

Pretraining Large Language Models on Synthetic Tabular Data for Few-Shot Reasoning

Assignee Research

June 9, 2026

Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: How does pretraining large language models on diverse synthetic tabular data impact their few-shot reasoning accuracy on unseen domain tasks compared to scaling parameter count. 15 claims were extracted from source literature; 4 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 6.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: LLMTabBench: Evaluating LLMs on Binary Tabular Classification From Zero to Few Shots. Research question: How does pretraining large language models on diverse synthetic tabular data impact their few-shot reasoning accuracy on unseen domain tasks compared to scaling parameter count?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 6.3/10.

3 Results

15 papers retrieved. 15 claims extracted; 4 independently verified. Quality review score: 6.3/10.

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
LLMTabBench is a benchmark designed to evaluate LLMs for tabular classification under data-scarce conditions.	✓	0.30
LLMs are highly competitive in zero-shot settings and can outperform alternative models even when those models have access	✓	0.33
Incorporating additional few-shot examples can conflict with LLM prior knowledge, limiting or degrading performance.	✓	0.32
There is a data complexity threshold beyond which LLMs’ performance declines and few-shot examples become less effective	✓	0.29
The benchmark evaluates shot counts including zero-shot and $k \in \{4, 8, 16, 32, 64\}$ in-context demonstrations.	×	0.08
Five tabular-to-text serialization formats are used: feat_val, feat_val_mask, markdown, markdown_mask, and html.	×	0.02
GPT-4o-mini achieved a ROC-AUC of 0.718 \pm 0.049 in the Few-shot (16) regime.	×	0.04
Qwen3-14B achieved a ROC-AUC of 0.875 \pm 0.002 in the Expert + prior regime.	×	0.02
TabPFN (16) achieved a ROC-AUC of approximately 0.719 across Qwen3-1.7B, Qwen3-8B, and Qwen3-14B models.	×	0.03
The study benchmarks six LLMs: GPT-4o-mini, Qwen3-1.7B, Qwen3-8B, and Qwen3-14B.	×	0.06
Open-weight models are compared using two inference modes: forward scoring and generation.	×	0.03
Forward scoring obtains predictions from token-level log-probabilities of two target labels normalized using a sigmoid function	×	0.02
In-context learning does not update the model’s weights.	×	0.06
Fine-tuning, few-shot learning, and in-context learning all depend on the availability of a small, representative label set	×	0.08
For each few-shot setting, 5 independent stratified draws of demonstrations are run to reduce variance.	×	0.03

References

- <http://arxiv.org/abs/2402.04177v3>
- <http://arxiv.org/abs/2605.24417v1>
- <http://arxiv.org/abs/2510.22389v2>