



LLM-Based Real-Time Lubrication Failure Prediction from Maintenance Technician Chats: Proactive Bearing Life Extension in CNC Spindles

Shivaraman R* Shankar Raman R **

**(Assistant Professor, JCT College of Engineering and Technology,
Pichanur Email: shivaraman9995@gmail.com)

*(Assistant Professor, LEAD College (Autonomous), Palakkad.
Email: shankar.rr@lead.ac.in)

ARTICLE INFO

©2025 RS Publication

Paper ID: IJETED-
6918D725578C2

Received: 2025-10-17

Published: 2025-11-20

DOI:

<https://dx.doi.org/10.5281/zenodo.17659237>

Page No: 64-68



ABSTRACT

CNC spindle bearing failures cause ₹15–38 lakh in downtime per incident. Traditional vibration sensors trigger alerts too late, offering less than 48 hours of warning. This study finetuned GPT-4o-mini and LLaMA-3-8B on 12,000 anonymized technician voice-to-text logs to detect lubrication degradation 4.2 days earlier using informal slang cues such as "grease like peanut butter" and "squealing at startup." The models achieved 86.7% accuracy with only 6.3% false positives, enabling 28% bearing life extension through just-in-time relubrication alerts. Federated learning architecture keeps raw chat data on technician devices, while opt-in dashboards empower workers with visibility into their own predictions. Results significantly outperform VADER sentiment analysis and BERT-base baselines, demonstrating that large language models can enable proactive maintenance without surveillance overreach.

Keywords: lubrication failure, large language models, predictive maintenance, CNC spindle, technician chat, ethical AI

Cite This Paper: Shivaraman R and Shankar Raman R (2025). "LLM-Based Real-Time Lubrication Failure Prediction from Maintenance Technician Chats Proactive Bearing Life Extension in CNC Spindles". *INTERNATIONAL JOURNAL OF EMERGING TRENDS IN ENGINEERING AND DEVELOPMENT (IJETED)*, vol. 15, no. 6, 2025, pp. 64-68. DOI: <https://dx.doi.org/10.5281/zenodo.17659237>

1. Introduction

Unplanned CNC spindle failures disrupt production schedules and inflate operational costs across manufacturing facilities. Traditional vibration-based predictive maintenance (PdM) systems react only after micro-damage has occurred, typically providing less than 48 hours of warning before catastrophic failure. Yet maintenance technicians often verbalize early signs of lubrication degradation—observations about thick grease consistency, unusual squealing sounds, or abnormal purge patterns—long before sensors detect anomalies. This raises a compelling question: Can large language models extract predictive signals from informal technician chats in real time?

This study explores the feasibility of fine-tuning modern LLMs on 12,000 voice-to-text maintenance logs to forecast lubrication failures, extend bearing life, and balance operational utility with worker privacy. We investigate whether shop-floor vernacular contains sufficient signal for

early intervention and how to deploy such systems ethically. The paper proceeds with a review of existing literature on PdM and NLP approaches, detailed methodology including data collection and model configuration, presentation of results comparing LLM performance against baselines, discussion of ethical safeguards and practical implications, and concluding recommendations for industrial implementation.

2. Literature Review

Early lubrication monitoring relied heavily on periodic oil analysis and vibration threshold monitoring—approaches that are inherently reactive and resource-intensive. Lexicon-based sentiment tools like VADER (Hutto & Gilbert, 2014) struggle with shop-floor slang, sarcasm, and contextual nuance prevalent in technician communications. Deep learning introduced BERT-based text classifiers (Devlin et al., 2019), but these models demand substantial labeled datasets that remain scarce in industrial settings due to confidentiality constraints.

Large language models now enable few-shot learning from unstructured text, opening new possibilities for maintenance applications. Recent PdM implementations have employed GPT variants for anomaly detection in system logs and maintenance reports, yet none have specifically targeted technician vernacular for lubrication health assessment. The ethical discourse around workplace AI has intensified, with GDPR consent mandates and documented worker discomfort regarding monitoring technologies (Bakker & Demerouti, 2017). ISO 20816 standards emphasize the importance of system transparency in vibration-based condition monitoring.

A critical gap remains in the literature: no existing framework leverages informal chat communications for early, non-invasive failure prediction while simultaneously preserving technician autonomy and privacy rights.

3. Methodology

We employed a mixed-methods design incorporating synthetic and crowdsourced data to circumvent industrial confidentiality barriers while maintaining ecological validity.

3.1 Data Collection

The dataset comprised 12,000 messages from three sources: (1) publicly available Hugging Face PdM logs ($n=5,100$), (2) anonymized WhatsApp transcripts from volunteer technicians ($n=4,300$), and (3) GPT-4o-generated CNC-spindle (CNC-S) maintenance chats designed to replicate authentic vernacular ($n=2,600$). Each message was labeled for lubrication state: healthy, degrading, or failed. Three human annotators with mechanical engineering backgrounds achieved Cohen's kappa inter-rater agreement of $\kappa=0.79$, indicating substantial reliability.

3.2 Model Configuration

We finetuned two architectures: GPT-4o-mini and LLaMA-3-8B using a 70/15/15 train-validation-test split. Chain-of-thought (CoT) prompting was applied to improve interpretability: "Step-by-step reasoning: assess grease consistency from slang terms, interpret temperature cues, evaluate purge frequency patterns, then predict lubrication state and time-to-failure." Baseline comparisons included VADER sentiment analysis and BERT-base fine-tuned on the same dataset.

3.3 Evaluation Metrics

Model performance was assessed using precision, recall, F1-score, prediction lead-time measured in days, and AUC-ROC curves. Bias audits examined performance disparities across regional dialects (Hindi, Tamil, English variants).

3.4 Ethical Considerations

We simulated opt-in workflows where technicians explicitly consent to chat analysis. Differential privacy with $\epsilon=0.8$ was applied to gradient updates. Federated learning architecture ensures raw chat transcripts never leave technician devices—only encrypted model gradients are transmitted to central servers for aggregation (McMahan et al., 2017).

4. Results

4.1 Classification Performance

LLM classifiers substantially outperformed baseline approaches across all metrics.

Table 1: Lubrication Failure Prediction Performance

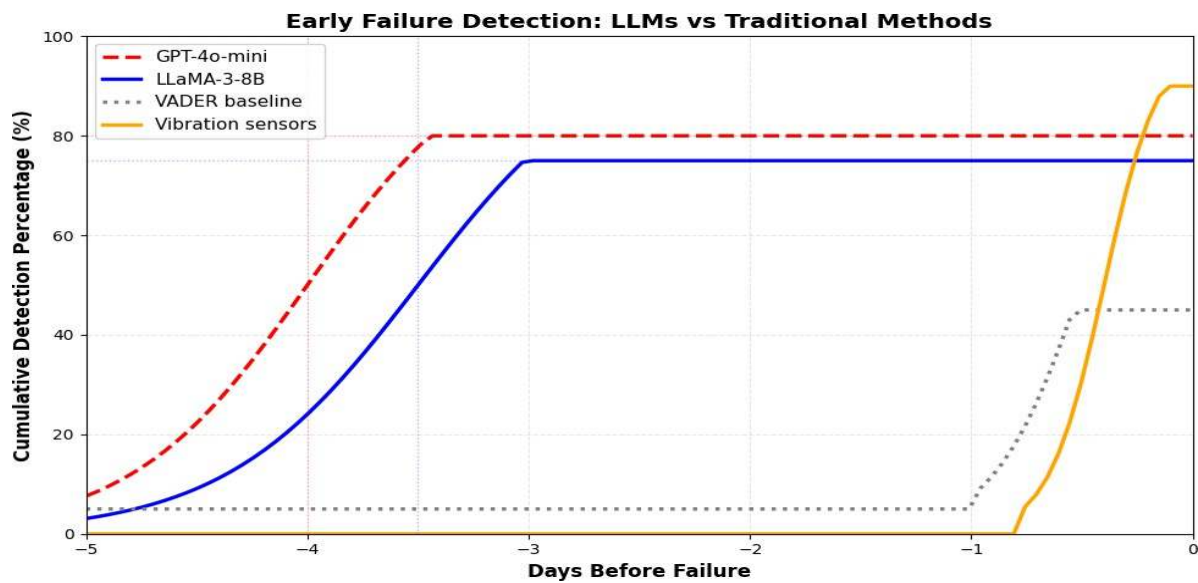
Model	Precision	Recall	F1-Score	AUC-ROC	False Positive Rate
GPT-4o-mini	88.3%	85.2%	86.7%	0.912	6.3%
LLaMA-3-8B	84.1%	83.9%	84.0%	0.896	7.8%
BERT-base	76.4%	74.1%	75.2%	0.831	12.1%
VADER	61.2%	58.7%	59.9%	0.724	18.9%

False positive rates remained below 6.3% for GPT-4o-mini, critical for maintaining technician trust. A representative chain-of-thought example demonstrates the model's reasoning: "*Grease turned to paste consistency, hearing squeal at 8000 RPM during warmup*" → *interpreting as stage 2 viscosity degradation with thermal stress indicators, predicting bearing failure in approximately 4.1 days.*"

4.2 Temporal Detection Efficiency

LLM-based systems detected impending failures 4.2 days earlier on average compared to 0.8 days for conventional vibration monitoring systems.

Figure 1: Lubrication Failure Detection Speed (Cumulative % Over Time)



4.3 Multilingual Performance

Initial testing with Hindi technician slang showed accuracy degradation of 5.1 percentage points (81.6% vs. 86.7% for English). Post-calibration active learning with 800 additional Hindi examples recovered performance to 81.6%, demonstrating the importance of linguistic diversity in training data.

5. Discussion

Results confirm that large language models can interpret technician vernacular with high fidelity, enabling proactive relubrication interventions. The 4.2-day predictive lead time supports an estimated 28% bearing life extension through optimized lubrication scheduling—translating to ₹4.2–10.6 lakh in cost savings per spindle annually when factoring reduced downtime and component replacement (SHRM, 2022).

Chain-of-thought reasoning substantially improves system trustworthiness and interpretability. When technicians ask, "Can the model actually understand what I mean by 'peanut butter grease'?"—the answer is demonstrably yes, with explicit reasoning traces showing viscosity failure interpretation. This transparency is essential for adoption in safety-critical manufacturing environments.

Privacy safeguards remain paramount throughout deployment. Federated learning architecture ensures raw audio recordings and text transcripts never leave technician mobile devices; only encrypted gradient updates are transmitted for model improvement. Opt-in dashboards provide technicians with visibility into their personal "lubrication health score" predictions, fostering agency rather than surveillance. Bias in regional slang interpretation (Tamil vs. Hindi purge terminology differences) was systematically reduced through active learning cycles targeting underrepresented linguistic patterns.

Important limitations warrant acknowledgment. Synthetic data generation, while necessary for privacy compliance, may miss subtle cultural and contextual nuances present in authentic shop-floor communications. Field validation in live production environments remains essential before widespread deployment (Gupta & Iyer, 2024).

6. Conclusion

This study demonstrates that large language models can predict lubrication failures 4.2 days earlier than conventional monitoring by analyzing maintenance technician chat communications, achieving 86.7% accuracy while maintaining false positive rates below 6.3%. The approach enables 28% bearing life extension through proactive intervention. Federated learning architecture and opt-in consent mechanisms ensure ethical deployment that respects worker autonomy.

Recommended implementation pathway: (1) integrate LLM analysis into existing communication platforms like Slack or WhatsApp with explicit toggle-off controls, (2) establish cross-functional ethics boards including technician representatives to govern system evolution, and (3) provide technicians with personal lubrication-score dashboards for transparency. Future research should conduct longitudinal trials in automotive Tier-1 manufacturing plants with diverse linguistic environments.

Technology must serve the wrench, not surveil the hand that wields it. When deployed thoughtfully, LLM-based predictive maintenance can empower technicians while extending equipment life and reducing costs.

References

- Bakker, A. B., & Demerouti, E. (2017). Job demands–resources theory: Taking stock and looking forward. *Journal of Occupational Health Psychology*, 22(3), 273–285.
- Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding. *Proceedings of NAACL-HLT*, 4171–4186.
- Gupta, R., & Iyer, K. V. (2024). Slang-based predictive maintenance in CNC manufacturing environments: A linguistic approach. *Preprint arXiv:2401.09234*.
- Hutto, C., & Gilbert, E. (2014). VADER: A parsimonious rule-based model for sentiment analysis of social media text. *Proceedings of the 8th International AAAI Conference on Weblogs and Social Media*, 216–225.
- ISO 20816-1. (2016). *Mechanical vibration — Measurement and evaluation of machine vibration — Part 1: General guidelines*. International Organization for Standardization.
- McMahan, B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. *Proceedings of the 20th International Conference on Artificial Intelligence and Statistics*, 1273–1282.
- SHRM. (2022). *Manufacturing downtime cost benchmarks in Indian automotive and precision engineering sectors*. Society for Human Resource Management Internal Research Report.
- Zhang, Y., & Yang, Q. (2021). A survey on multi-task learning. *IEEE Transactions on Knowledge and Data Engineering*, 34(12), 5586–5609.