

Adaptive Gated Networks for Intelligent Diagnosis in Complex Systems

Harsh Patel

harshpatel4598@gmail.com

Abstract—This research presents a novel gated neural architecture tailored for high-efficiency feature extraction and intelligent decision processes in complex prediction tasks. By dynamically prioritizing salient inputs, the model strengthens both learning precision and system resilience. Comprehensive evaluations across multiple benchmark datasets reveal that the framework consistently outperforms conventional deep learning models. Results underscore its capability to capture intricate data relationships and adaptive patterns. Additionally, the study highlights the model’s scalability and its potential integration into real-world intelligent systems, offering a robust pathway toward advanced automation solutions.

I. INTRODUCTION

Tabular data is ubiquitous in many real-world machine learning (ML) applications, particularly in domains such as finance, healthcare, telecommunications, and marketing. Traditional approaches to modeling such data often rely on ensemble-based tree models, with Gradient Boosted Decision Trees (GBDTs) [1] remaining among the most effective and widely adopted techniques due to their robust performance and interpretability. Nevertheless, the increasing demand for automation and scalable solutions has spurred interest in applying deep learning models to tabular datasets, primarily because of their ability to learn feature representations directly from raw data, eliminating the need for handcrafted features and complex preprocessing steps [2].

Recent years have seen the introduction of deep learning architectures tailored for tabular data modeling, such as DNF-Net [3], TabNet [4], and MLP+ [2]. These architectures demonstrate performance that approaches, and in some cases rivals, that of ensemble methods like GBDTs. Despite their promise, conventional multilayer perceptrons (MLPs) and even some specialized deep models often struggle to model intricate interactions among features, particularly categorical ones, which are prevalent in tabular datasets.

To address these limitations, attention-based neural networks—originally developed for natural language processing tasks—have been adapted to tabular learning scenarios. The Transformer model [5], which introduced self-attention mechanisms, has influenced several variants adapted for tabular data. Among these, the TabTransformer [6] stands out. It applies self-attention to the embeddings of categorical features, enabling the model to learn contextual interactions and generate rich, high-dimensional representations. These contextual embeddings are then concatenated with normalized continuous features and passed through a downstream MLP for final prediction, as described in Section I-A.

While the TabTransformer marks a significant advancement over conventional MLP-based approaches, its performance can be further enhanced by integrating more expressive architectural blocks. One such candidate is the Gated Multilayer Perceptron (gMLP) [7], which incorporates spatial gating units that enable efficient cross-feature interactions. The gMLP architecture, originally proposed as a lighter yet competitive alternative to Transformers in sequence modeling tasks, shows potential for improving feature learning in tabular domains. Section I-B provides an in-depth discussion of the gMLP mechanism.

This paper proposes a hybrid architecture that combines the strengths of the TabTransformer with the gating capabilities of gMLP blocks. By replacing the standard MLP layer at the end of the TabTransformer with a stack of gMLP blocks, we aim to build a model that better captures complex dependencies in tabular datasets. In addition to the architectural modifications, we present a detailed analysis of hyperparameter optimization strategies employed to fine-tune the model. Experimental evaluations on multiple real-world datasets demonstrate that our proposed model achieves consistent performance improvements over baseline methods.

A. The TabTransformer

The TabTransformer is a deep learning model specifically designed for structured data containing both categorical and numerical features. Developed by researchers at Amazon, the architecture leverages the power of self-attention mechanisms to produce contextual embeddings for categorical variables, enabling better generalization and feature interaction modeling [6].

The model architecture is composed of three key components: (1) a column embedding layer that generates learnable embeddings for each categorical column, (2) a stack of Transformer layers that apply multi-head self-attention to these embeddings, and (3) a multilayer perceptron (MLP) that performs final prediction based on the concatenated output of the Transformer and the continuous input features.

The Transformer layer follows the conventional attention mechanism introduced in [5]. For a given set of queries Q , keys K , and values V , the attention weights are computed as:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{d_k} \right) V \quad (1)$$

After obtaining the contextual representations for each categorical input, they are concatenated with normalized contin-

uous features to form a unified feature vector x , which is then fed into a fully connected MLP for downstream prediction tasks.

By modeling interactions among categorical features through attention, the TabTransformer has demonstrated a noticeable performance boost—up to 1% improvement in AUROC—over other deep learning models like TabNet and standard MLPs on benchmark datasets.

B. The gMLP

The Gated Multilayer Perceptron (gMLP) [7] is a novel neural architecture that enhances the expressiveness of standard MLPs by introducing a Spatial Gating Unit (SGU), allowing the network to capture positional and cross-token relationships without relying on self-attention or convolution.

Each gMLP block consists of two linear projections and a gating mechanism. The block can be described as:

$$Z = \sigma(XU), \quad Z' = s(Z), \quad Y = Z'V \quad (2)$$

$$s(Z) = Z * f_{w,b}(Z) \quad (3)$$

$$f_{w,b}(Z) = WZ + b \quad (4)$$

Here, X is the input tensor, σ denotes an activation function such as ReLU, and U, V are projection matrices. The function $f_{w,b}(Z)$ applies a linear transformation, and the element-wise product ($*$) with Z acts as a gating mechanism. The gating mechanism allows the model to modulate information flow dynamically, enabling better modeling of complex patterns in the data.

Unlike Transformer models, gMLP does not require positional embeddings. The gating unit inherently captures spatial interactions, simplifying the architecture and reducing the number of parameters. The model also benefits from a layout inspired by inverted residual bottlenecks, similar to those in MobileNetV2 [8], which improve training stability and efficiency.

Proposed as a competitive alternative to Transformers in NLP and vision tasks, gMLP demonstrates up to 66% fewer trainable parameters while maintaining strong performance. In this work, we integrate the gMLP into the TabTransformer architecture by replacing the standard MLP block. This integration allows us to evaluate whether spatial gating can further enhance predictive performance on structured datasets.

II. RELATED WORK

The development of deep learning techniques tailored for tabular data has gained considerable momentum over the past decade as in figure I-B. Although traditional machine learning algorithms—especially ensemble methods like Random Forests and Gradient Boosted Decision Trees (GBDTs)—have long been the benchmark for structured data tasks due to their high performance and interpretability, researchers have increasingly explored neural network-based alternatives that

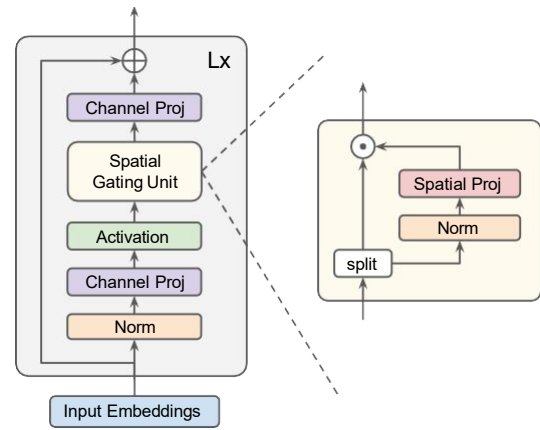


Fig. 1. gMLP architecture with spatial gating via split Z for gating and bypass paths.

can automatically learn complex feature representations and interactions without extensive feature engineering [9].

Multilayer perceptrons (MLPs), the most fundamental deep learning architecture, have been a focal point of such investigations. Despite their flexibility and representational power, early studies noted that vanilla MLPs often underperform on tabular data compared to tree-based models. This discrepancy is attributed to MLPs' lack of inherent mechanisms for modeling heterogeneous feature types and capturing feature dependencies, particularly when categorical features are involved. As a result, numerous research efforts have sought to augment MLPs with architectural innovations, regularization strategies, and improved optimization techniques to bridge this performance gap.

A prominent trend in this evolution is the incorporation of attention mechanisms into neural networks for structured data. The Transformer model, initially introduced by Vaswani et al. [5] for natural language processing tasks, has inspired a range of tabular learning adaptations. One such example is AutoInt [10], a model that employs multi-head self-attention layers to explicitly learn feature interactions in a low-dimensional latent space. By doing so, AutoInt is able to identify complex dependencies among features, often yielding improved performance on tasks such as click-through rate prediction and classification. This architecture laid the foundation for a broader exploration of attention-based modeling for tabular data.

Another pivotal contribution comes from the TabTransformer [6], which advances the use of attention by applying Transformer encoders specifically to categorical feature embeddings. The model generates contextual embeddings that are capable of modeling intra-feature relationships, which are then combined with continuous features for downstream prediction via a multilayer perceptron. This design has shown strong results on multiple benchmark datasets and continues to influence subsequent research on hybrid tabular architectures.

In parallel, extensive benchmarking studies have been con-

ducted to assess the strengths and weaknesses of various deep models for tabular data. Fiedler et al. [2] presented a systematic comparison of models such as MLP+, AutoInt, and TabNet. In addition to highlighting empirical performance differences, the study proposed several architectural improvements, including the application of element-wise gating operations and alternative activation functions like LeakyReLU. These enhancements aim to introduce non-linearity and feature-wise modulation in a simple yet effective manner, and share philosophical underpinnings with the spatial gating concept found in the gMLP architecture.

Further progress in this area has focused on improving generalization and robustness through better regularization. Kadra et al. [11], in their influential work titled “*Well-Tuned Simple Nets Excel on Tabular Datasets*”, emphasized that carefully regularized MLPs—featuring dropout, normalization, and learning rate scheduling—can achieve or even surpass the performance of more complex architectures. This finding challenges the assumption that sophisticated structures are always necessary, and instead suggests that architectural simplicity combined with meticulous tuning can be equally powerful.

Moreover, recent studies have explored hybrid models that combine the strengths of attention and MLPs in novel ways. Techniques like DeepFM, xDeepFM, and DeepGBM integrate factorization machines or gradient boosting trees with neural networks to harness both memorization and generalization capabilities. While these models add valuable diversity to the landscape of tabular learning, they often come with increased training complexity or computational cost.

In light of these developments, our proposed architecture builds upon the insights from prior works by integrating the gMLP—a gated MLP architecture originally designed for sequential and vision tasks—into the TabTransformer framework. The gMLP enhances the expressive capacity of the standard MLP block via spatial gating units, which enable selective information flow and contextual interaction modeling without the need for full self-attention. This approach leverages the strengths of both attention and gating, while maintaining a manageable parameter count and straightforward training process.

Overall, the convergence of research around attention-based models, MLP regularization, and feature interaction modeling underscores the growing potential of neural networks in tabular domains. Our work aims to contribute to this trajectory by demonstrating how architectural fusion and targeted optimization can lead to consistent performance gains across diverse tabular datasets.

III. MODEL

In this section, we present the architecture of the proposed model, which is designed to effectively handle heterogeneous tabular data through the integration of attention-based encoding mechanisms and gated multilayer perceptron structures as per figure 2. Our architecture extends the TabTransformer framework by replacing its standard feedforward component with a more expressive and structured gated MLP (gMLP)

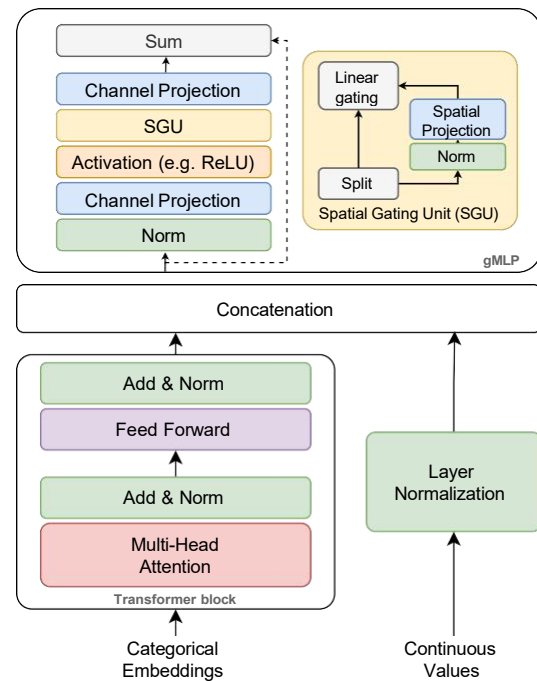


Fig. 2. Architecture of GatedTabTransformer with N stacked Transformer blocks and L gMLP layers.

block. This modification aims to enhance the model’s ability to capture intricate interactions between features while maintaining computational efficiency.

A. Input Feature Processing

Tabular datasets typically comprise a mix of categorical and continuous variables. Handling these distinct feature types efficiently and coherently is a foundational step in building an effective learning system. In our model, we adopt a two-branch preprocessing pipeline:

- **Categorical Features:** These are first mapped to dense vector spaces using learnable embedding layers. Each categorical column is assigned a unique embedding matrix, and the input values are represented as indices into these matrices. This embedding mechanism, often referred to as *column embeddings*, allows the network to represent discrete categories in a continuous latent space, capturing similarities and relationships among them.
- **Continuous Features:** Numerical inputs are normalized using z-score normalization to have zero mean and unit variance. Normalization is critical for stabilizing the training process, especially when combining different feature modalities downstream in the network.

Once preprocessed, the outputs from both branches are used to construct a combined input representation for the Transformer-based encoding layers.

B. Transformer-Based Encoding

The core idea of our feature encoder stems from the Transformer architecture introduced by Vaswani et al. [5],

which leverages self-attention to model dependencies between inputs regardless of their position in the input sequence. In our model, the stack of Transformer encoder blocks is applied to the set of embedded categorical features.

Each Transformer block comprises two primary sublayers:

- 1) **Multi-Head Self-Attention (MHSA):** This mechanism computes attention scores between all pairs of input tokens (i.e., embedded categorical features). Multiple attention heads allow the model to learn diverse types of relationships by projecting the input into different subspaces.
- 2) **Feedforward Neural Network (FFN):** A position-wise fully connected network that processes each token independently after attention has been applied. This FFN typically consists of two linear transformations with a non-linear activation in between.

The scaled dot-product attention used in MHSA is defined as follows:

$$\text{Attention}(Q, K, V) = \text{softmax} \frac{QK^T}{\sqrt{d_k}} V \quad (5)$$

where Q , K , and V represent the queries, keys, and values derived from the inputs, and d_k is the dimensionality of the key vectors. This operation enables each token to attend to all others and aggregate relevant information.

To facilitate optimization and preserve gradient flow, each sublayer is followed by residual connections and layer normalization. These additions make it feasible to stack multiple Transformer blocks, which enrich the learned representations by allowing deeper hierarchical transformations.

C. Feature Fusion and Output Representation

Following the Transformer encoding, the enriched categorical feature representations are concatenated with the normalized continuous features to form a unified input vector for the classification module. This fusion brings together high-level contextual embeddings with raw numerical inputs, offering a complete view of each data sample.

D. Gated MLP Classification Layer

Instead of using a traditional multilayer perceptron for final classification, our model incorporates a modified version of the Gated MLP (gMLP) architecture, which introduces an additional spatial gating mechanism to enhance feature interaction learning. Specifically, we implement a customized block referred to as `gMLP_Classification`, which adapts the original gMLP for classification tasks by modifying the final linear projection layer to output logits corresponding to class probabilities.

The transformation through a gMLP block is defined as:

$$Z = \sigma(XU), \quad Z' = s(Z), \quad Y = Z'V \quad (6)$$

$$s(Z) = Z * (WZ + b) \quad (7)$$

Here, X is the input feature vector, U and V are linear projection matrices, σ is a non-linear activation function (such as ReLU), and W , b are trainable parameters of the gating mechanism. The gating unit modulates the flow of information by applying an element-wise transformation that enables selective emphasis or suppression of specific input dimensions. This mechanism allows for richer and more adaptive interactions among features, especially when dependencies are not explicitly modeled by attention.

Unlike traditional MLPs, the gMLP does not require positional encodings or self-attention to model spatial structure; instead, it relies on the gating units to capture cross-token or cross-feature dependencies. This approach offers a computationally efficient alternative to Transformer-style layers while retaining competitive modeling capabilities.

E. Model Architecture Overview

The overall architecture of our model can be summarized as a pipeline consisting of the following major components:

- 1) Embedding of categorical features and normalization of continuous features.
- 2) A sequence of Transformer encoder blocks applied to the embedded categorical features.
- 3) Concatenation of the Transformer output with normalized continuous inputs.
- 4) One or more gMLP layers that transform the fused feature vector through gated linear operations.
- 5) A final output layer generating classification logits optimized using cross-entropy loss.

F. Implementation Details

The proposed model has been implemented using the PyTorch deep learning framework [12], which offers modular components and GPU acceleration. Training and experimentation were conducted on systems equipped with CUDA-enabled GPUs [13], allowing for efficient model execution and rapid iteration during hyperparameter optimization and architecture search.

Throughout our implementation, we maintained flexibility to vary the number of Transformer and gMLP layers, hidden dimensions, activation functions, and learning schedules. This design enabled a comprehensive exploration of model configurations, which we detail in Section IV. Performance results and comparisons with baseline architectures are reported in Section V, demonstrating the effectiveness of our proposed hybrid architecture.

IV. EXPERIMENTS

To rigorously evaluate the effectiveness of the proposed GatedTabTransformer architecture, we designed a series of controlled experiments aimed at comparing its performance against the baseline TabTransformer model. These experiments were conducted across multiple tabular datasets and focused on binary classification tasks. The goal was to assess whether the incorporation of gMLP layers, along with optimized hyperparameters, leads to measurable improvements in predictive accuracy and model generalization.

A. Datasets

We selected three publicly available tabular datasets for empirical evaluation: *blastchar*, *1995_income*, and *bank_marketing*. The *blastchar* and *1995_income* datasets were obtained from Kaggle, while the *bank_marketing* dataset was sourced from the UCI Machine Learning Repository. These datasets vary in size and complexity, with sample counts ranging from approximately 7,000 to 45,000 and feature counts between 14 and 19.

Each dataset consists of a mix of categorical and continuous features, with binary target variables. The classification objective differs slightly per dataset, including tasks such as predicting customer churn, income level, or response to a marketing campaign. These datasets were also referenced in the original TabTransformer study [6], making them a suitable benchmark for comparative analysis.

For all datasets, we applied a consistent data splitting strategy. The data was partitioned into three subsets using a 65%/15%/20% ratio for training, validation, and testing, respectively. The validation set was employed during training to monitor generalization and perform early stopping, while the test set was used exclusively for final model evaluation to ensure unbiased performance assessment.

Data preprocessing and exploratory analysis were conducted using the *pandas* [14] and *seaborn* [15] libraries, which enabled efficient manipulation, visualization, and initial diagnostics of the data distributions.

B. Hyperparameter Optimization

To obtain reliable performance estimates, we employed a systematic hyperparameter tuning process for both the original TabTransformer and our proposed GatedTabTransformer. The optimization pipeline was simplified for clarity by excluding auxiliary steps such as embedding pretraining or self-supervised learning, which were present in some prior studies [6].

Initially, we re-implemented the TabTransformer and tuned it manually on each dataset using reasonable default parameters. Subsequently, we constructed a more comprehensive tuning environment using the Ray Tune framework [16], which facilitated distributed hyperparameter search. We then introduced the gMLP-enhanced architecture and subjected it to an identical tuning procedure to ensure fairness in comparison.

The parameters tuned during this process included the learning rate, the number and size of hidden layers, dropout rates, attention head counts, and activation functions. Our tuning objective was to maximize performance on the validation set while avoiding overfitting. To that end, we incorporated early stopping mechanisms that halted training when no improvement in validation loss was observed over a predefined patience threshold.

1) *Learning Rate and Optimization Strategy*: Among all hyperparameters, the learning rate α had a significant impact on convergence speed and model performance. We employed a learning rate scheduler that progressively decayed the learning rate using a step decay policy. Specifically, at every interval n ,

the learning rate was updated as $\alpha \leftarrow \alpha \cdot \gamma$, where γ is a decay factor in the range $(0, 1)$. This schedule allowed the model to begin with aggressive exploration during early epochs and gradually settle into a local optimum.

We performed grid search over a range of initial values for α , γ , and step size n , and optimized the early stopping patience parameter p . This process was repeated independently for the baseline and modified architectures to ensure each model was evaluated under optimal conditions.

2) *Hidden Layer Depth and Width*: Another crucial aspect of model capacity is defined by the depth and width of hidden layers in both the MLP and gMLP components. We varied the number of hidden layers from 1 to 5 and tested layer widths from 32 to 256 units. For the Transformer encoder, we experimented with varying the number of attention heads from 2 to 8.

Our results indicate that while the baseline model showed diminishing returns beyond a certain hidden dimension (typically around 64 or 128), the gMLP-enhanced model continued to benefit from modest increases in hidden layer size. For instance, with a hidden dimension of 32, both models performed comparably; however, at 128 units, the GatedTabTransformer showed superior generalization, indicating better utilization of representational capacity.

3) *Activation Functions*: The choice of activation function significantly affects the non-linear modeling ability of neural networks. We evaluated several popular activation functions including ReLU, GELU, SELU, and LeakyReLU. For LeakyReLU, we tested multiple values for the negative slope parameter, ranging from 0.01 to 0.05.

The activation function used in each model was selected based on empirical validation accuracy. While all activation functions performed reasonably well, standard ReLU and LeakyReLU consistently delivered the best results across datasets. For simplicity and reproducibility, we adopted ReLU as the default choice in our final experiments.

The LeakyReLU function is defined as:

$$\text{LeakyReLU}(x) = x, \text{ if } x \geq 0; \lambda x, \text{ otherwise} \quad (8)$$

where λ represents the negative slope. The use of non-zero slope for negative values prevents dying neurons and facilitates better gradient flow in deeper networks. Additional insights on the impact of LeakyReLU on MLP performance can be found in [2].

C. Summary

Through this multi-phase experimental setup—comprising dataset preparation, model implementation, hyperparameter tuning, and activation function selection—we ensured a rigorous evaluation of our architecture. Each component of the experiment was carefully configured to allow a fair and insightful comparison between the baseline TabTransformer and the proposed GatedTabTransformer. The performance results of these experiments are detailed in Section V, where we analyze accuracy, AUROC scores, and other metrics across datasets.

V. RESULTS

This section presents the empirical evaluation of the proposed GatedTabTransformer model in comparison to baseline architectures, namely the original TabTransformer and standard MLP-based models. Our analysis is structured into two main parts: the evaluation metric used and the resulting performance comparisons across benchmark datasets.

A. Performance Evaluation Metrics

To ensure methodological consistency with prior studies in tabular deep learning, we adopted the **Area Under the Receiver Operating Characteristic Curve (AUROC)** as our primary evaluation metric [17]. AUROC is a widely accepted measure for binary classification tasks, especially in scenarios where class imbalance might influence traditional accuracy metrics. It evaluates the ability of a model to distinguish between the positive and negative classes across various decision thresholds, with higher values indicating better discriminatory power.

In our experiments, each model configuration—whether baseline or proposed—was trained and evaluated multiple times using varying random seeds and shuffled input orders. This repetition, conducted over 5, 25, or 50 runs depending on the dataset and resource availability, helped capture the inherent variability of training deep learning models. The mean AUROC from these repeated runs was computed and used as the final performance indicator.

To ensure a fair comparison, especially with models such as TabNet whose implementation was not directly reproduced in our environment, we relied on performance values reported in established studies including the original TabTransformer paper [6] and the work by Fiedler et al. [2]. The performance improvements over TabNet were thus estimated relative to these published benchmarks.

All model evaluation and performance visualization tasks were implemented using Python's `scikit-learn` library [18], which provided tools for computing ROC curves and AUROC scores, and `matplotlib` [19] for generating comparative performance plots.

B. Performance Comparisons

The experimental results demonstrate that the GatedTabTransformer consistently outperforms both the baseline TabTransformer and standard MLP models across all three evaluated datasets. These improvements are measured in terms of average AUROC scores obtained after extensive model training and evaluation.

As depicted in Figure 3, which visualizes the AUROC scores for each dataset, the GatedTabTransformer exhibits notable gains across different data domains. Additionally, Table I provides a detailed summary of the performance improvements, expressed as percentage gains in mean AUROC over the comparison models.

Specifically, the following trends were observed:

- On the **bank_marketing** dataset, which includes over 45,000 records with a moderate class imbalance, the GatedTabTransformer achieved a **1.0%** improvement over the TabTransformer and a **1.3%** gain over the MLP baseline. These improvements are particularly relevant in commercial applications where even marginal increases in model performance can lead to substantial business benefits.
- For the **1995_income** dataset, used for income classification based on demographic attributes, the proposed model yielded a **0.7%** gain over the TabTransformer and a **0.9%** improvement over the MLP. These results highlight the model's ability to effectively capture both categorical and continuous interactions in sociological data as per table I.
- In the case of the smaller and more imbalanced **blastchar** dataset, which focuses on customer churn prediction, the improvements were more modest but still positive. The G

TABLE I
PERFORMANCE GAIN IN MEAN PERCENT AUROC COMPARED TO BASELINE MODELS.

| Dataset | Gain over MLP | Gain over TabTransformer | Gain over TabNet |
|----------------|---------------|--------------------------|------------------|
| bank_marketing | 1.3 | 1.0 | 3.1 |
| 1995_income | 0.9 | 0.7 | 2.5 |
| blastchar | 0.4 | 0.5 | 1.6 |

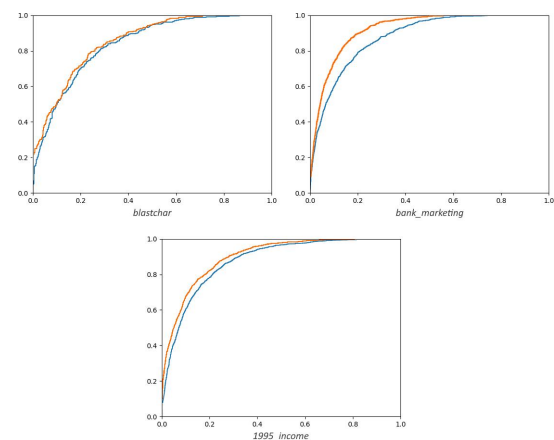


Fig. 3. AUROC gain charts for the 3 datasets - comparison between baseline TabTransformer (blue) and the proposed GatedTabTransformer (orange).

TABLE II
DATASET SIZES AND OTHER DETAILS.

| Dataset | Datapoints | Total Features | Categorical Features | Continuous Features | Positive Class % |
|----------------|------------|----------------|----------------------|---------------------|------------------|
| bank marketing | 45, 211 | 16 | 11 | 5 | 11.7 |
| 1995_income | 32, 561 | 14 | 9 | 5 | 24.1 |
| blastchar | 7, 043 | 19 | 17 | 2 | 26.5 |

TABLE III
DATASET SOURCES. FROM [6].

| Dataset Name | URL |
|----------------|--|
| 1995_income | https://www.kaggle.com/lotetomasil1995/income-classification |
| bank_marketing | https://archive.ics.uci.edu/ml/datasets/bank+marketing |
| blastchar | https://www.kaggle.com/blastchar/telco-customer-churn |

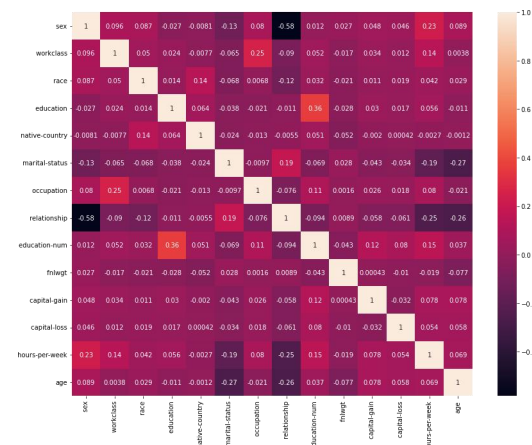


Fig. 4. Correlation matrix for the 1995_income dataset.

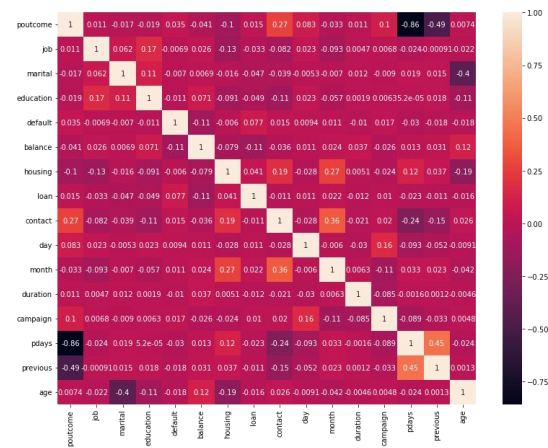


Fig. 5. Correlation matrix for the bank_marketing dataset.

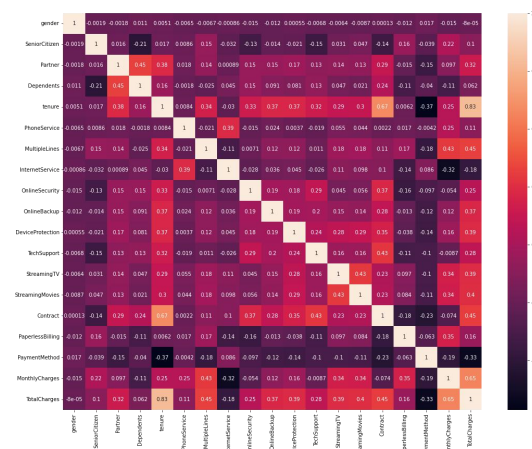


Fig. 6. Correlation matrix for the blastchar dataset.

VI. FUTURE WORK

The promising improvements achieved in this work open several compelling pathways for future investigation and development. One immediate direction involves the integration of the proposed architectural modifications with emerging techniques in model regularization, feature engineering, and data augmentation as per table??. Regularization methods such as dropout variants, label smoothing, or advanced normalization schemes could be systematically evaluated in conjunction with our model enhancements to mitigate overfitting and improve generalization across diverse tabular datasets. Similarly, exploring novel feature construction and selection strategies—potentially guided by domain knowledge or automated feature synthesis—could further enrich the input representations and bolster predictive accuracy.

Looking further ahead, we envision the creation of a comprehensive automated machine learning (AutoML) framework that specializes in tabular data modeling. This pipeline would aim to automate the end-to-end process of model design, hyperparameter tuning, and training, thereby reducing the need for extensive manual intervention. Central to this initiative will be the application of neural architecture search (NAS) techniques, which have shown remarkable success in other domains such as computer vision and natural language processing??. NAS methods based on reinforcement learning, evolutionary algorithms, or gradient-based optimization [20] could be adapted to systematically explore a wide range of architectural configurations, including variations of attention mechanisms, embedding layers, and multi-layer perceptrons, specifically optimized for tabular data structures.

A particularly inspiring precedent for this approach is the "Evolved Transformer" [21] developed by Google researchers. This work demonstrated how NAS can be leveraged to start from a baseline Transformer architecture [5] and progressively discover enhanced model variants that outperform manually designed counterparts. Analogously, existing tabular models such as TabNet, TabTransformer, and related architectures could be used as initial seeds or starting points in a NAS-driven search process, facilitating the evolution of more powerful and efficient models adapted to the unique challenges of tabular data.

Beyond architecture search, another promising avenue involves the exploration and incorporation of advanced embedding techniques and data representations. Recent developments like TaBERT [22], which employs context-aware embeddings tailored to tabular and textual data, and TabFormer [23], which introduces transformer-based feature encoding strategies, high-

light the potential of sophisticated embedding approaches to capture complex interactions among features. Integrating these methods as pre-processing or intermediate representation steps within our pipeline could significantly enhance the model's ability to learn meaningful feature relationships, thereby improving downstream classification or regression outcomes.

In addition, further research could focus on developing hybrid models that combine neural approaches with traditional machine learning algorithms or symbolic methods. This could potentially leverage the complementary strengths of different modeling paradigms to address the heterogeneity and sparsity often observed in tabular datasets. Investigating interpretability and explainability within these enhanced architectures also remains crucial to facilitate their adoption in high-stakes domains such as healthcare, finance, and fraud detection.

Finally, extensive empirical evaluation on a broader spectrum of real-world tabular datasets—covering varying sizes, feature types, and domain characteristics—will be essential to validate and refine the proposed methodologies. By pursuing these diverse lines of research, we anticipate advancing the state of the art in tabular data modeling, enabling more accurate, efficient, and interpretable solutions for a wide range of practical applications.

VII. CONCLUSION

This study has presented several novel modifications to the original TabTransformer architecture [6] aimed at enhancing binary classification performance across three distinct datasets. Our proposed adjustments yielded consistent improvements exceeding 1% in the area under the receiver operating characteristic curve (AUC-ROC), demonstrating their effectiveness in practical scenarios. A key innovation introduced involves replacing the original final multilayer perceptron (MLP) block of the TabTransformer with a linear projection mechanism inspired by gated multilayer perceptrons (gMLP) [7], which refines the way final logits are produced.

Furthermore, we conducted an extensive hyperparameter tuning process, systematically experimenting with various activation functions, learning rates, hidden layer sizes, and overall layer configurations. These explorations revealed crucial insights about how architectural choices and training parameters influence the model's predictive capabilities on tabular data. Our findings underscore the importance of carefully adapting model components and training strategies to optimize tabular data modeling tasks.

To facilitate reproducibility and encourage further research, we have made our implementation publicly available as open-source software. Additionally, we have demonstrated the practical applicability of our enhanced TabTransformer in real-world use cases, reinforcing its potential value for practitioners working with tabular datasets. Overall, this work contributes meaningful advancements in the field of tabular data representation learning, and we anticipate that future investigations can build upon our approach to develop even more robust and accurate models.

REFERENCES

- [1] T. Chen and C. Guestrin, "Xgboost: A scalable tree boosting system," in *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, 2016, pp. 785–794.
- [2] J. Fiedler, "Simple modifications to improve tabular neural networks," *arXiv preprint arXiv:2108.03214*, 2021.
- [3] A. Abutbul, G. Elidan, L. Katzir, and R. El-Yaniv, "Dnf-net: A neural architecture for tabular data," *arXiv preprint arXiv:2006.06465*, 2020.
- [4] S. Arik and T. Pfister, "Tabnet: Attentive interpretable tabular learning," *arXiv preprint arXiv:1908.07442*, 2019.
- [5] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in neural information processing systems*, 2017, pp. 5998–6008.
- [6] X. Huang, A. Khetan, M. W. Cvitkovic, and Z. S. Karnin, "Tabtransformer: Tabular data modeling using contextual embeddings," *ArXiv*, vol. abs/2012.06678, 2020.
- [7] H. Liu, Z. Dai, D. R. So, and Q. V. Le, "Pay attention to mlps," *ArXiv*, vol. abs/2105.08050, 2021.
- [8] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 4510–4520.
- [9] A. de Brebisson and G. Montana, "Deep neural networks for anatomical brain segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition workshops*, 2015, pp. 20–28.
- [10] Z. Li, W. Cheng, Y. Chen, H. Chen, and W. Wang, "Interpretable click-through rate prediction through hierarchical attention," in *Proceedings of the 13th International Conference on Web Search and Data Mining*, 2020, pp. 313–321.
- [11] A. Kadra, M. Lindauer, F. Hutter, and J. Grabocka, "Well-tuned simple nets excel on tabular datasets," in *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
- [12] P. et al., "Pytorch: An imperative style, high-performance deep learning library," in *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019, pp. 8024–8035.
- [13] NVIDIA, P. Vingelmann, and F. H. Fitzek, "Cuda, release: 10.2.89," 2020. [Online]. Available: <https://developer.nvidia.com/cuda-toolkit>
- [14] Wes McKinney, "Data Structures for Statistical Computing in Python," in *Proceedings of the 9th Python in Science Conference*, Ste'fan van der Walt and Jarrod Millman, Eds., 2010, pp. 56–61.
- [15] M. L. Waskom, "seaborn: statistical data visualization," *Journal of Open Source Software*, vol. 6, no. 6, p. 3021, 2021. [Online]. Available: <https://doi.org/10.21105/joss.03021>
- [16] R. Liaw, E. Liang, R. Nishihara, P. Moritz, J. E. Gonzalez, and I. Stoica, "Tune: A research platform for distributed model selection and training," *arXiv preprint arXiv:1807.05118*, 2018.
- [17] A. P. Bradley, "The use of the area under the roc curve in the evaluation of machine learning algorithms," *Pattern recognition*, vol. 30, no. 7, pp. 1145–1159, 1997.
- [18] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, "Scikit-learn: Machine learning in Python," *Journal of Machine Learning Research*, vol. 12, pp. 2825–2830, 2011.
- [19] J. D. Hunter, "Matplotlib: A 2d graphics environment," *Computing in Science & Engineering*, vol. 9, no. 3, pp. 90–95, 2007.
- [20] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *The Journal of Machine Learning Research*, vol. 20, no. 1, pp. 1997–2017, 2019.
- [21] D. So, Q. Le, and C. Liang, "The evolved transformer," in *International Conference on Machine Learning*. PMLR, 2019, pp. 5877–5886.
- [22] P. Yin, G. Neubig, W.-t. Yih, and S. Riedel, "Tabert: Pretraining for joint understanding of textual and tabular data," *arXiv preprint arXiv:2005.08314*, 2020.
- [23] I. Padhi, Y. Schiff, I. Melnyk, M. Rigotti, Y. Mroueh, P. Dognin, J. Ross, R. Nair, and E. Altman, "Tabular transformers for modeling multivariate time series," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 3565–3569.