

# Learning What and Where to Transfer

Yunhun Jang<sup>\*1,2</sup>, Hankook Lee<sup>\*1</sup>, Sung Ju Hwang<sup>3,4,5</sup>, Jinwoo Shin<sup>1,4,5</sup>

<sup>1</sup> School of Electrical Engineering, KAIST

<sup>2</sup> OMNIOUS

<sup>3</sup> School of Computing, KAIST

<sup>4</sup> Graduate School of AI, KAIST

<sup>5</sup> AITRICS

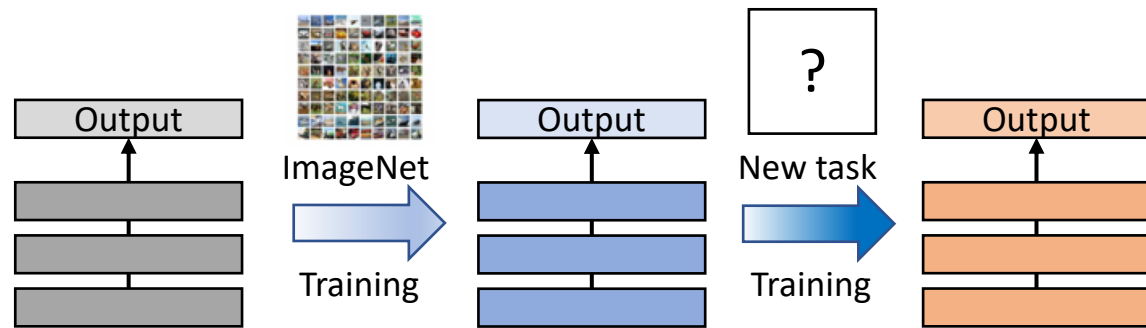
\* Equal contribution

# Transfer Learning

- DNNs require large labeled datasets to train
- *Transfer learning* is a popular method to mitigate the lack of samples
  - Improve the performance of a model *on a new task*
  - By utilizing the *knowledge* of pre-trained *source models*

# Transfer Learning

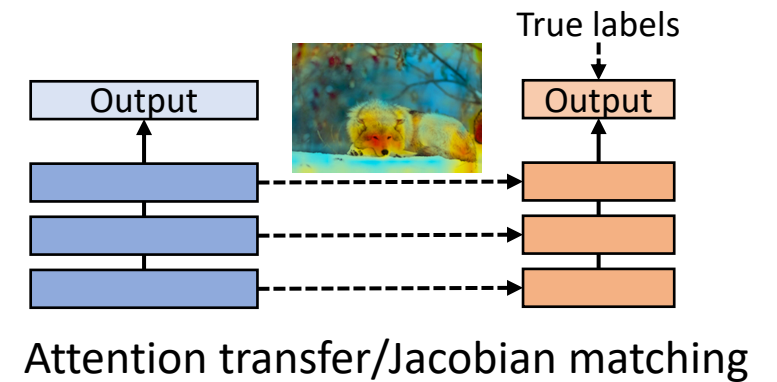
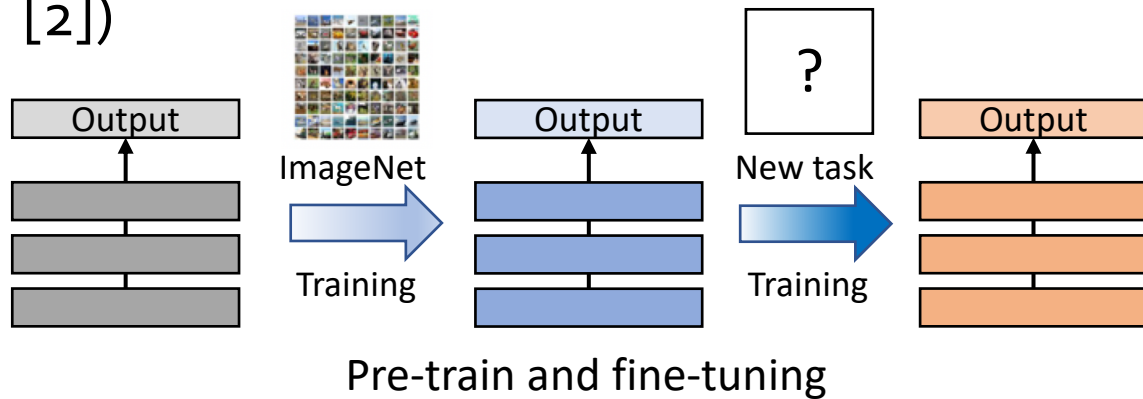
- DNNs require large labeled datasets to train
- *Transfer learning* is a popular method to mitigate the lack of samples
  - Improve the performance of a model *on a new task*
  - By utilizing the *knowledge* of pre-trained *source models*
- Limitations of previous methods
  - Require the same architecture between a source and target models (e.g., fine-tuning)



Pre-train and fine-tuning

# Transfer Learning

- DNNs require large labeled datasets to train
- *Transfer learning* is a popular method to mitigate the lack of samples
  - Improve the performance of a model *on a new task*
  - By utilizing the *knowledge* of pre-trained *source models*
- Limitations of previous methods
  - Require the same architecture between a source and target models (e.g., fine-tuning)
  - Require exhaustive hand-crafted tuning (e.g., attention transfer [1], Jacobian matching [2])



[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR 2017*

[2] Srinivas, S. and Fleuret, F. Knowledge transfer with Jacobian matching. In *Proceedings of the 35th International Conference on Machine Learning (ICML 2018)*, 2018.

# Learning What/Where to Transfer

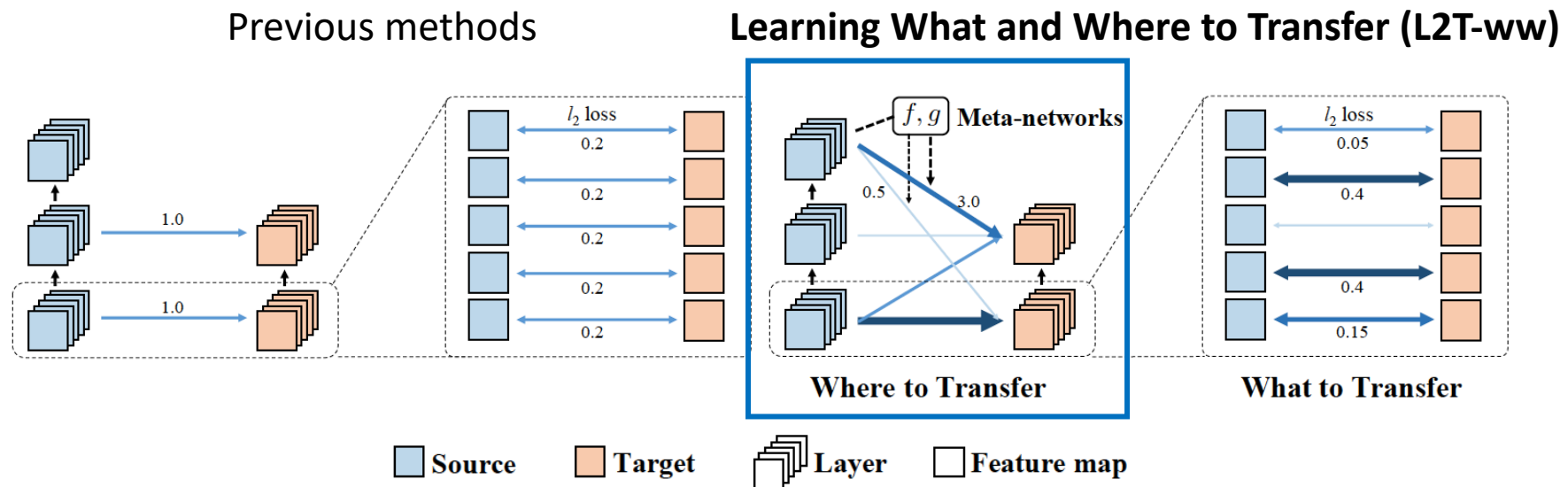
- Propose meta-networks  $f$  and  $g$ 
  - *Learn the learning rules* to transfer the source knowledge

# Learning What/Where to Transfer

- Propose meta-networks  $f$  and  $g$ 
  - *Learn the learning rules* to transfer the source knowledge

*Where* to transfer

- A meta-network  $g$  decides useful pairs of source/target layers to transfer



# Learning What/Where to Transfer

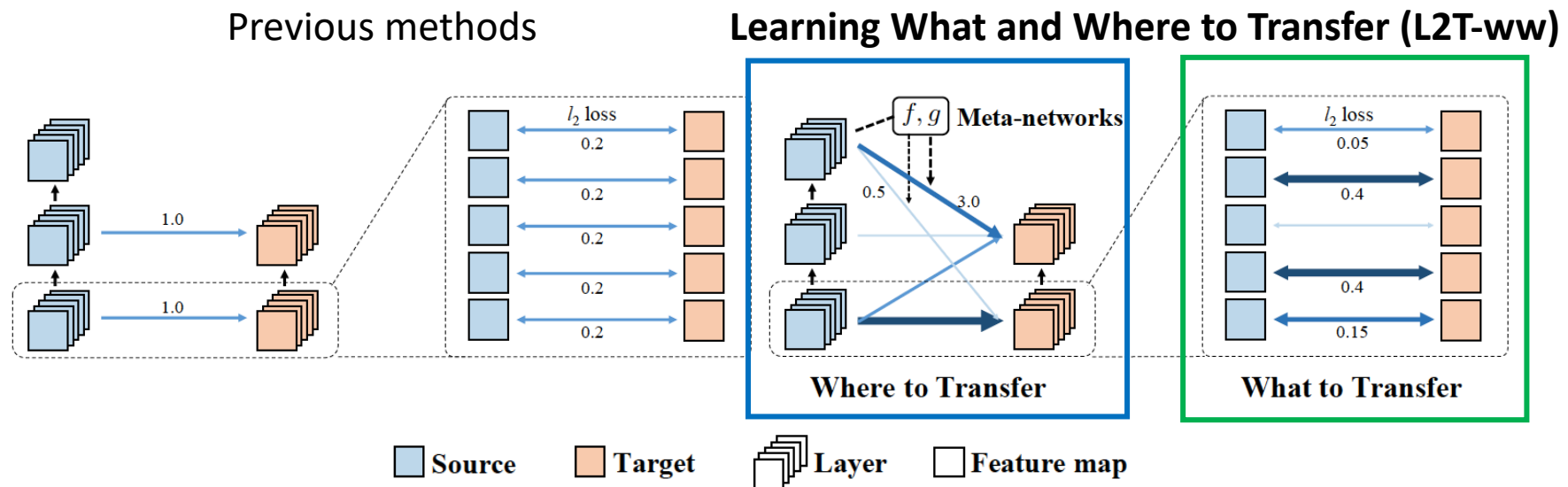
- Propose meta-networks  $f$  and  $g$  : Learning what/where to transfer (L2T-ww)
  - *Learn the learning rules* to transfer the source knowledge

*Where* to transfer

- A meta-network  $g$  decides useful pairs of source/target layers to transfer

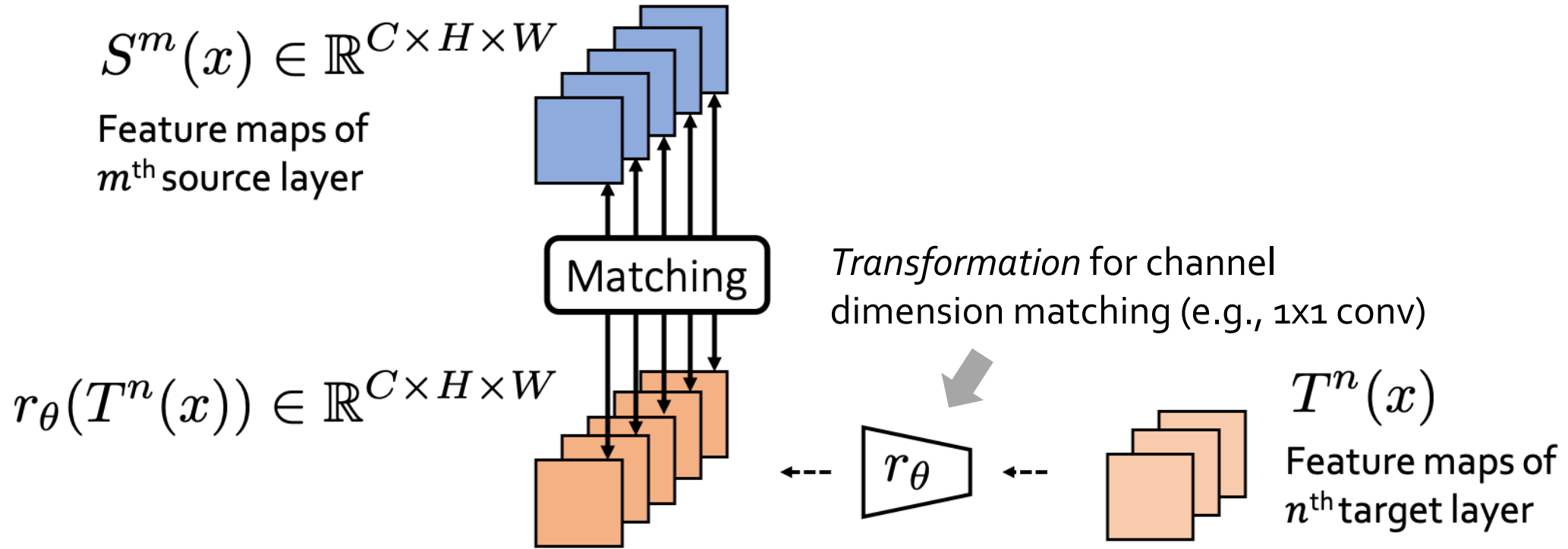
*What* to transfer

- A meta-network  $f$  decides useful channels to transfer



# L2T-ww: Learning *What* to Transfer

- Transfer by making target features similar to those of source [3]



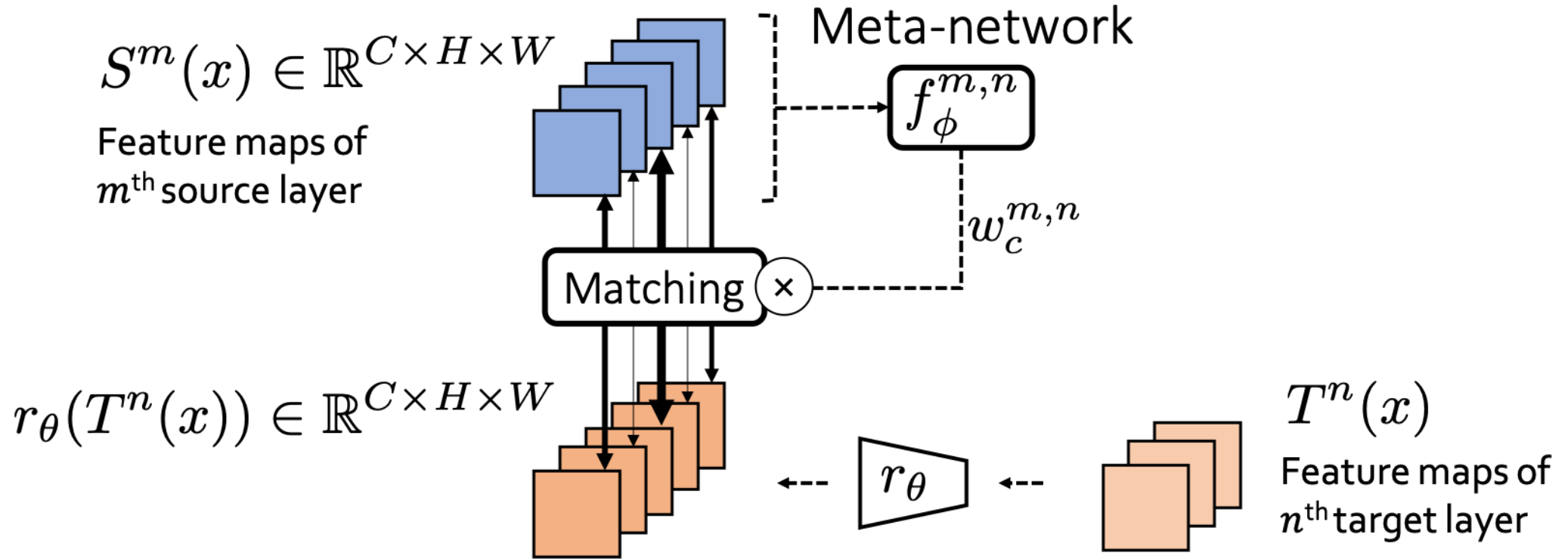
Feature  
Matching

$$\mathcal{L}_{\text{fm}}^{m,n}(\theta|x) = \frac{1}{CHW} \sum_{i,j} (r_\theta(T_\theta^n(x))_{c,i,j} - S^m(x)_{c,i,j})^2$$



# L2T-ww: Learning *What* to Transfer

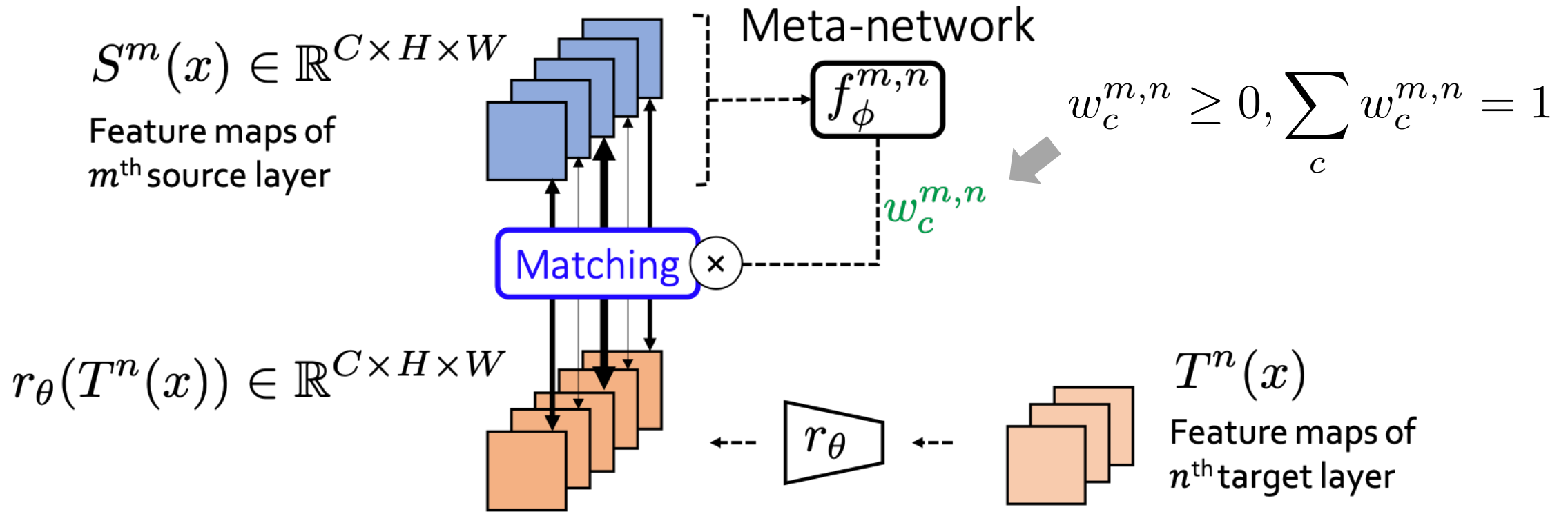
- Learn **what** to transfer



$$\mathcal{L}_{\text{wfm}}^{m,n}(\theta|x, w^{m,n}) = \frac{1}{HW} \sum_c w_c^{m,n} \sum_{i,j} (r_{\theta}(T_{\theta}^n(x))_{c,i,j} - S^m(x)_{c,i,j})^2$$

# L2T-ww: Learning *What* to Transfer

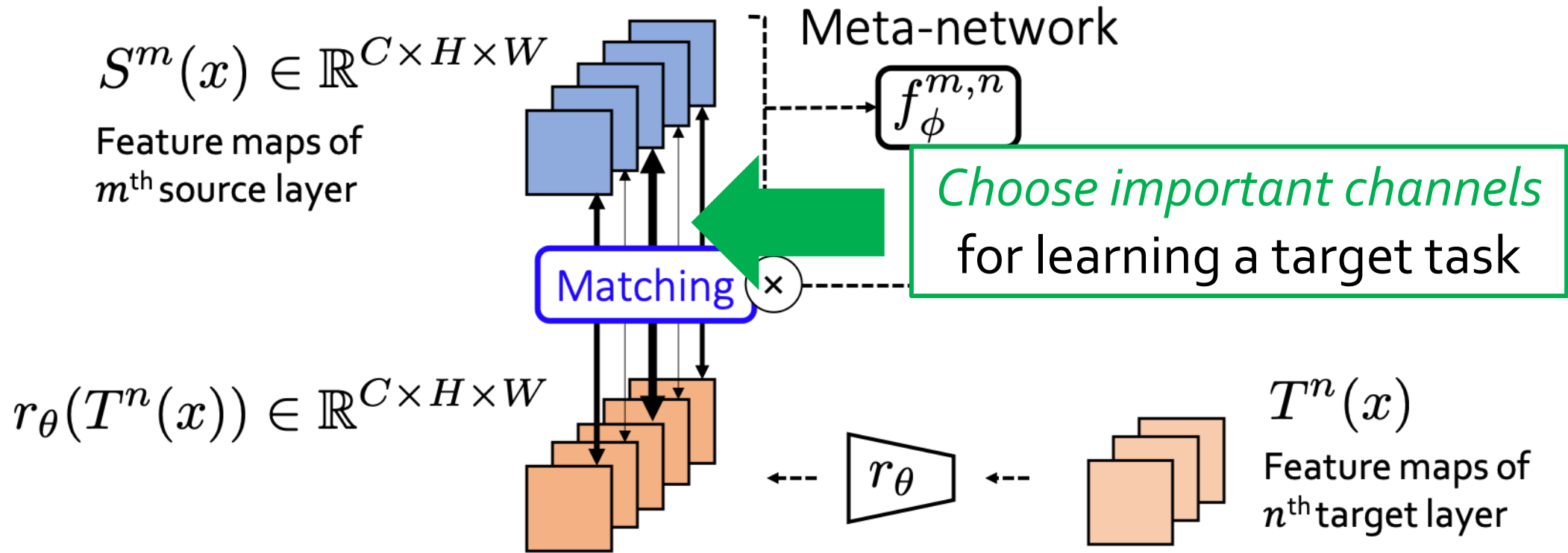
- Learn **what** to transfer



$$\mathcal{L}_{\text{wfm}}^{m,n}(\theta|x, w^{m,n}) = \frac{1}{HW} \sum_c w_c^{m,n} \sum_{i,j} (r_{\theta}(T_{\theta}^n(x))_{c,i,j} - S^m(x)_{c,i,j})^2$$

# L2T-ww: Learning *What* to Transfer

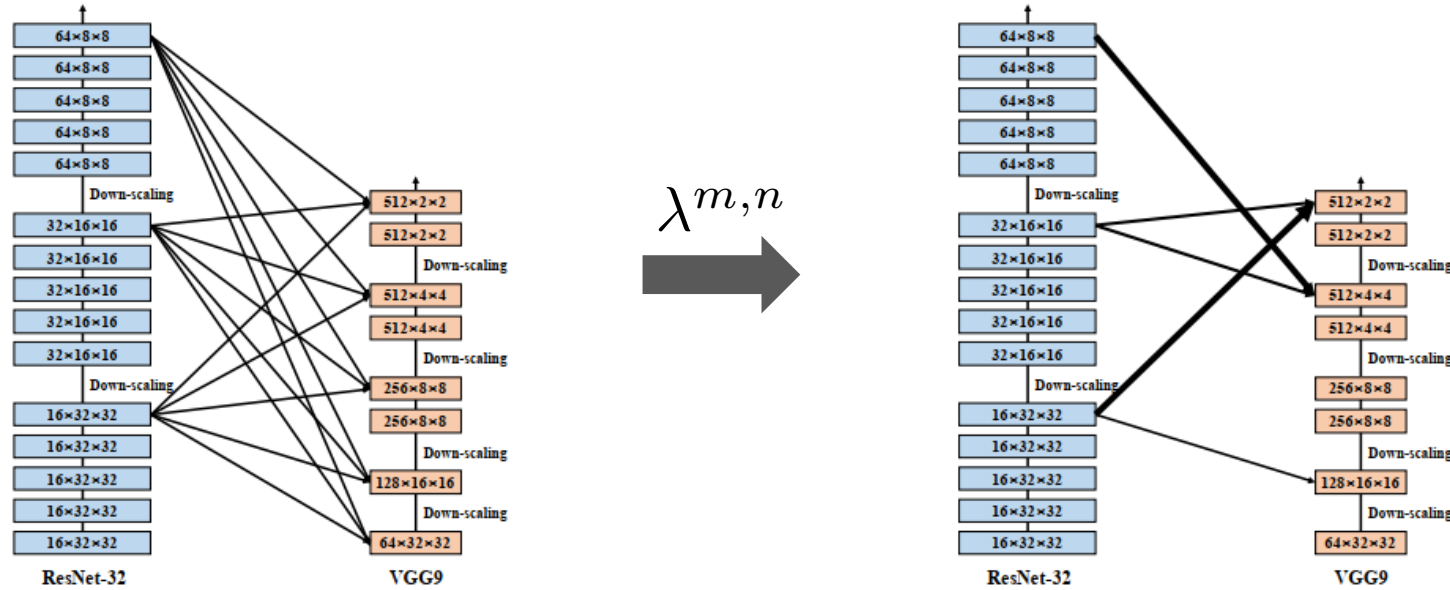
- Learn **what** to transfer



$$\mathcal{L}_{\text{wfm}}^{m,n}(\theta|x, w^{m,n}) = \frac{1}{HW} \sum_c w_c^{m,n} \sum_{i,j} (r_{\theta}(T^n(x))_{c,i,j} - S^m(x)_{c,i,j})^2$$

# L2T-ww: Learning *Where* to Transfer

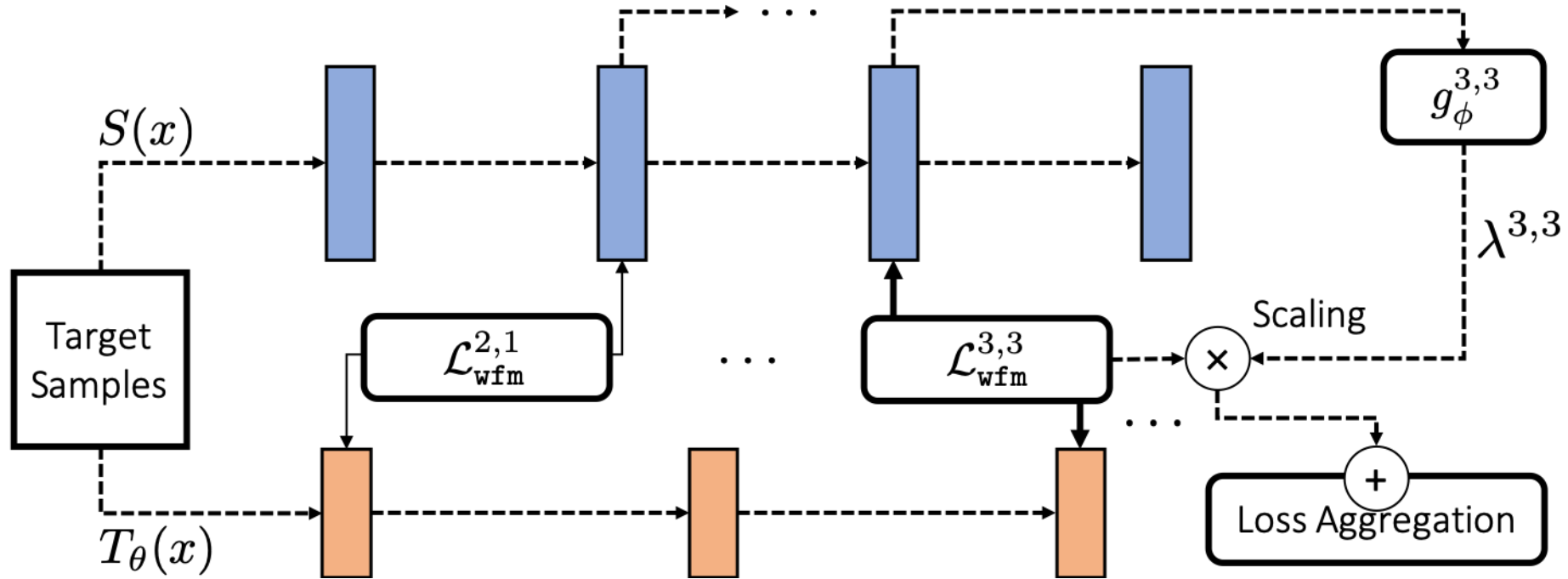
- Learn *where* to transfer



- Meta-networks choose important matching pairs to transfer
  - Given all possible candidate matching pairs  $\mathcal{C}$

# L2T-ww: Learning *Where* to Transfer

- Learn *where* to transfer

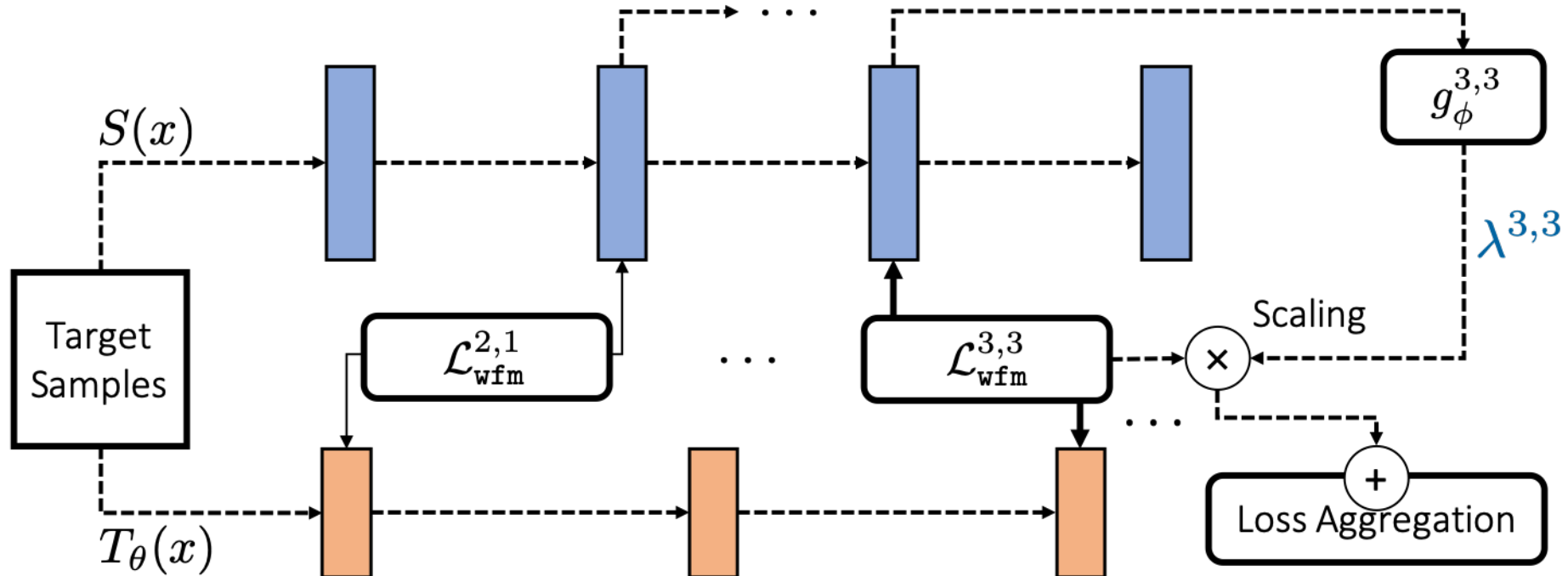


$$\mathcal{L}_{\text{wfm}}(\theta|x, \phi) = \sum_{(m,n) \in \mathcal{C}} \lambda^{m,n} \mathcal{L}_{\text{wfm}}^{m,n}(\theta|x, w^{m,n})$$

# L2T-ww: Learning *Where* to Transfer

7

- Learn *where* to transfer



$$\mathcal{L}_{\text{wfm}}(\theta|x, \phi) = \sum_{(m,n) \in \mathcal{C}} \lambda^{m,n} \mathcal{L}_{\text{wfm}}^{m,n}(\theta|x, w^{m,n})$$

- Choose pairs of feature-matched layers among all the possible pairs

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- A popular bilevel scheme [4,5] for training meta-parameters  $\phi$ :

[4] Colson, B., Marcotte, P., and Savard, G. An overview of bilevel optimization. *Annals of operations research*, 2007.

[5] Franceschi, L., Frasconi, P., Salzo, S., Grazzi, R., and Pontil, M. Bilevel programming for hyperparameter optimization and meta-learning. In *ICML*, 2018.

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- A popular bilevel scheme [4,5] for training meta-parameters  $\phi$ :

1. *Training simulation:* for  $t = 1, \dots, T$ ,

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}_{\text{total}}(\theta_t|x_t, y_t, \phi)$$

[4] Colson, B., Marcotte, P., and Savard, G. An overview of bilevel optimization. *Annals of operations research*, 2007.

[5] Franceschi, L., Frascioni, P., Salzo, S., Grazzi, R., and Pontil, M. Bilevel programming for hyperparameter optimization and meta-learning. In *ICML*, 2018.



# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- A popular bilevel scheme [4,5] for training meta-parameters  $\phi$ :

1. *Training simulation:* for  $t = 1, \dots, T$ ,

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}_{\text{total}}(\theta_t|x_t, y_t, \phi)$$

2. *Evaluation:*

$$\mathcal{L}_{\text{meta}}(\phi) = \mathcal{L}_{\text{org}}(\theta_{T+1}|x_{\text{val}}, y_{\text{val}})$$

[4] Colson, B., Marcotte, P., and Savard, G. An overview of bilevel optimization. *Annals of operations research*, 2007.

[5] Franceschi, L., Frascioni, P., Salzo, S., Grazzi, R., and Pontil, M. Bilevel programming for hyperparameter optimization and meta-learning. In *ICML*, 2018.

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- A popular bilevel scheme [4,5] for training meta-parameters  $\phi$ :

1. *Training simulation*: for  $t = 1, \dots, T$ ,

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}_{\text{total}}(\theta_t|x_t, y_t, \phi)$$

2. *Evaluation*:

$$\mathcal{L}_{\text{meta}}(\phi) = \mathcal{L}_{\text{org}}(\theta_{T+1}|x_{\text{val}}, y_{\text{val}})$$

3. Update  $\phi$  based on  $\nabla_{\phi} \mathcal{L}_{\text{meta}}(\phi)$  using second-order gradients

[4] Colson, B., Marcotte, P., and Savard, G. An overview of bilevel optimization. *Annals of operations research*, 2007.

[5] Franceschi, L., Frascioni, P., Salzo, S., Grazzi, R., and Pontil, M. Bilevel programming for hyperparameter optimization and meta-learning. In *ICML*, 2018.

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- A popular bilevel scheme [4,5] for training meta-parameters  $\phi$ :

1. *Training simulation*: for  $t = 1, \dots, T$ ,

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}_{\text{total}}(\theta_t|x_t, y_t, \phi)$$

2. *Evaluation*:

$$\mathcal{L}_{\text{meta}}(\phi) = \mathcal{L}_{\text{org}}(\theta_{T+1}|x_{\text{val}}, y_{\text{val}})$$

3. Update  $\phi$  based on  $\nabla_{\phi} \mathcal{L}_{\text{meta}}(\phi)$  using second-order gradients

- The transfer loss  $\mathcal{L}_{\text{wfm}}$  acts as a regularization
- A large number of steps  $T$  is required to obtain meaningful gradients
  - But it is time-consuming

[4] Colson, B., Marcotte, P., and Savard, G. An overview of bilevel optimization. *Annals of operations research*, 2007.

[5] Franceschi, L., Frascioni, P., Salzo, S., Grazzi, R., and Pontil, M. Bilevel programming for hyperparameter optimization and meta-learning. In *ICML*, 2018.

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- The proposed bilevel scheme for training meta-parameters  $\phi$ :

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- The proposed bilevel scheme for training meta-parameters  $\phi$ :

1. *Knowledge transfer*: for  $t = 1, \dots, T$ ,

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}_{\text{wfm}}(\theta_t | \mathbf{x}, \phi)$$

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- The proposed bilevel scheme for training meta-parameters  $\phi$ :

1. *Knowledge transfer*: for  $t = 1, \dots, T$ ,

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}_{\text{wfm}}(\theta_t | \mathbf{x}, \phi)$$

2. *One-step adaption*:

$$\theta_{T+2} = \theta_{T+1} - \alpha \nabla_{\theta} \mathcal{L}_{\text{org}}(\theta_{T+1} | \mathbf{x}, \mathbf{y})$$

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- The proposed bilevel scheme for training meta-parameters  $\phi$ :

1. *Knowledge transfer*: for  $t = 1, \dots, T$ ,

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}_{\text{wfm}}(\theta_t | \mathbf{x}, \phi)$$

2. *One-step adaption*:

$$\theta_{T+2} = \theta_{T+1} - \alpha \nabla_{\theta} \mathcal{L}_{\text{org}}(\theta_{T+1} | \mathbf{x}, \mathbf{y})$$

3. *Evaluation*:

$$\mathcal{L}_{\text{meta}}(\phi) = \mathcal{L}_{\text{org}}(\theta_{T+2} | \mathbf{x}, \mathbf{y}).$$

4. Update  $\phi$  based on  $\nabla_{\phi} \mathcal{L}_{\text{meta}}(\phi)$  using second-order gradients

# L2T-ww: Training Meta-Networks

- Total loss for target model:  $\mathcal{L}_{\text{total}}(\theta|x, y, \phi) = \mathcal{L}_{\text{org}}(\theta|x, y) + \beta\mathcal{L}_{\text{wfm}}(\theta|x, \phi)$
- The proposed bilevel scheme for training meta-parameters  $\phi$ :

1. *Knowledge transfer*: for  $t = 1, \dots, T$ ,

$$\theta_{t+1} = \theta_t - \alpha \nabla_{\theta} \mathcal{L}_{\text{wfm}}(\theta_t | \mathbf{x}, \phi)$$

2. *One-step adaption*:

$$\theta_{T+2} = \theta_{T+1} - \alpha \nabla_{\theta} \mathcal{L}_{\text{org}}(\theta_{T+1} | \mathbf{x}, \mathbf{y})$$

3. *Evaluation*:

$$\mathcal{L}_{\text{meta}}(\phi) = \mathcal{L}_{\text{org}}(\theta_{T+2} | \mathbf{x}, \mathbf{y}).$$

4. Update  $\phi$  based on  $\nabla_{\phi} \mathcal{L}_{\text{meta}}(\phi)$  using second-order gradients

- Ours is effective for learning  $\phi$  with a small number of steps  $T$
- Ours learns  $\theta$  and  $\phi$  jointly without separate meta-learning phase



# Experiments

- Learning **what** and **where** to transfer gives consistent improvements
  - Suggested method works well in various tasks and architectures

Source task	TinyImageNet			ImageNet		
Target task	CIFAR-100	STL-10	CUB200	MIT67	Stanford40	Stanford Dogs
Scratch	67.69 $\pm$ 0.22	65.18 $\pm$ 0.91	42.15 $\pm$ 0.75	48.91 $\pm$ 0.53	36.93 $\pm$ 0.68	58.08 $\pm$ 0.26
LwF <sup>[6]</sup>	69.23 $\pm$ 0.09	68.64 $\pm$ 0.58	45.52 $\pm$ 0.66	53.73 $\pm$ 2.14	39.73 $\pm$ 1.63	66.33 $\pm$ 0.45
AT <sup>[1]</sup> (one-to-one)	67.54 $\pm$ 0.40	74.19 $\pm$ 0.22	57.74 $\pm$ 1.17	59.18 $\pm$ 1.57	59.29 $\pm$ 0.91	69.70 $\pm$ 0.08
LwF <sup>[6]</sup> +AT <sup>[1]</sup> (one-to-one)	68.75 $\pm$ 0.09	75.06 $\pm$ 0.57	58.90 $\pm$ 1.32	61.42 $\pm$ 1.68	60.20 $\pm$ 1.34	72.67 $\pm$ 0.26
FM <sup>[3]</sup> (single)	69.40 $\pm$ 0.67	75.00 $\pm$ 0.34	47.60 $\pm$ 0.31	55.15 $\pm$ 0.93	42.93 $\pm$ 1.48	66.05 $\pm$ 0.76
FM <sup>[3]</sup> (one-to-one)	69.97 $\pm$ 0.24	76.38 $\pm$ 1.18	48.93 $\pm$ 0.40	54.88 $\pm$ 1.24	44.50 $\pm$ 0.96	67.25 $\pm$ 0.88
L2T-w (single)	70.27 $\pm$ 0.09	74.35 $\pm$ 0.92	51.95 $\pm$ 0.83	60.41 $\pm$ 0.37	46.25 $\pm$ 3.66	69.16 $\pm$ 0.70
L2T-w (one-to-one)	70.02 $\pm$ 0.19	76.42 $\pm$ 0.52	56.61 $\pm$ 0.20	59.78 $\pm$ 1.90	48.19 $\pm$ 1.42	69.84 $\pm$ 1.45
L2T-ww (all-to-all)	<b>70.96<math>\pm</math>0.61</b>	<b>78.31<math>\pm</math>0.21</b>	<b>65.05<math>\pm</math>1.19</b>	<b>64.85<math>\pm</math>2.75</b>	<b>63.08<math>\pm</math>0.88</b>	<b>78.08<math>\pm</math>0.96</b>

[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR 2017*

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR, 2015*.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

# Experiments

- Learning **what** and **where** to transfer gives consistent improvements
  - Suggested method works well in various tasks and architectures

Source task	TinyImageNet		ImageNet			
	CIFAR-100	STL-10	CUB200	MIT67	Stanford40	Stanford Dogs
Scratch	67.69 $\pm$ 0.22	65.18 $\pm$ 0.91	42.15 $\pm$ 0.75	48.91 $\pm$ 0.53	36.93 $\pm$ 0.68	58.08 $\pm$ 0.26
LwF <sup>[6]</sup>	69.23 $\pm$ 0.09	68.64 $\pm$ 0.58	45.52 $\pm$ 0.66	53.73 $\pm$ 2.14	39.73 $\pm$ 1.63	66.33 $\pm$ 0.45
AT <sup>[1]</sup> (one-to-one)	67.54 $\pm$ 0.40	74.19 $\pm$ 0.22	57.74 $\pm$ 1.17	59.18 $\pm$ 1.57	59.29 $\pm$ 0.91	69.70 $\pm$ 0.08
LwF <sup>[6]</sup> +AT <sup>[1]</sup> (one-to-one)	68.75 $\pm$ 0.09	75.06 $\pm$ 0.57	58.90 $\pm$ 1.32	61.42 $\pm$ 1.68	60.20 $\pm$ 1.34	72.67 $\pm$ 0.26
FM <sup>[3]</sup> (single)	69.40 $\pm$ 0.67	75.00 $\pm$ 0.34	47.60 $\pm$ 0.31	55.15 $\pm$ 0.93	42.93 $\pm$ 1.48	66.05 $\pm$ 0.76
FM <sup>[3]</sup> (one-to-one)	69.97 $\pm$ 0.24	76.38 $\pm$ 1.18	48.93 $\pm$ 0.40	54.88 $\pm$ 1.24	44.50 $\pm$ 0.96	67.25 $\pm$ 0.88
L2T-w (single)	70.27 $\pm$ 0.09	74.35 $\pm$ 0.92	51.95 $\pm$ 0.83	60.41 $\pm$ 0.37	46.25 $\pm$ 3.66	69.16 $\pm$ 0.70
L2T-w (one-to-one)	70.02 $\pm$ 0.19	76.42 $\pm$ 0.52	56.61 $\pm$ 0.20	59.78 $\pm$ 1.90	48.19 $\pm$ 1.42	69.84 $\pm$ 1.45
L2T-ww (all-to-all)	<b>70.96<math>\pm</math>0.61</b>	<b>78.31<math>\pm</math>0.21</b>	<b>65.05<math>\pm</math>1.19</b>	<b>64.85<math>\pm</math>2.75</b>	<b>63.08<math>\pm</math>0.88</b>	<b>78.08<math>\pm</math>0.96</b>

Maximum **+15%**  
relative improvements

- Learning **what** to transfer (channel importance) improves all the baselines

[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR* 2017

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

# Experiments

- Learning **what** and **where** to transfer gives consistent improvements
  - Suggested method works well in various tasks and architectures

Source task	TinyImageNet			ImageNet		
Target task	CIFAR-100	STL-10	CUB200	MIT67	Stanford40	Stanford Dogs
Scratch	67.69 $\pm$ 0.22	65.18 $\pm$ 0.91	42.15 $\pm$ 0.75	48.91 $\pm$ 0.53	36.93 $\pm$ 0.68	58.08 $\pm$ 0.26
LwF <sup>[6]</sup>	69.23 $\pm$ 0.09	68.64 $\pm$ 0.58	45.52 $\pm$ 0.66	53.73 $\pm$ 2.14	39.73 $\pm$ 1.63	66.33 $\pm$ 0.45
AT <sup>[1]</sup> (one-to-one)	67.54 $\pm$ 0.40	74.19 $\pm$ 0.22	57.74 $\pm$ 1.17	59.18 $\pm$ 1.57	59.29 $\pm$ 0.91	69.70 $\pm$ 0.08
LwF <sup>[6]</sup> +AT <sup>[1]</sup> (one-to-one)	68.75 $\pm$ 0.09	75.06 $\pm$ 0.57	58.90 $\pm$ 1.32	61.42 $\pm$ 1.68	60.20 $\pm$ 1.34	72.67 $\pm$ 0.26
FM <sup>[3]</sup> (single)	69.40 $\pm$ 0.67	75.00 $\pm$ 0.34	47.60 $\pm$ 0.31	55.15 $\pm$ 0.93	42.93 $\pm$ 1.48	66.05 $\pm$ 0.76
FM <sup>[3]</sup> (one-to-one)	69.97 $\pm$ 0.24	76.38 $\pm$ 1.18	48.93 $\pm$ 0.40	54.88 $\pm$ 1.24	44.50 $\pm$ 0.96	67.25 $\pm$ 0.88
L2T-w (single)	70.27 $\pm$ 0.09	74.35 $\pm$ 0.92	51.95 $\pm$ 0.83	60.41 $\pm$ 0.37	46.25 $\pm$ 3.66	69.16 $\pm$ 0.70
L2T-w (one-to-one)	70.02 $\pm$ 0.19	76.42 $\pm$ 0.52	56.61 $\pm$ 0.20	59.78 $\pm$ 1.90	48.19 $\pm$ 1.42	69.84 $\pm$ 1.45
L2T-ww (all-to-all)	<b>70.96<math>\pm</math>0.61</b>	<b>78.31<math>\pm</math>0.21</b>	<b>65.05<math>\pm</math>1.19</b>	<b>64.85<math>\pm</math>2.75</b>	<b>63.08<math>\pm</math>0.88</b>	<b>78.08<math>\pm</math>0.96</b>

Maximum **+15%**  
relative improvements

Maximum **+25%**  
relative improvements

- Learning **what** to transfer (channel importance) improves all the baselines
- Learning **where** to transfer (pair importance) gives more improvements on what to transfer

[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR* 2017

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

# Experiments

- Multi-source experiments

First source	TinyImageNet (ResNet32)			
Second source	None	TinyImageNet (ResNet20)	TinyImageNet (ResNet32)	CIFAR-10 (ResNet32)
Scratch	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91
LwF <sup>[6]</sup>	68.64 $\pm$ 0.58	68.56 $\pm$ 2.24	68.05 $\pm$ 2.12	69.51 $\pm$ 0.63
AT <sup>[1]</sup>	74.19 $\pm$ 0.22	73.24 $\pm$ 0.12	73.78 $\pm$ 1.16	73.99 $\pm$ 0.51
LwF <sup>[6]</sup> +AT <sup>[1]</sup>	75.06 $\pm$ 0.57	74.72 $\pm$ 0.46	74.77 $\pm$ 0.30	74.41 $\pm$ 1.51
FM <sup>[3]</sup> (single)	75.00 $\pm$ 0.34	75.83 $\pm$ 0.56	75.99 $\pm$ 0.11	74.60 $\pm$ 0.73
FM <sup>[3]</sup> (one-to-one)	76.38 $\pm$ 1.18	77.45 $\pm$ 0.48	77.69 $\pm$ 0.79	77.15 $\pm$ 0.41
L2T-ww (all-to-all)	<b>78.31<math>\pm</math>0.21</b>	<b>79.35<math>\pm</math>0.41</b>	<b>79.80<math>\pm</math>0.52</b>	<b>80.52<math>\pm</math>0.29</b>

[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR 2017*

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR, 2015*.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

# Experiments

- Multi-source experiments: Different *architectures*

First source	TinyImageNet (ResNet32)			
Second source	None	TinyImageNet (ResNet20)	TinyImageNet (ResNet32)	CIFAR-10 (ResNet32)
Scratch	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91
LwF <sup>[6]</sup>	68.64 $\pm$ 0.58	68.56 $\pm$ 2.24	68.05 $\pm$ 2.12	69.51 $\pm$ 0.63
AT <sup>[1]</sup>	74.19 $\pm$ 0.22	73.24 $\pm$ 0.12	73.78 $\pm$ 1.16	73.99 $\pm$ 0.51
LwF <sup>[6]</sup> +AT <sup>[1]</sup>	75.06 $\pm$ 0.57	74.72 $\pm$ 0.46	74.77 $\pm$ 0.30	74.41 $\pm$ 1.51
FM <sup>[3]</sup> (single)	75.00 $\pm$ 0.34	75.83 $\pm$ 0.56	75.99 $\pm$ 0.11	74.60 $\pm$ 0.73
FM <sup>[3]</sup> (one-to-one)	76.38 $\pm$ 1.18	77.45 $\pm$ 0.48	77.69 $\pm$ 0.79	77.15 $\pm$ 0.41
L2T-ww (all-to-all)	<b>78.31<math>\pm</math>0.21</b>	<b>79.35<math>\pm</math>0.41</b>	<b>79.80<math>\pm</math>0.52</b>	<b>80.52<math>\pm</math>0.29</b>

+2.45% relative improvements

[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR* 2017

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

# Experiments

- Multi-source experiments: Different *architectures*, *initialization*

First source	TinyImageNet (ResNet32)			
Second source	None	TinyImageNet (ResNet20)	TinyImageNet (ResNet32)	CIFAR-10 (ResNet32)
Scratch	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91
LwF <sup>[6]</sup>	68.64 $\pm$ 0.58	68.56 $\pm$ 2.24	68.05 $\pm$ 2.12	69.51 $\pm$ 0.63
AT <sup>[1]</sup>	74.19 $\pm$ 0.22	73.24 $\pm$ 0.12	73.78 $\pm$ 1.16	73.99 $\pm$ 0.51
LwF <sup>[6]</sup> +AT <sup>[1]</sup>	75.06 $\pm$ 0.57	74.72 $\pm$ 0.46	74.77 $\pm$ 0.30	74.41 $\pm$ 1.51
FM <sup>[3]</sup> (single)	75.00 $\pm$ 0.34	75.83 $\pm$ 0.56	75.99 $\pm$ 0.11	74.60 $\pm$ 0.73
FM <sup>[3]</sup> (one-to-one)	76.38 $\pm$ 1.18	77.45 $\pm$ 0.48	77.69 $\pm$ 0.79	77.15 $\pm$ 0.41
L2T-ww (all-to-all)	<b>78.31<math>\pm</math>0.21</b>	<b>79.35<math>\pm</math>0.41</b>	<b>79.80<math>\pm</math>0.52</b>	<b>80.52<math>\pm</math>0.29</b>

+2.45% relative improvements

+2.72% relative improvements

[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR* 2017

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

# Experiments

- Multi-source experiments: Different *architectures*, *initialization* and *datasets*

First source	TinyImageNet (ResNet32)			
Second source	None	TinyImageNet (ResNet20)	TinyImageNet (ResNet32)	CIFAR-10 (ResNet32)
Scratch	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91
LwF <sup>[6]</sup>	68.64 $\pm$ 0.58	68.56 $\pm$ 2.24	68.05 $\pm$ 2.12	69.51 $\pm$ 0.63
AT <sup>[1]</sup>	74.19 $\pm$ 0.22	73.24 $\pm$ 0.12	73.78 $\pm$ 1.16	73.99 $\pm$ 0.51
LwF <sup>[6]</sup> +AT <sup>[1]</sup>	75.06 $\pm$ 0.57	74.72 $\pm$ 0.46	74.77 $\pm$ 0.30	74.41 $\pm$ 1.51
FM <sup>[3]</sup> (single)	75.00 $\pm$ 0.34	75.83 $\pm$ 0.56	75.99 $\pm$ 0.11	74.60 $\pm$ 0.73
FM <sup>[3]</sup> (one-to-one)	76.38 $\pm$ 1.18	77.45 $\pm$ 0.48	77.69 $\pm$ 0.79	77.15 $\pm$ 0.41
L2T-ww (all-to-all)	<b>78.31<math>\pm</math>0.21</b>	<b>79.35<math>\pm</math>0.41</b>	<b>79.80<math>\pm</math>0.52</b>	<b>80.52<math>\pm</math>0.29</b>

+2.45% relative improvements      +2.72% relative improvements      +4.37% relative improvements

[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *ICLR* 2017

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

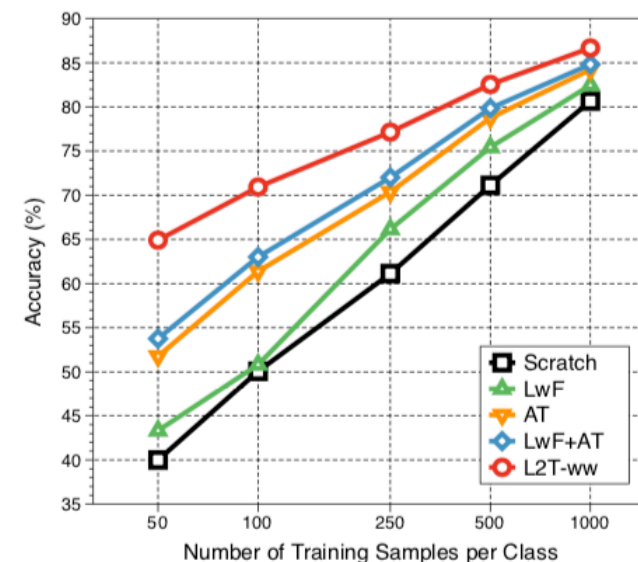
# Experiments

- Multi-source experiments: Different *architectures*, *initialization* and *datasets*

First source	TinyImageNet (ResNet32)			
Second source	None	TinyImageNet (ResNet20)	TinyImageNet (ResNet32)	CIFAR-10 (ResNet32)
Scratch	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91
LwF <sup>[6]</sup>	68.64 $\pm$ 0.58	68.56 $\pm$ 2.24	68.05 $\pm$ 2.12	69.51 $\pm$ 0.63
AT <sup>[1]</sup>	74.19 $\pm$ 0.22	73.24 $\pm$ 0.12	73.78 $\pm$ 1.16	73.99 $\pm$ 0.51
LwF <sup>[6]</sup> +AT <sup>[1]</sup>	75.06 $\pm$ 0.57	74.72 $\pm$ 0.46	74.77 $\pm$ 0.30	74.41 $\pm$ 1.51
FM <sup>[3]</sup> (single)	75.00 $\pm$ 0.34	75.83 $\pm$ 0.56	75.99 $\pm$ 0.11	74.60 $\pm$ 0.73
FM <sup>[3]</sup> (one-to-one)	76.38 $\pm$ 1.18	77.45 $\pm$ 0.48	77.69 $\pm$ 0.79	77.15 $\pm$ 0.41
L2T-ww (all-to-all)	<b>78.31<math>\pm</math>0.21</b>	<b>79.35<math>\pm</math>0.41</b>	<b>79.80<math>\pm</math>0.52</b>	<b>80.52<math>\pm</math>0.29</b>

- Limited data-regime experiments

- Smaller the volume of the target dataset  
→ More relative gain of ours
- Ours efficiently boosts up the performance of a target model



[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.



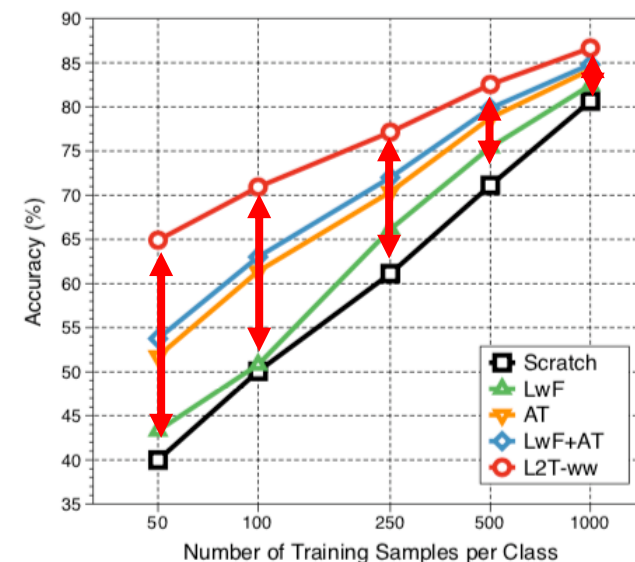
# Experiments

- Multi-source experiments: Different *architectures*, *initialization* and *datasets*

First source	TinyImageNet (ResNet32)			
Second source	None	TinyImageNet (ResNet20)	TinyImageNet (ResNet32)	CIFAR-10 (ResNet32)
Scratch	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91	65.18 $\pm$ 0.91
LwF <sup>[6]</sup>	68.64 $\pm$ 0.58	68.56 $\pm$ 2.24	68.05 $\pm$ 2.12	69.51 $\pm$ 0.63
AT <sup>[1]</sup>	74.19 $\pm$ 0.22	73.24 $\pm$ 0.12	73.78 $\pm$ 1.16	73.99 $\pm$ 0.51
LwF <sup>[6]</sup> +AT <sup>[1]</sup>	75.06 $\pm$ 0.57	74.72 $\pm$ 0.46	74.77 $\pm$ 0.30	74.41 $\pm$ 1.51
FM <sup>[3]</sup> (single)	75.00 $\pm$ 0.34	75.83 $\pm$ 0.56	75.99 $\pm$ 0.11	74.60 $\pm$ 0.73
FM <sup>[3]</sup> (one-to-one)	76.38 $\pm$ 1.18	77.45 $\pm$ 0.48	77.69 $\pm$ 0.79	77.15 $\pm$ 0.41
L2T-ww (all-to-all)	<b>78.31<math>\pm</math>0.21</b>	<b>79.35<math>\pm</math>0.41</b>	<b>79.80<math>\pm</math>0.52</b>	<b>80.52<math>\pm</math>0.29</b>

- Limited data-regime experiments

- Smaller the volume of the target dataset  
→ More relative gain of ours
- Ours efficiently boosts up the performance of a target model



[1] Zagoruyko, S. and Komodakis, N. Paying more attention to attention: Improving the performance of convolutional neural networks via

[3] Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., and Bengio, Y. Fitnets: Hints for thin deep nets. In *ICLR*, 2015.

[6] Li, Z. and Hoiem, D. Learning without forgetting. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

# Conclusion

12

- Meta-learning based transfer method
  - Selective transfer depending on a source and target task relation
  - Effective training scheme that learns meta-networks and target model jointly
  - Applicable between heterogeneous or/and multiple network architectures and tasks

**Poster #186**

**Thursday Jun 13<sup>th</sup> 6:30 – 9:00 PM**

**@ Pacific Ballroom**