# Inductive bias transfer between brains and machines

Fabian Sinz
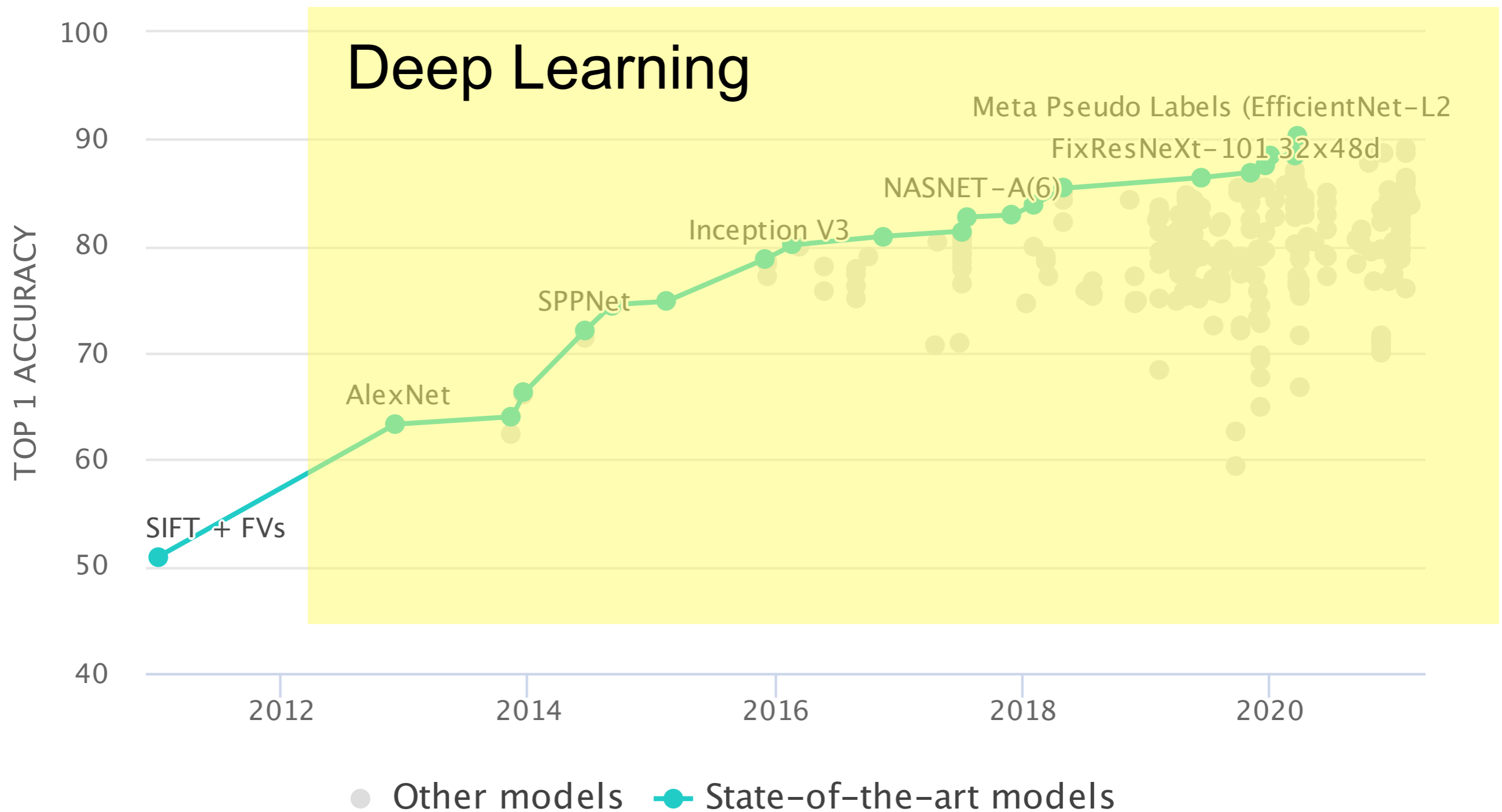
Neural Intelligence Group, Uni Tübingen
soon Uni Göttingen

@sinzlab
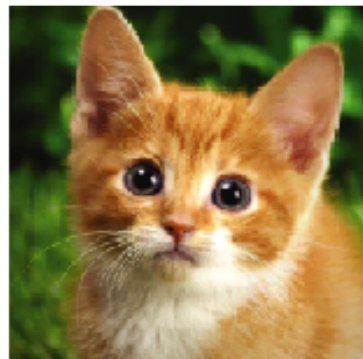
NICE 2021

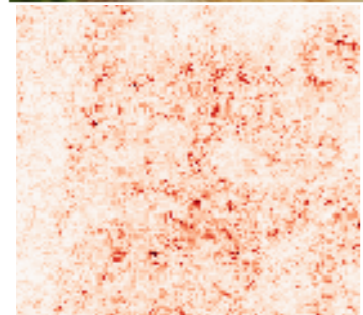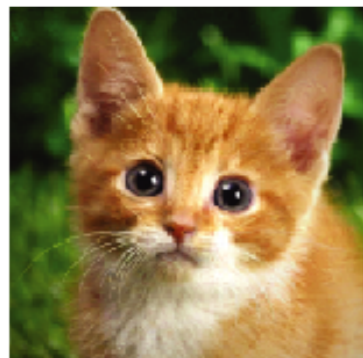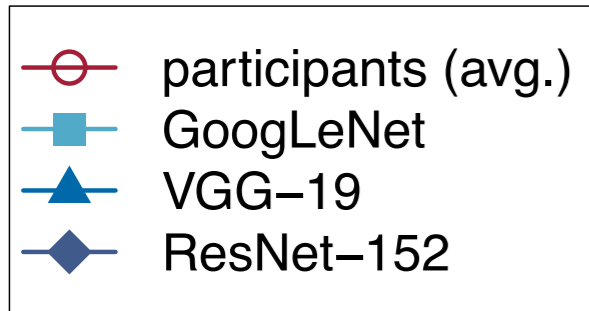# Success of deep learning
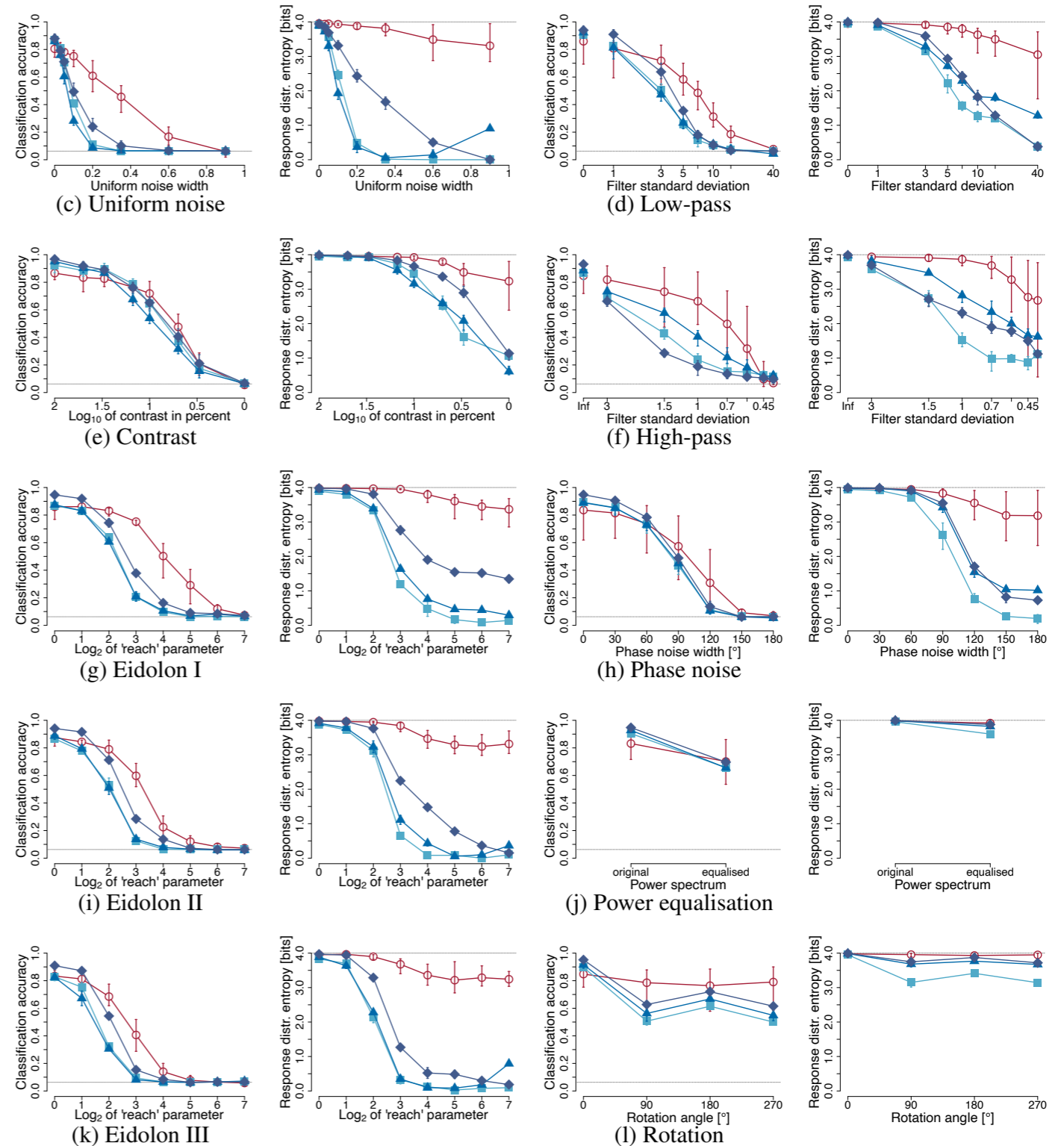
# Current state-of-the-art is brittle

"cat"

+

adversarial
perturbation

=

"moped"

original          texturised images

Geirhos et al. 2018. "Generalisation in Humans and Deep Neural Networks."
Szegedy et al. 2013. "Intriguing Properties of Neural Networks."
Sinz et al 2019. "Engineering a Less Artificial Intelligence." *Neuron.*

# Human visual system is robust

Legend:
- participants (avg.)
- GoogLeNet
- VGG−19
- ResNet−152

Axes: Classification accuracy, Response distr. entropy [bits], Uniform noise width, Filter standard deviation

colour    greyscale
Colour

# Human visual system is robust



Geirhos et al. 2018. "Generalisation in Humans and Deep Neural Networks."

EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

# Inductive Bias - Implicit Assumptions
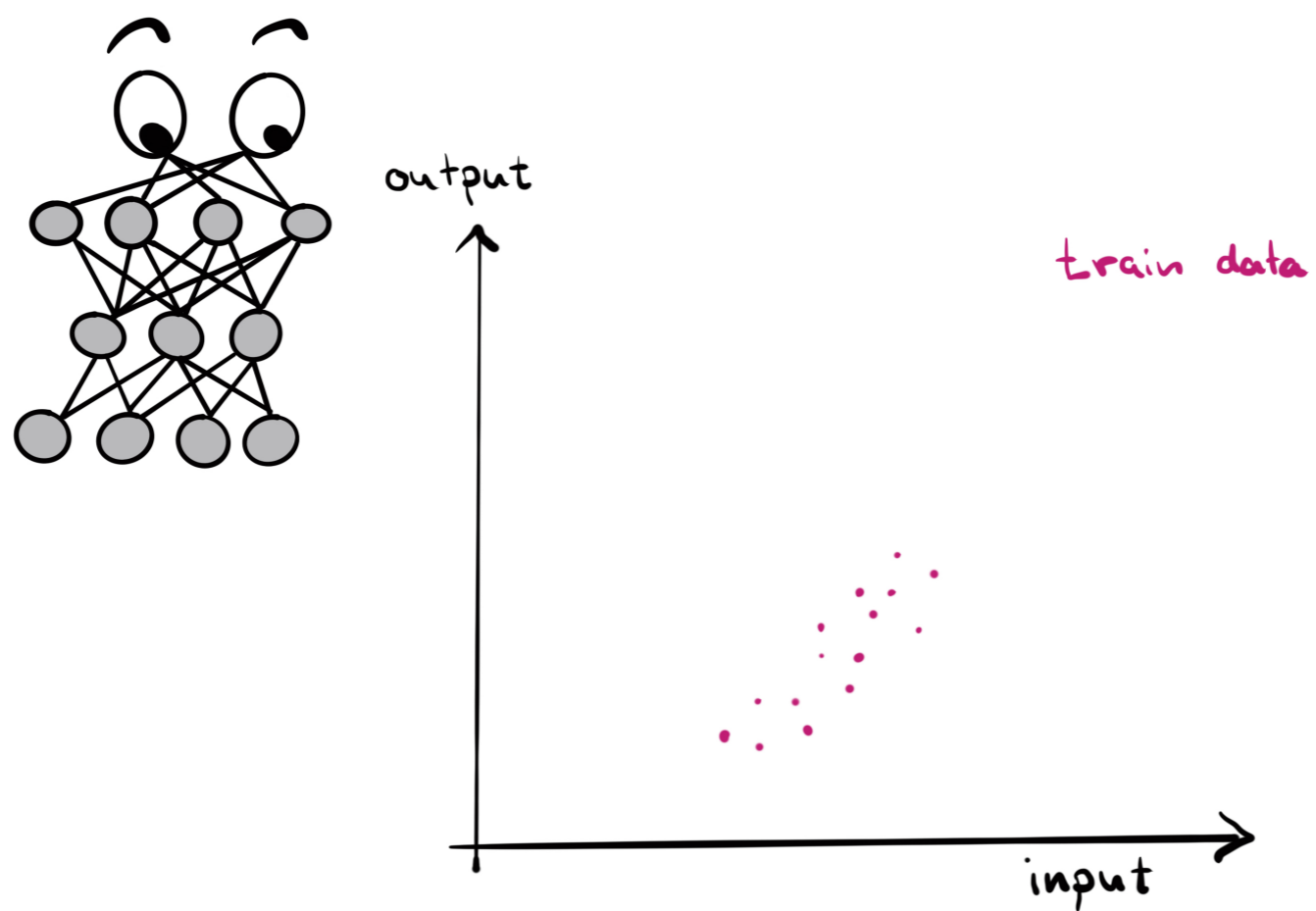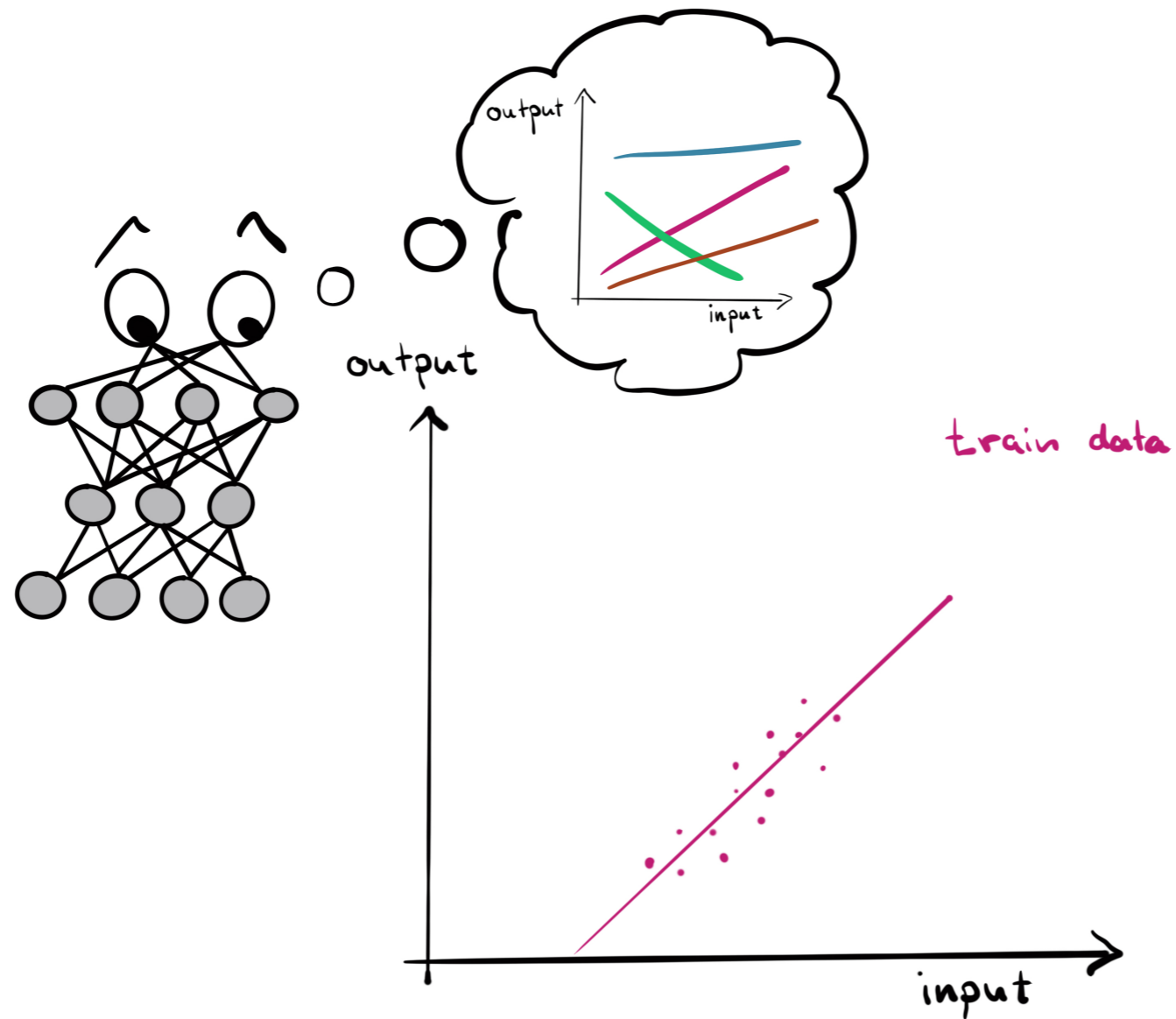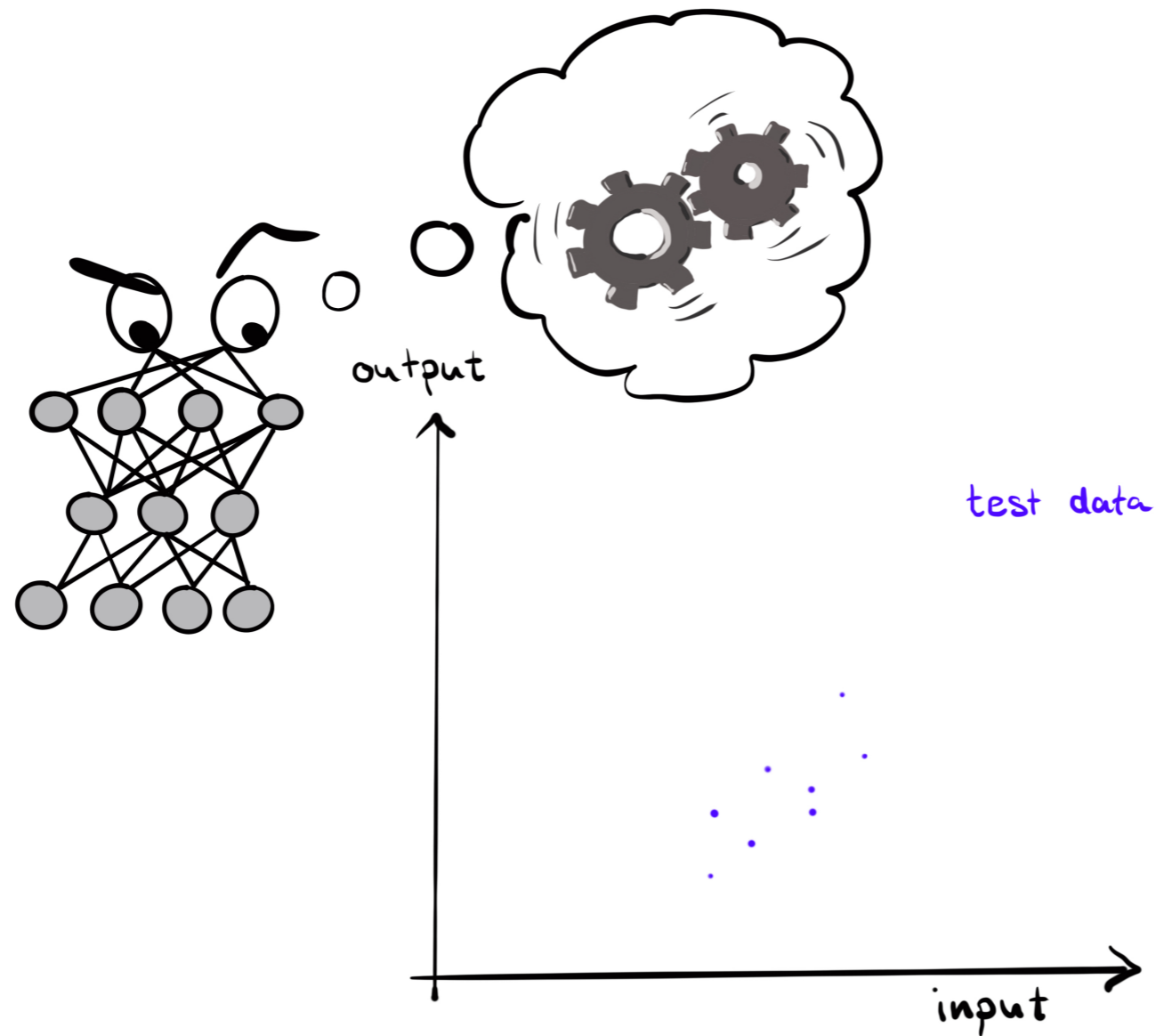
# Inductive Bias - Implicit Assumptions

# Inductive Bias - Implicit Assumptions
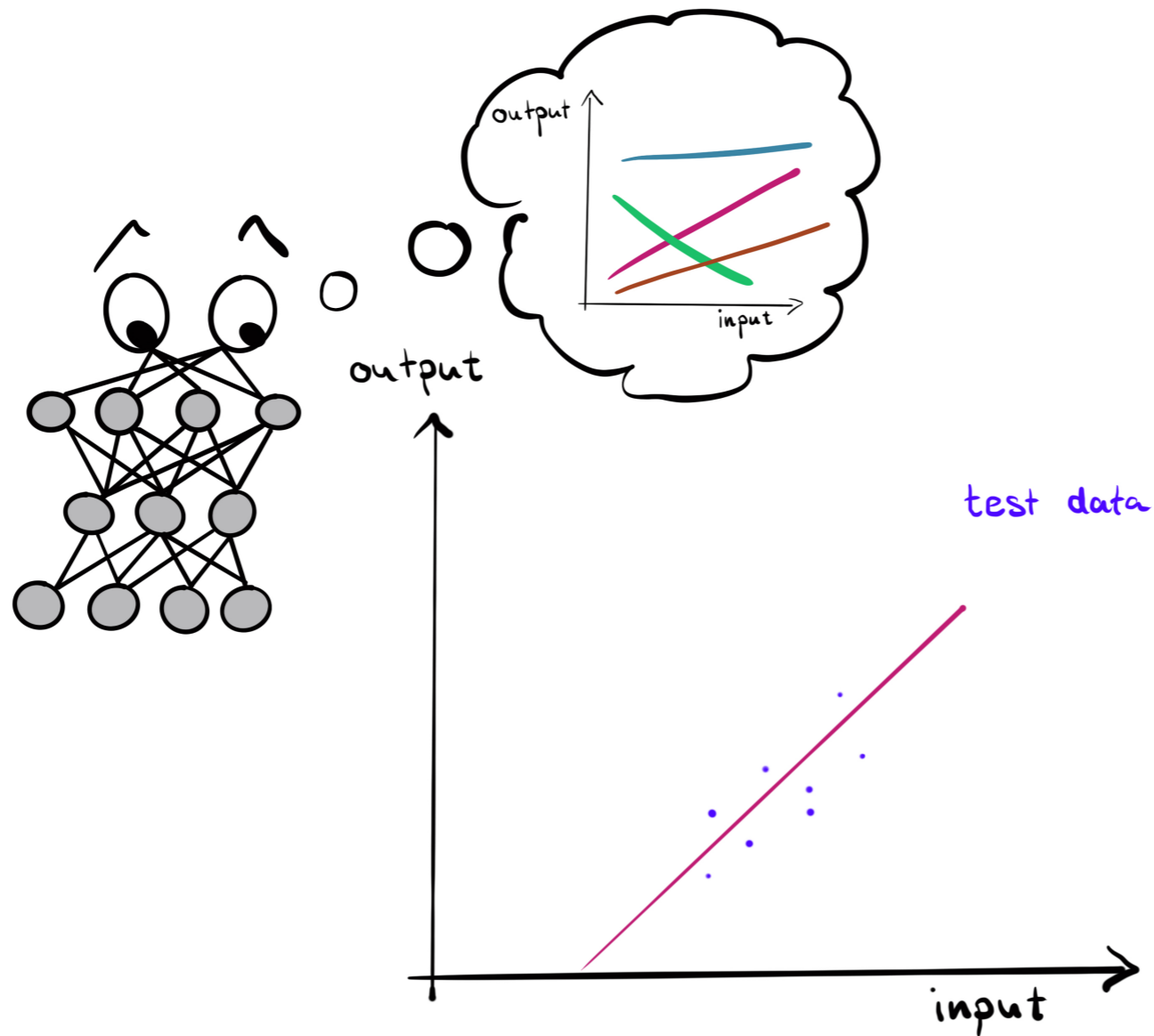
# Inductive Bias - Implicit Assumptions

# Inductive Bias - Implicit Assumptions
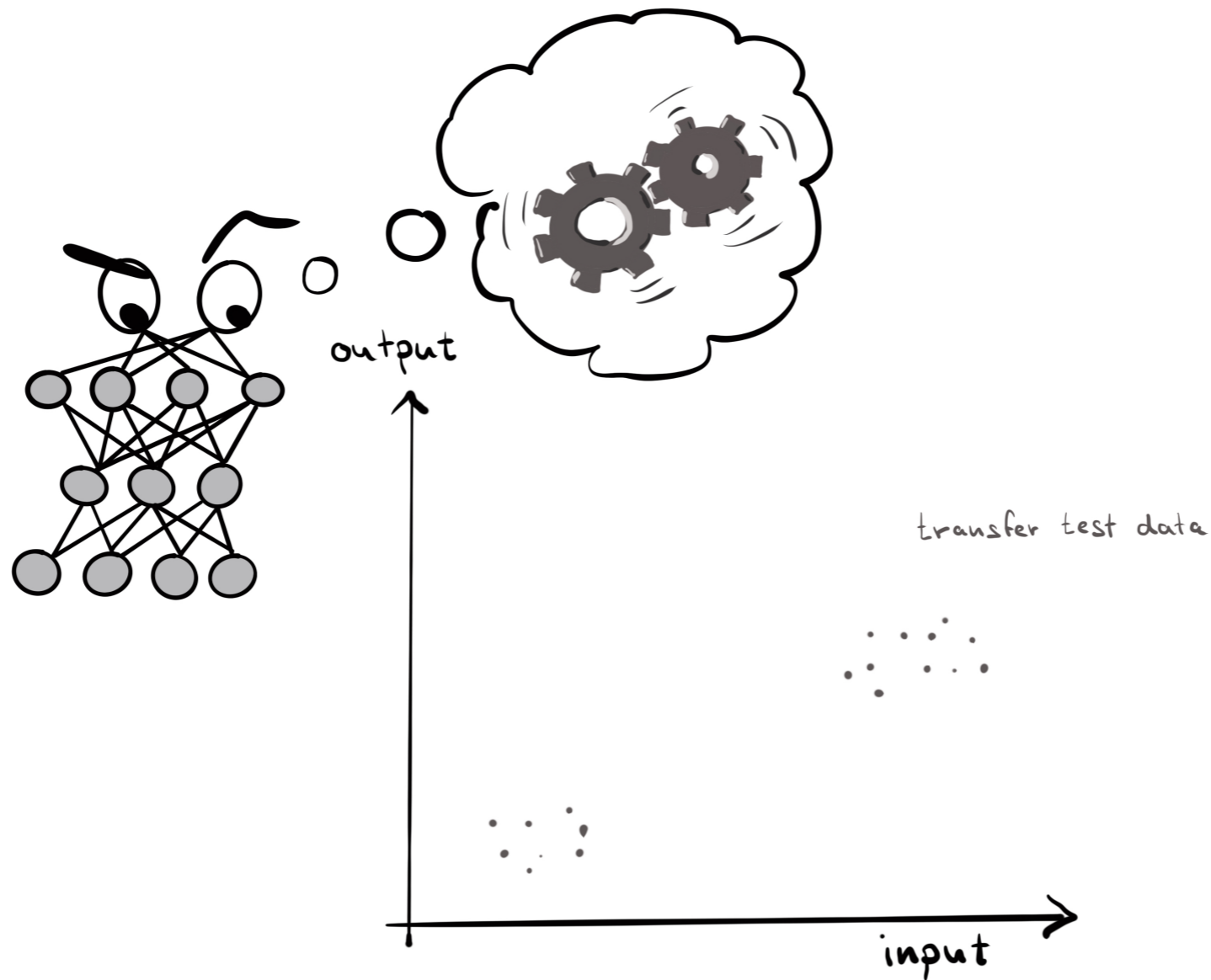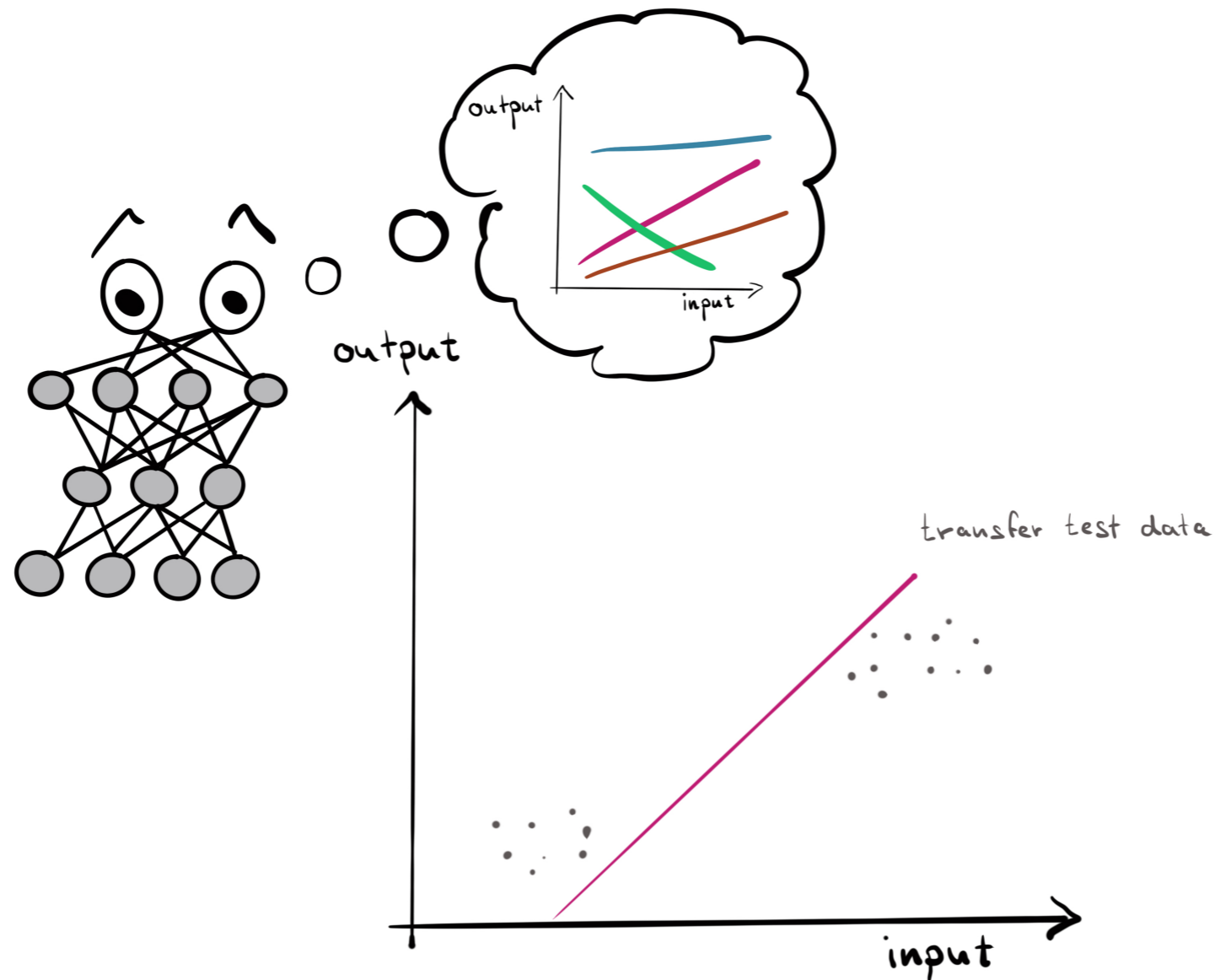
# Inductive Bias - Implicit Assumptions

# Inductive Bias - Implicit Assumptions
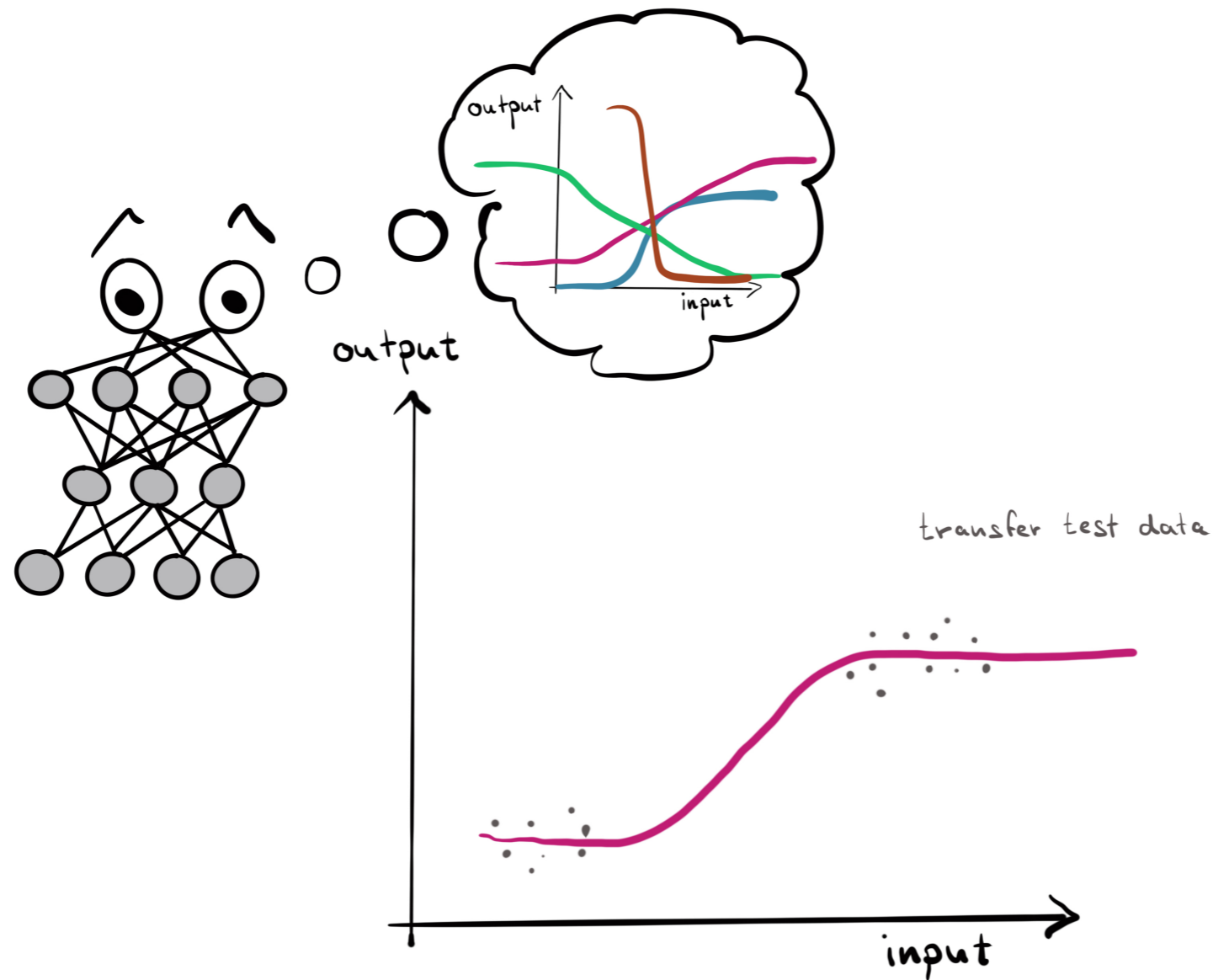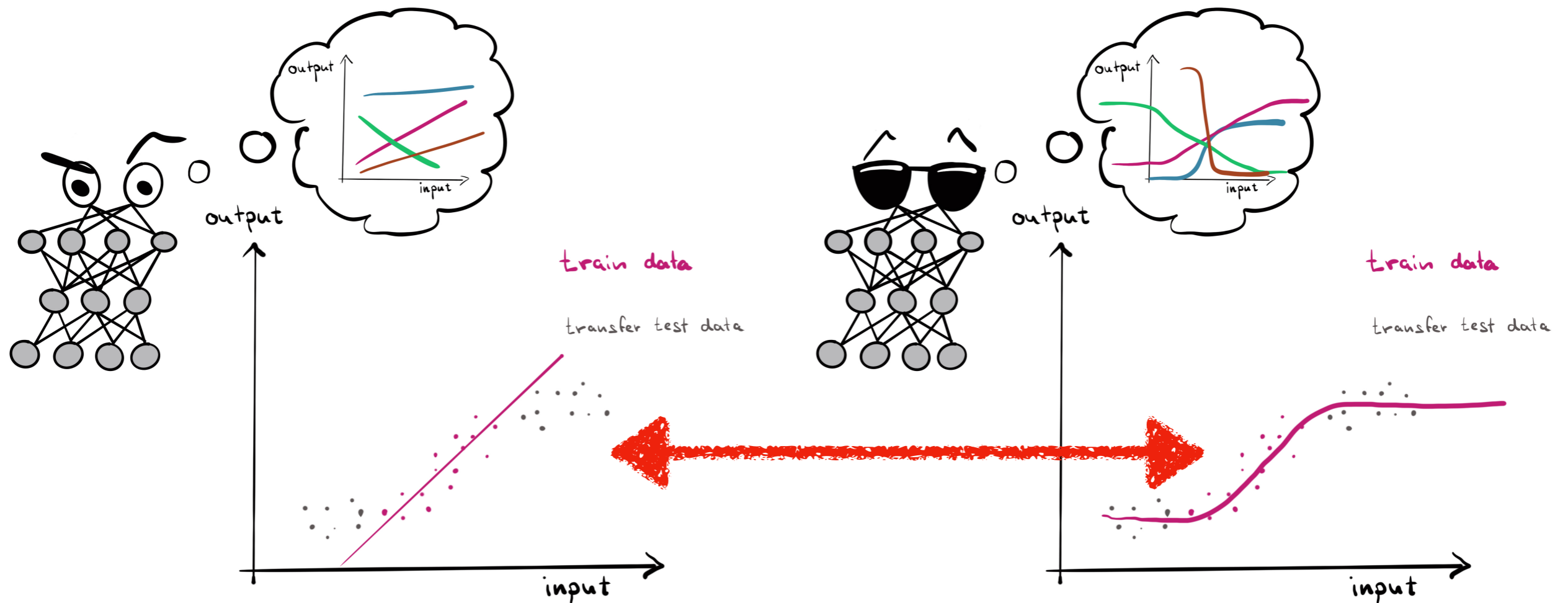
# Inductive Bias - Implicit Assumptions

# Inductive Bias - Implicit Assumptions

Differences in extrapolation between two algorithms given the **same** training data.

# Inductive Bias is Essential for Generalization

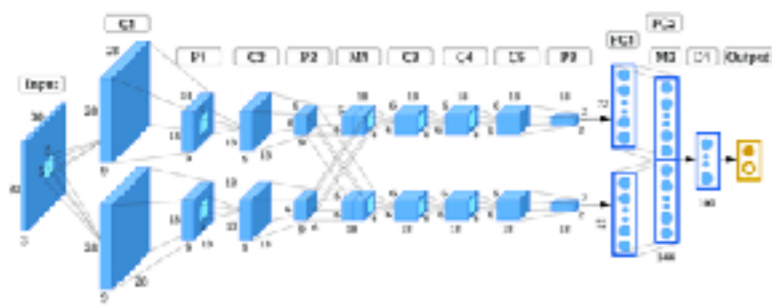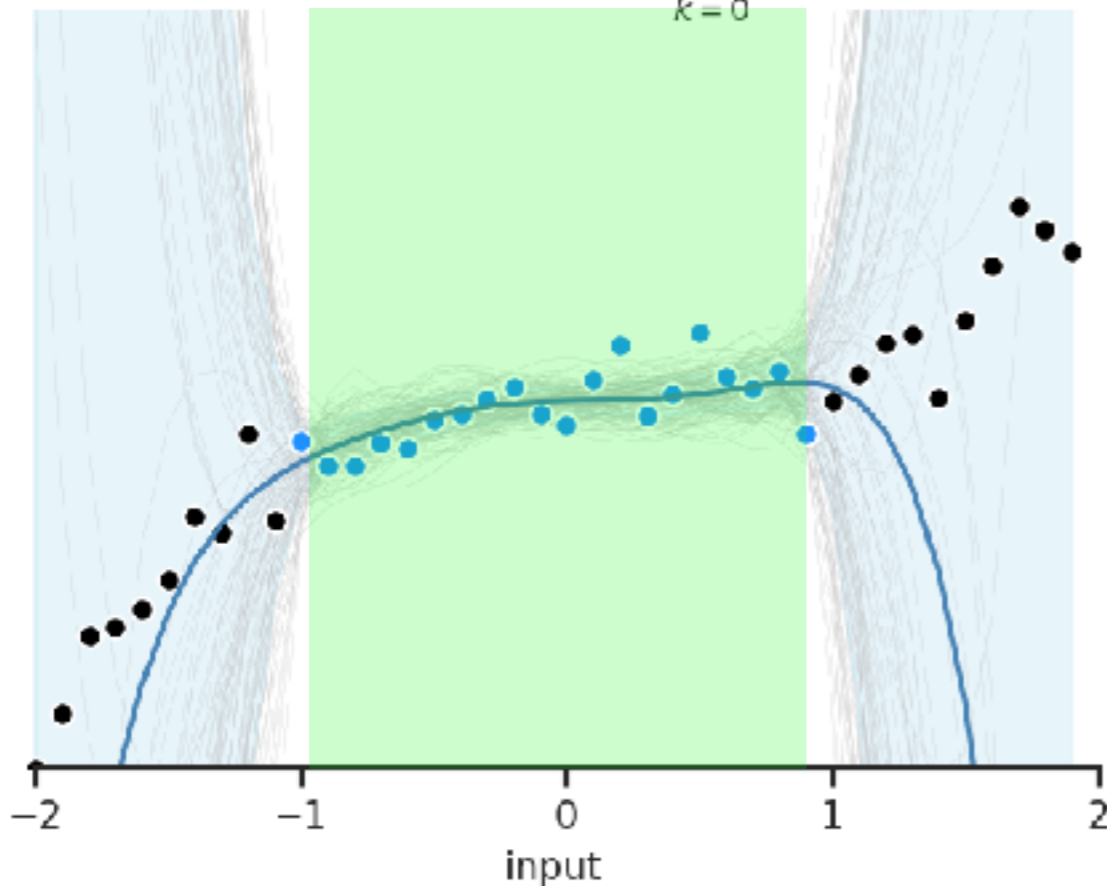No inductive bias
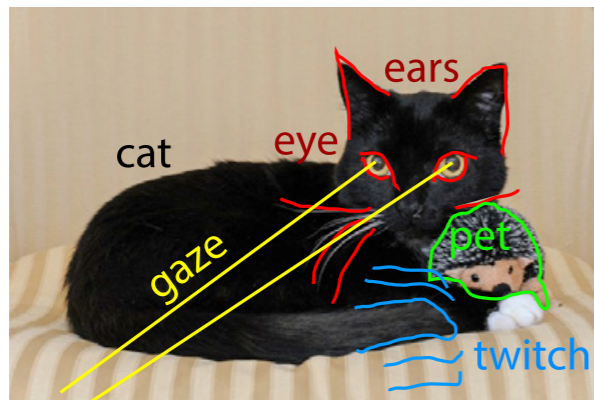
No free lunch

No generalization

# Inductive Bias



**less** ← bias → **more**

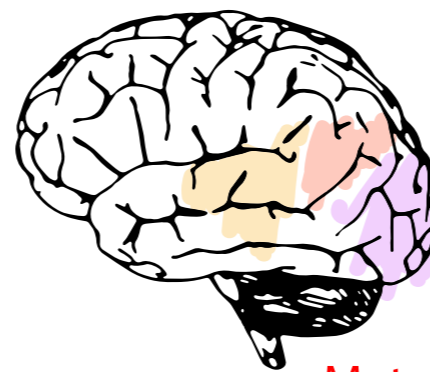unrestricted class $\sum_{k=0}^{6} a_k x^k$



input

# Levels of Inductive Bias Transfer
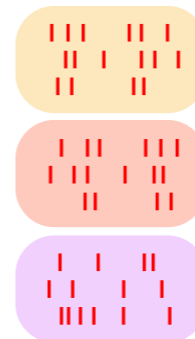


**Computational level**

Multi-task training:
act on causal variables
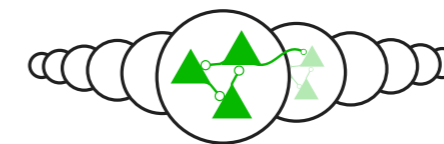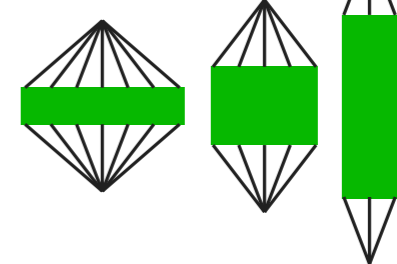
**Representational level**

Match neural
representations
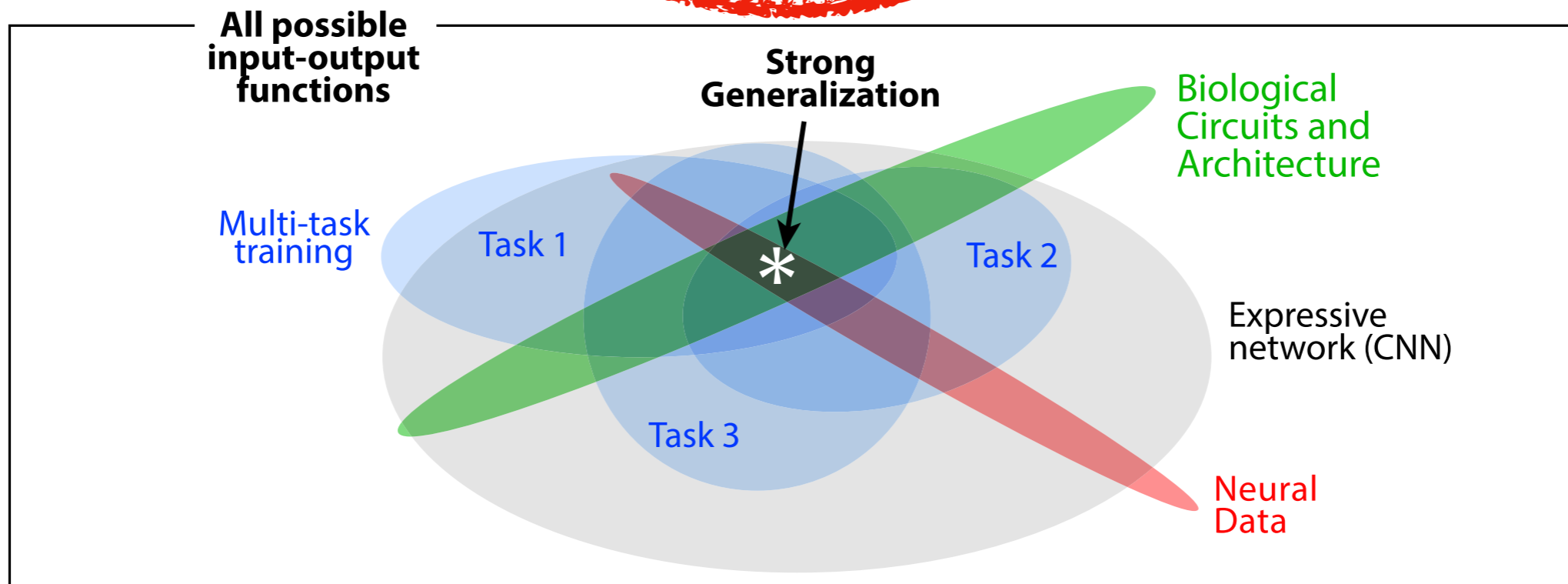of latent variables

**Implementational level**

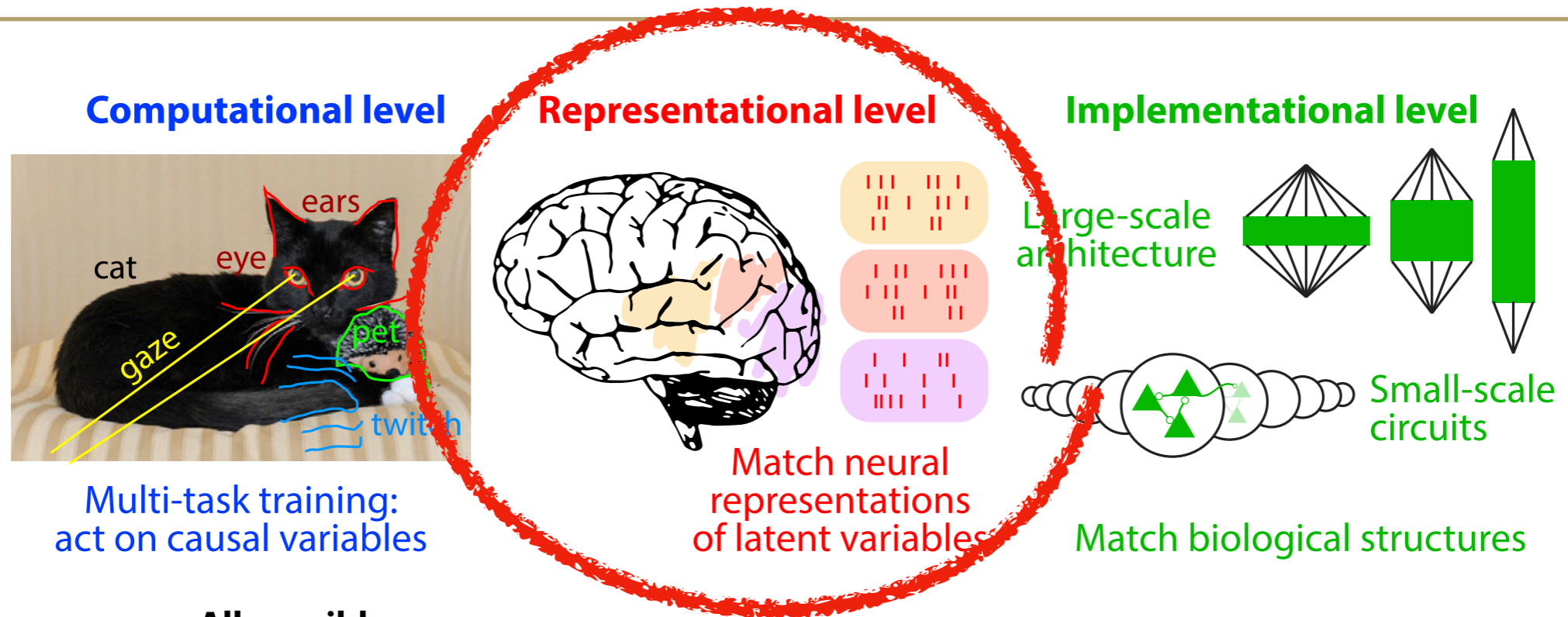Large-scale
architecture

Small-scale
circuits

Match biological structures

# How can we transfer good inductive biases?



Sinz et al 2019. "Engineering a Less Artificial Intelligence." *Neuron*.

# Neural co-training on monkey V1

Shahd Safarani

In collaboration with:

Arne Nix

Konstantin Willeke

Andreas Tolias, PhD
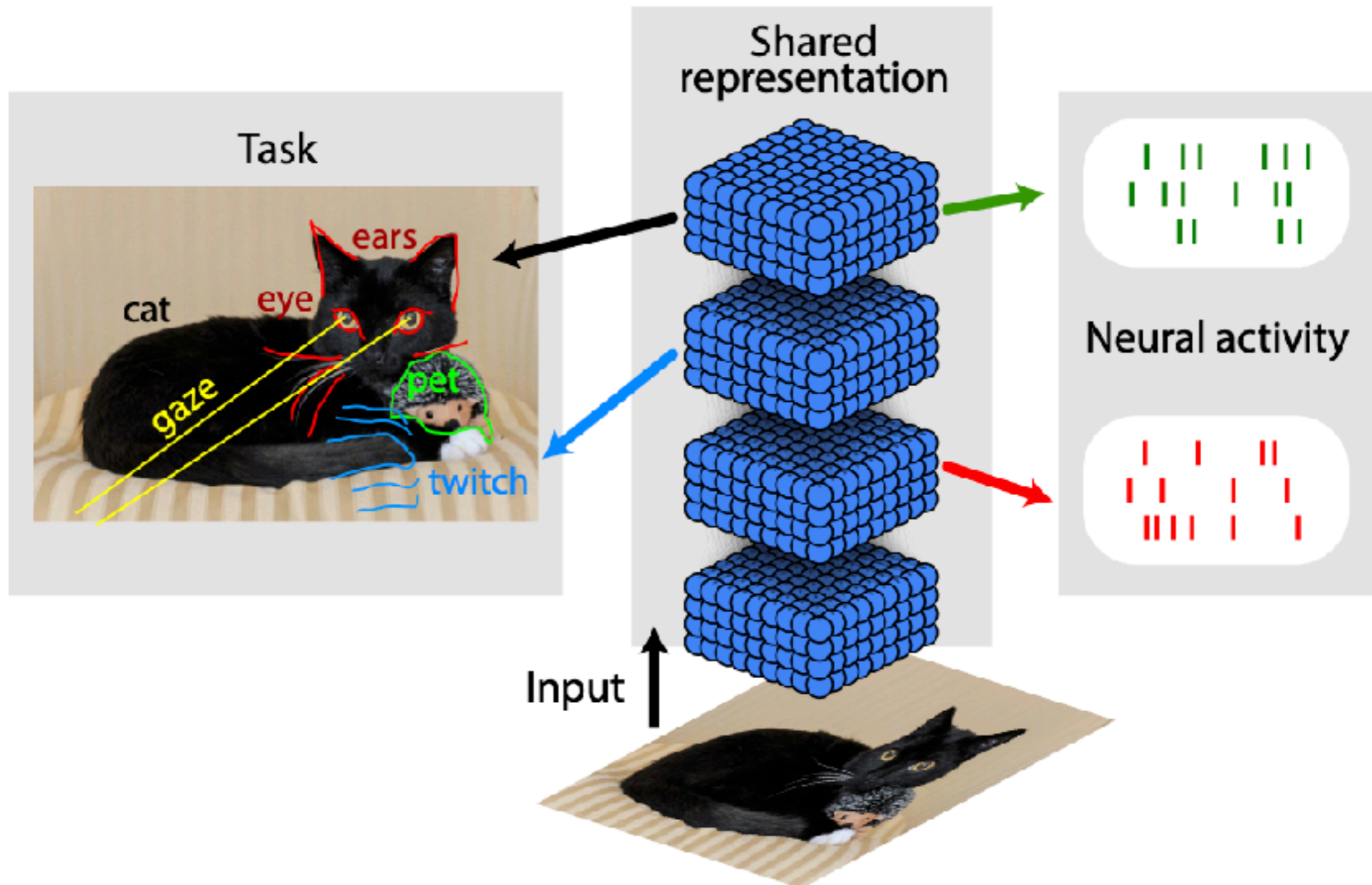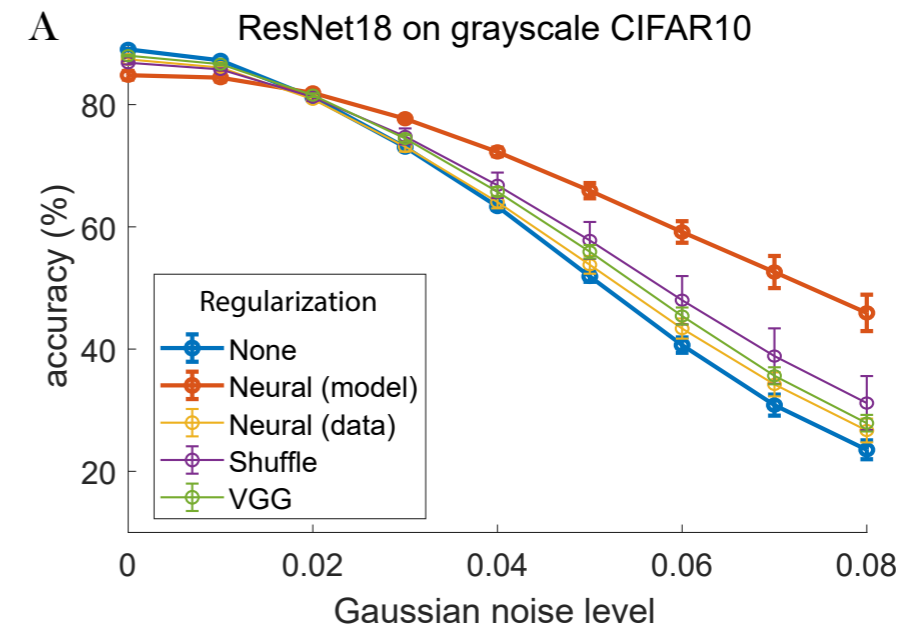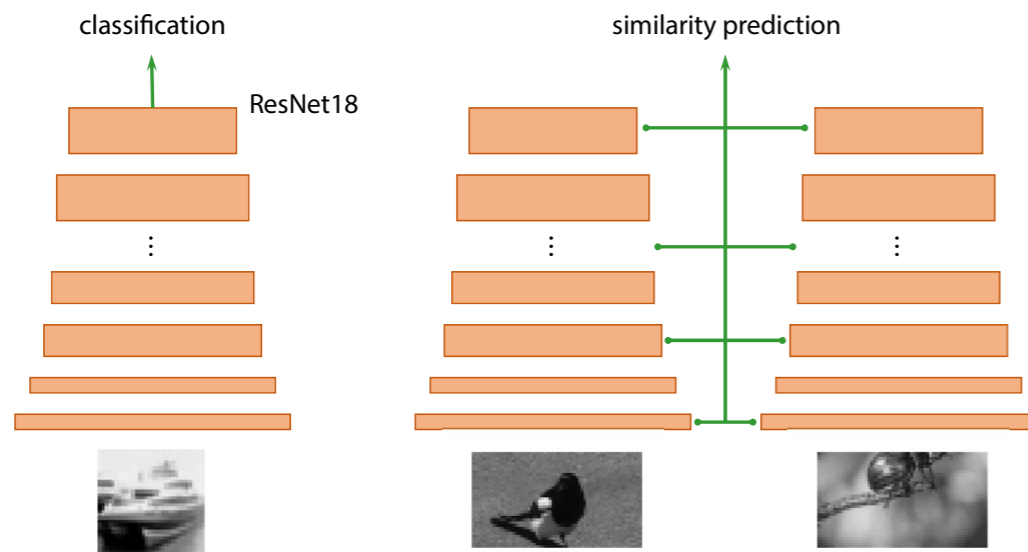
# Neural co-training hypothesis

# Do we expect it to work?

$$\text{loss} = \text{loss}_{\text{classification}} + \alpha\,\text{loss}_{\text{similarity}}$$

classification

ResNet18

similarity prediction

A

**ResNet18 on grayscale CIFAR10**

accuracy (%)

Regularization
- None
- Neural (model)
- Neural (data)
- Shuffle
- VGG

Gaussian noise level

Li et al. 2019. "Learning From Brains How to Regularize Machines." NeurIPS
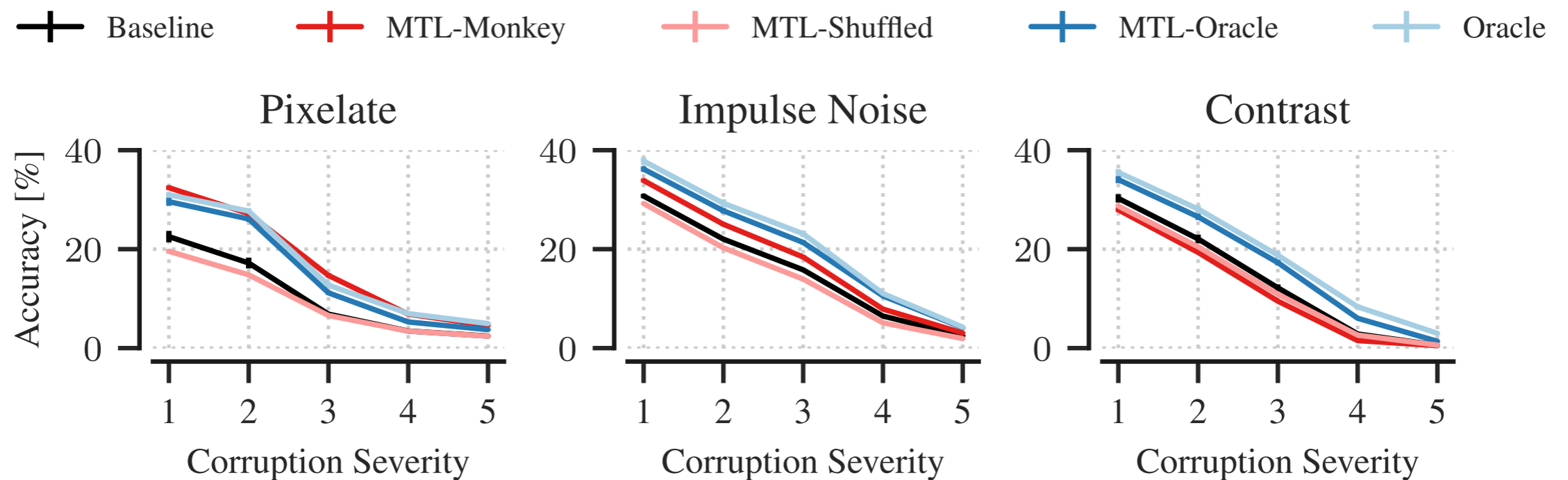
EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN

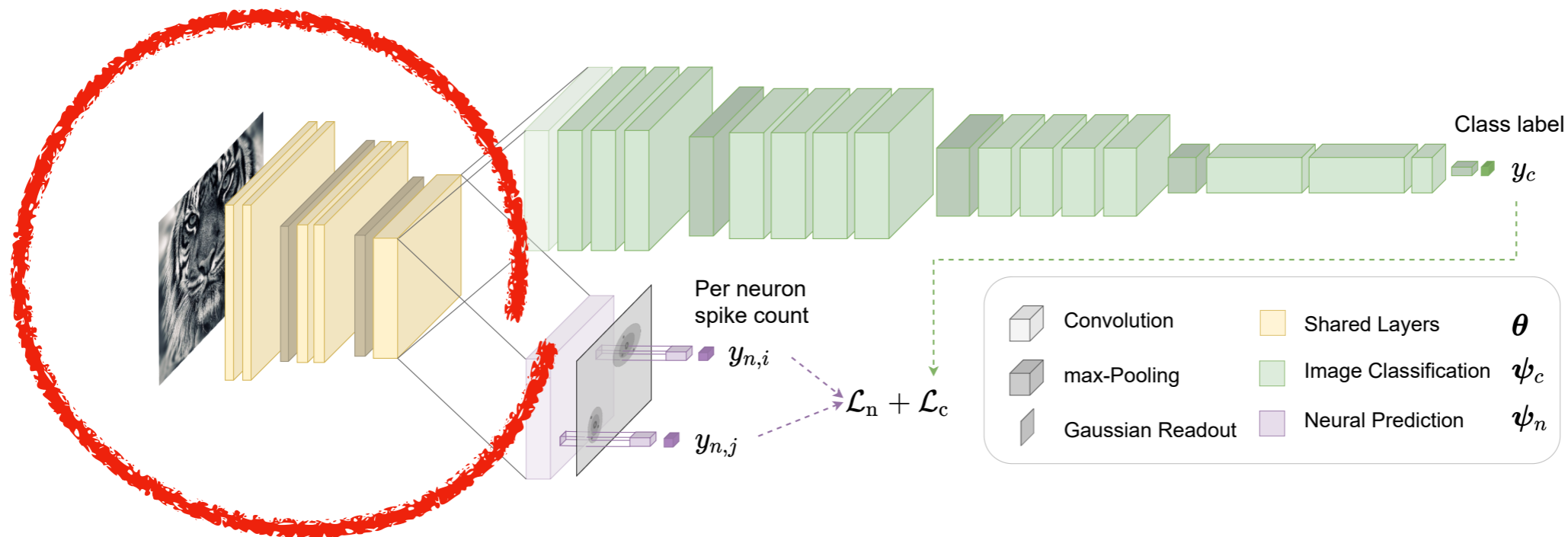# Multi-Task-Learning with monkey V1

# How do we test inductive bias ?

# V1 co-training yields benefits

# V1 co-training yields benefits

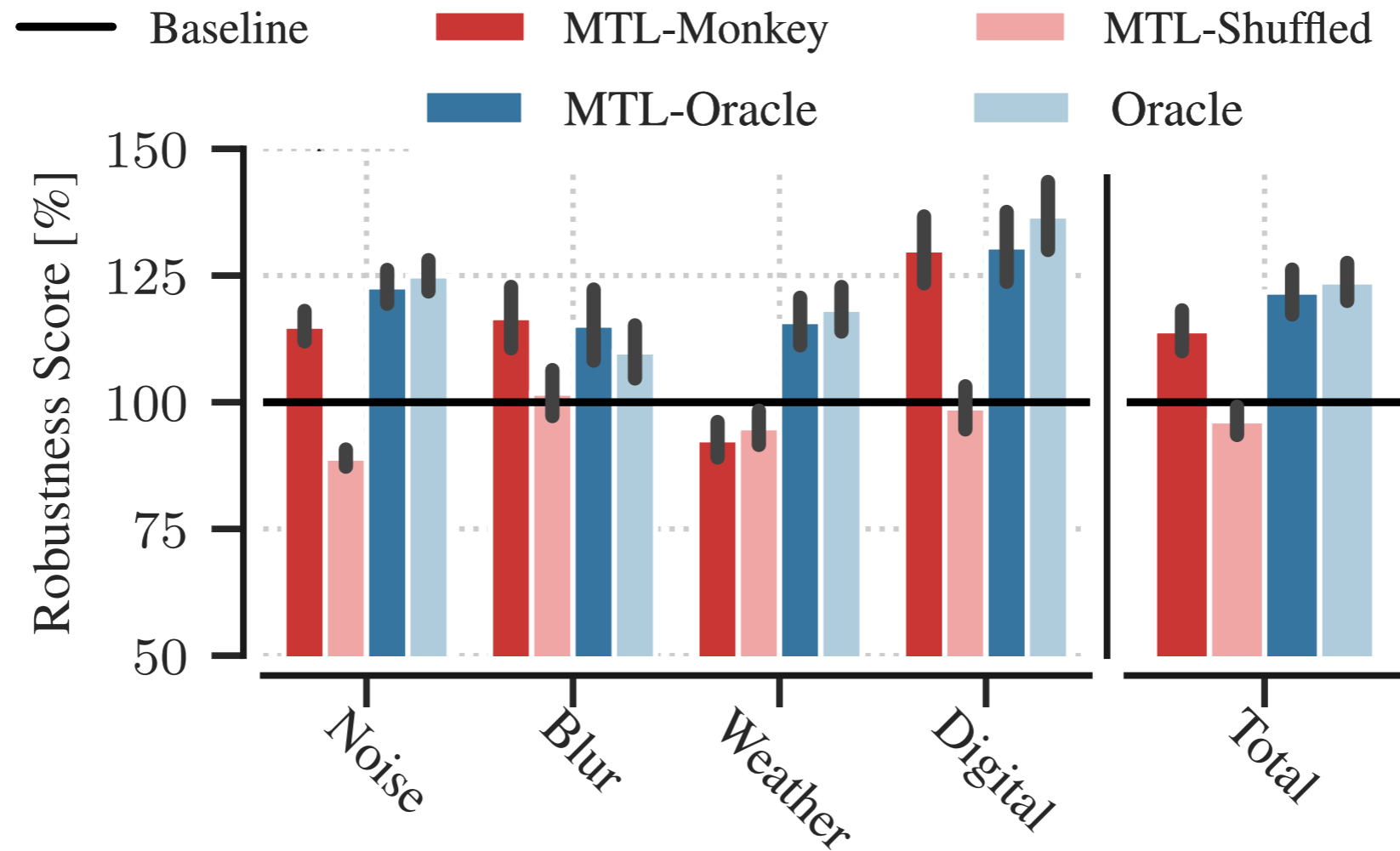# V1 co-training yields benefits

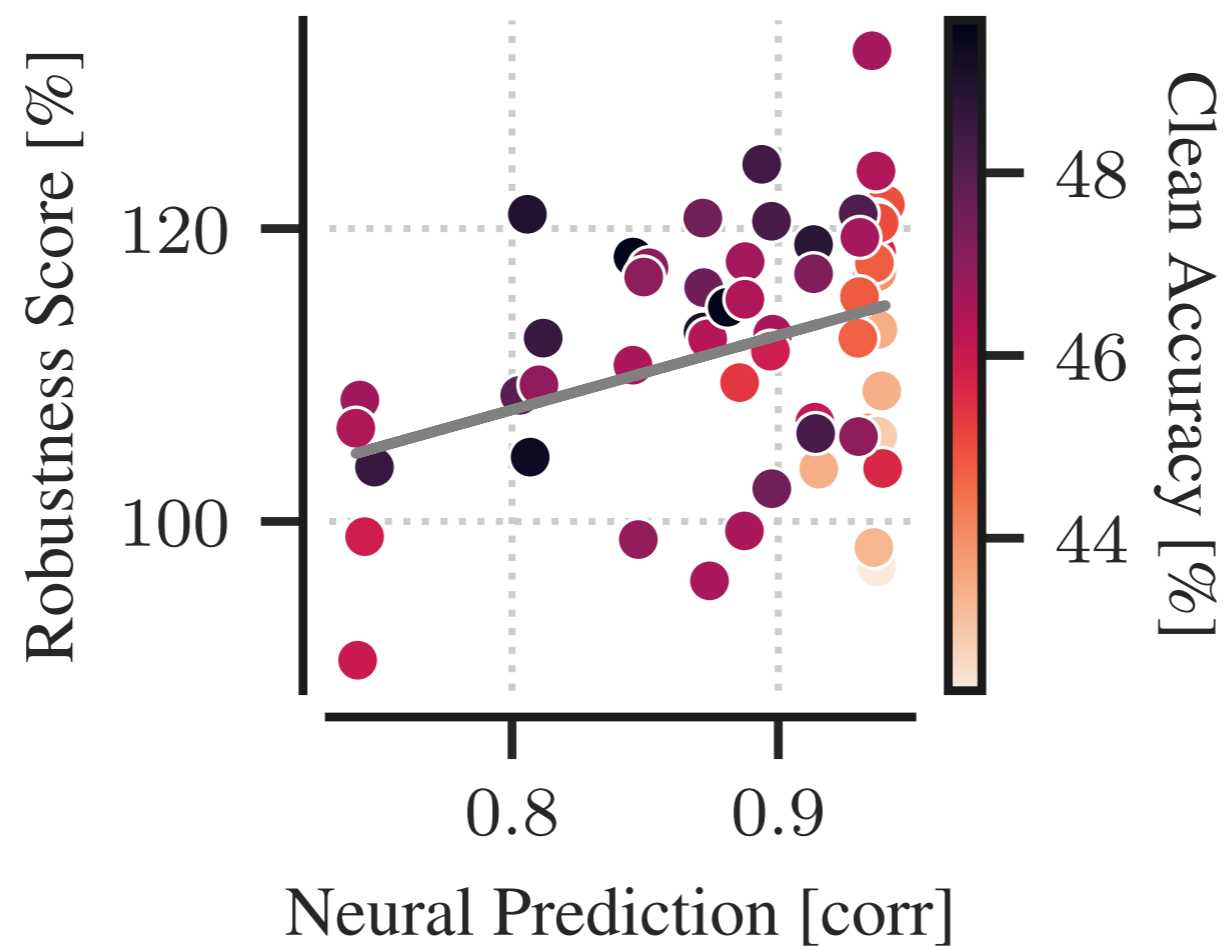# No all distortions show the same effect



$$\frac{1}{C}\sum_{c=1}^{C}\frac{A_c^{\textbf{robust}}}{A_c^{\textbf{baseline}}}$$

$$A_c = \frac{1}{5}\sum_{l,s=1}^{5}A_{l,c,s}$$

# Robustness correlates with "brain-likeness"

# Summary

- Mammalian visual systems have a better inductive bias than deep networks

- Multi task learning can be one avenue to improve inductive biases of models

- Co-training on monkey V1 yields improves robustness classification models

- Brain-likeness correlates with robustness

# Funding

# Thanks for listening! Questions?



We are looking for PhD students!

Check out: https://sinzlab.org/openpositions.html
or scan code



EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN