



Topics in AI (CPSC 532S): Multimodal Learning with Vision, Language and Sound

Lecture 14: Unsupervised Learning, Autoencoders [Part 2]

Logistics

- **Assignment 1 & 2** grades are posted
- **Assignment 3** handed in, solutions are give
- **Assignment 4** is out (due last day before the break) — Do Part 1!
- **Project pitches next week**
9 groups per class (~8 minutes / group, 5-6 min presentation + questions)

Review — Autoencoders

Self (i.e. self-encoding)

— Feed forward network intended to reproduce the input

— Encoder/Decoder architecture

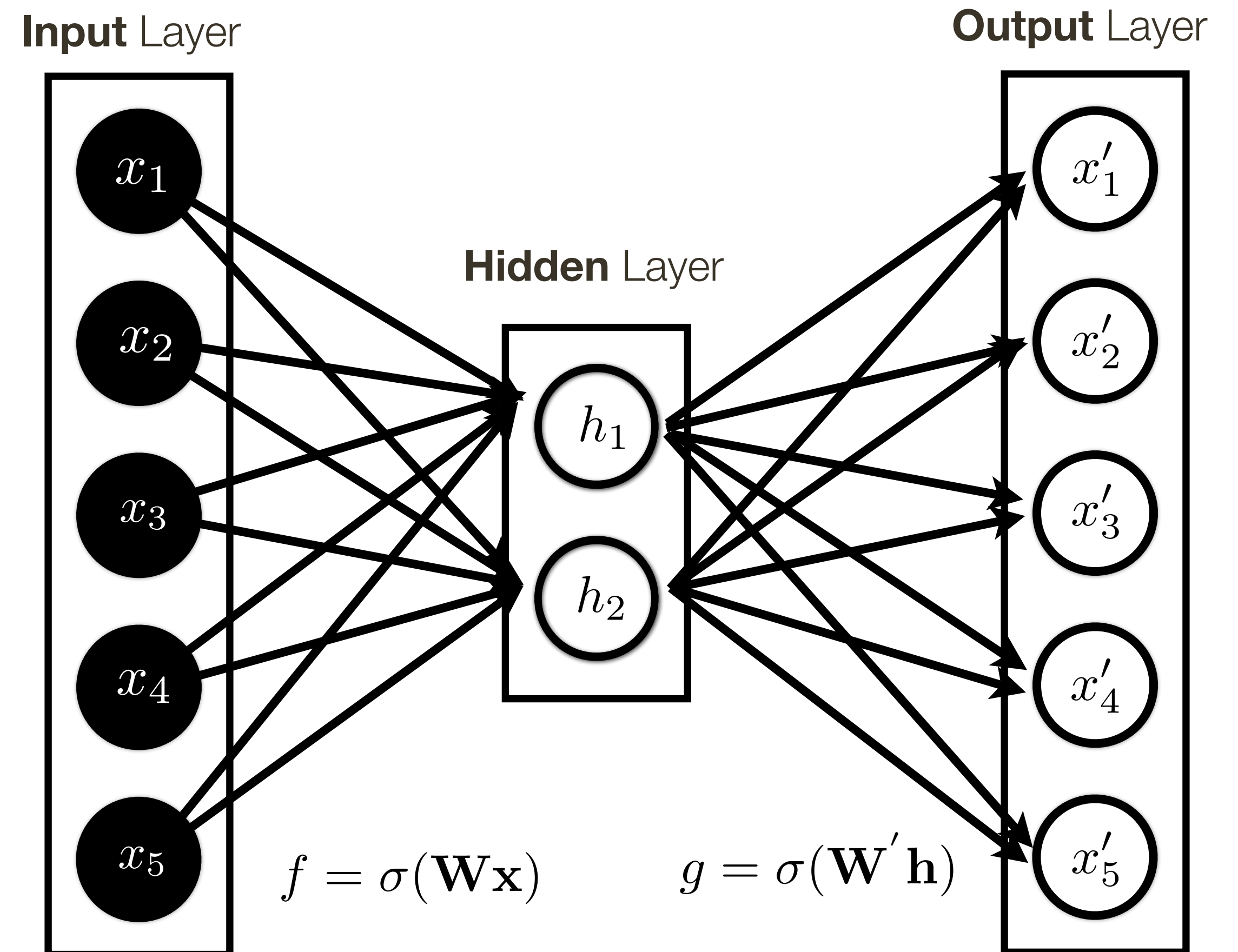
Encoder: $f = \sigma(\mathbf{W}\mathbf{x})$

Decoder: $g = \sigma(\mathbf{W}'\mathbf{h})$

— Score function

$$\mathbf{x}' = f(g(\mathbf{x}))$$

$$\mathcal{L}(\mathbf{x}', \mathbf{x})$$

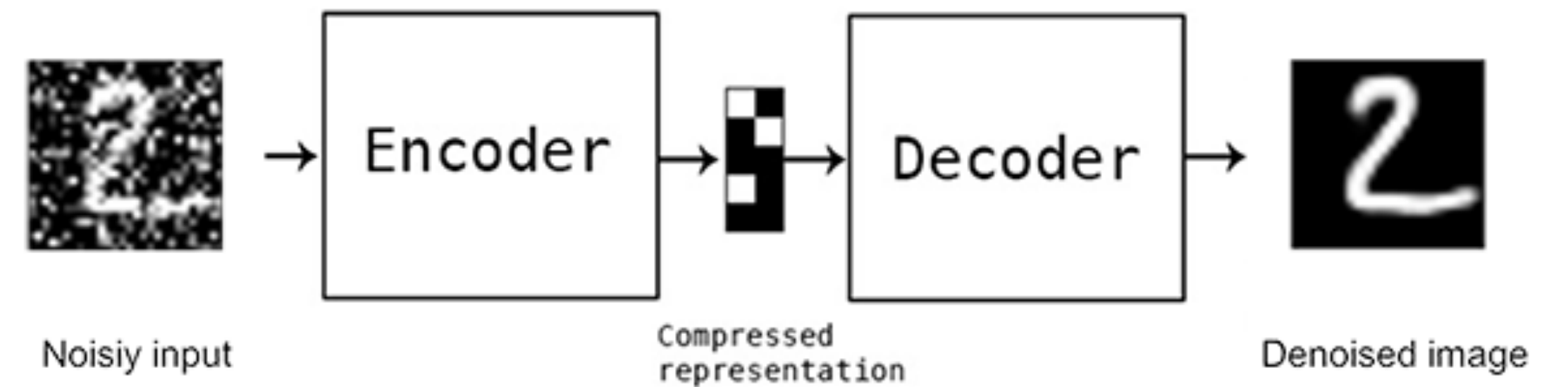


Review — De-noising Autoencoder

Idea: add noise to input but learn to reconstruct the original

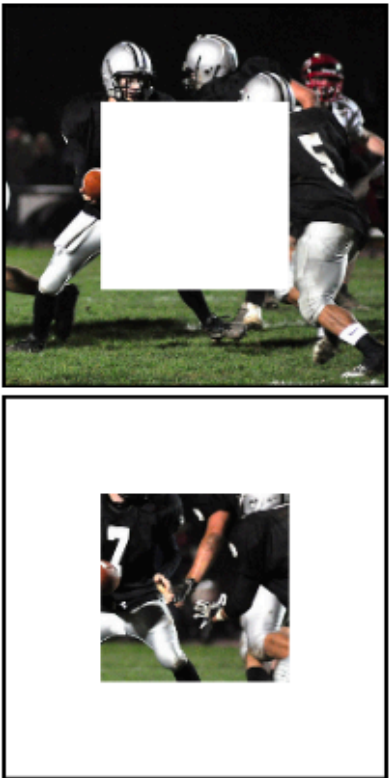
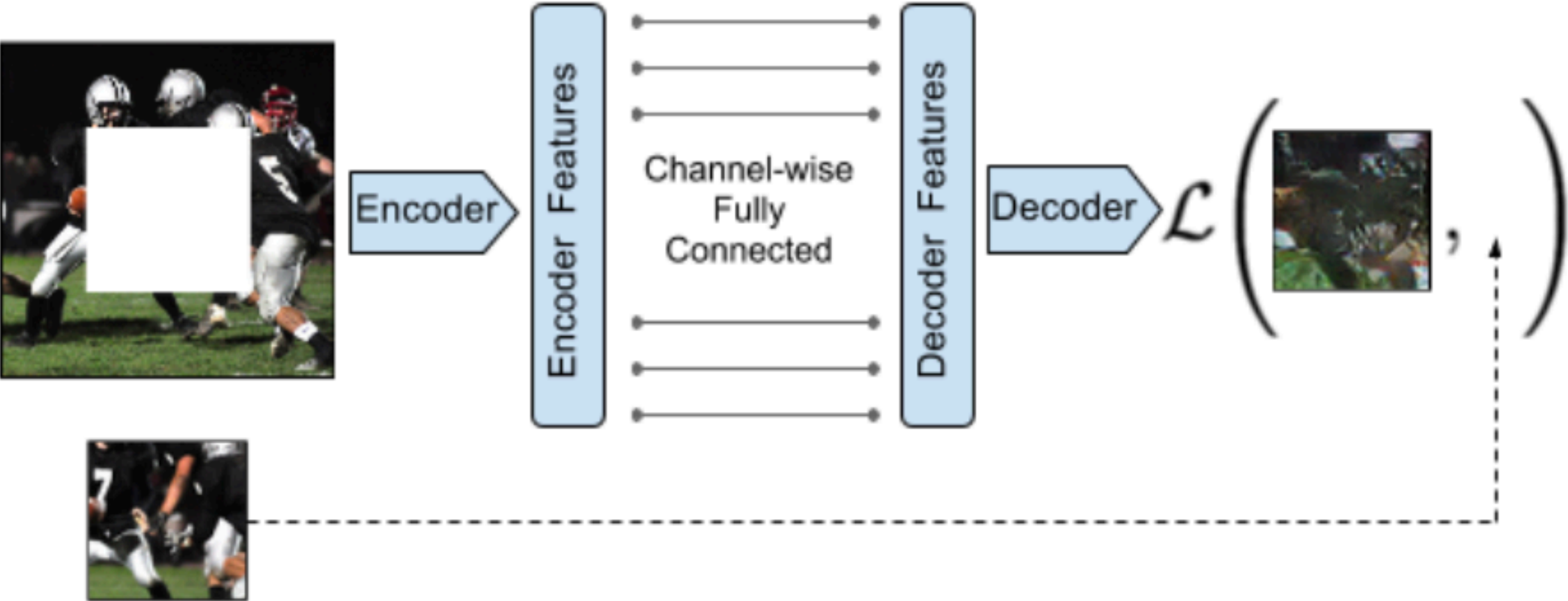
- Leads to better representations
- Prevents copying

Note: different noise is added during each epoch

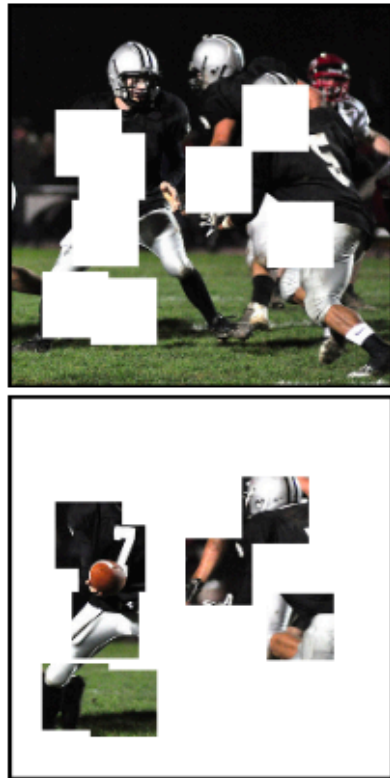


Review — Context Encoders

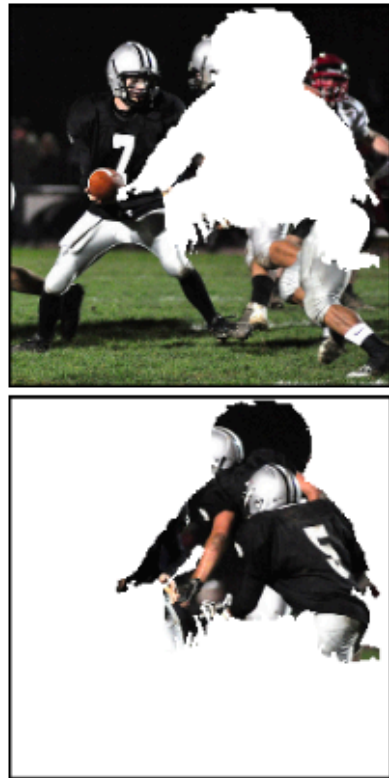
[Pathak et al., 2016]



(a) Central region



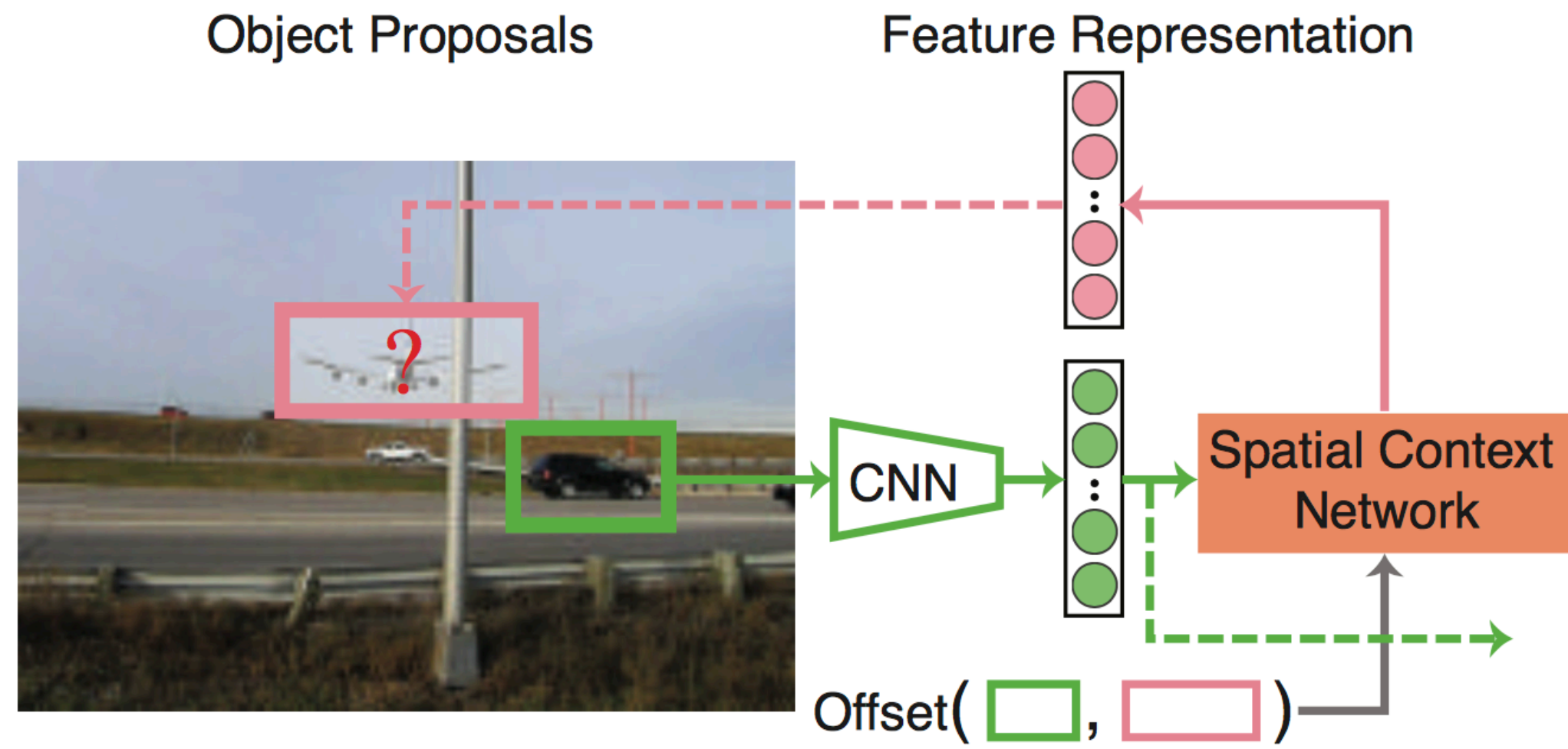
(b) Random block



(c) Random region

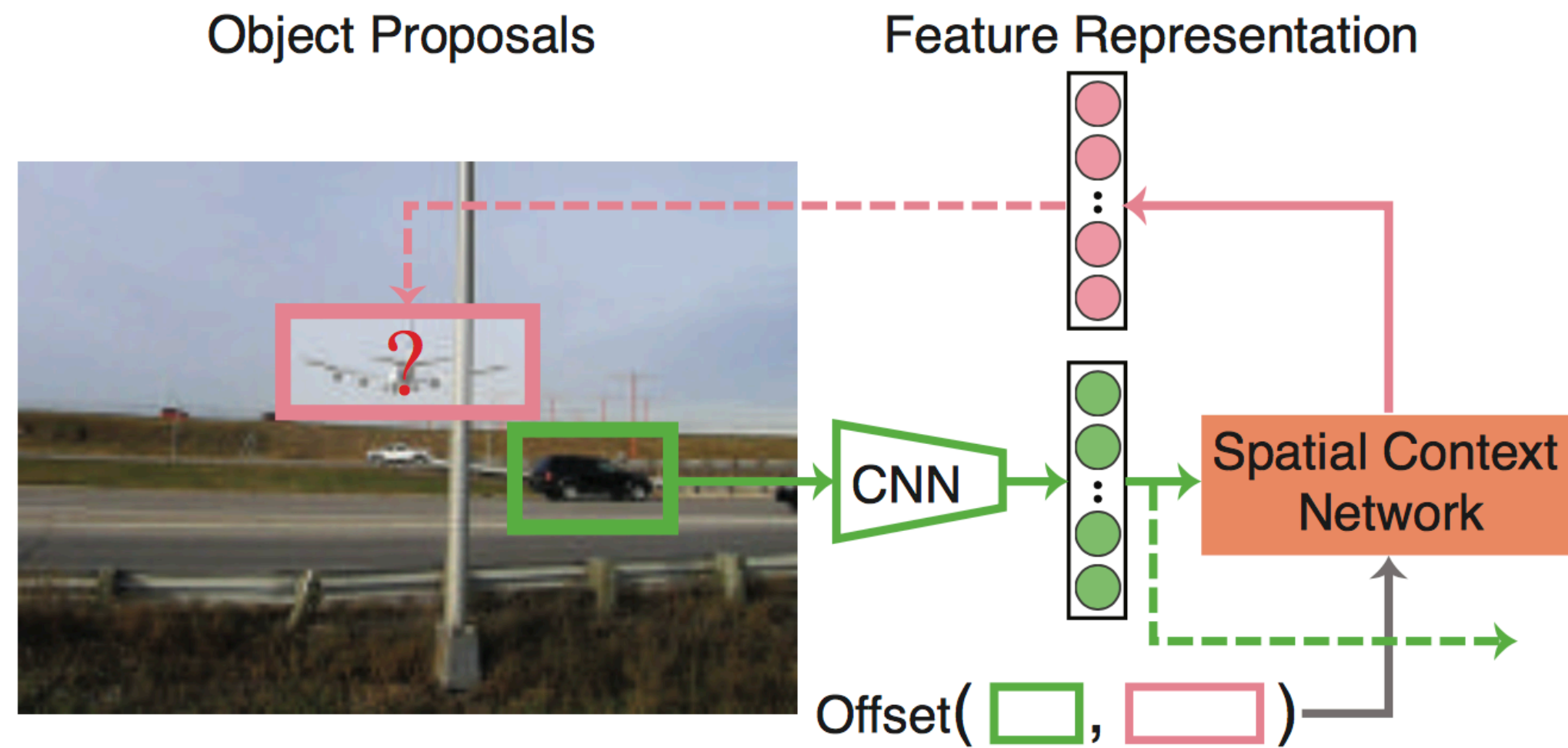
Spatial Context Networks

[Wu, Sigal, Davis, 2017]



Spatial Context Networks

[Wu, Sigal, Davis, 2017]

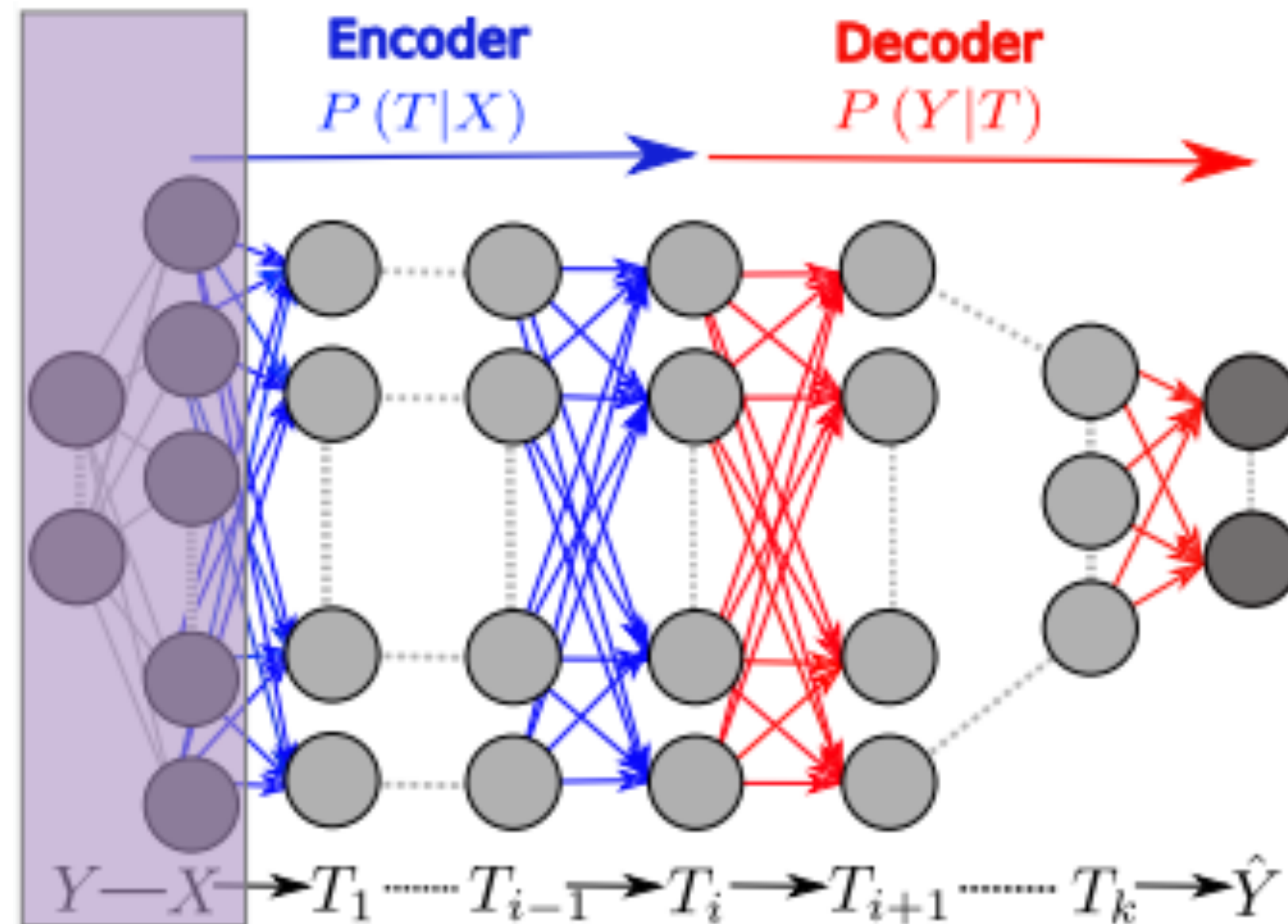


	Initialization	Supervision	Pretraining time	Classification	Detection
Random Gaussian	random	N/A	< 1 minute	53.3	43.4
Wang <i>et al.</i> [32]	random	motion	1 week	58.4	44.0
Doersch <i>et al.</i> [3]	random	context	4 weeks	55.3	46.6
*Doersch <i>et al.</i> [3]	1000 class labels	context	–	65.4	50.4
Pathak <i>et al.</i> [21]	random	context inpainting	14 hours	56.5	44.5
Zhang <i>et al.</i> [36]	random	color	–	65.6	46.9
ImageNet [21]	random	1000 class labels	3 days	78.2	56.8
*ImageNet	random	1000 class labels	3 days	76.9	58.7
SCN-EdgeBox	1000 class labels	context	10 hours	79.0	59.4

A Little Theory: Information Bottleneck [Tishbi et al., 1999]

Every layer could be treated as a random variable, then entire network is a Markov Chain

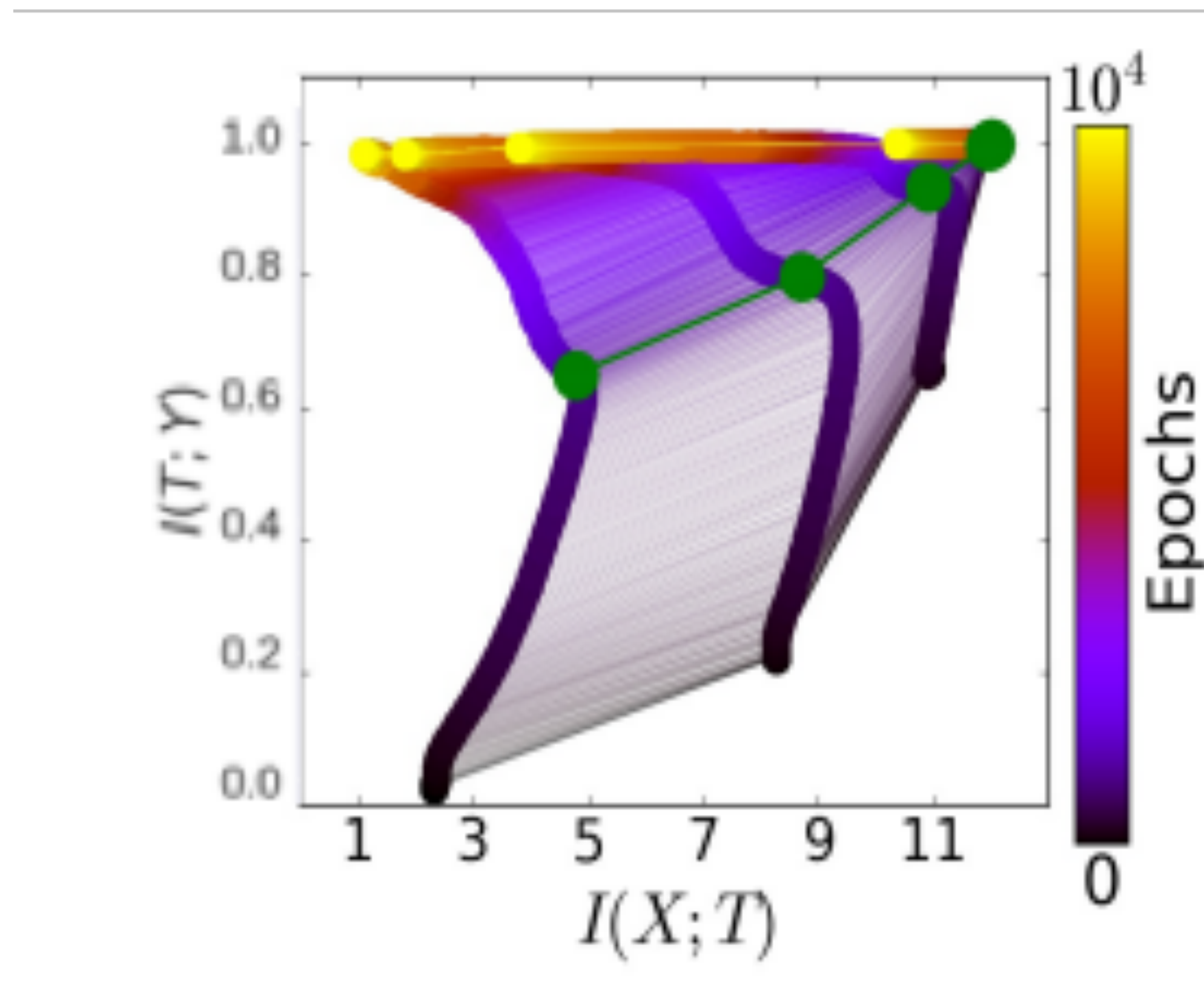
Data processing theorem: if the only connection between X and Y is through T , the information that Y gives about X cannot be bigger than the information that T gives about X .



$$I(X; Y) \leq I(T_1; Y) \leq I(T_2; Y) \leq \dots \leq I(\hat{Y}; Y)$$

A Little Theory: Information Bottleneck [Tishbi et al., 1999]

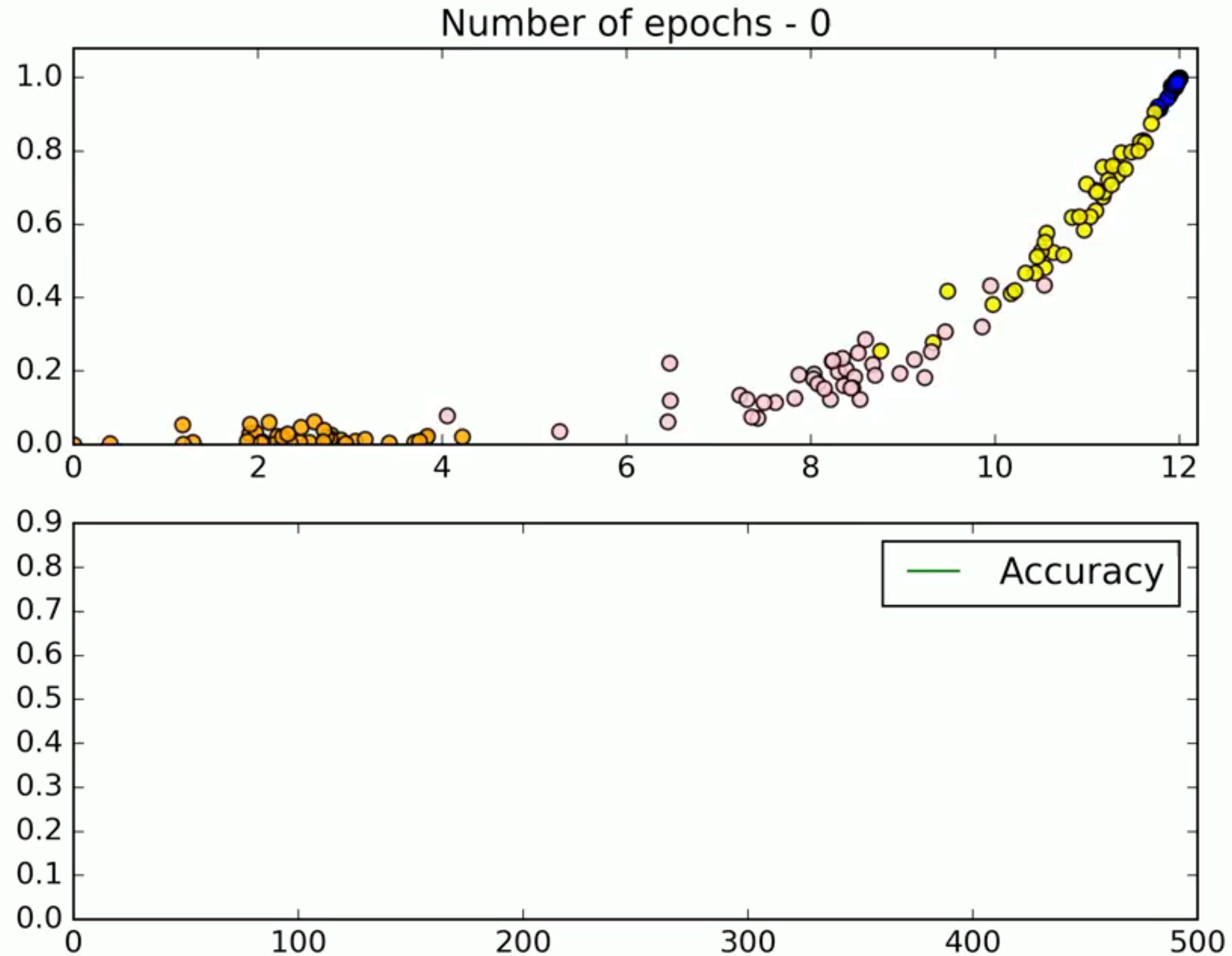
Observation: In the information plane layers first increase the mutual information between themselves and the output and then reduce information between themselves and the input (which leads to “forgetting” of irrelevant inputs and ultimately generalization)



A Little Theory: Information Bottleneck

[Tishbi et al., 1999]

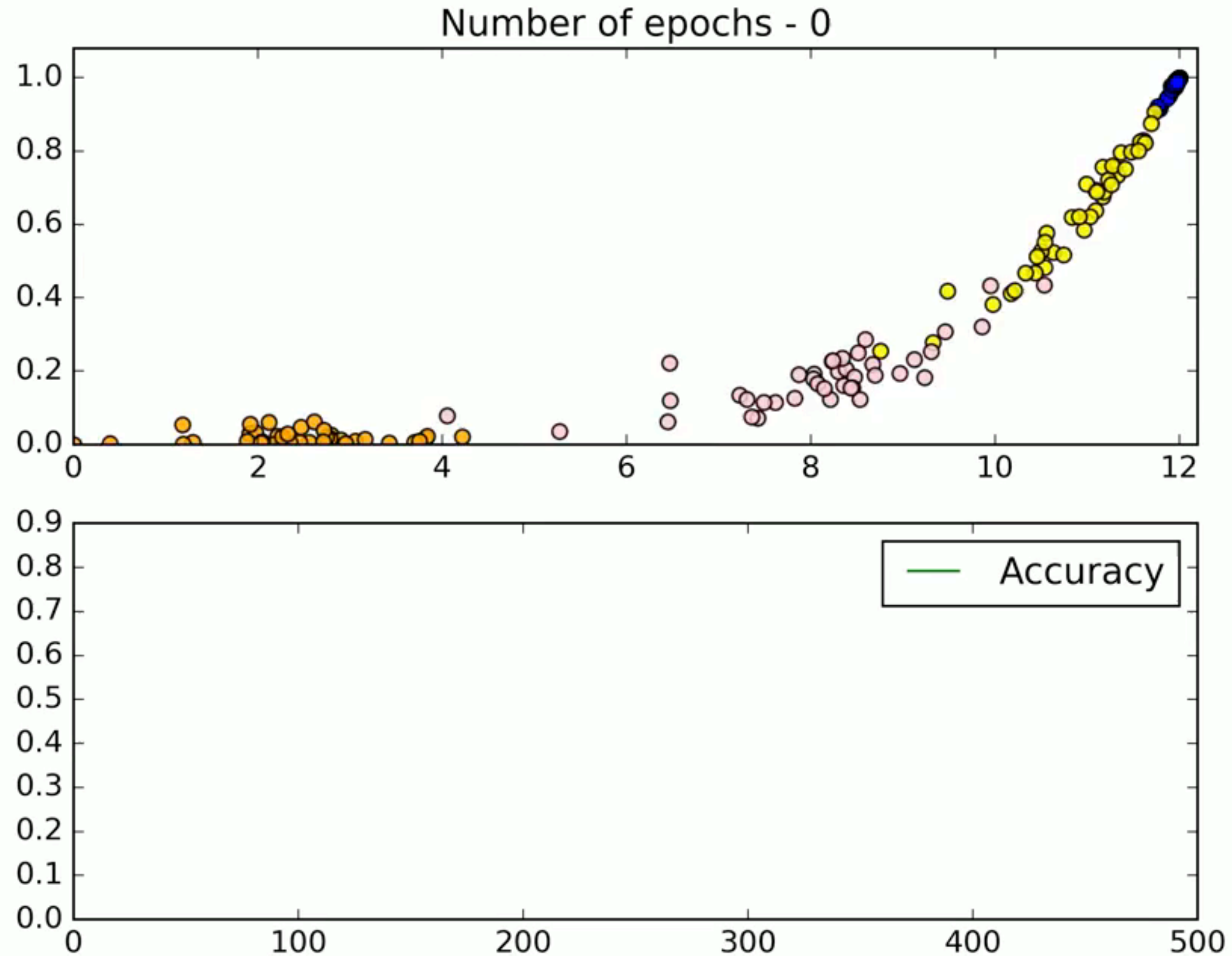
50 networks of same topology being optimized



A Little Theory: Information Bottleneck

[Tishbi et al., 1999]

50 networks of same topology being optimized



A Little Theory: Information Bottleneck [Tishbi et al., 1999]

Limitation: Does not seem to work for non-Tanh activations (e.g., ReLU)

