



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. Without skip-layer connections, GD gets trapped with constant probability ($\frac{1}{4} \sim \frac{3}{4}$);

A non-trival example provides new insights!

Two-layer Nonoverlapping CNNs

- Teacher Network Model:

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

- $\|\mathbf{w}^*\|_2 = 1$, $\mathbf{w} \in \mathbb{R}^p$, $\mathbf{a} \in \mathbb{R}^k$, $\sigma(\cdot) = \max\{\cdot, 0\}$.
- $\mathbf{Z} = [\mathbf{Z}_1, \dots, \mathbf{Z}_k]$ with \mathbf{Z}_j 's i.i.d. $N(\mathbf{0}, \mathbf{I})$.

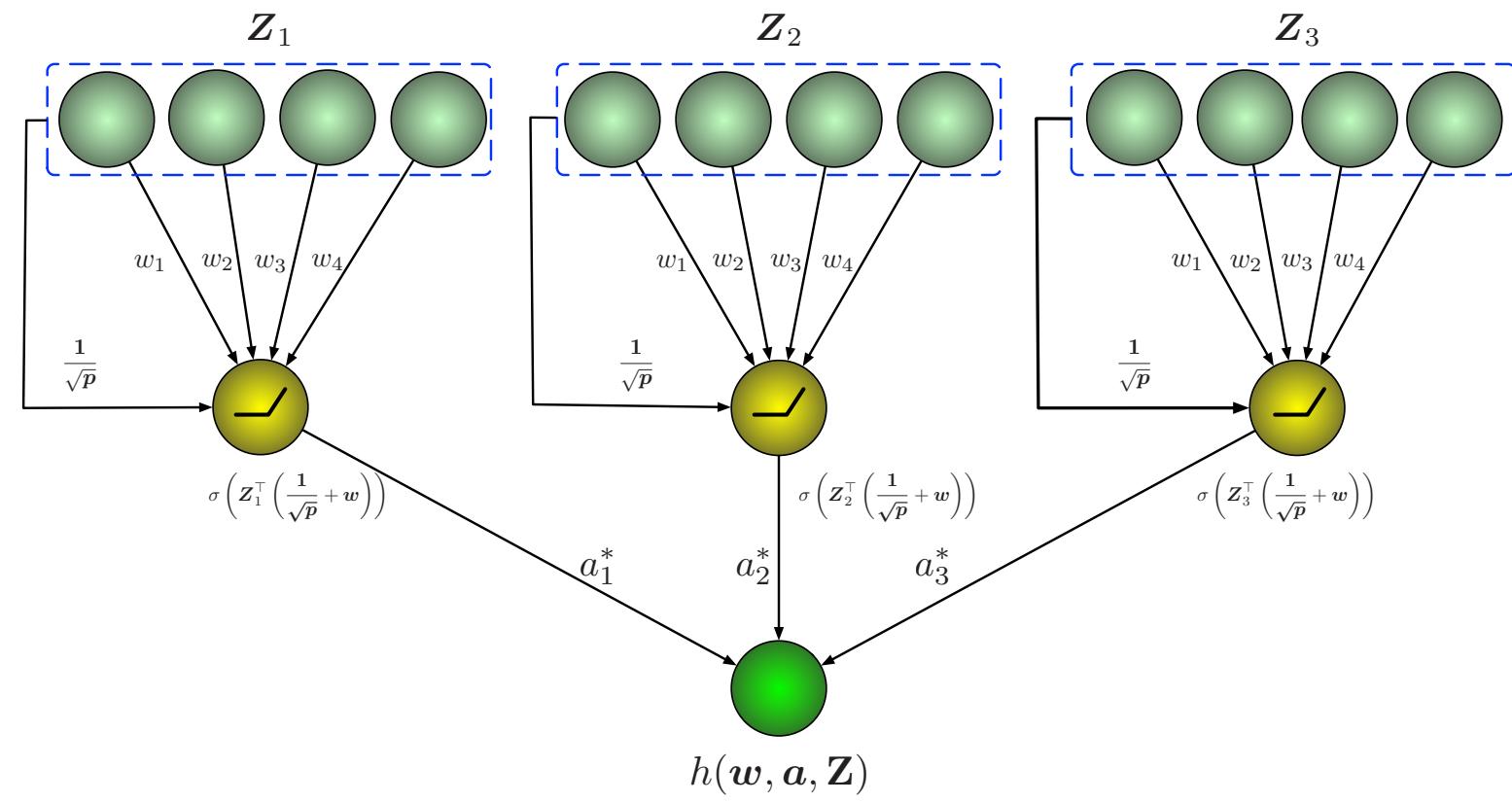
- Student Network with shortcut connection:

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

- Normalization to achieve identifiability:

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

Two-layer Nonoverlapping CNNs



- Nonconvex Optimization:

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

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

- (\mathbf{w}, \mathbf{a}) is a global optimum , if

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

- (\mathbf{w}, \mathbf{a}) is a spurious local optimum, if

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

Gradient Descent with Normalization

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

$$\begin{aligned} \tilde{\mathbf{w}}_{t+1} &= \mathbf{w}_t - \eta_w \nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}_t, \mathbf{a}_t), \\ \mathbf{w}_{t+1} &= \frac{\mathbf{1}/\sqrt{p} + \tilde{\mathbf{w}}_{t+1}}{\|\mathbf{1}/\sqrt{p} + \tilde{\mathbf{w}}_{t+1}\|_2} - \frac{\mathbf{1}}{\sqrt{p}}, \\ \mathbf{a}_{t+1} &= \mathbf{a}_t - \eta_a \nabla_{\mathbf{a}} \mathcal{L}(\mathbf{w}_t, \mathbf{a}_t). \end{aligned}$$

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

- Normalization ensures

$$\text{Var}\left(\mathbf{Z}_j^\top (\mathbf{1}/\sqrt{p} + \mathbf{w}_{t+1})\right) = 1,$$

\iff a population version of the batch normalization.

Skip-Layer Prior

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

- Supported by Existing Results:

- ILi 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 a constant $c \geq 0$,

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

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

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

- C1 holds \iff Improvement of \mathbf{a} .

- C2 does not hold, but \mathbf{w} will not move far away!

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

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

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

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

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) in Theorem 1, for any $\delta > 0$, choose

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

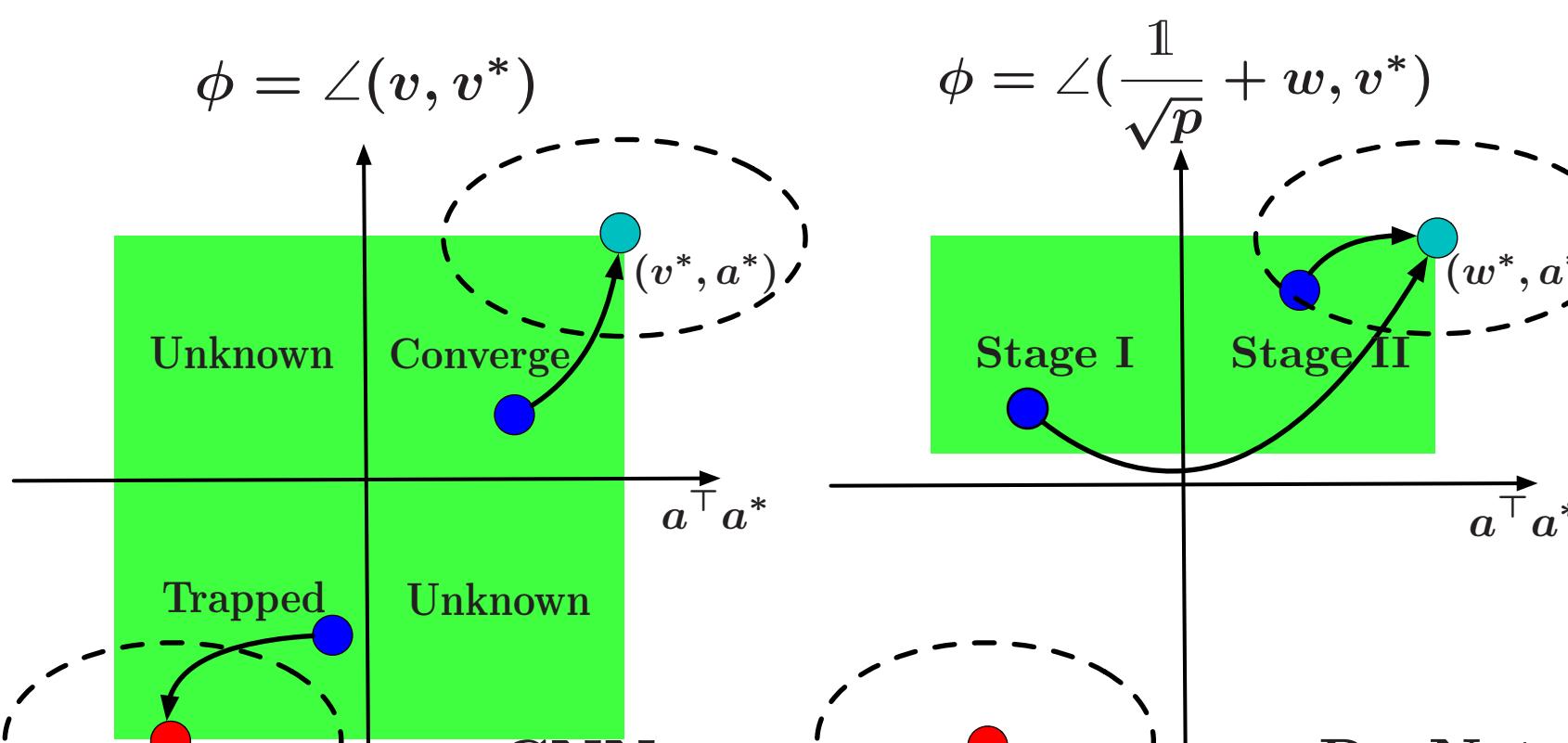
then we have

$$\|\mathbf{w}_t - \mathbf{w}^*\|_2^2 \leq \delta \quad \text{and} \quad \|\mathbf{a}_t - \mathbf{a}^*\|_2^2 \leq 5\delta$$

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

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

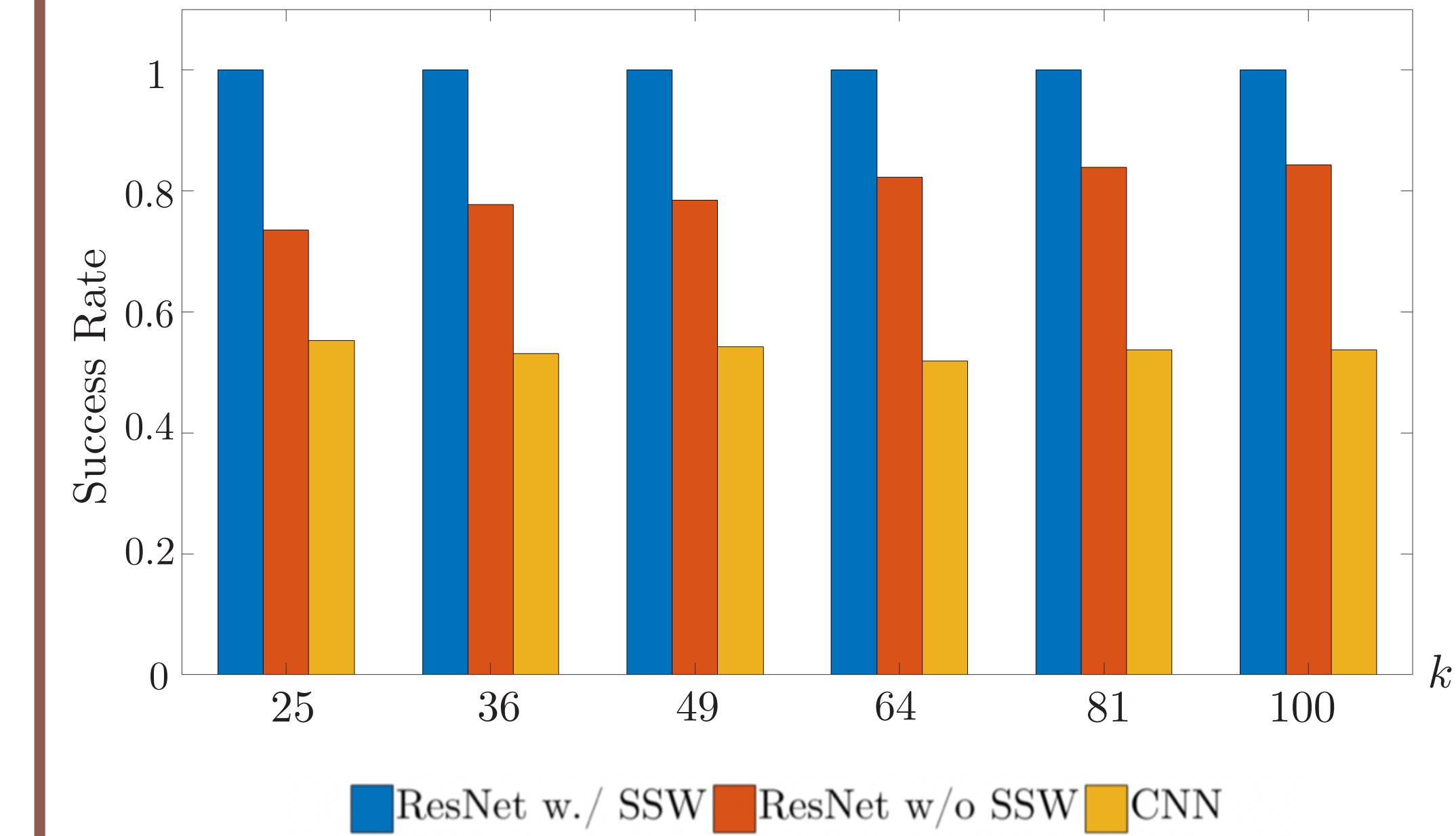
Comparison



Skip-layer prior helps avoid spurious local optima!

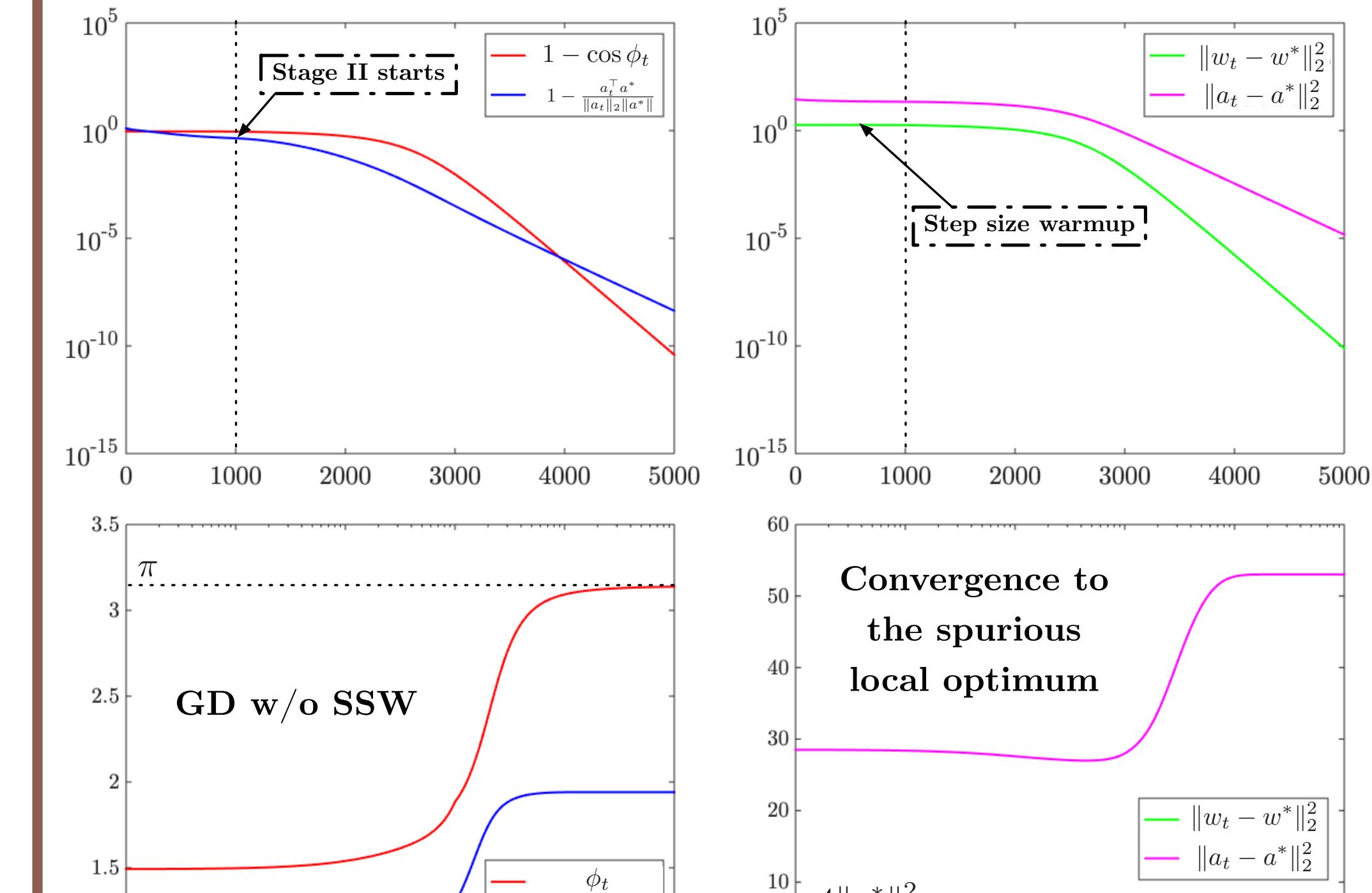
Experiments

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



- For ResNet, GD converge to the global optimum with **higher** probability.
- Step size warm-up makes ResNet even better (100% convergence).
- GD can get trapped with probability **50%**.

- Empirical Convergence:



- First Row: The algorithm has a phase transition.
- Second Row: GD w/o is trapped in the spurious local optima.

References

- [1] Du, S. S., Lee, J. D., Tian, Y., Poczos, B. and Singh, A. (2017). Gradient descent learns one-hidden- layer cnn: Don't be afraid of spurious local minima.In International Conference on Machine Learning 2018.
- [2] Li, Y. and Yuan, Y. (2017). Convergence analysis of two-layer neural networks with relu activation. In Advances in Neural Information Processing Systems.