

# In-Context Learning Accuracy Scaling in Tabular Foundation Models vs. Task-Specific Fine-Tuning

Assignee Research

June 9, 2026

## Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does the in-context learning accuracy of tabular foundation models scale with the diversity of training datasets compared to task-specific fine-tuning on standard tabular benchmarks. 10 claims were extracted from source literature; 1 was independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Causal Data Augmentation for Robust Fine-Tuning of Tabular Foundation Models. Research question: How does the in-context learning accuracy of tabular foundation models scale with the diversity of training datasets compared to task-specific fine-tuning on standard tabular benchmarks?.

## 2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

## 3 Results

12 papers retrieved. 10 claims extracted; 1 independently verified. Quality review score: 4.5/10.

## 4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

## 5 Extracted Claims

Claim	Verified	Confidence
CausalMixFT achieves the highest median improvement of $(+0.12 \pm 0.63)$ over the pre-trained model on 33 classification data	×	0.09
Default fine-tuning has a variability of $\pm 0.98$ , while CausalMixFT has a variability of $\pm 0.63$ , indicating greater instability	×	0.07
CausalMixFT ranks first overall in average ranks across datasets, followed by the default fine-tuning baseline, while pu	×	0.07
Early stopping based on limited validation data leads to significant validation set overfitting depending on the fine-tu	✓	0.16
The normalization strategy used to compare performance across different data generators is based on the zero-shot perfor	×	0.07
CausalMixFT extends the fine-tuning framework by mixing real and causally grounded synthetic samples into the fine-tunin	×	0.13
SCM-Based Synthetic Augmentation (CausalMixFT) uses SCMs fitted to the target dataset to generate synthetic data that re	×	0.12
The PC and FCI algorithms are used to estimate the structural relations between features in CausalMixFT.	×	0.02
DoWhy’s SCM framework with additive noise models is used to sample and fit DAGs in CausalMixFT.	×	0.03
Numerical features in CausalMixFT are modeled with regressors, and categorical features with classifiers.	×	0.04

## References

- <http://arxiv.org/abs/2601.04110v2>
- <http://arxiv.org/abs/2208.01009v2>
- <http://arxiv.org/abs/2602.09439v1>