

# TYTAN: Taylor-series based Non-Linear Activation Engine for Deep Learning Accelerators

Soham Pramanik<sup>1</sup>[0009-0005-1772-2901], Vimal William<sup>2</sup>[0000-0002-1069-9594], Arnab Raha<sup>3</sup>[0000-0002-8848-1069], Debayan Das<sup>4</sup>[0000-0003-1843-0124], Amitava Mukherjee<sup>5,6</sup>[0000-0003-1694-3911], and Janet L. Paluh<sup>6</sup>[0000-0002-5988-6075]

<sup>1</sup> Department of Electronics and Telecommunication Engineering, Jadavpur University, Kolkata, India

<sup>2</sup> SandLogic Technologies, Bangalore, India

<sup>3</sup> Intel Corporation, Santa Clara, CA, USA

<sup>4</sup> Department of Electronic Systems Engineering, Indian Institute of Science, Bangalore, India

<sup>5</sup> Department of Computer Science and Engineering, Amrita University, Amritapuri, Kollam, Kerala, India

<sup>6</sup> College of Nanoscale Science and Engineering, Nanobioscience, SUNY Polytechnic Institute, Albany, New York, USA

**Abstract.** The rapid advancement in AI architectures and the proliferation of AI-enabled systems have intensified the need for domain-specific architectures that enhance both the acceleration and energy efficiency of AI inference, particularly at the edge. This need arises from the significant resource constraints—such as computational cost and energy consumption—associated with deploying AI algorithms, which involve intensive mathematical operations across multiple layers. High-power-consuming operations, including General Matrix Multiplications (GEMMs) and activation functions, can be optimized to address these challenges. Optimization strategies for AI at the edge include algorithmic approaches like quantization and pruning, as well as hardware methodologies such as domain-specific accelerators. This paper proposes TYTAN: TaYlor-series based non-linear acTivAtion eNginE, which explores the development of a Generalized Non-linear Approximation Engine (GNAE). TYTAN targets the acceleration of non-linear activation functions while minimizing power consumption. The TYTAN integrates a reconfigurable hardware design with a specialized algorithm that dynamically estimates the necessary approximation for each activation function, aimed at achieving minimal deviation from baseline accuracy. The proposed system is validated through performance evaluations with state-of-the-art AI architectures, including Convolutional Neural Networks (CNNs) and Transformers. Results from system-level simulations using Silvaco’s FreePDK45 process node demonstrate TYTAN’s capability to operate at a clock frequency  $> 950MHz$ , showcasing its effectiveness in supporting accelerated, energy-efficient AI inference at the edge, which is  $\sim 2\times$  performance improvement, with  $\sim 56\%$  power reduction and  $\sim 35\times$  lower area compared to the baseline open-source NVIDIA Deep Learning Accelerator (NVDLA) implementation.

**Keywords:** Machine Learning · Domain-specific Accelerators · Algorithm/HW Co-design

## 1 Introduction

The continual advancements in Artificial Intelligence (AI) and the evolving demands for AI-assisted applications highlighted the limitations of existing infrastructures on processing such as data-intensive operations. These limitations significantly elevated the cost of accommodating modern-day AI-based applications and upgrading the existing infrastructure to support such applications could also be challenging. With the motivation of proliferating the benefits of AI, this paper focuses on optimizing the effective computation of non-linear activation functions employed in Neural Networks (NN) through state-of-the-art hardware-aware DNN algorithm optimization and algorithm-aware hardware design. The method aimed to improve the inference speed by accelerating the above-mentioned compute functions in the AI model. In addition, the methods focused on extending the support for AI-based solutions in resource-constrained environments.

### 1.1 Background

The rapid evolution of AI workloads demands optimized hardware infrastructures for modern deep learning models. A key bottleneck in AI accelerators is computing nonlinear functions like *softmax*, *sigmoid*, *GeLU*, and *layer normalization*. While multiply-accumulate (MAC) operations are well-optimized, efficient nonlinear operations implementation remains challenging.

Prior works address these challenges through approximate computation techniques. *ViTALiTy* [2] approximates *softmax* in linear attention layers using Taylor series expansion. *NN-LUT* [14] employs look-up-table (LUT)-based approximation for nonlinear functions, optimizing transformer models with precomputed values. However, LUT methods, while low-latency, suffer from high memory costs and limited scalability.

Alternative approximation techniques include *Fast Approximations of Activation Functions* [4] using *Posit arithmetic*, *TaylorNet* [10] with adaptive Taylor expansions reducing computational cost by up to 60%, and comparative studies of LUT-based, CORDIC, and Taylor series methods [9].

*UNO* [12, 13] introduces a MAC-based architecture unifying nonlinear function computation by mapping Taylor-approximated functions onto existing SIMD-based MAC arrays, eliminating dedicated processing units and enabling dynamic accuracy-energy scaling. However, UNO primarily uses low-order Taylor approximations, potentially compromising accuracy in complex workloads.

Building upon these advancements, this work integrates both algorithmic and hardware design optimizations to enhance nonlinear function mapping in AI accelerators. TYTAN achieves a better balance between accuracy, power, and area efficiency, ensuring minimal resource overhead while maintaining better performance and scalability than past attempts.

## 1.2 Motivation

A generalized and re-configurable activation engine for accelerating activation functions with minimal power consumption can be a promising method to handle the challenges addressed previously. The co-design of approximated computing and generalized hardware architecture for accelerating non-linear activation functions shows greater potential to enable ultra-low-powered AI inference at the edge. The proposed system combines re-configurable hardware design with a specialized algorithm for estimating the amount of approximation required per activation function, providing reduced power consumption with minimal deviation from the AI model’s baseline accuracy.

The mathematical limitations of the Taylor series toward the representation of discontinuous functions restrict the proposed hardware design from supporting piece-wise linear functions such as ReLU. However, ReLU requires only trivial hardware implementation (simple sign bit check), making it unsuitable for acceleration. TYTAN targets computationally intensive non-linear functions that involve expensive operations where hardware acceleration provides substantial benefits. The key contributions of this paper are summarized below:

- We propose a Taylor-series-based Non-Linear Activation Engine, TYTAN to accelerate a wide range of approximated non-linear activation functions.
- An iterative search-based algorithm is designed to estimate the amount of approximation required per activation function based on the layer placement and correlation with the AI model’s baseline accuracy.
- The proposed GNAE system design combines the TYTAN accelerator with add-ons to support a wide range of non-linear functions providing a modular and scalable solution.
- Using Silvaco’s FreePDK45 process node [8], we demonstrate a hardware accelerator for various non-linear activation functions operating at a clock frequency of  $950MHz$ , showing  $> 56\%$  power reduction,  $\sim 2\times$  performance improvement, and  $> 35\times$  lower area compared to the baseline open-source NVIDIA Deep Learning Accelerator NVDLA [6] implementation.

## 2 Methodology

The proposed architecture articulated in Fig. 1 is a software-hardware co-design for accelerating non-linear activation functions utilized in modern AI architectures. The algorithm is tailored to intake AI models from multiple frameworks and is dedicated to computing the approximation required by individual activation functions. Furthermore, the hardware can be programmed with the pre-calculated coefficients into its buffer to expedite the computing of the appropriate nonlinear functions. The detailed infrastructure of software and hardware sub-systems utilized in the TYTAN is described in the following sections.

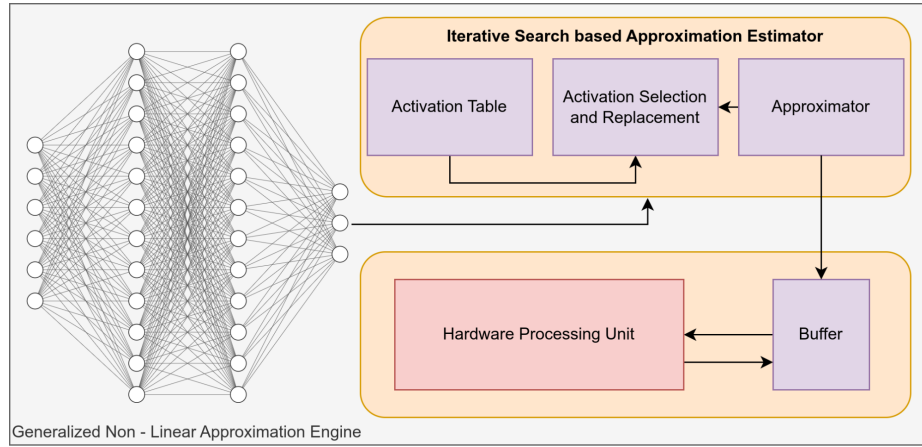


Fig. 1: Software/Hardware Co-design for Generalized Non-linear Approximation Engine

## 2.1 Algorithm Development

The software component of the co-design includes a well-structured, iterative search algorithm that determines the optimal number of terms,  $n$ , required in the Taylor series expansion to approximate the non-linear functions used in AI architectures. Additionally, the algorithm incorporates a simulated version of each approximated non-linear function, as described in Eqs. 10, 11, 12, 13, and 14, and replaces the original functions in the AI model with these approximations. Furthermore, the algorithm dynamically computes the optimal number of terms,  $n$ , based on the function's position within the model and a predefined accuracy budget, ensuring minimal degradation in baseline accuracy. The approximator described in Fig. 1 includes the above-illustrated search algorithm for identifying non-linear activation functions in the AI architecture. The activation selection and replacement process involves replacing the selected activation function with the pre-defined approximated activation from the activation table. The approximator triggers a counter at every activation selection and replacement, which computes the order of the Taylor series with the pre-defined search range. The search space for the approximation is determined through a set of brute-force experiments to understand the point of convergence of the approximated functions with the standard functions. The approximator also maintained a pre-defined deviation from the baseline accuracy by comparing the baseline accuracy to the estimated accuracy of the evaluated neural network at each iteration in computing the order of the Taylor series.

The algorithm is generic and can support various other non-linear activation functions similar to the above-described functions. The number of terms  $n$  in the approximated activation function of the Taylor series is directly proportional to the hardware resource requirements. The search space for approximation is constrained to balance power consumption and deviation from the AI

model’s baseline accuracy. Algorithm 1 presents the detailed infrastructure of the aforementioned algorithm used in TYTAN. The estimated approximation by the algorithm can be determined by factors such as activation placement in the AI architecture and its impact on the feature map to its preceding layers in the AI model. The approximation estimation is dynamic, and the depth of the Taylor series varies depending on the individual activation function based on the factors mentioned above. Utilizing pre-trained AI models as input can significantly lower the run-time cost of the algorithm, and it also inputs a minimal amount of dataset for accuracy estimation at every iteration. The algorithm’s run-time complexity depends on the complexity of the input AI architecture. The framework supports retraining AI models with the approximated non-linear activation functions by its pluggable infrastructure.

---

**Algorithm 1:** General Purpose Non-linear Approximation Algorithm

---

**Input:** Neural Network Model (*NN Model*), Test Data, Acceptable Accuracy Deviation (*Deviation*)

**Output:** Optimized Neural Network Model (*ApproxModel*)

```

1 L ← ActivationToBeApprox(NN Model);
2 BAcc ← Evaluate(NN Model);
3 ModelData ← [];
4 foreach Layer ∈ L do
5   | LayerData ← getLayer(NN Model, Layer);
6   | [nTerms, Acc] ← IterativeSearchBasedApprox(NN Model, Test Data);
7   | ModelData.append([nTerms, Acc]);
8   | if BAcc − Acc > Deviation then
9   |   | break;
10 ApproxModel ← Approximate(ModelData, NN Model);
11 DeltaAcc ← BAcc − Evaluate(ApproxModel);
12 if DeltaAcc > Deviation then
13   | call Approximator(ApproxModel, Test Data, Deviation);
14 return ApproxModel;

```

---

The above-discussed algorithm focused on improving the hardware acceleration for approximated non-linear activation function through dynamic approximation of every non-linear activation function utilized in the considered AI architecture. The algorithm can be scaled by considering hardware-aware factors, including support for handling higher depths of the Taylor series for better convergence with minimal utilization of system resources. Additionally, the dynamic lower and upper limits of the algorithms for Taylor-series computation enable it to work under different resource constraints, especially with limited hardware resources.

## 2.2 Hardware Modeling

The proposed hardware design, illustrated in Fig. 1, implements an optimal hardware/software co-design strategy to accelerate approximated non-linear activation functions in modern AI architectures. The architecture is reconfigurable and dynamic, enabling efficient power consumption and improved execution speed. The element-wise operator within the non-linear activation functions is mapped using TYTAN, which compute the approximated activation of the incoming tensor. The number of coefficients stored in the internal FIFO buffer determines the accuracy of approximation. The MAC units are modified from conventional MAC designs to implement the algorithm, as expressed in Eq. 3.

To establish the mathematical foundation of this approach, the  $e^x$  Taylor series expansion is considered:

$$\frac{x^n}{n!} = 1 + x + \frac{x^2}{2!} + \frac{x^3}{3!} + \frac{x^4}{4!} + \dots \quad (1)$$

For computational efficiency, this expansion is restructured as follows:

$$= 1 + x + \frac{x^2}{2!} + x^3 [c_3 + c_4(x)] \quad (2)$$

This reformulation  $T(x)$  enables efficient implementation in TYTAN using a nested multiplication approach, significantly reducing computational complexity:

$$T(x) = [c_0 + x [c_1 + x [c_2 + x [c_3 + c_4(x)]]]] \quad (3)$$

TYTAN processes four-dimensional tensor inputs as scalar batches using pre-computed Taylor series coefficients derived from the algorithm for efficient non-linear function computation. The Taylor series terms ( $n$ ) are precomputed by Algorithm 1 and stored in internal hardware buffers to control computational overhead. These variable terms enable flexible, adaptable architecture for efficiently mapping diverse non-linear functions.

The proposed approach can approximate various activation functions commonly used in neural networks. The standard forms of these functions are provided below:

$$\text{SELU}(x) = \lambda \begin{cases} x, & \text{if } x > 0 \\ \alpha(e^x - 1), & \text{if } x \leq 0 \end{cases} \quad (4)$$

$$\text{sigmoid}(x) = \frac{1}{1 + e^{-x}} \quad (5)$$

$$\text{Swish}(x) = x \cdot \text{sigmoid}(x) = \frac{x}{1 + e^{-x}} \quad (6)$$

$$\text{GELU}(x) \approx x \cdot \sigma(1.702x) = \frac{x}{1 + e^{-1.702x}} \quad (7)$$

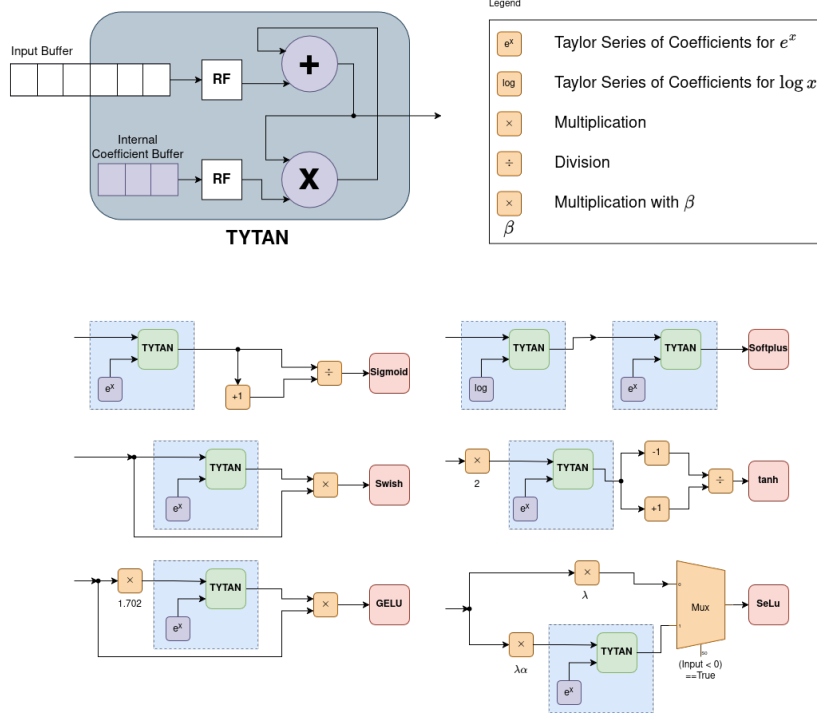


Fig. 2: TYTAN Architecture and Example Modes of Operations

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}} = \frac{e^{2x} - 1}{e^{2x} + 1} \quad (8)$$

$$\text{Softplus}(x) = \log(1 + e^x) \quad (9)$$

The key innovation in this approach lies in efficiently mapping these complex activation functions using the TYTAN hardware output. Equations 10 to 15 present the final formulations of the approximated non-linear functions to leverage hardware acceleration:

$$\text{SELU}(x) = \begin{cases} \lambda x, & \text{if } x > 0 \\ \lambda \cdot \alpha \cdot T_{e^x}(x), -1 & \text{if } x < 0 \end{cases} \quad (10)$$

$$\text{sigmoid}(x) = \frac{T_{e^x}(x)}{T_{e^x}(x) + 1} \quad (11)$$

$$\text{Swish}(x) = x \cdot T_{e^x}(x) \quad (12)$$

$$\text{GELU}(x) = x \cdot T_{e^x}(1.702x) \quad (13)$$

$$\tanh(x) = \frac{T_{e^x}(2x) - 1}{T_{e^x}(2x) + 1} \quad (14)$$

$$\text{Softplus}(x) = T_{\log}(T_{e^x}(x)) \quad (15)$$

The subscripts in each equation indicate the specific Taylor series coefficients utilized for each computation type. These reformulated equations are analogous to the six example TYTAN modes shown in Fig. 2

The proposed implementation offers significant flexibility, as any non-linear activation function can be efficiently mapped onto TYTAN by selecting appropriate coefficients. The arrangement of these activation functions within AI models determines the optimal number of terms  $n$  to be encoded into the hardware, dynamically adjusting based on accuracy requirements and power constraints. This adaptability allows the system to maintain high efficiency with minimal energy consumption while achieving acceptable accuracy in AI applications. Furthermore, the architecture provides a unified approach for accelerating diverse activation functions without requiring dedicated hardware for each function type.

### 3 HW-SW Co-Design & Results

The robust combination of domain-specific hardware and algorithm results in accelerated computation for non-linear activation functions involved in modern AI Algorithms. The hardware design proposed in section II was evaluated and compared with TensorFlow’s activation functions for correctness and accuracy. Detailed remarks on the algorithm’s performance and the hardware design are discussed in the following subsections.

Table 1: Performance Analysis of the Approximation Algorithm on the MobileViT Model

Model	Target	Baseline	Min.	Max.	Avg.	Est.	Deviation
		Acc. (%)	Length	Length	Length	Acc. (%)	
MobileViT	Swish	83.38	07	09	07.15	82.23	0.010
			07	19	08.81	82.95	0.005
			07	25	14.44	83.14	0.0025

#### 3.1 Software Setup

The experiment to evaluate the impact of the proposed algorithm for dynamic approximation utilizing the Taylor series is carried out by implementing approximated activation functions in the TensorFlow library [1]. The TensorFlow



simulations were used to compute the length of the Taylor series required by each approximated function to converge with the standard activation function. The approximation algorithm uses the point of convergence to restrict the search space and avoid unwanted computations beyond the convergence point. Further, the experiment focused on approximating the MobileViT [5], a state-of-the-art neural network with numerous compute-intensive activation functions. The approximation algorithm setup also includes T4 GPUs from Nvidia for extended hardware support. The algorithm utilizes the GPU acceleration to estimate the layer-wise length of the Taylor series for every activation function in the targeted neural network. MobileViT was initially trained on TFflower [11] with five classes and a training rate of 30 epochs.

The parameters, such as pre-trained MobileViT and a segment of test data, can be passed to the approximation algorithm for the layer-wise dynamic approximation. The algorithm begins with the baseline accuracy estimation for the targeted neural network, followed by the dynamic approximation for the targeted non-linear activation function. A search for the nominal length of the Taylor series is initiated by the algorithm on account of identifying the targeted non-linear activation function. The iterative search starts from the point of convergence to the lower limit, and the nominal length of the Taylor series is estimated for that particular activation function.

Table 1 demonstrates a clear correlation between the allowable deviation from the baseline accuracy and the length of the Taylor series utilized for approximation. The approximation algorithm operates for multiple deviation amounts from the baseline accuracy of the targeted neural network, and the length of the Taylor series for approximation is estimated to satisfy the constraints of the deviation from the baseline accuracy.

The algorithm targets about 32 swish non-linear activation functions at the MobileViT. The selection and replacement block described in Fig. 1 replaces the swish with the approximated swish for every function. Fig. 3 outlines the required length of the Taylor series for approximation at each targeted activation function at MobileViT for the different amounts of deviations from the baseline accuracy of the targeted neural network. Furthermore, Fig. 3 also illustrates that the intermediate layers of the considered model are sensitive and require significantly higher Taylor series terms to satisfy the accuracy constraints. It shows that the required length of the Taylor series is minimal for the higher deviation of estimated accuracy from the baseline accuracy of the target neural network and vice versa. The experiment also shows a tight requirement for compute power to elaborate the algorithm for neural networks with more non-linear activation functions. Similarly, the algorithm above and the experiment can be applied to next-generation state-space models like Mamba [3] to enable faster execution of state-of-the-art non-linear functions (Softplus) used within them.

### 3.2 Hardware Setup

This section provides a detailed elaboration of the hardware realization of the proposed model.

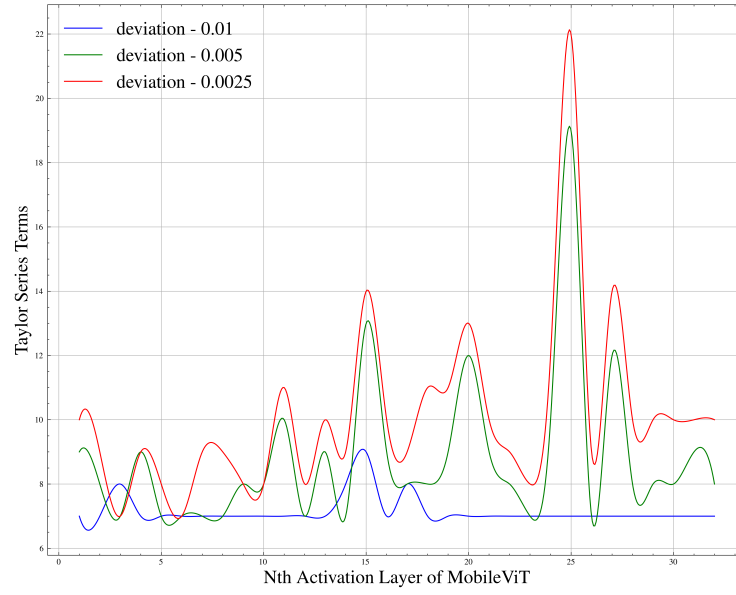


Fig. 3: Comparison of Nth Non-Linear Activation function vs. Order of the Taylor-series per Activation Layer (Swish) of the MobileViT

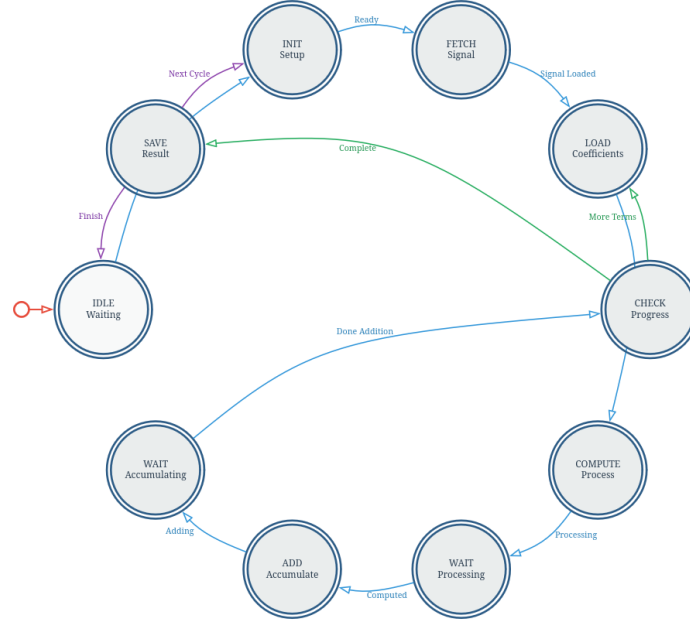


Fig. 4: Finite State Machine (FSM) diagram illustrating the operational flow of the TYTAN architecture

Table 2: TYTAN tanh computation latency with 30 Taylor coefficients

Operation	Clock Cycles Required	Execution Time (ms)*
Fill buffers (30 input values)	120	0.126
Generate output for every input	747	0.786
Full TYTAN operation (without Buffers)	22,474	23.656
Full TYTAN operation (with Buffers)	22,594	23.783

\*Calculated at 950MHz clock frequency (1.053 ns per clock cycle)

**TYTAN Operation** TYTAN efficiently computes activation functions using a multiply-accumulate (MAC) unit modified for polynomial approximation. Input data enters through a dedicated buffer, while Taylor series coefficients are stored in TYTAN’s internal buffer and reused for each input element.

The unit iteratively computes polynomial approximation through successive multiplications and additions (Eq. 3). The first coefficient is added to zero and passed to the MAC unit. Subsequent operations multiply the accumulated result with input and add the next coefficient, enabling efficient evaluation of non-linear functions.

Fig. 4 shows the finite state machine governing TYTAN operation. The FSM controls sequential processing through the computational pipeline, ensuring proper coefficient selection and accumulation at each polynomial evaluation step.

TYTAN output is mapped to activation functions such as SELU, Sigmoid, Swish, GELU, Tanh and Softplus, each expressed in terms of MAC operations (Eqs. 10 to 15). The architecture allows coefficient configuration through a dedicated port, providing flexibility to approximate diverse activation functions while maintaining scalability.

**Simulation and Synthesis** TYTAN is implemented in SystemVerilog HDL, with Taylor-series coefficients for the approximated activation functions derived from the software framework. A comprehensive verification testbench was developed using SystemVerilog to validate the functional correctness of TYTAN itself as well as with the SeLu, Sigmoid, and tanh add-ons. The testbench also counts the number of clock cycles it takes to finish the task across various modes of operations. A TYTAN is then simulated using both Synopsys VCS and Verilator, followed by synthesis using Synopsys Design Compiler V-2023.12-SP3 with the FreePDK45 process technology.

Table 3: Performance, Power, and Area (PPA) Comparison of SOTA Accelerators

Design	Process Node	Area (mm <sup>2</sup> )	Power (mW)	FMAX (MHz)
<b>TYTAN (This Work)</b>	FreePDK45	0.028 (base) 0.037 (w/ NL)	19.86 (base) 24.37 (w/ NL)	950
<b>NVDLA</b>	FreePDK45	1.002	35.165	450
<b>NN-LUT</b>	7 nm	0.001	0.059	N/A
<b>UNO</b>	TSMC 45 nm	0.283 (kernel)	66.5 (64 PE)	400
<b>ViTALiTy</b>	28 nm CMOS	5.22	1,460	500

NL = Non-Linear addons; PE = Processing Elements

### 3.3 Performance Evaluation

Figure 5 presents a detailed comparison of the accelerator’s accuracy across various activation functions, evaluated using different numbers of Taylor series terms. The results, benchmarked against TensorFlow’s implementations over an input range from -5 to +5, reveal that increasing the number of Taylor series coefficients consistently enhances the hardware’s accuracy.

Significantly, for every activation function, a threshold exists beyond which TYTAN’s output precisely matches TensorFlow’s reference. This finding implies that when the input range is limited to regions where the activation function is accurately represented with fewer coefficients, TYTAN can be optimally configured to use a reduced number of terms without sacrificing performance.

The computational latency of TYTAN is determined exclusively by the number of Taylor series coefficients employed, irrespective of the specific activation function being computed. Table 2 offers a comprehensive cycle analysis of the TYTAN hardware while calculating the tanh function using thirty Taylor series coefficients.

The current implementation of TYTAN operates in FP32 precision, with the total latency comprising contributions from individual module delays including buffers, multiplier, and adder operations. This overall computational latency can be further reduced through strategic optimization approaches.

By transitioning to lower precision formats or incorporating improved core hardware modules that is multiplier and adder circuits, the cumulative latency penalties from these core operations can be significantly minimized without changing rest of the system, thereby enhancing the accelerator’s overall performance efficiency.

### 3.4 Comparative Analysis

Table 3 compares neural network accelerator designs, with NVDLA serving as the primary baseline for evaluating TYTAN’s performance. The comparison focuses

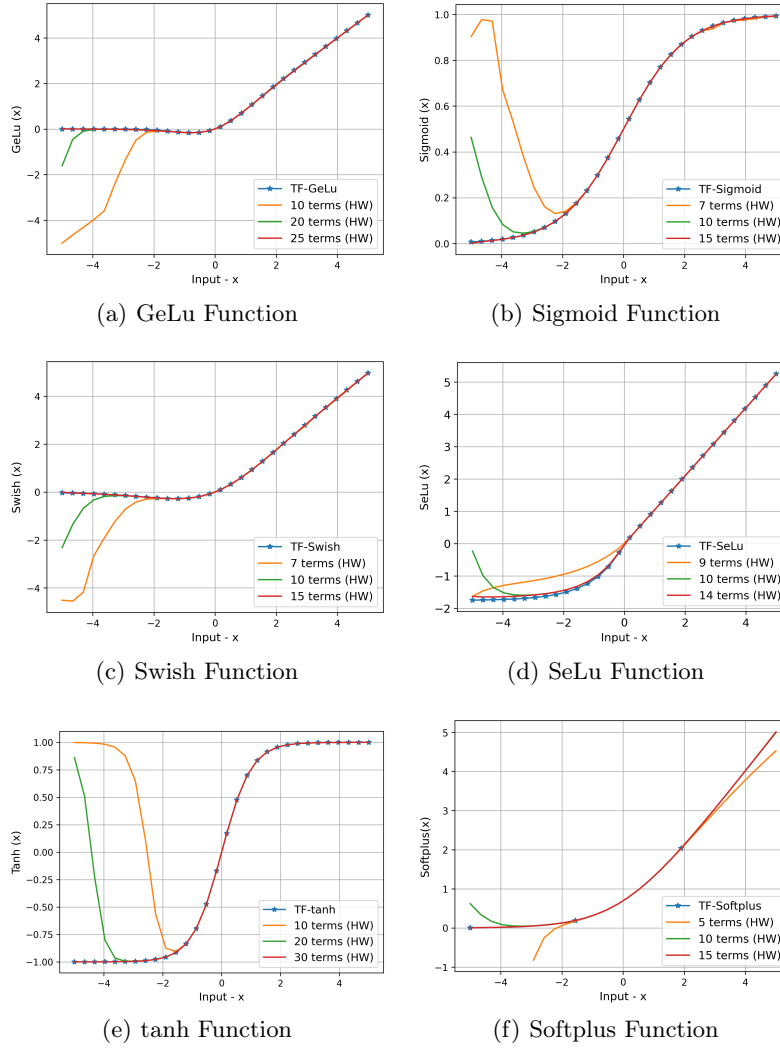


Fig. 5: TensorFlow Activation Functions vs. Hardware-Approximated Activation Results

on NVDLA's p-partition that includes the Single Data Point Processor (SDP), which uses lookup tables for linear and non-linear functions on individual data points.

TYTAN demonstrates exceptional performance improvements over NVDLA across all metrics using FreePDK45 technology. TYTAN achieves remarkable area efficiency with 35.785x improvement without add-ons and 27.081x with add-ons. Power consumption is reduced by 1.771x and 1.442x respectively, while

Table 4: Precision and Operation Support Comparison

Design	Precision Support	Supported Operations
<b>TYTAN (This Work)</b>	FP32 (configurable)	SELU, Sigmoid, Swish, GELU, Tanh, Softplus, etc
<b>NVDLA (baseline)</b>	INT8, INT16, FP16	ReLU, PReLU, Sigmoid, Tanh
<b>NN-LUT</b>	FP32, FP16, INT32	GELU, Softmax, LayerNorm
<b>UNO</b>	16-bit fixed-point	Division, exp, log, sigmoid, tanh
<b>ViTALiTy</b>	16-bit	Linear attention schemes

maximum operating frequency more than doubles at 2.11x improvement. These results highlight TYTAN’s scalability and modularity, efficiently handling additional capabilities like SeLu, Sigmoid, and tanh activation functions while maintaining substantial performance advantages.

The table includes other state-of-the-art accelerators for context. NN-LUT uses advanced 7nm technology with extremely low area and power but requires 2 clock cycle latency and lacks frequency specifications, while TYTAN achieves 950 MHz in FreePDK45. UNO, despite using TSMC 45nm, consumes more power (66.5 mW vs 19.86-24.37 mW) and operates at lower frequency (400 MHz vs 950 MHz). ViTALiTy, implemented in more advanced 28nm CMOS, exhibits significantly higher area (5.22 mm<sup>2</sup>) and power (1,460 mW) while operating at lower frequency (500 MHz) than TYTAN’s FreePDK45 implementation. These comparisons demonstrate TYTAN’s exceptional efficiency, achieving superior performance even in a less advanced process technology.

## 4 Conclusion

This paper presents TYTAN, a hardware-software co-design for energy-efficient non-linear activation function computation in AI architectures. TYTAN achieves  $\sim 2.11\times$  frequency improvement,  $\sim 1.77\times$  power reduction, and  $\sim 35.8\times$  area reduction compared to NVDLA baseline, while outperforming other state-of-the-art accelerators even when implemented in 45nm technology.

The modular and reconfigurable design ensures seamless integration with DNN accelerators, offering scalability from edge devices to high-performance systems. The parameterizable coefficient buffer optimizes precision-resource trade-offs, while supporting diverse AI workloads with configurable precision and a variety of activation functions. This approach delivers memory-optimized, accelerated activation computation with minimal resource constraints. The complete implementation is made available as open source [7] to enable reproducibility and foster further research in this domain.

## References

1. Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., Kudlur, M., Levenberg, J., Monga, R., Moore, S., Murray, D.G., Steiner, B., Tucker, P., Vasudevan, V., Warden, P., Wicke, M., Yu, Y., Zheng, X.: Tensorflow: A system for large-scale machine learning (2016), <https://arxiv.org/abs/1605.08695>
2. Dass, J., Wu, S., Shi, H., Li, C., Ye, Z., Wang, Z., Lin, Y.: Vitality: Unifying low-rank and sparse approximation for vision transformer acceleration with a linear taylor attention. In: 2023 IEEE International Symposium on High-Performance Computer Architecture (HPCA). pp. 415–428 (2023). <https://doi.org/10.1109/HPCA56546.2023.10071081>
3. Gu, A., Dao, T.: Mamba: Linear-time sequence modeling with selective state spaces (2024), <https://arxiv.org/abs/2312.00752>
4. Gustafson, J., Yonemoto, I.: Fast approximations of activation functions in deep neural networks using posit arithmetic. IEEE Transactions on Emerging Topics in Computing **9**(3), 1391–1402 (2021). <https://doi.org/10.1109/TETC.2020.2973431>
5. Mehta, S., Rastegari, M.: Mobilevit: Light-weight, general-purpose, and mobile-friendly vision transformer (2022), <https://arxiv.org/abs/2110.02178>
6. NVIDIA: Nvidia deep learning accelerator. [https://nvidia.org/hw/v2/integration\\_guide.html](https://nvidia.org/hw/v2/integration_guide.html) (2018), [https://nvidia.org/hw/v2/integration\\_guide.html](https://nvidia.org/hw/v2/integration_guide.html)
7. Pramanik, S.: Tytan: Taylor-series based non-linear activation engine for deep learning accelerators. <https://github.com/SoHam-56/GNAE> (2025), gitHub repository
8. Silicon Integration Initiative: Open cell and free PDK libraries. <https://si2.org/open-cell-and-free-pdk-libraries/> (2025), <https://si2.org/open-cell-and-free-pdk-libraries/>, accessed: March 16, 2025
9. Smith, A., Patel, R., Gupta, N.: Implementation of activation functions using various approximation methods. In: Proceedings of the IEEE Symposium on Computational Intelligence. pp. 1–6 (2022)
10. Taylor, J., Smith, K., Anderson, M.: Taylornet: A taylor-driven neural network approximation framework. Neural Processing Letters **52**, 1483–1497 (2020)
11. The TensorFlow Team: Flowers. [http://download.tensorflow.org/example\\_images/flower\\_photos.tgz](http://download.tensorflow.org/example_images/flower_photos.tgz) (Jan 2019), [http://download.tensorflow.org/example\\_images/flower\\_photos.tgz](http://download.tensorflow.org/example_images/flower_photos.tgz)
12. Wu, D., Li, J., Behroozi, S., Kim, Y., Miguel, J.S.: Uno: Virtualizing and unifying nonlinear operations for emerging neural networks. In: 2021 IEEE/ACM International Symposium on Low Power Electronics and Design (ISLPED). pp. 1–6 (2021). <https://doi.org/10.1109/ISLPED52811.2021.9502473>
13. Wu, D., Li, J., Behroozi, S., Kim, Y., San Miguel, J.: Uno: Virtualizing and unifying nonlinear operations for emerging neural networks. IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems **40**(12), 2690–2705 (2021). <https://doi.org/10.1109/TCAD.2021.3107632>
14. Yu, J., Park, J., Park, S., Kim, M., Lee, S., Lee, D.H., Choi, J.: Nn-lut: Neural approximation of non-linear operations for efficient transformer inference. In: Proceedings of the 59th ACM/IEEE Design Automation Conference. p. 577–582. DAC '22, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3489517.3530505>, <https://doi.org/10.1145/3489517.3530505>