

# ZeroS: Zero-Sum Linear Attention for Efficient Transformers



Jiecheng Lu<sup>1</sup> Xu Han<sup>2</sup> Yan Sun<sup>1</sup> Viresh Pati<sup>1</sup> Yubin Kim<sup>1</sup> Siddhartha Somani<sup>1</sup> Shihao Yang<sup>1</sup>

<sup>1</sup>Georgia Institute of Technology <sup>2</sup>AWS

## Motivation & Problem

- **Quadratic vs. Linear:** Softmax attention is accurate but  $O(N^2)$ . Linear attention variants achieve  $O(N)$  yet often underperform.
- **Two fundamental limitations (affecting linear and even softmax):**
  - *Convex-combination bottleneck:* per-head reads are convex averages, permitting only additive blending; single-layer contrast/differencing is impossible.
  - *Uniform weight bias / dilution:* a roughly uniform  $1/N$  component acts like persistent average pooling in long contexts, weakening focus.
- **Idea:** Remove the zero-order term behind that bias and reweight residuals to allow signed, zero-sum weights—enabling contrastive operations in one layer while preserving  $O(N)$  complexity.

## ZeroS in One Slide

**Zero-Sum Linear Attention (ZeroS)** removes the constant zero-order term and reweights the remaining softmax residuals to obtain *signed, zero-sum* attention weights.

$$\bar{s}_t = \frac{1}{t} \sum_{j=1}^t s_{t,j}, \quad \delta_{t,i} = s_{t,i} - \bar{s}_t,$$

$$\varepsilon_{t,i} = \frac{e^{s_{t,i}}}{\sum_{r=1}^t e^{s_{t,r}}} - \frac{1}{t} - \frac{\delta_{t,i}}{t}, \quad \sum_i \varepsilon_{t,i} = 0,$$

$$w_{t,i} = \underbrace{\sigma_t^1 \frac{\delta_{t,i}}{t}}_{1\text{st order}} + \underbrace{\sigma_t^h \varepsilon_{t,i}}_{2\text{nd+ order}}, \quad \sum_{i=1}^t w_{t,i} = 0.$$

Here  $\sigma_t^1, \sigma_t^h \in (0, 1)$  are learned gates (per step/channel) controlling linear vs. nonlinear competition components.

## Expressivity & Stability (intuition)

- **Beyond convex hull:** subtracting  $1/t$  unlocks signed, zero-sum weights. A single layer can *contrast* tokens (e.g., simple differencing), which convex mixing cannot represent in one step.
- **Stable by construction:** with bounded logits and  $\|v_i\| \leq B$ , zero-sum weights satisfy  $\max_i |w_{t,i}| = O(1/t)$ ; the head output is  $O(B)$ , independent of  $t$  (length-stable).
- **Residual alignment:** zero-sum produces centered updates ( $\sum_i w_{t,i} = 0$ ), matching decoder residual stream design.

## From Averaging to Contrast

- Keeping  $1/t$  in the *first* layer preserves the affine hull; subsequent ZeroS layers (zero-sum) add contrastive directions without shrinking the reachable set.
- Practically, we default to removing  $1/t$  in all layers; optionally retain it in the first layer (*minor impact* in our experiments).

## Radial-Angular Decoupling (with RoPE)

ZeroS combines a **radial** reweighted zero-sum softmax with an **angular** cosine term:

$$\mathbf{o}_t = \sum_{i=1}^t \underbrace{r_{t,i}}_{\text{signed, zero-sum}} \cdot \underbrace{\cos \theta_{t,i}}_{\text{direction}} \mathbf{v}_i, \quad \cos \theta_{t,i} = \hat{\mathbf{q}}_i^\top \hat{\mathbf{k}}_i \text{ or } \hat{\mathbf{q}}_i^\top \mathbf{R}_{t-i} \hat{\mathbf{k}}_i \text{ (RoPE)}$$

Unlike positive-only kernels, the sign flip of  $\cos \theta$  survives into the weights via  $r_{t,i}$ , restoring contrastive, distance-aware behavior at *linear* cost.

## Linear-Time Implementation (ZeroS-Lin)

We design  $O(N)$  logits and a prefix-sum scan:

$$\mathbf{u}_i = \mathbf{x}_i \mathbf{W}_u, \quad \bar{\mathbf{u}}_i = \frac{e^{\mu + \sum_{j=1}^i \mathbf{u}_j}}{e^{\mu + i}}, \quad s_i = -\frac{1}{\sqrt{d}} \mathbf{u}_i \bar{\mathbf{u}}_i^\top$$

Maintain prefix sums (per head/channel):

$$E_t = \sum_{i \leq t} e^{s_i}, \quad P_t = \sum_{i \leq t} s_i, \quad \mathbf{F}_t = \sum_{i \leq t} e^{s_i} \hat{\mathbf{k}}_i^\top \mathbf{v}_i, \quad \mathbf{G}_t = \sum_{i \leq t} s_i \hat{\mathbf{k}}_i^\top \mathbf{v}_i, \quad \mathbf{H}_t = \sum_{i \leq t} \hat{\mathbf{k}}_i^\top \mathbf{v}_i.$$

Compute scalars and output:

$$\alpha_t = \frac{\sigma_t^h}{E_t}, \quad \beta_t = \frac{\sigma_t^1 - \sigma_t^h}{t}, \quad \gamma_t = -\frac{\sigma_t^h + \beta_t P_t}{t}, \quad \mathbf{o}_t = \hat{\mathbf{q}}_t (\alpha_t \mathbf{F}_t + \beta_t \mathbf{G}_t + \gamma_t \mathbf{H}_t).$$

**Complexity:**  $O(d^2)$  per step,  $O(Nd^2)$  total; same hidden-state size as linear attention.

## ZeroS for Softmax Attention (ZeroS-SM)

Apply the same zero-order removal and residual reweighting to the *quadratic* softmax matrix:

$$\mathbf{S} = \frac{1}{\sqrt{d}} \mathbf{Q} \mathbf{K}^\top + \mathbf{M}, \quad \mathbf{A} = \text{softmax}(\mathbf{S}), \quad \mathbf{W} = (\mathbf{g}^1 \mathbf{1}^\top) \odot \Delta + (\mathbf{g}^h \mathbf{1}^\top) \odot \varepsilon, \quad \mathbf{O} = \mathbf{W} \mathbf{V}.$$

Empirically, ZeroS-SM *exceeds* vanilla softmax at matched budgets.

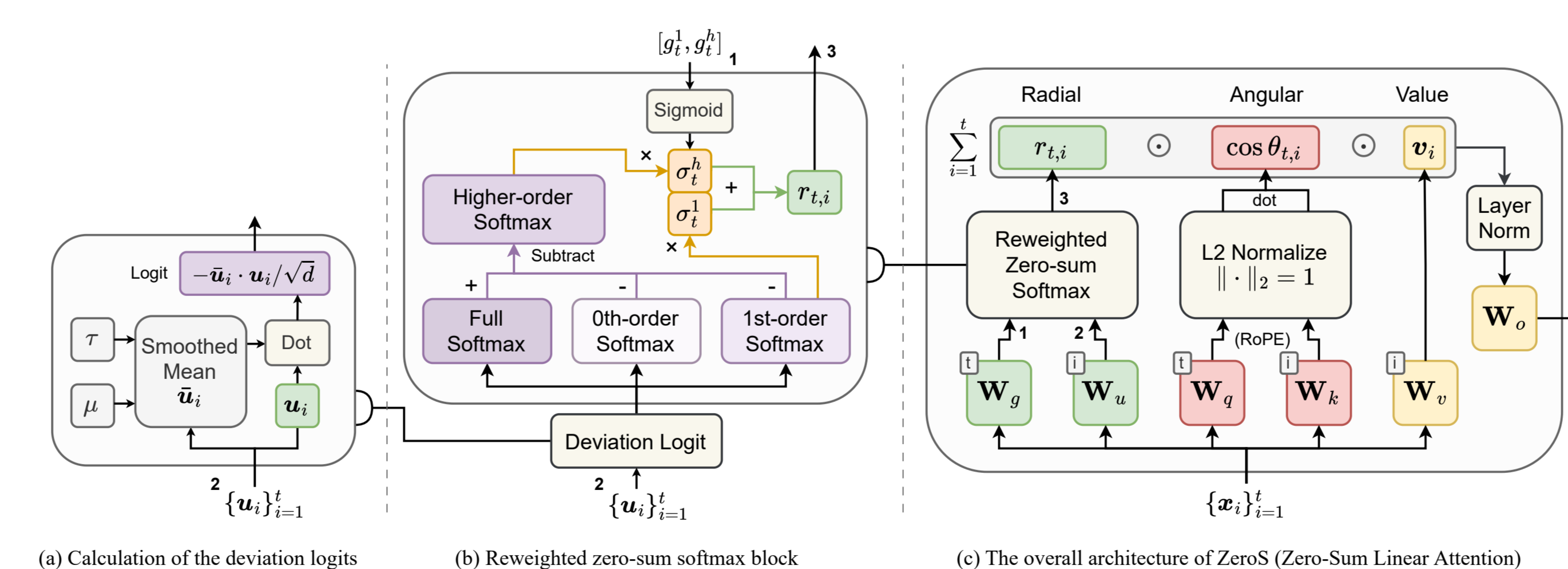


Figure 1. ZeroS block: deviation logits  $\{s_i\}$ , reweighted zero-sum softmax ( $r_{t,i}$ ), angular modulation (RoPE), linear-time scan.

## Empirical Highlights (matched parameter budgets)

**Benchmarks.** In-context learning (MAD, RegBench, MQAR), language modeling (WikiText-103, OWT2), ImageNet-1K (DeiT-Tiny), and multivariate TSF (ETT, Weather, Solar).

Model	Compress	Fuzzy	In-Ctx	Memorize	Noisy	Select	Avg
Transformer	51.6	29.8	94.1	85.2	86.8	99.6	74.5
GLA	38.8	6.9	80.8	63.3	81.6	88.6	60.0
Mamba	52.7	6.7	90.4	89.5	90.1	86.3	69.3
<b>ZeroS</b>	44.0	14.9	99.9	88.1	96.1	97.8	73.5
<b>ZeroS-SM</b>	45.2	28.0	<b>100</b>	84.3	96.6	98.5	<b>75.4</b>

## Ablations & Analysis

- **Zero-sum helps:** reintroducing the 0<sup>th</sup>-order term hurts In-Context Recall and WikiText PPL.
- **Residual reweighting matters:** replacing the reweighted zero-sum branch with vanilla softmax degrades performance, confirming the expressive gap vs. convex mixing.
- **Gating:** removing gates ( $\sigma^1, \sigma^h$ ) causes consistent but moderate drops—gates learn task-adaptive entropy/competition.
- **Normalization:** removing LayerNorm particularly harms algorithmic tasks (e.g., In-Context Recall), highlighting variance control in linear attention.

## Conclusion & Outlook

**ZeroS** removes the uniform zero-order term and reweights softmax residuals to produce *signed, zero-sum* linear attention:

- Enables single-layer *contrastive* operations and restores angular sign-flips at  $O(N)$  complexity.
- Achieves accuracy on par with or better than softmax and surpasses prior linear methods across NLP, vision, and time series.
- Maintains numerical stability (length-independent bounds) and integrates seamlessly with RoPE and existing linear-attention accelerators.