



Research Article

<https://doi.org/10.1631/jzus.A2500300>



PL-HLNet: a semi-supervised approach for tunnel boring machine disc cutter wear prediction

Zhaoyang LI, Wei TANG[✉], Xinyuan WANG, Huxiu XU, Huayong YANG, Jun ZOU[✉]

State Key Laboratory of Fluid Power and Mechatronic Systems, School of Mechanical Engineering, Zhejiang University, Hangzhou 310058, China

Abstract: Unanticipated wear of tunnel boring machine (TBM) disc cutters is a critical factor causing project delays and cost overruns in tunneling engineering. Accurate, real-time prediction of the cutter's wear state is therefore essential for enabling predictive maintenance. Data-driven methods, particularly deep learning, have shown promise for this task, but their performance is constrained by the scarcity of high-quality labeled data in practical industrial settings. To address this challenge, we propose a novel, decoupled semi-supervised framework called PL-HLNet. The first component of this framework is a multi-view pseudo-labeling (PL) module, which mines high-confidence supervisory signals from massive unlabeled data by leveraging heterogeneous views derived from feature engineering and diverse model architectures; it is followed by a consistency check to ensure label quality. This process effectively augments the training set while correcting for sampling bias. Subsequently, a specialized hierarchical hybrid attention network (HLNet) is used to make predictions. The HLNet organically integrates a temporal convolutional network (TCN) for local feature extraction, a bidirectional long short-term memory (Bi-LSTM) network for capturing temporal dynamics, and a custom attention mechanism for focusing on critical information. Experiments on a real-world tunneling dataset show that PL-HLNet significantly outperforms both supervised and mainstream semi-supervised baselines, such as the Mean Teacher and FixMatch. The framework's effectiveness is further substantiated by validations of its architectural design and data-driven selection of hyperparameters. Moreover, PL-HLNet has a high inference speed, showcasing its practicality for real-world scenarios. Our work provides an effective solution for machining equipment monitoring in data-scarce industrial environments.

Key words: Tunnel boring machine (TBM); Cutter wear; Semi-supervised learning; Pseudo-labeling (PL); Temporal convolutional network; Bidirectional long short-term memory (Bi-LSTM); Attention mechanism

1 Introduction

Tunnel boring machines (TBMs) serve as core equipment for modern tunneling and subsurface engineering projects. Renowned for their high efficiency, safety, and environmental friendliness, they have become the primary solution for constructing long-distance, deep-overburden tunnels. As the concept of "intelligent construction" has gained traction, advancing the intelligence of the TBM excavation process—particularly through predictive maintenance of critical components—has emerged as a key research objective.

Within a TBM's numerous subsystems, the cutterhead, equipped with its disc cutters, is the essential component that directly performs rock-breaking operations. However, it is also the most vulnerable component of a TBM, as it operates under the harshest conditions, being constantly exposed to complex and variable geology while withstanding immense compressive forces and impact vibrations. Unforeseen cutter failure can lead to costly downtime while waiting for replacement, significantly disrupt project schedules, and even introduce safety hazards. Consequently, developing accurate, real-time prediction models for cutter wear status holds high theoretical and practical significance for guaranteeing the continuity, economic viability, and safety of TBM operations.

Current research on TBM disc cutter wear prediction predominantly follows two paradigms: physics-based modeling and data-driven approaches. While

✉ Wei TANG, weitang@zju.edu.cn

Jun ZOU, junzou@zju.edu.cn

Jun ZOU, <https://orcid.org/0000-0003-2443-3516>

Received July 7, 2025; Revision accepted Sept. 16, 2025;
Crosschecked Feb. 3, 2026; Online first Apr. 1, 2026

© Zhejiang University Press 2026

mechanistic models, such as the classic Colorado School of Mines (CSM) model, offer a theoretical foundation for the wear process, they often struggle to adapt to the complex, variable, and unpredictable geological conditions encountered in real-world tunneling projects; this is due to their reliance on simplified assumptions and sensitive geomechanical parameters (Barzegari et al., 2021; Karami et al., 2021; Shen et al., 2022; Sun et al., 2023). In contrast, data-driven methods demonstrate greater adaptability and potential by learning patterns from historical data. Classic machine learning algorithms, such as support vector regression (SVR) and random forests, have been successfully applied to the task of cutter wear prediction (Mahmoodzadeh et al., 2021; Agrawal et al., 2022; Ghorbani and Yagiz et al., 2024; Shin et al., 2024). More recently, deep learning models, particularly long short-term memory (LSTM) and gated recurrent unit (GRU) networks, have emerged as mainstream approaches due to their powerful capabilities in capturing temporal features (Yu HG et al., 2021; Jia et al., 2022; Li et al., 2022; Pu et al., 2022; Mo et al., 2024; Bai et al., 2025). Building on this, bidirectional LSTM (Bi-LSTM) has demonstrated robust results by capturing contextual information from both past and future time steps. For instance, recent studies have shown that for tasks with high temporal correlation, such as tunnelling parameter prediction, Bi-LSTM can effectively mine long-term trends and outperforms traditional machine learning or unidirectional models (Lu et al., 2024). Furthermore, a large body of recent work combines LSTM with mechanisms such as attention or with convolutional neural networks (CNNs), confirming LSTM's robustness as a foundation for modeling temporal dependencies (Wang et al., 2023). However, the efficacy of all these data-driven methods is dependent upon the availability of large-scale, high-quality labeled datasets (Ding et al., 2022; Zhang et al., 2023). In practice, acquiring precise wear labels is a costly and time-consuming manual process, which has led to a scarcity of labeled samples. This constitutes a fundamental limitation for existing methods.

Semi-supervised learning (SSL) appears to be a promising solution to this limitation, leveraging massive amounts of unlabeled data to aid model training. Although classic SSL paradigms, such as co-training, have proven successful in various fields (Dou et al., 2024; Shi et al., 2024; Wu et al., 2025), they are often

susceptible to confirmation bias; this is where early misclassifications in pseudo-labels are reinforced during the iterative process, leading to error accumulation and performance degradation (Yu FX et al., 2021; Chen et al., 2022). Concurrently, significant advancements are being made in predictive model architectures. Temporal convolutional networks (TCNs) have shown distinct advantages in extracting local, structural features from time-series data (Hewage et al., 2020; Fan et al., 2023), while the attention mechanism offers a powerful solution for addressing long-term dependencies and focusing on salient information (Cheng et al., 2022; Wen and Li, 2023; Wang et al., 2025). Despite their individual strengths, creating a hierarchical fusion of these methods and integrating them within a robust semi-supervised framework that mitigates error propagation is a difficult task.

To address the aforementioned challenges, particularly the issue of labeled data scarcity in industrial settings, we propose a novel, decoupled semi-supervised predictive framework called PL-HLANe. The central principle of this framework is to decouple the generation of supervisory information from unlabeled data from the subsequent deep learning-based prediction, thereby mitigating the risk of error accumulation inherent in conventional methods. The framework first employs a multi-view pseudo-labeling (PL) module to generate high-confidence labels by leveraging heterogeneous views from both feature engineering and diverse model architectures, which is followed by a consistency check. Subsequently, a hierarchical hybrid attention network (HLANet) is utilized for prediction. This network organically integrates a TCN, a Bi-LSTM network, and a custom attention mechanism to perform hierarchical feature learning (from local patterns to global trends and finally to key moments) on TBM time-series data. The HLANet is trained with a combination of true labels and high-quality pseudo-labels. Additionally, we define the "instantaneous wear rate" as the core predictive indicator to enhance the model's real-time responsiveness.

The main contributions of this paper include the following:

(1) We propose a decoupled, multi-view PL framework that effectively addresses the challenge of label scarcity in industrial applications by mining valuable training data from massive unlabeled datasets.

(2) We design an HLANet that fuses TCN, Bi-LSTM, and an attention mechanism to achieve multi-scale feature extraction, from local-transient to global-temporal patterns, within complex TBM excavation data.

(3) We define and apply the “instantaneous wear rate” as a dynamic and responsive predictive target, which is better suited for the practical requirements of predictive maintenance than traditional cumulative wear metrics.

(4) We experimentally validate the method on a complex real-world dataset from a major national tunnel project, systematically demonstrating the superiority of the proposed framework and the necessity of its components through various comparisons and ablation studies.

The remainder of this paper is organized as follows. Section 2 details the research object and dataset construction. Section 3 describes the proposed PL-HLANet framework. Section 4 presents and analyzes the results of the comparative experiments. Section 5 showcases ablation studies that validate the roles of each module. Finally, Section 6 summarizes the findings and outlines future work.

2 Materials

This section describes the materials employed in this study. We begin by presenting an overview of the TBM project, establishing the engineering context, and introducing the research object—a disc cutter system. The chapter then details the pipeline for transforming raw excavation data into a structured dataset, a process that involves defining key performance indicators and selecting input features for the model.

2.1 Project description

The dataset for this study originates from a major tunnel project traversing the core of a plateau in central-western China. With a total length over 38 km and a maximum overburden of nearly 1700 m, the project is characterized by a confluence of formidable challenges: high altitude, high in situ stress, long-distance excavation, and complex geological conditions. Preliminary geological surveys indicate that the tunnel alignment is predominantly situated in hard rock strata, with over half of its length passing through rock masses of fair-to-poor stability (an overview of the project is depicted in Fig. 1). To address these conditions, a 10-m class, multi-support, open-type TBM was employed for the excavation. The TBM excavation system is illustrated in Fig. 2, and its main technical specifications are detailed in Table 1.

2.2 Research object and problem definition

The object of this study is the disc cutter system of the TBM cutterhead employed in this tunneling project. As illustrated in Fig. 3, the cutterhead is equipped with a total of 66 disc cutters, which comprise 70 individual cutting edges, systematically numbered from C1 to C70. To maximize the cutting coverage for higher excavation efficiency and ensure uniform load distribution for enhanced structural stability, each cutter is installed at a unique radial distance; this layout is presented in Fig. 3a. Based on their distribution areas, which are shown in Fig. 3b, the cutters are categorized into three distinct types: center cutters (C1–C8), which consist of four 17-inch twin-disc models (1 inch=0.0254 m); face cutters (C9–C56); gauge cutters (C57–C70), with the latter two types utilizing 19-inch single-disc cutters. To ensure a representative



Fig. 1 Overview of the TBM tunnel project

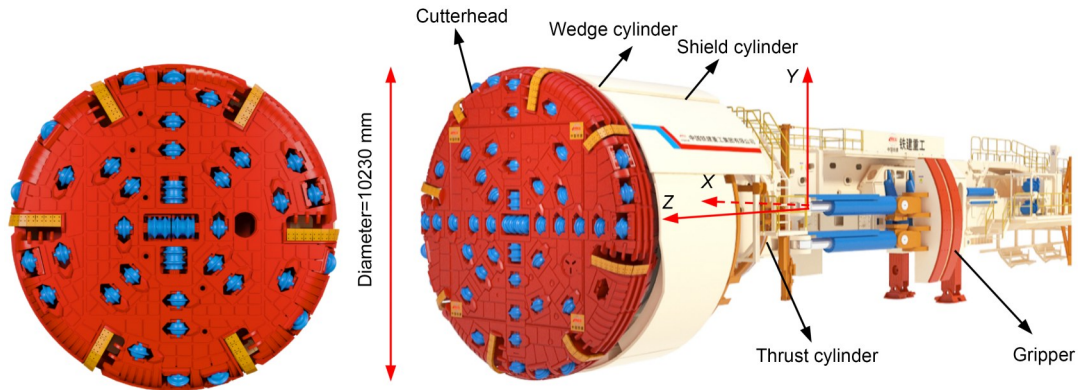


Fig. 2 TBM excavation system

Table 1 Principal TBM parameters

Parameter	Description
Machine stroke (mm)	21800
Maximum advance rate (mm/min)	100
Cutterhead drive rated torque (kN·m)	20000
Total thrust force (kN)	34000
Cutterhead type	Flat face
Number of cutters	66
Rock breaking capacity (MPa)	250

analysis, four cutters from different positions—C29 and C54 from the face region and C60 and C65 from the gauge region—were ultimately selected as the study objects.

During the TBM excavation process, the disc cutters achieve rock fragmentation by inducing fractures through the application of high compressive and shear stresses as the cutterhead rotates and thrusts forward. This process continuously subjects the cutters to severe impact loads, intense friction, and elevated temperatures. Consequently, disc cutters are the most critical and vulnerable components of the TBM and are highly susceptible to wear-induced failure.

The failure modes of these cutters encompass both normal wear and various forms of abnormal damage, such as eccentric wear, ring fracturing, and chipping. Analyzing on-site replacement logs reveals that normal wear is the predominant failure mode, accounting for over 85% of failure cases. Therefore, to address the most frequent cause of cutter replacement, we focus specifically on modeling and predicting the progression of the normal wear process.

2.3 Dataset construction and preprocessing

The predictive model developed in this study is built upon two primary data sources: multivariate

time-series data acquired in real time from TBM sensors and disc cutter wear measurements from on-site manual recordings. To transform these raw, heterogeneous data into structured samples suitable for machine learning, a comprehensive data processing pipeline was developed. To enable real-time assessment of performance and establish a consistent evaluation benchmark, we introduce the instantaneous wear rate as the core predictive indicator. It is defined as follows:

$$v_{\text{wear}} = \frac{\Delta W}{\Delta t}, \quad (1)$$

where v_{wear} , ΔW , and Δt represent the wear rate, change in wear value, and corresponding time for the cutter, respectively. Compared to cumulative wear, this indicator can more sensitively reflect the impact of current working conditions on wear, and it eliminates interference from differences in individual excavation cycle durations.

To comprehensively capture the key factors influencing wear, we selected 11 input features encompassing three categories: direct operational parameters such as the cutterhead water spray flow rate; machine state parameters such as the TBM attitude; comprehensive geological indices, namely, the force-penetration index (FPI) and the torque-penetration index (TPI). The four output targets are the wear rates of the individual cutters. A detailed description of all 11 input features and four output targets is provided in Eqs. (S1)–(S3) and Table S1 of the electronic supplementary materials (ESM).

Next, to construct sample units, excavation cycle segmentation and reconstruction were applied to the raw time-series data. Initial segmentation of the data

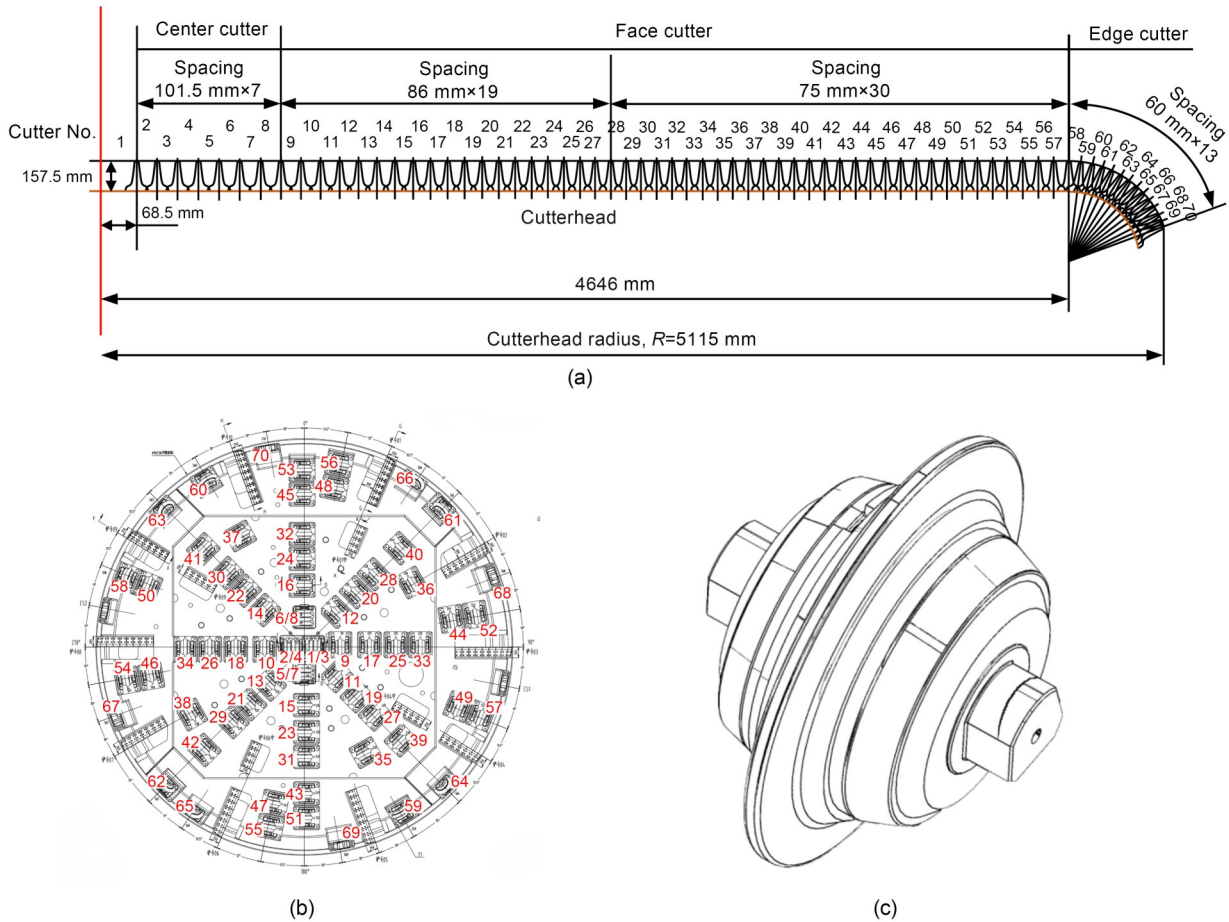


Fig. 3 TBM cutterhead and disc cutters: (a) radial distribution of the cutters; (b) cutter arrangement on the cutterhead face; (c) 19-inch single-disc cutter

stream involved identifying all continuous intervals where the thrust speed (v) exceeded zero. To address the issue of these segments being erroneously fragmented by data acquisition dropouts or brief operational pauses, a merging strategy based on segment duration was designed. A statistical analysis of these segment durations revealed a distinct bimodal distribution, as illustrated in Fig. 4. Based on this, a threshold of 300 s was established to identify and merge any pseudo-interrupted fragments shorter than this value. This reconstruction method yielded complete excavation cycles that demonstrated high consistency with on-site construction logs. Furthermore, to mitigate the impact of noise and anomalies on model learning, the data were cleaned by first removing outliers using the interquartile range (IQR) method and then performing smoothing and denoising using a combination of locally weighted scatterplot smoothing (LOWESS) and Savitzky–Golay filters.

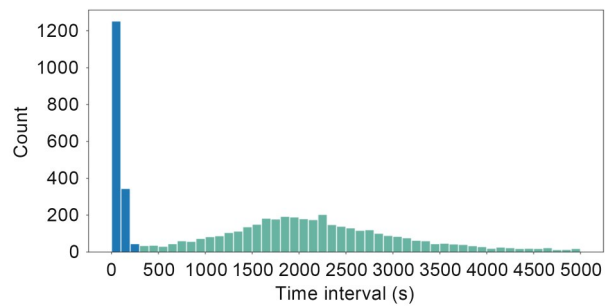


Fig. 4 Statistics of the time intervals

The data preparation pipeline culminates in a final dataset where the fundamental unit is a single excavation cycle. Each cycle is represented by 11 multivariate time-series variables that serve as the input features. Within this collection, cycles accompanied by wear measurements for the four target disc cutters are designated labeled samples, while all the others are unlabeled samples. This process yielded a dataset

of 3748 total cycles, comprised of 596 labeled and 3152 unlabeled samples, which serves as the foundation for the subsequent semi-supervised learning experiments.

3 Methods

This section presents the technical details of the proposed PL-HLANet framework, which is designed to address the challenges of TBM cutter wear prediction under data scarcity. An end-to-end flowchart of the methodology is provided in Fig. 5 to offer a high-level overview. The subsequent sections are structured in a top-down manner: we first introduce the architecture of the PL-HLANet framework, and then provide an explanation of its two primary components—the decoupled multi-view PL module and the hierarchical hybrid attention module.

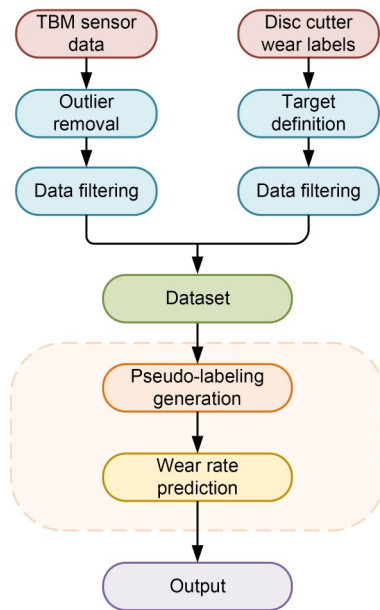


Fig. 5 Flowchart of the proposed method

3.1 Overall architecture

The core of our methodology is PL-HLANet, a decoupled, two-stage framework that separates the semi-supervised learning task from the final prediction task. Its detailed architecture is illustrated in Fig. 6, which shows its two synergistic components: (a) a decoupled multi-view PL module and (b) a hierarchical hybrid attention module (HLANet).

The first stage, the PL module, is designed to mine high-confidence supervisory signals from the massive pool of unlabeled data. It leverages a multi-view consistency checking mechanism to ensure the quality and reliability of the generated pseudo-labels, effectively augmenting the training set.

The second stage is the core predictive engine, HLANet, which is trained on the combination of original and pseudo-labeled data. Its internal architecture adheres to a “local-global-focus” philosophy, integrating a TCN, a Bi-LSTM, and a self-attention mechanism to accurately model the wear dynamics.

This decoupled design separates the semi-supervised data augmentation task from the final wear rate prediction task. Such an approach mitigates the risk of error accumulation inherent in traditional SSL methods and allows each specialized module to work in synergy, enabling more robust and accurate predictions.

3.2 Decoupled multi-view PL module

To address the challenge of acquiring sufficient ground-truth wear labels in industrial settings, we propose a decoupled, multi-view PL module. The central principle of this module is to separate the label generation process from the training of the downstream predictive model, thereby mitigating the risk of error accumulation inherent to conventional semi-supervised learning paradigms. As illustrated in Fig. 6a, the module operates by constructing and fusing information from multiple independent “views” to generate a high-quality set of pseudo-labels for the unlabeled data.

3.2.1 Multi-view construction

The effectiveness of multi-view learning relies on the complementarity among views to ensure learning sufficiency. To maximize the diversity between views, we employ a dual-level differentiation design at both the feature and model levels. At the feature level, we extracted three distinct feature subsets for each sample (detailed in Table 2), focusing on describing the signal’s energy and amplitude, impact characteristics, and statistical distribution shape. At the model level, each feature subset was then matched with a methodologically distinct regression model—Bayesian linear regression (BLR), SVR, and k -nearest neighbors (KNN) regression—to serve as a base learner for comprehensive and complementary feature learning.

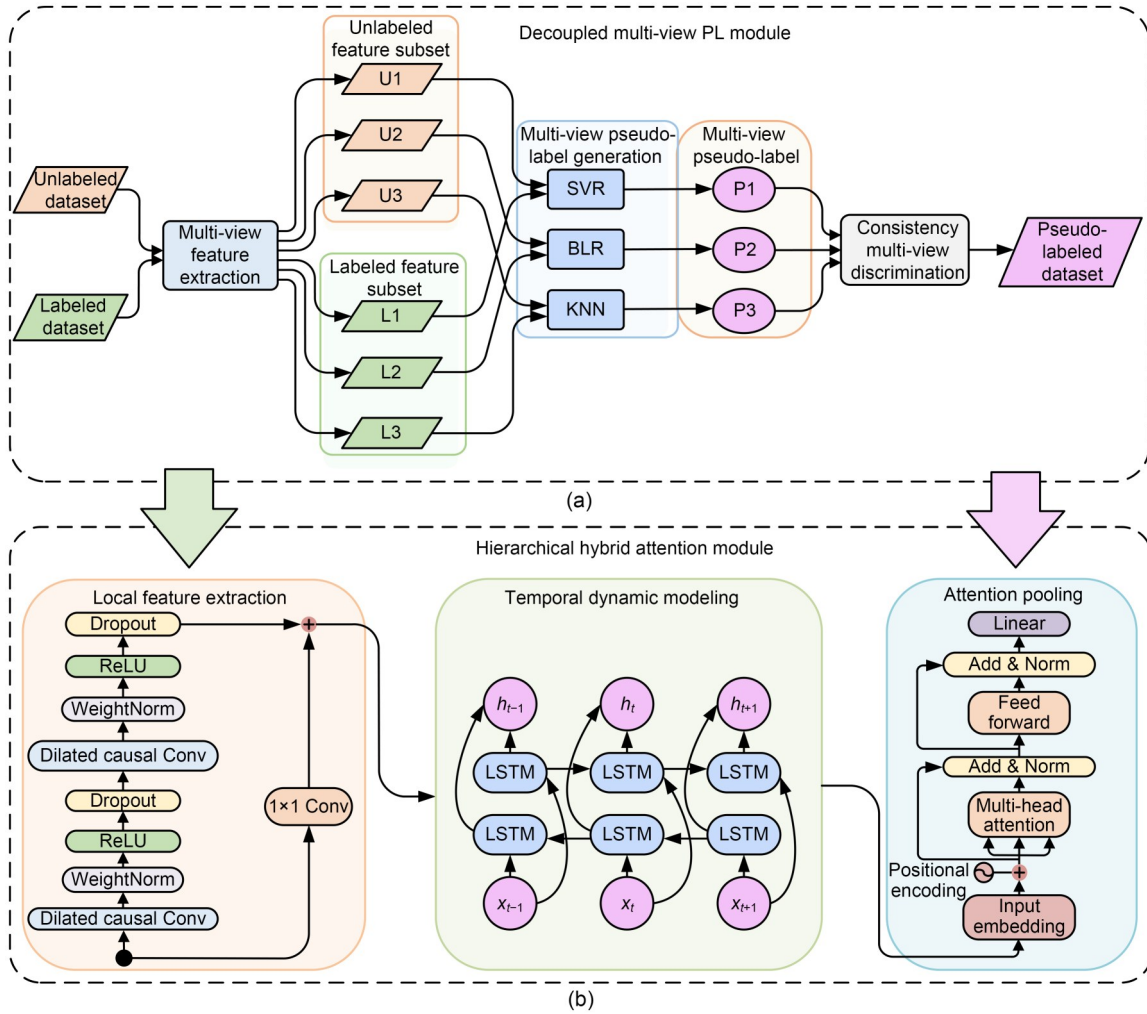


Fig. 6 Architecture of PL-HLNet: (a) decoupled multi-view PL module; (b) hierarchical hybrid attention module. ReLU: rectified linear unit; Conv: convolutional; WeightNorm: weight normalization; Add & Norm: addition and normalization; h_t : hidden state at time t ; x_t : input at time t

Table 2 Composition of feature subsets for multi-view learning

Feature subset	Feature
Subset 1	Mean, root mean square (RMS), root amplitude, and mean square
Subset 2	Mean, margin factor, crest factor, and impulse factor
Subset 3	Mean, standard deviation, skewness, and kurtosis

3.2.2 Sample reweighting for data imbalance

When training base learners on a small set of labeled data, a common challenge is an imbalanced distribution of label values. To address this, we designed an adaptive sample re-weighting strategy based on

kernel density estimation (KDE). The idea is to assign a higher training weight to samples located in low-density regions of the label distribution. To achieve this, we first estimate the probability density function $\hat{f}(y)$ of the observed labels using KDE, as defined by:

$$\hat{f}(y) = \frac{1}{n} \sum_{i=1}^n K_h(y - y_i) = \frac{1}{nh} \sum_{i=1}^n K\left(\frac{y - y_i}{h}\right), \quad (2)$$

where $\hat{f}(y)$ is the estimated probability density at label value y , n is the number of labeled training samples, y_i denotes the label value of the i th labeled sample, $K(\cdot)$ is the kernel function, and h is the bandwidth parameter that controls the smoothness of the density estimate.

Subsequently, a weighting function that is inversely proportional to the estimated density, $w(y)$, is employed to balance the model's focus between high- and low-density regions of the label space and thus mitigate prediction bias. This weight is then integrated into the loss function of the base learners to guide the optimization process. The final weighted loss L_w is calculated as follows:

$$w(y) = e^{-0.5(\hat{f}(y)-1)}, \quad (3)$$

$$L_w = \frac{1}{n} \sum_{i=1}^l (y_i - \hat{y}_i) w(y), \quad (4)$$

where y_i is the true label value, and \hat{y}_i is the predicted value.

The implementation of this strategy forms a complete logical chain, from quantifying label rarity to guiding model training. First, we employ KDE, as defined in Eq. (5), to quantify the rarity of each known label value y_i by calculating its probability density estimate, $\hat{f}(y_i)$. A high $\hat{f}(y_i)$ value signifies a common label in a dense region, while a low value indicates a rare one. This estimated density is then transformed into a training weight using our designed exponential function (Eq. (6)), which is inversely proportional to the density. This function amplifies the importance of rare samples (yielding a weight >1) and diminishes the influence of common samples (yielding a weight <1). Finally, this calculated weight is integrated into the training process via a weighted loss function (Eq. (7)), where each sample's squared error is multiplied by its corresponding weight. This mechanism compels the optimizer to prioritize the reduction of errors for these magnified rare samples, thereby promoting generalization across the entire, imbalanced label distribution.

3.2.3 Consistency-based filtering and label fusion

After the three base learners are trained, pseudo-labels can be generated and filtered for the unlabeled dataset. First, each unlabeled sample is input to the three independent base learners to generate three candidate pseudo-labels, $\{y_1^*, y_2^*, y_3^*\}$. Subsequently, this set of pseudo-labels is tested based on the consistency assumption of multi-view learning. We screen for consistent samples by comparing their average pairwise deviation against a threshold ε . Based on a parameter sensitivity analysis (later detailed in Section 5), we select an optimal threshold of $\varepsilon=0.15$ for our experiments.

$$\frac{|y_1^* - y_2^*| + |y_1^* - y_3^*| + |y_2^* - y_3^*|}{3} \leq \varepsilon. \quad (5)$$

Finally, for the samples that pass the consistency check, we adopt an ensemble learning strategy to fuse their pseudo-labels. The mean of the predicted values from the three views is used as the final, high-confidence pseudo-label for the sample; it is added to the pseudo-label dataset for training the downstream wear rate prediction model.

3.3 Hierarchical hybrid attention module

This module is the core engine for predicting cutter wear rates. Its internal architecture follows a hierarchical "local-global-focus" design philosophy, being divided into a local feature extraction layer, a temporal dynamic modeling layer, and an attention pooling layer. The goal of this design is to aid deep feature mining on high-dimensional TBM time-series data, so as to ultimately achieve accurate and robust prediction of cutter wear rates. The architecture details are illustrated in Fig. 6b.

3.3.1 Local feature extraction layer

The local feature extraction layer is designed to accurately capture key transient changes in the raw signal, such as impacts and vibrations. For this purpose, we employ a TCN, which is constructed from three stacked residual blocks, each utilizing a kernel size of 7. The effectiveness of the TCN stems from two core principles. First, its causal convolutions ensure that an output at a given time step depends only on past and present inputs, thus preserving temporal causality. Second, its dilated convolutions allow the network to achieve an exponentially large receptive field with minimal computational overhead. This operation is formally defined as:

$$y_s = \sum_{i=0}^{k-1} f(i) \cdot x_{s-d \cdot i}, \quad (6)$$

where y_s is the output at time step s , k is the kernel size, d is the dilation factor, $f(i)$ denotes the learnable kernel weight at position i , and $x_{s-d \cdot i}$ is the input at the corresponding historical time step.

3.3.2 Temporal dynamic modeling layer

The temporal dynamic modeling layer is designed to model the long-term, cumulative dynamics of the

wear process by processing the feature sequence generated by the TCN. A key aspect of this module's design is the nature of the connection between the TCN and the Bi-LSTM. Our TCN module utilizes causal convolutions with padding but intentionally omits any pooling layers or strided convolutions. The rationale behind this design is to ensure that the output feature sequence has the same temporal length as the input. Therefore, the Bi-LSTM module processes the complete feature sequence from the TCN. The data flow between layers and the corresponding tensor shape transformations are depicted in Fig. 7.

With this high-level feature sequence as input, we employ a single-layer Bi-LSTM network with a hidden dimension of 128 to model the evolutionary dynamics over time. Because of its unique gating mechanism and bidirectional structure, the Bi-LSTM is adept at capturing long-term dependencies while integrating contextual information from both the past and future, enabling a more comprehensive understanding of the cumulative dynamics of cutter wear throughout excavation. The Bi-LSTM combines the results from two LSTM modules that process the data in opposite directions, as shown in Eqs. (7) and (8).

$$h'_t, c'_t = \text{LSTM}(x_t, h'_{t-1}, c'_{t-1}), \quad (7)$$

$$h_t = \text{concatenate}[\vec{h}'_t; \vec{h}'_t], \quad (8)$$

where h'_t and c'_t represent the output and state of the LSTM, respectively, and h_t is the result of the Bi-LSTM combining the outputs from the two directions.

3.3.3 Attention pooling layer

Building upon the local and global understanding established by the preceding layers, the attention pooling layer filters and aggregates the most critical information from the entire time series to prepare for the prediction task. Recognizing that features at different time steps have varying importance, we introduce a self-attention pooling mechanism to this layer. This mechanism learns to assign a dynamic weight to each hidden state vector that is output by the Bi-LSTM, enabling the model to focus on the most discriminative feature segments. Ultimately, it aggregates the entire variable-length feature sequence into a fixed-dimension context vector. This concentrates the key information required for the decision and endows the model with a higher degree of interpretability. The working principle of this mechanism is shown in Eqs. (9)–(11):

$$e_t = \mathbf{h}_t^\top \mathbf{u}, \quad (9)$$

$$\alpha_t = \frac{\exp(e_t)}{\sum_{j=1}^T \exp(e_j)}, \quad (10)$$

$$\mathbf{c} = \sum_{t=1}^T \alpha_t \mathbf{h}_t, \quad (11)$$

where t is the index of the current time step, j is a dummy summation index over all time steps, T is the total sequence length, h_t is the hidden state at time step t , u is a learnable, globally shared context query vector, e_t is the raw attention score, α_t is the attention weight, and c is the final feature representation that condenses all time steps.

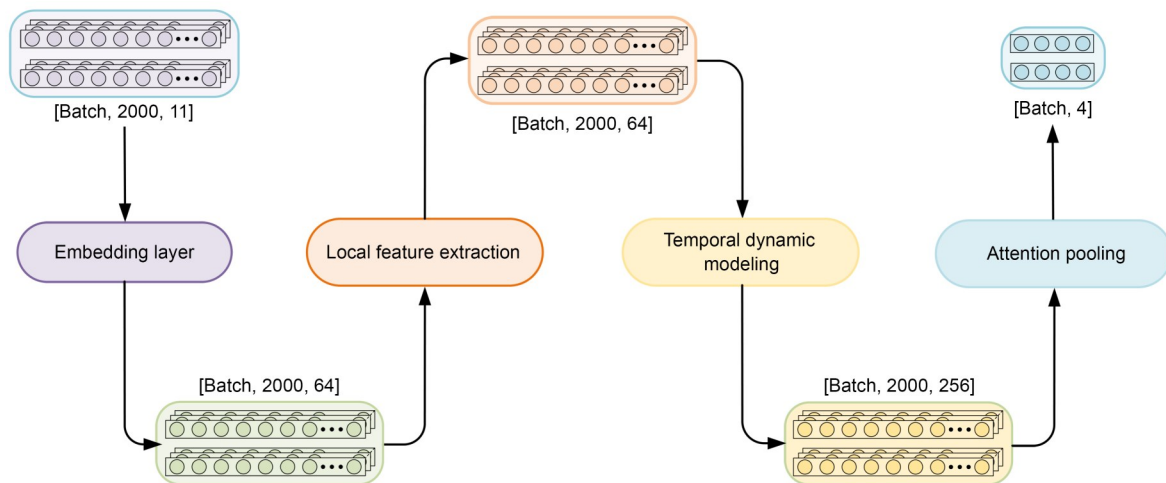


Fig. 7 Data flow in the hierarchical hybrid attention module

3.4 Training details

For experimental validation, the labeled dataset was partitioned, with 20% of its samples reserved as a test set to evaluate the final performance of the PL-HLNet model. The remaining 80% of the labeled samples, along with the entire unlabeled dataset, constituted the training set for our semi-supervised framework. The model was configured with an input dimension of 11, an output dimension of 4, and an embedding dimension of 64. The TCN module consisted of three layers with 64 channels, each using a kernel size of 7, while the Bi-LSTM layer had a hidden dimension of 128. A dropout rate of 0.2 was applied throughout the network. For training, the model was optimized using the Adam optimizer with an initial learning rate of 0.001. The model was trained by minimizing the mean squared error (MSE) between the predicted and actual wear rates. We set the maximum number of epochs to 300 and employed an early stopping mechanism with a patience of 5, monitoring the validation loss to prevent overfitting.

4 Results and analysis

4.1 Prediction performance

The performance of the proposed PL-HLNet was quantitatively evaluated on the test set, with detailed metrics presented in Table 3. The results show that our model achieves excellent prediction accuracy, with coefficient of determination (R^2) values for all four cutters exceeding 0.80 and the best-performing mean absolute error (MAE) and root mean square error (RMSE) values reaching 0.0243 and 0.0348, respectively. The strong agreement between the predicted and actual wear rates over time is further illustrated in the trend plots of Figs. 8a–8d. These results validate the effectiveness of the PL-HLNet framework for high-precision wear prediction, even with limited labeled data.

Furthermore, a discernible performance variation exists among the different cutters, with the gauge cutters (C60, C65) exhibiting higher prediction accuracy than the center cutters (C29, C54). This discrepancy can be attributed to their distinct operational characteristics. As center cutters, C29 and C54 have lower linear velocities and experience slower wear progression. In the harsh, hard-rock tunneling environment, this results in a lower signal-to-noise ratio (SNR) in the collected data, making it more challenging for the model to capture the underlying wear patterns. Conversely, the peripheral gauge cutters exhibit a faster wear rate, yielding a more pronounced and easily learnable pattern. The model's strong performance on even the low-SNR center cutters highlights its robustness.

To demonstrate the model's practical value for predictive maintenance, the predicted cumulative wear—obtained by integrating the instantaneous wear rate over time—was compared against the actual wear progression. As shown in Figs. 8e–8h, the predicted cumulative wear follows the actual measurements with high fidelity, even after 30 h of operation. This demonstrates the model's potential to provide reliable, long-term guidance for on-site personnel in scheduling optimal cutter replacement, thereby reducing downtime and operational costs. Its practicality is also demonstrated by its excellent computational efficiency; the model achieves a single-sample latency of just 7.37 ms and a batch throughput of nearly 6000 samples per second, making it fully viable for on-site deployment.

4.2 Performance comparisons

To evaluate the performance of the proposed PL-HLNet framework, we selected a diverse set of strong baselines for comparison. To ensure fairness and reproducibility, their key hyperparameters are detailed in Table S2 of the ESM. These models fall into two main categories: (1) purely supervised models, including SVR, LSTM, GRU, and a Transformer encoder; (2) mainstream semi-supervised learning frameworks, namely, Mean Teacher and FixMatch.

Table 3 Detailed performance metrics of the PL-HLNet model on the test set

Model	Cutter No.	R^2	MAE	RMSE	Latency (batch size=1) (ms)	Throughput (s^{-1})
PL-HLNet	29	0.8064	0.0305	0.0427	7.37	5987.15
PL-HLNet	54	0.8642	0.0339	0.0490	7.37	5987.15
PL-HLNet	60	0.9531	0.0342	0.0494	7.37	5987.15
PL-HLNet	65	0.9538	0.0243	0.0348	7.37	5987.15

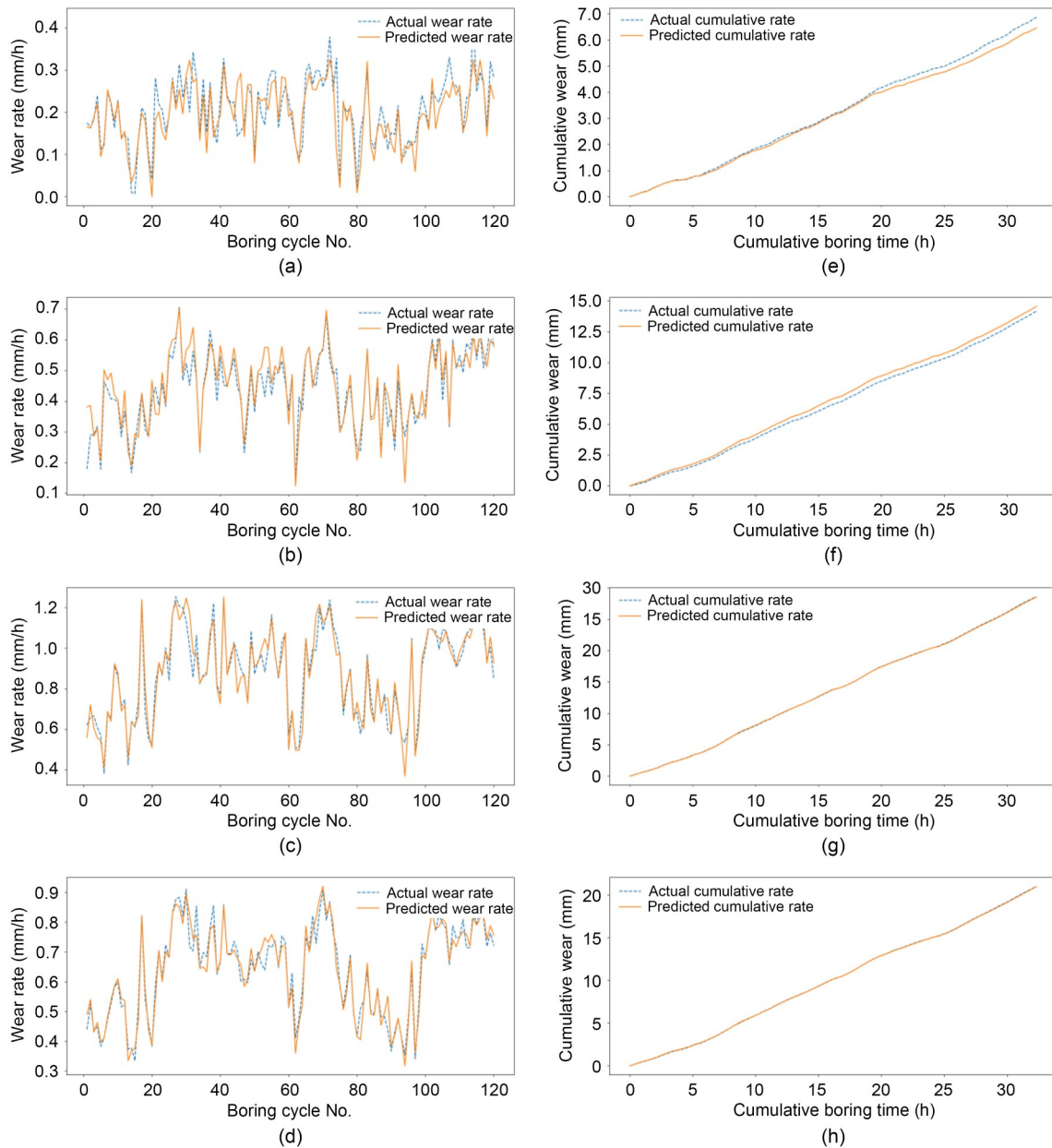


Fig. 8 Prediction performance for representative disc cutters: (a)–(d) instantaneous wear rate comparisons for cutters C29, C54, C60, and C65, respectively; (e)–(h) corresponding comparisons for cumulative wear

As indicated by the quantitative results in Table 4 and the visual evidence in Figs. 9 and 10, a clear performance hierarchy emerges. First, both the Mean Teacher and FixMatch significantly outperform the purely supervised baselines, confirming the value of leveraging unlabeled data. However, the proposed PL-HLNet consistently surpasses all baseline models, including the strong SSL competitors. Its predictions demonstrate the closest fit to the actual values in terms of both trend and magnitude.

An in-depth analysis reveals that the performance advantage of PL-HLNet is particularly apparent for the more challenging central cutters (especially C29), which have a lower SNR. The data in Table 4 show that while most baseline models struggle on the C29 cutter, our framework maintains high accuracy. For instance, while the strong FixMatch baseline improves the R^2 on C29 to a respectable 0.5914, our PL-HLNet elevates it further to 0.8064. This significant performance improvement on difficult samples demonstrates

Table 4 Predictive performance of different models

Model	Cutter No.	R^2	MAE	RMSE
LSTM	29	0.4919	0.0504	0.0631
	54	0.6695	0.0699	0.0884
	60	0.7587	0.1039	0.1394
	65	0.8107	0.0629	0.0875
GRU	29	0.4873	0.0511	0.0634
	54	0.6537	0.0719	0.0905
	60	0.7588	0.1035	0.1393
	65	0.8098	0.0626	0.0877
Transformer encoder	29	0.4541	0.0525	0.0654
	54	0.6579	0.0704	0.0900
	60	0.7683	0.1066	0.1366
	65	0.8238	0.0612	0.0844
SVR	29	0.1902	0.0639	0.0797
	54	0.5214	0.0801	0.1064
	60	0.6290	0.1270	0.1728
	65	0.6511	0.0843	0.1188
Mean Teacher	29	0.5587	0.0278	0.0361
	54	0.7026	0.0685	0.0893
	60	0.7579	0.1179	0.1549
	65	0.8216	0.0883	0.1158
FixMatch	29	0.5914	0.0317	0.0398
	54	0.7446	0.0572	0.0722
	60	0.7920	0.1086	0.1408
	65	0.8533	0.0820	0.1026
PL-HLNet	29	0.8064	0.0305	0.0427
	54	0.8642	0.0339	0.0490
	60	0.9531	0.0342	0.0494
	65	0.9538	0.0243	0.0348

that the strength of our model lies not only in its overall prediction accuracy but also, more importantly, in its robustness and generalizability under complex, noisy operating conditions. This is attributed to our framework's unique design: the multi-view PL module provides higher-quality supervisory signals than standard consistency-based methods; meanwhile, the hierarchical attention network more effectively mines weak wear signals from measurements with strong interference.

5 Ablation studies and discussion

In this section, we validate the effectiveness of the proposed PL-HLNet framework by analyzing the relative contributions of its components. We present a primary ablation study that quantifies the significant contribution of the semi-supervised PL module by comparing the full model's performance against its purely supervised counterpart.

For the sake of brevity in the main text, further in-depth validation studies—including a sensitivity analysis of the consistency threshold (ϵ) and an ablation study on the Bi-LSTM network depth—are provided in Sections S5 and S6 of the ESM. Collectively, these analyses provide data-driven evidence

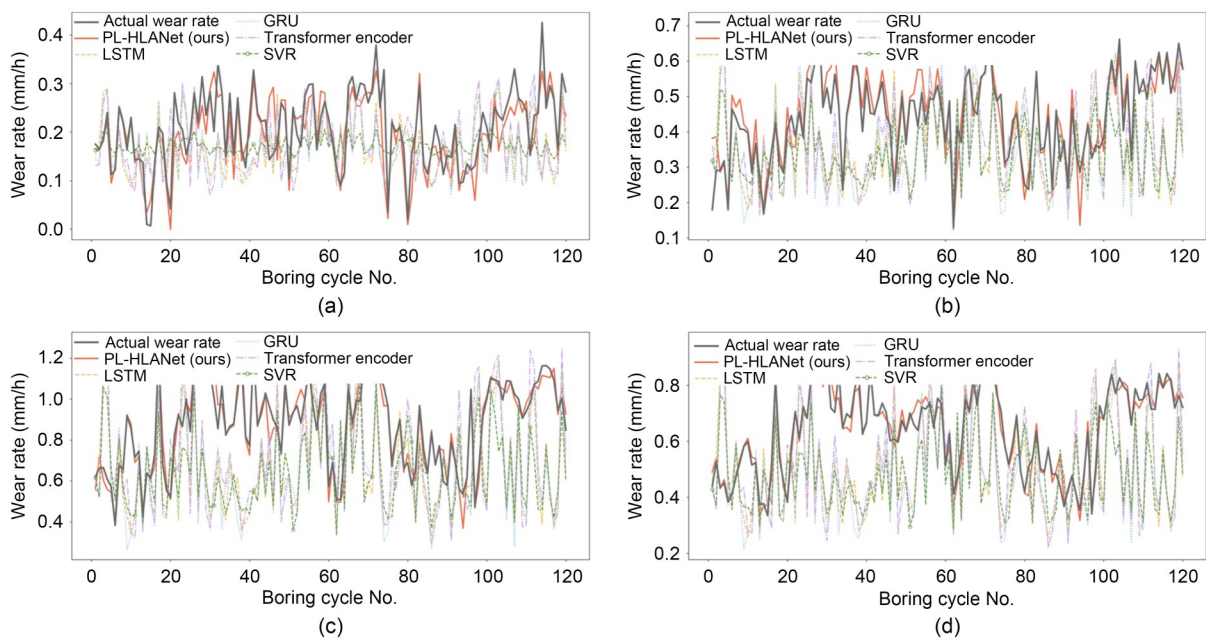


Fig. 9 Comparison of wear rate prediction performance across different models. Line plots compare the predicted instantaneous wear rates of each model with the actual values for four representative cutters: (a) C29; (b) C54; (c) C60; (d) C65

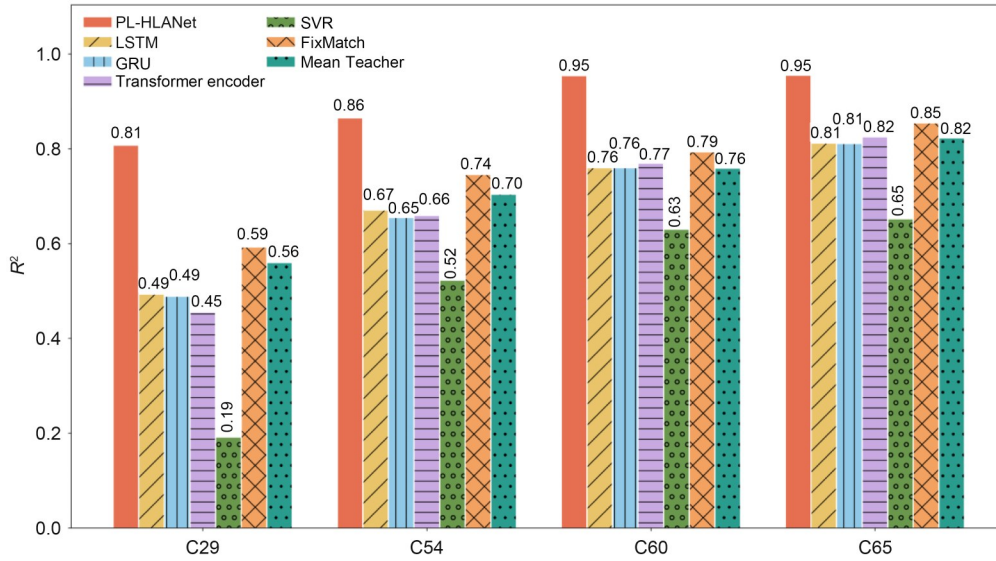


Fig. 10 Overall R^2 performance comparison across the different models. The bar chart displays the average R^2 score for each model, calculated across all four test cutters

for the framework’s effectiveness and its key design choices.

5.1 Ablation study on the PL module

To validate the effectiveness of the proposed decoupled PL framework, a core ablation study was conducted. We compared the performance of the complete PL-HLNet framework against that of a baseline HLNet model trained in a purely supervised manner. Specifically, the baseline model used an identical HLNet architecture but was only trained on the 596 labeled samples and not any unlabeled data. The performance on the test set is presented in Table 5.

Table 5 Ablation study results for the PL module

Model	Cutter No.	R^2	MAE	RMSE
HLNet	29	0.2653	0.0559	0.0700
	54	0.2288	0.0917	0.1168
	60	0.5109	0.1228	0.1596
	65	0.4942	0.0912	0.1150
PL-HLNet	29	0.8064	0.0305	0.0427
	54	0.8642	0.0339	0.0490
	60	0.9531	0.0342	0.0494
	65	0.9538	0.0243	0.0348

The metrics in Table 5 demonstrate that the model’s performance degrades sharply on each individual cutter without the supplementary supervisory information provided by the PL module. Specifically, by using the PL framework, for the more challenging

cutters C29 and C54, the R^2 scores dramatically improved from 0.2653 and 0.2288 to 0.8064 and 0.8642, respectively. Similarly, cutters C60 and C65 saw substantial gains, with their R^2 scores increasing from approximately 0.50 to 0.95. This numerically demonstrates the immense contribution of the PL module to the model’s performance.

Visual evidence for this performance gap is provided in Fig. 11, which illustrates the comparison on the most challenging central cutter, C29. As shown in the line plot (Fig. 11a), the supervised-only model does not track the curve of actual values as well as the full PL-HLNet framework. This performance disparity is even clearer in the scatter plot (Fig. 11b), where the predicted points from the supervised model are significantly more dispersed and exhibit systematic overestimation. The results for all the other cutters, which exhibit a similar pattern, are detailed in Fig. S1 of the ESM.

Overall, this ablation study showcases the effectiveness of the proposed PL module. By leveraging a large volume of unlabeled data, the module successfully corrects the distributional bias of the small training set and acts as a form of data augmentation for regularization, thereby significantly enhancing the model’s accuracy and generalizability.

5.2 Ablation study on the HLNet module

To validate the necessity and contributions of each component within the proposed HLNet, this

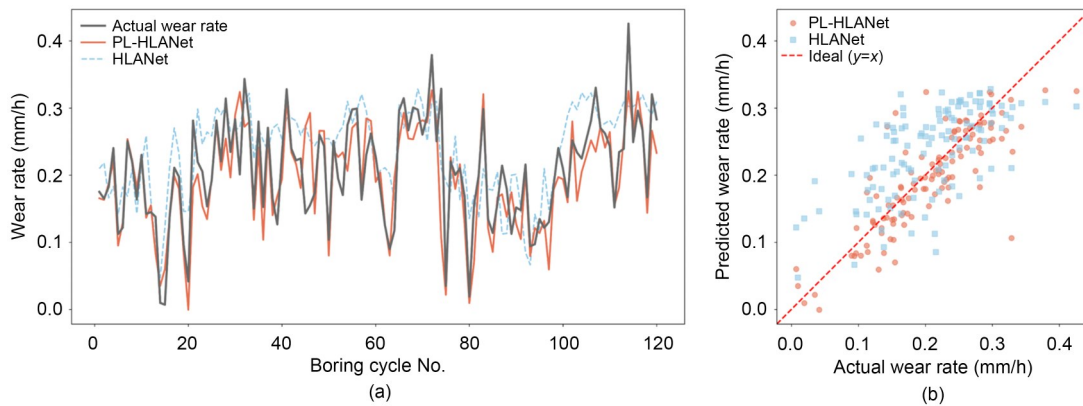


Fig. 11 Ablation study on the PL module for C29: (a) instantaneous wear rates; (b) cumulative wear rates

section presents a series of targeted ablation studies. We benchmarked the complete PL-HLANet against several ablated versions, where the TCN, Bi-LSTM, and attention modules were individually removed. All models were trained on the same dataset augmented with pseudo-labels, with the results presented in Table 6.

Table 6 Ablation study results for the HLANet module

Model	Cutter No.	R^2	MAE	RMSE
PL-HLANet (w/o TCN)	29	0.5326	0.0434	0.0558
	54	0.7379	0.0529	0.0681
	60	0.8698	0.0608	0.0824
	65	0.9111	0.0356	0.0482
PL-HLANet (w/o Bi-LSTM)	29	0.4091	0.0487	0.0628
	54	0.5911	0.0676	0.0851
	60	0.7438	0.0903	0.1155
	65	0.8101	0.0545	0.0705
PL-HLANet (w/o attention)	29	0.5922	0.0406	0.0522
	54	0.7616	0.0487	0.0650
	60	0.8990	0.0560	0.0726
	65	0.9134	0.0362	0.0476
PL-HLANet	29	0.8064	0.0305	0.0427
	54	0.8642	0.0339	0.0490
	60	0.9531	0.0342	0.0494
	65	0.9538	0.0243	0.0348

w/o indicates “without”

As observed from the average performance metrics in Table 6, the complete PL-HLANet makes the most accurate predictions. The removal of any core component from the network leads to reduced performance across all metrics. Fig. 12 provides further evidence for this conclusion by illustrating the ablation results on the most challenging central cutter, C29.

The plots for cutter C29 reveal the synergistic mechanism of the components. Removing the front-end TCN module leads to a significant decline in performance, indicating that directly feeding noisy, high-frequency raw signals into the Bi-LSTM is less effective than first extracting local, structured features. The contribution of the Bi-LSTM module is even more critical, as its removal causes the most severe performance degradation; the R^2 score on the C29 cutter decreases significantly from 0.8064 to 0.4091. This proves that it is essential to model the long-term time dependencies of such features. Finally, the ablation of the attention mechanism also validates the positive role of focusing on critical moments to optimize the final predictions. While the analysis here focused on C29, similar trends were observed across all cutters, with the complete results provided in Fig. S2 of the ESM.

In summary, the ablation studies systematically validate the rationale and necessity of each component in the HLANet architecture. The TCN serves as a local feature extractor, the LSTM as a temporal dynamic modeler, and the attention mechanism as an information aggregator, together forming a hierarchical information processing pipeline that progresses from local to global to key points. They work synergistically to yield strong performance on the disc cutter wear prediction task.

6 Conclusions

To address the dual challenge of labeled data scarcity and complex operating conditions in TBM disc cutter wear prediction for high-altitude, hard-rock

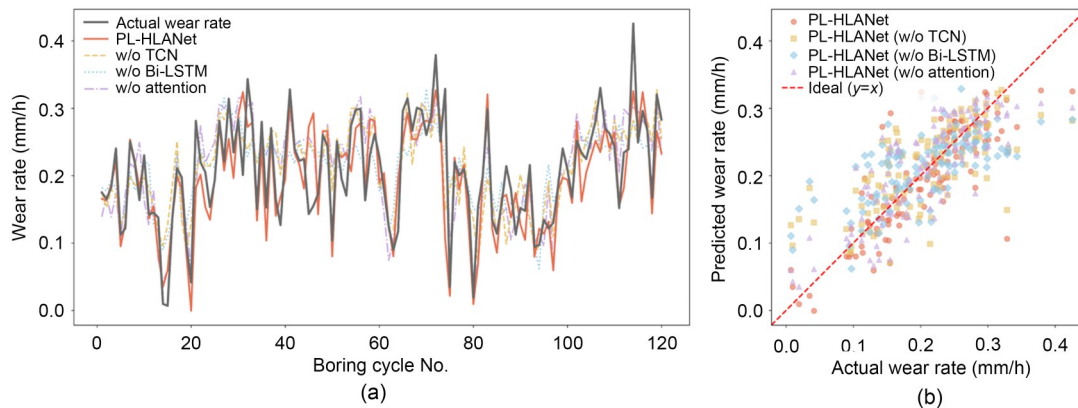


Fig. 12 Ablation study on the HLANet module for C29: (a) instantaneous wear rates; (b) cumulative wear rates

tunnels, we proposed a decoupled semi-supervised predictive framework termed PL-HLANet. The efficacy of the framework was demonstrated through systematic experiments and analysis on data from a real-world tunneling project.

The experimental results showed that the PL-HLANet consistently outperforms a comprehensive suite of baseline models, including classic supervised methods (SVR, LSTM) and advanced semi-supervised frameworks (Mean Teacher, FixMatch). Core ablation studies have shown that the decoupled PL module is a key driver of this performance enhancement, as it effectively corrects sampling bias and regularizes the model. In addition to its high accuracy, the model displays practicality with its high computational efficiency, exhibiting an average single-sample latency of just 7.37 ms.

Moreover, internal validation of the HLANet predictive module confirmed the superiority of its hierarchical TCN–Bi-LSTM–attention architecture. This superiority stems from fruitful design choices, such as the use of a single-layer Bi-LSTM and an empirically optimized consistency threshold, which were validated through the ablation studies and sensitivity analyses in Section 5. Additionally, the custom attention mechanism enables focusing on critical wear-related information, which leads to the model’s excellent robustness, particularly its ability to make accurate predictions on challenging, low-SNR samples (e.g., the C29 cutter).

In summary, we proposed and validated a predictive framework that organically combines a semi-supervised learning paradigm with a deep temporal network. Through a PL strategy, the framework effectively addresses the issue of labeled data scarcity in

industrial scenarios, while its built-in hybrid attention network demonstrates powerful feature learning capabilities. This work provides a reliable and effective approach for advancing intelligent TBM excavation and predictive maintenance.

Despite the promising results, this study has a few limitations that present avenues for future work. For instance, we focused exclusively on normal wear progression; future work could extend the model to identify and provide early warnings for abnormal failure modes, such as chipping and fracturing. In addition, more advanced self-supervised learning methods could be explored to further reduce the dependency on initially labeled data. Finally, deploying the model in a practical engineering context and enabling online learning for continuous optimization will be a key focus of our subsequent research.

Acknowledgments

This work is supported by the National Natural Science Foundation of China (No. 52305074) and the National Key Research and Development Program of China (No. 2021YFB3301603).

Author contributions

Wei TANG designed the research. Zhaoyang LI and Xinyuan WANG processed the corresponding data. Zhaoyang LI wrote the first draft of the manuscript. Jun ZOU and Huayong YANG helped to organize the manuscript. Zhaoyang LI, Huxiu XU, and Jun ZOU revised and edited the final version.

Conflict of interest

Huayong YANG is an Editor-in-Chief of this journal, and is NOT involved in the editorial review or the decision to publish this article. Zhaoyang LI, Wei TANG, Xinyuan WANG, Huxiu XU, Huayong YANG, and Jun ZOU declare that they have no conflict of interest.

References

- Agrawal AK, Murthy VMSR, Chattopadhyaya S, et al., 2022. Prediction of TBM disc cutter wear and penetration rate in tunneling through hard and abrasive rock using multi-layer shallow neural network and response surface methods. *Rock Mechanics and Rock Engineering*, 55(6):3489-3506.
<https://doi.org/10.1007/s00603-022-02834-7>
- Bai LP, Mo DY, Li HS, et al., 2025. Optimized data preprocessing and model selection for TBM cutter wear prediction. *Coatings*, 15(5):564.
<https://doi.org/10.3390/coatings15050564>
- Barzegari G, Khodayari J, Rostami J, 2021. Evaluation of TBM cutter wear in Naghadeh water conveyance tunnel and developing a new prediction model. *Rock Mechanics and Rock Engineering*, 54(12):6281-6297.
<https://doi.org/10.1007/s00603-021-02640-7>
- Chen BX, Jiang JG, Wang XM, et al., 2022. Debaised self-training for semi-supervised learning. Proceedings of the 36th International Conference on Neural Information Processing Systems, article 2349.
- Cheng Q, Chen YX, Xiao YT, et al., 2022. A dual-stage attention-based Bi-LSTM network for multivariate time series prediction. *The Journal of Supercomputing*, 78(14):16214-16235.
<https://doi.org/10.1007/s11227-022-04506-3>
- Ding XB, Xie YX, Xue HW, et al., 2022. A new approach for developing EPB-TBM disc cutter wear prediction equations in granite stratum using backpropagation neural network. *Tunnelling and Underground Space Technology*, 128:104654.
<https://doi.org/10.1016/j.tust.2022.104654>
- Dou YM, Li KW, Lv WJ, et al., 2024. ContrasInver: ultra-sparse label semi-supervised regression for multidimensional seismic inversion. *IEEE Transactions on Geoscience and Remote Sensing*, 62:5917613.
<https://doi.org/10.1109/TGRS.2024.3410022>
- Fan J, Zhang K, Huang Y, et al., 2023. Parallel spatio-temporal attention-based TCN for multivariate time series prediction. *Neural Computing and Applications*, 35(18):13109-13118.
<https://doi.org/10.1007/s00521-021-05958-z>
- Ghorbani E, Yagiz S, 2024. Predicting disc cutter wear using two optimized machine learning techniques. *Archives of Civil and Mechanical Engineering*, 24(2):106.
<https://doi.org/10.1007/s43452-024-00911-y>
- Hewage P, Behera A, Trovati M, et al., 2020. Temporal convolutional neural (TCN) network for an effective weather forecasting using time-series data from the local weather station. *Soft Computing*, 24(21):16453-16482.
<https://doi.org/10.1007/s00500-020-04954-0>
- Jia LX, Pu XB, Shang KD, et al., 2022. A deep learning combination model to predict TBM disc-cutter wear status. Proceedings of the IEEE 17th Conference on Industrial Electronics and Applications, p.1469-1474.
<https://doi.org/10.1109/ICIEA54703.2022.10006261>
- Karami M, Zare S, Rostami J, 2021. Introducing an empirical model for prediction of disc cutter life for TBM application in jointed rocks: case study, Kerman water conveyance tunnel. *Bulletin of Engineering Geology and the Environment*, 80(5):3853-3870.
<https://doi.org/10.1007/s10064-021-02166-w>
- Li L, Liu ZB, Zhou HY, et al., 2022. Prediction of TBM cutterhead speed and penetration rate for high-efficiency excavation of hard rock tunnel using CNN-LSTM model with construction big data. *Arabian Journal of Geosciences*, 15(3):280.
<https://doi.org/10.1007/s12517-022-09542-0>
- Lu DC, Liu YH, Kong FC, et al., 2024. A novel Bi-LSTM method fusing current and historical data for tunnelling parameters of shield tunnel. *Transportation Geotechnics*, 49:101402.
<https://doi.org/10.1016/j.trgeo.2024.101402>
- Mahmoodzadeh A, Mohammadi M, Hashim Ibrahim H, et al., 2021. Machine learning forecasting models of disc cutters life of tunnel boring machine. *Automation in Construction*, 128:103779.
<https://doi.org/10.1016/j.autcon.2021.103779>
- Mo DY, Bai LP, Huang WR, et al., 2024. TBM disc cutter wear prediction using stratal slicing and IPSO-LSTM in mixed weathered granite stratum. *Tunnelling and Underground Space Technology*, 148:105745.
<https://doi.org/10.1016/j.tust.2024.105745>
- Pu XB, Jia LX, Shang KD, et al., 2022. A new strategy for disc cutter wear status perception using vibration detection and machine learning. *Sensors*, 22(17):6686.
<https://doi.org/10.3390/s22176686>
- Shen X, Chen XS, Fu YB, et al., 2022. Prediction and analysis of slurry shield TBM disc cutter wear and its application in cutter change time. *Wear*, 498-499:204314.
<https://doi.org/10.1016/j.wear.2022.204314>
- Shi G, Qin CJ, Zhang ZN, et al., 2024. A novel decomposition and hybrid transfer learning-based method for multi-step cutterhead torque prediction of shield machine. *Mechanical Systems and Signal Processing*, 214:111362.
<https://doi.org/10.1016/j.ymssp.2024.111362>
- Shin YJ, Kwon K, Bae A, et al., 2024. Machine learning-based prediction model for disc cutter life in TBM excavation through hard rock formations. *Tunnelling and Underground Space Technology*, 150:105826.
<https://doi.org/10.1016/j.tust.2024.105826>
- Sun JD, Shang Y, Wang K, et al., 2023. A new prediction model for disc cutter wear based on Cerchar Abrasivity Index. *Wear*, 526-527:204927.
<https://doi.org/10.1016/j.wear.2023.204927>
- Wang HD, Qin CJ, Yu HG, et al., 2025. A real-time multi-head mixed attention mechanism-based prediction method for tunnel boring machine disc cutter wear. *Science China Technological Sciences*, 68(1):1120302.
<https://doi.org/10.1007/s11431-024-2794-6>
- Wang KY, Wu XG, Zhang LM, et al., 2023. Data-driven multi-step robust prediction of TBM attitude using a hybrid deep learning approach. *Advanced Engineering Informatics*, 55:101854.
<https://doi.org/10.1016/j.aei.2022.101854>

- Wen XY, Li WB, 2023. Time series prediction based on LSTM-attention-LSTM model. *IEEE Access*, 11:48322-48331.
<https://doi.org/10.1109/ACCESS.2023.3276628>
- Wu YH, Wu GQ, Lin JX, et al., 2025. Role exchange-based self-training semi-supervision framework for complex medical image segmentation. *IEEE Transactions on Neural Networks and Learning Systems*, 36(5):8372-8386.
<https://doi.org/10.1109/TNNLS.2024.3432877>
- Yu FX, Wang D, Chen YP, et al., 2021. Unsupervised domain adaptation for object detection via cross-domain semi-supervised learning. arXiv:1911.07158.
<https://doi.org/10.48550/arXiv.1911.07158>
- Yu HG, Tao JF, Huang S, et al., 2021. A field parameters-based method for real-time wear estimation of disc cutter on TBM cutterhead. *Automation in Construction*, 124: 103603.
<https://doi.org/10.1016/j.autcon.2021.103603>
- Zhang N, Shen SL, Zhou AN, 2023. A new index for cutter life evaluation and ensemble model for prediction of cutter wear. *Tunnelling and Underground Space Technology*, 131:104830.
<https://doi.org/10.1016/j.tust.2022.104830>

Electronic supplementary materials

Sections S1–S6, Tables S1–S5, Figs. S1–S3, Eqs. (S1)–(S3)