# Learning Cyclic Linear Groups in Neural Networks

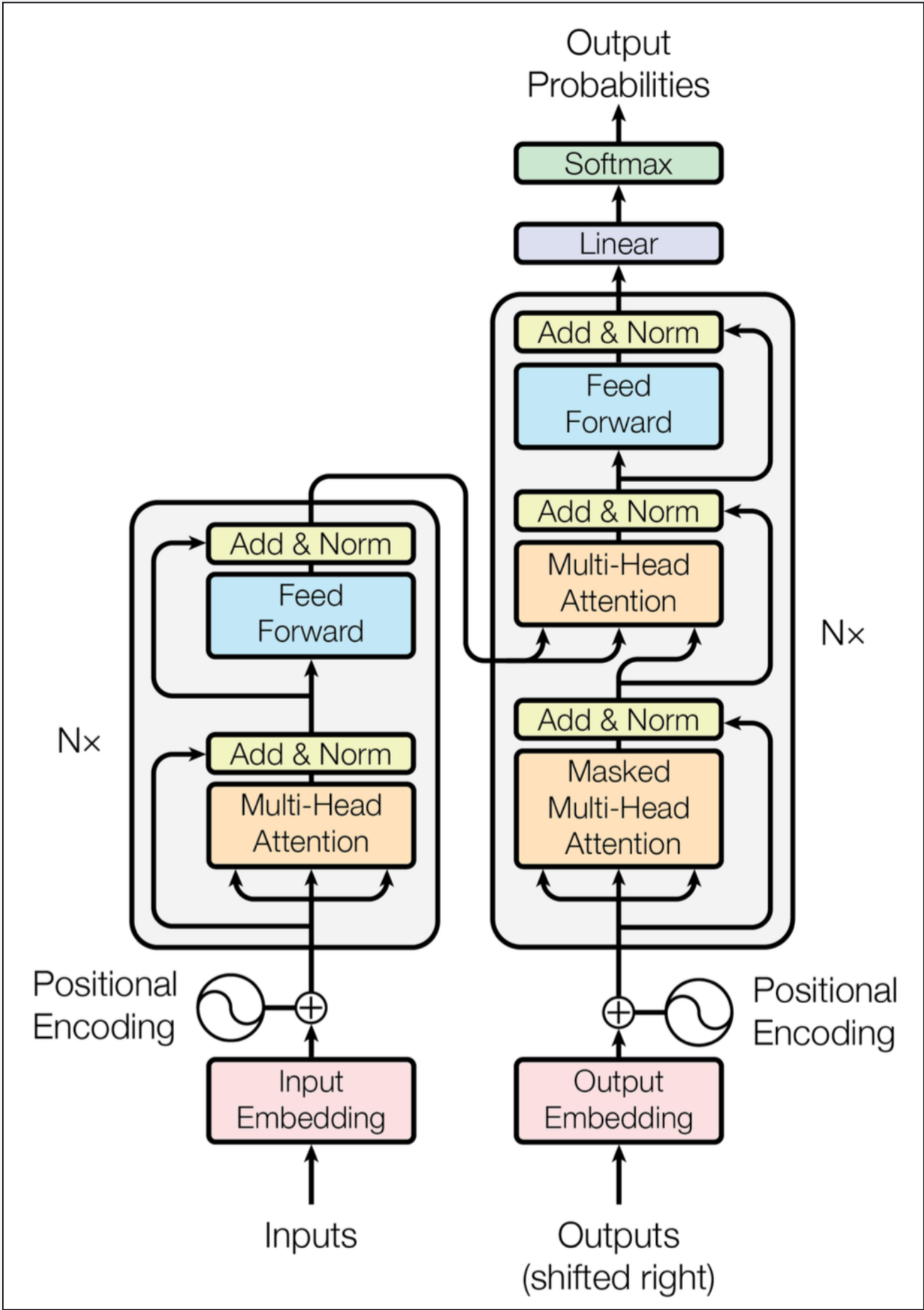Friday, September 8, 2023

Emmanouil Theodosis

✉ etheodosis@g.harvard.edu

🌐 manosth.github.io/

</> github.com/manosth/

# Deep learning is empirical…



**Transformer**

## Problems

- How to design?
- Functional understanding?
- Interpretability?

# ... and not very efficient

**ChatGPT**

- GPT-1 : 117M
- GPT-2 : 1.5B
- GPT-3 : 175B
- **GPT-4: 170T**

https://chat.openai.com

- LLaMA-2 7B    : 7B
- LLaMA-2 13B   : 13B
- **LLaMA-2 70B: 70B**

https://ai.meta.com/llama/

- LaMDA: 137B
- **PaLM : 540B**

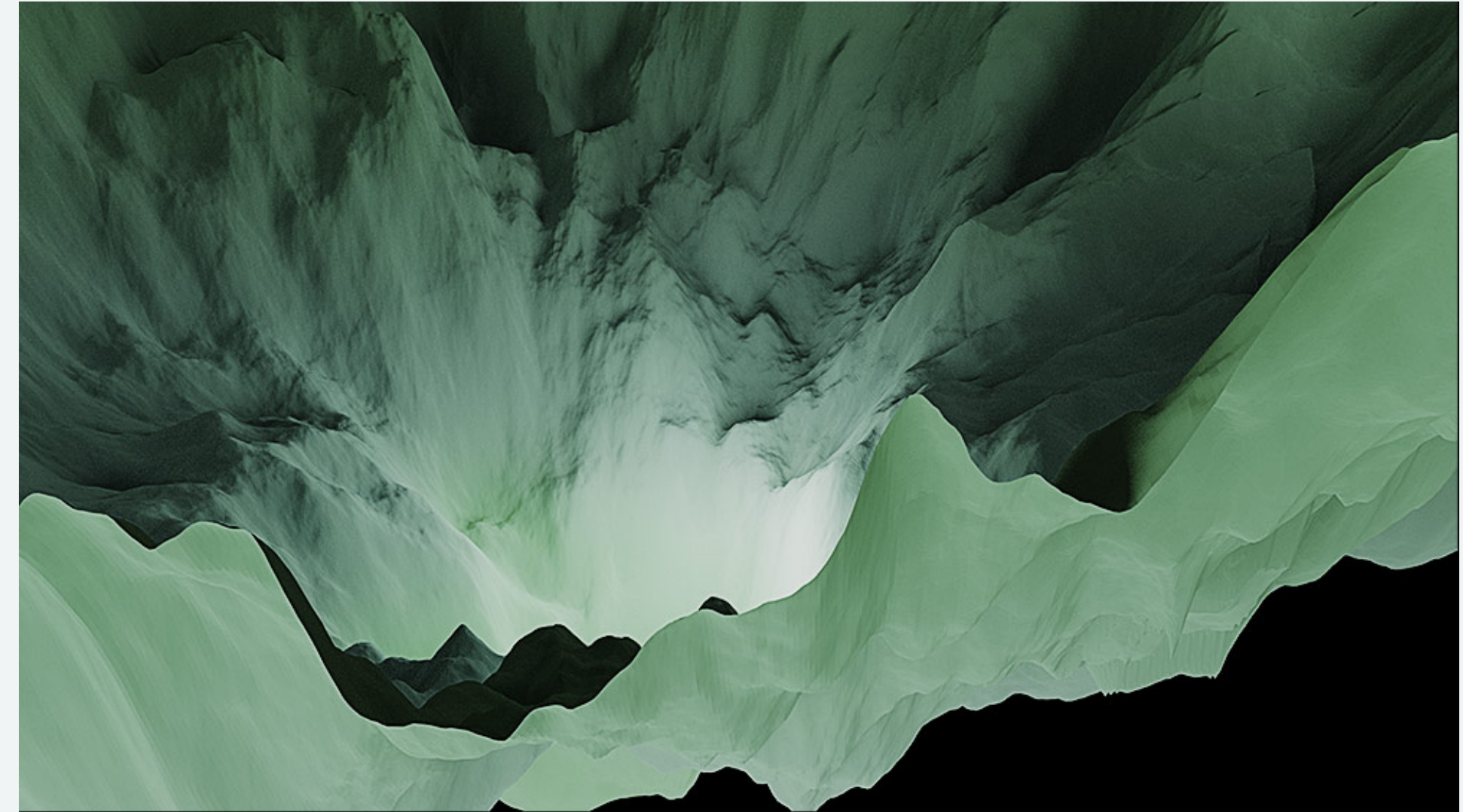https://bard.google.com

- **Claude 2: 175B**

https://claude.ai/

# So?

## Overparametrized models

- more parameters than data
- neurons memorize data points
- the rest interpolate

**result:** hope for the best.



## Use domain knowledge

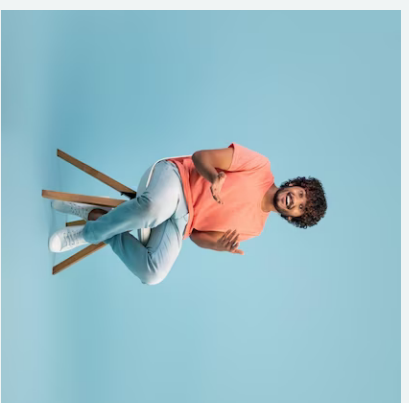- in a systematic way
- constrain the parameter space

**result:** similar (or better!) performance and more efficient.

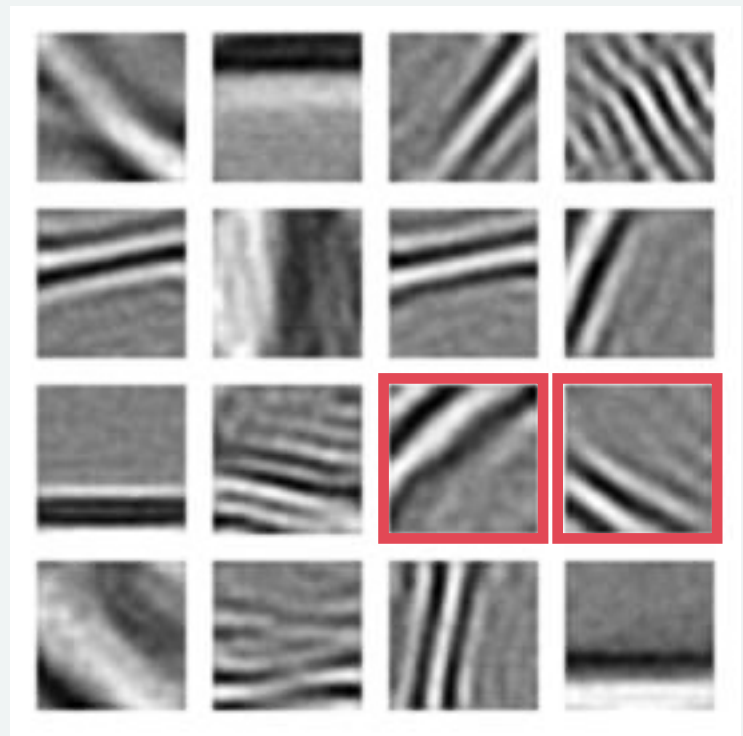# Enter equivariance

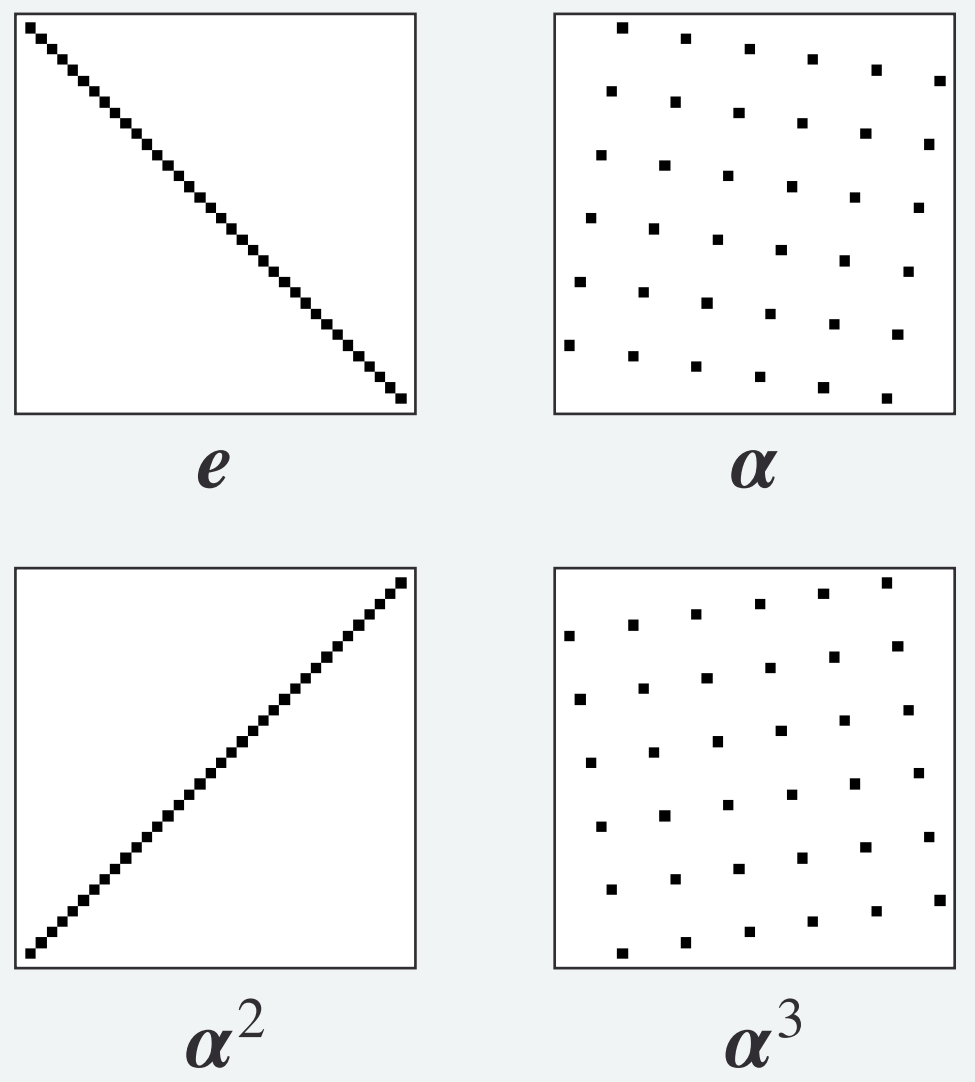## Classification has invariants


= Person on a chair


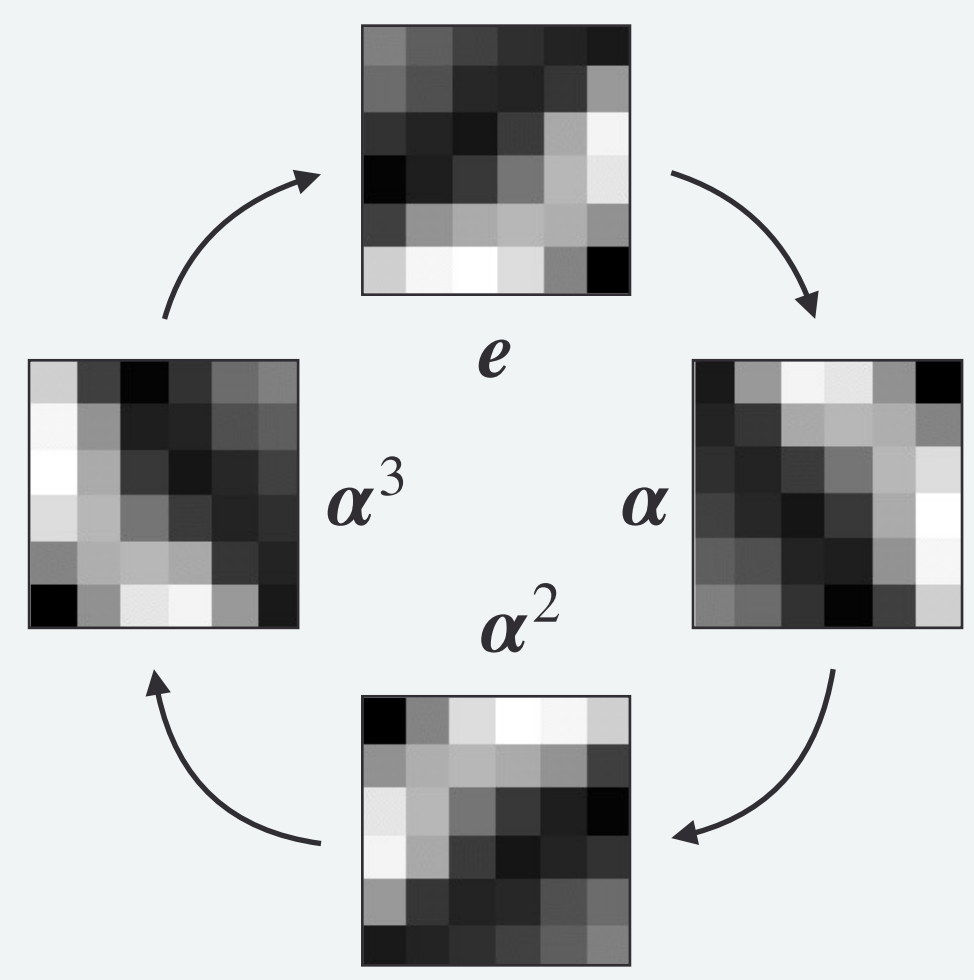= Person on a chair

## Filters have symmetries



<span style="color:red">different orientations</span>

## Group equivariant CNNs



$e$

$\alpha$

$\alpha^2$

$\alpha^3$



$e$    $\alpha$

$\alpha^2$    $\alpha^3$

# How do we generalize?



**code**

**privacy**

GCT
GCG
Alanine

···AAT**GCT**ACT···

···AAT**GCG**ACT···

**biology**

# Our work



$e$

$\boldsymbol{\alpha}^3$

$\boldsymbol{\alpha}$

$\boldsymbol{\alpha}^2$

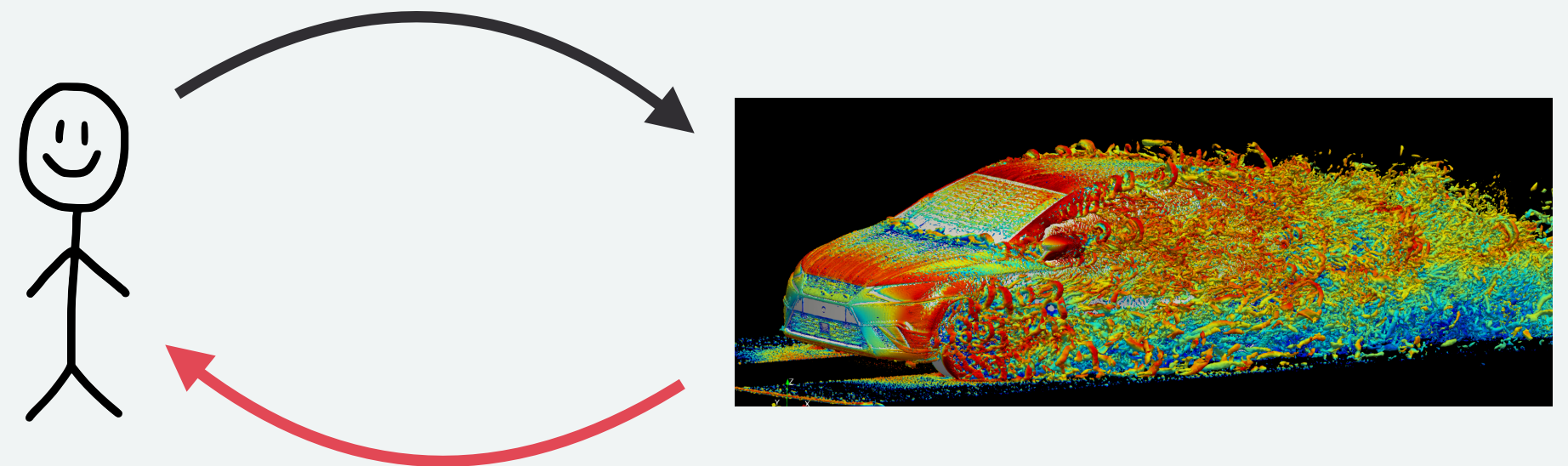other groups

# Our work

# Our work



other groups
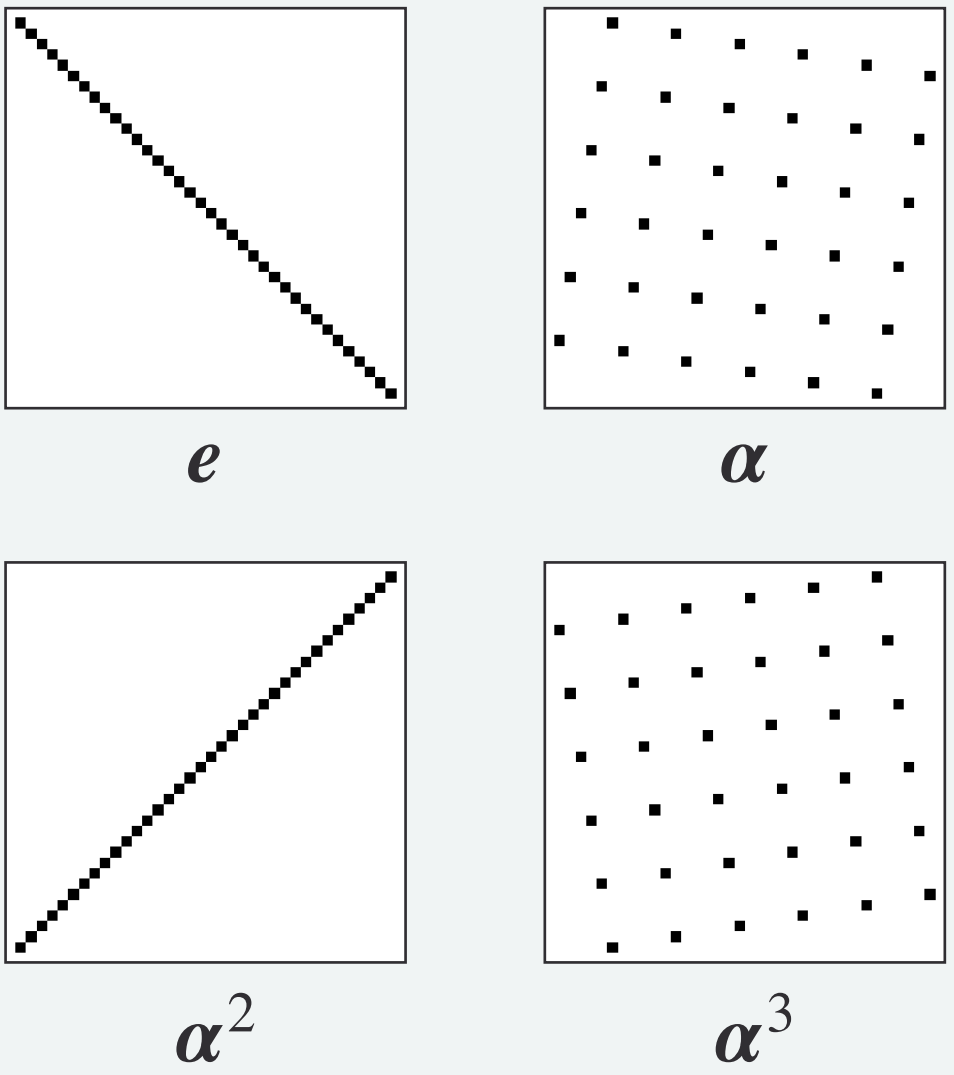
other actions

## Benefits



bidirectional information flow

efficient models

# Approach

# Elements of group theory



**Abstract group**
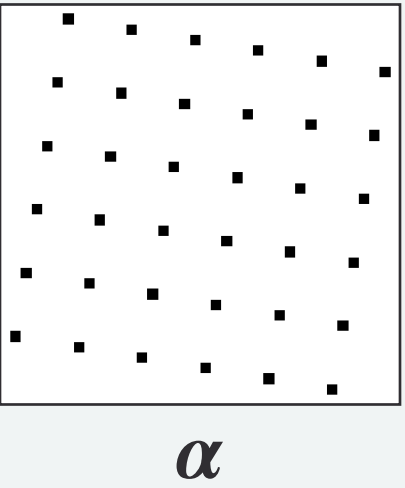High level interactions

$$\{e, \alpha, \alpha^2, \alpha^3\}$$

**Group representation**
"Implementation" matrices

$\alpha$

**Data space**
Where the group acts

# Elements of group theory

## Cyclic groups

**What are they?**

$$C_4 = \{e, \alpha, \alpha^2, \alpha^3\}$$

**Why do we care?**

- **short:** computation

- **long:** finite*, abelian, generator

$$T_{g_1}(\blacksquare) \quad T_{g_2}(\blacksquare) \quad \dots$$



$$\alpha^k \cdot \alpha^l = \alpha^l \cdot \alpha^k$$

$$T_g(\blacksquare) \quad T_g(\blacksquare) \quad \dots$$

## Equivariance

**What does it mean?**

$$f(\,\searrow \circ \,\blacksquare\,) = \searrow \circ f(\,\blacksquare\,)$$

(A bit more formally $f(T(x)) = T'(f(x))$.)

**Example:**

$$f(x) = x^2 \qquad \textbf{vs} \qquad f(\alpha \cdot x) = \alpha^2 \cdot x^2$$

Here $T(u) = \alpha \cdot u$ and $T'(u) = \alpha^2 \cdot u$.

So where do groups come in?

$$T = \boxed{\phantom{xxx}} \qquad\qquad T' = \boxed{\phantom{xxx}}$$
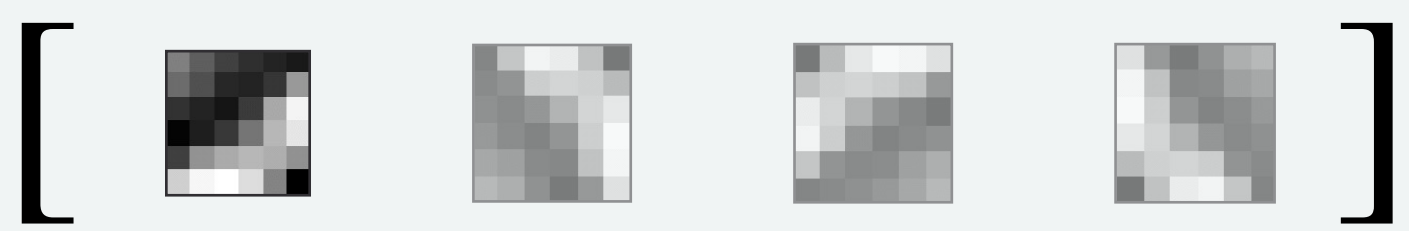
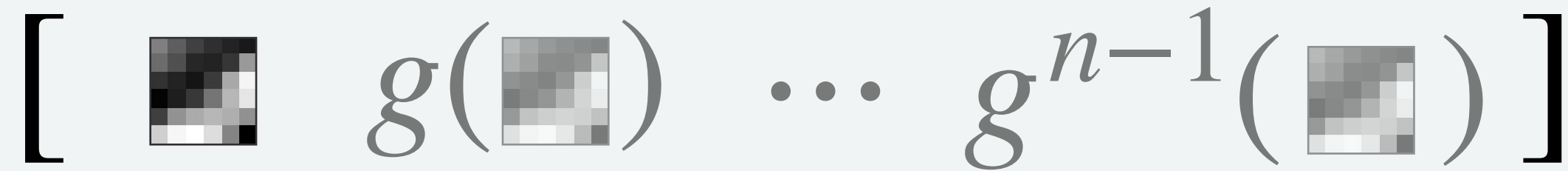Represent the "same" transformation!

# Main idea

**Theorem[1]:** For equivariance of a neural network with respect to a group $G$ each layer needs to be a $G$-convolution.

$$(f * g)(u) = \sum_{v \in G} f(uv^{-1})g(v) \qquad \left(\text{vs } (f * g)(u) = \sum_{v \in \mathbb{Z}} f(u-v)g(v)\right)$$
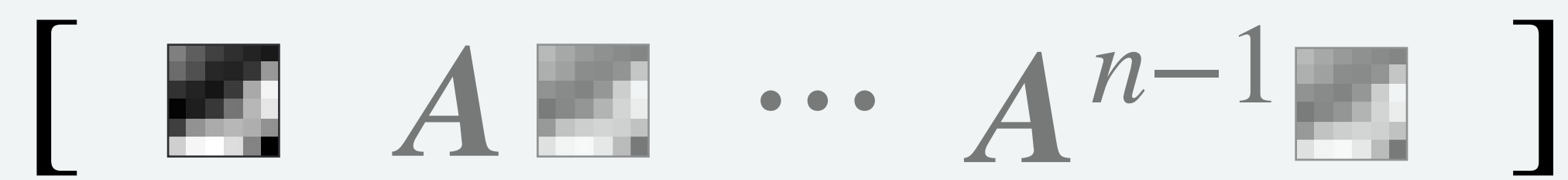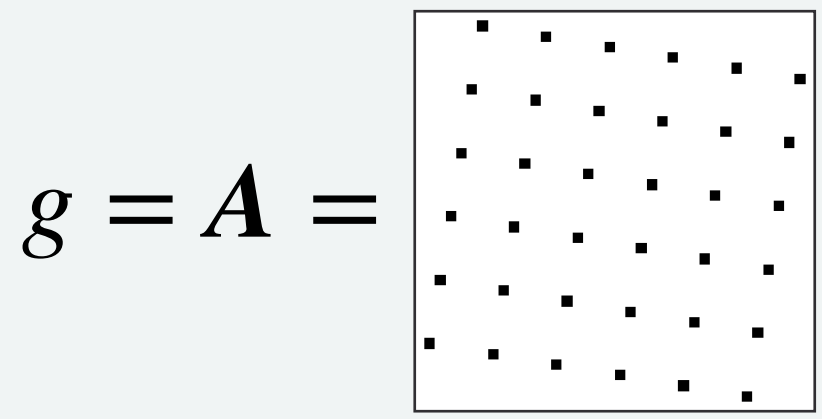
## For $C_4$
Convolution with the rotated set

$$\left[ \blacksquare \quad \blacksquare \quad \blacksquare \quad \blacksquare \right]$$

## What about $C_n$?
Cardinality and action change

$$\left[ \blacksquare \quad g(\blacksquare) \quad \cdots \quad g^{n-1}(\blacksquare) \right]$$

## How to represent $g$?
Use invertible matrices

$$\rho : G \to \mathrm{GL}_d(\mathbb{R}) \qquad\qquad g = A = \blacksquare \qquad\qquad \left[ \blacksquare \quad A\blacksquare \quad \cdots \quad A^{n-1}\blacksquare \right]$$

[1]: Risi Kondor and Shubhendu Trivedi. "On the generalization of equivariance and convolution in neural networks to the action of compact groups". In International Conference on Machine Learning, 2018.
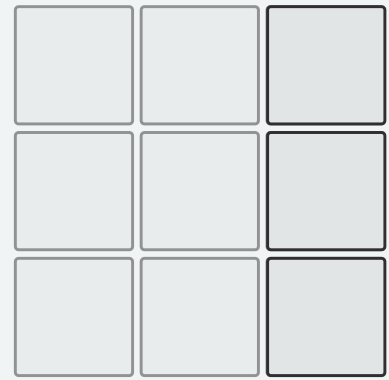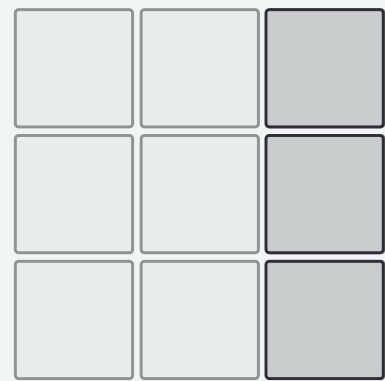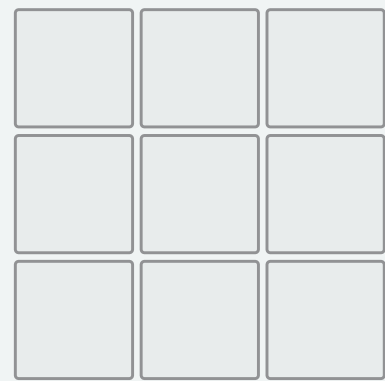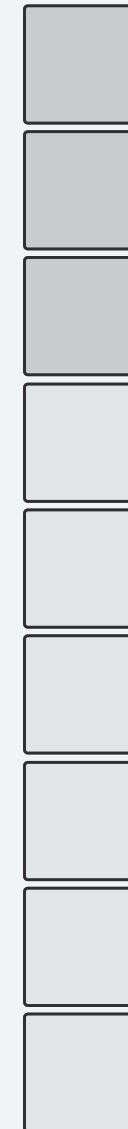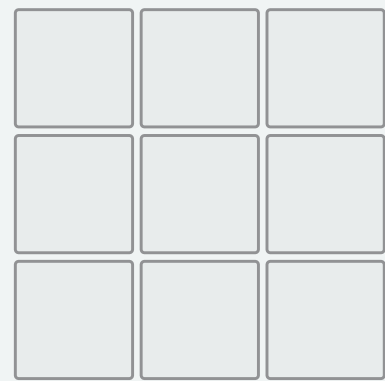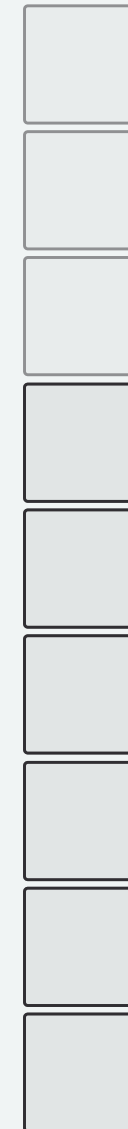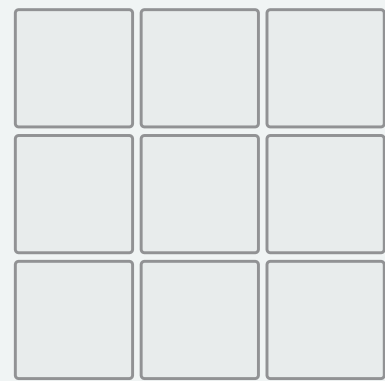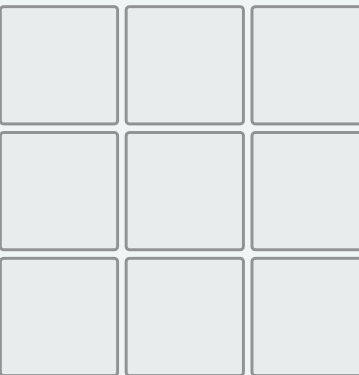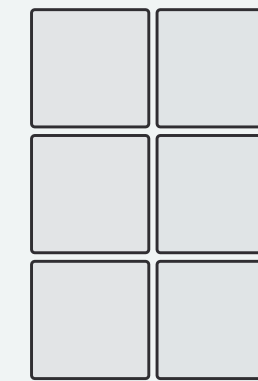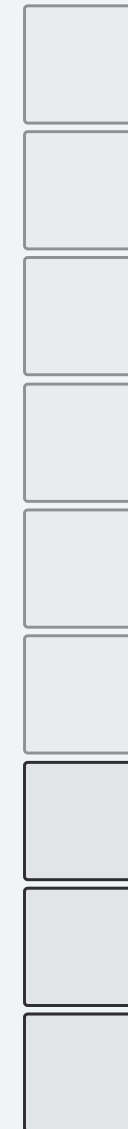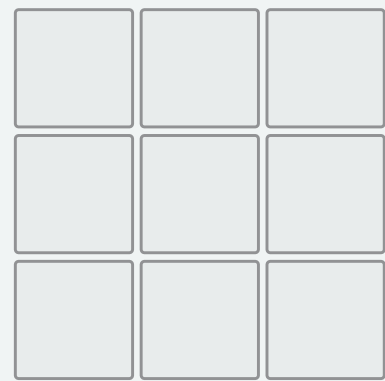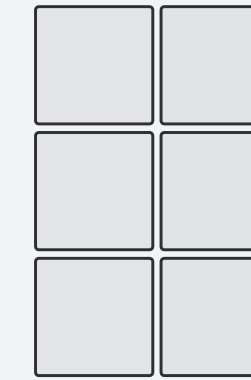
# Vectorization

# Vectorization

# Vectorization

# Vectorization

# Vectorization

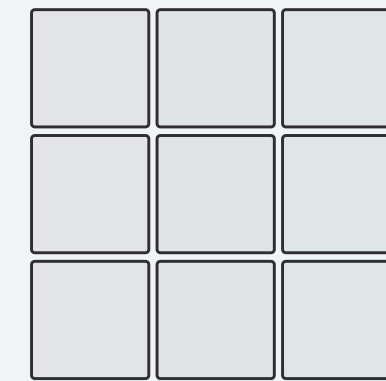# Vectorization

# Vectorization

# Vectorization
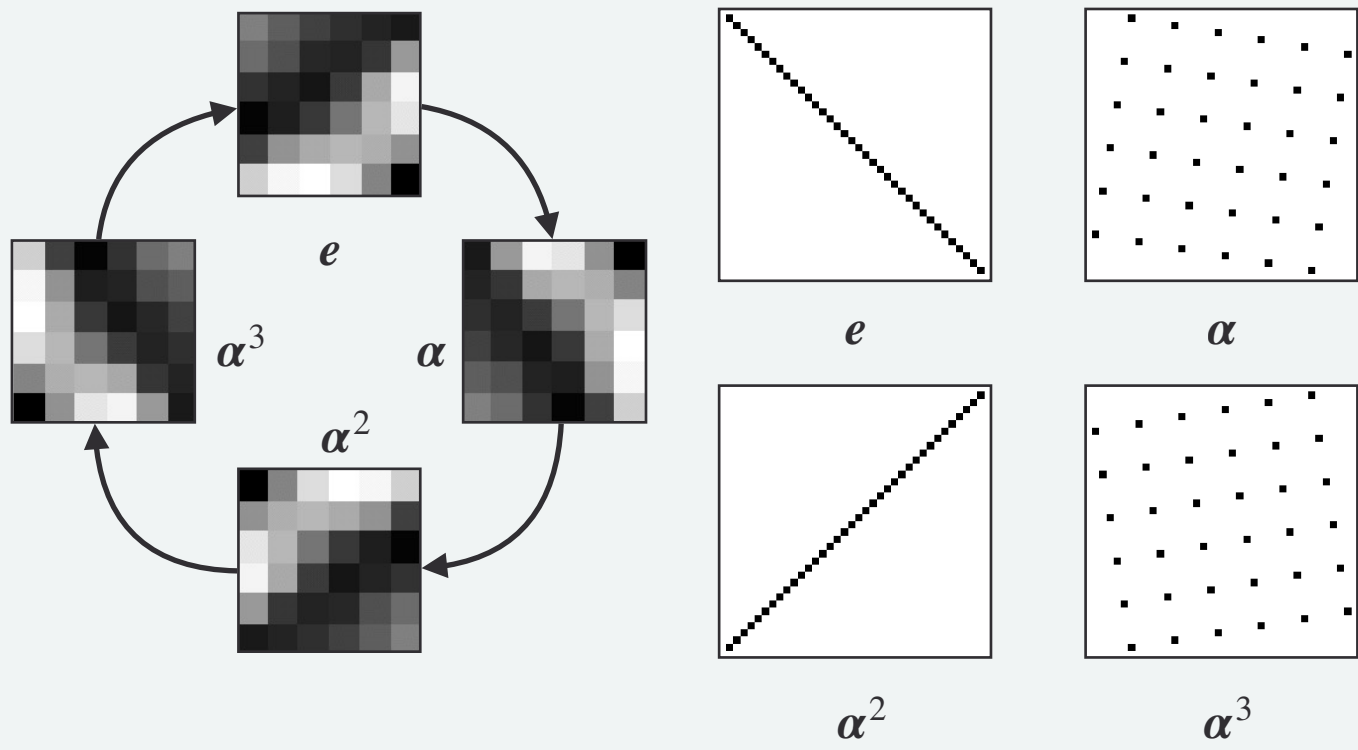
# Vectorization

# Vectorization

# Vectorization

# Vectorization

# Final model

**Filters**

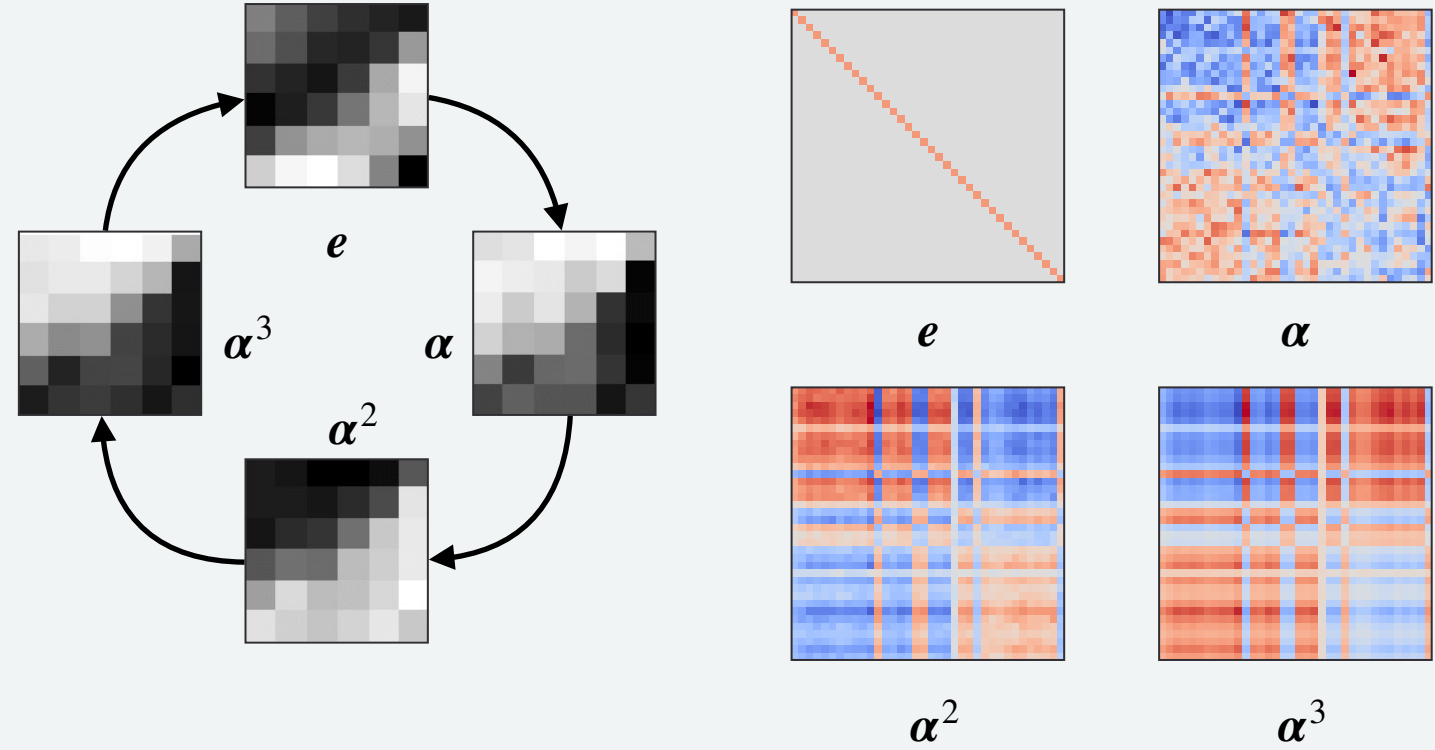$$\left[\ \blacksquare\ \phi_A(\blacksquare)\ \cdots\ \phi_A^{n-1}(\blacksquare)\ \right]$$

**learn with backprop**

$$\phi_A : \mathbb{R}^{n \times m} \to \mathbb{R}^{n \times m}$$
$$X \mapsto \text{vec}^{-1}(A\,\text{vec}(X))$$

**Now we go from this…**



**… to this!**



**Invertibility loss**

$A$ needs to be invertible
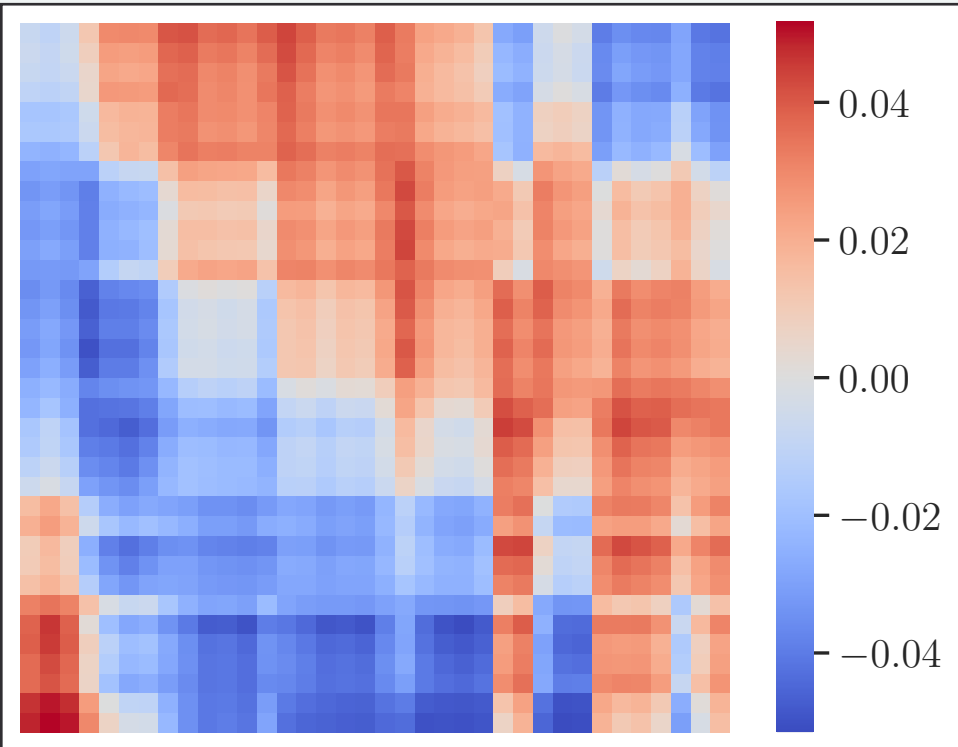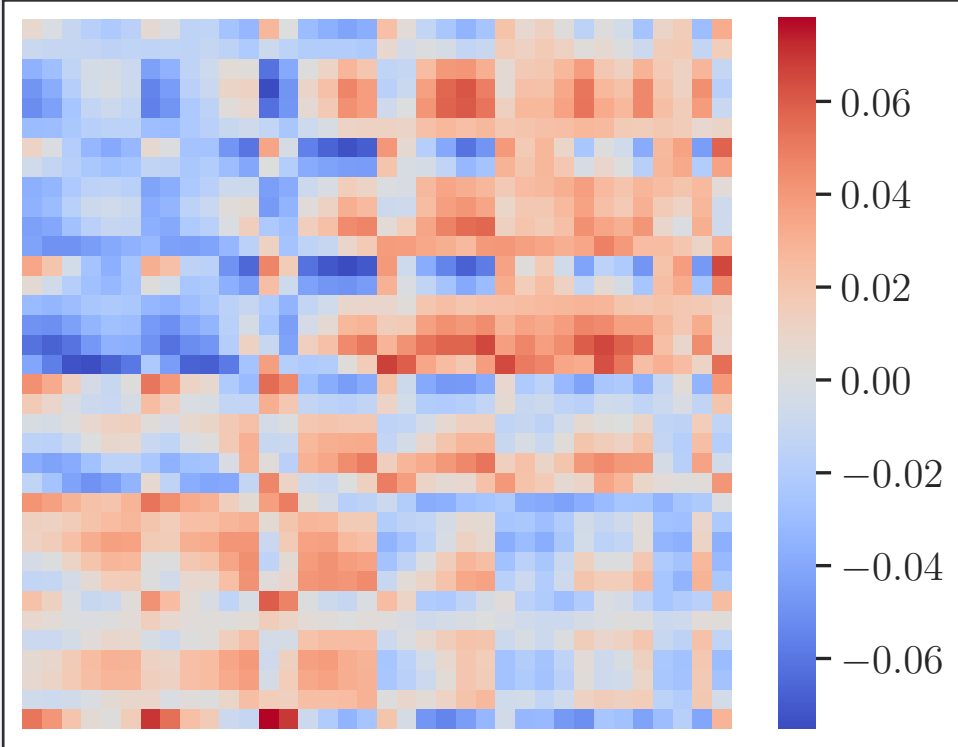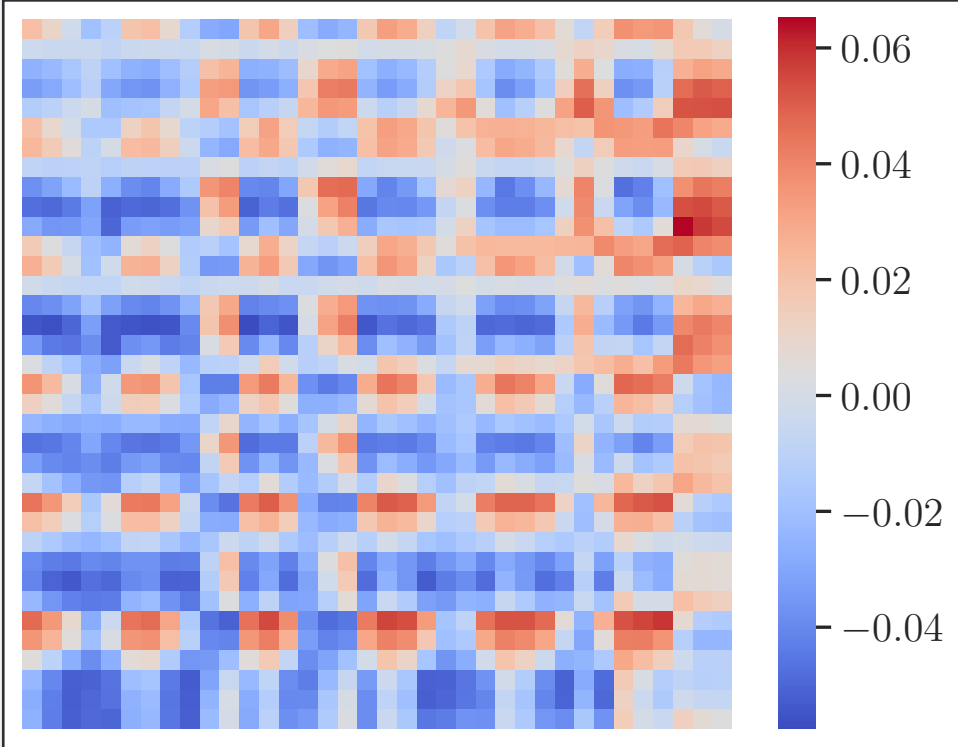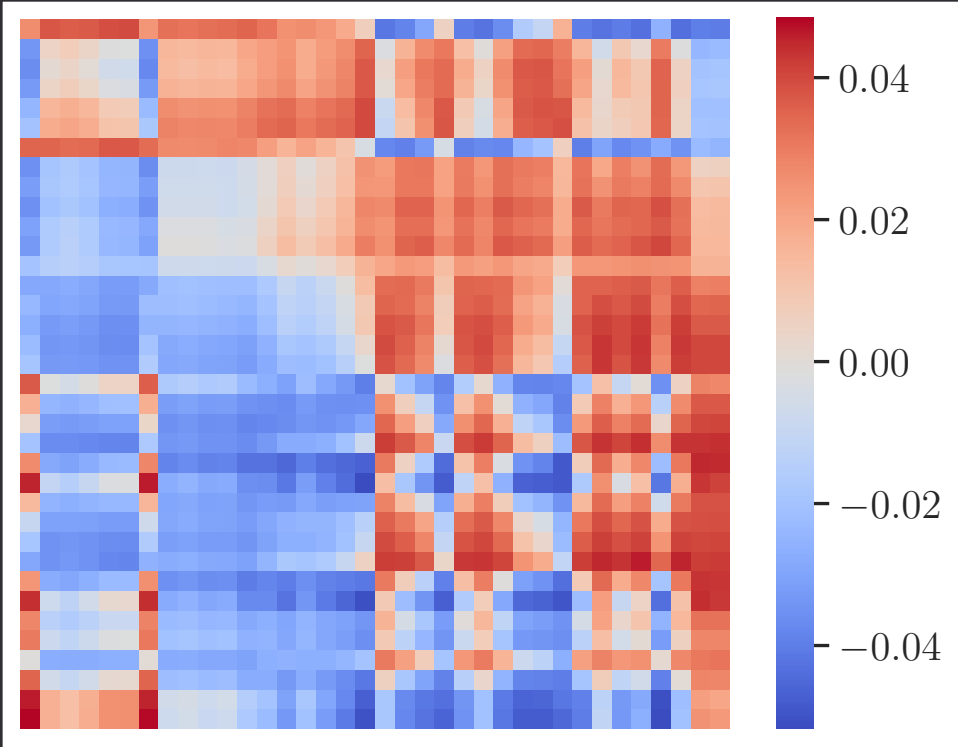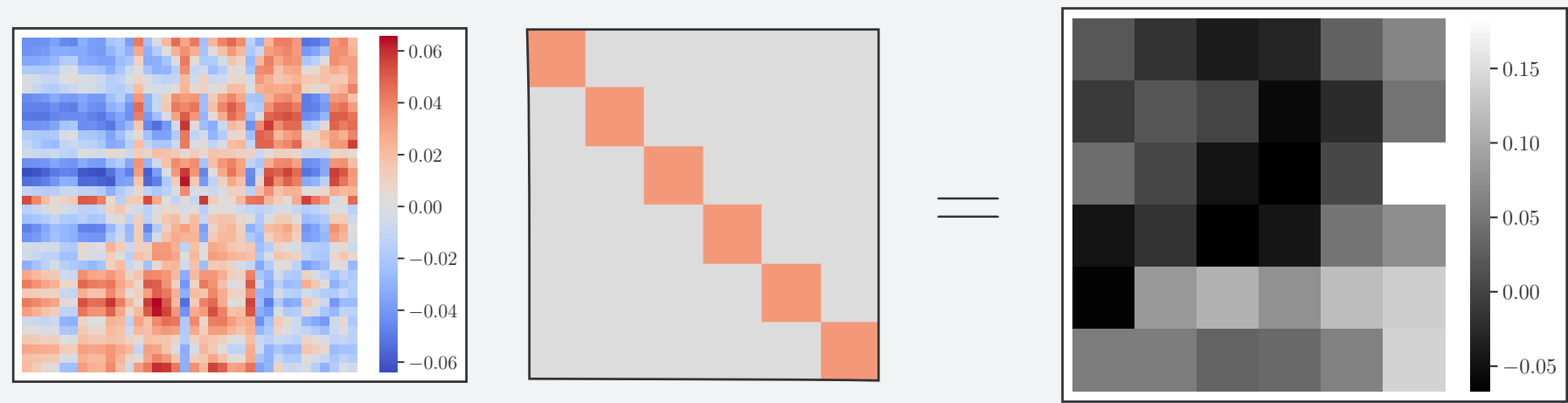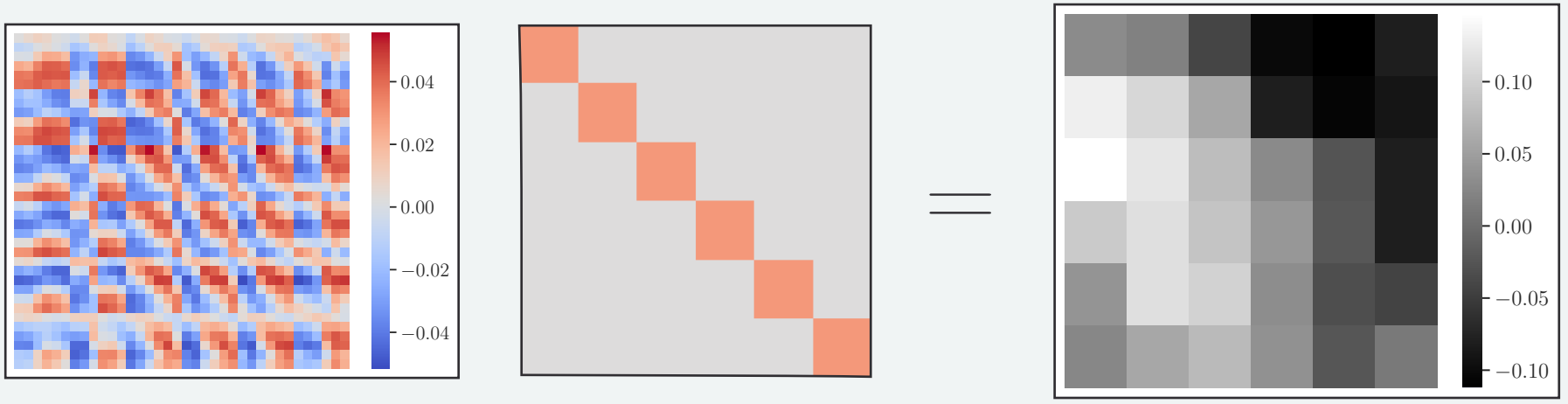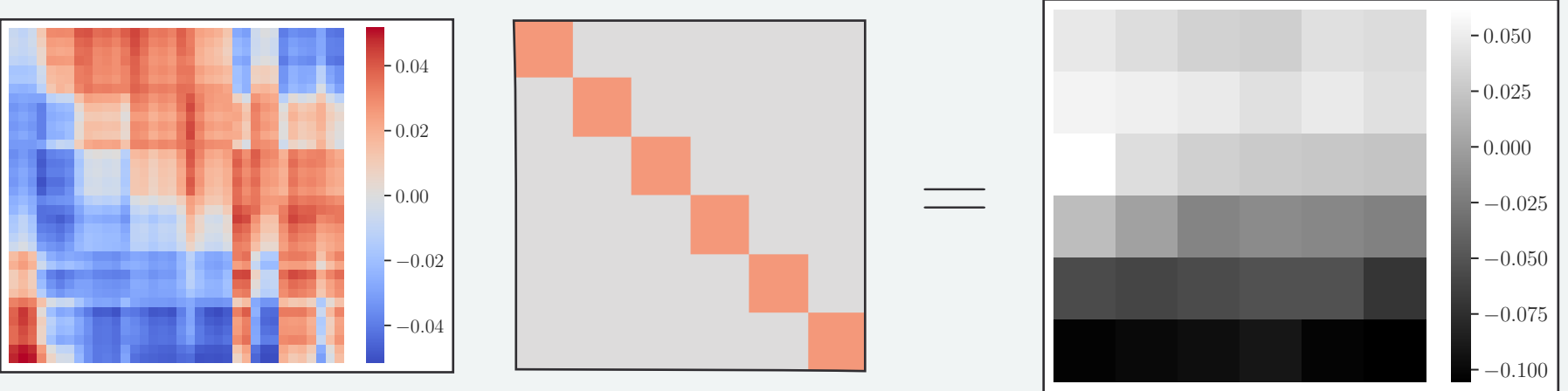
$$L = \mu \|A\tilde{A} - I\|_F$$

# Experiments

# Group structures



**Skew-symmetric**

**Toeplitz**

**Multi-scale**

# Interpreting the structures



**Skew-symmetric**

**Toeplitz**
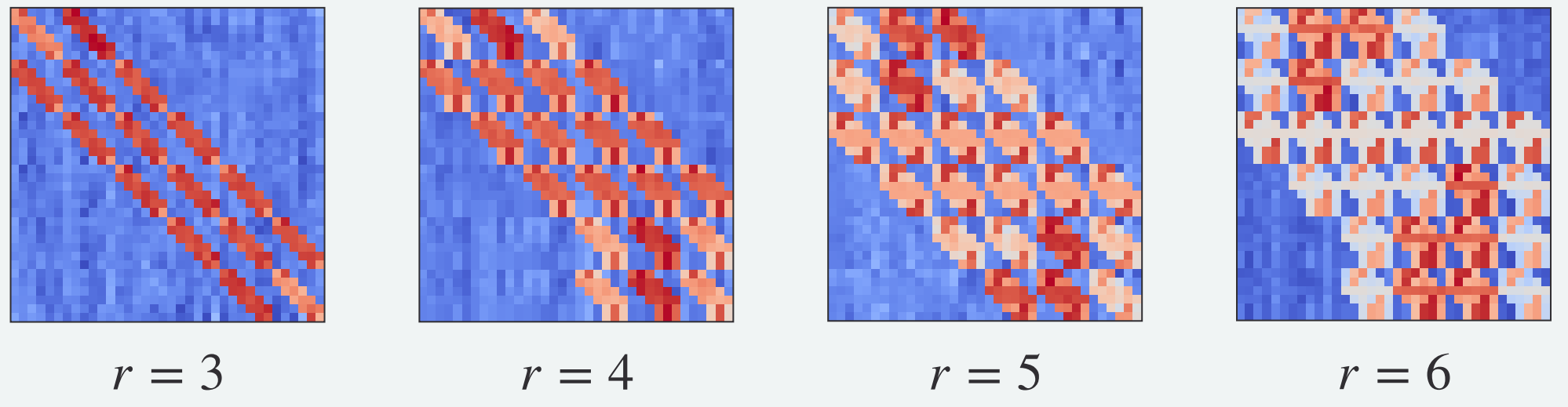
**Multi-scale**

# Interpreting the structures



**Rotations**

$\theta = 90°$     $\theta = 60°$     $\theta = 45°$     $\theta = 30°$

**Pooling**

$r = 3$     $r = 4$     $r = 5$     $r = 6$

**... both!**

$\theta = 60°$

$r = 4$     $r = 5$     $r = 6$
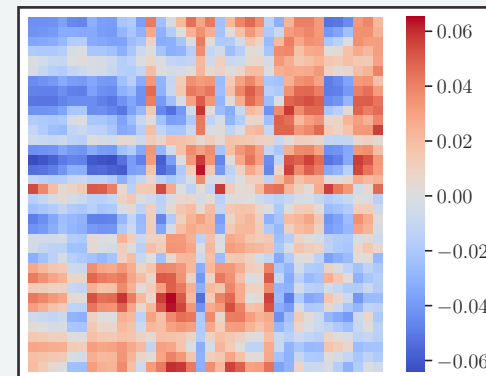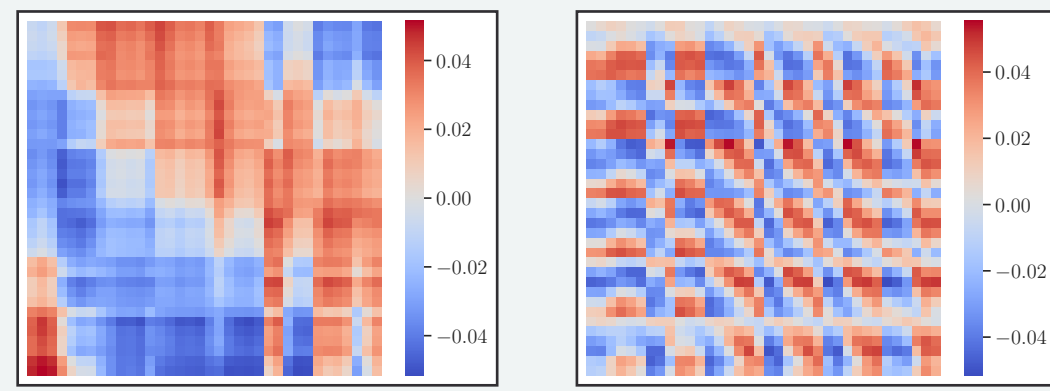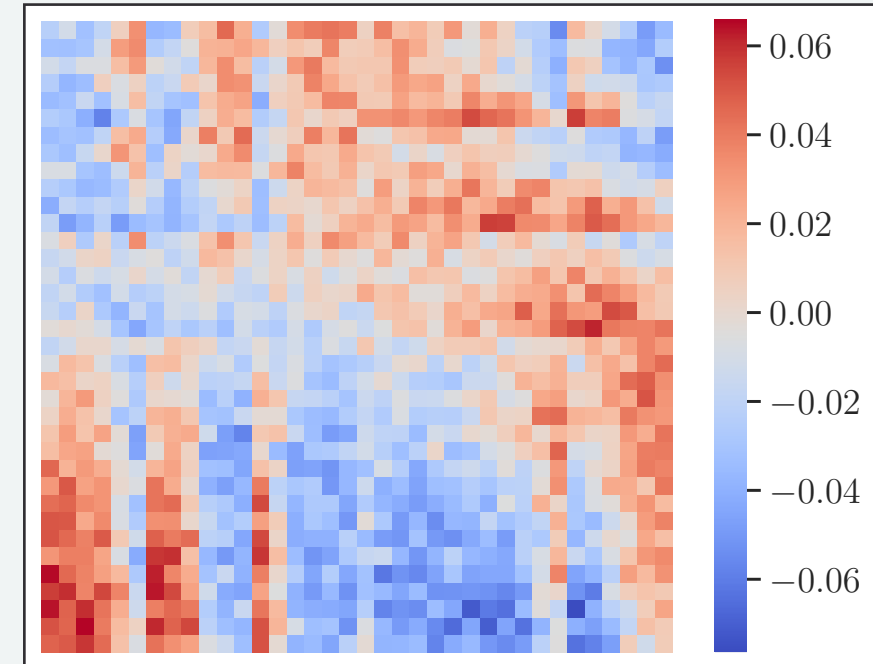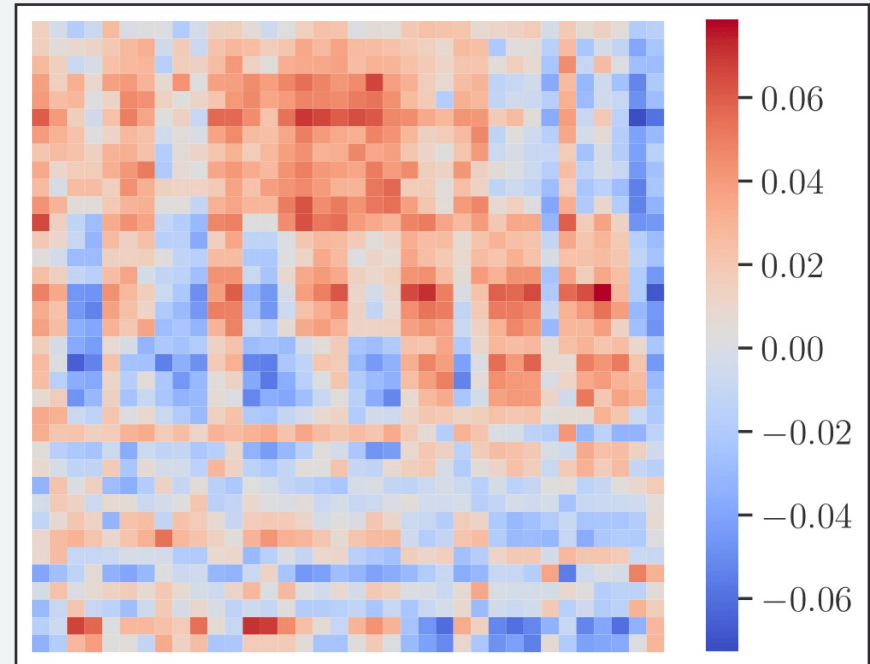
## Observations

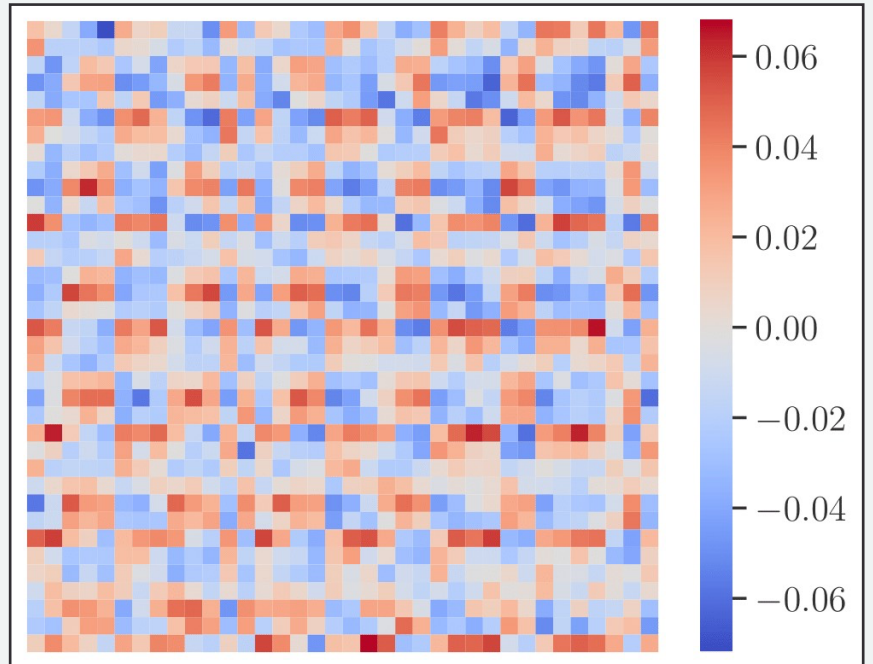- aligned on grid
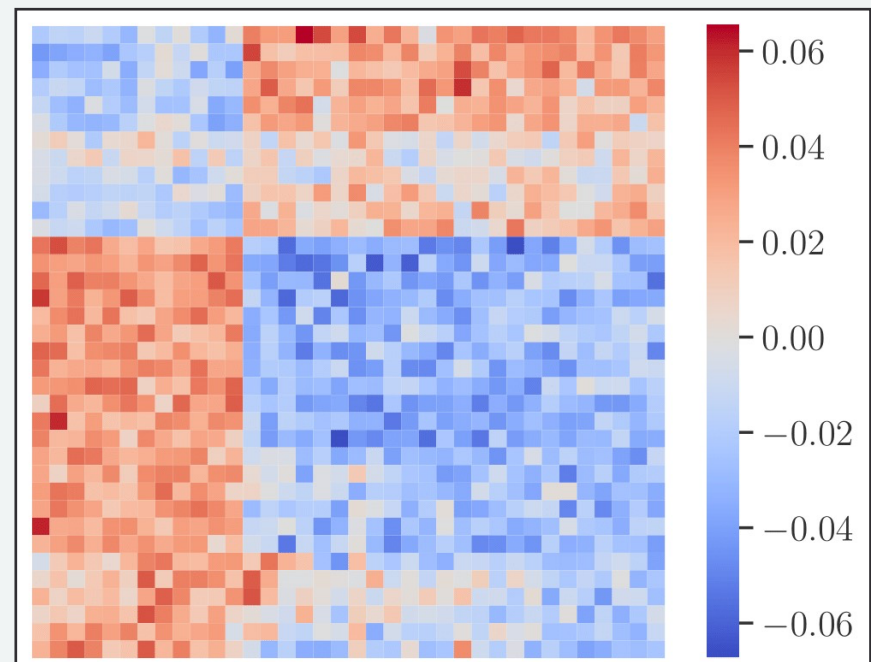- acts on many pixels
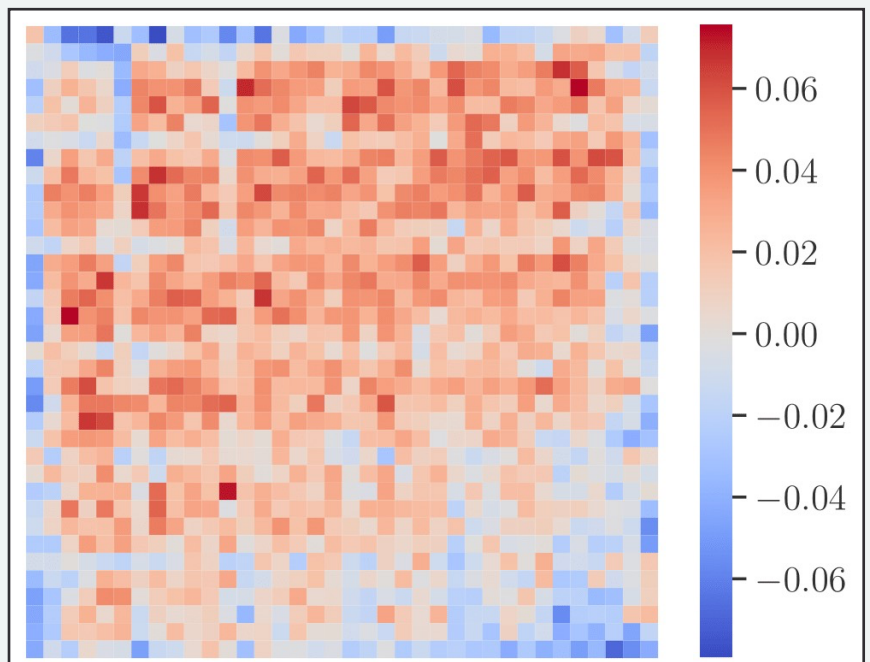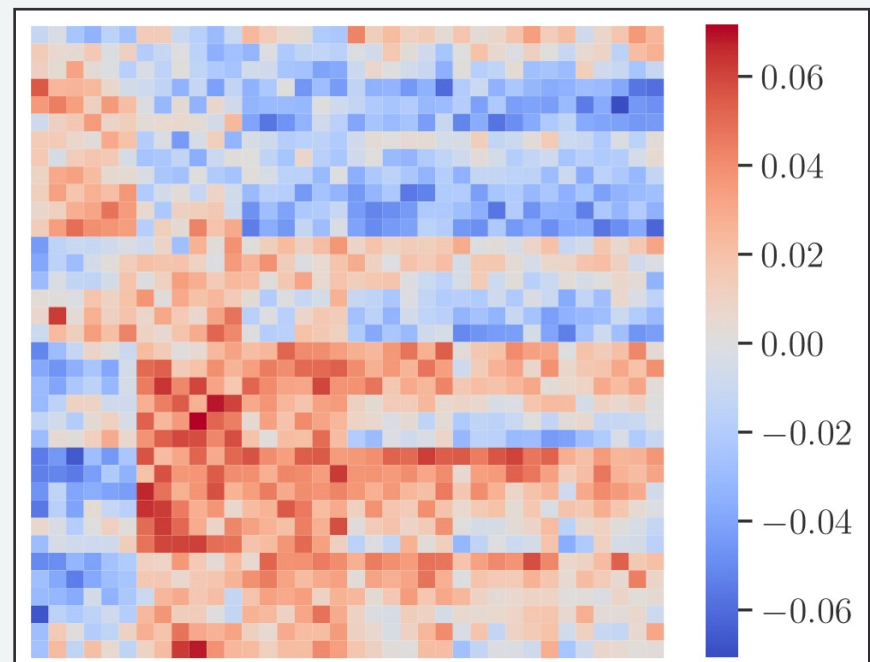
# More experiments

## Actions are consistent



CIFAR10

MNIST

rotMNIST

FashionMNIST

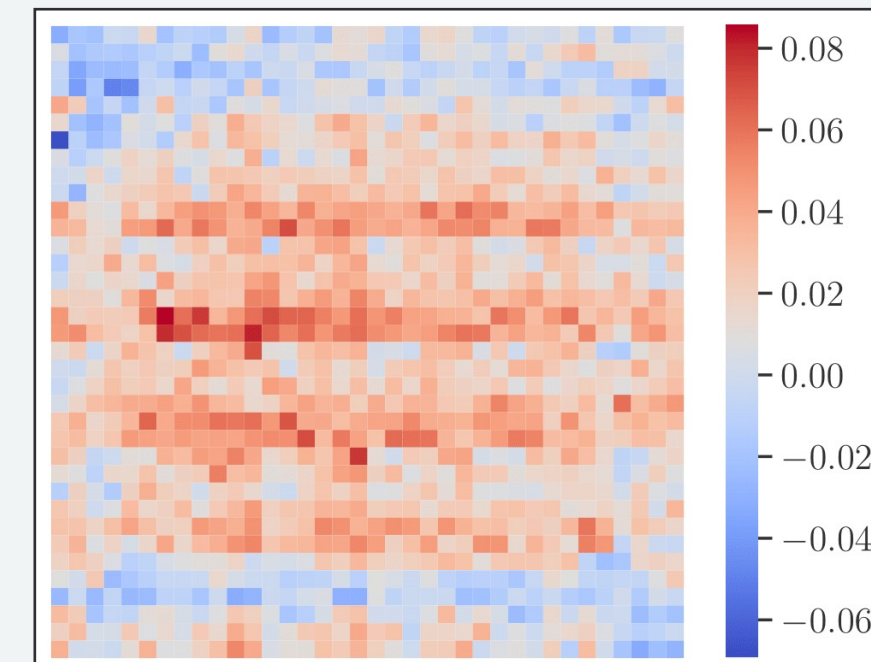## … and task dependent



**half** (CIFAR10)

**sim** (MNIST)

# Even more experiments

$$\left[ \begin{array}{cccc} \blacksquare & \phi_A(\blacksquare) & \cdots & \phi_A^{n-1}(\blacksquare) \end{array} \right] \quad \textbf{vs} \quad \left[ \begin{array}{c} \blacksquare \end{array} \right] \quad \textbf{vs} \quad \left[ \begin{array}{cccc} \blacksquare & \blacksquare & \cdots & \blacksquare \end{array} \right]$$

## Knowledge transfer

| MNIST | LGN | P4CNN | Single channel | Same channels |
|---|---|---|---|---|
| Accuracy | **98.86** | 98.68 | 98.52 | 97.42 |

| Fashion MNIST | LGN | P4CNN | Single channel | Same channels |
|---|---|---|---|---|
| Accuracy | 89.23 | 88.57 | 85.79 | **89.27** |

## CNN comparison

| | MNIST | rotMNIST | FashionMNIST | CIFAR10 | |
|---|---|---|---|---|---|
| LGCNN | 99.27 | **93.75** | **91.46** | **79.03** | (Based on ALL-CNN) |
| P4CNN | **99.35** | 92.96 | 91.44 | 75.26 | |

# Concluding

## Key takeaway



## What the future holds

- non-cyclic groups
- apply to other domains
- systematic interpration

# THANK YOU