

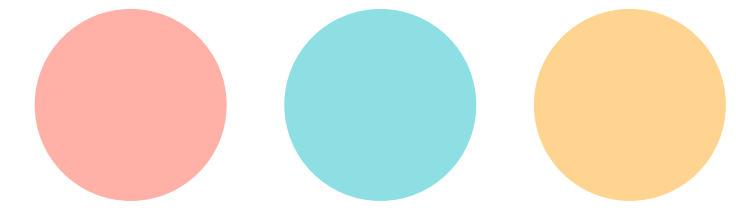


Stanford

Tripod

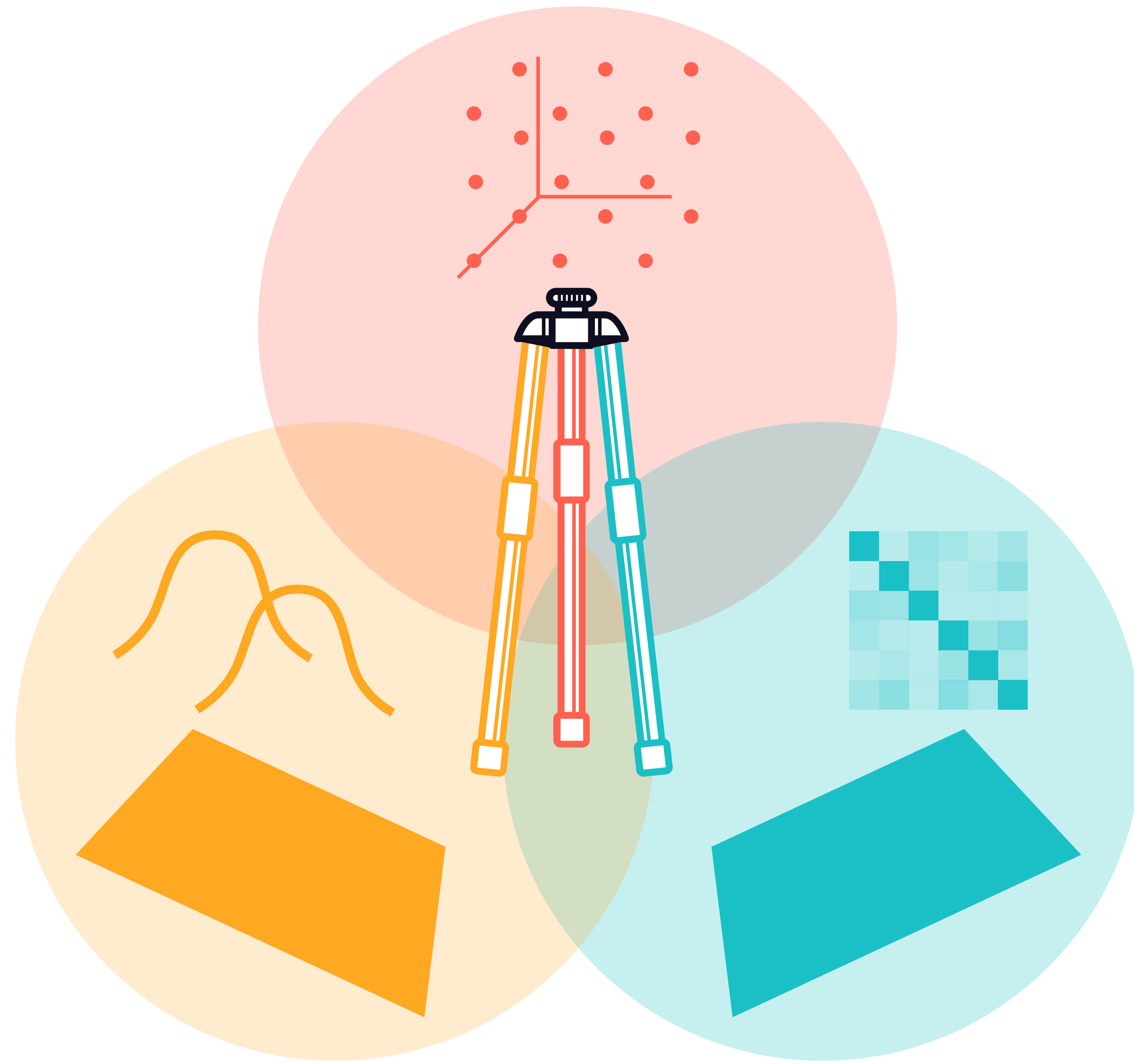
three complementary inductive biases for disentangled representation learning

Kyle Hsu\*, Jubayer Ibn Hamid\*, Kaylee Burns, Chelsea Finn, Jiajun Wu



**Inductive biases** facilitate disentanglement by paring down the solution space, but prior works largely propose and validate one new approach in isolation from others.

**latents** compressed and organized via quantization  
Hsu et al., 2023



**encoding** into independent latents  
Chen et al., 2018

**decoding** with small mixed derivatives  
Peebles et al., 2020

Tripod makes necessary adaptations to three **complementary** inductive biases to meld them into a state-of-the-art disentangling autoencoder.



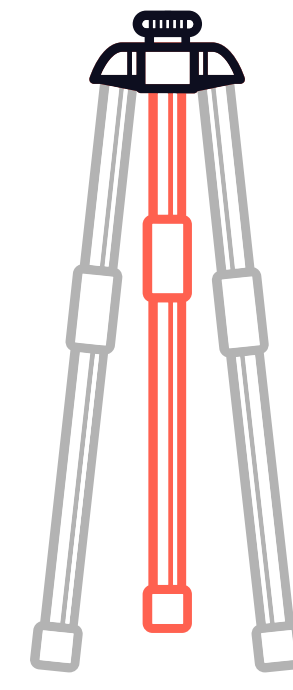
## preliminaries

- unlabelled data is generated noiselessly from independent sources

$$p(s) = \prod_{i=1}^{n_s} p(s_i) \quad x = g(s)$$

- goal: learn autoencoder whose latents recover sources

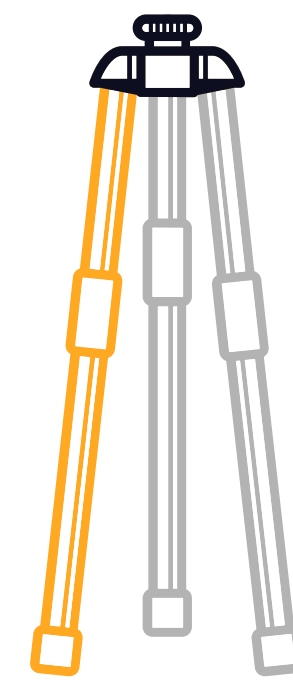
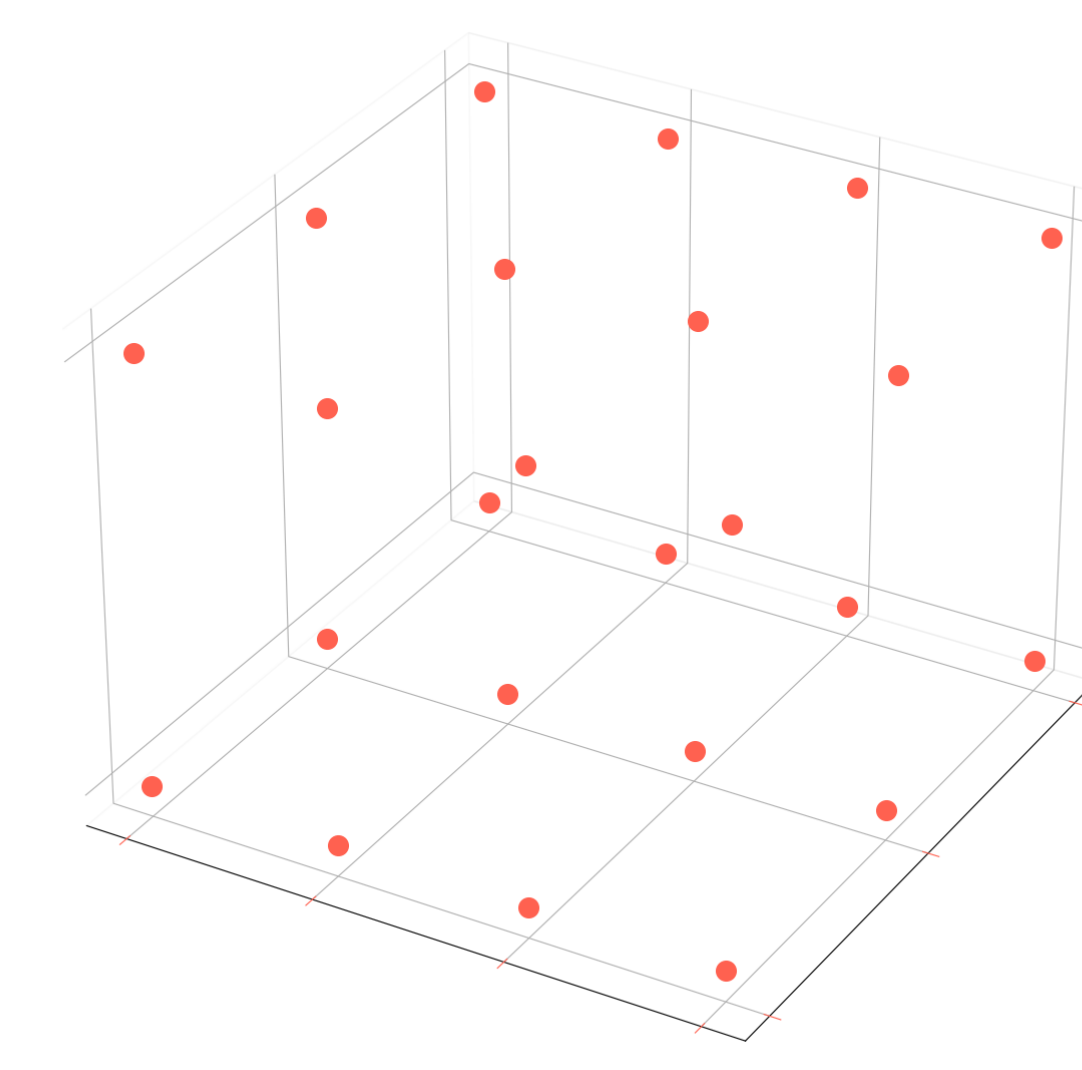
$$z = \hat{g}^{-1}(x) \quad \hat{x} = \hat{g}(z)$$



## finite scalar latent quantization

- motivation: true sources are highly compressed and organized
- goal: impose quantized, grid-like structure on latent space
- problem: dictionary learning destabilizes other components
- solution:** use finite scalar quantization (Mentzer et al., 2023)

$$z = \frac{2}{n_q - 1} \text{round} \left( \frac{n_q - 1}{2} (\tanh(\hat{g}^{-1}(x)) + 1) \right) - 1$$



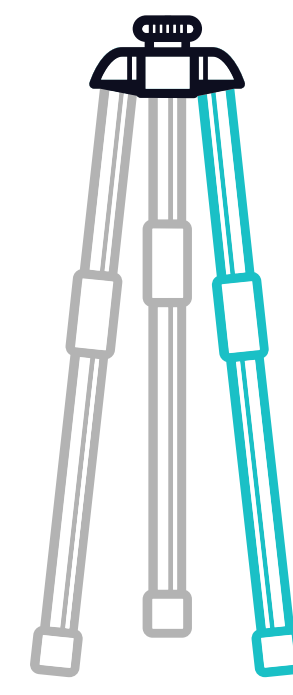
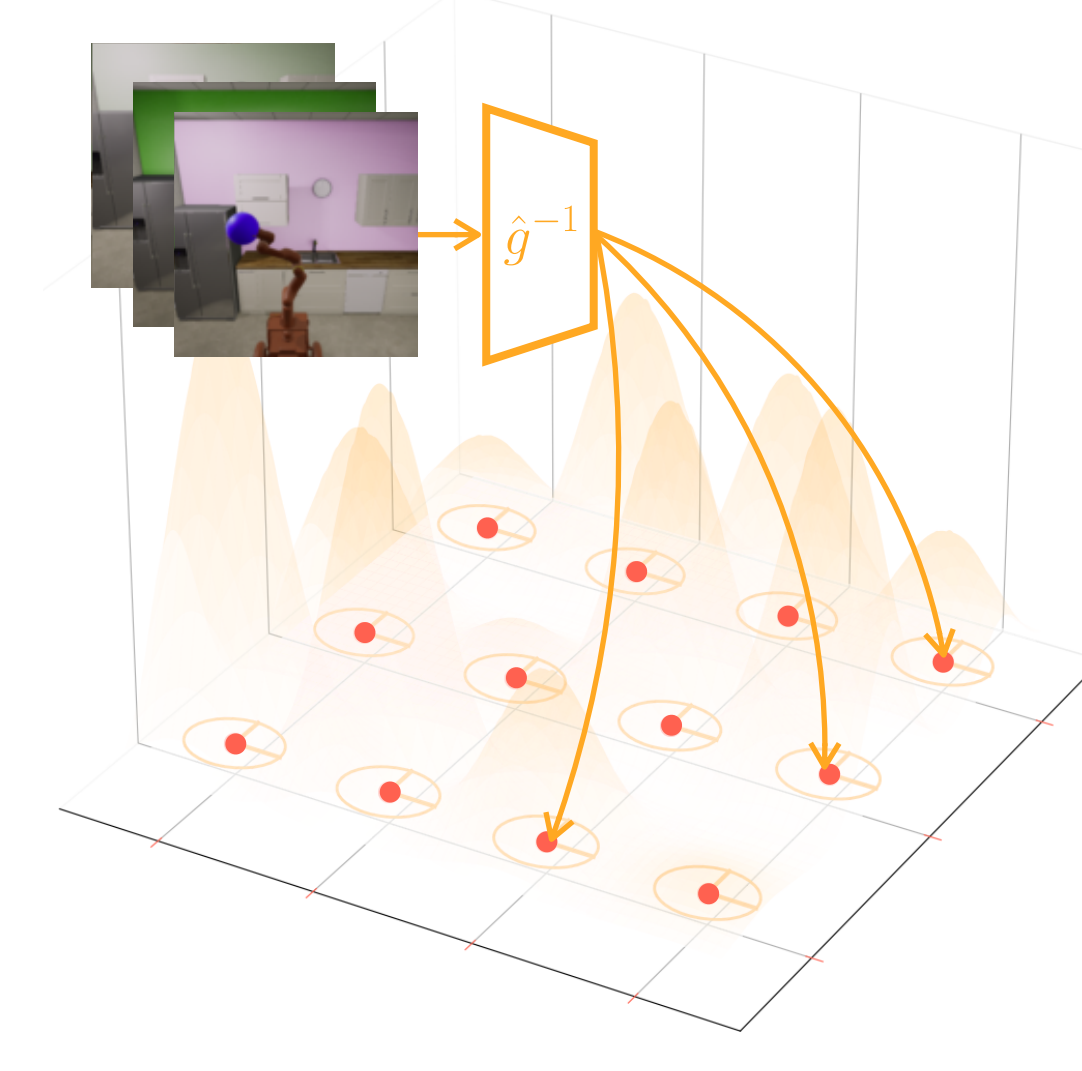
## kernel-based latent multiinformation

- motivation: true sources are collectively independent
- goal: regularize latent multiinformation

$$D_{\text{KL}} \left( q(z) \parallel \prod_{j=1}^{n_s} q(z_j) \right)$$

- problem: quantized latents aren't probabilistic
- solution:** use Gaussian kernel density estimation

$$q(z) \propto \sum_{i=1}^{n_b} \exp \left( -\frac{1}{2} (z - z^{(i)})^\top S^{-1} (z - z^{(i)}) \right)$$



## normalized Hessian penalty

- motivation: true sources interact minimally to generate data
- goal: regularize off-diagonal entries of decoder Hessians

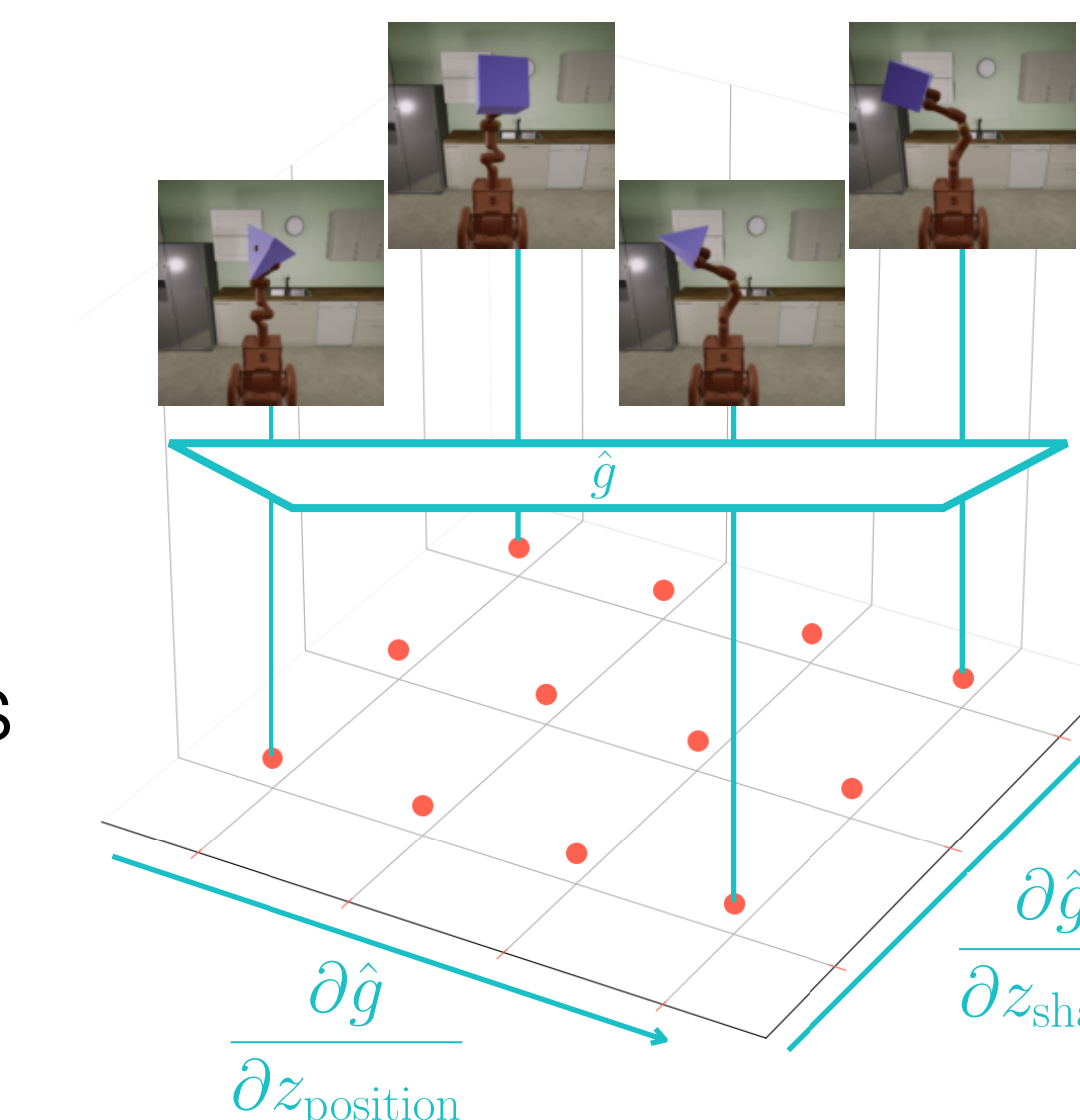
$$\sum_{j_1 \neq j_2} \left( H_{j_1 j_2}^{[k]} \right)^2 \quad H_{j_1 j_2}^{[k]} = \frac{\partial \hat{g}^{[k]}}{\partial z_{j_1} z_{j_2}}$$

- problem: sensitive to trivial rescalings of latents and activations
- solution:** replace with a normalized quantity

$$\frac{\sum_{j_1 \neq j_2} \left( H_{j_1 j_2}^{[k]} \sigma_{j_1} \sigma_{j_2} \right)^2}{\sum_{j_1, j_2} \left( H_{j_1 j_2}^{[k]} \sigma_{j_1} \sigma_{j_2} \right)^2}$$

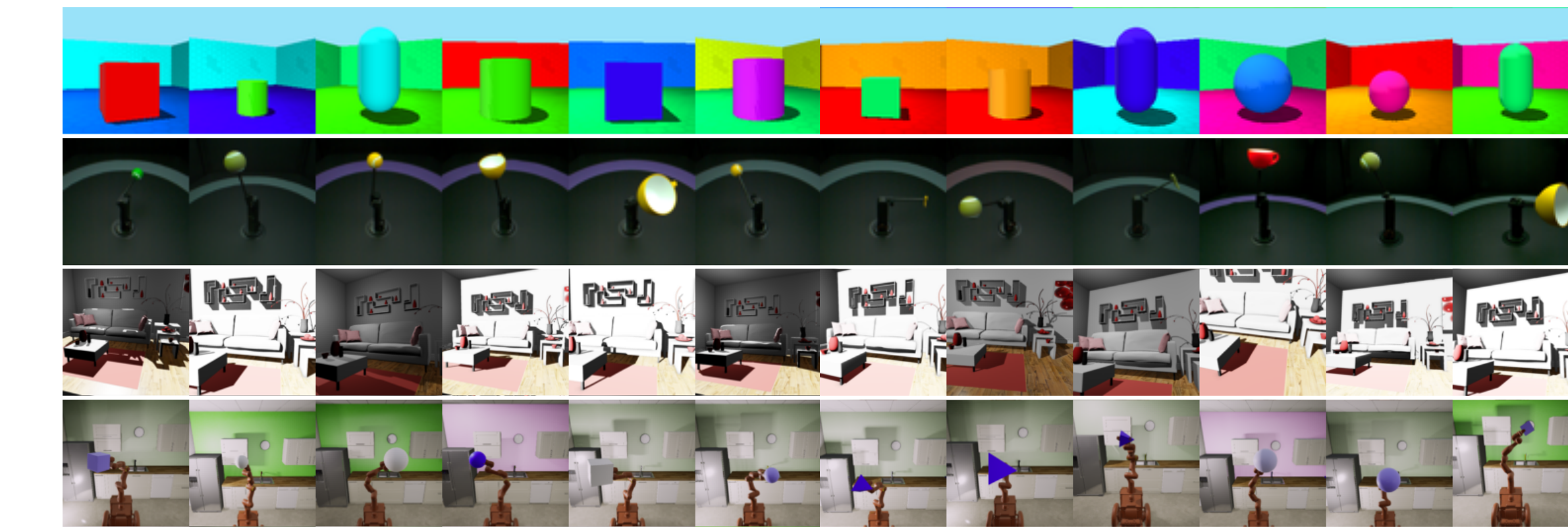
and associated estimator

$$\frac{\text{Var} [v^\top H^{[k]} v]}{\text{Var} [w^\top H^{[k]} w]} \quad v_j \sim \text{Rademacher}(\sigma_j) \\ w_j \sim \mathcal{N}(0, \sigma_j^2)$$



## experiments

- datasets



Shapes3D 6 sources  
MPI3D 7 sources  
Falcor3D 7 sources  
Isaac3D 9 sources

- Tripod greatly improves upon prior methods** that use only one of its legs

	InfoMEC: (InfoModularity, InfoCompactness, InfoExplicitness)									
	aggregated	Shapes3D		MPI3D		Falcor3D		Isaac3D		
$\beta$ -TCVAE	(0.68 0.43 0.88)	(0.86 0.45	<b>1.00</b> )	(0.52 0.46 0.75)	(0.62 0.39	<b>0.82</b> )	(0.72 0.42	<b>0.94</b> )		
QLAE	(0.62 <b>0.57</b> 0.77)	(0.68 0.55	<b>0.98</b> )	(0.45 0.42 0.61)	<b>(0.71 0.71</b> 0.72)	(0.65 0.61 0.78)				
Tripod (ours)	<b>(0.78 0.59 0.90)</b>	<b>(0.94 0.59 1.00)</b>		<b>(0.64 0.53 0.84)</b>	<b>(0.72 0.56 0.82)</b>	<b>(0.84 0.68 0.95)</b>				

- ablating each Tripod leg in turn shows that **all three legs are necessary for best performance**

	aggregated	Shapes3D		MPI3D		Falcor3D		Isaac3D		
Tripod (ours)	<b>(0.78 0.59 0.90)</b>	<b>(0.94 0.59 1.00)</b>		<b>(0.64 0.53 0.84)</b>	<b>(0.72 0.56 0.82)</b>	<b>(0.84 0.68 0.95)</b>				
Tripod w/o FSLQ	(0.56 0.46 <b>0.92</b> )	(0.69 0.48 <b>1.00</b> )		(0.43 0.40 <b>0.97</b> )	(0.54 0.41 <b>0.84</b> )	(0.57 0.54 0.87)				
Tripod w/o KLM	(0.73 0.50 <b>0.90</b> )	(0.89 <b>0.57 1.00</b> )		(0.57 0.50 0.80)	<b>(0.74 0.54 0.82)</b>	(0.72 0.38 <b>0.96</b> )				
Tripod w/o NHP	(0.70 0.48 0.89)	(0.85 0.46 <b>1.00</b> )		(0.60 0.50 0.81)	(0.59 0.40 0.81)	(0.75 0.57 0.93)				

- case study: naively combining Tripod's inductive biases fails to disentangle "robot x" and "robot y"

