



Towards Understanding the Importance of Shortcut Connections in Residual Networks

Tianyi Liu¹, Minshuo Chen¹, Mo Zhou², Simon S. Du³, Enlu Zhou¹, Tuo Zhao¹
¹ Georgia Institute of Technology ² Duke University
³ Institute for Advanced Study



Background

Success of Deep Neural Networks (DNNs):

- Speech and image recognition;
- Nature Language Processing;
- Recommendation Systems.

Among different types of networks, ResNet is a Milestone!

- Shortcut connections: skip layers in the forward step of an input.
- Success over CNNs: He et al.(2016a), He et al.(2016b), Srivastava et al.(2015), Huang et al.(2017).
- Our Empirical Observation:

# of Layers	≤ 30	≥ 30
CNN	Good	Bad
RNN	Good	Good

- **Shortcut connections helps training.**

Existing Results:

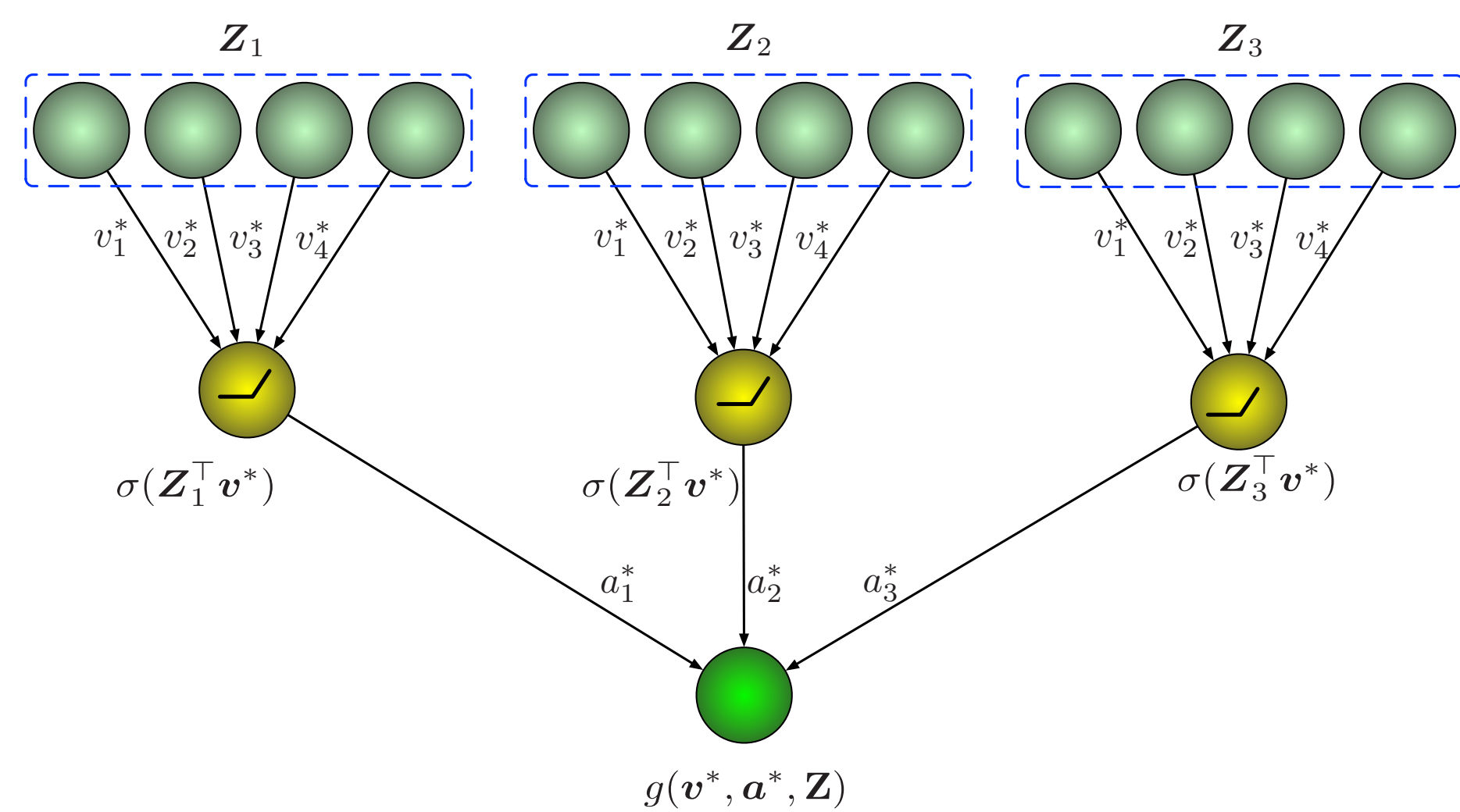
- Empirical: Veit et al. (2016), Balduzzi et al. (2017), Li et al. (2018).
- Hardt and Ma (2016): Linear ResNet has no spurious optima.
- Li and Yuan (2017): Two-layer ResNet with only one unknown layer has no spurious local optima and saddle points.

Question: How does the Shortcut Connection help training in the presence of bad optima?

- We Study: Two-Layer Nonoverlapping Convolutional NNs:
 1. A non-trivial spurious local optimum;
 2. GD gets trapped with constant probability ($\frac{1}{4} \sim \frac{3}{4}$);

A non-trivial example provides new insights!

Teacher Network



- Two-layer Nonoverlapping CNNs:

$$f(w^*, a^*, Z) = \sum_{j=1}^k a_j^* \sigma(Z_j^\top w^*),$$

- $\|w^*\|_2 = 1$, $w \in \mathbb{R}^p$, $a \in \mathbb{R}^k$, $\sigma(\cdot) = \max\{\cdot, 0\}$.
- $Z = [Z_1, \dots, Z_k]$ with Z_j 's i.i.d. $N(0, I)$,

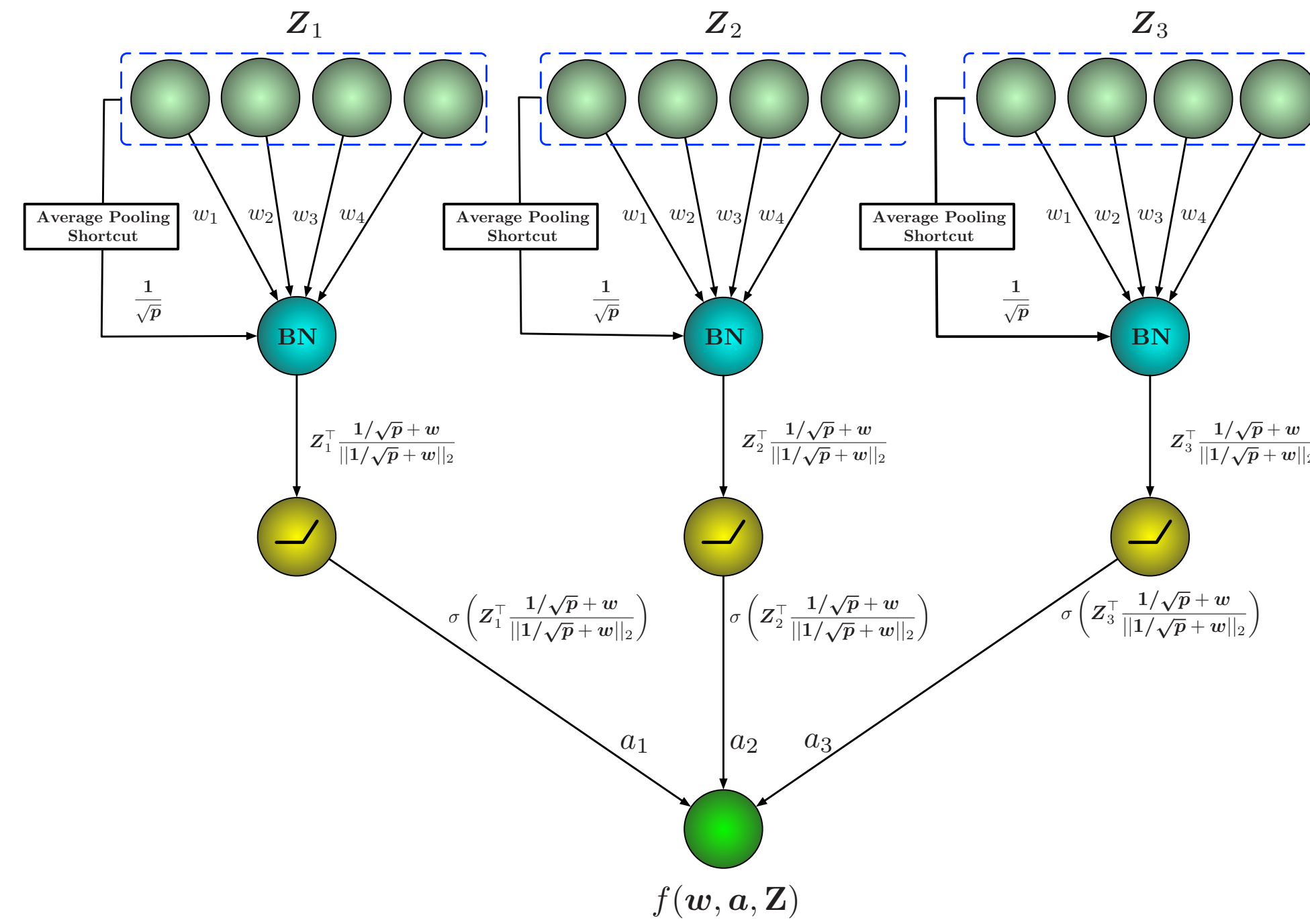
Student Network

- Student Network with shortcut connection:

$$h(w, a, Z) = \sum_{j=1}^k a_j \sigma \left(Z_j^\top \left(\frac{1}{\sqrt{p}} + w \right) \right)$$

- Normalization to achieve identifiability:

$$h(w, a, Z) = \sum_{j=1}^k a_j \sigma \left(Z_j^\top \frac{\frac{1}{\sqrt{p}} + w}{\|\frac{1}{\sqrt{p}} + w\|_2} \right).$$



- Nonconvex Optimization:

$$(\hat{w}, \hat{a}) = \operatorname{argmin}_{w, a} \mathcal{L}(w, a),$$

where $\mathcal{L}(w, a) = \mathbb{E}_Z(f(v^*, a^*, Z) - h(w, a, Z))^2$.

- (w, a) is a global optimum, if

$$\frac{1}{\sqrt{p}} + w = \alpha v^* \text{ and } a = a^*.$$

- (w, a) is a spurious local optimum, if

$$\bar{w} = -w^*, \bar{a} = (\mathbf{1}\mathbf{1}^\top + (\pi - 1)\mathbf{I})^{-1}(\mathbf{1}\mathbf{1}^\top - \mathbf{I})a^*.$$

Gradient Descent with Normalization

- Initialization: $a_0 \in \mathbb{B}_0(|\mathbf{1}^\top a^*|/\sqrt{k})$ and $w_0 = 0$.
- At the t -th iteration, we update w and a by

$$\begin{aligned} \tilde{w}_{t+1} &= w_t - \eta_w \nabla_w \mathcal{L}(w_t, a_t), \\ w_{t+1} &= \frac{\frac{1}{\sqrt{p}} + \tilde{w}_{t+1}}{\|\frac{1}{\sqrt{p}} + \tilde{w}_{t+1}\|_2} - \frac{1}{\sqrt{p}}, \\ a_{t+1} &= a_t - \eta_a \nabla_a \mathcal{L}(w_t, a_t). \end{aligned}$$

where $\mathcal{L}(w, a) = \mathbb{E}_Z(f(v^*, a^*, Z) - h(w, a, Z))^2$.

- Normalization ensures

$$\operatorname{Var} \left(Z_j^\top \left(\frac{1}{\sqrt{p}} + w_{t+1} \right) \right) = 1,$$

\iff a population version of the batch normalization.

Skip-Layer Prior

Assumption. There exists a w^* with $\|w^*\|_2 \leq 1$, such that $v^* = w^* + \mathbf{1}/\sqrt{p}$.

- Supported by Existing Results:

- Li et al. (2016) and Yu et al. (2018): The weight has a small and vanishing magnitude.
- Hardt and Ma (2016): For linear ResNet, the norm of the weight in each layer scales as $O(1/D)$ with D being the depth.
- Bartlett et al. (2018): The norm of the weight of order $O(\log D/D)$ is sufficient to express differentiable functions.

Convergence Analysis

Partial Dissipativity Condition: Given any $\delta \geq 0$ and $c \geq 0$,

$$C1: \langle -\nabla_w \mathcal{L}(w, a), w^* - w \rangle \geq c \|w - w^*\|_2^2 - \delta;$$

$$C2: \langle -\nabla_a \mathcal{L}(w, a), a^* - a \rangle \geq c \|a - a^*\|_2^2 - \delta;$$

- **Stage I: Avoid the spurious local optimum:**

- C1 holds \iff Improvement of a .
- C2 does not hold, but w **will not** move far away!

Theorem 1. Initialize with arbitrary $a_0 \in \mathbb{B}_0\left(\frac{|\mathbf{1}^\top a^*|}{\sqrt{k}}\right)$ and $w_0 = 0$. We choose step sizes

$$\eta_a = \frac{\pi}{20(k + \pi - 1)^2} = O\left(\frac{1}{k^2}\right), \eta_w = C \|a^*\|_2^2 \eta_a^2 = \tilde{O}(\eta_a^2)$$

for some constant $C > 0$. Then, we have

$$\phi_t \leq \frac{5\pi}{12} \text{ and } 0 \leq m \leq a_t^\top a^* \leq M, \quad (1)$$

for all $t \in [T_1, T]$, where $0 < m < M$ are some constants and

$$T_1 = \tilde{O}\left(\frac{1}{\eta_a}\right), \quad T = O\left(\frac{1}{\eta_a^2}\right).$$

- **Stage II: Converging to Global Optima:**

- C1, C2 jointly hold \iff Convergence!

Theorem 2. Given the output (1), for any $\delta > 0$, choose

$$\eta_a = \eta_w = \eta = \min \left\{ \frac{m}{2M^2}, \frac{5\pi^2}{4(k + \pi - 1)^2} \right\} = \tilde{O}\left(\frac{1}{k^2}\right),$$

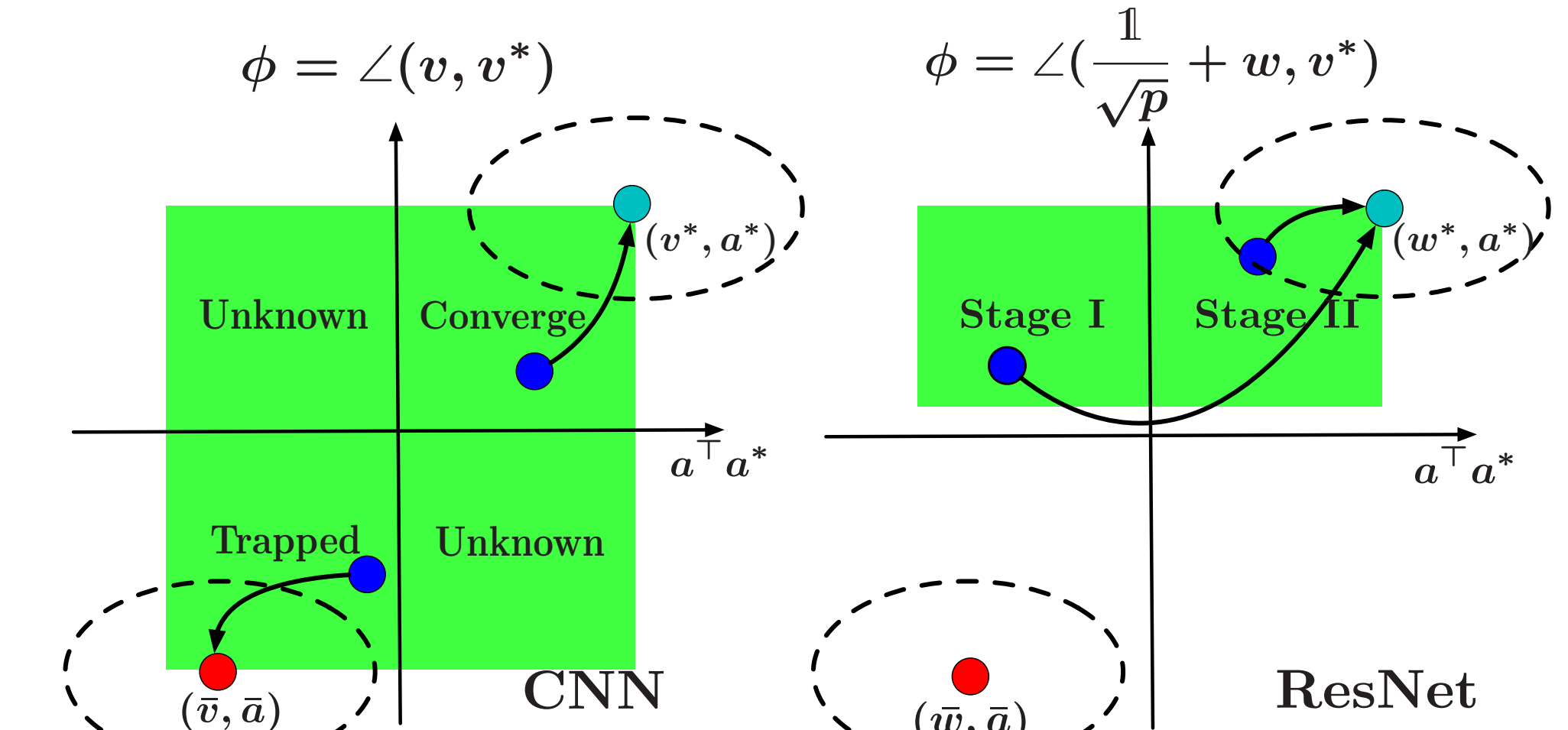
then we have

$$\|w_t - w^*\|_2^2 \leq \delta \text{ and } \|a_t - a^*\|_2^2 \leq 5\delta$$

for any $t \geq T_2 = \tilde{O}\left(\frac{1}{\eta} \log \frac{1}{\delta}\right)$.

- Remark: **Step Size Warm Up:** $\eta_w^1 < \eta_w^2$.

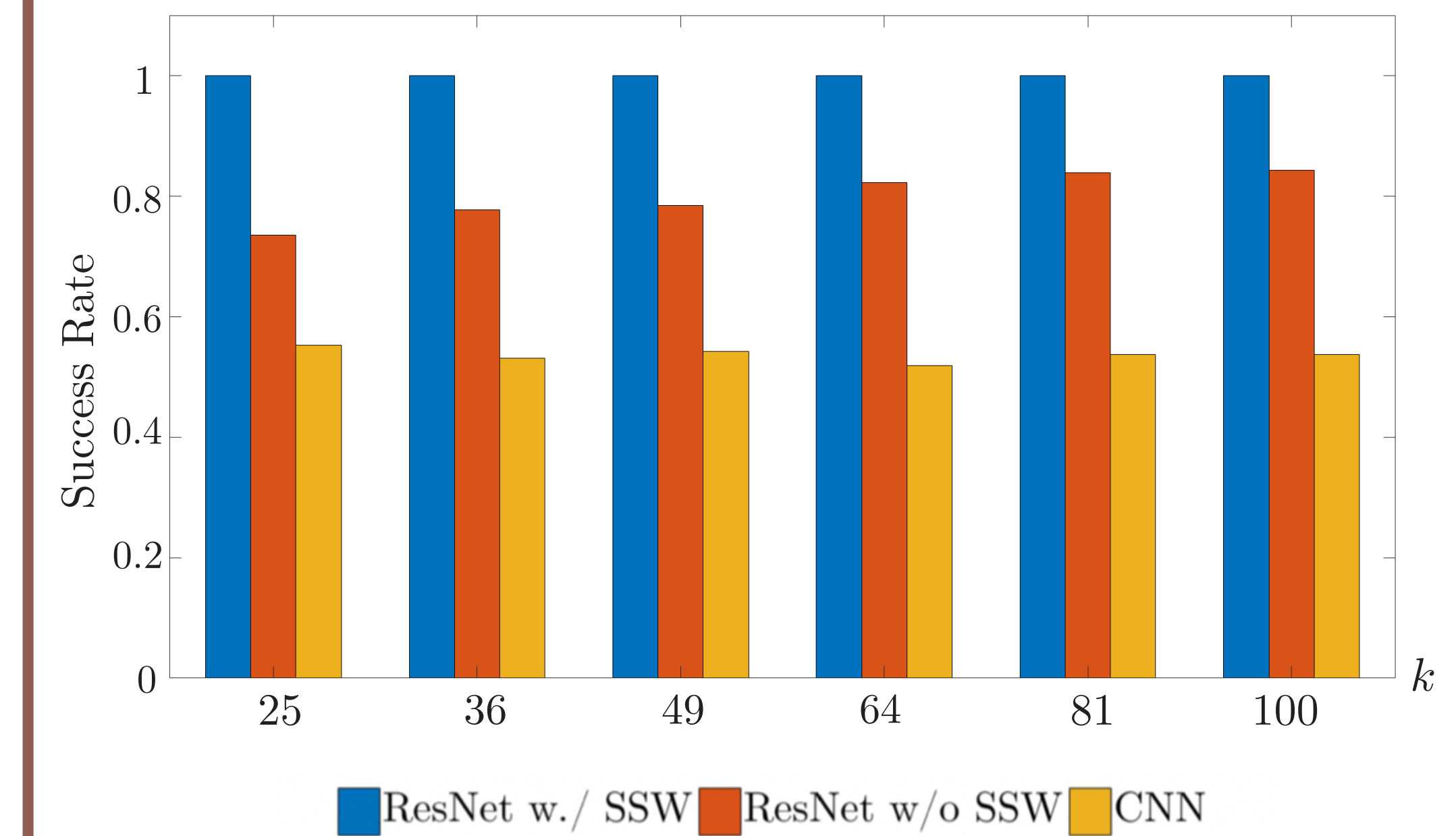
Comparison



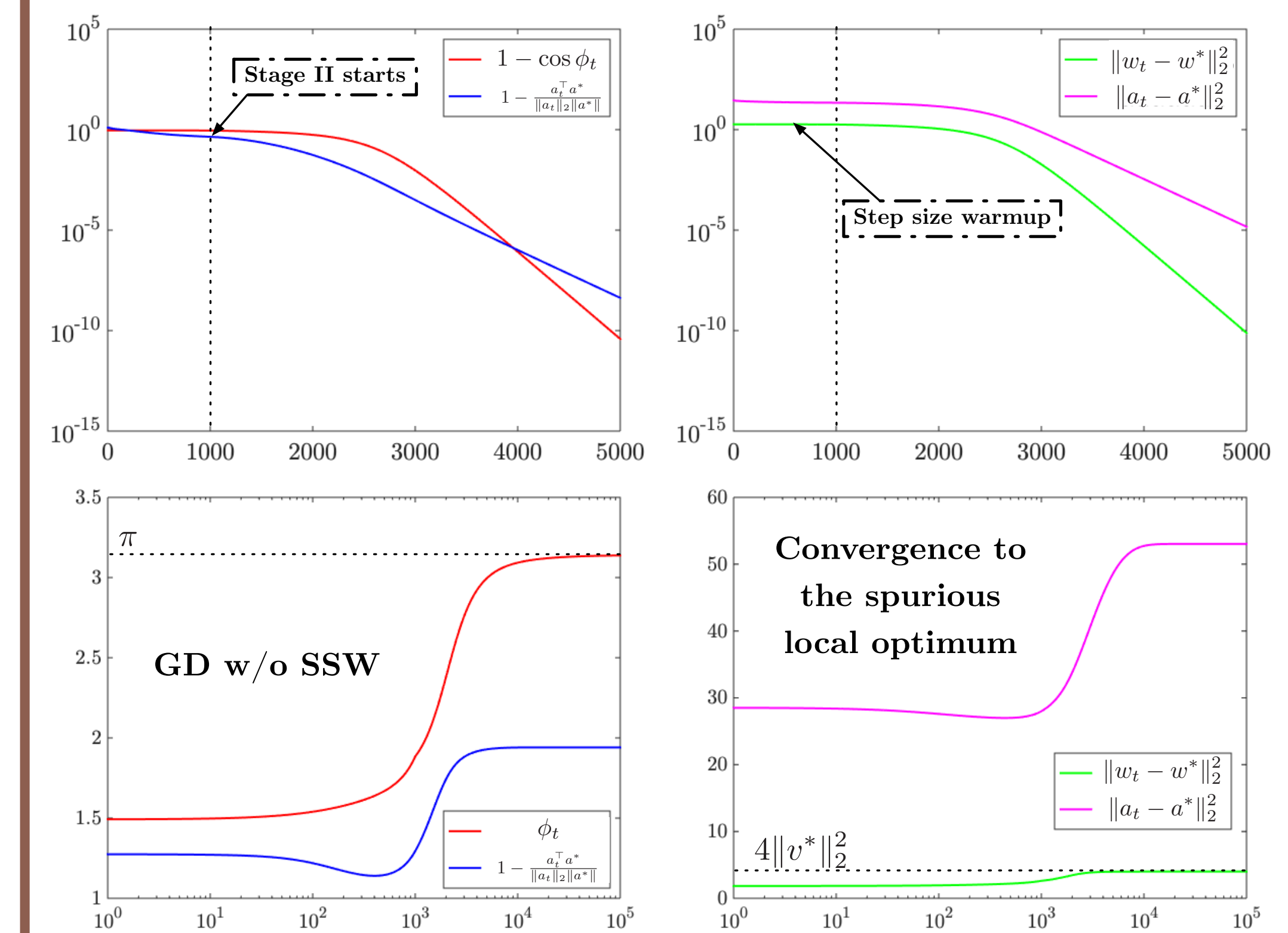
Skip-layer prior helps avoid spurious local optima!

Experiments

- Success Rates with $p = 8$ and Varing k .



- The skip-layer connection improves the success rate.
- Step size warm-up makes ResNet even better.
- Empirical Convergence:



- 1st Row: The algorithm has a phase transition.
- 2nd Row: GD w/o SSW is trapped in the spurious local optimum.