# Revisiting Residual Connections:
# Orthogonal Updates for Stable and Efficient Deep Networks

Giyeong Oh[1] Woohyun Cho[1] Siyeol Kim[1] Suhwan Choi[2] Youngjae Yu[3]

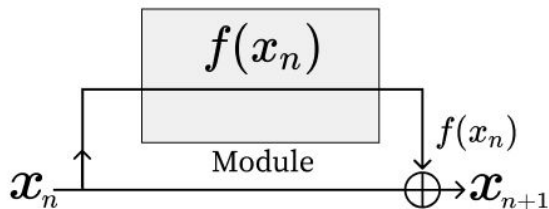Yonsei University[1] Maum.AI[2] Seoul National University[3]

# Overview & Achievement

- **Problem: Standard Residual Connection**

  - [Eq. 1] $x_{n+1} = x_n + f(x_n)$

  - [Eq. 2] $f(x_n) = f_{\parallel}(x_n) + f_{\perp}(x_n)$

- **Our Solution: Orthogonal Residual Update**

  - Update Only Orthogonal Part (New Representation)

  - [Eq. 3] $x_{n+1} = x_n + \boxed{f_{\perp}(x_n)}$

- **Result & Achievement**

  - **3.78 pp** Accuracy Gain ViT-B on ImageNet-1k

# Background

- **Residual Connection**



$$f(x_n)$$

$$f(x_n)$$

$$x_n \quad \text{Module} \quad \oplus \rightarrow x_{n+1}$$

  - [Fig. 1] Residual Connection, Pre-activation.

  - No Gradient Vanishing, enabling much deeper networks

  - Core Architecture of Deep Learning

# Background

- **Potential Problem**
  - We can decompose $f(x_n)$ into two parts.
  - Recall [Eq. 2] $f(x_n) = f_{\parallel}(x_n) + f_{\perp}(x_n)$
  - Parallel component *may* introduce redundant representation.
  - Orthogonal component can be considered as *new* representation.

# Orthogonal Residual Updates

- **Decomposition**
  - Gram-Schmidt Process
  - For $x_n, f(x_n) \in \mathbb{R}^d$
  - [Eq. 4]   $f_\|(x_n) := \dfrac{\langle x_n, f(x_n) \rangle}{\|x_n\|^2} x_n \,, \; f_\perp(x_n) := f(x_n) - \dfrac{\langle x_n, f(x_n) \rangle}{\|x_n\|^2} x_n$
  - [Eq. 5]   $f_\|(x_n) = \alpha x_n, \; \langle x_n, f_\perp(x_n) \rangle = \vec{0}$
  - $\langle \cdot, \cdot \rangle$ is inner product.

# Orthogonal Residual Updates

- **Numerical Stability**

  - For Preventing Zero Division, We adjust denominator.

$$\frac{\langle x_n, f(x_n) \rangle}{\|x_n\|^2} \rightarrow \frac{\langle x_n, f(x_n) \rangle}{\|x_n\|^2 + \epsilon} \quad \epsilon > 0$$

  - $\epsilon$ is stability constant. We use $\epsilon = 10^{-6}$.

# Orthogonal Residual Updates

- **Formulations**

  - [Eq. 6]   $s_n := \dfrac{\langle x_n, f(x_n) \rangle}{\|x_n\|^2 + \epsilon}$

  - Recall [Eq. 2]  $f(x_n) = f_{\parallel}(x_n) + f_{\perp}(x_n)$

  - [Eq. 7]   $f_{\parallel}(x_n) = s_n x_n$

  - [Eq. 8]   $f_{\perp}(x_n) = f(x_n) - s_n x_n$

  - Since $\epsilon$ exists, $\langle x_n, f_{\perp} \rangle \sim \vec{0}$. It is negligible.

# Orthogonal Residual Updates

- **What dimension to be orthogonal.**

    - Feature channel only.

        $$x \in \mathbb{R}^{B \times L \times \underline{d}}, \quad x \in \mathbb{R}^{B \times \underline{C} \times k \times k}$$

        – Apply orthogonalization *independently* to each feature vector.

    - Global, except batch dimension.

        $$x \in \mathbb{R}^{B \times (\underline{L \times d})}, \quad x \in \mathbb{R}^{B \times (\underline{C \times k \times k})}$$
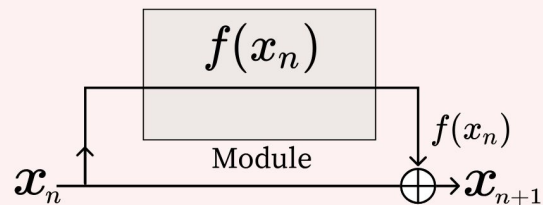
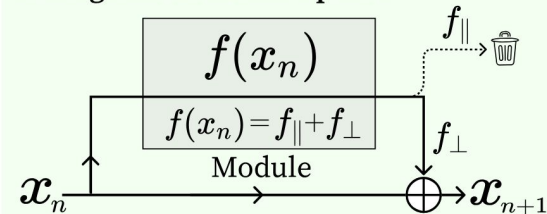        – Apply global orthogonalization to *the entire feature map*.

# Orthogonal Residual Updates

- ## Update Rule Comparison



Linear Residual Update

$f(x_n)$

Module

$x_n$   $f(x_n)$   $\oplus \rightarrow x_{n+1}$

Orthogonal Residual Update

$f(x_n)$

$f(x_n) = f_\parallel + f_\perp$

Module

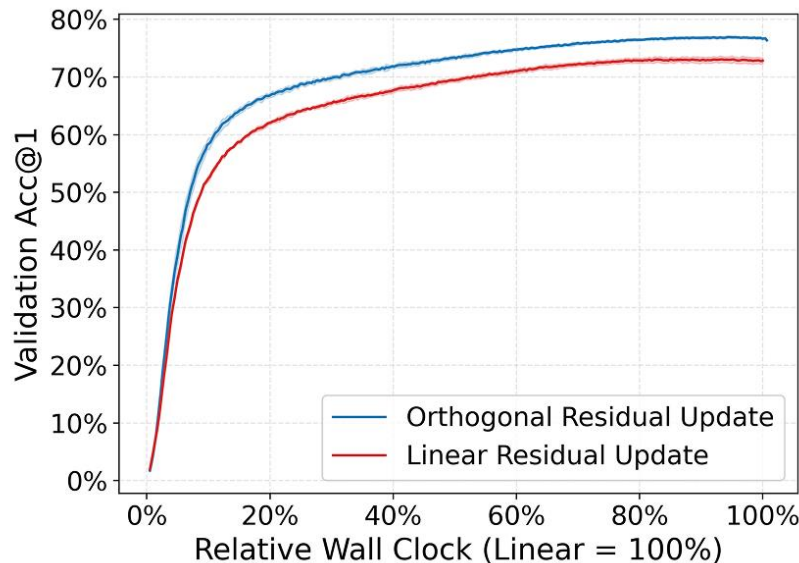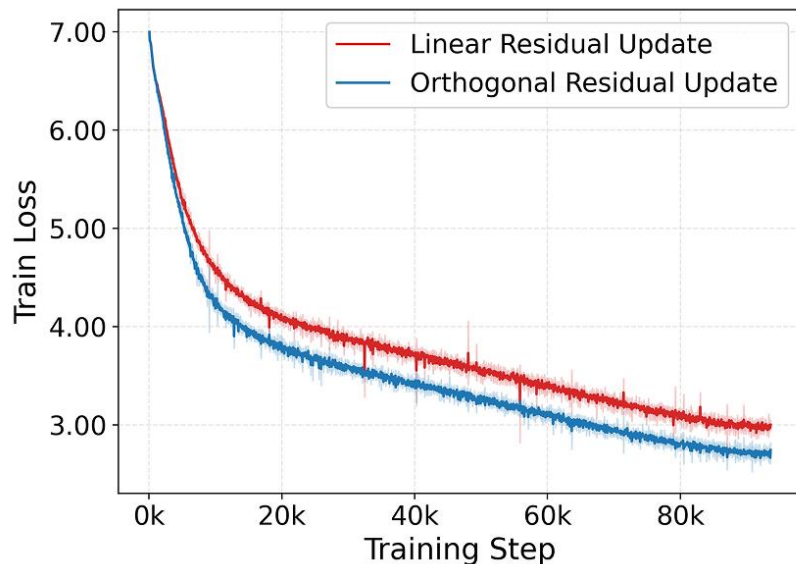$x_n$   $f_\parallel$   $f_\perp$   $\oplus \rightarrow x_{n+1}$

- Linear Residual Update (ResNet-V2)

  [Eq. 1]   $x_{n+1} = x_n + f(x_n)$

- Orthogonal Residual Update (Ours)

  [Eq. 3]   $x_{n+1} = x_n + f_\perp(x_n)$

# Experiments & Results

- **Image Classification: ViT-B, ImageNet-1K**



**Converges faster. Improves accuracy. Minimal overhead**

# Experiments & Results

- ## Image Classification: ViT, ResNet-V2

| Architecture | Connection | Dataset (Acc@1 % mean ± std.) | | | |
|---|---|---|---|---|---|
| | | CIFAR-10 | CIFAR-100 | TinyImageNet | ImageNet-1k |
| ViT-S | Linear | $89.82_{\pm 0.34}$ | $71.92_{\pm 0.24}$ | $51.30_{\pm 0.40}$ | $70.76_{\pm 0.26}$ |
| | Orthogonal-F | $\mathbf{90.61}_{\pm 0.21}$ | $\mathbf{73.86}_{\pm 0.31}$ | $\mathbf{52.57}_{\pm 0.71}$ | $\mathbf{72.53}_{\pm 0.49}$ |
| ViT-B | Linear | $87.28_{\pm 0.41}$ | $68.25_{\pm 0.88}$ | $55.29_{\pm 0.71}$ | $73.27_{\pm 0.58}$ |
| | Orthogonal-F | $\mathbf{88.73}_{\pm 6.06}$ | $\mathbf{75.07}_{\pm 0.43}$ | $\mathbf{57.87}_{\pm 0.37}$ | $\mathbf{77.05}_{\pm 0.21}$ |
| ResNetV2-18 | Linear | $95.06_{\pm 0.15}$ | $77.67_{\pm 0.28}$ | $62.04_{\pm 0.29}$ | — |
| | Orthogonal-F | $\mathbf{95.26}_{\pm 0.12}$ | $\mathbf{77.87}_{\pm 0.27}$ | $\mathbf{62.65}_{\pm 0.14}$ | |
| | Orthogonal-G | $95.25_{\pm 0.11}$ | $77.53_{\pm 0.19}$ | $62.32_{\pm 0.22}$ | |
| ResNetV2-34 | Linear | $95.49_{\pm 0.09}$ | $78.92_{\pm 0.31}$ | $64.61_{\pm 0.24}$ | — |
| | Orthogonal-F | $\mathbf{95.75}_{\pm 0.13}$ | $\mathbf{78.97}_{\pm 0.04}$ | $\mathbf{65.46}_{\pm 0.30}$ | |
| | Orthogonal-G | $95.53_{\pm 0.12}$ | $78.71_{\pm 0.24}$ | $65.38_{\pm 0.35}$ | |
| ResNetV2-50 | Linear | $\mathbf{94.75}_{\pm 0.09}$ | $\mathbf{77.90}_{\pm 0.24}$ | $63.74_{\pm 0.18}$ | — |
| | Orthogonal-F | $94.71_{\pm 0.11}$ | $77.43_{\pm 0.10}$ | $\mathbf{64.22}_{\pm 0.28}$ | |
| | Orthogonal-G | $\mathbf{94.75}_{\pm 0.10}$ | $77.56_{\pm 0.34}$ | $64.40_{\pm 0.36}$ | |
| ResNetV2-101 | Linear | $\mathbf{94.86}_{\pm 0.05}$ | $77.72_{\pm 0.33}$ | $63.77_{\pm 0.52}$ | — |
| | Orthogonal-F | $94.80_{\pm 0.13}$ | $\mathbf{78.50}_{\pm 0.26}$ | $65.78_{\pm 0.22}$ | |
| | Orthogonal-G | $94.75_{\pm 0.13}$ | $78.37_{\pm 0.19}$ | $\mathbf{65.87}_{\pm 0.23}$ | |

**Substantial gains on ViT**
- ViT, with *weaker inductive biases*, benefits more from our method's structural prior.

**Modest gains on ResNet-V2**
- ResNet's convolutional kernels are already a *powerful inductive bias*, resulting in smaller relative improvements.

**Linear**: Standard Residual Connection (pre-act)
**Orthogonal-F**: Feature-wise Orthogonalization
**Orthogonal-G**: Global Orthogonalization

# Experiments & Results

- **Image Classification: ViT-S, CIFAR-10/100.**

Table 3: ViT-S LR sweep on CIFAR-10/100: Val Acc@1 (mean±std over 3 runs).  ● **Stable Results**

| LR | CIFAR-10 | | CIFAR-100 | |
|---|---|---|---|---|
| | **Linear** | **Orthogonal** | **Linear** | **Orthogonal** |
| $5 \times 10^{-4}$ | $90.41_{\pm 0.15}$ | $\mathbf{90.56}_{\pm 0.32}$ | $71.46_{\pm 0.59}$ | $\mathbf{72.48}_{\pm 0.37}$ |
| $6 \times 10^{-4}$ | $90.70_{\pm 0.12}$ | $\mathbf{90.73}_{\pm 0.15}$ | $71.39_{\pm 0.37}$ | $\mathbf{72.95}_{\pm 0.14}$ |
| $7 \times 10^{-4}$ | $90.54_{\pm 0.23}$ | $\mathbf{91.01}_{\pm 0.22}$ | $71.55_{\pm 0.39}$ | $\mathbf{72.83}_{\pm 0.36}$ |
| $8 \times 10^{-4}$ | $90.45_{\pm 0.16}$ | $\mathbf{90.91}_{\pm 0.33}$ | $71.13_{\pm 0.92}$ | $\mathbf{72.99}_{\pm 0.29}$ |
| $9 \times 10^{-4}$ | $90.14_{\pm 0.23}$ | $\mathbf{90.57}_{\pm 0.27}$ | $70.99_{\pm 0.50}$ | $\mathbf{72.60}_{\pm 0.53}$ |
| $1 \times 10^{-3}$ | $89.95_{\pm 0.36}$ | $\mathbf{90.36}_{\pm 0.15}$ | $70.59_{\pm 0.80}$ | $\mathbf{73.11}_{\pm 0.39}$ |
| $2 \times 10^{-3}$ | $84.85_{\pm 1.07}$ | $\mathbf{87.38}_{\pm 0.42}$ | $62.17_{\pm 1.21}$ | $\mathbf{69.71}_{\pm 0.91}$ |
| $5 \times 10^{-3}$ | $66.09_{\pm 1.29}$ | $\mathbf{72.20}_{\pm 2.66}$ | $42.61_{\pm 2.82}$ | $\mathbf{48.24}_{\pm 2.50}$ |

# Experiments & Results

- **Adapting Connection Type**

**ViT-S, 22M params. 5 runs.**

| Dataset | Start Arch.→ End Arch. | | Acc@1 (%) | Acc@5 (%) |
|---|---|---|---|---|
| CIFAR-10 | Linear | → Linear | $89.82_{\pm0.34}$ | $99.65_{\pm0.03}$ |
| | Linear | → Orthogonal | $91.00_{\pm0.14}$ | $99.66_{\pm0.02}$ |
| | Orthogonal | → Linear | $\mathbf{93.18}_{\pm0.15}$ | $\mathbf{99.72}_{\pm0.03}$ |
| | Orthogonal | → Orthogonal | $90.61_{\pm0.21}$ | $99.69_{\pm0.03}$ |
| CIFAR-100 | Linear | → Linear | $71.92_{\pm0.24}$ | $92.11_{\pm0.18}$ |
| | Linear | → Orthogonal | $71.64_{\pm0.56}$ | $91.96_{\pm0.24}$ |
| | Orthogonal | → Linear | $\mathbf{74.14}_{\pm0.35}$ | $\mathbf{92.69}_{\pm0.19}$ |
| | Orthogonal | → Orthogonal | $73.86_{\pm0.31}$ | $92.23_{\pm0.26}$ |
| TinyImageNet | Linear | → Linear | $51.30_{\pm0.40}$ | $75.19_{\pm0.66}$ |
| | Linear | → Orthogonal | $50.78_{\pm0.42}$ | $73.91_{\pm0.32}$ |
| | Orthogonal | → Linear | $\mathbf{53.33}_{\pm0.62}$ | $\mathbf{76.06}_{\pm0.46}$ |
| | Orthogonal | → Orthogonal | $52.57_{\pm0.71}$ | $75.33_{\pm0.57}$ |

**Continuous Training Setup:**
- Train 150 epochs with **Start Arch.**
- Switch connection type → Continue 150 epochs with **End Arch.**
- Optimizer state is maintained throughout (no reset).

**Key Finding:**
- Starting with the Orthogonal update is crucial

**Interpretation:**
- Robust initial representations

# Experiments & Results

- **Orthogonal Connection Probability**
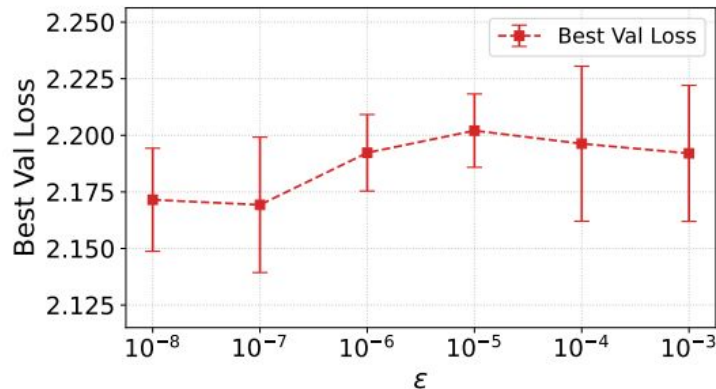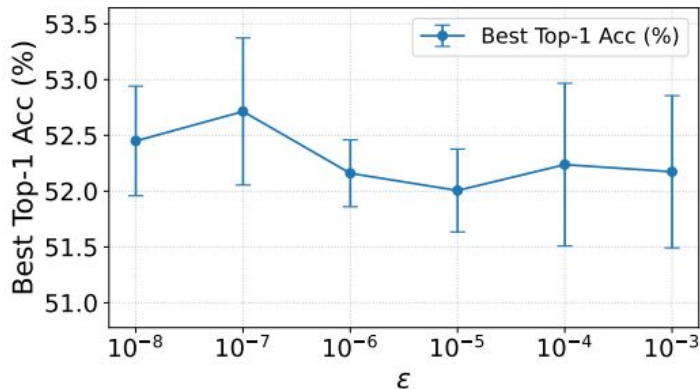


## Performance Scales with Application Probability

- Strong positive correlation observed between application probability (π) and accuracy.
- Suggests a fundamental improvement, not just a stochastic regularization effect.

# Experiments & Results

- ## Stability Constant



## Validating the Stability Constant (ε)

- Ensures numerical stability during orthogonal projection.
- Performance is robust across a wide range, with $\epsilon = 10^{-6}$ selected as default.

# Contribution Summary

🚀 **Converges faster and more stably.**

- **Enhanced Stability & Faster Convergence**: Provides a stable learning dynamic, reaching higher accuracy in less time.

🎯 **Improves accuracy across ResNets/ViTs.**

- **Broad Applicability & Performance Gains**: Demonstrates consistent improvements across diverse architectures like CNNs and Transformers.

⚡ **One-line Gram–Schmidt, negligible cost.**

- **High Efficiency**: Implemented as a simple Gram-Schmidt step with negligible computational overhead.

# Thesis Completion Plan

- **Quantitative Representation Analysis**

  - Effective Rank, CKA Similarity, … Representational Metric.

- **Practical Efficiency Analysis**

  - Measuring Practical Throughput.

- **Robustness to Hyperparameters**

  - LR Sweep

- **Expanding Discussion and Future Work**

  - Exploring Hybrid Updates (e.g., learnable balancing)

# Thank You!

# Backup Slides

Table 3: Mean ± std. of top-1 (Acc@1) and top-5 (Acc@5) accuracy (%) from 5 independent runs for ViT-S, evaluating adaptability to connection type changes. Models are trained for 300 epochs (**Start Arch.**) then another 300 epochs (**End Arch.**) on the same dataset, with connections (Linear '**L**' or Orthogonal '**O**') potentially switched. Compared are **L→L**, **L→O**, **O→L**, and **O→O** on CIFAR-10, CIFAR-100, and Tiny ImageNet. Results are averaged over the final epochs of the **End Arch.** phase.

| Dataset | Start Arch.→ End Arch. | | Acc@1 (%) | Acc@5 (%) |
|---|---|---|---|---|
| CIFAR-10 | Linear | → Linear | $92.78_{\pm0.06}$ | $99.74_{\pm0.03}$ |
| | Linear | → Orthogonal | $92.88_{\pm0.14}$ | $99.72_{\pm0.03}$ |
| | Orthogonal | → Linear | $93.89_{\pm0.12}$ | $\mathbf{99.75}_{\pm0.04}$ |
| | Orthogonal | → Orthogonal | $\mathbf{94.10}_{\pm0.12}$ | $\mathbf{99.73}_{\pm0.04}$ |
| CIFAR-100 | Linear | → Linear | $74.22_{\pm0.13}$ | $92.26_{\pm0.13}$ |
| | Linear | → Orthogonal | $74.02_{\pm0.24}$ | $91.96_{\pm0.17}$ |
| | Orthogonal | → Linear | $\mathbf{75.63}_{\pm0.17}$ | $\mathbf{92.91}_{\pm0.17}$ |
| | Orthogonal | → Orthogonal | $75.38_{\pm0.35}$ | $92.20_{\pm0.13}$ |
| TinyImageNet | Linear | → Linear | $53.24_{\pm0.13}$ | $75.25_{\pm0.21}$ |
| | Linear | → Orthogonal | $52.14_{\pm0.18}$ | $74.20_{\pm0.20}$ |
| | Orthogonal | → Linear | $\mathbf{54.58}_{\pm0.10}$ | $\mathbf{76.45}_{\pm0.24}$ |
| | Orthogonal | → Orthogonal | $53.88_{\pm0.29}$ | $75.34_{\pm0.23}$ |

# Backup Slides

(a) Approximate FLOPs per Transformer block. $s = n_{\text{seq}}$, $d = d_{\text{model}}$; FFN assumes a $4d$ expansion. Our *feature-wise* orthogonal connection introduces only $O(sd)$ FLOPs on top of the block.

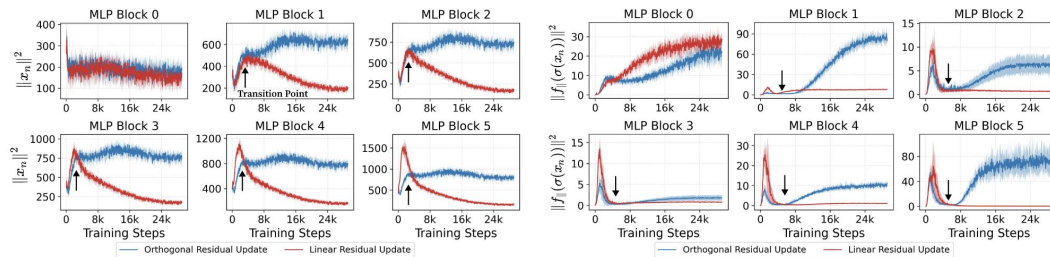| Module | Connection | Total FLOPs |
|---|---|---|
| Attention | Linear | $\approx 8sd^2 + 4s^2d + sd$ |
| | Orthogonal | $\approx 8sd^2 + 4s^2d + sd + \mathbf{6sd} + \mathbf{2s}$ |
| MLP (FFN) | Linear | $\approx 16sd^2 + sd$ |
| | Orthogonal | $\approx 16sd^2 + sd + \mathbf{6sd} + \mathbf{2s}$ |

(b) Training throughput (img/s) and overhead (%) of **Ortho-F** relative to the linear residual baseline.

| Arch. | Linear | Ortho-F | Overhead |
|---|---|---|---|
| ResNetV2-34 | 1737.2 | 1634.0 | 5.94% |
| ResNetV2-50 | 1002.8 | 876.7 | 12.58% |
| ViT-S | 3476.1 | 3466.3 | 0.28% |
| ViT-B | 1270.1 | 1246.2 | 1.88% |

Table 1: **Computation vs. practice.** Orthogonal projection adds $O(sd)$ FLOPs per block (bold in (a)); throughput in (b) is measured under identical condition.
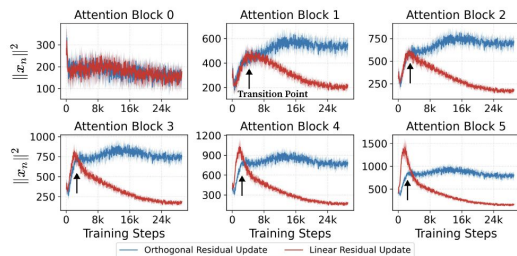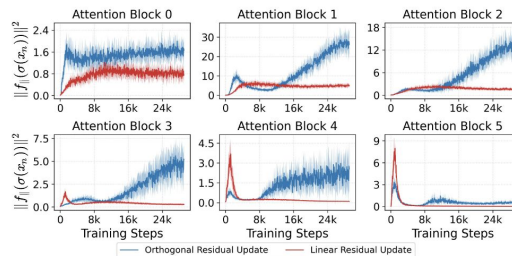
# Backup Slides



(a) Stream norm, MLP blocks.

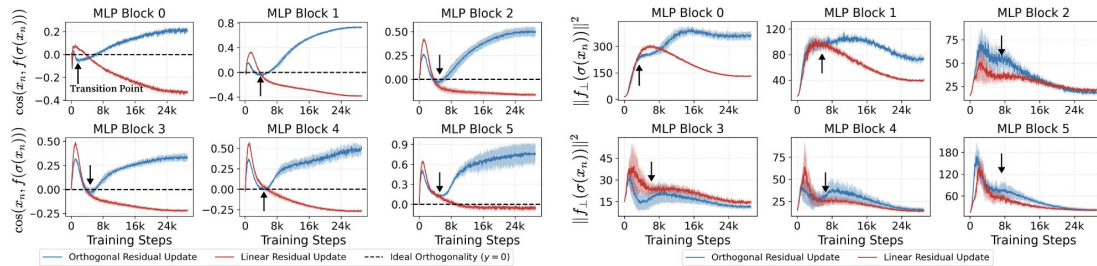(b) Parallel component norm, MLP blocks.

(c) Stream norm, Attention blocks.

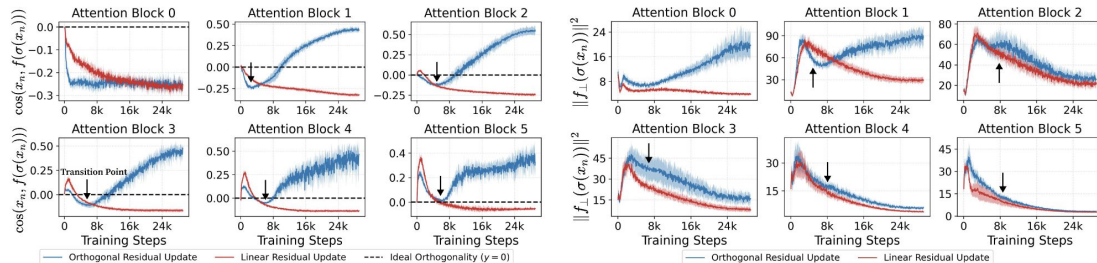(d) Parallel component norm, Attention blocks.

Figure 3: **Internal dynamics (ViT-S, TinyImageNet, 5 seeds).** Each subfigure shows blocks 0–5 (MLP top, Attention bottom). **Ours** denotes orthogonal updates; **Linear** denotes the standard residual. **(a,c)** After the *Transition Point*, orthogonal updates stabilize the stream norm $\|x_n\|$, whereas linear updates typically exhibit a post-transition decrease. **(b,d)** The parallel component energy $\|f_\|(\sigma(x_n))\|^2$ follows distinct layer-wise profiles for linear vs. orthogonal updates. Signed parallel coefficients and orthogonal-component traces are analyzed in the Appendix D.

# Backup Slides



(a) Cosine similarity (MLP blocks).
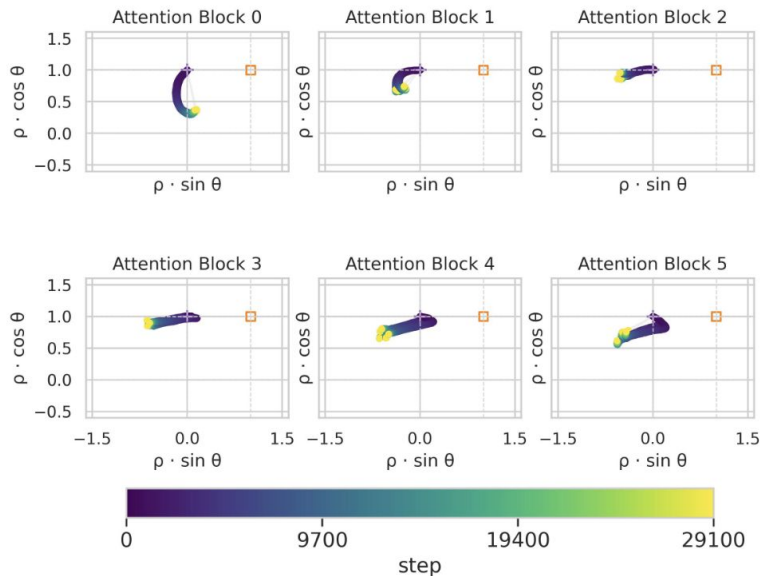
(b) Orthogonal component norm (MLP blocks).

(c) Cosine similarity (Attention blocks).

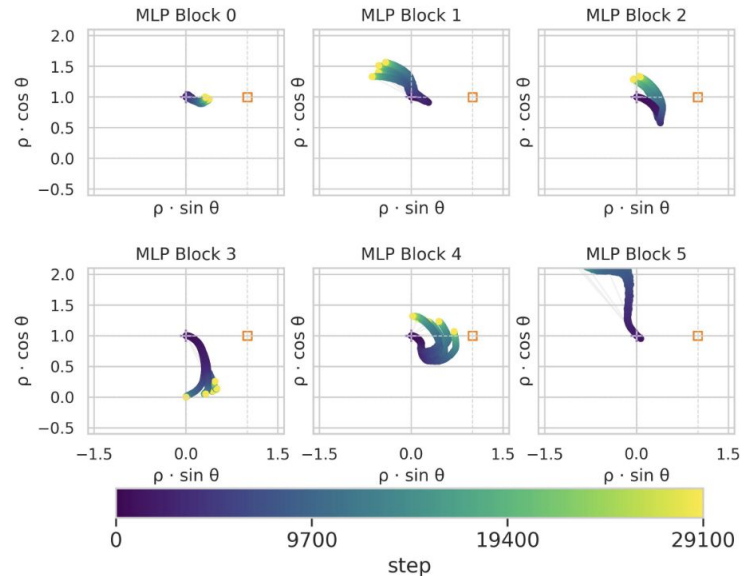(d) Orthogonal component norm (Attention blocks).

Figure 12: **Cosine similarity and orthogonal norm across layers (ViT–S, Tiny ImageNet, 5 seeds).** Each subfigure aggregates blocks 0–5 for MLP (top) and Attention (bottom). Ours (orthogonal updates) vs. Linear (standard residual). Orthogonal updates preserve the orthogonal component energy, while cosine trajectories diverge around the Transition Point. For stream norms and a broader view of alignment, see Fig. 11.

# Backup Slides

(a) Orthogonal start — Attention blocks.      (b) Orthogonal start — MLP blocks.

$$x_{n+1} = x_n + \rho_\ell\big(\sin\theta_\ell\, f_\parallel(x_n) + \cos\theta_\ell\, f_\perp(x_n)\big), \qquad \rho_\ell \geq 0,\ \theta_\ell \in \big[-\tfrac{\pi}{2}, \tfrac{\pi}{2}\big],$$

# Backup Slides



(c) Linear start — Attention blocks.

(d) Linear start — MLP blocks.

$$x_{n+1} = x_n + \rho_\ell \big( \sin \theta_\ell \, f_\parallel(x_n) + \cos \theta_\ell \, f_\perp(x_n) \big), \qquad \rho_\ell \geq 0, \; \theta_\ell \in \big[ -\tfrac{\pi}{2}, \tfrac{\pi}{2} \big],$$
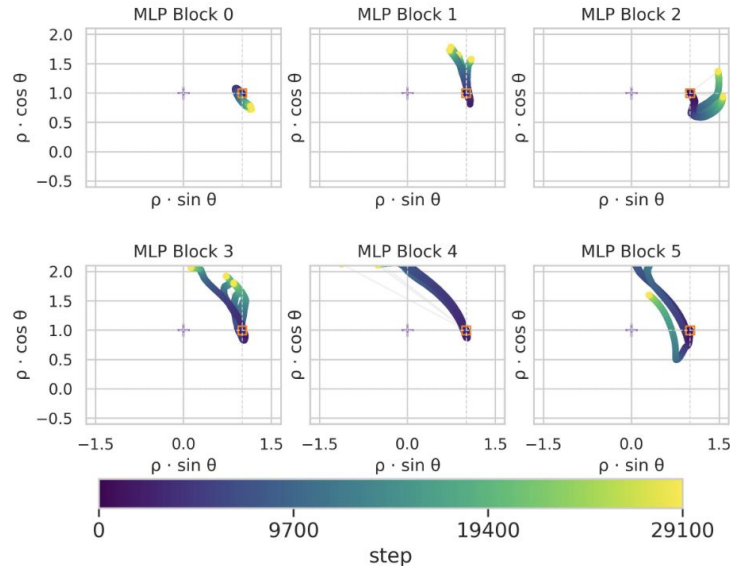
# Backup Slides

Table 4: **Representational metrics** on ViT–B (ImageNet–1k).

| Metric | Linear | Orthogonal | $\Delta$ |
|---|---|---|---|
| Effective Rank | 572.9 | **599.9** | +4.7% |
| Spectral Entropy | 6.512 | **6.539** | +0.41% |
| CKA (linear) | — | 0.546 | — |
| Feature Std. Dev. | 0.407 | **0.193** | $-52.5\%$ |

- **Effective Rank** is defined as $\exp(H)$, where $H := -\sum_i p_i \log p_i$ is the spectral entropy, and $p_i := \lambda_i / \sum_j \lambda_j$ are the normalized eigenvalues of the feature covariance matrix.
- **CKA (linear)** [28]: A similarity metric for two representations $X, Y$ defined as $\frac{\|X^\top Y\|_F^2}{\|X^\top X\|_F \|Y^\top Y\|_F}$.
- **Feature Std. Dev.**: The average per-feature standard deviation, $\frac{1}{d} \sum_{k=1}^{d} \mathrm{Std}(F_{\cdot k})$.

# Backup Slides

```python
def _orthogonal_channel(x: torch.Tensor, f_x: torch.Tensor, dim: int, eps:
    torch.Tensor) -> torch.Tensor:
    """
    Orthogonal residual connection (channel-wise).
    x   : residual stream tensor
    f_x : module output tensor (e.g., from Attention, MLP, or Conv if channel-wise)
    dim : dimension along which to compute orthogonality (e.g., channel dimension)
    eps : small epsilon tensor for numerical stability
    """
    # Ensure eps is on the same device as x if it's a tensor
    eps = eps.to(x.device)

    dot_product = (x * f_x).sum(dim, keepdim=True)
    norm_x_squared = (x * x).sum(dim, keepdim=True).float() + eps

    # Ensure scale is cast back to original dtype if x was float16/bfloat16
    scale_factor = (dot_product / norm_x_squared).to(dtype=x.dtype)

    projection_onto_x = scale_factor * x
    f_orthogonal = f_x - projection_onto_x

    return f_orthogonal
```