



## Interpretable predictions for science

One goal of machine learning for science:

- a model producing **interpretable predictions** that uses
- **learned concepts** close to the 'true' mechanism

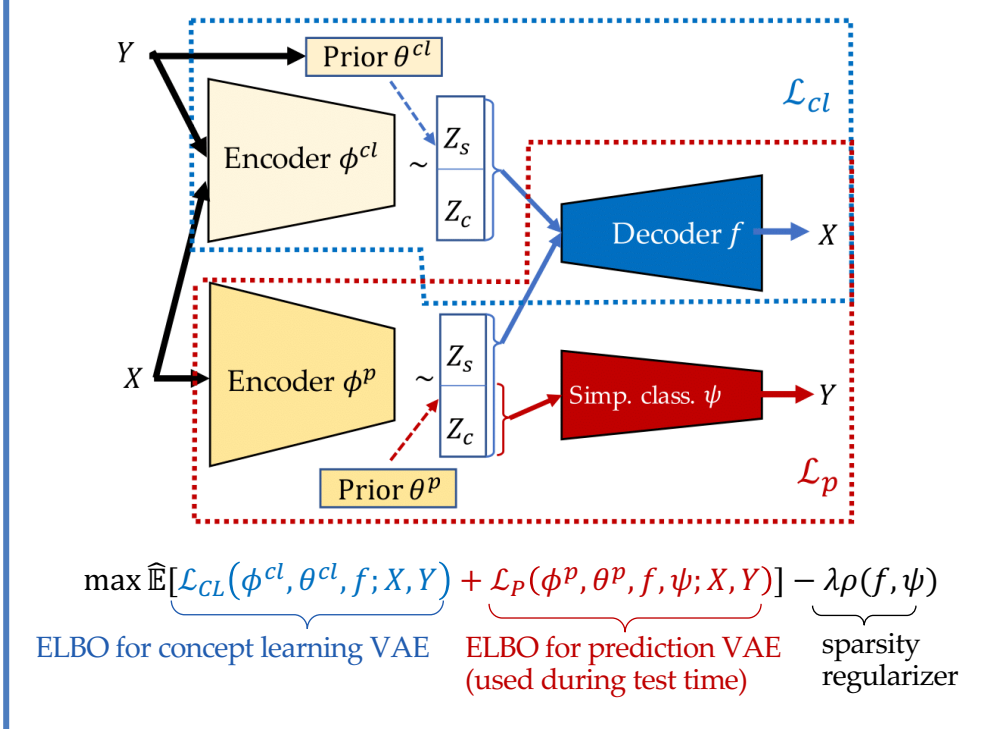
**Our contribution: CLAP, a VAE-based classifier**

- that provably learns unlabeled ground-truth concepts
- uses these learned features to obtain optimal predictions

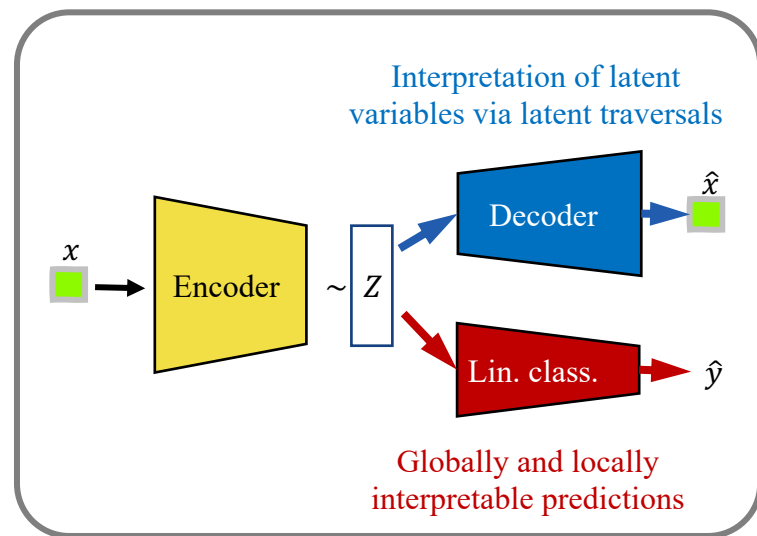
## Comparison with related work

Properties	Post-hoc explanations			Inherently interpretable	
	pixel attribution+ counterfactuals	pre-defined concepts	StyleGANs	existing VAEs/ autoencoders	CLAP
Learning visually distinct features	×	×	✓	✓*	✓
Global importance of predictive features	×	✓	×	×	✓
Guarantees: concept learning+prediction	×	×	×	×	✓

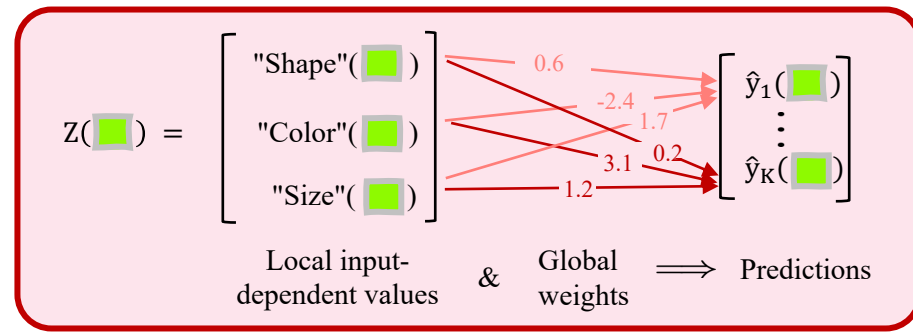
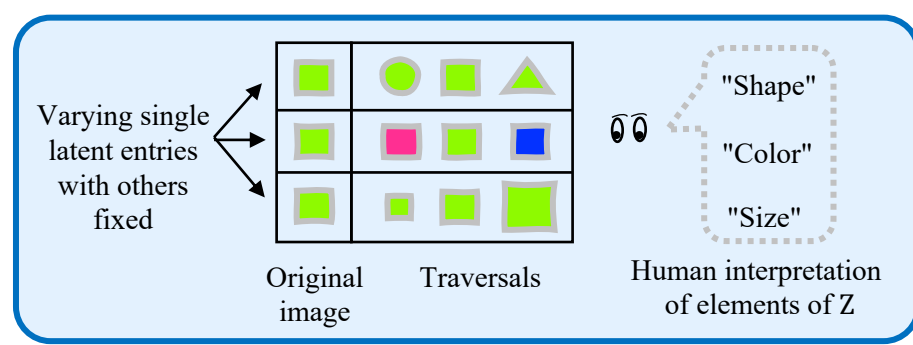
## Training CLAP



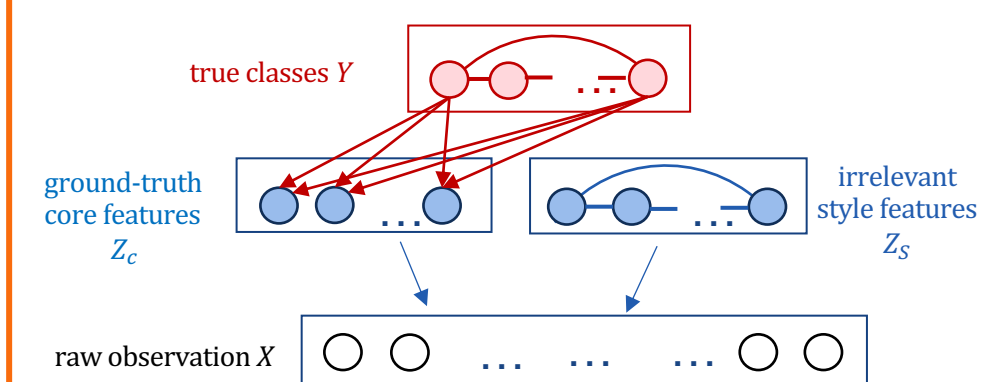
## C(oncept) L(earning) a(nd) P(rediction)



Trained predictive VAE



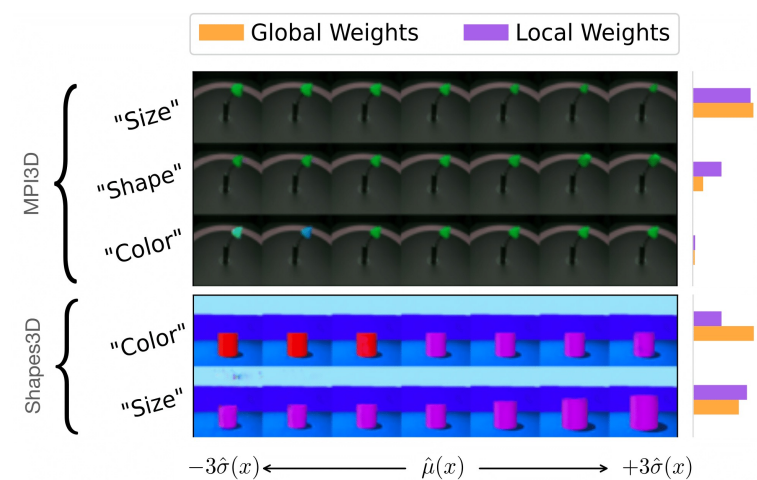
## Theoretical guarantees (informal)



Under above data-generating and further regularity and heterogeneity assumptions, in the infinite data limit:

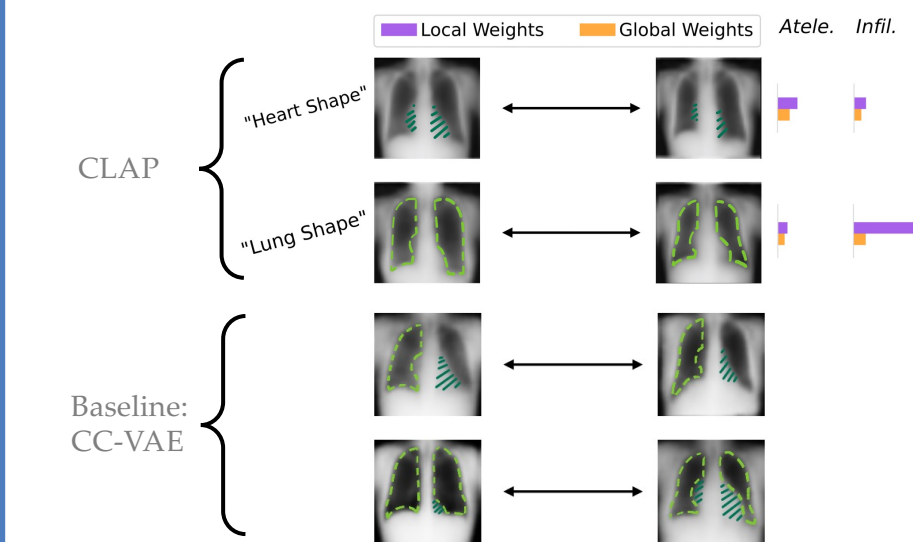
- CLAP identifies  $Z_c$  up to permutation and scaling
- CLAP learns the **optimal classifier that uses  $Z_c$**

## Example I: CLAP on toy datasets



Recovers ground-truth concepts + 99% class. accuracy!

## Example II: CLAP on the Chest-Xray dataset



At similar accuracy 90%, CLAP achieves better disentanglement & indicates global and local feature importance

## Work in progress

- Using GANs for sharper reconstructions
- Aggregate data across multiple sources
- Incorporating dependence between style features and  $Y$