  $W$  as wire width and  $H$  as wire height. Hence, we could naturally get:

$$r_{wire} = \frac{\rho}{A} \quad (1.2)$$

where  $r_{wire}$  is the unit length resistance of the wire of interest.

Besides, to model the interconnect capacitance, it could be first simplified as a parallel plate capacitor model with dielectric material filled in between the two conductors according to Figure 1.2:

$$C_{wire} = W \cdot L \cdot \frac{\varepsilon}{t} \quad (1.3)$$

where  $\varepsilon$  is the dielectric constant of the dielectric material usually as  $SiO_2$ , and  $t$  is the spacing between the interconnect and the silicon substrate. With the CMOS technology node keeps shrinking, the wire capacitance should decrease as the width  $W$  decreases so that this could somehow reduce the interconnect capacitance so as to compensate the total wire delay caused by RC product. However, as the dimensions of the wire

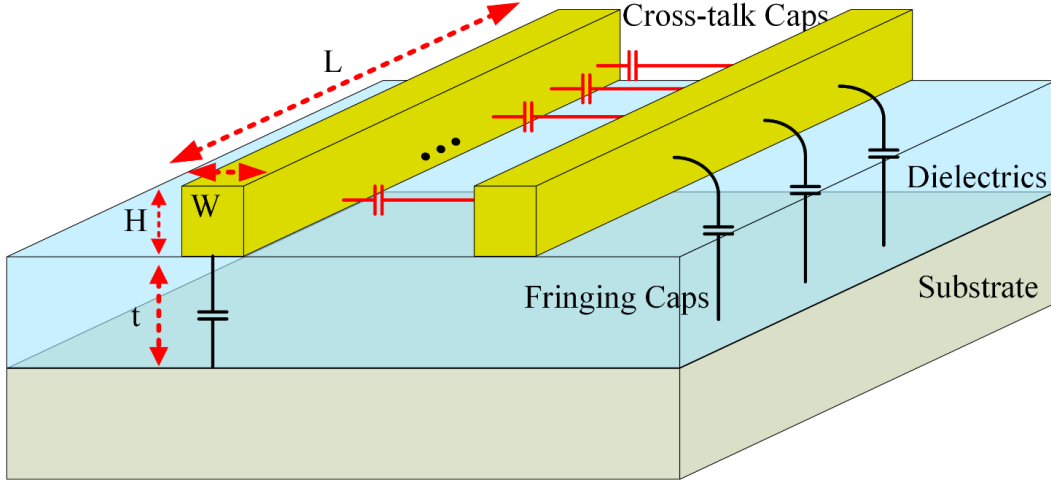


Figure 1.2: Simplified interconnect structures with wiring capacitance.

are getting smaller with the increasing wire density, the  $W/H$  ratio has already dropped below unity as shown in Figure 1.2. Under such circumstances, the original parallel plate capacitor is no longer accurate to model the interconnect capacitance. As Figure 1.2 shown below, this situation will lead to a sidewall effect which the side of any interconnect conductor will generate a fringing capacitance and this value will contribute to the overall wire capacitance. For the system working on high frequencies, those parasitic may severely reduce system performance in terms of signal integrity, propagation delay, etc. Cross-talks between interconnects will also be a troublesome problem because the sidewall area of a wire is essentially larger than the bottom or top surface area and hence producing the dominant contribution to wire capacitance, therefore generating interference and noise thus lead to erroneous signal transfer. The total normalised signal propagation delay within a conductor could then be modelled as:

$$T_{pd} = \tau r_{wire}(Lc_{wire} + C_{load}) = \tau^2 \frac{\rho}{H \cdot W} (c_{cross} + c_{fringing} + \frac{C_{load}}{L}) \quad (1.4)$$

where  $\tau$  is a conductor geometry-dependent constant,  $L$  is the total wire length,  $r_{wire}$  is the wire resistance per unit length,  $C_{load}$  is the load capacitance,  $c_{wire}$  is the total capacitance per unit length which consists of cross-talking capacitance  $c_{cross}$  and side-wall fringing capacitance  $c_{fringing}$  respectively. Using the above equation, we could consequently get the fact that the overall global signal propagation delay inside a metallic global interconnect keep increasing pseudo-exponentially as the geometry of the wire getting smaller. The above elaboration of the delay constant essentially presents the major reason for the increasing delay gap shown in Figure 1.3.

Besides, the increasing wire delay and chip density can also lead to larger global communication delay with the scaling up broadcast/multicast demands. The techniques utilised by modern many-core systems, such as cache-coherence protocol, will need to

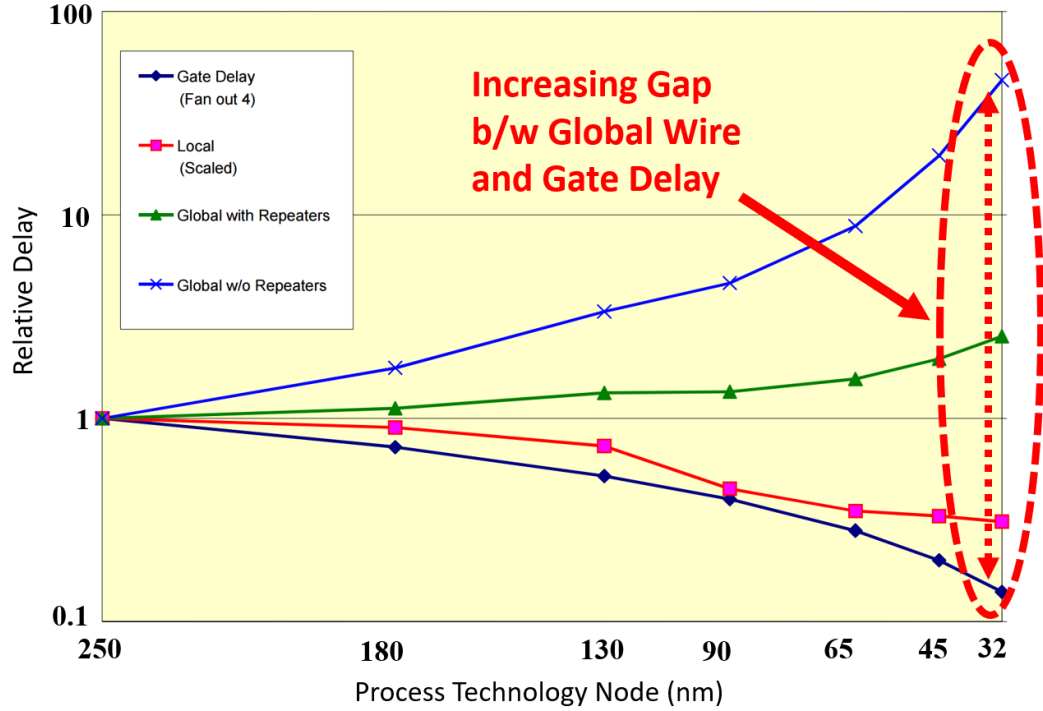


Figure 1.3: The propagation delay vs shrinking technology node for transistor gate delay, local communication delay and global communication delay, respectively, according to [1].

send calculated or updated data to the remaining nodes and cores to maintain the coherence of data in all cores. Therefore, this situation generates an increasing hop count which leads to large data traversal time/delay and power consumption if a conventional duplicated unicast scheme is still adopted.

To sum up the above existing problems, conventional interconnects seem to have reached their bottleneck in some particular applications, because of the reduced size and the increasing new communication requirements. To optimise this limitation, novel interconnect techniques are now emerging to meet the needs of new design targets and tackle out the existing drawbacks. In order to better incorporate with the existing CMOS process, most of these newly developed interconnects are mostly CMOS compatible which means they could be directly fabricated via standard CMOS developing procedures. Other novel interconnect techniques require extra on-chip components such as wave guides, forking points, etc. Further details related to the emerging interconnects and their application as a component of the CDN will be given in Chapter 2.

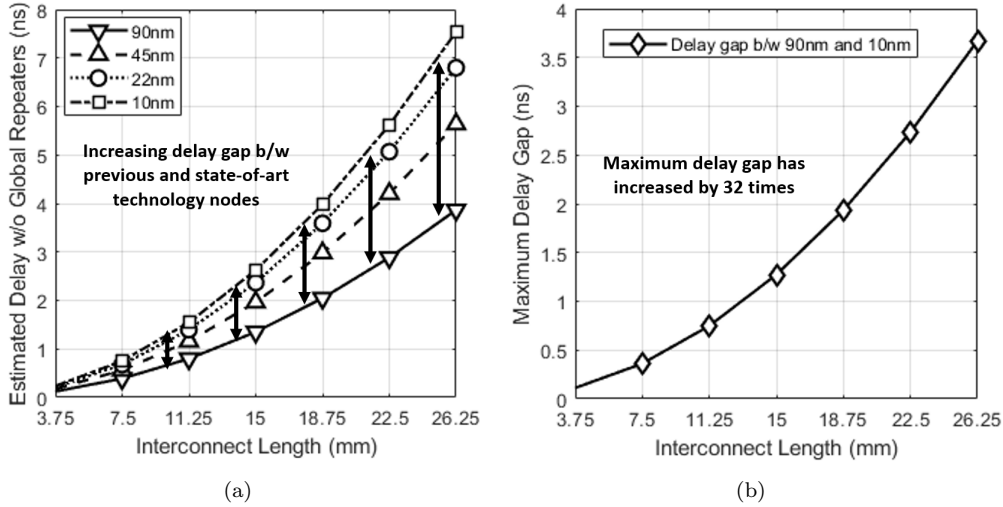


Figure 1.4: Increasing interconnect delay within a specific 4-level H-tree network with (a) interconnect delay under different technology nodes using the same geometry, and (b) Increasing delay gap between 90 nm and 10 nm process. ©2021 IEEE

## 1.2 Motivations for the Proposed Hybrid Clock Distribution Scheme

### 1.2.1 The Impact of Wire Delay

A balanced clock distribution network (CDN) is one of the most essential applications of interconnects to propagate clock signals with equally wire delay inside a sequential system. Conventional CDN incorporates metal base wires with different topologies to achieve the highest clock frequency as well as the smallest clock uncertainty including clock skew and clock jitter. Typical CDNs such as X-tree or H-tree network which is similar to a folded binary tree, which fanout to two new branches with similar length, at each end of a tree segment.

For an ideal H-tree network, each route from the clock source to the sink is equal in length, therefore generating minimal clock skew [16], [17]. However, as the load on a realistic chip would not be spread evenly, there could be difference in terms of route length and capacitive loadings among the root and sink pairs, hence generating different propagation delay and subsequently, clock skew. In order to avoid any setup or hold violation leading towards metastability, a synchronous system need to select the worst-case propagation delay as the clock period. The clock skew will possibly limit overall system performance, as it affects the system setup/hold margin.

As a global signal, clock tree should use clock buffers to improve driven strength and compensate different route delay and signal decay. A large number of clock buffers will be added onto the clock route to meet the user-specified design target during clock tree synthesis, thereby increasing power consumption and area occupation.

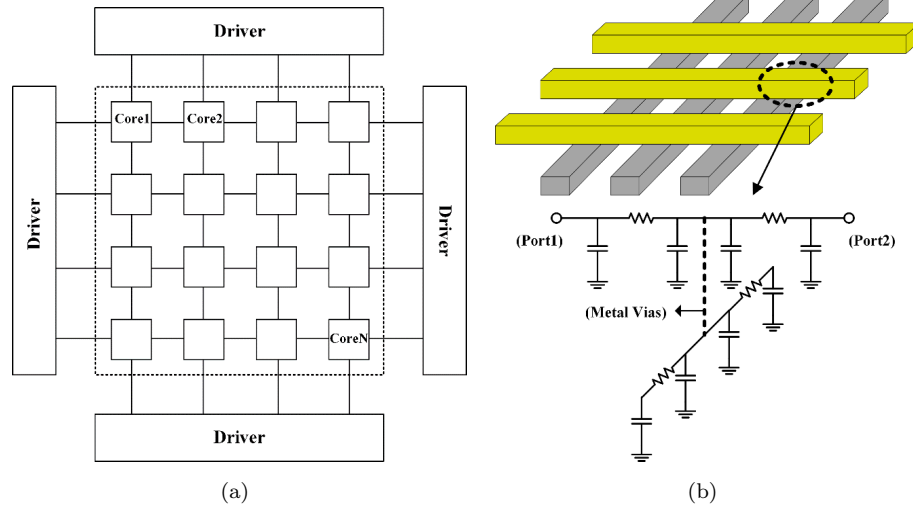


Figure 1.5: A typical structure of (a) plan view of a conventional mesh-based CDN and (b)  $\pi$ -model of the mesh wires.

Clock latency is another issue to be considered about, as the delay from clock root to an arbitrary clock sink is also a target to minimise during clock tree synthesis (CTS) stage. Clock routes with longer length lead to higher performance degradation caused by on-chip variation (OCV). Also, the clock uncertainties are proportional to the overall clock delay. Hence to design a CDN with reasonable performance, modern EDA tools are treating the clock latency as one of the most important goals during CTS stage. As an example, the global clock distribution delay has almost doubled its value when distributing a unit step signal in a designated H-tree network in 10 nm process compared with the same geometry using 90 nm process, which is presented in Figure 1.4. Therefore, the conventional interconnect seems to be reaching its limit for high-frequency global communication requirements such as broadcasting clock signals, because of the increasing delay.

As shown in Figure 1.4 from [18], a comprehensive clock delay comparison is given by using Elmore delay on a 4-level 16-fanout H-tree topology with its terminals tied to several 10 fF capacitive loadings. This model has a wirelength of the clock tree ranging from 3.75 to 26.25 mm and clock buffers allocated at the branching point of each level. Under the assumptions that this H-tree is implemented using top metals to maintain relatively small clock latency, based on the data from Predictive Technology Model (PTM), the RC product for global communication in 90 nm node has almost doubled its value in the more advanced technology node in 10 nm [18]. Since the Elmore delay of an interconnect segment is proportional to the product of the unit length resistance and capacitance, the clock latency from the clock source to an arbitrary clock tree fanout will also increase to twice of its original value. Besides, an increasing delay gap between 90 nm implementation and the 10 nm implementation can be found with the increase of

the CDN wirelength, hence request a fine-grain design of the CDN for applications with tight latency constraint.

Clock mesh/grid is another solution which produces less clock uncertainty for a high-performance system. Shown in Figure 1.5, as a simple but efficient topology, clock grid is easy for designers to use and generate as a mesh-based interconnect work is implemented for global clock transmission typically using the top two metal layers. However, because of the wiring requirements, the characteristic of a clock grid will lead to a high power and area consumption when compared with H-tree or X-tree when the same synchronisation area requested. The requirement of clock straps will also increase the demand for routing resources. Spine-based CDN is the balance between binary trees and clock grid. As the global clock CDN is in the form of a spine, local CDN still need to be customised with either comb or tree structures.

### 1.2.2 Synchronisation in Modern Many-Core Systems

Modern many-core systems adopt GALS structure to provide a solution to clocking and synchronisation problems. Based on the key techniques used, GALS system can be identified into 3 different categories, namely pausable clock, asynchronous interface and loosely synchronous.

First of all, pausable clock exists among asynchronous blocks with uncertain clock relationships. It performs certain types of handshaking and simply delays or pauses the sampling clock edge in a clock domain after the arrival of valid data signals from another clock domain, thus mitigating metastability [19]. Some challenges in designing a pausable clock GALS system are related to the design of a ring oscillator with robust and promising performance [20], [21]. As the clock signal and clock generators are delayed and restarted subject to clock domain crossing requirements, the varying cycle-to-cycle clock jitter generated by the ring oscillator will be accumulated and propagated along the down stream clock network and further impact the timing margin. Therefore, the pausable clock technique remains to be a niche technique.

A fully-asynchronous system means the blocks belong to different clock domains with irrelevant clock frequency and phase information. This kind of GALS system requires synchronisation circuits (synchronisers, asynchronous bridges) to mitigate the impact of clock domain crossing. One of the major concerns of using an asynchronous interface is related to the modelling, validation and delay analysis of synchronisation circuits. [19] indicates a rule of thumb design standard that at least 40 gates delay should be reserved for a valid and stable CDC signal. Based on this rule, asynchronous GALS style can be used in applications that can tolerate extra synchronisation delay, or that have low clock frequencies, which is not the first choice for high-performance GALS design.



Another GALS type contains loosely synchronous interfaces. This type of the CDC interface can offer a better performance because of the reduction of handshaking circuits, thus providing improved delay [19], [22]. The loosely synchronised GALS itself can be divided into three groups, namely *mesochronous*, *plesiochronous*, and *heterochronous* [23], based on different frequency and phase relations between each block in the same system.

*Mesochronous* yields same operation frequency and unknown phase information among blocks. *Plesiochronous* offers the same nominal operation frequency for all blocks with minor random frequency variation, thus creating and accumulating phase drift. *Heterochronous* operates at different nominal frequencies with periodic phase relationships between each block, which is similar to generated clock in a synchronous design.

To sum up, GALS based many-core system offers flexibility of synchronisation at the cost of extra interfaces between blocks, especially for pausable clock and asynchronous GALS, which lead to inevitable delay and impact system throughput. Loosely synchronous can benefit from the reduce of synchronisation circuits thus providing better performance, but still, it will require a global clock (with arbitrary clock latency and without tight skew constraint) or multiple PLLs under the same operation frequency. With the scale of the GALS system keeps growing, the wiring delay of conventional interconnects will eventually generate more challenges, which will impact system performance and power.

Alternatively, utilising the emerging interconnects as global distribution route of the CDN becomes available in recent decades. For instance, using optical or wireless interconnect including transmitter and receiver could effectively improve the performance and power consumption of CDN. As the radiated EM wave just need one hop from the source to the sink, overall signal propagation delay would be reduced rapidly. Moreover, taking advantage of the multicast feature, a wireless CDN becomes highly competitive comparing to conventional high-performance CDN with a much lower global power and global delay with moderate complexity overhead. However, similar to the spine or mesh-based CDN, the local network of wireless CDN also need to be customised manually. Further details of the novel interconnect and different clock distribution network would be given in Chapter 2.

Newly emerged interconnects in recent decades have significantly enlarged the area of application. Most of the developing interconnects are targeted at reducing multicast delay and enhancing fan-out architecture, which provides a positive potential to future many-core architecture that requires high speed and high-quality global communication. Hence, as shown in Figure 1.6, the proposed hybrid clock distribution network incorporates both a global-broadcasting-based interconnect with an integrated antenna inside the clock transmitter and receiver, and the conventional wire-based local network, such that the global delay and power can be potentially reduced significantly at the cost of moderate area and complexity overhead. This type of clock distribution can not only

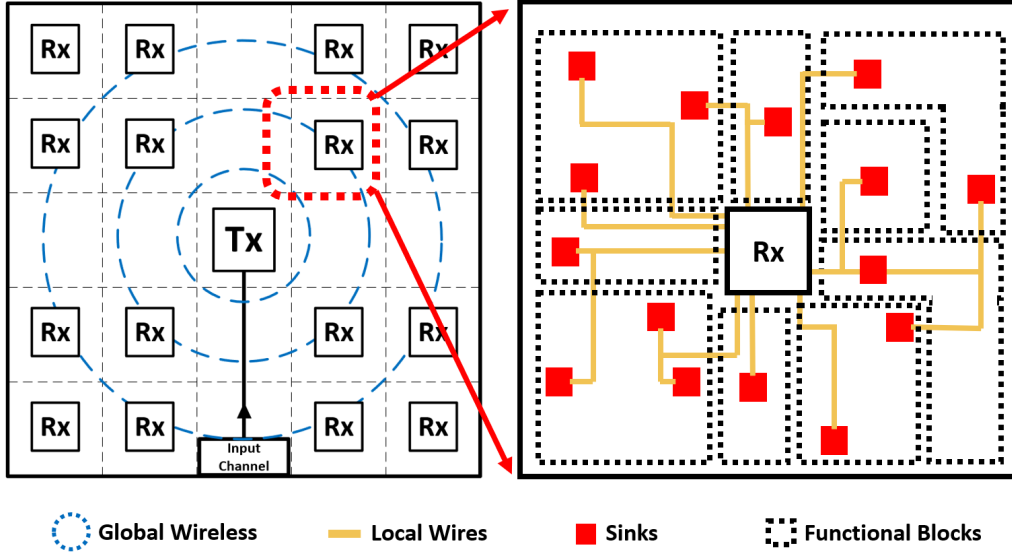


Figure 1.6: The proposed hybrid clock distribution network which combines both global wireless interconnect and local metallic wires.

be applied to fully-synchronous systems but also a GALS structure with mesochronous clocking. It provides a global clock signal with significant reduced delay and improved energy efficiency, which can benefit the synchronisation of a many-core system.

Under the assumption that the radio-frequency signal can propagate through the antenna following a omnidirectional radiation pattern, the latency of the wireless transmission is proportional to the distance between the receiver(Rx) and transmitter(Tx). Hence, wireless TRx pairs need to be allocated well such that the difference in communication distances can be minimised.

For general applications such as only part of the chip needs to be synchronised, or the synchronisation area has a non-nominal shape, the location of the Tx can then be adjusted accordingly to maintain the relatively balanced feature. Locations such as the centre of mass of a 2D shape or even place the clock Tx off-chip can be used, to satisfy certain DRC or design specifications. An off-chip wireless clock can benefit from the power reduction of the wireless clock Tx. Also, since the area overhead for the clock Tx is no longer to be a limitation, the design of the wireless Tx and antenna can be refined to have better SNR performance, therefore increase the frequency of the propagated baseband clock signal. One penalty of this implementation is that the distances between the Tx and Rx array are not identical/symmetrical, therefore the global clock skew will increase due to unbalanced structure. Other Tx allocations such as above die allocation (in a 3-D IC) exhibits different attributes. Although the Tx can be allocated in a different tier to provide more flexibility and to satisfy certain design specifications such as balanced latency, this implementation will bring challenges that the RF signal needs to penetrate the silicon substrate and the bounding bumps. Partial

RF energy will be wasted, absorbed or reflected by the metallic material, therefore mitigating overall wireless clock SNR [24].

This thesis will focus on the clock distribution in a nominal synchronisation area. For a simplified experimental model with a square topology as shown in Figure 1.6, in order to have a balanced structure for similar transmission latency, the Tx is located at the centre of the chip and hence the delay can then be quantified using the radius of the global wireless EM signals to the Tx.

As an essential attribute of a CDN, for a fully-synchronous system, ideally the clock skew should be reduced to the minimum level where possible, to mitigate the negative impact of clock uncertainties on timing closure. Based on certain CTS algorithms, it's possible to generate a zero-skew CDN. However, this comes with the cost of unnecessary resource waste such as excessive CDN wirelength, thereby occupying large chip area and generating large CDN capacitance. Taking advantage of the CAD support such as useful skew and concurrent clock and data optimisation (CCD) [25], [26], the constraint on clock skew has been relaxed in the recent decade. For high-performance VLSI design with high clock frequency, the global clock skew is set up to 10% of the overall system cycle time [27]. Hence, in this thesis, various implementations of the CDN are compared with a 10% skew upper bound, to observe the synchronicity of different solutions.

In conclusion, CDN using a novel interconnect instead of metal-based wires shows a promising prospect in terms of system performance improvement and power reduction. This research will investigate a hybrid wireless-wire architecture, which will adopt wireless interconnect during global clock distribution stage instead of wires.

### 1.3 Thesis Contributions and Structures

This thesis has proposed a low-delay and power efficient CDN design methodology, major contributions of which include:

1. A novel globally wireless CDN using wireless interconnect with efficient On-Off-Keying (OOK) modulation with the advantage of efficiency and performance was proposed and it provides significant power reduction of up to 32%.
2. A novel local CDN wired CDN design algorithm which provides low-local skew and wiring usage.
3. An improved circuit for the leakage-canceling modulator with 47.1% improvement in on-off isolation and a non-coherent noise-compensated demodulator were proposed.

4. Reductions of up to 23.9%, 28.8% and 11.4%, for power, delay and skew, respectively using mathematical modelling with MATLAB and Cadence co-simulations, were identified in our proposed novel industrial test case with Arm Mali-G77 GPU.

The basic structure of this thesis will show a throughout study of the proposed hybrid clock distribution schemes. Details of each chapter is given by:

Exploring the conventional and newly emerged interconnects for CDN design, reviewing the existing literature about clock networks using high-speed wireless and optical CDN will be given in Chapter 2.

Design of the proposed system architecture, mathematical model construction for local wired CDN will be shown in Chapter 3.

Design of the proposed system in circuit-level components for transmitter and receiver respectively, as well as the design of proposed two on-chip meandering antennas with structure and Electromagnetic simulation model setup, will be given in Chapter 4).

Experimental results and evaluation of the overall system performance of distributing clock signal will be shown in Chapter 5.

A comprehensive test case is then presented to give a full evaluation of the proposed hybrid CDN, and to test whether it can work on realistic scenarios. These materials will be given in Chapter 6.

Lastly, Chapter 7 will conclude this thesis towards the completion milestone and the potential future works to further explore and optimise the proposed design.

In addition, published and the on-going works would be included, for the purpose of academic evaluation. Walking towards the completion milestone, progress to-date includes:

1. "On High-speed Clock Distribution Network Using Hybrid Wire-Wireless Interconnects", Proceedings of Oxford Circuits and System Conference (OXCAS) 2017.
2. "Globally wireless locally wired (glowilow): A clock distribution network for many-core systems," in IEEE International Symposium on Circuits and Systems (IS-CAS), pp. 1-5, May 2018.
3. "Hybrid interconnect network for on-chip low-power clock distribution," IET Electronics Letters, vol. 55, no. 5, pp. 244-246, May 2019.
4. "An Active Silicon Interposer with Low-Power Hybrid Wireless-Wired Clock Distribution Network for Many-Core Systems", in IEEE Transactions on Very Large Scale Integration (VLSI) Systems, vol. 28, no. 9, pp. 2042-2054, September 2020.

5. “Energy-Efficient and Low-Skew 3D-KBP Algorithm for Clock Distribution Network Design in Many-Core 3-D Stacking Systems”, in IEEE Transactions on Computers, Special Issue on Communications for Many-Core Processors and Accelerators, October, 2020 (Major revision).



## Chapter 2

# Literature Review

This chapter will discuss and explore different types of clock distribution network (CDN) utilizing both conventional interconnect as well as newly emerged interconnects. These solutions will be compared through a qualitative analysis and the state-of-art interconnects will be used as the components of our proposed design.

The developing interconnect techniques have received great attention for the past few decades from researchers and engineers. For the consideration of high clock frequency and low clock uncertainty, various types of CDN using different topology have been developed. For conventional metal-based CDN, interference caused by crosstalk due to the scaling down technology, will consequently affect the quality of clock distribution, hence lead to unnecessary extra power and area dissipation. Thereby the mutual target for various newly emerged techniques is to reduce overall CDN power and area consumption caused by wiring issues [28], [29] and simultaneously remain rather low clock latency and uncertainties. These techniques are introduced as follows.

### 2.1 Reviewing of the Emerging Interconnect Techniques

First of all, various emerging interconnect techniques will be reviewed. The advantages and the drawbacks of these novel interconnects are then summarised to observe their potential of being used for on-chip global communication.

#### 2.1.1 Wireless Interconnect

Wireless interconnect is considered to be one of the most promising solutions to solve the above existing challenges for its developed background and CMOS compatibility. Under the help of matured RF technique applied to communication area as well as smaller CMOS technology nodes, chip designers could now design and implement a compact

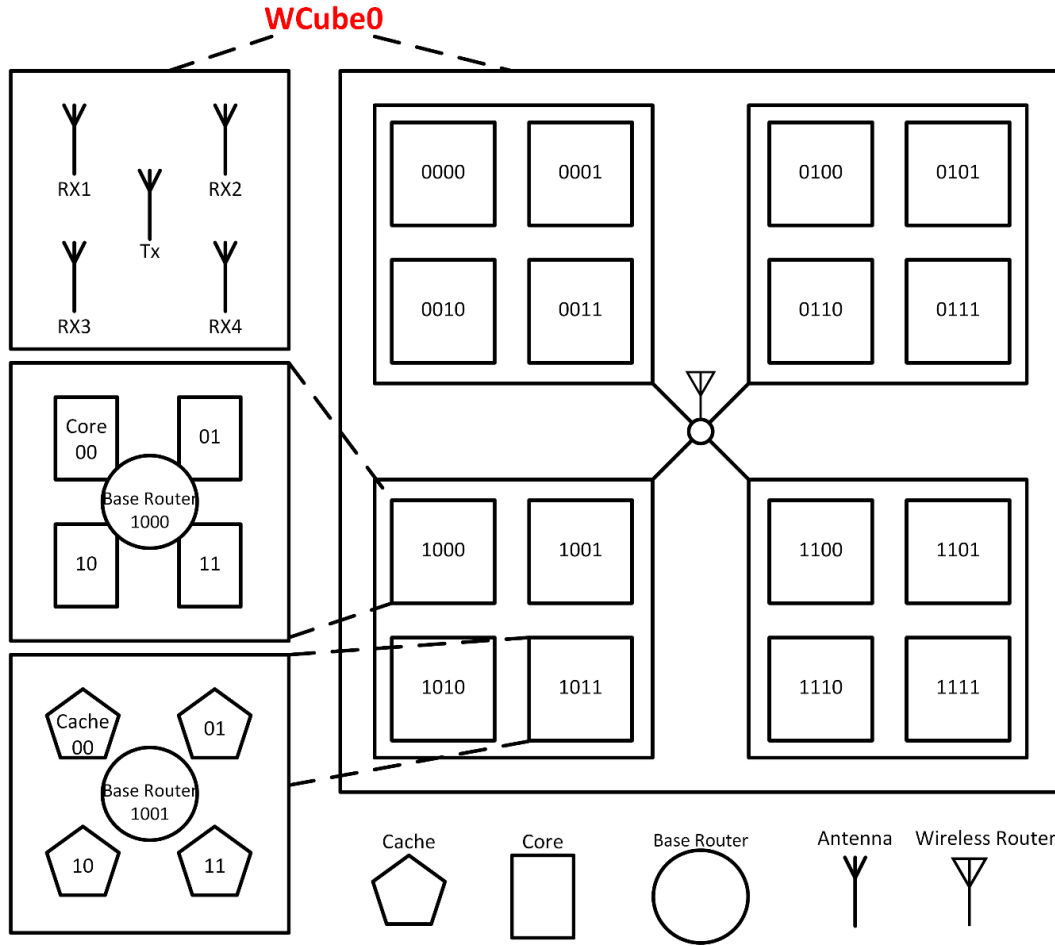


Figure 2.1: Wireless interconnect with WCube as application [2]. Wireless interconnect serves as the global communication channel.

on-chip wireless system with competitive performance when compared with conventional metallic interconnects. This cost-effective alternative aims to integrate on-chip transceivers with antennas so as to propagate electromagnetic radiation in between silicon substrate and bottom ground plane with a rather high speed. The surrounding nodes could therefore retrieve the decayed EM wave and then do the amplification/demodulation so that the originally transmitted data could be recovered accordingly. The inherent fan-out feature of the EM wave propagation significantly enhances the broadcast performance because the idealised propagation without considering any nearby interference could be regarded as a 360-degree circular pattern, hence generating less delay for higher data transmission rate. Receiving nodes could be arranged in a form of any geometry depends on the system performance requirements and load capacitance.

Ideally, wireless interconnect has the potential to take the place of conventional metallic interconnect in some particular areas, especially under the circumstances of large multicast demands or global signal transmission, however, the on-chip antenna integration currently becomes the bottleneck of this technique. Some on-chip antennas suffer from relatively low efficiency because of their dimension such that part of the power would not



be radiated out into substrate. This phenomenon will need extra energy to compensate the power loss in order to maintain adequate signal power level.

Another issue is the lossy silicon substrate. Radiated energy would essentially be dissipated for some amount because of the material property of the silicon [8]. Hence similarly, extra energy is required to mitigate this power decay, such as using power amplifiers at the transmitter front-end.

Apart from power loss, nearby electromagnetic interference (EMI) is another important issue that needed to be carefully considered. Analog signals are considered to be more "fragile" than common digital signals with only two logic levels. Since most of the wireless interconnects incorporate analogue techniques to local digital IPs, any signal distortion, for instance, spikes, would essentially cause severe error thereby generating erroneous data bits to downstream logic cells. The nearby flopping signal edge would affect analogue signals because of cross-talk existing between the analogue trace and digital components. Therefore, robustness is of paramount importance when considering adopting wireless interconnect for on-chip communications.

Many researchers have proposed the wireless interconnect utilised in network-on-chip design (NoC) for increasing data hop count [2, 30, 31] as a mature solution compared with conventional wire-based interconnects because of its compatibility with CMOS technology and broadcasting feature. The global transmitter is always placed at the centre of the system as shown in Figure 2.1 in order to achieve an average traversal time between any transmitter-receiver pair.

Besides, as the technology keep scaling down, smaller antenna size, as well as higher carrier frequency, could be realised, hence producing more bandwidth resources for higher data rates. Nevertheless, as mentioned above, the on-chip antenna needs to be carefully designed for the maximum power transfer. Trade-offs between antenna area and gains need to be made to balance the performance and transmission efficiency. Moreover, since wireless interconnect might contain power-hungry on-chip oscillating components to generate a high-frequency carrier wave, extra power consumption is another point that needs to be carefully considered about for chip designers.

### 2.1.2 Other Radio-frequency Based Interconnects

Using other media such as transmission line [32], [33] as well as a waveguided plane to propagate EM wave is similar to wireless interconnect. Among the radio frequency interconnect (RF-I), surface-wave interconnect (SW-I) is the newly emerged solution to multi-gigabit communication.

SW-I shows its superiority in multicast architecture similar to wireless interconnects [28], [34]. Since the EM wave is captured and propagated inside a waveguided plane,

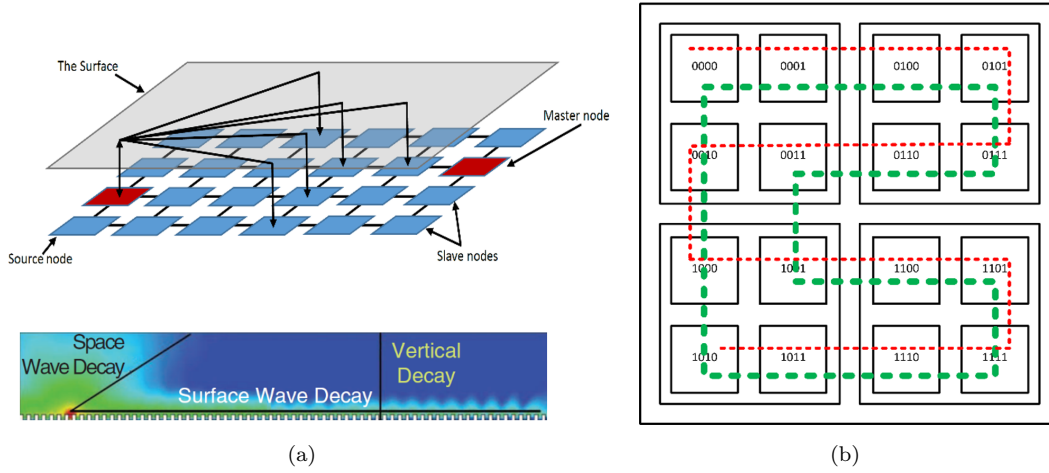


Figure 2.2: A typical structure of (a) surface-wave interconnect (SWI) and (b) transmission line interconnect (Tx-line) for global communication.

signal decay has been reduced remarkably, thus providing a desirable performance as well as better power efficiency, as shown in Figure 2.2(a). Nevertheless, extra waveguide plane and transducers need to be implemented, which could lead to an area and stack-up complexity overhead. Besides, a certain amount of 3D integration techniques for the link between RF transceivers and the surface or transducer is necessary, thus leading to thermal issues. Moreover, noise and interference need to be carefully considered during the design at integration stage of transceivers, transducers and the surface, as the surface-wave has the possibility to pick up unwanted noise from nearby circuits or ICs, which could consequently affect overall system performance. What's more, antenna multipath propagation is another critical problem. As EM wave would essentially travel through many media not only the designated surface, receivers may pickup multiple duplication of the transmitted signal, which increases the possibility of generating erroneous data bits.

Based on radio-frequency techniques, transmission line (Tx-line) is an alternative to wireless or pseudo-wireless interconnects such as SW-I. It adopts specified waveguide for global signal propagation as shown in Figure 2.2(b). Signal traces are allocated through every core in order to form a multicast architecture. Different from conventional wires, the TX-lines are specially designed for high-frequency signal transmission with adequate energy loss and cross-talk rejection. Hence, similar system structure or components in wireless interconnect could be adopted. The drawbacks of Tx-line are that extra area is occupied by the designated transmission lines [33], therefore it might not be an ideal solution to designs with tight area constraints.

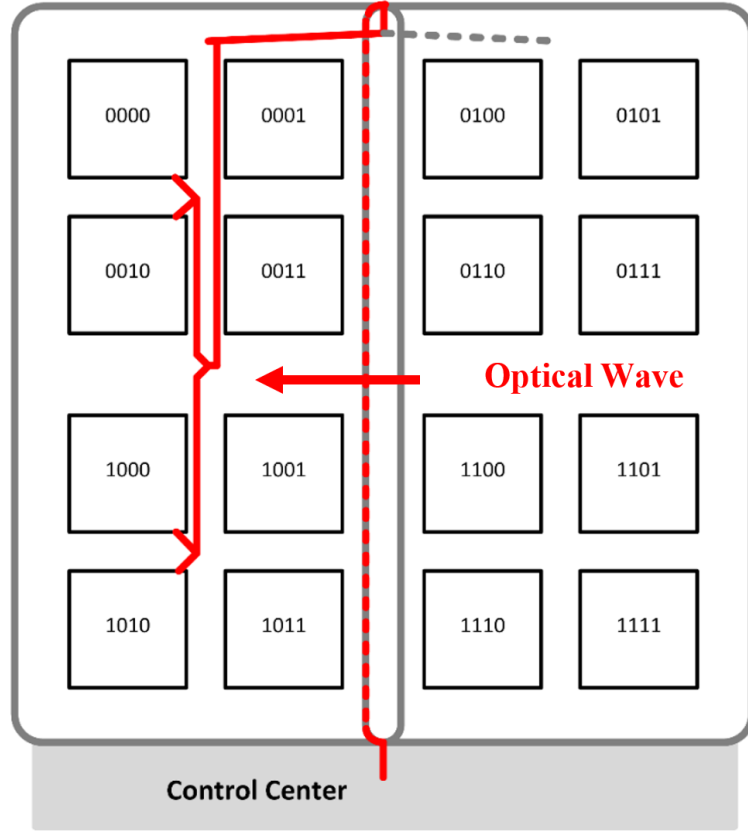


Figure 2.3: An global H-tree constructed using optical interconnect to distribute clock signals.

### 2.1.3 Optical Interconnect

Another novel approach is to use on-chip optical devices to build an optical interconnect network. Similar to transmission line interconnect, illustrated in Figure 2.3, optical interconnect could form a ring or tree-based architecture [35], [36]. As the signals are propagating at the speed of light inside the optical waveguide, the total transmission phase delay could be reduced remarkably. Moreover, the high bandwidth gives optical interconnect the potential to communicate even over terabit scale [37], [38]. Besides, optical interconnect has shown a natural immunity to nearby circuits interference caused by cross-talk, EMI, etc. This feature could be utilised to implement a noise-sensitive signal propagation route.

However, to implement the multicast architecture, extra components are necessary such as on-chip micro-lenses or laser source and so forth. These costly components are sometimes non-CMOS compatible and therefore produce extra area and power burden to overall design budget when mounting those components onto a chip. Besides, power-hungry components like on-chip laser source would burn a large amount of power to avoid signal attenuation inside the optical waveguide material, which could eventually become an intractable trouble to low power designs. In addition, the optical signal is

required to transform back to electronic signals before sending to the downstream logics, hence a potential energy loss might occur, and transmission efficiency would become a critical issue for on-chip signal transducers. What's more, energy loss might also appear at any forking point because of the reflection and refraction characteristics of light signal. Hence, optical interconnects would be challenging for clock distribution in near decades, despite its novel high electronic noise immunity and high-bandwidth features.

## 2.2 Conventional CDN with Broadcast Architecture

Given a fixed load allocation after cell placement procedures during IC physical design, for global clock distribution, clock signals need to travel from the clock generator (either on-chip or off-chip PLL) to local clock endpoints. Several types of pins/ports can be treated as clock endpoints, such as the clock pins of the sequential cells or macro cells, the root node of other CDNs, output ports or even non-clock input pins of sequential cells, among which consists two different types of clock endpoints, namely sink pins and ignored pins.

The detailed definition of the pins during CTS procedures are given as follows. A sink pin can be treated as a capacitive load and will be involved in the delay balancing and design rule constraint (DRC) violations fixing stage during clock tree synthesis (CTS), such as the clock input pin of an arbitrary sequential element (latch or flipflop). By contrast, an ignore pin will only contribute to DRC violations fixing such as non-clock input pin of a cell, or a clock-driven control signal pin in a tri-state gate.

A modern synchronous system might contain more than one clock signals, therefore increasing the complexity of designing a CDN with interclock skew constraints. The clocks (either master or generated) might have different frequencies (i.e. multiple times of the master clock period), but the phase relationships between each other might be fixed, hence these active areas can still be regarded as synchronous. If two or more clock signals are having both different frequencies and unpredictable phase information, the active areas with respective clocks are regarded as asynchronous areas, and therefore need asynchronous bridges (i.e. synchroniser or FIFO) to propagate data between two clock domains. To reduce the complexity, and simultaneously preserve the clock timing properties, this work is based on a synchronous system with a master clock solely. The generated clocks are regarded as floating pins [39] and the clock balancing calculation will trace through the intermediate sequential clock generation cells, until it reaches the ultimate clock sinks according to Figure 2.4.

The wiring, buffering and geometric methodologies all have a significant impact on the overall system performance so how to distribute clock signal with good signal quality remains to be a difficult task. For a common design rule, global signals such as the clock, power and ground nets ( $V_{dd}$ ), ground ( $V_{ss}$ ), etc. would adopt the upper thick metal

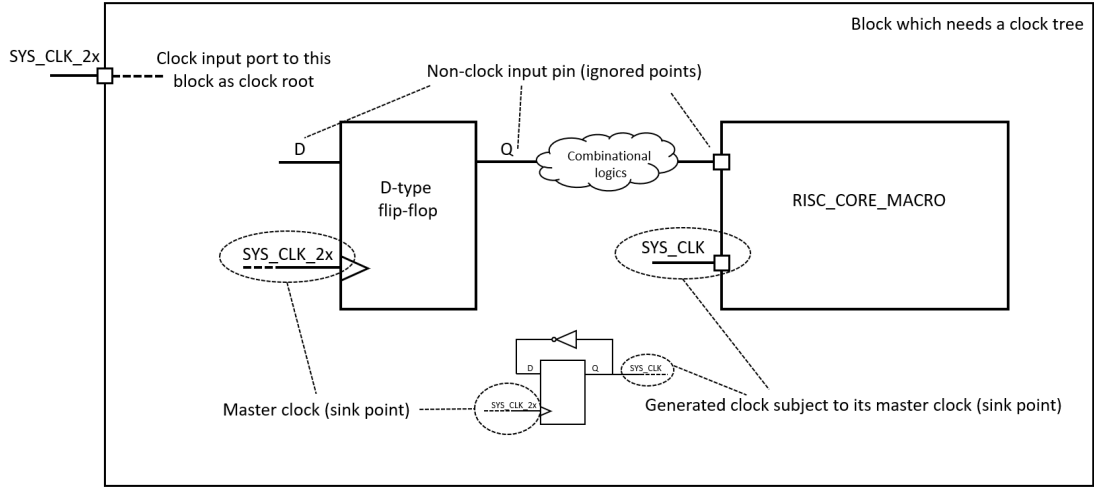


Figure 2.4: Identifying the start and the endpoints of a typical clock tree in a design.

layers or use non-default routing rules for global interconnect structures, to reduce wire resistance and IR drop across the global networks [40].

### 2.2.1 Combination of Binary-tree, H-tree, X-tree and Clock Grid

A tree-based structure is by far one of the most common methodologies to distribute clock inside a VLSI system. As a tree-structured interconnect network, the clock signal is entering the network via a root node and traverse until it reaches every single clock endpoint. The different clock tree branches are implemented, and clock signal could therefore be propagated through any branch buffer and passed to the next layer. At the end terminal of a clock tree, all the leaf nodes could consequently provide a clock to sequential cells to be synchronised. Figure 2.5 illustrates the basic structure of a typical clock tree. The ultimate goal of the clock tree synthesis (CTS) is to implement a clock tree with minimum clock skew and clock insertion delay, whilst preserving acceptable clock power and global congestion for the routing stage later. As a complex network, the property of the clock tree will have a substantial impact on the overall performance of the clock distribution networks, such as the topology, depth and balance of the tree.

A single buffer structure provides a single but powerful clock buffer/driver at the source of the CDN and there are no further buffers at the subsequent route of the network, as per Figure 2.5. Since the clock source directly drives the entire loads, the size of this buffer is designed to be large enough to provide driving strength so that the criteria such as overall clock slew rate, rise and fall time are acceptable.

On the contrary, a distributed buffer tree would have smaller clock buffers place on every clock route to enhance signal quality. The buffer tree would then converge on a bigger driver at the root node. This strategy provides more flexibility to designers to solve any timing violations by inserting more clock buffers onto the clock route with less delay

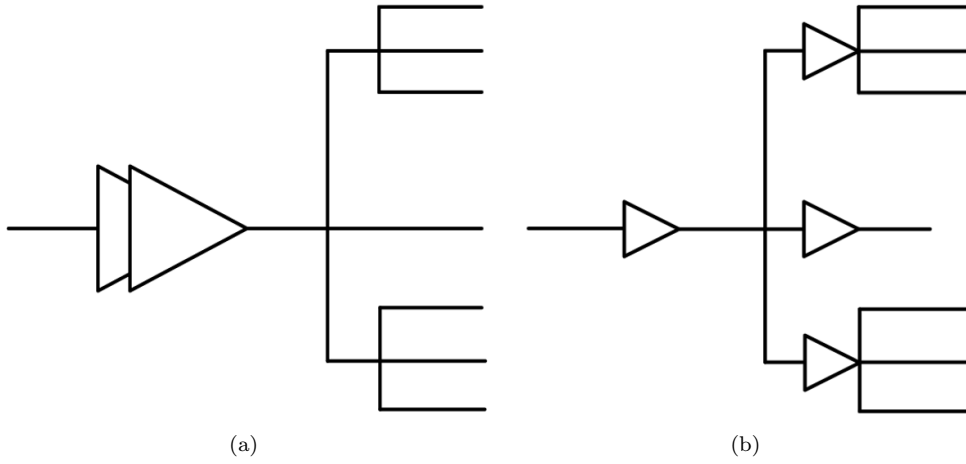


Figure 2.5: A typical structure of (a) single and big buffer/driver insertion strategy in a CDN and (b) the distributed and small buffer/driver insertion strategy in a CDN.

so that the overall time of arrival of clock signal is equalised and clock skew could be minimised. However, due to process variation during fabrication stage, any mismatch between different clock buffers would generate different propagation delay and essentially cause large clock skew.

A practical example of the distributed buffer tree is the DEC Alpha 21064 microprocessor [41] shown in Figure 2.6. It adopts a five-stage buffer tree with a minor modification at the intermediate layer. Orthogonal metals are connected to the original top metal interconnect to form the geometry of a clock grid, hence reducing the effective resistance seen by the previous clock buffers, which consequently mitigates overall clock propagation delay and draws the different time of arrival of the clock signal. As a counterpart and successor, DEC Alpha 21164 microprocessor in  $0.35\ \mu\text{m}$  process [42] is an example of the single buffer structure. It adopts a massive clock buffer to drive the entire loads around  $3.8\ \text{nF}$  over a dimension of around  $1.6\text{cm} \times 1.8\text{cm}$  physical die area. Due to the absence of active devices on any clock routes, any transistor mismatch or device noise could be neglected, therefore this design exhibits a rather low clock skew of around  $80\ \text{ps}$ . However, in order to drive the entire loads with only one clock driver, the driver needs to be large enough so that the entire CDN burns nearly 65% of the chip power around 50 Watts in the CDN. Other examples of CDN using binary-based global clock tree can be found in Intel Pentium series [43], which can provide robust clock signal up to  $3.8\ \text{GHz}$ .

The above CDNs adopt a combination of binary tree and clock mesh geometry for global clock distribution for its simplicity and energy efficiency. However, with system speed growing over multi-gigahertz, these wiring methodologies no longer seem to be capable of providing high speed clock signal with low clock power consumption. The H-tree/X-tree structure provides a symmetric alternative for global CDN design. Essentially, an H-tree or X-tree is a re-scheduled binary tree with a fully-symmetrical geometry with

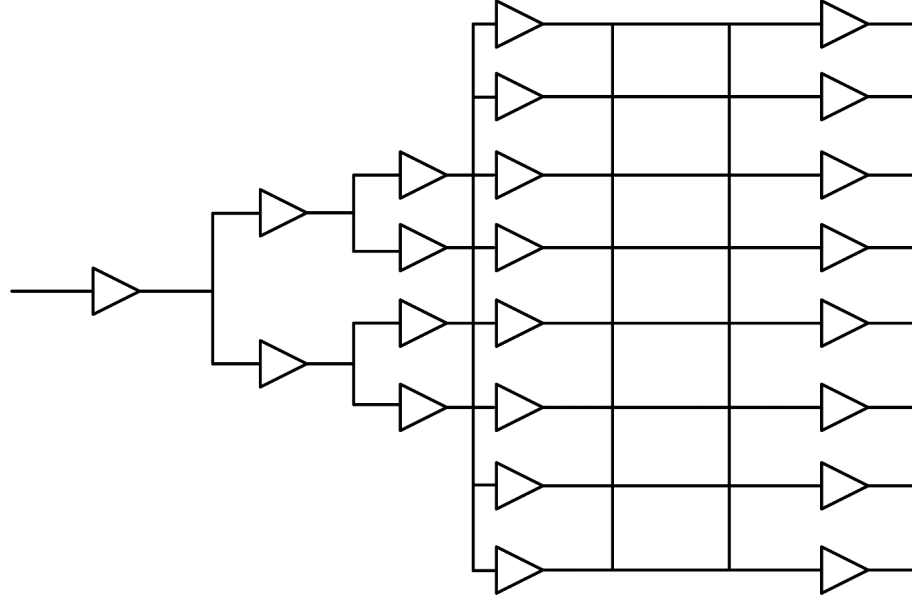


Figure 2.6: Buffered-tree with local clock grid for DEC Alpha 21064 clock distribution structure.

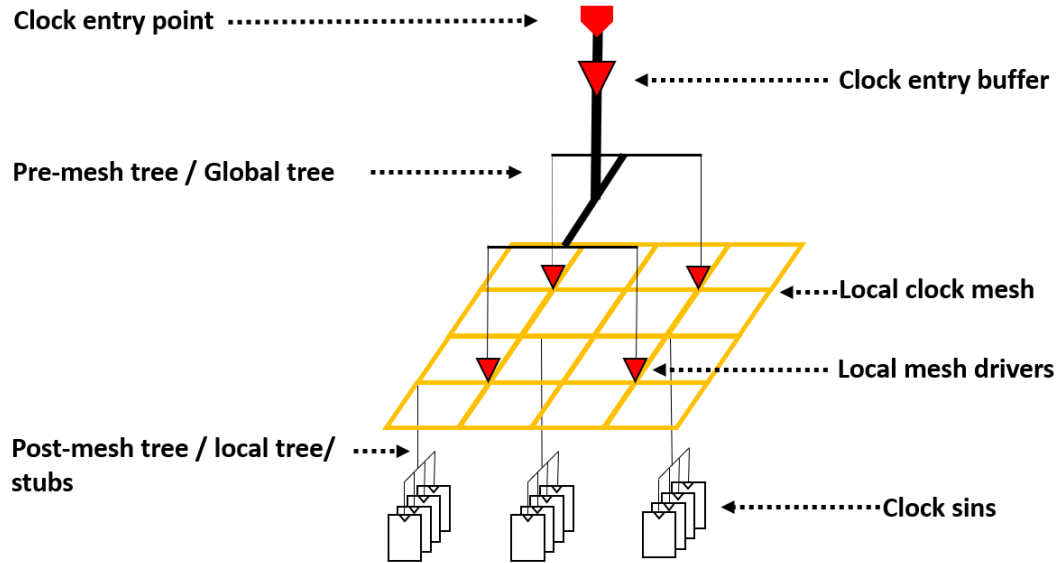


Figure 2.7: A typical structure of global tree and local mesh (TLM) clock distribution scheme.

clock signal distributed to adjacent sequential logics at the same time. As the total interconnect length of each branch is designed to be identical, this kind of CDN ideally has zero clock skew. To avoid energy reflection at any forking point, a tapered-branch is designed to have the successive conductor with half of the predecessors interconnect width. Therefore, it can generate twice of the original impedance seen by the forking point. As two identical branches are split at this joint position, effective parallel impedance of the downstream CDN seen by the forking point would stay the same to upper level CDN, hence reducing unwanted energy loss by a pseudo-impedance matching procedure [44],

[45].

However, fully balanced H-tree networks sometimes cost extra power and wiring budget on different applications. For some specific design with low chip density or unevenly distributed logics cells, H-tree structure always seems to be an overkill that the benefits of very low/zero clock skew could eventually be achieved by other clock routing architecture with less load capacitance, and hence, lower power consumption [45], [46]. Besides, a fully-symmetrical local load capacitance and displacement are always impractical, hence leading to a constant local clock skew caused by load-imbalance if local logics still strictly adopt H-tree as local CDN topology. To address this problem, a mesh-based topology can effectively reduce the significant difference in the length of different clock paths [47]. Some state-of-art and commercialised examples of using clock mesh as the global CDN include the AMD Zen microprocessor series in 14 nm process and the AMD Zen 2 microprocessor series [48] in 7 nm process, which can support the clock frequency up to around 4.7 GHz. Other variants such as Intel Core series [49] use a combination of global clock spine and mesh named Nehalem global clock distribution network, which efficiently attenuate the accumulation of clock skew along the clock routes.

Clock sinks are connected to the closest mesh segment, and vertical or horizontal metal lines is attached to the existing CDN, to further short different wires together as shown in Figure 2.7. [50] raised a local mesh topology to improve its local skew tolerance for various local loads as illustrated in Figure 2.6. [51] has also raised a similar opinion that a tolerable skew methodology is implemented with a relaxation on overall clock tree power consumption compared with H-tree/X-tree only.

Related research such as [52] and [53] provide an effective way of modelling local clock mesh, which points out that compared with a conventional full fan-out clock tree, a mesh structure can effectively reduce timing uncertainties. Besides, clock skew can be reduced by inserting ‘shorting’ segments, as suggested in [50], at the cost of increased power consumption. In fact, this power overhead becomes one of the major limitations of the clock mesh being adopted in the power-sensitive designs, especially for modern IoT devices.

### 2.2.2 Challenges of Existing Metallic Wire-based CDN

As mentioned earlier, global interconnect delay is increasing significantly because of the shrinking size of the conductor as well as the higher chip density and capacitive loads. As an essential signal inside a sequential system, a clock has several characteristics:

1. It has the highest fan-out number hence load capacitance.
2. It has the highest propagation distance across the chip area.



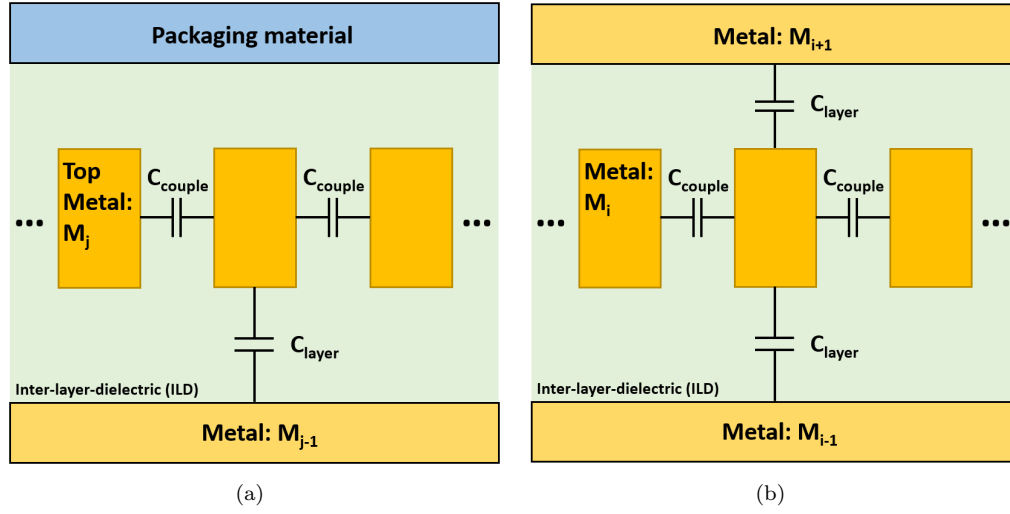


Figure 2.8: A typical structure of interconnect model from PTM with (a) global and (b) local wiring structure [3].

	Global	Global	Local	Local
Technology (nm)	R ( $\Omega/\text{mm}$ )	C (fF/mm)	R ( $\Omega/\text{mm}$ )	C (fF/mm)
180	21.9	232.5	174.6	203.8
130	30.6	258.1	244.4	228.7
90	36.7	259.7	488.9	186.1
65	40.7	259.6	1099.9	166.1
10	286.5	212.3	4583.3	201.4
7	733.3	213.5	4583.3	201.4

Table 2.1: Interconnect parameters for global and local wires from 180 nm process to 7 nm process in terms of both resistance and capacitance per unit length.

3. It has the highest speed among all on-chip digital signals.

These properties obviously have already caused a conflict between system performance requirements and the higher interconnect delay, hence to deliver high-speed clock signal across increasing chip area becomes a critical task. For a typical model and its data-set from Predictive Technology Model (PTM) [3] shown in Figure 2.8 and Table 2.1, CMOS process decreasing from 180 nm to the state-of-art 7 nm causes the unit length resistance of global traces almost double their original value.

From this perspective, interconnect delay would be affected by the increasing RC value as the CMOS process keeps shrinking. And this phenomenon would consequently affect clock distribution significantly as the interconnect delay might take several clock cycles from the clock source to different loads. Besides, the increasing RC characteristic would essentially cause a larger difference between the arrival time of the clock signal, consequently generating large clock skew. Cross-talk capacitance could no longer be neglected since the space between wires is getting smaller, and hence global signals like

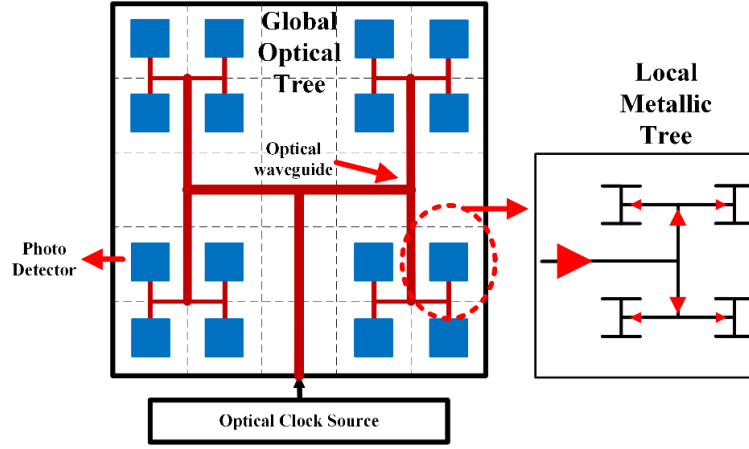


Figure 2.9: CDN architecture for global optical guided interconnect and local metallic H-tree.

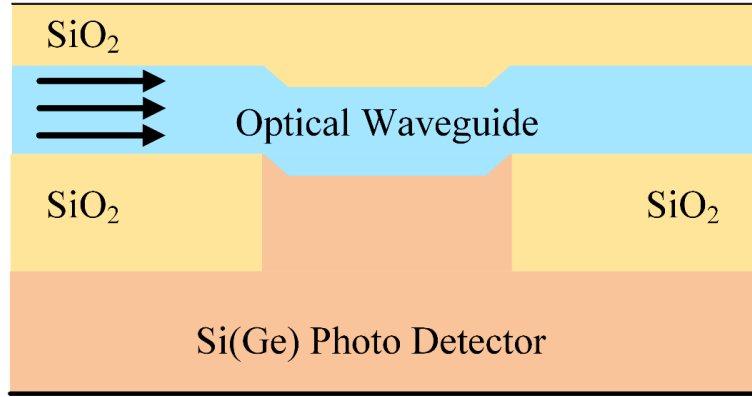


Figure 2.10: Stack-up schematic of the optical CDN proposed in [4].

clock are more interference-sensitive. Therefore, higher chances of glitches might occur due to clock propagation malfunction.

## 2.3 Optical CDN with Broadcast Architecture

Optical interconnect has the natural advantage of high bandwidth and data transmission speed. Taking advantage of the high-speed feature, using optical interconnect for global clock transfer could potentially bring huge merit to system performance.

### 2.3.1 Guided and Free-space Optical CDN

In General, two types of optical interconnect could be utilised as a part of the CDN. Firstly, a guided optical interconnect consists of optical fibres or integrated waveguides on the substrate could be regarded as the substitution of metal wires [37], [54]. As the dielectric waveguide could have very low loss compared to metals, distortion of the

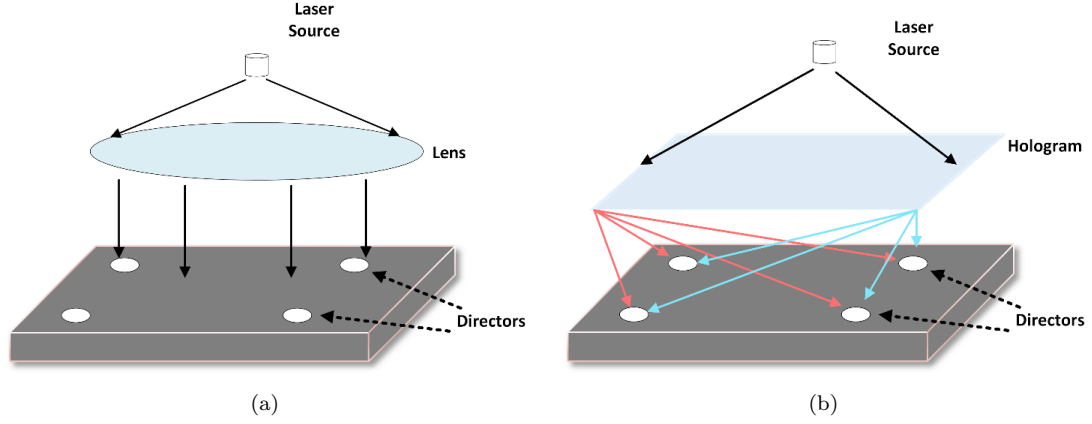


Figure 2.11: Free-space optical interconnect with (a) unfocused and (b) focused clock broadcast architecture [5].

electric signals could be restrained easily. In addition, overall propagation delay would also decline rapidly as the resistive load has been reduced and hence produce much less loss when comparing to conventional metal-based interconnects. EMAT Research Group has proposed a similar architecture which adopts an optical waveguide to construct a fully-balanced H-tree for global CDN, and uses conventional metallic interconnect to construct local H-tree as regional/local CDN [4] shown in Figure 2.9 and 2.10. This design has been successfully fabricated with an integrated on-chip splitter and an optical receiver (photon detector/ITA), which is capable of transmitting global clock signal to all 64 fan-out with a frequency of 10 GHz while keeping the heat generation at an acceptable level.

An alternative is that light signals are travelling in free-space through integrated lens or holograms and then captured by detectors located on a die [55], [56], which is shown in Figure 2.11. The literature has claimed that using on-chip optical devices, total performance would be massively promoted, as optical interconnect could produce large fan-in and fan-out feature by reducing the effect of capacitive loading. Besides, optical signals have a natural advantage of their extremely high signal frequency. As a carrier wave inside an interconnect could have its wavelength down to a few microns or even hundreds of nanometers in free space, optical interconnects begin to show their immunity against interference like signal reflection and cross-talk caused by frequency-depend reasons.

Both guided and unguided optical interconnect yield their superiority to the conventional metal-based CDN. As the effective transceiver data rate could even reach terabits per second, theoretical transmitted clock frequency has the potential to reach tens of GHz or even hundreds of GHz, which is highly sufficient for high-performance systems in the future. What's more, optical signals like a short pulse is capable of propagating clock signal directly as an effective timing edge [37] with outstanding precision, which could produce less skew and jitter subsequently.

### 2.3.2 Challenges of Existing Optical CDN

Despite of the high-performance feature of optical interconnect for CDN, the integration of power-hungry and sophisticated optical devices like optical signal sources and transducers, is still an issue needed to be carefully considered. Arguments between the drawbacks and the advantages of optical interconnect in terms of power efficiency, complexity and fan-out capability have been discussed in a variety of literatures. [54] and [36] have claimed that power-efficient optical interconnect could make full use of the reduction of resistive and capacitive loss thus providing higher energy efficiency with less J/bit. Yet the opposite opinion shows that regardless of the reduction of RC loss and delay, optical devices still face the challenge in power budget requirements unless interconnect length is relatively long [28], [37]. For instance, a certain amount of energy would be dissipated as signal from the source would need to be transferred into light first by the integrated transducer and then transferred back to electric signal at receiver end even if the communication distance is relatively short. There exists an amount of energy loss during transducer stage, which could eventually outweigh the benefits of optical CDN.

Besides, to implement the multicast architecture, forking nodes like splitters and combiners [28], [36] are necessary along the route. Incident light would face severe signal decay at the forking points, thus limiting fan-out capability and bringing new challenges to power management. Also, using a large percentage of non-CMOS devices would consequently bring more design complexity. For example, it would be relatively difficult to package an off-chip laser source for manufactures [57], [58]. On the other hand, silicon photonic devices are still under developing and being immature to eliminate CMOS compatibility issues at the moment [59]. Related opinions have been raised in [55] and [60] as well that extra power burden and complexity need to be balanced with technical merits.

Above all, optical interconnect for CDN still seems to be immature for at least the recent decade despite the promising prospect it shows. Currently, optical waveguide doesn't always seem to be a cost-effective solution compared with inexpensive metal wires.

## 2.4 Wireless CDN with Broadcast Architecture

When considering high-performance IC design, challenges have been raised for conventional CDN for years, as clock distribution is always a crucial aspect because of its significant impact on overall performance. As a single clock generator would drive the clock tree through several hierarchies of clock buffers to ensure drive strength, eventually, a large amount of phase delay would be generated and hence produce uncertainties like skew and jitter due to load unbalances.

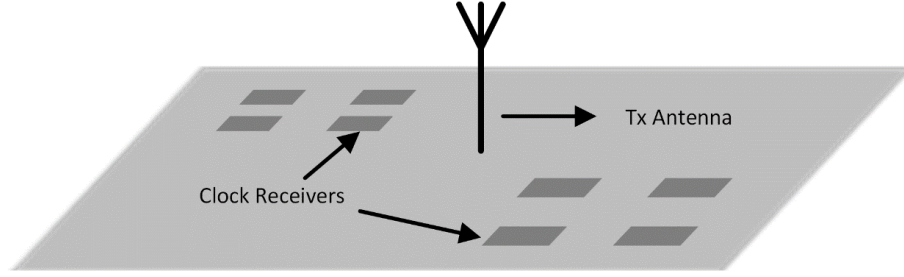


Figure 2.12: General architecture of a wireless CDN with integrated wireless clock transmitter and receivers.

Furthermore, for many-core systems, it would be excessively costly to use a conventional monolithic CDN to drive the entire die area because of the high latency and RC power caused by metal wires. Thereby, using wireless interconnect as CDN has shown its superiority in terms of power efficiency, outstanding fan-out feature and simplicity [61]. A typical structure of wireless CDN is shown in Figure 2.12. As electromagnetic (EM) wave propagated from on-chip antenna benefits from its omnidirectional propagation pattern, the global wireless clock Rx can receive the broadcasting clock signal with considerably small clock skew which is proportional to the communication distance between TRx pairs. Theoretically, the EM wave can travel near the speed of light, and therefore a global wireless CDN can produce extraordinary signal integrity and stability.

By freeing up the wiring resources, overall wire latency has been reduced almost to zero, therefore mitigating clock frequency dispersion and subsequently improving the clock frequency to be transmitted. Furthermore, the existing low power techniques such as power gating, clock gating or multiple supply voltage can be integrated with the proposed global wireless CDN and hence providing more flexibility. By allocating clock Rx to different synchronisation area, local clock receivers can independently manipulate the clock recovery procedures according to the clock domain enable/unable demand, hence forming the function of block-level clock gating to reduce more dynamic power when needed. For example, similar to conventional wired CDN, local logics can adopt circuits such as frequency divider or doubler to customise their own frequency of interest in a system with multiple clock domains without additional power-hungry clock generator like a phase-locked-loop (PLL) or delay-locked-loop (DLL) for clock generations with different frequency. Therefore, wireless interconnect seems to have great potential for clock distribution in future many-core systems, in terms of the natural broadcasting features and the support of the existing low-power techniques which can boost energy efficiency.

#### 2.4.1 Wireless CDN with VCO and Frequency Divider

Wireless interconnect has provided an alternative to skip conventional wires to propagate signals to destinations such as local registers. As technology keeps scaling down, higher

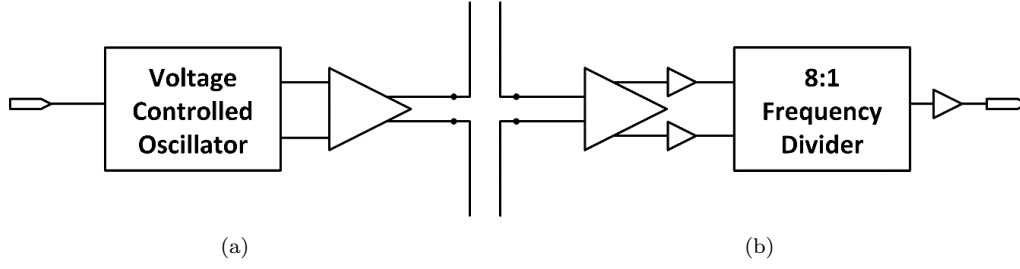


Figure 2.13: Block diagrams of wireless CDN [6] with (a) clock transmitter and (b) clock receiver.

carrier frequency with smaller integrated antenna size could be implemented, thus creating the potential of higher local clock frequency. As clock signals would generally jump over global clock branches and reach the corresponding synchronisation area within one hop [28], it becomes more power-saving and area-efficient because of the reduction of clock buffers and RC loss. Similar to board level radio frequency design, typical wireless interconnect implementation involves the design of a local oscillator, a frequency mixer, amplifiers and integrated antenna, which forms a complete RF transceiver [6], [60] shown in Figure 2.13 and 2.14.

Kenneth et al. has raised a solution to wireless clock distribution [7]. Using on-chip voltage control oscillator and antennas, a 20 GHz generated sinusoidal wave would be transmitted. At receiver front-end, an antenna would capture the radiated signal and amplify it through a low noise amplifier (LNA). Then the recovered signal would pass through a frequency divider to produce the required local clock signal given by:

$$f_{clk} = \frac{f_{carrier}}{\alpha} \quad (2.1)$$

where  $\alpha$  is the division ratio of the frequency-divider. Recovered clock signal would be buffered and then transmitted to downstream circuits accordingly. Compared with conventional metal-based CDN, this approach shows its superiority in simplicity without metal branches.

With the two integrated antennas placed 5.6 mm apart, clock signal generated at the receiver end could reach around 1.9 GHz, RC delay has been significantly reduced. Total power could be saved as well because of the reduction of wire resistance and load capacitance. Therefore, only a small amount of energy would be consumed by global clock transmission here.

#### 2.4.2 Wireless CDN with Digital Modulation Techniques

Related literatures such as [8], [2] and [62] have produced other kinds of wireless interconnect implementation shown in Figure 2.15. Although they were not directly to be

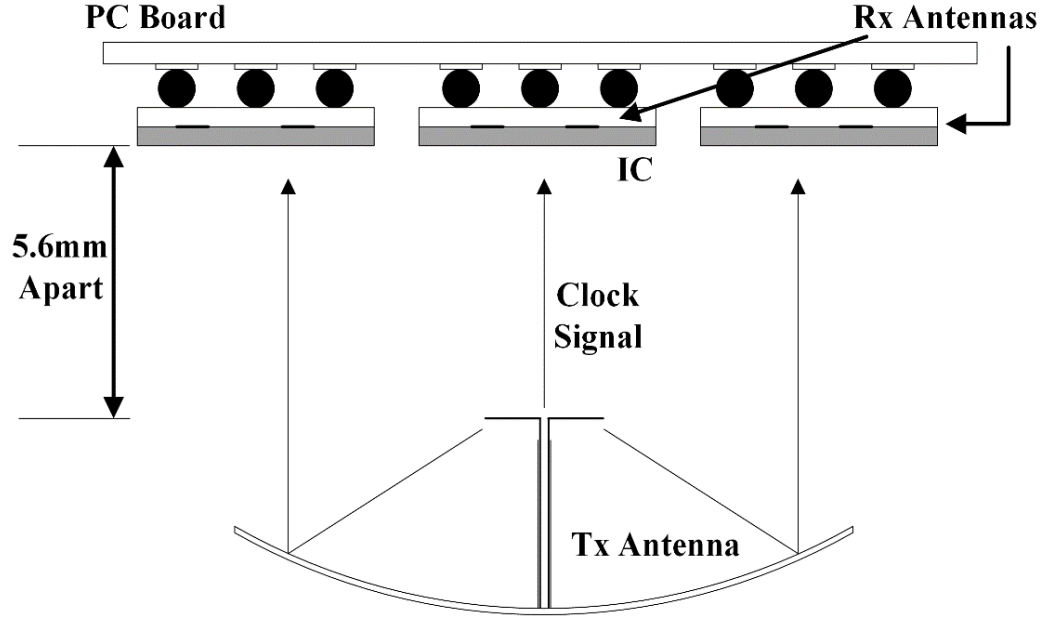


Figure 2.14: Wireless CDN transceiver architecture in [7] with 5.6mm propagation distance.

used as part of CDN, their broadcast architecture and fan-out capability also seem to be a good choice to form a CDN. These studies propose several throughout analysis of transmitting partial swing data bits on the on-chip/chip-to-chip basis. Hence, they are not directly suitable for global wireless clock distribution.

These techniques incorporate digital modulation [63] with high signal frequency band, therefore providing terrific communication quality when tested with pseudo-random bit series (PRBS) for near-field data transmission. Typical implementations using quadrature amplitude modulation (QAM), quadrature phase shift keying (QPSK) and amplitude shift keying (ASK) [64] to accommodate data stream onto a high-frequency carrier wave so as to reduce the disturbance of inter symbol interference (ISI). Experimental results indicate that within a communication distance less than one-meter, corresponding wireless interconnect/transceiver could reach an error-free (i.e. Bit Error Rate  $\leq 10^{-12}$  with  $2^{31} - 1$  PRBS) data rate over 40 Gbps [8], which provides a promising potential to multi-gigabits communication. With a stable communication quality using similar transceivers, clock frequency would be able to reach 10 GHz at a relatively low cost.

However, it is noticeable that the above wireless interconnects based on modern digital modulation techniques are not designated for clock distribution, such that the communication distance, system circuits complexity, etc. are quite different from the CDN performance or bandwidth requirements. Hence, to tailor wireless interconnect with modulation techniques, certain simplifications are necessary to meet the need for clock distribution. Specifically, power-hungry devices, such as power amplifiers and voltage-controlled oscillators which are commonly adopted inside modulation techniques and

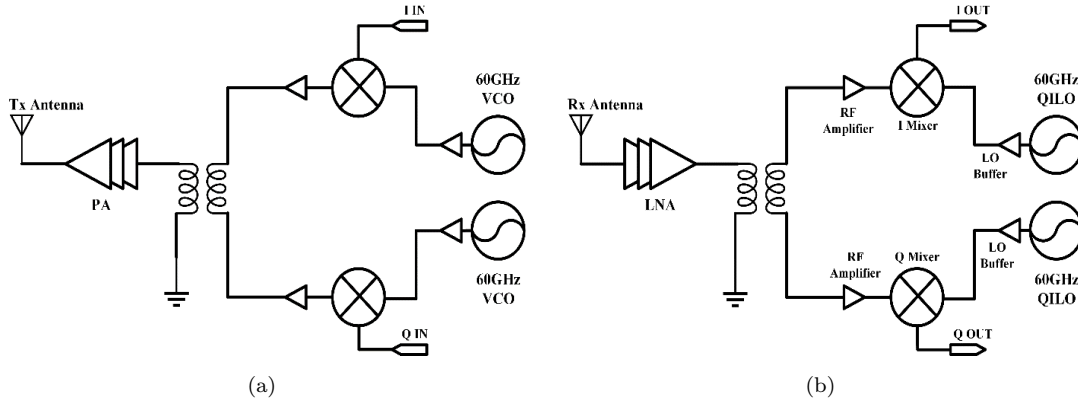


Figure 2.15: Simplified 64-QAM wireless interconnect architecture [8] with (a) transmitter and (b) receiver.

need to be considered whether they are necessary for wireless CDN, for the purposed of power efficiency.

### 2.4.3 Challenges of Existing Wireless CDN

Regardless of the advantage of wireless CDN, this emerging technique is still facing some challenges. As frequency divider based wireless CDN like [6] and [7] may directly transfer (down-convert) received signal without certain conditioning techniques, the recovered local clock signal might be highly dependent on the output signal frequency generated by VCO inside the transmitter. As VCO output frequency would eventually be affected by thermal and process variation, clock uncertainty would occur consequently if VCO output is unstable. Furthermore, to satisfy the output stability, a trade-off between performance and power efficiency needs to be made. The output of VCO could be phased locked to a phase-lock-loop (PLL), hence leading to extra circuit components which consume more power. Moreover, unless an adjustable frequency divider is adopted, the frequency of the clock to be transmitted would be limited to a fixed constant in design like [6], thus producing less flexibility.

On the other hand, wireless CDN based on modulators/demodulators in [30], [6] would be able to tackle this drawback, nevertheless, most researchers found it challenging to design a compact integrated antenna with high transmission gain ( $S_{21}$ ), wide frequency range and modest power consumption. The solution includes adding antennas with different operation frequency range to form an antenna cluster or array so as to cover wider bandwidth [65], while others adopt higher operation frequency even up to sub-millimetre range. However, these solutions either request a larger chip area or consume more energy on power-hungry local oscillators and passive components, which might consequently bring other drawbacks to the area and power budget during design stage [30], [66].



Besides, wireless CDN has the possibility of being affected by the nearby circuit as a microprocessor could be very noisy. Interference could be injected into the clock transmitter hence reducing the communication quality. Improving signal-to-noise-ratio (SNR) during clock transmission becomes a critical issue, which would request high-gain power amplifiers and low noise amplifiers thus cost more energy. In addition, transmission efficiency needs to be carefully considered to reduce antenna return loss caused by an impedance mismatch or chip packaging issues [66], thereby increases design complexity.

To sum up, the major challenges of global wireless clock distribution include: the recovered base-band clock frequency and uncertainties, the area and power overhead for wireless clock TRx pairs and the robustness and reliability of the wireless transmission. These points will be covered and discussed in Chapter 5 and Chapter 7.

The basic idea of other state-of-art clock distribution network (CDN) is to reduce the physical length or the resistance and capacitance (RC) product of the clock path, since clock skew is proportional to absolute interconnect latency [67]. Hence, the clock routes are designed to be off-die. A typical solution includes using the combination of a global tree and local mesh (or vice versa), so that the mesh segments can short different vertical spines together to compensate different arrival times of a clock signal [49].

A global mesh-based structure trades better clock distribution qualities with higher power consumption caused by extra wiring, hence literature such as [68] claim that it's not suitable to apply this structure to compact designs. Other mesh-based designs include using resonant loads to generate a standing wave on the clock route rather than propagate a clock using a conventional CDN [69, 70], which helps to deliver clock with ultra-low skew. However, resonators will occupy an extra area, and any fabrication mismatch among resonators would consequently generate more clock skew [70]. Other approaches such as [71] and [72] incorporate clock paths outside the die area either using a silicon interposer or using the package to deliver clock signals to the logic which requires a timing reference. Clock signals are sent back to the die by solder balls or bond wires. Nevertheless, the overall wire length, especially in the horizontal plane, does not exhibit significant reduction, thus becoming a potential limit to delay reduction for designs with large dimensions. Besides, bond wires and solder balls bring excess wiring parasitics to the CDN, and hence the cross-talking noise and power consumption.

## 2.5 Summary

Considering the emerging interconnect techniques introduced above, if the low latency feature can be incorporated with the local wires, to simultaneously provide increased bandwidth and the flexibility to route to each of the local sinks, the CDN performance could potentially be improved remarkably. Given a qualitative comparison table shown

Table 2.2:  
Qualitative analysis of the CDN using the conventional and emerging interconnect techniques.

Attributes	Metallic wires [41], [42]	Optical interconnect [37], [54], [4]	Wireless interconnect [6], [7]	Potential hybrid of wireless and wires
Signal delay	Interconnect delay proportional to the parasitics and length.	Very low latency with the signal speed near the speed of light.	Latency limited by carriers.	Latency limited by carriers and the length of local wires.
Bandwidth	Limited by the wire length and wire delay.	Very high bandwidth with typical values up to 500 Gbps.	Depends on transistor cut-off frequency and technology nodes. Typical value up to 100~200 Gbps.	Depends on transistor cut-off frequency and technology nodes, as well as the local CDN delay.
Power	Dynamic power depends on the capacitive loading of wires.	High power consumption.	High power consumption with multiple clock transmitters/receivers.	Moderate power consumption, global power depends on the number of clock transmitter and receivers. Local power depends on the local CDN wirelength and loads.
Complexity	Need repeaters/buffers for wires with long distance, yet it's the most common and simplest interconnect.	High complexity [28], as some devices might not be CMOS compatible. Requires: laser source, photon detectors, optical waveguide and lens.	Medium complexity as it requires: wireless clock transmitters/receivers, on-chip antennas.	Medium complexity as it incorporates both wireless interconnect and metallic interconnect.
Reliability	Cross-talk exist, need shielding wires for sensitive nets.	High signal integrity.	Noise interference caused by multi-path propagation.	Noise interference caused by multi-path propagation.

in Table 2.2, the novel interconnects with different attributes are summarised and analysed accordingly. From the comparison, wireless interconnect shows the advantage of higher bandwidth comparing to the conventional metallic wires whilst maintaining the moderate complexity, which would be a compelling solution to the global clock distribution. Besides, the omnidirectional fan-out feature of the wireless EM wave propagation can provide a natural support for wireless clock distribution. However, generating a fully wireless CDN would bring excess power consumption. To route the clock signal to an arbitrary clock endpoint, a fully wireless CDN would be impractical to implement as the power of wireless CDN heavily depends on the overall number of the clock transmitters and receivers.

Alternatively, if clock distribution could incorporate the inherent fan-out feature of the wireless interconnect and the efficiency of a conventional tree/mesh for local clock distribution, one can estimate a significant improvement in terms of the overall power and uncertainty decrease for the reduction of global wires and local clock receivers. Literature such as [73] suggests that with the CMOS fabrication technology scaling down, it's now possible to integrate wireless interconnect on-chip based on developed communication techniques. Although they are not naturally built for clock distribution, some of the compact modulation schemes can be adjusted and applied to broadcast global signals on-chip. Therefore, we propose a hybrid wireless-wire architecture, which is based on an efficient modulation scheme to accommodate to variable clock frequency. Details of the proposed hybrid architecture will be given in the following Chapter 3 and 4.



## Chapter 3

# Proposed Models and Algorithms for Local Wired CDN

According to the previous chapters, radio-frequency-based interconnect represented by wireless interconnect shows its superiority in terms of overall throughput as well as modest energy consumption. Hence the CDN using emerging wireless interconnects would naturally become a competitive alternative for future on-chip clock distribution. However, as multiple clock receivers are required, a global wireless CDN needs to be carefully designed to balance area, power and performance. Also, as the local CDN is still using wires, it is essential to define an effective and energy-efficient local network design methodology, so that the overall energy efficiency and clock performance can take advantage of the proposed wireless global CDN.

In this chapter, several conventional architectures and the proposed architecture would be discussed and compared. Two local CDN generation algorithms are given with tree and mesh topology, respectively, which will be used for the local clock network design. The conventional CDNs will be implemented as baseline architectures using the algorithms given in this chapter and will be compared to the proposed hybrid method in Chapter 5 and 6, respectively.

We shall first clarify the notations to be used in the remaining of this section, in the Nomenclature.

### 3.1 Delay Models of Conventional CDN

To observe the behaviour of the interconnect and quantify the delay factor inside a CDN, different interconnect models are adopted to make a comparison between each other. For a typical metallic interconnect illustrated in Figure 3.1, the interconnect segment could be modelled as an RC sub-circuit consists of resistance and capacitance. Elmore has

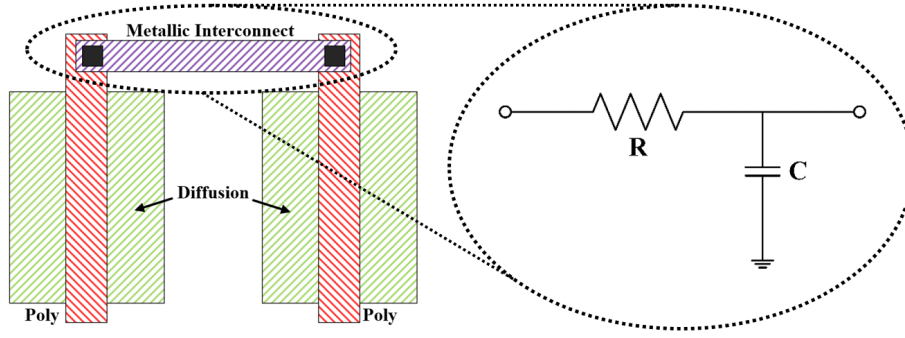


Figure 3.1: Conventional metallic interconnect structure and its lumped RC model.

developed an approximation model to describe wiring delay as the first-order moment of the system impulse response [74].

Given the transfer function of an interconnect segment normalised as:

$$g(s) = \frac{1 + a_1s + a_2s^2 + a_3s^3 + \dots + a_ns^n}{1 + b_1s + b_2s^2 + b_3s^3 + \dots + b_ms^m}, \forall (m, n) \in \{\mathbb{Z}^+, m \geq n\} \quad (3.1)$$

where the coefficients  $a_n, b_m$  in the numerator and denominator are real,  $m \geq n$ , in order to make the output of a linear time-invariant (LTI) system stable and converge, system poles must locate at the left side of complex s-plane. For a unit step input function, system transient response could be derived by:

$$U(s) = \frac{1}{s \cdot g(s)} \quad (3.2)$$

And hence the time domain step response could be given by the inverse Laplace transform of the s-domain product:

$$u(t) = \mathcal{L}^{-1} \left\{ \frac{1}{s} \cdot g(s) \right\} = \frac{1}{2\pi i} \lim_{x \rightarrow \infty} \int_{\gamma - iT}^{\gamma + iT} U(s) e^{st} dt \quad (3.3)$$

Given a unit-step input, the propagation delay can be commonly defined as the timing interval between the timing points that the input reaches 50% of its final input value and the output reaches 50% of its final output value [75], namely 50% logic transition delay. Since  $u(t)$  behaves several features like a monotonic increasing exponential function and has the final value of one (unit-step input), the derivative of  $u(t)$  which is the impulse response must be non-negative. Hence, the propagation delay could then be regarded as a probability distribution function (pdf) with a 50% logic transition delay of around the mean of this probability distribution given by:

$$T_d = \int_0^\infty t \cdot u'(t) dt \quad (3.4)$$

Meanwhile, this delay approximation has several mathematical equivalent expressions

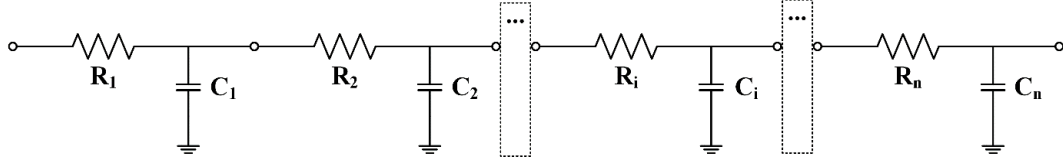


Figure 3.2: Distributed lumped circuit structure for improved accuracy of interconnect model.

such as the first moment of impulse response and the dominant pole of the system transfer function in complex  $s$ -domain. Given an example as a lumped circuit model. The voltage transfer function over capacitor  $C$  equals:

$$H(s) = \frac{Y(s)}{X(s)} = \frac{1}{sRC + 1} \quad (3.5)$$

where  $H(s)$  is the impulse response of this lumped RC network;  $Y(s)$  and  $X(s)$  are the output and input in  $s$ -domain, respectively. Time domain impulse response could then be derived by:

$$h(t) = \mathcal{L}^{-1}\{H(s)\} = \frac{e^{-\frac{t}{RC}}}{RC} \quad (3.6)$$

and therefore, Elmore delay could then be rewrite as:

$$T_D = \int_0^{+\infty} t \cdot h(t) dt = RC \quad (3.7)$$

which match the reciprocal of the dominant pole found inside system transfer function. Wyatt [76] raised a modified delay model which used the relationship that:

$$a_1 = \sum_{i=1}^n \frac{1}{z_i}, \quad b_1 = \sum_{i=1}^n \frac{1}{p_i}, \quad (3.8)$$

$$T_D = b_1 - a_1 \quad (3.9)$$

where  $a_i$  and  $b_i$  are the sum of the reciprocal of system zeros and poles respectively. Thus the time domain unit step response could then be given by the integral of impulse response:

$$u(t) = 1 - e^{-\frac{t}{T_D}} = 0.5 \quad (3.10)$$

and therefore the 50% logic transition delay is derived by:

$$t = \ln(2) \cdot T_D \quad (3.11)$$

which scale the delay of this single-stage RC network by a factor of  $\ln(2)$  rather than  $T_D$  originally raised by Elmore. This approximation seems suitable for system with an

outstanding dominant pole (significant different value than other poles) and no low-frequency zeros near it, that could be modelled as a single pole system. However, it becomes inaccurate as the total length of the interconnect increases. To have a more accurate interconnect model, for an interconnect wire with length  $L$ , distributed lumped model could be used to replace the single-segment model as shown in Figure 3.2. Where  $R_i$ ,  $C_i$  are the distributed resistance and capacitance respectively which equals to the resistance and capacitance per unit interconnect length  $r$  and  $c$ .  $R_m$ ,  $C_m$  are the total wire resistance and capacitance. For a distributed interconnect model with  $n$  identical segments based on the assumption that the wiring capacitance can be lump together, the effective time constant is now calculated as:

$$\tau_D = \frac{R_m C_m n}{n^2} + \frac{R_m C_m (n-1)}{n^2} + \frac{R_m C_m (n-2)}{n^2} + \dots + \frac{R_m C_m}{n^2} \tau_D = \frac{R_m C_m (n+1)}{2n} \quad (3.12)$$

Using this effective delay constant, given a unit step function input, the step response of a metallic interconnect in time domain could then be given by:

$$V_{out}(t) = u(t) \cdot (1 - e^{\frac{-t}{\tau_D}}) \quad (3.13)$$

And hence, for a 50% logic transition delay with a unit-step input, the distributed delay time can be given by:

$$T_D = \lim_{n \rightarrow \infty} \left[ -\ln(0.5) \cdot \frac{R_m C_m (n+1)}{2n} \right] = \frac{\ln(2) R_m C_m}{2} \quad (3.14)$$

From this equation, we could find that the propagation delay from source of the distribution RC network to the positive terminal of the capacitor  $C_n$  is further scaled by a factor of 2, which is more accurate and practical for modeling interconnect delays.

As most of the clock networks adopt H-tree as a basic implementation of CDN for its high balance and simplicity, an H-tree network is designed as the baseline architecture to evaluate the power and performance of the proposed design. H-tree network model is designed to have a clock signal injected at the centre of a square die area for global clock transmission. Tapered branch method is used in the H-tree model, which ensures that the reflected power at each branching node is compensated and minimised [44], as shown in Figure 3.3. More specifically, the baseline architecture is designed to have a branch width decreasing hierarchically with the propagation of clock signal towards lower levels. The impedance leaving the  $(k+1)$ -th level is set to be twice the impedance of the  $k$ -th level to get the desired performance [45]. Besides, 4 levels of H-tree are implemented to have a global branch with 16 leaf nodes. Lower levels of H-tree are neglected for overall system simplicity and replaced by load capacitors with different capacitance to mimic the impact of unbalanced loads.

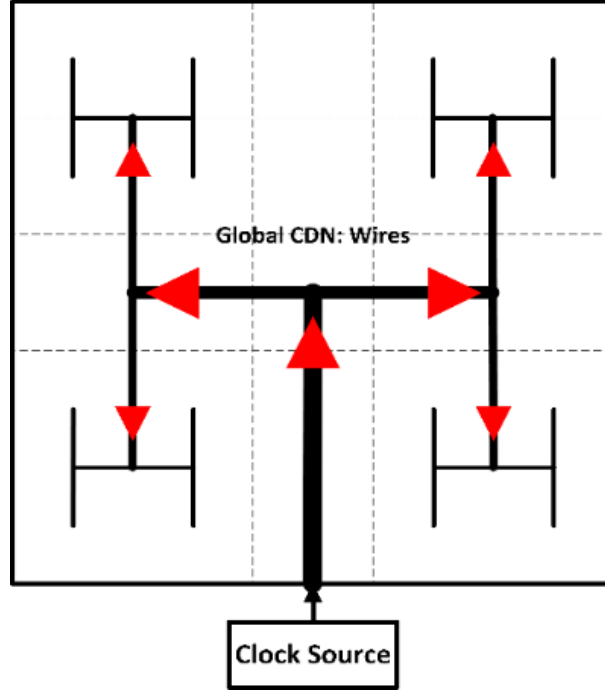


Figure 3.3: Conventional tapered H-tree network as baseline architecture with 16 fan-out nodes.

As per interconnect structure shown in Chapter 1 illustrated, wiring delay would essentially equals to the product of interconnect resistance and capacitance given by:

$$T_D = \tau L r_{wire} \cdot (L c_{wire} + C_{load}) = \tau L^2 \frac{\rho}{HW} (c_{cross} + c_{fringing} + \frac{C_{load}}{L}) \quad (3.15)$$

where the total capacitance consists of three components, cross-coupled capacitance  $c_{cross}$ , side-wall capacitance  $c_{fringing}$  and load capacitance  $C_{load}$ . The load capacitance could then be represented by:

$$C_{load} = C_i \quad (3.16)$$

where  $i$  is the fan-out index number of the  $i$ -level H-tree. For a fully symmetric tree structure, up to 16 fan-out nodes could be found. Besides, an H-tree could be naturally modelled as a binary tree composed of resistive and capacitive branches. The folded binary could have a structure illustrated in Figure 3.4.

Inside the binary tree, each of the wiring branch could be modelled as a sub-circuit segment consists of distributed resistance and capacitance. Two different sub-circuits could represent this concept, namely lumped model and  $\pi$ -model, respectively. The simple transfer function of a typical lumped model in s-domain could then be derived by:

$$G(s) = \frac{1}{sRC + 1} = \frac{1}{s\tau + 1} \quad (3.17)$$



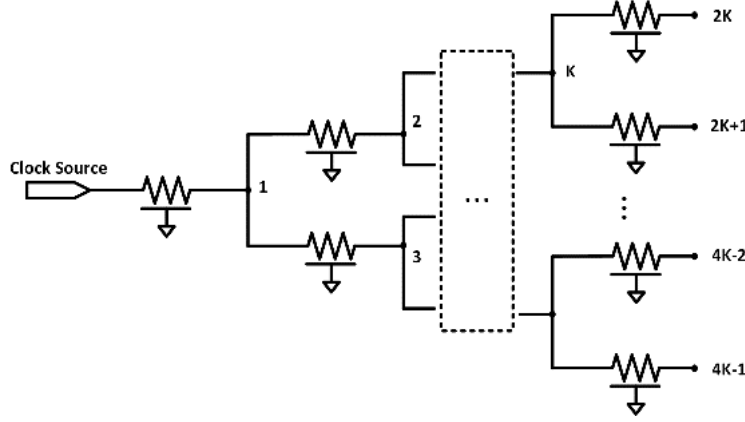


Figure 3.4: A  $k$ -level H-tree modelled as a folded binary RC tree with  $k^2$  fan-out.

where  $\tau$  is the pole of the transfer function and delay time constant composed of RC. Now if we treat an interconnect wire as a distributed network with  $n$  numbers of sub-circuits, the original delay constant could then be calculated as the combination of wiring resistance and capacitance.

For a 4-level clock tree structure, applying interconnect delay model on calculating the clock tree route delay, single branch unit length delay time constant at the  $k$ -th level using distributed sub-circuit model could then be modified as:

$$\tau_k = r_1(c_1 + c_2 + \dots + c_n + K_k) + r_2(c_2 + \dots + c_n + K_k) + \dots + r_n(c_n + C_k), \quad (3.18)$$

$$\tau_k = C_k \sum_{i=1}^n r_i + \sum_{i=1}^n r_i \cdot \sum_{i=1}^n c_i \quad (3.19)$$

where  $C_k$  is the load capacitance at the  $k$ -th level branching node.  $r_n, c_n$  are the distributed parameters for unit length resistance  $r$  and capacitance  $c$  respectively. If the distribution value  $n$  approaches to infinity, the unit length delay could then be derived as:

$$\tau_k = \lim_{n \rightarrow +\infty} \left[ C_{load} \sum_{i=1}^n r_i + \frac{n(n+1)r_i c_i}{2} \right] = r(C_{load} + \frac{c}{2}) \quad (3.20)$$

therefore the 50% logic transition delay of the  $k$ -th level H-tree branch responding to a unit step input is:

$$D(k) = \ln(2) \cdot L(k)r \left[ C_{load} + \frac{c}{2} \cdot L(k) \right], \forall k \geq 1 \quad (3.21)$$

where  $D(k)$  is the total delay at the  $k$ -th level.  $L(k)$  is the wire length between node  $k$  and its parent node  $k/2$  if  $k$  is a even node, and  $(k-1)/2$  if  $k$  is odd and larger than 1.

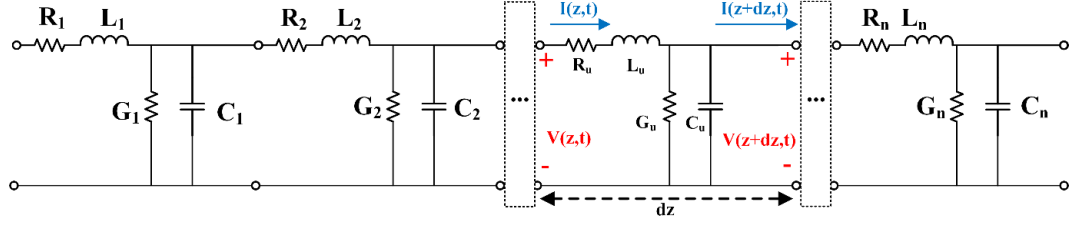


Figure 3.5: Transmission line model with distributed elements  $R_u$ ,  $L_u$ ,  $G_u$ , and  $C_u$ , respectively.

Hence for overall clock propagation delay from root node to any of the leaf node with  $m$  levels, it could be derived by:

$$C(m) = C(2m) + C(2m + 1) + c[L(2m) + L(2m + 1)], \quad (3.22)$$

$$D_{sum} = \sum_{i=1}^{2^{m+1}-1} D(i) + D_{buf}, \quad (3.23)$$

$$D(1) = \ln(2) \cdot R_{on} C(1) \quad (3.24)$$

where  $D_{buf}$  is the intrinsic delay of a driver cell,  $R_{on}$  is the on resistance of the clock buffer [77].  $C(m)$  is the lumped capacitance at the  $m$ -th node;  $D_{sum}$  is the overall delay from clock entry point to the  $m$ -th node.

Hence, under the assumption that the entire H-tree network could be modelled as a single pole system, Propagation delay of the clock signal could be calculated according to Equation 3.22 to 3.24.

### 3.2 Model of the Tree-based CDN

The above delay model only consists RC components, which could provide a quick estimation of the delay in an interconnect segment. However, with the wire dimensions getting smaller, the inductive effect also needs to be taken care of [78]. It is noticeable that, if a higher-order model is considered, the overall delay calculation would be more accurate. For example, for a system with an increasing frequency higher than multi-GHz, wiring inductance begins to show a significant impact on both signal propagation delay as well as power consumption because of the increasing system density, hence consequently coupling with nearby global interconnect. Besides, as the signal frequency keeps increasing, the clock signal would decay and delay more because of the transmission line effect [79]. Hence, to better accommodate CDN model to these emerging challenges, a second-order model containing the effect of inductive coupling is developed and studied.

According to the distributed transmission line model above shown in Figure 3.5, a single branch with single input buffer structure could be modelled as the combination of

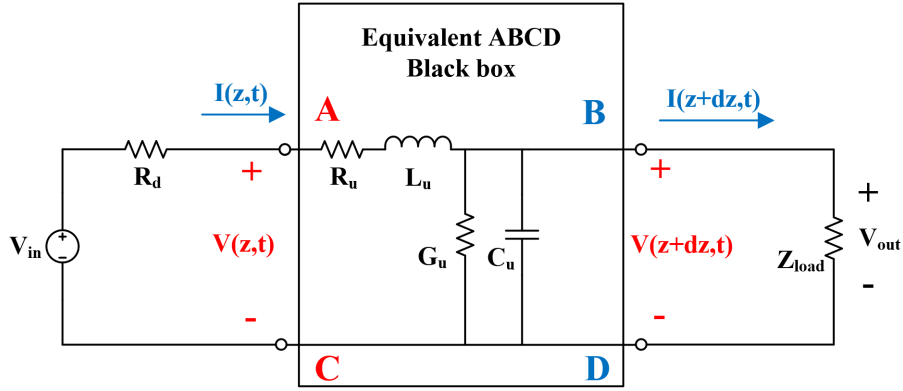


Figure 3.6: Equivalent ABCD representation of an interconnect in transmission line model.

lumped RLC elements. Given an interconnect with the length of  $l$ , unit length resistance, capacitance, admittance and inductance are given by  $R_u$ ,  $C_u$ ,  $G_u$  and  $L_u$  respectively. Each distributed segment has an interconnect length  $dz$ , a near end input voltage  $V(z, t)$ , and a far-end output voltage  $V(Z + dz, t)$ . Assume upper conductor is perfectly isolated with the ground line, that is,  $G_u$  is set to be zero. Hence, applying KVL and KCL to a single loop, its clear that:

$$V(z, t) = I(z, t) \left( R_u dz + s L_u dz + \frac{1}{C_u dz} \right), \quad (3.25)$$

$$I(z, t) = I(z + dz, t) + I_{C_u} \quad (3.26)$$

If  $dz$  is set to be approaching to zero, it could be derived that:

$$\frac{\partial^2 V(z, t)}{\partial z^2} = \gamma^2 V(z, t), \quad (3.27)$$

$$\frac{\partial^2 I(z, t)}{\partial z^2} = \gamma^2 I(z, t) \quad (3.28)$$

Also,

$$\gamma = \alpha + j\beta = \sqrt{(R_u + j\omega L_u)(G_u + j\omega C_u)}, \quad (3.29)$$

$$Z_o = \sqrt{\frac{R_u + j\omega L_u}{G_u + j\omega C_u}} \quad (3.30)$$

which is the well-known telegrapher equations. where  $\gamma$  is the propagation constant consists of distributed RLC values,  $Z_o$  is the characteristic impedance of the transmission determined by the dimension of the interconnect, hence essentially the distributed RLC values [78] and [80]. Solving differential equation set by using ordinary differential equation theory, the input impedance at the near end terminal could be given by:

$$Z_{in} = \frac{V(z, t)}{I(z, t)} = \frac{V^+ e^{j\omega t - \gamma z} + V^- e^{j\omega t + \gamma z}}{I^+ e^{j\omega t - \gamma z} + I^- e^{j\omega t + \gamma z}}, \quad (3.31)$$

$$Z_{load} = Z_o \frac{V^+ + V^-}{V^+ - V^-} = \frac{V^+ + V^-}{I^+ - I^-} \quad (3.32)$$

where  $Z_{load}$  is the load impedance,  $V^+, V^-, I^+, I^-$  are the insert and reflected voltage and current respectively. Substitute  $Z_{load}$  into  $Z_{in}$ , the input impedance at a distance of  $l$  away from the near end terminal could then be rewritten as:

$$Z_{in} = Z_o \cdot \frac{Z_{load} + jZ_o \tanh(\gamma l)}{Z_o + jZ_{load} \tanh(\gamma l)} \quad (3.33)$$

where  $l$  is the length of this interconnect. The input impedance of the single branch inside a H-tree is equivalent to its predecessors load impedance. According to transmission line theory, using two-port network model and ABCD matrix, this single branch also exhibits an in-out relationship that:

$$V(z, t) = A \cdot V(z + dz, t) + B \cdot I(z + dz, t), \quad (3.34)$$

$$I(z, t) = C \cdot V(z + dz, t) + D \cdot I(z + dz, t) \quad (3.35)$$

where  $V(z, t), I(z, t), V(z + dz, t), I(z + dz, t)$  are the injected and received voltage and current, respectively, according to Figure 3.6.

Substitute the above solutions to telegrapher equations to the above ABCD interconnect representation of Equation 3.34 and 3.35, an transmission line ABCD matrix can be derived by:

$$\begin{pmatrix} A & B \\ C & D \end{pmatrix} = \begin{pmatrix} \cosh(\gamma z) & Z_o \sinh(\gamma z) \\ \frac{1}{Z_o} \sinh(\gamma z) & \cosh(\gamma z) \end{pmatrix} \quad (3.36)$$

And therefore, for a system with input buffer impedance of  $R_d$  and capacitive loading with capacitance  $C_l$ , the single branch transfer function could then be derived by:

$$H(s) = \frac{V_{out}(s)}{V_{in}(s)} = \frac{V_{out}(s)}{A \cdot V_{out}(s) + B \cdot I_{out}(s) + C \cdot R_d V_{out}(s) + D \cdot R_d I_{out}(s)}, \quad (3.37)$$

$$H(s) = \frac{1}{\cosh(\gamma z) \cdot (1 + sC_l R_d) + \sinh(\gamma z) \cdot (Z_o sC_l + R_d Z_o^{-1})} \quad (3.38)$$

For a baseline tree architecture adopted in this section, the typical H-tree contains 4 levels of branches and 16 fanouts. From the above equations, load impedance can be obtained recursively, hence, by applying single branch transfer function onto the designated tree-based infrastructure, it is easy to get that the transfer function for H-tree is essentially the product of a series of single branch transfer functions along a specific route. Therefore, the universal transfer function from the root node  $N_0$  to any of the leaf nodes of a  $k$ -th level clock tree shown in Figure 3.3 and 3.4 can then be given

by:

$$H_k(s) = \frac{R_{origin}}{R_d + R_{origin}} \cdot \prod_{i=1}^k \frac{1}{\cosh(\gamma z)(1 + sC_{li}R_{di}) + \sinh(\gamma z)(Z_{oi}sC_{li} + \frac{R_{di}}{Z_{oi}})} \quad (3.39)$$

where  $k \geq 2$ ;  $R_{origin}$  is the input impedance at node  $N_0$ ;  $C_{li}$ ,  $R_{di}$ ,  $Z_{oi}$  and  $z_i$  are the load capacitance, input resistance, characteristic impedance at the  $i$ -th level and the interconnect length of the  $i$ -th branch, respectively.  $N$  stands for the set of nodes inside this fully-balanced H-tree RLC network. With this transfer function, it could be more accurate and convenient to quantify signal integrity including the time-domain output or the interconnect segment under different stimulations, by carrying an inverse Laplace transform after multiplying the transfer function of input and H-tree. Further detailed analysis of H-tree as a baseline architecture will be given in Chapter 5.

### 3.2.1 Balanced Bi-partitioning Algorithm for Clock Tree Synthesis

The notations of the algorithms proposed in this section are explained in Nomenclature. A balanced tree structure is always preferred when designing a conventional CDN for its skew-friendly features, such as an H-tree and X-tree. During the clock tree synthesis phase in a typical VLSI physical design flow, all the cells which need to get synchronised are defined as clock endpoints and are properly placed and optimised to avoid excess global congestion. These clock sinks with confirmed coordinates are then clustered into groups in which the number or the capacitance of the sinks can get balanced, following a top-down fashion, until each of the subsets only contains one clock sink. The sink subsets are then re-grouped pairwise in a bottom-up manner and the physical design CAD tool can then decide the length of the clock tree branch of each re-grouped sinks, based-on their capacitance and downstream delay, so that the overall delay to each clock endpoint can be equalised.

[9] gives a detailed study of a well-known heuristic algorithm named Balanced Bi-partitioning (BB) Algorithm, which controls the partitioning process based on the total loads of the two subsets. Two subsets are created when the difference of the capacitance from these two subsets are minimised. The BB algorithm also creates a boundary of all sinks within an octagon, as shown in Figure 3.7. The vertices (clock sinks) on the octagon are called the elements of reference sets. Half of the clock sinks within the octagon which are closer to the reference set, are divided into *subset A*. The rest of the sinks are grouped into *subset B*. By using an exhaustive search method, the reference set with minimum “merging cost” is then chosen, and these procedures are executed iteratively, until the last/leaf-level subset only contains one clock sink.

Each layer of the partition information is stored as a tuple, in which contains the  $2^{i-1}$  pairs of subsets in layer  $i$ . All tuples are retained from the root set to all leaf sets

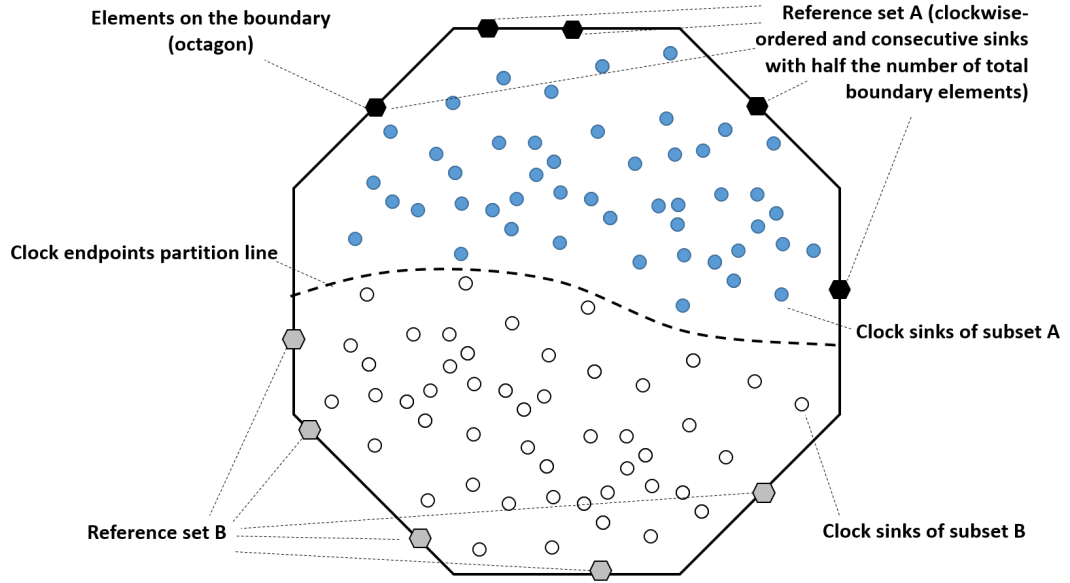


Figure 3.7: Top view of the Balanced Bi-partitioning (BB) Algorithm with half of the sinks on the boundary as a reference set. The clock sinks within the octagon that are closer to the reference set A are partitioned into subset A, the rest are grouped into subset B. Also, the summation of the capacitance of clock sinks are minimised.

for the subsequent merging stage, which is referred to as Deferred Merging Embedding (DME) Algorithm. The sinks are balanced and merged based on Manhattan geometry, as the metal layers for routing are either in the horizontal or vertical direction. All merging points are grouped together, which can satisfy that the delay between merging point and the two respective child merging points, according to the previous partition information, is equalised. The distance between the merging/balanced point and the clock sinks (named edges) are stored, which can be used in the following steps. An example of the merging scenario is shown in Figure 3.8, and conditions can be given that:

$$D_{MS_a} + D_{e_a} = D_{MS_b} + D_{e_b} \quad (3.40)$$

$$D_{MS_c} + D_{e_d} = D_{MS_c} + D_{e_d} \quad (3.41)$$

where  $D_{e_a}$ ,  $D_{e_b}$  represent the delay caused by the edge  $e_a$  and  $e_b$ , respectively. Similarly,  $D_{MS_a}$  and  $D_{MS_b}$  stands for the total delay from merging segment  $MS_a$  and  $MS_b$  to their leaf clock sinks, respectively. Using the delay models referring back to Equation 3.21 and 3.24, these procedures can ensure near-minimum clock skew.

As per this example, the collection of all possible merging points is called *merging segment*, which is a tilted line segment with a slope of either +1 or -1. The two merging segments can then be merged iteratively based on the partitioning information, until the merging process reaches the top/root sink set. The lower bound of merging cost is measured by the distance between two sink sets, which is defined as the minimum distance between two arbitrary sinks from the two sets, respectively. As the wire for

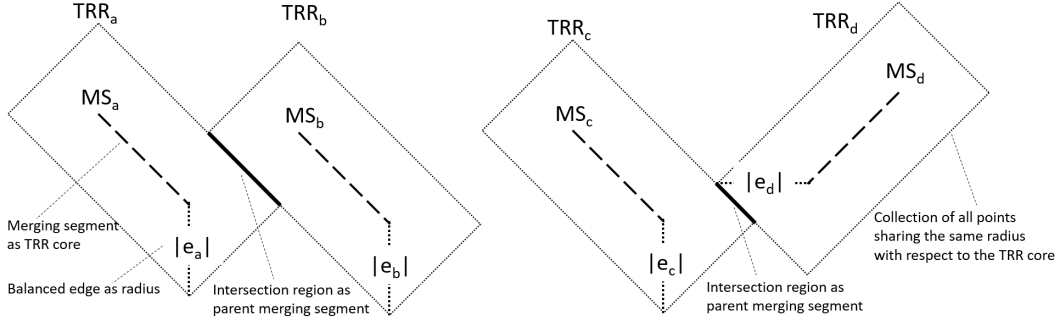


Figure 3.8: Example of the merging segment generation, with  $MS_a$ ,  $MS_b$  as a tuple element and  $MS_c$ ,  $MS_d$  as a tuple element, respectively. The generated parent merging segments are the intersection region of the two TRRs with the respective child merging segments as cores and calculated/balanced edges as radius.

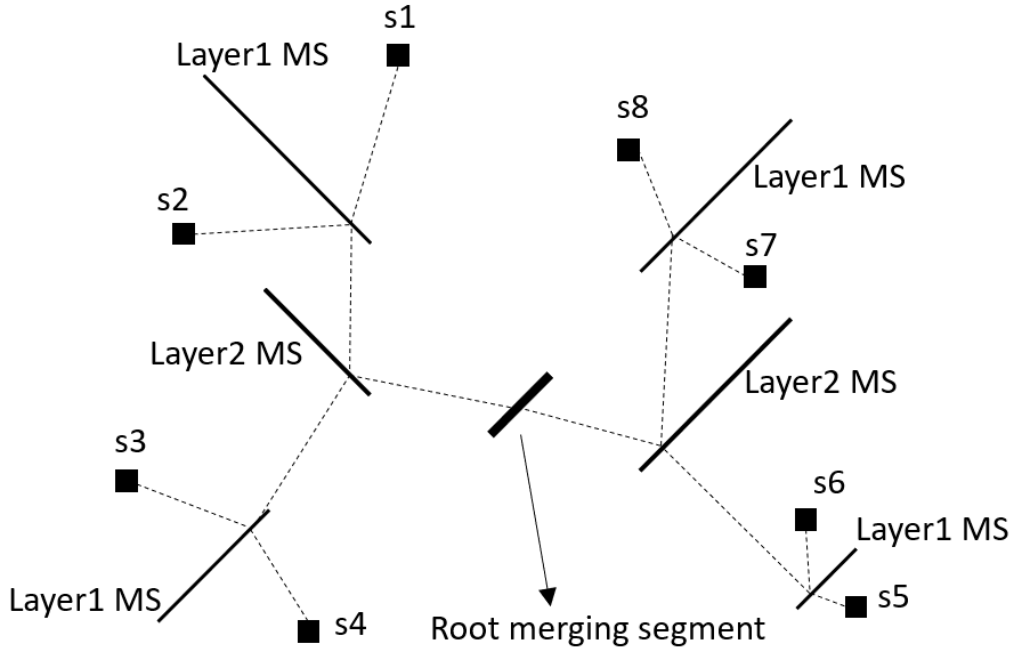


Figure 3.9: Example of the topology tree generation with a sink set containing 8 clock endpoints. The layer 1 MS are constructed in iteration 1 with tuple 1 (8 sinks) as input, the layer 2 MS are constructed in iteration 2 with tuple 2 (layer 1 MS) as input. The iterations are executed from a bottom-up fashion, until the root merging segment has been generated.

connecting two merging segments will contribute delay and capacitance, therefore, using the lower bound to merge two merging segments is always desirable, in light of improving energy efficiency and clock timing performance.

A *tilted rectangular region* (TRR) is then created based on the merging segment, which serves as the “core” of the TRR. The edge length which satisfies the balanced delay from the previous step is used as the “radius” of the TRR. The merging of the two merging segments is proven to be the intersection components of the two child merging segments

to be merged, as shown in Figure 3.9. Following these procedures, DME algorithm can then find out all possible locations of the merging points for delay balancing.

### 3.2.2 K-means Bi-partitioning (KBP) Algorithm for Clock Tree Synthesis

Based on the DME algorithm, we proposed a 2D-DME variant, which adapts the conventional algorithm to CDN generation. The exact locations of branching/merging points can then be determined based on the results of the proposed merging algorithm as well as following *find – exact – location* procedures in [9] and [81]. Detailed procedures for the proposed algorithm are given in the remaining of this section, the following equations are based on Manhattan geometry, unless otherwise specified.

Several related research [9], [82], [83], [84] have presented clustering algorithms such as the method of means and medians (MMM) or balanced bi-partitioning (BB) algorithms, which perform a near-balanced bi-partitioning to minimise the difference in the number of sinks or load capacitance in each subset. Thus, the clock skew in each subset is considered to be minimal and hence a ZST can be constructed based on the division results.

Energy efficiency is of paramount importance when considering clock distribution in a many-core system. According to [85], clock power is one of the major contributions in the overall power budget. With the dark-silicon phenomenon, only a subset of cores is allowed to run simultaneously with certain peak power and thermal constraints, energy efficiency has become one of the major limitations for boosting the clock speed to a higher level in a many-core system [86], [87]. And hence, it is vital to choose an energy-efficient CTS strategy during clock design phase.

In this section, we propose a new balanced-partitioning variant during CTS, which trades clock skew with overall energy efficiency by using less wires for delay balancing, for those applications whose power priority is higher than skew priority. Under the assumption that the clock power is proportional to the overall wire length [9] in a sink set with a specific topology, we first quantify the Manhattan distance between two sinks within a sink pair as:

$$dist(i, j) = |x_i - x_j| + |y_i - y_j|, \forall i, j \in k \quad (3.42)$$

where  $x_i, x_j, y_i, y_j$  are the x and y coordinates of the sink  $i$  and sink  $j$ , respectively. The distance between two sink sets is defined as the Manhattan distance between the center of mass of the two sets. Considering the capacitive loading of each sink in the given set, when calculating the center of mass of the given finite sink set, the coordinates of each sink is weighted by its capacitance accordingly. Let sink set  $a$  and  $b$  be the subsets of



set  $k$ , the mutual CDN interconnect wire length which connects a pair of sink sets (also as known as merging length) between set  $a$  and  $b$  can then be defined as follows:

$$WL_m(k) = |x_{cm}(a) - x_{cm}(b)| + |y_{cm}(a) - y_{cm}(b)| \quad (3.43)$$

$$x_{cm}(k) = \frac{\sum_{m=1}^n (x_m \cdot w_m)}{\sum_{m=1}^n (w_m)}, \forall k \in \theta \quad (3.44)$$

$$y_{cm}(k) = \frac{\sum_{m=1}^n (y_m \cdot w_m)}{\sum_{m=1}^n (w_m)}, \forall k \in \theta \quad (3.45)$$

where  $WL_m(k)$  yields the mutual distance between subset  $a$  and  $b$ ;  $x_m$ ,  $y_m$  and  $w_m$  are the x coordinates, y coordinates and the capacitance of sink  $m$ , respectively. Similarly,  $x_{cm}(k)$  and  $y_{cm}(k)$  are the x and y coordinates of the center of mass in the given sink set  $k$ .  $\theta$  is the set with all clock sinks.

In addition, we use the sum of the variance of two subsets containing finite vertices to yield the internal CDN interconnect wire length (fusion length) for a sink set  $k$  to be clustered, which can be defined as:

$$WL_s(k) = \sum_{i=a}^k \left\{ \sum_{j=1}^q [|x_j - x_{cm}(i)|^2 + |y_j - y_{cm}(i)|^2] \right\} \quad (3.46)$$

where  $q$  is the maximum number of sinks for an arbitrary subset  $a$  or  $b$ .

Based on the above notations, the total power cost  $PC(k)$  measured by wire length for sink set  $k$  with bi-partitioning is then given by the sum of mutual-CDN wire length (merging length) and self-CDN wire length (fusion length):

$$PC(k) = WL_m(k) + WL_s(k), k = a \cup b, \forall k \in \theta \quad (3.47)$$

where  $WL_m(k)$  and  $WL_s(k)$  stands for the merging length and fusion length, respectively. According to [88] and [89], a k-Means algorithm can naturally generate clusters with minimised sum of the intra-cluster variances. Hence, to reduce the overall wire usage during each partitioning process, topology tree generation using the proposed algorithm can follow the idea of updating the centroid of each subset at each iteration and forcing the type/number of the partitioning to be 2, to construct a balanced binary tree. Note that to reduce the impact of poor initial centroid selection, k-Means++ algorithm [90] is applied in the proposed KBP algorithm.

The detailed  $xy$ -cut procedures are given in Algorithm 1. If a sink set  $k$  to be partitioned contains  $n$  tiers, where  $n \geq 1$ , all sinks in different tiers will be projected in a  $xy$ -plane. The algorithm then partitions the projected sinks, as shown in Algorithm 1. The output  $index_a$  and  $index_b$  will then be mapped back to their original tier accordingly. The algorithm first of all randomly select the first initial centroid at the first iteration. Sinks with larger Manhattan distance will have larger probabilities to be chosen as the second

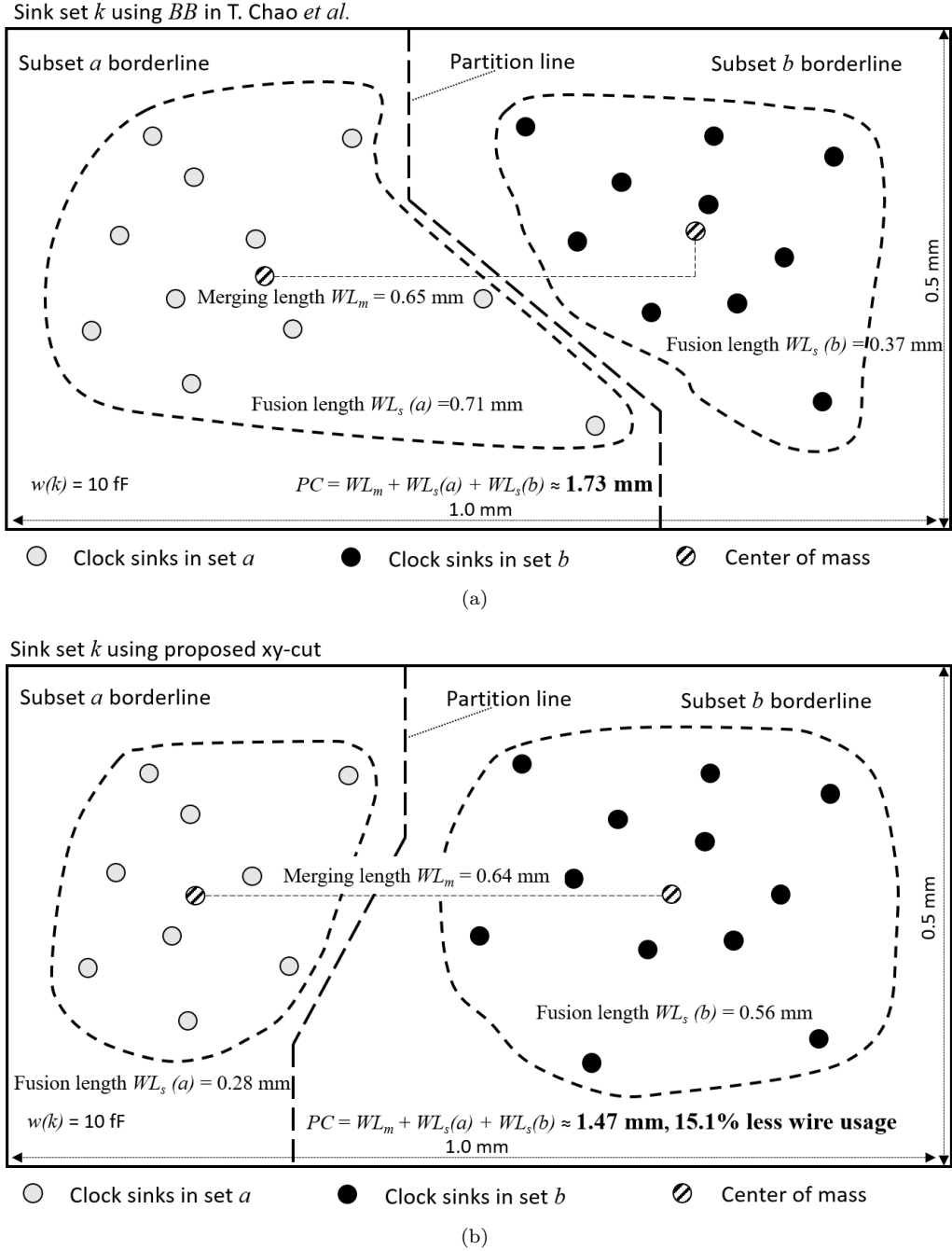


Figure 3.10: Illustration of the *xy-cut*, using (a) a conventional balanced bi-partition in [9], and (b) our proposed *xy-cut* algorithm for minimum wire usage. The total wire reduction between the proposed method and the conventional method is around 15.1% in this example, thus improving energy efficiency.

initial centroid. After the initialization process, all clock sinks will be classified to the corresponding group which has the minimal distance to either of the two initial centroids. The coordinates of two centroids will then be updated, respectively, among two groups of sinks and used to generate the distance vector in the next iteration. This updating process is performed iteratively until the coordinates of the centroids remain unchanged,

**Algorithm 1:** XY-cut for Low Wire Usage

---

Function( $x_m, y_m, w_m$ );

**Define Input:**

$x_m$  : x coordinate vector of clock sink set  $m$ ,  
 $y_m$  : y coordinate vector of clock sink set  $m$ ,  
 $w_m$  : capacitive loading vector of clock sink set  $m$ ;

**Define Output:**

$index_a$  : index vector of the sub-sink set  $a$ ,  
 $index_b$  : index vector of the sub-sink set  $b$ ;

**Define Parameters:**

$i_{max}, k, \delta$ ;

**Algorithm Procedures:**

Initialization:  $i = 2$ ; number of child sets = 2;  
 $cent_a(1) = \text{random index in set } m$ ;  
 $D = \text{dist}(m, cent_a(1))$ ;  
 $x = \text{random}(0 \text{ to } 1)$ ;

**for** ( $j = 1; j \leq \text{length}(m); j++$ )  
 $P(j) = D(j)^2 \cdot (\sum_m D(j)^2)^{-1}$   
**if**  $j > 1$   
 $Px(j, 1) = Px(j - 1, 2)$ ;  
 $Px(j, 2) = \sum_{k=1}^j P(k)$ ;  
**else**  
 $Px(j, 1) = 0; Px(j, 2) = P(j)$ ;  
**end if**  
**if**  $Px(j, 1) \leq x < Px(j, 2)$   
 $cent_b(1) = j$ ;  
**end if**  
**end for**

$cent = (cent_a(1), cent_b(1))$ ;

**while**  $|cent(i) - cent(i - 1)| \geq \delta$  **do**  
reset  $index_a$  and  $index_b$  to empty;  
**for** ( $j = 1; j \leq \text{length}(m); j++$ )  
**if**  $\text{dist}(j, cent_a(i)) \geq \text{dist}(j, cent_b(i))$   
 $index_b = index_b \cup j$ ;  
**else**  
 $index_a = index_a \cup j$ ;  
**end if**  
**end for**  
 $i = i + 1$ ;  
 $cent_a(i) = (x_{cm}(k), y_{cm}(k)), \forall k \in index_a$ ;  
 $cent_b(i) = (x_{cm}(k), y_{cm}(k)), \forall k \in index_b$ ;  
 $cent = (cent_a(i), cent_b(i))$ ;  
**if**  $i > i_{max}$ ;  
**break**  
**end if**  
**end while**

**return:**  $index_a, index_b$ ;

---

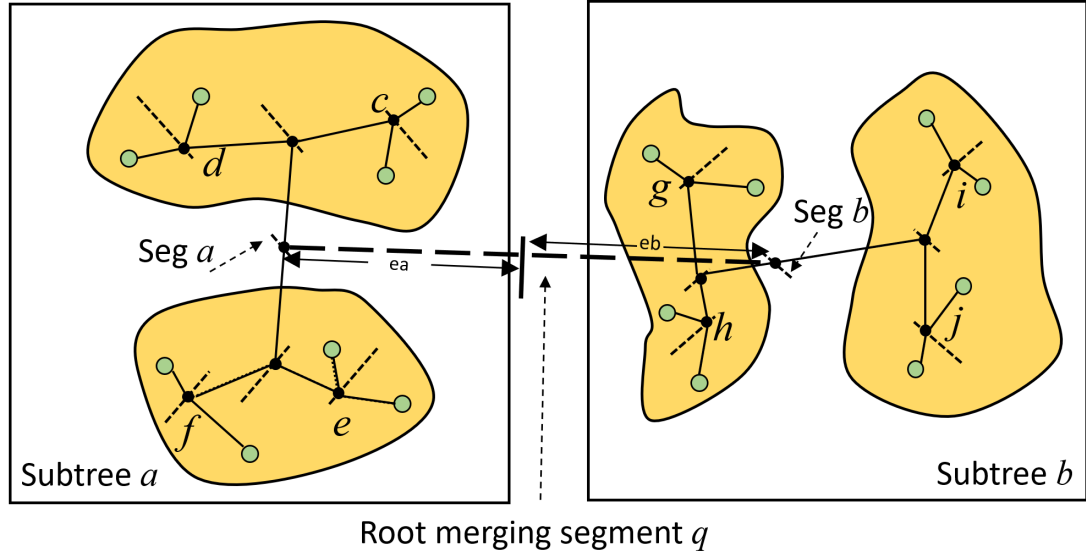


Figure 3.11: Example of the merging tree generation with two sink sets, each of which contains 8 clock endpoints (not to scale). The two subtrees are merged with root segment  $q$ .

quantified by a user-defined parameter  $\delta$ , and the max iteration number  $i_{max}$  to prevent the algorithm from running into deadlock. The output  $index_a$  and  $index_b$  contains  $2n$  subsets in total, if the input sink set  $k$  contains  $n$  tiers.

An example using both the conventional method and the proposed method for  $xy$ -cut is presented in Figure 3.10. As shown in Figure 3.10, comparing to conventional balanced bi-partitioning, the proposed  $xy$ -cut can have an up to 15.1% reduction in CDN wire length, based on the proposed metrics of power cost  $PC$ . Specifically, the proposed  $xy$ -cut can reduce total fusion length significantly, because of its low intra-cluster nature, which offers a promising solution to low power 3-D CDN design.

### 3.2.3 Merging Tree Generation

After the topology tree has been created, one can then use the DME algorithm to build the full merging tree, which contains all acceptable merging points for detailed routing/physical implementation of the clock tree. The DME algorithm [81] generates a tree of merging segments in a bottom-up fashion, which yields the possible placement of every merging/branching point in the binary CDN. Given a topology tree, the conventional DME algorithm can generate 2-D CDN connection with minimal wire length. By balancing the wire delay of each route using Elmore Delay Model [12], the total delay from leaf to root node can be minimised, and thus generating a CDN with minimal clock skew (ZST). The detailed merging algorithm is shown in Algorithm 2.

As shown in Figure 3.11, the example clock sinks are evenly spread across the rectangular area, which contains 8 sinks. The clock sinks will first be merged into merging segments

**Algorithm 2:** Merging tree generation

---

```

Function(topology(m));
Define Input:
topology(m) : partition matrix for sink set m;
Define Output:
layer :
    detailed merging tree matrix containing layer-
    -information for the tree of merging segments;
Define Parameters:
tE,  $\kappa$ , ea, eb, temp, seg, cap, delay;
Algorithm Procedures:
Initialization: i = 1,
inputvec = bottom layer of topology(m);
while i ≤ layers of topology tree do
    set j = 1;
    pair = pairs of sets to be merged in inputvec;
    while j ≤ pair do
        [a, b] = inputvec(j);
        if merging type == 2D
             $\kappa = \min(\text{dist}(a, b))$ ;
            solve ea and eb, satisfying con1 and con2;
            con1 :  $t_E(ea) + t_a = t_E(eb) + t_b$ ;
            con2 :  $ea + eb = \kappa$ ;
        end if
        temp = temp ∪ seg(ea, eb, cap, delay);
        j ++
    end while
    if j == pair;
        layer(i-th row) = temp;
        inputvec = temp;
        reset temp to empty;
    end if
    i ++;
end while
return: layer;

```

---

*c*, *d*, *e*, *f*, *g*, *h*, *i* and *j*, and the generated merging segments will be retained as the input vectors (tuples) for the next iteration of merging process. After three iterations of merging, every sink has been merged following the DME merging procedures. The two merging segments *a* and *b*, which are root merging segment of subtree *a* and subtree *b*, respectively, will need a root merging segment to link each other. Following the merging procedures, the two merging edges *ea* and *eb* are calculated, based on the delay balancing equation as presented in Algorithm 2. The generated merging information and edge information will be stored in the merging tree matrix, which includes details such as the vertices of the merging segment, lumped capacitance and delay at each layer of the 3-D CDN. Based on the output matrix *layer*, the exact coordinates for every merging point can then be located, following the conventional positioning method in [9].

Since the proposed algorithm will generate a fixed number of sub-sets at each iteration, the time complexity of the proposed algorithm for abstract tree generation using  $XY$ -cut is  $O(n)$ , where  $n$  is the number of the sink of an input set of clock endpoints. By comparison, the conventional method using BB-DME yields a worst-case complexity of  $O(n^3 \cdot \log n)$ . On the other hand, related works such as [82] and [91] adopt similar method of means and medians (MMM) methods for clock sink clustering, which present a reduced complexity of  $O(n \cdot \log n)$ . For the embedding phase, both the proposed work and [82] adopt an extended Deferred-Merging-Embedding algorithm for merging tree generation, which gives a similar complexity of  $O(n)$ . Therefore, the overall complexity reduction of the proposed method provides a promising solution to clock tree generation with a large number of data input. A detailed evaluation of the proposed tree generation algorithm is given in Chapter 5.

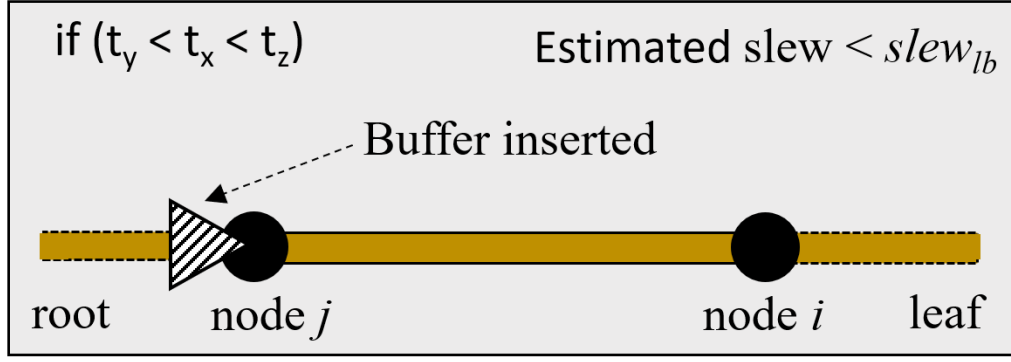
### 3.2.4 Buffer Insertion with Inverted Slew-estimation

To provide a reliable clock signal, a CDN should be capable of providing enough driving strength to drive an output, such as a clock gating cell or a flip-flop [92], [93], [94]. We therefore propose a buffer insertion algorithm, during the generation of the merging segment tree. The proposed buffer insertion algorithm monitors the clock slew at each merging point and then decide whether to add a buffer, depends on the slew estimation results. Therefore, clock buffers do not need to be uniformly located in the same layer of the CDN, thus saving clock power and offering higher flexibility to CDN design, compared with [85].

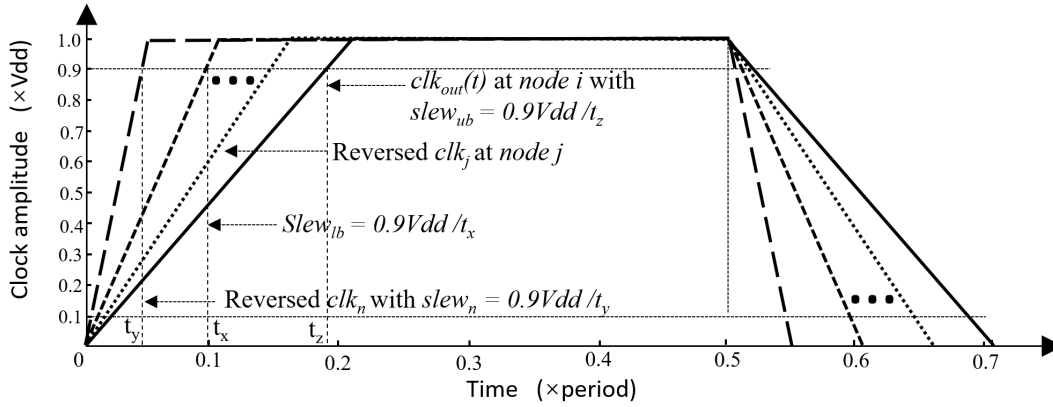
Since the slew estimation is executed at each iteration of the merging tree generation process, the extra delay and capacitance of clock buffers can be balanced by the subsequent merging procedures, thus having minimal impacts on the overall skew output. By contrast, conventional buffer insertion methods such as [95] will place clock buffers on a completely unbuffered tree. Considering a CDN generated by the proposed algorithms above, the unbuffered tree can already provide a very-low skew output. The conventional buffer insertion method will produce larger unbalanced delay to each clock sink, thus deteriorating the overall skew performance. Detailed procedures of the proposed buffer insertion method are given as follows.

Under the assumption that the slew constraint is around 25% of the total clock period  $T$ , the ideal clock signal will have a time interval of  $0.25T$  for signal attenuation along the CDN route. Considering the worst-case scenario, in which an arbitrary clock sink  $i$  have an output clock signal with  $0.25T$  clock slew, set the time-domain continuous clock output as  $clk_{out}(t)$ , the output signal in s-domain can be derived by running the Laplace transform given by:

$$clk_{out}(s) = \mathcal{L}\{clk_{out}(t)\} \quad (3.48)$$



(a)



(b)

Figure 3.12: Illustration of the (a) proposed buffer insertion method with inverted slew-estimation and (b) its time-domain behavior when the estimated slew is smaller than the slew lower bound of the driver. Without a clock driver, the child clock end point might violate the slew constraint.

If the precedent wire segment connecting the parent node  $j$  and the sink  $i$  has a transfer function of  $H_i(s)$ , the following relationship can be given:

$$H_i(s) = \frac{clk_{out}(s)}{clk_{in}(s)} = \frac{num_i(s)}{den_i(s)} \quad (3.49)$$

where  $num_i(s)$  and  $den_i(s)$  stands for the numerator and denominator for  $H_i(s)$ , respectively. Therefore, the time-domain input signal at node  $j$  can then be derived as:

$$clk_{in}(t) = \mathcal{L}^{-1} \left\{ \frac{\mathcal{L}\{clk_{out}(t)\} \cdot den_i(s)}{num_i(s)} \right\} \quad (3.50)$$

Therefore, the clock slew at the input node can be accurately measured based on the transfer function model proposed in Section II. The bottom level output signal  $clk_{out}(t)$  can be initialised as the clock signal with 25% slew rate, during the first iteration of the proposed 3-D CDN generation. By iteratively calculating the clock slew in a bottom-up order during the merging process, the clock slew at the input node will ultimately

approach to a slew lower bound  $slew_{lb}$ , which can be found when defining clock design rule constraints (max transition and max capacitance). If the calculated slew is less than  $slew_{lb}$ , a buffer can be inserted at the current merging point, providing enough driving strength to drive the very end of the CDN with a clock load, as shown in Figure 3.12.

### 3.3 Model of the Mesh-based CDN

Conventionally, modelling a clock mesh is considered to be more complicated than a balanced tree structure, because of its coupling and loop infrastructure. Related research [10] has suggested a method based on transient simulation in SPICE, namely the Sliding Window Scheme (SWS), which divides the large global clock mesh into smaller components called “Windows”. Since the attenuation of an input logic state transition is the exponential function of propagation distance between an arbitrary injection point and clock sinks, propagation delay will be calculated only inside the window region [10]. After the calculation is finished, the current window will move to a new position and this iteration will continue until all of the mesh units have been covered. Hence, complicated mesh modelling becomes a “divide and conquer” problem, which helps to reduce overall calculation complexity. Although this method could effectively provide an accurate estimation of timing information in a clock mesh, it is still time-consuming and memory-sensitive to finish all iterations [96] if the size of the mesh is increasing.

Thus, to provide an accurate and normalised solution to mesh modelling and simultaneously reduce unnecessary simulation time, two local mesh modelling methods are proposed in this chapter. Both the proposed methods adopt similar divide and conquer techniques to simplify the modelling process. However, instead of using SPICE-based modelling, transfer functions are adopted inside the region of interest. Details of both methods are listed below.

#### 3.3.1 Isolated Mesh Model

Similar to [97], a basic unit mesh cell of  $2 \times 2$  is chosen for overall simplicity. Moreover, loads and stubs are considered small enough and therefore could be lumped at the joint of mesh cells. The effective unit cell model is shown in Figure 3.13. According to Figure 3.13, mesh cells are considered to be isolated from each other. In other words, no coupling effect exists between each cell units, and no currents flow out of an isolated cell. Therefore, the analysis of the entire clock mesh simplifies to the analysis of isolated unit cells. Based on the transfer function of a single interconnect segment, the effective transfer function from buffer output to an arbitrary corner inside an isolated cell can be



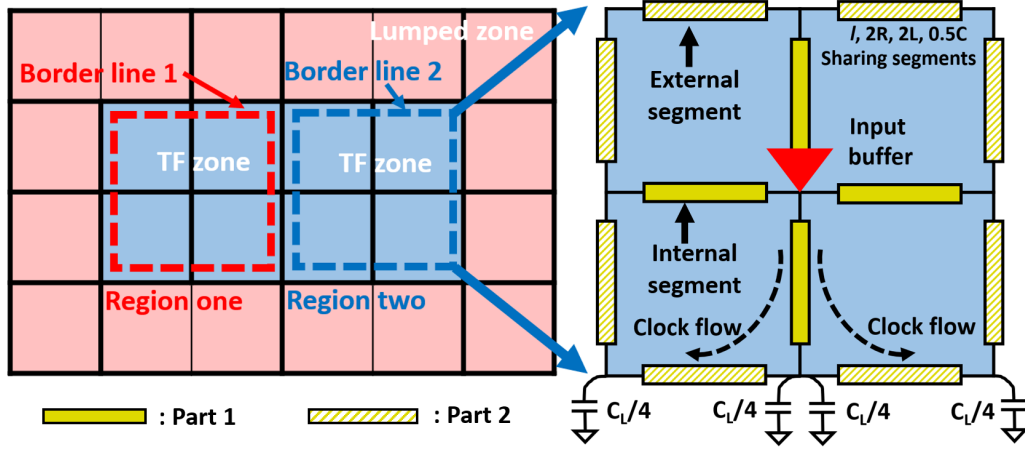


Figure 3.13: Partial top views of the proposed isolated model, including transfer function zone (TF zone) and lumped zone. Unit cells are considered independent of each other. ©2021 IEEE

derived as:

$$H_{unit}(s) = H_{int}(s) \cdot H_{ext}(s) \quad (3.51)$$

where  $H_{int}(s)$  and  $H_{ext}(s)$  are transfer functions for part one and part two respectively as shown in Figure 3.13. Load impedance of each mesh segment is modelled as the parallel combination of local capacitive loading and input impedance of successive mesh segments, which can be given by:

$$Z_{load}(s) = \frac{1}{sC_L} \parallel \left( Z_o \frac{Z_L + \tanh(\gamma l_x) \cdot Z_o}{Z_o + \tanh(\gamma l_x) \cdot Z_L} \right) \parallel \dots \quad (3.52)$$

where  $Z_L$  is the load impedance for successive mesh segments;  $Z_o$  and  $\gamma$  are the characteristic impedance and propagation constant of the segment of interest respectively given by:

$$Z_o = \sqrt{\frac{R_i + sL_i}{G_i + sC_i}} \approx \sqrt{\frac{R_i + sL_i}{sC_i}} \quad (3.53)$$

$$\gamma_i = \sqrt{(R_i + sL_i)(G_i + sC_i)} \approx \sqrt{(R_i + sL_i)(sC_i)} \quad (3.54)$$

From the above equations, load impedance can be obtained recursively. The internal transfer function in a unit cell could then be given by:

$$H_{int}(s) = \frac{1}{\cosh(\gamma l_x) + \frac{Z_o}{Z_{load_{iso1}}} \sinh(\gamma l_x)} \cdot \frac{Z_{in}}{Z_d + Z_{in}} \quad (3.55)$$

$$Z_{load_{iso1}} = Z_{CL} \parallel Z_{in1} \parallel Z_{in2} \quad (3.56)$$

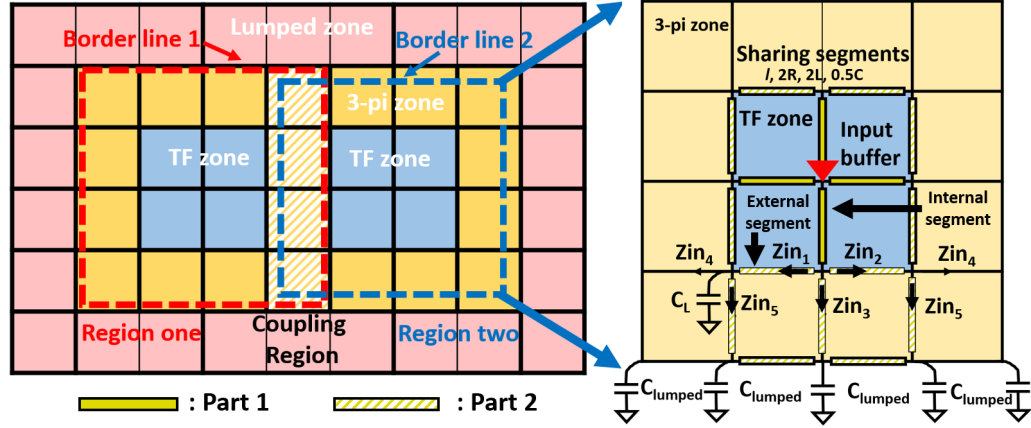


Figure 3.14: Partial top views of the proposed coupled model (not to scale), including transfer function zone (TF zone),  $3-\pi$  zone and lumped zone. Unit cells are considered coupled to each other shown in the overlapped region. ©2021 IEEE

Similarly, the external transfer function could be derived by:

$$H_{ext}(s) = \frac{1}{\cosh(\gamma l_x) + \left(\frac{Z_o}{Z_{load_{iso2}}}\right) \sinh(\gamma l_x)} \quad (3.57)$$

$$Z_{load_{iso2}} = Z_{C_L} || Z_{C_{lumped}} \quad (3.58)$$

where  $Z_{C_{lumped}}$  is the impedance of loads and wiring capacitance inside the lumped zone seen at the boundary between two TF zones, as in Figure 3.13.

### 3.3.2 Coupled Mesh Model

Since there will always be some difference in time of arrival for all clock buffers or Rx connected to local clock mesh, it is essential to evaluate the characteristic of averaging clock skew within the local mesh network. Therefore, coupled model take neighbour cells into consideration to accurately model the coupled section between different regions. according to Figure 3.14, an extra layer is added between the TF Zone and Lumped Zone, namely  $3-\pi$  Zone. The reason for this extra region is to provide a buffer layer to better mimic coupling effects with neighbour cells. Similar to the isolated method, transfer functions and lumped capacitors are used to model coupling cells inside the TF zones and Lumped zones respectively, However, a  $3-\pi$  model is adopted to trade slightly more calculation complexity with higher accuracy as shown in Figure 3.14. We thereby calculated the effective input impedance of the interconnect inside the  $3-\pi$  zone recursively as follows according to Figure 3.15:

$$Z_\alpha = (Z_{R_w} + Z_{L_w} + Z_{C'_w}) || Z_{C_w} \quad (3.59)$$

$$Z_\beta = (Z_{R_w} + Z_{L_w} + Z_\alpha) || Z_{C_w} \quad (3.60)$$

$$Z_\gamma = Z_{in_{3-\pi}} = (Z_{R_w} + Z_{L_w} + Z_\beta) || Z_{C'_w} \quad (3.61)$$

where  $R_w, L_w, C_w$  are the total resistance, inductance and capacitance (including load capacitance) of the segment of interest, respectively. Therefore, the internal transfer function for coupling cells can be modified as:

$$H_{int}(s) = \frac{1}{\cosh(\gamma l_x) + \frac{Z_o}{Z_{load_{cpl1}}} \sinh(\gamma l_x)} \cdot \frac{Z_{in}}{Z_d + Z_{in}} \quad (3.62)$$

$$Z_{load_{cpl1}} = Z_{C_L} || Z_{in1} || Z_{in2} || Z_{in3} \quad (3.63)$$

in which  $Z_{in3}$  represents the input impedance in 3-pi Zone. Similarly, the external transfer function is then modified as:

$$H_{ext}(s) = \frac{1}{\cosh(\gamma l_x) + (\frac{Z_o}{Z_{load_{cpl2}}}) \sinh(\gamma l_x)} \quad (3.64)$$

$$Z_{load_{cpl2}} = Z_{C_L} || Z_{C_{lumped}} || Z_{in4} || Z_{in5} \quad (3.65)$$

where  $Z_{in4}, Z_{in5}$  are the later stage input impedance seen at the corner of TF zone in a coupling cell. Taking advantage of the extra layer, the overall current output at a joint of interest is then given by the superposition of the response to different local buffers in 4 neighbour coupling cells. Hence, the transient response at an arbitrary joint  $e$  in the local CDN region is given by:

$$V_{out} = \sum_{k=1}^{n_b} (I_{cell_k}) \cdot Z_{load_e}, \forall e \in \{\theta_{m,n}\} \quad (3.66)$$

where  $I_{cell_k}$  is the current output generated by buffer  $k$  at joint  $e$ . Since  $I_{cell_k}$  is the superposition of all nearby buffer outputs, coupling cells take coupling effects into account. Therefore, this proposed method can model the local mesh better, even when input skew is injected into the local clock mesh from different neighbour buffers. The evaluation of both isolated cell and coupling cell are given in Section IV.

In summary, the overall interconnect delay from the central clock transmitter to an arbitrary clock sink (flip flop) would equal to the combination of global delay and local delay, which can be given by:

$$\begin{aligned} D_{total} &= D_{global} + D_{local} \\ &= v \cdot dist + (T_{V_{out}} - T_{V_{in}}) \end{aligned} \quad (3.67)$$

where  $V_{in}$  and  $V_{out}$  equals to  $0.5V_{dd}$ ;  $v$  is the velocity of the electromagnetic wave transmitted out from a clock Tx;  $V_{in}$  stands for the recovered clock signal output from a clock receiver.  $V_{out}$  is the actual clock signal received at an arbitrary clock sink. Assuming that the EM signal travels near the speed of light and local interconnect delay is relatively constant compared with conventional CDN, global interconnect delay during clock distribution could be reduced significantly, therefore allowing a larger range of frequency

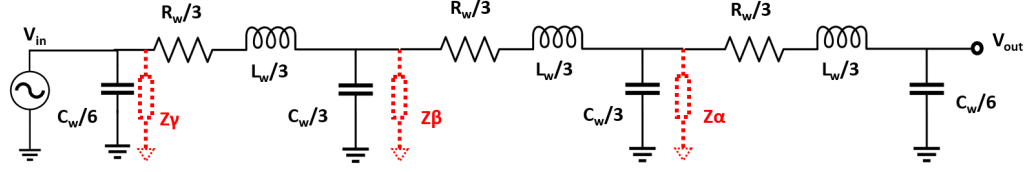


Figure 3.15: Typical structure of an RLC 3- $\pi$  model [10] which has been used in our proposed model for representing interconnect segments longer than 100  $\mu\text{m}$  inside 3- $\pi$  region.

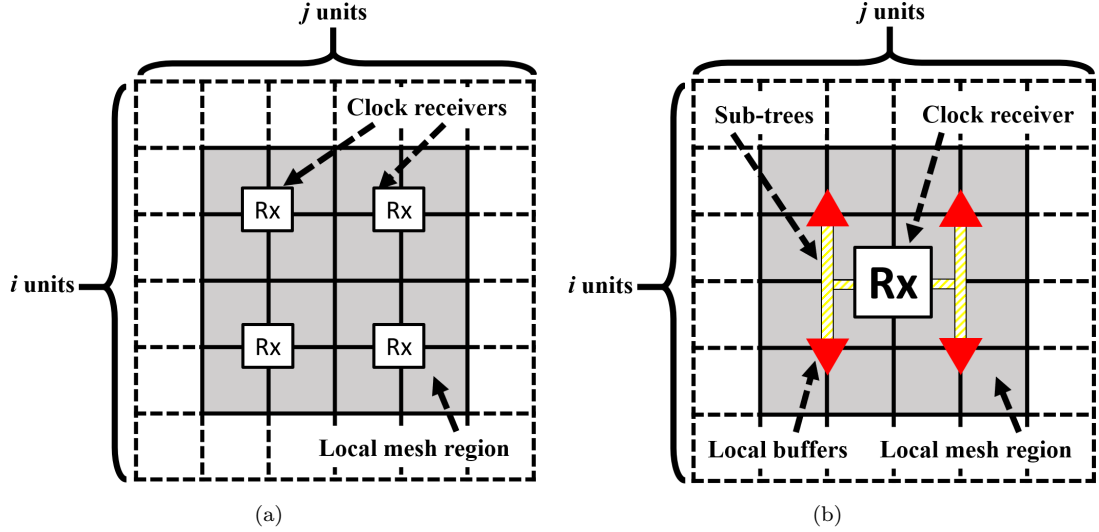


Figure 3.16: Top view of (a) distributed planning and (b) concentrated planning in an arbitrary  $i \times j$  local mesh network.

to be distributed. Besides, a local clock mesh can effectively mitigate output clock skew against a wide range of input skew, which helps to solve any potential timing violations.

### 3.3.3 Clock Receiver Planning

Clock receivers take the place of the conventional local clock buffers for distributing global clock signals to local CDN. It's important to quantify the area a clock receiver could cover. Obviously, clock Rx in close proximity to clock Tx experiences less delay, and by contrast, those located far from Tx would experience a higher delay, the delay difference between the shortest, and longest path from Tx to Rx being the global CDN skew. Although the local network has the potential of minimising uncertainties by applying the proposed framework, it is a more energy-efficient and resource-friendly choice to mitigate global skew at the first step by carefully allocating quantities and positions of Rx. This paper proposes two receiver insertion methods in terms of different topology. As shown in Figure 5.19(a), the first planning strategy is called the distributed method, in which a clock Rx is allocated at the centre of a  $2 \times 2$  unit mesh cell serving as a local clock buffer. As each Rx is assumed to drive its own mesh area, total clock

latency is minimised because of the very low global delay. However, due to the different placement of the clock Rx, the global skew will be higher.

An alternative Rx planning strategy is called the concentrated method is illustrated in Figure 5.19(b), in which a clock Rx drives local CDN with a series of local buffers, located in the centre of larger mesh area compared to the distributed method. A balanced sub-tree is utilised to deliver recovered clock signal with a very low global skew to local CDN. Thus, the concentrated method trades higher clock latency with better skew performance, compared with the distributed method.

Total power consumption for both distributed and concentrated planning are given as follows:

$$P_{distributed} = P_{mesh} + \left(\frac{n}{2}\right)^2 \cdot P_{Rx} \quad (3.68)$$

$$P_{concentrated} = P_{mesh} + P_{sub-tree} + P_{Rx} \quad (3.69)$$

where  $P_{mesh}$  is the power of local mesh CDN only;  $n$  is the number of mesh segment of interest, in this case  $n$  equals to 4;  $P_{Rx}$  is the power of a proposed clock Rx;  $P_{sub-tree}$  is the power dissipated in the sub clock tree connecting clock Rx with local mesh. It's worth noting that the total power of global clock Rx will become the major contribution of power dissipation in distributed Rx planning, as it will increase exponentially as a function of mesh segment numbers.

In summary, both planning methods can be applied to different scenarios. For delay-sensitive design constraint, distributed planning can fit in better for its low global latency feature. By contrast, if the application is more clock skew sensitive, concentrated planning shows a superior capability of skew reduction. Details of experimental results for both global and local CDN will be given in Chapter 5 and 6.

### 3.3.4 Local Mesh Design

In order to further reduce the skew in local CDN, it is essential to define a local distribution scheme based on different design constraints, such as distribution area, interconnect dimensions, power budget and upper skew bound. Researches such as [68] and [97] could be partially applied to local mesh planning, but still requires modification to define a boundary between wireless global CDN and wired local CDN.

From the evaluations of the local mesh, clock skew could be modelled as a function of the mesh size. Local skew decreases in a monotonous way with the increase in mesh numbers, as the size of each mesh, and hence the local interconnect length and number of clock endpoints in each mesh reduces. Also, to reduce the local network complexity, the stubs connecting clock endpoints to the clock mesh are using a comb structure, which directly tie the sinks to their nearest clock mesh segments where possible.

---

**Algorithm 3:** Local Mesh Design with Bounded Skew
 

---

 Function( $l_{ub}, sk_{ub}, P_{ub}, T$ );
**Define Input:**
 $l_{ub}$  : length upper bound,  $sk_{ub}$  : skew upper bound,

 $P_{ub}$  : power upper bound,  $T$  : clock period;
**Define Output:**
 $sk$  : overall clock skew,

 $P_{sum}$  : overall power consumption;
**Define Parameters:**
 $D_{global}, D_{local}, D_{i,j}, A_k, \tau, \theta_{m,n}, \mathbb{Z}^+$ ,

 $Z_{in}, P_i, num_x, num_y, I, I_{max}, C_t, V_{dd}$ ;
**Algorithm Procedures:**
 Initialization:  $sk = +\infty$ ;  $num_x = num_y = 0$ ;  $I = 0$ ;
**while**  $sk > sk_{ub}$  **do**
 $num_x = num_x + 2$ ;

 $num_y = num_y + 2$ ;

 $i$  : root of the CDN in  $I$ -th iteration;

 $j$  : index array of clock endpoints;

 $j = 0$ ;
**while**  $j < C_t$  **do**
 $D(i, j) = D_{global} + D_{local}, \forall i, j \in \theta_{m,n}$ 
 $sk = \max D(i, j) - \min D(i, j), \forall i, j \in \theta_{m,n}$ 
 $A_k = 2TV_{dd} \cdot (\tau k^2 \pi^2)^{-1}, \forall k \in \mathbb{Z}^+$ 
 $P_i = \sum_k^\infty A_k^2 \cdot Re(Z_{in}^{-1}), \forall k \in \{2\mathbb{Z}^+ - 1\}$ 
 $P_{sum} = P_{sum} + 0.5P_i$ 
 $j = j + 1$ 
**end while**
**if** ( $I > I_{max}$ )

 relax  $sk_{ub}$ ;

**break**;

**end if**
 $I = I + 1$ 
**end while**
**if** ( $P_{sum} > P_{ub}$ )

 relax  $P_{ub}$ ;

**end if**
**return:**  $sk, P_{sum}, num_x, num_y$ ;

---

The clock uncertainties such as skew decrease as a function of the segment length in a clock mesh, therefore, to reduce clock skew within the local wired CDN region, one can increase the mesh density recursively until the design constraint has been reached. Details of the procedures and pseudo-codes are given in Algorithm 3, where  $l_{ub}, sk_{ub}, P_{ub}$  are the upper bound of total segment length, skew and power constraint, respectively;  $num_x$  and  $num_y$  are the number of x-dimension segments and y-dimension segments, respectively.

An initial mesh segment number will be given at the very beginning of the algorithm,

which can generate a skew number which exceeds the skew upper bound. In order to reduce max clock skew, more metallic segments are inserted into the mesh cells, to provide better shorting effect [98]. After each iteration, the clock skew and power will be compared with the given specifications such as skew bound and power budget. If the current mesh size can satisfy such numbers, the algorithm will return the mesh size and other parameters such as clock skew and power estimation accordingly. The two proposed mesh cell models do not affect the overall iteration counts, but the algorithm with coupled model will execute slightly longer for better accuracy. The proposed algorithm shows a worst-case time complexity of around  $O(m \cdot n)$ , where  $m$  is the max iteration number to prevent the algorithm from a deadlock,  $n$  is the number of clock sinks given to be synchronised.

For a test circuit with over 360k FFs, the proposed algorithm can generate a  $64 \times 64$  local mesh network within 2 hour 45 minutes and around 1.2 GB peak memory usage on an Intel i7 processor [99] (algorithm run on a single-core). As a comparison, although using test machine with different specifications and testing circuits, the SPICE-based SWS implementation in [10] uses 6 hours and 48 minutes for a similar  $65 \times 65$  mesh design on a 300k FFs circuit with relatively high accuracy. Therefore, the proposed algorithm can provide a relatively efficient solution of generating a mesh network within a reasonable amount of time at the cost of slightly reduced accuracy, comparing to a full-SPICE analysis. Details of the test circuit and the analysis of the accuracy of the proposed local mesh design algorithm with two mesh models will be given in Chapter 5.

Overall power consumed in the local clock region can be calculated using Fourier series based method[78], which provides an accurate and comprehensive solution to calculating power consumption in an interconnect segment based on different clock parameters such as the logic transition time  $\tau$  and the overall clock period  $T$ , according to Algorithm 3. This algorithm would be repeated several times until one can get an optimal initial mesh design within acceptable criteria. Alternatively, this framework could also be used for power minimisation as well, if the design of interest is more power-sensitive.

In summary, the proposed mesh planning framework trades CDN timing performance with power consumption and wiring resources. One can easily adapt this framework to different design constraints or power budgets, thus provides more flexibility to initial CDN design and estimation. Further evaluations of this planning algorithm will be given in Chapter 5 and 6.

### 3.4 Model of the Proposed Hybrid Wireless-Wired CDN

As per Figure 3.17, the proposed architecture is designed to match the synchronisation area of the baseline architecture, which also has a 16-nodes global distribution network. A global clock transmitter is working as a master node transmitting clock signal out.

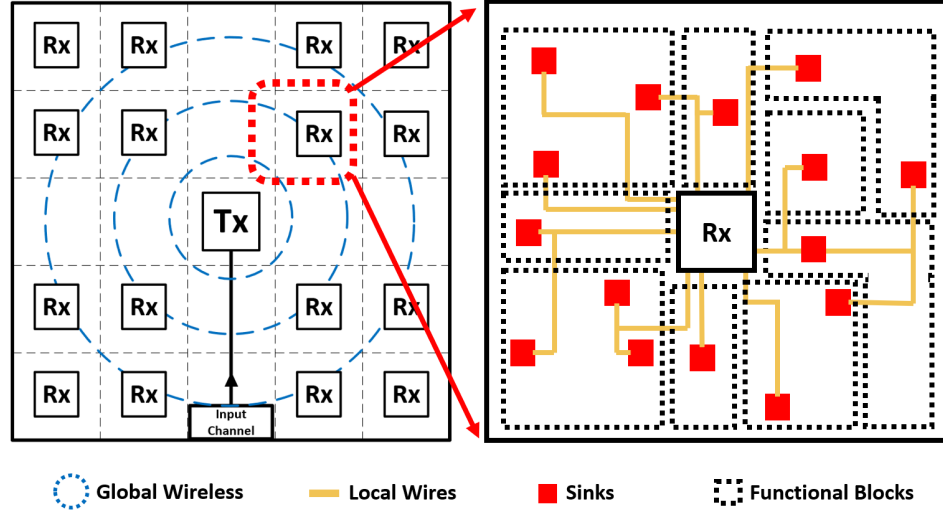


Figure 3.17: The proposed hybrid clock distribution network which combines both global RF-I and local metallic wires.

Besides, receivers are working as slave nodes which could listen to the transmitter to receive and recover clock signal accordingly. As the overall clock propagation delay equals to:

$$D_{hybrid} = D(i, j) = \omega(i, j) + T_{Tx} + T_{Rx} + T_{local}, \forall i, j \in \theta(m, n) \quad (3.70)$$

where  $\omega(i, j)$  stands for RF signal propagation time inside silicon substrate from node  $i$  to node  $j$ ,  $T_{Tx}$  and  $T_{Rx}$  represent to transmitter and receiver circuit propagation delay respectively,  $T_{local}$  equals to clock delay inside local wire CDN from the output of receivers to clocked registers.  $\theta(m, n)$  stands for the set of all nodes including clock source and clock sinks inside proposed  $m \times n$  distribution area. Clock skew is therefore defined by:

$$sk_{total} = \max D(i, j) - \min D(i, j), \forall i, j \in \theta(m, n) \quad (3.71)$$

The variations in clock TRx circuit and local wires will generate an arbitrary amount of delay which can impact the overall skew performance. Therefore, for two clock receivers under PVT variations, the clock skew need to be quantified as:

$$sk_{variation} = D_a - D_b = (T_{Rx_a} + T_{local_a} + \omega(i, a)) - (T_{Rx_b} + T_{local_b} + \omega(i, b)) \quad (3.72)$$

where  $D_a$  and  $D_b$  stand for the clock latency from Tx to Rx  $a$  and Rx  $b$ , respectively. [85] has suggested an inverter FO4 delay estimation method based on corners and effective transistor channel length. Based on [85], for a worst-case SSSS corner in NMOS, PMOS, wiring parasitics and supply voltage at a nominal temperature of 25 , an 27% speed variation can be given compared with a nominal TTTT corner. Hence, considering a worst-case scenario that one arbitrary receiver  $a$  is running at SSSS corner and receiver



b is running at FFFF corner, the Rx circuit can produce up to 54% delay variation, based on the above estimation method.

Since the TRx pairs can share one clock Tx in this model, Rx variation becomes the major issue to be considered. For a model length of 1cm with 16-receiver planning as shown in Figure 3.17, the largest communication distance among all radius can be given by approximately 2.8 mm, hence generating a delay difference  $\Delta\omega(i, j)$  of around 20 ps. Based on the estimation above and circuit simulation in TT corner in TSMC 65 nm pdk, the Rx circuit will generate a variation of around 100 ps, which is one order of magnitude higher than  $\Delta\omega(i, j)$ .

On the other hand, considering the worst-case corner which will impact local clock distribution, local CDN wires under  $RC_{max}$  corner will generate a delay variation compared with the wires in  $RC_{min}$  corner, which can be given by:

$$\Delta T_{local} = L_a R_a (0.5 L_a C_a + C_{load_a}) - L_b R_b (0.5 L_b C_b + C_{load_b}), \forall a, b \in \theta(m, n) \quad (3.73)$$

where  $L_a$ ,  $R_a$ ,  $C_a$ ,  $load_a$ ,  $L_b$ ,  $R_b$ ,  $C_b$ ,  $load_b$  are the effective clock wire length, unit resistance, unit capacitance and lumped capacitance for local synchronisation region  $a$  and  $b$ , respectively. For a specific example, based on the parasitic data from [18] and the model shown in Figure 3.17, under a 20% RC variation [100], the maximum delay variation can be given by approximately 40 ps. Using a non-default routing rule for local wire clock distribution can better tolerate the RC variations than using nominal wire dimensions.

To sum up, the PVT variations will have an impact on the overall performance (mostly related to clock latency and skew) in terms of the Rx circuit, wireless communication and local clock distribution. Among all impacting factors, clock Rx circuit contributes the largest clock latency variation under typical worst-case corner analysis, which remains to be an interesting topic for the future study about fine-grain receiver design that can better handle different corners.

To achieve a balanced global RF propagation delay  $\omega(i, j)$ , clock transmitter is designed to be located at the center of the proposed chip area broadcasting clock signal globally. In addition, clock receivers are uniformly placed over the chip area which replace the baseline global branches. Note that Tx node toward different receiver nodes with the same radius are assumed to have the same  $\omega(i, j)$ , hence the same skew. However, receivers located on inner circles naturally have smaller radius than those on the outer circles, thus lead to a small amount of skew between near and remote receivers. This skew is defined by:

$$sk_{global} = c \cdot (\max\{|a_i - a_j|\}), \forall i, j \in \theta(m, n) \quad (3.74)$$

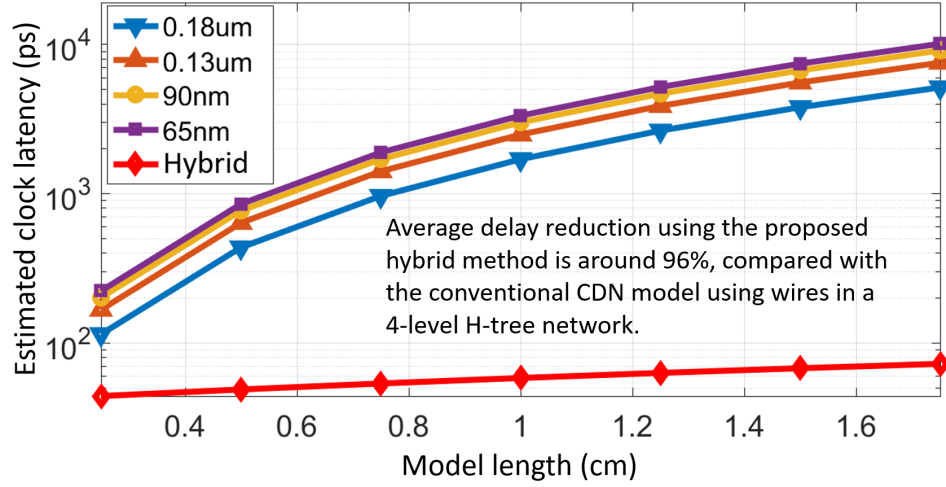


Figure 3.18: Estimated results using the proposed hybrid model and a conventional H-tree CDN with 96% average delay reduction.

where  $c$  stands for the EM wave propagation speed in silicon substrate,  $a_i$  and  $a_j$  stands for the radius around the transmitter. The maximum skew is limited to a proportion of the maximum difference of the radius between different transmitter-receiver pairs.

Besides, to reduce the impact of imbalanced capacitive loading, clock receivers (Rx) are placed following a ‘near-big’ and ‘far-small’ rule (NBFS) according to their load capacitance or, to be more specific, local buffer tree size. Transceiver pairs then transmit and receive a global clock signal via EM radiation without conventional wires. Finally, the recovered clocks are passed to a local CDN to synchronise digital logic in a different technology node mounted on top of the interposer. For overall system simplicity, the global wireless clock distribution adopts an On-Off Keying (OOK) based clock modulator and demodulator, which can modulate an input clock signal onto the carrier at clock Tx and recover RF signal back to baseband local clock at clock Rx.

Assuming that electromagnetic signals can travel near the speed of light in the dielectric layer, the transfer function for global CDN can be approximated as:

$$H_{global} = H_{Tx} \cdot H_{Rx} \cdot H_{dielectric} \quad (3.75)$$

$$= \frac{e^{-c \cdot dist \cdot s}}{s^2 \cdot R_T C_T R_R C_R + s \cdot (R_T C_T + R_R C_R) + 1} \quad (3.76)$$

where  $c$ ,  $dist$  denotes the speed of light and distance between clock Tx and Rx, respectively;  $R_T$ ,  $R_R$ ,  $C_T$ ,  $C_R$  represent for the effective resistance and capacitance for clock Tx and Rx, respectively. If we consider that propagation delay within Tx and Rx circuits is relatively constant for different stimulus input, the potential clock skew for hybrid architecture would merely depend on the allocation of mounted dies.

An estimated result in terms of the clock latency of the conventional CDN and the proposed hybrid CDN is given in Figure 3.18 shown above. Applying the two approaches

on the same test model with varying model side length and fixed clock loads, referring back to Figure 3.3, the maximum clock distribution delay from the clock root to the clock sinks is significantly reduced based on the assumptions and the models introduced above. An average delay reduction between the proposed hybrid model and the conventional 4-level H-tree CDN is around 96%, which yields a compelling delay performance.

### 3.5 Summary

To sum up, the wireless global clock Tx and Rx can reduce the clock latency significantly, as the conventional congested wires can be jumped. Therefore, the CDN performance using the proposed hybrid solution can be improved remarkably.

In addition, as the location of the clock Tx and clock Rx are fixed, the global clock uncertainties can be considered predictable and therefore providing more flexibility to the local CDN design. Also, the overall clock uncertainties can be reduced by adjusting the displacement for different logic with different loads.

Details of circuits for clock Tx and Rx are given in Chapter 4. Detailed evaluations for both the proposed CDN and the baseline architectures will be represented in Chapter 5 and 6, respectively.



## Chapter 4

# Proposed Clock Transmitter and Receiver for Global Wireless CDN

A typical RF transmitter involves components such as a local oscillator (VCO) to generate a carrier wave with the desired frequency, a modulator/mixer to up-convert baseband signal frequency, and an amplifier to amplify modulated signal power before or after being transmitted through an antenna. For clock delivering, some components for the purpose of long range communication can be removed to gain overall system simplicity. Hence, this chapter will examine different architectures of the proposed wireless clock transmitter and receivers.

Referring back to Chapter 2, a frequency divider based wireless approach will suffer from frequency distortions inside the VCO, hence produce unpredictable clock uncertainties. Our proposed architecture provides a solution to this potential drawback, that is, to build the system using a digital modulation scheme. Since the only signal to be transmitted for the one-to-all communication is the clock signal with a fixed signal pattern, therefore, an efficient and energy-efficient modulation scheme is necessary.

Hence, considering limited system power budget and simplicity, complicated system structures using modulation schemes such as Quadrature Phase Shift Keying (QPSK) and 16-Quadrature Amplitude Modulation (QAM) or even 64-QAM are neglected, as they are mainly designed for data communication to prevent issues like inter-symbol-interference (ISI) and simultaneously increase channel throughput. Thus, they will unnecessarily consume excess energy during clock distribution stage. Although they have been proved to be efficient, regardless of their outstanding bandwidth efficiency, the area and power overhead might be crucial for proposed clock distribution infrastructure with multiple receiver nodes.

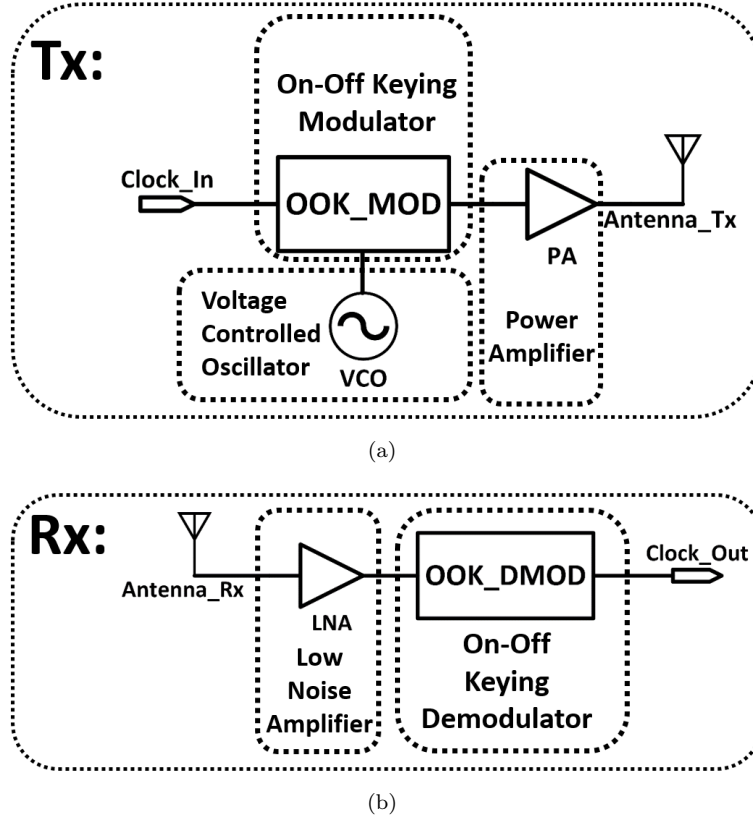


Figure 4.1: Block diagrams of the proposed (a) clock transmitter and (b) clock receiver.

Thus, a simplified and efficient modulation scheme is adopted for the purpose of simplicity and power-efficiency. Amplitude-Shift Keying (ASK) has been regarded as one of the most reliable and well-known digital modulation schemes. With signal logic level “1” represented by a positive voltage amplitude level, signal logic level “0” is modulated to another level, normally *gnd* level or a negative voltage amplitude, to distinguish the non-return-to-zero (NRZ) coding.

On-Off Keying (OOK) is an exceptional form of ASK. Just as its name implies, while a digital “1” is being modulated to a certain voltage (normally  $V_{dd}$ ), digital “0” would have zero voltage or at *gnd* level at the modulator output. Hence, the clock signal would only be transmitted during the period of logic “1” and would barely be transmitted during the period of logic “0”. Also, the dynamic switching power is much reduced as the overall system is at “on” state if the clock is at logic ‘1’ level, hence the power almost merely depends on the duty cycle of the clock signal. This symbolic feature makes it a power-saving alternative when comparing to an always-on modulation scheme.

However, as each symbol only contains one bit, OOK symbol rates is equal to its bit rate, thus making it less spectrum efficient compared with complex modulation schemes such as 16-QAM. This is a trade-off between bandwidth and power and it is an important characteristic for low power design. Considering the signal to be transmitted is a toggling

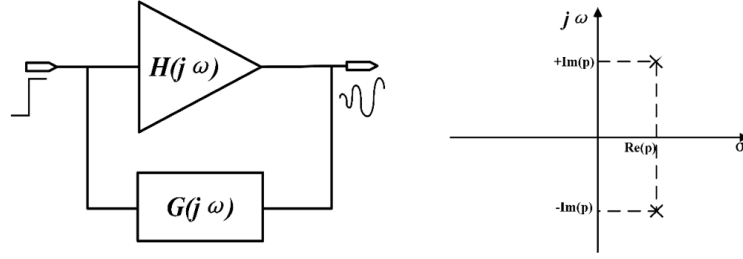


Figure 4.2: General oscillatory condition of feedback system.

clock, it is essential to adopt OOK for the purpose of low CDN power consumption. The proposed transmitter and receiver architecture are shown in Figure 4.1.

## 4.1 Clock Transmitter Design

### 4.1.1 Voltage-Controlled Oscillator Design

As clock signal needs to be modulated on a stable carrier wave with constant signal frequency oscillating in millimetre band to reduce the size of the integrated antennas, voltage control oscillator (VCO) design is quite essential as it would consume a large portion of DC power of the entire system [85], [101]. Besides, VCO output is set to be differential so as to alleviate cross talk inside transmitter, thus providing better performance and SNR.

Considering a general feedback system (i.e. an amplifier) with an open loop gain of  $H(j\omega)$  at frequency  $\omega$ , for the output track the input signal with finite time, it is necessary that system poles should be located on the left-half of the plane. By contrast, the oscillator is basically a non-stable close loop system which has its poles located on the right-half plane, so that the system could start to oscillate due to tiny disturbance at the input [102].

To get the VCO oscillating at its resonance frequency, certain feedback techniques could be attached on a general amplifier to satisfy certain criteria. For a simple amplifier, suppose the output is 180-degrees out of phase compared with input, for a negative feedback system, the output is enlarged by the subtraction and hence forms the functionality of regeneration allowing output swing becomes larger. To be more specific, the feedback system need to satisfy:

$$|\beta H(j\omega_0)| \geq 1, \quad (4.1)$$

$$\angle \beta H(j\omega_0) = 180^\circ \quad (4.2)$$

where  $\beta$  is the feedback coefficient. Named Barkhausen criteria, these two conditions are necessary, but not sufficient. While if total phase change along the loop equals to

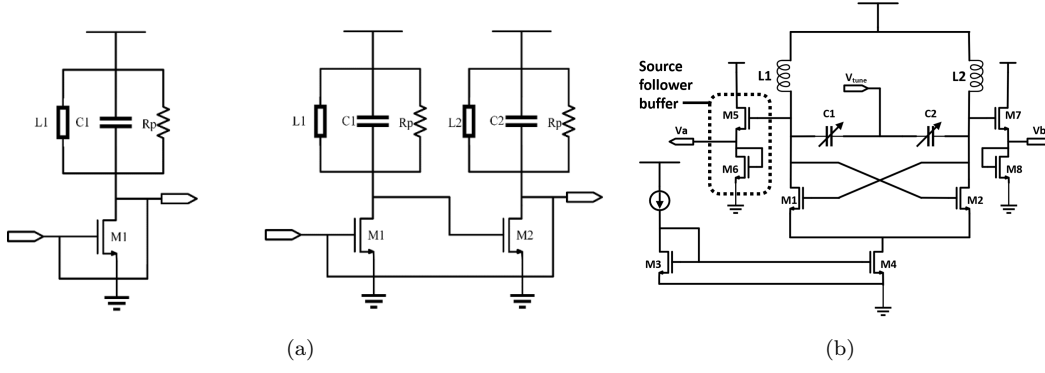


Figure 4.3: Circuit schematics of the proposed (a) tuned stage and (b) the proposed balanced VCO with cross-coupled NMOS pair and source follower buffers.

360 degrees and loop gain is larger than unity, this system will begin to oscillate at frequency  $\omega_0$  and  $\omega_0$  only.

To satisfy Barkhausen criteria while preventing positive feedback at DC region, a LC tank-based circuit was chosen as the frequency selector (resonator). While LC tank face very small impedance at both DC region, a common-source amplifier with LC tank load would have a DC phase change near zero and LC tank would have its phase approaching  $90^\circ$  because of the pure capacitance of the tank. In addition, LC tank would have its impedance equals to effective parallel resistance  $R_p$  at resonance frequency  $\omega_0$  while capacitive reactance and inductive reactance are cancelling each other. LC tank based circuit has similar features in high frequency region compared with DC region. The amplifier also has a frequency-dependent phase shift of  $90^\circ$  provided by its dominant pole at transistor drain terminal [103], hence this tuned stage could provide a total frequency-dependent phase change of  $180^\circ$ .

For a  $360^\circ$  total phase shift, simply cascade two tuned stage with a feedback connecting the output of second stage with the input of first stage so that total phase shift could be derived by:

$$\angle\beta H = \angle\beta H_{DC} + \angle\beta H_{freq} = 0 + 2 \times (90^\circ + 90^\circ) \quad (4.3)$$

where  $\angle\beta H_{DC}$  stands for DC phase shift and  $\angle\beta H_{freq}$  stands for frequency-depend phase shift. The oscillatory condition in terms of voltage gain:

$$|\beta H(j\omega_0)| = g_{m1}R_{p1} \times g_{m2}R_{p2} \geq 1 \quad (4.4)$$

where  $g_{m1}$  and  $g_{m2}$  are the transconductance of transistor  $M_1$  and  $M_2$  respectively. As long as Equation 4.1 to 4.4 are satisfied, this oscillator may oscillate at resonance.

In order to add immunity against load interference, two source followers as voltage buffer were added to the output of the cascade tuned-stage. Redrawing schematic of proposed



VCO, a differential output with two signals  $180^\circ$  out of phase is implemented shown in 4.3.

#### 4.1.2 On-Off Keying Modulator Design

For an ideal On-Off Keying (OOK) modulator, clock/data should be modulated onto the carrier signal according to its logic level. Carrier wave represented by a sinusoidal signal should occur on the output when the input baseband signal equals to logic “1”. On the contrary, no signal should be transmitted to the output when input is “0”, hence providing a high voltage difference as known as on-off isolation. In practice, because of the noise interference and transistor leakage, it is unavoidable to have some slight output during off stage, thereby reduce on-off isolation and increase the bit-error-rate (BER).

To simplify design while maintaining moderate on-off isolation, a MOS based sampling circuit is adopted to form the functionality of a switch. As the input clock signal is raised to  $V_{dd}$ , M1 starts to conduct hence the output tracks the input signal. When the input clock is low, M1 is shut down and the output would consequently hold the voltage of sample capacitor,  $C_{hold}$ . To get larger voltage swing with bigger on-off isolation at the output, a complementary device (PMOS) is adopted to form a transmission gate. This proposed architecture benefits from a better dynamic range as well as a faster response speed.

Assume body effect is neglected, the voltage output according to time could then be derived by:

$$V_{out} = V_{dd} - V_{th} - \frac{1}{0.5\mu_0 \frac{C_{ox}W_t}{C_{hold}L} + \frac{1}{V_{dd}-V_{th}}} \quad (4.5)$$

where  $C_{ox}$ ,  $W$  and  $L$  are effective transistor oxide capacitance, width and length respectively. As  $t$  gets larger, output voltage eventually approaches to  $V_{dd} - V_{th}$ .

Due to the cross-talk effect between clock signal to be transmitted and the output carrier signal, distortions might occur during the flipping edge of clock signal, thus reduce modulation quality. To alleviate distortions, a differential output scheme is proposed to enlarge the output swing and improve resistance against electromagnetic interference.

Considering the nature of differential signaling, an extra complementary clock is required as another control signal. To make the transistors of a transmission gate conducted and shut simultaneously to reduce risk and competition, a 1-to-2 clock buffer is constructed to provide the complementary clock pair shown in Figure 4.4. With the same delay through the buffer, complementary clock signal would turn off/on NMOS/PMOS at the same time, thus alleviate output ambiguity shown in Figure 4.5.

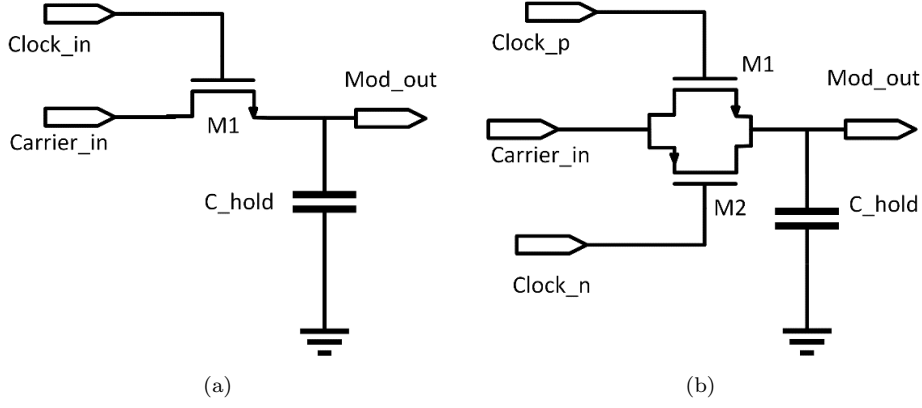


Figure 4.4: Typical structure of the (a) single and (b) complementary implementation of a MOS switch.

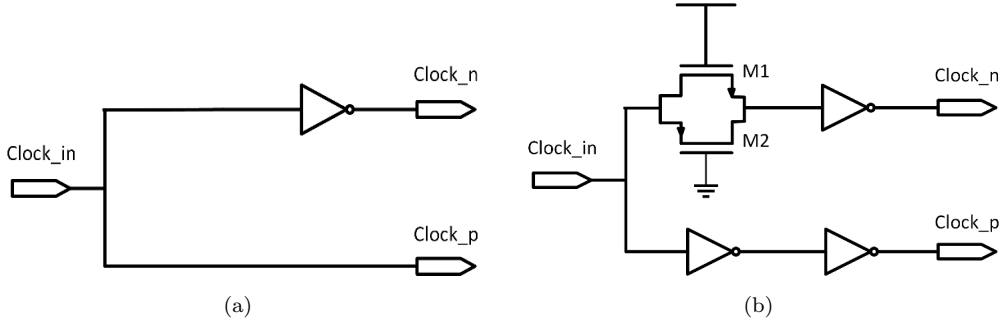


Figure 4.5: 1-to-2 Clock buffer with (a) unbalanced propagation delay and (b) balanced propagation delay.

The switching modulator design would consist two transmission gates presented in Figure 4.6 for differential signal input and output. Besides, transistor M5 is adopted to zero the difference of two outputs when input clock is “0”, hence further improve modulator on-off isolation. Out put voltage level of this modulator is given by:

$$V_{sw} = V_{dd} - V_{th} - \left( \frac{\mu_n C_{ox} W t}{2 C_{hold} L} + \frac{1}{V_{dd} - V_{th}} \right)^{-1} \quad (4.6)$$

where  $C_{ox}$ ,  $W$ ,  $L$  are the effective physical parameters of the CMOS switch. As the output sampling capacitor reaches its fully charged state, the output voltage will eventually approach  $V_{dd} - V_{th}$ .

This structure is adopted in our previous work [104] and [105] which shows an efficient switching characteristic. However, for carrier wave to be transmitted over 60 GHz, it is critical to block the entire signal out when the input clock signal is “1”. Since the input terminals are directly connected to the drain or source of the transistor pair according to Figure 4.6, any sub-threshold leakage would consequently transfer unwanted noise to subsequent power amplifier, therefore an enlarged noise would be transmitted out through transmitter antenna thus lead to potential receiver malfunction. As the core of

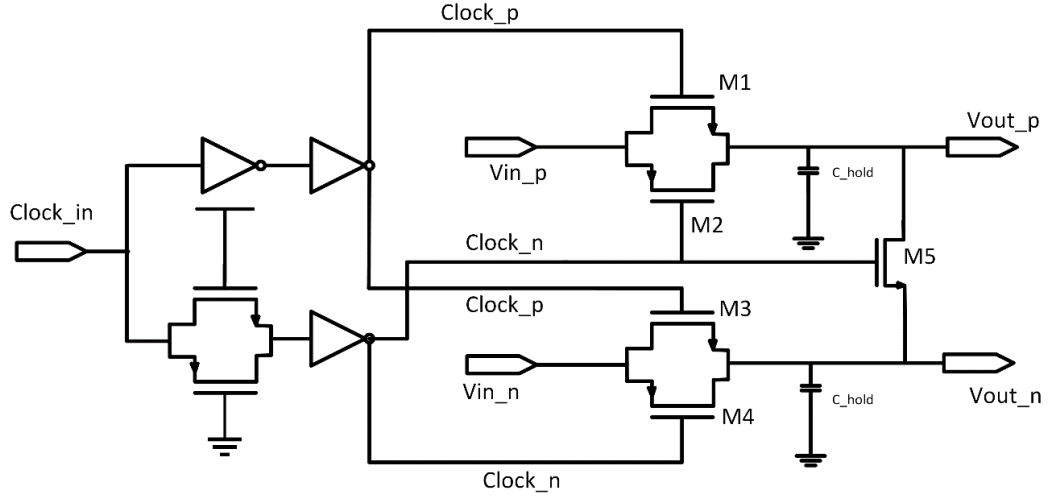


Figure 4.6: The schematic of MOS switch-based OOK modulator with differential signaling.

clock transmitter, it is of paramount importance to distinguish on state output from off state output while maintain good isolation so as to ensure an adequate signal-to-noise ratio (SNR) without generating possible error bits at receiver side. Hence, to address this challenge, a new novel architecture of on-off keying modulator is proposed.

Since differential signaling is adopted for common-mode noise rejection, taking advantage of the inherent characteristic that positive signal and negative signal are equal in magnitude and opposite in phase. Based on this property, any non-trivial signal/noise could be mitigated by its 180-degree out of phase counterpart. Therefore, by involving a differential signal pair onto output route, it is strait forward to eliminate the leakage signal by mathematically summing them together therefore cancelling each other out.

The new modulator structure is shown in Figure 4.7. As illustrated below, NMOS transistor M1 and M4 are always in “off” state. As the carrier attached to the source is over 60 GHz, leakage would always appear on the output terminal.  $V_a$  and  $V_b$  are differential carrier input respectively. Besides, depends on current state, input clock acts as a control signal applied on the gate of NMOS M2 and M3. Hence M2 and M3 start to conduct when the input clock is “1” and shut down when input clock is “0”. Therefore, M2 and M3 could be able to generate a leakage compensation signal during clock “off” state, which is 180-degree out of phase compared with the M4, M1 leakage. By attaching the drain of M2, M4 together as well as that of M1 and M3, leakage signals are self-cancelled at the cost of slight signal degradation during “on” state, therefor modulator on-off isolation and noise mitigation are significantly enhanced.

In addition, transistor M5 further works as a leakage rejection component, which exhibits a complimentary conducting timing compared with M2 and M3. Therefore, differential signal could be grounded to zero during “off” state thus further alleviate unwanted differential leakage.

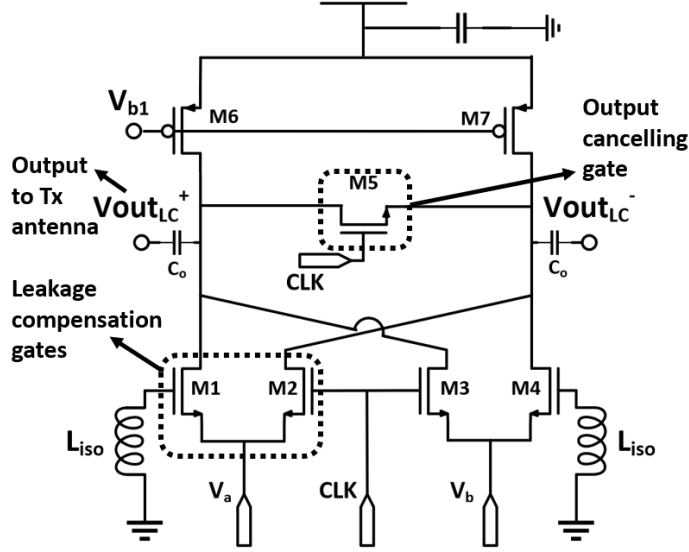


Figure 4.7: Schematic of the proposed OOK modulator with differential signaling and leakage compensation cross-coupling NMOS pairs.

Besides, load transistors M6, M7 and NMOS pair M2, M3 form a common gate amplifier structure which could enlarge signal level hence compensate the degradation due to phase cancelling effect, the new output signal voltage of effective half circuit could be given by:

$$V_{out_{lc}} = V_{dd} - \frac{1}{2} \mu_n C_{ox} \frac{W}{L} (V_b - V_{in} - V_{th})^2 R_d \quad (4.7)$$

where  $V_b$  is the voltage level of input clock signal,  $V_{in}$  is the input carrier signal and  $R_d$  is the load resistance. Compared with previous switching modulator, this new structure exhibits an output gain over original design given by:

$$\begin{aligned} Gain_{mod} &= \lim_{t \rightarrow \infty} \left[ \frac{V_{dd} - \frac{1}{2} \mu_n C_{ox} \frac{W}{L} (V_b - V_{in} - V_{th})^2 r_{op}}{V_{dd} - V_{th} - \left( \frac{\mu_n C_{ox} W t}{2 C_{hold} L} + \frac{1}{V_{dd} - V_{th}} \right)^{-1}} \right] \\ &= \mu_n C_{ox} \frac{W}{L} V_{th} r_{op} \end{aligned} \quad (4.8)$$

Hence, with this output gain, the leakage-compensated OOK modulator is capable of transmitting clock with higher clock frequency due to boosted on-off isolation compared to the previous switch-based design.

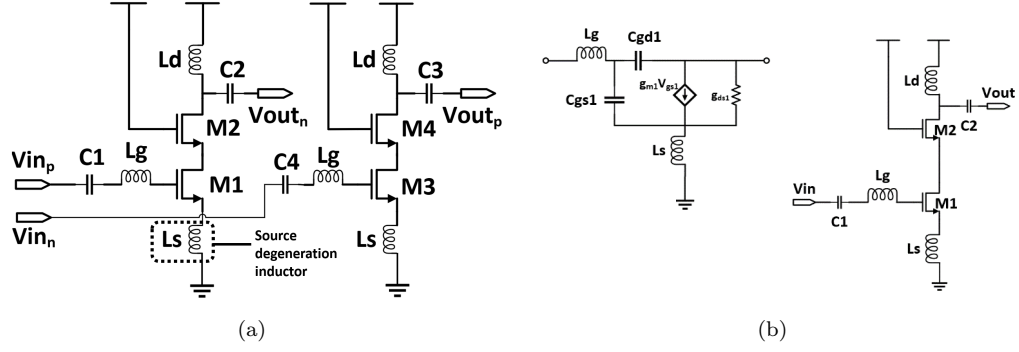


Figure 4.8: Schematics of the proposed (a) pseudo-differential Inductive degeneration LNA and (b) its high frequency single-stage small signal model.

## 4.2 Clock Receiver Design

### 4.2.1 Low Noise Amplifier Design

For the system simplicity, a single-stage low noise amplifier (LNA) is implemented with one cascoded amplifier. To be compatible with the subsequent circuits, a pseudo differential structure is adopted. Output matching network is designed to get maximum power transfer from the antenna. The input and output impedance of the amplifier is designed to be single-ended  $50\text{-}\Omega$  and differential  $100\text{-}\Omega$  seen into the input terminal at desired 69 GHz band for the purpose of impedance matching. Also, it helps to get low noise figure (NF) as LNA is the very first stage of the receiver circuit front-end.

The proposed LNA composes a single-stage cascoded structure to achieve high forward gain shown in Figure 4.8. A small inductor is adopted between input stage and cascode stage to resonate out the parasitic capacitance inside the internal node. The basic functionality of the LNA is to provide a large dynamic range with adequate linearity. In addition, shunt-peaking technique is implemented to improve LNA bandwidth and dynamic performance. However, single common source amplifier shows a relatively low isolation between output terminal and input terminal, thus makes it difficult to match the input/output impedance to predecessor stage and subsequent circuits, which requires fine tuning for matching network calibration therefore increases overall design complexity. In addition, low-isolation of common source structure means more noise sensitive. Any wiring parasitic would essentially affect impedance characteristics of the entire LNA, hence increase the potential of power reflection and receiver malfunction.

Therefore, a cascoded structure is adopted for higher isolation and robustness compared with original common-source stage. Given the small signal model of a cascode stage as shown in Figure 3.14, it is straight forward to get that it consists of a common-source stage(CS) and a common-gate stage(CG). Split the circuit into two portions, applying KVL, KCL to both CS and CG amplifier, it could be derived that for stage one (CS)

the small signal gain is:

$$Av_{cs} = -g_{m1}(R_{dx}||r_{o1}) \quad (4.9)$$

where  $g_{m1}$  and  $r_{o1}$  are the trans-conductance output impedance of NMOS M1 respectively,  $R_{dx}$  is the load impedance of M1, which is the input impedance of stage two (CG). Applying a voltage source with input current  $i_x$  onto the input terminal of CG, it is clear that:

$$v_x = v_{r_o} + i_x(R_d||R_L), \quad (4.10)$$

$$v_{r_o} = (i_x - g_{m2}v_x)r_{o2} \quad (4.11)$$

therefore, the input impedance of CG part could then be derived by:

$$R_{dx} = \frac{v_x}{i_x} = \frac{r_{o2} + R_d||R_L}{1 + g_{m2}r_{o2}} \quad (4.12)$$

where  $R_d$  and  $R_L$  are intrinsic load and outside load impedance respectively. Besides, the forward gain of CG stage could be given by:

$$Av_{cg} = \frac{v_{out}}{v_{in}} = \frac{i_x(R_d||R_L)}{v_x} = \frac{(1 + g_{m2}r_{o2})(R_d||R_L)}{r_{o2} + (R_d||R_L)} \quad (4.13)$$

For the cascode stage, effectively the total output gain could be regarded as the product of CS gain and CG gain, thus the overall voltage gain with resistive load could be given by:

$$Av_{total} = Av_{cs} \cdot Av_{cg} = \left\{ \left[ -g_{m1}(r_{o1}||\frac{r_{o2} + R_d||R_L}{1 + g_{m2}r_{o2}}) \right] \cdot \frac{(1 + g_{m2}r_{o2})(R_d||R_L)}{r_{o2} + (R_d||R_L)} \right\} \quad (4.14)$$

Compared with conventional CS stage, gain over forward transmission coefficient is derived by:

$$Gain_{fc} = \frac{Av_{total}}{Ac_{cs}} \approx g_{m2}r_{o2}||R_{dx} \quad (4.15)$$

Besides, assume that the LNA is perfectly matched to the impedance of input source such as voltage transmission line or antenna. The matching network could also provide an extra voltage gain to further enhance input signal level. The voltage across gate and source of the input transistor  $V_{gs}$  with respect to the input radio frequency signal is given by:

$$Gain_{matching} = \frac{v_{gs}}{v_{in}} = \frac{\sqrt{LC_{gg}^{-1}}}{R_{in}}, \quad (4.16)$$

$$Gain_{total} = Gain_{fc} \cdot Gain_{matching} \quad (4.17)$$

where  $L$  is the summation of inductance exists inside matching network and degeneration components,  $C_{gg}$  is the capacitance seen at the gate of input transistor,  $R_{in}$  is the input impedance under the assumption of perfect matching. Hence, cascode amplifier could essentially provide higher gain than conventional CS stage, thus becomes naturally suitable for LNA design inside this wireless CDN architecture.

Besides, LNA is designed to be the first stage of a clock receiver, thus decayed signal can pass through it to get enough energy level for subsequent circuit to recover the original clock without bringing too much noise into signal route. As one of the most important aspects of LNA design, for a typical two-port network such as amplifier, attenuator, etc., the noise figure can be modeled as:

$$NF = \frac{\frac{S}{N_{input}}}{\frac{S}{N_{output}}} \quad (4.18)$$

which generally represent the input signal-to-noise ratio (SNR) to output signal-to-noise ratio. Treating a typical RF receiver link as a black box [79], total output noise could be additive, and the receiver could be modeled as a noisy resistor. For system intrinsic noise, we could assume it follows the noise density function over frequency that:

$$N_{output} = N_{input} + N_{intrinsic}, \quad (4.19)$$

$$P_{system} = E^2 = 4kTR_{eff} \cdot \Delta f \quad (4.20)$$

where  $E$  is the RMS noise signal in volts (V),  $k$  is Boltzmann's constant,  $T$  is temperature in Kelvin (K) and  $R_{eff}$  is the effective resistance in ohms ( $\Omega$ ) and  $\Delta f$  is the frequency offset.

Assume the radio frequency signal is passing through a unit-gain buffer (output signal maintains the same energy level as input signal does), therefore, receiver noise figure could then be rewritten as:

$$NF = 10 \log \left( \frac{\frac{E^2}{4kTR_{in}\Delta f}}{\frac{E^2}{4kT(R_{in}+R_{eff})\Delta f}} \right) = 10 \log \left( 1 + \frac{R_{eff}}{R_{in}} \right) \quad (4.21)$$

which is independent of noise temperature and frequency/bandwidth. In terms of the input noise figure. For a typical input NMOS transistor [103], the effective resistance could then be modeled as:

$$R_{eff} = \frac{8}{3} \cdot \frac{kT\gamma}{g_m} \quad (4.22)$$

where the  $\gamma$  is a noise constant around 3 for 65 nm process and  $g_m$  is the transconductance of input transistor. Since the matching network has provided a matching gain to the overall voltage gain, to maintain the energy in and out of matching network at the

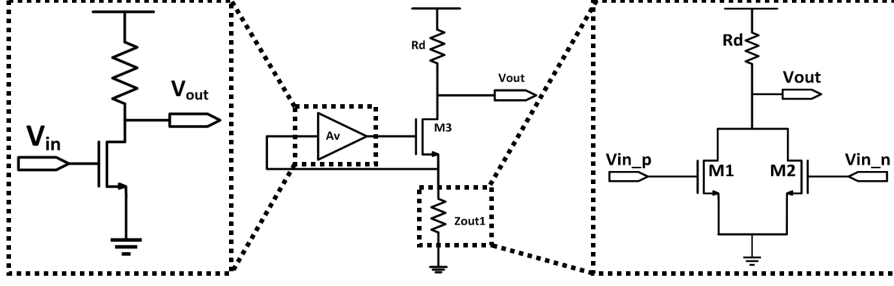


Figure 4.9: Schematic of the proposed gain-boosting load, with a conventional NMOS rectifier to rectify the input radio frequency signals.

same level, effective resistance is reduced by:

$$R'_{eff} = \frac{R_{eff}}{Gain_{matching}^2} \quad (4.23)$$

Therefore, the overall noise figure can be quantified as:

$$NF_{total} = 10 \log \cdot \left( 1 + \frac{\frac{8}{3} \cdot \frac{kT\gamma}{g_m}}{R_{in} \cdot Gain_{matching}^2} \right) \quad (4.24)$$

which can be minimised by adjusting transistor size and matching network parameters. Therefore, it's easy to find a balance point between the adequate noise figure and enough forward gain, considering overall design requirements.

#### 4.2.2 On-Off Keying Demodulator Design

As a conventional way of demodulating signals, coherent demodulator requires frequency mixer to down-convert carrier frequency to an intermediate frequency and then retrieve the original baseband signal accordingly. As carrier frequency (i.e. 69 GHz) is required to be removed from the modulated signal, an oscillator with operating frequency at 69 GHz is also necessary inside a clock receiver. Thus receiver array on-chip would consume a large amount of DC power due to high-frequency operation.

On the contrary, non-coherent demodulator would produce a one-step down convert which could alleviate overall power. Although non-coherent demodulator was originally considered lossy and inefficient [46], however, the reduction of local oscillator inside clock receivers would consequently lead to a desirable low-power feature, which is an essential criterion for clock receiver design.

In the light of power efficiency and output swing level, proposed non-coherent OOK demodulator would first rectify input RF signal to form the basic structure of the baseband signal. Similar to [106] and [107], a class AB-biased NMOS differential pair is adopted to rectify input signal accordingly. Differential pair output then feed to a gain-boosting



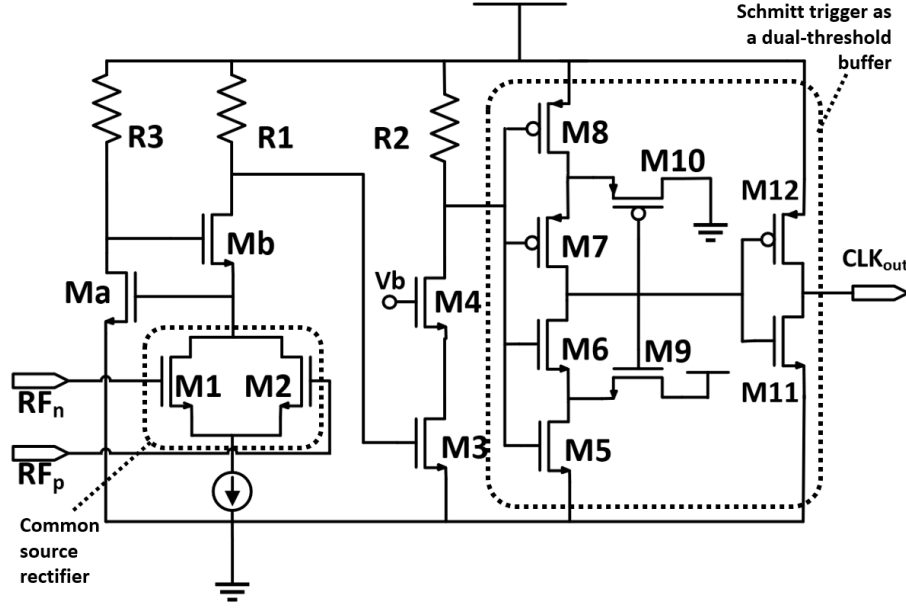


Figure 4.10: Schematic of the proposed Proposed OOK demodulator with gain-boosting technique and noise-compensation output buffer.

cascode device [103], [107] so as to enlarge the output impedance shown in Figure 4.9. Suppose the output impedance of NMOS differential pair approximately equals to:

$$Z_{out1} = r_{o1} || r_{o2} || R_D \quad (4.25)$$

where  $r_{o1}$  and  $r_{o2}$  are the on resistance of M1 and M2 respectively and  $R_D$  is the load resistance. The overall output impedance including gain-boosting device is then calculated by:

$$Z_{out_{total}} = Av g_{m3} r_{o3} Z_{out1} \quad (4.26)$$

where  $Av$  is the voltage gain of the amplifier inside feedback loop,  $g_{m3}$  and  $r_{o3}$  are the trans-conductance and the on-resistance of M3, respectively. Hence the output impedance of the overall rectifier could be substantially enlarged with a factor of  $Av g_{m3} r_{o3}$ , thus enhancing the overall gain.

The proposed OOK demodulator is shown in Figure 4.10, which first tracks the envelope of the input RF signal and then demodulates the signal as follows. Under the condition of class-AB biasing, the rectifier core first transfers the differential input into an inverted single-ended envelope. Besides, the rectifier output terminal (drain of M1 and M2) is exposed to the input capacitance of the downstream buffer, hence forming a low-pass filter preventing high-frequency noise from passing to subsequent circuits. The weak envelope then gets enlarged and inverted by a single-stage baseband amplifier. Assuming that all transistors are working in saturation mode, to provide enough driving strength

while reducing noise interference generated by nearby digital logic, a new Schmitt trigger-inverter buffer amplifier is designed to have a dual-switching threshold, therefore yields a better noise tolerance  $\delta$  around [108]:

$$\delta = \frac{V_{dd} + \sqrt{W_5/W_9} \cdot V_{tn}}{\sqrt{W_5/W_9} + 1} - \frac{\sqrt{W_8/W_{10}}(V_{dd} - |V_{tp}|)}{\sqrt{W_8/W_{10}} + 1} \quad (4.27)$$

where  $W_5$ ,  $W_8$ ,  $W_9$ ,  $W_{10}$  are the width of transistors in the proposed demodulator respectively according to Figure 4.10.  $V_{tn}$  and  $V_{tp}$  are the threshold voltage for NMOS and PMOS devices used in the buffer amplifier respectively. This property helps to provide a more robust output, therefore enhancing overall signal-to-noise ratio (SNR) performance significantly.

### 4.3 On-chip Antenna Design

The on-chip antenna is of paramount importance when considering wireless clock distribution. Due to the fact that the size of the antenna is pseudo-inversely proportional to the frequency of the transmitted signal, in order to maintain the area occupation at the minimum level, a 69 GHz carrier is generated and adopted with a wavelength of around 4mm. considering a conventional half-wavelength dipole antenna with a total length about 2.3mm excluding the antenna feedline, the size is still not practical for on-chip clock distribution [24], [109]. Hence, to address this problem, an area-efficient antenna geometry by wire-snaking is adopted in this design [110], [111].

The proposed Antenna shows a structure of meandering metal line with one antenna arm, hence is therefore formed a meandering monopole antenna (MMA). The proposed MMA has a dimension with 300  $\mu\text{m}$  total length and around 100  $\mu\text{m}$  width. Hence the overall area occupation is to the order of  $10^4 \mu\text{m}^2$ , which is two orders smaller than the conventional monopole/dipole with the effective area to the order of  $10^6 \mu\text{m}^2$ , which is naturally suitable for on-chip clock distribution.

The proposed MMA is implemented using the top Metal layer in any standard CMOS process shown in Figure 4.11 and 4.12. The vertical structure shows that the antenna is constructed over a device-free area, which aims to reduce the mutual interference between EM radiation and nearby passive devices and analog circuits, as they are fragile and sensitive to any noise, especially for on-chip inductors. If the magnetic flux through the coil inductor is affected by the EM wave radiated out from antenna, the induced current would therefore be generated and subsequently cause circuits malfunction. For the overall simulation simplicity, under the assumption of omnidirectional propagation, a single transmitter-receiver pair is constructed. Transmitted EM wave will traverse through multiple routes, including a 300  $\mu\text{m}$  silicon substrate, free-space and passivation/insulator layer as a surface wave. To mitigate the effect of multi-path propagation,

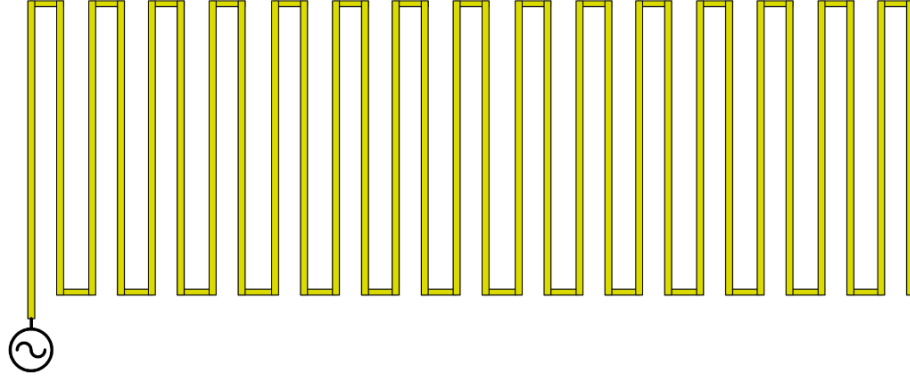


Figure 4.11: Top view of the proposed meander monopole antenna (MMA) structure with top copper layer.

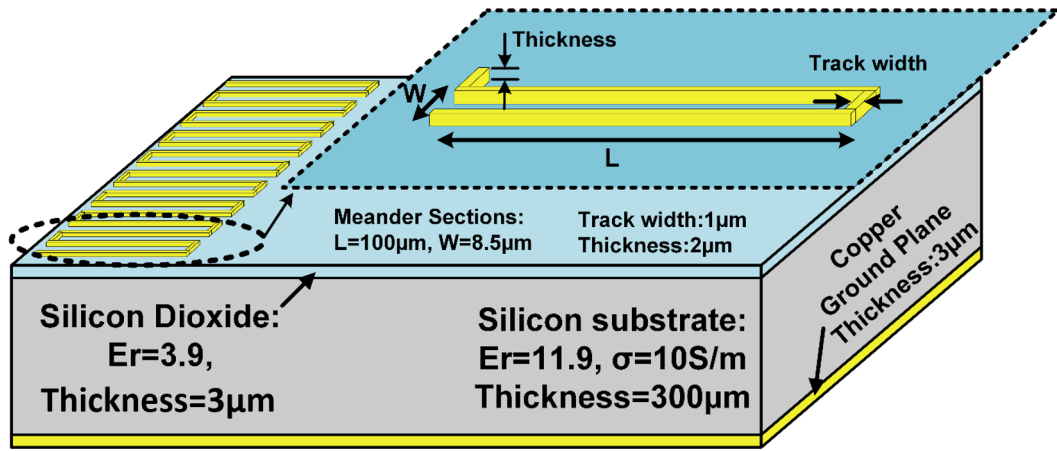


Figure 4.12: Proposed meandering monopole antenna (MMA) structure with EM simulation model setup (not to scale).

hence increase the potential of recovered clock jitter, metal shielding is placed between the transmitter and receiver antenna at the cost of slight degradation of antenna gain.

However, the MMA antenna requires extra balun to convert the fully-differential signal to a single-ended version as the modulator output does to suppress the nearby common mode digital noise. This coil based balun would consequently cause more area-occupation [46]. Hence, we proposed another antenna design which is a balanced meandering dipole antenna (MDA). The proposed MDA contains two meandering dipole arms while each of the arms consists of 6 meander sections. Each dipole arm is around  $440 \mu\text{m}$  in length and  $10 \mu\text{m}$  in width. Hence the total antenna physical length is around  $900 \mu\text{m}$  and the effective area still maintains to the order of  $10^4 \mu\text{m}^2$ , shown in Figure 4.13 and 4.14.

Besides, since the total effective radiation length of MDA is almost two times longer than that of MMA, the radiated energy with same modulator output is higher, hence the total radiation efficiency is increased. To sum up, both MMA and MDA antenna could be adopted as the on-chip antenna for wireless clock distribution, MMA has a

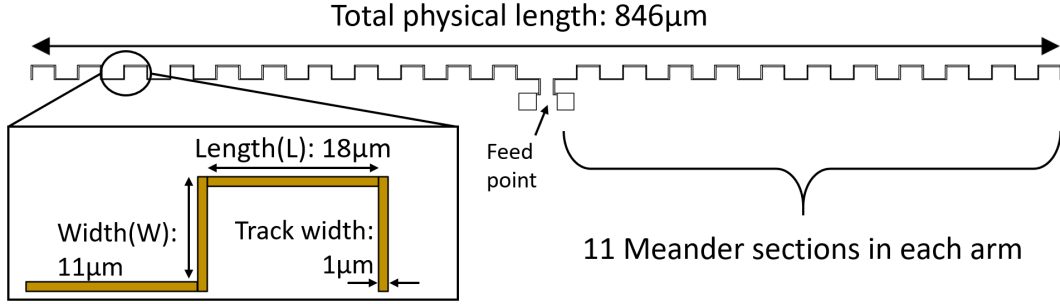


Figure 4.13: Top view of the proposed meander dipole antenna (MDA) structure with top copper layer.

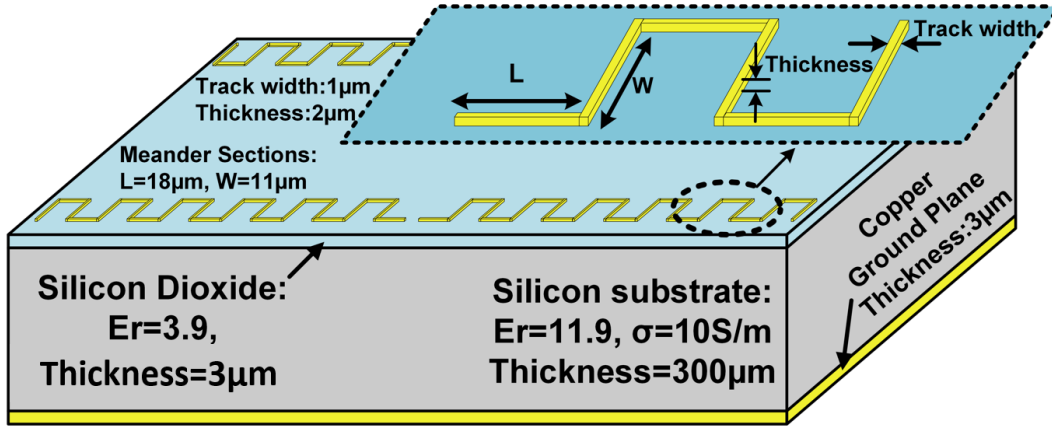


Figure 4.14: Proposed meandering dipole antenna (MDA) structure with EM simulation model setup (not to scale). ©2021 IEEE

smaller layout but less efficiency; On the contrary, MDA has a larger area occupation but higher radiation efficiency. Hence, these two on-chip antennas could be adopted depends on different application and system requirements. For example, for a 4-receiver architecture, both Tx and Rx could use MDA as on-chip antenna since the area constrain is not the highest concern and therefore the utilization of same antenna simply overall transceiver design. However, for a larger system requires more fanout number such as a 16-receiver architecture, MDA antenna could be used as the Tx antenna for its high efficiency. Meanwhile, MMA antenna could be utilised as the receiver antenna for its compact layout and both these two on-chip antennas shows the characteristic of horizontal polarization [112]. Also, the dark silicon phenomenon limits the number of cores to execute simultaneously, hence antenna sharing among multiple Rx becomes a possible option to mitigate the impact of area overhead of using global wireless CDN.

Last but not least, thermal control is an important task for on-chip antenna design in future many-core systems with variable switching rate and high density. The antenna itself is a passive device and hence does not generate a large amount of heat [113] compared with circuits, on-chip antennas can be used both as a EM wave radiator and a heat sink. Since the antenna is implemented using top metals and spanned over a

relatively large area, it is possible to gain some beneficial effect for spreading the heat generated by the circuit below [114].

To mitigate the impact of the presence of dark silicon phenomenon, thermal conductive materials are often added to the face/back of a designated core to cool down the heating circuit. This metallic object serves as a reflector which will have an impact on the antenna performance. Different studies have suggested that a near by metal object will shield the EM wave from certain propagation paths, therefore affecting overall radiation efficiency [7], [24]. The allocation of the heatsinks need to be carefully designed such that the performance degradation is minimised. [24] raised a solution of adopting a “metal-free” area above/below the integrated antenna to alleviate the performance attenuation, which is a contradiction when considering thermal control for many-core application. Other solutions include tuning the SNR of the radiated EM signal to a higher level such that the energy decay can be compensated, but this also leads towards an extra cost of higher Tx power dissipation. To sum up, finding the balance point between satisfying heating constraint and the antenna performance is an interesting and important topic for future study.

## 4.4 Summary

As a conclusion, both MMA and MDA exhibits an area-friendly layout and rather low directivity, which is naturally suitable for on-chip wireless clock distribution network. Developing CMOS technology makes it possible to mount an antenna over active die area. Although there are still some practical issues to be considered about such as antenna feeding, matching and lossy silicon substrate, MMA and MDA appear to be a promising solution to on-chip wireless communication.

The proposed design covers complete transmitter (Tx) components and simplified receiver (Rx) front end. The proposed OOK scheme helps to incorporate simplicity as well as low power consumption. In addition, differential signaling is adopted to alleviate common mode noise interference, thus produce better SNR. Besides, a gain boosting technique is utilised to form a non-coherent OOK demodulator with higher conversion gain and output swing which result in a low-power implementation.

However, there are certain limitations exist in the proposed design. First of all, the proposed design will produce an area overhead of around  $5.3 \text{ mm}^2$  for distributed Rx planning and  $1.6 \text{ mm}^2$  for concentrated Rx planning including the proposed MDA antenna, respectively. Hence, for applications which require small and compact chip layout, the proposed design might not be the straight forward choice. Besides, since the clock Tx and Rx will bring a power overhead, for applications with small synchronisation area, the conventional approach can provide better energy efficiency since the overall CDN wire length is limited. Another limitation is that the analog circuit, including the integrated

antenna, will suffer from the nearby interference. Although some approaches such as balanced antenna structure and differential signaling are adopted in the proposed design which can help to mitigate this reliability impact, the improvement of noise immunity of the clock TRx remains to be an essential task for future study. Variation in clock Rx becomes vital when considering different corners referring back to Chapter 3, and hence a robust circuit structure is also necessary for future improvements. Detailed discussion about the limitations will be given in Chapter 7. The proposed design has the potential to further improve system performance with moderate synchronisation area. A comprehensive test and simulation for global wireless clock distribution will be demonstrated in the following Chapter 5.



## Chapter 5

# Experimental Setup and Evaluations of the Proposed Hybrid CDN

The proposed clock Tx and Rx introduced previously are analog circuits, and hence, there exists some process dependencies such that the fabrication process may impact the overall performance. Although digital logic benefits from the ever-decreasing gate delay, certain attributes of the analog circuit are no longer available using the advanced process.

For example, voltage scaling is the first issue to be considered about when integrating analog circuit using digital processes. The voltage level of power supply has decreased from a common voltage of 1.8 V in 180  $\mu\text{m}$  process to around 0.7 V in some of the cells in 7 nm process, which is only less than 40% of the original value. The necessary voltage swing in analog circuit design has therefore been limited severely as there is no enough voltage headroom for the circuit operation. Furthermore, signal amplitude shrinks together with the voltage level, however, noise does not. The interference from nearby digital logics will generate an aggravate impact on the reliability of the sensitive analog components using the advanced process compared with the developed processes such as 65 nm or 45 nm nodes. Besides, as the threshold voltage now becomes around hundreds of millivolts, the analog device may always be in a weak “off” state and generate higher leakage current through the channel, therefore leading towards unacceptable leakage power [115].

Another issue to be considered about is the device mismatch, especially related to the threshold mismatch. According to [116] and [117], shrinking device size will generate an increasing variation/mismatch in device threshold, therefore reducing overall circuit performance. Also, the impact of parasitic and layout-dependant effects will get worse in each new technology nodes. Some solutions [118] have been raised such as using



ADC/DAC interfaces to convert analog signals to digital signals to mitigate the AMS design challenge in advanced nodes. Others [115] have suggested a functionality division, such as using an on-package die in a system in package (SiP) system, so as to increase yield and to be more AMS-friendly.

Some commercialised interposer technology such as [119] uses a developed process, such as 65 nm manufacturing process, to fabricate the silicon interposer, more advanced digital dies can be mounted on top of the silicon carrier, which can save chip area, cost and power. An example of using a 2.5-D stacking system will be given in Chapter 6, which helps to provide more flexibility by partitioning digital and analog devices when dealing with AMS design [120].

Hence, although the proposed design suffers from the above challenges, certain solutions can be applied to mitigate the impact of these drawbacks. This Chapter will focus on the evaluation of the proposed clock TRx circuit. To be more specific, simulations of the proposed circuits given in Chapter 4 based on TSMC 65 nm process will be analysed to verify the performance of the global wireless CDN.

To date, comprehensive evaluations of the proposed wireless CDN are implemented by using Cadence Virtuoso Analog Design Environment (ADE). An input clock signal with 50% nominal duty cycle is first feed into the clock Tx and then captured and recovered by the clock Rx. Input clock frequency of the clock TRx circuit can support a full-swing clock output without further amplification. A 2.5 GHz testing clock will be chosen as a nominal frequency through the remaining of this chapter.

Local wired CDN is evaluated using the proposed mathematical models and algorithms given in Chapter 3. The recovered global clock signals at the output of the clock receivers are sampled and used as the discrete input vectors of the proposed local CDN design algorithms implemented by using MATLAB. As the transfer function of local CDN contains hyperbolic functions which are difficult to directly modelled by analytical tools, thus they are first expanded by using Taylor series. By truncating the high-order poles and zeros of the transfer function, the overall local model can be much simplified and represented by its dominant poles and zeros only.

## 5.1 Results of the Proposed Antenna for Hybrid CDN

As the proposed CDN requires a compact, reliable and high-performance antenna, meandering dipole antenna (MDA) in the last chapter is adopted for on-chip wireless communication for its simplicity and area-efficiency. According to Chapter 4, a simulation model is first constructed under practical constraints. A single transmitter-receiver pair is proposed under the assumption of an omnidirectional propagation pattern. Since the propagation media is assumed to be the lossy silicon substrate, antenna efficiency is

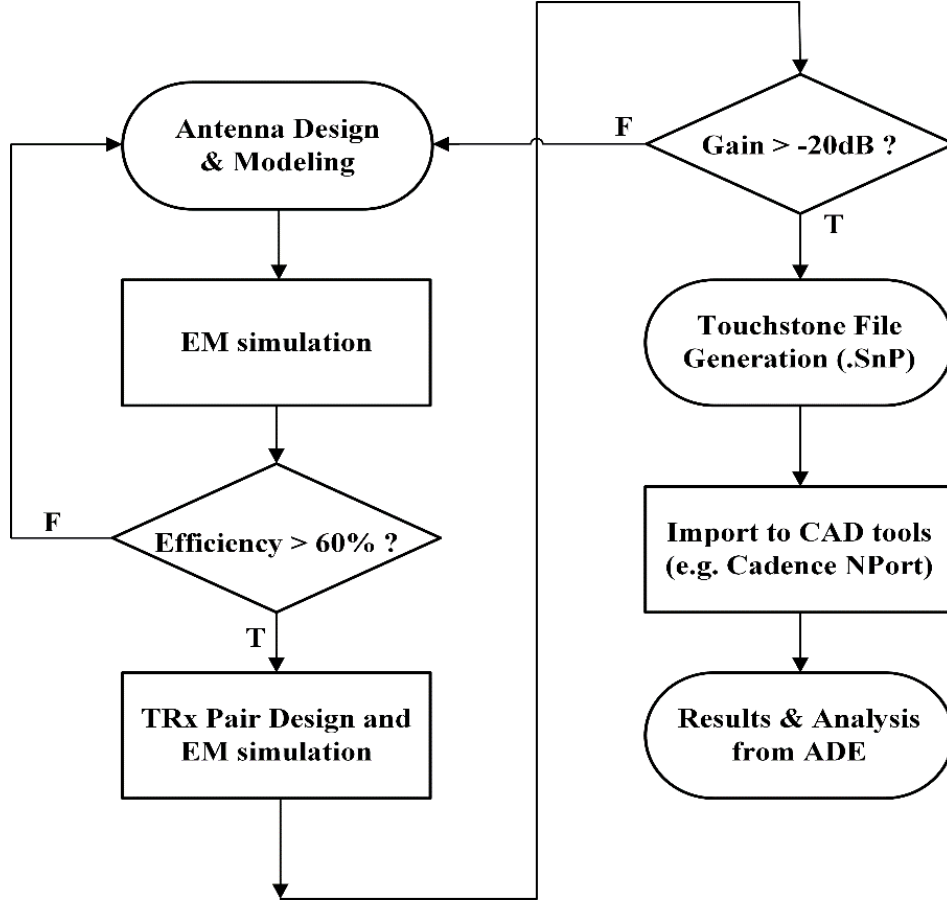


Figure 5.1: Complete antenna-circuits co-simulation flow diagram for global wireless CDN verification.

one of the most important considerations to determine how much radio frequency energy would be transmitted, which would consequently decide the performance of entire system. Hence to address that, a simple framework for antenna-circuit co-design is implemented for high radiation efficiency and relatively low directivity, according to Figure 5.1.

If a designated antenna has meet the requirements in terms of efficiency, directivity, gain, etc., a Touchstone file would be generated from EM simulation software such as CST Microwave Studio and Ansys HFSS by solving Maxwell Equations, containing all the ports information and system impulse response at the frequency band of interest. Taking advantage of Cadence NPort device in the default AnalogLib, given a transmitter output attached to this antenna, channel loss/interference could then be modelled as a black box with all necessary information for wireless communication. Similar characteristics could be applied on the connection between clock receivers and receiving antennas. Finally, received and recovered clock could then be observed at the receiver output terminals.

To provide robust clock quality with on-chip antennas, there are certain factors that need to be taken care of. As the density of the CDN circuits keeps growing, electromagnetic

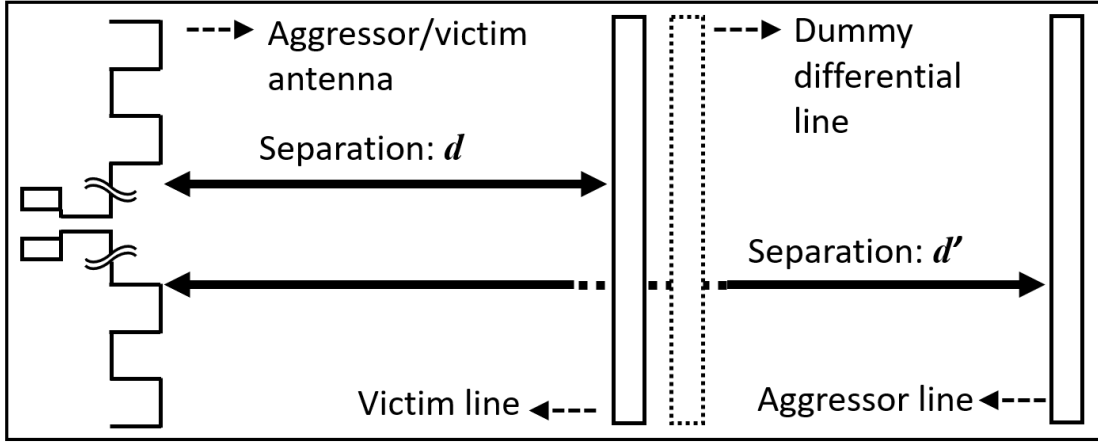
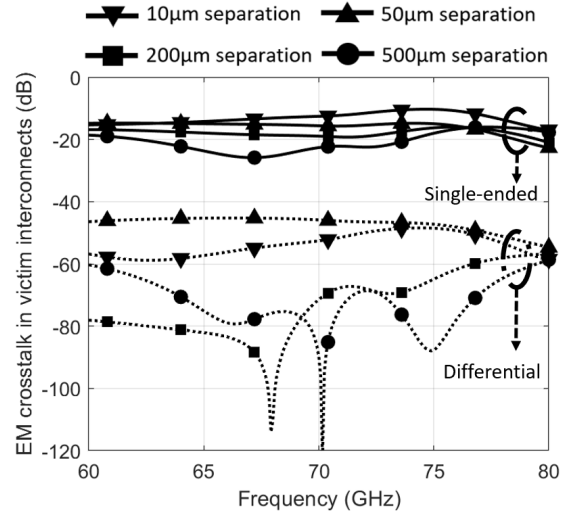


Figure 5.2: Top view of the test structure for the study of EM crosstalk between TRx antenna and the nearby interconnects with different separation distance. The antenna and nearby lines act as the victim of the crosstalk in the experiment, respectively. ©2021 IEEE

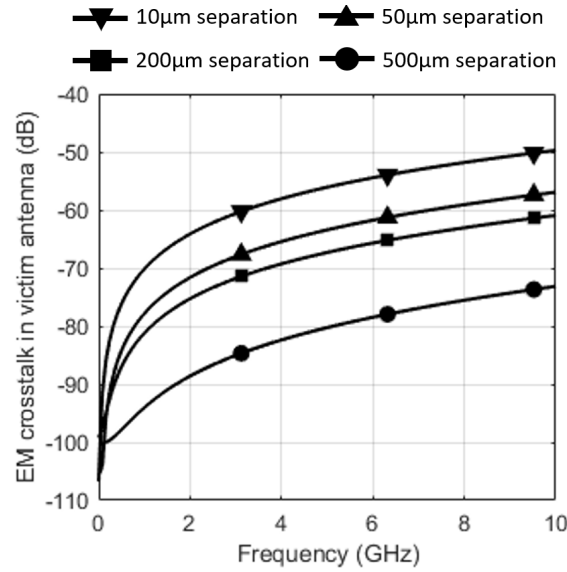
interference (EM crosstalk) will have an increasing impact on the overall clock quality. In this section, crosstalks on interconnects and circuits will be evaluated inside the active silicon interposer. The layout of the experiment is presented in Figure 5.2, the nearby interconnect and circuits will be modelled as victim/aggressor lines, with different separation distance between the test antenna and nearby interconnect ranging from  $10\ \mu\text{m}$  to  $500\ \mu\text{m}$ .

First of all, the Tx/Rx antenna serves as the aggressor and nearby interconnects are its victims. The crosstalk on the victims is measured in decibel and shown in Figure 5.3. The length of all victim lines is set to be  $300\ \mu\text{m}$ , and they share the same track width and height of the integrated antenna. As shown in the figure, the peak of the crosstalk occurs between 72-76 GHz, which is around -10 dB for the closest aggressor/victim pair with  $10\ \mu\text{m}$  separation. With the increase of signal frequency, a crosstalk drop will occur with the frequency larger than 76 GHz. Besides, the separation between the antenna and nearby interconnects is also a key factor which will impact crosstalks. There is an averagely 13.1 dB crosstalk drop between the  $10\ \mu\text{m}$  and  $500\ \mu\text{m}$  separation curves, according to Figure 5.3.

The  $10\ \mu\text{m}$  separation is set to be the worst-case as near 10% of the wireless power can be transferred to the nearby interconnects/circuits, thereby degrading SNR of the received clock signal. To alleviate the impact of EM crosstalk, differential signaling can be implemented for common mode noise rejection. According to Figure 5.2, a dummy element is implemented beside each of the nearby lines, with a separation of  $2.5\ \mu\text{m}$  for differential wires. As shown in Figure 5.3(a), by adopting this method, the crosstalk will have a significant drop averagely around 35 dB. The new worst-case has a peak of around -44 dB for  $50\ \mu\text{m}$  separation, which is 24 dB smaller than the -20 dB threshold [121]. Hence, it indicates a significant crosstalk rejection that the nearby interconnects/circuits



(a)



(b)

Figure 5.3: EM crosstalks versus the frequency of interest in (a) victim interconnects and (b) victim antenna, measured in  $S_{21}$ . The differential signaling can reduce the common-mode noise by 35 dB averagely, in the victim lines.

will only pick up 0.004% of the radiated wireless power at most, at the cost of increased interconnect usage. Hence the implementation of the clock Tx and Rx front-end circuit are all in differential mode, in order to reduce the impact of high-frequency crosstalk.

For the other scenario where the antenna is the victim suffering from an aggressor line, the crosstalks with local clock frequency up to 10 GHz are under -50 dB for all separation distances in Figure 5.3(b), which shows that the wired part has limited impact on the wireless transmission. Hence the EM crosstalk generated by aggressor lines is considered to be negligible.

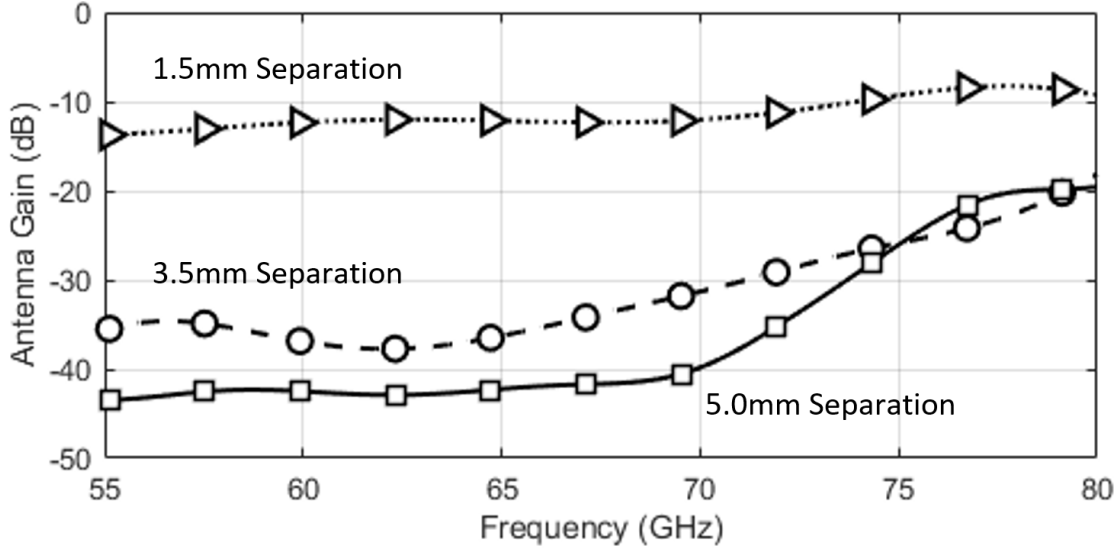


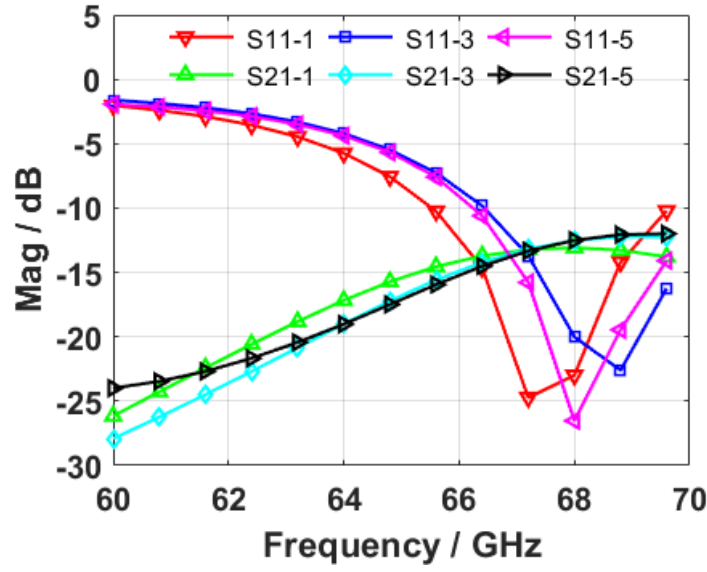
Figure 5.4: Antenna gain under matched condition with 3 separation distances ranging from 1.5 mm to 5 mm, which shows it's suitable for the proposed hybrid CDN. ©2021 IEEE

The proposed system adopts a new balanced structure for both circuits and antenna to increase common-mode noise rejection. The proposed meander dipole antenna (MDA) is shown in Figure 4.14. Although requiring another dipole arm, the size of the proposed antenna remains in the order of  $10^4 \mu\text{m}^2$  for thinner arm width around only  $10 \mu\text{m}$ , which shows a promising potential of being used in our hybrid architecture for its compact layout. Our proposed MDA contains 11 meander sections with a total physical length of  $840 \mu\text{m}$ . Since the physical length of the antenna has been increased and both of the arms are referred to each other as ground [122, 109], radiation efficiency has been boosted for higher energy efficiency compared to MMA. Each meander section has a track width of  $1 \mu\text{m}$  and a length-to-width ratio of around 1.68 to 1. A 3D EM simulation model has been setup using both CST Microwave Studio and Ansys HFSS. Considering a general test model of around  $300 \mu\text{m}$  thick silicon substrate with 6-layer metal stack-up and  $3 \mu\text{m}$  thick ground plane, our proposed MDA is stimulated by a  $50 \Omega$  horizontal lumped port and implemented with top metal M6 with  $2 \mu\text{m}$  wire thickness. For a matched antenna pair, according to [123], the gain of the antenna can be evaluated by the equation:

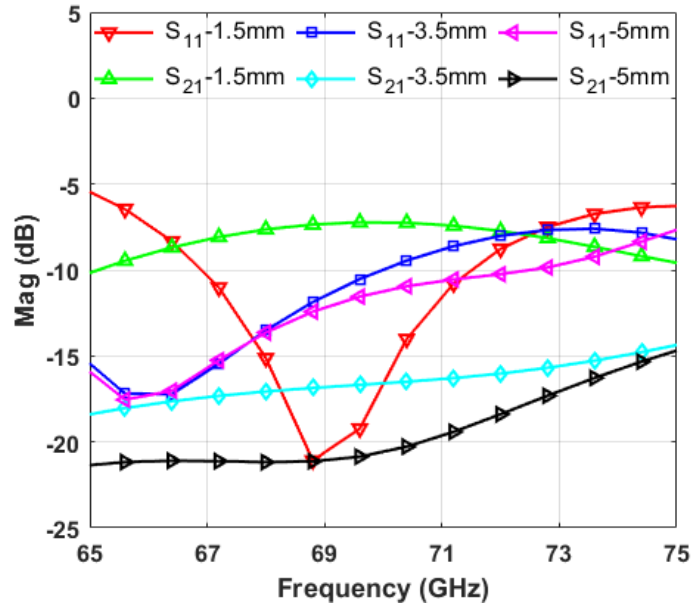
$$Gain_a = \frac{|S_{21}|^2}{(1 - |S_{11}|^2) \cdot (1 - |S_{22}|^2)} \quad (5.1)$$

where  $|S_{11}|$ ,  $|S_{22}|$  and  $|S_{21}|$  represent the magnitude of the s parameters, respectively. The antenna gain is then given in 5.4

According to Figure 5.5, given a  $50 \Omega$  horizontal lumped-port excitation, two proposed on-chip antenna both shows a good return and insertion loss at the frequency band of interest. Proposed MMA antenna exhibits a  $S_{11}$  between -12 dB and -23 dB depending on the use-case with different communication distance ranging from 1 mm to 5 mm.



(a)



(b)

Figure 5.5: EM simulation in terms of S-parameter results for (a) the proposed meander monopole antenna (MMA) antenna and (b) the proposed meander dipole antenna (MDA) antenna, with 3 separation distances ranging from 1.5 mm to 5 mm.

Besides, the insertion loss from 1 mm separation to 5 mm separation is ranging from -13 dB to -17 dB and with a total radiation efficiency around -10 dB. Meanwhile, for the proposed MDA antenna, return loss is ranging from -10 dB to -17 dB with the separation from 1 mm to 5 mm. Due to increased physical length, proposed MDA shows a high radiation efficiency up to -2.8 dB, which is 7 dB higher than MMA. And hence, MDA insertion loss has been boosted up to only -7 dB, which indicates a promising

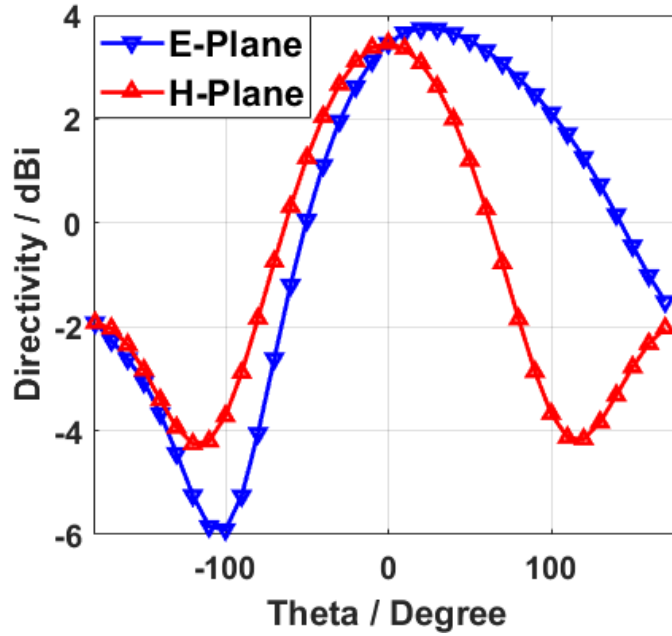


Figure 5.6: Directivity of the proposed meander dipole antenna (MDA) with 1.5 mm Tx and Rx separation, which shows it's suitable for our proposed hybrid CDN.

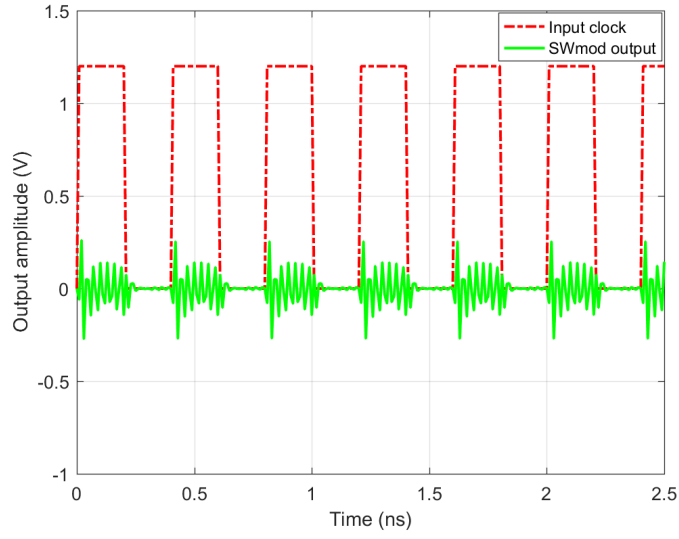
transmission gain for on-chip wireless clock distribution. Besides, comparing with MMA, the proposed MDA antenna shows a high total radiation efficiency of around 74% at resonance compared to that of 13% in [104] without losing significant power in metallic antenna fabrics, which improves overall energy efficiency. In addition, the moderate directivity of less than 4 dBi compared with a hypothetical isotropic antenna also enables this MDA as a potential antenna for the proposed hybrid CDN, as shown in Figure 5.6.

By importing the scattering parameters in Touchstone format to Cadence Spectre, SPICE-level simulation for TSMC 65 nm implementation can be performed in the remaining of this section.

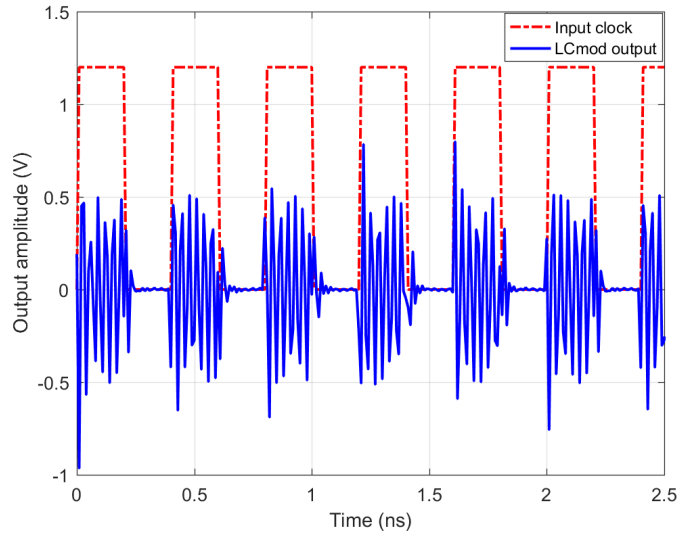
## 5.2 Results of the Proposed Clock Transmitter

As mentioned in the previous chapters, our proposed OOK modulator would work as analog switch to simplify circuit design and mitigate transmitter power consumption. Under 69 GHz input carrier signal, switching modulator has shown a clear characteristic of on-off isolation, which is beneficial to relieve inter-symbol-interference (ISI) caused by transmission loss and cross-talk.

As shown in Figure 5.7, previous switch-based modulator (SWM) follows the input clock as control signal, thereby flipping output accordingly. Simulated on-off isolation could reach 26.9 dB using 65 nm process, which is essentially due to the short channel length and therefore, higher leakage current in transistor sub-threshold region. Hence this



(a)



(b)

Figure 5.7: Output clock signal with an input 2.5 GHz nominal clock for the proposed (a) switch-based modulator (SWM) and (b) leakage-compensation modulator (LCM).

architecture is no longer adequate for modulators using smaller technologies. For the newly developed leakage compensation modulator (LCM), as the leaked carrier during off state is basically eliminated by phase cancelling, the on-off isolation has been boosted up to 47.8 dB with 65 nm process again which is 18 dB higher than the previous work [104] and indicates a 66% improvement compared to that of SWM, which indicates the potential of higher clock rate to be transmitted, according to Figure 5.8.

In addition, differential signaling facilitates the on-off isolation by further enlarge output



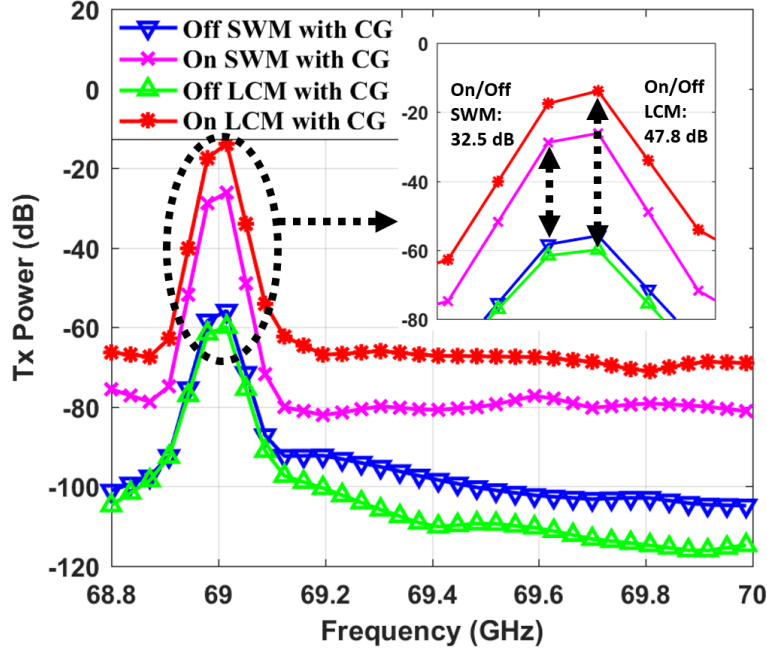


Figure 5.8: 65 nm leakage compensation modulator (LCM) and switching modulator (SWM) output in frequency domain with the proposed cancelling gate M5 and M6.

	[46]	[65]	This Work
Technology (nm)	90	90	65
Carrier frequency (GHz)	45	60	69
Modulation	ASK	OOK	OOK
Data rate (Gbps)	5	1	10
Clock frequency (GHz)	2.5	0.5	5
Clock amplitude (mV)	200	360	1200
On-off isolation (dB)	28.0	24.3	47.8
Power (mW)	61.6	183	20.1

Table 5.1: Comparisons between different implementation of the wireless transmitter in related researches and our proposed LCM-based clock transmitter.

swing by a factor of 2, thus provide better output quality than conventional single-ended design. Differential signaling could further effectively reduce the total common-mode noise generated by nearby digital logics, therefore increase noise immunity. Since the proposed wireless clock transmitter (Tx) is for delivering clock signals only, power amplifiers has been removed from a conventional implementation of RF transmitter. Hence, our proposed Tx exhibits a much improved energy efficiency compared to related designs. Also, the leakage-compensation modulator-based Tx shows its advantage in clock frequency to be transmitted, due to the boosted on/off isolation, as shown in Table 5.1.

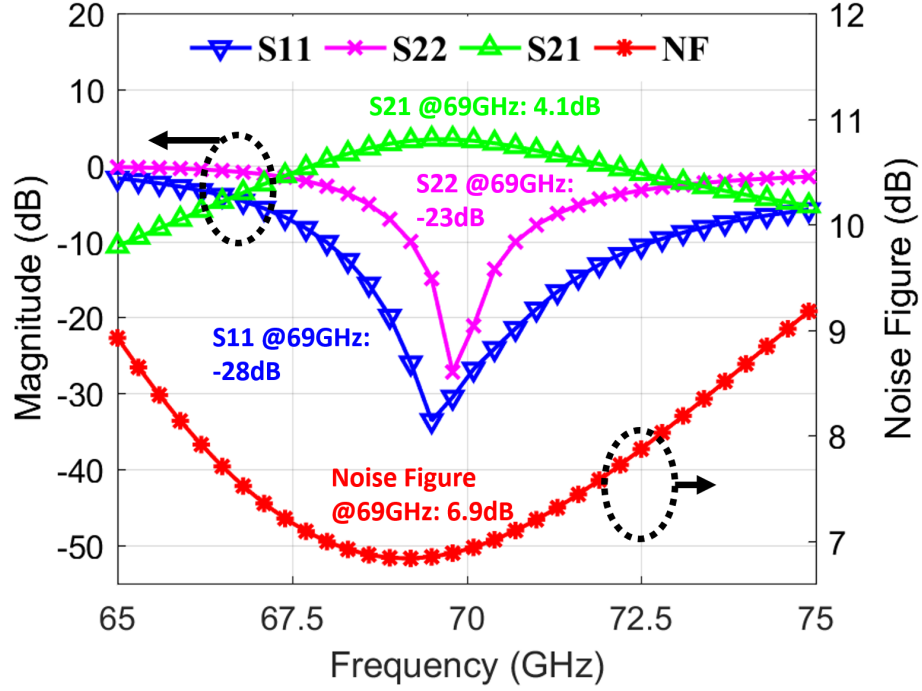


Figure 5.9: Rx front-end results in terms of the S-parameters of the proposed pseudo-differential low-noise amplifier (lna) and noise figure at 69 GHz.

### 5.3 Results of the Proposed Clock Receiver

Previously, with the assumption of zero wireless transmission loss, the proposed demodulator is chosen as the Rx front end and directly connected to the subsequent logics for the purpose of simplicity. In this chapter, channel loss caused by on-chip MDA antenna would be included and a pseudo-differential LNA is used for any signal decay compensation.

For the proposed two-stage cascoding LNA, according to Figure 5.9, the single-stage forward voltage gain  $Gain_{total}$  is around 2.7, which consequently mitigate the impact of signal energy loss inside lossy silicon substrate. Two identical LNAs form a pseudo differential structure to directly connect balanced MDA antenna with an input impedance around  $50\ \Omega$ , hence each input terminal would exhibit an input impedance of around  $25\ \Omega$ .

As shown in Figure 5.9, the pseudo-differential LNA compensates the energy loss in signal propagation with a total differential forward gain of around 7 dB, a well-matched return loss of -28 dB, which can recover the attenuated radio frequency signal to adequate voltage level that can be recognised by the subsequent rectifier inside the modulator. Besides, the proposed lna exhibits a minimum noise figure (NF) less than 7 dB without bringing too much device noise into enlarged RF signals at the frequency of interest, therefore providing the potential of delivering higher clock frequency for future improved designs.

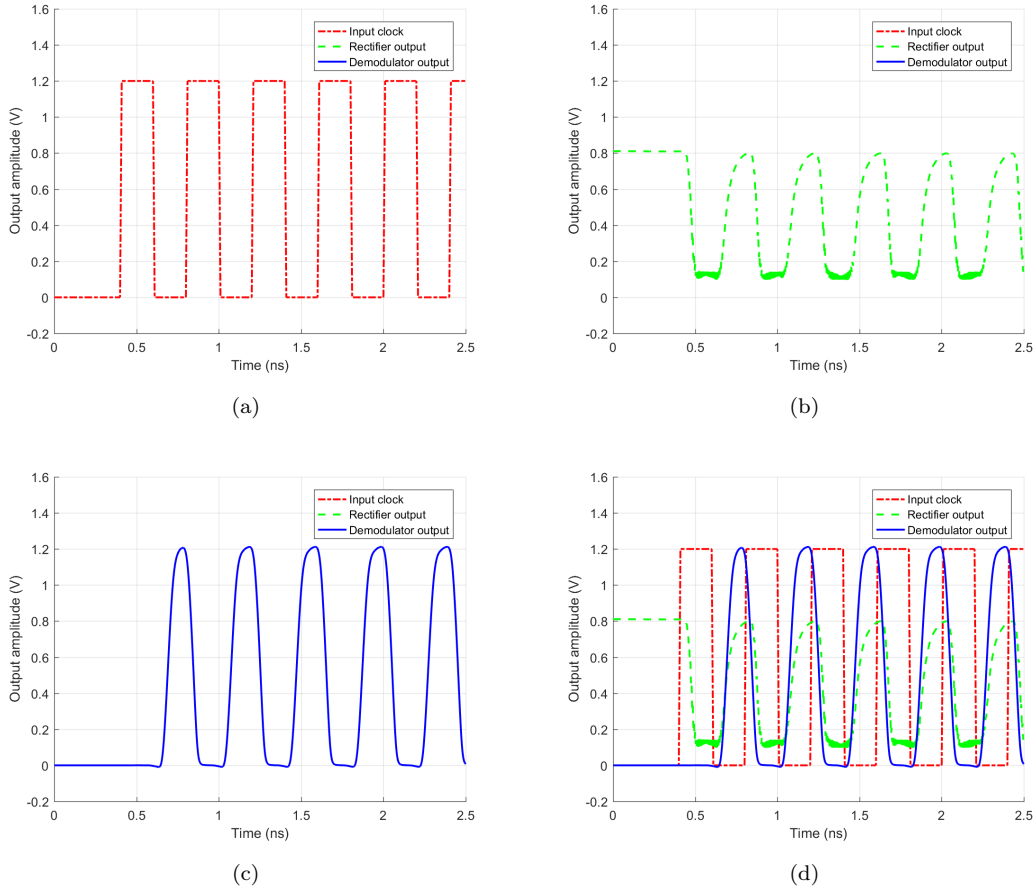


Figure 5.10: Waveform of the (a) input 2.5 GHz nominal clock signal, (b) rectifier output with inverted signal polarization, (c) the full-swing demodulator output signal at 1.2 V and (d) effective signal stack-up.

Demodulator would sense the output of the transmitter first, and then the small signal swing would be enlarged and filtered by subsequent baseband amplifiers. As illustrated in Figure 5.10, proposed gain-boost rectifier exhibits an output swing of 700.2 mV with a DC bias of 500 mV. Transmitted differential clock signal is then amplified and buffered by the subsequent dual-threshold Schmitt-trigger buffer, with a full-swing clock at 1.2 V ( $V_{dd}$ ) with a common mode level at around ground at 2.5 GHz nominal frequency as shown in Figure 5.10, therefore further amplification is reduced and energy efficiency is boosted.

Based on the experiments, the maximum frequency of the proposed circuit can be found around 5 GHz. For a full swing circuit, there's no need of any further amplification, hence the recovered clock can be directly used by downstream logics. However, with the frequency greater than 5 GHz, the amplitude of the clock signal begins to drop. As shown in Figure 5.11, the clock amplitude has been reduced by 20.8% with a 7 GHz input clock, because of the incomplete switching transistors. With the increasing clock frequency up to 10 GHz, the recovered clock amplitude will keep dropping to only

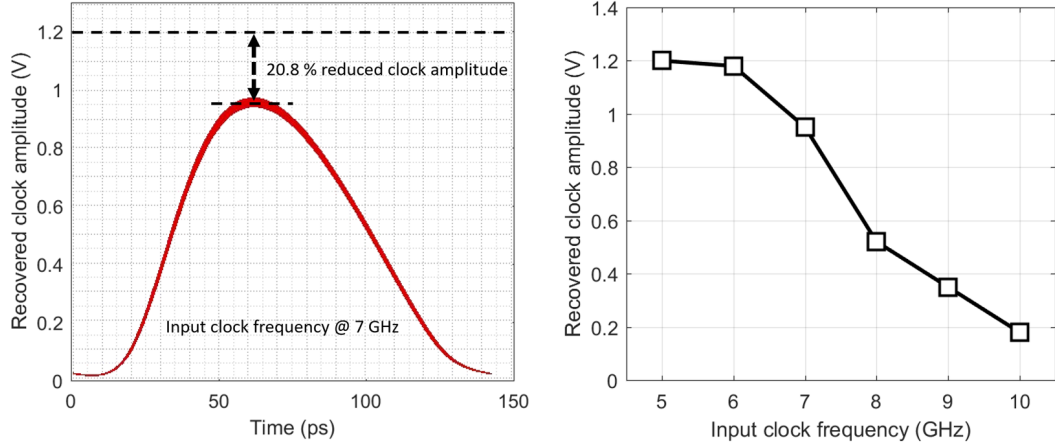


Figure 5.11: Example of 7 GHz recovered clock signal with 20.8% reduced signal swing. Clock amplitude will keep degrading with the increase of input clock frequency, transistors in clock Rx circuit cannot switch on/off completely in a reducing clock cycle, hence generating attenuated recovered clocks.

around 16% of the nominal V<sub>dd</sub> level at 1.2 V, which indicates that extra circuits such as level shifters or baseband amplifiers are necessary for generating effective clock edges for standard logic.

Besides, recovered clock signal shows a 10%-to-90% rise time of only 63 ps and 57 ps respectively at the nominal and maximum supported frequency of 2.5 GHz and 5 GHz. Due to faster switching frequency, this implementation shows a slew rate of 27.4 V/ns and 24.7 V/ns for 2.5 GHz and 5 GHz recovered clock respectively, which is highly adaptable compared with the input clock signal.

As the input clock frequency keeps increasing, transistors are more or less reaching the performance limit of speed, hence the recovered clock might show a decreased amplitude because of the incomplete on-off switching. A half-swing or a more general partial-swing clock might be more practical for future works with smaller technology nodes inside a potential application chip.

Measured from eye-diagrams generated by Cadence Spectre, in the presence of input random coupling noise with amplitude less than 50 mV, our proposed hybrid CDN exhibits a 19.3 dB SNR that is 2.1dB higher than TLM CDN with a peak-to-peak jitter reduction around 48%, which further shows the benefits of a hybrid solution in terms of noise rejection and lower uncertainties, according to Figure 5.12.

Figure 5.13 illustrates the global clock signal received at different locations. As the sixteen clock receivers are evenly located on four corners of the chip area, only the four receivers at the upright corner are studied as the other receivers are basically the duplication of this specific corner. Clock skew begin to occur because of the different distance from Tx to Rx. Experimental results indicate that a maximum skew of 24.4ps occurs between location 1 and location 4 for a 4-fanout architecture with uniform load

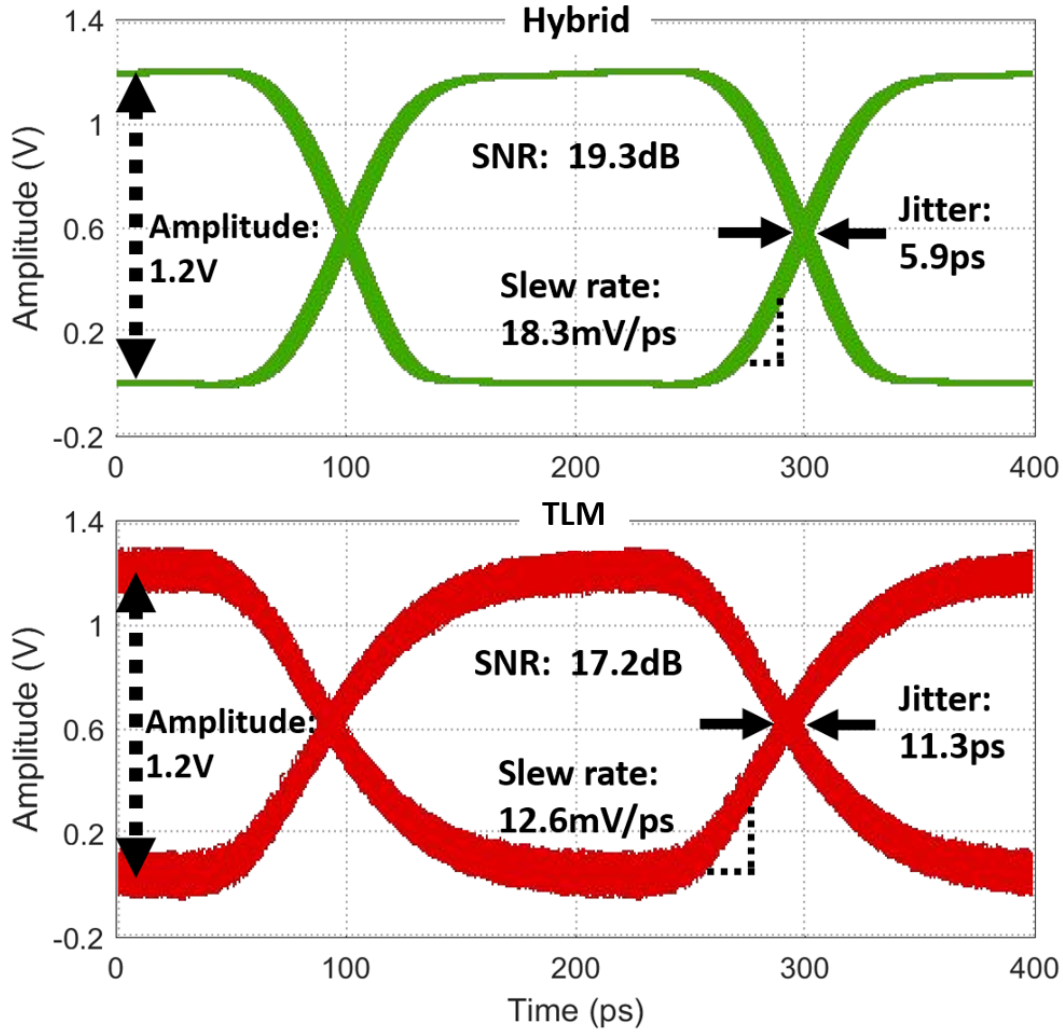


Figure 5.12: Comparisons between the measured eye diagrams of the recovered clock signal from conventional wired global tree local mesh (TLM) CDN and our proposed wireless global CDN at 2.5 GHz, with robust eye height/clock amplitude around 1.2 V with 50 mV additive noise at the output of clock buffers. ©2021 IEEE

capacitance, which shares the largest distance between each other. For a more critical situation, the impact of different values of load capacitance on clock skew performances must be taken into account.

To observe the delay performance of the proposed design, referring to Figure 5.13, a baseline architecture with 16 fanouts is adopted. Local load capacitances are modelled as capacitors shunt at the output terminal of each clock receiver. To maximumly model the impact of load unbalances, small load with only 10 fF and heavy loads with capacitance up to hundreds of fF is adopted and randomly distributed at different locations. Since the time-of-flight (ToF) of the EM wave would only depend on the distance between clock transmitter and receiver pairs, and it requires more time for receivers to charge up a load with bigger capacitance, naturally, it is beneficial for skew minimization to locate

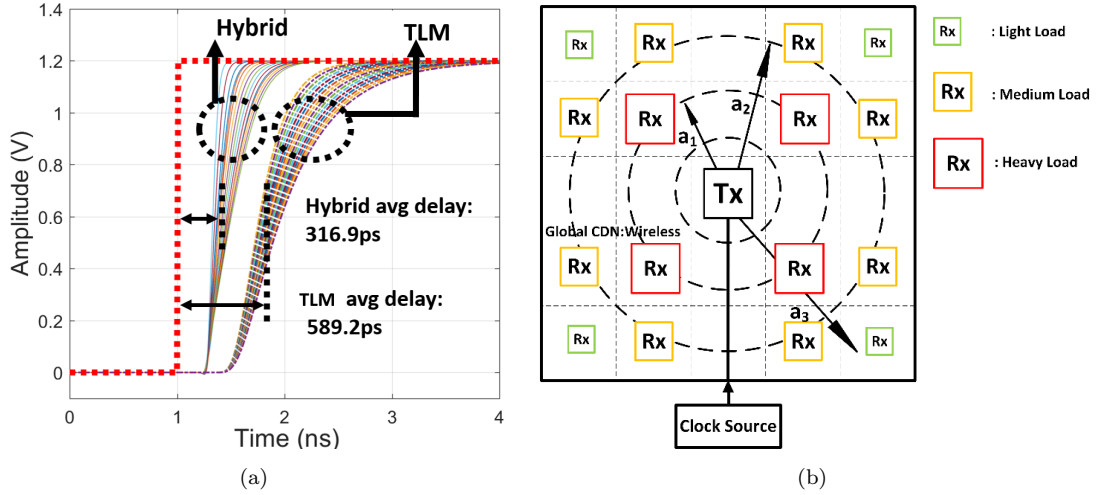


Figure 5.13: Clock skew with a pseudo-mesh based topology at 2.5 GHz using 65 nm Technology. (a) An overall 16 fan-out architecture is adopted with an average skew reduction of 45.5% under random unbalanced load ranging from 10 fF to 500 fF, and (b) unbalanced load allocation (NBFS) for skew minimisation.

	wired	wired	wired	wireless	wireless	Hybrid
Solutions	[43]	[49]	Tree & Mesh	[6]	[7]	This work
Technology (nm)	180	32	65	180	130	65
Efficiency	360	112	325	337.5	282.1	325
Power (mW)	14.4k	1k	619.2	Tx: 48 Rx: 40	Tx: 75 Rx: 69	<b>Tx: 24.3</b> <b>Rx: 17.5</b>
Distance (mm)	14.7	20.9	1-18	5.6	10-18	1-18
Max Frequency (GHz)	2	3.5	5	1.875	2.17	5
Clock amplitude (mV)	N/A	1000	1200	350	600	1200
Global jitter (ps)	35	N/A	11.3	6.6	12	4.6-5.9
Global skew (ps)	16	10	53.5	25	18	2.1-32.1

Table 5.2: Comparisons between different implementation of global CDNs using wire and wireless approaches

receivers with higher load capacitance at the inner ring with smaller radius, and locate receivers with smaller load capacitance at the outer ring with larger radius.

For global CDN timing performance, as shown in Figure 5.13, full-swing clock signals have been successfully recovered at 16 different positions within test model up to 5 GHz with different effective local loads from 100 fF to 2 pF for both conventional tree-grid (TLM) CDN and hybrid wireless-grid CDN. Benefit from the wireless global CDN, measured average 50% logic transition delay of our proposed hybrid solution with 1.2 V step input is approximately 316.9 ps, which is less than 50% of the delay of a conventional TLM solution.

Furthermore, the proposed hybrid architecture shows a competitive immunity against VCO phase noise. With a phase noise increasing from -100 dBc/Hz to -50 dBc/Hz at 1 MHz frequency offset, total peak-to-peak jitter has only increased for a limit value





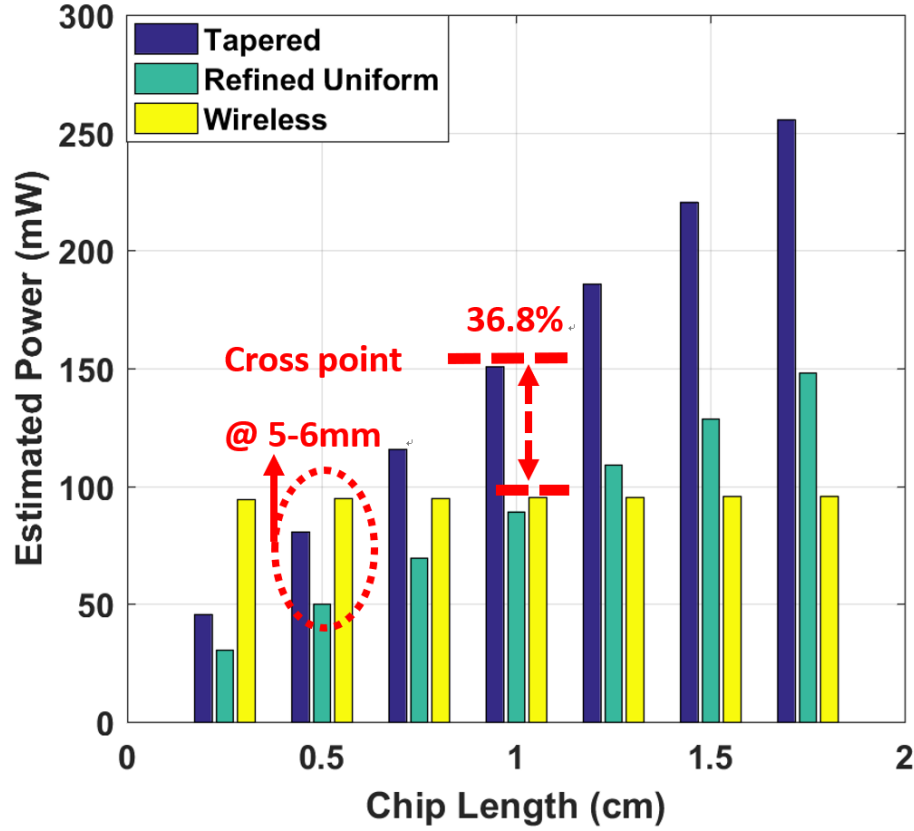


Figure 5.15: Power comparison between global H-tree and proposed global wireless CDN. ©2021 IET

from 0.3 cm to 1.7 cm. H-tree architecture shows a phase delay of 2.6 ns, which is almost 20 times larger than the proposed wireless CDN, thus might limit system performance significantly. Compared with delay model adopted in previous work [104] and [105], the delay prediction model adopted in this work utilizes distributed circuits elements, and therefore, its mathematically more accurate and practical than the previous version with a factor of  $\ln(2)$ , as shown in Figure 5.14.

Results shown in Figure 5.14 shows that using the new delay prediction model, proposed design still illustrates a significant delay reduction around 95%, which further verify the huge potential of the performance of wireless clock distribution.

Furthermore, if the overall power consumption of a H-tree could be calculated as [124] the switching power consumed by clock buffers, given by:

$$P_{tree} = a_v C_{total} f_{clk} \cdot V_{dd}^2 \quad (5.2)$$

where  $a_v$  is the activity factor indicating the fraction of the switching circuit components,  $C_{total}$  is the total capacitance seen at the source of global CDN,  $f_{clk}$  is the clock flipping frequency and  $V_{dd}$  is the supply voltage. For the clock distribution network, its effectively assumed that the CDN toggles twice per clock cycle hence  $a_v$  for H-tree equals to 1.



On the contrary, overall global power consumption of the proposed wireless CDN shows that:

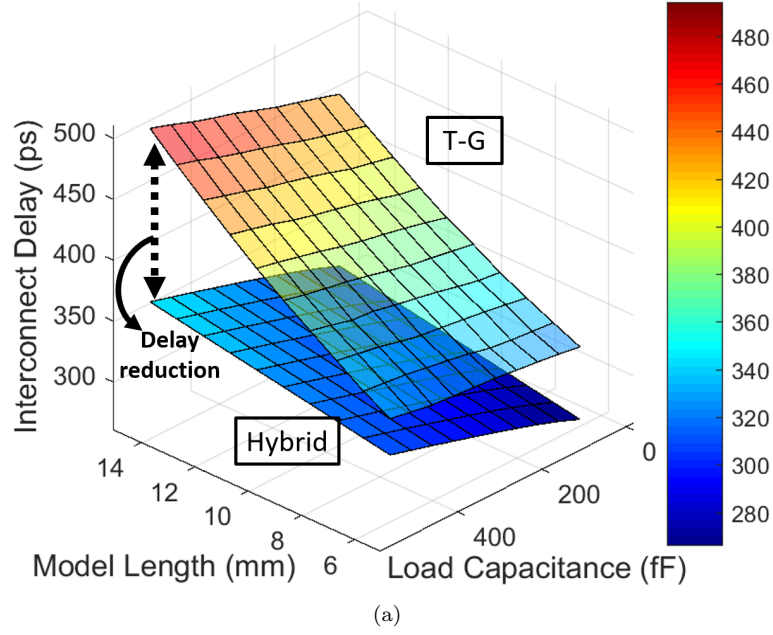
$$P_{wireless} = N_{Tx}P_{Tx} + N_{Rx}P_{Rx} \quad (5.3)$$

where  $N_{Tx}$ ,  $N_{Rx}$ ,  $P_{Tx}$  and  $P_{Rx}$  are the number of the transmitter, receiver and the power consumption of a single transmitter and receiver respectively. Given a 4-receiver global architecture, proposed design shows a cross point in terms of power consumption at 5 mm chip side around 95.5 mW as shown in Figure 5.15, which illustrates that when the size of the die is essentially small, conventional clock distribution network using metallic interconnect could be adopted for its power-efficiency. However, with the increase of the die size, its naturally suitable to utilize proposed clock distribution network for its power robustness and broadcasting feature. Estimated single transmitter and receiver power consumption using 65 nm process could be acquired through Cadence Spectre around 24.3 mW and 17.5 mW, respectively. For more receivers working in parallel, low power receiver design is required for overall power reduction.

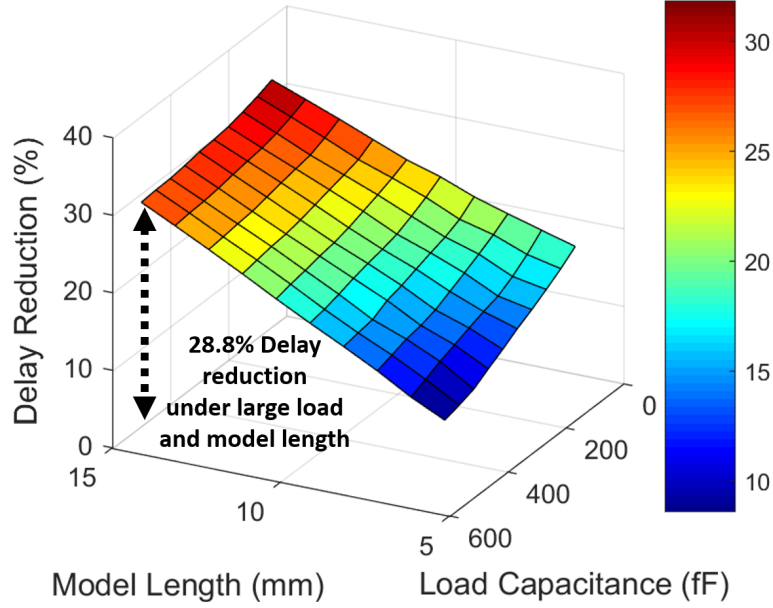
To have a more comprehensive evaluation about the impacts of local parameters, for instance, local clock frequency, system scale and loads, in this section, comparisons between hybrid and TLM solutions have been produced in terms of overall delay, skew and power consumption as per varying test model parameters to evaluate the impact on essential system performance criteria.

According to Figure 5.16(a), a conventional TLM solution exhibits a rapid increase of propagation delay in the presence of global repeaters (clock buffers) and is limited to a fixed input clock frequency. Metallic interconnects are more sensitive to total route length than local load capacitance. By contrast, our proposed hybrid solution shows a much-reduced interconnect delay. As shown in Figure 5.16(b), a normalised delay reduction of up to 28.8% is shown between the TLM and hybrid architecture under heavy load and large system size. Since the delay of our hybrid solution could be considered to be merely related to the distance of communication, under the assumption of signal propagation speed near the speed of light, our proposed hybrid solution would eventually exhibit a robust and competitive delay performance even with larger system size and local loads.

On the other hand, as shown in Figure 5.17(a), the TLM solution exhibits an ideal clock skew performance under small load imbalances because of its highly-symmetric structure. However, it's not always the case that the load is symmetric. With larger load imbalances between local nodes, the skew of TLM begins to increase in a pseudo-exponential way. Our hybrid solution, however, exhibits a complimentary skew tendency. A hybrid CDN would essentially get larger skew because of the increased distance between near-end and far-end nodes under uniform load assignment. Since clock propagation delay of a hybrid CDN only depends on the receiver displacement [105], skew could be effectively



(a)



(b)

Figure 5.16: Comparison between the TLM approach and the proposed hybrid architecture of (a) interconnect delay and (b) normalised delay reduction according to size and load variation with a maximum 28.8% decrease under a large model scale. ©2021 IEEE

mitigated via load compensation by allocating larger loads at near-end nodes and smaller loads at far-end nodes (NBFS rule). As shown in Figure 5.17(b), the normalised global skew of a hybrid CDN can be eventually reduced to less than 32% of that in a TLM CDN for large load imbalances, which indicates a promising feature for heterogeneous integration in 3D stacking systems.

Besides, for overall system power consumption, it is clear that power with fixed model size shown in Figure 5.18(a) indicates that TLM CDN exhibits a competitive efficiency

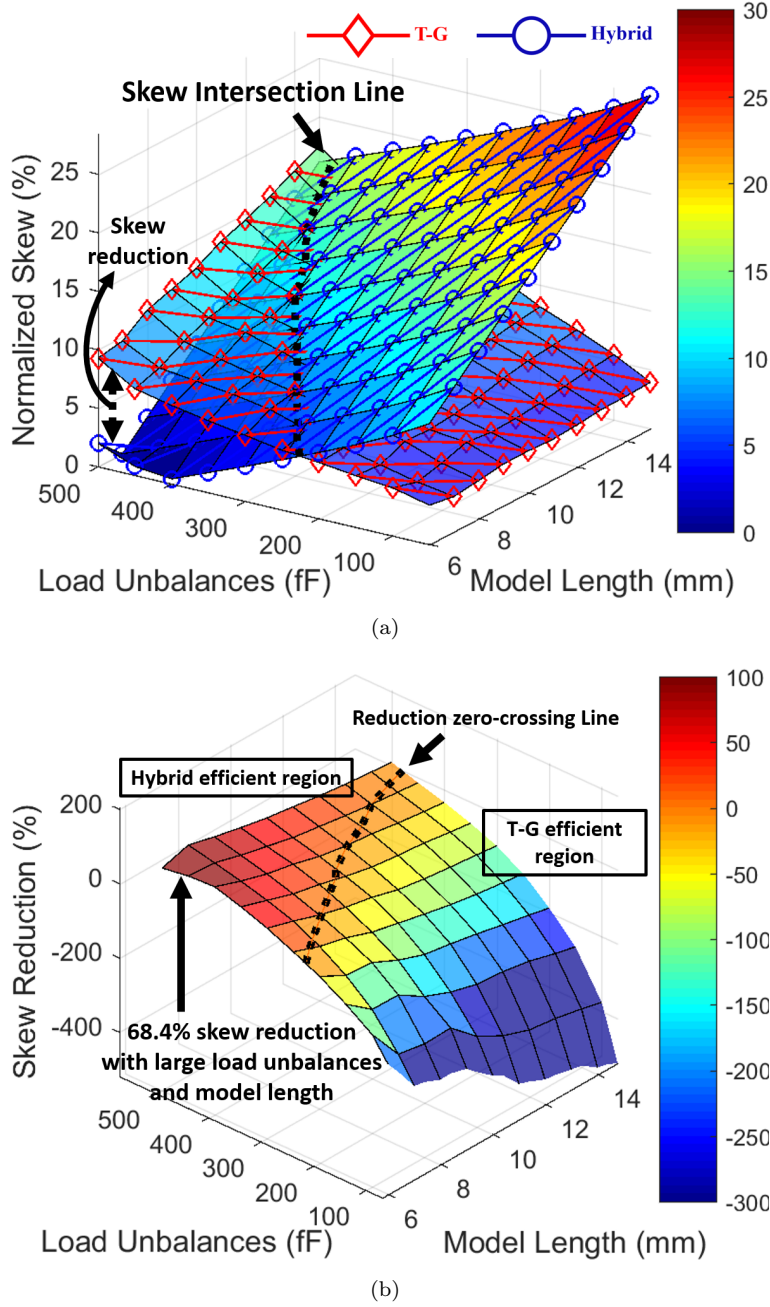
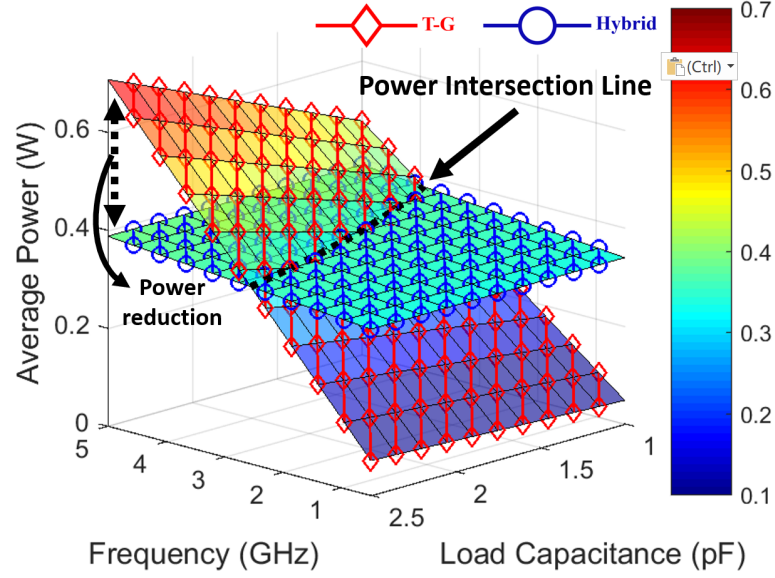
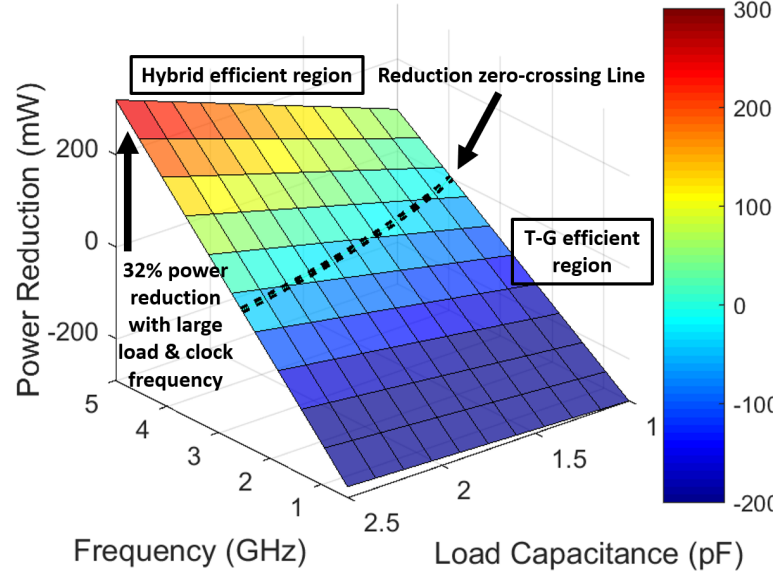


Figure 5.17: Comparison between the TLM and the proposed hybrid architecture of (a) global clock skew and (b) normalised skew reduction according to size and load variation with a maximum 68.1% decrease under large load unbalance, within one clock domain. ©2021 IEEE

with small local clock frequency under 2 GHz. However, power increases significantly with the increase of clock frequency and load capacitance, which becomes a major contribution to the overall power budget. Conversely, a hybrid CDN exhibits a less-competitive power consumption under low local frequency. Nevertheless, the overall power of clock Tx and Rx is almost independent of the clock frequency as the majority of the power is consumed when the Tx and Rx are in their on-state. Therefore, the power consumption



(a)



(b)

Figure 5.18: Comparisons between the TLM approach and the proposed hybrid architecture of (a) overall power consumption, (b) power reduction according to local clock frequency and load variation with a maximum 32% decrease under large local clock frequency. ©2021 IEEE

of our proposed hybrid CDN only depends on the duty cycle of the input clock signal and the number of Rx. From Figure 5.18(b), the proposed hybrid CDN shows a better power efficiency for higher input clock frequency and larger load capacitance on the left side of the reduction zero-crossing line. The power consumption of our proposed hybrid CDN exhibits a maximum 32% reduction compared to TLM under a maximum supported 5 GHz input clock and a 2.5 pF local capacitive loading.

Also, experimental results for concentrated and distributed Rx planning are given in

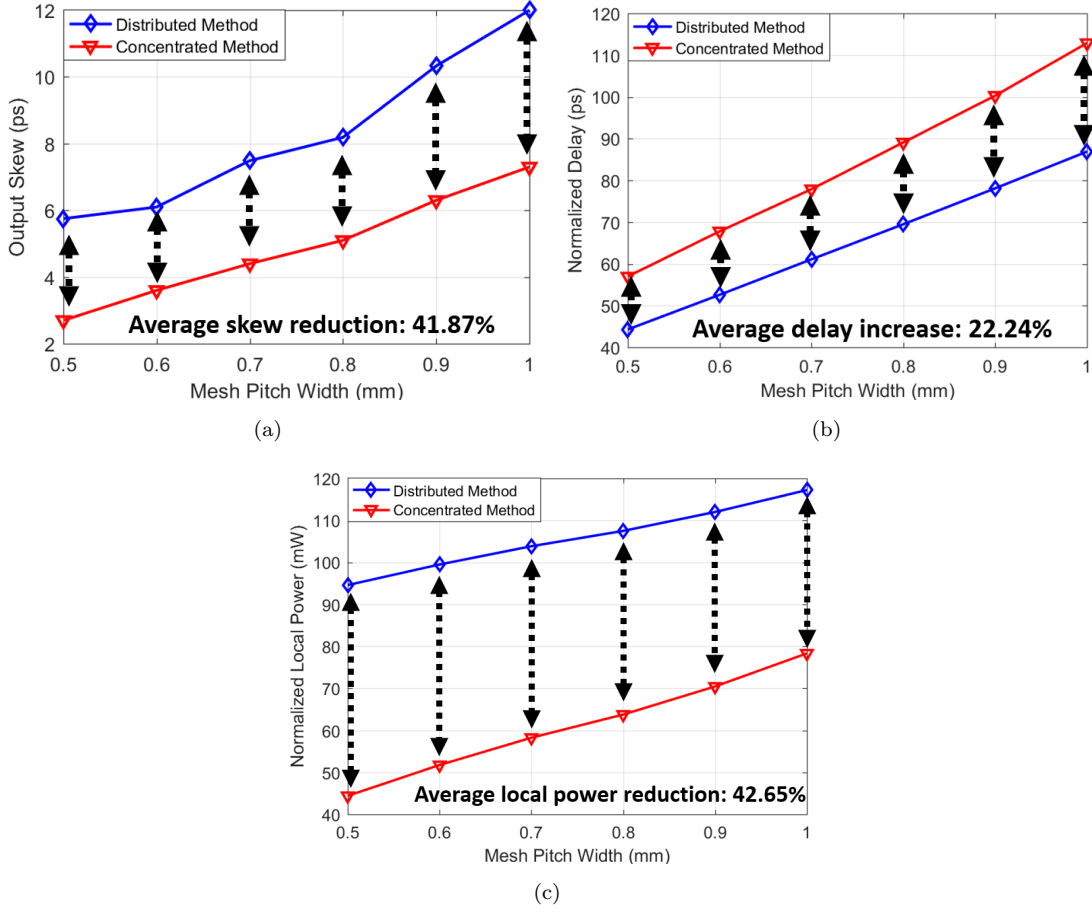


Figure 5.19: Different parameters in terms of (a) output skew, (b) clock latency and (c) power consumption in the test case for both distributed and concentrated planning. Concentrated planning trade higher propagation delay with better skew and power performance using coupled model.

Figure 5.19, based on a local clock distribution area of  $1.3 \times 1.3 \text{ mm}^2$  and 1fF to 4 fF random local load allocations. It is clear that concentrated planning shows a better skew performance with an average reduction of around 41.87%, with the mesh pitch width ranging from  $500 \mu\text{m}$  to 1 mm as shown in Figure 5.19(a) at the cost of 22.2% increased delay shown in Figure 5.19(b). Besides, concentrated planning shows an average power reduction of around 42.65% in the proposed test case, as shown in Figure 5.19(c).

The above results are based on global CDN, to verify the impact of the proposed global wireless CDN on the local clock performance, a local mesh structure is then selected using the local CDN generation algorithm in Chapter 3. Under the same testing conditions, a cost function is generated to find a local optimal operation point for a given test case input given as follows:

$$Cost_i = \frac{\sum_j^k (\alpha P_{j_i} + \beta A_{j_i} + \max\{|\gamma S_{j_i}|\})}{\max\{|\gamma S_{j_i}|\}}, \forall i, j \in \theta_{m,n} \quad (5.4)$$

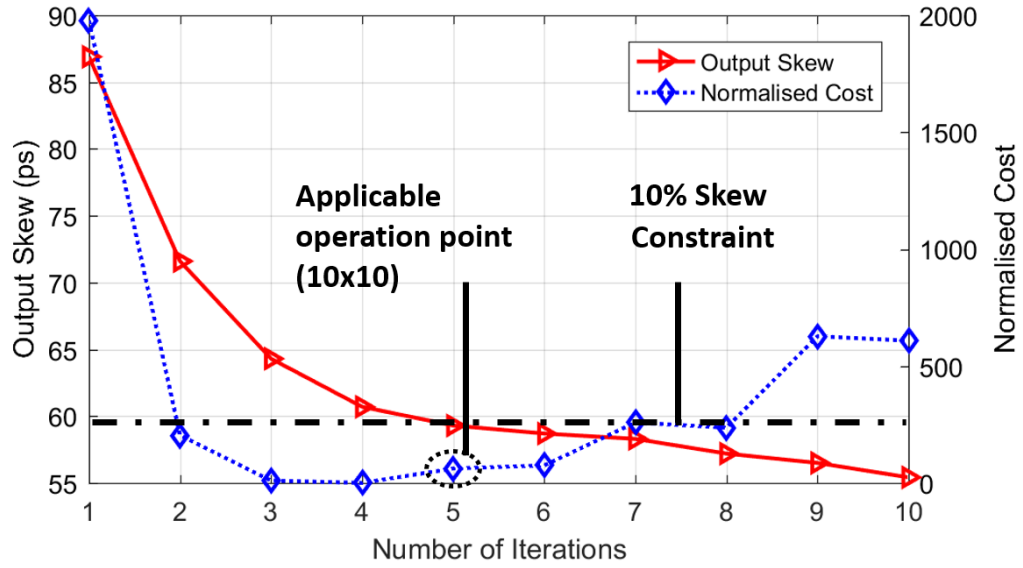


Figure 5.20: Total skew and normalised cost in the local region with increasing iterations. ©2021 IEEE

where  $P_{ji}$ ,  $A_{ji}$  and  $S_{ji}$  are the results of  $j_{th}$  section among all  $\theta_{mn}$  sections in the local area during the  $i_{th}$  planning iteration for power overhead, area overhead and skew reduction respectively. The weighting coefficients for the three inputs  $\alpha$ ,  $\beta$  and  $\gamma$  are equally set to be 1 for a balanced design constraint. However, they can be tuned for different design preferences to quantify the appropriate CDN topology [97]. Figure 5.20 shows that after the proposed CDN has reached below its initial skew target, an applicable operation point is found such that the size of the local mesh is fixed around  $10 \times 10$ , further iterations will consequently waste metal resources and increase power dissipation. Therefore, the introduction of the cost function can effectively reduce the unnecessary power waste through the proposed local mesh planning process.

To evaluate the accuracy of our proposed local mesh planning algorithm, an experiment with input global skew ranges from -30 ps to 30 ps is tested based on previous experimental conditions. Both proposed isolated model and coupled model indicate the capability of skew reduction under this framework. Compared with SPICE simulation results under the same topology and capacitive loadings, both proposed modelling methods indicate a slightly higher minimum clock skew while the absolute input skew is near zero.

However, both proposed methods show a stronger immunity against input skew especially for those larger than around 24 ps. The coupled model shows an average error reduction of around 55% compared with isolated cells under the condition of unbalanced input. Besides, the coupled model shows an average 5.7% error rate compared to the SPICE simulation, which indicates the promising accuracy of our proposed mesh planning algorithm shown in Figure 5.21. Since the clock output from clock Rx would not



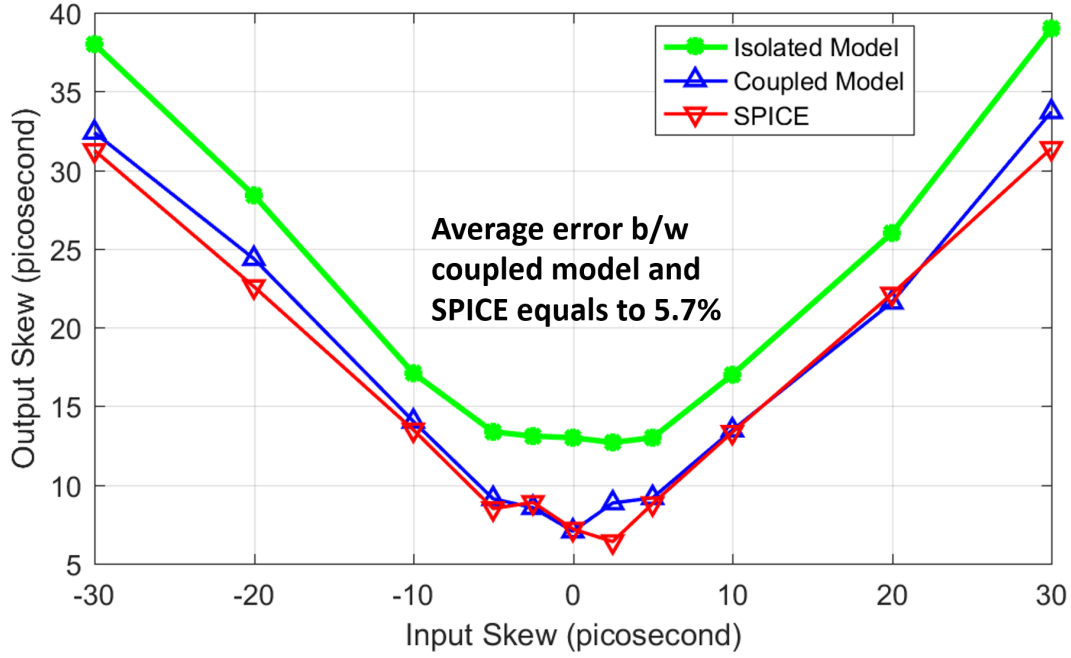


Figure 5.21: Output skew in local CDN under different input global skew caused by different communication distance or buffer mismatch. The coupled model shows an average 5.7% error comparing to SPICE simulation. ©2021 IEEE

arrive precisely at the same time, and consequently generates unpredictable input skew into local CDN, this method can accurately model most of the general cases for the local mesh.

## 5.5 Summary

In this chapter, the performance of the on-chip antenna is verified through EM simulation using CST Microwave Studio. Moreover, the performance of the proposed wireless CDN circuits have been evaluated using Cadence Virtuoso Analog Design Environment. In our previous work, a simplified simulation architecture is adopted to achieve system simplicity based on the assumption of zero antenna transmission loss and an omnidirectional propagation pattern. To observe the physical impact on the EM wave propagation in the EM simulation model, in this work, two different on-chip antennas are designed with compact layout and moderated radiation efficiency. Based on the scatter parameters obtained from EM simulations, cross-talks are evaluated for the on-chip antennas. Experimental results demonstrated in table III and IV indicates that CDN using wireless interconnect has the potential to transmit clock signals at a cost of relatively low propagation delay and low global clock uncertainties, thus improving performance of a many-core system with large synchronisation area such as CMP and NoC.

In addition, wireless CDN benefits from a natural fan-out feature, hence providing immunity against clock uncertainty caused by unbalanced load. Clock uncertainties in

terms of skew and jitter become independent of the length of conventional wires existing in local CDN. Besides, clock uncertainties are exclusively related to the displacement of clock receivers. By carefully consider Rx locations or the overall geometry of TRx architecture, clock uncertainty is now predictable, which could be further mitigated during design stage and provide more flexibility.

To verify the comprehensive CDN performance including both global and local CDN, two realistic test cases will be adopted to test the performance of the proposed hybrid wireless-wired CDN in [Chapter 6](#).





## Chapter 6

# Case Study: the Proposed CDN Verified by Testbench Circuits

In this chapter, two comprehensive applications have been proposed and the hybrid CDN has been applied to verify whether it could outweigh conventional clock distribution schemes in terms of both clock performance and energy efficiency. In both proposed test cases, clock sinks, such as clock gating cells or the clock-input pins for register sets are either represented by capacitive loadings or extracted from an actual chip. The physical locations of the clock sinks are modelled by a two-dimensional coordinate. Therefore, using the models defined in Chapter 3 and tools such as MATLAB, one can effectively quantify the time-domain response of each clock sink within the clock distribution area of interest, by solving the inverse-Laplace transform given the transfer functions of global and local CDN.

### 6.1 Architectures of the Proposed Test Case

Recent CDN implementations to reduce the high clock latency and power includes partitioning designs across multiple dies and stacking them, allowing vertically and physically closer placement of logic elements, therefore reducing interconnect delay. On the other hand, Interposer technology provides extra horizontal interconnect resources and bandwidth for die-to-die communications.

Taking advantage of the silicon interposer underneath flipped chips, 2.5D and 3D integration could effectively offer a high bandwidth connection between mounted dies such as microprocessor and memories [125]. Nevertheless, if the size of interposer keeps increasing, the latency of the silicon interposer will eventually become critical for long-distance on-chip communication between dies, compared with a monolithic IC [126]. On the other hand, parametric variations between dies will make it even more difficult to

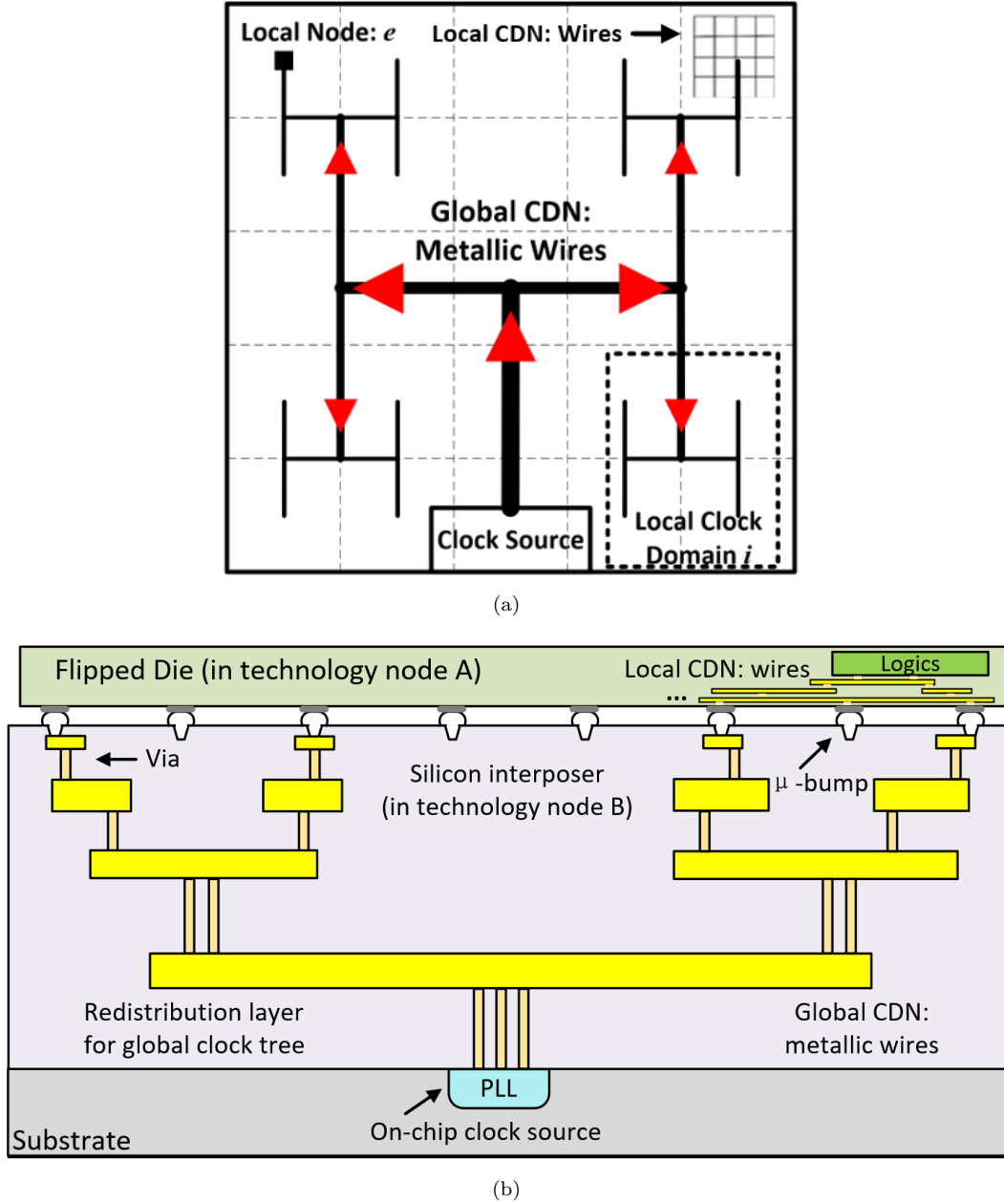


Figure 6.1: Models of a typical conventional global-tree local-mesh structure with (a) top view, (b) cross-sectional view, which consumes a large amount of metal. The flipped die could use a different technology node than the active silicon interposer, e.g. the flipped die is in 7 nm process and the interposer is in 65 nm process, which provides flexibility to the integration of different technologies.

deliver a clock signal with low skew and jitter inside interposer using conventional clock distribution solutions [127]. Different routes with different latency (skew) can cause potential timing violations in synchronous systems, therefore presenting a difficult task for synchronisation between dies within one clock domain mounted on interposer layer.

By incorporating the inherent fan-out feature of the wireless interconnect and the efficiency of a conventional tree/mesh for local clock distribution, one can estimate a

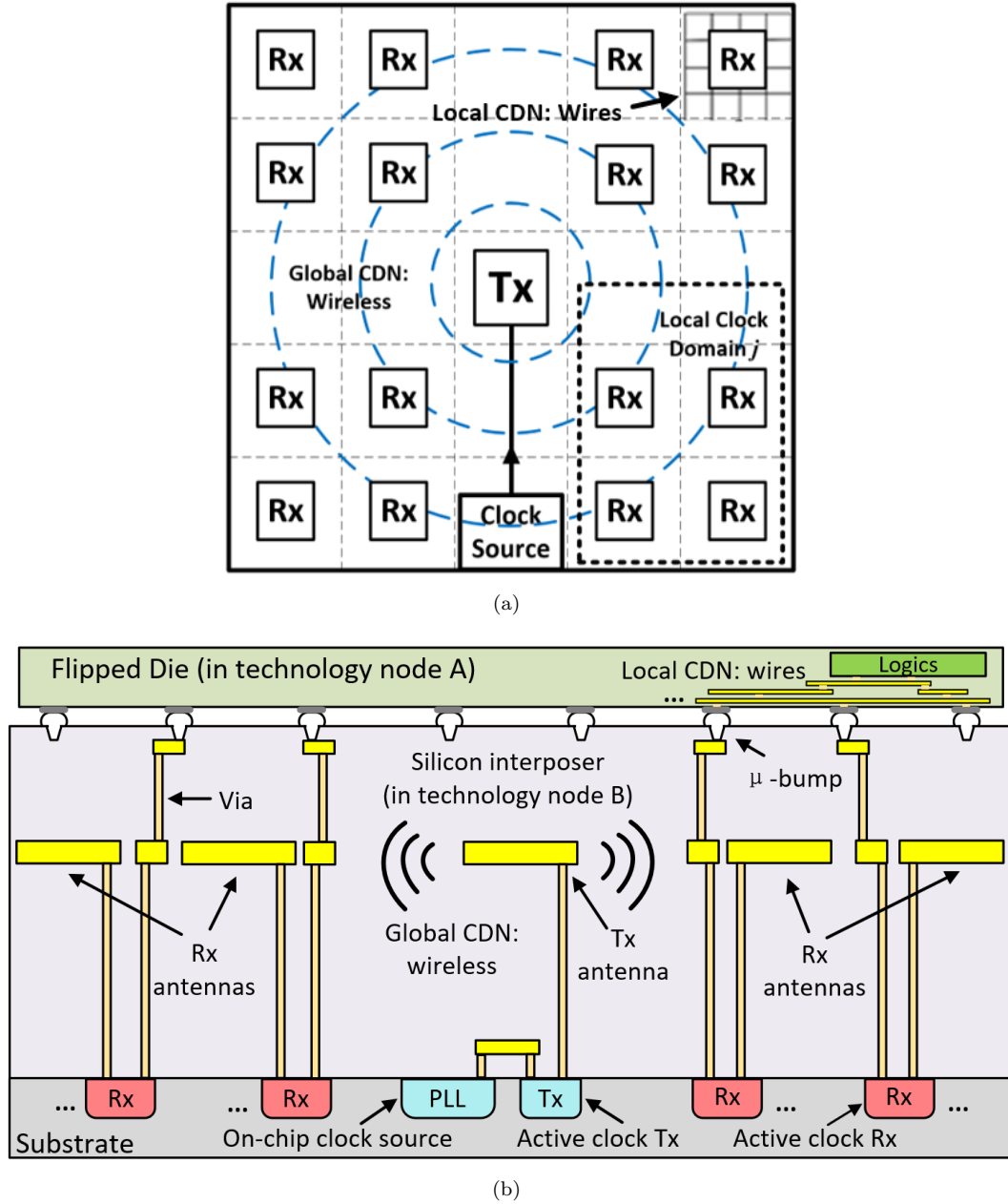


Figure 6.2: Models of our proposed hybrid clock distribution network, with (a) top view and (b) cross-sectional view, respectively. The flipped die could use a different technology node than the active silicon interposer, e.g. the flipped die is in 7 nm process and the interposer is in 65 nm process, which provides flexibility to the integration of different technologies.

significant improvement in terms of the overall power and uncertainty decrease for the reduction of global wires and local clock receivers [128]. Literature such as [73] suggests that with the CMOS fabrication technology scaling down, it's now possible to integrate wireless interconnect on-chip based on developed communication techniques. Although they are not naturally built for clock distribution, some of the compact modulation schemes can be adjusted and applied to broadcast global signals on-chip. Therefore,

we propose a hybrid wireless-wire test case, which is based on an efficient modulation scheme to accommodate to variable clock frequency by using an active silicon interposer. Moreover, unlike other solutions [6], the proposed solution allows a variable frequency clock to be distributed, which increases flexibility, hence making it naturally suitable for future many-core systems.

A conventional solution of a global tree and local mesh (TLM) clock distribution in a many-core system using silicon interposer is shown in Figure 6.1(a) and (b). Where the clock signal is transmitted via a global clock tree in global metal layers and then passed to local clock grids and local logics which require a timing reference. These conventional wires inside silicon interposer can also be modelled by using the transfer function of the interconnect [129], [130].

Our proposed global wireless CDN is shown in Figure 6.2(a) and (b). An integrated clock transmitter (Tx) is located at the centre of the clock distribution area inside an active silicon interposer, clock receivers (Rx) are placed following the ‘near-big’ and ‘far-small’ rule (NBFS) according to their load capacitance or, to be more specific, local buffer tree size.

Transceiver pairs then transmit and receive a global clock signal via EM radiation without conventional wires. Finally, the recovered clocks are passed to a local CDN to synchronise digital logic in a different technology node mounted on top of the interposer. For overall system simplicity, the global wireless clock distribution adopts an On-Off Keying (OOK) based clock modulator and demodulator, which can modulate an input clock signal onto the carrier at clock Tx and recover RF signal back to baseband local clock at clock Rx.

## 6.2 Test Case Circuit Generations and Experimental Setup

To evaluate the performance of our proposed CDN in realistic scenarios, clock signals with acceptable quality are delivered to an arbitrary local clock sink, such as a flip flop or a clock gating cell in the test case. Timing response of each clock sink will be monitored and stored. Also, the CDN power and skew performance will be compared with the conventional and state-of-art solutions, to better illustrate the attributes of our proposed work. For the specific test circuit in this chapter, only one clock Tx will be used and allocated at the center of the testing model to provide balanced global wireless clock latency, referring back to Chapter 3.

Furthermore, the test circuit will be regarded as a “core” in a many-core system, and will be expanded from a dual-core system to 64-core system, to test the adaptability of the proposed work. Clock sinks will be represented by capacitance of the clock sink pins defined in Chapter 2. According to their geometric locations, the proposed global

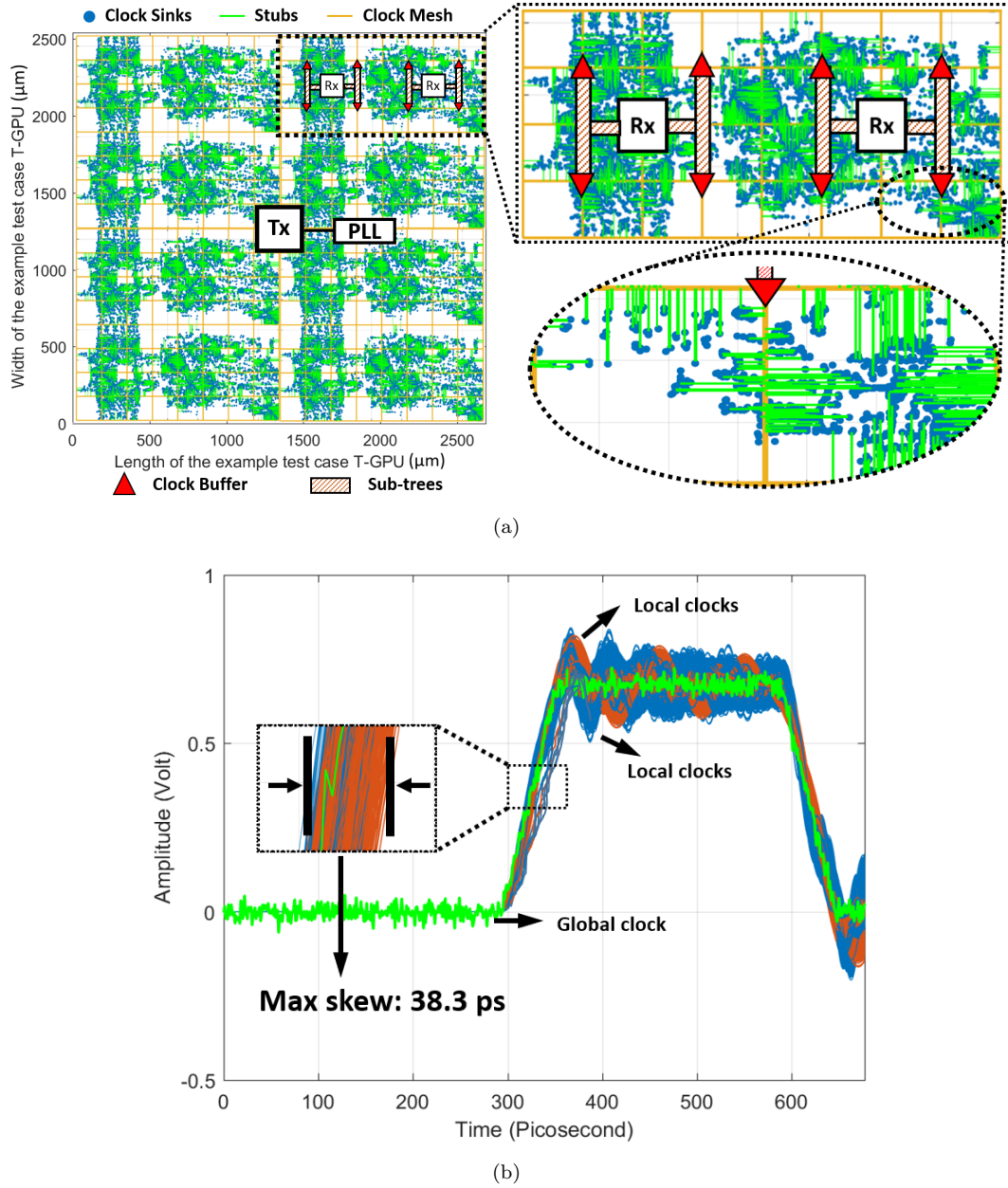


Figure 6.3: Unevenly distributed test case  $T_{GPU}$  extracted from Arm Mali-G77 GPU with (a) an example top view and (b) time-domain responses for  $T_{GPU}$  with 1.7 GHz recovered global clock input in an 8-core system. ©2021 IEEE

and local CDN design methodology introduced previously will be applied to implement a complete hybrid clock distribution network.

To test the timing performance under highly uneven loads allocation, we propose an industrial standard test case  $T_{GPU}$  with sink locations and capacitive loadings extracted of the Arm Mali-G77 GPU in 7 nm process [131] shown in Figure 6.3(a), which requires high power-efficiency and timing performance. Locations of the clock loads are extracted by using coordinates and capacitive loadings are ranging from 1.7 fF to 4 fF. Most of the capacitive loadings (over 90%) in the proposed test case are fixed at 1.77 fF, as clock

gating cells (ICGs) with similar size are commonly used in the GPU. Referring back to Figure 6.2(b), the 65 nm silicon interposer lies on the bottom of the proposed model, and 7 nm GPU cores are mounted on top of the interposer in another die.

For input signal in our experiments, a clock signal with 1.7 GHz frequency is chosen, which is the operation frequency of the actual fabricated Mali GPU with the power supply  $V_{dd}$  at 0.675 V. The skew upper bound is set to be 10% of the entire clock period. Interconnect parameters such as per unit length resistance, capacitance and inductance are set to be 286.5  $\Omega/\text{mm}$ , 189.6 fF/mm and 1.7 nH/mm, respectively, with a 0.16  $\mu\text{m}$  track width for modeling local clock distribution in 7 nm GPU cores; and 9.6  $\Omega/\text{mm}$ , 310 fF/mm and 1.6 nH/mm, respectively with a 2  $\mu\text{m}$  track width for modeling global tree in the conventional TLM solution in 65 nm technology, according to Predictive Technology Model [3], [18].

### 6.3 Test Case Results and Analysis

Time-domain response of our proposed test cases is shown in Figure 6.3(b). After several iterations of local mesh generation, the first applicable operation point is chosen to avoid excess metal waste as per Chapter 5. The  $10 \times 10$  local clock mesh can provide a 38.3 ps clock skew in local synchronisation region with a 1.7 GHz output clock, which satisfies the 10% skew constraint. To have a comprehensive view on the impact of the varying number of cores and size in  $T_{GPU}$ , results have been summarised in Table 6.1 and 6.2 from a basic dual-core architecture to a 64-core system. To compare the characteristics of the conventional approach and the proposed hybrid CDN, we also evaluated a baseline clock distribution solution and state-of-art implementation, which are the commonly adopted full fan-out tree (TLT) and Global Tree and Local Mesh (TLM) respectively. We apply TLT and TLM on  $T_{GPU}$  with the same topology, dimensions and buffer size to observe the impact of different CDNs on overall performance.

Given a  $8 \times 8$  mesh for each core using TLM method, comparing to a conventional full fan-out tree, the TLM CDN shows a robust skew performance against the increasing core counts and sink numbers, at the cost of higher clock power. By contrast, the conventional full fan-out tree CDN offers less power consumption because of the smaller load capacitance and wire length, at the cost of a significant rise of clock skew due to the increasing layers of tree structure and clock insertion delay.

Comparing the results in Table 6.1 and 6.2, our proposed approach shows a better power efficiency with a promising average reduction of around 31.9% and 16.0% for concentrated planning and distributed planning respectively, with one single clock transmitter, comparing to both TLT and TLM solutions. Average skew reduction, on the other hand, is around 20.4% and 2.4% for concentrated and distributed planning, respectively. However, with the increasing number of cores, global skew caused by unbalanced

Table 6.1: Overall CDN Performance Comparison for Test Case  $T_{GPU}$  using Concentrated Receiver Planning

Case (# of cores)	Full Fan-out Tree (TLT)		Tree & Local Mesh (TLM)		Hybrid (Concentrated)				
	Power (mW)	Skew (ps)	Power (mW)	Skew (ps)	Power (mW)	Skew (ps)	Aov (mm <sup>2</sup> )	Skew red. (%)	Power red. (%)
$T2(2)$	476.5	126.0	703.2	68.5	418.1	68.5	0.67	TLT: 45.6 TLM: 0	TLT: 12.3 TLM: 40.5
$T2(4)$	955.3	128.5	1408.7	69.1	815.2	69.5	0.97	TLT: 45.9 TLM: -0.6	TLT: 14.7 TLM: 42.1
$T2(8)$	1990.9	137.1	2816.8	71.3	1609.4	70.4	1.59	TLT: 48.6 TLM: 1.3	TLT: 19.1 TLM: 42.9
$T2(16)$	3882.6	149.5	5641.6	73.6	3197.8	78.5	2.82	TLT: 47.5 TLM: -6.6	TLT: 17.6 TLM: 43.3
$T2(32)$	8080.6	166.7	11303.2	78.8	6374.6	86.7	5.28	TLT: 48.0 TLM: -10.0	TLT: 21.1 TLM: 43.6
$T2(64)$	16163.4	187.5	22636.4	88.4	12728.2	104.5	10.21	TLT: 44.3 TLM: -18.2	TLT: 21.3 TLM: 43.8
Average	5258.2	149.2	7418.2	75.1	4190.6	79.6	3.59	<b>TLT: 46.7</b> <b>TLM: -5.7</b>	<b>TLT: 17.7</b> <b>TLM: 42.7</b>

Table 6.2: Overall CDN Performance Comparison for Test Case  $T_{GPU}$  using Distributed Receiver Planning

Case (# of cores)	Full Fan-out Tree (TLT)		Tree & Local Mesh (TLM)		Hybrid (Distributed)				
	Power (mW)	Skew (ps)	Power (mW)	Skew (ps)	Power (mW)	Skew (ps)	Aov (mm <sup>2</sup> )	Skew red. (%)	Power red. (%)
$T2(2)$	476.5	126.0	703.2	68.5	511.7	71.5	1.59	TLT: 43.3 TLM: -4.4	TLT: -7.3 TLM: 27.2
$T2(4)$	955.3	128.5	1408.7	69.1	1002.4	72.6	2.82	TLT: 43.5 TLM: -5.0	TLT: -4.9 TLM: 28.8
$T2(8)$	1990.9	137.1	2816.8	71.3	1983.8	75.8	5.28	TLT: 44.7 TLM: -6.3	TLT: 0.3 TLM: 29.6
$T2(16)$	3882.6	149.5	5641.6	73.6	3946.6	103.9	10.21	TLT: 30.5 TLM: -41.2	TLT: 0.9 TLM: 31.8
$T2(32)$	8080.6	166.7	11303.2	78.8	7872.2	118.9	20.07	TLT: 28.7 TLM: -50.9	TLT: 2.6 TLM: 30.4
$T2(64)$	16163.4	187.5	22636.4	88.4	15723.4	143.8	39.28	TLT: 23.3 TLM: -62.7	TLT: 2.7 TLM: 30.6
Average	5258.2	149.2	7418.2	75.1	5173.3	97.7	13.31	<b>TLT: 35.7</b> <b>TLM: -28.4</b>	<b>TLT: -0.96</b> <b>TLM: 29.7</b>

propagation distance will eventually generate larger local skew output at the leaf nodes for our proposed solution, especially compared with TLM. This drawback could be mitigated by introducing extra clock Tx's, and hence the different global wireless delays can be compensated.

The power savings, low delay and skew are at a cost of an estimated area overhead of around 5.3 mm<sup>2</sup> and 1.6 mm<sup>2</sup> for distributed and concentrated global planning respectively, including on-chip antennas. Besides, the power overhead of around 24.3 mW and 17.5 mW of clock Tx and Rx, respectively, are traded with the reduction of metallic global interconnect, which will in turn, be sensitive to any variation of the input clock frequency or model size. To sum up, from the experimental results, the concentrated receiver planning can provide a better and more balanced clock performance, whilst preserving high energy efficiency, comparing to the two baseline architectures.

In addition to varying number of cores, we also conduct an experiment to show the impact of changing size of T1 and T2. According to Figure 6.4, the power consumption



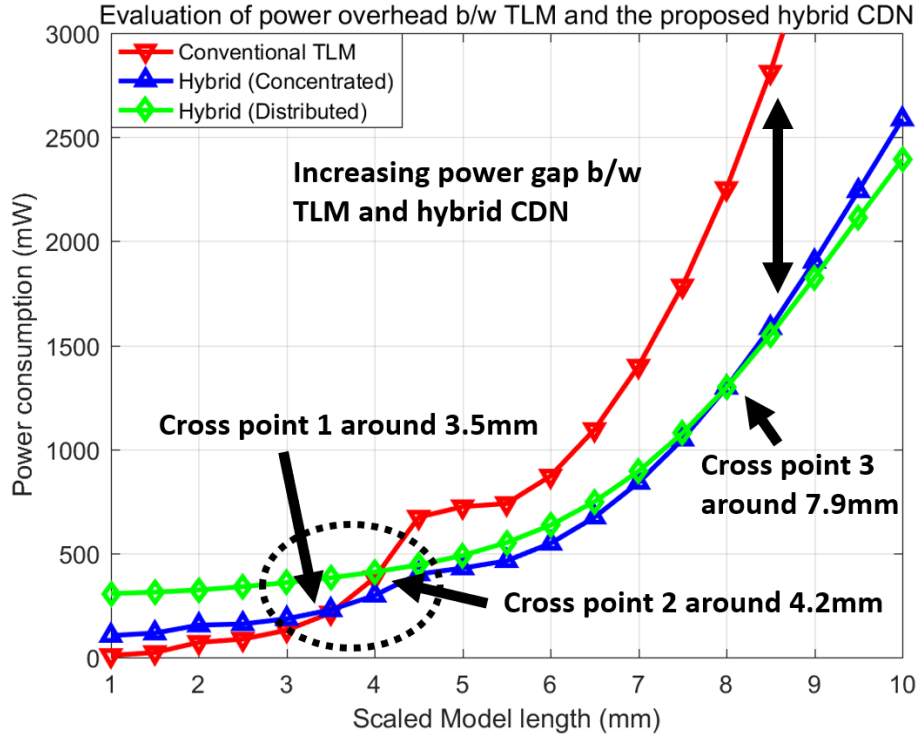


Figure 6.4: Power output in local region with incremental iterations in the proposed test case  $T_{GPU}$ , with an wireless/wire crossing point near 3.5 mm.

of the proposed hybrid CDN becomes more efficient than the conventional TLM after the crossing point near model length of 3.5 mm for a 4-global Rx planning. The hybrid CDN shows a constant and robust characteristic which only change less than 2% of the overall power with the increasing model length. Also, since the global clock Tx and Rx are based on OOK, global power consumption is almost independent of the input frequency, and only depends on the duty cycle of the clock. Hence the proposed CDN provides a viable solution to clock distribution in high performance applications.

## 6.4 Summary

As a conclusion, from the above experimental results, this chapter proposes a novel hybrid clock distribution architecture, consisting of global wireless links and local clock grids using embedded clock Tx and Rx. Comprehensive evaluations, including both global and local CDN, have been carried out to quantify the impacts of local parameters on overall CDN performance. Significant reductions in maximum global delay and power have been found around 28.8% and 42.7%, respectively, comparing to two baseline CDNs TLT and TLM, through the proposed Arm Mali-G77 test case. The skew performance of the proposed hybrid CDN in T2, however, is averagely -28.4% worse than TLM in the worst case scenario, because of the different signal arrival time of the wireless clock receivers. But still, the proposed hybrid CDN shows a up to 46.7% skew

reduction, compared to the data extracted from the fabricated TLT CDN in the test case. Therefore, our proposed hybrid CDN provides a balanced and promising solution to clock distribution with high power efficiency and low clock skew in fully-synchronous many-core systems.

Taking advantage of much-reduced physical interconnect length and the inherent fan-out feature of wireless clock distribution, congested wires could be skipped. The global clock could be distributed to several different dies wirelessly, thus providing an effective way of mitigating total signal propagation delay and wiring power loss. These attributes of our proposed design essentially indicate the potentials of distributing high-quality clock signal in future many-core systems and 3D-ICs.



## Chapter 7

# Conclusions and Future Works

It is worth noticing that the emerging interconnects can be beneficial to many of the conventional global communication schemes. Jumping out of the two-dimensional plane using optical or RF interconnects has offered a promising future for the on-chip and intra-chip global clock distribution. Besides, taking advantage of a broadcast architecture, global CDN would have less dependency on conventional wires, leading to a system with higher robustness.

On the other hand, as a modern digital system will be extremely noisy, challenges have been raised to integrated antenna designers to avoid interference of the noise generated by nearby circuits. In addition, a monolithic system with a compact layout would need a smaller antenna area and a higher carrier frequency eventually, thus requiring smaller process for lower gate delay and higher switching speed. This thesis has given a comprehensive introduction of the proposed work, a hybrid wired-wireless clock distribution network (CDN). The proposed work has been analysed in the following aspects.

### 7.1 Conclusions and Contributions of the Thesis

First, related works and techniques of generating state-of-art CDN solutions are discussed and compared. The challenges of the conventional CDN design methodology have limited the performance of modern synchronous systems in terms of both clock speed and energy efficiency.

Secondly, based on the assumption of using a hybrid CDN architecture, the design algorithms and models are given for local wired CDN design with low metal usage. The estimated results have indicated that the proposed work can have a substantial improvement with regard to global clock delay, and hence reducing the overall clock skew, referring back to Chapter 3.

In addition, the circuit design and the on-chip antenna with EM simulation model for global wireless clock TRx are given in Chapter 4. Based on the simplified OOK modulation and non-coherent demodulation, the proposed circuit has shown its improvement against conventional counterparts in terms of various aspects such as gain, on-off isolation and noise tolerance.

An electro-magnetic experimental model is given in Chapter 5, to quantify the attributes of the proposed design and its nearby logics in terms of noise tolerance and noise interference. Experimental scenarios include a: the clock TRx acts as the noise source and nearby interconnect acts as the victim who picks up common mode noise from the wireless CDN; and b: clock TRx works as a victim and nearby wires act as the aggressor. For scenario a, victim lines located ranging from  $10\ \mu\text{m}$  to  $500\ \mu\text{m}$  suffer from an average interference of around -20 dB, which indicates that almost 10% of the radiated energy will be received by the single-end wires. A solution is provided by using differential signalling, to reduce the common mode noise. Under the same experimental conditions with identical antenna-wire separation, differential wires indicate an average received interference of around -65 dB, which is three orders of magnitude smaller than the conventional single-end implementation. Hence improving overall signal integrity of the nearby wires. For scenario b, victim antennas will only pick up 0.06% of the radiated energy at the nominal operation frequencies of digital logics.

Besides, detailed evaluations of the proposed wireless TRx are presented in Chapter 5. Through SPICE-level simulations on a specific implementation of the proposed circuits in TSMC 65 nm PDK, the performance of the global clock TRx is verified. The proposed clock Tx offers a 47.8 dB on-off isolation of the transmitted clock signal, which provides higher energy efficiency and signal swing at the output of Tx terminal. On the other hand, clock signals with different frequencies can be recovered accordingly at the output of the Rx circuit. For a full swing clock signal, the proposed clock TRx can support a baseband clock signal up to 5 GHz, with 52.7% and 40.0% reduction compared with a conventional global H-tree with identical synchronisation area in global clock power and skew, respectively, taking advantage of wireless interconnect.

Furthermore, a novel industrial test circuit based on the physical extractions of a fabricated/commercialised many-core produce Arm Mali G77 GPU is introduced in Chapter 6. Through the comprehensive test circuit, given the actual chip operation parameters align to the G77 GPU, the proposed work is verified that significant reductions in average clock skew and power have been found around 17.7% and 46.7%, respectively, compared with the extracted baseline CDN TLT. Besides, a promising 42.7% power reduction is found comparing to TLM solution on the proposed test circuit. However, since TLM adopts a more balanced structure for global clock distribution, the TLM implementation exhibits a 5.7% better skew performance, compared with the proposed solution.

## 7.2 Challenges of the Thesis

Despite of the improvements in the hybrid CDN performance, there are still some limitations to be discussed.

First of all, the frequency of the transmitted clock signal is limited by the transistor switching speed. As shown in Figure 5.11 in Chapter 5, for input frequencies higher than 5 GHz (from experimental results), suppose all circuits and testing conditions remain the same, the output clock amplitude descends linearly with the input frequency. Hence, for applications which require full-swing clocks, extra components such as amplifiers and level shifters are necessary to compensate the amplitude decay. Other solutions for mitigating amplitude decay include implementing circuits with more advanced technology node, so that the transistor can share higher transition frequency ( $f_T$ ). However, for state-of-art solutions below 28 nm, certain challenges such as smaller transconductance [132] and mismatch need to be taken into consideration. Techniques such as cascoding, gain-boosting and bootstrapping can be employed to mitigate the negative impacts on the analog design using nanometer CMOS technologies. Therefore, the clock TRx circuits can incorporate the techniques mentioned above for higher maximum clock frequency, which is of utmost importance for high-performance many-core systems.

Secondly, MOSFET interference will increase for smaller processes, hence needs more attention to resolve its impact. As the thermal noise of a MOS transistor is inversely proportional to the transconductance, the thermal noise of single device will increase and thus creating more intrinsic noise for applications like a low-noise amplifier (LNA). FinFET technology however, can served as a solution to this challenge, as it can provide a relatively constant transconductance  $g_m$ . Besides, the increasing parasitics on both BEOL and MEOL wires will bring extra burden especially with the presence of resistance and capacitance mismatch and EM interference. Hence, it is vital to find a balanced point between higher  $f_T$  and less process impacts.

In addition, this work only describes a single clock structure, which transmit and recover a global clock signal across the entire synchronisation area, and is not sufficient for some modern many-core applications that contain several clock domains, either synchronous (including mesochronous) or asynchronous, to fit the working conditions of different integrated cores/blocks. However, this work can provide a robust timing reference, which can be combined with extra components such as frequency dividers or synthesisers. Hence, different blocks can generate multiple clocks, as per their own specifications, which provides more flexibility to synchronisation in many-core systems.

Last but not least, another limitation of the proposed design is the extra cost for the integrated antenna and clock TRx circuit. The estimated area overhead of a single clock Tx, Rx and the antenna can be given by  $0.37 \text{ mm}^2$ ,  $0.08 \text{ mm}^2$  and  $0.009 \text{ mm}^2$ , respectively. Also, an average power overhead is given by 24.3 mW and 17.5 mW, for

a single Tx and Rx, respectively. Hence, for certain applications with tight power and area constraints, the proposed solution will not be the most suitable option. However, for different applications with various load distribution, a cross-section point/curve can be found, to better quantify the efficiency of the proposed design, as per Chapter 5.

In general, the proposed hybrid CDN can offer a competitive performance in terms of clock latency, skew and power consumption, compared with conventional fully-synchronous solutions verified through an industrial many-core design. Besides, referring back to Chapter 1, the proposed design has shown its superiority in clock delay reduction, which is also a key aspect for Mesochronous based GALS many-core system [19]. As a conclusion, the proposed design has offered several attributes which can significantly boost the synchronisation and the performance of many-core systems. The proposed work has shown its potential in terms of both performance and stability. Our proposed hybrid CDN can provide a balanced and promising solution to clock distribution with high power efficiency and low clock skew in future many-core systems.

## 7.3 Potential Future Works

After the current stage, it is essential to have future research organised to keep polishing this topic to make it more practical. For example, CDN components such as the on-chip antenna needs to be optimised with appropriate size and performance. The presence of the local interconnects and logics will have an impact on the antenna performance, hence it is important to find out the antenna performance degradation and methods to compensate this situation.

Since clock receivers are designed to work simultaneously, low power design of the clock receiver is a necessary task in future improvements. By adopting transistors working not only in saturation region but also in sub-threshold region, the power consumption is likely to be reduced significantly, thus this design would become even more appealing and practical. For the continuation towards the detailed future research, short-term and long-term works are organised in the following sections.

### 7.3.1 Short Term Future Works

As per the discussion above, short term future works includes, but not limited to:

1. Continue to review literatures in terms of wireless CDN and wireless interconnect and applying emerging/novel architectures and techniques which could help improve the performance of proposed wireless CDN.
2. Refine on-chip antenna to form an appropriate receiver array with the high communication efficiency and quality.

3. Refine the transfer function model to provide a more accurate prediction of 50% point delay.
4. Optimise the boundary between the global wireless CDN and local wired CDN.
5. Optimise the circuit and algorithm, with the presence of PVT variation in circuits and RC parasitics.

### 7.3.2 Long Term Future Works

Potential long term future works includes, but not limited to:

1. Applying the hybrid CDN onto a multi-source clock structure.
2. Reduce the impact of multi-path propagation of the wireless CDN.
3. Design and test of the proposed hybrid CDN in 3D-IC, in terms of the performance as well as the EM characteristics of the wireless interconnect.
4. Design and test of a low power low-swing clock receiver with promising power reduction.
5. Potential prototype fabrication and testing.





# Bibliography

- [1] Semiconductor Industry Association, “International technology roadmap for semiconductors.” <http://www.itrs2.net/2012-itrs.html>, 2012. [Online; accessed 14 October 2019].
- [2] S.-B. Lee, S.-W. Tam, I. Pefkianakis, S. Lu, M. F. Chang, C. Guo, G. Reinman, C. Peng, M. Naik, L. Zhang, *et al.*, “A scalable micro wireless interconnect structure for cmps,” in *Proceedings of the 15th annual international conference on Mobile computing and networking*, pp. 217–228, ACM, 2009.
- [3] Arizona State University, “Predictive technology model (ptm).” <http://ptm.asu.edu/>, 2012. [Online; accessed 14 October 2019].
- [4] J.-F. Zheng, F. Robertson, E. Mohammad, I. Young, D. Ahn, K. Wada, J. Michel, and L. Kimerling, “On-chip optical clocking signal distribution,” in *Optics in Computing*, p. OWB3, Optical Society of America, 2003.
- [5] S. Chandran, *A Survey of Clock Distribution Techniques Including Optical and RF Networks*. PhD thesis, Auburn University, 2013.
- [6] B. A. Floyd, C.-M. Hung, *et al.*, “Intra-chip wireless interconnect for clock distribution implemented with integrated antennas, receivers, and transmitters,” *IEEE Journal of Solid-State Circuits*, vol. 37, no. 5, pp. 543–552, 2002.
- [7] R. Li, X. Guo, D.-J. Yang, *et al.*, “Wireless clock distribution system using an external antenna,” *IEEE Journal of Solid-State Circuits*, vol. 42, no. 10, pp. 2283–2292, 2007.
- [8] R. Wu, S. Kawai, Y. Seo, N. Fajri, K. Kimura, S. Sato, S. Kondo, T. Ueno, T. Siriburanon, S. Maki, *et al.*, “13.6 a 42gb/s 60ghz cmos transceiver for ieee 802.11 ay,” in *2016 IEEE International Solid-State Circuits Conference (ISSCC)*, pp. 248–249, IEEE, 2016.
- [9] K. D. Boese and A. B. Kahng, “Zero-skew clock routing trees with minimum wirelength,” in *[1992] Proceedings. Fifth Annual IEEE International ASIC Conference and Exhibit*, pp. 17–21, 1992.

- [10] H. Chen *et al.*, “A sliding window scheme for accurate clock mesh analysis,” in *ICCAD-2005. IEEE/ACM International Conference on Computer-Aided Design, 2005.*, pp. 939–946, Nov 2005.
- [11] Semiconductor Industry Association, “International technology roadmap for semiconductors.” <http://www.itrs2.net/itrs-reports.html>, 2014. [Online; accessed 14 October 2019].
- [12] S. S. Sapatnekar, “Rc interconnect optimization under the elmore delay model,” in *31st Design Automation Conference*, pp. 387–391, IEEE, 1994.
- [13] A. I. Abou-Seido, B. Nowak, and C. Chu, “Fitted elmore delay: a simple and accurate interconnect delay model,” in *Proceedings. IEEE International Conference on Computer Design: VLSI in Computers and Processors*, pp. 422–427, Sep. 2002.
- [14] R. Ho, K. W. Mai, and M. A. Horowitz, “The future of wires,” *Proceedings of the IEEE*, vol. 89, no. 4, pp. 490–504, 2001.
- [15] P. Kapur, J. P. McVittie, and K. C. Saraswat, “Technology and reliability constrained future copper interconnects. i. resistance modeling,” *IEEE Transactions on Electron Devices*, vol. 49, no. 4, pp. 590–597, 2002.
- [16] E. G. Friedman, “Clock distribution networks in synchronous digital integrated circuits,” *Proceedings of the IEEE*, vol. 89, no. 5, pp. 665–692, 2001.
- [17] B. Ravelo and A. Jastrzebski, “Modelling of symmetrical distributed clock rc h-tree,” in *International Symposium on Electromagnetic Compatibility-EMC EUROPE*, pp. 1–6, IEEE, 2012.
- [18] Q. Ding, G. Knight, and T. Mak, “An active silicon interposer with low-power hybrid wireless-wired clock distribution network for many-core systems,” *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 28, no. 9, pp. 2042–2054, 2020.
- [19] P. Teehan, M. Greenstreet, and G. Lemieux, “A survey and taxonomy of gals design styles,” *IEEE Design Test of Computers*, vol. 24, no. 5, pp. 418–428, 2007.
- [20] D. M. Chapiro, “Globally-asynchronous locally-synchronous systems,” tech. rep., Stanford Univ CA Dept of Computer Science, 1984.
- [21] K. Y. Yun and R. P. Donohue, “Pausible clocking: a first step toward heterogeneous systems,” in *Proceedings International Conference on Computer Design. VLSI in Computers and Processors*, pp. 118–123, 1996.
- [22] M. R. Greenstreet, “Implementing a stari chip,” in *Proceedings of ICCD ’95 International Conference on Computer Design. VLSI in Computers and Processors*, pp. 38–43, 1995.

- [23] A. Chakraborty and M. R. Greenstreet, "Efficient self-timed interfaces for crossing clock domains," in *Ninth International Symposium on Asynchronous Circuits and Systems, 2003. Proceedings.*, pp. 78–88, 2003.
- [24] A. Rashid, N. Sultana, M. R. Khan, and T. Kikkawa, "Efficient design of integrated antennas on si for on-chip wireless interconnects in multi-layer metal process," *Japanese Journal of Applied Physics*, vol. 44, no. 4S, p. 2756, 2005.
- [25] K. Ravindran, A. Kuehlmann, and E. Sentovich, "Multi-domain clock skew scheduling," in *ICCAD-2003. International Conference on Computer Aided Design (IEEE Cat. No. 03CH37486)*, pp. 801–808, IEEE, 2003.
- [26] P. Cunningham, M. Swinnen, and S. Wilcox, "Clock concurrent optimization," *Azuro white paper*, 2009.
- [27] C. Lenzen, T. Locher, and R. Wattenhofer, "Clock synchronization with bounded global and local skew," in *2008 49th Annual IEEE Symposium on Foundations of Computer Science*, pp. 509–518, 2008.
- [28] A. Karkar, T. Mak, K.-F. Tong, and A. Yakovlev, "A survey of emerging interconnects for on-chip efficient multicast and broadcast in many-cores," *IEEE Circuits and Systems Magazine*, vol. 16, no. 1, pp. 58–72, 2016.
- [29] D. C. Keezer and V. K. Jain, "Clock distribution strategies for wsi: A critical survey," in *1991 Proceedings, International Conference on Wafer Scale Integration*, pp. 277–283, IEEE, 1991.
- [30] S. Deb, A. Ganguly, P. P. Pande, B. Belzer, and D. Heo, "Wireless noc as interconnection backbone for multicore chips: Promises and challenges," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 2, no. 2, pp. 228–239, 2012.
- [31] D. Tse and P. Viswanath, *Fundamentals of wireless communication*. Cambridge university press, 2005.
- [32] H.-M. Hsu, T.-H. Lee, and C.-J. Hsu, "Millimeter-wave transmission line in 90-nm cmos technology," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 2, no. 2, pp. 194–199, 2012.
- [33] A. Carpenter, J. Hu, J. Xu, M. Huang, H. Wu, and P. Liu, "Using transmission lines for global on-chip communication," *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, vol. 2, no. 2, pp. 183–193, 2012.
- [34] A. Karkar, N. Dahir, K. Tong, T. Mak, A. Yakovlev, *et al.*, "Hybrid wire-surface wave architecture for one-to-many communication in networks-on-chip," in *2014 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1–4, IEEE, 2014.

- [35] R. Morris, E. Jolley, and A. K. Kodi, "Extending the performance and energy-efficiency of shared memory multicores with nanophotonic technology," *IEEE Transactions on Parallel and Distributed Systems*, vol. 25, no. 1, pp. 83–92, 2013.
- [36] P. Dong, Y.-K. Chen, T. Gu, L. L. Buhl, D. T. Neilson, and J. H. Sinsky, "Re-configurable 100 gb/s silicon photonic network-on-chip," *IEEE/OSA Journal of Optical Communications and Networking*, vol. 7, no. 1, pp. A37–A43, 2015.
- [37] D. A. Miller, "Device requirements for optical interconnects to silicon chips," *Proceedings of the IEEE*, vol. 97, no. 7, pp. 1166–1185, 2009.
- [38] D. Chang, F. Yu, Z. Xiao, Y. Li, N. Stojanovic, C. Xie, X. Shi, X. Xu, and Q. Xiong, "Fpga verification of a single qc-ldpc code for 100 gb/s optical systems without error floor down to ber of 10<sup>-15</sup>," in *Optical Fiber Communication Conference*, p. OTuN2, Optical Society of America, 2011.
- [39] H. B. Kommuru and H. Mahmoodi, "Asic design flow tutorial using synopsys tools," *Nano-Electronics & Computing Research Lab, School of Engineering, San Francisco State University San Francisco, CA, Spring*, 2009.
- [40] J. Lu and B. Taskin, "From rtl to gdsii: An asic design course development using synopsys® university program," in *2011 IEEE International Conference on Microelectronic Systems Education*, pp. 72–75, IEEE, 2011.
- [41] E. McLellan, "The alpha axp architecture and 21064 processor," *IEEE Micro*, vol. 13, no. 3, pp. 36–47, 1993.
- [42] E. Anderson, J. Brooks, C. Grassl, and S. Scott, "Performance of the cray t3e multiprocessor," in *Proceedings of the 1997 ACM/IEEE conference on Supercomputing*, pp. 1–17, ACM, 1997.
- [43] N. A. Kurd *et al.*, "A multigigahertz clocking scheme for the pentium(r) 4 microprocessor," *IEEE Journal of Solid-State Circuits*, vol. 36, pp. 1647–1653, Nov 2001.
- [44] M. A. El-Moursy and E. G. Friedman, "Exponentially tapered h-tree clock distribution networks," *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 13, no. 8, pp. 971–975, 2005.
- [45] K. Matsumaru, "Reflection coefficient of e-plane tapered waveguides," *IRE transactions on Microwave theory and techniques*, vol. 6, no. 2, pp. 143–149, 1958.
- [46] F. Zhu, W. Hong, W.-F. Liang, J.-X. Chen, X. Jiang, P.-P. Yan, and K. Wu, "A low-power low-cost 45-ghz ook transceiver system in 90-nm cmos for multi-gb/s transmission," *IEEE Transactions on Microwave Theory and Techniques*, vol. 62, no. 9, pp. 2105–2117, 2014.

- [47] G. R. Wilke and R. Murgai, "Design and analysis of" tree+ local meshes" clock architecture," in *8th International Symposium on Quality Electronic Design (ISQED'07)*, pp. 165–170, IEEE, 2007.
- [48] T. Singh, S. Rangarajan, D. John, R. Schreiber, S. Oliver, R. Seahra, and A. Schaefer, "2.1 zen 2: The amd 7nm energy-efficient high-performance x86-64 microprocessor core," in *2020 IEEE International Solid- State Circuits Conference - (ISSCC)*, pp. 42–44, 2020.
- [49] G. Shamanna *et al.*, "Scalable, sub-1w, sub-10ps clock skew, global clock distribution architecture for intelcorei7/i5/i3 microprocessors," in *2010 Symposium on VLSI Circuits*, pp. 83–84, June 2010.
- [50] S. Rusu and G. Singer, "The first ia-64 microprocessor," *IEEE Journal of Solid-State Circuits*, vol. 35, no. 11, pp. 1539–1544, 2000.
- [51] D. W. Dobberpuhl, R. T. Witek, R. Allmon, R. Anglin, D. Bertucci, S. Britton, L. Chao, R. A. Conrad, D. E. Dever, B. Gieseke, *et al.*, "A 200-mhz 64-b dual-issue cmos microprocessor," *IEEE Journal of Solid-State Circuits*, vol. 27, no. 11, pp. 1555–1567, 1992.
- [52] W. Liu, G. Chen, Y. Wang, and H. Yang, "Modeling and optimization of low power resonant clock mesh," in *The 20th Asia and South Pacific Design Automation Conference*, pp. 478–483, IEEE, 2015.
- [53] A. Farshidi, *Power and Timing Driven Optimal Gate, Clock Buffer and Clock Wire Sizing in High Performance Digital Integrated Circuits*. PhD thesis, University of Calgary, 2016.
- [54] Z. Chen, H. Gu, Y. Yang, L. Bai, and H. Li, "A power efficient and compact optical interconnect for network-on-chip," *IEEE Computer Architecture Letters*, vol. 13, no. 1, pp. 5–8, 2013.
- [55] J. Xue, A. Garg, B. Ciftcioglu, J. Hu, S. Wang, I. Savidis, M. Jain, R. Berman, P. Liu, M. Huang, *et al.*, "An intra-chip free-space optical interconnect," in *ACM SIGARCH Computer Architecture News*, vol. 38, pp. 94–105, ACM, 2010.
- [56] J. M. Miller, N. De Beaucoudrey, P. Chavel, J. Turunen, and E. Cambril, "Design and fabrication of binary slanted surface-relief gratings for a planar optical interconnection," *Applied optics*, vol. 36, no. 23, pp. 5717–5727, 1997.
- [57] D. Huang, T. Sze, A. Landin, R. Lytel, and H. L. Davidson, "Optical interconnects: out of the box forever?," *IEEE Journal of Selected Topics in Quantum Electronics*, vol. 9, no. 2, pp. 614–623, 2003.
- [58] Semiconductor Industry Association, "International technology roadmap for semiconductors." <http://www.itrs2.net/2013-its.html>, 2013. [Online; accessed 14 October 2019].

- [59] R. Meade, J. S. Orcutt, K. Mehta, O. Tehar-Zahav, D. Miller, M. Georgas, B. Moss, C. Sun, Y.-H. Chen, J. Shainline, *et al.*, “Integration of silicon photonics in bulk cmos,” in *2014 Symposium on VLSI Technology (VLSI-Technology): Digest of Technical Papers*, pp. 1–2, IEEE, 2014.
- [60] R. K. Dokania and A. B. Apsel, “Analysis of challenges for on-chip optical interconnects,” in *Proceedings of the 19th ACM Great Lakes symposium on VLSI*, pp. 275–280, ACM, 2009.
- [61] K. Kim, B. A. Floyd, J. L. Mehta, H. Yoon, C.-M. Hung, D. Bravo, T. O. Dickson, X. Guo, R. Li, N. Trichy, *et al.*, “On-chip antennas in silicon ics and their application,” *IEEE Transactions on Electron Devices*, vol. 52, no. 7, pp. 1312–1323, 2005.
- [62] K. Okada, N. Li, K. Matsushita, K. Bunsen, R. Murakami, A. Musa, T. Sato, H. Asada, N. Takayama, S. Ito, *et al.*, “A 60-ghz 16qam/8psk/qpsk/bpsk direct-conversion transceiver for ieee802. 15.3 c,” *IEEE Journal of Solid-State Circuits*, vol. 46, no. 12, pp. 2988–3004, 2011.
- [63] S. G. Wilson, *Digital Modulation and Coding*. Upper Saddle River, NJ, USA: Prentice-Hall, Inc., 1996.
- [64] T. S. Rappaport *et al.*, *Wireless communications: principles and practice*, vol. 2. prentice hall PTR New Jersey, 1996.
- [65] J. Lee, Y. Chen, and Y. Huang, “A low-power low-cost fully-integrated 60-ghz transceiver system with ook modulation and on-board antenna assembly,” *IEEE Journal of Solid-State Circuits*, vol. 45, no. 2, pp. 264–275, 2010.
- [66] A. Ganguly, K. Chang, S. Deb, P. P. Pande, B. Belzer, and C. Teuscher, “Scalable hybrid wireless network-on-chip architectures for multicore systems,” *IEEE Transactions on Computers*, vol. 60, no. 10, pp. 1485–1502, 2010.
- [67] F. Minami and M. Takano, “Clock tree synthesis based on rc delay balancing,” in *1992 Proceedings of the IEEE Custom Integrated Circuits Conference*, pp. 28.3.1–28.3.4, May 1992.
- [68] A. Rajaram and D. Z. Pan, “Meshworks: An efficient framework for planning, synthesis and optimization of clock mesh networks,” in *2008 Asia and South Pacific Design Automation Conference*, pp. 250–257, March 2008.
- [69] D. Shan *et al.*, “Resonant clock mega-mesh for the ibm z13tm,” in *2015 Symposium on VLSI Circuits (VLSI Circuits)*, pp. C322–C323, June 2015.
- [70] X. Hu and M. R. Guthaus, “Distributed lc resonant clock grid synthesis,” *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 59, pp. 2749–2760, Nov 2012.

- [71] C. Ryu *et al.*, “A three-dimensional stacked-chip star-wiring interconnection for a digital noise-free and low-jitter i/o clock distribution network,” *IEEE Microwave and Wireless Components Letters*, vol. 16, pp. 651–653, Dec 2006.
- [72] D. Chung *et al.*, “Chip-package hybrid clock distribution network and dll for low jitter clock delivery,” *IEEE Journal of Solid-State Circuits*, vol. 41, pp. 274–286, Jan 2006.
- [73] A. Karkar *et al.*, “A survey of emerging interconnects for on-chip efficient multicast and broadcast in many-cores,” *IEEE Circuits and Systems Magazine*, vol. 16, pp. 58–72, Firstquarter 2016.
- [74] W. C. Elmore, “The transient response of damped linear networks with particular regard to wideband amplifiers,” *Journal of applied physics*, vol. 19, no. 1, pp. 55–63, 1948.
- [75] M. Balch, *Complete digital design: a comprehensive guide to digital electronics and computer system architecture*. McGraw-Hill Education, 2003.
- [76] J. L. Wyatt, “Circuit analysis, simulation and design,” 1987.
- [77] F. Minami and M. Takano, “Clock tree synthesis based on rc delay balancing,” in *Proc. IEEE Custom Integrated Circuits Conf*, pp. 28–3, 1992.
- [78] Guoqing Chen and E. G. Friedman, “An rlc interconnect model based on fourier analysis,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 24, pp. 170–183, Feb 2005.
- [79] N. Ida, *Engineering electromagnetics*, vol. 2. Springer, 2000.
- [80] A. B. K. et al., “An analytical delay model for rlc interconnects,” *IEEE Trans. Comput. Aided Des. Integr. Circuits Syst.*, vol. 16, no. 12, pp. 1507–1514, 1997.
- [81] Ting-Hai Chao, Yu-Chin Hsu, Jan-Ming Ho, and A. B. Kahng, “Zero skew clock routing with minimum wirelength,” *IEEE Transactions on Circuits and Systems II: Analog and Digital Signal Processing*, vol. 39, no. 11, pp. 799–814, 1992.
- [82] X. Zhao, J. Minz, and S. K. Lim, “Low-power and reliable clock network design for through-silicon via (tsv) based 3d ics,” *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 1, no. 2, pp. 247–259, 2011.
- [83] K. D. Boese and A. B. Kahng, “Zero-skew clock routing trees with minimum wirelength,” in *[1992] Proceedings. Fifth Annual IEEE International ASIC Conference and Exhibit*, pp. 17–21, IEEE, 1992.
- [84] M. A. Jackson, A. Srinivasan, and E. S. Kuh, “Clock routing for high-performance ics,” in *27th ACM/IEEE Design Automation Conference*, pp. 573–579, IEEE, 1990.



- [85] N. H. Weste and K. Eshraghian, "Principles of cmos vlsi design: a systems perspective," *NASA STI/Recon Technical Report A*, vol. 85, 1985.
- [86] A. Rahmani, M. Haghbayan, A. Kanduri, A. Y. Weldezion, P. Liljeberg, J. Plosila, A. Jantsch, and H. Tenhunen, "Dynamic power management for many-core platforms in the dark silicon era: A multi-objective control approach," in *2015 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED)*, pp. 219–224, 2015.
- [87] M. Haghbayan, A. Rahmani, A. Y. Weldezion, P. Liljeberg, J. Plosila, A. Jantsch, and H. Tenhunen, "Dark silicon aware power management for manycore systems under dynamic workloads," in *2014 IEEE 32nd International Conference on Computer Design (ICCD)*, pp. 509–512, 2014.
- [88] S. Lu, H. Yu, X. Wang, Q. Zhang, F. Li, Z. Liu, and F. Ning, "Clustering method of raw meal composition based on pca and kmeans," in *2018 37th Chinese Control Conference (CCC)*, pp. 9007–9010, 2018.
- [89] S. Zahra, M. A. Ghazanfar, A. Khalid, M. A. Azam, U. Naeem, and A. Prugel-Bennett, "Novel centroid selection approaches for kmeans-clustering based recommender systems," *Information sciences*, vol. 320, pp. 156–189, 2015.
- [90] T. Kanungo, D. M. Mount, N. S. Netanyahu, C. D. Piatko, R. Silverman, and A. Y. Wu, "An efficient k-means clustering algorithm: analysis and implementation," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 881–892, 2002.
- [91] X. Zhao and S. K. Lim, "Power and slew-aware clock network design for through-silicon-via (tsv) based 3d ics," in *2010 15th Asia and South Pacific Design Automation Conference (ASP-DAC)*, pp. 175–180, 2010.
- [92] G. E. Tellez and M. Sarrafzadeh, "Minimal buffer insertion in clock trees with skew and slew rate constraints," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 16, no. 4, pp. 333–342, 1997.
- [93] S. Hu, C. J. Alpert, J. Hu, S. K. Karandikar, Z. Li, W. Shi, and C. N. Sze, "Fast algorithms for slew-constrained minimum cost buffering," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 26, no. 11, pp. 2009–2022, 2007.
- [94] X. Zhao, D. L. Lewis, H.-H. S. Lee, and S. K. Lim, "Pre-bond testable low-power clock tree design for 3d stacked ics," in *2009 IEEE/ACM International Conference on Computer-Aided Design-Digest of Technical Papers*, pp. 184–190, IEEE, 2009.
- [95] S. Saini, A. M. Kumar, S. Veeramachaneni, and M. B. Srinivas, "An alternative approach to buffer insertion for delay and power reduction in vlsi interconnects," in *2010 23rd International Conference on VLSI Design*, pp. 411–416, 2010.

- [96] S. M. Reddy, G. R. Wilke, and R. Murgai, "Analyzing timing uncertainty in mesh-based clock architectures," in *Proceedings of the Design Automation Test in Europe Conference*, vol. 1, pp. 1–6, March 2006.
- [97] W. Liu *et al.*, "Modeling and optimization of low power resonant clock mesh," in *The 20th Asia and South Pacific Design Automation Conference*, pp. 478–483, Jan 2015.
- [98] A. Rajaram and D. Z. Pan, "Meshworks: a comprehensive framework for optimized clock mesh network synthesis," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 29, no. 12, pp. 1945–1958, 2010.
- [99] P. Hammarlund, A. J. Martinez, A. A. Bajwa, D. L. Hill, E. Hallnor, H. Jiang, M. Dixon, M. Derr, M. Hunsaker, R. Kumar, R. B. Osborne, R. Rajwar, R. Singhal, R. D'Sa, R. Chappell, S. Kaushik, S. Chennupaty, S. Jourdan, S. Gunther, T. Piazza, and T. Burton, "Haswell: The fourth-generation intel core processor," *IEEE Micro*, vol. 34, no. 2, pp. 6–20, 2014.
- [100] K. Yamada and N. Oda, "Statistical corner conditions of interconnect delay (corner lpe specifications)," in *Asia and South Pacific Conference on Design Automation, 2006.*, pp. 6 pp.–, 2006.
- [101] C. H. Doan, S. Emami, A. M. Niknejad, and R. W. Brodersen, "Millimeter-wave cmos design," *IEEE Journal of solid-state circuits*, vol. 40, no. 1, pp. 144–155, 2005.
- [102] B. Razavi, "Cmos transceivers for the 60-ghz band," *Proc. IEEE RFIC*, pp. 11–13, 2006.
- [103] B. Razavi, *Design of Analog CMOS Integrated Circuits*. New York, NY, USA: McGraw-Hill, Inc., 1 ed., 2001.
- [104] Q. Ding, B. J. Fletcher, and T. Mak, "Globally wireless locally wired (glowilow): A clock distribution network for many-core systems," in *2018 IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 1–5, May 2018.
- [105] Q. Ding and T. Mak, "Hybrid interconnect network for on-chip low-power clock distribution," *Electronics Letters*, vol. 55, no. 5, pp. 244–246, 2019.
- [106] X. Yu, H. Rashtian, S. Mirabbasi, P. P. Pande, and D. Heo, "An 18.7-gb/s 60-ghz ook demodulator in 65-nm cmos for wireless network-on-chip," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 62, no. 3, pp. 799–806, 2015.
- [107] C. W. Byeon, J. J. Lee, K. C. Eun, and C. S. Park, "A 60 ghz 5 gb/s gain-boosting ook demodulator in 0.13  $\mu\text{m}$  cmos," *IEEE Microwave and Wireless Components Letters*, vol. 21, no. 2, pp. 101–103, 2011.

- [108] S. Chen and M. Ker, "A new schmitt trigger circuit in a 0.13-/spl mu/m 1/2.5-v cmos process to receive 3.3-v input signals," *IEEE Transactions on Circuits and Systems II: Express Briefs*, vol. 52, pp. 361–365, July 2005.
- [109] O. O. Olaode, W. D. Palmer, and W. T. Joines, "Effects of meandering on dipole antenna resonant frequency," *IEEE Antennas and Wireless Propagation Letters*, vol. 11, pp. 122–125, 2012.
- [110] S.-L. Zuo, Z.-Y. Zhang, and J.-W. Yang, "Planar meander monopole antenna with parasitic strips and sleeve feed for dvb-h/lte/gsm850/900 operation in the mobile phone," *IEEE Antennas and Wireless Propagation Letters*, vol. 12, pp. 27–30, 2012.
- [111] O. O. Olaode, W. D. Palmer, and W. T. Joines, "Characterization of meander dipole antennas with a geometry-based, frequency-independent lumped element model," *IEEE Antennas and Wireless Propagation Letters*, vol. 11, pp. 346–349, 2012.
- [112] J. Montero-de Paz, E. Ugarte-Muñoz, L. E. García-Muñoz, I. C. Mayorga, and D. Segovia-Vargas, "Meander dipole antenna to increase cw thz photomixing emitted power," *IEEE Transactions on Antennas and Propagation*, vol. 62, no. 9, pp. 4868–4872, 2014.
- [113] L. Covert and J. Lin, "Simulation and measurement of a heatsink antenna: a dual-function structure," *IEEE Transactions on Antennas and Propagation*, vol. 54, no. 4, pp. 1342–1349, 2006.
- [114] J. Park, D. Choi, and W. Hong, "Millimeter-wave phased-array antenna-in-package (aip) using stamped metal process for enhanced heat dissipation," *IEEE Antennas and Wireless Propagation Letters*, vol. 18, no. 11, pp. 2355–2359, 2019.
- [115] A. L. Loke, D. Yang, T. T. Wee, J. L. Holland, P. Isakanian, K. Rim, S. Yang, J. S. Schneider, G. Nallapati, S. Dundigal, *et al.*, "Analog/mixed-signal design challenges in 7-nm cmos and beyond," in *2018 IEEE Custom Integrated Circuits Conference (CICC)*, pp. 1–8, IEEE, 2018.
- [116] H. . P. Wong, D. J. Frank, P. M. Solomon, C. H. J. Wann, and J. J. Welser, "Nanoscale cmos," *Proceedings of the IEEE*, vol. 87, no. 4, pp. 537–570, 1999.
- [117] B. Greskamp, S. R. Sarangi, and J. Torrellas, "Threshold voltage variation effects on aging-related hard failure rates," in *2007 IEEE International Symposium on Circuits and Systems*, pp. 1261–1264, 2007.
- [118] M. Li, K. Khalaf, C. Li, V. Vojkan, M. Ingels, A. Bourdoux, P. Wambacq, J. Craninckx, and L. Van Der Perre, "Signal processing challenges for emerging digital intensive and digitally assisted transceivers with deeply scaled technology (invited)," in *SiPS 2013 Proceedings*, pp. 324–329, 2013.

- [119] S. K. Goel, S. Adham, M. Wang, J. Chen, T. Huang, A. Mehta, F. Lee, V. Chickermane, B. Keller, T. Valind, S. Mukherjee, N. Sood, J. Cho, H. H. Lee, J. Choi, and S. Kim, "Test and debug strategy for tsmc cowos stacking process based heterogeneous 3d ic: A silicon case study," in *2013 IEEE International Test Conference (ITC)*, pp. 1–10, 2013.
- [120] J. Jayabalan, V. Chidambaram, S. L. P. Siang, W. Xiangyu, J. Ming Ching, and S. Bhattacharya, "Active through-silicon interposer based 2.5d ic design, fabrication, assembly and test," in *2019 IEEE 69th Electronic Components and Technology Conference (ECTC)*, pp. 587–593, 2019.
- [121] A. B. M. H. Rashid, S. Watanabe, and T. Kikkawa, "Crosstalk isolation of monopole integrated antenna on si for ulsi wireless interconnect," in *Proceedings of the IEEE 2003 International Interconnect Technology Conference*, June 2003.
- [122] C. Hu, B. Tai, and A. Yang, "Meander-line folded monopole design for umts-hsdap-based data-card applications," *IEEE Antennas and Wireless Propagation Letters*, vol. 7, pp. 279–282, 2008.
- [123] B. A. Floyd, Chih-Ming Hung, and K. K. O, "Intra-chip wireless interconnect for clock distribution implemented with integrated antennas, receivers, and transmitters," *IEEE Journal of Solid-State Circuits*, vol. 37, pp. 543–552, May 2002.
- [124] M. Donno, A. Ivaldi, L. Benini, and E. Macii, "Clock-tree power optimization based on rtl clock-gating," in *Proceedings of the 40th annual Design Automation Conference*, pp. 622–627, ACM, 2003.
- [125] V. Sukumaran *et al.*, "Low-cost thin glass interposers as a superior alternative to silicon and organic interposers for packaging of 3-d ics," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 2, pp. 1426–1433, Sep. 2012.
- [126] A. Kannan *et al.*, "Enabling interposer-based disintegration of multi-core processors," in *Proceedings of the 48th International Symposium on Microarchitecture, MICRO-48*, (New York, NY, USA), pp. 546–558, ACM, 2015.
- [127] G. H. Loh *et al.*, "Interconnect-memory challenges for multi-chip, silicon interposer systems," in *Proceedings of the 2015 International Symposium on Memory Systems, MEMSYS '15*, (New York, NY, USA), pp. 3–10, ACM, 2015.
- [128] J. Song, S. Park, S. Kim, J. J. Kim, and J. Kim, "Active silicon interposer design for interposer-level wireless power transfer technology for high-density 2.5-d and 3-d ics," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 6, no. 8, pp. 1148–1161, 2016.

- 
- [129] Q. Zou and Y. Xie, “Compact models and model standard for 2.5 d and 3d integration,” in *Proceedings of SLIP (System Level Interconnect Prediction) on System Level Interconnect Prediction Workshop*, pp. 1–7, ACM, 2014.
  - [130] S.-Y. Huang and C.-C. Zheng, “Die-to-die clock skew characterization and tuning for 2.5 d ics,” in *2016 IEEE 25th Asian Test Symposium (ATS)*, pp. 221–226, IEEE, 2016.
  - [131] Arm, “The mali-g77 graphics processors.” <https://www.arm.com/products/silicon-ip-multimedia/gpu/mali-g77>, 2019. [Online; accessed 18 October 2019].
  - [132] W. Sansen, “Analog design challenges in nanometer cmos technologies,” in *2007 IEEE Asian Solid-State Circuits Conference*, pp. 5–9, 2007.