

A Hybrid Sparse-Semantic Image Classification Framework to Support Marine and Coastal Monitoring and Management

Ying Shi ^a, Lan Song^{b, c *}, Luo Cheng Wu ^d, Huanhuan Li ^e

^aDepartment of Basic Courses, Suzhou City University, Suzhou, China

^bInstitute of Environmental Science (CML), Leiden University, The Netherlands

^c College of Engineering, Eastern Institute of Technology, Ningbo, China

^d Institute of Sound and Vibration Research, University of Southampton, Southampton, UK

^e School of Engineering, University of Southampton, Southampton, UK

* Corresponding author: lansong@eitech.edu.cn

Abstract: Effective ocean and coastal management increasingly relies on timely, reliable, and interpretable information derived from large volumes of visual data collected by satellites, Unmanned Aerial Vehicles (UAVs), and shipborne sensing systems. Such data underpin a wide range of management and governance functions, including environmental monitoring, maritime surveillance, port operations, and coastal planning. However, the operational use of automated image classification in marine and coastal contexts remains constrained by complex visual conditions and the persistent scarcity of labelled data, limiting its integration into routine decision-support **workflows**. This study develops a hybrid sparse-semantic image classification framework designed to enhance the robustness and applicability of marine and coastal image understanding in management-oriented settings. The framework integrates structured local representations, which capture geometric and spatial characteristics of marine scenes through a sparsity-driven encoding strategy with locality, non-negativity, and semantic consistency constraints, with global semantic features learned via self-supervised Vision Transformers (ViTs). Specifically, semantic information extracted from the Classification (CLS) token of a pretrained self-Distillation with No Labels (DINO) model is fused with sparse local descriptors to form interpretable and discriminative image representations that remain effective under limited labelled data conditions. Experiments conducted on four public benchmark datasets and three maritime image datasets demonstrate consistent performance improvements over representative state-of-the-art approaches. Beyond classification accuracy, the proposed framework provides a transferable and expandable analytical tool to support marine environmental monitoring, maritime traffic surveillance, and coastal management practices **by enhancing upstream image interpretation within monitoring and decision-support workflows in data-constrained marine and coastal contexts**.

Keywords: Marine and coastal management; Image classification; Self-supervised learning; Maritime surveillance; Decision support

1. Introduction

Effective ocean and coastal management increasingly depends on timely, reliable, and actionable information to support planning, regulation, and **decision-support workflows** across complex marine and coastal systems. Advances in Earth observation technologies, Unmanned Aerial Vehicles (UAVs), autonomous platforms, and shipborne sensing systems have led to a rapid expansion in the availability of visual data covering marine environments, coastal zones, and maritime activities. These data streams play a central role in management-oriented

1 applications such as marine environmental monitoring (Galgani et al., 2024; Ghaban et al., 2025;
2 Puskic et al., 2026), maritime traffic surveillance (Dong et al., 2024; Teixeira et al., 2025; Yang
3 et al., 2024), port and terminal operations management (Abu Bakar et al., 2025; He et al., 2025;
4 Yorulmaz and Susoy, 2025), habitat assessment and conservation planning (Agrillo et al., 2025;
5 Wang et al., 2025; Xiao et al., 2025), and maritime safety and security governance (Adetunji et
6 al., 2025; Wang et al., 2024). While the volume and diversity of marine and coastal imagery
7 continue to grow, transforming these data into management-relevant knowledge remains a
8 persistent challenge, particularly in large-scale and heterogeneous operational contexts.

9 Recent studies have increasingly adopted data-driven approaches to support maritime and
10 coastal management. For instance, predictive modelling has been applied to forecast shipping-
11 related emissions and identify key driving factors, thereby informing policy development and
12 operational decision-making within the shipping sector (Xu et al., 2025c). In parallel,
13 optimisation-based methods have been used to analyse trade-offs in infrastructure planning and
14 regulatory design, including shore-to-ship electricity deployment and port competition
15 strategies (Xu et al., 2025a). Additional research has investigated the low-carbon efficiency of
16 ship-port systems (Zhang et al., 2025) and assessed the environmental impacts of fuel policies
17 (J. Shi et al., 2025), highlighting the need for robust monitoring and evaluation frameworks to
18 support evidence-based governance. Furthermore, large-scale initiatives such as the Maritime
19 Silk Road underscore the increasing complexity of cross-regional carbon management and
20 policy coordination (Xu et al., 2025b). Collectively, these studies reveal an increasing demand
21 for robust, data-driven analytical tools capable of supporting decision-relevant processes in
22 complex and data-constrained marine and coastal environments.

23 Image classification constitutes a foundational analytical component in many marine and
24 coastal monitoring pipelines, enabling raw visual data to be translated into semantically
25 meaningful categories that can inform downstream decision-support systems, regulatory
26 oversight, and management interventions (Rangel-Buitrago et al., 2026; Vehmaa et al., 2024;
27 Yang et al., 2021). However, compared with terrestrial environments, marine and maritime
28 imagery is typically acquired under highly variable and adverse conditions, including complex
29 illumination, strong reflections, low contrast, scale variability, background clutter, and partial
30 occlusion. These characteristics substantially increase the difficulty of extracting reliable and
31 discriminative image representations, often leading to degraded classification performance
32 when methods are deployed in real-world marine and coastal settings.

33 In parallel, labelled marine imagery remains scarce, expensive to obtain, and unevenly
34 distributed across application domains, posing a major constraint for the adoption of supervised
35 data-driven approaches in routine management and governance processes.

36 From a methodological perspective, sparse coding-based image classification techniques
37 have been widely explored due to their robustness, interpretability, and ability to capture local
38 geometric and structural characteristics of visual scenes (Cao et al., 2024; Cheng et al., 2023;
39 Shu, 2026). By representing visual patterns as sparse combinations of dictionary atoms, sparse
40 coding can effectively capture local geometric and structural characteristics while suppressing

1 redundant information. When combined with handcrafted local descriptors such as the Scale-
2 Invariant Feature Transform (SIFT) (Lowe, 2004), sparse coding has demonstrated resilience
3 to scale and illumination variations, making it attractive for complex visual environments,
4 including marine scenes. Nevertheless, such approaches mainly rely on low-level handcrafted
5 features and have limited capability to capture high-level semantic information, which is
6 increasingly important for distinguishing visually similar marine categories.

7 More recently, self-supervised deep learning models, particularly Vision Transformers (ViTs)
8 (Dosovitskiy et al., 2021), trained using frameworks such as self-Distillation with No Labels
9 (DINO) (Caron et al., 2021), have shown strong capability in learning transferable global
10 semantic representations without requiring large quantities of labelled data. The Classification
11 (CLS) token produced by a pretrained DINO model encodes holistic image-level semantics and
12 has demonstrated strong generalisation across diverse visual recognition tasks. This property is
13 especially appealing for marine and coastal applications, where annotated datasets are often
14 limited. However, deep semantic features alone typically lack explicit modelling of local
15 geometric structures and offer limited interpretability, which may reduce robustness and
16 transparency in challenging visual conditions.

17 These observations highlight the need for image representation frameworks that can
18 simultaneously preserve local geometric structure, capture global semantic meaning, and
19 remain robust and interpretable under limited labelled data conditions. Such frameworks are
20 particularly important for marine and coastal management, where analytical outputs must
21 support operational decision making, risk assessment, and governance processes across diverse
22 environmental and institutional contexts. Sparse coding and self-supervised ViT models offer
23 complementary strengths: sparse coding provides structured and interpretable local
24 representations, while ViTs deliver powerful global semantic understanding. Leveraging this
25 complementarity, this study proposes a hybrid sparse-semantic image classification framework
26 tailored to the analytical needs of marine and coastal monitoring and management.

27 Within the proposed framework, local SIFT descriptors are encoded using a Histogram
28 intersection and Semantic information-based Locality Non-negative Laplacian Sparse Coding
29 (HS-LNLSC) scheme to preserve neighbourhood similarity, non-negativity, and structural
30 consistency. In parallel, global semantic representations are extracted from the CLS token of a
31 pretrained DINO ViT. The two complementary feature types are fused at the feature level to
32 construct discriminative image descriptors, which are subsequently classified using a linear
33 multi-class Support Vector Machine (SVM). By combining structured local representations
34 with deep semantic features, the proposed approach enhances robustness, interpretability, and
35 classification performance in complex marine and maritime scenarios.

36 **The HS-LNLSC-DINO framework generates robust and discriminative image-level**
37 **representations by integrating structured local sparse features with self-supervised global**
38 **semantic features. This combination enables the simultaneous encoding of fine-grained spatial**
39 **characteristics, such as vessel contours and coastline geometries, and higher-level contextual**
40 **information, including scene categories and environmental conditions. As such, the framework**

1 provides a strong analytical basis for applications in marine environmental monitoring, vessel
2 traffic analysis, and coastal infrastructure assessment.

3 It should be noted, however, that the present study evaluates performance at the image-
4 classification level. The extent to which these representations can be translated into actionable
5 insights for downstream applications, such as risk prioritisation or policy evaluation, has not
6 been empirically examined and remains an important direction for future work. Accordingly,
7 the proposed framework is positioned as an analytical module within broader monitoring and
8 decision-support workflows, rather than as a standalone decision-making tool. Its primary
9 contribution is to improve the reliability and interpretability of upstream image analysis,
10 thereby enhancing the quality of information available to monitoring systems. Issues related to
11 domain adaptation and deployment in operational settings are beyond the scope of this study.

12 The main contributions of this study are summarised as follows:

13 (1) A hybrid sparse-semantic image classification framework is developed to support robust
14 marine and coastal image analysis under complex visual conditions and limited labelled data
15 constraints relevant to management practice.

16 (2) A structured local feature encoding strategy is designed to preserve geometric and spatial
17 characteristics of marine scenes through locality, non-negativity, and semantic consistency
18 constraints.

19 (3) Self-supervised Vision Transformer features are integrated to enhance global semantic
20 understanding and cross-domain generalisation in data-scarce marine and coastal monitoring
21 contexts.

22 (4) Extensive experiments on four public benchmark datasets and three maritime image
23 datasets demonstrate consistent performance improvements over representative state-of-the-art
24 methods, highlighting the potential of the proposed framework as a transferable analytical tool
25 for marine environmental monitoring, maritime surveillance, and coastal management
26 applications.

27 The HS-LNLSC-DINO framework contributes a robust and interpretable approach to image-
28 based analysis in marine and coastal contexts. By strengthening information extraction under
29 low-sample and visually complex conditions, it enhances the quality of data inputs for
30 monitoring and assessment processes, thereby supporting indirectly decision-relevant
31 workflows without explicitly performing policy or operational decision-making.

32 The remainder of this paper is organised as follows. Section 2 reviews related literature.
33 Section 3 details the proposed framework. Section 4 presents experimental results. Section 5
34 discusses management implications and practical considerations for integrating the proposed
35 framework into marine and coastal monitoring workflows. Section 6 concludes the paper and
36 discusses future research directions.

37 **2. Literature review**

38 The automated interpretation of marine and coastal imagery has become an increasingly
39 important analytical component of ocean and coastal management systems, supporting
40 environmental monitoring, maritime surveillance, port operations, and coastal governance.

1 Visual data collected from satellites, UAVs, and shipborne sensing platforms are now routinely
2 used to inform situational awareness, risk assessment, and **decision-support workflows** across
3 marine and coastal domains. Image classification plays a central role in these workflows by
4 transforming large volumes of heterogeneous visual data into semantically meaningful
5 categories that can be integrated into monitoring dashboards, early-warning systems, and
6 management decision-support tools.

7 Existing research on image classification for complex natural environments can be broadly
8 organised into three interrelated methodological streams that are particularly relevant to marine
9 and coastal applications: (i) structured local representation learning based on sparse coding, (ii)
10 semantic-aware and joint representation-classification models, and (iii) deep and self-
11 supervised representation learning.

12 This section reviews representative studies within each stream and critically examines their
13 strengths and limitations from the perspective of robustness, interpretability, and applicability
14 to marine and coastal monitoring and management.

15 **2.1. Sparse coding and structured local representations**

16 Sparse coding has been widely investigated for image classification due to its ability to model
17 local structural characteristics and produce interpretable representations (Cao et al., 2024;
18 Cheng et al., 2023; Shu, 2026). By expressing visual features as sparse linear combinations of
19 dictionary atoms, sparse coding highlights salient local patterns while suppressing redundancy,
20 which is advantageous for scenes with cluttered backgrounds and variable appearance. To
21 enhance descriptive capability, Yang et al. incorporated Spatial Pyramid Matching (SPM) into
22 Sparse Coding (SCSPM), enabling multi-scale spatial layout modelling (Yang et al., 2009).
23 Related spatially aware extensions have also demonstrated effectiveness in capturing spatial
24 context for image categorisation (Lazebnik et al., 2006; Perronnin et al., 2010).

25 However, conventional sparse coding often suffers from encoding instability, whereby
26 visually similar local features may be assigned to different codewords, leading to inconsistent
27 representations. To address this issue, Gao et al. introduced Laplacian Sparse Coding (LSC),
28 which incorporates Laplacian regularisation to encourage similar local descriptors to have
29 similar codes (Gao et al., 2010). Wang et al. proposed Locality-constrained Linear Coding
30 (LLC), enforcing locality constraints so that nearby descriptors contribute more strongly to
31 reconstruction (Wang et al., 2010). Min et al. further integrated Laplacian regularisation into
32 LLC to strengthen neighbourhood uniformity (Min et al., 2016). Subsequent studies explored
33 adaptive neighbourhood selection and manifold-preserving constraints to further improve
34 locality-aware coding (Lin et al., 2011; Mairal et al., 2008).

35 Non-negativity constraints have also been introduced to enhance stability and interpretability.
36 Lee and Seung proposed Matrix Factorization (NMF) to learn part-based representations (Lee
37 and Seung, 1999), followed by extensions such as Robust NMF (RNMF) (Zhang et al., 2011).
38 Hoyer integrated NMF with sparse coding principles to derive non-negative sparse coding
39 (Hoyer, 2002), while Cai et al. introduced graph-constrained NMF to preserve intrinsic data
40 structures (Cai et al., 2011). Han et al. proposed Non-negative Laplacian Sparse Coding (Lap-

1 NMF-SPM), integrating Laplacian regularisation, NMF, and sparse coding to maintain
2 dependencies among local features (Han et al., 2015).

3 Most of these methods rely on Euclidean distance to measure feature similarity, which is not
4 always suitable for histogram-based descriptors. To overcome this limitation, Wu et al.
5 proposed the Histogram Intersection Kernel (HIK) as a more appropriate similarity measure
6 (Wu and Rehg, 2009), which has been shown to be effective for distribution-based features
7 (Vedaldi and Zisserman, 2012). Chen et al. incorporated the intersection of histograms into LLC
8 for improved environmental categorisation (Chen et al., 2018), while Wan et al. developed
9 Elastic-net and Histogram intersection-based Non-negative Local Sparse Coding (EH-NLSC),
10 jointly modelling sparsity, locality, and non-negativity (Wan et al., 2019).

11 These structured local representation methods are attractive for marine imagery because they
12 preserve fine-grained geometric patterns that are important for recognising vessels, coastlines,
13 and sea-surface structures. Nevertheless, they mainly rely on handcrafted descriptors and have
14 limited capability to capture high-level semantic information.

15 **2.2. Semantic-aware and joint representation-classification models**

16 Traditional image classification pipelines usually treat feature encoding and classifier
17 learning as separate stages, which may result in semantic information loss. To address this
18 limitation, semantic representation learning has been introduced into image representation
19 frameworks (Zhang et al., 2013). From a generative perspective, Rasiwasia et al. formulated
20 compact semantic spaces based on Gaussian mixture modelling for Scene categorization and
21 retrieval (Rasiwasia and Vasconcelos, 2008a, 2008b). Discriminative approaches incorporate
22 semantic cues more explicitly. Zhang et al. proposed a joint image representation and
23 classification framework by integrating sparse coding with a Random Semantic Space (RSS)
24 (Zhang et al., 2015). Shen et al. combined Global visual features and associated label semantics
25 within a unified segmentation and classification framework (Shen and Zeng, 2019). Additional
26 studies have explored semantic-aware sparse and joint representation-classification models to
27 improve discriminative power (Ruan et al., 2025; Xie et al., 2025).

28 Moreover, semantic-aware and joint representation-classification models have attracted
29 increasing attention in the maritime and marine engineering domain, where visual data are
30 typically affected by complex backgrounds, dynamic illumination conditions, and a scarcity of
31 annotated samples. By jointly learning visual representations and high-level semantic
32 information, these approaches mitigate semantic ambiguity and enhance feature
33 discriminability, thereby improving robustness in challenging marine environments.
34 Consequently, semantic-aware modelling strategies have demonstrated promising performance
35 in a variety of maritime vision tasks, including maritime object recognition, ship detection, and
36 coastal scene understanding. In particular, incorporating contextual and semantic information
37 into feature learning has been shown to significantly improve the interpretation of marine and
38 coastal imagery under adverse sea-surface and atmospheric conditions (Kanjir et al., 2018; Yang
39 et al., 2025).

1 Although these approaches partially bridge the gap between low-level features and semantic
2 categories, they remain largely dependent on handcrafted descriptors and local modelling
3 strategies, which limit their ability to capture complex, high-level semantics commonly
4 encountered in marine and coastal imagery. To address these limitations, the proposed HS-
5 LNLSC formulation explicitly integrates locality, non-negativity, and semantic consistency
6 constraints within a unified optimisation framework. These constraints are designed to preserve
7 neighbourhood structure, ensure physically meaningful non-negative representations, and
8 enforce consistency across semantically related local patches. Specifically, the locality
9 constraint encourages similar local descriptors to be represented using nearby dictionary atoms,
10 thereby preserving intrinsic geometric relationships. The non-negativity constraint enhances
11 interpretability by ensuring additive and physically meaningful feature contributions. In
12 addition, the semantic consistency constraint promotes coherence among local features
13 belonging to similar semantic regions, improving discriminative capability under complex
14 visual conditions.

15 Compared with conventional sparse coding variants, the proposed HS-LNLSC method
16 provides more stable and interpretable local feature representations, particularly in low-sample
17 and visually complex maritime scenarios. By jointly preserving geometric structure and
18 incorporating semantic consistency, the method improves discriminability in challenging
19 conditions such as low contrast, background clutter, and small target objects, which are
20 commonly encountered in marine and coastal imagery.

21 **2.3. Deep and self-supervised representation learning for marine imagery**

22 Deep learning has substantially advanced image representation learning, particularly through
23 Convolutional Neural Networks (CNNs) (Krizhevsky et al., 2012; Simonyan and Zisserman,
24 2015; Smirnov et al., 2014), recurrent architectures such as LSTM (Graves, 2012), and
25 Transformer-based models (Vaswani et al., 2017). ViTs model images as sequences of patches
26 and have demonstrated strong performance in image classification (Dosovitskiy et al., 2021),
27 with subsequent works improving data efficiency and architectural design (Liu et al., 2021;
28 Touvron et al., 2021).

29 Self-supervised learning approaches, such as DINO, further enhance ViT by enabling robust
30 representation learning without large labelled datasets (Caron et al., 2021). The CLS token
31 produced by pretrained DINO models encodes global semantic information and shows strong
32 transferability, making such models attractive for marine and coastal applications where
33 labelled data are scarce.

34 However, ViT-based representations primarily focus on global semantic modelling and
35 typically lack explicit mechanisms for preserving local geometric structures. Consequently,
36 they are often treated as black-box features with limited interpretability and may exhibit
37 reduced robustness in visually complex environments.

38 Recent marine remote sensing studies highlight these challenges. Bakirci et al. conducted an
39 extensive review of deep learning-based approaches for satellite-driven ship detection and
40 ocean monitoring, noting that background clutter and scale variation remain major sources of

error (Bakirci, 2025). Wang et al. proposed Ship-YOLO for ship detection in remote sensing imagery, demonstrating improvements for small targets and complex backgrounds, yet still relying on deep features alone (W. Shi et al., 2025). Beyond ship detection, deep learning-based visual analysis methods have been increasingly applied to a wide range of maritime tasks, including mooring line monitoring and oil spill detection. For instance, Khatri et al. proposed an optimised Faster R-CNN framework for mooring line detection in marine images. By refining the feature extraction strategy of ResNet-50 and reallocating high-level features to the region proposal and classification stages, the method improves localisation accuracy and robustness in complex marine scenes (Khatri et al., 2026). To address the inefficiency and limited accuracy of traditional manual inspection and radar-based approaches for oil spill detection, He et al. proposed a YOLOv12-based framework for accurate oil spill detection and boundary delineation in UAV imagery. By constructing a diverse remote sensing dataset and incorporating high-resolution inputs, pretrained initialisation, and cosine annealing optimisation, the proposed method significantly enhances detection accuracy and robustness in complex marine environments, supporting timely ecological risk assessment and response (He et al., 2026). These studies indicate that while deep learning offers strong semantic modelling capability, complementary local structural modelling remains important for robust marine image understanding.

According to the above discussion, the referenced studies primarily focus on four commonly adopted approaches, denoted as SCSPM, LSC, ViT, and DINO. For a clearer comparison, Table 1 summarises the advantages, disadvantages, and application scenarios of these four methods.

Based on Table 1, existing approaches can be broadly categorised into traditional sparse coding methods and deep learning-based models. Sparse coding techniques, such as SCSPM and LSC, emphasise local geometric structure through sparsity constraints, offering strong noise robustness and interpretability, but often suffer from high computational cost and sensitivity to parameter tuning, which limits scalability. In contrast, deep learning-based models, including ViT and DINO, excel at learning global semantic representations and long-range dependencies, particularly under self-supervised settings, albeit at the expense of substantial data and computational requirements. The complementary characteristics highlighted in Table 1 suggest that integrating sparse coding with deep semantic learning has the potential to yield more robust and efficient image representation frameworks.

Table 1. Comparative analysis of four methods.

Methods	Advantages	Disadvantages	Applications
SCSPM	(1) Efficient and compact feature encoding; (2) Feature selection and dimensionality reduction; (3) Robustness to noise and redundant data; (4) Nonlinear mapping.	(1) Higher computational complexity; (2) Sensitivity to parameters; (3) Risk of overfitting.	(1) Image processing: Image compression, denoising, and feature extraction; (2) Signal processing; (3) Machine learning; (4) Neuroscience.
LSC	(1) Preservation of local structural information; (2) Enforcement of data non-negativity;	(1) Higher computational complexity; (2) Dependency on parameters.	(1) Image compression, denoising, and feature representation; (2) Signal processing;

	(3) Addressing the instability of encoding.		(3) Bioinformatics.
ViT	(1) Strong capability to capture long-range dependencies in images; (2) Flexibility to handle images of varying sizes with patch embeddings; (3) High performance on large-scale image classification tasks.	(1) High dependence on large-scale training data for optimal performance; (2) High computational and memory costs compared to CNNs; (3) Less effective on small datasets without pretraining.	(1) Image classification and recognition; (2) Object detection and segmentation; (3) Medical image analysis and remote sensing.
DINO	(1) Self-supervised learning: does not require labelled data; (2) Generation of high-quality visual representations suitable for downstream tasks; (3) Robust to data augmentations and variations.	(1) Training can be unstable without careful hyperparameter tuning; (2) Sensitive to batch size and learning rate schedules; (3) High computational resource requirement for large-scale training.	(1) Unsupervised feature learning for image classification; (2) Object detection and segmentation pretraining; (3) Image retrieval and clustering.

2.4. Research gaps

Despite substantial progress, several critical gaps remain in the context of marine and coastal image classification:

(1) Limited incorporation of robust similarity measures and semantic cues within structured local encoding.

Most sparse coding-based methods focus on locality, non-negativity, and manifold consistency but rely on Euclidean distance and perform encoding independently of semantic information. A unified structured local encoding framework that simultaneously integrates robust similarity measurement and semantic guidance remains underexplored.

(2) Insufficient integration of structured local representations with deep semantic features.

Existing approaches typically emphasise either handcrafted local modelling or deep semantic learning. Few studies provide principled integration of these two complementary paradigms within a unified framework.

(3) Lack of feature-level fusion strategies that preserve both interpretability and semantic discrimination.

Many hybrid methods adopt heuristic fusion schemes that do not explicitly balance fine-grained local details and global semantic context, leading to suboptimal robustness and transparency.

(4) Limited validation in complex marine and maritime scenarios.

Most methods are evaluated primarily on generic benchmark datasets, with insufficient large-scale validation on maritime imagery characterised by severe illumination changes, background clutter, and limited labelled data.

In summary, these gaps motivate the development of a unified hybrid sparse-semantic image representation framework that combines structured local modelling and self-supervised deep semantic features to support robust marine and coastal image classification. The proposed HS-LNLSC-DINO framework is designed to address these challenges in a complementary and application-oriented manner.

3. Methodology

This section introduces the proposed hybrid sparse-semantic image classification framework, referred to as HS-LNLSC-DINO. The framework combines structured sparse representations with global semantic features extracted from DINO, with the goal of leveraging both local geometric structures and high-level semantic information for robust image classification. An overview of the framework is provided first, followed by detailed descriptions of each component.

3.1. The proposed framework

The proposed framework consists of five main stages, as illustrated in Fig. 1. First, local SIFT features are computed from each image to capture invariant local patterns. Second, these descriptors are encoded using Histogram Intersection-based Local Non-negative LSC (HI-LNLSC) to generate discriminative image representations. Third, semantic information is incorporated into the HI-LNLSC framework to jointly model image representation and classification, resulting in the HS-LNLSC feature representation. Through these three steps, a structured sparse representation is obtained, preserving local geometric relationships. Fourth, global semantic features are extracted from the CLS token of a pretrained DINO model, providing image-level semantic representations learned in a self-supervised manner. Simultaneously, the sparse and semantic features are fused at the feature level to form the final feature representation. Finally, classification is performed using a linear multi-class SVM.

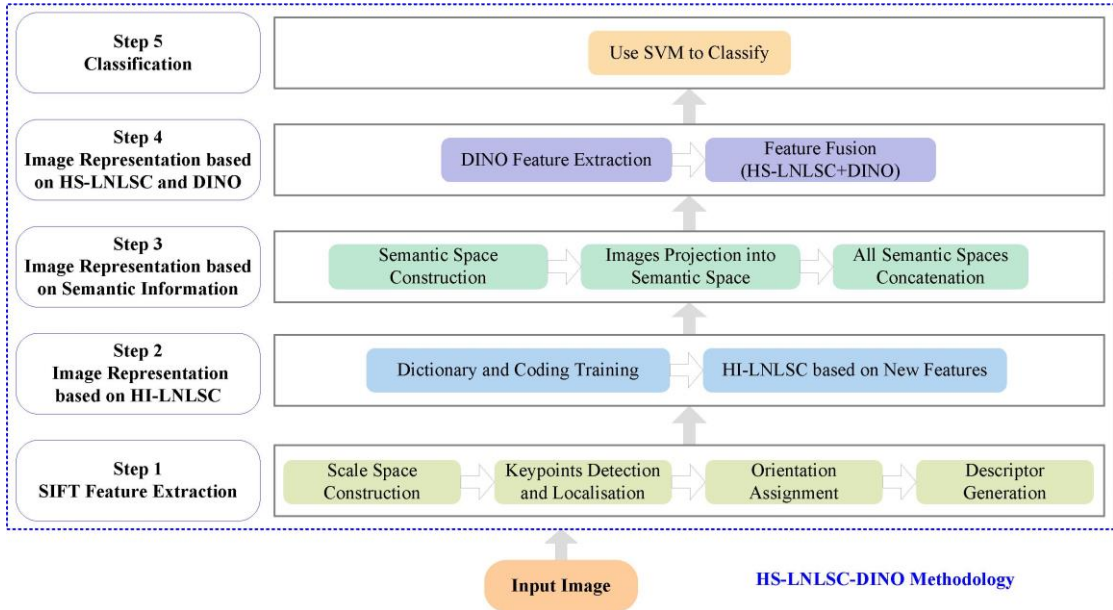


Fig. 1. An overview of the HS-LNLSC-DINO framework.

3.2. HS-LNLSC

This subsection introduces the LNLSC approach, which extends the LSC model by jointly enforcing locality and non-negativity constraints. To further improve locality modelling, the distance metric of the optimisation framework is enhanced by integrating histogram intersection. Building on this formulation, the HS-LNLSC method is developed through combining image representation and classification to leverage semantic information.

1 In the first stage, the HI-LNLSC method is used to encode local image features, with Max
 2 Pooling (MP) applied to generate the initial image representations. In the second stage, a
 3 randomly sampled subset of these image representations is used to build the semantic space. A
 4 trained classifier then projects all training images into this space, producing the image feature
 5 representations that incorporate semantic information.

6 By jointly incorporating locality and non-negativity constraints, histogram intersection-based
 7 similarity, and semantic information, this framework substantially improves the effectiveness
 8 of image representation and the overall classification performance of the HS-LNLSC method.

9 3.2.1. HI-LNLSC-based image representation

10 This section explains how the HI-LNLSC method generates the original image
 11 representations. It achieves this by embedding histogram intersection into the LNLSC
 12 optimisation framework, which redefines the similarity measure between feature vectors and
 13 dictionary atoms, thereby providing a more effective way to quantify their similarity.

14 (1) Dictionary and corresponding coding training

15 As a result of the prohibitive complexity of constructing local adaptors (Wang et al., 2010)
 16 and the Laplacian matrix for the entire set of local features, a subset of template features is
 17 randomly sampled and employed to train the dictionary and coding. The original formulation
 18 of the HI-LNLSC approach is described below. Let X represent the non-negative input feature
 19 matrix, while D and S denote the non-negative dictionary and the associated non-negative
 20 sparse representation, respectively, with the definitions of $X = [x_1, x_2, \dots, x_N] \in R^{M \times N}$,

21 $D = [d_1, d_2, \dots, d_Q] \in R^{M \times Q}$, $V = [v_1, v_2, \dots, v_N] \in R^{Q \times N}$. The optimisation problem is formulated by
 22 integrating locality and non-negativity constraints into the LSC framework as follows:

$$\begin{cases} \min_{D, S} \sum_{i=1}^N (\|x_i - Ds_i\|_2^2 + \lambda \|b_i \odot v_i\|_2^2) + \beta \text{tr}(VLV^T) \\ \text{s.t.} \quad \text{sparseness}(d_j) = S_d \\ \quad \quad \|d_j\|_2^2 \leq 1, D \geq 0, V \geq 0, \forall j \end{cases} \quad (1)$$

23 where λ , β and S_d are predefined constant parameters. The sparseness measure is S_d defined
 24 by jointly incorporating the l_2 -norm and l_1 -norm, as expressed below:

$$\text{sparseness}(d_j) = \frac{\sqrt{M} - \|d_j\|_1 / \|d_j\|_2}{\sqrt{M} - 1} \quad (2)$$

25 where M represents the dimension of d_j , namely, $d_j \in R^{M \times 1}$.

26 In the present work, the conventional Euclidean distance used in the LLC framework to
 27 evaluate the similarity between feature vectors and dictionary atoms is replaced by a histogram
 28 intersection-based similarity measure. $b_i \in R^Q$ represents a local weighting operator, which is
 29 given by:

$$\begin{aligned} b_i &= \sigma / \text{dist}(x_i, D) \\ \text{dist}(x_i, D) &= [\text{dist}(x_i, d_1), \dots, \text{dist}(x_i, d_M)]^T \end{aligned} \quad (3)$$

1 Here, σ is a parameter that regulates the strength of weight regularisation, and $dist(x_i, d_j)$
 2 denotes the distance between x_i and d_j , evaluated based on histogram intersection. The
 3 computation can be expressed as:

$$4 \quad dist(x_i, d_j) = \sum_{k=1}^Q \min(x_{ik}, d_{jk}) \quad (4)$$

5 Here, Q denotes the histogram dimension, which corresponds to the number of dictionary
 6 atoms. Additionally, x_{ik} and d_{jk} denote the k -th components of the feature vector x_i and the
 7 dictionary atom d_j , respectively.

8 To solve Eq. (1), an alternating optimisation strategy is adopted, in which D and V are
 9 iteratively updated by fixing one variable while optimising the other. First, with X and D held
 10 fixed, V is optimised. By applying the Method of Lagrange Multiplier (MLM) along with the
 11 Karush-Kuhn-Tucker (KKT) theorem, the corresponding update strategy for V is derived as
 12 follows:

$$13 \quad v_{ij} = v_{ij} \frac{(D^T X + \beta V W)_{ij}}{(D^T D V + \beta V D + \lambda \text{diag}(g_i) \bar{V})_{ij}}, \quad (5)$$

$$14 \quad \forall i = 1, 2, \dots, M, \quad j = 1, 2, \dots, N$$

15 where $\text{diag}(g_i)$ denotes a diagonal matrix whose diagonal entries are given by g_i ,
 16 $g_i = [b_{i1}^2, b_{i2}^2, \dots, b_{iN}^2]$ is an N -dimensional vector represented as a row vector, and
 17 $b_{ij}^2 = (\sigma / \text{dist}(x_j, d_i))^2$ and $\bar{V} = [\text{diag}(v_1), \text{diag}(v_2), \dots, \text{diag}(v_N)]^T$ represent vectors of ones and
 18 zeros, respectively.

19 Next, with X and V fixed, the optimisation is performed with respect to D . For this problem,
 20 the diagonal matrix Λ is derived by solving the corresponding Lagrange dual formulation
 21 using the conjugate gradient method. The result Λ is then substituted into the following
 22 expression to compute D :

$$23 \quad D = (XV^T)(VV^T + \Lambda)^{-1} \quad (6)$$

24 (2) HI-LNLSC based on new features

25 A set of template features, X , is randomly selected in Section 3.2.1 (1) to train the matrices D
 26 and V . For an unseen local feature matrix Y , the learned matrices D and V are directly used for
 27 HI-LNLSC encoding, leading to the following optimisation formulation:

$$28 \quad \begin{cases} \min_V \|Y - DS\|_F^2 + \lambda \|b \odot S\|_F^2 + \frac{\beta}{2} \sum_{ji} \|s_j - v_i\|_2^2 w_{ji} \\ \text{s.t. } s_{ij} \geq 0, \forall i, j \end{cases} \quad (7)$$

29 Here, v_i denotes the i -th element of the column vector in V , and S denotes the sparse
 30 representation of the updated feature matrix Y . The similarity matrix W is constructed using a
 31 K -nearest neighbour strategy based on histogram intersection. If a template feature x_i is among
 32 the K -nearest neighbours of a new feature y_j , then w_{ji} is set to 1; otherwise, it is set to 0. This
 33 similarity metric is consistent with Eq. (4).

34 The update rule for S can be derived by applying the MLM and the KKT conditions, as
 35 outlined below.

$$s_{ij} = s_{ij} \frac{(D^T X + \beta V W^T)_{ij}}{(D^T D S + \frac{1}{2} \beta S A + \lambda \text{diag}(g_i) \bar{S})_{ij}} \quad (8)$$

where A is a diagonal matrix with $a_{jj} = \sum_i w_{ji}$ as its diagonal elements, $\text{diag}(g_i)$ follows the definition provided in Eq. (5), and $\bar{S} = [\text{diag}(s_1), \text{diag}(s_2), \dots, \text{diag}(s_N)]^T$ denotes a vector of ones. After training D and V on the template features, Eq. (8) is applied to encode new features using HI-LNLSC.

To achieve feature fusion, the MP method is utilised in this paper. The detailed procedure is outlined below.

$$f_l = \max\{|s_{1l}|, |s_{2l}|, \dots, |s_{Nl}|\}, l = 1, 2, \dots, Q \quad (9)$$

Here, the l -th component of the sparse representation s_N is denoted by s_{Nl} , and f_l corresponds to the l -th component of the vector f . Consequently, each spatial pyramid region is represented by a Q -dimensional descriptor f , as defined by $f = [f_1, f_2, \dots, f_Q]^T$.

After computing the region-level representations, the SPM method is applied to derive the final image descriptor.

3.2.2. Semantic information-based image representation

To address the limitations of the HI-LNLSC method, specifically its relatively independent processing, as well as the semantic gap between visual features and human interpretation, the HS-LNLSC approach is introduced. This approach aims to effectively combine both visual cues and semantic information, thereby modelling interactions between semantic entities and their contextual surroundings. First, the HI-LNLSC approach is used to perform encoding of local image features, resulting in an initial image representation. Next, a semantic space is created to refine and derive the resulting visual representation. Ultimately, a SVM is utilised to categorise the images based on these semantic-aware representations.

(1) Semantic space construction.

The semantic space is formulated as a unified representation domain constructed from various visual representations, providing a basis for semantic modelling of images. Each semantic space is constructed via classifier training on a randomly sampled image subset.

Let X denote the initial image feature representations of H training samples obtained utilising the HI-LNLSC approach, where the samples belong to E classes with corresponding labels denoted by z . A subset of P ($P \leq H$) images is randomly selected from these representations to build the semantic space, and the above procedure is iteratively performed R times. The resulting sets from each iteration are denoted as

$\{(u^{1,1}, z^{1,1}), \dots, (u^{P,1}, z^{P,1})\}, \dots, \{(u^{1,R}, z^{1,R}), \dots, (u^{P,R}, z^{P,R})\}$. For the r -th random selection

$\{(u^{1,r}, z^{1,r}), \dots, (u^{P,r}, z^{P,r})\}$, a SVM is employed to construct the corresponding semantic space,

as follows:

$$f_e^r(z^{p,r}) = \bar{z}_e^{p,r} = w_e^r u^{p,r} + b_e^r, \quad p=1,2,\dots,P, \quad e=1,2,\dots,E \quad (10)$$

An optimisation problem based on the hinge loss function is subsequently formulated as follows:

$$\min_{w_e^r} \|w_e^r\|^2 + \alpha \sum_{p=1}^P p(\bar{z}_e^{p,r}, z^{p,r}) \quad (11)$$

Solving Eq. (11) yields the corresponding values of w_e^r and b_e^r . Each component of the resulting semantic space is associated with an individual classifier, which is trained using a randomly chosen subset of samples. Given E image classes, the resulting semantic space consists of E dimensions.

(2) Image projection into semantic space.

After training the SVM, all training images are projected into the constructed semantic space as follows:

$$u_{r,k}^{ss} = (f_1^r(u^h); f_2^r(u^h); \dots; f_E^r(u^h)), \quad r=1,2,\dots,R; \quad h=1,2,\dots,H \quad (12)$$

Here, the superscript ‘ss’ is used to indicate the ‘semantic space’.

(3) All semantic spaces concatenation.

Once all spaces are constructed, all training images are mapped into each of the R learned spaces. The final image representation is then obtained by concatenating the corresponding feature vectors across these spaces, as follows:

$$u_k^{ss} = (u_{1,h}^{ss}; u_{2,h}^{ss}; \dots; u_{R,h}^{ss}), \quad h=1,2,\dots,H \quad (13)$$

Thus, the locally structured image representation is obtained via Eq. (13).

3.3. HS-LNLSC-DINO

While HS-LNLSC effectively captures local geometric structures and preserves discriminative local relationships through locality, non-negativity, and histogram intersection constraints, it primarily relies on handcrafted local descriptors and lacks explicit modelling of high-level semantic information. In contrast, DINO excels at learning rich global semantic representations but does not explicitly encode local geometric structures, often treating these features as black-box representations. Motivated by the complementary strengths of these two approaches, this section introduces a hybrid sparse-semantic image representation framework, termed HS-LNLSC-DINO, which integrates structured sparse coding with self-supervised semantic features.

3.3.1. DINO feature extraction

In parallel with local sparse encoding, global semantic features are extracted using DINO. Given the same input image, it is divided into fixed-size, non-overlapping patches, which are then processed by the Transformer encoder. Let $f_{CLS} \in R^{Dino}$ denote the output CLS token of DINO. This token encodes holistic, image-level semantic information learned through self-distillation, without the need for labelled supervision, and has demonstrated strong transferability across various visual tasks.

1 Unlike local sparse features, f_{CLS} captures long-range contextual dependencies and high-
2 level semantics, making it particularly well-suited for complex marine and maritime scenes,
3 where global context plays a crucial role.

4 3.3.2. Sparse-semantic feature fusion

5 To leverage the complementary strengths of structured sparse representations and deep
6 semantic features, the HS-LNLSC and DINO features are fused at the feature level. The local
7 structural features of the image, obtained through Eq. (13), are denoted by $f_{HS-LNLSC}^T$, and the
8 final image representation is constructed as follows:

$$f = [f_{HS-LNLSC}^T, \alpha f_{CLS}^T]^T \quad (14)$$

9 where α is a scaling factor that balances the contributions of the local sparse features and global
10 semantic features.

11 This fusion strategy preserves fine-grained local geometric details through HS-LNLSC,
12 while simultaneously incorporating global semantic understanding via the DINO CLS token.
13 As a result, the fused representation achieves enhanced robustness, improved discriminative
14 capability, and better interpretability compared to purely sparse or purely deep representations.

15 Finally, the fused feature vector f is fed into a multi-class linear SVM for classification. The
16 use of a linear classifier further emphasises the discriminative power of the proposed
17 representation without introducing additional model complexity.

18 Rather than relying on simple feature concatenation, the proposed framework adopts a
19 principled feature-level integration of structured local descriptors derived from HS-LNLSC and
20 global semantic representations obtained from DINO CLS tokens. This hybrid design
21 effectively combines fine-grained geometric information with high-level contextual semantics,
22 thereby enhancing discriminative capability and robustness under data-limited conditions
23 typical of maritime image analysis. By leveraging the complementary strengths of local and
24 global representations, the framework achieves more reliable and interpretable classification
25 performance than either component used in isolation.

26 3.4. Algorithmic description

27 Once the formulations of matrices D and V are established, the dictionary and corresponding
28 sparse representations of the template features are obtained via a two-stage procedure. First,
29 Algorithm 2 iteratively updates the diagonal matrix D associated with the Lagrange multipliers.
30 Then, Algorithm 3 identifies the optimal approximation that enforces the desired sparsity, with
31 D updated accordingly. The specific algorithmic procedure is provided in Appendix A and B.

32 Algorithm 1 is then employed to learn both the dictionary for HS-LNLSC-DINO and its
33 associated sparse codes. This algorithm iteratively updates both D and V until a predefined
34 stopping criterion is reached. Notably, Algorithm 1 consolidates the operations of Algorithms 2
35 and 3 in the Appendix.

36 In conclusion, the overall procedure is summarised below.

Algorithm 1 (HS-LNLSC-DINO)

Input: Non-negative feature matrix X ; Initial dictionary D ; Initial sparse coding V ; Laplacian matrix L ; Parameters $\lambda, \beta, \sigma, S_d$; Training samples H ; Training times R .

Output: Updated dictionary D ; Sparse coding V ; Category annotations.

Step 1. Initialise the variables X and V , namely, $X = X / \max(X(:))$, $V = V / \|V\|_1$.

Step 2. Convergence Check: While convergence criteria are not met, repeat the following steps:

Step 3. Update sparse coding V with Eq. (5).

Step 4. Normalise the dictionary D and sparse coding V , namely, $d_{ij} = d_{ij} / \sqrt{\sum_i v_{ij}}$, $v_{ij} = v_{ij} / \sqrt{\sum_i v_{ij}}$.

Step 5. Update Lagrange dual matrix Λ according to **Algorithm 2**.

Step 6. Each column vector of matrix D is projected individually according to **Algorithm 3** to acquire the updated values. Let D^* represent the optimal dictionary and V represent the corresponding sparse coding.

Step 7. Check the convergence condition for Step 8 in **Algorithm 2**.

Step 8. If the convergence condition is satisfied, proceed to the next step.

Step 9. If not, return to Step 3 and repeat the process.

Step 10. Post-Processing: After updating D and V , the sparse codes S for the new features are computed using Eq. (8).

Step 11. Apply SPM via Eq. (9) to implement MP for the encoded features and acquire the initial image representation.

Step 12. Semantic Space Construction: Eq. (10) is employed to build the semantic space, after which w_e^r and b_e^r are obtained using Eq. (11).

Step 13. All images in the training set are embedded within the semantic space in accordance with Eq. (12).

Step 14. Integrate all semantic spaces to form the final image representation using Eq. (13).

Step 15. Fuse the HS-LNLSC and DINO features to obtain HS-LNLSC-DINO features according to Eq. (14).

Step 16. Implement a SVM to categorise the image based on its semantic representation.

4. Experiments

The section evaluates the proposed hybrid sparse-semantic framework in terms of accuracy, robustness, and practical suitability for marine and coastal monitoring tasks. The experiments are organised as follows. Section 4.1 introduces the datasets and implementation settings, with particular attention to the characteristics that reflect operational marine imaging conditions (e.g., clutter, scale variation, and limited labelled data). Section 4.2 presents ablation studies to quantify the contribution of each component and to examine design choices relevant to deployment (feature selection, Transformer layer depth, and training data availability). Section 4.3 benchmarks the proposed approach against representative methods to assess its competitiveness under standard evaluation protocols. Section 4.4 investigates generalisation and stability to evaluate reliability under repeated sampling and across diverse datasets. **Section 4.5 presents a discussion of the computational complexity of the proposed algorithm.**

4.1. Datasets and experimental settings

4.1.1. Datasets description

To evaluate performance across both general-purpose and domain-specific conditions, experiments are conducted on four widely used benchmark datasets (Corel-10, Scene-15, Caltech-101, and Caltech-256) and three maritime image datasets representing realistic marine visual complexity. Table 2 summarises the benchmark datasets, and Fig. 2 provides example images from Corel-10. These benchmark datasets are included to ensure comparability with established baselines and to verify that the proposed framework is not domain-specific.

1 For marine and coastal relevance, three maritime datasets are considered: Open Seaship,
 2 Seaship-trimming, and SMD-trimming. The Open Seaship dataset comprises 31,455 images
 3 across seven common ship categories (passenger ships, mixed types, container ships, bulk
 4 carriers, ore carriers, general cargo ships, and fishing vessels), as summarised in Table 3. **In
 5 addition, the Seaship dataset represents the complete collection of maritime images, whereas
 6 Seaship-trimming refers to a pre-processed subset specifically employed for ablation studies or
 7 small-sample experiments. This distinction has been applied consistently throughout the
 8 manuscript to ensure clarity and avoid potential confusion.** The Singapore Maritime Dataset
 9 (SMD) includes imagery captured in diverse operational settings (on-board, near-infrared, and
 10 on-shore) with substantial variability in background, illumination, and target scale, as
 11 summarised in Table 4. Since the original SMD was designed for detection rather than
 12 classification, experiments employ the SMD-trimming dataset to enable controlled evaluation
 13 within an image classification pipeline. Example images from Seaship are shown in Fig. 3.

14 Overall, the selected maritime datasets reflect key challenges encountered in practice,
 15 including background interference, small targets, viewpoint changes, and limited labelled
 16 samples, conditions that directly affect the reliability of monitoring outputs in marine and
 17 coastal management settings.

18 Table 2. Four benchmark image datasets.

Datasets	Classes	Images Per Class	Total Images
Corel-10	10	100	1000
Scene-15	15	200~400	4485
Caltech-101	101	31~800	9144
Caltech-256	256	≥ 80	29,780

16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

bus



elephant



flower



building



1
2

Fig. 2. Sample images from the Corel-10 dataset.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

fishing boat



container ship



ore carrier



passenger ship



Fig. 3. Some pictures of the Seaship.

Table 3. The Seaship dataset.

Ship Class	Images	Percentage
Passenger ship	3171	0.1008
Mixed type	3440	0.1094
Container ship	3657	0.1163
Bulk cargo carrier	5067	0.1610

Ore carrier	5126	0.1630
General cargo ship	5342	0.1698
Fishing boat	5652	0.1797
Total	31455	1

Table 4. The SMD dataset.

Subdataset	Videos (Annotated)	Labelled Frames	Number of Labels
VIS on-board	11 (4)	2400	3173
NIR	30 (23)	11,286	83,174
VIS on-shore	40 (36)	17,967	154,495
Total	81 (63)	31,653	240,842

4.1.2. Experimental settings

The experimental pipeline follows a hybrid representation strategy combining structured local encoding and self-supervised global semantic representations.

Structured local feature extraction and coding. Dense SIFT descriptors are extracted using a sliding window (step size: 8 pixels) with a fixed window size. Each local descriptor has dimensionality 128. A dictionary of size 1024 is learned for encoding. **To clarify the selection of key parameters in the HS-LNLSC-DINO framework, we provide the following rationale.**

The locality regularisation parameter (λ), Laplacian smoothness (β), and histogram intersection weight (σ) were determined through a combination of empirical stability analysis across multiple datasets and guidance from previous studies (Han et al., 2015; Wang et al., 2010). Preliminary experiments demonstrated that the model’s performance is relatively insensitive to small variations around these values, indicating that the chosen default settings achieve consistent and robust performance without the need for dataset-specific tuning. The fusion parameter α , which governs the relative contribution of local and global features, was set to 0.4 to ensure a stable and balanced integration of HS-LNLSC and DINO representations. Similarly, the sparsity parameter S_d was set to 0.4, providing sufficiently discriminative sparse coding while preventing excessively sparse representations that could compromise model performance.

Semantic space construction. A semantic space is constructed using a fixed proportion of images ($P=0.3H$) and repeated sampling ($R=30$) to reduce sensitivity to data selection and to improve robustness.

DINO feature extraction. Images are resized to a fixed resolution prior to feeding into a pretrained DINO ViT with patch size 16. The network weights are frozen during feature extraction to reduce overfitting risk, which is particularly relevant for operational marine datasets where labelled data may be limited. The CLS token is used as the global semantic representation. Key configurations are summarised in Table 5.

Classification and evaluation. For all experiments, classification is performed using a multi-class linear SVM to ensure consistent evaluation and to isolate the effect of representation learning from classifier complexity. Unless otherwise specified, accuracy is reported as mean \pm standard deviation across repeated runs.

To evaluate model performance under realistic operational constraints, training protocols with varying numbers of samples per class (50, 100, 150, 200, and 250) are adopted, simulating low-data conditions common in maritime monitoring scenarios. The rationale for these choices is explicitly described to clarify that the protocols are designed to assess robustness under limited labelled data rather than arbitrary selection.

Table 5. Parameter settings used in each stage of the experiment.

Image Datasets	Corel-10	Scene-15	Caltech-101	Caltech-256	Seaship	Seaship-trimming	SMD-trimming
The SIFT feature extraction stage							
Approaches	SIFT	SIFT	SIFT	SIFT	SIFT	SIFT	SIFT
Step size	8	8	8	8	8	8	8
Window size	16×16	16×16	16×16	16×16	16×16	16×16	16×16
The HS-LNLSC stage							
Dictionary size	1024	1024	1024	1024	1024	1024	1024
λ	0.4	0.4	0.4	0.4	0.4	0.4	0.4
β	0.2	0.2	0.2	0.2	0.2	0.2	0.2
σ	100	100	100	100	100	100	100
Sparseness							
S_d	0.4	0.4	0.4	0.4	0.4	0.4	0.4
The construction of the semantic space stage							
The number of images utilised to form the semantic space (P)	$0.3H$	$0.3H$	$0.3H$	$0.3H$	$0.3H$	$0.3H$	$0.3H$
The times of repetition selected (R)	30	30	30	30	30	30	30
The DINO feature extraction stage							
The size of the input image	320×320	320×320	320×320	320×320	320×320	320×320	320×320
Patch size	16×16	16×16	16×16	16×16	16×16	16×16	16×16
α	0.4	0.4	0.4	0.4	0.4	0.4	0.4

4.2. Ablation studies

To identify the components that most critically influence performance under realistic marine imaging conditions, a comprehensive set of ablation studies is conducted from three complementary perspectives: (i) feature source and fusion strategy, (ii) the choice of Transformer layers for semantic representation learning, and (iii) sensitivity to the size of the training dataset. Collectively, these analyses provide actionable insights for practical deployment, particularly with respect to computational efficiency, annotation burden, and robustness in data-scarce scenarios.

4.2.1. Selection of feature combinations

This experiment investigates the relative contributions of different feature sources and their combinations to classification performance, with the objective of identifying the most effective

1 representation strategy within the proposed framework. Three categories of features are
 2 considered: (1) structured sparse representations generated by HS-LNLSC, (2) patch-level
 3 features derived from intermediate DINO representations (denoted as DINO-Patch), and (3)
 4 global semantic features extracted from the CLS token of DINO (denoted as DINO-CLS).
 5 Based on these feature types, five representative configurations are evaluated: HS-LNLSC,
 6 DINO-Patch, DINO-CLS, HS-LNLSC+DINO-Patch, and HS-LNLSC+DINO-CLS.

7 Experiments are conducted on two representative datasets, **Seaship** and Scene-15, under
 8 identical experimental settings. For each class, 50 images are randomly selected for training,
 9 with the remaining samples reserved for testing. To further investigate the influence of semantic
 10 feature depth, DINO representations are extracted from multiple Transformer layers (L2, L4,
 11 L6, L8, L10, and L12). A multi-class linear SVM classifier is employed for all configurations
 12 to ensure a fair and consistent comparison.

13 The quantitative results are reported in Tables 6 and 7, with corresponding performance
 14 trends illustrated in Fig. 4. Several important observations emerge. First, HS-LNLSC alone
 15 yields stable and competitive performance across all Transformer layers on both datasets,
 16 indicating strong robustness and minimal sensitivity to semantic depth. This behaviour is
 17 expected, as HS-LNLSC is built upon handcrafted local descriptors and structured sparse
 18 coding, which are independent of Transformer-based semantic representations.

19 Second, DINO-CLS consistently outperforms DINO-Patch on both datasets, and its
 20 performance improves markedly as deeper Transformer layers are used. This finding suggests
 21 that the CLS token progressively encodes more discriminative global semantic information at
 22 deeper layers, whereas patch-level features mainly emphasise local visual patterns that are less
 23 effective when used in isolation for classification.

24 Third, the combination of HS-LNLSC and DINO-CLS achieves the best performance across
 25 nearly all layers and datasets. HS-LNLSC+DINO-CLS consistently surpasses both individual
 26 components, demonstrating strong complementarity between structured local sparse
 27 representations and global semantic features. By contrast, integrating HS-LNLSC with DINO-
 28 Patch yields only marginal gains or even performance degradation, implying that patch-level
 29 semantic features provide limited additional information beyond that already captured by HS-
 30 LNLSC.

31 Overall, these results indicate that HS-LNLSC+DINO-CLS constitutes the most effective
 32 feature combination. HS-LNLSC contributes stable and interpretable local geometric modelling,
 33 while DINO-CLS supplies high-level semantic discrimination. Their integration produces a
 34 unified representation that is both robust and semantically expressive. Consequently, HS-
 35 LNLSC+DINO-CLS is selected as the optimal feature configuration and adopted in all
 36 subsequent experiments.

37 Table 6. Classification accuracy on the Seaship dataset with different feature combinations (%).

	Layers	L2	L4	L6	L8	L10	L12
Methods							
	HS-LNLSC	78.73	78.73	78.73	78.73	78.73	78.73
	DINO-Patch	55.94	57.18	56.54	54.86	56.05	53.08
	DINO-CLS	41.13	55.97	77.53	84.64	89.48	90.88

HS-LNLSC+DINO-Patch	56.57	57.98	57.28	55.39	55.51	55.15
HS-LNLSC+DINO-CLS	79.52	80.93	81.12	87.14	90.92	92.76

Table 7. Classification accuracy on the Scene-15 dataset with different feature combinations (%).

Methods	Layers					
	L2	L4	L6	L8	L10	L12
HS-LNLSC	87.53	87.53	87.53	87.53	87.53	87.53
DINO-Patch	68.64	67.90	67.35	66.92	67.98	68.21
DINO-CLS	60.74	82.68	90.45	94.39	96.28	96.27
HS-LNLSC+DINO-Patch	67.95	68.95	69.21	67.45	69.25	68.15
HS-LNLSC+DINO-CLS	88.46	90.01	93.55	95.19	96.59	97.51

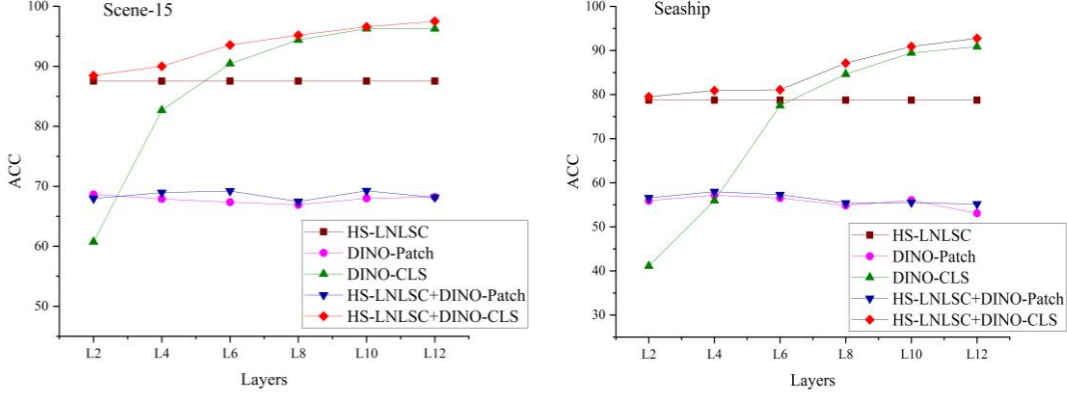


Fig. 4. Classification accuracy for five different methods in two image datasets.

4.2.2. Influence of Transformer layer selection

Following the identification of HS-LNLSC+DINO-CLS as the most effective feature combination in Section 4.2.1, the impact of Transformer layer selection on overall classification performance is further examined. Because different layers of a ViT encode semantic information at progressively higher levels of abstraction, determining the layer that provides the most discriminative and complementary semantic representation when fused with HS-LNLSC is of practical importance.

In this experiment, CLS-token features are extracted from multiple Transformer layers of the DINO model, namely L2, L4, L6, L8, L10, and L12, while the HS-LNLSC component is kept unchanged. **To further investigate the robustness of the HS-LNLSC-DINO framework under more constrained training conditions, we reduce the number of training samples per class to 30, while maintaining the same protocol as in Section 4.2.1. This allows assessment of relative performance across ablation variants without affecting comparability.**

The quantitative results are reported in Table 8, and the corresponding trends are illustrated in Fig. 5. A clear performance increase is observed as the Transformer depth grows, indicating that deeper layers encode increasingly discriminative semantic information. In contrast, shallow layers (L2 and L4) mainly capture low-level visual cues such as edges and textures, which offer limited additional complementarity to HS-LNLSC and therefore lead to smaller performance gains.

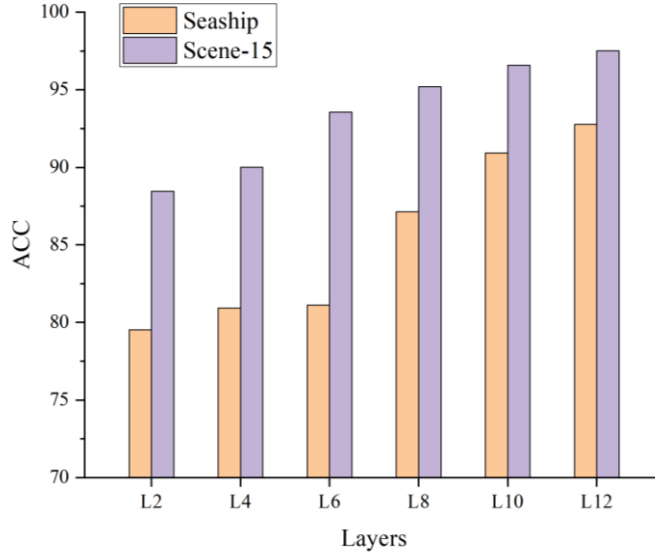
Mid-to-deep layers (L6 and L8) yield substantial improvements, suggesting a favourable trade-off between semantic abstraction and generalisation capability. Among these, L8 consistently demonstrates strong performance across both datasets, indicating that features at

1 this depth effectively complement the structured sparse representations produced by HS-
 2 LNLSC.

3 The highest classification accuracy is consistently obtained using features from L12. This
 4 outcome indicates that semantic representations extracted from the deepest Transformer layer
 5 are the most beneficial when integrated with HS-LNLSC. Consequently, Layer 12 (L12) is
 6 selected as the optimal semantic feature extraction layer and is adopted in all subsequent
 7 experiments, including analyses of training-sample sensitivity and generalisation performance.

8 Table 8. Classification accuracy for different Transformer layers on the two image datasets (%).

Layers Datasets	L2	L4	L6	L8	L10	L12
Seaship	79.52	80.93	81.12	87.14	90.92	92.76
Scene-15	88.46	90.01	93.55	95.19	96.59	97.51



10 Fig. 5. Classification accuracy for different Transformer layers on the two image datasets.

11 4.2.3. Impact of the number of training samples

12 This subsection evaluates the robustness of the proposed HS-LNLSC-DINO framework
 13 under varying training sample sizes, a critical consideration for real-world marine and maritime
 14 applications where annotated data are often scarce. The objective is to characterise the evolution
 15 of classification performance as the amount of training data increases and to assess the data
 16 efficiency of the proposed approach.

17 Experiments are conducted on three datasets with distinct characteristics: **Seaship**, SMD-
 18 trimming, and Scene-15. For the **Seaship** and SMD-trimming datasets, the number of training
 19 samples per class is set to 50, 100, 150, 200, and 250, whereas for the Scene-15 dataset, 50,
 20 100, 150, and 200 samples per class are considered. In all cases, the remaining samples are
 21 reserved for testing. The optimal configuration identified in the preceding subsections, namely,
 22 the HS-LNLSC+DINO-CLS feature fusion strategy with semantic features extracted from
 23 Layer 12 (L12), is adopted throughout.

The quantitative results are reported in Table 9 and visualised in Fig. 6. As expected, classification accuracy increases monotonically with the number of training samples, reflecting typical learning behaviour. More importantly, strong performance is observed even in low-data regimes. With only 50 training samples per class, the HS-LNLSC-DINO framework achieves competitive accuracy across all three datasets, demonstrating high data efficiency.

This favourable behaviour can be attributed to the complementary nature of the hybrid representation. The structured sparse coding component captures stable and discriminative local geometric patterns that generalise well with limited supervision, while the self-supervised DINO features contribute transferable semantic knowledge learned from large-scale unlabeled data. As additional training samples become available, the classifier is able to further exploit these rich representations, leading to consistent performance gains. Collectively, these results indicate that the proposed framework not only achieves high accuracy with sufficient training data but also demonstrates strong robustness and generalisation in low-sample scenarios, underscoring its suitability for practical marine and maritime image classification tasks.

Table 9. Classification accuracy for different Training images across the three image datasets (%).

Datasets	Training Images				
	50	100	150	200	250
Seaship	92.76	96.38	97.85	98.12	98.75
SMD-trimming	94.90	97.35	98.65	99.13	99.68
Scene-15	97.51	98.50	99.26	99.61	/

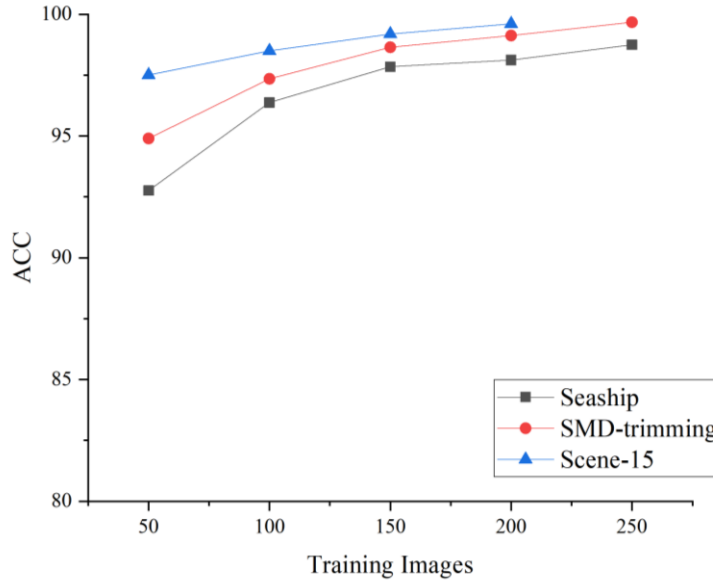


Fig. 6. Classification accuracy for different Training images on the three image datasets.

4.3. Comparison with state-of-the-art methods

To comprehensively assess the effectiveness of the proposed HS-LNLSC-DINO framework, extensive comparisons are conducted with a broad spectrum of representative state-of-the-art image classification methods. The selected methods encompass diverse image representation paradigms, such as SCSPM, LLC, LSC, Lap-NMF-SPM, EH-NLSC, RSS, and HS-NLLSC (Y. Shi et al., 2025) approaches.

1 Regarding the four benchmark datasets, different training-testing splits are adopted according
2 to commonly used evaluation protocols. In particular, with respect to the Corel-10 and Scene-
3 15 datasets, 50 and 100 images from each class are randomly chosen to train the model, while
4 the leftover images serve as the testing set. For the Caltech-101 dataset, experiments are
5 conducted using 15 and 30 training images per category, while all remaining images are
6 reserved for testing. In the case of Caltech-256, 15, 30, 45, and 60 images per class are randomly
7 selected for training, and the remaining images are used for testing. For the three maritime
8 datasets, the number of training samples per category is varied to assess performance under
9 different data availability conditions. Specifically, 50, 100, 150, 200, and 250 images from each
10 class are randomly chosen for training, while the remaining images are allocated for testing.

11 Table 10 presents the comparative results between the proposed HS-LNLSC-DINO method
12 and several representative cutting-edge approaches over four benchmark datasets. The
13 evaluated methods include classical sparse coding models (SCSPM, LLC, LSC), structured
14 sparse representation techniques (Lap-NMF-SPM, EH-NLSC), semantic-based models (RSS),
15 HS-NLLSC, and deep feature-based baselines. Across all datasets and experimental settings,
16 HS-LNLSC-DINO consistently achieves the highest classification accuracy, demonstrating
17 clear and robust performance advantages.

18 **In addition, statistical significance analyses were performed to evaluate the reliability of HS-
19 LNLSC-DINO’s performance gains relative to the second-best baseline across all datasets and
20 training sample sizes. As shown in Table 10, the reported P-values indicate that most
21 improvements are statistically significant ($P < 0.05$), suggesting that the observed advantages
22 are unlikely to arise from random variation. These findings further support the robustness and
23 generalizability of the proposed framework across diverse datasets and varying sample
24 conditions.**

25 The superior performance of HS-LNLSC-DINO can be attributed to its integrated sparse-
26 semantic representation mechanism. The HS-LNLSC component incorporates locality, non-
27 negativity constraints, Laplacian regularisation, and histogram intersection similarity,
28 facilitating the stable and consistent encoding of local SIFT features while preserving spatial
29 geometric relationships. This design effectively addresses the coding instability and feature
30 cancellation issues commonly encountered in traditional sparse coding frameworks. In parallel,
31 the incorporation of global semantic features extracted from DINO introduces high-level
32 contextual information that is challenging to capture using handcrafted descriptors alone.

33 By contrast, traditional sparse coding approaches, such as SCSPM, LSC, and Lap-NMF-SPM,
34 primarily rely on Euclidean-distance-based encoding and operate independently of semantic
35 information, which limits their discriminative capability in complex scenes. While EH-NLSC
36 enhances similarity measurement through histogram intersection, the absence of Laplacian
37 regularisation reduces its ability to effectively model spatial structure. Semantic-based methods
38 like RSS incorporate label information at the representation level but still rely on unstable
39 sparse coding during the encoding process, resulting in limited performance improvements.

Through the joint exploitation of structured local representations and globally discriminative semantic features, HS-LNLSC-DINO effectively bridges the gap between handcrafted sparse modelling and deep representation learning. This complementary integration results in more robust, informative, and semantically consistent image representations, leading to significant performance improvements over existing methods.

Table 10. Classification accuracy (mean \pm standard deviation) across four benchmark image datasets (%); P-value versus second-best.

Methods	Corel-10	Scene-15	Caltech-101 (15)	Caltech-101 (30)	Caltech-256 (15)	Caltech-256 (30)	Caltech-256 (45)	Caltech-256 (60)
SCSPM	86.76 \pm 1.18	81.12 \pm 0.45	66.87 \pm 0.45	72.10 \pm 1.14	27.53 \pm 0.42	33.86 \pm 0.55	37.35 \pm 1.64	40.08 \pm 0.79
LLC	87.83 \pm 1.03	81.53 \pm 0.87	68.57 \pm 0.88	72.54 \pm 0.71	31.27 \pm 0.85	34.17 \pm 0.33	35.93 \pm 0.51	37.58 \pm 0.49
LSC	88.43 \pm 0.75	89.65 \pm 0.41	70.32 \pm 1.35	74.86 \pm 0.53	29.88 \pm 0.15	35.67 \pm 0.33	38.37 \pm 0.46	40.35 \pm 0.24
Lap-NMF-SPM	91.24 \pm 0.95	90.46 \pm 0.87	74.35 \pm 0.94	76.81 \pm 0.49	35.24 \pm 0.83	37.46 \pm 0.32	39.87 \pm 0.75	41.35 \pm 0.72
EH-NLSC	93.64 \pm 0.78	91.82 \pm 0.67	73.34 \pm 0.62	78.89 \pm 0.39	35.89 \pm 0.56	38.87 \pm 0.59	41.65 \pm 0.53	43.61 \pm 0.47
RSS	95.72 \pm 0.78	92.45 \pm 0.93	77.63 \pm 0.89	82.91 \pm 0.22	40.16 \pm 0.53	44.96 \pm 0.85	48.25 \pm 0.47	51.32 \pm 0.41
HS-NLLSC	98.86 \pm 0.23	97.56 \pm 0.64	81.54 \pm 0.48	87.73 \pm 0.67	46.34 \pm 0.45	49.86 \pm 0.65	53.60 \pm 0.78	57.68 \pm 0.12
HS-LNLSC-DINO	99.84\pm0.17	98.50\pm0.12	97.63\pm0.11	99.21\pm0.16	87.24\pm0.19	92.07\pm0.16	94.86\pm0.26	96.86\pm0.11
P-value (vs second-best)	5.95×10^{-3}	1.21	1.37×10^{-10}	2.95×10^{-8}	6.66×10^{-14}	7.11×10^{-13}	4.44×10^{-12}	0.00
Significant (P < 5%)	*	*	*	*	*	*	*	*

For clearer visualisation, Fig. 7 presents the mean classification accuracy and corresponding standard deviation across seven key methods over four benchmark datasets. HS-LNLSC-DINO consistently achieves the best performance. Compared with HS-NLLSC and other competing approaches, accuracy gains range from approximately 1% to 16% across the datasets.

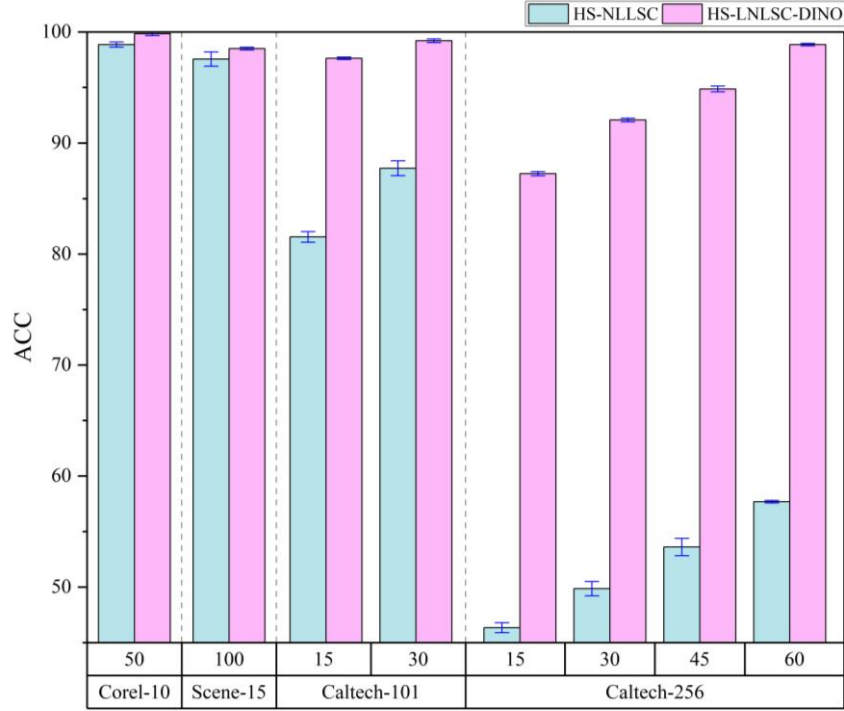


Fig. 7. Classification accuracy (mean \pm standard deviation) across four benchmark image datasets.

With particular emphasis on ship classification, performance is evaluated on the **Seaship**, **Seaship-trimming**, and **SMD-trimming** datasets. The average classification accuracy across these three maritime datasets is reported in Table 11. For improved visual interpretation and ease of comparison, the corresponding numerical results are further depicted in Fig. 8.

Table 11. Classification accuracy (mean \pm standard deviation) across the three maritime datasets (%).

Training Images		50	100	150	200	250
HS-NLLSC	Seaship	78.62 \pm 0.72	89.28 \pm 0.26	92.99 \pm 0.13	94.62 \pm 0.23	94.61 \pm 0.43
	Seaship-trimming	87.59 \pm 0.67	92.16 \pm 0.62	95.33 \pm 0.23	96.64 \pm 0.16	97.63 \pm 0.23
	SMD-trimming	93.22 \pm 0.36	96.41 \pm 0.41	97.73 \pm 0.29	98.97 \pm 0.15	99.28 \pm 0.08
HS-LNLSC-DINO	Seaship	92.76\pm0.21	96.38\pm0.18	97.85\pm0.31	98.12\pm0.24	98.75\pm0.11
	Seaship-trimming	95.78\pm0.48	98.04\pm0.27	98.62\pm0.22	99.23\pm0.16	99.31\pm0.11
	SMD-trimming	94.90\pm0.42	97.35\pm0.28	98.65\pm0.17	99.13\pm0.17	99.68\pm0.07

As presented in Table 11 and illustrated in Fig. 8, the proposed HS-LNLSC-DINO approach demonstrates consistently strong classification performance across all three maritime datasets. Compared to HS-NLLSC, the HS-LNLSC-DINO method achieves an accuracy improvement of approximately 1% to 15% as the number of training samples increases.

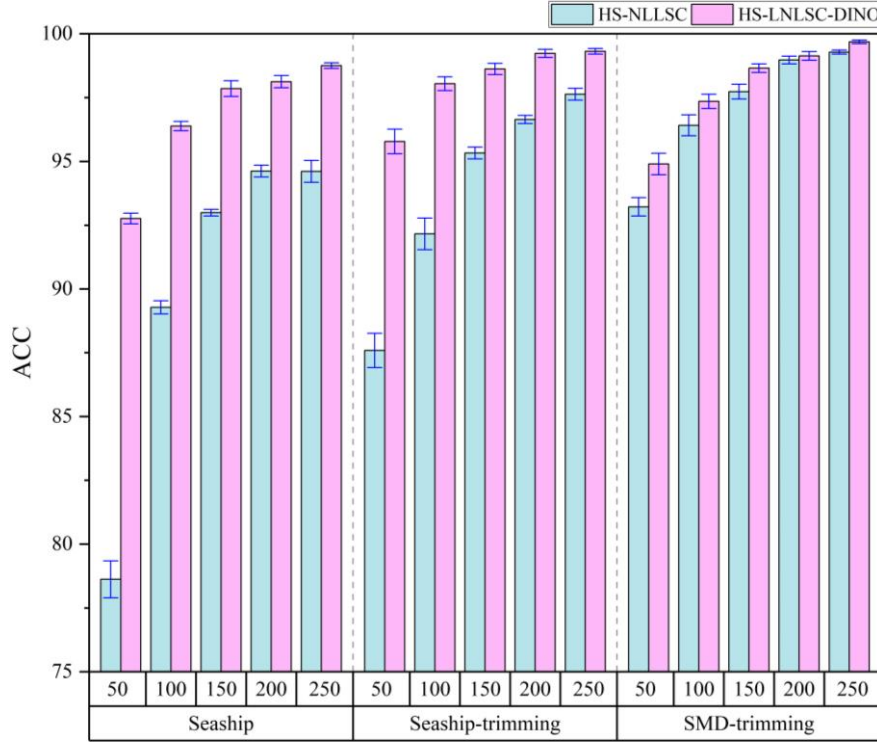


Fig. 8. Classification accuracy (mean \pm standard deviation) across the three maritime datasets.

4.4. Generalisation and stability analysis

To further evaluate the robustness and generalisation capability of the proposed HS-LNLSC-DINO method, additional experiments are conducted from two perspectives: cross-dataset generalisation performance and random repeatability analysis.

4.4.1. Generalisation evaluation across datasets

Based on the results of the ablation studies in Section 4.2, the optimal configuration, specifically, the fusion of HS-LNLSC with DINO CLS features extracted from the L12 Transformer layer, is selected for all subsequent experiments. This configuration is consistently applied across all seven datasets, including four benchmark datasets and three maritime image datasets. The classification accuracies achieved on each dataset are presented in Table 12.

As presented in Table 12, the proposed HS-LNLSC-DINO framework consistently achieves strong performance across all datasets. Particularly notable advantages are observed on the maritime datasets, which are characterised by complex backgrounds, illumination variations, and significant scale changes, highlighting the robustness of the proposed framework in challenging visual environments. At the same time, competitive performance is maintained on standard datasets, indicating that the proposed hybrid representation is not limited to a specific domain. These results validate the effectiveness and broad applicability of integrating structured sparse representations with self-supervised semantic features.

To provide a more intuitive illustration of the results, Fig. 9 presents the classification accuracy of three methods, HS-LNLSC, DINO-CLS, and HS-LNLSC-DINO, across seven datasets. As shown in Fig. 9, the proposed HS-LNLSC-DINO method outperforms the other

two methods on all seven datasets, further demonstrating the robustness and generalisation capability of the proposed approach.

Table 12. Classification accuracy on the seven image datasets (%).

Datasets	Training Images			
	15	30	45	60
Corel-10	98.85	99.60	99.82	99.85
Scene-15	94.07	95.89	96.95	97.36
Caltech-101	97.63	99.21	/	/
Caltech-256	87.24	92.07	94.86	96.86
Seaship	77.50	87.82	91.87	93.04
Seaship-trimming	90.38	94.88	96.04	96.73
SMD-trimming	88.92	93.51	95.64	95.53

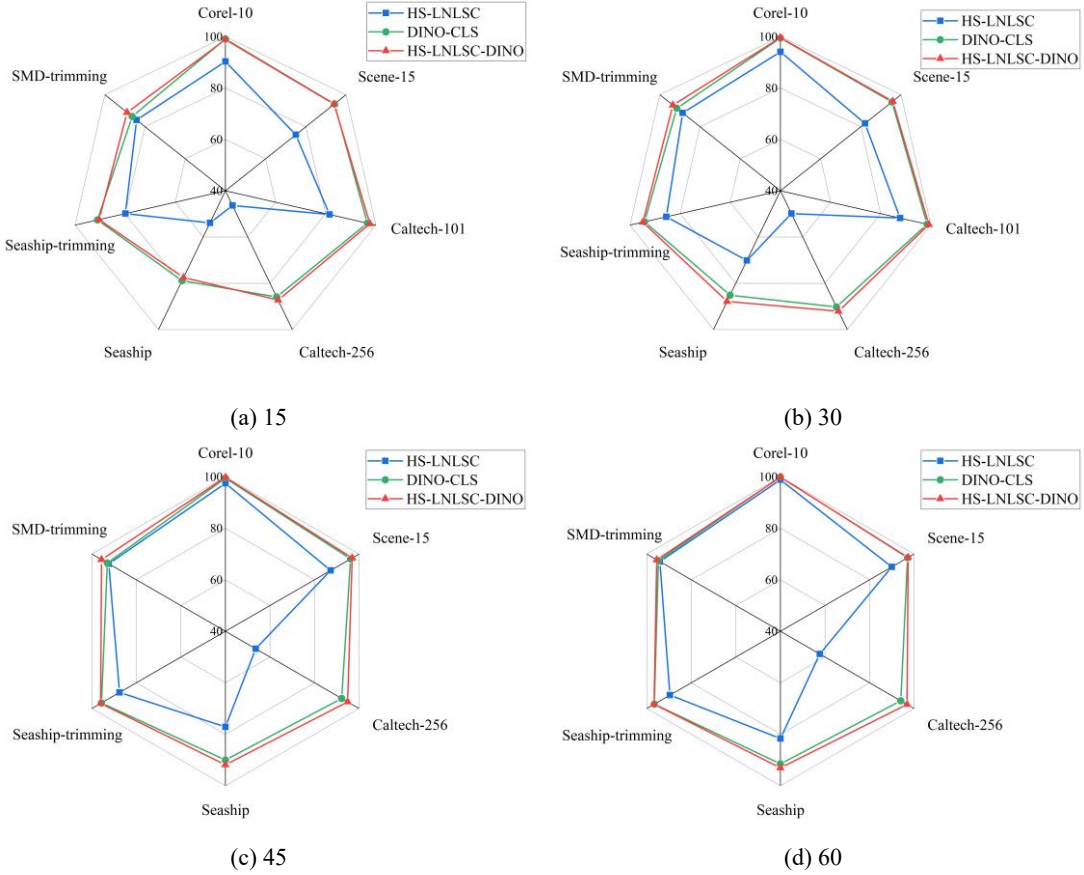


Fig. 9. Classification accuracy of the three methods on the seven image datasets.

4.4.2. Random repeatability and stability analysis

Beyond generalisation assessment, the stability of the proposed method is further evaluated through random repeatability experiments. Specifically, for the seven image datasets, 30 images from each category are randomly selected as training samples, with the remaining images used for testing. The average classification accuracy and the corresponding standard deviation are then calculated to assess the robustness of the method against data sampling variability. The experimental results from five random runs are shown in Table 13 and depicted in Fig. 10. Additionally, Fig. 11 displays the classification accuracy, along with the mean and standard deviation across the seven image datasets.

As shown in Table 13 and Figures 10 and 11, the proposed HS-LNLSC-DINO framework achieves not only high average classification accuracy but also relatively low standard deviation across repeated runs. This result demonstrates strong stability and limited sensitivity to random fluctuations in training sample selection. Taken together with the generalisation results, these findings confirm that HS-LNLSC-DINO exhibits robust and reliable performance. Through the effective integration of local geometric modelling and global semantic representation, the framework delivers consistent results across diverse datasets and experimental conditions, further supporting its suitability for practical image classification applications, particularly in marine and maritime environments.

While the Seaship dataset exhibits a slightly higher variance in classification accuracy (SD = 1.39%), this is largely attributable to the inherent characteristics of the dataset. Specifically, the images are relatively large, whereas the target objects occupy only a small portion of each scene, and the majority of pixels correspond to background or non-target regions. Consequently, forced resizing and feature extraction can lead to partial loss of target-specific features, which increases variability across experimental runs.

These observations are discussed here to highlight dataset-specific sensitivity, while simultaneously confirming the overall robustness and stability of the proposed HS-LNLSC-DINO framework across diverse maritime datasets.

Table 13. Stability analysis over five random runs on the seven image datasets (%).

Datasets	Rounds					Mean Accuracy	Standard Deviation
	1	2	3	4	5		
Corel-10	99.71	99.57	99.86	99.29	99.57	99.60	0.21
Scene-15	96.08	96.15	95.61	96.06	95.57	95.89	0.28
Caltech-101	99.01	99.15	99.17	99.36	99.39	99.21	0.16
Caltech-256	92.13	92.01	91.98	92.31	91.91	92.07	0.16
Seaship	89.45	86.51	86.57	89.12	87.49	87.82	1.39
Seaship-trimming	95.07	94.23	94.60	95.61	94.91	94.88	0.52
SMD-trimming	94.39	92.66	93.77	93.52	93.21	93.51	0.64

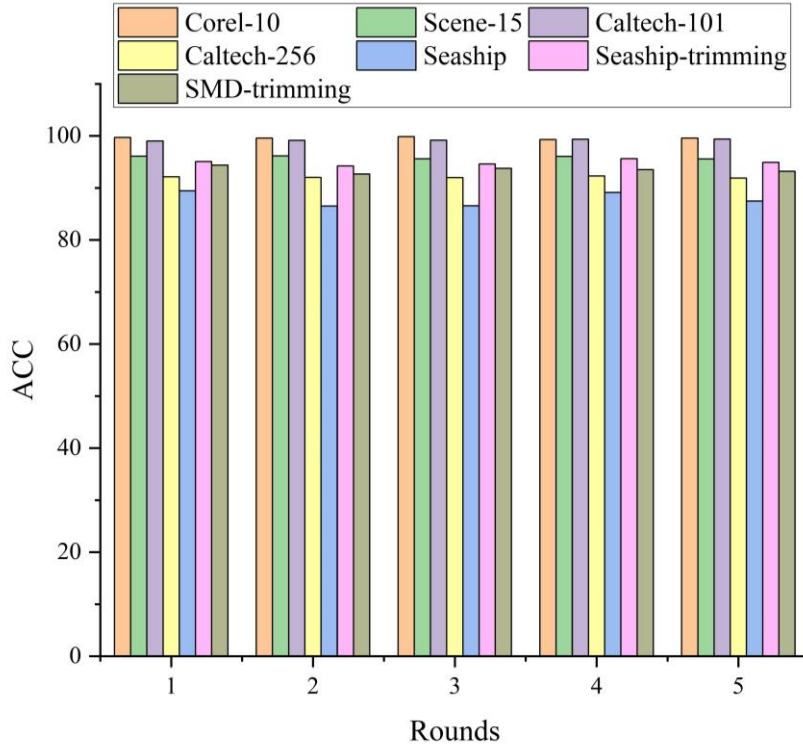


Fig. 10. Stability analysis over five random runs on the seven maritime image datasets.

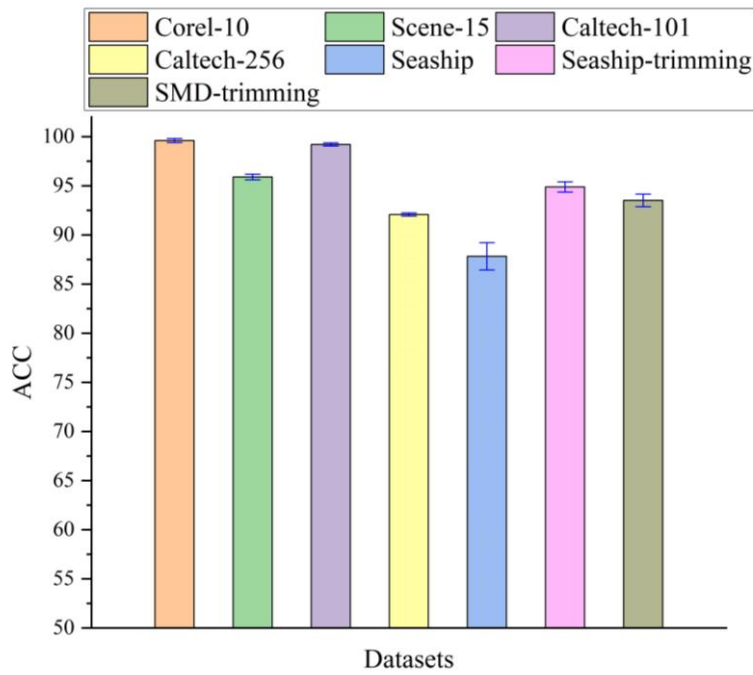


Fig. 11. Classification accuracy (mean \pm standard deviation) on the seven image datasets.

4.5. Algorithm Complexity Analysis

Let n denote the number of local features per image, m the number of template features, and k the dictionary size. The computational complexity of the coding stage in LSC can be expressed as $O(n \times (m + k))$. After incorporating histogram intersection and locality constraints,

1 the complexity of the HI-LNLSC coding stage becomes $O(k^2)$. In the MP stage, given l
2 pyramid levels and b histogram bins, the computational complexity is $O(n+l \times b)$. Accordingly,
3 the total computational complexity of the HI-LNLSC stage can be summarised as
4 $O(n+(m \times k)+k^2+l \times b)$. For the semantic information stage, let nR denote the number of
5 cross-validation folds and C the number of categories in the dataset for SVM classification. The
6 computational complexity of this stage can be expressed as $O(nR \times C)$. Accordingly, the
7 overall computational complexity of the HS-LNLSC algorithm is given by
8 $O(n+(m \times k)+k^2+l \times b+nR \times C)$.

9 For the DINO feature extraction stage, let p denote the number of image patches, d the token
10 dimension, and t the number of Transformer layers. Since the pretrained DINO model is used
11 in a frozen manner and only forward propagation is required, the computational complexity of
12 this stage is given by $O(tp^2d)$.

13 The feature fusion stage consists of a simple weighted concatenation of HS-LNLSC and
14 DINO features, whose computational cost is negligible compared with the feature extraction
15 stages. So the complexity is $O(d_1+d_2)$. Finally, given H training samples, the linear SVM
16 classification stage has a complexity of $O(H(d_1+d_2))$, where d_1 and d_2 denote the feature
17 dimensions of HS-LNLSC and DINO representations, respectively. In summary, the overall
18 computational complexity of the proposed HS-LNLSC-DINO framework can be expressed as
19 $O(n+(m \times k)+k^2+l \times b+nR \times C+tp^2d+(d_1+d_2)+H(d_1+d_2))$.

20 Table 14 presents the complexity analysis of three different methods, including HS-LNLSC,
21 DINO, and HS-LNLSC-DINO. As shown in Table 14, the proposed HS-LNLSC-DINO
22 framework integrates both local sparse representations and global semantic features, and its
23 overall complexity reflects the combination of these components. So the complexity is higher.

24 Although the hybrid framework incurs higher computational complexity compared with
25 individual components, the additional cost is mainly concentrated in the feature extraction stage.
26 The use of a frozen DINO backbone and a linear classifier ensures that the classification stage
27 remains computationally efficient. Therefore, the overall computational burden remains
28 manageable while enabling improved representation capability and classification performance
29 in complex marine and coastal scenarios.

30 Table 14. The complexity analysis of three different methods.

Methods	Complexity
HS-LNLSC	$O(n+(m \times k)+k^2+l \times b+nR \times C)$
DINO	$O(tp^2d)$
HS-LNLSC-DINO	$O(n+(m \times k)+k^2+l \times b+nR \times C+tp^2d+(d_1+d_2)+H(d_1+d_2))$

5. Management and policy implications, and practical considerations

The findings of this study extend beyond methodological performance and offer important insights for the design, implementation, and governance of marine and coastal monitoring systems. By improving the reliability, interpretability, and scalability of image-based classification under visually complex and data-limited conditions, the proposed hybrid sparse-semantic framework provides practical support for evidence-based ocean and coastal management across diverse institutional and operational contexts.

5.1. Supporting evidence-based monitoring under data and capacity constraints

HS-LNLSC-DINO exhibits robust performance in low-sample and visually complex maritime scenarios, delivering stable and interpretable data outputs that enhance the reliability and consistency of routine marine monitoring operations. From both methodological and management perspectives, HS-LNLSC-DINO demonstrates robust performance in low-sample and visually complex maritime scenarios, maintaining stable and interpretable local feature representations. These enhancements not only strengthen the reliability of upstream image-based analyses but also improve the quality of data inputs for routine monitoring systems, particularly under data-constrained conditions where labelled samples are limited, thereby supporting more consistent and trustworthy operational monitoring. A persistent challenge in marine and coastal management is the uneven availability of labelled data across regions, activities, and monitoring objectives, often compounded by limited financial and technical resources. Many management authorities must prioritise surveillance and environmental monitoring under strict budgetary and operational constraints. The demonstrated effectiveness of the proposed framework when trained on small datasets suggests that automated image classification can be more feasibly integrated into routine monitoring programmes without requiring extensive and continuous annotation efforts. This capability enables management agencies to extract actionable information from existing image data streams, including satellite and UAV imagery, even in data-poor or newly monitored areas.

From a policy perspective, reduced dependence on labelled data is particularly relevant for adaptive management frameworks, where monitoring priorities may evolve in response to emerging environmental risks, regulatory changes, or shifts in maritime activity patterns. Faster deployment and updating of monitoring systems can support more timely regulatory responses and flexible implementation of monitoring mandates at local, national, and regional scales.

5.2. Enhancing robustness for regulatory and governance applications

The hybrid integration of structured local descriptors with global DINO features enables consistent, reliable identification of vessel types and environmental features, thereby supporting maritime traffic supervision and compliance-oriented monitoring workflows. In regulatory and operational contexts, the hybrid integration of structured local sparse descriptors and global DINO features provides robust, accurate, and consistent classification performance, enabling reliable identification of maritime targets, including vessel types and environmental features, and thereby supporting surveillance and compliance-related applications. Marine and coastal imagery increasingly underpins regulatory functions related to maritime traffic management,

1 port operations, environmental protection, and safety enforcement. However, classification
2 errors caused by reflections, haze, cluttered backgrounds, and small targets can undermine
3 regulatory confidence in automated monitoring outputs. The improved robustness observed
4 across challenging maritime datasets indicates that combining structured local encoding with
5 global self-supervised semantic representations can mitigate these sources of error, supporting
6 more reliable surveillance outputs for governance-relevant tasks such as vessel-type monitoring,
7 nearshore activity assessment, and long-term evaluation of human pressures on coastal and
8 marine environments.

9 More consistent and reliable classification results can contribute to fairer and more
10 transparent enforcement practices by reducing both false positives that impose unnecessary
11 regulatory burdens and false negatives that weaken compliance objectives. In this way, the
12 proposed framework aligns technical performance improvements with broader governance
13 goals.

14 **5.3. Strengthening transparency, accountability, and institutional trust**

15 **Beyond classification accuracy, the interpretability of HS-LNLSC-DINO outputs provides**
16 **traceable and transparent feature representations, thereby enhancing auditability and fostering**
17 **institutional trust in marine and coastal governance. These practical validations underscore the**
18 **added value of the hybrid sparse-semantic framework in operationally relevant management**
19 **contexts.** Transparency and accountability are central principles of contemporary ocean and
20 coastal governance, particularly where automated systems inform regulatory actions, risk
21 assessments, or public reporting. Compared with purely deep learning-based approaches, the
22 structured sparse encoding component of the proposed framework provides interpretable
23 representations of local feature contributions, improving the traceability of classification
24 outcomes. Such interpretability supports internal audits, post-event evaluations, and external
25 scrutiny when monitoring outputs are used to justify management interventions or regulatory
26 decisions.

27 From an institutional standpoint, enhanced transparency can strengthen trust in automated
28 monitoring systems among regulators, port authorities, coastal managers, and other
29 stakeholders. This is especially important in governance settings where decisions must be
30 communicated across organisational boundaries or defended in legal, administrative, or public
31 forums.

32 **5.4. Practical deployment considerations and future policy-relevant directions**

33 The proposed framework adopts a modular design that combines frozen self-supervised
34 features with a lightweight classifier, limiting training complexity and reducing the risk of
35 overfitting. This design makes the framework suitable for deployment in operational settings
36 where computational resources and technical capacity may be constrained. In practice, the
37 feature extraction and classification pipeline can be integrated as a modular component within
38 broader monitoring systems, including satellite and UAV data ingestion, image categorisation
39 workflows, and event summarisation platforms.

1 Nevertheless, several considerations remain relevant for policy- and management-oriented
2 deployment. In practical monitoring workflows, geographic variability, seasonal changes, and
3 sensor heterogeneity can affect classification performance. While HS-LNLSC-DINO
4 demonstrates robustness across multiple datasets, systematic evaluation under these domain-
5 shift conditions remains necessary. The current study does not include cross-region transfer,
6 seasonal robustness, or sensor-shift experiments, nor does it address open-set recognition or
7 uncertainty-aware prediction. In addition, the current evaluation focuses primarily on image-
8 level classification rather than downstream decision outcomes, such as risk prioritisation or
9 regulatory effectiveness. Finally, operational decision making often requires explicit
10 representation of uncertainty, particularly in risk-sensitive contexts related to safety,
11 environmental protection, and enforcement.

12 Addressing these limitations represents an important direction for future research.
13 Incorporating cross-region adaptation, open-set recognition, and uncertainty-aware
14 classification could further enhance the policy relevance of automated image analysis,
15 supporting risk-informed decision making and precautionary management approaches. Such
16 developments would strengthen the alignment of image-based monitoring tools with the
17 principles of adaptive, transparent, and accountable ocean and coastal governance.

18 Overall, the proposed framework contributes to marine and coastal management in an indirect
19 yet structured manner by reinforcing the upstream data interpretation layer within monitoring
20 systems. This enhanced data foundation supports downstream processes, including indicator
21 construction, surveillance analysis, and decision-support activities, thereby establishing a
22 coherent pathway from image analysis to governance applications. It should be noted, however,
23 that the current evaluation primarily focuses on classification performance at the image level,
24 rather than on downstream decision-making outcomes such as risk prioritisation, regulatory
25 effectiveness, or operational decisions.

26 **6. Conclusions**

27 This paper developed a hybrid sparse-semantic image classification framework (HS-LNLSC-
28 DINO) to support robust marine and coastal image interpretation in operationally realistic and
29 data-constrained settings. By integrating structured local representations derived from HS-
30 LNLSC sparse coding with global semantic features extracted from a self-supervised DINO
31 Vision Transformer, the proposed framework bridges fine-grained geometric modelling and
32 high-level semantic understanding. This unified representation is well-suited to the visual
33 complexity commonly encountered in marine and coastal environments, including variable
34 illumination, background clutter, and limited availability of labelled data.

35 From a methodological standpoint, the proposed framework enhances the stability and
36 interpretability of sparse coding by jointly preserving locality, non-negativity, and
37 neighbourhood consistency, while introducing semantic guidance to improve discriminative
38 capability. In parallel, self-supervised DINO features provide transferable global semantics
39 without reliance on extensive labelled datasets. The feature-level fusion of these
40 complementary representations, combined with a lightweight linear classifier, results in a

1 practical and computationally efficient analytical pipeline. Such characteristics are particularly
2 relevant for marine and coastal monitoring systems that must operate under resource constraints
3 and integrate with existing observational infrastructures.

4 Extensive evaluation on four public benchmark datasets and three maritime image datasets
5 demonstrates consistent and robust performance improvements compared with representative
6 sparse coding-based, semantic-aware, and deep-feature baselines. The proposed framework
7 achieves stable gains across diverse conditions, including challenging maritime scenarios, and
8 remains effective under limited training data. Ablation analyses further confirm the
9 complementary value of structured local modelling and self-supervised semantic
10 representations, highlighting the importance of combining interpretability and semantic
11 richness in complex marine image classification tasks.

12 Beyond methodological performance, the findings have direct implications for marine and
13 coastal management practice. More reliable and interpretable image classification can
14 strengthen environmental monitoring, maritime traffic observation, and coastal surveillance by
15 improving the consistency and transparency of automated analytical outputs. In management
16 and governance contexts where monitoring results inform regulatory oversight, operational
17 decision making, or stakeholder communication, the ability to balance robustness, data
18 efficiency, and interpretability is particularly important. The proposed framework therefore
19 contributes to the broader objective of enabling evidence-based ocean and coastal management
20 through improved utilisation of visual data.

21 Several directions for future research emerge from this work. **It should be emphasised that**
22 **the primary contribution of this study lies in enhancing the reliability and interpretability of**
23 **upstream image-based analyses within marine and coastal monitoring systems. While these**
24 **improvements strengthen the foundation for downstream decision-support activities, evaluation**
25 **of decision-level outcomes-such as risk prioritisation, regulatory effectiveness, or decision-**
26 **support workflows-remains beyond the scope of the current work and is identified as a key**
27 **focus for future research.** First, extending the framework towards end-to-end optimisation while
28 retaining interpretability could further enhance operational performance. Second, incorporating
29 adaptive and uncertainty-aware fusion strategies would better support **upstream information**
30 **extraction within monitoring and decision-support workflows** under domain shifts associated
31 with seasonal variation, geographic heterogeneity, and sensor diversity. **Future research will**
32 **further investigate cross-region transfer, seasonal generalisation, sensor-shift adaptation, open-**
33 **set recognition, and uncertainty-aware prediction to enhance robustness and reliability in**
34 **heterogeneous, real-world monitoring environments.** Finally, future studies should expand
35 evaluation beyond image-level classification to management-oriented applications, such as
36 fine-grained vessel categorisation, maritime scene understanding, and cross-region or cross-
37 season generalisation within real-world monitoring workflows. Advancing along these
38 directions will further strengthen the contribution of image-based analytics to resilient,
39 transparent, and adaptive marine and coastal management.

1 **Declaration of competing interest**

2 The authors declare that they have no known competing financial interests or personal
3 relationships that could have influenced the work reported in this paper.

4 **Acknowledgements**

5 This work presented in this study is financially supported by the National Natural Science
6 Foundation of China (NSFC) under Grant No. 42407114 and No. 52571407.

7 **Data availability**

8 <https://github.com/shiying820212/HS-LNLSC-DINO.git>.

9 **References**

10 Abu Bakar, N.N., Bazmohammadi, N., Vasquez, J.C., Guerrero, J.M., 2025. Seaside port
11 operation optimization and energy management system with integrated seaport microgrid and
12 cold ironing. *Next Energy* 9, 100439. <https://doi.org/10.1016/j.nxener.2025.100439>

13 Adetunji, A.S., Vasanthan, C., Glomsrud, J.A., Galeazzi, R., Rokseth, B., 2025. Safety
14 Assurance for Autonomous Ships Using Contract-Based Design and Simulation-Based Testing.
15 *IFAC-PapersOnLine*, 16th IFAC Conference on Control Applications in Marine Systems,
16 Robotics and Vehicles CAMS 2025 59, 722–727. <https://doi.org/10.1016/j.ifacol.2025.11.720>

17 Agrillo, E., Alessi, N., Angelini, P., Buffi, F., Carli, E., Casella, L., Cutini, M., Filipponi, F.,
18 Fratarcangeli, C., Massimi, M., Mercatini, A., Pezzarossa, A., Pretto, F., Sarmati, S., Tartaglione,
19 N., Attorre, F., 2025. Enhancing Natura 2000 habitat monitoring: A framework for biodiversity
20 conservation assessment. *Ecological Indicators* 181, 114403.
21 <https://doi.org/10.1016/j.ecolind.2025.114403>

22 Bakirci, M., 2025. Advanced ship detection and ocean monitoring with satellite imagery and
23 deep learning for marine science applications. *Regional Studies in Marine Science* 81, 103975.
24 <https://doi.org/10.1016/j.rsma.2024.103975>

25 Cai, D., He, X., Han, J., Huang, T.S., 2011. Graph Regularized Nonnegative Matrix
26 Factorization for Data Representation. *IEEE Transactions on Pattern Analysis and Machine*
27 *Intelligence* 33, 1548–1560. <https://doi.org/10.1109/TPAMI.2010.231>

28 Cao, C., Yi, H., Xiang, H., He, P., Hu, J., Xiao, F., Gao, X., 2024. Accelerated Sparse-Coding-
29 Inspired Feedback Neural Architecture Search for Hyperspectral Image Classification. *IEEE*
30 *Transactions on Geoscience and Remote Sensing* 62, 1–14.
31 <https://doi.org/10.1109/TGRS.2024.3363777>

32 Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A., 2021.
33 Emerging Properties in Self-Supervised Vision Transformers. Presented at the Proceedings of
34 the IEEE/CVF International Conference on Computer Vision, pp. 9650–9660.

35 Chen, H., Xie, K., Wang, H., Zhao, C., 2018. Scene image classification using locality-
36 constrained linear coding based on histogram intersection. *Multimed Tools Appl* 77, 4081–4092.
37 <https://doi.org/10.1007/s11042-017-4830-7>

38 Cheng, C., Peng, J., Cui, W., 2023. A Two-Stage Convolutional Sparse Coding Network for
39 Hyperspectral Image Classification. *IEEE Geoscience and Remote Sensing Letters* 20, 1–5.
40 <https://doi.org/10.1109/LGRS.2023.3245210>

- 1 Dong, S., Feng, J., Fang, D., 2024. A Novel Multiscale Contrastive Learning Network for
2 Fine-Grained Ocean Ship Classification. *IEEE Journal of Selected Topics in Applied Earth
3 Observations and Remote Sensing* 17, 9989–10005.
4 <https://doi.org/10.1109/JSTARS.2024.3399310>
- 5 Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T.,
6 Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N., 2021. An Image
7 is Worth 16x16 Words: Transformers for Image Recognition at Scale.
8 <https://doi.org/10.48550/arXiv.2010.11929>
- 9 Galgani, F., Lusher, A.L., Strand, J., Haarr, M.L., Vinci, M., Molina Jack, E., Kagi, R., Aliani,
10 S., Herzke, D., Nikiforov, V., Primpke, S., Schmidt, N., Fabres, J., De Witte, B., Solbakken,
11 V.S., van Bavel, B., 2024. Revisiting the strategy for marine litter monitoring within the
12 european marine strategy framework directive (MSFD). *Ocean & Coastal Management* 255,
13 107254. <https://doi.org/10.1016/j.ocecoaman.2024.107254>
- 14 Gao, S., Tsang, I.W.-H., Chia, L.-T., Zhao, P., 2010. Local features are not lonely – Laplacian
15 sparse coding for image classification, in: 2010 IEEE Computer Society Conference on
16 Computer Vision and Pattern Recognition. Presented at the 2010 IEEE Computer Society
17 Conference on Computer Vision and Pattern Recognition, pp. 3555–3561.
18 <https://doi.org/10.1109/CVPR.2010.5539943>
- 19 Ghaban, W., Ahmad, J., Siddique, A.A., Alshehri, M.S., Saghir, A., Saeed, F., Ghaleb, B.,
20 Rehman, M.U., 2025. Sustainable Environmental Monitoring: Multistage Fusion Algorithm for
21 Remotely Sensed Underwater Super-Resolution Image Enhancement and Classification. *IEEE
22 Journal of Selected Topics in Applied Earth Observations and Remote Sensing* 18, 3640–3653.
23 <https://doi.org/10.1109/JSTARS.2024.3522202>
- 24 Graves, A., 2012. Long Short-Term Memory, in: *Supervised Sequence Labelling with
25 Recurrent Neural Networks, Studies in Computational Intelligence*. Springer Berlin Heidelberg,
26 Berlin, Heidelberg, pp. 37–45. https://doi.org/10.1007/978-3-642-24797-2_4
- 27 Han, H., Liu, S., Gan, L., 2015. Non-negativity and dependence constrained sparse coding
28 for image classification. *Journal of Visual Communication and Image Representation* 26, 247–
29 254. <https://doi.org/10.1016/j.jvcir.2014.12.002>
- 30 He, J., Huang, M., Wang, Y., Zhang, Y., Yu, H., 2025. A holistic approach to resilient port
31 management: Synergizing space template and crane deployment via benders decomposition
32 algorithm. *Ocean & Coastal Management* 269, 107806.
33 <https://doi.org/10.1016/j.ocecoaman.2025.107806>
- 34 He, L., Zhou, Y., Yang, H., Su, L., Ma, J., 2026. A deep learning-based method for marine oil
35 spill detection and its application in UAV imagery. *Marine Pollution Bulletin* 222, 118889.
36 <https://doi.org/10.1016/j.marpolbul.2025.118889>
- 37 Hoyer, P.O., 2002. Non-negative sparse coding, in: *Proceedings of the 12th IEEE Workshop
38 on Neural Networks for Signal Processing*. Presented at the the 12th IEEE Workshop on Neural
39 Networks for Signal Processing, pp. 557–565. <https://doi.org/10.1109/NNSP.2002.1030067>

- 1 Kanjir, U., Greidanus, H., Oštir, K., 2018. Vessel detection and classification from spaceborne
2 optical images: A literature survey. *Remote Sensing of Environment* 207, 1–26.
3 <https://doi.org/10.1016/j.rse.2017.12.033>
- 4 Khatri, T.K., Wei, K.T., Sharif, K.Y., 2026. An optimized deep learning framework for
5 detecting and localizing mooring lines in marine images using Faster R-CNN. *Ocean*
6 *Engineering* 343, 123559. <https://doi.org/10.1016/j.oceaneng.2025.123559>
- 7 Krizhevsky, A., Sutskever, I., Hinton, G.E., 2012. ImageNet Classification with Deep
8 Convolutional Neural Networks, in: *Advances in Neural Information Processing Systems*.
9 Curran Associates, Inc.
- 10 Lazebnik, S., Schmid, C., Ponce, J., 2006. Beyond Bags of Features: Spatial Pyramid
11 Matching for Recognizing Natural Scene Categories, in: *2006 IEEE Computer Society*
12 *Conference on Computer Vision and Pattern Recognition (CVPR'06)*. Presented at the 2006
13 *IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'06)*,
14 pp. 2169–2178. <https://doi.org/10.1109/CVPR.2006.68>
- 15 Lee, D.D., Seung, H.S., 1999. Learning the parts of objects by non-negative matrix
16 factorization. *Nature* 401, 788–791. <https://doi.org/10.1038/44565>
- 17 Lin, Z., Liu, R., Su, Z., 2011. Linearized Alternating Direction Method with Adaptive Penalty
18 for Low-Rank Representation, in: *Advances in Neural Information Processing Systems*. Curran
19 Associates, Inc.
- 20 Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B., 2021. Swin Transformer:
21 Hierarchical Vision Transformer Using Shifted Windows. Presented at the Proceedings of the
22 *IEEE/CVF International Conference on Computer Vision*, pp. 10012–10022.
- 23 Lowe, D.G., 2004. Distinctive Image Features from Scale-Invariant Keypoints. *International*
24 *Journal of Computer Vision* 60, 91–110. <https://doi.org/10.1023/B:VISI.0000029664.99615.94>
- 25 Mairal, J., Ponce, J., Sapiro, G., Zisserman, A., Bach, F., 2008. Supervised Dictionary
26 Learning, in: *Advances in Neural Information Processing Systems*. Curran Associates, Inc.
- 27 Min, H., Liang, M., Luo, R., Zhu, J., 2016. Laplacian regularized locality-constrained coding
28 for image classification. *Neurocomputing* 171, 1486–1495.
29 <https://doi.org/10.1016/j.neucom.2015.07.084>
- 30 Perronnin, F., Sánchez, J., Mensink, T., 2010. Improving the Fisher Kernel for Large-Scale
31 Image Classification, in: Daniilidis, K., Maragos, P., Paragios, N. (Eds.), *Computer Vision –*
32 *ECCV 2010*. Springer, Berlin, Heidelberg, pp. 143–156. [https://doi.org/10.1007/978-3-642-](https://doi.org/10.1007/978-3-642-15561-1_11)
33 [15561-1_11](https://doi.org/10.1007/978-3-642-15561-1_11)
- 34 Puskic, P.S., Cramer, I., Church, E., Deery, E., Egger, M., Fox, N., de Haan, W.P., Lebreton,
35 L., Liconti, A., Sanchez-Vidal, A., Wolter, H., 2026. Highly engaged marine users can help
36 monitor marine plastic pollution in under accessed environments. *Ocean & Coastal*
37 *Management* 272, 107983. <https://doi.org/10.1016/j.ocecoaman.2025.107983>
- 38 Rangel-Buitrago, N., Giarrizzo, T., Brabo, L., Silva Filho, F.J., Cooper, J.A.G., Neal, W.J.,
39 2026. Updating coastal beach classification: A cluster-based typology for contemporary human

1 use and management. *Ocean & Coastal Management* 273, 108064.
2 <https://doi.org/10.1016/j.ocecoaman.2025.108064>

3 Rasiwasia, N., Vasconcelos, N., 2008a. Scene classification with low-dimensional semantic
4 spaces and weak supervision, in: 2008 IEEE Conference on Computer Vision and Pattern
5 Recognition. Presented at the 2008 IEEE Conference on Computer Vision and Pattern
6 Recognition, pp. 1–6. <https://doi.org/10.1109/CVPR.2008.4587372>

7 Rasiwasia, N., Vasconcelos, N., 2008b. Image retrieval using query by contextual example,
8 in: Proceedings of the 1st ACM International Conference on Multimedia Information Retrieval,
9 MIR '08. Association for Computing Machinery, New York, NY, USA, pp. 164–171.
10 <https://doi.org/10.1145/1460096.1460124>

11 Ruan, H., Xu, Z., Yang, Z., Lu, Y., Qin, J., Chen, T., 2025. Learning Semantic-aware
12 Representation in Visual-Language Models for Multi-label Recognition with Partial Labels.
13 *ACM Trans. Multimedia Comput. Commun. Appl.* 21, 79:1-79:19.
14 <https://doi.org/10.1145/3708991>

15 Shen, F., Zeng, G., 2019. Semantic image segmentation via guidance of image classification.
16 *Neurocomputing* 330, 259–266. <https://doi.org/10.1016/j.neucom.2018.11.027>

17 Shi, J., Chen, J., Wan, Z., Zhou, S., Jun, Y., Shu, Y., 2025. The impact of low-sulfur marine
18 fuel policy on air pollution in global coastal cities. *Sustainable Horizons* 14, 100130.
19 <https://doi.org/10.1016/j.horiz.2024.100130>

20 Shi, W., Zheng, W., Xu, Z., 2025. Ship-Yolo: A Deep Learning Approach for Ship Detection
21 in Remote Sensing Images. *Journal of Marine Science and Engineering* 13.
22 <https://doi.org/10.3390/jmse13040737>

23 Shi, Y., Wan, Y., Wang, X., Li, H., 2025. Incorporation of Histogram Intersection and
24 Semantic Information into Non-Negative Local Laplacian Sparse Coding for Image
25 Classification. *Mathematics* 13, 219. <https://doi.org/10.3390/math13020219>

26 Shu, Z., 2026. Hybrid quantum sparse coding and dynamic convolution capsule network for
27 enhanced image classification. *Pattern Recognition* 169, 111974.
28 <https://doi.org/10.1016/j.patcog.2025.111974>

29 Simonyan, K., Zisserman, A., 2015. Very Deep Convolutional Networks for Large-Scale
30 Image Recognition. <https://doi.org/10.48550/arXiv.1409.1556>

31 Smirnov, E.A., Timoshenko, D.M., Andrianov, S.N., 2014. Comparison of Regularization
32 Methods for ImageNet Classification with Deep Convolutional Neural Networks. *AASRI*
33 *Procedia*, 2nd AASRI Conference on Computational Intelligence and Bioinformatics 6, 89–94.
34 <https://doi.org/10.1016/j.aasri.2014.05.013>

35 Teixeira, E.H., Mafra, S.B., De Figueiredo, F.A.P., 2025. InaTechShips: A validation study of
36 a novel ship dataset through deep learning-based classification and detection models for
37 maritime applications. *Ocean Engineering* 326, 120823.
38 <https://doi.org/10.1016/j.oceaneng.2025.120823>

39 Touvron, H., Cord, M., Douze, M., Massa, F., Sablayrolles, A., Jegou, H., 2021. Training
40 data-efficient image transformers & distillation through attention, in: Proceedings of the 38th

1 International Conference on Machine Learning. Presented at the International Conference on
2 Machine Learning, PMLR, pp. 10347–10357.

3 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł. ukasz,
4 Polosukhin, I., 2017. Attention is All you Need, in: Advances in Neural Information Processing
5 Systems. Curran Associates, Inc.

6 Vedaldi, A., Zisserman, A., 2012. Efficient Additive Kernels via Explicit Feature Maps. IEEE
7 Transactions on Pattern Analysis and Machine Intelligence 34, 480–492.
8 <https://doi.org/10.1109/TPAMI.2011.153>

9 Vehmaa, A., Lanari, M., Jutila, H., Mussaari, M., Pätsch, R., Telenius, A., Banta, G., Eklöf,
10 J., Jensen, K., Krause-Jensen, D., Quintana, C.O., von Numers, M., Boström, C., 2024.
11 Harmonization of Nordic coastal marsh habitat classification benefits conservation and
12 management. Ocean & Coastal Management 252, 107104.
13 <https://doi.org/10.1016/j.ocecoaman.2024.107104>

14 Wang, H., Han, B., Shu, C., Ouyang, Z., 2025. An integrated assessment of urban habitat
15 quality based on the InVEST–IUEMS model. Ecological Frontiers.
16 <https://doi.org/10.1016/j.ecofro.2025.10.011>

17 Wang, J., Yang, J., Yu, K., Lv, F., Huang, T., Gong, Y., 2010. Locality-constrained Linear
18 Coding for image classification, in: 2010 IEEE Computer Society Conference on Computer
19 Vision and Pattern Recognition. Presented at the 2010 IEEE Computer Society Conference on
20 Computer Vision and Pattern Recognition, pp. 3360–3367.
21 <https://doi.org/10.1109/CVPR.2010.5540018>

22 Wang, Q., Zhang, H., Xi, S., 2024. China’s law and policy framework for maritime safety
23 regulation of alternative fuel ships in the decarbonization transition. Marine Policy 163, 106142.
24 <https://doi.org/10.1016/j.marpol.2024.106142>

25 Wu, J., Rehg, J.M., 2009. Beyond the Euclidean distance: Creating effective visual
26 codebooks using the Histogram Intersection Kernel, in: 2009 IEEE 12th International
27 Conference on Computer Vision. Presented at the 2009 IEEE 12th International Conference on
28 Computer Vision, pp. 630–637. <https://doi.org/10.1109/ICCV.2009.5459178>

29 Xiao, X., Wang, P., Ge, Y., Luo, J., Chen, H., He, Y., Zhang, D., Li, Y., Fang, C., Lin, H.,
30 2025. GeoKG-HSA: A framework for habitat suitability assessment with geospatial knowledge
31 graphs. International Journal of Applied Earth Observation and Geoinformation 144, 104921.
32 <https://doi.org/10.1016/j.jag.2025.104921>

33 Xie, R.-D., He, Z.-F., Li, B., Liu, B., Hu, J.-Y., 2025. Semantic-Aware Representation
34 Learning via Conditional Transport for Multi-Label Image Classification.
35 <https://doi.org/10.48550/arXiv.2507.14918>

36 Xu, L., Li, X., Yan, R., Chen, J., 2025a. How to support shore-to-ship electricity constructions:
37 Tradeoff between government subsidy and port competition. Transportation Research Part E:
38 Logistics and Transportation Review 201, 104258. <https://doi.org/10.1016/j.tre.2025.104258>

1 Xu, L., Shen, C., Chen, J., 2025b. The impact of the Maritime Silk Road Initiative on the
2 carbon intensity of the participating countries. *Maritime Economics & Logistics* 27, 331–349.
3 <https://doi.org/10.1057/s41278-024-00295-z>

4 Xu, L., Wu, J., Yan, R., Chen, J., Fu, S., 2025c. Who predicts better? A comparison of
5 machine learning and econometrics in forecasting CO2 emissions from global shipping. *Energy*
6 338, 138967. <https://doi.org/10.1016/j.energy.2025.138967>

7 Yang, J., Rao, Y., Fan, H., Dong, J., Yu, H., 2025. Learning Semantic-Aware Point-Line
8 Features for Localization and Reconstruction. *IEEE Transactions on Circuits and Systems for*
9 *Video Technology* 35, 11783–11796. <https://doi.org/10.1109/TCSVT.2025.3578998>

10 Yang, J., Yu, K., Gong, Y., Huang, T., 2009. Linear spatial pyramid matching using sparse
11 coding for image classification, in: 2009 IEEE Conference on Computer Vision and Pattern
12 Recognition. Presented at the 2009 IEEE Conference on Computer Vision and Pattern
13 Recognition, pp. 1794–1801. <https://doi.org/10.1109/CVPR.2009.5206757>

14 Yang, K., Yang, T., Yao, Y., Fan, S., 2021. A transfer learning-based convolutional neural
15 network and its novel application in ship spare-parts classification. *Ocean & Coastal*
16 *Management* 215, 105971. <https://doi.org/10.1016/j.ocecoaman.2021.105971>

17 Yang, X., Zeng, Z., Yang, D., 2024. Adaptive Mid-Level Feature Attention Learning for Fine-
18 Grained Ship Classification in Optical Remote Sensing Images. *IEEE Transactions on*
19 *Geoscience and Remote Sensing* 62, 1–10. <https://doi.org/10.1109/TGRS.2024.3351874>

20 Yorulmaz, M., Susoy, M., 2025. Risk analysis and management for STS operations in ports
21 using an integrated hybrid method. *Ocean Engineering* 316, 120019.
22 <https://doi.org/10.1016/j.oceaneng.2024.120019>

23 Yuan, W. a. N., Jinghui, Z., Zhiping, C., Xiaojing, M., 2019. Non-negative local sparse coding
24 algorithm based on elastic net and histogram intersection. *Journal of Computer Applications* 39,
25 706. <https://doi.org/10.11772/j.issn.1001-9081.2018071483>

26 Zhang, C., Liu, J., Tian, Q., Liang, C., Huang, Q., 2013. Beyond visual features: A weak
27 semantic image representation using exemplar classifiers for classification. *Neurocomputing,*
28 *Image Feature Detection and Description* 120, 318–324.
29 <https://doi.org/10.1016/j.neucom.2012.07.056>

30 Zhang, C., Zhu, X., Li, L., Zhang, Y., Liu, J., Huang, Q., Tian, Q., 2015. Joint image
31 representation and classification in random semantic spaces. *Neurocomputing* 156, 79–85.
32 <https://doi.org/10.1016/j.neucom.2014.12.083>

33 Zhang, L., Chen, Z., Zheng, M., He, X., 2011. Robust non-negative matrix factorization.
34 *Front. Electr. Electron. Eng. China* 6, 192–200. <https://doi.org/10.1007/s11460-011-0128-0>

35 Zhang, W., Chen, W., Wu, X., Chen, J., Zhang, Z., 2025. Evaluation and spatiotemporal
36 analysis of low-carbon efficiency of the “ship-port” system. *Maritime Policy & Management*
37 1–33. <https://doi.org/10.1080/03088839.2025.2514765>

39 **Appendix A**

40 **Algorithm 2** (Apply an iterative process to Λ in order to compute D)

41 **Input:** Non-negative feature matrix X ; Sparse coding S ; Precision δ .

Output: Diagonal matrix Λ ; Dictionary D .

Step 1. Initialise: Set Λ^0 , $i = 0$, and δ as the initial precision.

Step 2. Convergence Check: While convergence criteria are not met, proceed with the following steps:

Step 3. Set $h_0 = \nabla f(\Lambda^0)$, if $\|h_0\|_2 \leq \delta$ then

Step 4. Λ^0 denotes the desired extremum.

Step 5. else

Step 6. $d^0 = -h_0$.

Step 7. end if

Step 8. Let $h_{i+1} = \nabla f(\Lambda^{i+1})$, if $\|h_{i+1}\|_2 \leq \delta$ or $\|\Lambda^{i+1} - \Lambda^i\|_F \leq \delta \|\Lambda^1 - \Lambda^0\|_F$ then

Step 9. Λ^{i+1} denotes the desired extremum.

Step 10. else

Step 11. $d^{i+1} = -h_{i+1} + \beta_i d^i$, $\beta_i = \|h_{i+1}\|_2^2 / \|h_i\|_2^2$.

Step 12. end if

Step 13. Identify the most suitable step size ξ_i through a rough search method along with one dimension, i.e., $f(\Lambda^i + \xi_i d^i) = \min_{\xi} f(\Lambda^i + \xi d^i)$.

Step 14. $\Lambda^{i+1} = \Lambda^i + \xi_i d^i$.

Step 15. Let $i = i + 1$, go back to Step 8.

Step 16. end while

Step 17. Go back to Λ , and acquire dictionary D using Eq. (6).

1 Appendix B

Algorithm 3 (The most suitable approximation D^* for appropriate sparseness of D)

Input: A randomly sampled column vector d from matrix D .

Output: The closest non-negative vector d^* within D^* .

Step 1. Calculate the sparseness S_{d^*} for the column vector d^* using Eq.

$sparseness(d^*) = \frac{\sqrt{M} - \|d^*\|_1 / \|d^*\|_2}{\sqrt{M} - 1}$ and $l_1 = \|d^*\|_1 / \|d^*\|_2 = \sqrt{M} - (\sqrt{M} - 1)S_{d^*}$. Here, M denotes the

dimension of d^* .

Step 2. Project the vector d into the space defined by the l_1 constraint, $d_k^* = d_k + \frac{l_1 - \|d\|_1}{M}$, $\forall k$, i.e., let

$$\|d^*\|_1 = l_1.$$

Step 3. Initialise a set $W = \{ \}$ containing the indices of negative elements within d .

Step 4. While the vector d^* does not meet the non-negativity constraint and the required sparseness condition, namely, $\|d^*\|_2^2 = \|d\|_2^2$, continue the iteration process:

Step 5. Compute the midpoint $m_k = \begin{cases} \frac{l_1}{M - \text{length}(W)}, & i \notin W \\ 0, & k \in W \end{cases}$ within the non-negative space defined by

the l_1 constraint.

Step 6. Solve the quadratic Eq. $\|d^*\|_2^2 = \|m + \alpha(d^* - m)\|_2^2$ to get the non-negative solution α , and update d^* according to Eq. $d^* = m + \alpha(d^* - m)$.

Step 7. If all elements of d^* are non-negative proceed to Step 8. Otherwise, go to Step 9.

Step 8. Go back to d^* .

Step 9. If any elements of d^* remain negative:

Step 10. For each negative element, set it to zero via $W = W \cup \{k : d_k^* < 0\}$, and set $d_k^* = 0, \forall k \in W$.

1 **Step 11.** Recalculate the projection of d^* , ensuring d^* invariant within the non-negative space defined
2 by the l_1 constraint, i.e., $d_k^* = d_k^* + (l_1 - \|d^*\|_1) / (M - \text{length}(W))$.

3 **Step 12.** Repeat the process by returning to Step 5.

4 **Step 13.** Continue iterating until the vector d^* satisfies both the non-negativity and sparseness
5 conditions.
6

7 1
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65