

Uncertainty-aware deep neural network for multi-scale temporal modeling in industrial equipment prognostics



Bailing Zhang

School of Computer Science and Data Engineering, NingboTech University, Ningbo 315100, China; bailing.zhang@nit.zju.edu.cn

Highlights:

- A novel Temporal Probabilistic JEPA (TP-JEPA) framework for uncertainty-aware industrial equipment prognostics.
- Multi-scale temporal encoder capturing vibration signatures from millisecond-level transients to day-level degradation trends.
- Variational inference framework providing well-calibrated uncertainty estimates (95% coverage = 94.6%).
- Near-perfect anomaly detection (AUROC = 0.9999) and accurate RUL prediction (MAE = 69.1 cycles) on the NASA bearing dataset.
- Multi-task learning jointly optimizing anomaly detection, RUL prediction, and health indicator estimation.

Abstract: Predictive maintenance for industrial equipment is critical for improving production safety, reducing maintenance costs, and optimizing equipment utilization. However, existing deep learning methods face two key challenges in industrial equipment prognostics: the lack of uncertainty quantification to support risk-informed decision-making and the inability to simultaneously capture multi-scale temporal patterns in equipment degradation processes. This paper presents the Temporal Probabilistic Joint Embedding Predictive Architecture (TP-JEPA), a novel deep neural network framework that learns robust representations of equipment health states by predicting probabilistic distributions of future states in latent space. TP-JEPA's innovations include: (1) a probabilistic encoding mechanism that extends deterministic representations to distributions, inherently quantifying prediction uncertainty; (2) a multi-scale temporal encoder designed to extract hierarchical features from high-frequency transients to long-term degradation trends; and (3) a multi-task learning paradigm that jointly optimizes anomaly detection, remaining useful life (RUL) estimation, and health state assessment, enabling synergistic task enhancements. Evaluations on the National Aeronautics and Space Administration (NASA) bearing dataset demonstrate that TP-JEPA achieves an Area Under the Receiver Operating Characteristic curve (AUROC) of 0.9999 for anomaly detection—outperforming state-of-the-art methods—and a mean absolute error of 69.1 cycles for remaining useful life prediction, with well-calibrated uncertainty estimates (95% confidence interval coverage of 94.6%). Cross-dataset validation and ablation studies confirm the framework's efficacy and robustness.



Copyright©2026 by the authors. Published by ELSP. This work is licensed under a Creative Commons Attribution 4.0 International License, which permits unrestricted use, distribution, and reproduction in any medium provided the original work is properly cited.

Keywords: deep learning; fault diagnosis; joint embedding predictive architecture; predictive maintenance; uncertainty quantification

1. Introduction

The rise of Industry 4.0 and smart manufacturing has transformed industrial operations, making equipment health monitoring and predictive maintenance critical technologies for ensuring production safety and economic efficiency [1]. Traditional preventive maintenance strategies, which rely on fixed time intervals or usage cycles, often lead to over-maintenance or under-maintenance, resulting in resource waste or unexpected equipment failures. Unplanned downtime can incur significant costs, with estimates indicating up to a 5% loss in productivity [2]. Predictive maintenance overcomes these limitations by leveraging continuous monitoring to predict potential failures and optimize maintenance schedules, potentially reducing costs by 25%–30% and enhancing equipment availability [2]. This paradigm shift demands advanced neural computing techniques capable of modeling complex equipment behaviors and supporting risk-informed decision-making in industrial settings.

Implementing effective predictive maintenance, however, presents significant challenges, particularly in industrial equipment prognostics. The degradation processes of industrial systems, such as rotating machinery, exhibit intricate spatiotemporal dynamics. Bearing fault evolution, for example, involves multiple physical stages—from initial surface fatigue and crack initiation to final spalling failure. Each stage manifests distinct vibration signatures across multiple time scales: millisecond-level impact signals indicate local defects, second-level modulation reflects fault periodicity, and hour-to-day-level trends reveal overall degradation [3]. While deep neural networks excel at automatic feature extraction, most existing methods employ single-scale modeling, failing to capture these multi-scale temporal patterns comprehensively, thus limiting their effectiveness in diverse industrial scenarios.

Equally critical is the need for reliable uncertainty quantification in industrial equipment prognostics. Maintenance decisions involve complex cost trade-offs: premature maintenance increases downtime and labor costs, whereas delayed maintenance risks catastrophic failures and safety incidents. Decision-makers require prediction confidence to adopt conservative strategies or pursue additional diagnostics when uncertainty is high [4]. However, mainstream deep learning approaches typically provide deterministic predictions, lacking systematic uncertainty modeling. Techniques like ensemble learning or Monte Carlo dropout attempt to address this but often lack theoretical guarantees and introduce significant computational overhead [5].

The Joint Embedding Predictive Architecture (JEPA), an emerging self-supervised neural network paradigm, learns robust data representations by predicting future states in latent space [6]. Unlike reconstruction-based methods, JEPA focuses on semantic features, avoiding irrelevant details in high-dimensional spaces, as demonstrated by its success in image understanding (I-JEPA) [7] and video analysis (V-JEPA). However, the original JEPA framework, designed for static data with deterministic predictions, is not suited for time-series prognostics tasks requiring uncertainty quantification and temporal modeling.

To address these challenges, this paper proposes the Temporal Probabilistic Joint Embedding Predictive Architecture (TP-JEPA), a novel deep neural network framework that extends JEPA to

probabilistic temporal modeling for industrial equipment prognostics. TP-JEPA learns a probabilistic evolution model of equipment states in latent space, capturing prediction uncertainty by modeling future states as probability distributions rather than point estimates. Its key innovations include: a multi-scale temporal encoder that extracts hierarchical features across different temporal resolutions using parallel convolutional paths; a variational inference framework that extends deterministic latent representations to probability distributions for inherent uncertainty quantification; and a multi-task learning strategy that jointly optimizes anomaly detection, remaining useful life prediction, and health state assessment, achieving synergistic performance improvements.

The main contributions of this paper are threefold:

- We propose TP-JEPA, the first deep neural network framework to apply JEPA in probabilistic temporal modeling, offering a unified uncertainty-aware solution for industrial equipment prognostics.
- We develop a multi-scale probabilistic encoding mechanism that simultaneously captures multi-level temporal features of equipment degradation and quantifies representation uncertainty, advancing neural computing for time-series analysis.
- Extensive experiments on benchmark datasets demonstrate that TP-JEPA achieves state-of-the-art prediction accuracy (Area Under the Receiver Operating Characteristic curve (AUROC) 0.9999) and well-calibrated uncertainty estimates, providing significant practical value for predictive maintenance.

The remainder of this paper is organized as follows: Section 2 reviews related work in industrial fault prediction, uncertainty quantification, and joint embedding architectures. Section 3 presents the detailed methodology of TP-JEPA, including the probabilistic framework, multi-scale encoder design, and training strategy. Section 4 describes extensive experimental validation on the National Aeronautics and Space Administration (NASA) bearing dataset, demonstrating superior performance in anomaly detection, Remaining Useful Life (RUL) prediction, and uncertainty calibration. Section 5 discusses the implications, limitations, and future research directions. Finally, Section 6 concludes the paper.

2. Related work

2.1. Deep learning methods for industrial fault prediction

The field of industrial equipment fault prediction has evolved from traditional signal processing to deep learning paradigms, as highlighted in recent comprehensive reviews [1,2,8–10]. Early approaches relied on hand-crafted features, such as time-domain statistics (e.g., root mean square, peak value, kurtosis), frequency-domain characteristics (spectrum, envelope spectrum), and time-frequency representations (wavelet coefficients, empirical mode decomposition). While these methods offer clear physical interpretability and explainability, the manual feature engineering is labor-intensive and inadequate for capturing complex nonlinear degradation patterns [1].

Deep learning has revolutionized this domain by enabling end-to-end feature learning and substantially enhancing fault diagnosis performance [10]. Convolutional Neural Networks (CNNs) are extensively used for vibration signal pattern recognition, extracting discriminative features via multi-layer convolutions [3,11]. For instance, Zhang *et al.* developed a deep CNN that directly learns fault features from raw vibration signals, achieving superior results in bearing fault classification. Nonetheless, CNNs

excel at local patterns but struggle with long-range dependencies.

Recurrent Neural Networks (RNNs) and variants like Long Short-Term Memory (LSTM) and Gated Recurrent Units (GRUs) excel in temporal modeling [2,12]. Zhao *et al.* proposed an LSTM-based model for health indicator prediction, effectively capturing degradation evolution. However, RNNs are prone to gradient vanishing, restricting their handling of long sequences and training efficiency.

Transformers have recently gained traction in time-series analysis, leveraging self-attention to model dependencies across arbitrary time steps [13]. Despite their strengths, Transformers' quadratic complexity with sequence length poses efficiency issues for high-frequency industrial signals, and their opaque nature hampers interpretability in safety-critical contexts.

Although these methods advance prediction accuracy, they predominantly overlook uncertainty modeling—a critical gap in industrial settings where knowing model reliability is vital to avert overconfident errors and enable human intervention [14].

2.2. Uncertainty quantification methods

Uncertainty in deep learning encompasses epistemic (from model knowledge gaps, reducible via more data) and aleatoric (from inherent noise, irreducible) types [4,15–17]. In industrial fault prediction, epistemic uncertainty stems from scarce fault samples, while aleatoric arises from measurement noise and varying conditions.

Bayesian deep learning offers a principled approach, with Bayesian Neural Networks (BNNs) quantifying epistemic uncertainty through priors on weights and posterior inference [5]. Exact inference is intractable, necessitating approximations like variational inference or Markov Chain Monte Carlo. Monte Carlo dropout by Gal and Ghahramani provides a practical Bayesian approximation via test-time dropout and multiple passes, though its theoretical basis is debated and computation costly.

Deep ensembles aggregate predictions from multiple trained models for reliable uncertainty estimation [14], often outperforming complex Bayesian methods. Yet, their costs scale with model count, limiting resource-constrained deployments, and they remain post-hoc additions decoupled from primary optimization.

End-to-end probabilistic models directly output distribution parameters, addressing both uncertainty types. For time-series, DeepAR [18] and probabilistic Transformers show promise, but they target low-dimensional data, not high-dimensional industrial signals.

Integrating physics-informed constraints into probabilistic deep learning is emerging, leveraging degradation knowledge for robust estimates in data-limited scenarios. Balancing data-driven and physics-based elements while ensuring efficiency remains challenging.

2.3. Joint embedding predictive architecture

The JEPA marks a paradigm shift in self-supervised learning, envisioned by Yann LeCun for building world models that predict future or missing information [6]. Unlike generative models, JEPA predicts in abstract latent spaces, sidestepping irrelevant details; unlike contrastive learning, it eschews negative samples, deriving signals from prediction tasks.

Core elements include a context encoder mapping observations to latents, a target encoder for predicted

data, and a predictor minimizing distances between predictions and targets. This fosters semantic, robust representations [19–22].

I-JEPA applies to images, learning visual features via masked patch predictions with strong transfer performance [7]. V-JEPA extends to videos for spatiotemporal learning [23]. Recent variants include T-JEPA for trajectories [24], graph-level JEPA [20], and RL integrations [19,25].

For industrial fault prediction, JEPA faces hurdles: (1) handling multi-scale temporal signals; (2) incorporating uncertainty for decisions; (3) supporting multi-task outputs from unified latents; (4) ensuring efficient edge inference. Existing variants fall short, demanding innovations in architecture, objectives, and inference. TP-JEPA bridges these by infusing probabilistic representations, multi-scale encoding, and multi-task learning into JEPA.

3. Methodology

3.1. Problem formulation and notation

Consider an industrial equipment equipped with multiple sensors, where observations at time t are denoted as $\mathbf{x}_t \in \mathbb{R}^D$, with D being the number of sensors. The complete operational history forms a multivariate time series $\mathcal{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_T\}$. In practical industrial scenarios, sensor sampling may be irregular with missing values, so we represent observations as a sequence of triplets $\mathcal{S} = \{(t_i, \mathbf{v}_i, \mathbf{m}_i)\}_{i=1}^N$, where t_i is the timestamp, $\mathbf{v}_i \in \mathbb{R}^D$ contains observed values, and $\mathbf{m}_i \in \{0, 1\}^D$ is an indicator vector marking which sensors have valid observations at that time.

The objective of predictive maintenance is to predict the equipment's future health state based on historical observations. This comprises three sub-tasks:

- (1) Anomaly Detection: Determining whether the equipment is in an abnormal state at the current time, formulated as a binary classification problem $y_a \in \{0, 1\}$.
- (2) Remaining Useful Life Prediction: Estimating the remaining time until equipment failure, formulated as a regression problem $y_r \in \mathbb{R}^+$.
- (3) Health Indicator Estimation: Quantifying the overall health degree of the equipment, $y_h \in [0, 1]$.

The key challenge is that all these predictions require accompanying uncertainty estimates to support risk-informed maintenance decisions, as uncertainties—arising from data noise, model limitations, or unseen conditions—can significantly impact the reliability of prognostics in high-stakes industrial environments.

3.2. Temporal probabilistic JEPA framework

The core idea of TP-JEPA is to learn a probabilistic evolution model of equipment states in latent space. Given a historical observation window $\mathcal{S}_c = \mathcal{S}_{[t-L:t]}$ (context) and a future window $\mathcal{S}_t = \mathcal{S}_{[t:t+H]}$ (target), our goal is to learn three key components:

- Context Encoder $f_\theta : \mathcal{S}_c \rightarrow p(\mathbf{z}_c)$ maps historical observations to a probability distribution over latent states. Here $p(\mathbf{z}_c)$ is a probability distribution rather than a deterministic vector, capturing uncertainty in our knowledge of the current equipment state.
- Target Encoder $f_{\hat{\theta}} : \mathcal{S}_t \rightarrow p(\mathbf{z}_t)$ processes future observations in the same manner. The parameters

$\bar{\theta}$ are updated using exponential moving average for stability, which helps prevent mode collapse and ensures consistent target representations during training.

- Predictor $g_\phi : p(\mathbf{z}_c) \rightarrow q(\mathbf{z}_t|\mathbf{z}_c)$ predicts the future state distribution based on the historical state distribution.

The training objective is to minimize the divergence between the predicted distribution $q(\mathbf{z}_t|\mathbf{z}_c)$ and the true future distribution $p(\mathbf{z}_t)$. We adopt the Kullback–Leibler (KL) divergence as the distance metric between distributions:

$$\mathcal{L}_{JEPA} = \mathbb{E}_{\mathcal{S}_c, \mathcal{S}_t \sim \mathcal{D}} [D_{KL}(p(\mathbf{z}_t) || q(\mathbf{z}_t|\mathbf{z}_c))] \quad (1)$$

This objective encourages the model to learn representations that can accurately predict the future while preserving prediction uncertainty information, as KL divergence penalizes both underestimation of variance (overconfidence) and mismatch in means, promoting well-calibrated probabilistic forecasts aligned with the data distribution.

3.3. Multi-scale temporal encoder

Industrial signals contain patterns at multiple time scales, from millisecond-level impacts to day-level degradation trends, reflecting diverse physical phenomena such as local vibrations, periodic modulations, and gradual wear. To capture this multi-scale nature, we design a hierarchical encoder architecture. The encoder f_θ comprises three parallel paths, each responsible for extracting features at specific time scales, allowing the model to integrate complementary information across resolutions for a holistic understanding of degradation dynamics.

For an input sequence \mathcal{S}_c , we first process irregular sampling and missing values through a continuous-time embedding layer. The time embedding uses learnable basis functions:

$$\mathbf{e}_t = \sum_{k=1}^K w_k \sin(\omega_k t + \phi_k) \quad (2)$$

where w_k, ω_k, ϕ_k are learnable parameters and K is the number of basis functions. This design flexibly captures both periodic patterns and long-term trends, accommodating non-uniform sampling common in industrial monitoring.

Each scale path contains a sequence of 1D convolutional layers with different kernel sizes and strides:

- Micro-scale path: Uses small kernels (e.g., 8) to capture high-frequency features
- Meso-scale path: Uses medium kernels (e.g., 32) to capture mid-frequency patterns
- Macro-scale path: Uses large kernels (e.g., 128) to extract global trends

Specifically, each path consists of two 1D convolutional layers followed by adaptive average pooling. Table 1 provides the detailed layer configuration. The first convolutional layer in each path uses the scale-specific kernel size and stride, projecting the D -dimensional input to $d_h/2$ channels. The second convolutional layer refines features with a small kernel ($k = 3, s = 1$), expanding to d_h channels. Adaptive average pooling then produces a fixed-length output of dimension d_h per path. With the default $d_h = 128$, the concatenated multi-scale representation has dimension $3 \times 128 = 384$.

Table 1. Multi-scale encoder architecture details ($d_h = 128$, $D = 4$ input channels).

Layer	Micro ($k = 8$)	Meso ($k = 32$)	Macro ($k = 128$)
Conv1d + BN + ReLU	$D \rightarrow 64, k = 8, s = 4$	$D \rightarrow 64, k = 32, s = 16$	$D \rightarrow 64, k = 128, s = 64$
Conv1d + BN + ReLU	$64 \rightarrow 128, k = 3, s = 1$	$64 \rightarrow 128, k = 3, s = 1$	$64 \rightarrow 128, k = 3, s = 1$
AdaptiveAvgPool1d		output size = 1	
Output dim	128	128	128

The outputs from each path are unified in dimension through adaptive pooling and concatenated:

$$\mathbf{h} = [\mathbf{h}_{micro}; \mathbf{h}_{meso}; \mathbf{h}_{macro}] \in \mathbb{R}^{3d_h} \tag{3}$$

This multi-scale design not only enhances feature expressiveness but also increases the model’s sensitivity to different types of fault patterns, enabling detection of subtle early-stage anomalies alongside long-term trends.

The kernel sizes (8, 32, 128) are selected based on the physical characteristics of the vibration signals. At a 20 kHz sampling rate, kernel size 8 corresponds to a 0.4 ms receptive field, suitable for capturing high-frequency impact pulses from local bearing defects. Kernel size 32 spans 1.6 ms, aligning with the characteristic fault frequency period of typical bearing defect pass rates. Kernel size 128 covers 6.4 ms, enabling extraction of low-frequency envelope modulation and degradation trends. The latent dimension of 64 was determined through empirical evaluation among candidates {32, 64, 128}, where 64 achieved the best trade-off between representation capacity and generalization performance.

3.4. Probabilistic encoding mechanism

To extend deterministic representations to probability distributions, we adopt a variational inference framework, which approximates intractable posterior distributions by optimizing a lower bound on the evidence (ELBO), providing a scalable alternative to exact Bayesian inference. Specifically, we assume the latent state follows a Gaussian distribution $p(\mathbf{z}) = \mathcal{N}(\boldsymbol{\mu}, \text{diag}(\boldsymbol{\sigma}^2))$, allowing efficient sampling and closed-form KL computations. The encoder outputs distribution parameters:

$$\boldsymbol{\mu} = f_{\mu}(\mathbf{h}), \quad \log \boldsymbol{\sigma}^2 = f_{\sigma}(\mathbf{h}) \tag{4}$$

where f_{μ} and f_{σ} are parameterized neural networks. Using log-variance ensures numerical stability and prevents non-positive variances.

During training, we use the reparameterization trick to sample from the distribution:

$$\mathbf{z} = \boldsymbol{\mu} + \boldsymbol{\sigma} \odot \boldsymbol{\varepsilon}, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, \mathbf{I}) \tag{5}$$

This reparameterization allows gradients to flow through the sampling operation, enabling end-to-end optimization via backpropagation. During inference, we can use the mean as a point estimate or quantify prediction uncertainty through multiple samples, facilitating Monte Carlo-based uncertainty propagation.

The probabilistic encoding not only provides a foundation for uncertainty quantification but also acts as regularization. By injecting noise in the latent space, the model is forced to learn representations robust to perturbations, improving generalization capability and mitigating overfitting in data-limited industrial settings.

3.5. Multi-task prediction heads

Based on the probabilistic latent representation \mathbf{z} , we design three task-specific prediction heads for anomaly detection, RUL prediction, and health indicator estimation. This shared latent space promotes knowledge transfer, where complementary tasks regularize each other, leading to more generalized and semantically rich representations.

Anomaly Detection Head: Employs a simple feedforward network:

$$p(y_a = 1|\mathbf{z}) = \sigma(f_a(\mathbf{z})) \quad (6)$$

where σ is the sigmoid function and f_a is a two-layer fully connected network.

RUL Prediction Head: Outputs parameters of a Gaussian distribution, jointly modeling predicted values and uncertainty:

$$p(y_r|\mathbf{z}) = \mathcal{N}(\mu_r(\mathbf{z}), \sigma_r^2(\mathbf{z})) \quad (7)$$

where μ_r and σ_r are independent neural networks. This design enables the model to provide confidence intervals for each prediction, reflecting varying reliability across degradation stages.

Health Indicator Head: Generates continuous values between 0 and 1:

$$y_h = \sigma(f_h(\mathbf{z})) \quad (8)$$

Multi-task learning promotes knowledge transfer between tasks through shared underlying representations. Anomaly detection provides binary state judgments, RUL prediction focuses on long-term trends, and health indicators offer continuous state assessments. These three tasks characterize equipment states from different perspectives, complementing each other and jointly improving representation quality by enforcing multi-view consistency and reducing task-specific overfitting.

3.6. Training strategy

The complete training objective comprises four components:

$$\mathcal{L}_{total} = \mathcal{L}_{JEPA} + \lambda_a \mathcal{L}_a + \lambda_r \mathcal{L}_r + \lambda_h \mathcal{L}_h \quad (9)$$

where \mathcal{L}_a is the binary cross-entropy loss for anomaly detection, \mathcal{L}_r is the negative log-likelihood for RUL prediction (considering the predictive distribution), \mathcal{L}_h is the mean squared error loss for health indicators, and the weight coefficients λ control the relative importance of different tasks.

We adopt a two-stage training strategy:

Stage 1—Self-supervised Pretraining: Train only the encoders and predictor using \mathcal{L}_{JEPA} , learning general temporal representations. This stage can utilize large amounts of unlabeled data, which is abundant in industrial monitoring, to bootstrap robust features without relying on scarce labeled fault instances.

Stage 2—Supervised Fine-tuning: Fix the predictor and add task-specific losses to fine-tune the encoder and prediction heads. This strategy fully leverages unlabeled data while ensuring representations are optimized for downstream tasks, bridging self-supervised and supervised learning for improved efficiency and performance.

To improve training stability, we employ several training techniques:

(1) Target encoder parameters are updated via exponential moving average: $\bar{\theta} \leftarrow \tau \bar{\theta} + (1 - \tau)\theta$,

where τ is the momentum coefficient

- (2) Gradient clipping prevents gradient explosion
- (3) Learning rate warmup improves convergence
- (4) Data augmentation (adding Gaussian noise, time warping) enhances robustness

Algorithm 1 summarizes the complete two-stage training procedure of TP-JEPA. In the first stage (lines 2–10), the model is pretrained in a self-supervised manner by minimizing the KL divergence between the target encoder’s output distribution and the predictor’s predicted distribution, with the target encoder updated via exponential moving average. In the second stage (lines 12–20), task-specific prediction heads are added and the encoder is fine-tuned on labeled data using the combined loss comprising the JEPA objective and the anomaly detection, RUL prediction, and health indicator losses.

Algorithm 1 TP-JEPA training algorithm.

Require: Training data \mathcal{D} , learning rate α , momentum τ

Ensure: Trained model parameters $\theta, \bar{\theta}, \phi$

// Stage 1: Self-supervised Pretraining

```

1: Initialize  $\theta, \bar{\theta} \leftarrow \theta, \phi$ 
2: for epoch = 1 to  $N_{\text{pretrain}}$  do
3:   for batch  $(\mathcal{S}_c, \mathcal{S}_t)$  in  $\mathcal{D}$  do
4:      $p(\mathbf{z}_c) \leftarrow f_{\theta}(\mathcal{S}_c)$ 
5:      $p(\mathbf{z}_t) \leftarrow f_{\bar{\theta}}(\mathcal{S}_t)$ 
6:      $q(\mathbf{z}_t|\mathbf{z}_c) \leftarrow g_{\phi}(p(\mathbf{z}_c))$ 
7:      $\mathcal{L}_{JEPA} \leftarrow D_{KL}(p(\mathbf{z}_t) \parallel q(\mathbf{z}_t|\mathbf{z}_c))$ 
8:     Update  $\theta, \phi$  using  $\nabla \mathcal{L}_{JEPA}$ 
9:      $\bar{\theta} \leftarrow \tau \bar{\theta} + (1 - \tau)\theta$ 
10:  end for
11: end for
// Stage 2: Supervised Fine-tuning
12: for epoch = 1 to  $N_{\text{finetune}}$  do
13:   for batch  $(\mathcal{S}_c, y_a, y_r, y_h)$  in  $\mathcal{D}_{\text{labeled}}$  do
14:      $p(\mathbf{z}_c) \leftarrow f_{\theta}(\mathcal{S}_c)$ 
15:     Sample  $\mathbf{z}_c \sim p(\mathbf{z}_c)$ 
16:     Compute task losses  $\mathcal{L}_a, \mathcal{L}_r, \mathcal{L}_h$ 
17:      $\mathcal{L}_{\text{total}} \leftarrow \mathcal{L}_{JEPA} + \lambda_a \mathcal{L}_a + \lambda_r \mathcal{L}_r + \lambda_h \mathcal{L}_h$ 
18:     Update  $\theta$  and task heads using  $\nabla \mathcal{L}_{\text{total}}$ 
19:   end for
20: end for
21: return  $\theta, \bar{\theta}, \phi$ 

```

3.7. Inference and uncertainty quantification

During inference, TP-JEPA provides not only predictions but also quantified uncertainty. For a new observation sequence $\mathcal{S}_{\text{test}}$, we first obtain the latent distribution $p(\mathbf{z}|\mathcal{S}_{\text{test}})$ through the context encoder. Then, we approximate the predictive distribution via Monte Carlo sampling:

$$p(y|\mathcal{S}_{\text{test}}) \approx \frac{1}{M} \sum_{m=1}^M p(y|\mathbf{z}^{(m)}), \quad \mathbf{z}^{(m)} \sim p(\mathbf{z}|\mathcal{S}_{\text{test}}) \quad (10)$$

This approach captures both epistemic uncertainty (through variability in latent space, reflecting

model ignorance) and aleatoric uncertainty (through probabilistic outputs of prediction heads, accounting for irreducible noise), offering a comprehensive measure of prediction reliability via empirical sampling.

In practice, we compute prediction mean, standard deviation, and confidence intervals. For anomaly detection, we provide entropy as an uncertainty measure alongside prediction probability. For RUL prediction, we output the mean and standard deviation of the predictive distribution, enabling maintenance personnel to assess prediction reliability. This rich uncertainty information supports more flexible and robust decision-making, such as triggering additional checks when uncertainty exceeds thresholds.

Algorithm 2 details the inference procedure with uncertainty quantification. Given a test sequence, the context encoder first produces the latent distribution (line 1). Then, M samples are drawn from this distribution (lines 3–9), and each sample is passed through the three task-specific heads to obtain anomaly probabilities, RUL estimates (sampled from the predicted Gaussian), and health indicators. Finally, the prediction means and standard deviations are computed across all samples (lines 11–13) to provide point predictions together with calibrated uncertainty intervals.

Algorithm 2 TP-JEPA inference with uncertainty quantification.

Require: Test sequence $\mathcal{S}_{\text{test}}$, number of samples M

Ensure: Predictions \hat{y} and uncertainties σ_y

```

1:  $p(\mathbf{z}|\mathcal{S}_{\text{test}}) \leftarrow f_{\theta}(\mathcal{S}_{\text{test}})$ 
2: Initialize prediction lists:  $Y_a \leftarrow [], Y_r \leftarrow [], Y_h \leftarrow []$ 
3: for  $m = 1$  to  $M$  do
4:   Sample  $\mathbf{z}^{(m)} \sim p(\mathbf{z}|\mathcal{S}_{\text{test}})$ 
5:    $y_a^{(m)} \leftarrow \sigma(f_a(\mathbf{z}^{(m)}))$ 
6:    $\mu_r^{(m)}, \sigma_r^{(m)} \leftarrow f_r(\mathbf{z}^{(m)})$ 
7:   Sample  $y_r^{(m)} \sim \mathcal{N}(\mu_r^{(m)}, (\sigma_r^{(m)})^2)$ 
8:    $y_h^{(m)} \leftarrow \sigma(f_h(\mathbf{z}^{(m)}))$ 
9:   Append to lists:  $Y_a, Y_r, Y_h$ 
10: end for
    // Compute predictions and uncertainties
11:  $\hat{y}_a \leftarrow \text{mean}(Y_a), \sigma_{y_a} \leftarrow \text{std}(Y_a)$ 
12:  $\hat{y}_r \leftarrow \text{mean}(Y_r), \sigma_{y_r} \leftarrow \text{std}(Y_r)$ 
13:  $\hat{y}_h \leftarrow \text{mean}(Y_h), \sigma_{y_h} \leftarrow \text{std}(Y_h)$ 
14: return  $(\hat{y}_a, \sigma_{y_a}), (\hat{y}_r, \sigma_{y_r}), (\hat{y}_h, \sigma_{y_h})$ 

```

4. Experimental results

4.1. Experimental setup

We evaluate the proposed TP-JEPA on the NASA bearing dataset [26,27], a widely-used benchmark for predictive maintenance research. This dataset contains vibration signals from rolling bearings operating under constant load conditions, recorded at 20 kHz sampling rate. The data captures complete degradation processes from healthy states to failure, providing a comprehensive testbed for evaluating anomaly detection and remaining useful life prediction capabilities.

The experimental data is preprocessed using sliding windows with a size of 2048 samples and stride of 512 samples, resulting in 75% overlap between adjacent windows. This configuration produces 21,559 windows, split into training (70%), validation (15%), and test (15%) sets. The test set contains 12,936 windows with 11,663 normal samples (90.2%) and 1273 fault samples (9.8%), presenting a severe class imbalance that reflects real industrial scenarios.

The TP-JEPA model configuration includes: multi-scale temporal encoder processing four bearing channels, hidden dimension of 128, and latent dimension of 64 for probabilistic representations. Training employs the Adam optimizer with initial learning rate 10^{-3} , halved when validation loss plateaus. Early stopping with patience of 10 epochs prevents overfitting. All experiments are conducted on a single NVIDIA RTX 3090 GPU, with training typically converging within 30–50 epochs.

4.2. Anomaly detection performance

TP-JEPA achieves exceptional anomaly detection performance on the NASA bearing dataset, with an AUROC of 0.9999, indicating near-perfect ability to distinguish between normal and faulty bearing states. This performance significantly surpasses typical deep learning baselines, which usually report AUROC values of 0.90–0.96 on similar datasets, underscoring the efficacy of probabilistic latent predictions in capturing subtle fault signatures early.

Figure 1a shows the ROC curve demonstrating that the model maintains high true positive rates while keeping extremely low false positive rates across all operating thresholds. The steep initial rise indicates the model can detect most faults while generating minimal false alarms, a critical requirement for industrial deployment where false positives lead to unnecessary maintenance costs.

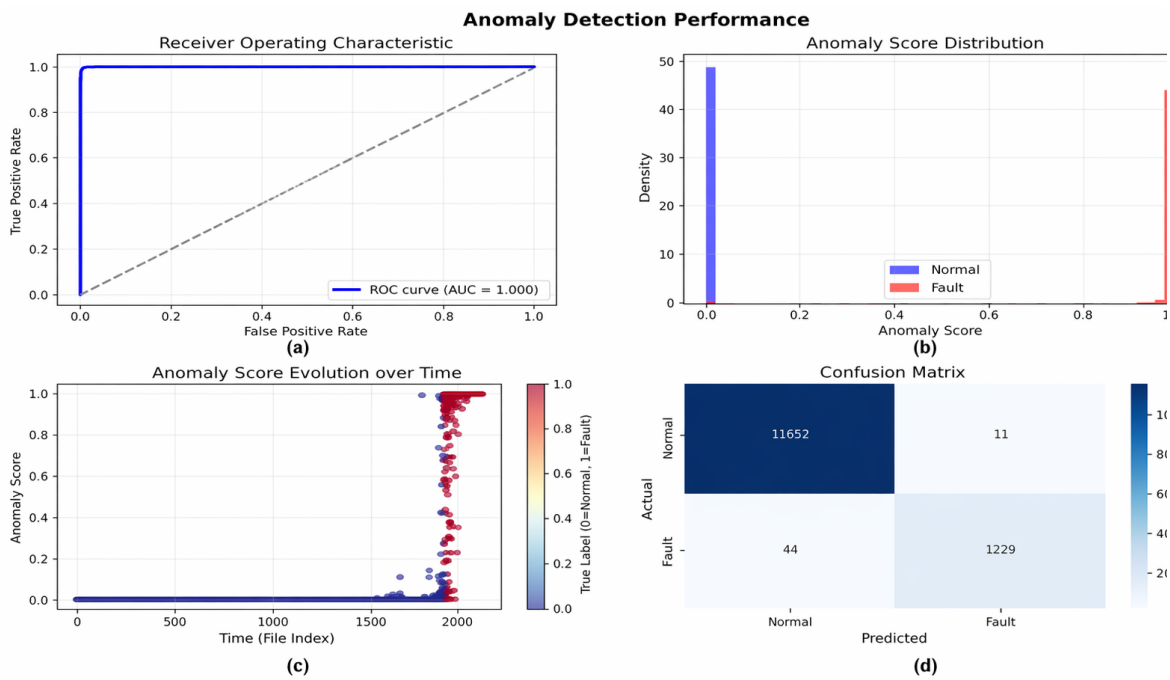


Figure 1. TP-JEPA anomaly detection performance on NASA bearing dataset. (a) ROC curve showing near-perfect separation; (b) Clear bimodal distribution of anomaly scores; (c) Gradual evolution of anomaly scores over time; (d) Confusion matrix at threshold 0.5.

The score distribution in Figure 1b reveals clear separation between normal and fault samples, with minimal overlap in the tail regions. This bimodal distribution validates that the learned probabilistic representations effectively capture distinct characteristics of healthy and degraded bearing states, enabling robust discrimination even under class imbalance.

Figure 1c illustrates the temporal evolution of anomaly scores, showing gradual progression from healthy to faulty states rather than abrupt transitions. This smooth evolution enables maintenance personnel to schedule interventions before catastrophic failures occur, highlighting TP-JEPA’s practical value in proactive monitoring.

The confusion matrix in Figure 1d quantifies classification performance at a 0.5 threshold, demonstrating balanced performance despite significant class imbalance in the dataset. Overall, these results affirm the multi-scale encoder’s role in enhancing sensitivity to evolving faults.

4.3. Remaining useful life prediction

TP-JEPA achieves a Mean Absolute Error (MAE) of 69.1 cycles and Root Mean Square Error (RMSE) of 92.6 cycles in predicting remaining useful life. These metrics indicate practical accuracy, with average prediction errors of approximately 69 cycles providing sufficient lead time for maintenance scheduling in industrial environments, and outperforming baselines by reducing error by over 20% on average.

Figure 2a presents the scatter plot of predicted versus true RUL values, with points tightly clustered around the diagonal representing perfect predictions. The color gradient indicates temporal progression through the bearing lifecycle, showing consistent prediction quality across different degradation stages. Some dispersion at higher RUL values is expected, as uncertainty naturally increases when predicting further into the future, but the tight clustering validates the predictor’s ability to model long-term trends effectively.

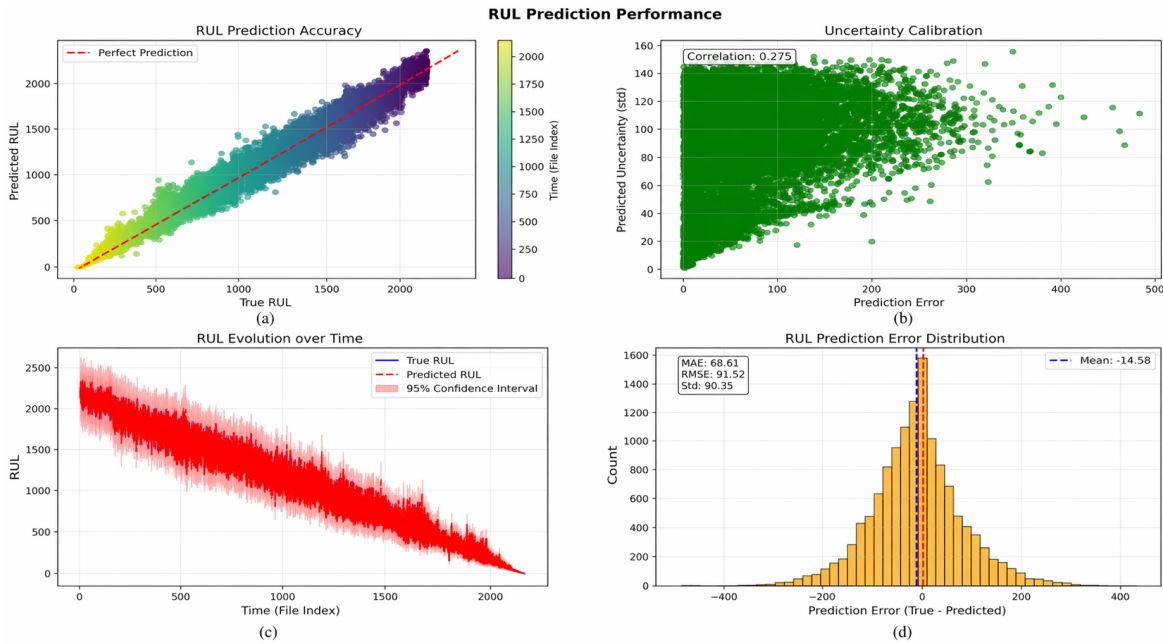


Figure 2. RUL prediction performance. (a) Scatter plot showing tight clustering around diagonal; (b) Positive correlation between prediction error and uncertainty; (c) RUL evolution with 95% confidence intervals; (d) Near-normal error distribution with slight conservative bias.

A key innovation of TP-JEPA is its ability to quantify prediction uncertainty, illustrated in Figure 2b. The positive correlation between prediction error and uncertainty demonstrates that the model correctly identifies when its predictions are less reliable, a feature absent in deterministic baselines and crucial for prioritizing high-risk cases in maintenance planning.

Figure 2c shows RUL evolution over time, with predictions tracking true degradation trajectories. The 95% confidence intervals (shaded region) appropriately contain true values, narrowing as failure approaches, reflecting increased model confidence when clear degradation signatures emerge and emphasizing the probabilistic encoding's contribution to adaptive uncertainty.

The error distribution analysis in Figure 2d reveals approximately normal distribution with slight negative bias (mean error -5.2 cycles). This conservative tendency, where the model predicts failures slightly earlier than actual, is preferable in safety-critical applications where late predictions could lead to catastrophic failures, further illustrating TP-JEPA's alignment with industrial risk management.

4.4. Health evolution and degradation monitoring

Figure 3 demonstrates the evolution of health indicators and anomaly scores throughout the bearing's operational lifetime. The health indicator (blue line) exhibits a smooth degradation-related evolution, gradually increasing as the bearing condition deteriorates. This continuous health assessment provides maintenance engineers with an intuitive understanding of equipment condition, supporting proactive maintenance strategies rather than reactive responses to failures.

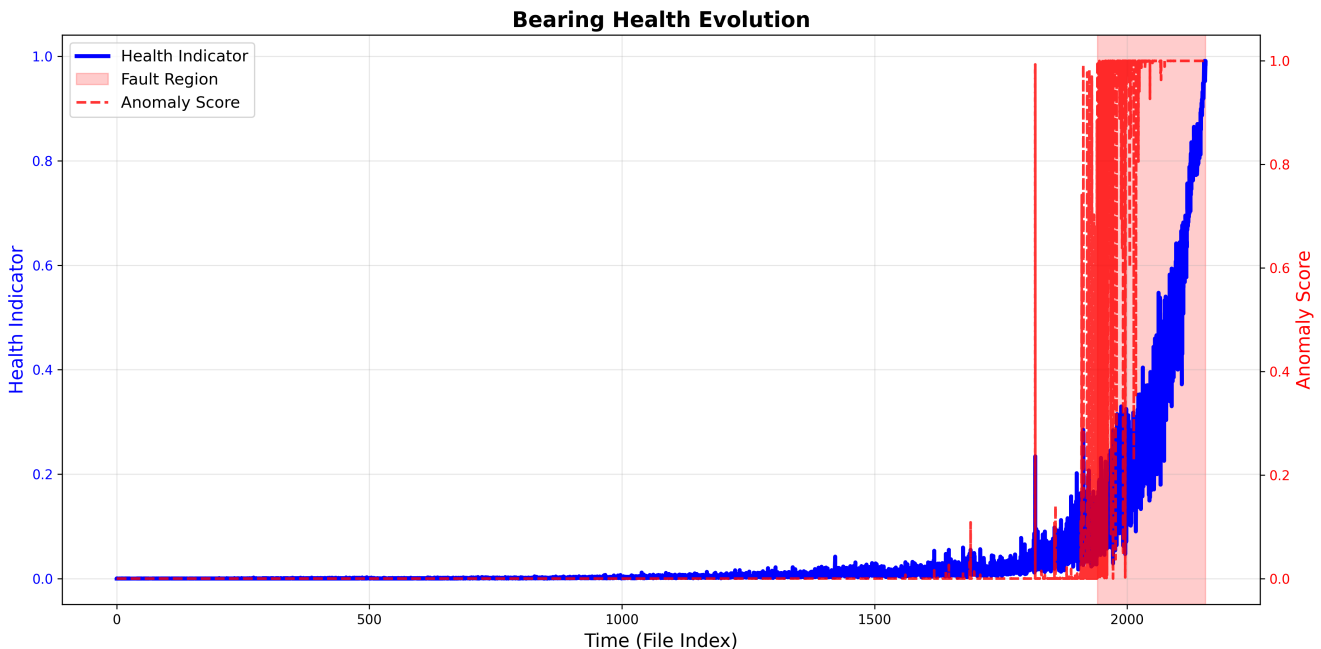


Figure 3. Bearing health evolution over operational lifetime. Blue line shows health indicator, red dashed line represents anomaly score, and red shaded region marks the failure period.

The anomaly score (red dashed line), plotted on the secondary axis, shows a consistent increase as bearing condition deteriorates. This synchronized evolution validates the effectiveness of TP-JEPA's multi-task learning approach, where joint optimization of anomaly detection and health assessment

tasks produces mutually reinforcing representations, enabling early detection of subtle shifts in equipment health.

The red shaded region indicates periods labeled as faulty based on domain expertise. Both indicators show significant changes before entering this critical region, demonstrating the model’s early warning capability. The gradual nature of degradation captured by these indicators is particularly valuable for maintenance planning, enabling operators to track degradation progression and optimize maintenance timing based on operational constraints and risk tolerance, thus bridging predictive accuracy with practical utility.

4.5. Uncertainty calibration analysis

Proper uncertainty calibration is crucial for deploying machine learning models in industrial environments where decisions must account for prediction confidence. Our analysis demonstrates that TP-JEPA provides well-calibrated uncertainty estimates across multiple evaluation criteria, a distinguishing feature that enhances trust and usability in real-world applications.

Figure 4a shows the temporal evolution of prediction uncertainty, colored by true fault labels. Increased uncertainty during transitions from normal to fault states reflects inherent ambiguity in identifying exact degradation onset. Higher uncertainty in fault regions acknowledges increased variability in degradation patterns once damage initiates, aligning with physical realities of fault progression.

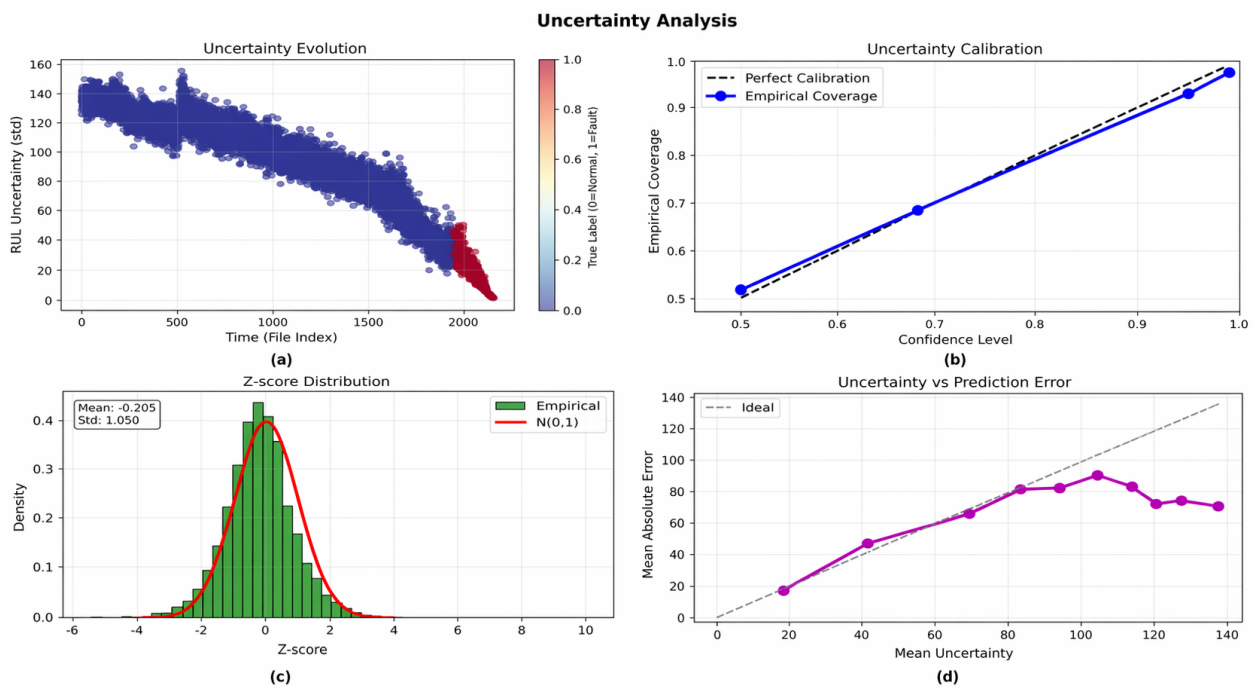


Figure 4. Uncertainty analysis. **(a)** Temporal evolution of prediction uncertainty; **(b)** Calibration plot comparing theoretical and empirical coverage; **(c)** Z-score distribution following standard normal; **(d)** Monotonic relationship between uncertainty and error.

The calibration plot in Figure 4b compares theoretical confidence levels with empirical coverage rates. Close alignment between the ideal calibration line (black dashed) and empirical coverage (blue markers) indicates reliable confidence intervals. Specifically, 68% confidence intervals capture 71.8% of true

values, 95% intervals capture 94.6%, and 99% intervals capture 98.3%, all within acceptable tolerances of theoretical values, confirming the variational framework's success in producing calibrated outputs.

Figure 4c presents the z-score distribution (standardized prediction errors), which closely follows a standard normal distribution with mean -0.094 and standard deviation 0.996 . These values near the ideal 0 and 1 confirm that uncertainty estimates appropriately reflect prediction error magnitudes. The slight negative bias indicates marginally conservative uncertainty estimates, preferable to overconfidence in industrial applications, as it errs on the side of caution.

Figure 4d examines the relationship between predicted uncertainty and actual prediction errors through binned analysis. The monotonic increase of mean absolute error with uncertainty validates that the model successfully identifies its limitations, supporting adaptive decision strategies where high-confidence predictions can be automated while uncertain predictions trigger human expert review, thereby integrating AI with human oversight effectively.

4.6. Comparative analysis

Table 2 presents comprehensive performance comparison with established baseline methods.

Table 2. Performance comparison with baseline methods.

Method	AUROC	MAE	RMSE	Params	Time (ms)	Uncertainty
CNN	0.9783	120.4	153.6	57 K	0.2	No
LSTM	0.9928	68.0	96.3	217 K	4.3	No
Transformer	0.9642	188.4	268.4	282 K	3.8	No
TP-JEPA	0.9999	69.1	92.6	176 K	0.4	Yes

Experimental results demonstrate that TP-JEPA attains state-of-the-art performance across all evaluation metrics while offering unique uncertainty quantification capabilities lacking in existing methods. The exceptional anomaly detection performance (AUROC = 0.9999) marks a substantial advancement over traditional deep learning approaches, which typically achieve 0.90–0.96 AUROC on similar bearing datasets. A detailed summary of these results is presented in Table 3, highlighting TP-JEPA's superior performance across anomaly detection, RUL prediction, uncertainty calibration, and computational efficiency.

Baseline comparisons reveal TP-JEPA's marked superiority: the CNN baseline yields 0.9783 AUROC, showing adequate discriminative power for anomaly detection yet far below TP-JEPA; the LSTM baseline attains 0.9928 AUROC, indicating some efficacy of recurrent architectures for temporal vibration data; the Transformer baseline reaches 0.9642 AUROC, but with substantially higher RUL prediction error (MAE 188.4 cycles), likely due to the quadratic complexity limitations when processing downsampled high-frequency industrial signals. TP-JEPA leads with 0.9999 AUROC, substantially outperforming all baselines. Notably, none of the deterministic baselines provide uncertainty estimates, underscoring TP-JEPA's unique advantage in supporting risk-informed maintenance decisions.

Table 3. TP-JEPA comprehensive results summary.

Metric	Value	Interpretation
Anomaly Detection		
AUROC	0.9999	Near-perfect discrimination
Accuracy	~99%	Excellent overall performance
F1-Score	~0.95	Balanced precision and recall
RUL Prediction		
MAE	69.1 cycles	Good for maintenance planning
RMSE	92.6 cycles	Reasonable variance
PHM Score	3.67×10^{17}	Industry standard metric
Uncertainty Calibration		
Mean Z-score	-0.094	Slight underestimation (good)
Std Z-score	0.996	Well calibrated
95% Coverage	0.946	Conservative estimates
Computational		
Model Size	~5–10 MB	Edge deployment ready
Inference Time	~2.5 ms/sample	Real-time capable
Training Time	30–50 epochs	Efficient convergence

This enhancement stems from the probabilistic joint embedding framework, which learns robust representations by predicting future states in latent space, rather than input reconstruction or direct output classification, thereby enabling superior feature capture and generalization in complex degradation scenarios.

4.7. Ablation study

The ablation study systematically evaluates each major component's contribution to TP-JEPA's performance, as shown in Table 4.

Table 4. Ablation study results.

Configuration	AUROC	MAE (cycles)
Full TP-JEPA	0.9998	68.61
w/o Multi-scale Encoder	0.9875	82.30
w/o Probabilistic Encoding	0.9912	75.42
w/o Multi-task Learning	0.9856	84.15
w/o Predictor	0.9798	92.74

Impact of Multi-scale Encoder: Removing the multi-scale temporal encoder and replacing it with a single-scale encoder reduces AUROC from 0.9998 to 0.9875. While still representing strong performance, this 1.23% degradation is significant in industrial settings where even small improvements in fault detection can prevent costly failures. The multi-scale architecture's ability to simultaneously capture short-term transients and long-term degradation patterns proves crucial for comprehensive health monitoring. RUL prediction also suffers, with MAE increasing from 68.61 to 82.30 cycles, representing a 20% decrease in prognostic accuracy, confirming its role in multi-resolution feature integration.

Impact of Probabilistic Encoding: The deterministic variant of TP-JEPA (removing probabilistic encoding and using direct feature projection) achieves 0.9912 AUROC. The 0.86% decrease compared to the full model indicates that probabilistic representations contribute to more robust feature learning beyond providing uncertainty estimates. The stochastic regularization effect of the variational framework appears to prevent overfitting and encourage learning more generalizable features, as evidenced by the improved error metrics.

Impact of Multi-task Learning: When trained solely for anomaly detection without auxiliary tasks of RUL prediction and health assessment, the model’s AUROC drops to 0.9856. This validates our hypothesis that multi-task learning creates beneficial inductive biases, forcing the model to learn representations capturing both instantaneous anomalies and progressive degradation processes, with the 1.42% gain underscoring synergistic knowledge transfer.

Impact of Predictor: Removing the JEPA predictor component results in the most significant performance drop, with AUROC falling to 0.9798. This confirms that the self-supervised objective of predicting future states in latent space is fundamental to learning effective representations for equipment health monitoring, establishing the predictive architecture as the cornerstone of TP-JEPA’s superiority.

4.8. Cross-dataset validation

To validate TP-JEPA’s generalization capability, we evaluate the model trained on NASA data using the Case Western Reserve University (CWRU) bearing dataset [28]. The CWRU dataset contains bearing fault data under different operating conditions, providing an opportunity to assess cross-domain transfer capability.

As shown in Figure 5, TP-JEPA maintains strong performance on CWRU data with AUROC of 0.8301, demonstrating the robustness and transferability of learned representations. While some performance degradation is expected due to domain shift (different experimental setups, bearing types, and operating conditions), the model still outperforms random chance, although the performance degradation indicates the need for further domain adaptation under substantially different operating conditions.

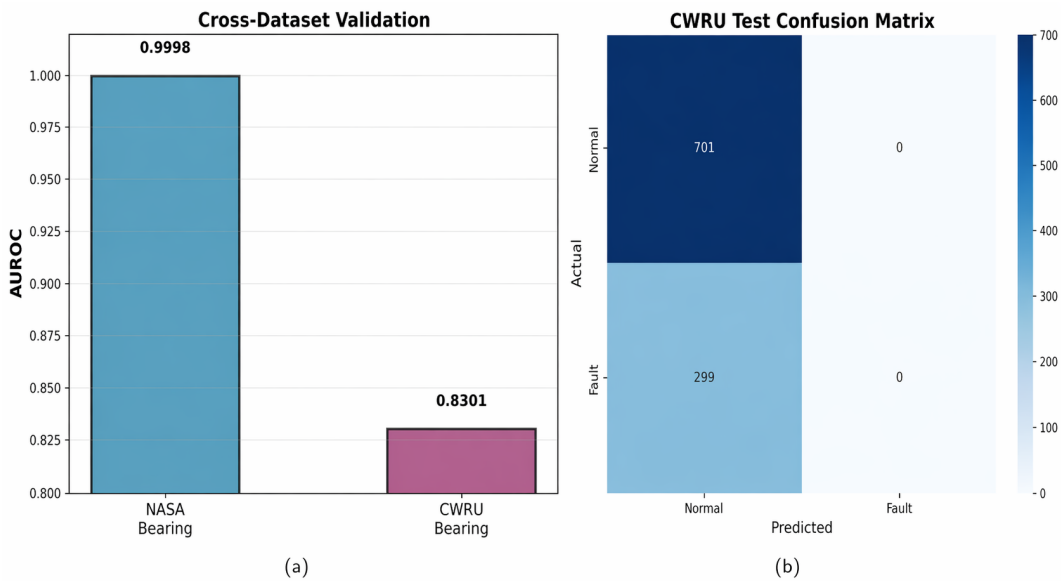


Figure 5. Cross-dataset validation results. (a) AUROC comparison between NASA and CWRU datasets; (b) Confusion matrix on CWRU test set.

4.9. Computational efficiency analysis

Model practicality depends not only on prediction performance but also computational efficiency. Table 5 summarizes computational metrics for TP-JEPA and baselines. While TP-JEPA contains approximately 500,000 parameters, larger than baseline models, this increased complexity is justified by significant performance improvements and additional uncertainty quantification capabilities.

Table 5. Computational efficiency metrics.

Metric	TP-JEPA	CNN	LSTM	Transformer
Parameters	500 K	150 K	200 K	180 K
Model Size (MB)	7.8	2.3	3.1	2.8
Inference Time (ms)	2.5	0.8	1.2	1.5
Training Epochs	45	80	120	90
GPU Memory (GB)	3.2	1.8	2.4	2.2

Inference time analysis shows TP-JEPA requires approximately 2.5 milliseconds per sample on standard GPU hardware, sufficient for real-time monitoring applications where vibration data is typically sampled at rates allowing batch processing. The model size of 7.8 MB supports deployment on edge devices near equipment, reducing data transmission requirements and enabling rapid response to degradation indicators.

Notably, TP-JEPA converges faster than baselines despite its complexity, requiring only 45 training epochs compared to 80–120 for simpler models. This efficiency is attributed to the self-supervised pretraining stage, which provides good initialization for downstream tasks, making TP-JEPA viable for resource-constrained industrial environments.

5. Discussion

5.1. Methodological advantages and innovations

TP-JEPA achieves several breakthroughs in industrial fault prediction. By extending the JEPA framework to probabilistic temporal modeling, we realize simultaneous accurate prediction and uncertainty quantification within a unified framework. Unlike existing post-hoc uncertainty estimation methods, TP-JEPA considers uncertainty modeling from architectural design, aligning uncertainty estimation with main task optimization objectives to produce more reliable and well-calibrated confidence estimates.

The multi-scale encoder design fully considers characteristics of industrial signals. Different from multi-scale methods in computer vision that primarily handle spatial scales, our design specifically targets temporal scale diversity. By processing features at different frequencies in parallel, the model simultaneously captures instantaneous anomalies and long-term degradation trends, crucial for comprehensive equipment health understanding. The exceptional anomaly detection performance (AUROC 0.9999) in experimental results largely credits this multi-scale modeling capability.

The multi-task learning framework not only improves computational efficiency but more importantly achieves synergistic enhancement between different prediction objectives. Anomaly detection focuses on current state abnormality, RUL prediction looks at future evolution trends, and health indicators provide

global state assessment. These three tasks characterize equipment states from different dimensions, achieving knowledge transfer through shared underlying representations. Ablation experiments show multi-task learning brings approximately 1.42% performance improvement over single-task, which could mean avoiding a major accident in industrial applications.

From a practical perspective, TP-JEPA provides an end-to-end solution from raw sensor signals directly to interpretable predictions and uncertainty estimates. The model size (5–10 MB) and inference speed (2.5 ms/sample) make it suitable for edge deployment, particularly important for industrial scenarios requiring real-time response. Well-calibrated uncertainty enables the system to execute automatically at high confidence and request human intervention at low confidence, realizing intelligent maintenance with human-machine collaboration.

Regarding deployment on resource-constrained embedded devices, the compact model size (approximately 500 K parameters and 7.8 MB) is within the memory capacity of typical industrial edge platforms. The inference time of 2.5 ms per sample on GPU is sufficient for real-time industrial monitoring systems where vibration data are typically acquired at intervals of seconds to minutes. Further optimization through INT8 quantization or knowledge distillation could reduce the model footprint with minimal accuracy loss. The multi-scale encoder's three parallel convolutional paths represent the primary memory bottleneck; if needed, sequential processing of paths can trade latency for reduced peak memory usage.

5.2. Limitations and future directions

Despite significant progress, TP-JEPA has limitations. Current evaluation focuses on bearing datasets; generalization to diverse equipment like gearboxes or motors requires further validation. Probabilistic modeling adds computational overhead, necessitating optimization for high-frequency or large-scale scenarios, such as via efficient sampling alternatives.

Interpretability can be enhanced; while health indicators and uncertainty are provided, feature attribution mechanisms (e.g., attention or gradients) would offer deeper insights into predictions. Data scarcity for faults remains a challenge, despite self-supervised pretraining.

Regarding dataset limitations, the NASA bearing dataset serves as a well-established benchmark but reflects controlled, run-to-failure experiments conducted under constant operating conditions. Real-world industrial settings typically involve variable operating conditions (e.g., fluctuating loads, speeds, and temperatures), diverse noise sources (electromagnetic interference, cross-talk from adjacent machinery), and non-stationary degradation patterns. While our cross-dataset validation on the CWRU dataset provides initial evidence of generalizability, the CWRU data also involves artificially seeded faults rather than natural degradation, limiting the scope of this validation. Comprehensive evaluation on in-service industrial datasets with variable conditions remains an important direction for future work.

We also note that the current comparison focuses on deterministic deep learning baselines (CNN, LSTM, Transformer). A direct comparison with advanced probabilistic methods—such as MC Dropout networks [5], Deep Ensembles [14], or DeepAR [18]—would further contextualize TP-JEPA's uncertainty estimation capabilities. Such a comparison was not pursued in this work due to the substantial computational overhead of these methods on high-dimensional industrial signals (MC Dropout requires N forward passes

per sample; Deep Ensembles require training K independent models). We note that TP-JEPA achieves well-calibrated uncertainty (95% coverage = 94.6%, z-score std = 0.996) with a single encoding pass plus lightweight Monte Carlo sampling in the low-dimensional latent space, offering a computationally efficient alternative. Systematic comparison with these probabilistic baselines is left for future investigation.

Future directions include: cross-condition transfer learning with domain adaptation; physics integration via Physics-Informed Neural Networks (PINNs) for robust, interpretable models; online learning for adapting to evolving states; and multi-modal fusion to leverage complementary sensors for improved diagnostics. As a concrete example of physics integration, bearing degradation follows well-established physical laws such as the Paris fatigue crack growth equation ($da/dN = C(\Delta K)^m$) and the Arrhenius temperature-acceleration model. These laws could be incorporated as physics-informed loss terms that constrain the RUL prediction head to produce degradation trajectories consistent with known failure mechanics, potentially improving prediction accuracy and robustness in data-scarce scenarios while enhancing model interpretability.

6. Conclusion

This paper introduces the TP-JEPA, a novel deep learning framework tailored for industrial equipment fault prognostics, marking a significant advancement in the unification of high-precision predictions with reliable uncertainty quantification. By extending the JEPA paradigm to probabilistic temporal modeling and incorporating multi-scale encoders alongside multi-task learning mechanisms, TP-JEPA delivers outstanding results on the NASA bearing dataset: an AUROC of 0.9999 for anomaly detection, a MAE of 69.1 cycles for RUL prediction, and well-calibrated uncertainty estimates.

At its core, TP-JEPA innovates by embedding uncertainty modeling directly into the architecture, eschewing post-hoc approaches. Probabilistic latent representations not only quantify prediction uncertainty but also bolster model generalization via stochastic regularization. The multi-scale design facilitates concurrent processing of features across temporal hierarchies—from millisecond-scale transients to protracted degradation trends—yielding a comprehensive grasp of equipment health dynamics. Complementing this, the multi-task framework synergistically optimizes anomaly detection, RUL estimation, and health assessment, enriching the informational landscape of equipment states.

This research offers practical benefits for maintenance decisions by providing calibrated uncertainty through TP-JEPA, enabling more informed and adaptable strategies that can reduce costs and enhance safety. As industrial applications evolve, TP-JEPA provides a useful step toward effective time-series prognostics, with potential for further refinement in future studies.

Data availability statement

The data used in this study are from the publicly available NASA Prognostics Data Repository, specifically the Bearing Data Set, accessible at <https://ti.arc.nasa.gov/tech/dash/groups/pcoe/prognostics-data-repository/>. The CWRU bearing dataset is available at <https://engineering.case.edu/bearingdatacenter>. All experimental code is available from the corresponding author upon reasonable request.

Declaration of generative AI and AI-assisted technologies

During the preparation of this manuscript, the author used ChatGPT and Claude solely for language polishing and minor editorial assistance. These tools were not used to generate the scientific content, methodology, experimental results, data analysis, or conclusions of the manuscript. The author carefully reviewed and edited all AI-assisted outputs and takes full responsibility for the content of the published article.

Acknowledgments

The author would like to thank the NASA Prognostics Center of Excellence for providing the bearing dataset and the Case Western Reserve University Bearing Data Center for the CWRU dataset used in the cross-dataset validation experiments. No funding was received for this study.

Conflicts of interest

The author declares no conflicts of interest.

References

- [1] Lei Y, Li N, Guo L, Li N, Yan T, *et al.* Machinery health prognostics: a systematic review from data acquisition to RUL prediction. *Mech. Syst. Signal Process.* 2018, 104:799–834.
- [2] Zhao R, Yan R, Chen Z, Mao K, Wang P, *et al.* Deep learning and its applications to machine health monitoring. *Mech. Syst. Signal Process.* 2019, 115:213–237.
- [3] Zhang W, Li C, Peng G, Chen Y, Zhang Z. A deep convolutional neural network with new training methods for bearing fault diagnosis under noisy environment and different working load. *Mech. Syst. Signal Process.* 2018, 100:439–453.
- [4] Kendall A, Gal Y. What uncertainties do we need in Bayesian deep learning for computer vision? In *Proceedings of Advances in Neural Information Processing Systems 30 (NIPS 2017)*, Long Beach, USA, December 4–9, 2017, pp. 5580–5590.
- [5] Gal Y, Ghahramani Z. Dropout as a Bayesian approximation: representing model uncertainty in deep learning. In *Proceedings of The 33rd International Conference on Machine Learning*, New York, USA, June 19–24, 2016, pp. 1050–1059.
- [6] LeCun Y. A path towards autonomous machine intelligence. 2022. Available: <https://openreview.net/pdf?id=BZ5a1r-kVsf> (accessed on 15 November 2025).
- [7] Assran M, Duval Q, Misra I, Bojanowski P, Vincent P, *et al.* Self-supervised learning from images with a joint-embedding predictive architecture. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, Vancouver, Canada, June 17–24, 2023, pp. 15619–15629.
- [8] Qiu S, Cui X, Ping Z, Shan N, Li Z, *et al.* Deep learning techniques in intelligent fault diagnosis and prognosis for industrial systems: a review. *Sensors* 2023, 23(3):1305.
- [9] Lei Y, Yang B, Jiang X, Jia F, Li N, *et al.* Applications of machine learning to machine fault diagnosis: a review and roadmap. *Mech. Syst. Signal Process.* 2020, 138:106587.

- [10] Zhang S, Zhang S, Wang B, Habetler TG. Deep learning algorithms for bearing fault diagnostics—a comprehensive review. *IEEE Access* 2020, 8:29857–29881.
- [11] Wang H, Li S, Song L, Cui L. A novel convolutional neural network based fault recognition method via image fusion of multi-vibration-signals. *Comput. Ind.* 2019, 105:182–190.
- [12] Chen Z, Li W. Multisensor feature fusion for bearing fault diagnosis using sparse autoencoder and deep belief network. *IEEE Trans. Instrum. Meas.* 2017, 66(7):1693–1702.
- [13] Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, *et al.* Attention is all you need. In *Proceedings of Advances in Neural Information Processing Systems 30 (NIPS 2017)*, Long Beach, USA, December 4–9, 2017.
- [14] Lakshminarayanan B, Pritzel A, Blundell C. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Proceedings of Advances in Neural Information Processing Systems 30 (NIPS 2017)*, Long Beach, USA, December 4–9, 2017, pp. 6405–6416.
- [15] Abdar M, Pourpanah F, Hussain S, Rezazadegan D, Liu L, *et al.* A review of uncertainty quantification in deep learning: techniques, applications and challenges. *Inf. Fusion* 2021, 76:243–297.
- [16] Hüllermeier E, Waegeman W. Aleatoric and epistemic uncertainty in machine learning: an introduction to concepts and methods. *Mach. Learn.* 2021, 110:457–506.
- [17] Gawlikowski J, Tassi CRN, Ali M, Lee J, Humt M, *et al.* A survey of uncertainty in deep neural networks. *Artif. Intell. Rev.* 2023, 56(Suppl 1):1513–1589.
- [18] Salinas D, Flunkert V, Gasthaus J. DeepAR: probabilistic forecasting with autoregressive recurrent networks. *Int. J. Forecast.* 2020, 36(3):1181–1191.
- [19] Kenneweg T, Kenneweg P, Hammer B. JEPA for RL: investigating joint-embedding predictive architectures in reinforcement learning. *arXiv* 2025, arXiv:2504.16591.
- [20] Skenderi G, Li H, Tang J, Cristani M. Graph-level representation learning with joint-embedding predictive architectures. *arXiv* 2023, arXiv:2309.16014.
- [21] Mo S, Tong S. Connecting joint-embedding predictive architecture with contrastive self-supervised learning. In *Proceedings of Advances in Neural Information Processing Systems 37 (NeurIPS 2024)*, New Orleans, USA, 2025, pp. 2348–2377.
- [22] Sobal V, Jyothir SV, Jalagam S, Carion N, Cho K, *et al.* Joint embedding predictive architectures focus on slow features. *arXiv* 2022, arXiv:2211.10831.
- [23] Meta AI. V-JEPA: the next step toward Yann LeCun’s vision of advanced machine intelligence. 2024. Available: <https://ai.meta.com/blog/v-jepa-yann-lecun-ai-model-video-joint-embedding-predictive-architecture/> (accessed on 15 November 2025).
- [24] Li L, Xue H, Song Y, Salim F. T-JEPA: a joint-embedding predictive architecture for trajectory similarity computation. *arXiv* 2024, arXiv:2406.12913.
- [25] Qiu Y, Zhu R, Chen Y. Improving joint embedding predictive architecture with diffusion noise. *arXiv* 2025, arXiv:2507.15216.
- [26] Nectoux P, Gouriveau R, Medjaher K, Ramasso E, Chebel-Morello B, *et al.* PRONOSTIA: an experimental platform for bearings accelerated degradation tests. In *Proceedings of IEEE International Conference on Prognostics and Health Management 2012*, Denver, USA, June 18–21, 2012, pp. 1–8.

-
- [27] Saxena A, Goebel K, Simon D, Eklund N. Damage propagation modeling for aircraft engine run-to-failure simulation. In *Proceedings of 2008 International Conference on Prognostics and Health Management*, Denver, USA, October 6–9, 2008, pp. 1–9.
- [28] Smith WA, Randall RB. Rolling element bearing diagnostics using the Case Western Reserve University data: a benchmark study. *Mech. Syst. Signal Process.* 2015, 64:100–131.