

Appendix: Hybrid Autoencoders for Tabular Data: Leveraging Model-Based Augmentation in Low-Label Settings

Extending the spectral analysis in the main paper (Section 6), we present additional visualizations across more datasets to further examine how the gating mechanisms in TANDEM shape the frequency content of the input. For each dataset, we focus on a single class and compute the unnormalized discrete Fourier transform (NUDFT) over the 50 most variant features. We compare the spectra of the original input x ; the neural-gated input $\hat{x}^{\text{NN}} = x \odot g^{\text{NN}}(x)$, as produced by both TANDEM and SS-AE with gating; and the tree-gated input $\hat{x}^{\text{OSDT}} = x \odot \bar{g}^{\text{OSDT}}(x)$, where \bar{g}^{OSDT} is the average gating mask across all trees and depths.

These extended plots confirm the patterns observed in the main paper: namely, that neural gating acts as a strong low-pass filter, suppressing high-frequency components, while tree-based gating preserves more high-frequency variation, which reinforces their complementary inductive biases in shaping representations.

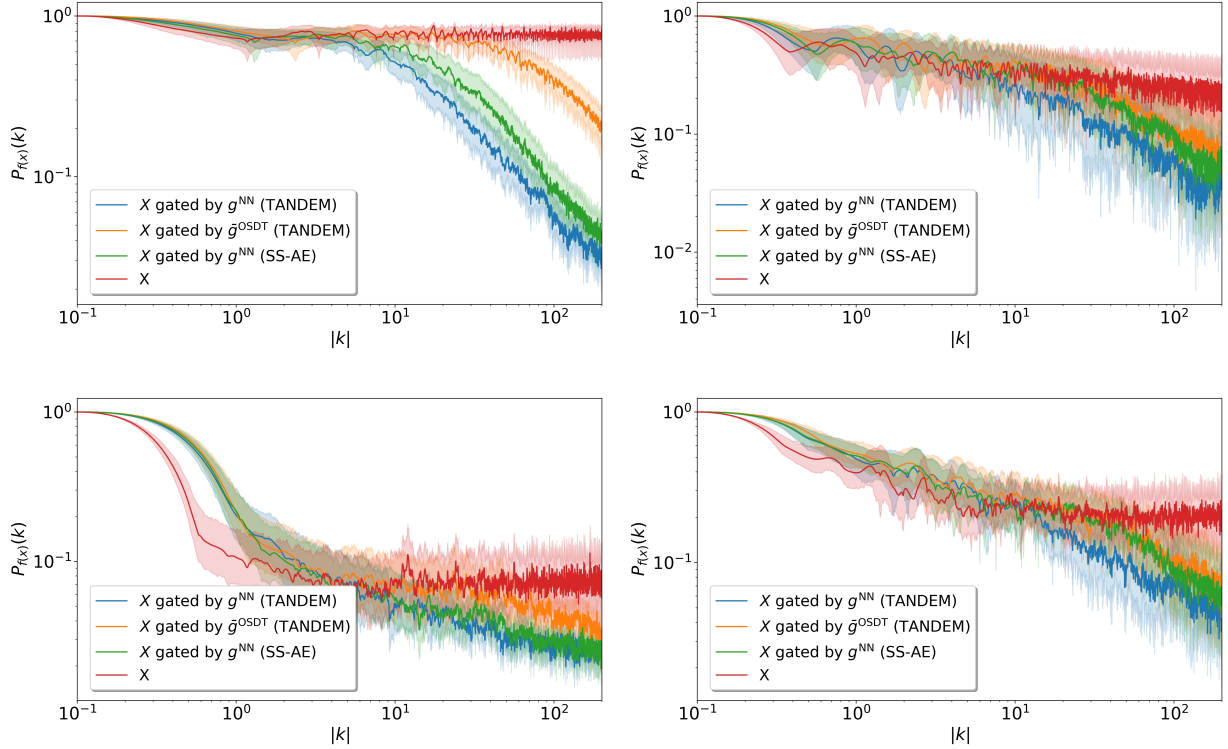


Figure 1: Frequency spectra of gated inputs for the NN and OSDT encoders. Visualized datasets include OG, CP, VO, and CC. NN gating results in stronger suppression of high-frequency components compared to tree-based gating.

A Comparative Performance and Ablation Analysis

This section presents the same tables shown in the main paper, with added standard deviation error bars.

Dataset	LogReg	DeepTLF	TabM	TabPFN	XGBoost	CatBoost	MLP	SS AE	TANDEM
CP	0.4927 \pm 0.0138	0.5438 \pm 0.0390	0.5827 \pm 0.0283	0.5998 \pm 0.0294	0.5822 \pm 0.0298	0.5401 \pm 0.0114	0.5792 \pm 0.0291	0.6505 \pm 0.0303	0.6779 \pm 0.0325
MT	0.7521 \pm 0.0252	0.6757 \pm 0.0529	0.7720 \pm 0.0252	0.8139 \pm 0.0216	0.7868 \pm 0.0215	0.7929 \pm 0.0166	0.7798 \pm 0.0213	0.7996 \pm 0.0255	0.8180 \pm 0.0228
OG	0.6228 \pm 0.0134	0.3712 \pm 0.1913	0.6228 \pm 0.0132	0.6514 \pm 0.0125	0.6136 \pm 0.0115	0.6363 \pm 0.0116	0.6321 \pm 0.0125	0.6500 \pm 0.0118	0.6870 \pm 0.0123
PW	0.9373 \pm 0.0315	0.9281 \pm 0.0161	0.9384 \pm 0.0152	0.9477 \pm 0.0136	0.9353 \pm 0.0137	0.9327 \pm 0.0296	0.9315 \pm 0.0152	0.9353 \pm 0.0173	0.9618 \pm 0.0143
AD	0.7992 \pm 0.0221	0.7645 \pm 0.0544	0.8115 \pm 0.0207	0.8199 \pm 0.0200	0.7875 \pm 0.0216	0.8019 \pm 0.0177	0.8006 \pm 0.0260	0.8202 \pm 0.0220	0.8200 \pm 0.0217
ALB	0.6065 \pm 0.0265	0.5512 \pm 0.0349	0.6196 \pm 0.0228	0.6494 \pm 0.0324	0.6096 \pm 0.0234	0.6354 \pm 0.0248	0.5929 \pm 0.0347	0.6497 \pm 0.0322	0.7038 \pm 0.0283
BM	0.8325 \pm 0.0155	0.6306 \pm 0.1093	0.8132 \pm 0.0292	0.8241 \pm 0.0287	0.7991 \pm 0.0300	0.8171 \pm 0.0246	0.7913 \pm 0.0301	0.8141 \pm 0.0218	0.8233 \pm 0.0178
CO	0.4624 \pm 0.0086	0.4881 \pm 0.0250	0.5049 \pm 0.0172	0.5485 \pm 0.0160	0.5326 \pm 0.0167	0.4975 \pm 0.0111	0.4963 \pm 0.0143	0.5007 \pm 0.0154	0.5491 \pm 0.0150
CC	0.6169 \pm 0.0249	0.6302 \pm 0.0370	0.6218 \pm 0.0394	0.6451 \pm 0.0374	0.6780 \pm 0.0362	0.6690 \pm 0.0240	0.7190 \pm 0.0358	0.6833 \pm 0.0248	0.7331 \pm 0.0251
EL	0.6617 \pm 0.0313	0.6424 \pm 0.0274	0.6610 \pm 0.0315	0.7723 \pm 0.0279	0.7668 \pm 0.0292	0.7654 \pm 0.0249	0.7525 \pm 0.0296	0.7429 \pm 0.0309	0.6940 \pm 0.0240
HE	0.4593 \pm 0.0119	—	—	—	0.4590 \pm 0.0173	0.4853 \pm 0.0111	0.4448 \pm 0.0184	0.5250 \pm 0.0132	0.5462 \pm 0.0140
HI	0.5454 \pm 0.0257	0.4965 \pm 0.0284	0.5532 \pm 0.0259	0.6499 \pm 0.0248	0.6096 \pm 0.0266	0.6069 \pm 0.0309	0.6035 \pm 0.0254	0.6179 \pm 0.0326	0.6459 \pm 0.0333
JA	0.5105 \pm 0.0166	0.4649 \pm 0.0292	0.5743 \pm 0.0248	0.5986 \pm 0.0254	0.5409 \pm 0.0262	0.5413 \pm 0.0159	0.5417 \pm 0.0249	0.5168 \pm 0.0195	0.5660 \pm 0.0212
NU	0.4833 \pm 0.0221	0.4932 \pm 0.0153	0.5221 \pm 0.0144	0.4333 \pm 0.0170	0.5308 \pm 0.0141	0.5231 \pm 0.0177	0.5517 \pm 0.0151	0.6092 \pm 0.0320	0.6545 \pm 0.0355
RS	0.6992 \pm 0.0186	0.6590 \pm 0.0308	0.7886 \pm 0.0272	0.7554 \pm 0.0283	0.7057 \pm 0.0289	0.7328 \pm 0.0169	0.7243 \pm 0.0273	0.7016 \pm 0.0253	0.7576 \pm 0.0289
VO	0.4624 \pm 0.0086	0.3966 \pm 0.0340	0.4790 \pm 0.0306	0.5082 \pm 0.0313	0.4736 \pm 0.0304	0.4975 \pm 0.0111	0.5400 \pm 0.0298	0.5154 \pm 0.0128	0.5220 \pm 0.0138
Mean Accuracy	0.6257	0.5824	0.6577	0.6812	0.6507	0.6612	0.6551	0.6708	0.6975
Mean Rank	6.94	8.53	5.13	2.93	5.38	4.81	5.31	3.63	1.81

Table A.1: Accuracy results for 8 baseline models and TANDEM, across 19 datasets. Bold indicates the best result per dataset.

Dataset	SS AE + Gating	OSDT AE + Gating	TANDEM (no gate)	TANDEM (no LRS + Alignment)	TANDEM
CP	0.6602 \pm 0.0318	0.6321 \pm 0.0303	0.6740 \pm 0.0324	0.6740 \pm 0.0324	0.6779 \pm 0.0325
MT	0.8096 \pm 0.0255	0.7639 \pm 0.0240	0.8117 \pm 0.0214	0.8014 \pm 0.0206	0.8180 \pm 0.0228
OG	0.6740 \pm 0.0120	0.6601 \pm 0.0356	0.6659 \pm 0.0152	0.6671 \pm 0.0132	0.6870 \pm 0.0123
PW	0.9353 \pm 0.0181	0.8259 \pm 0.0310	0.9464 \pm 0.0129	0.9485 \pm 0.0141	0.9618 \pm 0.0143
AD	0.8085 \pm 0.0232	0.7688 \pm 0.0585	0.8111 \pm 0.0223	0.7990 \pm 0.0216	0.8200 \pm 0.0217
ALB	0.6797 \pm 0.0320	0.6652 \pm 0.0284	0.6763 \pm 0.0284	0.6800 \pm 0.0277	0.7038 \pm 0.0283
BM	0.8141 \pm 0.0237	0.7552 \pm 0.0199	0.8138 \pm 0.0214	0.8059 \pm 0.0231	0.8233 \pm 0.0178
CO	0.5207 \pm 0.0170	0.5373 \pm 0.0499	0.5220 \pm 0.0148	0.5164 \pm 0.0163	0.5491 \pm 0.0150
CC	0.6900 \pm 0.0284	0.6892 \pm 0.0523	0.7434 \pm 0.0254	0.7188 \pm 0.0265	0.7331 \pm 0.0251
EL	0.7429 \pm 0.0332	0.6149 \pm 0.0444	0.7134 \pm 0.0252	0.7001 \pm 0.0276	0.6940 \pm 0.0240
HE	0.5350 \pm 0.0136	0.4416 \pm 0.0656	0.5315 \pm 0.0139	0.5193 \pm 0.0142	0.5462 \pm 0.0140
HI	0.6279 \pm 0.0330	0.6621 \pm 0.0528	0.6300 \pm 0.0299	0.6350 \pm 0.0274	0.6459 \pm 0.0333
JA	0.5468 \pm 0.0234	0.5012 \pm 0.0465	0.5529 \pm 0.0175	0.5541 \pm 0.0204	0.5660 \pm 0.0212
NU	0.6092 \pm 0.0342	0.6140 \pm 0.0393	0.6325 \pm 0.0318	0.6377 \pm 0.0307	0.6545 \pm 0.0355
RS	0.7016 \pm 0.0260	0.7252 \pm 0.0274	0.7311 \pm 0.0252	0.7335 \pm 0.0270	0.7576 \pm 0.0289
VO	0.5254 \pm 0.0142	0.4334 \pm 0.0231	0.4741 \pm 0.0126	0.4983 \pm 0.0140	0.5220 \pm 0.0138
Mean Accuracy	0.6801	0.6431	0.6831	0.6768	0.6975
Mean Rank	2.75	3.56	2.38	2.88	1.31

Table A.2: Accuracy results for TANDEM with encoder and gating variants across 19 datasets. Bold indicates the best result per dataset.

Dolan–Moré curves: Figure 4 in the main paper presents Dolan–Moré performance profiles, which show the proportion of datasets where each model achieves performance within a factor τ of the best-performing model. Unlike rank-based comparisons, these curves capture both accuracy and robustness, providing a more complete view of model consistency across diverse tasks. Higher curves indicate models that maintain strong performance across a larger share of datasets, even if not ranked first.

B Dataset Details

Table B.1 provides an overview of the datasets used in our experiments, including the number of samples, input features, and target classes. This supplements the dataset summary introduced in Section 5.1 of the main paper.

Dataset (Acronym)	#Samples (N)	#Features (F)	#Classes (C)
Click_prediction_small (CP)	9200	27	2
MagicTelescope (MT)	19020	11	2
Otto-Group-Product-Classification-Challenge (OG)	16400	93	5
PhishingWebsites (PW)	9200	30	2
adult (AD)	9200	14	2
albert_categorical (ALB)	9200	25	2
bank-marketing (BM)	45200	16	2
covertime (CO)	283300	54	6
default-of-credit-card-clients_categorical (CC)	5200	23	2
electricity (EL)	19240	11	2
helen (HE)	36260	27	13
higgs (HI)	470080	28	2
jannis (JA)	83600	54	4
numera128.6 (NU)	9200	119	2
road-safety_categorical (RS)	363240	67	2
volkret (VO)	14800	181	8
pol (PO)	15000	48	2
california (CA)	20634	8	2
eye_movements (EM)	10935	27	3

Table B.1: Dataset statistics: number of samples (N), features (F), and target classes (C).

C Experimental Setup

We follow the training and evaluation protocol described in Section 5.2 of the main paper. All experiments use 2,000 unlabeled and 400 labeled examples per class, with 100 pretraining epochs and batch size 128. During supervised training, encoders were frozen for the first 25 epochs and then fine-tuned for an additional 25 epochs. Gating modules were frozen throughout.

All experiments were conducted on a local machine equipped with an Intel i7 CPU, 16GB RAM, and an NVIDIA RTX 3060 GPU (12GB).

D Model Architectures and Training

This section outlines the architectural components used in our experiments, following the design described in Sections 4.1 and 4.3 of the main paper. The final embedding dimensionality is determined by the depth parameter L of the OSDT encoder, which was selected separately for each dataset based on validation performance.

- **Neural Encoder:** 4-layer MLP with BatchNorm and Leaky ReLU activations. Hidden dimensions are chosen to match the embedding size dictated by the OSDT encoder.
- **OSDT Encoder:** Ensemble of oblivious soft decision trees with fixed depth L ; each tree outputs a soft assignment over 2^L leaves. The mean-aggregated output defines the embedding dimension.
- **Neural Decoder:** Mirrors the architecture and dimensionality of the neural encoder.
- **Gating Network:** 2-layer MLP with tanh activation and hard-sigmoid output; used for per-sample feature selection.
- **Fine-tuning MLP:** A single-layer fully connected classifier trained on top of the encoder for downstream supervised evaluation.

E Hyperparameter Tuning

Table E.1 summarizes the hyperparameter search space used in Optuna-based tuning (50 trials per model, per dataset). This supports the configuration described in Section 5.2 of the main paper. Gating components, when applicable, were tuned separately.

Table E.1: Hyperparameter search space for all models. Each model was optimized using Optuna with 50 trials. Gating components (when applicable) were tuned separately.

Model	Hyperparameter Search Space	Notes
TabNet	learning rate: $[10^{-4}, 10^{-1}]$ mask type: {sparsemax, entmax} scheduler step size: [5–20], gamma: [0.8–0.95] batch size: {512, 1024, 2048} virtual batch size: {64, 128, 256} max epochs: [10–100], patience: [5–20]	Pretrained with TabNetPretrainer before supervised fine-tuning
MLP	learning rate: $[10^{-4}, 10^{-2}]$ (log scale) hidden sizes: [64–256] batch size: {32, 64, 128} epochs: [10–50]	4-layer MLP with ReLU activations
XGBoost	max depth: [3–10] gamma: $[10^{-8}, 1]$ subsample, colsample bytree: [0.5, 1.0] reg alpha, reg lambda: $[10^{-8}, 10]$	<code>use_label_encoder=False</code>
CatBoost	depth: [2–10] learning rate: $[10^{-3}, 0.3]$ iterations: [100–300]	<code>verbose=0</code>
Logistic Regression	C : $[10^{-4}, 10]$ (log-uniform)	<code>max_iter=5000</code>
DeepTFL	n estimators: [10–100] max depth: [2–7] dropout: [0.0–0.3] n layers: [1–3]	Internal tree-based layers
TabM	n blocks: [2–4] d block: [64–256] dropout: [0.0–0.3] k : [8–64] learning rate: $[10^{-4}, 5 \times 10^{-3}]$ (log)	–
TabPFN	n estimators: {1, 2, 4, 8} softmax temperature: [0.5, 1.0] balance probabilities: {True, False} average before softmax: {True, False}	Pretrained; limited tuning allowed
SS-AE	learning rate: $[10^{-4}, 10^{-2}]$ weight decay: $[10^{-5}, 10^{-2}]$ depth: [3–10]	Encoder uses powers-of-two hidden sizes. Gating tuned separately.
TANDEM	learning rate: $[10^{-4}, 10^{-2}]$ weight decay: $[10^{-5}, 10^{-2}]$ depth: [3–10] num trees: [2–16] gating optimizer, activation, learning rate	Shared decoder. Gating network tuned separately.
Gating Network	hidden size: [32–128] learning rate: $[10^{-4}, 10^{-2}]$ weight decay: $[10^{-5}, 10^{-2}]$	Used in SS-AE and TANDEM for per-sample masking.