# Neural Similarity Learning

**Weiyang Liu[1]\*, Zhen Liu[2]\*, James M. Rehg[1], Le Song[1]**

**1. Georgia Tech    2. Mila, University of Montreal    \* Equal Contribution**

## Background

- Convolution operator contains two components:
  - Learnable template (Kernel)
  - Similarity measure (inner product)
- Learning (modifying) the shape of kernel:
  - Dilated (atrous) convolution
  - Deformable convolution, Active convolution
- Learning (modifying) the similarity measure:
  - Hyperspherical convolution
  - Decoupled convolution
- Our work aims to generalize the current convolution operator by jointly learning both kernel shape and similarity measure.
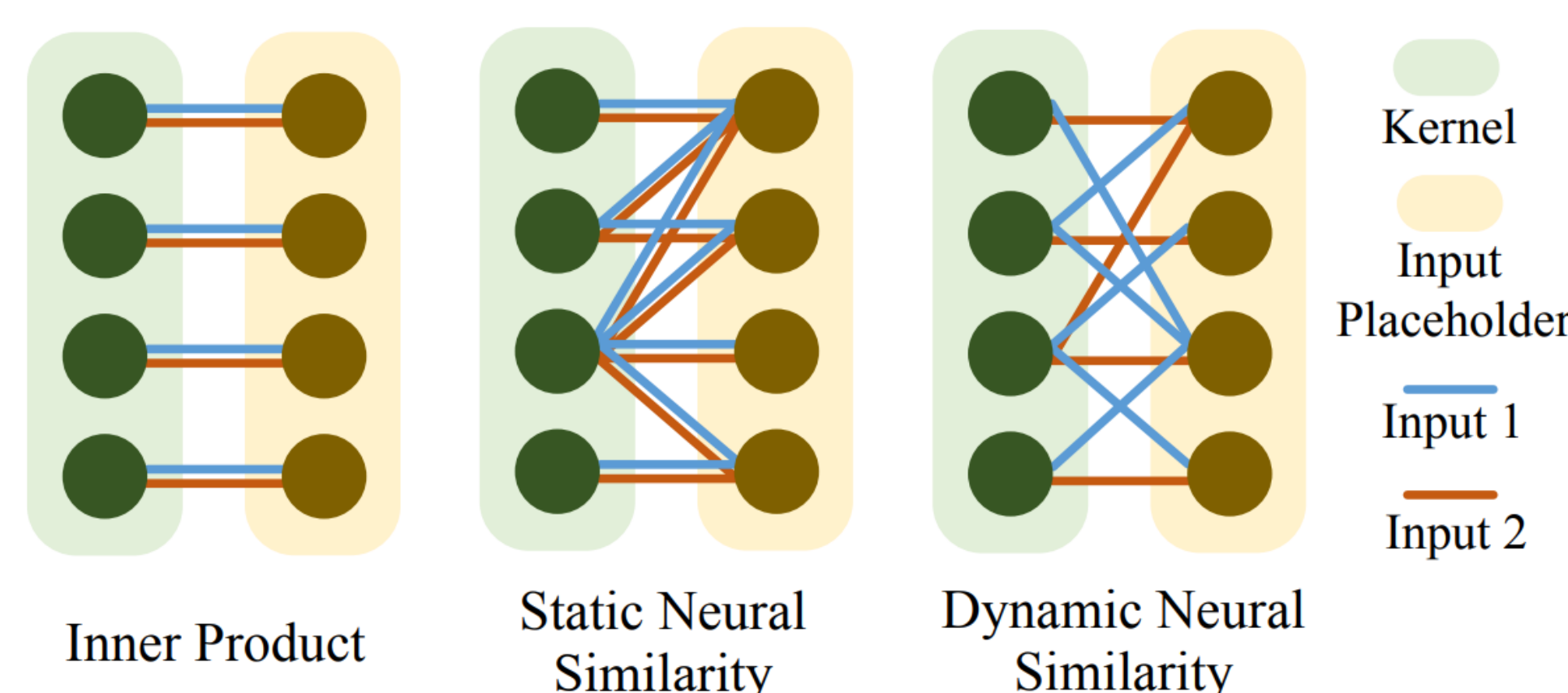
## Motivation

- Hand-designed inner-product based convolution is unlikely to be optimal for every task.
- Optimizing an underdetermined quadratic objective over a matrix $W$ with gradient descent on a factorization of this matrix leads to an implicit regularization for the solution

## Main Contribution

- **Neural similarity** generalizes the inner product via bilinear similarity.
- **Neural similarity network** stacks convolution layers with neural similarity.
- **Static** and **dynamic** learning strategies for the neural similarity.
- Significant performance gain in visual recognition and few-shot learning.

## High-level Comparison with Inner Product



Inner Product    Static Neural Similarity    Dynamic Neural Similarity

Kernel
Input Placeholder
Input 1
Input 2

- A line represents a multiplication operation and a circle denotes an element in a vector. Green color denotes kernel and yellow denotes input.

## Neural Similarity Learning

- Notation:
  $\tilde{W}$ : a convolution kernel with size $C \times H \times V$.
  $W = \{\tilde{W}^F_{1,:,:}, \tilde{W}^F_{2,:,:}, \cdots, \tilde{W}^F_{C,:,:}\} \in \mathbb{R}^{CHV}$ : a flatten kernel.
  $X$ : a flatten input patch.

- Generalizing convolution with bilinear similarity:

$$f_M(W, X) = W^\top M X$$

where $M \in \mathbb{R}^{CHV \times CHV}$ denotes the bilinear similarity matrix.
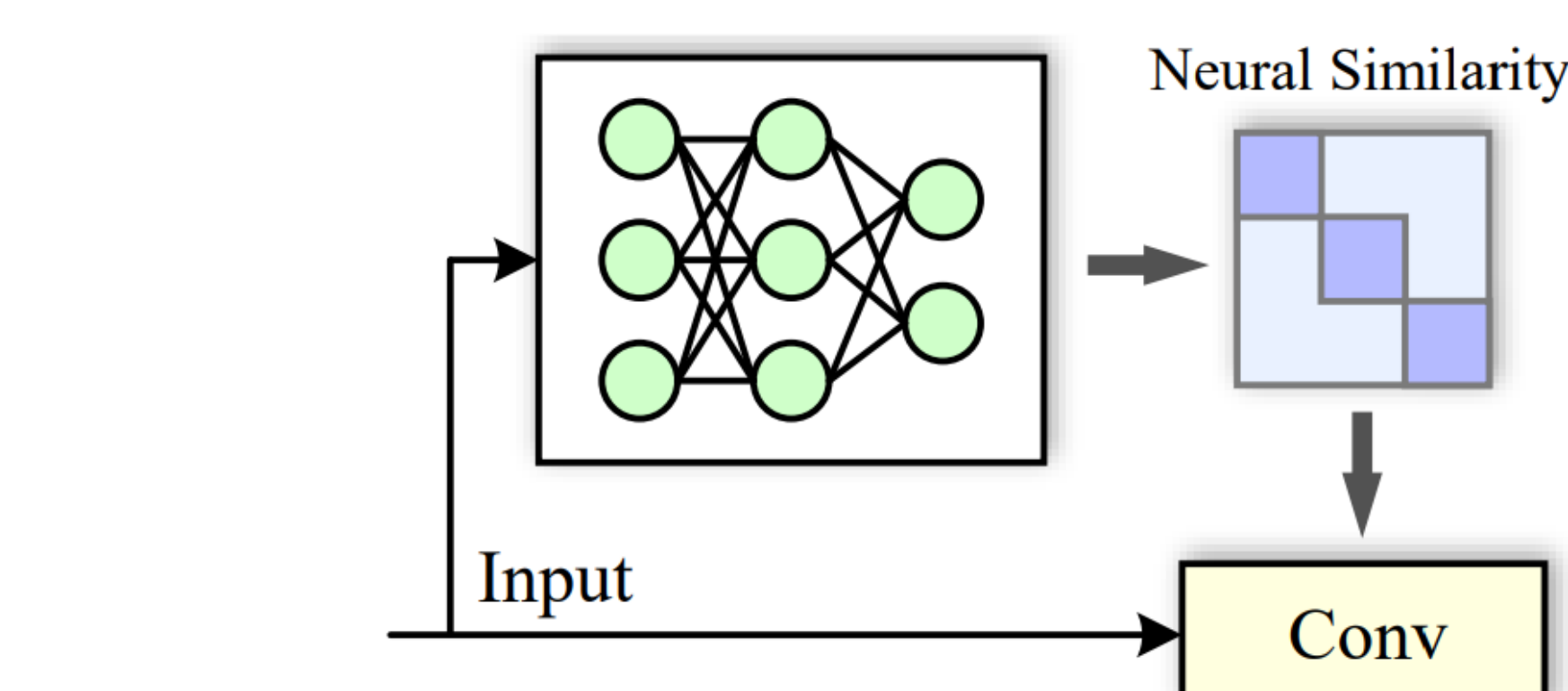
- Constraining M to be block-diagonal:

$$f_M(W, X) = W^\top \begin{bmatrix} M_s & & \\ & \ddots & \\ & & M_s \end{bmatrix} X$$

where $M = \mathrm{diag}(M_s, \cdots, M_s)$ and $M_s$ is of size $HV \times HV$. Note that, hyperspherical convolution becomes a special case of this bilinear formulation if $M$ is a diagonal matrix with diagonal being $\frac{1}{\|W\|\|X\|}$.

### Learning Static Neural Similarity
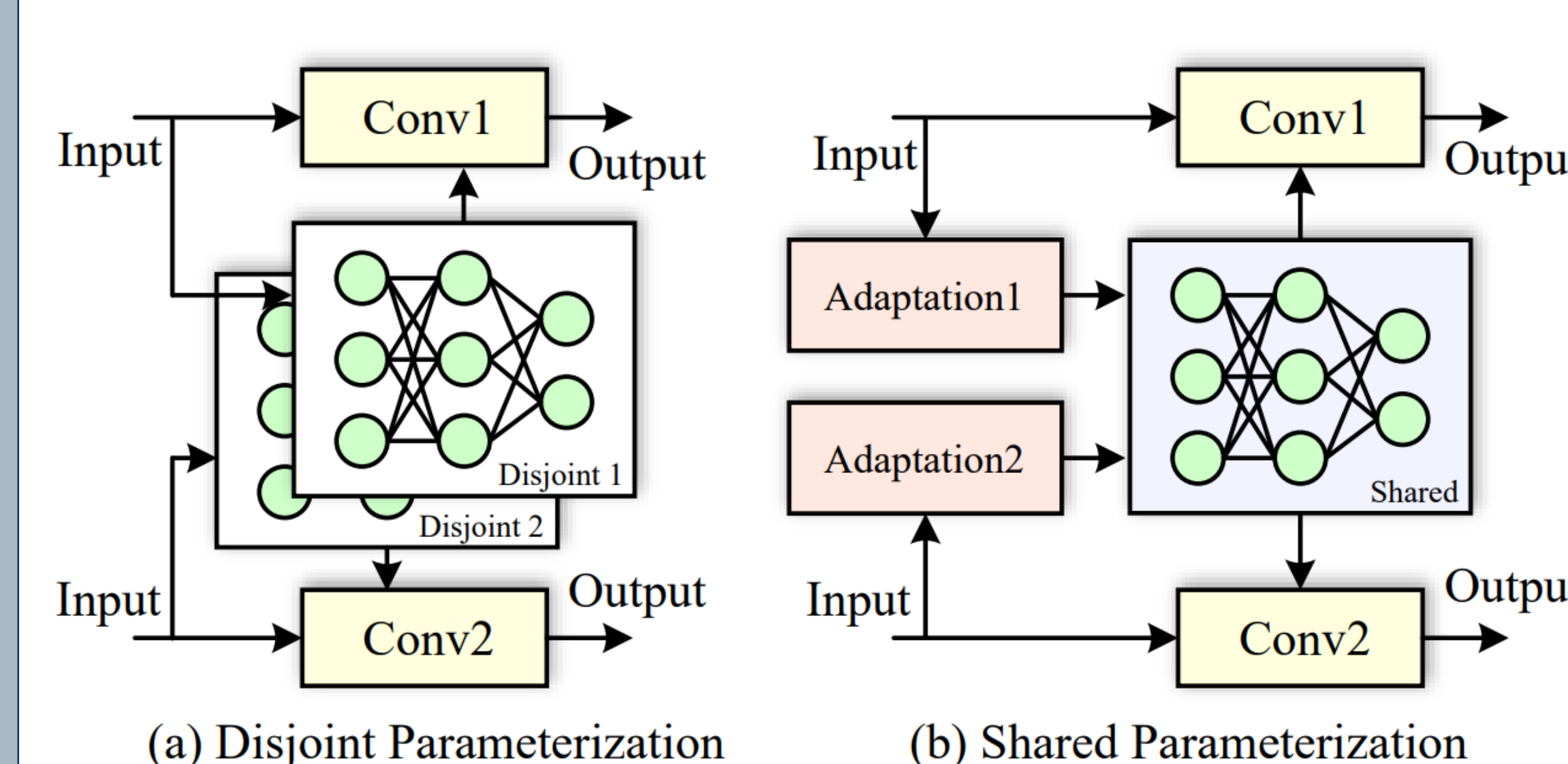


Neural Similarity
Input
Conv

- We learn the matrix $M$ jointly with the convolution kernel via back-propagation.
- Learning static neural similarity can be viewed as a factorized learning of neurons.
- Recent theories suggest that such factorization tends to give minimum nuclear norm solution.

### Learning Dynamic Neural Similarity



Neural Similarity
Input
Conv

- We use a neural network to predict the neural similarity.
- Such neural similarity is dynamic in the sense that it is dependent on the input and dynamically determines the neural similarity during inference.
- It is equivalent to a **dynamic neural network**.

## Disjoint and Shared Parameterization



(a) Disjoint Parameterization    (b) Shared Parameterization

## Learning Both Kernel Shape and Similarity

$$f_M(W, X) = W^\top \begin{bmatrix} DR & & \\ & \ddots & \\ & & DR \end{bmatrix} X$$

$$= W^\top \cdot \begin{bmatrix} D & & \\ & \ddots & \\ & & D \end{bmatrix} \cdot \begin{bmatrix} R & & \\ & \ddots & \\ & & R \end{bmatrix} \cdot X$$

Kernel Shape    Similarity Measure

where $D = \mathrm{diag}(d_1, \cdots, d_{HV})$ and $d_i \in \{0, 1\}, \forall i$.

## Theoretical Insights

- **Implicit regularization induced by NSL**: NSL can be viewed as a form of matrix multiplication where the weight matrix $W$ is factorized as $M^\top W'$.

- Such factorization form not only provides more modeling and regularization flexibility, but it also introduces an **implicit regularization** (in gradient descent).

- Comparison of gradient flow:

  **Standard derivative**

  $$\dot{W}_t = \sum_i X_i (y_i - W_t^\top X_i)^\top = \sum_i X_i (r_t^i)^\top$$

  **NSL derivative**

  $$\dot{W}_t = \dot{M}_t^\top \dot{W}'_t + \dot{M}_t^\top W'_t$$
  $$= M_t^\top M_t \sum_i X_i (r_t^i)^\top + \sum_i X_i (r_t^i)^\top {W'_t}^\top W'_t$$

- Connection to dynamic neural unit (DNU): an isolated DNU is given by a differential equation:

  $$\dot{x}(t) = -\alpha x(t) + f(w, x(t), u), \quad y(t) = g(x(t))$$

- Different from DNU, dynamic NSN does not have the state feedback and self-recurrence.

## Generic Image Recognition

| Method | Error (%) |
|---|---|
| Baseline CNN | 7.78 |
| Dynamic NSN (Shared) | 7.20 |
| Dynamic NSN (Disjoint) | **6.85** |

Error of different parameterization on CIFAR-100

- Shared parameterization has better generalizability than disjoint parameterization.

| Method | CIFAR-10 | CIFAR-100 |
|---|---|---|
| Baseline CNN | 7.78 | 28.95 |
| Baseline CNN++ | 7.29 | 28.70 |
| Static NSN w/ DNS | 7.15 | 28.35 |
| Static NSN w/ UNS | 7.38 | 28.11 |
| Dynamic NSN w/ DNS | 6.85 | **27.81** |
| Dynamic NSN w/ UNS | **6.5** | 28.02 |

Testing error on CIFAR-10 and CIFAR-100

| Method | Top-1 | Top-5 | # params |
|---|---|---|---|
| Baseline CNN | 42.72 | 19.11 | 8.90M |
| Baseline CNN++ | 42.11 | 18.98 | 9.71M |
| Dynamic NSN w/ DNS | **40.61** | **18.04** | 9.61M |

Testing error on ImageNet-2012

- NSL generally yields **better generalization power**.
- NSL has **better parameter efficiency**.
- NSL **does not affect the inference speed** and has the same inference speed as its CNN counterpart.

## Few-shot Image Recognition

| Method | Backbone | 5-shot Accuracy |
|---|---|---|
| Finetuning Baseline | CNN-4 | $49.79 \pm 0.79$ |
| Nearest Neightbor Baseline | CNN-4 | $51.04 \pm 0.65$ |
| MatchingNet | CNN-4 | $55.31 \pm 0.73$ |
| ProtoNet | CNN-4 | $68.20 \pm 0.66$ |
| MAML | CNN-4 | $63.15 \pm 0.91$ |
| RelationNet | CNN-4 | $65.32 \pm 0.70$ |
| Static NSN (ours) | CNN-4 | $65.74 \pm 0.68$ |
| Meta-learned static NSN (ours) | CNN-4 | $66.21 \pm 0.69$ |
| Dynamic NSN (ours) | CNN-4 | $\mathbf{71.26 \pm 0.65}$ |
| Discriminative k-shot | ResNet-34 | $73.90 \pm 0.30$ |
| Tadam | ResNet-12 | $76.7 \pm 0.3$ |
| LEO | ResNet-28 | $\mathbf{77.59 \pm 0.12}$ |
| Dynamic NSN (ours) | CNN-9 | $77.44 \pm 0.63$ |

Few-shot classification on Mini-ImageNet test set

- Meta-learned static NSN is to meta-learn the neural similarity matrix $M$ during training.
- NSL generally has **better generalization power** on few-shot learning.
- Dynamic NSL performs the best and also outperforms the variant where $M$ is meta-learned instead of being learned by a neural network.