

1 **Physics-Informed Neural Networks for Passive Scalar Emission**
2 **and Transport**

3 Joshua Ian Rawden,* Christina Vanderwel, and Sean Symon

4 *Aerodynamics & Flight Mechanics Research Group,*

5 *University of Southampton, Southampton SO17 1BJ, UK*

6 (Dated: January 19, 2026)

Abstract

Accurate modeling of harmful pollutant concentrations is an important field of interest for protecting public health and the environment. In this study, physics-informed neural networks (PINNs) are applied to a low Reynolds number, time-averaged cylinder wake, which interacts with a variety of different passive scalar regimes. The PINN reconstructs both the time-averaged velocity and scalar concentration from limited measurements. In addition to satisfying the incompressible Reynolds-averaged Navier-Stokes (RANS) equations, the reconstructed fields must also obey the time-averaged advection-diffusion equation. This was done to extend the applicability of mean field reconstruction and to lay the foundations for PINNs to be used in more complex passive scalar modeling in future studies, particularly the prediction of pollutant behavior in urban environments. It was found that the PINN could successfully reconstruct the flow fields on a macro-scale in almost all scenarios, having considerable success with first-order quantities and managing to accurately infer the spatial structure of the unknown closure term of the advection-diffusion equation in all cases. When reconstructing scalar source characteristics, the PINN could identify the source location and size in many cases. The results allowed basic guidelines for optimal sensor placement to be described to inform future studies, and suggests ways in which this technique could be developed to contribute to future air quality forecasting models.

* jr1g20@soton.ac.uk

7 I. INTRODUCTION

8 A passive scalar is any concentration that does not influence the physical behavior of a
9 fluid, meaning the velocity fields of a flow with and without a passive scalar are identical.
10 In recent years, the importance of understanding how passive scalars are transported and
11 dispersed in urban environments has increased tremendously. This includes the way we
12 model the transfer of heat, particulates, and disease within closed spaces, as well as outdoor
13 public health in ever-growing cities as the world continues to industrialize. According to the
14 World Health Organization (WHO), 99% of the population breathe air that exceeds WHO
15 guideline limits [1]. The first step toward reducing this number is developing the ability
16 to reliably model the evolution of pollutant plumes. To do so, accurate reconstruction of
17 the passive scalar properties at/near the source(s) is critical, as an understanding of the
18 origin and typical paths of harmful particulates and gases in the atmosphere allows public
19 health measures to be enacted by authorities. These would aim to minimize the impact on
20 the public and appropriately direct the enforcement of environmental protection laws when
21 necessary.

22 Experiments and high-fidelity computational fluid dynamics (CFD) simulations have been
23 performed to improve understanding of passive scalars in urban pollution, such as the sim-
24 ulations of Refs. [2–4] and the scaled experiments of Refs. [5–7]. While experiments pro-
25 vide physically justified measurements at a discrete number of sampling points, and CFD
26 methods give a modeled solution on a dense mesh, both approaches are expensive and time-
27 consuming to set up while having their own sets of limitations. Sparse field measurements
28 and the complexity of large-scale urban geometries present challenges due to the sharp spa-
29 tial variations of pollutant concentrations, making optimal placement of sensors difficult
30 [8, 9]. This presents an opportunity for data-driven flow reconstruction methods to demon-
31 strate their effectiveness where the use of CFD or experimentation alone is challenging. The
32 most established data-driven method used in mean-flows is variational data assimilation,
33 which uses observations from experimental/high-fidelity CFD to inform and improve the re-
34 sults of a cheap, low-fidelity CFD simulation (e.g. RANS), often reconstructing features and
35 multi-scale structures not present in either dataset separately [10, 11]. Data assimilation is
36 a fundamental component of modern weather models, and is employed in a variety of roles
37 from large-scale storm monitoring to regional-scale pollution modeling [12, 13]. Data assim-

38 ilation has also been a useful tool for reconstructing passive scalar fields on smaller scales
39 [9, 14–17], but its effectiveness is highly dependent on the locations of individual sensors and
40 the monitoring paradigm being used, especially when reconstructing source characteristics
41 as shown by Ref. [9].

42 Physics-informed neural networks (PINNs) are a recent development in the flow recon-
43 struction community, leveraging the rapid progress in machine learning technologies to create
44 a new class of neural network that has a physical understanding of the problem it is solv-
45 ing embedded within its objective function. PINNs were formally introduced by Ref. [18],
46 where they were used in a variety of physical problems to demonstrate their novelty and
47 flexibility. More recently, Ref. [19] used PINNs to reconstruct the mean fields and turbu-
48 lence statistics around a low-Reynolds number cylinder flow. Reference [20] introduced the
49 Spalart-Allmaras turbulence model to the PINN when solving a periodic hill flow problem
50 where it was found to outperform traditional data assimilation methods. Reference [21]
51 applied PINNs to DNS data of an unsteady, buoyancy-driven Rayleigh-Bénard flow and had
52 success with reconstructing the temperature field. Reference [22] apply PINNs to an urban
53 environment, but did not consider any passive scalar. Several other studies have been per-
54 formed that also attempt to solve fluid-based problems using PINNs, suggesting that there
55 is an appetite in this technology to push the field of fluid dynamics forward [23–30].

56 The novelty of this study is to identify the size and location of unknown passive scalar
57 sources by reconstructing the mean velocity and concentration fields with limited measure-
58 ments using PINNs. This is performed on a well-researched, canonical flow case under a
59 wide range of passive scalar regimes defined using different boundary conditions (BCs) and
60 Schmidt numbers. The PINN’s effectiveness will be determined by its ability to reconstruct
61 the mean velocities, the mean scalar concentration, the unknown closure field of the advec-
62 tion diffusion equation, and the properties of the scalar source. If the PINN is successful
63 in completing these objectives, it demonstrates that the methodology may be applicable to
64 more complex problems in the near future. It will also suggest that PINNs can complement
65 with the established variational data assimilation techniques for tackling these problems and
66 may play a foundational role in future air quality forecasting frameworks [31].

67 The rest of the paper is organized as follows. Section II introduces the physical problem
68 and equations, as well as the mathematics and principles that govern the PINN. Section
69 III describes the generation of the CFD reference data used to train the PINN. Section

70 IV describes the PINN models in more detail, listing the architectures and training proce-
 71 dures used. Section IV D summarizes the variations in the independent variables for the
 72 cases, followed by Sec. V which presents the results. Finally, Sec. VI concludes the study,
 73 summarizing the reconstructive qualities of PINNs in passive scalar problems.

74 II. PHYSICS & PROBLEM DEFINITION

75 This section will introduce the physical problem, beginning with introducing the govern-
 76 ing equations in Sec. II A, the time-averaged versions of which are presented and described
 77 more rigorously for the passive scalar and velocity in Secs. II B and II C, respectively, fol-
 78 lowed by an explanation of the PINN’s loss function in Sec. II D.

79 A. Governing Equations

80 We start with the vector form of the advection-diffusion equation and the full incom-
 81 pressible Navier-Stokes equations

$$\frac{\partial c}{\partial t} + \mathbf{u} \cdot \nabla c - \frac{1}{Re \cdot Sc} \cdot \nabla^2 c = 0, \quad (1a)$$

$$\nabla \cdot \mathbf{u} = 0, \quad (1b)$$

$$\frac{\partial \mathbf{u}}{\partial t} + \mathbf{u} \cdot \nabla \mathbf{u} + \nabla p - \frac{1}{Re} \cdot \nabla^2 \mathbf{u} = 0, \quad (1c)$$

84 where \mathbf{u} is velocity, p is pressure, c is passive scalar concentration, $Re = \frac{UL}{\nu}$, and $Sc = \frac{\nu}{\alpha}$. U
 85 is the characteristic velocity, L is a characteristic length scale, ν is the kinematic viscosity
 86 of the fluid, and α is the molecular diffusivity of the passive scalar. The Schmidt number
 87 represents the ratio between momentum diffusivity and molecular diffusivity. Certain flows
 88 concerning small heat fluctuations may also be considered passive scalar flows, as buoyancy
 89 effects are negligible. In this context, the Prandtl number can be used instead of the Schmidt
 90 number. A flow with a high Sc will generally result in a field with higher gradients of scalar
 91 concentration, as it follows the velocity field more strictly. In contrast, a low Sc allows the
 92 scalar to spread regardless of the advective fluxes, resulting in a generally more diffuse and
 93 homogeneous concentration field. For some physical intuition, CO_2 has a molecular Schmidt
 94 number in air of 0.94, and air itself has a Prandtl number of 0.708 [32]. Naphthalene (an

95 organic compound used to make plastics and pesticides) has a higher molecular Schmidt
 96 number of 2.35 [33]. These values inform the ranges used in this study.

97 B. Time-Averaged Advection-Diffusion Equation

98 To transform the system into a set of time-averaged equations, we use the Reynolds
 99 decomposition

$$\mathbf{u} = \bar{\mathbf{u}} + \mathbf{u}', \quad (2)$$

100 applied to u , v , p and c , where u is the streamwise velocity and v is the transverse velocity.
 101 Here, the $\bar{[\]}$ terms are the mean flow quantities and the $[\]'$ terms are the fluctuations from
 102 the mean flow.

103 Applying Eq. 2 to Eq. 1a and averaging yields the time-averaged advection-diffusion
 104 equation

$$\mathcal{P}_1 = \bar{u} \cdot \frac{\partial \bar{c}}{\partial x} + \bar{v} \cdot \frac{\partial \bar{c}}{\partial y} - \frac{1}{Re \cdot Sc} \cdot \left(\frac{\partial^2 \bar{c}}{\partial x^2} + \frac{\partial^2 \bar{c}}{\partial y^2} \right) + g, \quad (3a)$$

105

$$g = \nabla \cdot (\overline{\mathbf{u}'c'}) = \frac{\partial \overline{u'c'}}{\partial x} + \frac{\partial \overline{v'c'}}{\partial y}, \quad (3b)$$

106 where g is an unknown closure term, as it requires knowledge of the fluctuations and is
 107 defined as the divergence of the mean turbulent scalar flux. Here, we assume a 2D problem
 108 (explained in Sec. III), hence only 2 components of velocity are present. In Eq. 3a, the first
 109 two terms are the scalar advection terms, and the 3rd is the diffusion term. Since g depends
 110 on the turbulent scalar fluxes $\overline{u'c'}$ and $\overline{v'c'}$ (which are absent from the mean flow data),
 111 any reconstruction of the macro-level structure of g indicates that the PINN is inferring
 112 the divergence of the fluxes from the mean concentration fields, which is a useful property
 113 for turbulence modeling. The notation \mathcal{P}_i indicates that the equation will be reduced to a
 114 residual which the PINN uses to inform its gradient descent algorithm (see Sec. IID).

115 C. RANS Equations with Solenoidal Forcing

116 Similarly to Eq. 3a, applying Eq. 2 to Eqs. 1b and 1c then averaging transforms them
 117 into

$$\mathcal{P}_2 = \frac{\partial \bar{u}}{\partial x} + \frac{\partial \bar{v}}{\partial y}, \quad (4a)$$

$$\mathcal{P}_3 = \bar{u} \cdot \frac{\partial \bar{u}}{\partial x} + \bar{v} \cdot \frac{\partial \bar{u}}{\partial y} + \frac{\partial(\bar{p} - \phi)}{\partial x} - \frac{1}{Re} \cdot \left(\frac{\partial^2 \bar{u}}{\partial x^2} + \frac{\partial^2 \bar{u}}{\partial y^2} \right) + f_{su}, \quad (4b)$$

$$\mathcal{P}_4 = \bar{u} \cdot \frac{\partial \bar{v}}{\partial x} + \bar{v} \cdot \frac{\partial \bar{v}}{\partial y} + \frac{\partial(\bar{p} - \phi)}{\partial y} - \frac{1}{Re} \cdot \left(\frac{\partial^2 \bar{v}}{\partial x^2} + \frac{\partial^2 \bar{v}}{\partial y^2} \right) + f_{sv}. \quad (4c)$$

In addition to the mean-flow transformation, new variables ϕ , f_{su} and f_{sv} have been introduced. The terms f_{su} and f_{sv} are derived from the vector $\mathbf{f} = \nabla \cdot \mathbf{R}$ (which emerges after Reynolds-averaging) through the process described below. Tensor \mathbf{R} contains the Reynolds stress components $R_{ij} = \overline{u'_i u'_j}$, and so taking Eqs. 4 and 3a, we have 4 equations for 9 unknowns: \bar{c} , \bar{u} , \bar{v} , \bar{p} , $\overline{u'u'}$, $\overline{v'v'}$, $\overline{u'v'}$, $\overline{u'c'}$ and $\overline{v'c'}$. This is an underdetermined system which cannot be solved for a unique solution.

To resolve this, a modified formulation of the RANS equations is to be used, which was first proposed by Ref. [10] for use in variational data assimilation algorithms, and has since been used successfully in studies of mean flow with PINNs [19, 20]. This modification performs a Helmholtz decomposition on \mathbf{f} , separating it into a potential component and a solenoidal (divergence-free) component:

$$\mathbf{f} = \nabla \phi + \mathbf{f}_s, \quad (5a)$$

$$\nabla \cdot \mathbf{f}_s = 0. \quad (5b)$$

The potential forcing ϕ is combined with the pressure gradient term of equations 4b and 4c. As described in detail in Sec. 2.6 of Ref. [10], this results in the true pressure field being unrecoverable after reconstruction; therefore, the reconstructed pressure fields will be omitted from the results.

The addition of Eq. 5b to the roster as:

$$\mathcal{P}_5 = \frac{\partial f_{su}}{\partial x} + \frac{\partial f_{sv}}{\partial y}, \quad (6)$$

increases the count to 5 equations for 9 unknowns, which is still under determined. One way to make the PDE system tractable is to provide the \bar{u} , \bar{v} and \bar{c} fields on a sufficiently resolved grid, allowing the velocity data to fully satisfy Eq. 4a. This reduces the system to 4 equations (Eqs. 4b, 4c, 5b and 3a) for 5 unknowns: $\bar{p} - \phi$, f_{su} , f_{sv} , $\overline{u'c'}$ and $\overline{v'c'}$. As discussed

141 in [19] and [20], a sufficiently dense velocity data grid will close the RANS equations and
 142 result in a unique solution; the expectation is that this principle can also be extended to \bar{c} for
 143 the advection-diffusion system. However, this is infeasible in realistic applications and many
 144 experimental techniques. Therefore, we rely on the deterministic nature of the idealized
 145 velocity and passive scalar fields to allow the amount of necessary data to be reduced to a
 146 practical amount while still producing an accurate reconstruction.

147 D. PINN Problem

148 With some generalization, neural networks aim to “learn” a relationship between pairs
 149 of input/output data with or without training data serving as a reference. For PINNs, this
 150 is effectively the minimization of an objective loss function:

$$\mathcal{L}(\boldsymbol{\theta}; \mathbf{y}) = w_D \mathcal{L}_D(\boldsymbol{\theta}; \mathbf{y}) + w_B \mathcal{L}_B(\boldsymbol{\theta}; \mathbf{y}) + w_P \mathcal{L}_P(\boldsymbol{\theta}; \mathbf{y}), \quad (7a)$$

$$\mathcal{L}_D(\boldsymbol{\theta}; \mathbf{y}) = \sum_{i=1}^N \left[\frac{1}{|\mathbf{x}^D|} \sum_{j=1}^{|\mathbf{x}^D|} [y_i(\mathbf{x}^D_j) - y_{i,j}^D]^2 \right], \quad (7b)$$

$$\mathcal{L}_B(\boldsymbol{\theta}; \mathbf{y}) = \sum_{i=1}^N \left[\frac{1}{|\mathbf{x}^B|} \sum_{j=1}^{|\mathbf{x}^B|} [y_i(\mathbf{x}^B_j) - y_{i,j}^B]^2 \right], \quad (7c)$$

$$\mathcal{L}_P(\boldsymbol{\theta}; \mathbf{y}) = \sum_{i=1}^N \left[\frac{1}{|\mathbf{x}^P|} \sum_{j=1}^{|\mathbf{x}^P|} [\mathcal{P}_i|_{\mathbf{x}_j^P}]^2 \right], \quad (7d)$$

154 with respect to the set of network parameters $\boldsymbol{\theta}$, where \mathbf{y} is the vector of N reconstructed
 155 fields $\mathbf{y} = [y_1(\mathbf{x}) \ y_2(\mathbf{x}) \ \dots \ y_N(\mathbf{x})]^T$ at each domain coordinate $\mathbf{x} = [x \ y]^T$, and \mathbf{x} can
 156 represent the set of data point, boundary point or collocation point coordinates (\mathbf{x}^D , \mathbf{x}^B ,
 157 \mathbf{x}^P) depending on which loss term the PINN is solving for. Here, $\mathcal{L}_D(\boldsymbol{\theta}; \mathbf{y})$ represents the
 158 data loss, i.e. the summation of L_2 errors between the training and reconstructed values
 159 across each data point within the PINN domain [19, 34] where y^D is the reference data field.
 160 Next, $\mathcal{L}_B(\boldsymbol{\theta}; \mathbf{y})$ is the boundary loss with y^B as the boundary condition values for each field,
 161 which sums the L_2 errors across the domain boundary to enforce the specified boundary
 162 conditions. Finally, $\mathcal{L}_P(\boldsymbol{\theta}; \mathbf{y})$ is the physics loss, with \mathcal{P}_i as the residual for each governing
 163 equation evaluated at each collocation point within the domain. $\mathcal{L}_P(\boldsymbol{\theta}; \mathbf{y})$ is the primary
 164 means by which the PINN is able to enforce the physics of the problem when producing

165 a solution. w_D , w_B and w_P are the respective data, boundary and physics loss weights,
 166 which can be specified by the user to bias the PINN to conform more toward the data
 167 or the equations. Training to minimize $\mathcal{L}(\boldsymbol{\theta}; \mathbf{y})$ is performed iteratively through numerical
 168 optimization of the parameter space.

169 III. GENERATION OF CFD REFERENCE DATA

170 OpenFOAM [35, 36] was used to simulate a canonical 2D cylinder flow. This flow is
 171 well described in literature and serves as the foundation onto which the novel PINN is ap-
 172 plied. Three different flow cases were considered, the boundary conditions/initial conditions
 173 (BCs/ICs) for which are given in Sec. III B: a cylinder with passive scalar emission from the
 174 surface (hereafter referred to as the *Surface* case), a cylinder with passive scalar emission
 175 from a point source upstream of the cylinder (hereafter referred to as the *Upstream* case),
 176 and a cylinder with passive scalar emission from a point source within the recirculation
 177 region of the cylinder wake (hereafter referred to as the *Wake* case). From these three cases,
 178 datasets were generated for each case at Schmidt numbers of 0.1, 0.3, 1.0, 3.0, and 10.0 to
 179 give variation in the diffusive behavior of the scalar.

180 A. Domain & Meshing

181 The CFD domain extends from $-15 < \frac{x}{d} < 40$, $-10 < \frac{y}{d} < 10$ with a cylinder obstacle of
 182 diameter $d = 1$ located at $(0, 0)$. The z dimensionality of the domain is superficial, as the
 183 mesh only contains a single cell along the z axis, and the front/back boundary conditions are
 184 set to empty. ‘Empty’ is a type of boundary condition in OpenFOAM reserved for reduced-
 185 dimension cases, as OpenFOAM can only operate on 3D domains [37]. This boundary
 186 condition type, along with the single cell in the z -direction, enforces uniformity in the
 187 spanwise solution and makes the domain effectively 2D. The mesh was generated using
 188 BlockMesh within OpenFOAM. The mesh is composed of 6 individual blocks that were
 189 refined independently: the pre-block, 4 obstacle blocks surrounding the cylinder, and the
 190 post-block. The obstacle blocks were the only blocks to be refined, and were refined using
 191 the simpleGrading function with a refinement factor of 80 approaching the cylinder surface.
 192 This resulted in a total cell count of 51,000. BlockMesh’s internal mesh checking function

193 verified that the cell aspect ratios and skewnesses are satisfactory.

194 For the surface case, the passive scalar is injected into the domain by specifying a Dirichlet
195 boundary condition of magnitude 1.0 on the cylinder surface (see Sec. III B). The point
196 source upstream and wake cases use an additional OpenFOAM class called a cellZone to
197 apply special properties to a group of cells without imposing strict boundary conditions that
198 influence other properties of the flow. These cells are assigned a constant value of 5.0 for
199 the scalar concentration which is maintained throughout the simulation. For the upstream
200 case, the cellZone was placed at (-1.6, 0.3) and applied to 2 cells. For the wake case, the
201 cellZone was placed at (1.0, 0.3) and applied to 3 cells. This ensures that the point source
202 cases produce scalar concentration fields with sufficient magnitude to be reconstructed by
203 the PINN, since the total number of cells emitting passive scalar is far lower than in the
204 surface emission case. The difference in source strength does not affect the results since the
205 investigation that compares these cases is concerned with the impact of specifying BCs and
206 not the strength of the source.

207 B. Boundary/Initial Conditions & Solver

208 The solver and solution parameters for all three cases were the same except for details
209 of the scalar injection. The pisoFoam solver within OpenFOAM was used to solve the
210 unsteady Navier-Stokes equations transiently without a turbulence model. The solver used
211 a second-order backward time differentiation scheme, a Gauss linear gradient scheme, and
212 a Gauss linear corrected Laplacian scheme. For all cases, a uniform inlet is prescribed. No-
213 slip Dirichlet boundary conditions are applied to the cylinder surface, symmetric boundary
214 conditions at the upper and lower boundaries, and advective conditions to the outlet. The
215 Reynolds number based on the diameter was $Re = 150$.

216 The solution was advanced in time up to 500 units with $\Delta t = 0.001$ to give 500,000
217 individual time steps. At a freestream velocity of 1, this corresponded to a maximum
218 CFL number of around 0.15 throughout the simulation which ensured temporal stability
219 and justified the use of the backward time scheme. In the surface case, the scalar was
220 emitted from the cylinder surface at a constant rate throughout the entire simulation. In
221 the upstream and wake cases, the point-source emission cellZone was switched on at $t = 100$
222 to allow the velocity fields to settle into a periodic vortex shedding state before the scalar

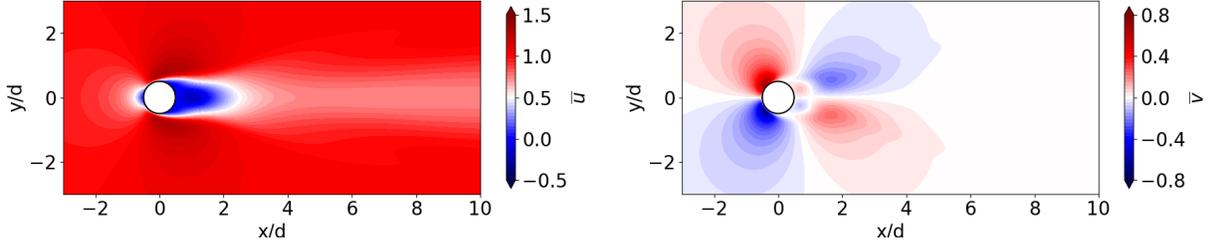


FIG. 1: Reference fields of \bar{u} and \bar{v} .

223 was introduced.

224 Following the simulations, the data were time-averaged to provide mean fields for the
 225 PINN. Time-averaging was performed after the fluid had reached periodic stability. This
 226 was found to be at $t \approx 150$; therefore, averaging was performed for 350 snapshots at intervals
 227 of 1 between $150 < t < 500$. The g fields were computed using the turbulent scalar fluxes
 228 taken from the transient data from this same period, and the results were then averaged to
 229 give g as per Eq. 3b. To validate the CFD data, the frequency of shedding was calculated
 230 to be $St = \frac{fL}{U} = 0.186$, which is the expected frequency for $Re = 150$ [38]. The mean fields
 231 were confirmed to be in agreement with those of other studies [10, 19], and are presented
 232 in Figs. 1 and 2. Note that only \bar{u} , \bar{v} , and \bar{c} are explicitly given to the PINN, as g is an
 233 inferred term and will be reconstructed by the PINN.

234 IV. DETAILS OF PINN MODELS & EXPERIMENTS

235 This section introduces the procedure used to create the PINN models used throughout
 236 the investigation, giving their architectures and hyperparameters, the training protocol, and
 237 the independent variables that were varied over the course of the study.

238 A. PINN Model Generation Program

239 The program used to generate the PINN models was written in Python and was developed
 240 from the code made publicly available in Ref. [19]. The program makes use of the DeepXDE
 241 package [34] - a library that interfaces with commonly used deep-learning packages.

242 To provide the PINN with the reference CFD training data, the data at each cell within
 243 the CFD mesh were linearly interpolated to a uniform grid with spacing h , which indicates

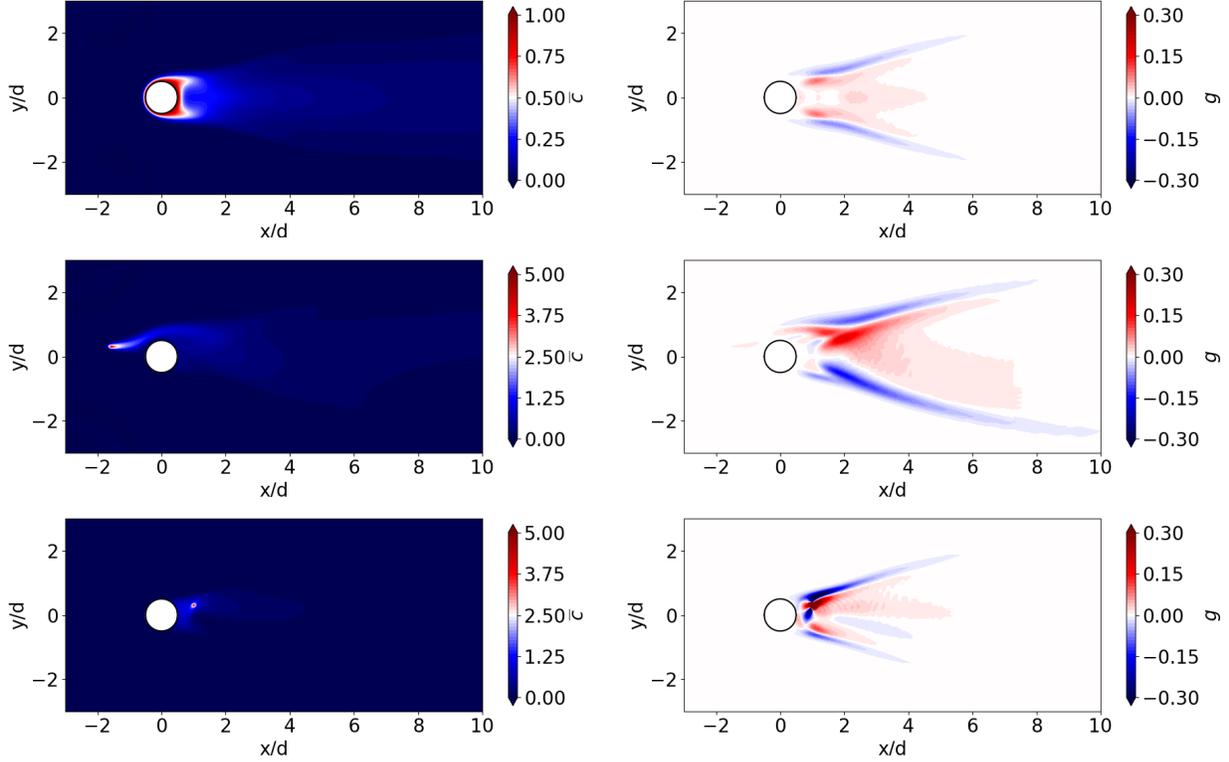


FIG. 2: Surface, upstream and wake case CFD fields at $Sc = 1$; \bar{c} provided as a reference field to the PINN and g to be reconstructed by the PINN.

244 the number of points per cylinder diameter, e.g. $h = 2$ represents a data point spacing
 245 of $d/2$. h with no subscript indicates the velocity and scalar data were supplied on the
 246 same grid. The PINN models used the following BCs for the cylinder surface for all cases:
 247 Dirichlet no-slip conditions for $\bar{\mathbf{u}}$ and \mathbf{f}_s , Neumann conditions for g , Dirichlet $\bar{c} = 1$ in the
 248 surface case, and Neumann conditions for \bar{c} in the upstream and wake cases. All models used
 249 the same 5 residual equations within their partial differential equation (PDE) models with
 250 partial derivatives computed using automatic differentiation (AD). This allows the PINN to
 251 be continuous and differentiable, meaning that no underlying grid/mesh is necessary.

252 B. PINN Architecture & Parameters

253 All PINN models used the same architecture: an input layer of width 2, 7 hidden layers of
 254 width 100, and an output layer of width 7, which maps $\begin{bmatrix} x & y \end{bmatrix}^T \rightarrow \begin{bmatrix} \bar{u} & \bar{v} & \bar{p} - \phi & \bar{c} & f_{su} & f_{sv} & \bar{g} \end{bmatrix}^T$.
 255 The impact of varying the number of nodes per layer and number of layers is presented in

256 the appendix, but the network size was the same for all models considered in the results
 257 section. Each layer used the *tanh* activation function and the *Glorot uniform* initialiser
 258 function. The relative loss weights were the same as in Ref. [20], and were set as: $w_D = 10$,
 259 $w_B = 10$, $w_P = 1$.

260 The full domain collocation points were defined using a custom function to cluster points
 261 near the cylinder surface and along the flow centerline. This was done using a lognormal
 262 distribution in x and a normal distribution in y . Figure 3 shows the extent of the full
 263 PINN domain and the collocation point distribution described above, with the dimensions
 264 normalized by the cylinder diameter d . The custom distribution is expected to improve
 265 convergence as the PINN can more effectively resolve the PDEs in regions of the domain that
 266 experience high gradients. The number of domain collocation points was around 58,000, and
 267 the number of boundary collocation points was set to 1500. 50 testing points were randomly
 268 selected to update the loss function each iteration to give a fairly uniform distribution
 269 throughout the domain without hindering the optimizer loop with excessive passes. The
 270 testing points do not influence the loss function seen by the optimizer and therefore do
 271 not inform the PINN training process. Including testing collocation points is useful for
 272 monitoring the testing loss to determine whether the PINN model is overfitting to the
 273 training data and failing to adhere to the physics, which was not an issue for any models
 274 used in this study.

275 C. PINN Training Protocol

276 All models were trained in two stages: an initial stage using the *Adam* optimizer [39] which
 277 is based on stochastic gradient descent, followed by a secondary longer training stage using
 278 the *L-BFGS-B* optimizer [40, 41]. This order was decided due to *Adam*'s ability to explore
 279 and smoothen the PINN's rough initial parameter space, which prepares the model for a
 280 longer training period with *L-BFGS-B* to significantly reduce the losses [34]. This procedure
 281 is consistent with other studies that investigated similar fluid cases [19, 20, 24, 25, 42]. The
 282 PINN was trained using *Adam* for 25,000 epochs at a learning rate of 0.001, followed by
 283 *L-BFGS-B* up to a maximum of 60,000 total epochs, or until the optimizer converged.
 284 The convergence criterion is set as the DeepXDE default of $\leq 1 \times 10^{-8}$ for the maximum
 285 component of the projected descent gradient. A typical loss graph is given in Fig. 4. Training

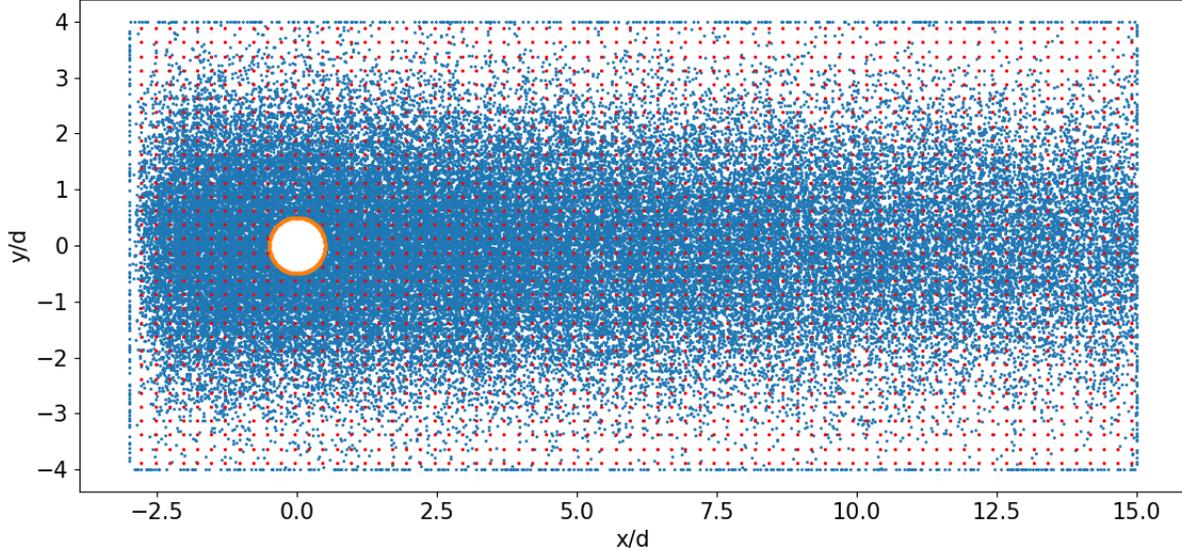


FIG. 3: Collocation and typical data point distribution: domain points (blue), boundary points (orange), data points (red).

286 was performed using the *Iridis X* High-Performance Computing cluster at the University of
 287 Southampton, utilizing NVIDIA A100 enterprise GPUs. The total training time was heavily
 288 dependent on the number of data points presented to the PINN, but a typical training cycle
 289 for a single model was on the order of 2-4 hours. In total, the time taken to train all models
 290 used in this investigation was around 7-9 days. All investigations required many models to
 291 be trained with slight variations to the network conditions. This process was automated
 292 using a Python loop.

293 The code used to train the PINN models and some example models used in this study
 294 can be found at <https://github.com/roardon/PassiveScalarPINNs> or at <https://doi.org/10.5258/SOTON/D3810>.
 295

296 D. Experimental Conditions & Independent Variables

297 The investigations performed can be divided into two categories: investigation of the
 298 full domain and investigation of the cropped domain as seen in Fig. 5. The separation
 299 of the investigations allows reconstruction of the macro-level structure of the fields to be
 300 considered independently of the identification of the source characteristics. The full domain
 301 investigations were performed on all three flow cases and consider a smaller region of the
 302 CFD domain by cropping unimportant (freestream) parts of the flow. This was done to

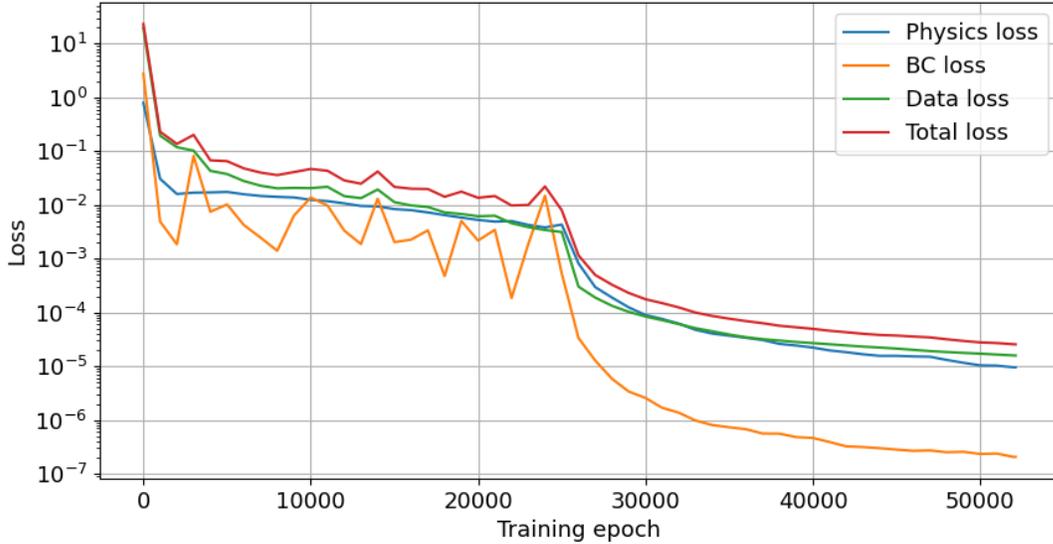


FIG. 4: Typical loss throughout the training cycle; optimiser switches from *Adam* to *L-BFGS-B* at 25,000 epochs.

303 allow the data/collocation points to be placed within the more sensitive areas of the PINN
 304 domain to make the most efficient use of computational resources. The full domain for the
 305 PINN extends from $-3 < \frac{x}{d} < 15$, $-4 < \frac{y}{d} < 4$ (green in Fig. 5). The investigations in
 306 the cropped domain were performed on the two point source cases (upstream & wake) and
 307 reduce the reconstruction region to a smaller window of size 1×1 , with its location dependent
 308 on the flow case but ensuring that the passive scalar point source was located at the center
 309 of this window (blue in Fig. 5). In the cropped domain investigations, the scalar data were
 310 supplied within the aforementioned 1×1 window, while the velocity data were supplied
 311 within a larger window of size 6×6 from $-3 < \frac{x}{d} < 3$, $-3 < \frac{y}{d} < 3$ which includes the entire
 312 cylinder obstacle (see orange in Fig. 5). Focus was placed on identification of the source
 313 characteristics (size, location) because they represent important properties influencing the
 314 structure and evolution of a passive scalar plume. The reduced domain size was chosen
 315 because it allows for a more efficient use of computational resources, which means that more
 316 models could be trained in the same amount of time, resulting in a more robust investigation
 317 where a wider range of values for h could be tested. We found that reducing the size of
 318 the domain had no impact on the conclusions of this experiment compared to using the
 319 full domain. Passive scalar flows are highly sensitive to changes in the source properties;
 320 therefore, this knowledge is critical when attempting to accurately predict dispersion over

321 large, complex domains or long time horizons [9]. Within these two categories, investigations
 322 were performed to assess the impact of adjusting a single PINN parameter or flow constant
 323 on the quality of the reconstruction. Details of the cases are presented in Table I.

Case	Variable range	Variables fixed
Full-domain		
Data sparsity	$1 \leq h \leq 10$	$Sc = 1.0$
Schmidt number	$0.1 \leq Sc \leq 10$	$h = 4$
Cropped-domain		
Data sparsity	$3 \leq h_c \leq 100$	$Sc = 1.0, h_u = 6$
Schmidt number	$0.1 \leq Sc \leq 10$	$h = 10$

TABLE I: Case details.

324

325

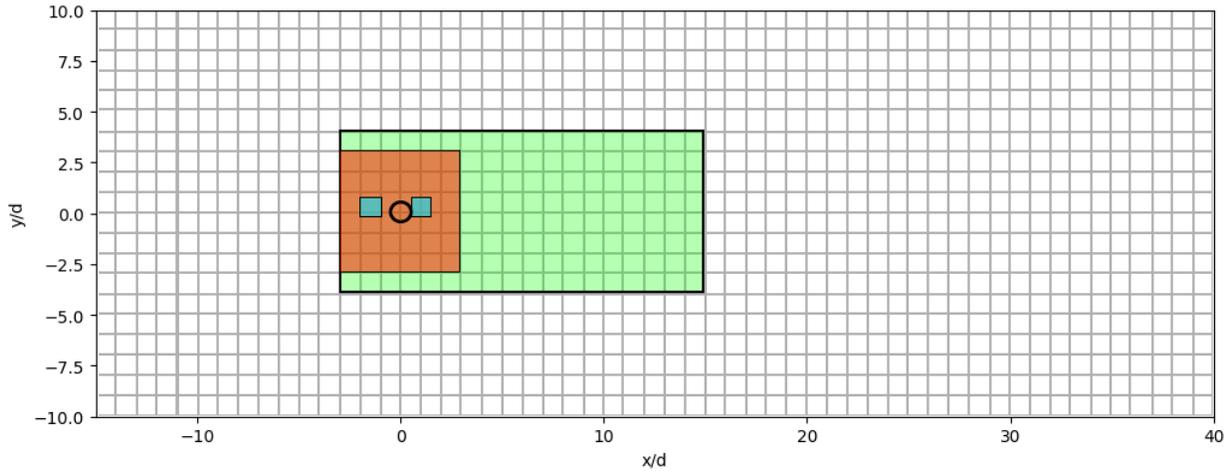


FIG. 5: Subdomains of CFD domain. Green: full PINN domain. Orange: cropped-domain velocity data range. Blue: cropped-domain scalar data ranges.

326 In the cropped-domain data sparsity investigation, the data were supplied on separate
 327 velocity and scalar grids whose densities could be varied independently, given by h_u and
 328 h_c , respectively. The ‘variable range’ for this investigation refers to the range of h_c used,
 329 while the velocity grid is kept constant at $h_u = 6$. An investigation was performed for the
 330 inverse case (fixed h_c with varying h_u), but this did not have any noticeable impact on scalar
 331 reconstruction, and so it will be omitted from the results. A PINN model was trained for

332 the combinations of parameters listed in Table I, giving a total number of models used in
 333 the results of 116.

334 V. RESULTS

335 The results will be presented as ordered in Table IV D, starting with the full-domain cases
 336 in Sec. V A, followed by the cropped-domain investigations in Sec. V B.

337 A. Full-Domain Investigations

338 Within the full-domain investigations, a volume-weighted L_2 -norm error metric is used
 339 as the primary measure for assessing the accuracy of the reconstructed fields relative to the
 340 baselines. This metric is defined as:

$$\varepsilon_2 = \sqrt{\sum_{i \in \Omega} [\omega_i \times (y_i(\mathbf{x}_i) - y_i^D)^2]}, \quad (8)$$

341 where ω_i is an individual cell volume inside the domain Ω . This metric gives a single scalar
 342 value as a measure of the reconstructed field’s accuracy, and is not normalized with respect
 343 to the scale of the original data. Therefore L_2 norm values for different fields (e.g. \bar{u} , \bar{c}) are
 344 not intended to be compared to one another.

345 1. Velocity Reconstruction & Influence of Passive Scalar Data

346 Since the velocity fields are identical among all cases, the influence of data sparsity on
 347 these fields is presented first in Fig. 6. Reference [19] performed a similar analysis, so this
 348 section aims in part to validate the the current PINN framework as well as to extend the
 349 application to passive scalar cases, as per the investigation objectives.

350 Figure 6 shows that a data grid density of $h = 1$ is inadequate for reconstructing the
 351 velocity fields around the cylinder as the PINN does not have enough information in the
 352 regions of importance to produce a unique solution. These ‘regions of importance’ are
 353 typically identified to be areas of the domain that exhibit high spatial gradients of the
 354 reconstructed fields. Furthermore, we have found by comparing this model to an additional
 355 model with a specified inlet BC of $\bar{u} = 1$, $\bar{v} = 0$, that this results in a substantial improvement

356 to the reconstruction both upstream of the cylinder and within the near-wake when $h = 1$.
 357 Despite this, the influence of specifying the inlet BCs was found to rapidly diminish beyond
 358 $h \geq 2$, suggesting that under these conditions the training data and governing equations
 359 are sufficient to produce a unique solution, a similar result to Ref. [19]. For these reasons,
 360 we believe that inlet BCs are unnecessary for most of the cases in our results and rely on
 361 external knowledge about the flow case that may not be available in practical settings, but
 362 we maintain their importance when using extremely sparse training data.

363 Increasing toward $h = 10$, we see greatly improved quality of reconstruction at the expense
 364 of requiring a higher measurement window resolution and increased computational resources.
 365 The velocity absolute error fields' structures and magnitudes are very similar to those for
 366 the first-order velocity statistics presented in Ref. [19] for $h = 50$. This suggests that
 367 reconstruction accuracy saturates with data at around $h = 10$, as the error reduces minimally
 368 beyond this point. We see that there is reduction in error of 2 orders of magnitude between
 369 just $1 < h < 2$. This result is consistent with Ref. [19], and we believe that this rapid
 370 increase in accuracy is due to both the convergence to a unique flow solution in the presence
 371 of additional data, as well as the PINN more accurately reconstructing the size and shape
 372 of the cylinder recirculation bubble.

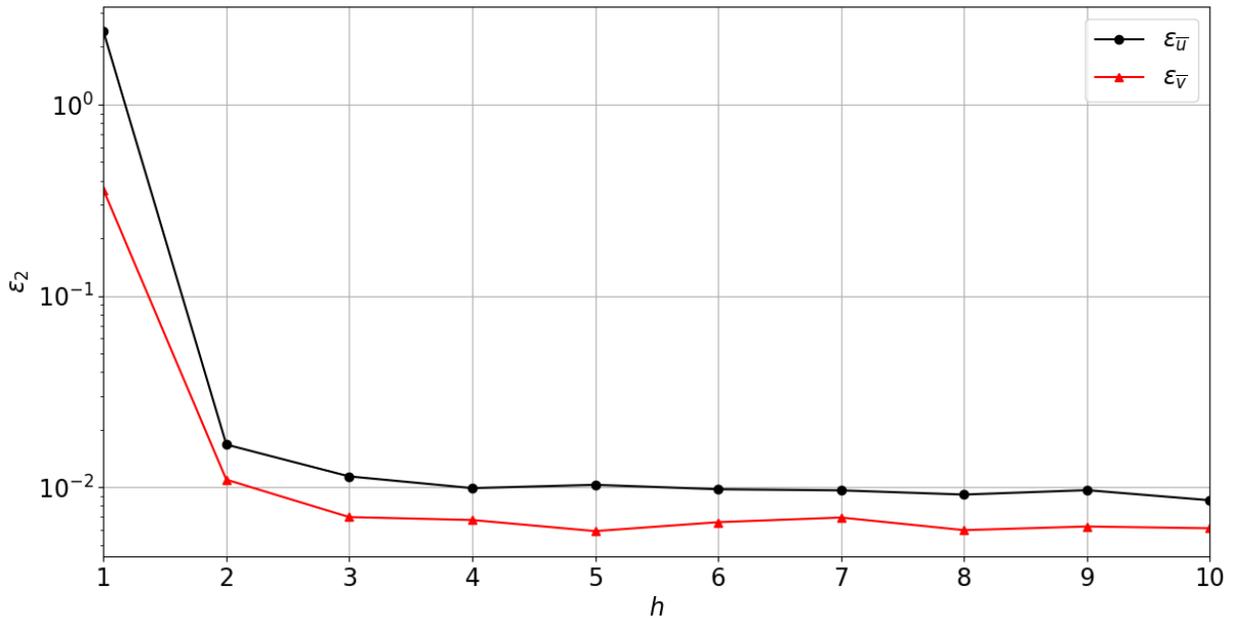


FIG. 6: L_2 norm versus h for velocity fields \bar{u} and \bar{v} .

373 Next, the effects of introducing the passive scalar transport system are investigated when

374 applied to the velocity PINN model. This was done to determine whether knowledge of
 375 the passive scalar distribution within the domain would feed back into the RANS equa-
 376 tion system via the advection-diffusion equation to improve the velocity reconstruction. In
 377 this investigation, the velocity data grid is kept constant at $h_u = 1$, while the scalar data
 378 ($Sc = 1.0$) is supplied on a separate grid that varies in density between $1 < h_c < 10$. These
 379 results help to satisfy the objectives of accurately reconstructing the passive scalar fields,
 380 and are shown in Fig. 7, where it is clear that the addition of the advection-diffusion equa-
 381 tion system and uniformly distributed scalar data produces more accurate reconstructions
 382 of the velocity fields. This effect is most pronounced in the surface case, where the passive
 383 scalar is specified through a boundary condition, but minor improvements can also be seen
 384 for the upstream and wake cases. It is likely that the surface case benefits the most from
 385 this change since accurate information about the passive scalar is provided on the cylinder
 386 surface - the region of the domain in which the velocity error is highest. Within the flow
 387 reconstruction framework, information that improves the solution of one equation can im-
 388 plicitly influence the solutions of other equations in the system. It is probable that solving
 389 the advection-diffusion equation (Eq. 3a) near the cylinder surface resulted in more accurate
 390 reconstructions of \bar{u} and \bar{v} by constraining the range of values that these quantities can take,
 391 since \bar{u} and \bar{v} appear in Eq. 3a, but \bar{c} does not appear in Eq. 4. It is also noted that the
 392 error of \bar{u} reconstruction was affected more than \bar{v} by this change. This result illustrates
 393 that knowledge of passive scalar boundary conditions can greatly improve reconstructions
 394 of local velocities.

395 2. Effect of Data Sparsity on Full-Domain Scalar Reconstruction

396 The remainder of the results will focus only on reconstruction of the scalar quantities \bar{c}
 397 and g . Reconstructed fields of \bar{c} and g for the 3 flow cases at $Sc = 1$ are displayed in Fig.
 398 8. For the sake of clarity and conciseness, only the dense data grid reconstructions ($h = 10$)
 399 are shown as the sparser reconstructions showed similar trends to the velocity fields (Sec.
 400 V A 1). Unlike \mathbf{f} , we expect that the reconstructions of g will resemble the CFD data since
 401 the term has not been modeled or decomposed in any way. Figure 8 confirms that this is
 402 the case since the reconstructed g fields match closely with the CFD data presented in Fig.
 403 2 for all 3 cases throughout the majority of the domain. This is a significant result, as it

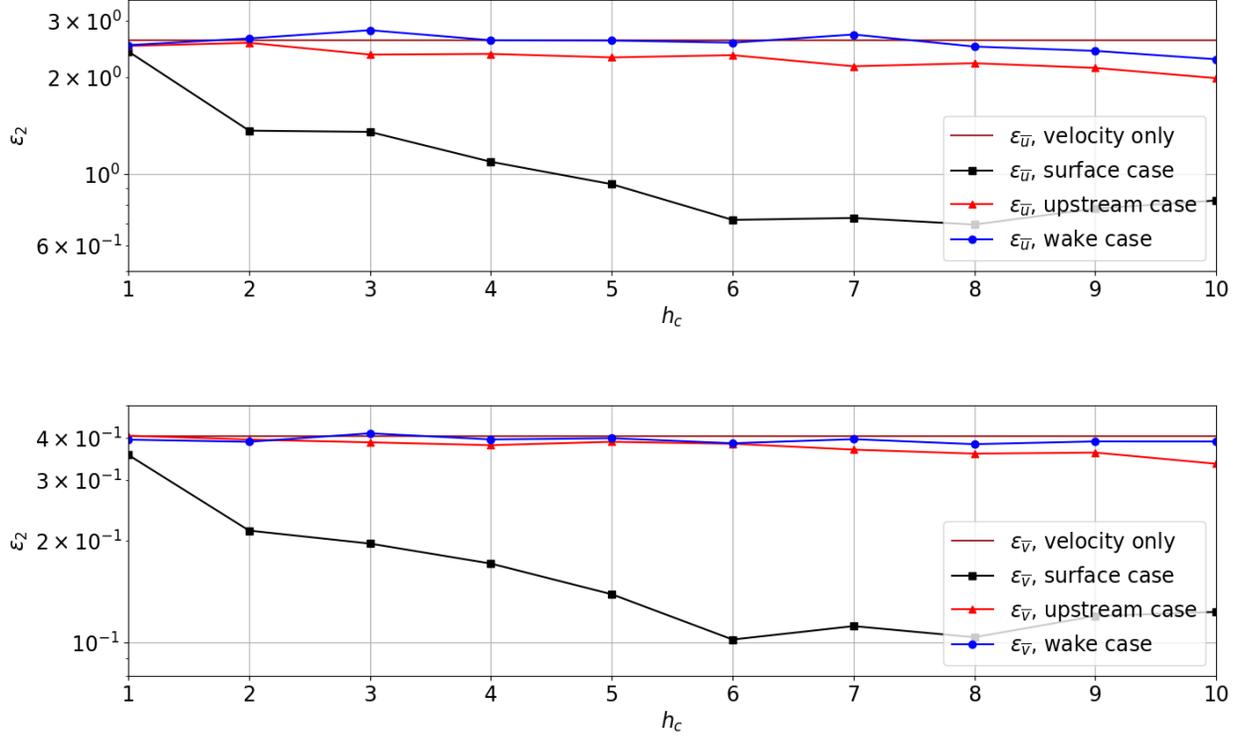


FIG. 7: L_2 norm versus h_c for velocity fields (top) \bar{u} and (bottom) \bar{v} .

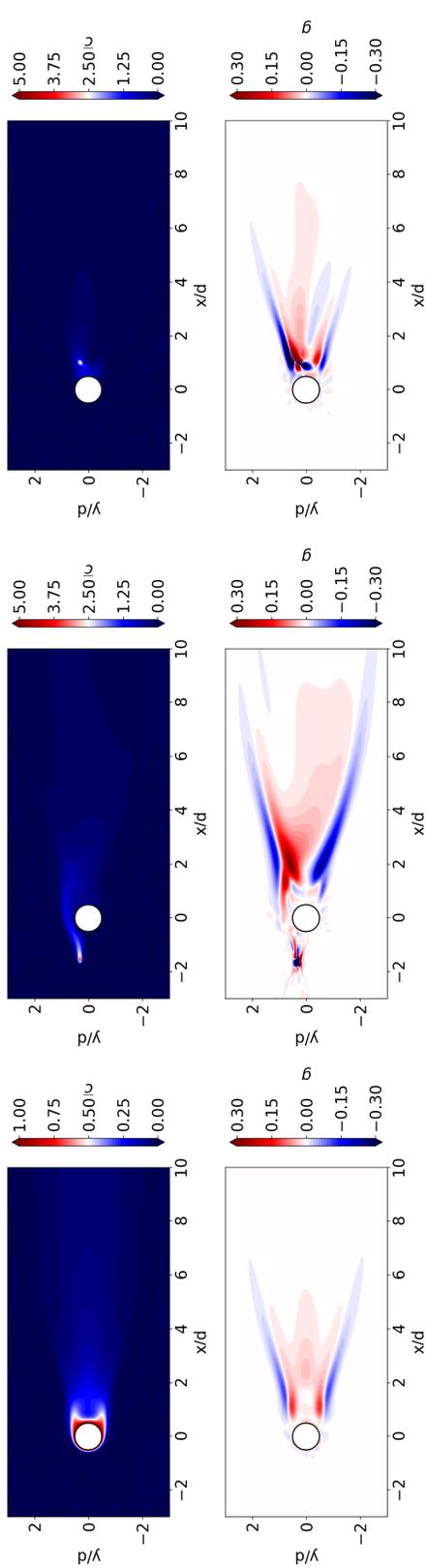
404 demonstrates that the PINN is inferring the true solution to the time-averaged advection-
405 diffusion equation in the presence of limited data. The areas that exhibit the highest absolute
406 error are found near the passive scalar sources, particularly in the upstream and wake cases
407 in which the scalar source was intentionally not given as a boundary condition to study the
408 PINN's response. In these cases and in the source region, the absolute error of g reaches
409 magnitudes of > 10 , suggesting that the PINN struggles with the high spatial gradients of
410 \bar{c} present in this part of the domain. As discussed in Sec. II B, g depends on the spatial
411 gradients of the turbulent scalar fluxes $\overline{u'c'}$ and $\overline{v'c'}$, which are at a maximum at the point-
412 source due to the step-function nature of the \bar{c} field at this point. This causes the PINN
413 to overestimate the magnitude of g at the source, and this behavior is only seen in the two
414 point-source cases since in the surface case (where \bar{c} is given as a boundary condition on the
415 cylinder surface), the scalar gradients and fluxes are well-defined which results in an absolute
416 error of near-zero. However, the regions of high absolute error in Fig. 8 seen downstream
417 of the cylinder surface for the surface case is due to rapid dissipation of the passive scalar.
418 The large decrease in magnitude of \bar{c} challenges the PINN's optimization of the data loss in

419 this region, resulting in a less accurate reconstruction.

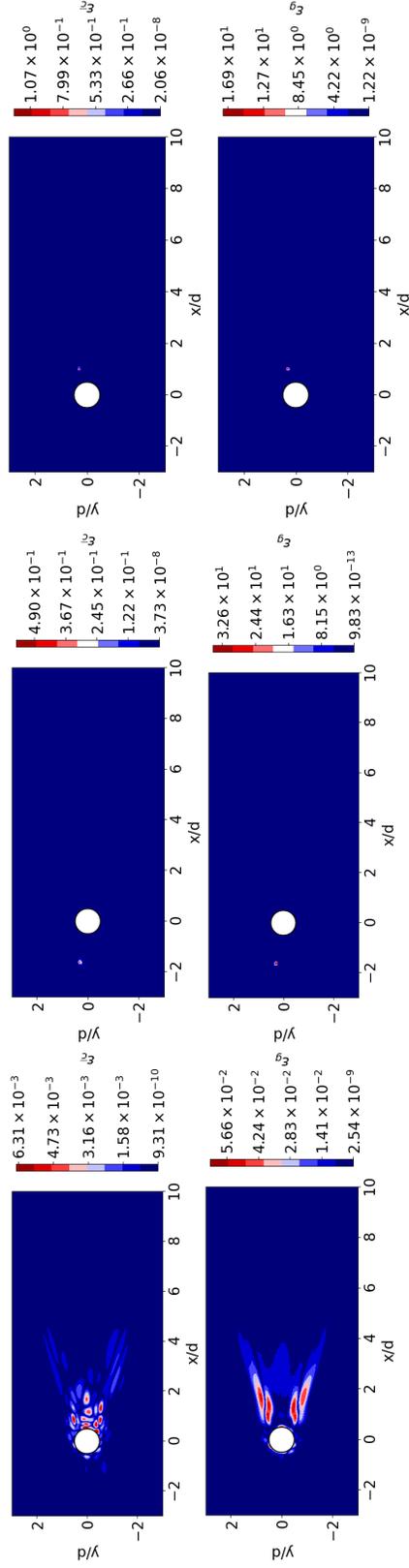
420 For the point-source cases, the overestimation behavior could be remedied through a
421 relaxation function or constraint placed within the PINN. The reconstructed \bar{c} fields show
422 more predictable behavior, where the absolute error is highest around the scalar sources
423 and far lower in the rest of the domain for similar gradient-dependent reasons as mentioned
424 previously. Once again, the absolute error is the lowest in the surface case. The reader may
425 notice asymmetry about the streamwise centerline in the absolute error fields of the surface
426 case in Fig. 8, as well as figures in the previous section. This may be unexpected, since
427 the RANS solution for this case should be entirely symmetric, yet the PINN reconstructions
428 exhibit some asymmetries. This feature is a residual symptom of the PINN’s random initial-
429 ization (see Sec. IV B). The architecture of the PINN and its parameters do not have any
430 direct relationship to the structure of the flow domain, meaning that the fully-connected
431 network architecture has no mechanism by which to identify or respect spatial symmetries
432 in the flow. Any learned symmetrical structures are purely a result of the PDE system
433 and training data. Enforcing symmetry in the PINN solution is possible through applying
434 Neumann BCs at the flow centerline, but we have chosen to avoid this approach in this in-
435 vestigation as it relies on a priori knowledge of the flow conditions - information that would
436 be unavailable in practical applications.

437 The L_2 norms for all 3 cases are presented in Fig. 9, which also illustrates the differences
438 between supplying \bar{c} as a boundary condition versus only as data, as well as the differences
439 between the ability of the PINN to reconstruct \bar{c} and g . Figure 9 shows that the surface case
440 is the most responsive to changes in h , where the reconstructions of \bar{c} and g reach plateaus at
441 $h \approx 5$ to 6. This is due to the additional reinforcement of the emissive boundary conditions
442 on the cylinder surface. The cylinder surface is a region of importance (as discussed in Sec.
443 V A 1) due to the high spatial gradients present, which is made even more critical in the
444 surface case due to passive scalar emission from this boundary. These factors increase the
445 sensitivity of the PINN to additional flow information and result in rapid improvement in
446 reconstruction accuracy.

447 In contrast, the upstream and wake cases are slower to benefit from additional data.
448 Nevertheless, the L_2 errors of the \bar{c} fields are seen to improve by a factor of 10 between
449 $1 < h < 10$, and show no signs of plateauing or saturating unlike the surface case or the
450 velocity fields analyzed previously. On the other hand, $\varepsilon_2(g)$ is seen to diverge as the data



(a) Top: reconstructed fields of \bar{c} . Bottom: reconstructed fields of g . Left: surface case. Middle: upstream case. Right: wake case.



(b) Top: absolute error fields of \bar{c} . Bottom: absolute error fields of g .

FIG. 8: Reconstructions of \bar{c} and g at $Sc = 1$ and $h = 10$.

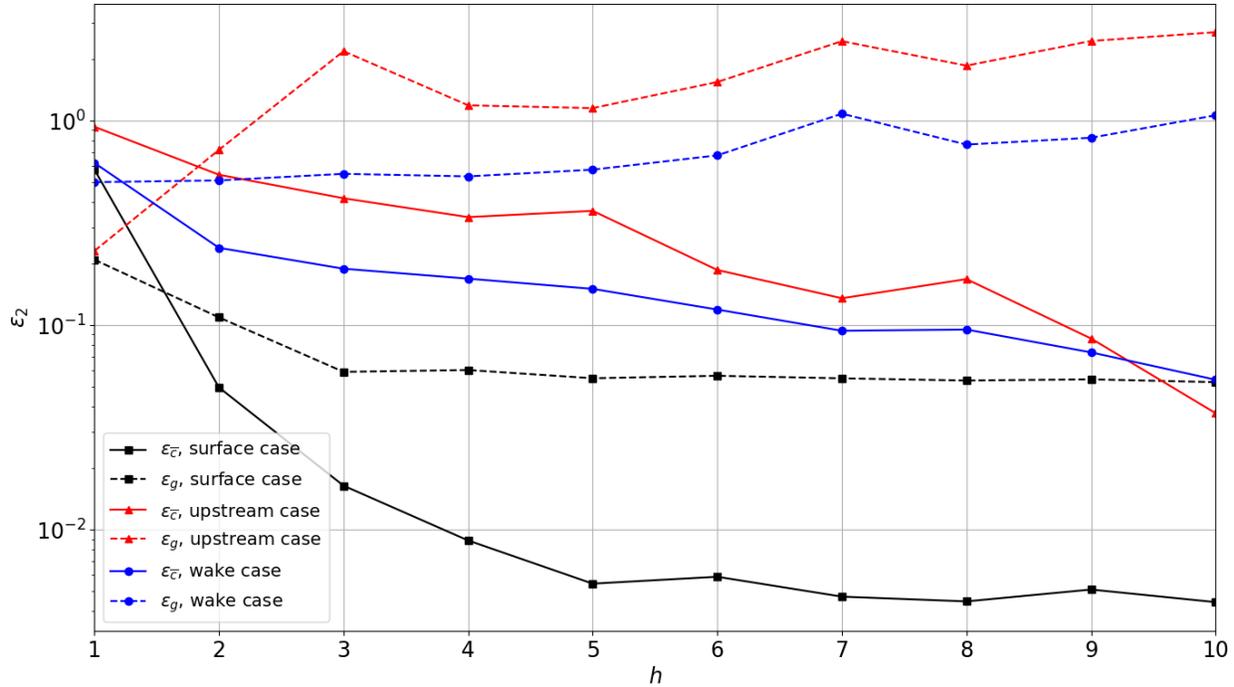


FIG. 9: L_2 norm versus h for scalar quantities \bar{c} and g for all 3 flow cases at $Sc = 1$.

451 grid density increases in the upstream and wake cases. When observing Fig. 8b, it is clear
 452 again that the source region contributes almost entirely to this behavior. The reasons for
 453 this were previously discussed and it can be seen that this does not affect the reconstruction
 454 of \bar{c} to the same degree, which responds positively to additional data. This may be due to
 455 differences in the PINN's reconstruction of directly accessible fields (\bar{u} , \bar{v} , \bar{c}) versus inferred
 456 fields (f_{su} , f_{sv} , g). Despite these limitations, the reconstructions of g remain remarkably
 457 accurate and well-structured away from the locations of the point-sources (see Fig. 2 to
 458 compare), demonstrating that the PINN can successfully reconstruct g on a macro scale
 459 even with relatively sparse data.

460 3. Effect of Schmidt Number on Full-Domain Scalar Reconstruction

461 In this investigation, the data point grid density was set as $h = 4$ and the Schmidt number
 462 of the flow was varied according to Table IV D. This data point grid density was chosen based
 463 on Fig. 9, ensuring the reconstructed fields would be sufficiently accurate while providing
 464 room for improvement without risk of overfitting to the different Sc datasets. Different
 465 Schmidt numbers were considered to determine how generalizable the PINN is to different

466 passive scalar regimes, and to identify any bias it may have to one regime versus another.

467 The L_2 norms for each Schmidt number are shown in Fig. 10, allowing comparison
 468 between the PINNs' reconstructive capabilities for the three flow cases at the five simulated
 469 Schmidt numbers. The reconstruction of \bar{c} is seen to behave uniquely in all three cases. For
 470 the surface case, $\varepsilon_2(\bar{c})$ is lowest at $Sc = 0.1$ and reaches a maximum between $3 < Sc < 10$.
 471 For the upstream case, $\varepsilon_2(\bar{c})$ once again rises as Sc increases, but at a slower, steadier rate
 472 than for the surface case. For the wake case, $\varepsilon_2(\bar{c})$ remains near constant and appears to be
 473 reconstructed to the same accuracy independently of the Schmidt number.

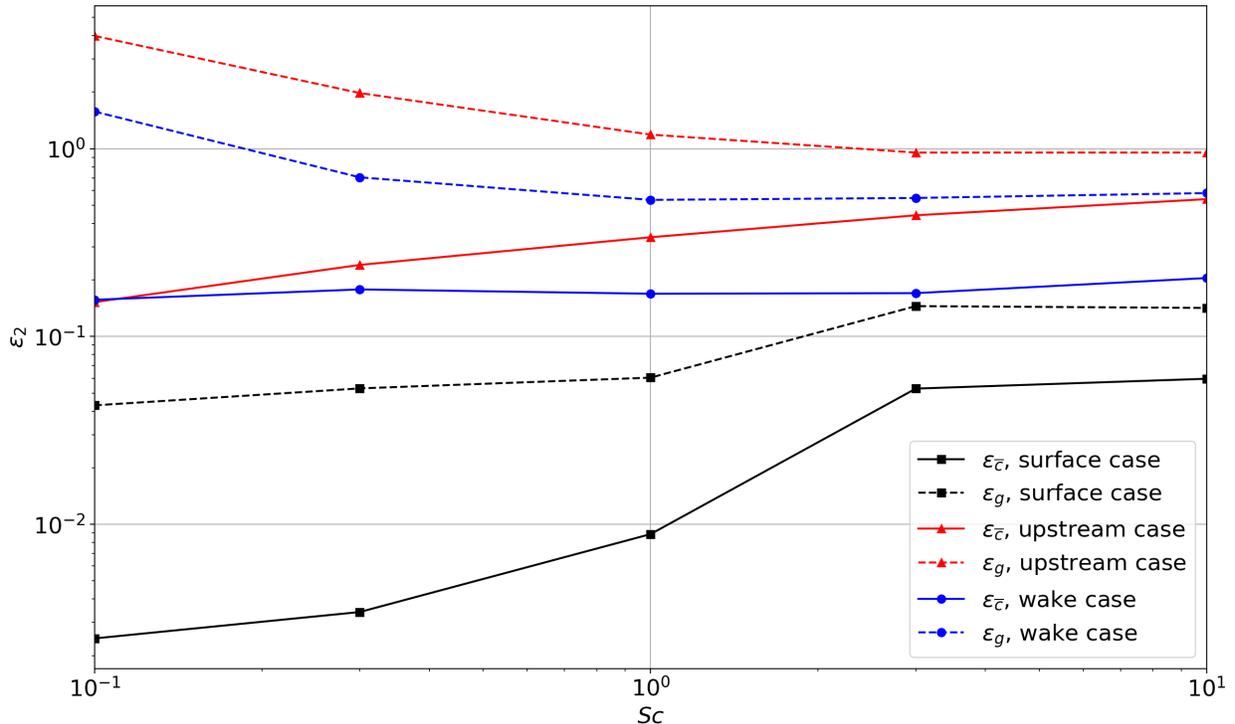


FIG. 10: L_2 norm versus Schmidt number for scalar quantities \bar{c} and g for all 3 flow cases at $h = 4$. Solid lines: $\varepsilon_2(\bar{c})$. Dashed lines: $\varepsilon_2(g)$. Black: surface case. Red/triangle: upstream case. Blue/circle: wake case.

474 In the surface case, the increased diffusivity of the passive scalar at a low Sc comple-
 475 ments the boundary conditions and helps the PINN compute smooth concentration gradients
 476 around the cylinder region. This advantage is lost as Sc increases, as the steeper gradients
 477 do not benefit as much from the boundary conditions and require more data to be resolved
 478 adequately since most of the diffusive activity of the scalar takes place over a very small
 479 region of the domain. A similar explanation can be applied to the upstream case, albeit

480 with a less severe loss of accuracy as this case never had any \bar{c} boundary conditions to use
481 as a basis for the PINN reconstruction.

482 For the wake case, it is possible that a different mechanism results in a Sc -independent
483 error field. The point source for this case is located at (1.0, 0.3), which is near the center
484 of the cylinder recirculation bubble. The oscillatory vortex shedding behind the cylinder,
485 when averaged, results in values of \bar{u} and \bar{v} of near zero (see Fig. 1), which results in
486 advective fluxes (\overline{uc} , \overline{vc}) which are also near zero in magnitude. This means the transport
487 of the passive scalar is almost entirely dominated by diffusion (see Eq. 1a), resulting in
488 a circular plume. We can infer from Fig. 10 that the PINN can reliably model purely
489 diffusive behavior of \bar{c} , and that the error fields are dependent on the balance of advection
490 to diffusion within the flow. The results suggest that the most challenging case for the
491 PINN’s reconstruction of \bar{c} is an advection-dominated flow (high advective fluxes and a
492 high Schmidt number corresponding to low diffusivity). This result is intuitive, as accurate
493 modeling of the advective term of Eq. 1a is conditional on a well-resolved grid of both \bar{c}
494 as well as \bar{u} and \bar{v} data, meaning that error in the PINN’s estimation of either field will be
495 compounded when solving for the residual of Eq. 3a.

496 Interestingly, as Sc increases, $\varepsilon_2(g)$ is seen to decrease for the point-source cases, and
497 increase for the surface case. The point-source region is responsible for the majority of the
498 accumulated error. It is possible that a reduction in the effective source size as Sc increases
499 results in a smaller region for $|g|$ to be incorrectly overestimated by the PINN, which reduces
500 the total error of the reconstruction. The behavior of $\varepsilon_2(g)$ in the surface case is likely due
501 to the steepening of scalar concentration gradients near the cylinder surface, as previously
502 discussed with respect to $\varepsilon_2(\bar{c})$ in this case.

503 B. Source Identification

504 The cropped-domain investigations focus only on the upstream and wake cases in which
505 \bar{c} was not given to the PINN as a boundary condition, allowing the study of the PINNs’
506 reconstructions when presented solely with sampled measurement data. Within these inves-
507 tigation, we instead use a more specific ‘source metric’ to quantify the PINN’s accuracy
508 when reconstructing the passive scalar source properties in order to satisfy the objectives.
509 This metric is based on a *modified scalar concentration*, c^* proposed in Ref. [9] and applies

510 the following scaling to the scalar concentration field:

$$c^* = \begin{cases} \log_{10}(10c) & c \geq 0.0001 \\ c & c < 0.0001 \end{cases}. \quad (9)$$

511 This scaling prevents the high concentrations at the source from overpowering the plots
 512 and helps to preserve the relations between small values of \bar{c} while also eliminating the
 513 hassle of performing logarithms on values that approach machine precision. The quantity \bar{c}^*
 514 is used to generate isocontours of the scalar concentration field localized around the source,
 515 which allows the dispersion of the source to be quantified as well as the location of highest
 516 concentration. An isocontour was produced at levels of $\bar{c}^* = 1.0$ and $\bar{c}^* = 1.2$ for the data
 517 sparsity and Schmidt number investigations, respectively. These levels were chosen to give
 518 the PINN a flexibility margin in the event that it underestimated the source strength of
 519 $\bar{c} = 5$. Similarly, a value of \bar{c}^* too low would produce an isocontour that encompasses all
 520 of the source domain, and would therefore be impractical for assessing the PINN's source
 521 prediction capacity. The source properties are then compared to the reference fields. The
 522 dispersion is calculated by measuring the distance between the two furthest points on the
 523 isocontour in x and y to give the measures L_x and L_y which are normalized relative to the
 524 corresponding case's CFD source size. The location of highest concentration is simply taken
 525 as the location of the maximum value within the \bar{c}^* array and is given as a distance from
 526 the normalized true source location as Δx , Δy .

527 1. *Effect of Data Sparsity on Source Reconstruction*

528 This section presents the results for the cropped-domain data sparsity investigation, in
 529 which the grid density of the \bar{c} data points was varied between $3 < h_c < 100$, while the
 530 velocity data (supplied on a separate grid as described in Sec. IV D) remained at a constant
 531 density of $h_u = 6$. The Schmidt number was set as $Sc = 1$ since it is within the range of
 532 typical experimental Schmidt numbers. The results illustrating how L_x , L_y , Δx and Δy
 533 vary as a function of the density of \bar{c} data are shown in Fig. 11.

534 Figure 11a shows the convergence trend of L_x and L_y to the CFD reference values as
 535 h_c increases. The manner in which this occurs is different between the cases. For the
 536 upstream case, sparse \bar{c} data is shown to result in an underestimate for both L_x and L_y ,

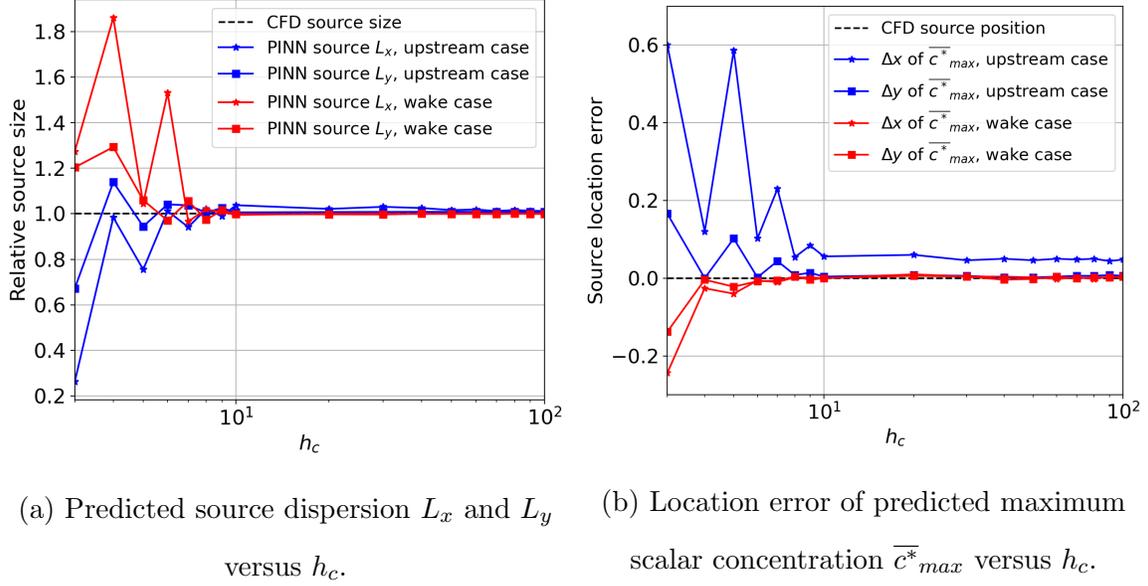


FIG. 11: Impact of h_c on reconstructed source properties for upstream and wake cases at $Sc = 1$.

537 while the inverse is true in the wake case which can also be seen for Δx and Δy in Fig.
538 11b. Furthermore, the upstream case (which is advection-dominant) exhibits a small steady-
539 state error in L_x beyond a data point grid density of around 10 points per unit, in which the
540 PINN overestimates the source size by a small margin. This trend is also visible in Fig. 11b,
541 showing that the PINN predicts the source location to be around 0.05 units downstream
542 of the true location. This error is only present for L_x in the advective case, which is the
543 regime under which the corresponding advective flux $\bar{u}\bar{c}$ is significantly higher than in other
544 cases. Supplying additional data to the PINN does not improve the discrepancy, suggesting
545 once again that high advective fluxes pose difficulties when solving Eq. 3a and accurately
546 resolving regions of high advective transport. It is very possible that providing the PINN
547 with a finer grid of the velocities (particularly \bar{u}) would remedy this. For all other cases in
548 Fig. 11, the source properties converge to their respective CFD reference values before or
549 around 10 points per unit. These results demonstrate that the PINN is capable of identifying
550 the correct source properties from a density of data points similar to that which is needed to
551 accurately reconstruct the entire flow field. Performance may be improved further by scaling
552 the scalar concentration data to reveal smaller spatial variations in the concentration field,
553 which could push the PINN toward identifying the correct source location under advective

554 regimes. Finally, Fig. 11 illustrates a ‘zig-zagging’ pattern of convergence for all cases, which
555 is a consequence of the distribution of points within the data point grid. It is clear when
556 considering Fig. 11 that the quality of the reconstructed field can vary significantly between
557 neighboring models, and that there is an element of chance of where high-importance data
558 points will be placed on the reconstructed field.

559 2. *Effect of Schmidt Number on Source Reconstruction*

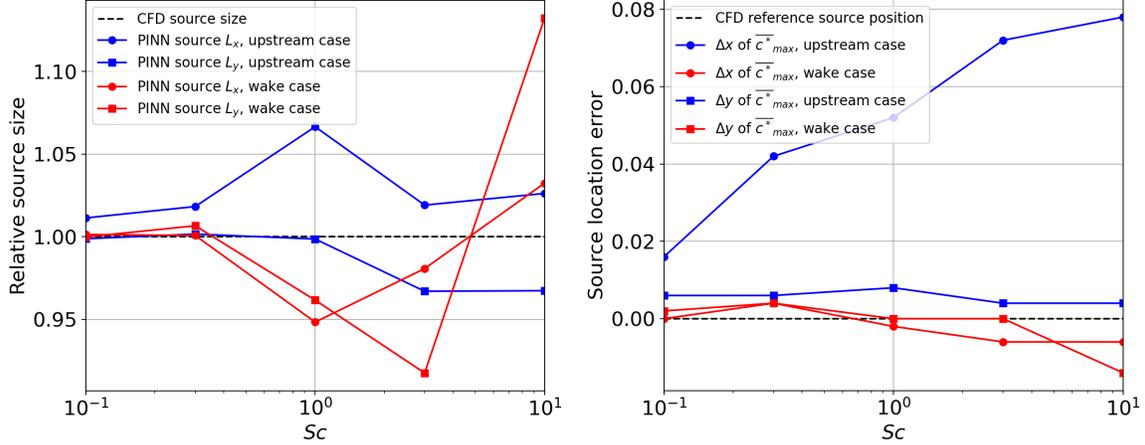
560 The final investigation focused on the influence of varying the Schmidt number between
561 $0.1 < Sc < 10$ on the source reconstruction characteristics. The \bar{c} , \bar{u} and \bar{v} data grids had a
562 coincident, uniform grid density of $h = 10$.

563 Figure 12 presents the relative reconstructed source sizes and locations for all Schmidt
564 numbers. A noticeable difference between Fig. 12 and Fig. 11 is that the adjustment of
565 the Schmidt number appears to be much less influential on source reconstruction than the
566 density of the data point grid. The maximum discrepancy of the source size in Fig. 11 was
567 around 90% (reducing until $h \approx 8$ to 9), while in Fig. 12 the error in the source size does not
568 exceed around 14%, occurring at the extreme Schmidt number of $Sc = 10$. This suggests
569 that data point placement is a hyperparameter that (when set optimally) can mitigate the
570 uncertainty in performance when other properties of the physical system are adjusted.

571 Figure 12a shows that the PINN has difficulty reconstructing the source size at high
572 Schmidt numbers, as seen by the more erratic behaviors of L_x and L_y beyond $Sc = 0.3$.
573 Divergent behavior is seen in Fig. 12b for Δx of the upstream case, which is shown to
574 increase linearly with Sc up to a maximum of $\Delta x = 0.08$, and is representative of the
575 PINN’s bias toward downstream prediction of source location in advection-dominated flows
576 - a phenomenon also seen in Fig. 11b. All other cases in Fig. 12b do not stray beyond
577 $\Delta = 0.01$, suggesting that the Schmidt number does not influence the prediction of the
578 source location in diffusive regimes.

579 VI. CONCLUSIONS

580 This article has introduced passive scalar transport systems as another application of
581 PINNs. This was done by investigating the mean fields of three classic cylinder flow cases



(a) Predicted source dispersion L_x and L_y versus Sc . (b) Location error of predicted maximum scalar concentration $\overline{c^*_{max}}$ versus Sc .

FIG. 12: Impact of Schmidt number on reconstructed source properties for upstream and wake cases at $h = 10$.

582 at a Reynolds number of 150 and Schmidt numbers of $0.1 < Sc < 10$, with a passive scalar
 583 source active in a different part of the domain and under different boundary conditions
 584 for each case. The focus was on reconstructing the terms of the time-averaged advection-
 585 diffusion equation, the mean velocities, and identifying the source characteristics. This was
 586 done to assess the effectiveness of PINNs for solving the inverse problem of modeling passive
 587 scalar dispersion using sparse field measurements.

588 A PINN model trained using sparse velocity measurements was provided with passive
 589 scalar data on an increasingly dense uniform grid, where it was found that the additional
 590 scalar data resulted in improvements to the velocity reconstructions. These improvements
 591 were amplified when the PINN was given specified boundary conditions for the scalar. This
 592 trend was present throughout the first half of the results, where the implementation of
 593 Dirichlet scalar boundary conditions was found to produce very accurate reconstructed fields
 594 that performed better than the two point-source flows in all cases. This strongly suggests
 595 that any knowledge of scalar boundary conditions is greatly beneficial to PINN accuracy, and
 596 this information should be sought after in real-world applications to make the most effective
 597 use of PINNs. This result has implications for future applications of PINNs, particularly
 598 when training data for certain fields have differing degrees of availability.

599 The PINN was found to be able to reliably infer and reconstruct the unknown closure term
600 g of the advection-diffusion equation, producing reconstructed fields that are remarkably
601 accurate across the domain as a whole. However, g was found to be very sensitive in regions
602 of high turbulent fluxes, requiring scalar boundary conditions to be reconstructed accurately
603 in these areas. Changing the Schmidt number of the flow demonstrated that the PINN can be
604 used at a wide range of Schmidt numbers, but generally performs better at low-intermediate
605 Schmidt numbers ($0.3 < Sc < 3$) where diffusion dominates. The reconstruction of g was
606 found to improve for the two point-source cases as Sc increases, but reduce for the surface
607 case. The L_2 errors for the surface case remained lower than the point source cases at all
608 Schmidt numbers. The PINN’s ability to infer g from mean-flow measurements of scalar
609 concentration and velocity, while simultaneously reconstructing and refining these fields,
610 suggests advantages over simpler interpolation-based methods for modeling this type of
611 problem.

612 When focusing on reconstruction of the scalar source properties (source size/location), the
613 PINN reached peak accuracy at $h_c \approx 10$ and showed little improvement beyond this point.
614 The source size and location predictions were found to exhibit a small steady-state error for
615 properties parallel to the freestream flow direction, representing uncertainty in the sources’
616 sizes and locations along the streamlines. The oscillatory nature of the reconstructed proper-
617 ties highlighted the importance of considering data point placement in high-sensitivity areas
618 of the domain to ensure that gradients and fluxes can be resolved adequately. Nevertheless,
619 the PINN was broadly successful in identifying the scalar source properties when provided
620 with an amount of data similar to that which is necessary to reconstruct the entire flow
621 field. This result is useful for applications regarding urban pollution control and modeling,
622 and should be pursued further.

623 The use of PINNs should also be extended to more passive scalar transport cases, in-
624 cluding complex geometries, 3-dimensional cases, as well as the use of experimental data
625 [5, 6] instead of idealized CFD simulations. Studies could be performed using transient data
626 instead of mean-flow fields, as it is possible that a time-resolved PINN solver is more likely
627 to produce a unique solution than one trained on mean fields [21]. Other neural network
628 architectures could be employed to solve this problem, such as variational autoencoders for
629 model order reduction, or transformers to leverage long/short-term network attention [43].

630 **Appendix: Network Hyperparameter Sensitivity Analysis**

631 Due to the increased complexity of the physical systems in this study compared to pre-
 632 vious work [19, 24, 42], a sensitivity study was conducted to determine the optimal network
 633 size for the physical system being investigated. As described in Sec. IV B, the hidden
 634 network size used in this article was 7 layers of width 100 ($N_L = 7$, $N_N = 100$). In this
 635 sensitivity study, additional PINNs were trained with varying increases in hidden layers and
 636 widths as presented in Table II.

Network size	Network 1	Network 2	Network 3	Network 4	Network 5
N_L	7	8	9	8	9
N_N	100	100	100	120	120

TABLE II: Network architectures for sensitivity investigation.

637 The PINN models reconstructed the \bar{u} , \bar{v} , \bar{c} and g fields for the two point source flow
 638 cases (upstream and wake) at $h = 4$ and $Sc = 1.0$. The L_2 errors for the reconstructions
 639 were computed as shown in Fig. 13. It can be seen that increasing the network size through
 640 both an increase in N_L and N_N provides different improvements to the reconstructed fields.
 641 The \bar{c} and g fields showed improvements of $< 3\%$ as the network size was increased. For \bar{u}
 642 and \bar{v} on the other hand, more tangible improvements can be found of at least 20% for both
 643 cases with a maximum of $\approx 46\%$ for \bar{v} in the wake case between the smallest and largest
 644 networks. While this may appear significant, it is worth noting that these improvements are
 645 calculated relative to the original (smallest) model, where the L_2 errors were already of the
 646 order $O(10^{-3})$ prior to increasing the network size. Whether improvements of this magnitude
 647 are detectable and useful in practice remains to be seen, since the PINN’s accuracy in this
 648 application is currently limited by its passive scalar modeling abilities. Nevertheless, this
 649 outcome may indicate that a larger PINN network architecture would be more suitable in
 650 use cases where training data is sparse in order to make the most effective use of available
 651 data.

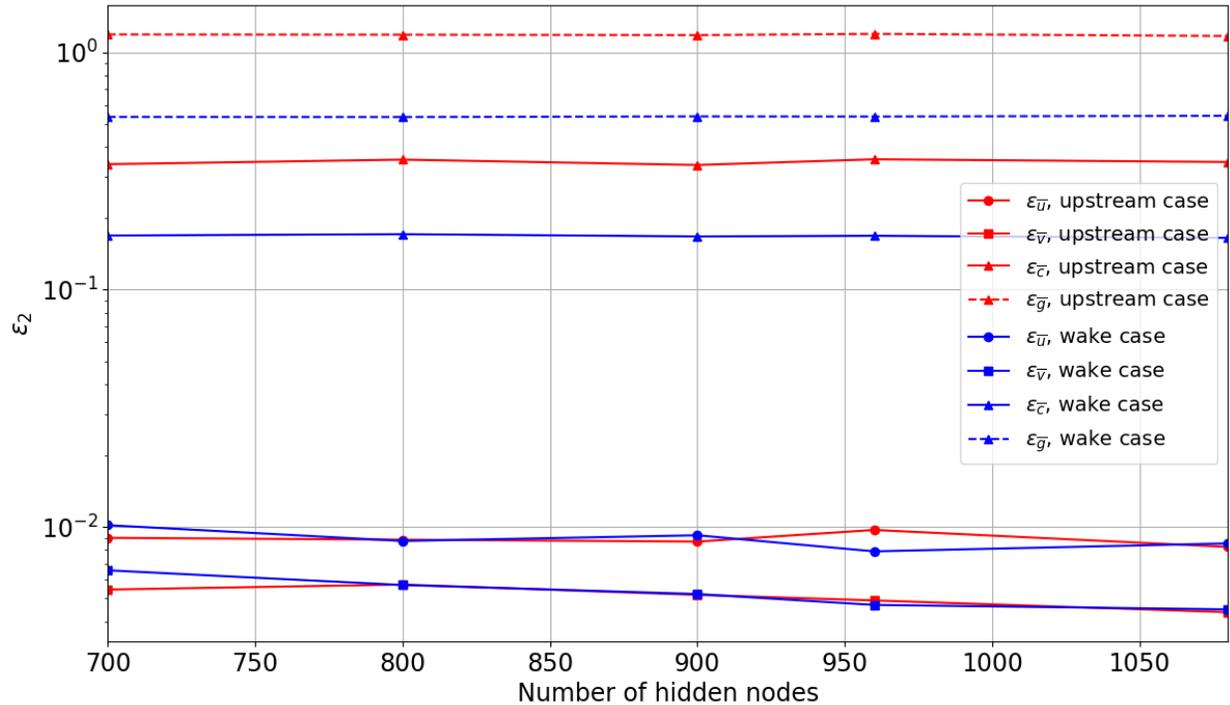


FIG. 13: Network size investigation.

652 **ACKNOWLEDGMENTS**

653 The authors would like to acknowledge the use of the IRIDIS High Performance Com-
 654 puting Facility, and associated support services at the University of Southampton, in the
 655 completion of this work.

656 J. I. Rawden gratefully acknowledges the financial support of the Antony Wright Schol-
 657 arship, provided by the Department of Aeronautics and Astronautics at the University of
 658 Southampton, and the School of Engineering.

659 **DATA AVAILABILITY**

660 Data supporting the findings of this article are openly available at [44].

661 [1] WHO, Air pollution, https://www.who.int/health-topics/air-pollution#tab=tab_1,
 662 [Accessed 17-03-2025].

- 663 [2] V. Fuka, Z.-T. Xie, I. P. Castro, P. Hayden, M. Carpentieri, and A. G. Robins, Scalar fluxes
664 near a tall building in an aligned array of rectangular buildings, *Bound.-Layer Meteorol.* **167**,
665 53 (2018).
- 666 [3] C. Górlé, J. van Beeck, and P. Rambaud, Dispersion in the wake of a rectangular build-
667 ing: validation of two reynolds-averaged navier–stokes modelling approaches, *Boundary-layer*
668 *meteorology* **137**, 115 (2010).
- 669 [4] I. P. Castro, I. R. Cowan, and A. G. Robins, Simulations of flow and dispersion around
670 buildings, *Journal of Aerospace Engineering* **12**, 145 (1999).
- 671 [5] H. Lim, D. Hertwig, T. Grylls, H. Gough, M. v. Reeuwijk, S. Grimmond, and C. Vanderwel,
672 Pollutant dispersion by tall buildings: laboratory experiments and large-eddy simulation, *Exp.*
673 *Fluids* **63**, 92 (2022).
- 674 [6] T. Rich and C. Vanderwel, Pollutant dispersion around a single tall building, *Bound.-Layer*
675 *Meteorol.* **190**, 34 (2024).
- 676 [7] P. Salizzoni, S. Fellini, H. Gamel, M. Marro, and L. Soulhac, Atmospheric dispersion down-
677 stream a two-dimensional obstacle: experimental evaluation of turbulence closure models,
678 *Boundary-Layer Meteorology* **191**, 15 (2025).
- 679 [8] J. S. Apte and C. Manchanda, High-resolution urban air pollution mapping, *Science* **385**, 380
680 (2024).
- 681 [9] V. Mons, L. Margheri, J.-C. Chassaing, and P. Sagaut, Data assimilation-based reconstruction
682 of urban pollutant release characteristics, *J. Wind Eng. Ind. Aerodyn.* **169**, 232 (2017).
- 683 [10] D. P. Foures, N. Dovetta, D. Sipp, and P. J. Schmid, A data-assimilation method for Reynolds-
684 averaged Navier–Stokes-driven mean flow reconstruction, *J. Fluid Mech.* **759**, 404 (2014).
- 685 [11] S. Symon, N. Dovetta, B. J. McKeon, D. Sipp, and P. J. Schmid, Data assimilation of mean
686 velocity from 2D PIV measurements of flow over an idealized airfoil, *Exp. Fluids* **58**, 1 (2017).
- 687 [12] D. Walters, A. J. Baran, I. Boutle, M. Brooks, P. Earnshaw, J. Edwards, K. Furtado, P. Hill,
688 A. Lock, J. Manners, *et al.*, The met office unified model global atmosphere 7.0/7.1 and jules
689 global land 7.0 configurations, *Geoscientific Model Development* **12**, 1909 (2019).
- 690 [13] A. Fraser, J. Abbott, and R. Rose, Application of data assimilation to the uk air quality
691 forecast, *Air Pollution Modeling and its Application XXIII* , 617 (2014).
- 692 [14] U. Cadambi Padmanaban, B. Ganapathisubramani, C. Vanderwel, and S. Symon, Towards
693 passive scalar reconstruction using data assimilation, 14th UK Conference on Wind Engineer-

- 694 ing (2024).
- 695 [15] P. Zille, T. Corpetti, L. Shao, and C. Xu, Super-resolution of turbulent passive scalar images
696 using data assimilation, *Exp. Fluids* **57**, 21 (2016).
- 697 [16] S. Li, W. Zhou, H. J. Sung, and Y. Liu, On the origin of counter-gradient transport in
698 turbulent scalar flux: physics interpretation and adjoint data assimilation, *J. Fluid Mech.*
699 **999**, A81 (2024).
- 700 [17] R. Wang, B. Chen, S. Qiu, Z. Zhu, and X. Qiu, Data assimilation in air contaminant dispersion
701 using a particle filter and expectation-maximization algorithm, *Atmosphere* **8**, 170 (2017).
- 702 [18] M. Raissi, P. Perdikaris, and G. E. Karniadakis, Physics-informed neural networks: A deep
703 learning framework for solving forward and inverse problems involving nonlinear partial dif-
704 ferential equations, *J. Comput. Phys.* **378**, 686 (2019).
- 705 [19] L. Sliwinski and G. Rigas, Mean flow reconstruction of unsteady flows using physics-informed
706 neural networks, *Data-Centric Eng.* **4**, e4 (2023).
- 707 [20] Y. Patel, V. Mons, O. Marquet, and G. Rigas, Turbulence model augmented physics-informed
708 neural networks for mean-flow reconstruction, *Phys. Rev. Fluids* **9**, 034605 (2024).
- 709 [21] M.-C. Volk, A. Sergent, D. Lucor, M. Mommert, C. Bauer, and C. Wagner, A PINN method-
710 ology for temperature field reconstruction in the PIV measurement plane: Case of Rayleigh-
711 Bènard convection, arXiv preprint arXiv:2503.23801 (2025).
- 712 [22] Y. Wang, B. Zhang, C. Hu, H. Jia, C. Wei, and H. Kikumoto, Flow field reconstruction and
713 wind pressure estimation from sparse measurements using physics-informed neural networks:
714 Application to two-dimensional street canyon flow, *Building and Environment* , 113616 (2025).
- 715 [23] C. Hu, Y. Cui, W. Zhang, F. Qian, H. Wang, Q. Wang, and C. Zhao, Solution of conservative-
716 form transport equations with physics-informed neural network, *Int. J. Heat Mass Trans.* **216**,
717 124546 (2023).
- 718 [24] H. Eivazi, M. Tahani, P. Schlatter, and R. Vinuesa, Physics-informed neural networks for
719 solving Reynolds-averaged Navier–Stokes equations, *Phys. Fluids* **34** (2022).
- 720 [25] X. Jin, S. Cai, H. Li, and G. E. Karniadakis, NSFnets (Navier-Stokes flow nets): Physics-
721 informed neural networks for the incompressible Navier-Stokes equations, *J. Comput. Phys.*
722 **426**, 109951 (2021).
- 723 [26] M. Raissi, A. Yazdani, and G. E. Karniadakis, Hidden fluid mechanics: Learning velocity and
724 pressure fields from flow visualizations, *Science* **367**, 1026 (2020).

- 725 [27] R. Vinuesa, S. L. Brunton, and B. J. McKeon, The transformative potential of machine learn-
726 ing for experiments in fluid mechanics, *Nat. Rev. Phys.* **5**, 536 (2023).
- 727 [28] J. G. von Saldern, J. M. Reumschüssel, T. L. Kaiser, M. Sieber, and K. Oberleithner, Mean
728 flow data assimilation based on physics-informed neural networks, *Phys. Fluids* **34** (2022).
- 729 [29] A. Villié, S. Schmitter, J. G. von Saldern, S. Demange, and K. Oberleithner, Physics-informed
730 neural networks for enhancing medical flow magnetic resonance imaging: Artifact correction
731 and mean pressure and Reynolds stresses assimilation, *Phys. Fluids* **37** (2025).
- 732 [30] H. Wang, Y. Liu, and S. Wang, Dense velocity reconstruction from particle image velocime-
733 try/particle tracking velocimetry using a physics-informed neural network, *Phys. Fluids* **34**
734 (2022).
- 735 [31] A. Baklanov and Y. Zhang, Advances in air quality modeling and forecasting, *Global Transi-*
736 *tions* **2**, 261 (2020).
- 737 [32] I. P. Castro and C. Vanderwel, *Turbulent flows: an introduction* (IOP Publishing, 2021).
- 738 [33] C. Keumnam, T. F. Irvine Jr, and J. Karni, Measurement of the diffusion coefficient of naph-
739 thalene into air, *International journal of heat and mass transfer* **35**, 957 (1992).
- 740 [34] L. Lu, X. Meng, Z. Mao, and G. E. Karniadakis, DeepXDE: A deep learning library for solving
741 differential equations, *SIAM Rev.* **63**, 208 (2021).
- 742 [35] H. G. Weller, G. Tabor, H. Jasak, and C. Fureby, A tensorial approach to computational
743 continuum mechanics using object-oriented techniques, *Comput. Phys.* **12**, 620 (1998).
- 744 [36] H. Jasak, OpenFOAM: Open source CFD in research and industry, *Int. J. Nav. Archit. Ocean*
745 *Eng.* **1**, 89 (2009).
- 746 [37] OpenFOAM, Openfoam user guide, [https://www.openfoam.com/documentation/
747 user-guide/a-reference/a.4-standard-boundary-conditions](https://www.openfoam.com/documentation/user-guide/a-reference/a.4-standard-boundary-conditions), [Accessed 11/08/2025].
- 748 [38] C. Williamson, Three-dimensional wake transition, *J. Fluid Mech.* **328**, 345 (1996).
- 749 [39] D. P. Kingma and J. L. Ba, Adam: A method for stochastic gradient descent, in *Int. Conf.*
750 *Learn. Represent.* (ICLR US., 2015) pp. 1–15.
- 751 [40] D. C. Liu and J. Nocedal, On the limited memory BFGS method for large scale optimization,
752 *Math. Program.* **45**, 503 (1989).
- 753 [41] R. H. Byrd, P. Lu, J. Nocedal, and C. Zhu, A limited memory algorithm for bound constrained
754 optimization, *SIAM J. Sci. Comput.* **16**, 1190 (1995).

- 755 [42] G. Hasanuzzaman, H. Eivazi, S. Merbold, C. Egbers, and R. Vinuesa, Enhancement of PIV
756 measurements via physics-informed neural networks, *Meas. Sci. Technol.* **34**, 044002 (2023).
- 757 [43] A. Solera-Rico, C. Sanmiguel Vila, M. Gómez-López, Y. Wang, A. Almashjary, S. T. Dawson,
758 and R. Vinuesa, β -variational autoencoders and transformers for reduced-order modelling of
759 fluid flows, *Nat. Commun.* **15**, 1361 (2024).
- 760 [44] <https://doi.org/10.5258/SOTON/D3810>.