Value Residual Learning

Anonymous ACL submission

Abstract

While Transformer models have achieved remarkable success in various domains, the effectiveness of information propagation through deep networks remains a critical challenge. Standard hidden state residuals often fail to adequately preserve initial token-level information in deeper layers. This paper introduces ResFormer, a novel architecture that enhances information flow by incorporating value residual connections in addition to hidden state residuals. And a variant is the SVFormer, where all layers share the first layer's value embedding. Comprehensive empirical evidence demonstrates ResFormer achieves equivalent validation loss with 13.3% fewer model parameters and 15.4% less training data compared to Transformer, while maintaining similar memory usage and computational cost. Besides, SV-Former reduces KV cache size by nearly half with only a small performance penalty and can be integrated with other KV-efficient methods, yielding further reductions in KV cache, with performance influenced by sequence length and cumulative learning rate.

1 Introduction

The Transformer (Vaswani et al., 2017) model has become one of the leading architectures in recent years, excelling in both language modeling (Devlin et al., 2019; Lan et al., 2020; Brown et al., 2020) and computer vision tasks (Dosovitskiy et al., 2021). Among its variants, decoder-only architectures have become the most prominent (Kaplan et al., 2020; Dubey et al., 2024). The discovery of scaling laws (Hoffmann et al., 2022; Kaplan et al., 2020) has driven the pursuit of larger Transformer models by increasing network depth and width.

In a standard decoder-only transformer, initial token embeddings contain localized information, which rapidly evolves into abstract semantic features through early attention layers (Sun et al., 2024b; Clark et al., 2019). As Transformer deepens, a critical question arises: How effectively is the initial information propagated to deeper layers? One common answer is that residual connections of hidden states ensure access to initial information throughout the network. However, some studies (Zhou et al., 2021; Shi et al., 2022) have identified that the smoothing effect of attention mechanisms leads to over-smoothing, where token representations become increasingly similar as the network deepens. This indicates that in deeper layers, sequence-level features become dominant, while token-level features are diluted. DenseFormer (Pagliardini et al., 2024) applied the idea of learnable dense connections from DenseNet (Huang et al., 2016) to transformers, and the learned connection coefficients shows that deeper layers indeed require larger attention to initial embeddings. Given the low similarity between initial token embeddings and deeper hidden states (Sun et al., 2024b), their directly summation may significantly impact the modeling of attention distribution for abstract semantic information in later layers. NeuTRENO (Nguyen et al., 2023) alleviates over-smoothing from the view of regularizers by considering the difference between value vectors of the first and current layers.

In this paper, we propose ResFormer, which enhances the propagation of initial local information by introducing value residual connections in addition to the standard hidden residual connections. Specifically, ResFormer applies a residual connection between the value vectors of the current layer and the first layer before the attention operation. In other words, both value states share the existing attention matrix of the current layer. The value states of the first attention layer and the preceding hidden states differ only by a linear transformation along the channel dimension, both representing token-level raw information. We hypothesize that introducing residual connections for values has



Figure 1: Simplified illustration of the vanilla Transformer, NeuTRENO, DenseFormer, ResFormer, and SVFormer, with only three-layer structures and no operations other than attention. \mathbf{A}^i , \mathbf{V}^i , and \mathbf{H}^i denote the attention matrix, value vectors, and attention outputs at the *i*-th layer, respectively. \oplus , \ominus , and \otimes represent standard matrix addition, subtraction, and multiplication, respectively.

a less impact on modeling attention distributions for sequence-level semantic information in higher layers and complements the original hidden state residual. Fig. 1 illustrates a comparison of the extra skip connections introduced by different models.

During inference, deep networks require substantial KV cache, severely impacting model deployment (Xiao et al., 2024). Existing KV-efficient methods often process keys and values simultaneously. Building on ResFormer, we decouple the value from the attention operation and propose a new kind of Transformer (SVFormer) where all layers share a single value states.

We experiment on a 20B SlimPajama subsampled dataset, using settings similar to popular large language models (Wei et al., 2023; Dubey et al., 2024; Kaplan et al., 2020). We compare different models based on the valid loss against the vanilla Transformer. Results show that ResFormer outperforms the vanilla Transformer, DenseFormer, and NeuTRENO. ResFormer achieves equivalent validation loss with 13.3% fewer model parameters and 15.4% less training data compared to Transformer, while maintaining similar memory usage and computational cost. Besides, SVFormer, while reducing the KV-cache by nearly half, requires a 12.2% increase in parameters to achieve the same validation loss as Transformer. And SVFormer performs better when the training sequence length is longer. It further reduces the KV cache when integrated with GQA (Ainslie et al., 2023).

2 Related Work

2.1 Shortcut Connections

Deep learning models often consist of multiple layers, posing a challenge to minimize information loss during transmission. ResNet (He et al., 2016) mitigates the vanishing gradient problem with identity connections. Stochastic Depth (Huang et al., 2016) enhances training by randomly dropping layers. DenseNet (Huang et al., 2017) allows subsequent layers to directly access the hidden states of all preceding layers. These two methods further enhance the information flow after ResNet.

Related research indicates that, although increasing depth continues to yield performance improvements in language modeling tasks, the gains become less significant with further increases (Petty et al., 2024). Furthermore, (Zhou et al., 2021) illustrates that a 32-layer ViT underperforms a 24-layer ViT. DenseFormer (Pagliardini et al., 2024) integrates weighted fusion of outputs from all preceding layers after each layer. To explore why increasing depth in Transformers does not yield expected gains, (Wang et al., 2022) finds that self-attention acts as a low-pass filter, smoothing token representations in ViTs. Additionally, (Shi et al., 2022) investigates over-smoothing from a graph perspective in BERT-based language modeling tasks. Neu-TRENO (Nguyen et al., 2023) adds the difference between the value vectors of the first and current layers to each layer's attention output and signifi144

cantly alleviates the over-smoothing problem.

2.2 KV cache compressing

The KV cache significantly impacts the efficiency of long-text model inference, attracting extensive research. One category of Transformer-based methods addresses this by employing parameter or activation value sharing techniques. The most representative works include Multi-Query Attention (Shazeer, 2019) and Grouped-Query Attention (Ainslie et al., 2023) which suggest to share key and value across a group of queries. Besides, CLA (Brandon et al., 2024) and LISA (Mu et al., 2024) respectively point out that we can reuse keys, values, or the attention matrix across layers to reduce redundancy between layers. While these methods typically process both key and value simultaneously, SVFormer is the first approach to decouple value from query and key during attention.

3 Method

3.1 Notations

Let $\mathbf{H}_n \in \mathbb{R}^{l \times d}$ be the input hidden state of the *n*-th layer, where *l* denotes the sequence length and *d* is the dimension size. For each layer, the hidden state \mathbf{H}_{n-1} will be firstly projected into $\mathbf{Q}_n, \mathbf{K}_n, \mathbf{V}_n \in \mathbb{R}^{l \times d}$ through three linear projections $\mathbf{W}_n^{\mathbf{Q}}, \mathbf{W}_n^{\mathbf{K}}, \mathbf{W}_n^{\mathbf{V}} \in \mathbb{R}^{d \times d}$ respectively. After these projections, the attention operation (Attn), output projection ($\mathbf{W}_n^{\mathbf{O}} \in \mathbb{R}^{d \times d}$), and Multi-Layer-Perceptron (Mlp) are applied sequentially:

$$\mathbf{U}_n = \operatorname{Attn}(\mathbf{Q}_n, \mathbf{K}_n, \mathbf{V}_n).$$
(1)

$$\mathbf{H}_n = \mathrm{Mlp}(\mathbf{U}_n \mathbf{W}_n^{\mathbf{O}}). \tag{2}$$

3.2 NeuTRENO and DenseFormer

After Eqn. 1, NeuTRENO adds the difference between the first and current layer's value:

$$\mathbf{U}_n = \operatorname{Attn}(\mathbf{Q}_n, \mathbf{K}_n, \mathbf{V}_n) + \boldsymbol{\lambda}_n(\mathbf{V}_1 - \mathbf{V}_n).$$
 (3)

After Eqn. 2, DenseFormer performs a weighted average between all previous hidden states:

$$\mathbf{H}_{n} = \boldsymbol{\lambda}_{n,n} \operatorname{Mlp}(\mathbf{U}_{n} \mathbf{W}_{n}^{\mathbf{O}}) + \sum_{i=1}^{n-1} \boldsymbol{\lambda}_{n,i} \mathbf{H}_{i}.$$
 (4)

where $\mathbf{H}_0 = \text{Embedding}(\mathbf{X})$ for the input \mathbf{X} . λ_n in Eqn. 3 and $\{\lambda_n\}$ in Eqn. 4 are new parameters.

3.3 ResFormer

In contrast, before Eqn. 1, ResFormer introduces a skip connection from the first layer's value $\mathbf{V}_0 = \mathbf{H}_0 \mathbf{W}_1^{\mathbf{V}}$ to current layer's value $\mathbf{V}_n = \mathbf{H}_n \mathbf{W}^{\mathbf{V}}$:

$$\mathbf{V}_n = \boldsymbol{\lambda}_{n,1} \mathbf{V}_1 + \boldsymbol{\lambda}_{n,2} \mathbf{H}_n \mathbf{W}^{\mathbf{V}}.$$
 (5)

where $\lambda_{n,1}$ and $\lambda_{n,2}$ are flexible scalars.

When all λ_1 and λ_2 are predetermined constants, it is termed **Constant-ResFormer**. If $\{\lambda_{n,1} = \lambda_{n,2}\}_{n=1}^N$, where *N* is the total number of layers, the model is called **Identity-ResFormer**. Another variant where some layers have $\lambda_1 = 0$ is referred to as **Sparse-ResFormer**. Besides, if λ_1 and λ_2 are trainable parameters, the model is termed **Learnable-ResFormer**. Unless otherwise specified, λ_1 and λ_2 are initialized to 0.5 for Learnable-ResFormer and are predetermined as 0.5 for Identity-ResFormer.

A more general form is the **Dense-ResFormer**, defined as $\mathbf{V}_n = \boldsymbol{\lambda}_{n,n} \mathbf{H}_n \mathbf{W}^{\mathbf{V}} + \sum_{i=1}^{n-1} \boldsymbol{\lambda}_{n,i} \mathbf{V}_i$ for $n \ge 2$, where $\{\boldsymbol{\lambda}_n\}$ are constants or trainable scalars. Unless noted, all $\boldsymbol{\lambda}$ are initialized to 1 here.

3.4 SVFormer

Shared Parts	-	values	keys	keys & values
Valid Loss	2.739	2.743	2.753	2.776

Table 1: Results of sharing different parts every 2 layers.

Beyond ResFormer, SVFormer adopts standard attention in the first layer and obtain the attention output $\mathbf{U_n}$ for *n*-th layer where $n \ge 2$ through $\mathbf{U_n} = \mathbf{A_nV_1}$. Its main advantage is that it only requires computing and storing the value vectors for the first layer, saving nearly half of the KVcache during inference. Similar methods like CLA reduce KV cache by sharing both of the key and value vectors every two layers. However, the results in Table 1 show that sharing values has less negative impact compared with sharing keys.

4 **Experiments**

4.1 Setting

Training Details Following (Brandon et al., 2024), we choose the Llama-like architecture and SlimPajama (Soboleva et al., 2023) data for main experiments. Specifically, the architecture includes pre-normalization, SwiGLU activations (Shazeer, 2020), rotary position embedding (Su et al., 2024), and no dropout. For slimpajama, we randomly sample nearly 20B tokens based to the original data

distribution of seven domains during training and adopt tokenizer used for "RedPajama-INCITE-7B-Base". See Table 8 in Appendix for data details.

Unless otherwise noted, we train all models using AdamW optimizer with 0.1 weight decay, $\beta_1 = 0.9, \beta_2 = 0.95$ and the max grad norm 1.0. The batch size is set to be around 2M tokens (Zhang et al., 2024) with a sequence length of 2,048 and the total steps is fixed 10,000 steps (Kaplan et al., 2020). We adopt linear learning rate warmup for the first 1,200 steps with the initial learning rate and the peak learning rate to be 1e-7 and 6e-4 respectively. The cosine decay schedule gradually decays to 10% of the peak learning rate by the end of training (Zhou et al., 2024; Wei et al., 2023). The detailed hyperparameters for models of various sizes and different training sequence lengths in Appendix A.3. Moreover, All models are trained with 8 Nvidia A100 80G GPUs using mixed-precision training in FP16. We adopt deepspeed zero-2 optimizer and flash attention mechanism.

Evaluation For model evaluation and comparison, we primarily utilized the average validation loss across seven domains, computed on the entire SlimPajama validation split. Additionally, we randomly selected a fixed set of 1,000 sample sequences for subsequent visualization analysis.

4.2 Scaling Analysis of ResFormer



Figure 2: (Left) Validation loss as model size scales from 82M to 468M parameters. (Right) Validation loss for the 468M parameter model evaluated every 2B tokens. ResFormer achieves approximately 13.3%-15.4% reduction in both model parameters and training data.

We analyzed how ResFormer and Transformer scale at model size and data size. ResFormer and Transformer are trained on similar experiment setting. On the one hand, we trained model with 82M, 180M, 320M and 468M parameters on 20B training tokens and evaluated them on a separate validation set. As shown in Fig.2 (Left), ResFormer achieves equivalent validation loss to the Transformer while utilizing 13.3% fewer model parameters. On the other hand, we evaluated the 468M models every 2B tokens and ResFormer needs 15.4% fewer training tokens to achieve the same loss as Transformer.

4.3 ResFormer vs. NeuTRENO, DenseFormer

Model	Initial Form	Loss
Vanilla Transformer	-	2.739
DenseFormer*	$\mathbf{H}_n' = 1 imes \mathbf{H}_i + \sum_{i=1}^{n-1} 0 imes \mathbf{H}_i$	2.722
NeuTRENO	$\mathbf{U}_n' = 0.4(\mathbf{V}_1 - \mathbf{V}_n) + \mathbf{U}_n$	2.72
Identity-ResFormer	$\mathbf{V}_n'=\mathbf{0.5V}_1+\mathbf{0.5V}_n$	2.712
Dense-ResFormer*	$\mathbf{V}'_n = \sum_{i=1}^n 1 imes \mathbf{V}_i$	2.709
Learnable-ResFormer*	$\mathbf{V}_n'=\mathbf{0.5V}_1+\mathbf{0.5V}_n$	2.705
Constant-ResFormer	$\mathbf{V}_n'=\mathbf{2V}_1+\mathbf{0.5V}_n$	2.7
Sparse-ResFormer	$\begin{cases} \mathbf{V}_n' = \mathbf{V}_n, 1 \le n \le 5 \\ \mathbf{V}_n' = \mathbf{V}_1, 6 \le n \le 8 \end{cases}$	2.696
Sparse-ResFormer	$\begin{cases} \mathbf{V}'_n = \mathbf{V}_n, 1 \le n \le 5\\ \mathbf{V}'_n = 5\mathbf{V}_1 + 0.5\mathbf{V}_n,\\ 6 \le n \le 8 \end{cases}$	2.687

Table 2: Average valid loss for 8-layer, 82M-parameter models. "Initial form" shows deviations from vanilla transformer. Red numbers are the λ values from Eqn. 3, Eqn. 4, and Eqn. 5. For models marked with "*", λ is learnable, and the red numbers indicate the initial value; otherwise, red numbers are fixed constants.

We compared the average validation loss of different models with 82M parameters trained on 20B tokens. In our experimental setup, NeuTRENO performed best with $\lambda = 0.4$, achieving comparable results to DenseFormer. All ResFormer variants, including the simplest Identity-ResFormer, demonstrated significant performance improvements. The λ values for Constant-ResFormer and Sparse-ResFormer were optimized through multiple experiments. While Sparse-ResFormer achieved the best overall performance, it is challenging to pre-determine the optimal layers for V_1 connections and their corresponding λ_1 values in more general scenarios. Although, Learnable-ResFormers are struggle to identify the optimal value residual pattern, they outperform Identity-ResFormer. Interestingly, the second to last row shows that Sparse-ResFormer achieved better performance despite having three fewer $\mathbf{W}^{\mathbf{V}}$.

Both Constant-ResFormer and NeuTRENO rely on predetermined λ constants. Fig. 3 shows the performance curves of these models against varying λ . Results indicate that Constant-ResFormer significantly outperforms NeuTRENO and demonstrates greater robustness across a wider range of λ values, achieving optimal performance at $\lambda = 2$.



Figure 3: The impact of varying λ values on 82M 8layer Constant-ResFormer and NeuTRENO.

4.4 Truly Better or Just Faster?



Figure 4: (Left) Average gradient norms of model outputs with respect to parameter matrices across different layers in Transformer and ResFormer. (Right) Comparison of Transformer and ResFormer performance across various learning rates during training.

To verify that ResFormer's performance improvements are not solely due to accelerated training from its shortcuts, we examined model performance across different learning rates. We compared Identity ResFormer and Vanilla Transformer under five learning rate settings. As shown in Fig. 4 (Right), both models achieved optimal results around a learning rate of 0.003, with Identity ResFormer significantly outperforming Vanilla Transformer across all rates.

Analysis of the grad norm for the four parameter matrices ($\mathbf{W}^{\mathbf{Q}}$, $\mathbf{W}^{\mathbf{K}}$, $\mathbf{W}^{\mathbf{V}}$, $\mathbf{W}^{\mathbf{O}}$) in each layer's attention module revealed that Identity Res-Former's output had approximately twice the grad norm for $\mathbf{W}_{1}^{\mathbf{V}}$ and half for $\mathbf{W}_{1}^{\mathbf{O}}$ in the first layer compared to Vanilla Transformer. This indicates that a portion of the gradient originally propagated to \mathbf{V}_{1} through \mathbf{H}_{1} is now transmitted via the value residual directly for Identity ResFormer.

In this way, we conducted the other two ablation experiments on Vanilla Transformer: doubling the learning rate for only the first layer, and doubling it exclusively for $\mathbf{W}_1^{\mathbf{V}}$ in the first layer. Neither modification yielded significant improvements. This further demonstrates that the performance improvements brought by ResFormer are unrelated to the changes in gradient magnitude.

4.5 Ablation Study of Value Residual



Figure 5: (Left) Impact of value skip connections source from different layers on model performance, where all connections are identity connections and $\lambda = 1$ in Dense-ResFormer. (Right) Average validation loss of various Sparse-ResFormer configurations, which retain only single or multiple skip connections from V₁.

Where From, Where To? We analyzed which value skip-connections are necessary for the vanilla transformer. For an 8-layer transformer, we added various pre-defined value skip-connections (with constant λ) and evaluated the resulting validation loss. As shown in Fig. 5 (Left), we first examined the impact of skip-connections from different sources. Our findings indicate that only skipconnections originating from the first layer's value (\mathbf{V}_1) yield significant performance improvements. Skip-connections from the second layer's value (V_2) offer no significant benefit to subsequent layers. Skip-connections from later layers, occurring only in the final few layers, even lead to performance degradation. Both of the two special cases in Fig. 5 (Left) include V_1 skip-connections. However, when these connections occur only between adjacent layers, the information in V_1 fails to effectively reach the final layers. Conversely, dense value skip-connections dilute the impact of V_1 with information from other sources.

Furthermore, we investigated spare ResFormer, a variant of identity ResFormer where the value residual connection $\mathbf{V}'_n = 0.5\mathbf{V}_1 + 0.5\mathbf{V}_n$ is applied selectively to specific layers. As shown in Fig. 5 (Right), for an 8-layer model, when limited to a single layer, applying the residual connection to the 7th layer yields the most significant improvement. When applied to multiple layers, the greatest benefit is observed when incorporating layers 6 to 8. Extending the residual connection to earlier layers, such as the 5th, diminishes the overall effect. It suggests that the model's final few layers benefit most from the first layer's value information.

We further trained 8-layer and 24-layer Learnable-ResFormers, as well as an 8-layer Learnable-Dense-ResFormer, and visualized the



Figure 6: (Left) Visualization of λ_1/λ_2 across different layers in the 82M and 468M Learnable-ResFormer. (Right) Heatmap visualization of learned λ across different layers in the 468M Dense-Learnable-ResFormer.

learned λ values. As shown in Fig. 6, the later layers tend to require more value residual connections from V_1 , which aligns with the findings in Fig. 5. Fortunately, the Learnable-ResFormer can, to some extent, identify similar sparse residual patterns to those of the best performing Sparse-ResFormer in Table 2. Notably, the Learnable-Dense-ResFormer learns value residual patterns that closely resemble those of the Learnable-ResFormer.

Residual Type	Initial Form	Valid Loss
value residual	$\mathbf{V}_n' = \mathbf{0.5V}_1 + \mathbf{0.5V}_n$	2.7389 2.7049
hidden residual value residual hidden residual	$\mathbf{H}'_n = \mathbf{0.5H}_0 + \mathbf{0.5H}_n$ $\mathbf{V}'_n = 0 imes \mathbf{V}_1 + 1 \mathbf{V}_n$ $\mathbf{H}'_n = 0 imes \mathbf{H}_0 + 1 \mathbf{H}_n$	2.7812 2.7298 2.7216

Table 3: Comparison of additional value residual (to V_1) and hidden residual (to H_0) connections against the default hidden residual, under various λ initializations. Trainable λ parameters are highlighted in red.

Why needed beyond hidden residual? Our experiments revealed that V_1 information provides additional benefits to later network layers, despite both H_1 and V_1 containing initial, unfused token information. H_0 is propagated through default hidden residual connections, but it may be diluted by subsequent information, hindering its effective utilization in later layers. To test this hypothesis, we introduced an additional skip connection to \mathbf{H}_0 : $\mathbf{H}'_n = \lambda_1 \mathbf{H}_1 + \lambda_2 \mathbf{H}_n$, where λ is learnable. We conducted experiments with two λ initialization settings and compared them to value residual.

Results showed that when $\lambda_1 = \lambda_2$ initially, the extra hidden residual had adverse effects. However, initializing $\lambda_1 = 0$ yielded some improvements, suggesting possible dilution of H_0 information. Nevertheless, these gains were smaller than those from value residual connections, which consistently outperformed vanilla transformers across different initializations. Actually, the connection $\mathbf{H}_0: \mathbf{H}'_n = \lambda_1 \mathbf{H}_1 + \lambda_2 \mathbf{H}_n$, is similar to applying residuals to queries, keys, and values at the same

time, may disrupt attention distributions and hinder higher-level semantic information fusing. The reduction performance brought by identity residuals of queries or keys shown in Table 5 can support it.

Hidden Residual Starts Place	Value Residual Target Place	Valid Loss
H ₀ (Default)	-	2.7389
H ₀ (Default)	V_1	2.7117
H ₀ (Default)	V_2	2.7375
H_1	-	2.7802
H_1	V_2	2.784
H_2	-	2.8196
H_2	V_2	2.7873
H_2	-	2.8196
H_2	V_3	2.8333
H_3	-	3.0568
H_3	V_3	2.8832

Table 4: Comparison of performance across different value residual target and varying hidden residual start settings. "value residual target place" V_i indicates the earliest value accessible to subsequent layers, while "hidden residual starts place" denotes the earliest hidden state available, without prior residual connections.



Figure 7: (Left) The relative training loss curve between different cross layer residual and vanilla hidden residual. (Right) Layer-to-layer hidden states similarity.

Why V_1 instead of V_2 ? In Fig. 5 (Left), connections to V_1 show significant improvement, while those to V_2 yield minimal gains. This likely occurs because the original hidden residual propagates information from H_1 to the network $(\mathbf{V}_2 = \mathbf{H}_1 \mathbf{W}_2^{\mathbf{V}})$. To verify, we adjusted residual connections, introducing them at different points. For example, starting from \mathbf{H}_1 , we use $\mathbf{H}_1 =$ $Layer_1(\mathbf{H}_0)$ instead of $\mathbf{H}_1 = Layer_1(\mathbf{H}_0) + \mathbf{H}_0$.

Table 4 results show that when residual connections begin from H_0 or H_1 , allowing H_2 and subsequent layers access to \mathbf{H}_1 , $\mathbf{V}_2 + \mathbf{V}_n$ offers no improvement. However, starting from H_2 , skip connections from V_2 provide substantial benefits. Regarding the disparity in information propagation between V_2 (via H_1) and V_1 (via H_0), we posit that after the first layer's integration, H_1 contains higher-level semantic information more similar to subsequent hidden states, see Fig. 7 (Right). This may ensure that the attention distribution remains relatively undisturbed when connecting to H_1 . Besides, Fig. 7 (Left) shows that Dense-ResFormer performs better than ResFormer when there is no cross layer hidden residual.

Superior to other residual For vanilla transformers, to better propagate information from

Residual Type	Valid Loss	Residual Mapping	Valid Loss
Query Key Attention Value	2.739 2.742 2.746 2.757 2.712	Identity Mapping Cross Layer Attention Current Attention	2.739 3.137 2.729 2.712

Table 5: The impact of various residual types, ces when adding V_1 to where all residual connec- U_n , with "Current Attions adopt a form similar tention" corresponding to

Table 6: Comparison of different mapping matrito $\mathbf{V}'_n = 0.5 \mathbf{V}_1 + 0.5 \mathbf{V}_n$. Identity-ResFormer.

the first layer, new residual connections can be introduced at various points in addition to the existing hidden residual: query states Q, key states K, value states V, and post-softmax attention matrix A. Results in Table 5 indicate that only the value residual connection improves performance. When connecting V_1 and V_n , three approaches free of extra parameters are possible: (1) the proposed residual connection, directly summing the two and then sharing an attention matrix; (2) cross layer attention (Softmax $(\mathbf{Q}_n \operatorname{Concat}(\mathbf{K}_n, \mathbf{K}_1)^T) \operatorname{Concat}(\mathbf{V}_n, \mathbf{V}_1)),$ recomputing an attention matrix for V_1 based on \mathbf{K}_1 and \mathbf{Q}_n ; and (3) directly adding \mathbf{V}_1 to \mathbf{U}_n in Eqn. 1, equivalent to using an identity mapping as V_1 's attention matrix in layer N. The second approach significantly increases computational cost. Results in Table 6 demonstrate that sharing the attention matrix yields the best performance.

Post-Analysis of ResFormer 4.6



Figure 8: (Left) Token-to-token similarity across sequence for value states at different place. (Right) Similarity between first layer values and other layers' values.

How value residual works? We performed postanalysis on trained ResFormer and vanilla Transformer models to understand value residual learning. Fig. 8 (Left) shows cosine similarities between value states at different layers and the first layer, averaged across token positions. For ResFormer, we calculated this before and after applying value residual. Results show that in vanilla Transformers, the first layer's value has low similarity with other

layers. In contrast, ResFormer maintains high similarity between the first layer's value and the postresidual values in subsequent layers due to value residual connections. Notably, in layers where Res-Former relies more heavily on the first layer's value (see Fig. 6), the pre-residual value exhibits lower similarity with the first layer's value, indicating that $\mathbf{W}^{\mathbf{V}}$ in these layers is learning the value residual.

For ResFormer, we also examined the average pairwise similarity between tokens' values before and after the residual connection. The results Fig. 8 (Right) reveal that with value residual connections, the learned values (before the value residual) from each layer become increasingly similar as the network deepens. We hypothesize that this is because, given the default hidden residual and value residual, each layer learns a ΔV , with the magnitude of necessary adjustments decreasing in later layers. This phenomenon is unique to ResFormer and not observed in vanilla Transformers.



Figure 9: (Left) The change in test loss as model modules are progressively removed, starting from the back to front while keeping the first layer intact. (Right) The number of core features in each layer's hidden state after PCA dimensionality reduction, where core features represent the minimum number of principal components required to explain 99% of the variance.

Representation and Module Analysis We analyzed the overall network changes, focusing on the hidden state representation capabilities and the contributions of different modules. (Tyukin et al., 2024) suggests that removing Attention in Transformers has a significantly smaller impact than removing Mlp. We progressively removed attention or MLP layers, starting from the last layer while retaining the first layer. Fig. 9 (Left) demonstrates that for ResFormer, the impact of removing Attention is more comparable to that of removing Mlp, in contrast to vanilla Transformers. This indicates that the Attention in ResFormer, with value residual, contribute more significantly to each layer's hidden states than in vanilla Transformers.

Furthermore, we performed PCA dimensionality reduction on the hidden states of each layer in both ResFormer and vanilla Transformer models. We determined the minimum number of principal components required to explain 99% of the variance. Fig. 9 reveals that ResFormer, starting from the second layer where value residual connections are introduced, consistently produces hidden states with a higher minimum number of principal components compared to vanilla Transformers. This suggests that ResFormer generates hidden states with higher information density.

4.7 SVFormer vs. GQA,CLA



Figure 10: The relative training loss for SVFormer and other KV efficient model compared with vanilla attention. The numbers in parentheses represent the training sequence length. Left: Model with nearly $1/2 \ KV$ cache. Right: Model with nearly $1/8 \ KV$ cache.



Figure 11: Left: The relative training loss for SVFormer under different sequence lengths with a fixed batch size of 2M tokens. Right: Analysis of critical point, and we predict it for length 64,000 using linear regression with the last 1,000 data points.

In the Fig. 10, at a training sequence length of 64,000, SVFormer demonstrates lower final loss compared to existing KV-efficient methods such as CLA and GQA. Moreover, it can be used concurrently with GQA to enhance KV efficiency further. However, we observed that with a training sequence length of 2,048, SVFormer underperforms compared to GQA. The results indicate that sequence length significantly affects SVFormer's performance. Thus, we conducted more comprehensive experiments on sequence length.

Effects of sequence length. Results in Fig. 11 (Left) demonstrate that SVFormer will always be gradually surpassed by vanilla attention during training while its training speed is faster than vanilla Transformer at the early stage. However, as the training sequence length increases, the SV-Former model performs better. In this way, we



Figure 12: The relative training loss for SVFormer under different hyper-parameter setting and the validation loss as model size scales from 82M to 468M parameters.

focus on the critical point, defined as the number of training steps exceeded. Fig. 11 (Right) illustrates that the relationship between the critical point and sequence length exhibits an exponential trend. We argue that it's due to the challenge deep models face in fully optimizing the increasingly larger first-layer value matrix as the sequence length grows.

Other Factors. Fig. 12a and Fig. 12b show SV-Former benefits more from smaller learning rates than from warmup. This aligns with performance correlating to total summed learning rate (Kaplan et al., 2020). Larger models, requiring smaller learning rates, suit SVFormer better. Fig. 12c indicates the SVFormer-Transformer difference is not architecture-sensitive. Compared with Transformer, SVFormer requires a 12.2% increase in parameters to achieve the same loss while reducing the KV-cache by nearly half.

5 Conclusion

In this paper, we demonstrate the inadequacy of existing hidden residual connections in propagating information from the initial token-level to deeper layers. To address this limitation, we propose Res-Former, which incorporates a residual connection between the value vectors of the current layer and those of the first layer prior to the attention operation. Furthermore, we introduce SVFormer, an extension of ResFormer, which achieves a nearly 50% reduction in the KV cache. We conducted extensive experiments on language modeling tasks to evaluate the efficacy of these two Transformer variants across diverse scenarios.

547 548 549

The proposed learnable ResFormer, still falls short of identifying the optimal λ setting through current training, instead converging on a relative optimum. This limitation suggests that further refinement of initialization strategies and learning algorithms may be necessary. Due to computational constraints, we were unable to conduct experimental validation on larger-scale models at this time.

Ethics Statement

Limitations

On the one hand, the data employed in this paper is sourced from publicly available datasets provided by the company, which have undergone a certain level of filtering. On the other hand, the models trained in our study are solely utilized for experimental analysis and will not be publicly deployed.

References

- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebrón, and Sumit Sanghai. 2023. Gqa: Training generalized multi-query transformer models from multi-head checkpoints. *arXiv preprint arXiv:2305.13245*.
- Stella Biderman, Hailey Schoelkopf, Quentin Gregory Anthony, Herbie Bradley, Kyle O'Brien, Eric Hallahan, Mohammad Aflah Khan, Shivanshu Purohit, USVSN Sai Prashanth, Edward Raff, et al. 2023. Pythia: A suite for analyzing large language models across training and scaling. In *International Conference on Machine Learning*, pages 2397–2430. PMLR.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. 2020. PIQA: reasoning about physical commonsense in natural language. In *AAAI*, pages 7432–7439. AAAI Press.
- William Brandon, Mayank Mishra, Aniruddha Nrusimha, Rameswar Panda, and Jonathan Ragan Kelly. 2024. Reducing transformer key-value cache size with cross-layer attention. *arXiv preprint arXiv:2405.12981.*
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners. In *NeurIPS*.

- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D. Manning. 2019. What does BERT look at? an analysis of bert's attention. In *BlackboxNLP@ACL*, pages 276–286. Association for Computational Linguistics.
- 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. *CoRR*, abs/1803.05457.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, pages 4171–4186. Association for Computational Linguistics.
- Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. 2021. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*. OpenReview.net.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The Ilama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Tianyu Guo, Druv Pai, Yu Bai, Jiantao Jiao, Michael I Jordan, and Song Mei. 2024a. Active-dormant attention heads: Mechanistically demystifying extreme-token phenomena in llms. *arXiv preprint arXiv:2410.13835*.
- Zhiyu Guo, Hidetaka Kamigaito, and Taro Watanabe. 2024b. Attention score is not all you need for token importance indicator in kv cache reduction: Value also matters. *arXiv preprint arXiv:2406.12335*.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 770– 778.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, et al. 2022. Training compute-optimal large language models. *arXiv preprint arXiv:2203.15556*.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4700–4708.
- Gao Huang, Yu Sun, Zhuang Liu, Daniel Sedra, and Kilian Q Weinberger. 2016. Deep networks with stochastic depth. In *Computer Vision–ECCV 2016:*

14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14, pages 646–661. Springer.

- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut.
 2020. ALBERT: A lite BERT for self-supervised learning of language representations. In *International Conference on Learning Representations*. OpenReview.net.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. 2018. Can a suit of armor conduct electricity? a new dataset for open book question answering. *Preprint*, arXiv:1809.02789.
- Yongyu Mu, Yuzhang Wu, Yuchun Fan, Chenglong Wang, Hengyu Li, Qiaozhi He, Murun Yang, Tong Xiao, and Jingbo Zhu. 2024. Cross-layer attention sharing for large language models. *arXiv preprint arXiv:2408.01890*.
- Tam Nguyen, Tan Nguyen, and Richard Baraniuk. 2023. Mitigating over-smoothing in transformers via regularized nonlocal functionals. *Advances in Neural Information Processing Systems*, 36:80233–80256.
- Matteo Pagliardini, Amirkeivan Mohtashami, Francois Fleuret, and Martin Jaggi. 2024. Denseformer: Enhancing information flow in transformers via depth weighted averaging. *arXiv preprint arXiv:2402.02622*.
- Jackson Petty, Sjoerd Steenkiste, Ishita Dasgupta, Fei Sha, Dan Garrette, and Tal Linzen. 2024. The impact of depth on compositional generalization in transformer language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 7232–7245.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. 2019. Winogrande: An adversarial winograd schema challenge at scale. *Preprint*, arXiv:1907.10641.
- Noam Shazeer. 2019. Fast transformer decoding: One write-head is all you need. *CoRR*, abs/1911.02150.
- Noam Shazeer. 2020. GLU variants improve transformer. *CoRR*, abs/2002.05202.
- Han Shi, Jiahui Gao, Hang Xu, Xiaodan Liang, Zhenguo Li, Lingpeng Kong, Stephen Lee, and James T Kwok.
 2022. Revisiting over-smoothing in bert from the perspective of graph. arXiv preprint arXiv:2202.08625.

- 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.
- Mingjie Sun, Xinlei Chen, J Zico Kolter, and Zhuang Liu. 2024a. Massive activations in large language models. *arXiv preprint arXiv:2402.17762*.
- Qi Sun, Marc Pickett, Aakash Kumar Nain, and Llion Jones. 2024b. Transformer layers as painters. *CoRR*, abs/2407.09298.
- Georgy Tyukin, Gbètondji J.-S. Dovonon, Jean Kaddour, and Pasquale Minervini. 2024. Attention is all you need but you don't need all of it for inference of large language models. *CoRR*, abs/2407.15516.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need.(nips), 2017. *arXiv preprint arXiv:1706.03762*, 10:S0140525X16001837.
- Peihao Wang, Wenqing Zheng, Tianlong Chen, and Zhangyang Wang. 2022. Anti-oversmoothing in deep vision transformers via the fourier domain analysis: From theory to practice. *arXiv preprint arXiv:2203.05962*.
- Tianwen Wei, Liang Zhao, Lichang Zhang, Bo Zhu, Lijie Wang, Haihua Yang, Biye Li, Cheng Cheng, Weiwei Lü, Rui Hu, et al. 2023. Skywork: A more open bilingual foundation model. *arXiv preprint arXiv:2310.19341*.
- Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024. Efficient streaming language models with attention sinks. In *International Conference on Learning Representations*. OpenReview.net.
- Sang Michael Xie, Hieu Pham, Xuanyi Dong, Nan Du, Hanxiao Liu, Yifeng Lu, Percy S Liang, Quoc V Le, Tengyu Ma, and Adams Wei Yu. 2024. Doremi: Optimizing data mixtures speeds up language model pretraining. *Advances in Neural Information Processing Systems*, 36.
- Rowan Zellers, Ari Holtzman, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. Hellaswag: Can a machine really finish your sentence? In *ACL* (1), pages 4791–4800. Association for Computational Linguistics.
- Peiyuan Zhang, Guangtao Zeng, Tianduo Wang, and Wei Lu. 2024. Tinyllama: An open-source small language model. *arXiv preprint arXiv:2401.02385*.

- Chunting Zhou, Pengfei Liu, Puxin Xu, Srinivasan Iyer, Jiao Sun, Yuning Mao, Xuezhe Ma, Avia Efrat, Ping Yu, Lili Yu, et al. 2024. Lima: Less is more for alignment. *Advances in Neural Information Processing Systems*, 36.
- Daquan Zhou, Bingyi Kang, Xiaojie Jin, Linjie Yang, Xiaochen Lian, Zihang Jiang, Qibin Hou, and Jiashi Feng. 2021. Deepvit: Towards deeper vision transformer. arXiv preprint arXiv:2103.11886.

A Appendix

A.1 Downstream Evaluations

We compare the different models on several classical reasoning tasks following (Zhang et al., 2024) in a zero-shot way. The tasks include Hellaswag (Zellers et al., 2019), OpenBookQA (Mihaylov et al., 2018), WinoGrande (Sakaguchi et al., 2019), ARC-Easy and ARC-Challenge (Clark et al., 2018) and PIQA (Bisk et al., 2020). The results in Table 7 show that ResFormer achieved an average accuracy improvement of nearly 3% compared to the vanilla Transformer.



(c) Norms of value states.

(d) Norms of hidden states.

Figure 13: The token importance (Xiao et al., 2024), value-state norms (Guo et al., 2024b), and hidden-state norms (Sun et al., 2024a) of the first token across layers of 468M models. "Attention Entropy" refers to the entropy of token importance across each sequence.

Attention Pattern Analysis Given the attention matrix $\mathbf{A} \in \mathbb{R}^{l \times l}$ at one layer, we use entropy e to represent its concentration effect. To obtain entropy E, calculate the importance vector $\mathbf{a} = \frac{1}{l} \sum_{j=1}^{l} A_{ij}$ firstly where \mathbf{A} is a lower triangular matrix. The entropy can be formulated as follows: $e = -\sum_{i=1}^{l} a'_i \log a'_i$, where $a'^i_i = a_i / (\sum_{i=1}^{l} a_i)$ for i = 1, 2, ..., l and the higher the entropy e, the greater the degree of clustering

in a, i.e., attention matrix A is more likely to focus on several specific tokens.

Fig. 13a illustrates that the clustering effect of attention increases significantly with the number of layers for the vanilla Transformer, whereas the clustering effect is relatively less pronounced for the ResFormer. We further visualize the attention weights, value-state norms $\|v\|_2$, and hidden-state norms $\|h\|_2$ of tokens at different layers and positions. Given that attention clustering often occurs on the first token, we primarily show its results in Fig. 13. The results indicate that using ResFormer significantly mitigates attention sinks (Xiao et al., 2024), value-state drains (Guo et al., 2024b) and residual-state peaks (Sun et al., 2024a). (Guo et al., 2024a) attributes these phenomena to the mutual reinforcement mechanism of model and we suggest that the value shortcut disrupts this mechanism by alleviating value-state drains. Specifically, for tokens lacking semantic information like start tokens, a large value state magnitude can adversely affect the prediction of subsequent tokens if they are overly attended to. When there is no value-state drains, models will reduce attention clustering to these tokens to minimize loss.

A.2 Pre-train Dataset

Based on the equation $D \ge 5000 \cdot N^{0.74}$ (Kaplan et al., 2020) where D is data size and N is the number of non-embedding parameters, we need to collect at least 17.5B for model has N = 700M non-embedding parameters (corresponding to complete 1B model with 2,048 hidden size, 50,277 vocab size and 2,048 sequence length) to avoid over-fitting. Besides, (Xie et al., 2024) indicates that the mixture proportions of pre-training data domains significantly affects the training results. In this way, we sampled 20B tokens data from original 627B data based on the original data proportions shown in the Table 8.

A.3 Training Details

Section 4.1 introduces the main experimental hyperparameters used in the paper. This section further details the training parameters for various model sizes and training sequence lengths. Table 10 demonstrates the differences among models of various sizes. The configurations for the number of layers, attention heads, hidden dimensions, and FFN dimensions are based on (Biderman et al., 2023). Additionally, the λ in Eqn. 3 is set to be 0.4 for NeuTRENO. Moreover, as reported in Table 9,

Model	Max Length	HellaSwag	Obqa	WinoGrande	ARC-c	ARC-e	PIQA	Avg
Transformer	2,048	0.263	0.142	0.492	0.199	0.331	0.572	0.333
ResFormer	2,048	0.273	0.148	0.512	0.182	0.414	0.604	0.355
Transformer	64,000	0.267	0.142	0.485	0.179	0.322	0.570	0.328
ResFormer	64,000	0.274	0.136	0.513	0.184	0.407	0.588	0.350

Table 7: Zero-shot accuracy on commonsense reasoning tasks.

Data source	proportions	Tokens
Commoncrawl	50%	10 B
C4	20%	4 B
GitHub	10%	2 B
Books	5%	1 B
ArXiv	5%	1 B
Wikpedia	5%	1 B
StackExchange	5%	1 B

Table 8: The details of pre-train dataset.

the batch size that a single GPU can accommodate varies depending on the length of the training sequences. Note that the total number of tokens in each batch is consistently 2 million.

839 840

Max Sequence Length	512	2,048	8,192	32,000	64,000
Total Batch Size	4,096	1,024	256	64	32
Per-GPU Batch Size	128	32	8	2	1
Gradient Accumulation Step			32		
GPUs			8		

Table 9: Training details for training dataset with different sequence length.

Model Size	2M	82M	180M	468M	
Layers	4	8	12	24	
Attention Heads	2	8	12	16	
Hidden Dimension	16	512	768	1,024	
FFN Dimension	56	1,792	2,688	3,584	
Tie Word Embedding		Fa	alse		
(Peak Learning Rate, Final Learning Rate)		(6e - 4)	(6e-5))	
Learning Rate Schedule		Cosine Decay			
Vocabulary Size	50,277				
Activation Function	vation Function SwiGLU				
Position Embedding		RoPE (θ	= 10,000	0)	
Batch Size		2M 1	tokens		
Data Size		20B	tokens		
(Warmup Steps, Training Steps)		(120, 10,000)			
Adam β	(0.9, 0.95)				
Dropout		().0		
Weight Decay		().1		

Table 10: Training details for models with different size.