# Multicoated Supermasks Enhance Hidden Networks

Yasuyuki Okoshi*, Ángel López García-Arias*, Kazutoshi Hirose, Kota Ando, Kazushi Kawamura, Thiem Van Chu, Masato Motomura, Jaehoon Yu*
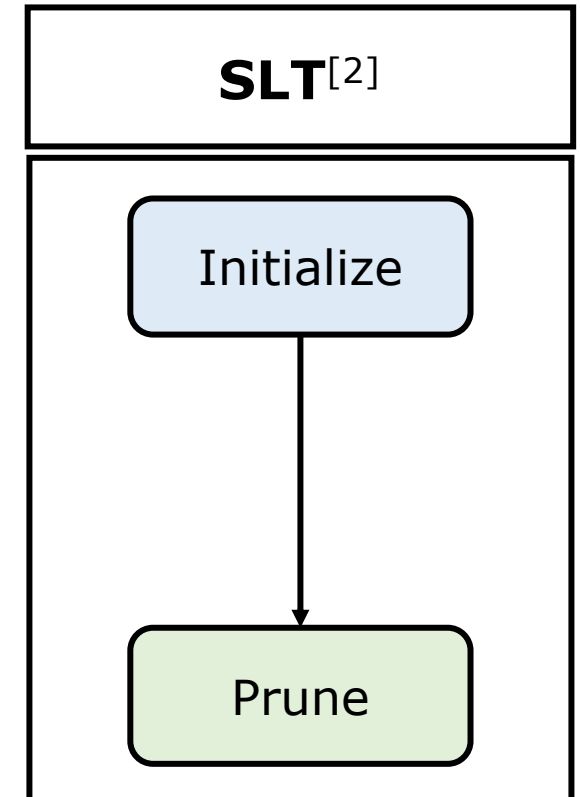
*Equally-Credited Authors

Code available

1

# New Pruning Scheme : Strong Lottery Ticket (SLT)

## SLT is a neural network obtained by learning only connections instead of weights
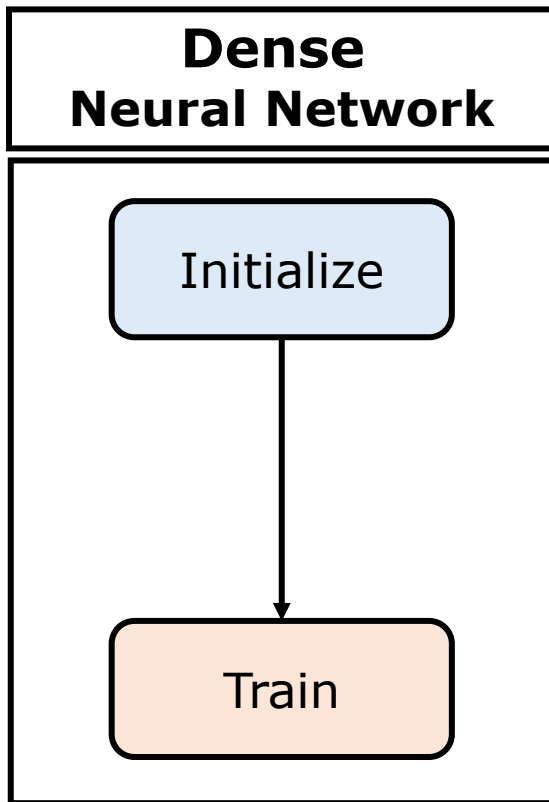


**SLT**[2]

Initialize

Prune

Find sparse NN
without training

[1] J.Frankle, and C.Michael. "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks." ICLR. 2018.
[2] E.Malach, et al. "Proving the lottery ticket hypothesis: Pruning is all you need." ICML. 2020.

# New Pruning Scheme : Strong Lottery Ticket (SLT)

**SLT is a neural network obtained by learning only connections instead of weights**



Dense Neural Network: Initialize → Train
Overparametrized NN

SLT[2]: Initialize → Prune
Find sparse NN without training

[1] J.Frankle, and C.Michael. "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks." ICLR. 2018.
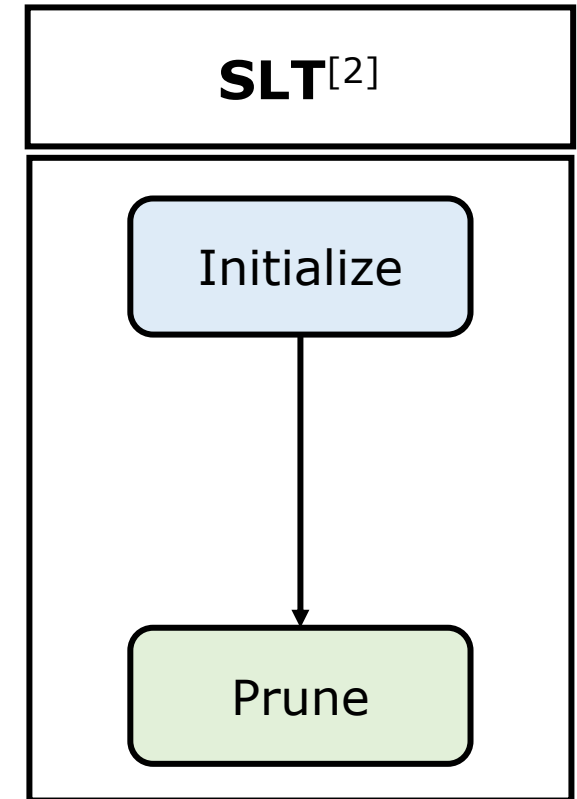[2] E.Malach, et al. "Proving the lottery ticket hypothesis: Pruning is all you need." ICML. 2020.

# New Pruning Scheme : Strong Lottery Ticket (SLT)

## SLT is a neural network obtained by learning only connections instead of weights



**Dense Neural Network**

Initialize → Train

Overparametrized NN

⇨

**Post-training pruning**

Initialize → Train ⇄ Prune

Find sparse NN after training
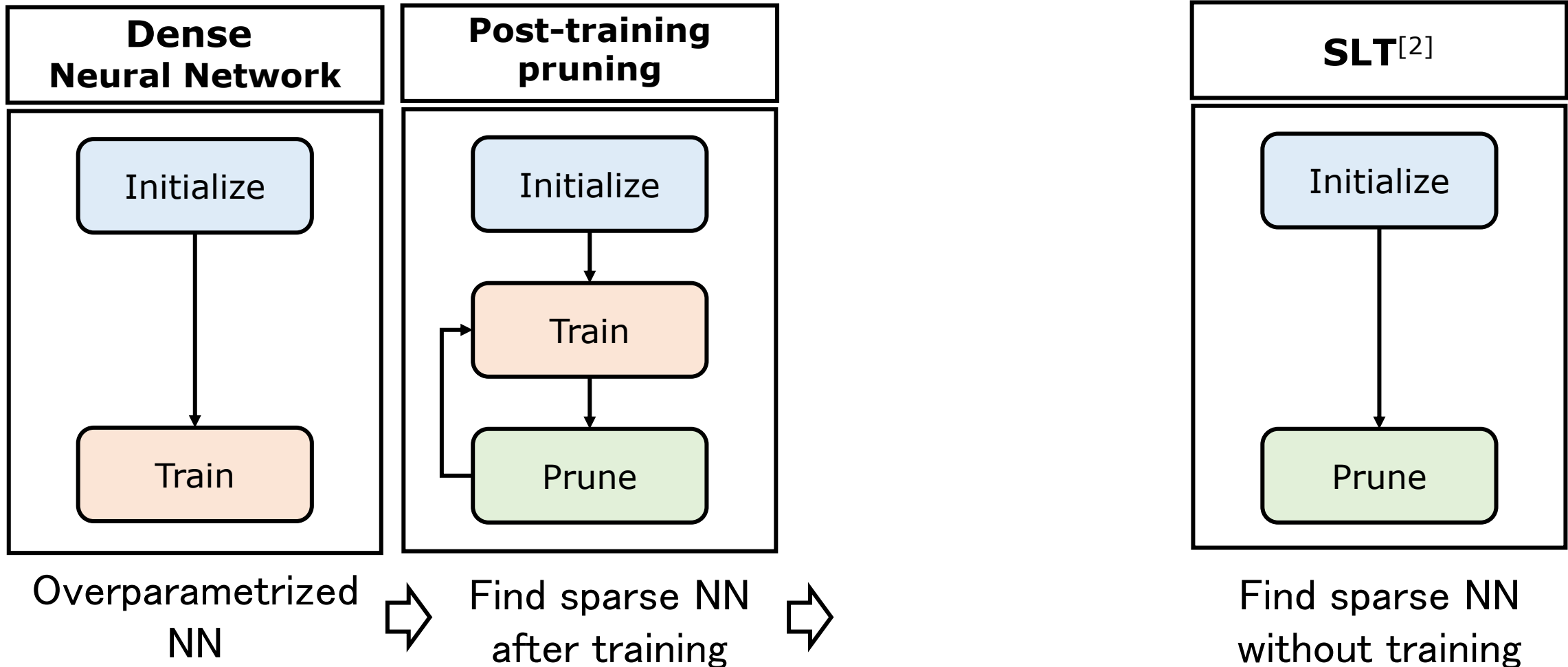
⇨

**SLT[2]**

Initialize → Prune

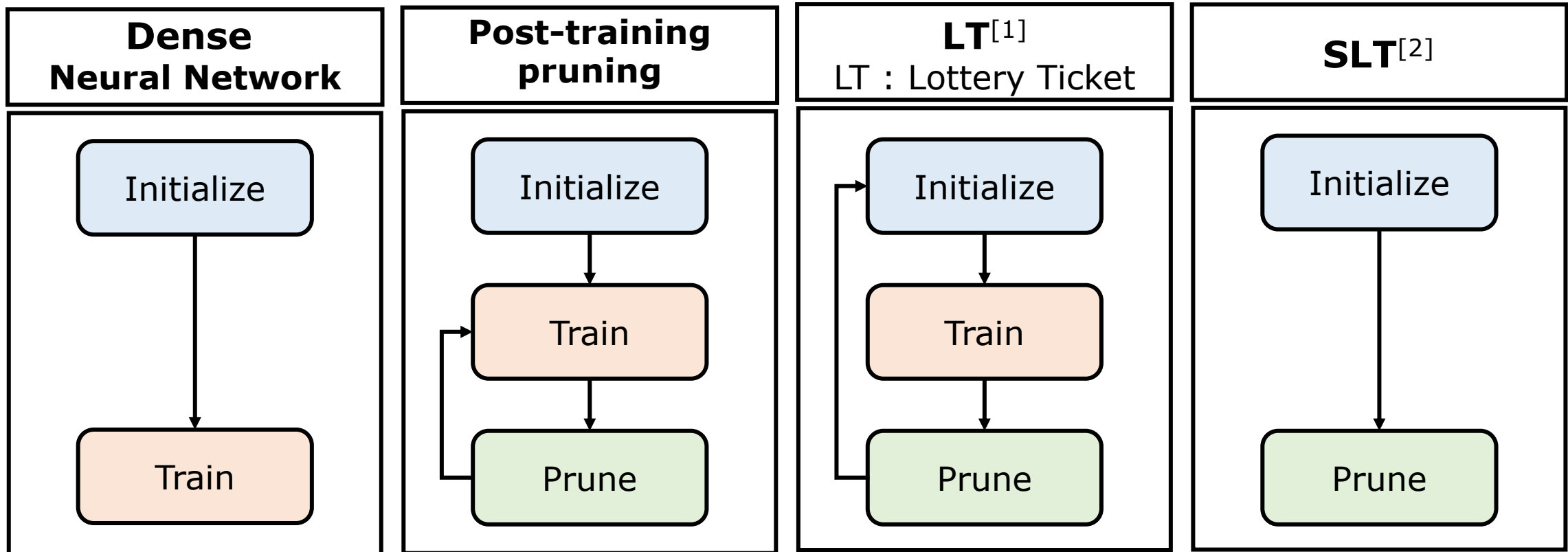Find sparse NN without training

[1] J.Frankle, and C.Michael. "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks." ICLR. 2018.
[2] E.Malach, et al. "Proving the lottery ticket hypothesis: Pruning is all you need." ICML. 2020.

# New Pruning Scheme : Strong Lottery Ticket (SLT)

**SLT is a neural network obtained by learning only connections instead of weights**



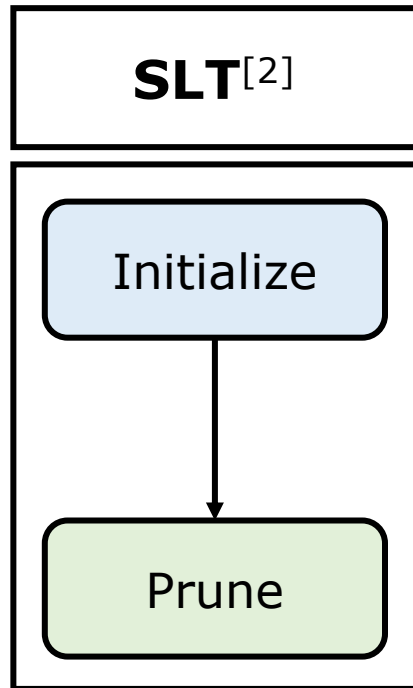| Dense Neural Network | Post-training pruning | LT[1] LT : Lottery Ticket | SLT[2] |
|---|---|---|---|
| Initialize → Train | Initialize → Train → Prune | Initialize → Train → Prune | Initialize → Prune |
| Overparametrized NN | Find sparse NN after training | Find trainable sparse NN | Find sparse NN without training |

[1] J.Frankle, and C.Michael. "The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks." ICLR. 2018.
[2] E.Malach, et al. "Proving the lottery ticket hypothesis: Pruning is all you need." ICML. 2020.
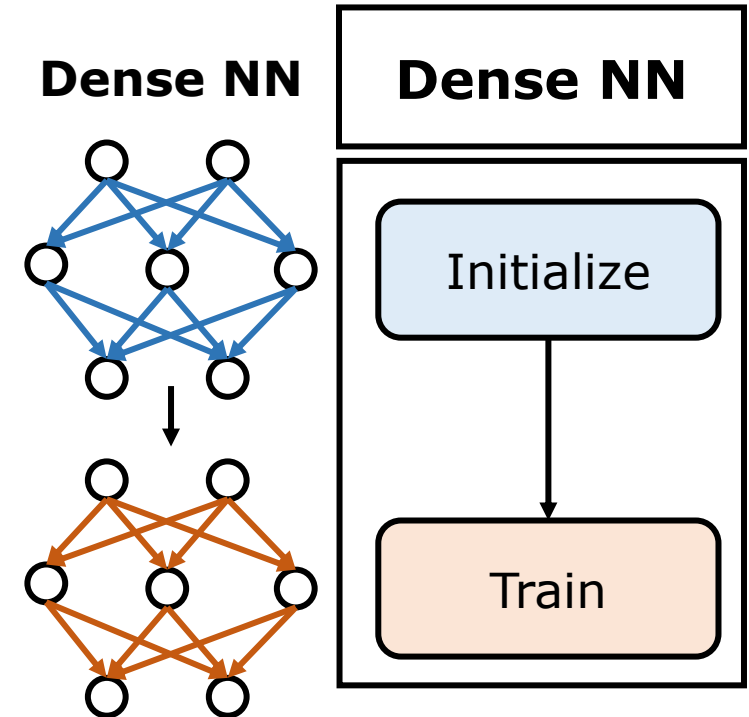
# Comparison of SLT and Dense NN



SLT[2]

Initialize → Prune

**Pros of SLT**
Small model size

**Model Information**

| Connectivity (Supermask) | Learned weights |
| --- | --- |

Dense NN

Dense NN

Initialize → Train

[3] V.Ramanujan, et al. "What's hidden in a randomly weighted neural network?." CVPR. 2020.

International Conference on Machine Learning

*Multicoated Supermasks Enhance Hidden Networks*

6

# Comparison of SLT and Dense NN



SLT[2]

Initialize

Prune

Hidden Networks [3]

Pruning with **edge-popup**

Subnetwork

**Pros of SLT**
Small model size

**Model Information**

| Connectivity (Supermask) | Learned weights |
|---|---|

Dense NN

Dense NN

Initialize

Train

[3] V.Ramanujan, et al. "What's hidden in a randomly weighted neural network?." CVPR. 2020.

*Multicoated Supermasks Enhance Hidden Networks*

# Comparison of SLT and Dense NN



**SLT**[2]

Initialize

↓

Prune

Pruning with **edge-popup**

**Hidden Networks** [3]

Subnetwork

**Pros of SLT**
Small model size

**Model Information**

| Connectivity (Supermask) | Learned weights |

**Cons of SLT**
Lower accuracy

**Search Space**

| Supermask | Weights |

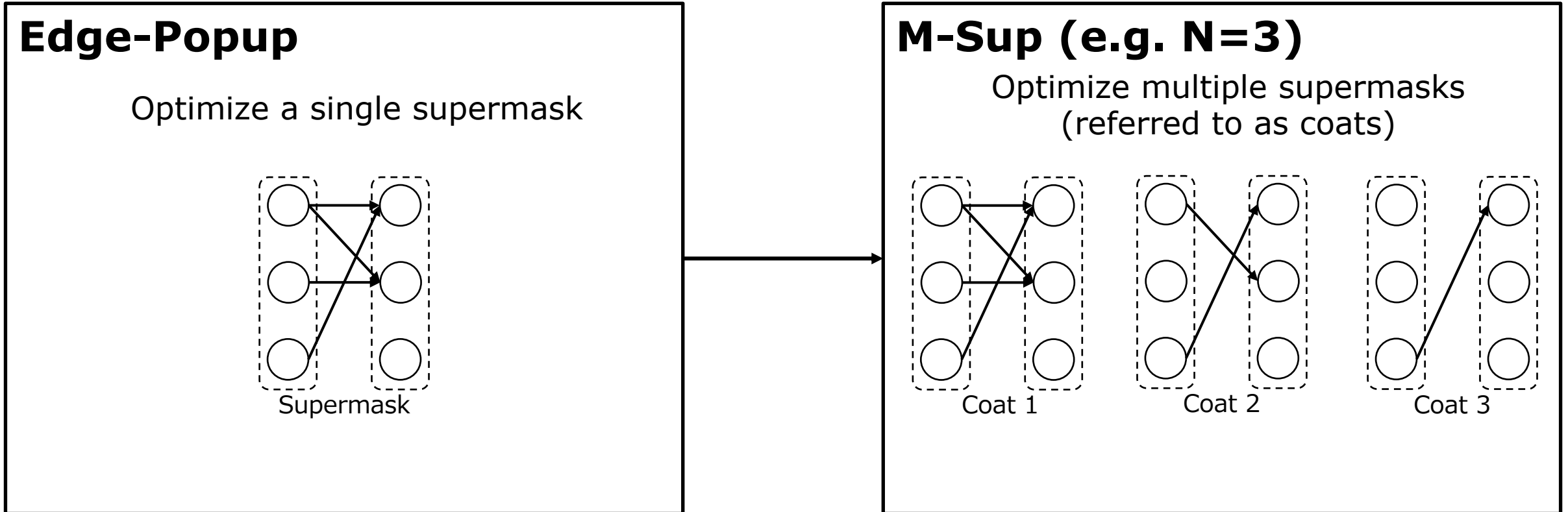**Dense NN**

**Dense NN**

Initialize

↓

Train

⇩

***Multicoated Supermask* extend edge-popup to use multiple supermasks**

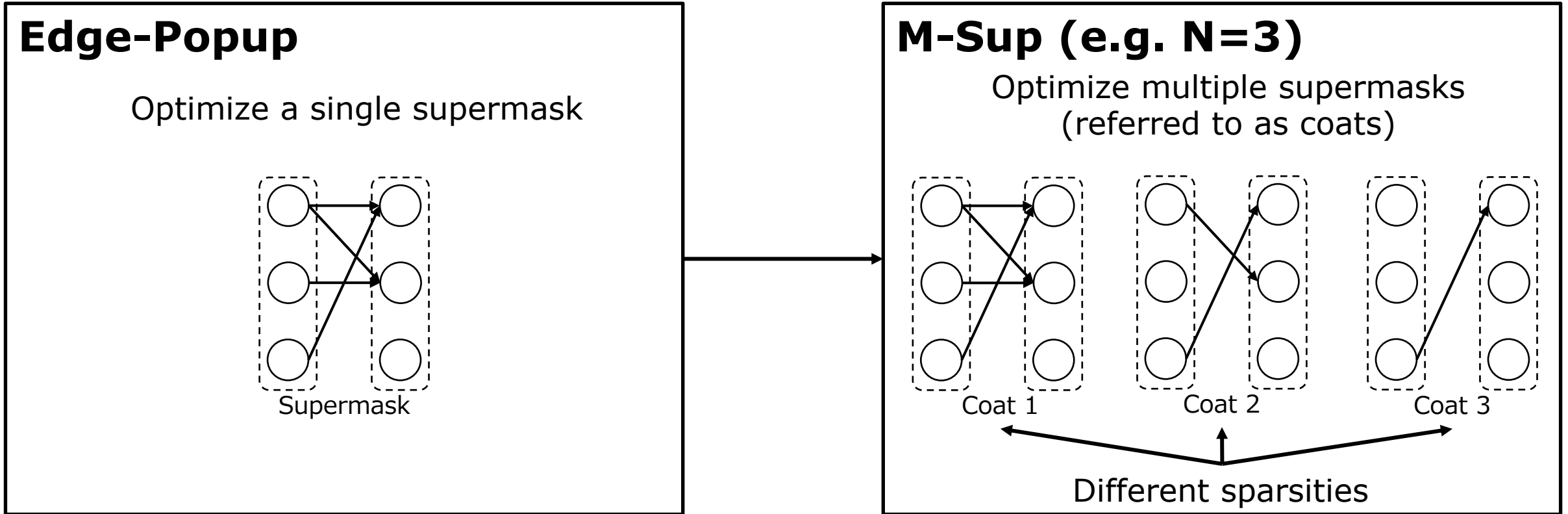[3] V.Ramanujan, et al. "What's hidden in a randomly weighted neural network?." CVPR. 2020.

*Multicoated Supermasks Enhance Hidden Networks*

# Multicoated Supermasks (M-Sup)

**_Multicoated Supermasks_ optimize multiple supermasks simultaneously**

International Conference on Machine Learning          *Multicoated Supermasks Enhance Hidden Networks*

# Multicoated Supermasks (M-Sup)

***Multicoated Supermasks* optimize multiple supermasks simultaneously**



**Edge-Popup**

Optimize a single supermask

Supermask

**M-Sup (e.g. N=3)**

Optimize multiple supermasks
(referred to as coats)

Coat 1     Coat 2     Coat 3

Different sparsities

# Multicoated Supermasks (M-Sup)

**_Multicoated Supermasks_ optimize multiple supermasks simultaneously**

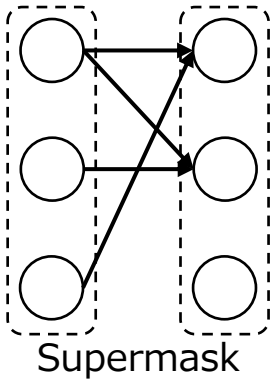**International Conference on Machine Learning**      *Multicoated Supermasks Enhance Hidden Networks*

# Multicoated Supermasks (M-Sup)

**Compute subnetwork from random weights and supermasks**

**Edge-Popup**

Optimize a single supermask



Supermask

**M-Sup (e.g. N=3)**

Optimize multiple supermasks
(referred to as coats)



M-Sup

International Conference on Machine Learning

*Multicoated Supermasks Enhance Hidden Networks*

12

# Multicoated Supermasks (M-Sup)

## Compute subnetwork from random weights and supermasks



**Edge-Popup**

Optimize a single supermask

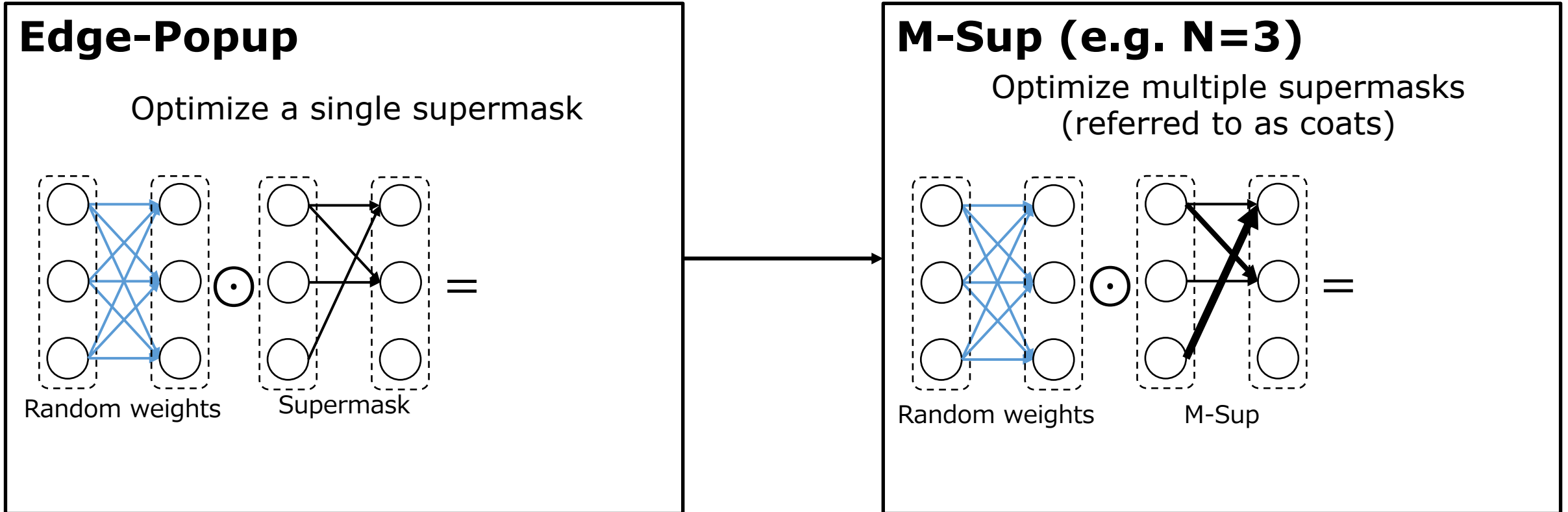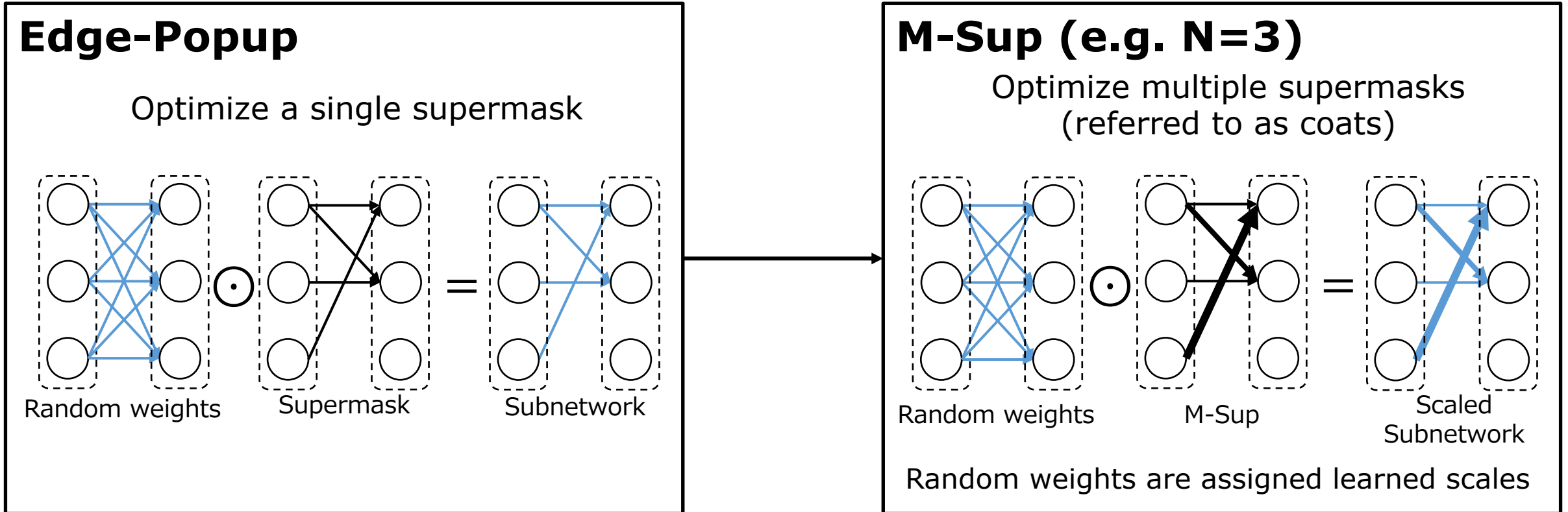Random weights    Supermask

**M-Sup (e.g. N=3)**

Optimize multiple supermasks
(referred to as coats)

Random weights    M-Sup

International Conference on Machine Learning                    *Multicoated Supermasks Enhance Hidden Networks*

# Multicoated Supermasks (M-Sup)

**Compute subnetwork from random weights and supermasks**

## Edge-Popup

Optimize a single supermask



Random weights        Supermask        Subnetwork

## M-Sup (e.g. N=3)

Optimize multiple supermasks
(referred to as coats)



Random weights        M-Sup        Scaled
Subnetwork

Random weights are assigned learned scales

*Multicoated Supermasks Enhance Hidden Networks*

14

# Multicoated Supermasks (M-Sup)

**Compute subnetwork from random weights and supermasks**



**Edge-Popup**

Optimize a single supermask

Random weights  ⊙  Supermask  =  Subnetwork

**M-Sup (e.g. N=3)**

Optimize multiple supermasks
(referred to as coats)

Random weights  ⊙  M-Sup  =  Scaled Subnetwork
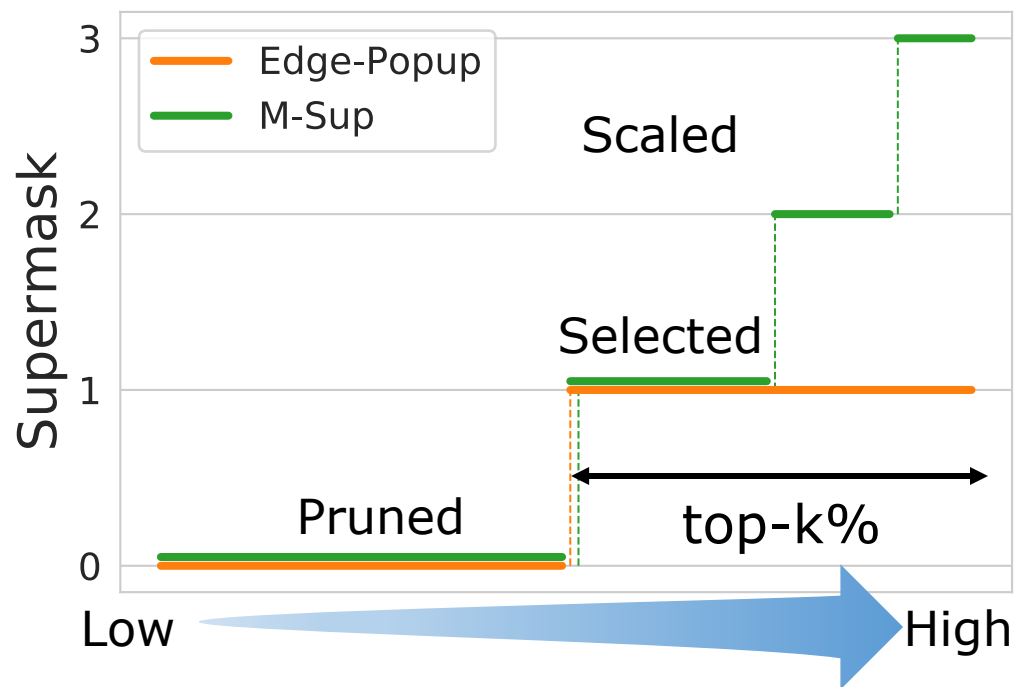
Random weights are assigned learned scales

***Multicoated Supermasks* expand search space with additional coats**
- Training connections of neural network
- Training scales of random weights

*Multicoated Supermasks Enhance Hidden Networks*

15

# Edge-Popup vs M-Sup

## Minimizing the additional cost and additional model size of our method

Score information
used more efficiently

**International Conference on Machine Learning**                    *Multicoated Supermasks Enhance Hidden Networks*

# Edge-Popup vs M-Sup

**Minimizing the additional cost and additional model size of our method**

No additional pruning cost

Score information used more efficiently

# Edge-Popup vs M-Sup

**Minimizing the additional cost and additional model size of our method**
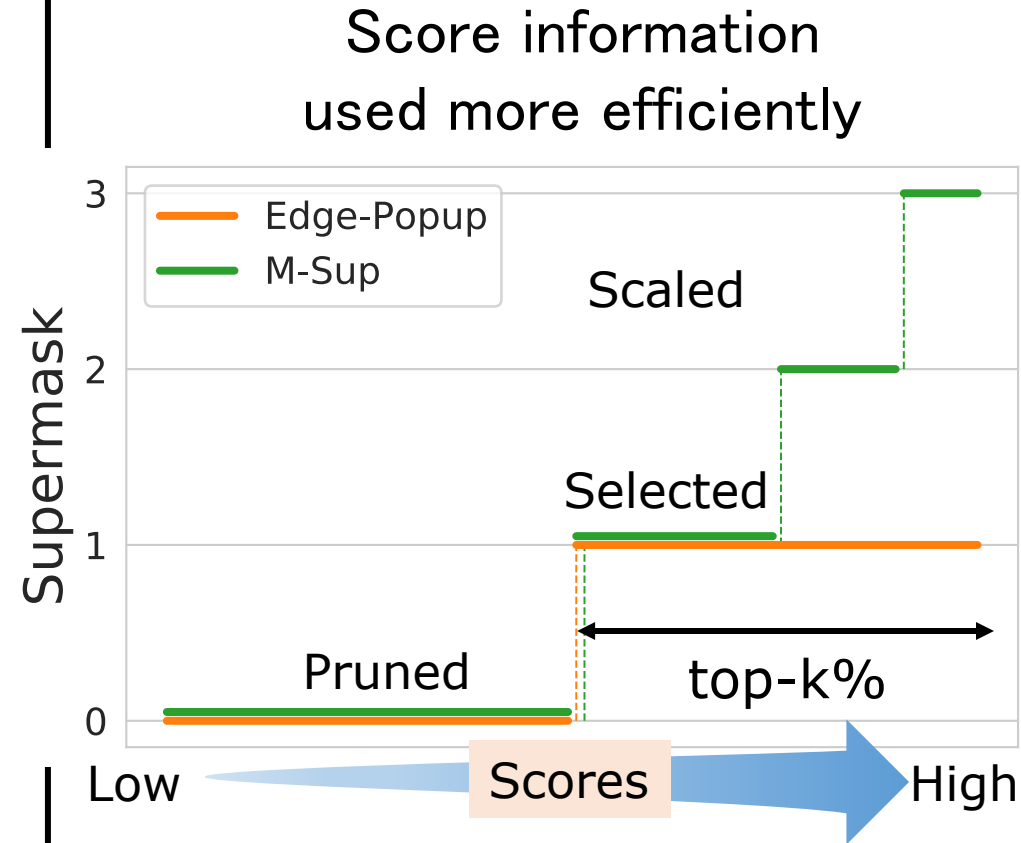
**No additional pruning cost**

Share pruning scores

Score information used more efficiently

*Multicoated Supermasks Enhance Hidden Networks*

# Edge-Popup vs M-Sup

## Minimizing the additional cost and additional model size of our method

**No additional pruning cost**

Share pruning scores

↓

No additional parameters

↓

No additional Training cost

Score information used more efficiently

**International Conference on Machine Learning**                    *Multicoated Supermasks Enhance Hidden Networks*

# Edge-Popup vs M-Sup

**Minimizing the additional cost and additional model size of our method**
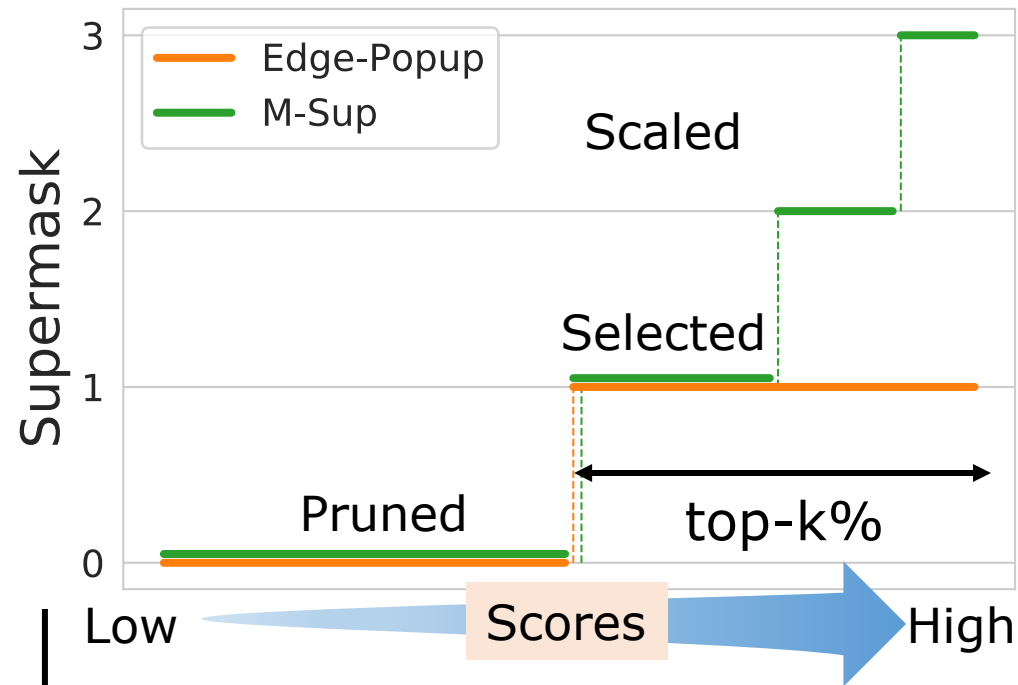
**No additional pruning cost**

Share pruning scores

↓

No additional parameters

↓

No additional Training cost

Score information used more efficiently

- Edge-Popup
- M-Sup

Supermask

3

2

1

0

Scaled

Selected

Pruned

top-k%

Low — Scores — High

**Small increase in model size**

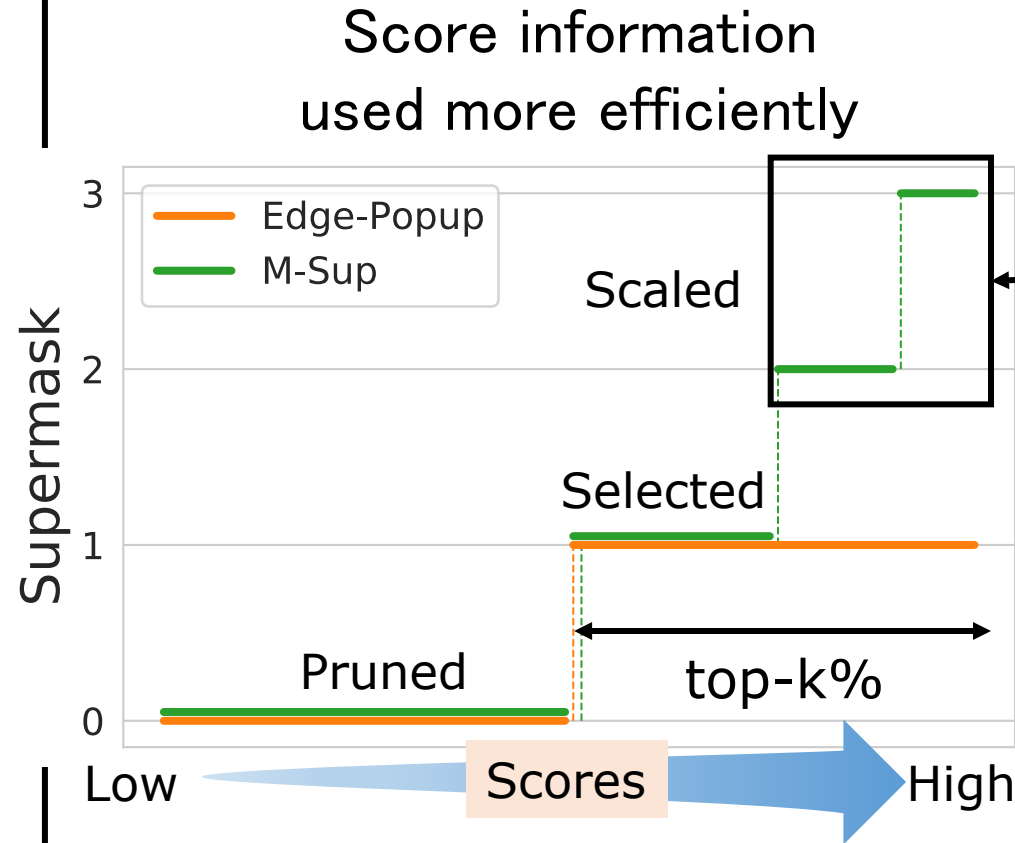*Multicoated Supermasks Enhance Hidden Networks*

# Edge-Popup vs M-Sup

**Minimizing the additional cost and additional model size of our method**

**No additional pruning cost**

Share pruning scores

No additional parameters

No additional Training cost

Score information used more efficiently



**Small increase in model size**

Additional masks are increasingly sparser

# Edge-Popup vs M-Sup

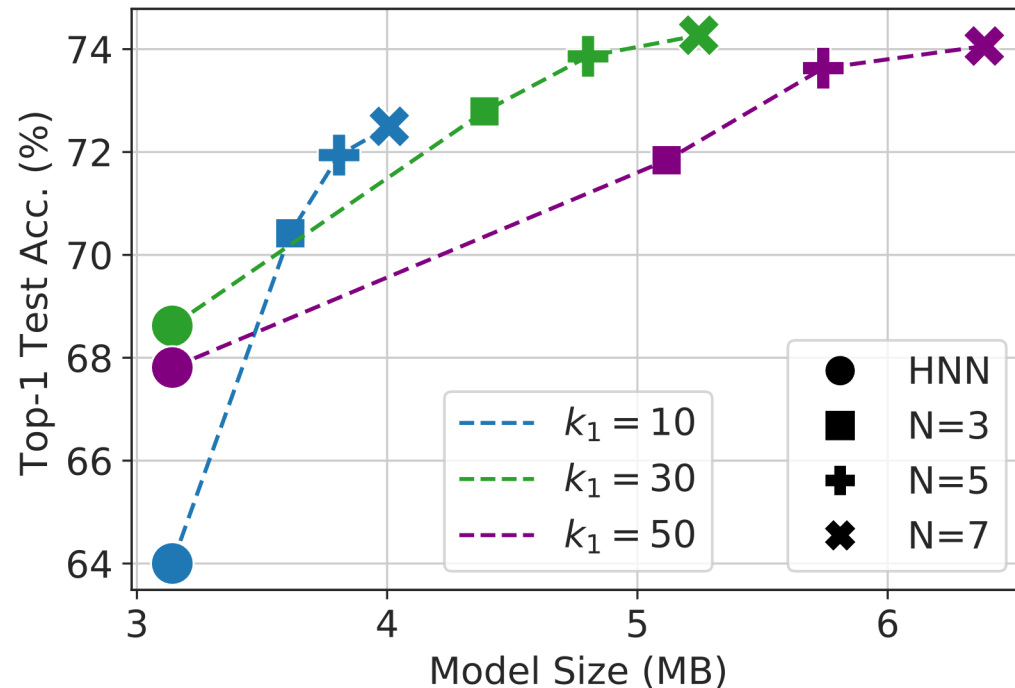## Minimizing the additional cost and additional model size of our method

**No additional pruning cost**

Share pruning scores

↓

No additional parameters

↓

No additional Training cost

Score information used more efficiently



Supermask

Edge-Popup
M-Sup

Scaled

Selected

Pruned

top-k%

Low ← Scores → High

**Small increase in model size**

Additional masks are increasingly sparser

↓

Entropy encoding

Larger mask →

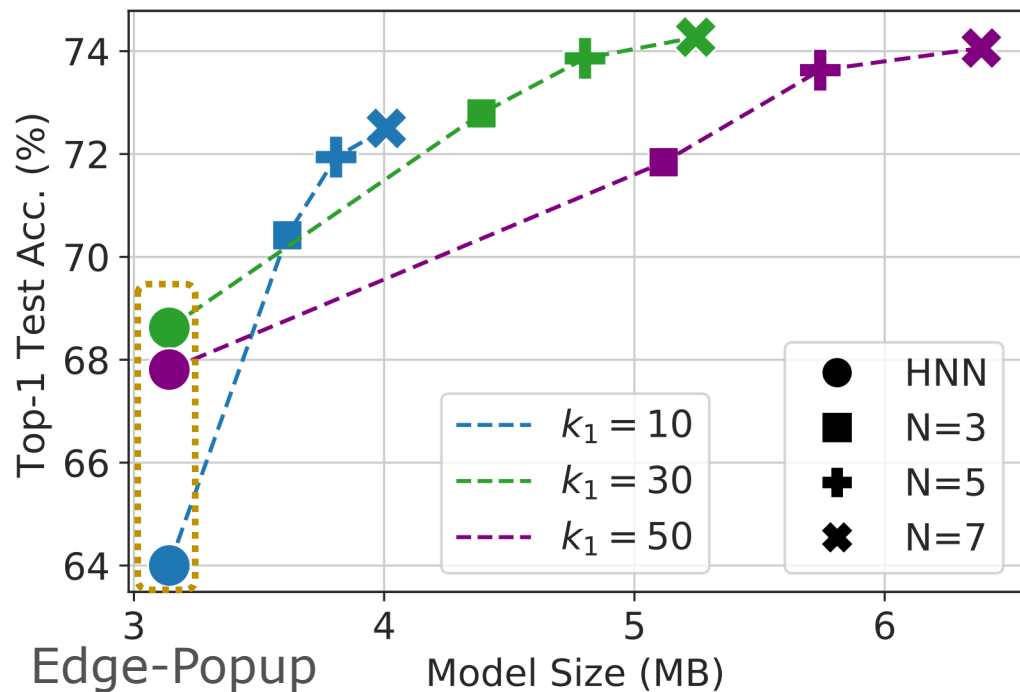| Mask | 0 | 1 | 2 | 3 |
|------|---|----|-----|-----|
| Code | 0 | 10 | 110 | 111 |

Longer Length →

# Comparison – Model Size VS. Accuracy (ImageNet)

**M-Sup achieve competitive results on ImageNet**

- Comparison of #Coats
  - ResNet-50

- Comparison of Model Size

International Conference on Machine Learning          *Multicoated Supermasks Enhance Hidden Networks*
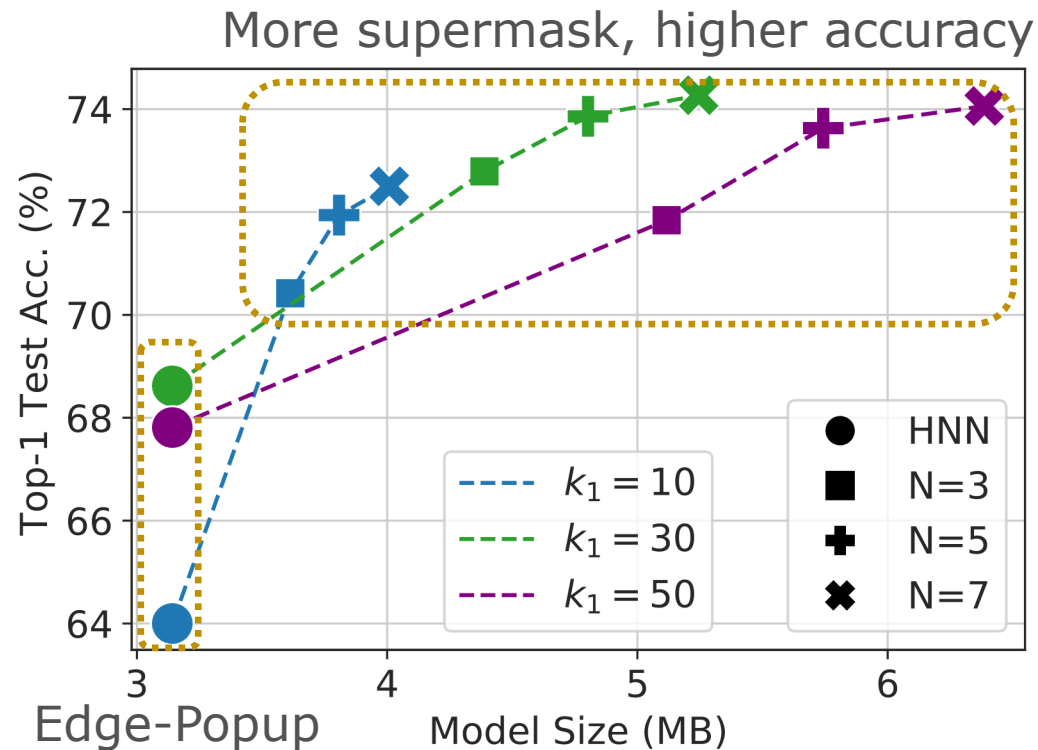
# Comparison – Model Size VS. Accuracy (ImageNet)

**M-Sup achieve competitive results on ImageNet**

- Comparison of #Coats
  - ResNet-50
- Comparison of Model Size

*Multicoated Supermasks Enhance Hidden Networks*

# Comparison – Model Size VS. Accuracy (ImageNet)

**M-Sup achieve competitive results on ImageNet**

- Comparison of #Coats
  - ResNet-50

- Comparison of Model Size

**International Conference on Machine Learning**                    *Multicoated Supermasks Enhance Hidden Networks*
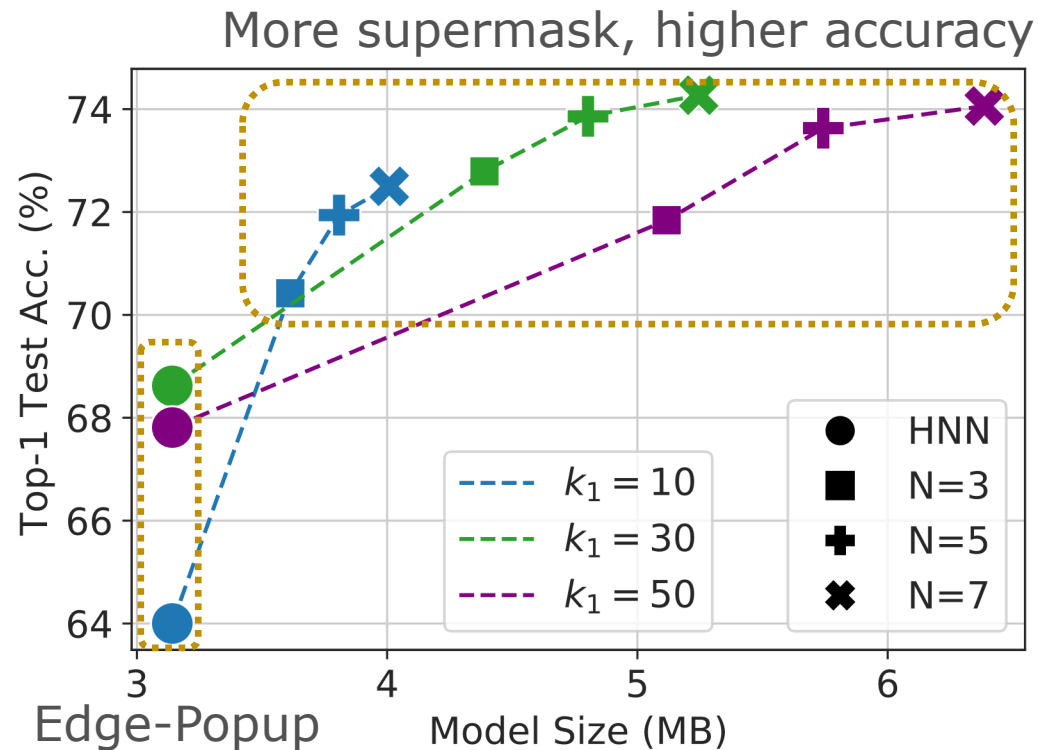
# Comparison – Model Size VS. Accuracy (ImageNet)

**M-Sup achieve competitive results on ImageNet**

- Comparison of #Coats
  - ResNet-50

- Comparison of Model Size

*Multicoated Supermasks Enhance Hidden Networks*

# Conclusion

- First work discussing multiple supermasks for trainable scaling
    - Expanded search space finds a network with high accuracy

# Conclusion

- First work discussing multiple supermasks for trainable scaling
  - Expanded search space finds a network with high accuracy

- Multicoated supermasks achieve
  - **+5% Accuracy** on ImageNet w.r.t. Edge-Popup (ResNet-50)
  - **10x Smaller** size than dense model

**International Conference on Machine Learning**

*Multicoated Supermasks Enhance Hidden Networks*

# Conclusion

- First work discussing multiple supermasks for trainable scaling
  - Expanded search space finds a network with high accuracy

- Multicoated supermasks achieve
  - **+5% Accuracy** on ImageNet w.r.t. Edge-Popup (ResNet-50)
  - **10x Smaller** size than dense model

- The Combination of pruning, quantization, and random weights achieves accurate, highly compressed models