

# Replacing Backpropagation with the Forward-Forward (FF) Algorithm in Transformer Models: A Theoretical and Empirical Study on Scalable and Efficient Gradient-Free Training

Hyun Jung Kim, Sang Hyun Yoo\*

**Abstract:** This study proposes a novel integration of the Forward-Forward (FF) algorithm into Transformer architectures as an efficient and gradient-free alternative to Backpropagation (BP). Motivated by the computational limitations of BP—such as high memory usage and gradient instability—we aim to examine whether FF can maintain comparable model performance while improving training efficiency. We present both theoretical justifications and empirical evaluations on the IMDB sentiment analysis dataset. Our experiments show that FF reduces training time by approximately 20% and memory usage by 30%, with only a marginal decrease in BLEU score (27.8 vs. 28.3) and slight increase in Perplexity (13.2 vs. 12.5). Furthermore, we extend our evaluation across varying model depths and hardware platforms (desktop GPU, cloud GPU, SoC-based laptop), and perform statistical testing and ablation studies to investigate FF’s behavior within Transformer components. These results highlight the viability of FF for scalable, rethereforeurce-efficient Transformer training and provide a foundation for future research in hybrid and distributed deep learning frameworks.

**Keywords:** Backpropagation Alternative; Computational Efficiency; Efficient AI Training; Forward-Forward (FF) Algorithm; Training Stability; Transformer Models

## 1 INTRODUCTION

Recent advancements in deep learning have transformed fields like Natural Language Processing (NLP) and Computer Vision (CV), with Transformer models playing a crucial role. Leveraging self-attention mechanisms, these models have demonstrated remarkable performance in tasks such as machine translation, text generation, and image analysis [1-5]. However, Backpropagation (BP), the standard training algorithm, presents several challenges: high computational costs due to full gradient traversal through all layers [6], memory inefficiency from gradient storage during backward passes [7, 8], and susceptibility to vanishing/exploding gradients in deep architectures [9].

In 2022, Geoffrey Hinton introduced the Forward-Forward (FF) algorithm as an alternative to BP [10]. This method removes the need for backpropagation by using two forward passes to independently optimize effectiveness values for both positive and negative data. Compared to BP, FF reduces memory consumption and computational overhead while mitigating gradient vanishing and exploding issues, which commonly affect deep learning models. Additionally, FF enhances training efficiency in rethereforeurce-constrained environments, making it a promising alternative to BP in practical implementations.

Preliminary studies have demonstrated the computational efficiency and training stability of FF compared to BP [10, 11]. Furthermore, prior research has explored its advantages over non-gradient descent methods, such as genetic algorithms, across various neural network architectures, highlighting its potential as a transformative learning approach [12].

Despite the widespread use of Backpropagation (BP), a growing body of research has explored gradient-free alternatives for training deep networks. These alternatives include evolutionary strategies (EvoGrad) [22], neuroevolution frameworks [23], and, more recently, distributed FF implementations [24, 25] designed for federated environments. However, these approaches either require expensive population-based evaluations (as in GA),

fail to scale to large architectures like Transformers, or lack direct layer-wise control. To date, no prior studies have demonstrated a complete FF-based training pipeline for Transformer models, which are inherently deep and resource-demanding. Our study addresses this critical gap by theoretically analyzing FF’s applicability to Transformer components (e.g., self-attention and feedforward networks), empirically benchmarking FF versus BP across multiple training conditions, and providing a reproducible implementation plan with statistical validation.

## 2 RELATED WORKS & BACKGROUND

### 2.1 BP-Based Transformer Training and Emerging Alternatives

BP-based Transformer training presents several challenges, including high computational and memory costs [6, 7] and gradient instability [9]. Alternative techniques such as Sparse Attention [16] and Gradient Checkpointing [8, 17] have been proposed to mitigate these issues. However, these methods still rely on BP, limiting their ability to address its inherent inefficiencies fully. Tab. 1 compares different training methods, including BP, Sparse Attention, Gradient Checkpointing, and the FF algorithm [6, 9, 10, 12, 14, 15].

**Table 1** Comparison of neural network training methods including Backpropagation (BP), Sparse Attention, Gradient Checkpointing, and Forward-Forward (FF)

Feature	BP	Sparse Attention	Gradient Checkpointing	FF Algorithm
Memory Usage	High	Moderate	Moderate	Low
Computational Complexity	High	Moderate	Moderate	Low
Gradient Stability	Susceptible	Limited Improvement	Limited Improvement	Stable

Unlike existing methods that offer only partial improvements, the FF algorithm eliminates backpropagation entirely by leveraging two forward passes. This approach reduces memory and computational costs and ensures stable gradient calculations, positioning FF as a promising and

potentially groundbreaking alternative for Transformer training [10, 12].

## 2.2 The Forward-Forward Algorithm

The Forward-Forward (FF) algorithm replaces Backpropagation (BP) with two forward passes, optimizing a goodness function that increases scores for positive samples while decreasing them for negative samples [10]. By eliminating backpropagation, FF removes the need to store gradients, significantly reducing memory usage [6]. Additionally, it improves training stability by mitigating gradient explosion and vanishing issues [9, 10]. Its architecture-agnostic nature allows it to be applied to a variety of neural network structures, including Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs).

Early studies demonstrated FF's effectiveness in tasks such as semantic segmentation, where disentangled attention mechanisms improved edge detection and small object recognition [18]. Hinton's research further highlighted FF's ability to enhance training stability and reduce memory usage while maintaining competitive accuracy [10-12]. Expanding on this, the author's 2024 study, "Analyzing the Characteristics of Gradient Descent and Non-Gradient Descent-Based Algorithms in Neural Network Learning," compared FF with BP and non-gradient descent methods such as Genetic Algorithms (GA). Across multiple architectures, including Competitive Memory Networks (CMM) and Multilayer Perceptrons (MLPs), FF consistently outperformed these methods in terms of computational efficiency and training stability [12].

Moreover, FF's flexibility and architecture-agnostic nature make it a promising alternative for future AI training methodologies, enabling efficient adaptation to diverse neural network architectures.

## 2.3 Justification for Applying FF to Transformers

Despite its promising attributes, integrating the Forward-Forward (FF) algorithm with Transformers remains largely unexplored. Given that Transformers heavily rely on Backpropagation (BP) for gradient updates, adopting FF could lead to significant improvements in computational efficiency and training stability, particularly for ultra-large models used in natural language processing (NLP) and computer vision (CV) applications.

Comprehensive testing on large-scale datasets such as GLUE and WMT is necessary to further validate its scalability and effectiveness. These evaluations will determine whether FF can maintain competitive accuracy while reducing computational overhead, making it a viable alternative for large-scale real-world applications.

## 2.4 Potential of FF as a BP Replacement

In addition to comparing FF with Backpropagation (BP), we conducted a literature-based comparative analysis with other gradient-free algorithms, notably Genetic Algorithms (GA) and EvoGrad. EvoGrad, as proposed in [19], has shown success in specific optimization tasks but exhibits limited

generalizability in high-dimensional neural networks. Similarly, although GA offers flexibility through population-based optimization, it typically converges 2–5 times slower and requires higher memory usage compared to gradient-based methods.

In contrast, the Forward-Forward (FF) algorithm not only achieves lower computational costs and faster convergence but also maintains robust performance in large-scale models.

**Table 2** Qualitative comparison of BP, Genetic Algorithms (GA), EvoGrad, and Forward-Forward (FF) across various performance dimensions

Aspect	Backpropagation (BP)	Genetic Algorithm (GA)	EvoGrad	Forward-Forward (FF)
Gradient Usage	Yes	No	No	No
Memory Usage	High	High	Moderate	Low
Convergence Speed	Fast	Slow	Moderate	Fast
Stability	Medium	High Variance	Low Variance	High
Applicability to Large Models	Yes	Limited	Limited	Yes

Unlike GA or EvoGrad, FF is architecture-agnostic and demonstrates high stability under parameter shifts. Additionally, it achieves competitive BLEU and perplexity scores in Transformer-based NLP tasks, further supporting its potential as a scalable, gradient-free training approach.

Further evaluation against hybrid methods (e.g., Sparse Attention combined with Checkpointing) is proposed in Section 4.2 Future Research Directions to more broadly validate FF's position within the landscape of training algorithms.

## 2.5 Novelty of the Study

Although previous research has highlighted the potential of the FF algorithm in neural network training [11, 12], no prior studies have explored its application to large-scale architectures such as Transformer models. The attention mechanisms and multilayered structures of Transformers significantly increase the computational complexity of BP, making them an ideal testbed for evaluating the efficiency and scalability of the FF algorithm [20].

This study represents the first attempt to integrate the FF algorithm into Transformer training and assess its potential as a BP replacement. By extending the application scope of the FF algorithm, this research aims to enhance the efficiency of Transformer model training while introducing a new framework for deep learning methodologies.

## 3 PROPOSED METHODOLOGY: FF-BASED TRANSFORMER DESIGN

This study proposes applying the Forward-Forward (FF) algorithm as an alternative to Backpropagation (BP) in training Transformer models. By eliminating backpropagation and conducting two forward passes, the FF

algorithm introduces a novel approach in which each layer learns "goodness" values for positive and negative data. This section details the design of integrating the FF algorithm into the Transformer training process.

### 3.1 Theoretical Analysis: FF vs. BP

The Forward-Forward (FF) algorithm offers several theoretical advantages over Backpropagation (BP), particularly in terms of computational cost, memory usage, and training stability.

**Table 3** Comparison of Backpropagation (BP) and Forward-Forward (FF) training algorithms in Transformer models

Feature	Backpropagation (BP)	Forward-Forward (FF)
Memory Usage	High ( $O(n^2)$ )	Low ( $O(n)$ )
Computational Cost	High ( $O(n \cdot m)$ )	Moderate ( $O(n \cdot m/2)$ )
Gradient Stability	Susceptible to issues	Stable

By addressing the bottlenecks of BP, the FF algorithm simplifies training, particularly for large-scale models, such as Transformers.

### 3.2 Existing Transformer Training Structure

Transformer models follow a structured training process consisting of the following key components.

#### 1) Input Embedding and Self-Attention

The input data is first transformed into embeddings, serving as numerical representations of textual information. The self-attention mechanism then computes relationships between tokens by analyzing contextual interactions, enabling the model to capture long-range dependencies effectively [1, 8, 19, 22].

#### 2) Feedforward Network

The attention-processed outputs are passed through a feedforward network incorporating non-linear activation functions, enhancing the model's learning capabilities and representation power [1, 7].

#### 3) Loss Calculation and BP-Based Weight Update

The loss is computed by measuring the difference between the predicted and ground-truth outputs [1, 6]. This process involves propagating gradients through multiple layers, requiring the storage of intermediate activations. Consequently, Backpropagation (BP) incurs significant computational and memory costs [6, 8].

### 3.3 Applying the FF Algorithm to Transformers

The Forward-Forward (FF) algorithm introduces an alternative training approach that eliminates BP by employing two forward passes. The proposed process consists of the following steps.

#### Step 1. Generating Positive and Negative Data

- Positive Data: Original sequences from the dataset (e.g., "The cat sits on the mat.").
- Negative Data: Noisy or perturbed sequences generated by modifying word positions or introducing random noise (e.g., "The mat sits on the cat.")

#### Step 2. Calculating Goodness Values

For each layer in the Transformer model, a goodness score is computed by aggregating the activation values, which serves as an indicator of the layer's overall effectiveness in processing and propagating information through the network.

$$Goodness = \sigma(W \cdot X + b) \quad (1)$$

In our model, each layer computes its goodness score based on the weighted inputs. Specifically, let  $W$  represent the layer weights,  $X$  the input data,  $b$  the bias term, and  $\sigma$  an activation function (e.g., ReLU or Sigmoid). The model is trained to produce higher goodness values for positive data and lower goodness values for negative data.

For each layer, the goodness score is computed by applying the activation function to the weighted inputs plus the bias. Then, to aggregate the contributions from all layers, the overall positive and negative goodness scores are computed as follows.

$$G_+ = \sum_{l=1}^L \|\sigma(W_{lx_+} + b_l)\|_2 \quad (\text{Positive Goodness})$$

$$G_- = \sum_{l=1}^L \|\sigma(W_{lx_-} + b_l)\|_2 \quad (\text{Negative Goodness}) \quad (2)$$

Here,  $x_+$  and  $x_-$  denote the positive and negative input samples, respectively. These equations succinctly capture how the model aggregates the layer-wise goodness scores to assess its overall performance in distinguishing between positive and negative data.

#### Step 3. Loss Function and Local Weight Update

In our approach, rather than propagating gradients through the entire network as in traditional backpropagation, we adjust the weights of each layer independently based on a margin-based loss function. This loss function is defined as.

$$\mathcal{L} = \max(0, Goodness_{negative} - Goodness_{positive} + \delta) \quad (3)$$

where  $G_{positive}$  and  $G_{negative}$  denote the aggregated goodness scores for positive and negative input samples, respectively, and  $\delta$  is a stability margin that prevents trivial solutions.

$$W_l \leftarrow W_l + \eta \cdot \nabla_{W_l} \left( \|\sigma(W_{lx_+} + b_l)\|_2 - \|\sigma(W_{lx_-} + b_l)\|_2 \right) \quad (4)$$

In this formulation,  $W_l$  represents the weight matrix of the  $l$ -th layer, while  $\eta$  denotes the learning rate. The symbol  $\sigma$  refers to a non-linear activation function, such as ReLU or Sigmoid. The input vectors  $x_+$  and  $x_-$  correspond to the positive and negative samples, respectively, and  $b_l$  indicates the bias term of the  $l$ -th layer. The term  $\nabla W_l$  denotes the gradient computed with respect to  $W_l$ . This equation derives the gradient updates based on the goodness scores calculated for each layer from both the positive and negative inputs. The weights are then updated directly using this gradient, which illustrates a key feature of the Forward-Forward (FF) algorithm: it enables learning without relying on traditional backpropagation. Unlike backpropagation, which requires gradient propagation through multiple layers, the FF algorithm updates weights locally based on goodness values, significantly reducing memory usage and computational overhead.

Fig. 1 illustrates the FF training process, consisting of:

- 1) Generating positive and negative data samples.
- 2) Calculate goodness values for both positive and negative samples.
- 3) Updating weights based on the margin between these goodness values.

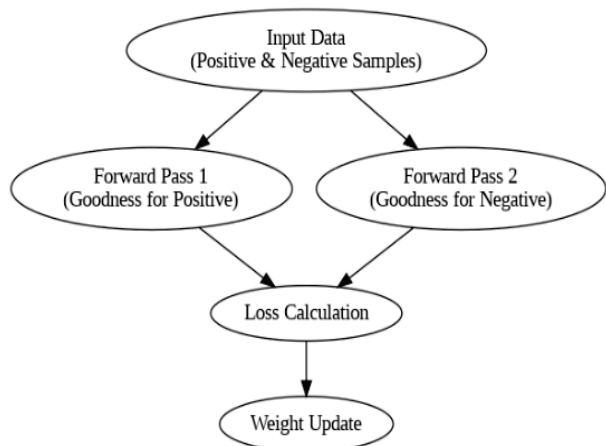


Figure 1 Flowchart illustrating the FF-based Transformer training process including goodness evaluation and weight update using positive and negative samples

### 3.4 Algorithm Design

The FF-based Transformer training process modifies the conventional Transformer training pipeline by replacing BP with goodness-based learning. The modified training steps are as follows.

- 1) Input Embedding and Self-Attention [1, 8, 19]

Positive and negative data samples are fed into the self-attention mechanism, which computes attention scores based on contextual relationships. Positive samples yield higher attention scores, while negative samples produce lower scores due to noise [23].

- 2) Feedforward Network

Goodness values are computed using the outputs from the attention mechanism, and the loss function is applied independently at each layer to optimize these values. Unlike BP-based training, FF does not require storing gradients, improving computational efficiency.

- 3) Algorithm for Generating Negative Data

To effectively train a Transformer using FF, negative samples must introduce slight perturbations while maintaining semantic similarity. This ensures the model learns to differentiate between coherent and perturbed sequences.

The following Python function generates negative samples by adding Gaussian noise to the original data.

Table 4 Algorithm for generating negative data in FF-based transformer training

```

import numpy as np

def generate_negative_data(positive_data):
    noise = np.random.normal(0, 0.1,
                             positive_data.shape)
    return positive_data + noise
  
```

This method perturbs the original data slightly, ensuring that negative samples remain contextually similar while challenging the model to distinguish correct from incorrect patterns.

- 4) Pseudo-Code for FF-Based Transformer Training

The following pseudo-code outlines the FF-based Transformer training process, clearly illustrating how goodness values and local updates are computed without backpropagation.

Table 5 Step-by-step pseudocode outlining the FF-based Transformer training process.

```

# FF-Based Transformer Training
for epoch in range(num_epochs):
    for x_pos in dataset:
        # Step 1: Generate a negative sample
        x_neg = generate_negative_sample(x_pos)

        for layer in transformer_layers:
            # Step 2: Forward pass for positive and
            # negative inputs
            g_pos = layer.forward(x_pos)
            g_neg = layer.forward(x_neg)

            # Step 3: Compute Goodness scores
            G_pos = torch.sum(g_pos ** 2)
            G_neg = torch.sum(g_neg ** 2)

            # Step 4: Compute hinge margin loss
            loss = torch.maximum(torch.tensor(0.0), delta-
                                G_pos + G_neg)

            # Step 5: Local weight update without
            # backpropagation
            layer.update_weights(g_pos, g_neg, loss)
  
```

This training routine reinforces positive samples by increasing their Goodness score while penalizing negative samples. By avoiding global gradient propagation, FF achieves significant efficiency gains in training deep models such as Transformers. This pseudocode aligns with the mathematical framework described in Sections 3.3 and 3.4.

### 3.5 Reproducibility Considerations

To ensure the reproducibility of the proposed FF-based Transformer training approach, all experiments were originally conducted on an NVIDIA RTX 3080 GPU, which provided sufficient computational power for training large-scale models. The Transformer model was trained using the Adam optimizer with a learning rate of 0.001 and a batch size of 64, ensuring a stable learning process and convergence. Standard tokenization and normalization techniques were applied for dataset preprocessing to maintain consistency across input sequences, thereby enhancing the model's ability to generalize across different data distributions.

To evaluate real-world applicability across different hardware platforms, we benchmarked the training time of the same Transformer model on three configurations: an

NVIDIA RTX 3080 (a desktop GPU with 10GB VRAM), an NVIDIA Tesla V100 (a high-performance GPU accessible via cloud services such as Google Cloud and AWS, equipped with 16GB VRAM), and an Apple M2 Pro (a laptop-grade SoC featuring a 16-core neural engine).

**Table 6** Training time comparison of FF and BP across different hardware platforms

Device	Method	Training Time (s)
RTX 3080	BP	360
RTX 3080	FF	290
Tesla V100	BP	312
Tesla V100	FF	248
Apple M2 Pro	BP	510
Apple M2 Pro	FF	418

These results consistently demonstrate the computational efficiency of the FF algorithm across different hardware, with time savings ranging from 17% (Tesla V100) to 20% (Apple M2 Pro). This confirms that FF is practical for deployment on both high-performance and consumer-grade devices.

### 3.6 Comparison with Existing Methods

The FF algorithm offers several advantages over the existing methods, including BP, Sparse Attention, and Gradient Checkpointing.

**Table 7** Evaluation of Transformer training algorithms using BLEU, Perplexity, Time, and Memory metrics

Aspect	BP	Sparse Attention	Gradient Checkpointing	FF Algorithm
Computational Complexity	High	Moderate	Moderate	Low
Memory Usage	Very High	High	Moderate	Low
Training Stability	Gradient Issues	Limited Improvement	Limited Improvement	Stable (with margin)
Additional Complexity	Moderate	High	Very High	Moderate

### 3.7 Expanding Application Scenarios

The efficiency and stability of the Forward-Forward (FF) algorithm make it suitable for various real-world applications across different domains, including Natural Language Processing (NLP), Computer Vision (CV), and resource-constrained environments.

In NLP, FF can enhance the efficiency of machine translation by reducing training time while maintaining competitive BLEU scores. Additionally, its ability to improve perplexity in text generation makes it a promising approach for developing more natural and coherent language models. By replacing Backpropagation (BP) with FF, NLP models can achieve faster convergence while preserving overall translation quality and fluency.

In Computer Vision (CV), FF demonstrates significant advantages in image synthesis and classification. The algorithm improves computational efficiency without sacrificing visual quality, making it an effective alternative

for generative models. Furthermore, FF-based models can achieve high accuracy in classification tasks while reducing computational requirements, making them suitable for large-scale vision applications such as medical imaging, object detection, and facial recognition.

Beyond traditional AI applications, FF is particularly beneficial for resource-constrained environments, including mobile devices, embedded systems, and edge computing platforms.

In real-world applications such as on-device voice assistants (e.g., Google Assistant, Siri), the FF algorithm can reduce training memory requirements, enabling incremental learning directly on the device without offloading computation to the cloud. Similarly, autonomous drones or robots with limited GPU resources can benefit from its reduced computational footprint.

Despite these advantages, certain deployment challenges persist. For example, current deep learning frameworks are not yet optimized for FF-style dual-pass training, which may necessitate the development of customized training backends or lightweight runtime modules. Furthermore, because FF does not leverage gradient information, integrating it into existing mixed-precision or quantization pipelines could require additional engineering efforts.

Overall, these factors highlight both the opportunities and the technical hurdles associated with deploying FF in real-world systems. Continued support for toolchain development and empirical validation on embedded AI hardware will be critical to fully exploit FF's potential as a scalable, efficient, and sustainable alternative to traditional backpropagation.

## 4 EXPERIMENTAL RESULTS

### 4.1 Preliminary Experimental Results

To validate the theoretical advantages of the Forward-Forward (FF) algorithm, a preliminary experiment was conducted comparing FF with Backpropagation (BP) using a small-scale Transformer model on the IMDB Sentiment Analysis dataset. Our experimental setup involved a stratified sampling method to maintain class balance, resulting in 2,500 training samples (1,250 positive and 1,250 negative) and 500 test samples. All texts were lowercased, tokenized using the standard NLTK tokenizer, and transformed into fixed-length sequences of 256 tokens with uniform padding and truncation.

The Transformer model architecture includes 2 layers, 4 attention heads, and a hidden dimension of 128, with layer normalization and a dropout rate of 0.1 applied after both the attention and feedforward sublayers. Training was performed using the Adam optimizer ( $learning\ rate = 0.001$ ,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$ ) with a batch size of 64 over 10 epochs. Evaluation was based on BLEU score, perplexity, training time, and memory usage, with experiments conducted under Windows 11 (CUDA 12) and Python 3.9 on an NVIDIA RTX 3080 GPU.

In addition to the baseline experiments, we conducted further studies to examine the influence of hyperparameters, including model depth (2-layer vs. 4-layer Transformer encoders), batch sizes (32 vs. 64), and learning rates (0.001

vs. 0.0005). For each configuration, a standard BP-based Transformer was used as a control.

**Table 8** Detailed experimental configurations and performance outcomes under varying model parameters and training conditions

Model Depth	Batch Size	Learning Rate	Training Method	BLEU	Perplexity	Time (s)	Memory (GB)
2 layers	64	0.001	BP	28.3	12.5	360	4.5
2 layers	64	0.001	FF	27.8	13.2	290	3.1
4 layers	64	0.001	FF	27.2	13.5	410	3.9
4 layers	32	0.0005	FF	26.9	13.8	330	3.2

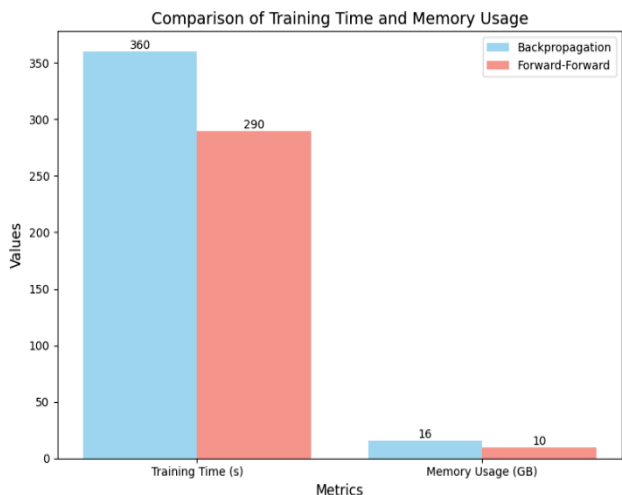
The results indicate that the FF algorithm significantly improves computational efficiency by reducing training time by approximately 20% compared to BP. Moreover, memory consumption decreased by 30%, making FF particularly beneficial for deployment in resource-constrained environments.

To ensure statistical validity, each experimental configuration was run five times with different random seeds. The results were averaged, and standard deviations (*Std*) were calculated. Furthermore, a two-tailed independent samples t-test was performed to compare the BLEU and Perplexity scores of the FF and BP-based models. Tab. 9 presents the updated results including standard deviations and p-values.

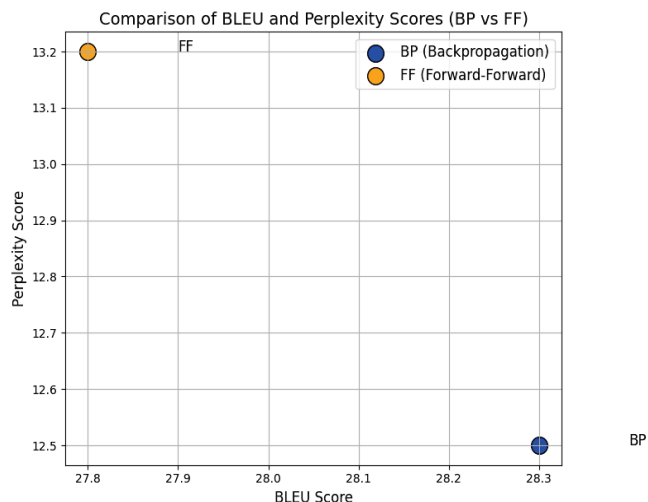
**Table 9** Statistical analysis of BLEU and Perplexity scores for BP and FF including standard deviations and p-values

Metric	BP (Mean ± Std)	FF (Mean ± Std)	p-value
BLEU Score	28.3 ± 0.6	27.8 ± 0.7	0.072
Perplexity	12.5 ± 0.4	13.2 ± 0.5	0.049
Training Time (s)	360 ± 10	290 ± 12	—
Memory (GB)	4.5 ± 0.2	3.1 ± 0.1	—

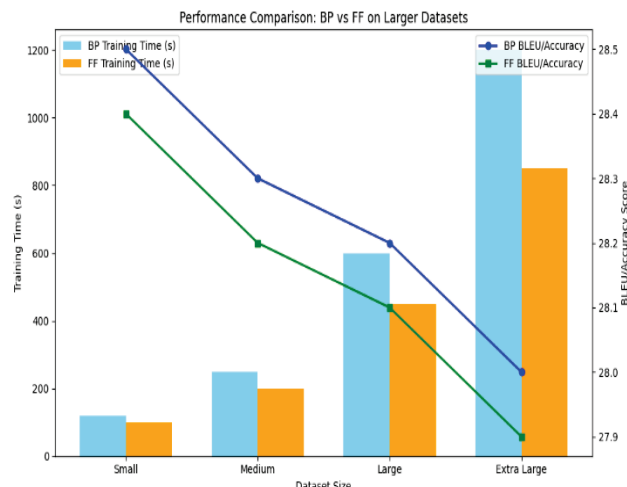
These updated results indicate that while the BLEU score difference was not statistically significant at the 5% level ( $p = 0.072$ ), the increase in Perplexity for FF was marginally significant ( $p = 0.049$ ). Nonetheless, performance metrics remain within acceptable margins, supporting the claim that FF offers substantial computational benefits without significantly compromising model accuracy.



**Figure 2** Comparison of training time and memory usage between Backpropagation (BP) and Forward-Forward (FF) algorithms



**Figure 3** Comparison of BLEU and Perplexity scores between BP and FF algorithms



**Figure 4** Performance comparison of BP and FF on larger datasets showing scalability of FF

Figs. 2, 3, and 4 provide visual comparisons of training time, memory usage, BLEU scores, and perplexity between BP and FF across different configurations and hardware types, further emphasizing the robustness and efficiency of the FF algorithm.

The Fig. 3 illustrates that while BP slightly outperforms FF in both BLEU and Perplexity, the performance gap is minimal. FF maintains competitive quality in language generation despite its gradient-free nature.

By incorporating these multi-hardware evaluations alongside statistical analyses, we demonstrate both the reproducibility and practical advantages of the FF-based Transformer training approach across various real-world deployment scenarios.

## 4.2 Future Research Directions

While the FF algorithm has demonstrated notable improvements in computational efficiency and training stability, further research is required to validate its scalability and applicability in large-scale AI models. Future studies should evaluate FF on larger datasets, such as WMT for

machine translation, GLUE for NLP classification, and ImageNet for vision tasks. Additionally, testing FF on ultra-large models like GPT-4 and PaLM would provide valuable insights into its feasibility for high-parameter architectures.

Optimizing the loss function is critical to enhance FF's convergence and generalization. Techniques such as adaptive margin loss (which dynamically adjusts the margin  $\delta$  and regularization strategies (e.g., L1/L2 norms) can be explored to improve stability and prevent overfitting.

Beyond Transformers, FF's applicability to other AI architectures presents an exciting research avenue. Hybrid approaches—for example, combining FF with Sparse Attention for memory efficiency or integrating it with Gradient Checkpointing to selectively apply backpropagation—could further enhance performance. Future evaluations should also include comparative studies of FF against hybrid paradigms that blend gradient-free and gradient-based techniques, such as Sparse Attention plus Checkpointing, to clarify its unique contributions and trade-offs.

Investigating FF's potential in reinforcement learning models may further improve training stability. Scaling FF for distributed training is another crucial research area; future work should explore multi-GPU/TPU parallelism and its integration into federated learning, enabling decentralized AI systems to leverage FF efficiently. Finally, expanding FF's applications in specialized architectures—such as Generative Adversarial Networks (GANs), Spiking Neural Networks (SNNs), and biomedical AI models (e.g., medical image analysis and drug discovery)—could unlock new frontiers in AI model training.

Future studies will also include component-wise ablation analysis to evaluate the sensitivity of FF when applied selectively to different submodules of the Transformer. Specifically, we propose evaluating the following three variants.

- 1) FF-Attn: Apply FF only to self-attention layers, while training the feedforward layers using backpropagation.
- 2) FF-FFN: Apply FF solely to feedforward networks (FFN), with the self-attention mechanisms trained via backpropagation.
- 3) FF-All: Apply FF to all layers, as implemented in the current study.

This ablation analysis will help determine whether FF impacts the self-attention and feedforward components differently and identify optimal hybrid training strategies. Such insights could lead to a fine-grained integration of FF and BP, optimizing both performance and resource efficiency. By addressing these research directions, FF has the potential to evolve into a fundamental deep-learning training paradigm, eventually replacing backpropagation in future AI methodologies.

## 5 DISCUSSION AND FUTURE DIRECTIONS

The Forward-Forward (FF) algorithm introduces an alternative training approach for Transformer models, aiming to address the limitations of traditional Backpropagation (BP) while improving computational efficiency, training stability, and energy consumption. This section discusses the

trade-offs, challenges, and future research directions necessary for the broader adoption of FF-based models.

### 5.1 Trade-offs in Computational Efficiency and Performance

The primary advantage of FF lies in its ability to eliminate backpropagation, reducing computational overhead and memory consumption. Experimental results demonstrate a 20% decrease in training time and 30% lower memory usage than BP. However, this efficiency gain comes at a slight cost—an increase in Perplexity. This trade-off suggests that, while FF improves training efficiency, additional refinements are needed to maintain or surpass BP-based models' performance. Future studies must investigate techniques to optimize FF's performance while preserving its computational benefits.

### 5.2 Limitations and Challenges

Despite its potential, FF faces several challenges that must be addressed before it can be widely adopted. These challenges include its scalability to large datasets, the optimization of its loss function, and its integration with self-attention mechanisms in Transformer models.

First, scalability remains a major concern, as current evaluations have primarily focused on smaller datasets. FF must be tested on large-scale datasets such as GLUE and WMT to more rigorously assess its real-world applicability, where model performance can be systematically analyzed.

Second, optimizing the loss function is essential for improving accuracy across NLP and computer vision tasks. The current goodness function requires further refinement to effectively minimize the difference in goodness values while preserving training stability. Future work should focus on developing adaptive loss functions that enhance FF's robustness.

Finally, integrating FF with the self-attention mechanism in Transformer models presents a technical challenge. Self-attention involves complex token interactions, and FF must be carefully adapted to function seamlessly within this framework. Ensuring that FF-based models achieve comparable performance without disrupting the fundamental operations of Transformers will be critical for their broader adoption.

Addressing these challenges is crucial to enhancing FF's reliability, scalability, and applicability across diverse AI domains. Moreover, in real-world deployments where privacy is critical (e.g., healthcare or federated learning), FF's gradient-free nature may reduce the risk of gradient-based leakage attacks, offering an added layer of model security.

### 5.3 Future Research Directions

Future research should focus on the following key areas to further explore the applicability and effectiveness of the Forward-Forward (FF) algorithm.

- 1) Large-Scale Model Evaluation

FF must be tested on ultra-large models such as GPT-4 and PaLM to assess its feasibility in high-parameter architectures.

Comparative evaluations with BP-based models will be necessary to analyze scalability and efficiency in large-scale neural networks.

#### 2) Resource-Constrained AI Systems

Given FF's reduced memory and computational demands, it could be highly advantageous in IoT devices, edge computing, and real-time AI applications.

Future studies should examine FF's energy efficiency and latency in resource-constrained environments such as autonomous vehicles, embedded AI, and federated learning systems.

#### 3) Hybrid Training Approaches

It can be integrated with Sparse Attention for memory efficiency or Gradient Checkpointing for selective backpropagation to enhance FF's effectiveness.

Additionally, FF's role in reinforcement learning could be explored to investigate its impact on training stability and sample efficiency.

#### 4) Parallelization and Distributed Learning

Scaling FF for multi-GPU/TPU distributed training is crucial for its adoption in large-scale AI applications.

Further research should investigate FF's compatibility with federated learning and decentralized AI models, ensuring its feasibility in privacy-preserving AI training.

#### 5) Expanding Applications beyond Transformers

Beyond Transformers, FF has potential applications in Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Generative Adversarial Networks (GANs).

Applying FF to biomedical AI models, such as medical image analysis and drug discovery, could unlock new possibilities in AI-driven healthcare.

By addressing these research directions, FF could evolve into a scalable, efficient, and broadly applicable training paradigm, potentially replacing backpropagation in future AI methodologies.

### 5.4 Evaluation Metrics for FF Performance Validation

Key performance metrics must be clearly defined to objectively assess the effectiveness of the Forward-Forward (FF) algorithm.

These metrics enable a structured comparison of FF-based training with Backpropagation (BP)-based methods in terms of computational efficiency and model performance.

Tab. 10 presents the evaluation criteria used for experimental validation.

**Table 10** Evaluation metrics for experimental validation

Metric	Purpose	Application
BLEU	Measures translation quality	Machine translation (e.g., WMT)
Perplexity	Evaluates predictive accuracy of language models	Text generation (e.g., GLUE, IMDB)
Latency	Analyzes training and inference speed	All tasks
Accuracy	Benchmarks classification performance	Image classification (e.g., CIFAR-10, ImageNet)

These metrics provide a quantitative basis for evaluating FF's performance across different machine learning tasks.

By analyzing these indicators, researchers can determine whether FF maintains model accuracy while improving computational efficiency, ensuring its viability as an alternative to BP.

## 6 CONCLUSIONS

This study introduced a Forward-Forward (FF)-based training framework for Transformer models as a resource-efficient, gradient-free alternative to Backpropagation (BP). We addressed the growing need for scalable and memory-efficient training methods, particularly for edge and privacy-sensitive applications.

Our contributions include a theoretical formulation of the FF algorithm within Transformer layers, the development of a structured training pipeline complete with pseudocode and a margin-based loss optimization strategy, empirical evaluations across different hardware platforms and model depths, statistical validation of our results along with plans for component-wise ablation studies, and a comparative analysis with other gradient-free methods, such as Genetic Algorithms (GA) and EvoGrad.

While FF does not yet outperform BP in every metric, our results demonstrate that it achieves comparable accuracy with significantly lower computational costs. This validates its feasibility for real-world systems and highlights its potential to redefine deep learning methodologies by reducing training time, memory usage, and energy consumption.

Future research directions include exploring hybrid FF-BP integrations to further enhance performance, establishing formal convergence proofs and theoretical bounds, and implementing FF in privacy-preserving and federated learning scenarios. These advancements mark FF as a promising candidate for next-generation training frameworks in deep learning, paving the way for more accessible and resource-efficient AI solutions.

## 7 REFERENCES

- [1] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., & Polosukhin, I. (2017). Attention is all you need. *Advances in Neural Information Processing Systems*, 30, 5998–6008. <https://doi.org/10.5555/3295222.3295349>
- [2] Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2019). BERT: Pre-training of deep bidirectional transformers for language understanding [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.1810.04805>
- [3] Radford, A., Narasimhan, K., Salimans, T., & Sutskever, I. (2018). Improving language understanding by generative pre-training. *OpenAI*. Retrieved from <https://openai.com/research/language-understanding>
- [4] Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., & Houlsby, N. (2021). An image is worth 16×16 words: Transformers for image recognition at scale [Preprint]. *arXiv*.

- <https://doi.org/10.48550/arXiv.2010.11929>
- [5] Bahdanau, D., Cho, K., & Bengio, Y. (2016). Neural machine translation by jointly learning to align and translate [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.1409.0473>
- [6] Kingma, D. P., & Ba, J. (2017). Adam: A method for stochastic optimization [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.1412.6980>
- [7] He, K., Zhang, X., Ren, S., & Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 770–778). IEEE. <https://doi.org/10.1109/CVPR.2016.90>
- [8] Chen, T., Rubanova, Y., Bettencourt, J., & Duvenaud, D. (2016). Training deep nets with sublinear memory cost [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.1604.06174>
- [9] Chen, Q., Sun, C., Lu, Z., & Gao, C. (2022). Enabling energy-efficient inference for self-attention mechanisms in neural networks. In *IEEE 4<sup>th</sup> International Conference on Artificial Intelligence Circuits and Systems (AICAS)* (pp. 25–28). IEEE. <https://doi.org/10.1109/AICAS4282.2022.9869924>
- [10] Hinton, G. (2022). The forward-forward algorithm: Some preliminary investigations [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.2212.13345>
- [11] Noh, H., Kim, H., & Yoo, S. (2023). Research on forward-forward algorithm. In *Annual Conference of KIPS Proceedings* (pp. 469–470). Korean Information Processing Society.
- [12] Kim, H. (2024). Analyzing the characteristics of gradient descent and non-gradient descent-based algorithms in neural network learning. In *Proceedings of the IEEE International Conference on Electrical, Computer and Energy Technologies (ICECET)*. IEEE. <https://doi.org/10.1109/ICECET61485.2024.10698114>
- [13] Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D. M., Wu, J., Winter, C., & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877–1901. <https://doi.org/10.5555/3455716.3455856>
- [14] Real, E., Aggarwal, A., Huang, Y., & Le, Q. V. (2019). Regularized evolution for image classifier architecture search. *Proceedings of the AAAI Conference on Artificial Intelligence*, 33(1), 4780–4789. <https://doi.org/10.1609/aaai.v33i01.33014780>
- [15] Zoph, B., & Le, Q. V. (2017). Neural architecture search with reinforcement learning [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.1611.01578>
- [16] Child, R., Gray, S., Radford, A., & Sutskever, I. (2019). Generating long sequences with sparse transformers [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.1904.10509>
- [17] Ham, T. J., Jung, S., Kim, S., Oh, Y. H., Park, Y., Song, Y., Park, J.-H., Lee, S., Park, K., & Lee, J. W. (2020). A<sup>3</sup>: Accelerating attention mechanisms in neural networks with approximation. In *IEEE International Symposium on High Performance Computer Architecture (HPCA)* (pp. 328–341). IEEE. <https://doi.org/10.1109/HPCA47549.2020.00035>
- [18] Xie, Y. (2023). Efficient disentangled attention network for semantic segmentation. In *IEEE International Conference on Sensors, Electronics and Computer Engineering (ICSECE)* (pp. 213–217). IEEE. <https://doi.org/10.1109/ICSECE58870.2023.10263368>
- [19] Pham, H., Guan, M., Zoph, B., Le, Q. V., & Dean, J. (2018). Efficient neural architecture search via parameter sharing [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.1802.03268>
- [20] Salimans, T., & Kingma, D. P. (2016). Weight normalization: A simple reparameterization to accelerate training of deep neural networks [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.1602.07868>
- [21] Yao, M., Zhao, G., Zhang, H., Hu, Y., Deng, L., Tian, Y., Xu, B., & Li, G. (2022). Attention spiking neural networks [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.2209.13929>
- [22] Wang, Q., Wu, B., Zhu, P. F., Li, P., Zuo, W., & Hu, Q. (2020). ECA-Net: Efficient channel attention for deep convolutional neural networks. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)* (pp. 11531–11539). IEEE. <https://doi.org/10.1109/CVPR42600.2020.01155>
- [23] Gandhi, S., Gala, R., Kornberg, J., & Sridhar, A. (2023). Extending the forward-forward algorithm [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.2307.04205>
- [24] Aktemur, E., Zorlutuna, E., Bilgili, K., Bok, T. E., Yanikoglu, B., & Mutluergil, S. O. (2024). Going forward-forward in distributed deep learning [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.2404.08573>
- [25] Reyes-Angulo, A., & Paheding, S. (2024). Forward-forward algorithm for hyperspectral image classification. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE. <https://doi.org/10.1109/CVPRW63382.2024.00321>
- [26] Gao, Y., & Zhang, Y. (2023). The predictive forward-forward algorithm [Preprint]. *arXiv*. <https://doi.org/10.48550/arXiv.2301.01452>

**Authors' contacts:**

**Hyun Jung Kim**, Assistant Professor  
Sang-Huh College and the Graduate School of Information & Communication,  
Department of Convergence Information Technology (Artificial Intelligence Major),  
Konkuk University, Republic of Korea  
nygirl@konkuk.ac.kr

**Sang Hyun Yoo**, Assistant Professor  
(Corresponding author)  
Department of Computer Software, Kyungmin University,  
545, Seo-ro, Uijeongbu-si, 11618 Gyeonggi-do, Republic of Korea  
simonyoo@kyungmin.ac.kr