



Matteo Rovere

Giovanni Paolini

Glen Mbeng

Stefano Soatto

Information in the Weights and Emergent Properties of Deep Neural Networks

Alessandro Achille

University of California, Los Angeles

Menu

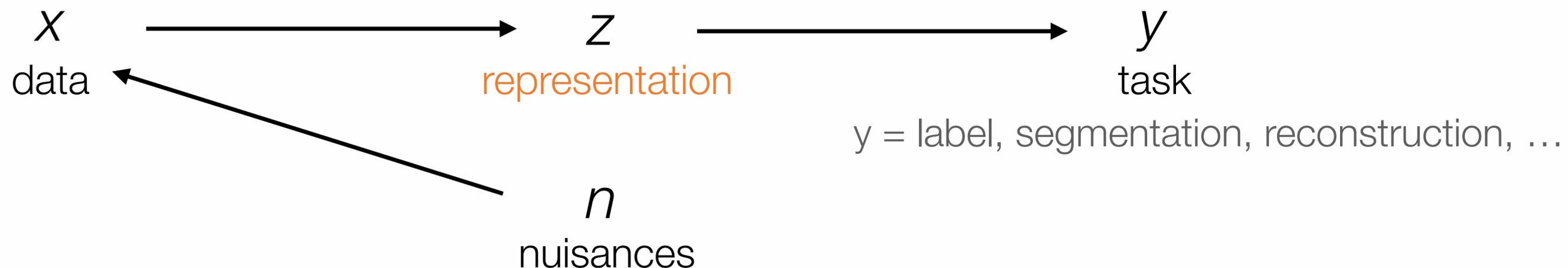
Two representations are important in machine learning:

- The representation of the test image (**activations**)
- The representation of the train dataset (**weights**)

In a DNN, the optimality of one implies optimality of the other (**Emergence Bound**)

Information gives an **asymmetric distance** distance on the space of learning tasks, which allows prediction of transfer learning.

What is an optimal representation (of test data)?



Sufficient

$$I(z; y) = I(x; y)$$

Nuisance invariance

$$n \perp y \Rightarrow I(n; z) = 0$$

Minimal

$$I(x; z) = \text{minimal}$$

Compositional

Minimal component correlation?

A Variational Principle for representation learning: The Information Bottleneck principle

A minimal sufficient representation \mathbf{z} of the data \mathbf{x} for the task \mathbf{y} is the solution to:

$$\begin{aligned} & \text{minimize}_{p(z|x)} \quad I(z; x) \\ & \text{s.t.} \quad I(z; y) = I(x; y) \end{aligned}$$

Information Bottleneck Lagrangian: (Tishby et al., 1999)

$$\mathcal{L}(p(z|x)) = \underbrace{H(y|z)}_{\text{cross-entropy}} + \beta \underbrace{I(z;x)}_{\text{regularizer}}$$

Invariant if and only if minimal

Recall. A representation z is minimal for the task y if it minimizes $I(z; x)$ among the sufficient representations.

Theorem (A., Soatto) Let z be a sufficient representation and n a nuisance. Then,

$$I(z; n) \leq I(z; x) - I(x; y)$$

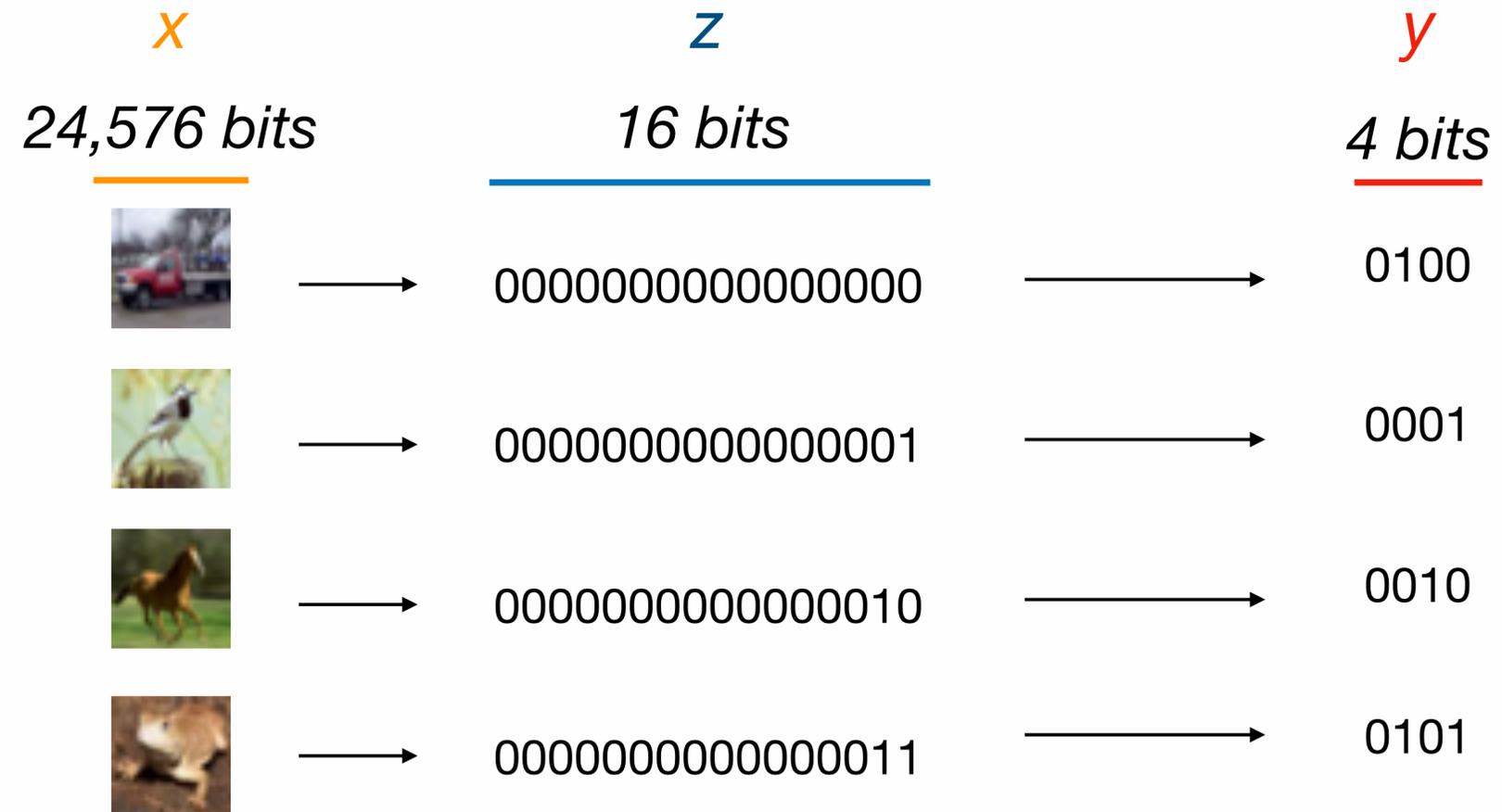
invariance minimality constant

Moreover, there exists a nuisance n for which equality holds.

Corollary: A representation is maximally invariant if and only if it is minimal

The catch

What if we just represent an image by its index in the training set (or by a unique hash)?



It is a sufficient representation and it is close to minimal.

This Information Bottleneck is wishful thinking

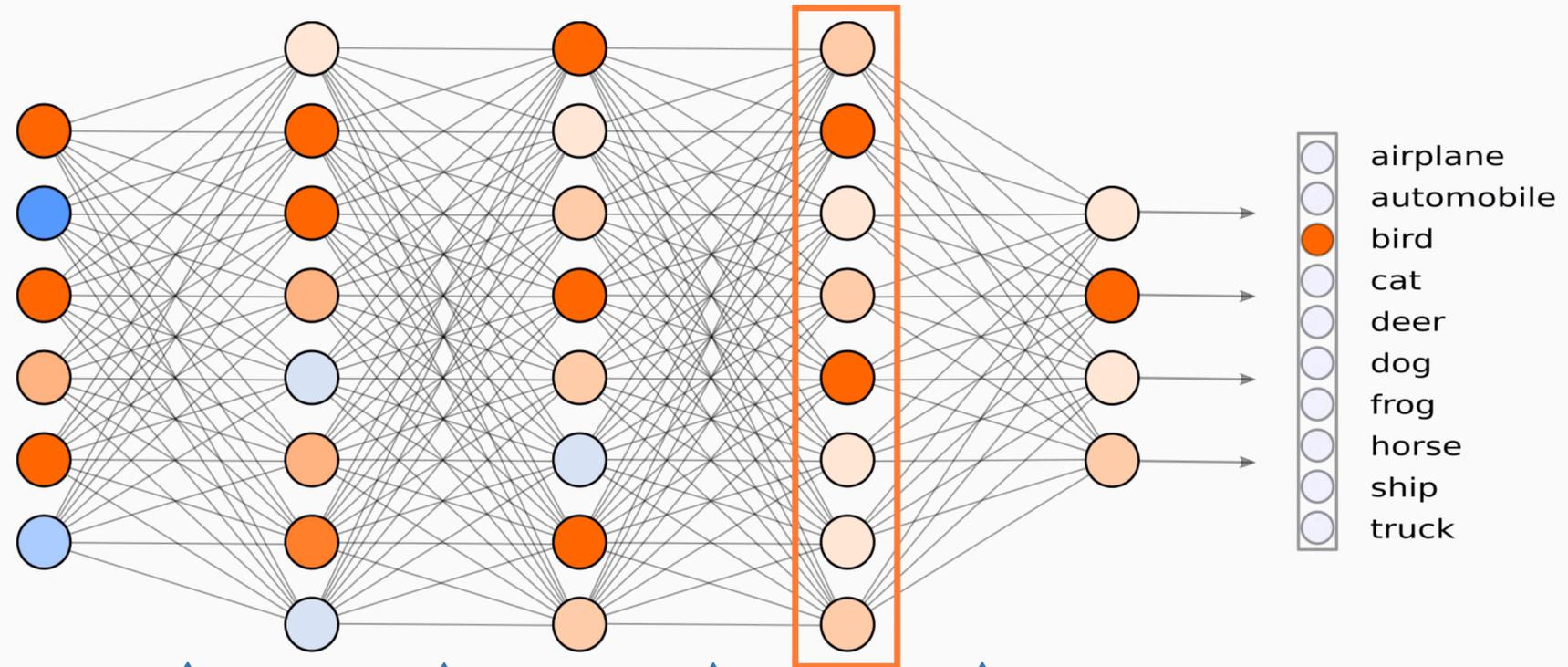
The IB is a **statement of desire** for future data we do not have:

$$\min_{q(z|x)} \mathcal{L} = H_{p,q}(y|z) + \beta I(z; x)$$

What we have is the data collected in the past.

What is the best way to use the past data in view of future tasks?

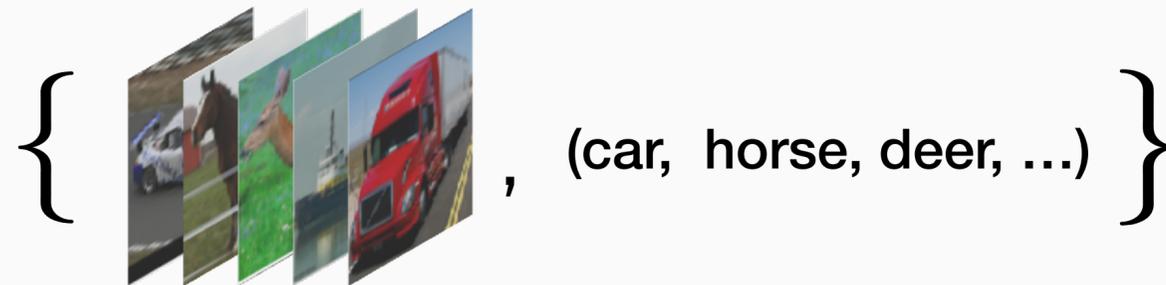
Test Image



Weights

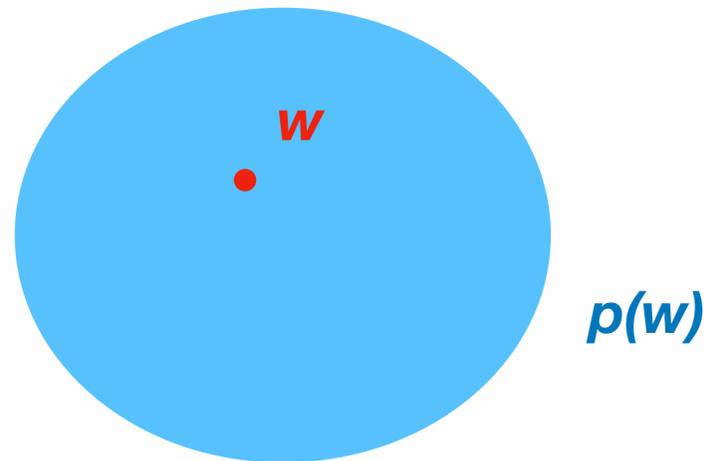
Representation of past data

Training Set



How do we measure the information in the weights of a DNN?

Problem. Assume some prior distribution $p(w)$ over the weights. Codifying a particular set of weights as real numbers requires *infinite information*.



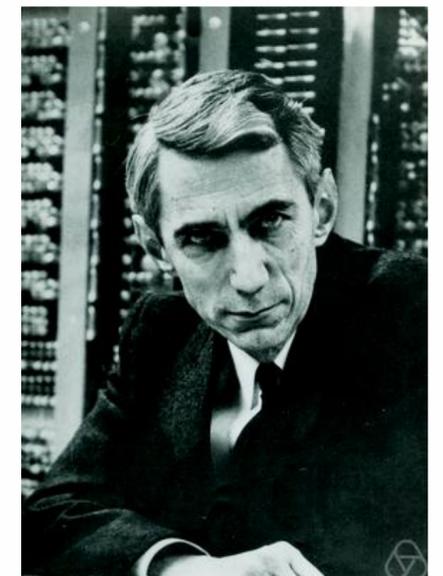
Example: Measuring Information by Adding Noise

Estimate the amount of information by adding noise and measuring the decrease in performance.

Prediction and Entropy of Printed English

By C. E. SHANNON

(Manuscript Received Sept. 15, 1950)



Example: Shannon (1951) estimates the information content of the English language by corrupting random letters and measuring the reconstruction error of English speakers.

“Thif is a veyy moisy party” → “This is a very noisy party”

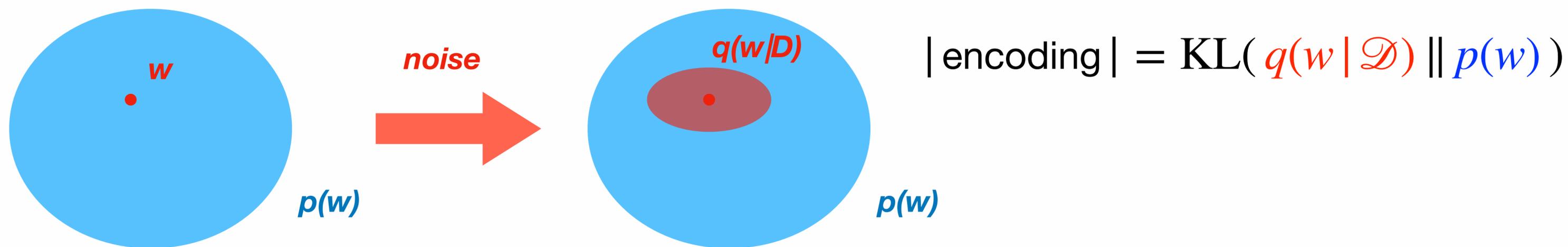
Reducing the description length by adding noise

Keeping Neural Networks Simple by Minimizing the Description Length of the Weights

Geoffrey E. Hinton and Drew van Camp
Department of Computer Science
University of Toronto



Idea: Add noise to the weights to encode with a finite amount of information



The Information in the Weights

We want to measure the trade-off between the amount $\text{KL}(q(w | \mathcal{D}) \parallel p(w))$ of noise added and the accuracy of the network.

$$S(t) = \min \text{KL}(q(w | \mathcal{D}) \parallel p(w)) \quad \leftarrow \text{Minimum information in weights...}$$
$$\text{s.t. } \mathbb{E}_{w \sim q(w | \mathcal{D})} [L_{\mathcal{D}}(w)] < t \quad \leftarrow \text{... such that the expected training loss is less than } t$$

Or, equivalently the Lagrangian:

$$\mathcal{L} = \mathbb{E}_{w \sim q(w | \mathcal{D})} [L_{\mathcal{D}}(w)] + \beta \text{KL}(q(w | \mathcal{D}) \parallel p(w))$$

Expected loss **Information in the Weights**

optimal noise fixed prior

Grounding: The PAC-Bayes generalization bound

PAC-Bayes bound on the test error: (Catoni, 2007; McAllester 2013)

$$L_{\text{test}} \leq \frac{1}{1 - \frac{1}{2\beta}} \left[\mathbb{E}_w [L_{\mathcal{D}}(w)] + \beta \text{KL}(q(w|\mathcal{D}) \parallel p(w)) \right]$$

What matters for generalization is not the number of weights, but the information they contain.

This gives **non-vacuous** generalization bounds.

The sharpest PAC-Bayes bound is the IB Lagrangian

Which encoding prior $p(w)$ gives the tightest PAC-Bayes bound?

Proposition. On expectation over the sampling of the dataset \mathcal{D} , the sharpest PAC-Bayes bound is obtained when $p(w) = \mathbb{E}_{\mathcal{D}}[q(w | \mathcal{D})]$, in which case $\mathbb{E}_{\mathcal{D}}[\text{KL}(q(w | \mathcal{D}) || p(w))] = I(w; \mathcal{D})$.

The Weight Lagrangian then becomes:

$$\mathcal{L}(q(w | \mathcal{D})) = \mathbb{E}_{\mathcal{D}}[H(\mathcal{D} | w)] + \beta I(w; \mathcal{D})$$

IB Lagrangian for the weights

The Information in a Deep Neural Network

$$L(w) = H_{p,q}(\mathcal{D}|w) + \beta \text{KL}(\underbrace{q(w|\mathcal{D})}_{\text{output of training}} \parallel \underbrace{p(w)}_{\text{encoding prior}})$$

Shannon Information: using the encoding prior $p(w) := \mathbb{E}_{\mathcal{D}}[q(w|\mathcal{D})]$

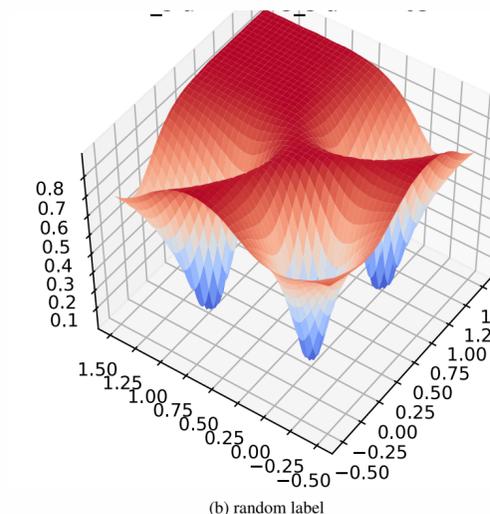
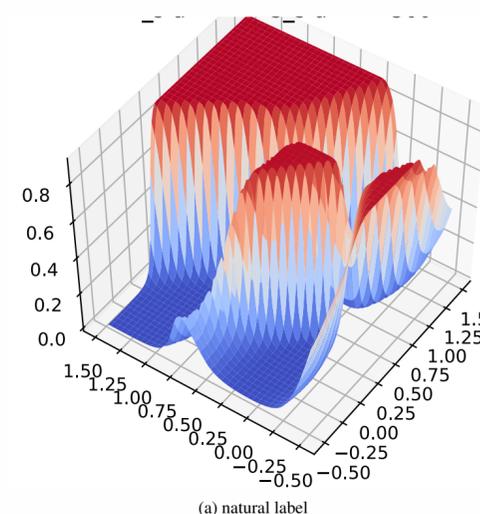
$$\mathbb{E}_{\mathcal{D}}[\text{KL}] = I(w; \mathcal{D})$$

⇒ Minimizes the PAC-Bayes generalization bound

Fisher Information: $p(w) =$ uniform encoding prior

$$\text{KL} = -\log |F|$$

⇒ Implicitly minimized by **SGD**



SGD connects Fisher and Shannon

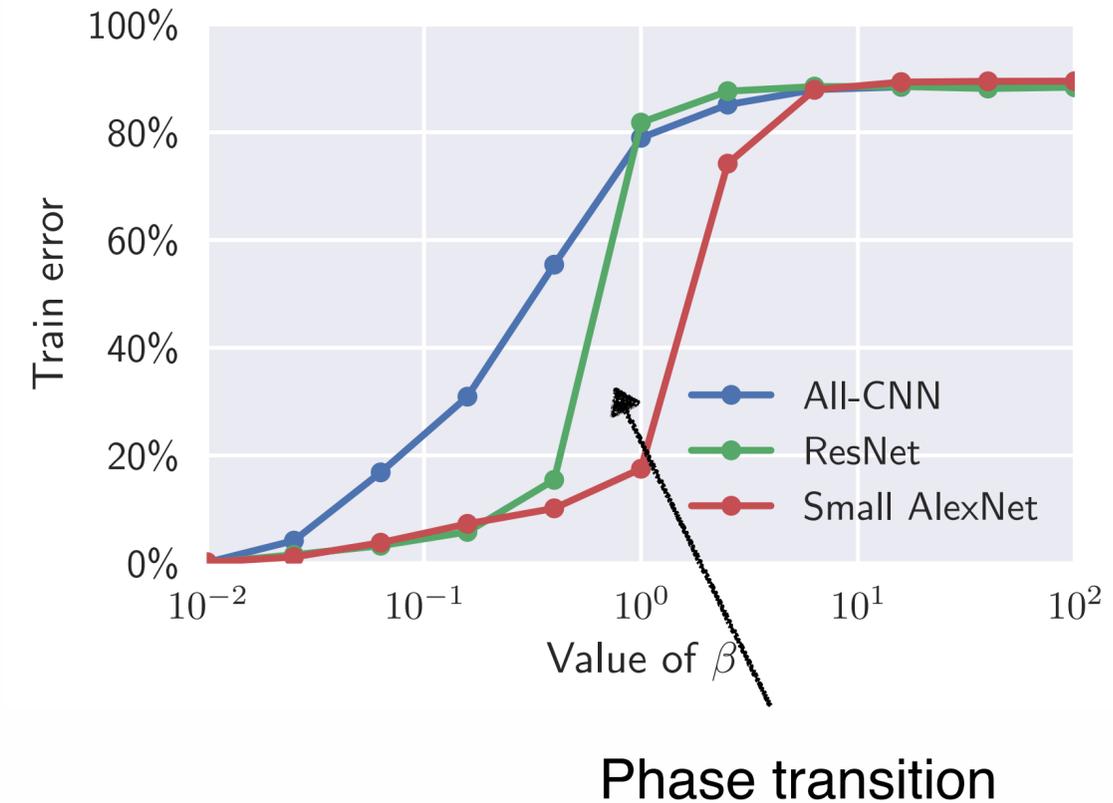
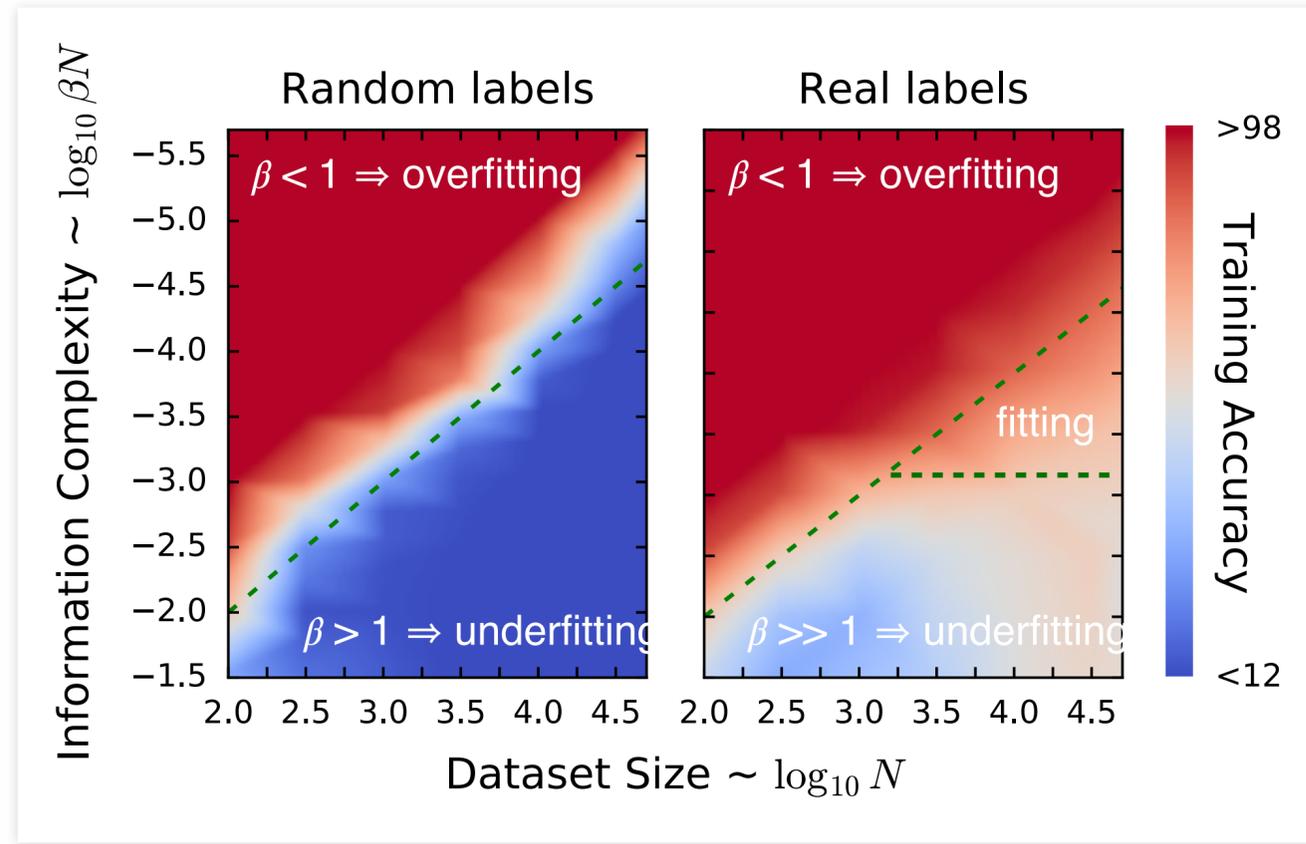
SGD minimizes the **Fisher Information** of the Weights. However, generalization is governed by the **Shannon Information**.

Proposition. Assuming the dataset is parametrized in a differentiable way, we have:

$$\underbrace{I(w; \mathcal{D})}_{\text{Shannon Information}} \approx H(\mathcal{D}) - \mathbb{E} \left[\log \left(\frac{(2\pi e)^k}{\underbrace{|\nabla_{\mathcal{D}} w^* F(w^*) \nabla_{\mathcal{D}} w^{*T}|}_{\text{Fisher Information}}} \right) \right]$$

Where $w^* = w^*(D)$ is the result of running SGD on dataset D and $F(w)$ is the Fisher Information Matrix in w .

Empirical verification: Phase transition



Using the regularized loss:

$$L(w) = H_{p,q}(\mathcal{D}|w) + \beta KL(q(w|\mathcal{D})||p(w))$$

For random labels there is a transition between over- and under-fitting at $\beta = 1$.

Two Bottlenecks

Activations IB
Invariance



$$\min_{q(z|x)} \mathcal{L} = H_{p,q}(y|z) + \beta I(z; x)$$

Weights IB
Generalization



$$\min_w \mathcal{L} = H_{p,q_w}(y|z) + \beta I(\mathcal{D}; w)$$

The Emergence Bound: connecting the two bottlenecks

Let $z = f_w(x)$ be a layer of a network, and let z_n be the representation obtained by adding noise to the weights. We define the **effective information** as $I_{\text{eff}}(x; z) = I(x; z_n)$

Proposition. To the first order the information in the activations is given by:

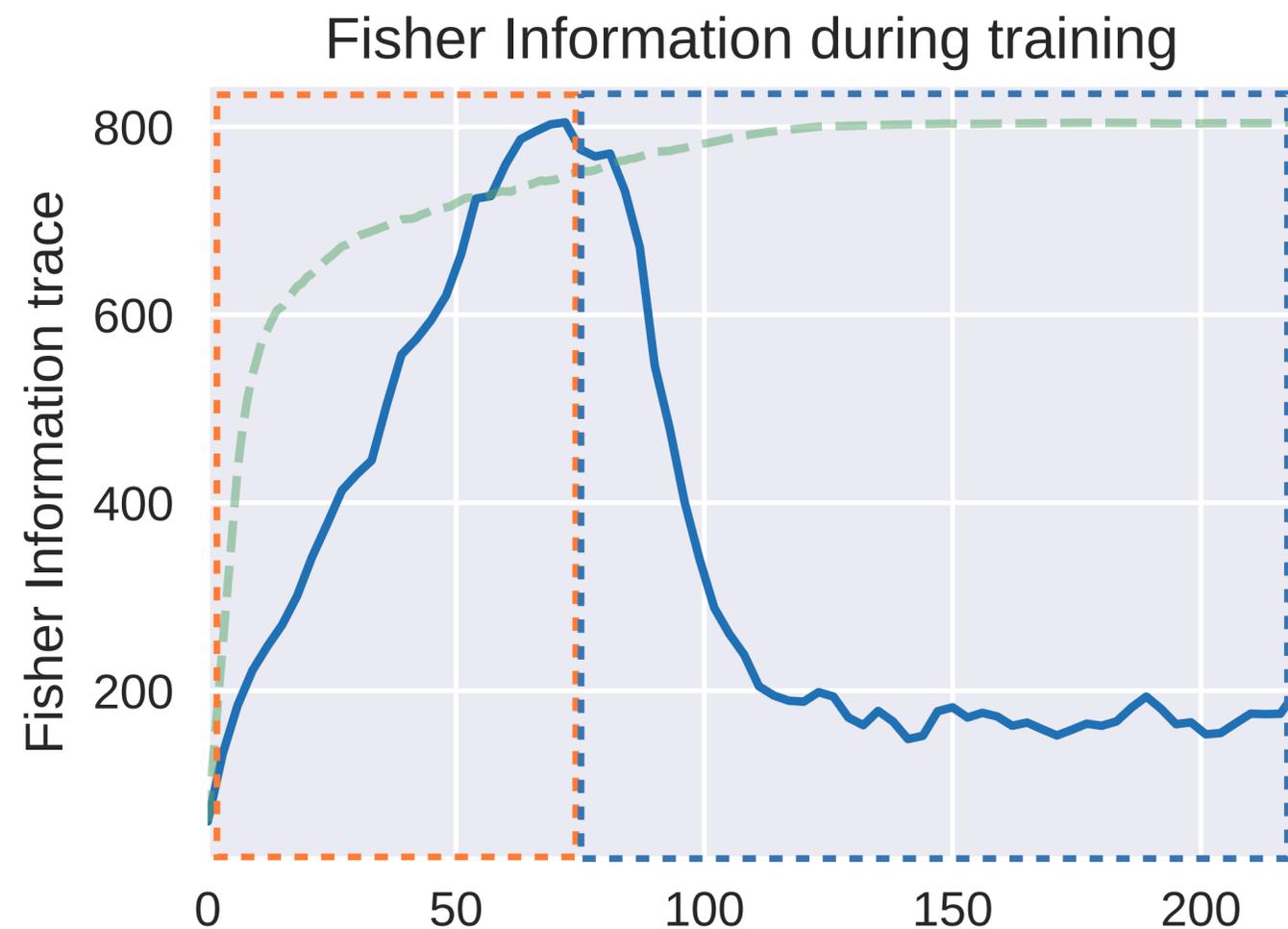
$$I_{\text{eff}}(x; z) \approx H(x) - \log \left(\frac{(2\pi e)^k}{|\nabla_x f_w(x)^t J_f^t \boxed{F(w)} J_f \nabla_x f_w(x)|} \right)$$

Information in activations Fisher Information in weights

Take-away: Reducing information in the weights reduces information in the activations, hence it promotes invariant classifiers.

Compression of weights leads to compression of activations

We empirically observe that the Fisher Information decreases later during training.

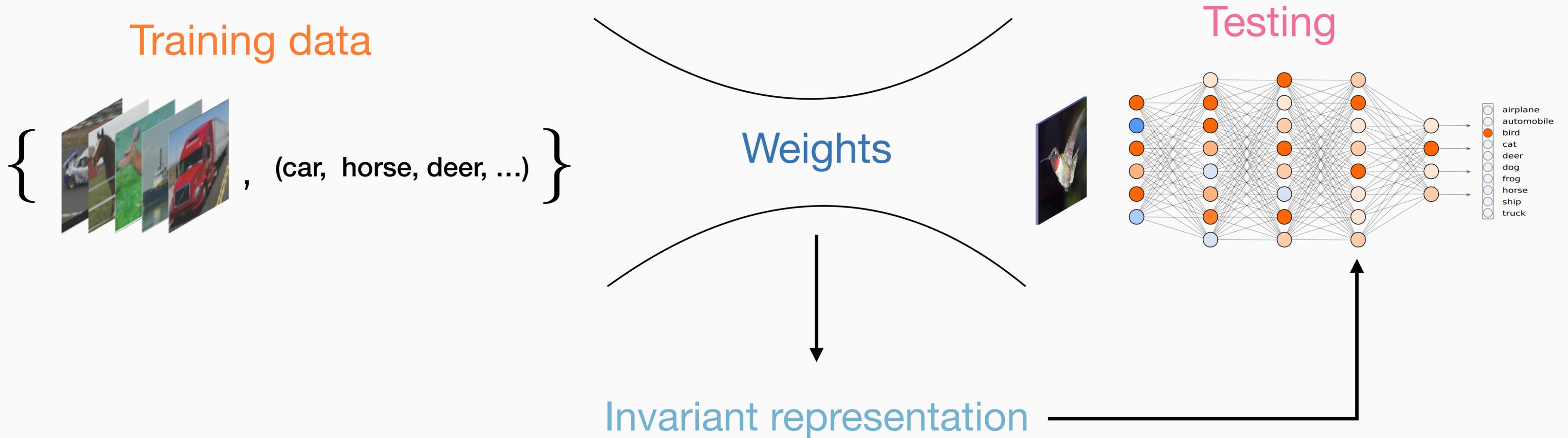


The (Fisher) information in the **weights** decreases later in training



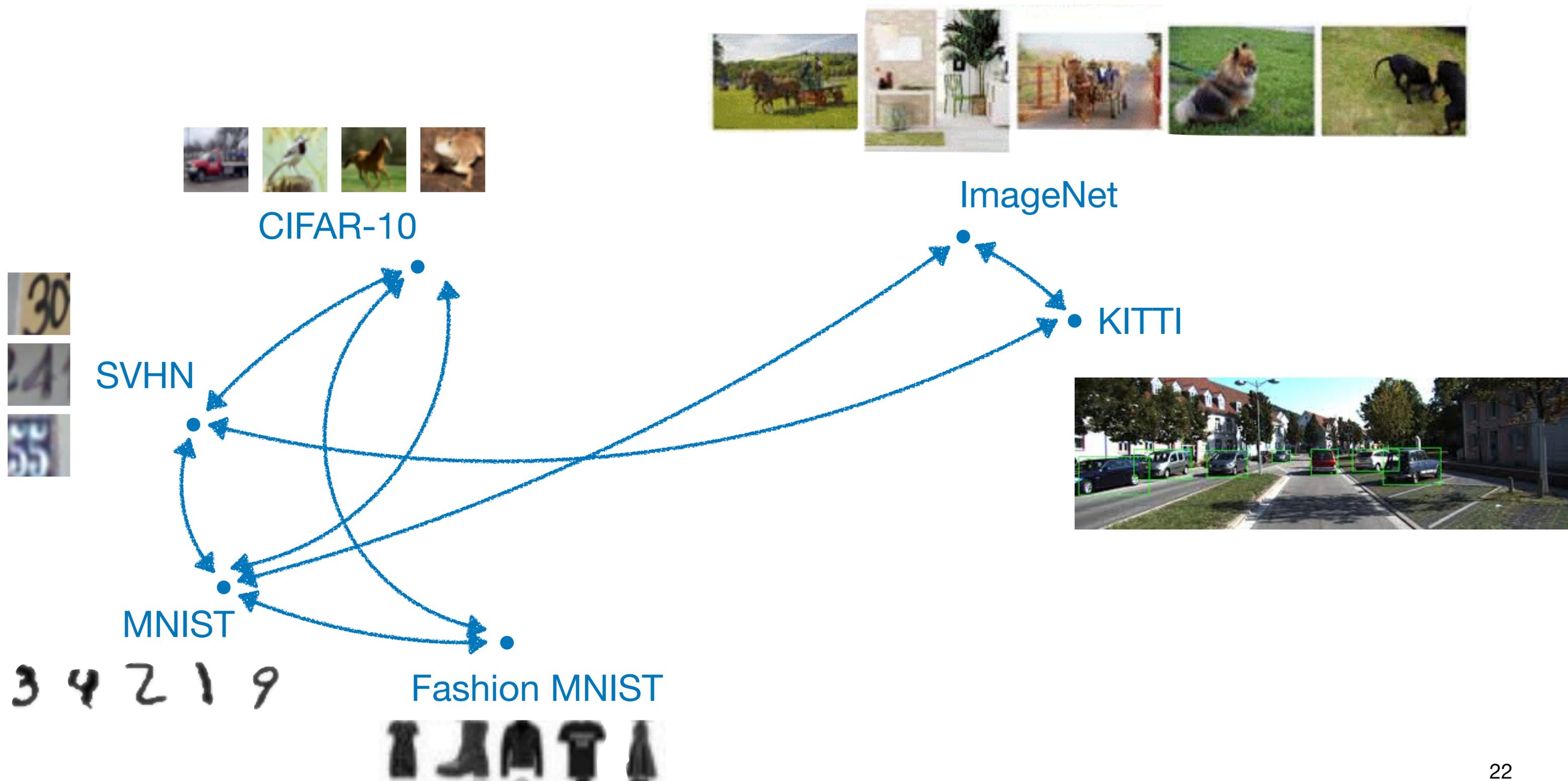
The effective information in the **activations** decreases

Compression of the **weights** biases toward invariant and disentangled **representations**.



PAST	FUTURE
Weights	Activations
Generalization (PAC-Bayes)	Invariance (Emergence)
Minimality (Shannon)	Minimality (Fisher)

What is the distance between two tasks?



A Topology on the Space of Tasks

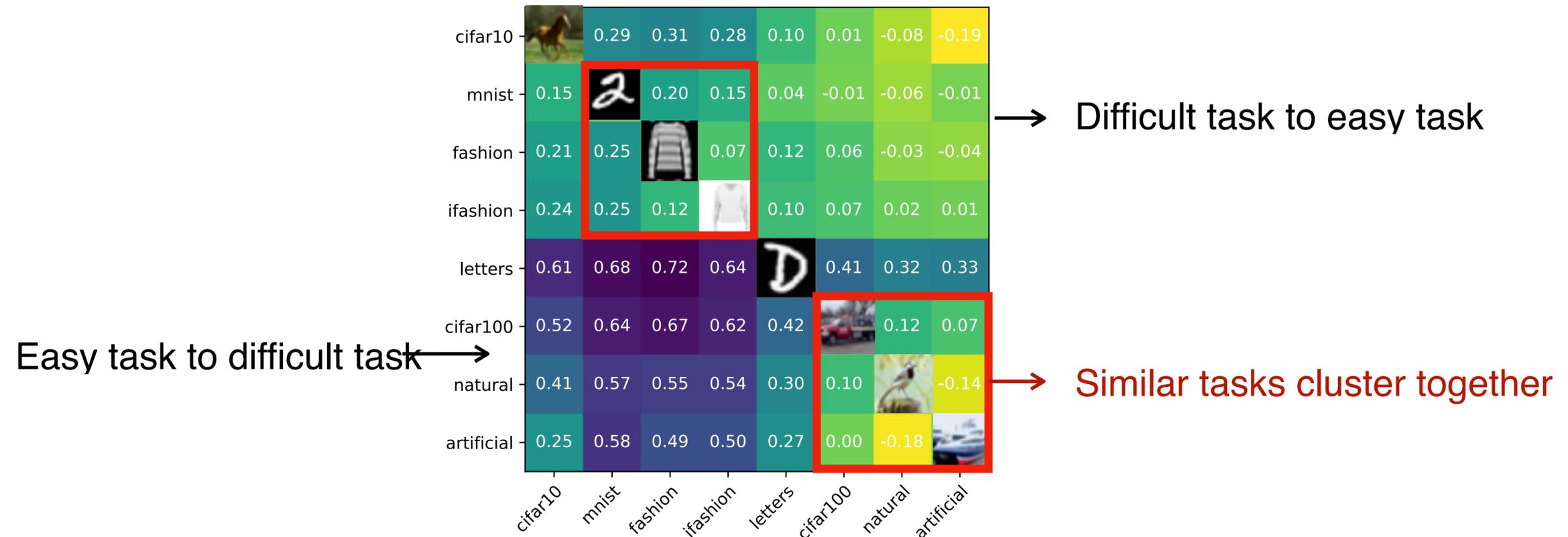
Distance between tasks:

$$d(\mathcal{D}_1 \rightarrow \mathcal{D}_2) = I(\mathcal{D}_1 \mathcal{D}_2; w) - I(\mathcal{D}_1; w)$$

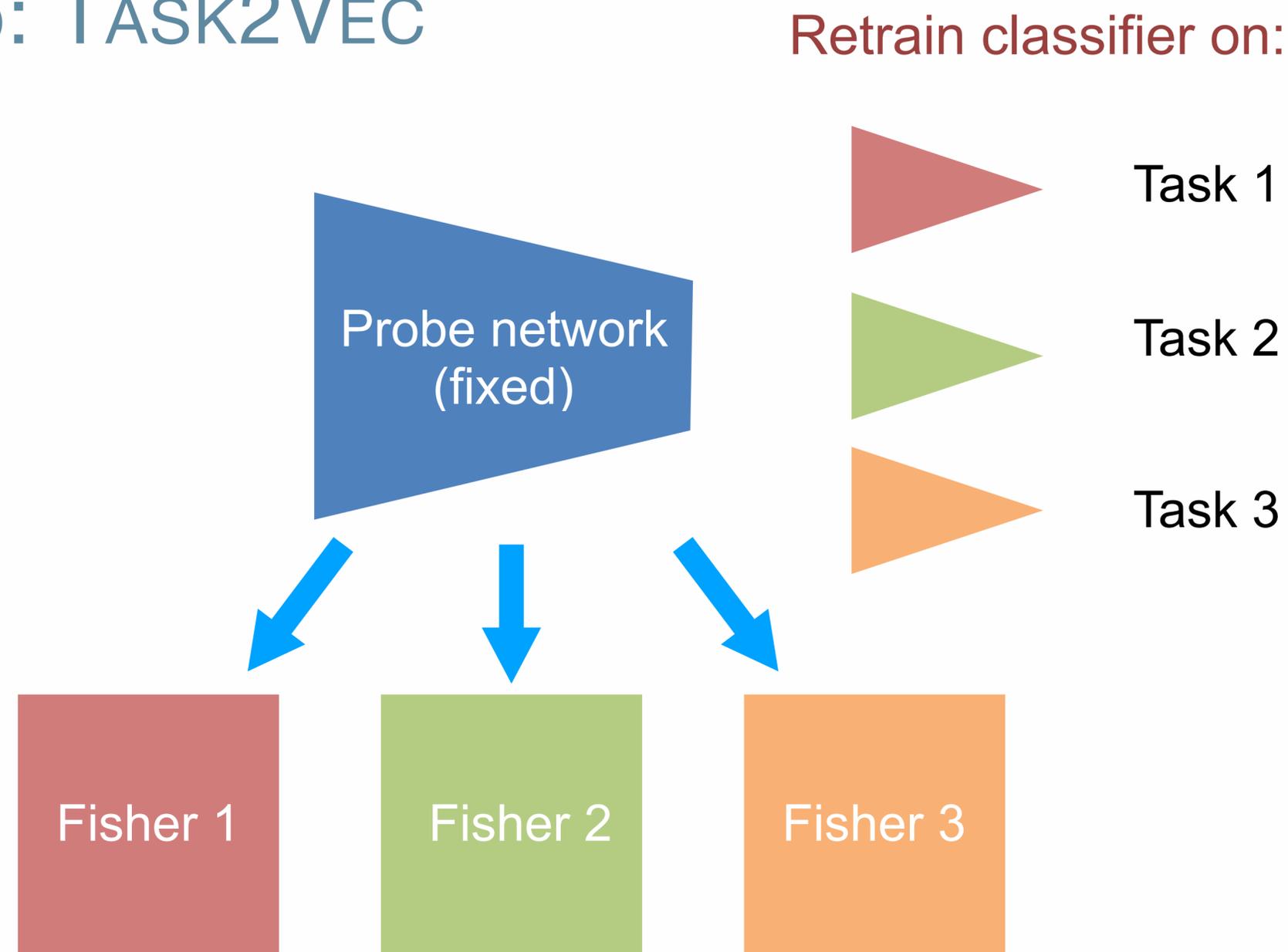
Complexity of
learning together

Complexity of
learning one

That is, how much more information do we need to learn?



Scaling this up: TASK2VEC



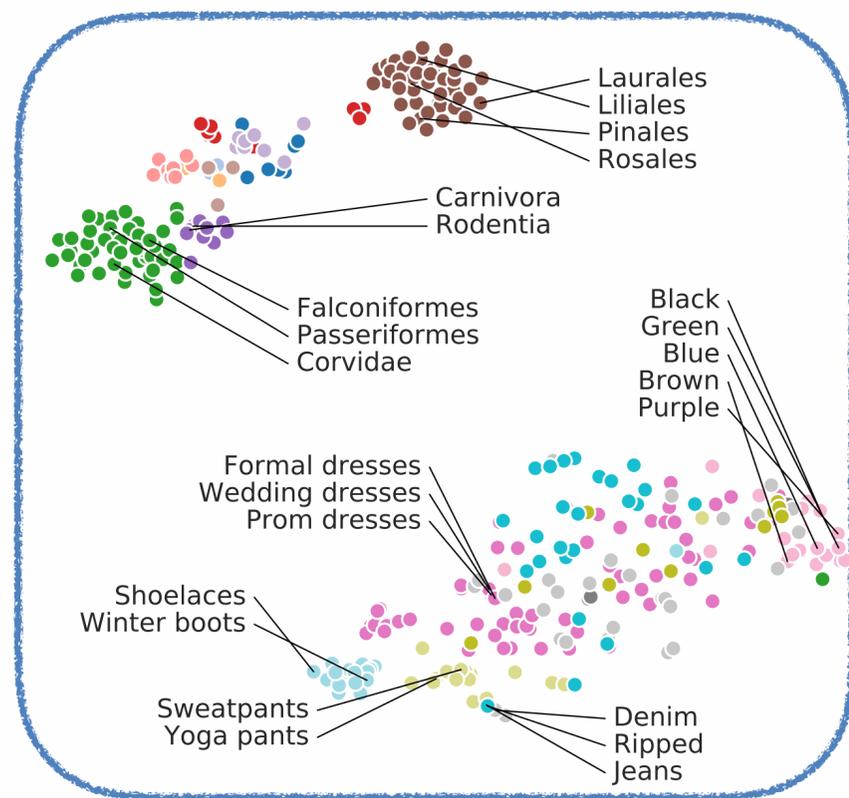
We compute the Fisher Information for the task using a pretrained set of weights and measure the distance between the compute the cosine distance between the Fisher:

$$d(\mathcal{D}_i, \mathcal{D}_j) := \cos(F_i, F_j)$$

TASK2VEC: Embedding tasks in a metric space

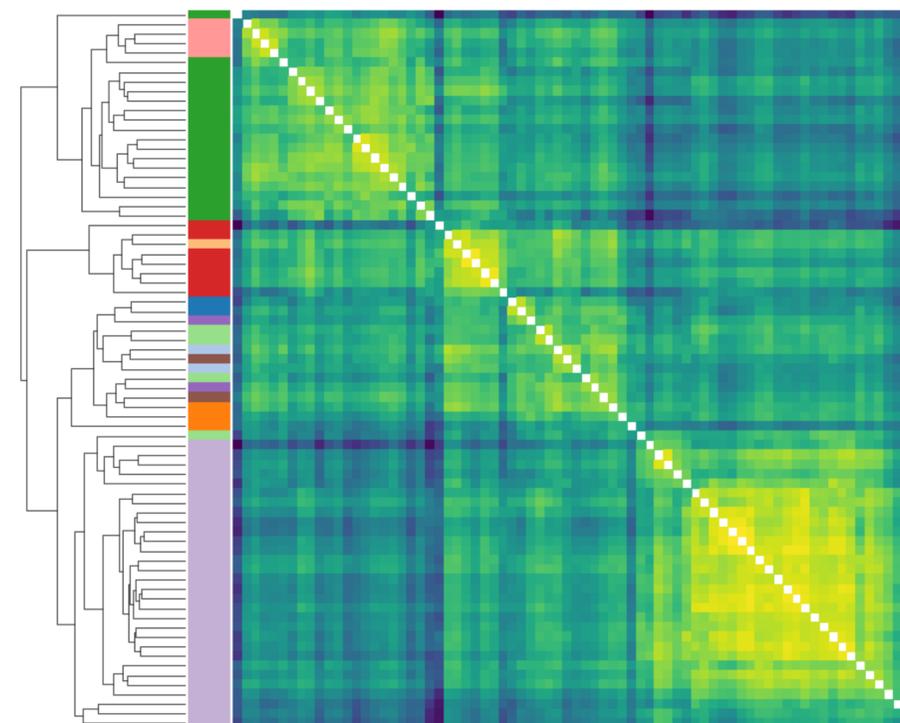
Idea: Represent each tasks in a metric space using its Fisher Information Matrix diagonal.

- Actinopterygii (n)
- Amphibia (n)
- Arachnida (n)
- Aves (n)
- Fungi (n)
- Insecta (n)
- Mammalia (n)
- Mollusca (n)
- Plantae (n)
- Protozoa (n)
- Reptilia (n)
- Category (m)
- Color (m)
- Gender (m)
- Material (m)
- Neckline (m)
- Pants (m)
- Pattern (m)
- Shoes (m)

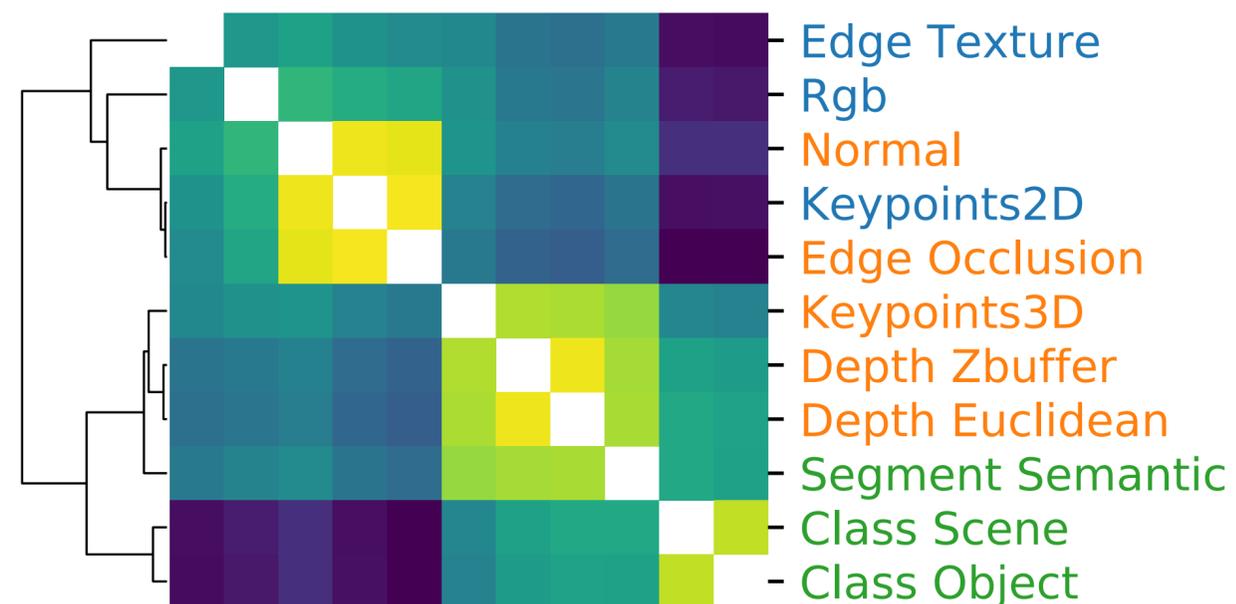
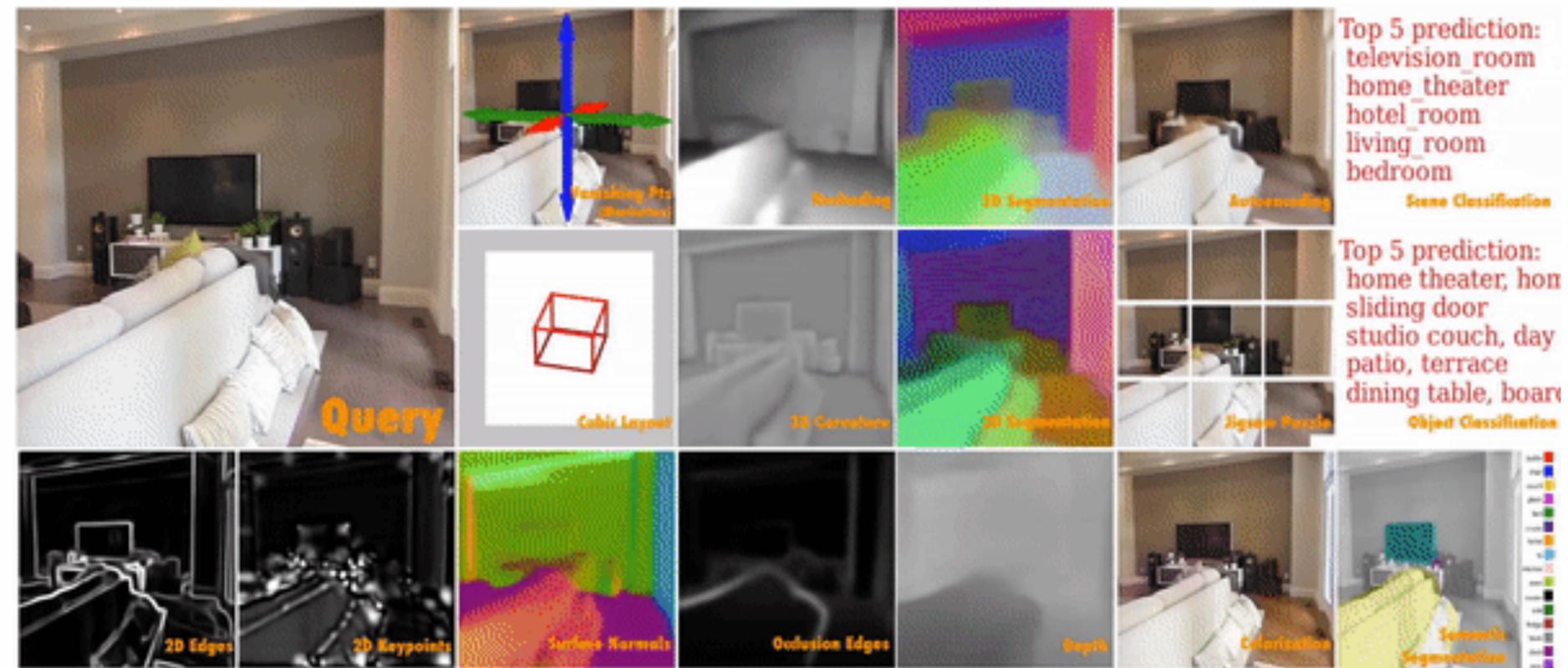


Recovers a meaningful topology
on hundred of tasks

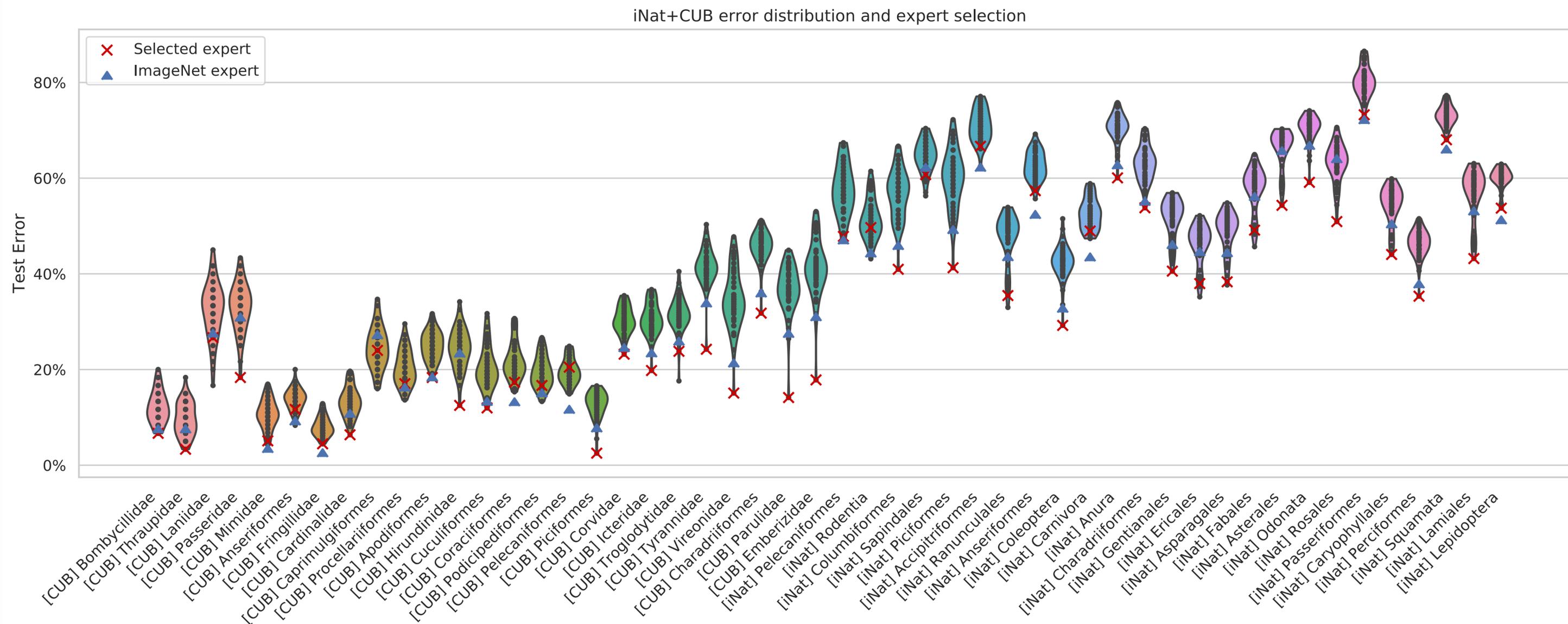
Recovers species taxonomy on
iNaturalist



Distance between different tasks on the same inputs

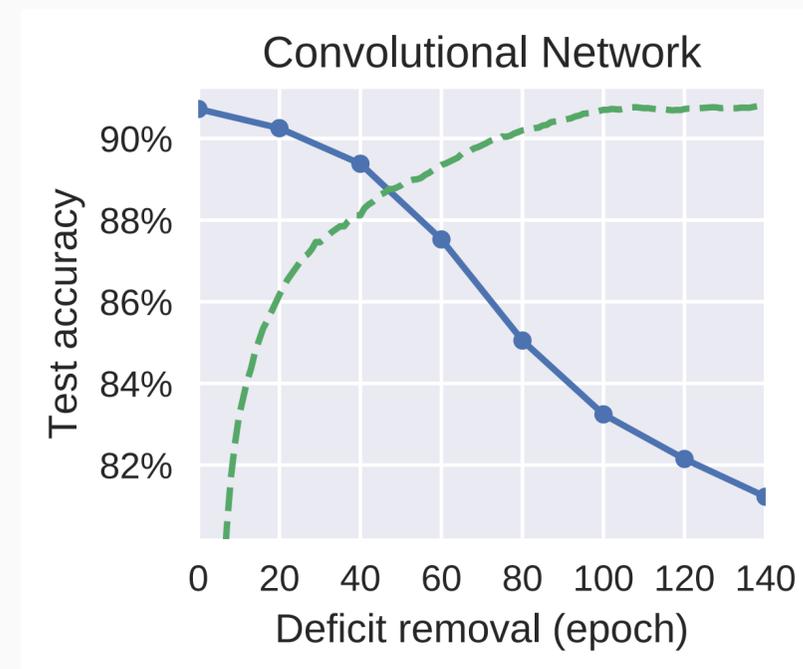
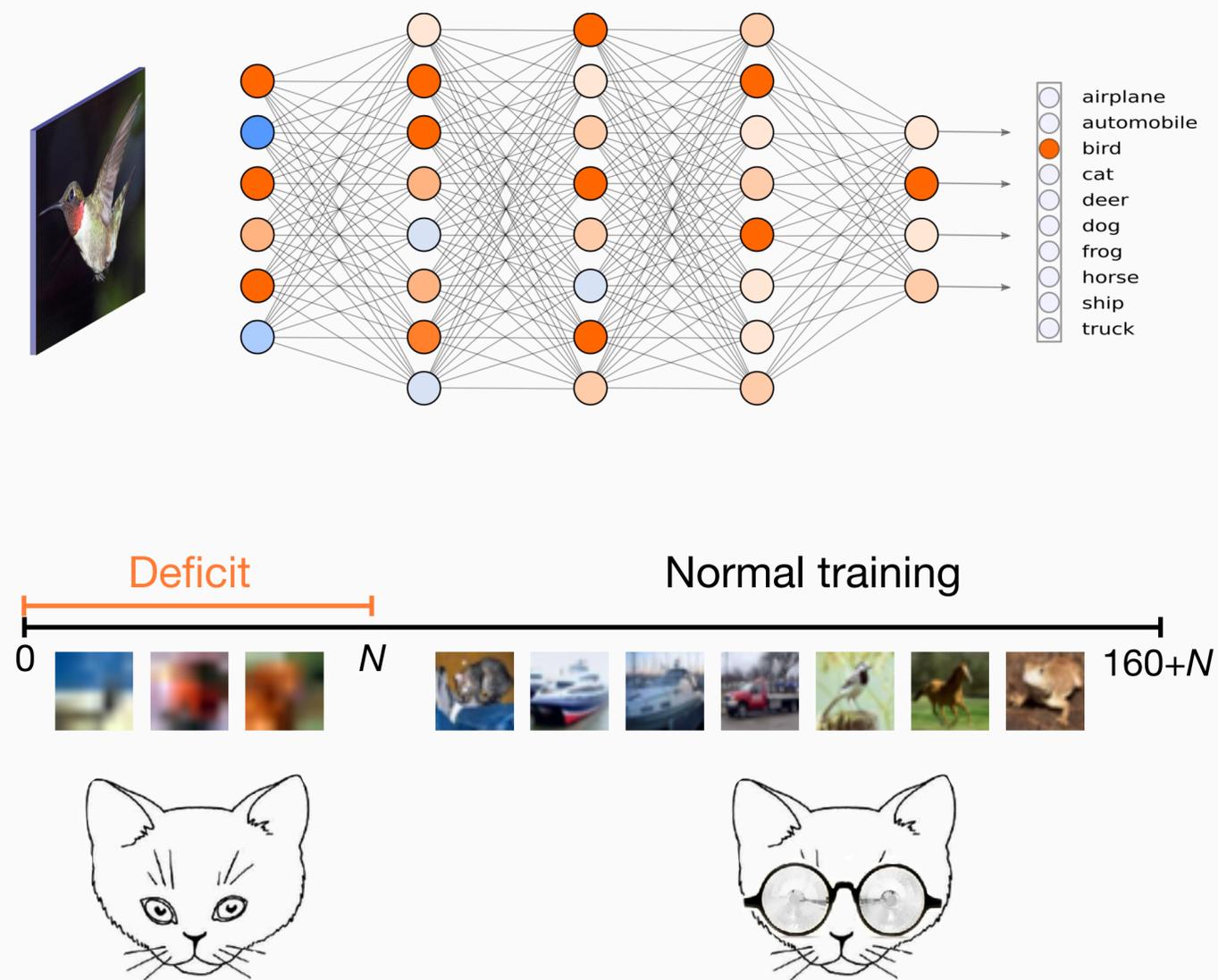


Proposing an optimal expert for the task



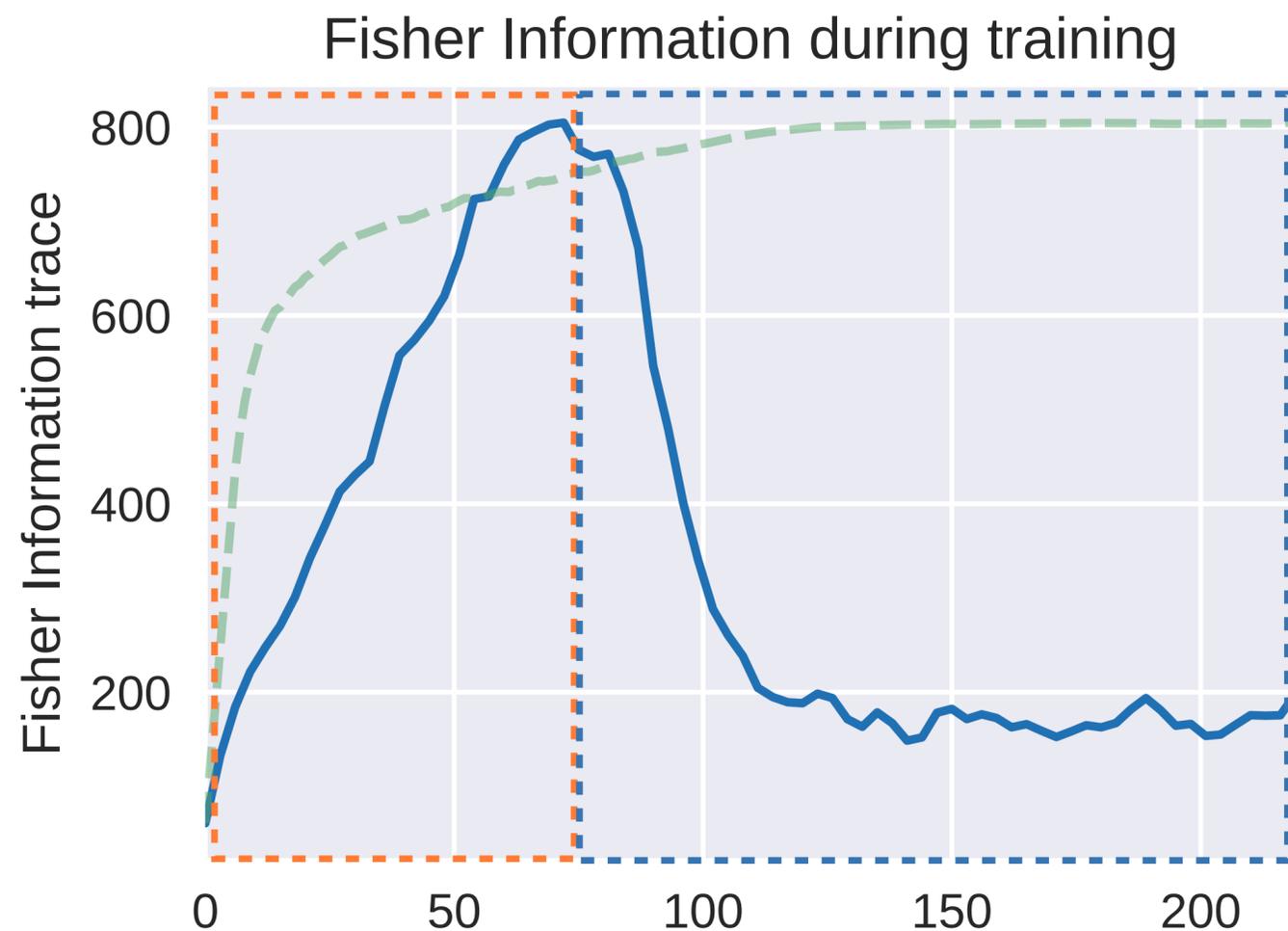
Allows to select the best expert to solve a task and substantially reduce error and training time.

A snag: Critical Learning Periods in Deep Networks



Compression of weights leads to compression of activations

Information acquisition behaves in a non-trivial way during training.



The (Fisher) information in the **weights** decreases later in training



The effective information in the **activations** decreases

THANKS!

- *Where is the Information in a Deep Neural Network?*, arXiv1:905.12213 (2019)
- *The Information Complexity of Learning Tasks, their Structure and their Distance*, arXiv:1904.03292 (2019)
- *TASK2VEC: Task embedding for meta-learning*, ICCV (2019)
- *Critical Learning Periods in Deep Neural Networks*, ICLR (2019)
- *Emergence of Invariance and Disentangling in Deep Representations*, JMLR (2018)



Matteo Rovere



Giovanni Paolini



Glen Mbeng



Stefano Soatto