

## INTRODUCTION

### Motivation:

- Equivariance is an effective inductive bias
- Convolution networks achieve shift equivariance through parameter-sharing
- Can we discover equivariance from data?

### Contribution:

- Equivariance discovery by learning how to share model parameters from data
- Analysis on the benefit of learning a parameter-sharing scheme
- Empirical results on recovering shift and permutation equivariance

## RUNNING EXAMPLE

### Convolution:

$$y[k] = \sum_j x[j+k]\theta[j]$$

and let  $\mathbf{x} \in \mathbb{R}^3$  and  $\mathbf{k} = [2, 1]$ , we equivalently write convolution as  $\mathbf{y} = \mathbf{K}\mathbf{x}$ , where

$$\mathbf{K} = \begin{bmatrix} 2 & 1 & 0 \\ 0 & 2 & 1 \end{bmatrix}$$

## APPROACH

### Parameterizing Parameter-Sharing:

$$\begin{bmatrix} 2 \\ 1 \\ 0 \\ 0 \\ 2 \\ 1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 0 \\ 2 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

Flatten( $\mathbf{K}$ ) =

$\mathbf{A}$

$\psi$

### Learning parameter-sharing as a bi-level optimization:

$$\min_{\mathbf{A}} \underbrace{\mathcal{L}(\mathbf{A}\psi^*(\mathbf{A}), \mathcal{V})}_{\text{upper-level task}} \text{ s.t. } \psi^*(\mathbf{A}) = \arg \min_{\psi} \underbrace{\mathcal{L}(\mathbf{A}\psi, \mathcal{T})}_{\text{lower-level task}}$$

$$\mathbf{A} \in \{0, 1\}^{K \times K}, \sum_j \mathbf{A}_{ij} = 1 \quad \forall i.$$

### Practical Considerations:

- Relax  $\mathbf{A} \in [0, 1]^{K \times K}$  and optimize using gradient based methods
- Penalty term 1: Entropy  $H(\mathbf{A})$  encourages  $\mathbf{A}_{ij}$  to be closer to 0 or 1
- Penalty term 2: Nuclear norm  $\|\mathbf{A}\|_*$  encourages  $\mathbf{A}$  to be low-rank

### Analysis on Gaussian Data:

$$\mathbf{y} \sim \mathcal{N}(\mathbf{A}_{\text{gt}}\psi_{\text{gt}}, \sigma^2 \mathbf{I}). \quad (1)$$

- *i.i.d.* Gaussian with shared means across dimensions
- $\theta_{\text{gt}} = \mathbf{A}_{\text{gt}}\psi_{\text{gt}}$ : Ground-truth mean
- $\hat{\theta}_{\text{val}} = \mathbf{A}_{\text{val}}\hat{\psi}$ : Estimated  $\psi$  with  $\mathbf{A}_{\text{val}}$  from our approach
- $\hat{\theta}_{\text{gt}} = \mathbf{A}_{\text{gt}}\hat{\psi}$ : Estimated  $\psi$  with  $\mathbf{A}_{\text{gt}}$
- MSE Gap:  $\text{MSE}(\hat{\theta}_{\text{val}}) - \text{MSE}(\hat{\theta}_{\text{gt}})$

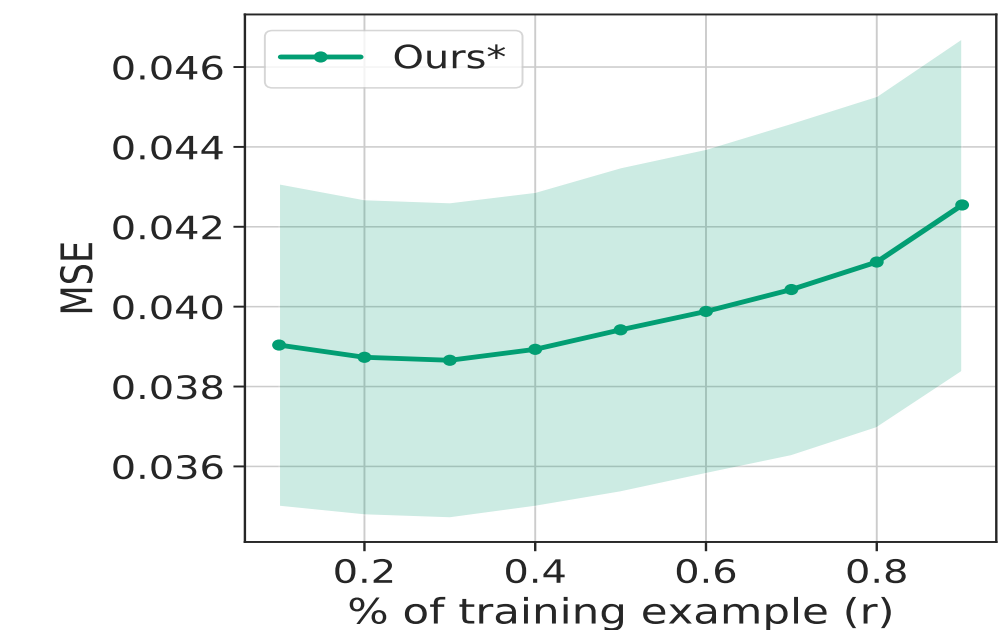
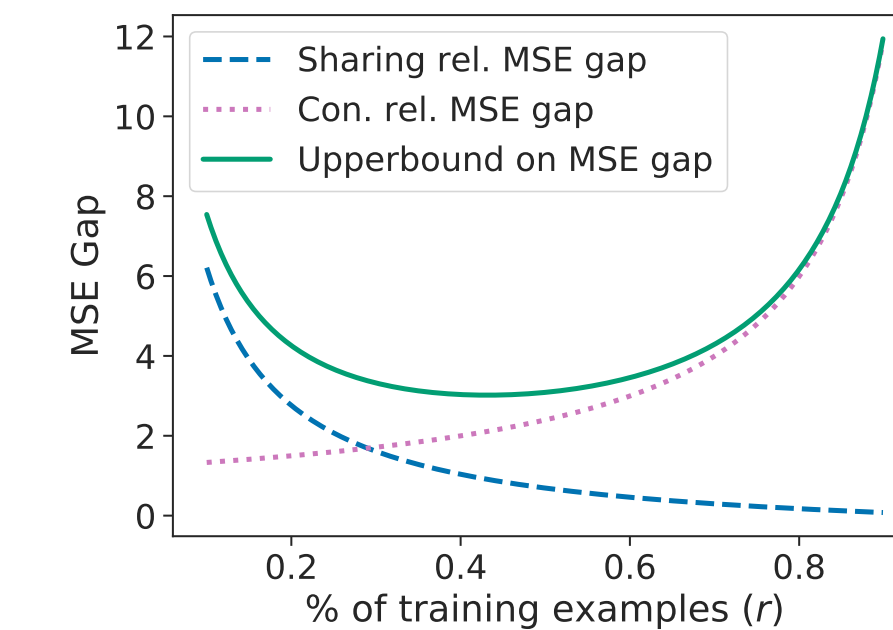
- A selection matrix  $\mathbf{A}$  to represent a sharing scheme
- $\mathbf{A}_{ij} \in \{0, 1\}, \sum_j \mathbf{A}_{ij} = 1$
- $\mathbf{A}$  selects the model parameters from  $\psi$

## RESULTS

**Claim:** Given data following Eq. (1), with probability  $1 - \alpha$  and  $\alpha < \exp \frac{-K}{10}$ , the MSE gap is upper bounded by

$$\sigma^2 \left( \underbrace{\frac{1-r}{r|\mathcal{D}|} (\text{rk}(\mathbf{A}_{\text{gt}}) - 1)}_{\text{sharing rel. MSE gap}} - \underbrace{\frac{40 \ln(\alpha)}{(1-r)|\mathcal{D}|}}_{\text{con. rel. MSE gap}} \right),$$

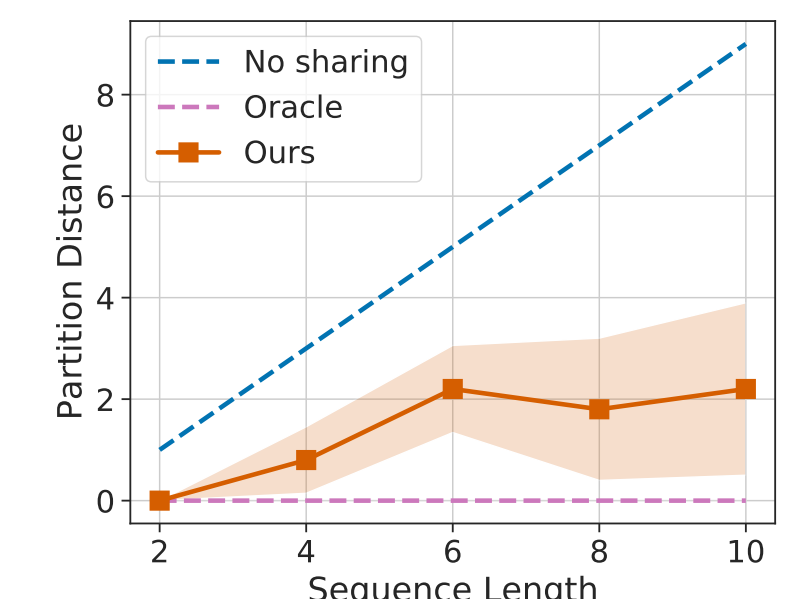
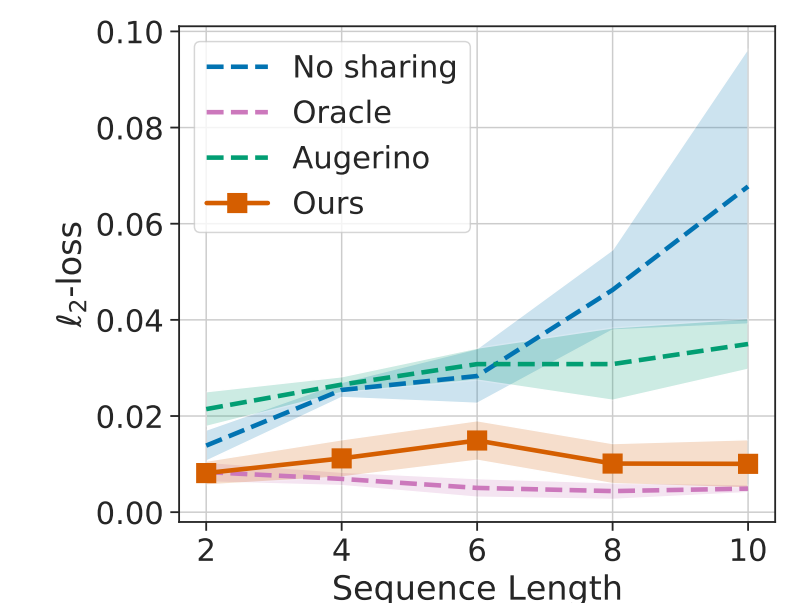
where  $r = \frac{|\mathcal{T}|}{|\mathcal{D}|}$  denotes the ratio between the size of training and overall dataset.



### Experiments (Sum of Numbers):

The task is to regress to the sum of a sequence of numbers provided in text format.

- *E.g.*, given the input (‘one, five’) the model should output 6
- $\ell_2$ -loss between the predicted and ground-truth
- Partition distance between the recovered  $\mathbf{A}_{\text{val}}$  and  $\mathbf{A}_{\text{gt}}$
- Partition distance measures the number of assignments that must be changed for one sharing scheme to be identical to the other.



Please find additional experiments in the paper