

# Attention Residuals: A Memoir

Jianlin Su

2026-03-19

This article introduces our latest work, Attention Residuals (AttnRes); as the name suggests, it uses the idea of Attention to improve Residuals.

Many readers will have heard of the Pre Norm / Post Norm debate, but ultimately this is just “infighting” within Residuals itself, and many later variants of Normalization are the same. A more interesting change is HC, which began down the route of expanding the residual stream, but perhaps because of unstable results it did not attract much response. The later story is probably known to everyone: at the end of last year DeepSeek’s mHC improved HC and verified its effectiveness on a larger scale.

Rather than further expanding the residual stream, we chose another radical route: directly performing Attention between layers to replace Residuals. Of course, getting the full pipeline to work involved many details and efforts; here we simply recall the related journey.

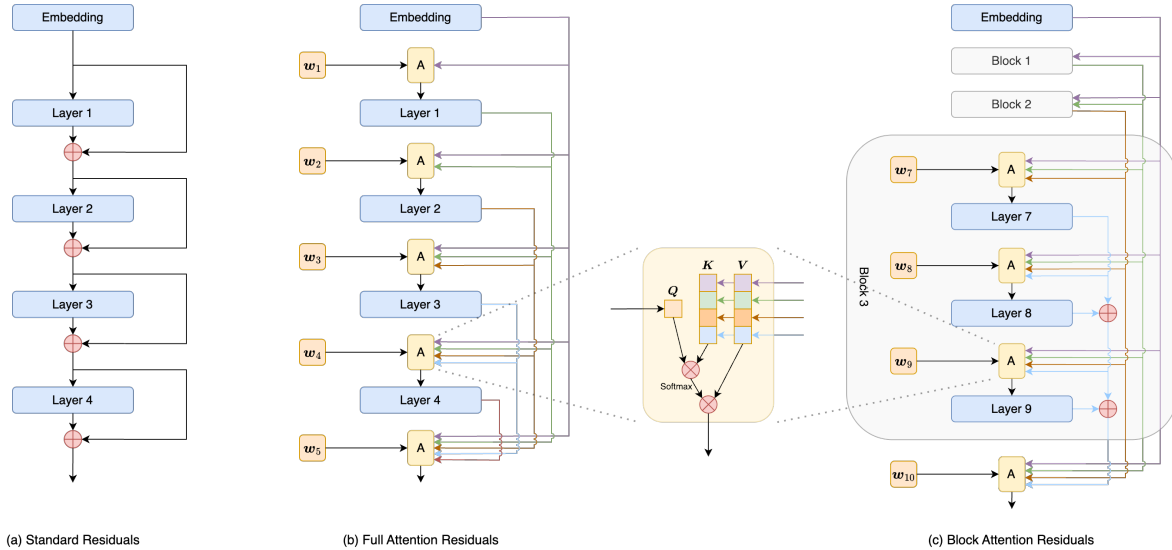


Figure 1: AttnRes schematic diagram

## 1 Inter-layer Attention

As usual, we start from Residuals, which everyone knows by heart. Its form is

$$\mathbf{x}_t = \mathbf{x}_{t-1} + \mathbf{f}_t(\mathbf{x}_{t-1}) \quad (1)$$

Here we switch to another notation that lets us see something deeper. Let  $\mathbf{y}_t = \mathbf{f}_t(\mathbf{x}_{t-1})$ , then  $\mathbf{x}_t = \mathbf{x}_{t-1} + \mathbf{y}_t$ , with the convention  $\mathbf{y}_0 = \mathbf{x}_0$ , so we easily obtain  $\mathbf{x}_t = \mathbf{y}_0 + \mathbf{y}_1 + \dots + \mathbf{y}_t$ , and it

can equivalently be written as

$$\mathbf{y}_{t+1} = \mathbf{f}_{t+1}(\mathbf{y}_0 + \mathbf{y}_1 + \cdots + \mathbf{y}_t) \quad (2)$$

That is, from the  $\mathbf{y}$  perspective, Residuals takes the equal-weight sum of  $\mathbf{y}_0, \mathbf{y}_1, \dots, \mathbf{y}_t$  as the input to  $\mathbf{f}_{t+1}$  to obtain  $\mathbf{y}_{t+1}$ . A natural generalization is to replace it with a weighted sum:

$$\mathbf{y}_{t+1} = \mathbf{f}_{t+1} \left( \sum_{s=0}^t a_{t+1,s} \mathbf{y}_s \right) \quad \text{where} \quad a_{t,s} \geq 0, \quad \sum_{s=0}^t a_{t+1,s} = 1 \quad (3)$$

This is the germ of AttnRes. The above formula adds two extra constraints on  $a_{t,s}$ ; let us discuss their necessity:

1. The constraint  $a_{t,s} \geq 0$  ensures that the same  $\mathbf{y}_s$  always contributes in the same direction to different layers, avoiding the inconsistency where one layer wants to increase  $\mathbf{y}_s$  while another wants to shrink it, which is intuitively more friendly to model learning;
2. The  $\mathbf{f}$  we use includes In Norm and will first apply RMSNorm to the input. Since  $\text{RMSNorm}(\mathbf{x}) = \text{RMSNorm}(c\mathbf{x})$  holds identically for all  $c > 0$ , weighted averaging and weighted summation are completely equivalent, so the constraint  $\sum_{s=0}^t a_{t,s} = 1$  does not reduce expressiveness.

## 2 Hyper-Connections

Before launching into AttnRes, let us briefly review HC (Hyper-Connections) and show that it too can be understood as inter-layer Attention, thereby demonstrating that inter-layer Attention is indeed a more fundamental route. HC changes Residuals to

$$\mathbf{X}_t = \mathbf{H}_t^{\text{res}} \mathbf{X}_{t-1} + \mathbf{H}_t^{\text{post}} \mathbf{f}_t(\mathbf{H}_t^{\text{pre}} \mathbf{X}_{t-1}) \quad (4)$$

where  $\mathbf{X} \in \mathbb{R}^{k \times d}$ ,  $\mathbf{H}^{\text{res}} \in \mathbb{R}^{k \times k}$ ,  $\mathbf{H}^{\text{pre}} \in \mathbb{R}^{1 \times k}$ ,  $\mathbf{H}^{\text{post}} \in \mathbb{R}^{k \times 1}$ , and the classic choice is  $k = 4$ . Simply put, the state variable is expanded  $k$  times; before input to  $\mathbf{f}_t$  a matrix  $\mathbf{H}_t^{\text{pre}}$  projects it back to one time; after output a matrix  $\mathbf{H}_t^{\text{post}}$  expands it back to  $k$  times, and finally it is added to the  $\mathbf{x}_{t-1}$  adjusted by  $\mathbf{H}_t^{\text{res}}$ . If we do not restrict the forms of  $\mathbf{H}_t^{\text{res}}$ ,  $\mathbf{H}_t^{\text{pre}}$ ,  $\mathbf{H}_t^{\text{post}}$ , then Post Norm, Highway, and others are all special cases of HC.

Similarly let  $\mathbf{y}_t = \mathbf{f}_t(\mathbf{H}_t^{\text{pre}} \mathbf{X}_{t-1})$ , then  $\mathbf{X}_t = \mathbf{H}_t^{\text{res}} \mathbf{X}_{t-1} + \mathbf{H}_t^{\text{post}} \mathbf{y}_t$ , with the convention  $\mathbf{X}_0 = \mathbf{H}_0^{\text{post}} \mathbf{y}_0$ , so it can also be expanded as  $\mathbf{X}_t = \mathbf{H}_{t \leftarrow 1}^{\text{res}} \mathbf{H}_0^{\text{post}} \mathbf{y}_0 + \mathbf{H}_{t \leftarrow 2}^{\text{res}} \mathbf{H}_1^{\text{post}} \mathbf{y}_1 + \cdots + \mathbf{H}_{t \leftarrow t}^{\text{res}} \mathbf{H}_{t-1}^{\text{post}} \mathbf{y}_{t-1} + \mathbf{H}_t^{\text{post}} \mathbf{y}_t$ , where  $\mathbf{H}_{t \leftarrow s}^{\text{res}}$  is defined as  $\mathbf{H}_t^{\text{res}} \mathbf{H}_{t-1}^{\text{res}} \cdots \mathbf{H}_{s+1}^{\text{res}} \mathbf{H}_s^{\text{res}}$ . Further adopting the convention  $\mathbf{H}_{t \leftarrow t+1}^{\text{res}} = \mathbf{I}$ , we can write

$$\mathbf{y}_{t+1} = \mathbf{f}_{t+1}(\mathbf{H}_{t+1}^{\text{pre}} \mathbf{x}_t) = \mathbf{f}_{t+1} \left( \sum_{s=0}^t \underbrace{\mathbf{H}_{t+1}^{\text{pre}} \mathbf{H}_{t \leftarrow s+1}^{\text{res}} \mathbf{H}_s^{\text{post}}}_{a_{t+1,s}} \mathbf{y}_s \right) \quad (5)$$

Note that each  $\mathbf{H}_{t+1}^{\text{pre}} \mathbf{H}_{t \leftarrow s+1}^{\text{res}} \mathbf{H}_s^{\text{post}}$  is a  $1 \times 1$  matrix, equivalent to a scalar, so it too is of the inter-layer Attention form of Eq. (3). Readers familiar with linear attention should quickly understand this result: HC is essentially a DeltaNet “rotated by 90 degrees”. In practice, the three  $\mathbf{H}$  matrices are computed from simple linear layers with tanh activation, which causes the chained product  $\mathbf{H}_{t \leftarrow s}^{\text{res}}$  to risk explosion or collapse, and also cannot guarantee the non-negativity of  $a_{t+1,s}$ .

Later mHC made improvements: first, all three  $\mathbf{H}$  were changed to Sigmoid activation to guarantee  $a_{t+1,s} \geq 0$ ; then it alternately normalized  $\mathbf{H}_t^{res}$  to make it doubly stochastic, and by the closure of doubly stochastic matrices under multiplication it guaranteed the stability of  $\mathbf{H}_{t \leftarrow s}^{res}$ ; finally experiments verified the effectiveness of these changes. However, some new experiments such as “Your deepseek mHC may not need ‘m’” showed that directly setting  $\mathbf{H}_t^{res}$  to the identity matrix is good enough.

### 3 Teamwork

Let us return to AttnRes. After realizing the feasibility of AttnRes, the next question is: what form should  $a_{t+1,s}$  take? A very natural idea is to follow the standard Scaled Dot-Product Attention, but at the time the author wanted a quick first try, so chose a simpler form

$$a_{t+1,s} \propto \exp(\mathbf{w}_{t+1} \cdot \mathbf{y}_s) \tag{6}$$

where  $\mathbf{w}_t$  is a learnable vector parameter, i.e., directly using a data-independent static vector as Q while K and V are both  $\mathbf{y}_s$  to perform Softmax Attention. This is the first version of AttnRes. Surprisingly, with such a simple design the improvement over Residuals is already very significant!

When the author shared the preliminary experimental results of AttnRes within the group, @Zhang Yu and @Guang Yu showed great interest and joined in to start validation on larger-scale models, and the results were all encouraging. During this period, we also tried some more complex designs and found that most were not as good as this simple version; only adding an extra RMSNorm operation to K could obtain relatively stable gains, which constitutes the final form of AttnRes

$$a_{t+1,s} \propto \exp(\mathbf{w}_{t+1} \cdot \text{RMSNorm}(\mathbf{y}_s)) \tag{7}$$

However, AttnRes is after all an intensive inter-layer connection scheme. Is training and inference feasible on K2 or even larger scales? Excitingly, @V via an elegant analysis first affirmed the feasibility of inference, and the “finishing touch” was precisely the convenient static-Q design! This allows us to compute the attention  $a_{t,s}$  for  $t > s$  immediately after computing  $\mathbf{y}_s$ , giving Infra enough room to maneuver.

But unfortunately, the training side, such as @Wang, after careful analysis judged that AttnRes was still not feasible in our current training environment (to put it bluntly, we are still poor), and needed a further scheme to reduce communication and memory; thus came the Block version below; correspondingly, the previous version is called the Full version.

### 4 Block Version

From Full AttnRes to Block AttnRes is analogous to the past process of linearizing quadratic Attention; various existing Efficient Attention ideas can be tried, such as SWA (Sliding Window Attention), which was our first attempt, but the actual effect was very poor, even worse than Residuals.

After reflection, the author thinks it can be understood this way: Residuals itself is already a very strong baseline, corresponding to the equal-weight sum of all state vectors. Any new design that wants to surpass it must at least formally be able to cover it. The Full AttnRes obviously satisfies this condition, but adding SWA does not, because it discards part of the state and cannot cover the special case of “equal-weight sum of all state vectors”.

Thus we realized that for AttnRes, “compression” may be more effective than “sparsity”, and the compression does not need to be too fine-grained; simple weighted summation may suffice. After





*If reproducing, please include this URL:* <https://kexue.fm/archives/11664>

*For more detailed reproduction matters please refer to:* “Science Space FAQ”