

# Diversity in Synthetic Pretraining Distributions and Zero-Shot Transfer in Tabular Foundation Models

Assignee Research

June 9, 2026

## Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: How does increasing the diversity of synthetic pretraining distributions impact the zero-shot transfer accuracy of tabular foundation models on out-of-domain TabBench tasks. 10 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.2/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Shaping the Prior: How Synthetic Task Distributions Determine Tabular Foundation Model Quality. Research question: How does increasing the diversity of synthetic pretraining distributions impact the zero-shot transfer accuracy of tabular foundation models on out-of-domain TabBench tasks?.

## 2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.2/10.

## 3 Results

14 papers retrieved. 10 claims extracted; 0 independently verified. Quality review score: 4.2/10.

## 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
The central methodological commitment of this work is isolation: holding the model architecture, optimizer, training bud	×	0.10
NANOTABPFN is used as the base TFM, a lightweight, fully open-source reimplementa-tion of the TABPFN V2 architecture desi	×	0.05
All experiments use the TFM-Playground training protocol, which provides a unified interface for loading synthetic prior	×	0.05
NANOTABPFN at reduced scale requires only tens of thousands of synthetic datasets for pre-training rather than millions,	×	0.04
All models are trained under an identical budget: 40,000 synthetic datasets per prior, each of size $T \in [512, 1024]$ rows	×	0.04
TabPFN v1 generator provides functional diversity but no realism perturbations, no structured missingness, and no distri	×	0.08
TabICL-v1 generator extends the MLP-based SCM prior with tree-based structural equations (XGBoost), producing tasks whos	×	0.04
TabICL-v2 generator augments TabICL-v1 with additional SCM diversity and a richer set of feature-level transforms, inclu	×	0.05
TabICL-v2 broadens functional coverage while remaining free of explicit realism perturbations and shift mechanisms.	×	0.05
Nine variants of O’PRIOR are constructed and organized into four groups following the sequential pipeline of Section 2.	×	0.02

## References

- <http://arxiv.org/abs/2605.18971v1>
- <http://arxiv.org/abs/2512.03307v1>
- <http://arxiv.org/abs/2601.04110v2>