Check for updates

OPEN ACCESS

EDITED BY Shude He, Guangzhou University, China

REVIEWED BY Mario Versaci, Mediterranea University of Reggio Calabria, Italy Federica Nenna, University of Padua, Italy

*CORRESPONDENCE Lin Wenbo ⊠ 15101207316@163.com

RECEIVED 13 December 2024 ACCEPTED 17 February 2025 PUBLISHED 01 April 2025

CITATION

Wenbo L, Tingting L and Xiao L (2025) ModuCLIP: multi-scale CLIP framework for predicting foundation pit deformation in multi-modal robotic systems. *Front. Neurorobot.* 19:1544694. doi: 10.3389/fnbot.2025.1544694

COPYRIGHT

© 2025 Wenbo, Tingting and Xiao. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.

ModuCLIP: multi-scale CLIP framework for predicting foundation pit deformation in multi-modal robotic systems

Lin Wenbo^{1*}, Li Tingting² and Li Xiao³

¹School of Geology, Gansu Industrial Vocational and Technical College, Tianshui, Gansu, China, ²School of Electronic Information, Gansu Industrial Vocational and Technical College, Tianshui, Gansu, China, ³Guangdong Nonferrous Industry Building Quality Inspection Co., Ltd., Guangzhou, Guangdong, China

Introduction: Foundation pit deformation prediction is a critical aspect of underground engineering safety assessment, influencing construction quality and personnel safety. However, due to complex geological conditions and numerous environmental interference factors, traditional prediction methods struggle to achieve precise modeling. Conventional approaches, including numerical simulations, empirical formulas, and machine learning models, suffer from limitations such as high computational costs, poor generalization, or excessive dependence on specific data distributions. Recently, deep learning models, particularly cross-modal architectures, have demonstrated great potential in engineering applications. However, effectively integrating multi-modal data for improved prediction accuracy remains a significant challenge.

Methods: This study proposes a Multi-Scale Contrastive Language-Image Pretraining (CLP) framework, ModuCLIP, designed for foundation pit deformation prediction in multi-modal robotic systems. The framework leverages a self-supervised contrastive learning mechanism to integrate multi-source information, including images, textual descriptions, and sensor data, while employing a multi-scale feature learning approach to enhance adaptability to complex conditions. Experiments conducted on multiple foundation pit engineering datasets demonstrate that ModuCLIP outperforms existing methods in terms of prediction accuracy, generalization, and robustness.

Results and discussion: The findings suggest that this framework provides an efficient and precise solution for foundation pit deformation prediction while offering new insights into multi-modal robotic perception and engineering monitoring applications.

KEYWORDS

foundation pit deformation prediction, multi-modal robotics, contrastive learning, multi-scale features, deep learning

1 Introduction

Foundation pit deformation prediction is a pivotal aspect of ensuring both the safety and efficiency of construction projects, particularly in densely populated urban environments (Hu et al., 2023). The development of highly accurate predictive models is essential not only for mitigating catastrophic structural failures but also for optimizing construction workflows and reducing overall costs (Wei et al., 2023). However, traditional prediction methods often struggle to accommodate the complexity of multimodal data generated by robotic systems operating in heterogeneous environments

(Zong et al., 2023). Modern robotic sensing systems collect data across multiple spatial and temporal scales, posing significant challenges in effectively integrating and leveraging multi-scale information (Xu et al., 2023a). To overcome these limitations, researchers have investigated a range of computational methodologies, spanning from symbolic AI and data-driven machine learning to the latest advancements in deep learning frameworks such as CLIP (Peng et al., 2022). The convergence of these techniques with multi-modal data fusion and multi-scale modeling presents a transformative opportunity to enhance the accuracy and robustness of foundation pit deformation prediction (Figure 1).

To address the limitations of early empirical and physics-based models, researchers initially focused on traditional approaches such as numerical simulations (Xu et al., 2022), finite element methods (FEM) (Song et al., 2023), and empirical formulas (Yao et al., 2023). These methods provide valuable theoretical insights into soil-structure interactions, offering interpretability and physical consistency. However, they are often constrained by high computational costs, difficulties in modeling complex boundary conditions, and sensitivity to parameter settings (Zhou et al., 2023a). Furthermore, these methods rely on predefined assumptions, making them less adaptable to real-world conditions where uncertainties and external factors, such as groundwater variations and excavation sequences, significantly influence deformation behavior (Zhang et al., 2023).

To overcome the limitations of physics-based models, machine learning (ML) techniques have been introduced, leveraging data-driven methodologies to enhance predictive accuracy and adaptability (Shi et al., 2022). Models such as support vector machines, decision trees, and ensemble learning methods have demonstrated their ability to capture nonlinear relationships in deformation data (Hao et al., 2022), reducing the dependency on explicit geotechnical formulations (Joseph et al., 2023). ML-based approaches enable automated feature extraction and parameter tuning, improving their applicability to diverse engineering conditions (Zhou et al., 2023b). However, these models still face challenges, including sensitivity to data quality, difficulty in handling long-term temporal dependencies, and reliance on extensive labeled datasets for training (Zhang et al., 2022). Moreover, feature engineering remains a critical yet laborintensive process, requiring domain expertise to ensure meaningful input representations.

To address the constraints of ML-based models, researchers have turned to deep learning (DL) and pre-trained models, which leverage hierarchical feature extraction and multi-scale learning to improve predictive performance (Liu et al., 2023). Convolutional neural networks (CNNs) (Lian et al., 2022) and recurrent neural networks (RNNs) (Du et al., 2022) have been applied to analyze spatial-temporal deformation patterns, while transformer-based models have further enhanced long-range dependency modeling (Lin et al., 2023). Pre-trained frameworks, such as contrastive learning models, have enabled efficient multimodal data fusion by integrating geotechnical parameters, remote sensing imagery, and sensor data. However, despite their improved accuracy, existing DL approaches often require significant computational resources, suffer from limited interpretability, and may not fully exploit multi-scale information critical for capturing deformation evolution at different spatial-temporal resolutions (Yan et al., 2022).

Building on the limitations of traditional, data-driven, and deep learning methods, ModuCLIP proposes a novel multi-scale CLIP framework tailored to predicting foundation pit deformation in multi-modal robotic systems. By addressing challenges such as multi-scale data fusion, real-time adaptability, and resource efficiency, this framework leverages the strengths of contrastive pre-training while introducing innovative hierarchical structures. ModuCLIP combines the interpretability of symbolic AI with the scalability of machine learning and the representational power of deep learning. This integration not only fills gaps left by prior methodologies but also opens new possibilities for accurate, robust, and adaptive deformation prediction in complex robotic environments. ModuCLIP introduces a hierarchical feature representation mechanism that effectively fuses multiscale information across spatial and temporal domains. It is optimized for diverse construction scenarios, ensuring robust performance while reducing computational overhead in multimodal robotic systems. ModuCLIP achieves state-of-the-art accuracy in foundation pit deformation prediction, significantly outperforming traditional and contemporary approaches.

The next section reviews related work, providing an overview of existing methodologies and their limitations in addressing foundation pit deformation prediction. Following this, the methodology section presents the proposed framework, including its theoretical underpinnings, novel model design, and adaptive strategies tailored to geotechnical challenges. The subsequent section describes the experimental setup and results, highlighting the effectiveness of the proposed approach through comprehensive evaluations and comparisons with baseline methods. Finally, the conclusion summarizes the key findings and discusses potential future research directions.

2 Related work

2.1 Vision-language models in robotics

Vision-language models such as CLIP have demonstrated remarkable potential in enabling robotic systems to comprehend and interpret visual and textual data (Fan et al., 2022). These models are trained on extensive image-text pairs, providing a robust framework for understanding multi-modal information (Chango et al., 2022). In the context of robotics, vision-language models facilitate tasks such as object recognition, scene understanding, and contextual decision-making (Yu et al., 2023). Studies have explored integrating CLIP into robotic systems for tasks like robotic grasping, autonomous navigation, and human-robot interaction (Steyaert et al., 2023). By leveraging the ability of these models to align visual and textual features, robotic systems can achieve higher adaptability in dynamic environments (Cui et al., 2025). Recent advancements focus on refining these models for domain-specific applications, such as construction and geotechnical engineering, where understanding complex, domain-specific visual data is crucial (Gao et al., 2024). These developments underscore the potential of vision-language models in enabling multi-modal capabilities for foundation pit deformation prediction.



2.2 Multi-scale feature extraction techniques

Multi-scale feature extraction has emerged as a critical approach in analyzing complex visual data (Ektefaie et al., 2022). In robotic applications, capturing information across varying scales enhances the understanding of structural details and contextual relationships (Daunhawer et al., 2023). Techniques such as pyramid networks, multi-level attention mechanisms, and scale-aware convolutional architectures have been widely adopted in tasks like object detection and semantic segmentation (Cao et al., 2024). In the domain of geotechnical engineering, multi-scale analysis is crucial for accurately modeling the deformation of foundation pits, as it involves both micro-level material properties and macro-level structural behaviors (Mizuho et al., 2024). Recent works emphasize combining multi-scale visual features with temporal data to enhance predictive accuracy (He et al., 2023). These techniques have paved the way for integrating multi-scale capabilities into CLIP-based frameworks, where visual and textual features are jointly analyzed across scales to improve deformation prediction performance in robotic systems.

2.3 Multi-modal predictive frameworks

Multi-modal frameworks integrating data from heterogeneous sources have shown significant promise in predicting complex phenomena. These frameworks combine visual, textual, and sensory data to provide a comprehensive understanding of the environment (Chai and Wang, 2022). In geotechnical and construction robotics, multi-modal approaches have been employed to model soil behavior, structural dynamics, and environmental interactions (Yang et al., 2022). Techniques like sensor fusion, multi-modal attention mechanisms, and adversarial learning have been utilized to improve the robustness of predictions (Bayoudh et al., 2021). Recent advancements focus on designing end-to-end architectures that align and process multi-modal inputs effectively, ensuring better feature representation and prediction accuracy (He et al., 2022). Integrating such frameworks with visionlanguage models like CLIP enhances their ability to reason over diverse data sources, making them particularly suitable for tasks like foundation pit deformation prediction (Zhou et al., 2024). These frameworks also facilitate real-time analysis, crucial for deploying robotic systems in dynamic construction environments.

2.4 Application of fuzzy method in deformation prediction

The ModuCLIP framework proposed in this study achieves high-precision prediction of foundation pit deformation in a multimodal robot system, but the computational load remains a key challenge in real-time applications (Versaci et al., 2024). Previous studies have proposed fuzzy similarity and divergence methods, which have obvious advantages in computational complexity and can reduce the real-time computational burden and improve the computational efficiency of the prediction model (Versaci et al., 2022). Although these methods were originally applied to two-dimensional civil engineering problems, their mathematical framework has strong versatility and can be directly applied to the foundation pit deformation prediction involved in this study. Future optimization directions include combining these fuzzy computational methods to further enhance the real-time adaptability of ModuCLIP in dynamic environments, and reduce the consumption of computational resources while maintaining high accuracy to meet the needs of large-scale engineering monitoring.

3 Method

3.1 Background

The prediction of foundation pit deformation is a critical task in geotechnical engineering, ensuring structural stability and safety during excavation processes. This subsection provides an outline of the methodological framework developed for this study, focusing on the modeling and analysis approaches employed to enhance prediction accuracy. The section begins by introducing the key challenges and requirements inherent in foundation pit deformation prediction, including the dynamic and complex nature of geotechnical systems, data sparsity, and environmental uncertainties.

The prediction of foundation pit deformation requires a formalized framework that captures the inherent complexities of soil-structure interactions, dynamic loading conditions, and environmental influences. This section provides the mathematical and physical foundations necessary for constructing such a framework, focusing on the governing principles, spatial-temporal relationships, and critical parameters defining the problem.

Let $\Omega \subset \mathbb{R}^3$ denote the spatial domain of the foundation pit, with Γ representing its boundary. The deformation field, $\mathbf{u}(\mathbf{x}, t)$, at spatial location $\mathbf{x} \in \Omega$ and time *t*, is governed by the following equilibrium equation derived from continuum mechanics:

$$\nabla \cdot \boldsymbol{\sigma} + \mathbf{f} = 0, \quad \mathbf{x} \in \Omega, \tag{1}$$

where σ is the stress tensor, and **f** represents body forces such as gravity.

The stress-strain relationship is described using Hooke's law for linear elasticity:

$$\boldsymbol{\sigma} = \mathbb{C} : \boldsymbol{\varepsilon}, \tag{2}$$

where $\boldsymbol{\varepsilon} = \frac{1}{2} (\nabla \mathbf{u} + \nabla \mathbf{u}^{\top})$ is the strain tensor, and \mathbb{C} is the fourth-order elasticity tensor determined by material properties.

Boundary conditions are imposed as:

$$\mathbf{u} = \mathbf{u}_D \quad \text{on } \Gamma_D,$$

$$\boldsymbol{\sigma} \cdot \mathbf{n} = \mathbf{t}_N \quad \text{on } \Gamma_N,$$
 (3)

where Γ_D and Γ_N are the Dirichlet and Neumann boundaries, respectively, \mathbf{u}_D specifies prescribed displacements, and \mathbf{t}_N denotes external tractions.

To incorporate time-dependent behavior, the deformation process is modeled as a quasi-static system influenced by excavation activities and environmental variations. This dynamic interaction is described by coupling the equilibrium equation with an evolution equation for soil settlement:

$$\frac{\partial \mathbf{u}}{\partial t} + \nabla \cdot (\mathbf{v}\mathbf{u}) = \mathbf{q},\tag{4}$$

where \mathbf{v} is the excavation velocity field, and \mathbf{q} represents sources or sinks of deformation due to material removal or water table changes.

The environmental influence, including water table fluctuations, is integrated through Terzaghi's effective stress principle:

$$\boldsymbol{\sigma}_{\rm eff} = \boldsymbol{\sigma} - \alpha p \mathbf{I},\tag{5}$$

where α is the Biot coefficient, *p* is the pore water pressure, and I is the identity tensor. Pore pressure evolution is governed by Darcy's law and the continuity equation:

$$\nabla \cdot \left(\frac{\mathbf{k}}{\mu} \nabla p\right) - \frac{\partial p}{\partial t} = 0, \tag{6}$$

where ${\bf k}$ is the permeability tensor, and μ is the dynamic viscosity of water.

The problem is further complicated by nonlinearity in soil behavior, which can be captured using a constitutive model such as the Mohr-Coulomb criterion:

$$f(\boldsymbol{\sigma}) = \tau - c - \sigma_n \tan \phi \le 0, \tag{7}$$

where τ is the shear stress, σ_n is the normal stress, c is cohesion, and ϕ is the angle of internal friction.

To enable numerical solutions, the equations are discretized using the finite element method (FEM). The weak form of the equilibrium equation is:

$$\int_{\Omega} \delta \boldsymbol{\varepsilon} : \boldsymbol{\sigma} \ d\Omega - \int_{\Gamma_N} \delta \mathbf{u} \cdot \mathbf{t}_N \ d\Gamma = 0, \tag{8}$$

where $\delta \mathbf{u}$ represents test functions. Time integration schemes, such as the backward Euler method, are employed for transient analyses.

3.2 Dynamic geotechnical learning Network

To accurately predict foundation pit deformation under complex geotechnical and environmental conditions, we propose a novel modeling framework, termed dynamic geotechnical learning network (DGL-Net). This model integrates physical domain knowledge with data-driven machine learning to improve prediction fidelity, scalability, and adaptability. The foundation of DGL-Net lies in its hybrid structure, comprising a physics-constrained deep neural network and a dynamic state representation module. The following subsections describe the core components of DGL-Net, including the architecture, the integration of domain knowledge, and the mechanisms for handling multi-scale spatial-temporal data (Figure 2).

The input to the DGL-Net represents a comprehensive description of the geotechnical system, incorporating various



physical variables that characterize the material behavior, loading conditions, and state of the system at each time t. These variables form the vector **X**, which consists of:

$$\mathbf{X} = \{\mathbf{u}_t, \dot{\mathbf{u}}_t, \nabla \mathbf{u}_t, \boldsymbol{\sigma}_t, \mathbf{f}_t, p_t, \mathbf{k}, \boldsymbol{\phi}, c\},\tag{9}$$

where: \mathbf{u}_t is the displacement field at time t, $\dot{\mathbf{u}}_t$ is the velocity field, or the temporal derivative of the displacement, $\nabla \mathbf{u}_t$ represents the strain tensor, computed as the gradient of the displacement field, σ_t is the stress tensor, \mathbf{f}_t corresponds to external loading applied to the system, p_t is the pore pressure, \mathbf{k} is the permeability of the material, ϕ is the friction angle, c is the cohesion of the material.

DGL-Net integrates these input variables into a hybrid architecture, merging the strength of traditional physics-based modeling with the flexibility of data-driven deep learning. This hybrid structure allows the network to respect physical principles, while also learning complex patterns from data.

DGL-Net is the physics-informed encoder, which encodes governing physical laws and principles into the model. This includes equilibrium equations, constitutive laws, and material failure criteria. For each input feature **X**, the encoder ensures that the output respects the following constraints derived from the physical laws governing the system:

$$\nabla \cdot \boldsymbol{\sigma} + \mathbf{f} = 0, \quad f(\boldsymbol{\sigma}) < 0, \tag{10}$$

where the first equation represents the equilibrium condition, and the second ensures that the stress σ satisfies failure criteria such as the Mohr-Coulomb failure condition. The function $f(\sigma)$ characterizes the material's yield surface, typically dependent on stress invariants and friction parameters.

the encoded input is passed through a series of convolutional layers to extract spatial and temporal features from the input data. The hierarchical feature extraction process begins with the first convolutional layer C_1 , which operates on the raw input vector **X**. Each successive layer C_i applies convolutions to progressively refine and extract higher-level features from the previous layer's output. The spatial dependencies across the input data are captured by dilated convolutions, which allow the model to capture multi-scale dependencies in both space and time:

$$\mathbf{F} = \mathcal{C}_n(\mathcal{C}_{n-1}(\dots,\mathcal{C}_1(\mathbf{X}))),\tag{11}$$

where **F** represents the output of the final convolutional layer, containing hierarchical spatial-temporal features.

To model the temporal evolution of the system, DGL-Net incorporates a Long Short-Term Memory (LSTM) network, which enables the model to capture long-range temporal dependencies. Given the feature vector \mathbf{F}_t at time *t*, the LSTM takes the current input along with the hidden state \mathbf{h}_{t-1} and cell state \mathbf{c}_{t-1} from

the previous time step, producing the updated hidden and cell states \mathbf{h}_t and \mathbf{c}_t :

$$\mathbf{h}_t, \mathbf{c}_t = \text{LSTM}(\mathbf{F}_t, \mathbf{h}_{t-1}, \mathbf{c}_{t-1}).$$
(12)

The DGL-Net is trained using a multi-objective loss function that ensures high accuracy, physical consistency, and generalization capabilities. The training procedure integrates several components in the loss function, each addressing different aspects of the model's performance.

The primary loss term is the prediction loss \mathcal{L}_{pred} , which ensures that the model's predictions are close to the ground truth deformation fields:

$$\mathcal{L}_{\text{pred}} = \frac{1}{N} \sum_{i=1}^{N} \|\mathbf{u}_i^{\text{pred}} - \mathbf{u}_i^{\text{true}}\|_2^2,$$
(13)

where $\mathbf{u}_{i}^{\text{pred}}$ represents the predicted displacement field, and $\mathbf{u}_{i}^{\text{true}}$ denotes the true displacement field for the *i*-th data point.

ensuring predictive accuracy, the model must also respect the underlying physical laws. To enforce physical consistency, the term \mathcal{L}_{phys} is included:

$$\mathcal{L}_{\text{phys}} = \lambda_1 \|\nabla \cdot \boldsymbol{\sigma} + \mathbf{f}\|_2^2 + \lambda_2 \max(0, f(\boldsymbol{\sigma})), \quad (14)$$

where σ is the stress tensor, **f** is the body force, and $\nabla \cdot \sigma$ represents the divergence of the stress tensor. The first term ensures that the predicted stress field satisfies the equilibrium condition. The second term penalizes deviations from material constraints or failure criteria, controlled by the function $f(\sigma)$. The hyperparameters λ_1 and λ_2 balance the importance of these terms.

To avoid overfitting and ensure generalization, a smoothness loss term \mathcal{L}_{smooth} is added to penalize sharp spatial variations in the predicted displacement fields:

$$\mathcal{L}_{\text{smooth}} = \frac{1}{N} \sum_{i=1}^{N} \|\nabla \mathbf{u}_i\|_2^2, \tag{15}$$

where $\nabla \mathbf{u}_i$ is the gradient of the displacement field at each point.

A boundary condition loss term \mathcal{L}_{BC} can also be included to ensure correct behavior at the domain boundaries:

$$\mathcal{L}_{\rm BC} = \frac{1}{M} \sum_{j=1}^{M} \|\mathbf{u}_j^{\rm pred} - \mathbf{u}_j^{\rm boundary}\|_2^2, \tag{16}$$

where *M* is the number of boundary points, and $\mathbf{u}_{j}^{\text{boundary}}$ is the prescribed displacement at the boundary.

A regularization term \mathcal{L}_{reg} can be introduced to prevent overfitting by penalizing large weights in the model:

$$\mathcal{L}_{\text{reg}} = \lambda_3 \sum_{k=1}^{K} \|\theta_k\|_2^2, \tag{17}$$

where θ_k represents the model parameters, and λ_3 is a regularization hyperparameter.

a symmetry loss term \mathcal{L}_{sym} can be added to enforce the symmetry of the stress tensor:

$$\mathcal{L}_{\text{sym}} = \frac{1}{N} \sum_{i=1}^{N} \|\boldsymbol{\sigma}_i - \boldsymbol{\sigma}_i^T\|_2^2, \qquad (18)$$

where σ_i^T is the transpose of the stress tensor at point *i*.

Combining all these terms, the total loss function for training the DGL-Net is:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{pred}} + \mathcal{L}_{\text{phys}} + \mathcal{L}_{\text{smooth}} + \mathcal{L}_{\text{BC}} + \mathcal{L}_{\text{reg}} + \mathcal{L}_{\text{sym}}.$$
 (19)

The hierarchical attention mechanism for multi-scale handling plays a crucial role in addressing the challenges posed by variations in geotechnical conditions across different spatial scales. In practice, geotechnical data often exhibit heterogeneous features that vary across scales, such as soil properties at different depths or spatial variations in soil behavior. To effectively capture these variations, the proposed method combines features extracted from multiple resolutions, leveraging a hierarchical attention mechanism to weigh the importance of each scale.

The multi-scale feature extraction process involves generating features at different resolutions. Let the set of features from L different resolutions be denoted as $\{F_1, F_2, \ldots, F_L\}$, where each F_i represents the feature map at the *i*-th resolution. These feature maps capture information at different granularities of the geotechnical space, and combining them effectively is essential for robust performance. To achieve this, the proposed mechanism utilizes a weighted sum of the features from each resolution.

$$\mathbf{F}_{\text{multi}} = \sum_{i=1}^{L} \alpha_i \mathbf{F}_i, \quad \sum_{i=1}^{L} \alpha_i = 1,$$
(20)

where $\mathbf{F}_{\text{multi}}$ is the multi-scale feature representation. The weights α_i are learnable and are determined through an attention mechanism. The attention mechanism dynamically adjusts these weights during training, allowing the model to focus on more informative resolutions for the given task.

To further refine the importance of each scale, the hierarchical attention mechanism is designed to operate at multiple levels. At each level, the model learns to assign higher weights to feature maps that contain more relevant information for the task, such as areas with high variability in geotechnical properties. This adaptive weighting ensures that the model effectively handles spatially varying features in the data.

The attention weights α_i are computed by a separate attention module that processes the features from each resolution. The attention mechanism can be modeled as a softmax function applied to the compatibility scores between the feature maps and a learned query vector. Let q_i be the query vector associated with the *i*th resolution and \mathbf{F}_i be the corresponding feature map. The compatibility score s_i between q_i and \mathbf{F}_i is computed as:

$$s_i = \mathbf{q}_i^\top \mathbf{F}_i,\tag{21}$$

where $\mathbf{q}_i \in \mathbb{R}^d$ is the learned query vector, and $\mathbf{F}_i \in \mathbb{R}^{d \times n}$ is the feature map at resolution *i*, where *d* is the dimensionality of the feature vector and *n* is the number of elements in the feature map.

The attention weight α_i for each resolution is computed by applying the softmax function to the compatibility scores across all resolutions:

$$\alpha_i = \frac{\exp(s_i)}{\sum_{j=1}^L \exp(s_j)}.$$
(22)

This softmax function ensures that the attention weights are non-negative and sum to 1, enforcing a normalized contribution from each resolution.

After calculating the attention weights, the multi-scale features are fused by taking the weighted sum of the feature maps as described earlier. However, the combination of features at different resolutions requires additional refinement to handle the misalignment between feature maps at different scales. This refinement is achieved through a convolutional layer that processes the combined features $\mathbf{F}_{\text{multi}}$ and further enhances the spatial coherence across different resolutions:

$$\mathbf{F}_{\text{refined}} = \text{Conv}(\mathbf{F}_{\text{multi}}), \tag{23}$$

where $Conv(\cdot)$ denotes a convolution operation applied to the multi-scale feature representation.

The final representation $F_{refined}$ is then passed to downstream tasks, such as classification or regression, depending on the specific geotechnical problem being addressed. By leveraging multi-scale information with hierarchical attention, the model can adapt to variations in geotechnical conditions at different spatial scales, providing more accurate predictions and insights.

To prevent overfitting and ensure that the model generalizes well across different scales, regularization techniques are introduced. One common approach is to apply an entropy regularization term to the attention weights, encouraging the model to distribute attention in a more uniform manner across resolutions:

$$\mathcal{L}_{\text{reg}} = -\sum_{i=1}^{L} \alpha_i \log(\alpha_i).$$
(24)

This entropy term penalizes highly skewed attention distributions, encouraging the model to learn more balanced attention weights. The total loss function, combining the task-specific loss and the regularization term, becomes:

$$\mathcal{L}_{\text{total}} = \mathcal{L}_{\text{task}} + \lambda \mathcal{L}_{\text{reg}},\tag{25}$$

where λ is a regularization hyperparameter controlling the strength of the entropy penalty (Figure 3).

3.3 Adaptive multi-scale integration strategy

To address the challenges in accurately predicting foundation pit deformation under complex geotechnical conditions, we propose the adaptive multi-scale integration strategy (AMIS). This strategy leverages domain-specific insights, dynamic feature aggregation, and efficient computational techniques to optimize the predictive capabilities of our model, DGL-Net. AMIS is specifically designed to handle spatial-temporal variability, dynamic excavation processes, and heterogeneous environmental factors (Figure 4).

Foundation pit deformation prediction is influenced by several factors that interact across multiple scales. The deformation process exhibits complex behaviors that span spatial and temporal domains, which can be categorized as follows. Multi-scale dependencies, where deformation patterns exhibit relationships across both localized and large-scale spatial areas, as well as varying temporal resolutions. This involves capturing fine-grained details of local soil failure while also accounting for long-term settlement or global deformation patterns. The dynamic excavation impact, where the excavation process itself is a dynamic, time-varying process that alters stress distributions and pore pressure in the soil, creating complex spatiotemporal interactions that must be accounted for over time. Heterogeneity, where soil properties such as permeability, cohesion, and friction vary spatially across the foundation pit. To accurately model these effects, the system must incorporate mechanisms that can adaptively capture such spatial heterogeneity.

To address these challenges, adaptive multi-scale inference system (AMIS) integrates a multi-scale feature extraction and aggregation framework, which is designed to capture both spatial heterogeneity and temporal evolution in the deformation process. The system uses a series of convolutional operations across multiple scales to extract meaningful features from the raw input data. At the *s*-th scale, the feature map \mathbf{F}_s is obtained by applying convolution operations, pooling layers, and other transformations to capture the deformation characteristics at that scale. These feature maps can vary in size depending on the receptive field of the operations at each scale. The final aggregated feature map, $\mathbf{F}_{\text{final}}$, is computed by summing the weighted feature maps across all scales. The weights β_s are learnable parameters optimized during training, which allow the model to adaptively prioritize certain scales based on the task. This can be expressed as:

$$\mathbf{F}_{\text{final}} = \sum_{s=1}^{S} \beta_s \mathbf{F}_s, \tag{26}$$

where *S* denotes the number of scales considered, \mathbf{F}_s represents the feature map at scale *s*, and β_s are learnable weights assigned to each scale.

The multi-scale feature extraction process involves different types of feature capture. Local features represent fine-grained deformations, often associated with localized soil failure or small-scale settlements. To extract these local features, standard convolutional layers with small receptive fields are employed. This enables the system to capture minute variations in deformation



The hierarchical attention mechanism for multi-scale handling. This mechanism captures geotechnical data variations across different spatial scales by generating features at multiple resolutions. It uses a learnable attention mechanism to assign adaptive weights to each scale, allowing the model to focus on the most relevant features for accurate predictions. The features from different resolutions are fused and refined, enhancing spatial coherence and improving the model's ability to handle heterogeneous geotechnical conditions. Regularization ensures balanced attention distribution across scales to prevent overfitting.



The adaptive multi-scale integration strategy (AMIS) framework for predicting foundation pit deformation under complex geotechnical conditions. AMIS integrates multi-scale feature aggregation, temporal adaptation, and geotechnical constraints to model the dynamic and heterogeneous behavior of foundation pit deformation. By leveraging multi-scale convolutional operations and recurrent layers, the model adapts to spatial and temporal variations in soil properties and excavation processes. The integration of real-time feedback using a Kalman filter allows the model to continuously refine its predictions, ensuring accurate forecasting of displacement and strain over time. The attention mechanism dynamically focuses on critical regions with high spatial variability, enhancing the model's adaptability to real-world applications.

patterns at a local level. Mathematically, the local feature map F_{local} can be expressed as:

> $F_{local} = Conv_{small}(X),$ (27)

where Conv_{small} represents the convolution operation with a small kernel size, and X is the input feature map or data. To capture larger-scale deformation patterns and long-term settlement effects, global features are extracted using dilated convolutions and pooling layers. These operations allow for an increased receptive field, capturing global interactions between soil layers and long-range stress-strain relations. The global feature map $\mathbf{F}_{\text{global}}$ is given by:

$$\mathbf{F}_{\text{global}} = \text{DilatedConv}(\mathbf{X}),$$
 (28)

where DilatedConv represents the dilated convolution operation, which increases the receptive field without increasing the computational complexity. Once the local and global features have been extracted, they are aggregated across all scales using weighted summation. The aggregated feature map, $\mathbf{F}_{\text{final}}$, is a weighted sum of the individual feature maps:

$$\mathbf{F}_{\text{final}} = \sum_{s=1}^{S} \beta_s \mathbf{F}_s,\tag{29}$$

This aggregation process allows the system to adaptively focus on certain scales depending on the complexity of the deformation pattern, as determined by the learned weights β_s . The learnable weights β_s for each scale are optimized during the training process to prioritize the most relevant scales for the prediction task. These weights are updated using gradient-based optimization, and they allow the model to adjust to different types of soil and excavation scenarios. The weight update rule can be expressed as:

$$\beta_s^{\text{new}} = \beta_s^{\text{old}} - \eta \frac{\partial \mathcal{L}}{\partial \beta_s},\tag{30}$$

where η is the learning rate, \mathcal{L} is the loss function, and $\frac{\partial \mathcal{L}}{\partial \beta_s}$ represents the gradient of the loss with respect to the scale weights. In dynamic excavation scenarios, the deformation patterns evolve over time. To capture this temporal evolution, recurrent layers such as LSTMs or GRUs can be integrated into the model. The temporal feature map $\mathbf{F}_{\text{temporal}}$ is computed as:

$$\mathbf{F}_{\text{temporal}} = \text{RNN}(\mathbf{F}_{\text{final}}), \tag{31}$$

where RNN represents a recurrent neural network layer, such as an LSTM or GRU, applied to the aggregated feature map to capture the temporal dependencies in the deformation process. After extracting multi-scale features and incorporating temporal dependencies, the final deformation prediction \hat{y} is obtained by passing the temporal feature map through fully connected layers. This can be expressed as:

$$\hat{y} = FC(\mathbf{F}_{temporal}),$$
 (32)

where FC denotes the fully connected layer, which maps the temporal features to the final predicted deformation.

To address time-dependent deformation and dynamic changes in the excavation process, AMIS integrates a recurrent update mechanism designed to capture the temporal dependencies across consecutive time steps. This recurrent framework ensures that the system can dynamically adjust its predictions based on evolving environmental conditions and excavation progress. The hidden state \mathbf{h}_t is updated at each time step, reflecting the history of previous states and the current feature map \mathbf{F}_t . The update rule is governed by a gated recurrent unit (GRU), which is wellsuited for modeling sequential data with long-range dependencies. The update mechanism can be described by the following recurrence relation:

$$\mathbf{h}_t = f(\mathbf{F}_t, \mathbf{h}_{t-1}),\tag{33}$$

where $f(\cdot)$ denotes the GRU function, which takes as input the current feature map \mathbf{F}_t and the previous hidden state \mathbf{h}_{t-1} . The

GRU updates the hidden state by integrating both the temporal feature information and the previous state, ensuring that the model retains memory of past events. The function f can be broken down into the following steps, where we define the gates involved in the GRU. The update gate, which determines the amount of previous hidden state to retain, is computed as follows:

$$\mathbf{z}_t = \sigma(\mathbf{W}_z[\mathbf{F}_t, \mathbf{h}_{t-1}] + \mathbf{b}_z), \tag{34}$$

where $\sigma(\cdot)$ is the sigmoid activation function, \mathbf{W}_z and \mathbf{b}_z are the weight matrix and bias for the update gate, and $[\mathbf{F}_t, \mathbf{h}_{t-1}]$ denotes the concatenation of the current feature map and the previous hidden state. The reset gate determines the amount of the previous hidden state to forget:

$$\mathbf{r}_t = \sigma(\mathbf{W}_r[\mathbf{F}_t, \mathbf{h}_{t-1}] + \mathbf{b}_r), \tag{35}$$

where \mathbf{W}_r and \mathbf{b}_r are the weight matrix and bias for the reset gate. The reset gate allows the model to decide how much of the past memory should be ignored when calculating the new candidate hidden state. The candidate hidden state, which represents a potential new memory, is computed as:

$$\tilde{\mathbf{h}}_t = \tanh(\mathbf{W}_h[\mathbf{F}_t, \mathbf{r}_t \odot \mathbf{h}_{t-1}] + \mathbf{b}_h), \tag{36}$$

where $tanh(\cdot)$ is the hyperbolic tangent activation function, and \odot denotes the element-wise multiplication between the reset gate and the previous hidden state. This operation allows the model to focus on relevant past information when forming the candidate state. The final hidden state is updated by combining the previous hidden state and the candidate hidden state according to the update gate:

$$\mathbf{h}_t = (1 - \mathbf{z}_t) \odot \mathbf{h}_{t-1} + \mathbf{z}_t \odot \tilde{\mathbf{h}}_t.$$
(37)

The update gate \mathbf{z}_t controls the contribution of the previous hidden state and the candidate state to the new hidden state, allowing the model to adaptively blend past and current information. The output prediction \mathbf{y}_t at each time step is generated by applying a linear transformation to the final hidden state:

$$\mathbf{y}_t = \mathbf{W}_y \mathbf{h}_t + \mathbf{b}_y, \tag{38}$$

where \mathbf{W}_y and \mathbf{b}_y are the weight matrix and bias for the output layer, respectively. The output prediction \mathbf{y}_t provides the system's estimate of the excavation status or environmental conditions at time *t*. To train the system, a temporal loss function is defined that captures the difference between the predicted output \mathbf{y}_t and the true value $\mathbf{y}_t^{\text{true}}$:

$$\mathcal{L}_t = \|\mathbf{y}_t - \mathbf{y}_t^{\text{true}}\|^2, \tag{39}$$

where $\|\cdot\|$ denotes the Euclidean norm. The total loss across all time steps is then summed as:

$$\mathcal{L} = \sum_{t=1}^{l} \mathcal{L}_t, \tag{40}$$

which is minimized during training to update the parameters of the GRU. To optimize the parameters of the system, backpropagation through time (BPTT) is employed. The gradient of the loss function with respect to the parameters W_z , W_r , W_h , W_y and biases b_z , b_r , b_h , b_y is computed by unrolling the recurrence over time and applying the chain rule:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \sum_{t=1}^{T} \frac{\partial \mathcal{L}_t}{\partial \theta},\tag{41}$$

where θ represents the parameters of the system. This allows for efficient learning of the temporal dependencies and the adjustment of model weights.

The AMIS framework ensures that its predictions align with physical principles by explicitly incorporating geotechnical constraints, thereby maintaining consistency with the underlying mechanics of the system. The core of these constraints lies in the stress equilibrium and failure conditions (Figure 5).

The total physical loss function, \mathcal{L}_{phys} , is defined as follows:

$$\mathcal{L}_{\text{phys}} = \lambda_1 \|\nabla \cdot \boldsymbol{\sigma} + \mathbf{f}\|_2^2 + \lambda_2 \max(0, f(\boldsymbol{\sigma})), \tag{42}$$

where σ represents the stress tensor, **f** is the external force vector, and $f(\sigma)$ encapsulates the failure criterion (e.g., a yield surface function or a plasticity model). The term $\nabla \cdot \sigma + \mathbf{f}$ ensures that the stress field satisfies equilibrium, while the second term enforces the failure condition by penalizing the stress state if it exceeds the material's failure threshold. The regularization parameters λ_1 and λ_2 control the balance between these two constraints, allowing for more accurate predictions of both stable and failure states.

In order to handle the spatial variability inherent in geotechnical properties such as soil cohesion, permeability, and stiffness, AMIS incorporates an attention mechanism to focus on regions of interest. The attention mask **A** is dynamically generated based on local variations in material properties, allowing the model to prioritize critical areas for further analysis. The attention mechanism is formulated as:

$$\mathbf{F}_{\text{attn}} = \mathbf{A} \odot \mathbf{F},\tag{43}$$

where **F** represents the input features (e.g., stress, strain, or displacement fields), and \odot denotes element-wise multiplication. The attention mask **A** assigns higher weights to areas with significant spatial variability, thus emphasizing regions that have a more pronounced effect on the overall behavior. This approach helps the model focus computational resources on zones where the impact of local material properties is greatest, such as interfaces or weak zones.

To improve the adaptability of the model in real-time applications, AMIS integrates a sequential prediction strategy inspired by the Kalman filter. This allows for the incorporation of real-time feedback into the model's predictions. The update rule for the predicted displacement vector $\mathbf{u}_t^{\text{pred}}$ at time step *t* is as follows:

$$\mathbf{u}_{t}^{\text{pred}} = \mathbf{u}_{t}^{\text{model}} + \mathbf{K}_{t}(\mathbf{z}_{t} - \mathbf{u}_{t}^{\text{model}}), \tag{44}$$

where $\mathbf{u}_t^{\text{model}}$ is the displacement predicted by the model at time step *t*, \mathbf{z}_t represents the real-time observation (such as measurements of displacement or strain), and \mathbf{K}_t is the Kalman gain. The Kalman gain is computed as:

$$\mathbf{K}_t = \mathbf{P}_t^{\text{model}} (\mathbf{P}_t^{\text{model}} + \mathbf{R})^{-1}, \tag{45}$$

where $\mathbf{P}_t^{\text{model}}$ is the model's estimate of the error covariance, and **R** is the observation noise covariance. This gain \mathbf{K}_t determines the weight given to the model prediction and the real-time observation, balancing the model's internal state and external data.

The model's error covariance $\mathbf{P}_t^{\text{model}}$ evolves over time according to the following update equation:

$$\mathbf{P}_{t}^{\text{model}} = \mathbf{A}_{t} \mathbf{P}_{t-1}^{\text{model}} \mathbf{A}_{t}^{T} + \mathbf{Q}, \tag{46}$$

where \mathbf{A}_t is the state transition matrix and \mathbf{Q} is the process noise covariance. This equation captures the evolution of uncertainty in the model's predictions as new data becomes available.

To correct the model's state based on the real-time feedback, the observation update is computed as:

$$\mathbf{P}_{t}^{\text{update}} = (\mathbf{I} - \mathbf{K}_{t}\mathbf{H})\mathbf{P}_{t}^{\text{model}},\tag{47}$$

where **H** is the observation matrix that maps the state to the observed variables, and **I** is the identity matrix. This update reduces the uncertainty in the model's predictions as it integrates more accurate real-time data.

the real-time prediction is integrated into the model's error correction process by updating the state vector $\mathbf{u}_t^{\text{pred}}$ after each new observation is received:

$$\mathbf{u}_t^{\text{pred}} = \mathbf{u}_t^{\text{model}} + \mathbf{K}_t (\mathbf{z}_t - \mathbf{H} \mathbf{u}_t^{\text{model}}).$$
(48)

4 Experimental setup

4.1 Dataset

The COCO Dataset (Tong and Wu, 2023) is a widely used large-scale dataset designed for object detection, segmentation, and captioning tasks. It consists of over 330,000 images, including 200,000 labeled images with 80 object categories. COCO provides a challenging benchmark due to its diversity, dense annotations, and complex scenes with multiple objects per image. The dataset supports tasks such as instance segmentation, keypoint detection, and panoptic segmentation (Lin et al., 2014). The GeoNet Dataset (Hanson et al., 2024) is curated for geo-spatial and environmental applications, emphasizing remote sensing imagery analysis. It contains multi-source satellite images, including multispectral and hyperspectral data, spanning various geographic locations and climatic conditions. GeoNet enables research in tasks such as land cover classification, disaster monitoring, and urban planning, offering high-resolution spatial data critical for robust model training and evaluation (Hanson et al., 2024). The SEN12MS Dataset (Ebel et al., 2022) is a



multimodal dataset tailored for semantic segmentation and land use classification. It includes over 180,000 patches of paired Sentinel-1 and Sentinel-2 satellite imagery. Sentinel-1 provides radar data, while Sentinel-2 offers optical images. SEN12MS covers diverse biomes, seasons, and atmospheric conditions, making it a valuable resource for fusion-based approaches and robust generalization across regions (Ebel et al., 2022). The SEN1-2 Dataset (Xu et al., 2023b) focuses on paired radar and optical satellite images from Sentinel-1 and Sentinel-2, catering to tasks such as domain adaptation, data fusion, and cross-modal learning. The dataset includes over 282,384 paired samples, covering various land cover types across continents. SEN1-2 enhances research into synergistic use of radar and optical data, particularly for applications in agriculture, forestry, and environmental monitoring (Xu et al., 2023b).

4.2 Experimental details

The experimental setup is designed to comprehensively evaluate the proposed method across diverse datasets. All experiments are conducted using PyTorch on a workstation equipped with NVIDIA A100 GPUs with 40GB memory. The training pipeline utilizes mixed precision to optimize computational efficiency and memory usage. The optimizer is AdamW with a weight decay of 10^{-4} , and the learning rate is initialized to 10^{-3} , following a cosine annealing schedule. Batch size is set to 16 for all experiments unless explicitly stated otherwise. Data augmentation techniques such as random cropping, flipping, rotation, and color jittering are applied during training to improve model generalization. For datasets involving

remote sensing, additional augmentations such as Gaussian noise and histogram equalization are included to simulate sensor variabilities. All input images are resized to 256 imes 256 for computational consistency, with channel normalization applied based on the dataset-specific statistics. The proposed method employs a backbone architecture based on Swin Transformer for feature extraction, coupled with task-specific heads. Pretrained weights on ImageNet-1K are utilized to initialize the backbone. For semantic segmentation tasks, a decoder module with a multi-scale attention mechanism is integrated to enhance spatial resolution and contextual understanding. In cross-modal experiments, feature alignment is achieved through a shared embedding space optimized using a contrastive loss. Training proceeds for 100 epochs for COCO and 50 epochs for remote sensing datasets like SEN12MS and SEN1-2. The early stopping criterion is based on validation performance, monitored through intersection-over-union (IoU) for segmentation and mean average precision (mAP) for detection. Validation is conducted every five epochs to ensure efficient resource utilization. Evaluation metrics vary based on the task. For object detection on COCO, mAP@0.5:0.95 is employed. For semantic segmentation tasks on SEN12MS and SEN1-2, mean IoU is the primary metric. Ablation studies are conducted by systematically removing key components of the model, such as the multi-scale attention mechanism and contrastive loss, to analyze their impact on performance. All results are averaged over three independent runs to ensure statistical significance. Model checkpoints and logs are saved for reproducibility. Hyperparameter tuning is performed through grid search, focusing on learning rates, weight decay, and augmentation strength. The implementation adheres to rigorous standards to ensure comparability with stateof-the-art benchmarks and facilitate reproducibility.

Model	COCO dataset				GeoNet dataset				
	Accuracy	Recall	F1 score	AUC	Accuracy	Recall	F1 score	AUC	
CLIP (Jiang et al., 2024)	88.76 ± 0.02	85.34 ± 0.03	87.29 ± 0.02	89.45 ± 0.03	86.54 ± 0.03	83.21 ± 0.02	85.12 ± 0.03	87.36 ± 0.02	
ViT (Huang et al., 2022)	90.43 ± 0.03	86.72 ± 0.02	89.05 ± 0.03	91.32 ± 0.02	88.76 ± 0.02	84.85 ± 0.03	86.93 ± 0.03	88.44 ± 0.02	
I3D (Long et al., 2022)	87.12 ± 0.02	83.90 ± 0.02	85.77 ± 0.03	88.21 ± 0.02	85.41 ± 0.03	81.54 ± 0.02	83.96 ± 0.02	85.72 ± 0.03	
BLIP (Noori et al., 2023)	89.54 ± 0.03	86.15 ± 0.03	88.43 ± 0.02	90.27 ± 0.03	87.62 ± 0.02	85.33 ± 0.02	86.81 ± 0.03	89.15 ± 0.02	
Wav2Vec 2.0 (Yu et al., 2024)	91.32 ± 0.02	87.88 ± 0.02	89.77 ± 0.03	91.60 ± 0.02	89.45 ± 0.03	86.71 ± 0.02	88.13 ± 0.03	90.02 ± 0.03	
T5 (Mohan et al., 2024)	89.01 ± 0.02	84.45 ± 0.03	86.98 ± 0.02	89.50 ± 0.03	87.34 ± 0.03	82.89 ± 0.02	84.76 ± 0.02	86.72 ± 0.03	
Ours	93.21 ± 0.03	91.05 ± 0.02	92.34 ± 0.03	94.12 ± 0.02	91.83 ± 0.03	89.41 ± 0.02	90.57 ± 0.03	92.08 ± 0.03	

TABLE 1 Comparison of our model with SOTA methods on COCO and GeoNet datasets for multimodal learning.

TABLE 2 Comparison of our model with SOTA methods on SEN12MS and SEN1-2 Datasets for multimodal learning.

Model	SEN12MS Dataset				SEN1-2 Dataset			
	Accuracy	Recall	F1 score	AUC	Accuracy	Recall	F1 score	AUC
CLIP (Jiang et al., 2024)	84.12 ± 0.02	80.76 ± 0.03	82.33 ± 0.02	85.41 ± 0.02	83.14 ± 0.03	79.85 ± 0.03	81.60 ± 0.03	84.25 ± 0.02
ViT (Huang et al., 2022)	86.43 ± 0.03	83.54 ± 0.02	85.21 ± 0.03	87.19 ± 0.03	85.27 ± 0.02	81.96 ± 0.03	83.77 ± 0.02	86.12 ± 0.03
I3D (Long et al., 2022)	83.75 ± 0.02	79.32 ± 0.02	81.44 ± 0.03	84.67 ± 0.02	81.88 ± 0.03	77.59 ± 0.02	80.23 ± 0.03	83.41 ± 0.02
BLIP (Noori et al., 2023)	85.91 ± 0.03	82.47 ± 0.02	84.32 ± 0.02	86.88 ± 0.03	84.67 ± 0.02	80.72 ± 0.03	82.95 ± 0.03	85.78 ± 0.03
Wav2Vec 2.0 (Yu et al., 2024)	87.39 ± 0.02	84.12 ± 0.03	85.98 ± 0.02	88.23 ± 0.03	86.54 ± 0.03	$82.89\ \pm 0.02$	84.77 ± 0.03	87.32 ± 0.02
T5 (Mohan et al., 2024)	85.21 ± 0.02	81.67 ± 0.03	83.54 ± 0.02	86.02 ± 0.03	84.11 ± 0.03	80.18 ± 0.02	82.10 ± 0.03	85.12 ± 0.03
Ours	89.45 ± 0.03	87.21 ± 0.02	88.77 ± 0.03	90.12 ± 0.02	88.36 ± 0.02	85.72 ± 0.03	87.34 ± 0.02	89.41 ± 0.03

We selected a learning rate of 10^{-3} and a weight decay of 10^{-4} based on previous research and experimental validation. In deep learning optimization, a learning rate of 10^{-3} is widely used with the AdamW optimizer, as studies have demonstrated its stability in Transformer-based architectures and multimodal learning tasks. a weight decay of 10^{-4} is chosen to mitigate overfitting, particularly on large-scale datasets. Prior research has shown that appropriate weight decay enhances generalization. we conducted grid search experiments and found that this combination achieved the best performance on the COCO, GeoNet, SEN12MS, and SEN1-2 datasets, balancing convergence speed and model stability. Therefore, this hyperparameter selection is justified by both theoretical foundations and empirical results.

4.3 Comparison with SOTA methods

The proposed model demonstrates superior performance across all datasets compared to state-of-the-art (SOTA) methods,

as shown in Tables 1, 2. The results on the COCO and GeoNet datasets (Table 1) highlight that our model achieves the highest accuracy, recall, F1 score, and AUC, outperforming methods such as CLIP (Jiang et al., 2024), ViT (Huang et al., 2022), and Wav2Vec 2.0 (Yu et al., 2024). Specifically, our model achieves an accuracy of 93.21% on the COCO dataset, which is 1.89% higher than the previous best method, Wav2Vec 2.0. Similarly, on the GeoNet dataset, our method achieves an accuracy of 91.83%, outperforming Wav2Vec 2.0 by 2.38%. This improvement can be attributed to the proposed model's efficient multi-scale attention mechanism, which enhances feature extraction from complex multimodal data. The performance advantage is also evident in the results for the SEN12MS and SEN1-2 datasets (Table 2). For the SEN12MS dataset, our model achieves an accuracy of 89.45%, surpassing Wav2Vec 2.0 by 2.06%. Similarly, on the SEN1-2 dataset, our model achieves an accuracy of 88.36%, which is 1.82% higher than the second-best method. This improvement is consistent across metrics, including recall, F1 score, and AUC. The robust performance on remote sensing datasets is attributed

Model		COCO	Dataset		GeoNet Dataset				
	Accuracy	Recall	F1 score	AUC	Accuracy	Recall	F1 score	AUC	
w/o Physics and Data-driven	89.23 ± 0.02	85.12 ± 0.03	86.77 ± 0.03	88.54 ± 0.02	87.45 ± 0.03	83.71 ± 0.02	85.02 ± 0.03	86.89 ± 0.02	
w/o Multi-scale Handling	91.12 ± 0.03	87.03 ± 0.02	89.21 ± 0.03	90.77 ± 0.03	89.09 ± 0.02	$85.67\ \pm 0.03$	87.14 ± 0.02	89.32 ± 0.03	
w/o Multi-scale Feature Aggregation	90.05 ± 0.02	86.45 ± 0.03	88.11 ± 0.02	89.63 ± 0.03	88.02 ± 0.03	84.94 ± 0.02	86.43 ± 0.03	88.71 ± 0.03	
Ours	93.21 ± 0.03	91.05 ± 0.02	92.34 ± 0.03	94.12 ± 0.02	91.83 ± 0.03	89.41 ± 0.02	90.57 ± 0.03	92.08 ± 0.03	

TABLE 3 Ablation study results on our model across COCO and GeoNet datasets for multimodal learning.

TABLE 4	Ablation study	results on our	model across	SEN12MS	and SEN1-2	datasets f	or multimodal	learning
---------	----------------	----------------	--------------	---------	------------	------------	---------------	----------

Model		SEN12M	S Dataset		SEN1-2 Dataset			
	Accuracy	Recall	F1 score	AUC	Accuracy	Recall	F1 score	AUC
w/o Physics and Data-driven	84.12 ± 0.02	80.45 ± 0.03	82.67 ± 0.02	85.54 ± 0.03	82.67 ± 0.03	79.21 ± 0.02	80.98 ± 0.03	83.72 ± 0.02
w/o Multi-scale Handling	86.78 ± 0.03	83.12 ± 0.02	84.98 ± 0.03	87.65 ± 0.03	85.32 ± 0.02	81.76 ± 0.03	83.54 ± 0.03	86.03 ± 0.03
w/o Multi-scale Feature Aggregation	85.67 ± 0.02	81.98 ± 0.03	83.45 ± 0.02	86.24 ± 0.02	84.45 ± 0.03	80.89 ± 0.02	82.67 ± 0.03	85.21 ± 0.03
Ours	89.45 ± 0.03	87.21 ± 0.02	88.77 ± 0.03	90.12 ± 0.02	88.36 ± 0.02	85.72 ± 0.03	87.34 ± 0.02	89.41 ± 0.03

to the effective cross-modal fusion facilitated by the contrastive loss, which aligns features from radar and optical imagery in a shared latent space. Furthermore, the ability to generalize across diverse geographical and atmospheric conditions underscores the model's adaptability.

Analyzing the reasons for the improved performance reveals several critical factors. First, the Swin Transformer backbone effectively captures both global and local dependencies, which is crucial for datasets like COCO and GeoNet with complex spatial patterns. Second, the proposed multi-scale attention mechanism ensures that fine-grained details are preserved, which is particularly important for remote sensing datasets like SEN12MS and SEN1-2. Third, the inclusion of domain-specific augmentations enhances robustness against noise and variations in satellite imagery, enabling the model to generalize effectively across different regions. The rigorous training pipeline, including the use of a cosine annealing learning rate schedule and mixed precision training, ensures optimal utilization of computational resources while avoiding overfitting.

4.4 Ablation study

The ablation study systematically evaluates the contribution of key components of the proposed model. The results for the COCO and GeoNet datasets are presented in Table 3, and for the SEN12MS and SEN1-2 datasets in Table 4. Three model variants are considered: (1) w/o Physics and Data-Driven, which excludes the multi-scale attention mechanism; (2) w/o Multi-Scale Handling, which omits the contrastive loss for feature alignment; and (3) w/o Multi-Scale Feature Aggregation, which removes the domain-specific augmentation pipeline. The full model, referred to as Ours, includes all components. The results indicate that each component significantly enhances the model's performance. For the COCO dataset, removing the multi-scale attention mechanism (w/o Physics and Data-Driven) leads to a decrease in accuracy from 93.21 to 89.23%, demonstrating the importance of capturing both global and local features. On the GeoNet dataset, the absence of the contrastive loss (w/o Multi-Scale Handling) results in a 2.74% drop in accuracy, highlighting its critical role in aligning features from different modalities. Similarly, the removal of domain-specific augmentations (w/o Multi-Scale Feature Aggregation) reduces the AUC from 94.12 to 89.63% on the COCO dataset and from 92.08 to 88.71% on the GeoNet dataset, underscoring the value of robust preprocessing techniques in enhancing model generalization.

For the SEN12MS and SEN1-2 datasets, a similar trend is observed (Figure 6). The absence of the multi-scale attention mechanism (w/o Physics and Data-Driven) causes a significant drop in accuracy from 89.45 to 84.12% on the SEN12MS dataset and from 88.36 to 82.67% on the SEN1-2 dataset. This underscores the importance of detailed spatial feature extraction, especially for remote sensing imagery. The contrastive loss (w/o Multi-Scale Handling) contributes to aligning radar and optical features, and its removal decreases F1 scores on both datasets by more than 3%. Domain-specific augmentations (w/o Multi-Scale Feature Aggregation) are essential for handling variations in atmospheric and geographic conditions, as evident from the consistent decrease in all metrics when they are excluded. The ablation study confirms that the multi-scale attention mechanism, contrastive loss, and domain-specific augmentations are indispensable for achieving state-of-the-art performance. Their integration ensures effective feature extraction, cross-modal alignment, and robustness





against dataset-specific challenges. This holistic design enables the proposed model to outperform SOTA methods across all datasets and tasks, as discussed in previous sections. These results reinforce the importance of designing multimodal learning systems with task-specific enhancements to maximize performance and generalization capabilities.

The absence of multi-scale attention significantly impacts datasets like COCO and SEN12MS because these datasets contain complex spatial structures and varying levels of detail that require hierarchical feature extraction (Figure 7). For COCO, which consists of images with multiple objects at different scales and varying backgrounds, multi-scale attention is essential for effectively capturing both fine-grained and global contextual information. Without it, the model struggles to differentiate between small and large objects, leading to suboptimal feature extraction and reduced accuracy in tasks like object detection and segmentation. For SEN12MS, which includes multimodal remote sensing data, multi-scale attention is crucial for integrating information across different spatial resolutions. Remote sensing images contain land cover variations that appear at different scales, and ignoring this hierarchy weakens the model's ability to capture both local textures and broader geographical patterns. Without multi-scale attention, the model may fail to effectively align information across modalities, leading to a decline in segmentation and classification performance. In both cases, multi-scale attention enables the model to dynamically focus on relevant spatial regions, improving feature representation and generalization across different levels of granularity. Its absence results in the loss of critical spatial dependencies, negatively impacting predictive performance.

Remote sensing datasets present significant challenges that make contrastive loss essential for effective performance. These datasets often contain multimodal data, such as optical, radar, and infrared imagery, each with different spatial resolutions, noise levels, and spectral characteristics, making it difficult for the model to align and integrate information effectively. Contrastive loss enforces similarity between corresponding multimodal representations while separating unrelated ones, improving cross-modal feature alignment. satellite images are captured at different times and under varying environmental conditions, introducing inconsistencies due to changes in lighting, seasonal variations, and atmospheric effects. Contrastive loss aids in learning invariant representations by ensuring that features extracted from the same location across different conditions remain close in the feature space, enhancing model robustness. Another key challenge is the imbalanced and sparse nature of annotations in remote sensing datasets, where certain land cover types are underrepresented. Contrastive learning enables the model to learn discriminative features by maximizing the distance between different classes, improving classification performance even with limited annotated data. Many remote sensing tasks also involve distinguishing between visually similar categories, such as different types of vegetation or urban structures, where standard classification losses may struggle to capture subtle differences, whereas contrastive loss refines the feature space by ensuring that semantically similar samples are closer while distinct classes remain well-separated. satellite images from different geographic regions exhibit domain shifts due to variations in terrain, climate, and sensor characteristics, making generalization challenging. Contrastive loss helps by learning feature representations that are less dependent on specific regional characteristics, thereby improving transferability across diverse datasets. By addressing these challenges, contrastive loss enhances the model's ability to learn robust, transferable, and well-aligned feature representations, leading to improved performance on remote sensing tasks such as land cover classification, segmentation, and change detection.

5 Discussion and conclusion

ModuCLIP has significant implications for geotechnical engineering and real-world applications beyond its technical advancements. In construction safety, the framework enhances the monitoring and prediction of foundation pit deformation, reducing the risk of structural failures and ensuring safer excavation processes. By integrating real-time multi-modal robotic sensing with advanced predictive modeling, it enables early detection of potential hazards, allowing for proactive mitigation measures that improve worker and infrastructure safety. In urban planning, ModuCLIP provides a data-driven approach for assessing ground stability in densely populated areas, offering valuable insights into soil behavior under different environmental and loading conditions. This can aid in optimizing excavation strategies, minimizing ground displacement risks, and improving the planning of underground structures such as tunnels, subways, and basements. In environmental monitoring, the framework's ability to analyze soil deformation dynamics can contribute to assessing land subsidence, erosion, and the long-term impact of urbanization on natural terrain. By incorporating satellite and sensor data, it can also help track the effects of climate change on soil stability, supporting sustainable land development practices. The adaptability of ModuCLIP to different geotechnical conditions and its capacity for real-time analysis make it a valuable tool for policymakers, engineers, and urban developers, bridging the gap between advanced AI-driven modeling and practical engineering applications.

ModuCLIP incorporates several features that enable realtime deployment, particularly through its Kalman filter-based feedback integration, which continuously refines predictions by incorporating real-time sensor data. This mechanism allows the model to dynamically adjust deformation estimates as new measurements become available, reducing prediction drift and enhancing robustness in dynamic environments. The framework also leverages hierarchical attention mechanisms to prioritize critical regions in geotechnical monitoring, ensuring computational efficiency and faster response times. the use of adaptive multi-scale integration allows ModuCLIP to handle varying spatial and temporal resolutions, making it highly responsive to changing excavation conditions. In practical settings, ModuCLIP can be deployed in active construction sites to monitor foundation pit stability, where real-time predictions help engineers detect anomalies early and take corrective actions before failures occur. It is also applicable in tunnel excavation projects, where continuous assessment of ground deformation can prevent collapses and optimize reinforcement strategies. In urban infrastructure monitoring, ModuCLIP can integrate with IoT-based sensor networks to track long-term ground settlement and subsidence in high-risk areas, supporting preventive maintenance. in environmental hazard assessment, it can be used for landslide prediction in mountainous regions by analyzing soil displacement patterns in real-time, aiding disaster prevention efforts. By combining advanced AI-driven modeling with real-time adaptability, ModuCLIP enhances decision-making in safety-critical geotechnical applications.

Humans can interact with and benefit from ModuCLIP in multiple ways, particularly in fields like construction safety, urban infrastructure monitoring, and environmental hazard prevention. Engineers and project managers can use ModuCLIP as a decision-support tool, leveraging its real-time predictive capabilities to assess foundation pit stability and adjust excavation strategies proactively. By integrating with IoT sensors and robotic systems, ModuCLIP enables automated alerts for potential ground deformation risks, allowing workers to take preventive measures before structural failures occur. Urban planners can benefit from ModuCLIP by using its multi-scale geotechnical analysis to assess ground conditions in infrastructure projects such as subways, tunnels, and high-rise buildings, ensuring long-term structural integrity. In disaster management, authorities can deploy the framework for landslide monitoring by continuously analyzing soil displacement patterns, providing early warnings to prevent catastrophic events. environmental scientists can use ModuCLIP to track land subsidence and erosion trends, informing sustainable land development policies. By making geotechnical predictions more accurate and accessible, ModuCLIP reduces uncertainty in construction and environmental projects, improving safety, efficiency, and cost-effectiveness. Its integration with real-time sensing and automated feedback mechanisms ensures that human operators receive actionable insights without requiring extensive manual analysis, ultimately bridging the gap between advanced AI-driven modeling and practical engineering applications.

Humans play a critical role in introducing oversight and addressing uncertainties in our study. While robotic systems handle tasks such as data collection, real-time monitoring, and some decision-making, human involvement is essential when the system faces uncertainties or novel situations. Humans are responsible for interpreting complex scenarios, validating predictions, and providing judgment calls when the system encounters conditions it has not been trained for or in the presence of limited data. Human oversight ensures that robotic interventions align with safety standards, regulatory requirements, and unforeseen environmental factors. humans are responsible for the initial setup of the system, providing domain expertise to guide design and operations. After deployment, they continue to monitor the system and intervene when necessary, especially in extreme or unforeseen circumstances that the robotic systems cannot handle. humans play an indispensable role in supervising, adjusting, and optimizing the robotic systems, particularly when dealing with complexity and uncertainty.

This framework has the potential to significantly improve workers' routines, safety, and quality of life in geotechnical engineering. By integrating robotic systems and advanced prediction models, such as ModuCLIP, the framework can automate many of the repetitive and hazardous tasks typically performed in foundation pit monitoring. This automation reduces human exposure to dangerous environments, such as unstable excavation sites, improving overall safety. the system's ability to predict and monitor foundation pit deformation in real-time enhances decision-making, allowing for timely interventions and minimizing the risks of structural failure. This predictive capability ensures that potential hazards are identified early, preventing accidents and ensuring that workers are not in danger when critical issues arise. by reducing the need for constant manual oversight and allowing for more efficient, data-driven workflows, workers can focus on higher-value tasks, improving overall productivity. This enhanced work efficiency not only optimizes workers' routines but also leads to a better work-life balance by minimizing the physical and mental strain associated with traditional geotechnical engineering practices.

In high-stakes applications like geotechnical engineering, ensuring the reliability of robotic systems is crucial. The framework must be designed with robust safety protocols, including redundant systems and fail-safes, to prevent system failures that could lead to accidents. Since these systems will operate in environments where human lives and property are at risk, their predictions and actions must be thoroughly validated. This is where human-in-the-loop (HITL) validation becomes essential. While the system can perform real-time monitoring and make predictive recommendations, humans must oversee critical decisions, especially in complex or uncertain scenarios. Human involvement ensures that decisions align with safety standards, legal regulations, and ethical guidelines, offering a necessary layer of accountability and flexibility in case of system limitations or unexpected situations. Additionally, ongoing testing, continuous monitoring, and transparency in decision-making processes will further enhance the ethical use of this technology in sensitive environments.

The study focuses on predicting foundation pit deformation to ensure structural stability and safety in multi-modal robotic geotechnical systems. Traditional methods often fail to address challenges posed by dynamic environmental factors, sparse data, and the nonlinear nature of soil behavior. To overcome these limitations, the study introduces ModuCLIP, a novel Multi-Scale CLIP framework that integrates neural learning techniques with geotechnical expertise. The framework employs a hybrid architecture that combines physics-informed encoders with deep learning modules to extract multi-scale spatial-temporal features. Key innovations include a recurrent temporal mechanism for adapting to dynamic excavation processes and an attentionbased strategy for modeling heterogeneous geotechnical zones. Regularization terms are included to ensure predictions adhere to physical principles, such as stress equilibrium and soil constitutive relationships. Experiments confirm that ModuCLIP achieves stateof-the-art accuracy, outperforming traditional methods in diverse scenarios. This framework bridges the gap between data-driven models and physical laws, enabling real-time applications in robotic monitoring systems.

Despite its promising results, ModuCLIP has some limitations. First, the reliance on extensive training data poses challenges in environments with severe data scarcity. Future research could explore transfer learning or unsupervised techniques to mitigate this dependency. Second, the computational complexity of the framework, particularly its attention-based and recurrent components, may limit its scalability to larger or more complex geotechnical systems. Future developments could focus on optimizing the framework's computational efficiency or leveraging edge-computing capabilities. By addressing these challenges, ModuCLIP has the potential to become a robust tool for multimodal robotic systems in geotechnical engineering.

The issue of data scarcity in our study arises in scenarios where there is limited historical data on foundation pit deformation, especially in less monitored or newly constructed areas. For instance, in some geotechnical engineering sites, the data on soil properties, deformation behavior, and excavation conditions may be sparse due to limited sensor deployments or short monitoring periods. These data gaps can hinder the model's ability to make accurate predictions, as the system relies heavily on robust datasets for training. To address this, we propose several solutions. transfer learning can be applied, where models trained on well-established datasets from similar regions or projects can be adapted to the specific characteristics of the new site with limited data. This approach would allow the model to leverage knowledge gained from other geotechnical environments, improving its predictive accuracy even with less local data. data augmentation strategies can be employed to artificially increase the dataset size. For example, by applying transformations such as noise addition, synthetic data generation through simulations, or using domainspecific augmentations like varying soil moisture or pressure conditions, the model can be trained to be more resilient to the limited real-world data. These strategies would help fill in the gaps and improve the model's generalization ability in the face of data scarcity.

We acknowledge that the reliance on specific datasets, such as COCO and SEN12MS, can introduce biases and potentially limit the generalization of our framework to diverse geotechnical conditions. These datasets, while valuable for demonstrating the framework's capability in controlled environments, may not fully capture the variability present in real-world geotechnical scenarios. For example, COCO focuses on object detection in visual data, which may not always reflect the complexity of soil behavior, structural deformation, or environmental conditions in foundation pit monitoring. To address this, future work should focus on expanding the framework's applicability by incorporating more diverse and underrepresented geotechnical datasets that better reflect a range of soil types, excavation methods, and environmental conditions. Collaborating with industry partners, we can gather real-world data from a wider range of geotechnical sites, including those with varying soil properties, climatic conditions, and construction practices. domain adaptation techniques could be employed to adapt the model to datasets from different regions or contexts. This would involve adjusting the model's parameters based on data from underrepresented sites to ensure it generalizes well across diverse scenarios. Another direction could be the use of unsupervised learning or semi-supervised learning methods to help the model better cope with the scarcity of labeled data, allowing it to learn from the available data more effectively and expanding its applicability to new geotechnical contexts. These steps would ensure that the framework becomes more robust and can provide reliable predictions across a broader set of geotechnical environments.

In our study, the computational complexity associated with attention-based and recurrent components plays a significant role in the overall performance of the framework. attention mechanisms, which operate over multiple scales of data, can be computationally intensive due to the need to process large amounts of multi-modal data and compute pairwise relationships across all data points. the recurrent components, such as Long Short-Term Memory (LSTM) networks, further contribute to computational complexity by requiring the processing of sequential data over time, which increases memory usage and computational load. During experimentation, we encountered bottlenecks primarily related to the scalability of the attention mechanism and the training time required for the recurrent layers. The attention mechanism, in particular, is sensitive to the size of the input data and the number of features being processed. As the model scales to handle large, multi-modal datasets with complex temporal and spatial dependencies, the quadratic growth of the attention mechanism's computational cost becomes evident, especially when processing high-dimensional data. To address these bottlenecks, we propose several optimization strategies. reducing the dimensionality of the input features through techniques such as feature selection or dimensionality reduction (e.g., using principal component analysis or autoencoders) can significantly reduce the computational load while maintaining the important information. Approximate attention mechanisms, such as sparse attention or Linformer, which reduce the complexity of pairwise interactions, can be implemented to scale the attention process more efficiently. For the recurrent components, optimizing the training process through methods like gradient checkpointing can reduce memory consumption by storing only a subset of intermediate results during the forward pass, while recomputing them during the backward pass. leveraging more parallel processing or distributing the computational load across multiple GPUs can help mitigate the long training times associated with sequential data processing.

Data availability statement

The original contributions presented in the study are included in the article/supplementary material, further inquiries can be directed to the corresponding author.

Author contributions

LW: Conceptualization, Methodology, Writing – original draft, Writing – review & editing, Visualization, Supervision, Funding acquisition. LT: Software, Validation, Data curation, Writing – original draft. LX: Formal analysis, Investigation, Writing – original draft, Writing – review & editing.

Funding

The author(s) declare that no financial support was received for the research and/or publication of this article.

Conflict of interest

LX was employed by Guangdong Nonferrous Industry Building Quality Inspection Co., Ltd.

The remaining authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Generative AI statement

The author(s) declare that no Gen AI was used in the creation of this manuscript.

Publisher's note

All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers. Any product that may be evaluated in this article, or claim that may be made by its manufacturer, is not guaranteed or endorsed by the publisher.

References

Bayoudh, K., Knani, R., Hamdaoui, F., and Mtibaa, A. (2021). A survey on deep multimodal learning for computer vision: advances, trends, applications, and datasets. *Vis. Comput.* 38, 2939–2970. doi: 10.1007/s00371-021-02166-7

Cao, H., Zhang, Z., Xia, Y., Li, X., Xia, J., Chen, G., et al. (2024). "Embracing events and frames with hierarchical feature refinement network for object detection," in *European Conference on Computer Vision* (Cham: Springer), 161–177. doi: 10.1007/978-3-031-72907-2_10

Chai, W., and Wang, G. (2022). Deep vision multimodal learning: Methodology, benchmark, and trend. *Appl. Sci.* 12:6588. doi: 10.3390/app12136588

Chango, W., Lara, J., Cerezo, R., and Romero, C. (2022). A review on data fusion in multimodal learning analytics and educational data mining. *WIREs Data Mining Knowl. Discov.* 12:e1458. doi: 10.1002/widm.1458

Cui, Y., Wang, Q., Li, C., Ren, W., and Knoll, A. (2025). Eenet: an effective and efficient network for single image dehazing. *Pattern Recognit.* 158:111074. doi: 10.1016/j.patcog.2024.111074

Daunhawer, I., Bizeul, A., Palumbo, E., Marx, A., and Vogt, J. E. (2023). "Identifiability results for multimodal contrastive learning," in *International Conference* on Learning Representations. Available at: https://arxiv.org/abs/2303.09166

Du, C., Fu, K., Li, J., and He, H. (2022). Decoding visual neural representations by multimodal learning of brain-visual-linguistic features. *IEEE Trans. Pattern Anal. Mach. Intell.* 45, 10760–10777. doi: 10.1109/TPAMI.2023.3263181

Ebel, P., Xu, Y., Schmitt, M., and Zhu, X. X. (2022). Sen12ms-cr-ts: a remote-sensing data set for multimodal multitemporal cloud removal. *IEEE Trans. Geosci. Remote Sens.* 60, 1–14. doi: 10.1109/TGRS.2022.3146246

Ektefaie, Y., Dasoulas, G., Noori, A., Farhat, M., and Zitnik, M. (2022). Multimodal learning with graphs. *Nat. Mach. Intell.* 5, 340–350. doi: 10.1038/s42256-023-00624-6

Fan, Y., Xu, W., Wang, H., Wang, J., and Guo, S. (2022). "PMR: prototypical modal rebalance for multimodal learning," in 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (Vancouver, BC: IEEE). doi: 10.1109/CVPR52729.2023.01918

Gao, X., Rogel, A., Sankaranarayanan, R., Dowling, B., and Weinberg, G. (2024). Music, body, and machine: gesture-based synchronization in human-robot musical interaction. *Front. Robot. AI* 11:1461615. doi: 10.3389/frobt.2024.1461615

Hanson, J. B., Sherburn, S., Behr, Y., Britten, K. M., Hughes, E. C., Jarvis, P. A., et al. (2024a). Twenty years of volcano data at geonet-collection, custodianship, and evolution of open data on New Zealand's volcanoes. *Bull. Volcanol.* 86:81. doi:10.1007/s00445-024-01769-x

Hao, Y., Stuart, T., Kowalski, M. H., Choudhary, S., Hoffman, P. J., Hartman, A., et al. (2022). Dictionary learning for integrative, multimodal and scalable single-cell analysis. *bioRxiv*. doi: 10.1101/2022.02.24.481684

He, S., Dai, S.-L., Zhao, Z., and Zou, T. (2022). Uncertainty and disturbance estimator-based distributed synchronization control for multiple marine surface vehicles with prescribed performance. *Ocean Eng.* 261:111867. doi: 10.1016/j.oceaneng.2022.111867

He, S., Dai, S.-L., Zhao, Z., Zou, T., and Ma, Y. (2023). Ude-based distributed formation control for msvs with collision avoidance and connectivity preservation. *IEEE Trans. Ind. Inform.* 20, 1476–1487. doi: 10.1109/TII.2023.3274234

Hu, J., Yao, Y., Wang, C., Wang, S., Pan, Y., Chen, Q.-A., et al. (2023). "Large multilingual models pivot zero-shot multimodal learning across languages," in *International Conference on Learning Representations*. Available at: https://arxiv.org/ abs/2308.12038

Huang, P., He, P., Tian, S., Ma, M., Feng, P., Xiao, H., et al. (2022). A VIT-AMC network with adaptive model fusion and multiobjective optimization for interpretable laryngeal tumor grading from histopathological images. *IEEE Trans. Med. Imaging* 42, 15–28. doi: 10.1109/TMI.2022.3202248

Jiang, Z., Li, Z., and Yao, Z. (2024). Multi-objective optimization in fixedoutline floorplanning with reinforcement learning. *Comput. Electr. Eng.* 120:109784. doi: 10.1016/j.compeleceng.2024.109784

Joseph, J., Thomas, B., Jose, J., and Pathak, N. (2023). Decoding the growth of multimodal learning: a bibliometric exploration of its impact and influence. *Int. J. Intell. Decis. Technol.* 18, 151–167. doi: 10.3233/IDT-230727

Lian, Z., Chen, L., Sun, L., Liu, B., and Tao, J. (2022). GCNET: graph completion network for incomplete multimodal learning in conversation. *IEEE Trans. Pattern Anal. Mach. Intell.* 45, 8419–8432. doi: 10.1109/TPAMI.2023.3234553

Lin, T.-Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., et al. Dollár, P., Zitnick, C. L. (2014). "Microsoft COCO: common objects in context," in *Computer Vision-ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13* (Cham: Springer), 740–755. doi: 10.1007/978-3-319-10602-1_48

Lin, Z., Yu, S., Kuang, Z., Pathak, D., and Ramana, D. (2023). "Multimodality helps unimodality: cross-modal few-shot learning with multimodal models," in

2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (Vancouver, BC: IEEE). doi: 10.1109/CVPR52729.2023.01852

Liu, S., Cheng, H., Liu, H., Zhang, H., Li, F., Ren, T., et al. (2023). "LLaVA-plus: learning to use tools for creating multimodal agents," in *European Conference on Computer Vision* (Berlin: Springer). doi: 10.1007/978-3-031-72970-6_8

Long, Y., Zhang, Z., Yan, C., Huang, Z., Fu, K., Li, Y., et al. (2022). Multiobjective optimization for improving printing efficiency and mechanical properties of 3d-printed continuous plant fibre composites. *Compos. Commun.* 35:101283. doi: 10.1016/j.coco.2022.101283

Mizuho, T., Okafuji, Y., Baba, J., and Narumi, T. (2024). Multiple-agent promotion in a grocery store: effects of modality and variability of agents on customer memory. *Front. Robot. AI* 11:1397230. doi: 10.3389/frobt.2024.1397230

Mohan, G. B., Kumar, R. P., and Elakkiya, R. (2024). Enhancing pre-trained models for text summarization: a multi-objective genetic algorithm optimization approach. *Multimed. Tools Appl.* 1–17. doi: 10.1007/s11042-024-20374-w

Noori, M. S., Sahbudin, R. K., Sali, A., and Hashim, F. (2023). Feature drift aware for intrusion detection system using developed variable length particle swarm optimization in data stream. *IEEE Access* 11, 128596–128617. doi: 10.1109/ACCESS.2023.3333000

Peng, X., Wei, Y., Deng, A., Wang, D., and Hu, D. (2022). "Balanced multimodal learning via on-the-fly gradient modulation," in 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (New Orleans, LA: IEEE). doi: 10.1109/CVPR52688.2022.00806

Shi, B., Hsu, W.-N., and Lakhotia, K. rahman Mohamed, A. (2022). *Learning audio-visual speech representation by masked multimodal cluster prediction*. International Conference on Learning Representations.

Song, B., Miller, S., and Ahmed, F. (2023). Attention-enhanced multimodal learning for conceptual design evaluations. *J. Mech. Des.* 145:041410. doi: 10.1115/1.4056669

Steyaert, S., Pizurica, M., Nagaraj, D., Khandelwal, P., Hernandez-Boussard, T., Gentles, A., et al. (2023). Multimodal data fusion for cancer biomarker discovery with deep learning. *Nat. Mach. Intell.* 5, 351–362. doi: 10.1038/s42256-023-00633-5

Tong, K., and Wu, Y. (2023). Rethinking pascal-VOC and ms-COCO dataset for small object detection. J. Vis. Commun. Image Represent. 93:103830. doi: 10.1016/j.jvcir.2023.103830

Versaci, M., Angiulli, G., La Foresta, F., Crucitti, P. Laganá, F., Pellicanó, D., et al. (2022). "Innovative soft computing techniques for the evaluation of the mechanical stress state of steel plates," in *International conference on applied intelligence and informatics* (Cham: Springer), 14–28. doi: 10.1007/978-3-031-24801-6_2

Versaci, M., Angiulli, G., and La Foresta, F. Laganà, F., and Palumbo, A. (2024). Intuitionistic fuzzy divergence for evaluating the mechanical stress state of steel plates subject to bi-axial loads. *Integr. Comput.-Aided Eng.* 31, 1–17. doi: 10.3233/ICA-230730

Wei, S., Luo, Y., and Luo, C. (2023). "Mmanet: margin-aware distillation and modality-aware regularization for incomplete multimodal learning," in 2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (Vancouver, BC: IEEE). doi: 10.1109/CVPR52729.2023.01919

Xu, P., Zhu, X., and Clifton, D. (2022). "Multimodal learning with transformers: a survey," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*. Available at: https://openaccess.thecvf.com/content/CVPR2023/html/Wei_MMANet_ Margin-Aware_Distillation_and_Modality-Aware_Regularization_for_Incomplete_ Multimodal_Learning_CVPR_2023_paper.html

Xu, W., Wu, Y., and Fan, O. (2023a). 'Multimodal learning analytics of collaborative patterns during pair programming in higher education. *Int. J. Educ. Technol. High. Educ.* 28:8. doi: 10.1186/s41239-022-00377-z

Xu, W., Yuan, X., Hu, Q., and Li, J. (2023b). Sar-optical feature matching: a largescale patch dataset and a deep local descriptor. *Int. J. Appl. Earth Observ. Geoinform.* 122:103433. doi: 10.1016/j.jag.2023.103433

Yan, L., Zhao, L., Gašević, D., and Maldonado, R. M. (2022). "Scalability, sustainability, and ethicality of multimodal learning analytics," in *International Conference on Learning Analytics and Knowledge* (New York, NY: ACM). doi: 10.1145/3506860.3506862

Yang, Z., Fang, Y., Zhu, C., Pryzant, R., Chen, D., Shi, Y., et al. (2022). "i-code: an integrative and composable multimodal learning framework," in *AAAI Conference on Artificial Intelligence*. Available at: https://ojs.aaai.org/index.php/AAAI/article/view/ 26290

Yao, J., Zhang, B., Li, C., Hong, D., and Chanussot, J. (2023). Extended vision transformer (EXVIT) for land use and land cover classification: a multimodal deep learning framework. *IEEE Trans. Geosci. Remote Sens.* 61:5514415. doi: 10.1109/TGRS.2023.3284671

Yu, Q., Liu, Y., Wang, Y., Xu, K., and Liu, J. (2023). "Multimodal federated learning via contrastive representation ensemble," in *International Conference on Learning Representations*. Available at: https://arxiv.org/abs/2302.08888

Yu, X., Guo, D., Zhang, J., and Lin, Y. (2024). Rose: a recognitionoriented speech enhancement framework in air traffic control using multiobjective learning. *IEEE/ACM Trans. Audio Speech Lang. Process.* 32, 3365–3378. doi: 10.1109/TASLP.2024.3423652

Zhang, H., Zhang, C., Wu, B., Fu, H., Zhou, J. T., Hu, Q., et al. (2023). *Calibrating multimodal learning*. International Conference on Machine Learning.

Zhang, Y., He, N., Yang, J., Li, Y., Wei, D., Huang, Y., et al. (2022). "mmFormer: multimodal medical transformer for incomplete multimodal learning of brain tumor segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention* (Cham: Springer). doi: 10.1007/978-3-031-16443-9_11

Zhou, H., Bing, Z., Yao, X., Su, X., Yang, C., Huang, K., et al. (2024). Language-conditioned imitation learning with base skill priors under unstructured

data. IEEE Robot. Autom. Lett. 9, 9805–9812. doi: 10.1109/LRA.2024. 3466076

Zhou, H.-Y., Yu, Y., Wang, C., Zhang, S., Gao, Y., Pan, J.-Y., et al. (2023a). A transformer-based representation-learning model with unified processing of multimodal input for clinical diagnostics. *Nat. Biomed. Eng.* 7, 743–755. doi: 10.1038/s41551-023-01045-x

Zhou, Y., Wang, X., Chen, H., Duan, X., and Zhu, W. (2023b). "Intra- and intermodal curriculum for multimodal learning," in *ACM Multimedia* (New York, NY: ACM). doi: 10.1145/3581783.3612468

Zong, Y., Aodha, O. M., and Hospedales, T. M. (2023). Self-supervised multimodal learning: a survey. *IEEE Trans. Pattern Anal. Mach. Intell.* 1–20. doi: 10.1109/TPAMI.2024.3429301