

# Beyond Shannon Applications in Machine Learning

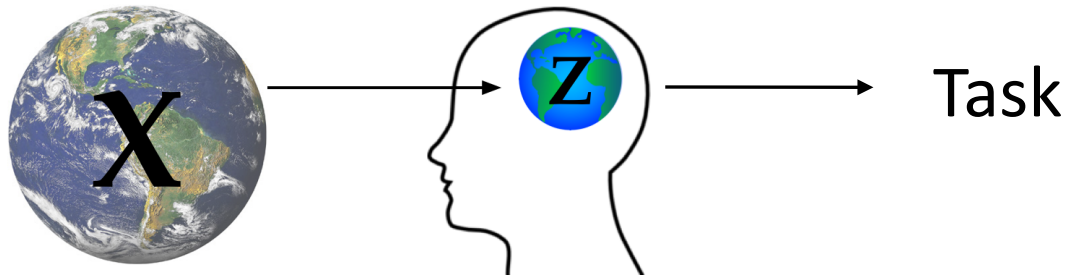
Greg Ver Steeg

Beyond Shannon Workshop, 2019



*Information Sciences I*





Z is a learned  
*representation\** of X

$$z = J\theta$$

E.g., parametrized by a  
neural network

Outline:

“Disentangling”

“Fair” representation learning

Lossy Compression in representation learning

- Gaussian vs Echo compression
- Higher order interactions

Slides with all paper links: <http://bit.ly/ST>

\* Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. *Transactions on pattern analysis and machine intelligence*, 35(8): 1798–1828, 2013.

# KL is King

In machine learning, Kullback-Leibler divergence and related quantities are dominant: *Maximum likelihood, cross entropy loss, variational inference, Mutual information*



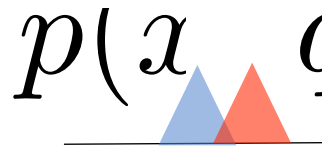
But there are practical issues for machine learning:

- Estimation is hard
- Often rely on loose bounds
- Usage is incorrect / inappropriate / poorly motivated
- Highly non-convex optimization
- Underflows lead to numerical errors with rare events,  $\log 0 = -\infty$

# Other information measures in ML?

Smoother optimization (than KL divergence) with Wasserstein distance (earth mover distance for distributions)

- Martin Arjovsky, Soumith Chintala, Léon Bottou, Wasserstein GAN, <https://arxiv.org/abs/1701.07875>



Other (in)dependence measures: Hilbert-Schmidt Independence Criterion (HSIC) / Maximum-Mean Discrepancy (MMD) / kernel trick

- Introduced: Gretton et al 2007
- Example usage: InfoVAE, [arXiv:1706.02262](https://arxiv.org/abs/1706.02262)

Almost never in ML: Renyi, info decomposition, synergy, intersection, common info

- Wyner common information: [arXiv:1606.02307](https://arxiv.org/abs/1606.02307)
- Synergy: [arXiv:1710.03839](https://arxiv.org/abs/1710.03839)
- Hierarchical decomposition of total correlation: [arXiv:1410.7404](https://arxiv.org/abs/1410.7404)

“Disentangling” – encourage  $Z_j$  to match  $x$  (without

labels)

, Tian Qi, et al.

Identifying sources of

entanglement in

neural

encoders." *NIPS* 2018

? Good old statistical

independence



$$D(Z) \equiv \mathbb{E}_p \left[ \log \frac{p(\approx 1)}{\prod_j} \right]$$

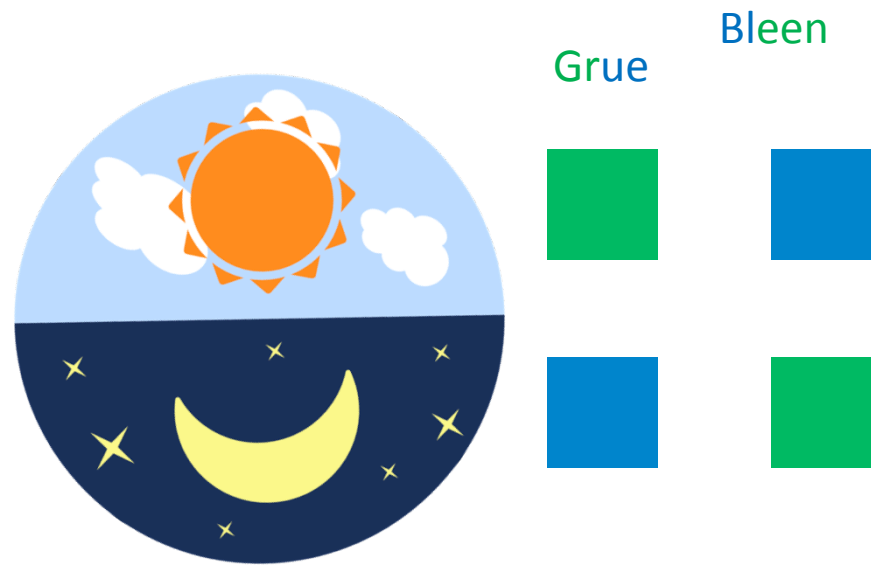
papers working on this angle... with dubious results, see review:

Luca Locatello, Stefan Bauer, Mario Lucic, Sylvain Gelly, Bernhard Scholkopf, and Olivier Bachem. Challenging conditions in the unsupervised learning of disentangled representations. arXiv preprint arXiv:1811.12359, 2018.

The Grue Language doesn't have words for "blue" or "green"

**Grue** = **G**reen during the **day**, **bl**ue at *night*

**Bleen** = **B**lue during the **day**, **gr**een at *night*



# Needlessly complicated?

English is needlessly complicated because:

bleen = Grue during day and bleen at night

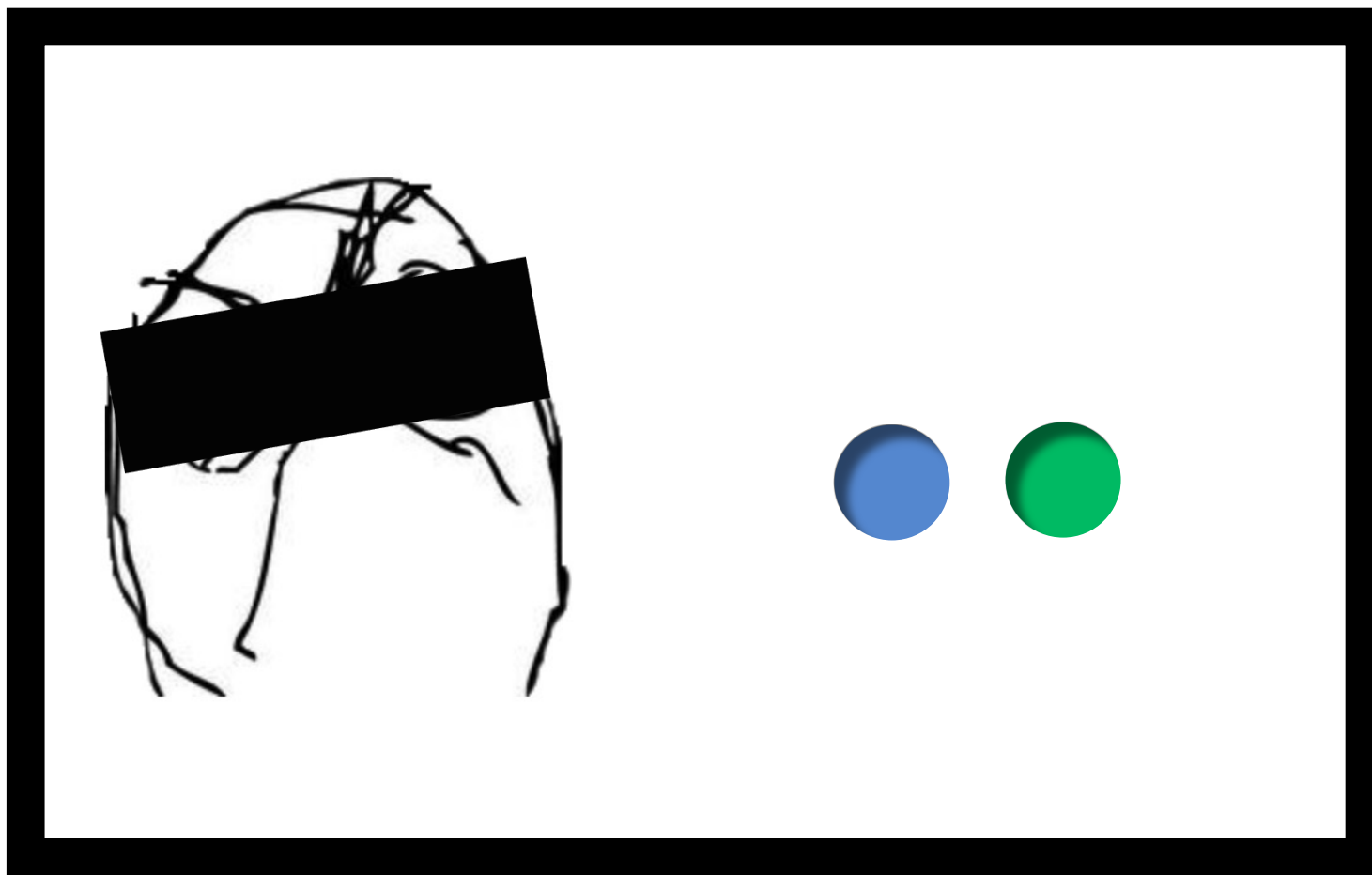
grue = Bleen during day and grue at night

Grue language  
speakers



(Also, what is deal with “the”  
and “a”? Why you have these  
useless words?)

Push the bleen button and we let you go...



Elaborate windowless ca  
←

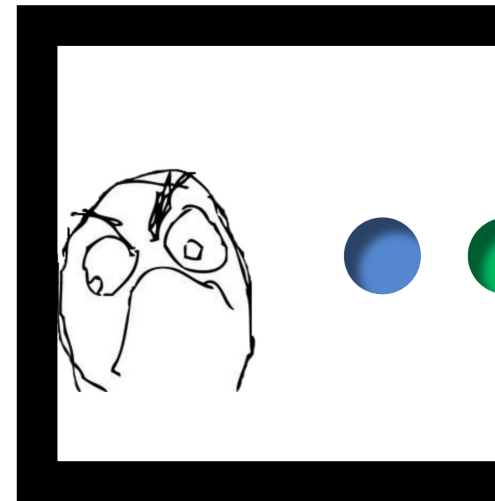


# Grue language is synergistic



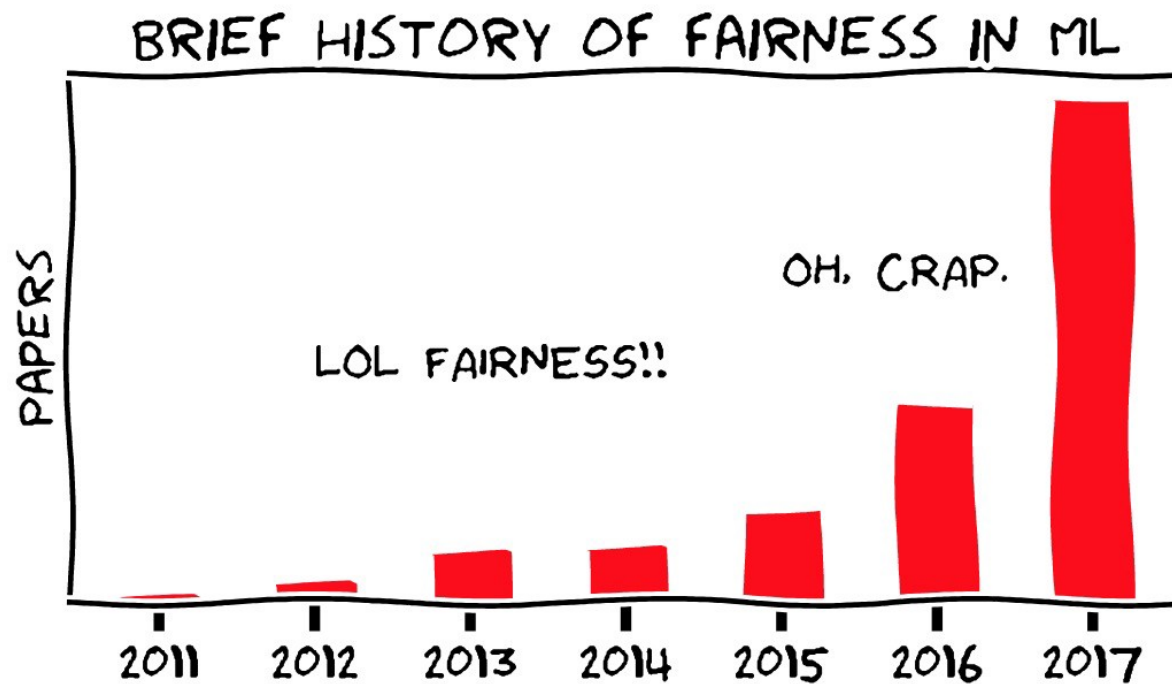
time and grue > time only *or* grue only  
(complete info) (no info)

*predicting what you will see*



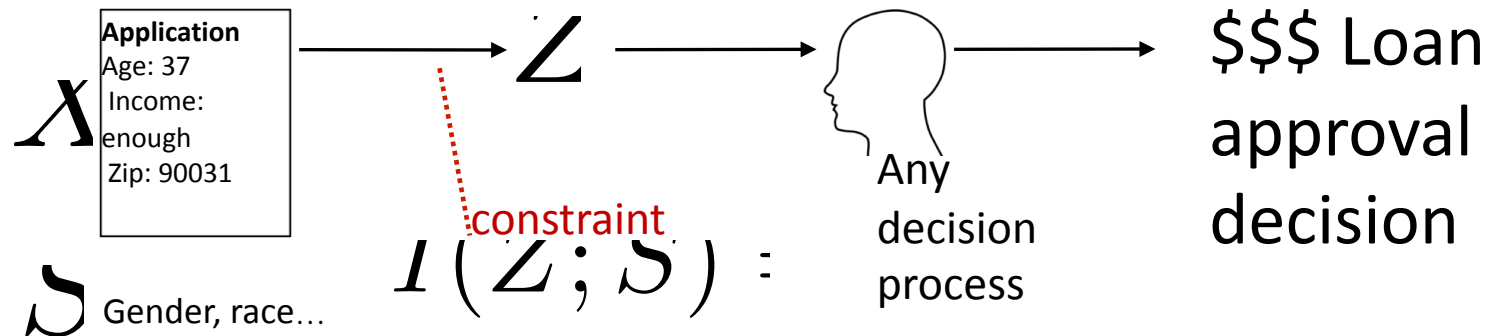
**Is there hope for a plausible information-theoretic principle of “disentangling”? And one that can be optimized for learning?**

# Fair and Invariant Representations



<https://towardsdatascience.com/a-tutorial-on-fairness-in-machine-learning-3ff8ba1040cb>

# An information-theoretic notion of fairness



Example: your task is to approve loans

To be fair, you don't want to discriminate based on protected variables

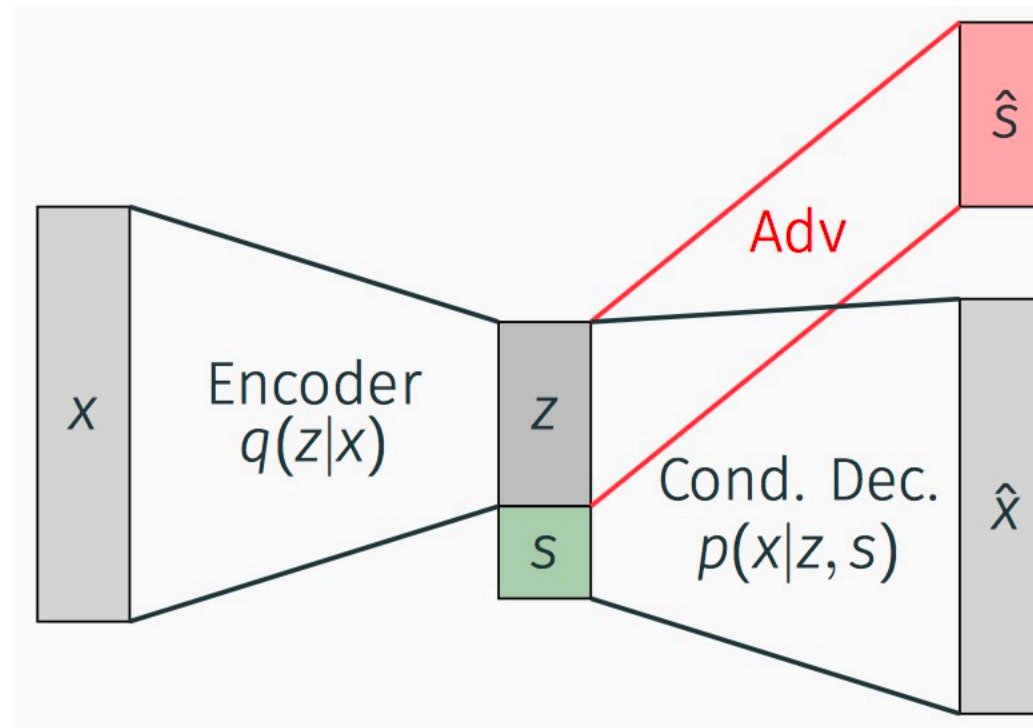
For historical reasons,  $S$  may be correlated with  $X$

By removing information about  $S$  from the representation  $Z$  that we use for decision making, we make it impossible to be biased against  $S$

# Adversarial approach

Make sure an adversary cannot reconstruct  $S$  from  $Z$   
This doesn't guarantee that  $I(Z;S)$  is small! Adversaries only provide *lower* bounds on mutual information.

(Nevertheless adversarial learning revolutionized some problems)



C. Louizos, K. Swersky, Y. Li, M. Welling, and R. Zemel. variational fair autoencoder. *arXiv preprint arXiv:1511.00830*, 2015.

Q. Xie, Z. Dai, Y. Du, E. Hovy, and G. Neubig. Controllable invariance through adversarial feature learning. *NeurIPS*, 2017.

# A direct, information-theoretic approach

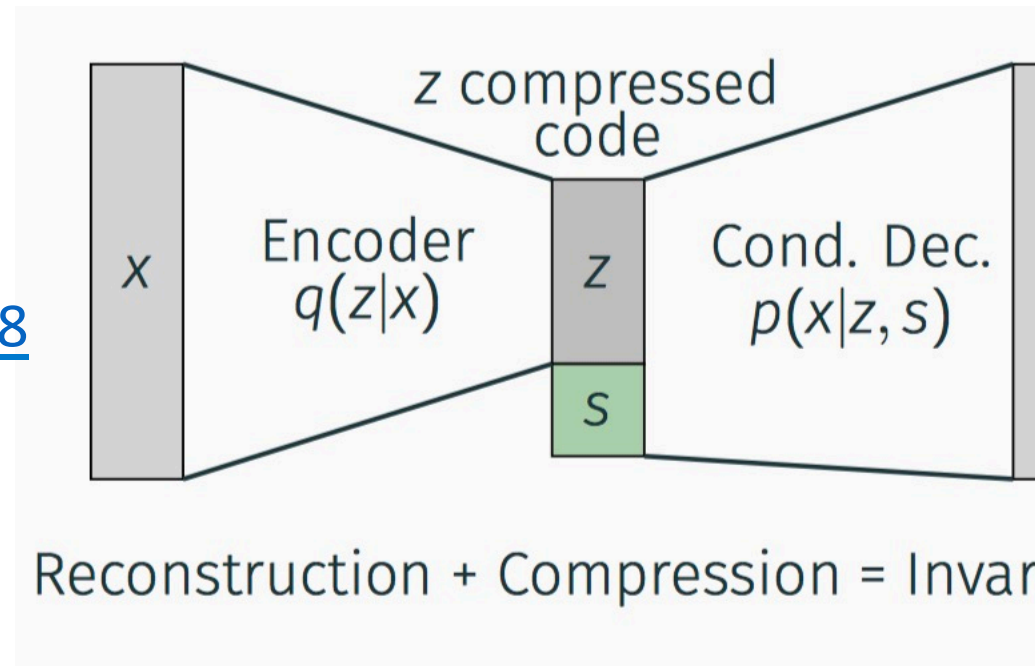
We derive an *upper* bound on  $I(Z;S)$   
that can be tractably optimized

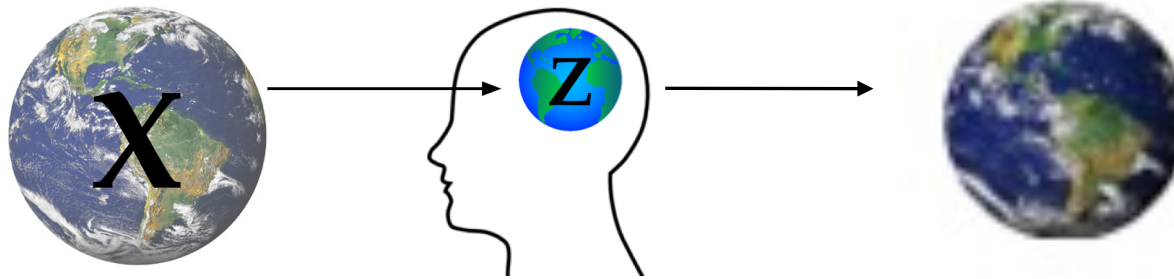
$$I(Z;S) \leq \underbrace{-\mathbb{E}_q [\log p(x|z,s)]}_{\text{Conditional reconstruction}} + \underbrace{I_q(Z;S)}_{\text{Constant}}$$

Main idea in NIPS 2018, [arXiv:1805.09458](https://arxiv.org/abs/1805.09458)

Application to fMRI: [arXiv:1904.05375](https://arxiv.org/abs/1904.05375)

**When and how to remove information?**





Z is a compressed *representation* of X,  
that should be useful for reconstruction

# Lossy compression (and VAE's)

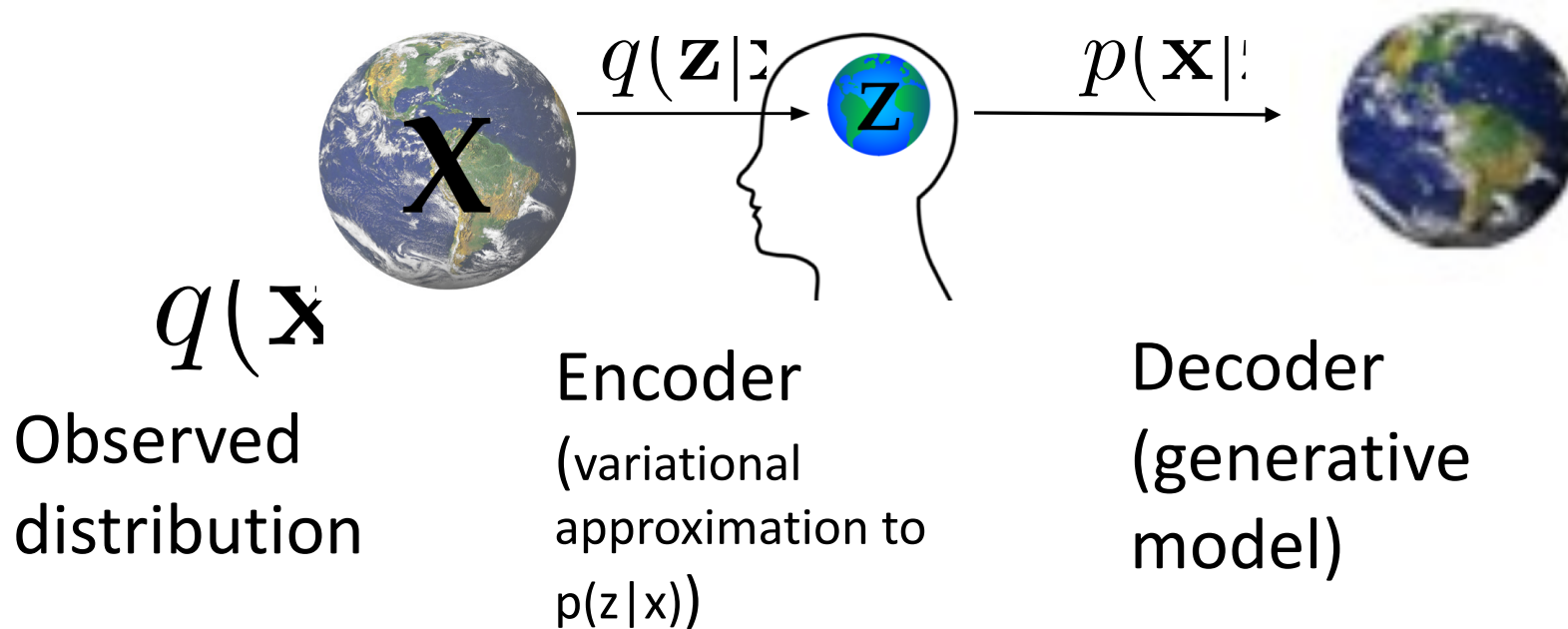
Variational auto-encoders (VAEs) are our motivating example. They have an interpretation in terms of generative models as well as with rate-distortion, as discussed in several papers:

Delellis, A., Poole, B., Fischer, I., Dillon, J., Saurous, R. A., and Murphy, K. Fixing a broken ELBO. In *International Conference on Machine Learning*, pp. 159–168, 2018.

Rezende, D. J. and Viola, F. Taming VAEs. *arXiv preprint arXiv:1810.00597*, 2018.

Gregor, K., van de Walle, M., Moyer, Galstyan, and Ver Steeg. "Exact Rate-Distortion in Autoencoders via Echo Noise." *arXiv preprint arXiv:1904.07199* (2019).

# Variational Auto-Encoder (VAE)



Encoder and decoder parametrized by neural nets

Goal is to maximize likelihood of observations under generative model

# Variational Evidence Lower Bound (ELBO)

$$\mathbb{E}_q \log p_\theta(\mathbf{x}) - \mathbb{E}_q \log p_\theta(\mathbf{x}|\mathbf{z}) - D_{KL}[q(\mathbf{z}) \| p(\mathbf{z})]$$

$$= \mathbb{E}_q \log p_\theta(\mathbf{x}|\mathbf{z}) - D_{KL}[q(\mathbf{z}) \| p(\mathbf{z})]$$

True for any  $p(\mathbf{z})$ . The choice with the tightest bound is  $p(\mathbf{z})=q(\mathbf{z})$

$$\mathbb{E}_{q(\mathbf{x})} \log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{x}|\mathbf{z})].$$

$$\max_{\theta, \phi} \mathbb{E}_{q_\phi} [\log p_\theta(\mathbf{x}|\mathbf{z})] -$$

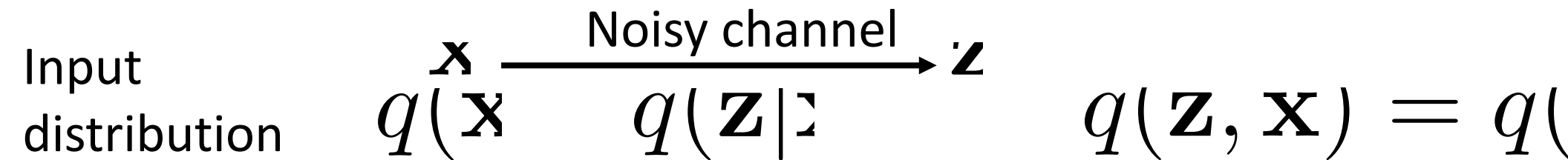
Distortion /  
reconstruction loss

Compression



Compression: bounds and a  
new exact approach with echo  
noise

# Information in a noisy channel



Mutual information

$$= D_{KL} [q(\mathbf{z}) | q(\mathbf{z})]$$

# Problem

$$q(\mathbf{z}) = \int d\mathbf{x} q(\mathbf{z} |$$

High-d  
integral

Could be complex  
(images, audio,  
gene expression...)

Approximate as Gaussian for (MaxEnt)  
Upper bound

$$Z_i = X_i + \epsilon_i \quad \epsilon_i \sim \mathcal{N}(0, \sigma_i^2)$$

Approximate  $q(z)$  as Normal:  $p(z_i) = \mathcal{N}(\mathbb{E}(x_i), \text{Var}(x_i))$

$$I(Z; X) \leq \sum_i \frac{1}{2} \log \left( 1 + \frac{\text{Var}(x_i)}{\sigma_i^2} \right)$$

Only tight if the *input is Gaussian*, and *each channel is independent*

# Echo noise

By choosing a more flexible noise channel, we can exactly specify information rates for arbitrary inputs



# Echo noise: make the noise look like the signal

Why?

- For correlated Gaussian noise, optimal signal is correlated in same basis
- **Key property** for analytic mutual information under arbitrary input

# Echo noise: make the noise look like the signal

How do we make the noise look like the signal?

$$= f(\mathbf{x}) -$$

$$= f(\mathbf{x}^{(0)}) + s f(\mathbf{x}^{(1)}) + s^2 f(\mathbf{x}^{(2)}) \dots$$

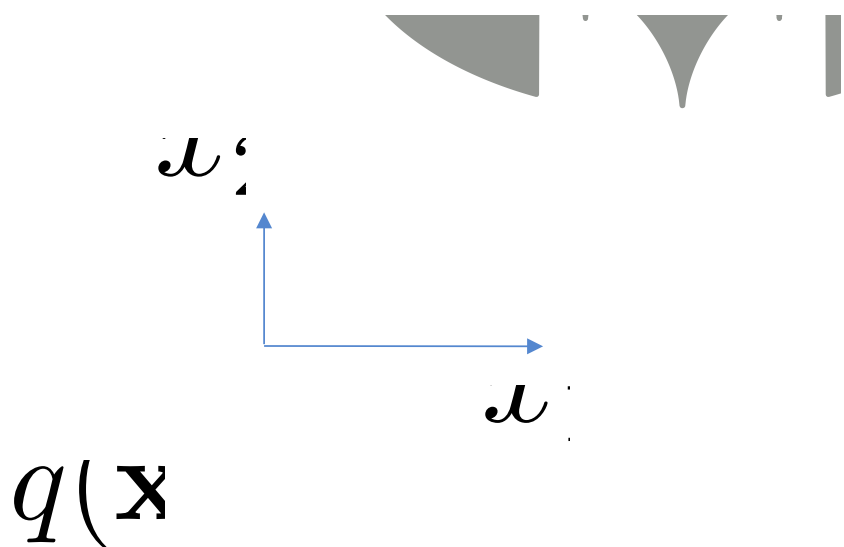
$$= f(\mathbf{x}) + s \left( f(\mathbf{x}^{(0)}) + s f(\mathbf{x}^{(1)}) + \dots \right)$$

Multiply out and re-label iid samples....  
these are the same (in distribution)!



# Example: a non-Gaussian input distribution

A uniform distribution in  $\mathbb{R}^2$  with a shape that strikes fear into the hearts of villains and Gaussians





Echo example,  $z = x + s \varepsilon$ , with  $s = 1/2$



$$q(x) = x^{(0)} + \frac{1}{2} x^{(1)} + \frac{1}{4} x^{(2)} + \dots$$



$$z = x + \varepsilon$$



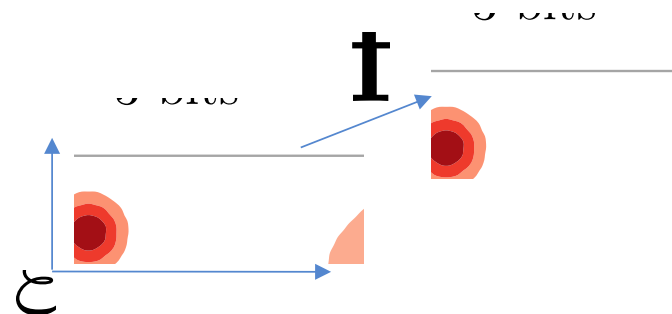
# Sample-dependent noise scaling

$$\mathbf{z} = f(\mathbf{x}) + \varepsilon$$

$$= \mathbb{E}_{\mathbf{x}} H(f(\mathbf{x}) + S(\mathbf{x})\varepsilon)$$

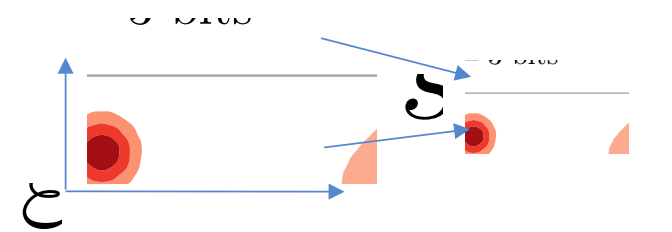
$$= \mathbb{E}_{\mathbf{x}} H(S(\mathbf{x})\varepsilon) =$$

$$= H(\varepsilon) + \mathbb{E} \log |d\varepsilon|$$



$$H(\zeta + t) = H(\zeta)$$

**Translation invariance**



$$H(s\zeta) = H(\zeta) - \dots$$

**Scale property**

# Mutual information for the echo noise channel

$$H(Z) = H(\mathcal{E})$$

Sample-dependent noise s

$$I(Z; X) = H(Z) - H(Z |$$

Self-similar noise (Echo)

Mutual information decom

$$I(Z; X) = -\mathbb{E} \log | C$$

+



# Mutual information for the echo noise channel

Works for any input (sampling noise requires samples of input)

Set  $S(x) = s$  to get a simple, exact  $MI = -\log s$

But  $S(x)$  is controllable (e.g. specify with a neural net) – a powerful way to get more flexible noise models

$$(Z; X) = -\mathbb{E} \log |C| + \frac{\dots}{\dots}$$

Echo results: [arXiv:1904.07199](https://arxiv.org/abs/1904.07199)

Exact mutual information, rather than bound

Better log likelihood bounds

Better rate-distortion trade-offs

Simpler than other state-of-the-art methods that require a complicated “autoregressive flow” model to parametrize non-Gaussianity

**Bigger question: How should we compress?**

**Simple noise (Gaussians, dropout) may not be optimal.**

# Variational Evidence Lower Bound (ELBO)

$$\max_{\theta, \phi} \mathbb{E}_{q_{\phi}} [\log p_{\theta}(\mathbf{x}|\mathbf{z})] -$$

Distortion /  
reconstruction loss

Compression

Reconstruction and higher  
order interactions

# Higher order dependence

Total Correlation / multivariate mutual information (Watanabe, 1960)

$$I \cup (\Lambda) = \sum_i I(\Lambda_i)$$

- Maximized when variables are redundant

“Cohesion” measure of order  $k$

$$C^{(n)}(X) = \frac{1}{\binom{n-1}{k-1}} \sum_{X_A \in \mathcal{E}_k} H(X_A)$$

- Papers discussing: (Fujishige, 1978), Nihat Ay information geometry 2007/2011
- Maximizers correspond to error correcting codes: [arXiv:1811.10839](https://arxiv.org/abs/1811.10839)



# ELBO for higher order interactions

Normal ELBO (can be derived from TC)

$$\mathbb{E} \log p(x) \leq \sum_i \mathbb{E} \log p(x_i | z),$$

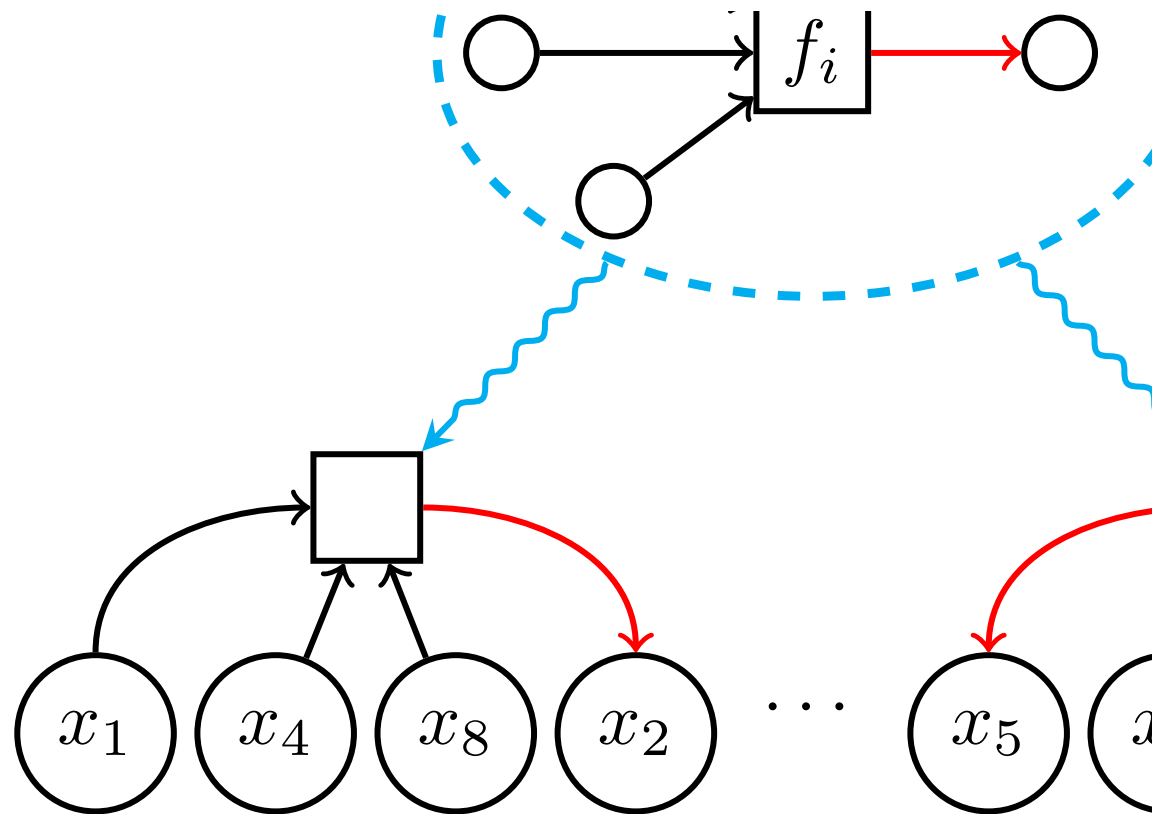
New ELBO based on  $\mathcal{C}^k$

$$\mathbb{E} \log p(x) \geq \frac{1}{\binom{n-1}{k-1}} \sum_{X_A \in \mathcal{E}_k} \mathbb{E} \log p(x)$$

Can we use it to optimize VAEs to detect higher order interactions?  
(preliminary results: maybe!)

# Detect embedded synergies? (challenge idea?)

$x_1$   $x_2$   $x_3$   $x$

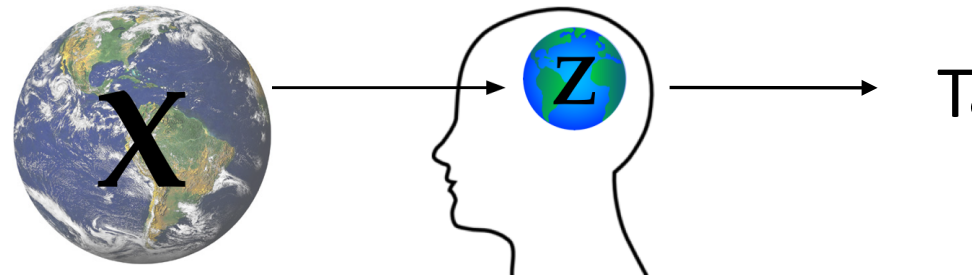


E.g.

J.C: "fraudulent white noise"

Kristian: information lost in high order correlations

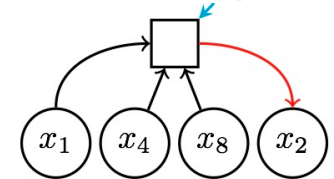
# Conclusion/questions



## Manipulating information in (high-d) representation learning

- Does it matter how we compress? (E.g. Gaussian vs echo)
- What do we want to reconstruct? (marginals versus higher order interactions)
- Fair / invariant representations – How to omit sensitive info?
- Is there an information principle to disentangle?
- Puzzles about the success of adversarial learning
- Do information bottlenecks improve generalization?

## Synergy



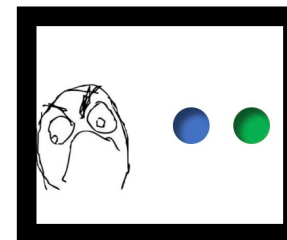
## Practical issues:

- Appropriate measures, *that also have*:
- Efficient estimates or bounds, *and*
- Differentiable / smoothly optimizable

Contact: [gregv@isi.edu](mailto:gregv@isi.edu)

Get these slides: [http://bit.ly/STEEG\\_BS](http://bit.ly/STEEG_BS)

## Grue

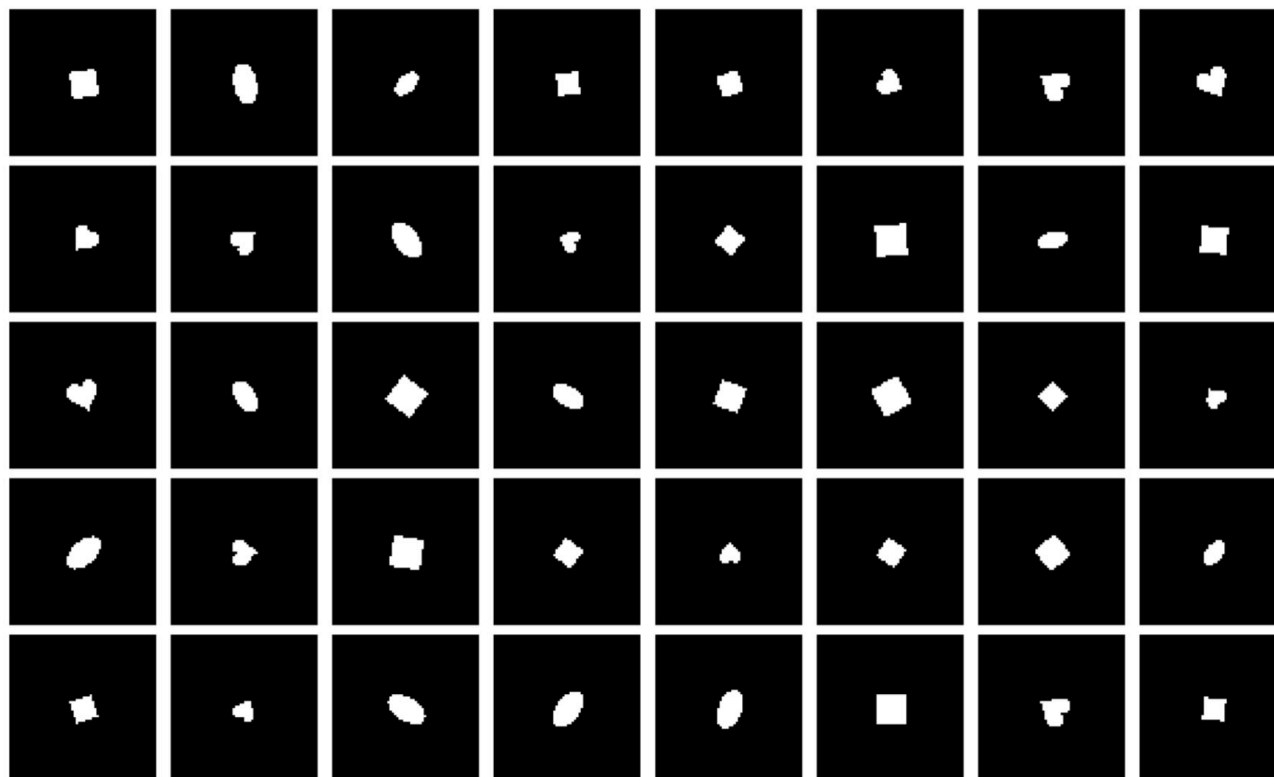


## \$\$ Fair

<b>Application</b>
Age: 37
Income: enough
Zip: 90031

# Disentangling toy data: dSprites

Suppose we observe the following dataset:



Factors:

- Scale
- Shape
- Orientation

This is part of a common artificial test set you'll see in the literature. Generally, we need a test set with known ground truth factors for evaluation

# Adversarial learning and Jensen Shannon Divergence

Examples generated by StyleGAN [arXiv:1809.02926](https://arxiv.org/abs/1809.02926)

Generative Adversarial Networks are the state-of-the-art way to generate realistic images.

Jensen-Shannon Divergence (JSD) tells us how easily we can distinguish two distributions: in this case “real images” and “generated images”.

Instead of minimizing JSD, adversaries that try to distinguish provide a **lower bound** on JSD which is then **minimized**



Figure 2. Uncurated set of images produced by our style-generator (config F) with the FFHQ dataset. Here we used a

# Neural nets as information bottlenecks

inputs) – T ( representation) – Y (labels to predict)

$$\max_T I(X; Y) - \rho I(X; T)$$

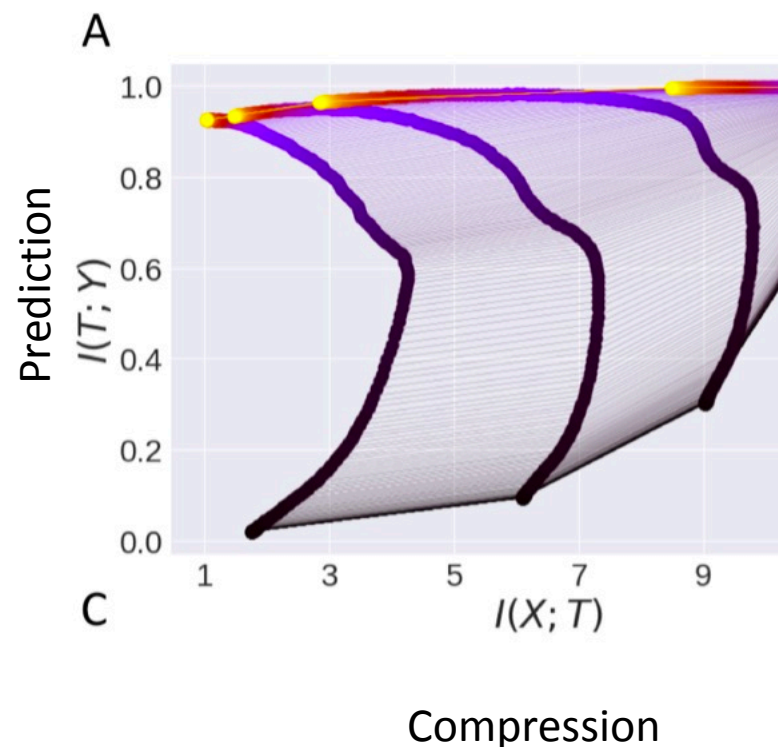
Shwartz-Ziv and N. Tishby. Opening the black box of deep neural networks via information theory. arXiv:1703.00810, 2017.

Neural networks function as bottlenecks with a “fitting phase” and “compression phase”

Compression aids generalization

Shwartz-Ziv et al. “On the information bottleneck theory of deep learning”, ICLR 2018.

Counter-examples for each of Tishby’s claims



# One problem for Tishby's analysis is estimating mutual information

For continuous random variables, if  $z=f(x)$ ,  $I(Z;X) \rightarrow \infty$

Tishby discretizes to estimate the mutual information...

If you discretize differently (or apply invertible transformations with the same discretization), you get different results

Alternatives:

- Add noise to bound mutual information of continuous variables
- Other estimators?

# Use neural nets to estimate mutual information

## Mixture of Gaussians

- Kolchinsky and Tracey, [arXiv:1706.02419](https://arxiv.org/abs/1706.02419)

## Donsker-Varadhan (lower bound on KL/mutual information)

- Belghazi, Mohamed Ishmael, et al. "Mine: mutual information neural estimation." arXiv preprint arXiv:1801.04062 (2018)
- Poole, Ben, Sherjil Ozair, Aaron van den Oord, Alexander A. Alemi, and George Tucker. "On variational lower bounds of mutual information." In NeurIPS Workshop on Bayesian Deep Learning. 2018.

## Impossible with few samples?

- McAllester, David, and Karl Statos. "Formal Limitations on the Measurement of Mutual Information." *arXiv preprint arXiv:1811.04251* (2018).



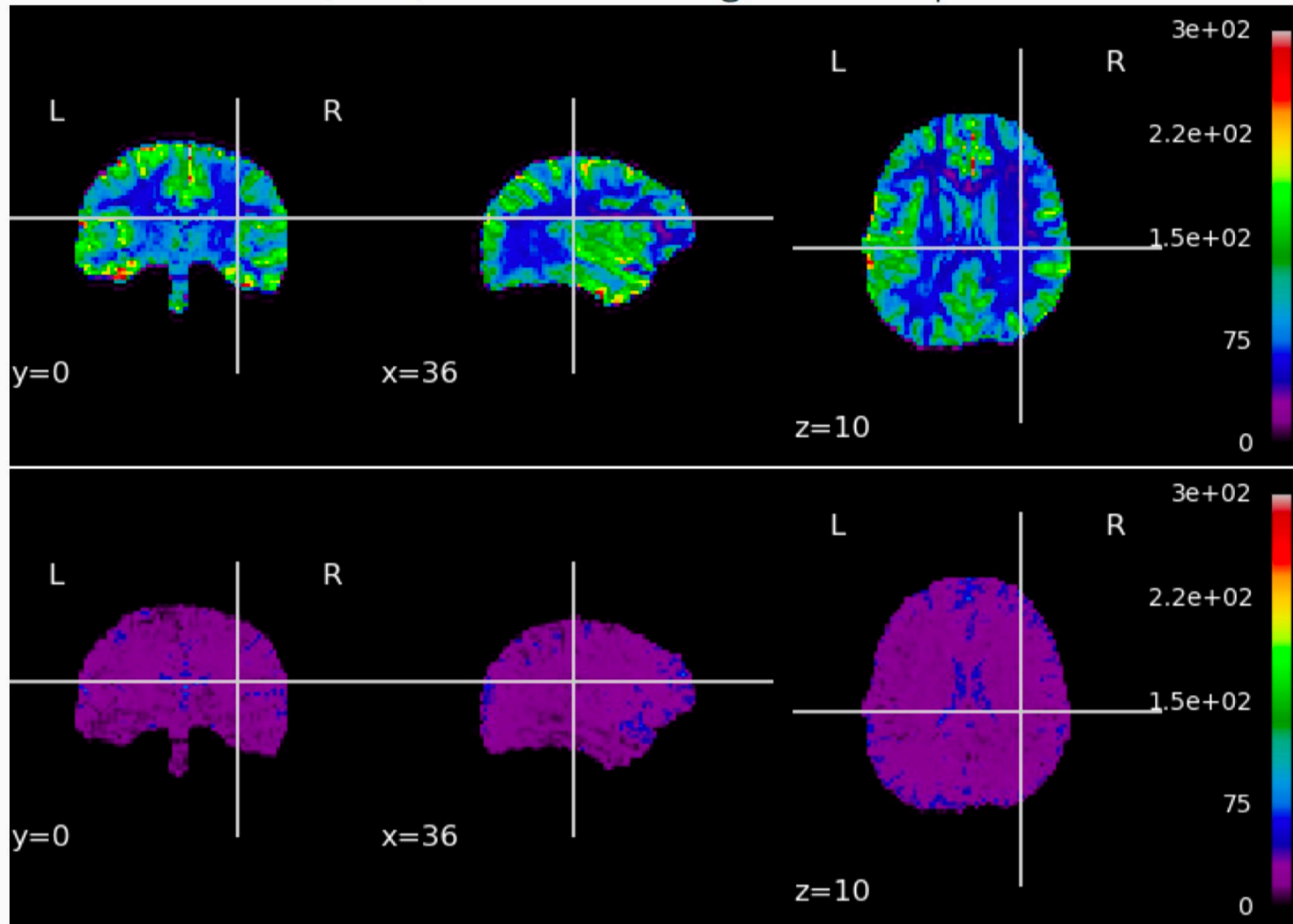
# Other papers I forgot

The Information Complexity of Learning Tasks, their Structure and their Distance. Alessandro Achille \* & Giovanni Paolini † & Glen Mbeng ‡ & Stefano Soatto \*

Information-Theoretic Lower Bounds on BayesRisk in Decentralized Estimation. Aolin Xu and Maxim Raginsky Senior Member, IEEE

# Error Plots – Prisma 30 to Connectom 30

Mirzaalian et al. (2018) Baseline vs Single-site Proposed Method

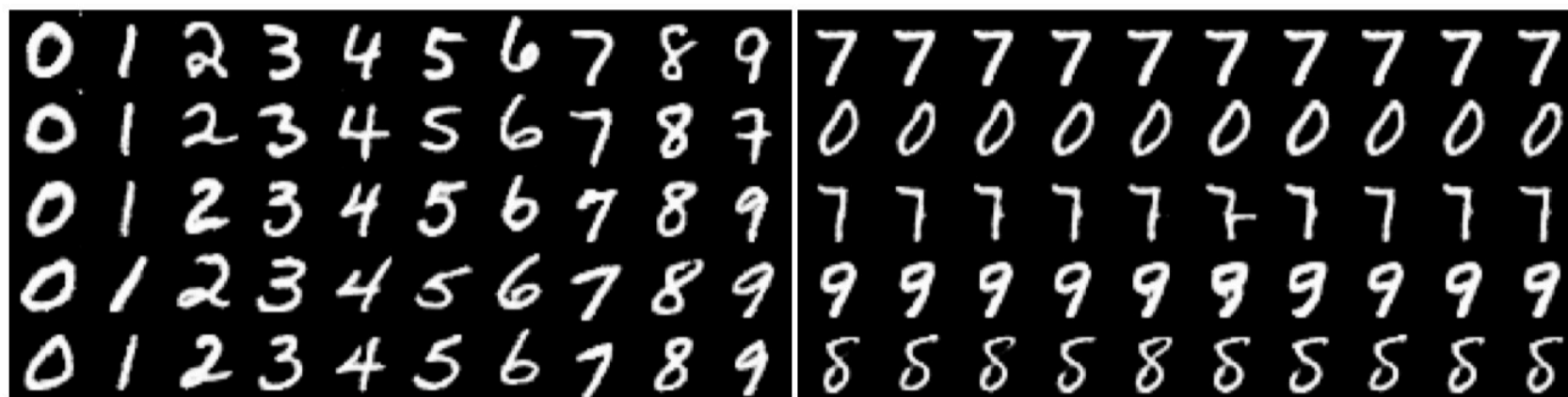


# Interpretability through mutual information regularization

Chen, Yan Duan, Rein Houthoofd, John Schulman, Ilya Sutskever, Pieter Abbeel (2016) InfoGAN: Interpretable Representation Learning by Information Maximizing Generative Adversarial Nets

Guide most information into a special neuron,  $c_1$

$$\text{Max Loss} + I(c_1; x)$$



(a) Varying  $c_1$  on InfoGAN (Digit type)

(b) Varying  $c_1$  on regular GAN (No clear meaning)

Another example of this phenomenon: [arXiv:1802.05822](https://arxiv.org/abs/1802.05822)