

IJET: Efficient Numerical Transformer via Implicit Iterative Euler Method

Xinyu Liu^{1*}, Bei Li^{2*†}, Jiahao Liu², Junhao Ruan¹,
Kechen Jiao³, Hongyin Tang², Jingang Wang², Xiao Tong^{1,4†}, Jingbo Zhu^{1,4}

¹ NLP Lab, School of Computer Science and Engineering, Northeastern University, Shenyang, China

² Meituan Inc. ³ Tsinghua University ⁴ NiuTrans Research, Shenyang, China

lxy1051493182@gmail.com, libei17@meituan.com

{xiaotong, zhujingbo}@mail.neu.edu.com

Abstract

High-order numerical methods enhance Transformer performance in tasks like NLP and CV, but introduce a performance-efficiency trade-off due to increased computational overhead. Our analysis reveals that conventional efficiency techniques, such as distillation, can be detrimental to the performance of these models, exemplified by PCformer. To explore more optimizable ODE-based Transformer architectures, we propose the Iterative Implicit Euler Transformer (IIET), which simplifies high-order methods using an iterative implicit Euler approach. This simplification not only leads to superior performance but also facilitates model compression compared to PCformer. To enhance inference efficiency, we introduce Iteration Influence-Aware Distillation (IIAD). Through a flexible threshold, IIAD allows users to effectively balance the performance-efficiency trade-off. On lm-evaluation-harness, IIET boosts average accuracy by 2.65% over vanilla Transformers and 0.8% over PCformer. Its efficient variant, E-IIET, significantly cuts inference overhead by 55% while retaining 99.4% of the original task accuracy. Moreover, the most efficient IIET variant achieves an average performance gain exceeding 1.6% over vanilla Transformer with comparable speed.

1 Introduction

The integration of advanced numerical ordinary differential equation (ODE) solvers into Transformer architectures (Vaswani, 2017) has spurred significant progress in natural language processing (NLP) (Li et al., 2022, 2024; Tong et al., 2025) and image synthesis (Ho et al., 2020; Lu et al., 2022a,b; Zheng et al., 2024). Leveraging high-order methods, particularly Predictor-Corrector (PC) schemes, within Transformer residual connections has demonstrated the capacity to enhance model learning without

increasing parameter counts, offering a pathway to both performance and parameter efficiency (Li et al., 2022, 2024).

However, the promise of high-order PCformer (Li et al., 2024) is often constrained by deployment inefficiencies. The inherent linear dependency in nested computations across layers during inference poses critical inference latency. A straightforward approach to mitigating this deployment bottleneck is Knowledge Distillation (Hinton, 2015; Kim and Rush, 2016). However, our preliminary experiments demonstrate that the inherent architectural discrepancy between the predictor and corrector within PCformer impedes effective knowledge transfer via distillation. Our empirical investigations reveal an obvious 54% loss in performance advantage for distilled student models, even for those initialized with PCformer parameters.

Confronted with these deployment bottlenecks, we pivot towards architectural innovations grounded in numerical method principles. A naive yet seemingly logical initial approach might be to pursue uniformity in numerical methods between predictor and corrector, such as pairing explicit and backward Euler schemes. Similar attempts have been validated in previous studies (Li et al., 2024; Zhao et al., 2024), where a high-order predictor combined with a single-step backward Euler method demonstrated promising results, particularly on smaller datasets. However, ensuring solution precision inherently requires iterative solvers to obtain the final solution, a process that shares the same merits as high-order methods. Building on this insight, we take a step further to explore whether an iterative corrector mechanism is equally critical for achieving both superior solution fidelity and unlocking genuine efficiency gains.

To this end, we introduce the Iterative Implicit Euler Transformer (IIET). Concretely, in IIET, each iteration represents a computational step within an implicit Euler iterative solver, where multiple

*These authors contributed equally to this work.

†Corresponding authors.

corrections to the initial prediction are made to ensure output precision. To further strengthen numerical stability, we also employ linear multi-step methods during each correction step. This architecture, detailed in Figure 1d, is designed not only to achieve superior performance that scales with increasing iterations, exhibiting competitive results against PCformer, but also to be inherently compressible due to its iterative nature. Notably, our top-performing IIET models demonstrate a growing advantage over equivalent vanilla Transformers as they scale, achieving improvements of 2.4% (340M), 2.9% (740M), and a more substantial 5.6% for the 1.3B model.

In this way, we can effectively accelerate the inference of IIET via distillation techniques. Here, we further propose Iteration Influence-Aware Distillation (IIAD), a method inspired by structured pruning techniques (Men et al., 2024; Chen et al., 2024), to reduce dispensable iterations. Specifically, IIAD first assesses “iteration influence” by calculating input-output similarity for each iteration. The optimal number of iterations per layer is then determined according to a predefined influence threshold. Subsequently, a continued pre-training phase is employed to restore the model’s capabilities. This process enables users to tailor the iterative correction steps of the IIET model according to their computational budget, yielding efficient IIET variants. Experiments demonstrate that our efficient variant, E-IIET, reduces the inference computational cost of IIET by over 60% while retaining 98.7% of its performance. The lower bound of our efficient IIET variants not only outperforms the vanilla Transformer by an average of 1.6 points but also matches its inference efficiency, showcasing a significant advancement in both performance and deployment efficiency.

2 Background

We begin by establishing the connection between residual connections and the Euler method, and then discuss Transformer optimization strategies informed by advanced explicit and implicit numerical solutions of ODEs. Our work builds upon the standard Transformer architecture (Vaswani, 2017), which comprises a stack of identical layers. For language modeling, each layer typically comprises a causal attention (CA) block and a feed-forward network (FFN) block. With residual connections, the output of each block can be formu-

lated as $y_{n+1} = y_n + \mathcal{F}(y_n, \theta_n)$, where $\mathcal{F}(y_n, \theta_n)$ represents the transformation performed by either the CA or FFN block with parameters θ_n .

2.1 Euler Method in Residual Networks

The Euler method provides a linear approximation for first-order ODEs, defined as $y'(t) = f(y(t), t)$ with an initial value $y(t_0) = y_0$. Given a step size h where $t_{n+1} = t_n + h$, the method computes the subsequent value y_{n+1} as:

$$y_{n+1} = y_n + hf(y_n, t_n) \quad (1)$$

where $f(y_n, t_n)$ represents the rate of change of y , determined by its current value and time t . Notably, this formulation shares a structural similarity with residual networks, where a trainable function, $\mathcal{F}(\cdot)$, approximates these changes. Consequently, from an ODE perspective, residual connections can be interpreted as a first-order discretization of the Euler method. Although the success of residual connections highlights the benefits of the Euler method, its first-order nature introduces significant truncation errors (Li et al., 2022, 2024), limiting the precision of y_{n+1} . Fortunately, more advanced numerical methods exist and have been successfully applied to neural networks.

2.2 Advanced Numerical Transformers

To improve the precision of y_{n+1} , the Runge-Kutta (RK) method offers a more accurate alternative. Inspired by the o -order RK method, the ODE Transformer (Li et al., 2022) replaces residual connections with a RK process:

$$y_{n+1} = y_n + \sum_{i=1}^o \gamma_i \mathcal{F}_i \quad (2)$$

$$\mathcal{F}_1 = \mathcal{F}(y_n, \theta_n) \quad (3)$$

$$\mathcal{F}_i = \mathcal{F}(y_n + \sum_{j=1}^{i-1} \beta_{ij} \mathcal{F}_j, \theta_n) \quad (4)$$

where \mathcal{F}_i represents the i^{th} order results computed by a shared transformer block $\mathcal{F}(*, \theta_n)$. The coefficients γ_i, β_{ij} are learnable parameters. This architecture effectively mitigates truncation error, leading to significant performance gains in generation tasks such as machine translation.

Compared to explicit numerical methods, implicit numerical methods typically offer higher precision and stability. The Predictor-Corrector (PC) method, using an explicit predictor for initial estimates and an implicit corrector for refinement, is a classic example. Recent work has demonstrated the

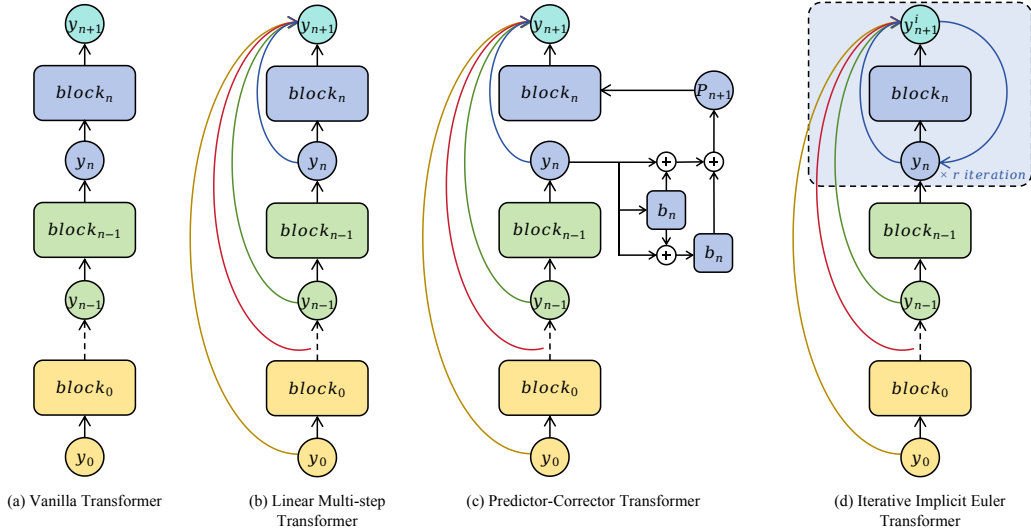


Figure 1: Architectural comparison: (a) Vanilla Transformer; (b) Linear multistep-enhanced Transformer; (c) PCformer with 2nd-order Runge-Kutta predictor and 1st-order Euler corrector; (d) Our proposed Iterative Implicit Euler Transformer (IIET). The iteration steps r in IIET is configurable, with experimental validation determining $r = 3$ as the optimal setting in this work. All blocks follow an identical computational procedure as the $block_n$.

benefits of integrating PC components into neural network architecture. PCformer (Li et al., 2024) employs an o -order RK predictor and a linear multi-step (Wang et al., 2019) corrector, defined as:

$$y_p = y_n + \sum_{i=1}^o \gamma(1-\gamma)^{o-i} \mathcal{F}_i \quad (5)$$

$$y_{n+1} = y_n + \alpha \mathcal{F}(y_p, \theta_n) + \sum_{i=n-2}^n \beta \tilde{\mathcal{F}}_i \quad (6)$$

where \mathcal{F}_i shares the same meaning as in Eq. 2 and $\tilde{\mathcal{F}}_i$ denotes the outputs of previous blocks. α, β , and γ are learnable coefficients. Specifically, PCformer’s predictor incorporates an Exponential Moving Average (EMA) to weight the contributions of different orders, while the corrector integrates previous block outputs for increased precision. PCformer achieves superior performance over the ODE Transformer and, to some extent, unifies structural paradigms for Transformers improved with implicit numerical methods. Our IIET can be interpreted as a specific instance within the PC paradigm, with a particular emphasis on the iterative corrector component.

3 Iterative Implicit Euler Transformer

In this section, we detail the theoretical foundation and core architectural design of the Iterative Implicit Euler Transformer (IIET). Our approach leverages the inherent stability of the implicit Euler method, a cornerstone of numerical analysis, to address key challenges in deep sequence modeling.

3.1 Iterative Implicit Euler Method

The implicit Euler method, also known as the backward Euler method, is a foundational first-order implicit numerical technique celebrated for its robust stability properties, particularly advantageous in handling stiff systems (LeVeque, 2007). Unlike its explicit counterparts, the implicit Euler method employs a backward difference quotient, formulated as:

$$y_{n+1} = y_n + hf(y_{n+1}, t_{n+1}). \quad (7)$$

The implicit nature of Eq. 7, where the computation of y_{n+1} depends on its value at the same time step t_{n+1} , inherently requires iterative solvers from numerical analysis to obtain a solution. Specifically, in traditional numerical methods for solving such implicit equations, Newton’s iteration is frequently employed due to its quadratic convergence rate and robustness (Zhang et al., 2017; Shen et al., 2020; Kim et al., 2024). However, within the context of neural sequence modeling, where computational efficiency and architectural simplicity are often prioritized, we propose to investigate the efficacy of a simpler alternative: fixed-point iteration (Rhoades, 1976). While prior works like Li et al. (2024) have utilized explicit methods for initial approximations followed by a single backward Euler correction, the potential of iterative refinement within the implicit corrector remains largely unexplored.

Thus, challenging the implicit assumption that a strong predictor is sufficient for high precision

(Li et al., 2024), we propose the central hypothesis that iterative refinement inside the implicit corrector constitutes a pivotal mechanism for enhancing solution fidelity. We argue that a single-step correction inherently limits the achievable accuracy, particularly when modeling intricate sequence dynamics and seeking high-fidelity representations of y_{n+1} . Consequently, this work rigorously investigates whether leveraging iterative solutions within the implicit corrector can translate to demonstrable gains in downstream model performance.

Intriguingly, our empirical findings reveal that computationally efficient fixed-point iteration yields surprisingly high precision, particularly within our neural sequence modeling framework. Our proposed Iterative Implicit Euler (IIE) method commences with an initial approximation, y_{n+1}^0 , derived from an explicit Euler step. This initial estimate is then iteratively refined through r fixed-point iterations as defined below:

$$y_{n+1}^0 = y_n + hf(y_n, t_n) \quad (8)$$

$$y_{n+1}^i = y_n + hf(y_{n+1}^{i-1}, t_{n+1}), \quad i \in [1..r]. \quad (9)$$

The final approximation y_{n+1} is thus given by y_{n+1}^r , representing the output of the r^{th} iteration.

The IIE method, while formally retaining its first-order numerical accuracy, achieves a significant enhancement in the approximation of y_{n+1} through iterative refinement. This iterative process engenders a structured form of nested computations that superficially resemble higher-order methods, albeit through a fundamentally distinct mechanism rooted in repeated fixed-point iterations. Acknowledging the increased computational cost, the inherent structural regularity of IIE, predicated solely on the preceding iteration’s output, emerges as a crucial enabler for inference efficiency optimizations, as detailed in Section 4. This carefully engineered balance between iteratively enhanced precision and structural simplicity underpins the design philosophy of the IJET architecture.

3.2 Model Architecture

Building on the IIE method, we propose the Iterative Implicit Euler Transformer (IJET) as a foundational architecture for sequence modeling, particularly for large language models. Adopting the LLaMA architecture (Touvron et al., 2023b) (Transformer++), IJET consists of N stacked transformer decoder layers. Each layer comprises a causal attention module followed by a feedforward module, and employs rotary positional encoding (Su

Algorithm 1 Iterative Implicit Euler Paradigm

```

1: procedure IJET BLOCK( $y_n, \theta_n, \mathbf{L}, r$ )
2:    $\mathbf{L}$  is the global stack for historical hidden states
3:    $f_n^0 \leftarrow \mathcal{F}(y_n, \theta_n)$   $\triangleright$  Compute initial value
4:    $\mathbf{L}.append(f_n^0)$   $\triangleright$  Store the initial context
5:   for  $i \leftarrow 0$  to  $r - 1$  do  $\triangleright$  Refinement loop
6:     Compute  $y_{n+1}^{i+1}$  using  $\mathbf{L}$  via the update rule
7:      $f_n^{i+1} \leftarrow \mathcal{F}(y_{n+1}^{i+1}, \theta_n)$   $\triangleright$  Compute refined value
8:      $\mathbf{L}.update(f_n^i \rightarrow f_n^{i+1})$   $\triangleright$  Update the context  $\mathbf{L}$ 
9:   end for
10:  Compute  $y_{n+1}^r$  using the final context  $\mathbf{L}$ 
11:  return  $y_{n+1}^r$   $\triangleright$  Return the final layer output
12: end procedure

```

et al., 2024), SiLU activation (Shazeer, 2020), and RMS normalization (Zhang and Sennrich, 2019). Given an input sequence $x = x_1, \dots, x_L$ of length L , the initial input embeddings are represented as $X^0 = [x_1, \dots, x_L] \in \mathbb{R}^{L \times d_{\text{model}}}$, where d_{model} is the hidden dimension. The output of each subsequent layer is then computed as $X^n = \text{Decoder}(X^{n-1})$, for $n \in [1, N]$.

The key distinction between IJET and Transformer++ lies in IJET’s integration of the IIE method within each decoder layer (Figure 1). Unlike Transformer++, which directly computes the layer’s output using a single Euler step (standard residual), IJET employs an iterative refinement process. Specifically, IJET first estimates an initial value, y_{n+1}^0 , via a single Euler step (Eq. 8):

$$y_{n+1}^0 = y_n + \mathcal{F}(y_n, \theta_n) \quad (10)$$

where $\mathcal{F}(*, \theta_n)$ represents the n^{th} transformer layer with parameters θ_n . This initial estimate in IJET corresponds to the direct output of each layer in Transformer++.

In the subsequent iterations, our preliminary experiments suggest that incorporating outputs from previous layers, similar to Transformer-DLCL (Wang et al., 2019), can enhance the performance. We thus modify Eq. 9 as follows:

$$y_{n+1}^i = y_n + \alpha_n \mathcal{F}(y_{n+1}^{i-1}, \theta_n) + \sum_{j=0}^{n-1} \alpha_j \tilde{\mathcal{F}}_j \quad (11)$$

where $i \in [1..r]$ denotes the iteration step, $\tilde{\mathcal{F}}_j$ represents the output of the previous layers j , and α represents learnable layer merge coefficients.

Algorithm 1 further details the computational flow within a single IJET layer. Specifically, the matrix \mathbf{L} stores the hidden states from previously computed layers, thereby providing the necessary historical context. Within a single block, an initial

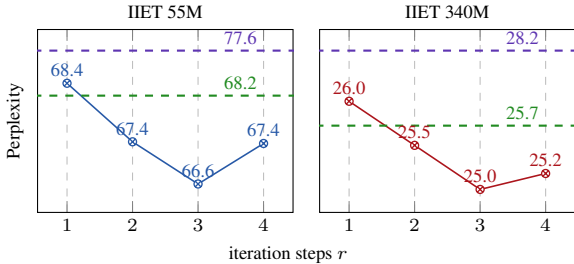


Figure 2: PPL on the Wikitext test set for 55M and 340M IJET across varying iteration steps r . Dashed lines indicate Transformer++ and PCformer performance at corresponding parameter scales. Note that IJET’s FLOPs is nearly $r + 1$ times of Transformer++.

estimate of the output, denoted as y_{n+1}^0 , is generated and then iteratively refined. Each iteration i updates this estimate to y_{n+1}^i by leveraging the function \mathcal{F} and the context \mathbf{L} . This fixed-point iteration process ensures that the hidden state progressively converges to a more precise and stable final output, y_{n+1}^r .

3.3 Experimental Setups

Baselines. We evaluate IJET’s performance against two strong baselines: Transformer++ (Touvron et al., 2023a) and PCformer (Li et al., 2024). Transformer++ adopts the LLaMA architecture. PCformer employs a 2nd-order Runge-Kutta predictor and a linear multi-step corrector¹. We train all models from scratch at three parameter scales: 340M, 740M, and 1.3B. All models are trained on the same dataset with identical token counts to ensure controlled comparison. Detailed training hyperparameter settings can be found in Appendix A.1.

Datasets and Evaluation Metrics. Our models are pre-trained on SlimPajama (Soboleva et al., 2023) and tokenized using the LLaMA2 tokenizer (Touvron et al., 2023a). From the original 627B-token dataset, we sample 16B, 30B and 100B tokens for training the 340M, 740M and 1.3B parameter models, respectively. For comprehensive evaluation, we assess perplexity (PPL) on Wikitext (Wiki.) (Merity et al., 2016) and consider several downstream tasks covering common-sense reasoning and question answering: LAMBADA (LMB.) (Paperno et al., 2016), PiQA (Bisk et al., 2020), HellaSwag (Hella.) (Zellers et al., 2019), WinoGrande (Wino.) (Sakaguchi et al., 2021),

¹We also explored a 4th-order Runge-Kutta predictor and more complex correctors, but these increased training costs without substantially improving performance.

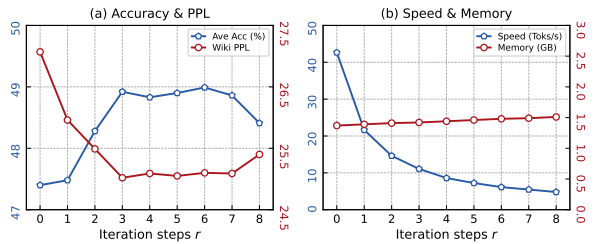


Figure 3: Ablation study on iteration steps r : (a) Impact on model performance. (b) Corresponding effects on inference speed and VRAM utilization.

ARC-Challenge (ARC-c) (Clark et al., 2018), and SCIQ (Welbl et al., 2017). We report PPL on Wikitext and LAMBADA; length-normalized accuracy on HellaSwag, ARC-Challenge, and PiQA; and standard accuracy on the remaining tasks. All evaluations are conducted using the lm-evaluation-harness (Gao et al., 2021). In addition to language modeling, we also conduct experiments on machine translation and summarization tasks, with detailed results in Appendix B.

Ablation Study on Iteration Steps.

3.4 Experimental Results

Iteration Steps. To identify the optimal iteration steps r , we first apply varying r values to the 340M IJET model and a smaller 55M parameter variant (detailed in Appendix A.1). All models were evaluated on Wikitext test set. As illustrated in Figure 2, which showcases the benefit of iterative correction, IJET’s performance exceeds PCformer at $r = 2$ and achieves its peak at $r = 3$. Therefore, we adopt $r = 3$ in this work.

Results. IJET’s advantages are clearly demonstrated on LLM evaluation benchmarks. As shown in Table 1 (*Pre-training Phase*), IJET consistently surpasses Transformer++ at comparable parameter scales. At the 340M scale, IJET achieves a mean accuracy 2.4 points higher than Transformer++ across six challenging subtasks. This performance gap widens with model size, reaching 2.9 points at 740M and a substantial 5.6 points at 1.3B parameters. Moreover, IJET’s performance also matches or exceeds that of PCformer across all tested scales, demonstrating the advantages of the iterative correction paradigm. This strong scaling behavior, consistent with the findings in Li et al. (2024)’s work, confirms the robust scalability of IJET and showcases its potential with larger models and datasets.

Scale	Model	Wiki. ppl ↓	LMB. ppl ↓	LMB. acc ↑	PiQA acc_norm ↑	Hella. acc_norm ↑	SCIQ acc ↑	ARC-c acc_norm ↑	Wino. acc ↑	Avg. ↑
<i>Pre-training Phase</i>										
<i>340M Params</i> <i>16B Tokens</i>	Transformer++	28.2	78.3	28.9	64.3	34.2	76.0	23.6	51.9	46.5
	PCformer	25.7	47.0	33.1	64.9	36.3	77.5	24.7	53.3	48.3
	IJET	25.0	30.5	37.1	65.2	36.9	79.4	23.9	51.0	48.9
<i>740M Params</i> <i>30B Tokens</i>	Transformer++	23.3	34.8	36.1	66.4	38.4	78.6	24.5	50.2	49.0
	PCformer	21.2	22.0	41.0	66.3	41.3	82.0	23.3	51.2	50.9
	IJET	20.7	21.1	41.2	68.9	42.5	82.1	23.8	53.1	51.9
<i>1.3B Params</i> <i>100B Tokens</i>	Transformer++	16.3	11.8	51.6	71.0	51.7	86.7	28.1	54.6	57.2
	PCformer	14.0	7.9	59.6	73.8	60.0	90.7	31.7	61.7	62.9
	IJET	14.0	7.8	59.8	73.7	60.5	88.6	32.3	61.6	62.8
<i>Iteration Influence-Aware Distillation Phase</i>										
<i>340M Params</i> <i>5B Tokens</i>	Distil PCformer	27.2	50.4	32.2	64.6	34.9	78.0	24.7	51.3	47.6
	Lower Bound	27.0	34.6	36.1	64.0	35.0	80.7	23.0	51.5	48.4
	E-IJET	25.7	30.9	37.4	64.4	35.8	80.4	23.5	52.1	48.9
<i>740M Params</i> <i>10B Tokens</i>	Distil PCformer	22.5	29.5	37.4	66.8	39.2	80.0	23.2	50.9	49.6
	Lower Bound	23.0	29.9	37.6	67.4	38.7	79.7	25.2	53.0	50.3
	E-IJET	21.2	24.2	40.1	68.5	41.0	81.0	24.6	52.4	51.3
<i>1.3B Params</i> <i>30B Tokens</i>	Lower Bound	15.9	9.3	57.0	72.7	56.7	87.7	29.6	59.1	60.5
	E-IJET	14.8	8.9	57.5	73.2	58.4	88.6	30.7	59.4	61.3

Table 1: Comparison of results between our models and baselines in the *Pre-training Phase* and *Iteration Influence-Aware Distillation Phase*. The individual task performance is via zero-shot. We report the main results on the same set of tasks reported by Gu and Dao (2023). The last column shows the average over all benchmarks that use (normalized) accuracy as the metric. **Bold** values represent the best results in each set.

3.5 Analysis

A key question concerning IJET is whether its performance improves monotonically with an increasing number of iterative correction steps. To investigate this, we conducted an ablation study on the 340M IJET model, varying the number of iteration r .² As illustrated in Figure 3a, performance initially improves with increasing r . However, beyond a certain threshold, further increases in r lead to a plateau in performance gains. This suggests that the iterative refinement process guides the final representation towards a more precise ODE solution, but with diminishing returns after optimal convergence. Detailed downstream results can be found in Appendix C. Furthermore, to assess the impact of r on inference efficiency, we measured the autoregressive generation throughput of IJET variants on a single A100 GPU. Figure 3b shows that while IJET’s inference speed substantially declines with increasing r , its VRAM footprint remains largely unaffected as it incurs no extra parameters.

Comparison with Equal FLOPs. Given that IJET’s iterative correction adds FLOPs (to approximately four times that of Transformer++ when

²In the case where $r = 0$, IJET is structurally the same as the DLCL Transformer.

Model	LMB.	PiQA	Hella.	SCIQ	ARC-c	Wino.	Avg.
IJET	37.1	65.2	36.9	79.4	23.9	51.0	48.9
Trans WS	30.7	63.1	34.4	75.7	23.2	50.4	46.3
Trans 1.3B	37.3	65.7	37.6	78.6	23.7	51.5	49.0

Table 2: Performance comparison of models with FLOPs comparable to the 340M IJET.

$r = 3$), we aimed for a performance comparison under equivalent computational budgets. Thus, we trained a 1.3B Transformer++ model on identical 16B training data. The results in Table 2 show that IJET performs comparably to the much larger Transformer++ but with substantially fewer parameters, thereby reducing memory and training overhead. To ensure a fair comparison based on model size, we compressed the 1.3B Transformer++ model using pruning and quantization to align with IJET’s storage requirements. The results of this evaluation are detailed in Appendix D. Moreover, models with exactly matched parameter scale and FLOPs were benchmarked. Since IJET’s architecture closely resembles weight-sharing methods, we established a naive weight-sharing baseline: the Transformer++ model’s depth was quadrupled, with weights shared every four layers, namely Trans WS. As shown in Table 2, this simple weight-sharing approach alone does not yield per-

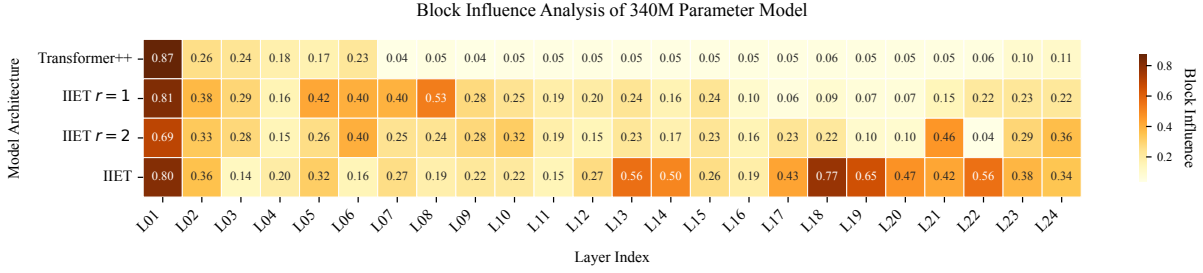


Figure 4: Distribution of Block Influence (BI) for Transformer++ and IJET models with varying iteration steps r . Higher BI values indicate lower model redundancy.

formance gains, highlighting the crucial contribution of IJET’s implicit iterative solver-based design to its enhanced performance.

Parameter Redundancy of IJET. We hypothesize that the iterative correction process of IJET enhances learning efficiency and reduces parameter redundancy. To investigate this, we used Block Influence (BI) (Men et al., 2024) to measure layer redundancy in IJET and Transformer++. BI assesses the influence of each model block on the hidden state by measuring the similarity between its input and output; lower similarity indicates a higher influence. Specifically, the BI of a Transformer block is calculated as:

$$BI_i = 1 - \frac{\mathbb{E}_{\mathbf{H},t} \mathbf{H}_{i,t}^T \mathbf{H}_{i+1,t}}{\|\mathbf{H}_{i,t}\|_2 \|\mathbf{H}_{i+1,t}\|_2} \quad (12)$$

where $\mathbf{H}_{i,t}$ represents the t^{th} row of the i^{th} layer’s input hidden states. We randomly sampled 5,000 text segments from Wikitext to calculate the BI of each model. As shown in Figure 4, the influence of IJET’s blocks increases significantly with iteration steps, demonstrating higher layer utilization. This also indicates that the learning potential of existing large-scale language models remains under-exploited.

4 Iteration Influence-Aware Distillation

While IJET achieves strong downstream task performance, its iterative structure introduces computational overhead that curtails inference speed. This added latency is particularly non-negligible for autoregressive generation in large language models. To enhance IJET’s inference efficiency without performance loss, we explore whether continuous pre-training combined with distillation can enable fewer forward passes, ideally a single one, to yield outputs equivalent to those from the complete, multi-step iterative correction process. To

this end, we analyze the impact of each iterative correction step on the hidden state within each block. Surprisingly, Figure 5 shows that not all layers require the same number of iteration steps to achieve accurate output, with deeper layers benefiting more from additional iterative corrections, which is potentially due to the varying roles layers play in the Transformer’s representation-building process.

4.1 Methodology

In this section, we propose Iteration Influence-Aware Distillation (IIAD). IIAD first analyzes the iterative process of a pre-trained IJET, identifying and eliminating non-essential iterative computations to yield an efficient variant, E-IJET. Subsequently, a layer-wise self-distillation phase restores the performance of E-IJET.

Iteration Influence. Iteration influence employs a computational methodology similar to block influence; however, its calculation is performed specifically within individual IJET blocks. For a given n^{th} block, we consider its input y_n and the output y_{n+1}^i of each internal iteration i . The pairwise differences between these representations are calculated using Eq. 12 to obtain the iteration influence values. Based on these values and a specified computational budget, users can determine the number of iteration steps to retain per block.

In this work, we primarily investigate two designs for efficient IJET variants: **1 Lower Bound:** Each layer performs only a single forward pass, establishing a performance lower bound for efficient IJET. **2 E-IJET:** This variant establishes a threshold using the minimum of the initial iteration influence values computed in each layer. Consequently, iteration steps with influence scores below this threshold are omitted, preserving each layer’s initial computation and essential iteration steps. Specifically, E-IJET reduces the number of iteration steps from a baseline of 72 to 15 in the 340M

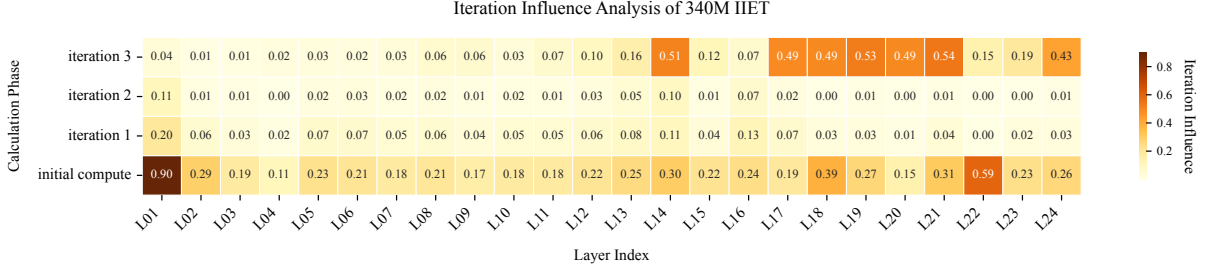


Figure 5: **Iteration Influence** within each layer of the 340M IIET model. Deeper colors indicate larger hidden state changes after this iteration. The 740M IIET results are presented in Appendix E due to space constraints.

variant and 23 in the 740M variant.

Iteration Influence-Aware Distillation. In the continuous pre-training stage, we employ a warm-start initialization strategy, directly inheriting parameters from the pre-trained IIET model to retain knowledge acquired during its initial pre-training phase. To enable efficient IIET variants (e.g., E-IIET) to approximate the precise output representations of the full IIET, we utilize a fine-grained, block-specific knowledge distillation framework incorporating two complementary losses: **1) Mean Squared Error (MSE) Loss:** For each block, an MSE loss encourages E-IIET to mimic the refined hidden states produced by the full IIET. This loss is computed as:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{n} \sum_{i=1}^n \|\mathbf{h}_i^{\text{IIET}} - \mathbf{h}_i^{\text{E-IIET}}\|_2^2 \quad (13)$$

where \mathbf{h}_i are the hidden state outputs of the i^{th} block. **2) Kullback-Leibler (KL) Loss:** To further align prediction behavior, we compute the KL divergence between the final output probability distributions of the full IIET and E-IIET:

$$\mathcal{L}_{\text{KL}} = D_{\text{KL}}(p(\mathbf{z}^{\text{IIET}}/\tau) \| p(\mathbf{z}^{\text{E-IIET}}/\tau)) \quad (14)$$

where \mathbf{z} represent the output logits and τ is the distillation temperature. By combining these two distillation losses with Cross-Entropy loss, we train E-IIET to effectively capture the knowledge embedded within the full IIET’s iterative refinement process. The final training objective for this continuous pre-training stage is thus:

$$\mathcal{L}_{\text{E-IIET}} = \mathcal{L}_{\text{CE}} + \alpha \mathcal{L}_{\text{MSE}} + \beta \mathcal{L}_{\text{KL}} \quad (15)$$

4.2 Experiments and Results

Setups. To train efficient IIET variants, we sample one-third of the original pre-training tokens (see Appendix A.2 for detailed training settings). For

Model	340M			740M		
	Spd.	FLOPs	VRAM	Spd.	FLOPs	VRAM
Transformer++	49.97	0.38	1.37	48.91	0.80	2.80
PCformer	14.14	1.06	1.41	14.38	2.30	2.86
IIET	11.07	1.40	1.42	10.95	3.05	2.89
Lower Bound	42.66	0.38	1.37	42.03	0.80	2.80
E-IIET	25.95	0.60	1.38	22.12	1.52	2.83

Table 3: A comparison of inference speed (tokens per second), FLOPs (T) and VRAM (GB) for baseline models, PCformer, and efficient IIET variants.

performance comparison against E-IIET, we also prepare two key baselines: a *Lower Bound* variant, which omits all iterative corrections, and a distilled version of PCformer. All models are trained following the method outlined in Section 4.1.

Main Results. Table 1 presents the main results for IIAD. As a baseline, directly distilling PCformer into a standard Euler architecture (namely *Distil PCformer*) leads to substantial performance degradation, highlighting the importance of the sophisticated numerical solvers employed by higher-order methods to achieve their accuracy. In contrast, E-IIET, compared to the full IIET model, retains the vast majority of its performance while reducing the average iterative correction overhead by about 55%. Importantly, even the *Lower Bound* efficient IIET variant achieves performance on par with PCformer, demonstrating IIET’s strength in balancing efficiency with strong performance.

Inference Efficiency. We analyze the inference speed, FLOPs and VRAM usage of our main models. As Table 3 indicates, E-IIET achieves over a 2x speedup compared to full IIET, while largely maintaining IIET’s performance advantage (E-IIET vs. full IIET scores: 48.9/48.9 for 340M and 51.3/51.9 for 740M model). However, due to the FLOPs incurred by its remaining iteration steps, E-IIET still exhibits nearly twice the inference latency of Trans-

former++. A key characteristic of these efficient IIET variants is the inverse relationship between performance and efficiency: fewer iterations lead to lower performance but higher efficiency. Notably, Table 3 shows that the maximum efficiency attained by these variants (i.e., *Lower Bound*) is close to that of the Transformer++, with their average performance surpassing it by 1.6 points. This adaptability makes E-IIET a flexible solution for practical deployment, as users can select the iteration steps based on their resource constraints (e.g., reducing iterations to maximize inference speed).

5 Related Work

Ordinary Differential Equations in Deep Learning The conceptual link between Ordinary Differential Equations (ODEs) and residual networks, first established by Weinan (2017), has catalyzed the development of numerous ODE-inspired neural architectures. In computer vision, this perspective give rise to models such as PolyNet (Zhang et al., 2017), FractalNet (Larsson et al., 2016), Multi-stepNet (Lu et al., 2018), and Momentum Residual Networks (Sander et al., 2021). The ODE viewpoint has also been highly influential in generative modeling, particularly for diffusion models (Ho et al., 2020). For example, Liu et al. (2022) reframed Denoising Diffusion Probabilistic Models (DDPMs) as a process of solving differential equations on manifolds, introducing a pseudo linear multi-step method to improve performance. Building on this, DPM-Solver (Lu et al., 2022a,b) significantly accelerated the sampling process by employing exact ODE solutions and higher-order numerical methods. More recently, ODE principles have been leveraged to enhance Transformer architectures for sequence modeling and generation (Lu et al., 2019; Wang et al., 2019). Dutta et al. (2021) redesigned the Transformer as a more efficient multi-particle dynamic system, while Tong et al. (2025) proposed high-order methods to mitigate error accumulation in first-order ODE blocks. Further advancing this line of work, PCformer (Li et al., 2024) introduced a predictor-corrector framework to boost performance, and they first show the potential of such numerical design in large language model literature. In our work, we build upon this foundation by employing the implicit Euler method, aiming to enhance language modeling performance while achieving a superior trade-off between accuracy and computational efficiency.

Implicit ODE Method Existing approaches that utilize implicit ODE solvers can be broadly categorized into two paradigms. The first, exemplified by seminal works like Neural ODEs (Chen et al., 2018; Zhang et al., 2021) and Deep Equilibrium Models (DEQs) (Bai et al., 2019), represents a significant shift towards implicit deep learning. Neural ODEs model the network as a continuous transformation by using a neural network to parameterize the state’s derivative, which is solved with a numerical ODE solver. In contrast, DEQs define network layers implicitly through an equilibrium point, denoted as $z^* = f(z^*, x)$, which is found iteratively, with gradients computed via implicit differentiation. While subsequent research has explored these models for applications like time-series forecasting, their computational cost remains a significant barrier for large-scale language modeling; for instance, a DEQ model can require over five times the computation to match the performance of a standard Transformer-XL model (240M). More recently, the Neural ODE Transformer (Tong et al., 2025) has achieved performance surpassing that of the vanilla Transformer on language modeling tasks. Another class of methods utilizes advanced implicit numerical solvers to optimize mainstream model architectures (Li et al., 2024, 2020; Shen et al., 2020; Kim et al., 2024). For instance, IE-Skips (Li et al., 2020) modifies the original skip connection in ResNet for robustness. Similarly, IM-BERT (Kim et al., 2024) introduced the use of the implicit Euler method to enhance the adversarial robustness of BERT. Our proposed IIET distinctively focuses on autoregressive generation, and the core formulation is quite different, that we adopted an iterative refinement schema began from a fixed point. To our knowledge we are the first to apply iterative implicit Euler into LLMs paradigm.

6 Conclusions

We propose a novel Transformer architecture, the Iterative Implicit Euler Transformer (IIET), designed for enhanced language modeling performance. IIET leverages the iterative implicit Euler method, providing substantial improvements over vanilla Transformers with a simplified architecture compared to PCformer. Furthermore, we introduce an inference acceleration technique for IIET, which uses self-distillation to prune the iterative process, allowing users to adjust inference efficiency based on their budget.

Limitations

Although the IIAD method is designed to produce efficient IIET variants for inference, the IIAD process itself introduces notable computational overhead during its application. Future research will focus on integrating the determination of layer-specific iteration requirements directly into the pre-training stage. This could facilitate the direct training of inherently efficient IIET models, potentially bypassing a separate, resource-intensive distillation phase.

Beyond optimizing IIET’s per-token efficiency, we also identify a promising avenue for broader application. Current large reasoning models often achieve high performance by generating substantially more tokens than are present in the final answer, leading to significant inference latency. IIET, on the contrary, enhances per token representational power through depth-wise iterative refinement, albeit at an increased per-token computational cost. We hypothesize that this trade-off could be ultimately advantageous in multi-step reasoning tasks: IIET’s more precise computation per token might enable it to generate complete and correct answers in fewer overall autoregressive steps, thereby reducing the total token count and potentially overall latency. Validating this hypothesis, however, necessitates training and evaluating IIET at larger model and data scales, which remains a key direction for future investigation.

Acknowledgments

This work was supported in part by the National Science Foundation of China (Nos. 62276056 and U24A20334), the Yunnan Fundamental Research Projects (No.202401BC070021), the Yunnan Science and Technology Major Project (No. 202502AD080014), and the Program of Introducing Talents of Discipline to Universities, Plan 111 (No.B16009).

References

Shaojie Bai, J Zico Kolter, and Vladlen Koltun. 2019. Deep equilibrium models. *Advances in neural information processing systems*, 32.

Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, and 1 others. 2020. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439.

Ricky TQ Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. 2018. Neural ordinary differential equations. *Advances in neural information processing systems*, 31.

Xiaodong Chen, Yuxuan Hu, and Jing Zhang. 2024. Compressing large language models by streamlining the unimportant layer. *arXiv preprint arXiv:2403.19135*.

Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. 2018. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*.

Subhabrata Dutta, Tanya Gautam, Soumen Chakrabarti, and Tanmoy Chakraborty. 2021. Redesigning the transformer architecture with insights from multi-particle dynamical systems. *Advances in Neural Information Processing Systems*, 34:5531–5544.

Leo Gao, Jonathan Tow, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Kyle McDonell, Niklas Muennighoff, and 1 others. 2021. A framework for few-shot language model evaluation. *Version v0. 0.1. Sept*, 10:8–9.

Albert Gu and Tri Dao. 2023. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*.

Geoffrey Hinton. 2015. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.

Jonathan Ho, Ajay Jain, and Pieter Abbeel. 2020. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851.

Shengding Hu, Yuge Tu, Xu Han, Chaoqun He, Ganqu Cui, Xiang Long, Zhi Zheng, Yewei Fang, Yuxiang Huang, Weilin Zhao, and 1 others. 2024. Minicpm: Unveiling the potential of small language models with scalable training strategies. *arXiv preprint arXiv:2404.06395*.

Mihyeon Kim, Juhyoung Park, and Youngbin Kim. 2024. Im-bert: Enhancing robustness of bert through the implicit euler method. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 16217–16229.

Yoon Kim and Alexander M Rush. 2016. Sequence-level knowledge distillation. *arXiv preprint arXiv:1606.07947*.

Gustav Larsson, Michael Maire, and Gregory Shakhnarovich. 2016. Fractalnet: Ultra-deep neural networks without residuals. *arXiv preprint arXiv:1605.07648*.

Randall J. LeVeque. 2007. *Finite difference methods for ordinary and partial differential equations - steady-state and time-dependent problems*. SIAM.

- Bei Li, Quan Du, Tao Zhou, Yi Jing, Shuhan Zhou, Xin Zeng, Tong Xiao, Jingbo Zhu, Xuebo Liu, and Min Zhang. 2022. Ode transformer: An ordinary differential equation-inspired model for sequence generation. *arXiv preprint arXiv:2203.09176*.
- Bei Li, Tong Zheng, Rui Wang, Jiahao Liu, Qingyan Guo, Junliang Guo, Xu Tan, Tong Xiao, Jingbo Zhu, Jingang Wang, and 1 others. 2024. Predictor-corrector enhanced transformers with exponential moving average coefficient learning. *arXiv preprint arXiv:2411.03042*.
- Mingjie Li, Lingshen He, and Zhouchen Lin. 2020. Implicit euler skip connections: Enhancing adversarial robustness via numerical stability. In *International Conference on Machine Learning*, pages 5874–5883. PMLR.
- Aixin Liu, Bei Feng, Bing Xue, Bingxuan Wang, Bochao Wu, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, and 1 others. 2024. Deepseek-v3 technical report. *arXiv preprint arXiv:2412.19437*.
- Luping Liu, Yi Ren, Zhijie Lin, and Zhou Zhao. 2022. Pseudo numerical methods for diffusion models on manifolds. *arXiv preprint arXiv:2202.09778*.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. 2022a. Dpm-solver: A fast ode solver for diffusion probabilistic model sampling in around 10 steps. *Advances in Neural Information Processing Systems*, 35:5775–5787.
- Cheng Lu, Yuhao Zhou, Fan Bao, Jianfei Chen, Chongxuan Li, and Jun Zhu. 2022b. Dpm-solver++: Fast solver for guided sampling of diffusion probabilistic models. *arXiv preprint arXiv:2211.01095*.
- Yiping Lu, Zhuohan Li, Di He, Zhiqing Sun, Bin Dong, Tao Qin, Liwei Wang, and Tie-Yan Liu. 2019. Understanding and improving transformer from a multi-particle dynamic system point of view. *arXiv preprint arXiv:1906.02762*.
- Yiping Lu, Aoxiao Zhong, Quanzheng Li, and Bin Dong. 2018. Beyond finite layer neural networks: Bridging deep architectures and numerical differential equations. In *International Conference on Machine Learning*, pages 3276–3285. PMLR.
- Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. 2024. Shortgpt: Layers in large language models are more redundant than you expect. *arXiv preprint arXiv:2403.03853*.
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. 2016. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*.
- Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. 2016. The lambda dataset: Word prediction requiring a broad discourse context. *arXiv preprint arXiv:1606.06031*.
- BE Rhoades. 1976. Comments on two fixed point iteration methods. *Journal of Mathematical Analysis and Applications*, 56(3):741–750.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2021. Winogrande: An adversarial winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106.
- Michael E Sander, Pierre Ablin, Mathieu Blondel, and Gabriel Peyré. 2021. Momentum residual neural networks. In *International Conference on Machine Learning*, pages 9276–9287. PMLR.
- Noam Shazeer. 2020. Glu variants improve transformer. *arXiv preprint arXiv:2002.05202*.
- Jiawei Shen, Zhuoyan Li, Lei Yu, Gui-Song Xia, and Wen Yang. 2020. Implicit euler ode networks for single-image dehazing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 218–219.
- Daria Soboleva, Faisal Al-Khateeb, Robert Myers, Jacob R Steeves, Joel Hestness, and Nolan Dey. 2023. Slimpajama: A 627b token cleaned and deduplicated version of redpajama.
- Jianlin Su, Murtadha Ahmed, Yu Lu, Shengfeng Pan, Wen Bo, and Yunfeng Liu. 2024. Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing*, 568:127063.
- Anh Tong, Thanh Nguyen-Tang, Dongeun Lee, Duc Nguyen, Toan Tran, David Leo Wright Hall, Cheongwoong Kang, and Jassik Choi. 2025. Neural ode transformers: Analyzing internal dynamics and adaptive fine-tuning. In *ICT.R.2025 Poster*. Unpublished.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, and 1 others. 2023a. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.
- Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, and 1 others. 2023b. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*.
- A Vaswani. 2017. Attention is all you need. *Advances in Neural Information Processing Systems*.
- Qiang Wang, Bei Li, Tong Xiao, Jingbo Zhu, Changliang Li, Derek F Wong, and Lidia S Chao. 2019. Learning deep transformer models for machine translation. *arXiv preprint arXiv:1906.01787*.

- Ee Weinan. 2017. A proposal on machine learning via dynamical systems. *Communications in Mathematics and Statistics*, 1(5):1–11.
- Johannes Welbl, Nelson F Liu, and Matt Gardner. 2017. Crowdsourcing multiple choice science questions. *arXiv preprint arXiv:1707.06209*.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? *arXiv preprint arXiv:1905.07830*.
- Biao Zhang and Rico Sennrich. 2019. Root mean square layer normalization. *Advances in Neural Information Processing Systems*, 32.
- Jing Zhang, Peng Zhang, Baiwen Kong, Junqiu Wei, and Xin Jiang. 2021. Continuous self-attention models with neural ode networks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pages 14393–14401.
- Xingcheng Zhang, Zhizhong Li, Chen Change Loy, and Dahua Lin. 2017. Polynet: A pursuit of structural diversity in very deep networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 718–726.
- Wenliang Zhao, Lujia Bai, Yongming Rao, Jie Zhou, and Jiwen Lu. 2024. Unipc: A unified predictor-corrector framework for fast sampling of diffusion models. *Advances in Neural Information Processing Systems*, 36.
- Kaiwen Zheng, Cheng Lu, Jianfei Chen, and Jun Zhu. 2024. Dpm-solver-v3: Improved diffusion ode solver with empirical model statistics. *Advances in Neural Information Processing Systems*, 36.

A Training Settings

A.1 Pre-training Phase

To evaluate IIET’s performance across different model sizes, we train models from scratch at three parameter scales: 340M, 740M, and 1.3B. For all training runs, we employ the AdamW optimizer with a maximum learning rate of $3e-4$. A batch size of 0.5M tokens is used for the 340M model, while 1M tokens are used for the 740M and 1.3B models. We apply a cosine learning rate schedule to all model scales, which includes a 0.01 warmup ratio, 0.01 weight decay, and gradient clipping at 1.0. Furthermore, a smaller 55M parameter IIET variant is trained to determine the optimal iteration count, r . The complete hyperparameter details for the pre-training phase are provided in Table 4.

A.2 Iteration Influence-Aware Distillation Phase

To train efficient IIET variants, we sample one-third of the total pre-training tokens for each configuration (e.g., 5 billion tokens for 340M models, 10 billion tokens for 740M models and 30 billion tokens for 1.3B models). Users can customize the corrective iteration process for these variants based on their computational budget. In this study, we focus on two main types of efficient IIETs: a ‘lower bound’ configuration that removes all iterative steps, and E-IIET, which utilizes a threshold for iteration selection. For training, all efficient IIET variants use the full IIET as a teacher model and are trained with the fine-grained supervision method detailed in Section 4.1. We apply a cosine decay learning rate schedule with an initial value of $2e-4$, while other pre-training hyperparameters are kept consistent. Furthermore, for comparison purposes, we train Distil PCformer, a self-distilled version of PCformer using the same methodology. To ensure a fair comparison, we use the same evaluation dataset and metrics as described in Section 3.3.

A.3 Two-Stage Training

We conduct the two-stage training experiment focusing on a model with a 740 million parameter scale. The specific model configurations are provided in Table 4. In the initial stage, we train a standard Transformer architecture with a linear learning rate warmup followed by a constant rate. In the second stage, we adapt the model to the IIET structure with $r = 3$. The detailed hyperparameters for both training phases are summarized in Table 5.

Hyperparameters	55M	340M	740M	1.3B
model_type	llama	llama	llama	llama
hidden_act	silu	silu	silu	silu
initializer_range	0.02	0.02	0.02	0.02
hidden_size	512	1024	1536	2048
intermediate_size	1408	2816	4224	5504
max_position_embeddings	2048	2048	2048	2048
num_attention_heads	4	8	8	16
num_hidden_layers	12	24	24	24
num_key_value_heads	4	8	8	16
pretraining_tp	1	1	1	1
rms_norm_eps	1.00×10^{-6}	1.00×10^{-6}	1.00×10^{-6}	1.00×10^{-6}
tie_word_embeddings	True	True	True	False
torch_dtype	float16	float16	float16	float16
vocab_size	32000	32000	32000	32000
training_len	2048	2048	2048	2048
total_batch_size	128	256	512	512
learning_rate	0.0004	0.0003	0.0003	0.0002
max_steps	5000	30000	30000	100000
warm_up	0.05	0.05	0.01	0.01

Table 4: The key hyperparameters for both the model architecture and the training process.

Hyperparameters	Value
<i>Foundational Pretraining</i>	
batch_size	512
learning_rate	3e-4
lr_scheduler_type	constant_with_warmup
warmup_steps	300
training_len	2048
total_steps	22,000
<i>Architectural Transition</i>	
batch_size	512
learning_rate	3e-4
lr_scheduler_type	cosine
warmup_steps	0
training_len	2048
total_steps	8,000

Table 5: Hyperparameter configuration for two-stage training.

B Experimental Results on Other Tasks

To validate the generalizability of our proposed IJET architecture, we present experimental results on machine translation and summarization. For a fair and direct comparison, our experimental setup strictly follows the one established in PCformer (Li et al., 2024). Specifically, we train IJET in a standard big configuration (6-layer encoder and 6-layer decoder), with three iterations of our iterative implicit Euler method applied to the encoder, which is consistent with the methodology in our main exper-

Model	En-De BLEU	En-Fr BLEU	OPUS SacreBLEU	Summarization Rouge-1/2/L
Transformer++	29.2	42.9	30.8 (34.0/27.6)	40.5/17.7/37.3
PCformer	30.9	43.9	32.6 (36.0/29.1)	42.0/19.0/38.7
IJET	30.7	43.8	32.1 (35.7/28.6)	42.0/19.0/38.7

Table 6: Performance comparison of IJET against baseline models on machine translation (WMT14) and summarization tasks.

iments. Results for all baseline models, including PCformer, are taken directly from the original publication. For the OPUS benchmark, we report the average score of the X-En and En-X translation directions.

The results in Table 6 confirm that IJET achieves performance competitive with the state-of-the-art PCformer model while substantially outperforming the strong Transformer++ baseline. However, we posit that IJET’s primary practical advantage lies in its superior adaptability for low-FLOPs deployment in resource-constrained environments. This adaptability contrasts with models like PCformer, where the inherent predictor-corrector discrepancy can hinder effective compression. Furthermore, our training logs reveal an increase in validation perplexity during the final training stages, which we attribute to overfitting. This observation suggests that even stronger performance may be achievable

Model	Iter.	LMB.	PiQA	Hella.	SCIQ	ARC-c	Wino.	Avg.
Transformer++	-	28.9	64.3	34.3	76.0	23.6	51.9	46.5
PCformer	-	33.1	64.9	36.3	77.5	24.7	53.3	48.3
IJET	0	32.4	65.1	34.8	78.3	23.5	50.4	47.4
IJET	1	34.4	64.7	36.1	76.3	23.3	50.1	47.5
IJET	2	34.6	65.0	36.8	77.2	24.2	51.9	48.3
IJET	3	37.1	65.2	36.9	79.4	23.9	51.0	48.9
IJET	4	36.8	64.3	37.3	78.1	22.8	53.8	48.8
IJET	5	36.3	64.5	37.3	78.6	23.4	53.2	48.9
IJET	6	35.8	64.7	37.4	79.0	24.1	52.9	49.0
IJET	7	35.5	65.2	37.0	79.2	23.3	53.0	48.9
IJET	8	35.2	65.6	36.5	79.4	22.6	51.1	48.4

Table 7: Performance comparison of IJET with varying iteration steps at 340 million parameters.

through simple hyperparameter tuning, such as increasing the dropout rate. We leave this exploration for future work.

C IJET with Varying Iteration Steps

We evaluated the downstream task performance of our 340M model across iteration steps $r = 0$ to $r = 8$, as detailed in Section 3.3. Table 7 shows that as the number of iterations increases, IJET’s performance on downstream tasks initially improves progressively before these gains begin to plateau. Although performance slightly degrades at $r = 8$, IJET still surpasses both Transformer++ and PCformer. Notably, with $r = 2$ iterations, IJET achieves performance comparable to PCformer with its per-block forward pass count is also similar to PCformer’s. This demonstrates that our proposed iterative implicit Euler (IJET) architecture, despite its simpler design, offers representation refinement capabilities that are close to those of higher-order methods. Finally, using identical training data, IJET exhibited superior data-fitting ability over the other models, as indicated by its perplexity scores.

D Comparison with Compressed Vanilla Models

To ensure a fair comparison based on model size, we compressed the 1.3B Transformer++ model (trained on 16B tokens) using pruning and quantization. These compressed baselines are further compared with the 340M IJET to demonstrate the structural advantages of IJET.

Quantization We employ the AutoGPTQ library³ to perform 4-bit quantization on the 1.3B Transformer++. This process produces Trans quant, a model with a parameter storage footprint comparable to the 340M IJET model in bfloat16 format.

³<https://github.com/AutoGPTQ/AutoGPTQ>

Model	Para.	LMB.	PiQA	Hella.	SCIQ	ARC-c	Wino.	Avg.
IJET	340M	37.1	65.2	36.9	79.4	23.9	51.0	48.9
Transformer++	1.3B	37.3	65.7	37.6	78.6	23.7	51.5	49.0
Quantization (4-bit)	1.3B	35.9	65.3	36.5	77.3	22.5	49.3	47.8
Prun3	1.1B	23.3	62.5	34.2	72.8	23.0	50.0	44.3
Prun6	975M	11.5	59.4	32.5	67.7	22.7	50.4	40.7
Prun12	650M	2.9	53.3	28.9	52.5	22.1	50.9	35.1

Table 8: Performance comparison between IJET and compressed variants of the vanilla Transformer. The Transformer variants were created using pruning and quantization to match the storage footprint of IJET.

The quantization is calibrated on a set of 512 samples from the C4 dataset. We apply the quantization using a group size of 128, a damping percentage of 0.01, and the standard activation order. Table 8 shows that Transformer++ (4-bit) experiences a slight performance degradation compared to 1.3B Transformer++ and also underperforms the 340M IJET model (47.8 vs 48.9). The results suggest that IJET offers structural advantages beyond what post-hoc compression can achieve. Moreover, we do not apply quantization to the activations, which is in line with common practice. As a result, while 4-bit quantization reduces the parameter storage footprint, it does not lead to a significant reduction in FLOPs, since activations are still processed at 16-bit precision.

Pruning We use the method from Short-GPT (Men et al., 2024) to perform layer-wise pruning experiments on the vanilla Transformer. We calculate the block influence for each layer and create three pruned models (Prun3, Prun6, and Prun12) by removing the 3, 6, and 12 layers with the lowest influence scores, respectively. However, we find that pruning has a significant impact on the performance of Transformer++. As shown in Table 8, compared to the original model, the accuracy of Prun3 decreases by 4.6 points, while Prun6 shows a more substantial drop of 8.2 points. When half of the layers are pruned (Prun12), the model’s performance on all test sets was nearly random, so we did not further increase the pruning ratio. Notably, all pruned models performed worse than the 340M IJET model.

E Iteration Influence of 740M IJET

Figure 6 displays the Iteration Influence of the 740M IJET model. By selecting the minimum initial computation of each layer as the threshold, we can reduce the number of corrective iterations from 72 to 23.

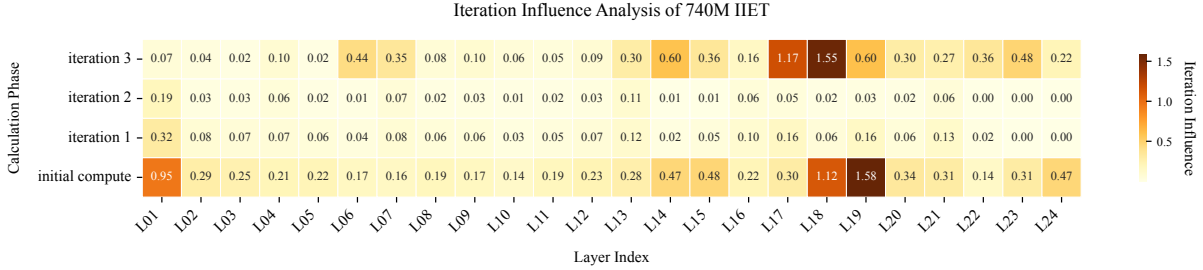


Figure 6: Impact of different iteration stages on the hidden state within each layer of the 740M IJET model, which we term **iteration influence**. Deeper colors indicate larger hidden state changes after this iteration.

F Two-Stage Training for IJET

Motivated by the significant computational expense of pre-training large language models from scratch, we investigated a two-stage training paradigm for transitioning a vanilla Transformer to an IJET. Our methodology draws inspiration from the learning rate strategies of recent large models like MiniCPM (Hu et al., 2024) and Deepseek-v3 (Liu et al., 2024), which characteristically maintain a constant learning rate after an initial warmup period during pre-training. A key advantage of this stable training phase is that it facilitates dynamic adjustments to the data curriculum. In this work, we extend this principle of in-training adaptation from the data to the model itself by transitioning the model architecture from a vanilla Transformer to an IJET during the stable learning rate phase.

Our training methodology is partitioned into two distinct phases: **① Foundational Pretraining ($\approx 75\%$ of tokens)**. We first pre-train a standard Transformer++ using a constant learning rate of $3e-4$, preceded by a 300-step linear warmup. This stage efficiently establishes a robust feature foundation while avoiding the IJET’s computational overhead. **② Architectural Transition ($\approx 25\%$ of tokens)**. The model’s architecture is then transitioned to an IJET for the remaining training. The learning rate decays from $3e-4$ via a cosine schedule, leveraging the IJET’s dynamics for final model refinement. The experiments are conducted at the 740M parameter scale, see Appendix A.3 for detailed experimental settings.

As presented in Table 9, our two-stage IJET (denoted “Two-stage”) demonstrates impressive performance. A significant performance gain is observed during the second training phase, ultimately achieving results comparable to an IJET trained from scratch (51.5 vs. 51.9). Furthermore, this approach significantly outperforms the vanilla Transformer (Trans++) trained on the same data. By

Model	LMB.	PiQA	Hella.	SCIQ	ARC-c	Wino.	Avg.
Transformer++	36.1	66.4	38.4	78.6	24.5	50.2	49.0
PCformer	41.0	66.3	41.3	82.0	23.3	51.2	50.9
IJET	41.2	68.9	42.5	82.1	23.8	53.1	51.9
Two-stage	39.3	67.5	41.3	81.9	25.3	53.7	51.5

Table 9: Performance comparison of the two-stage IJET and baseline models at the 740M parameter scale.

introducing IJET only during the final quarter of the training process, the associated computational overhead is reduced by 75%. This finding offers a promising path for applying IJET to models at much larger scales, and we anticipate that future work will continue to explore this direction.