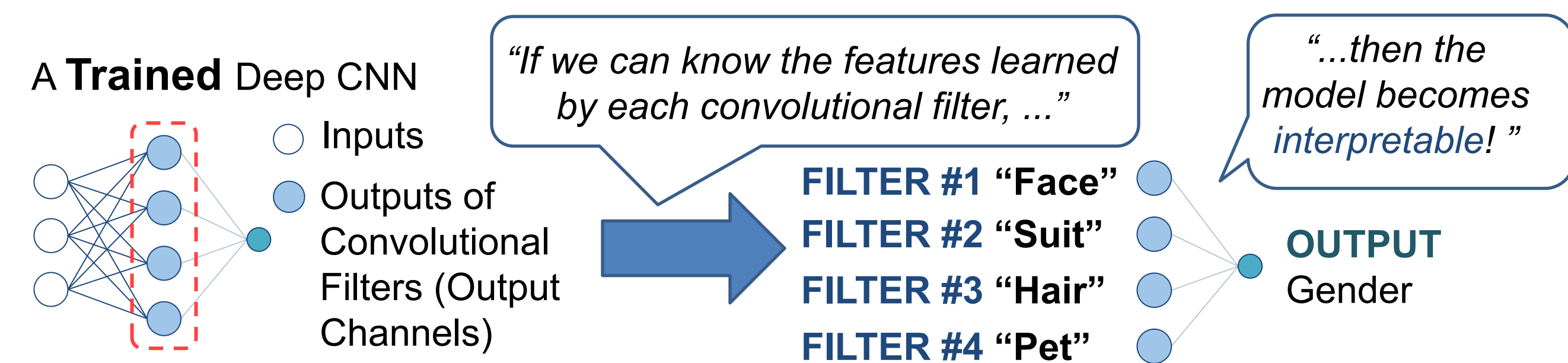
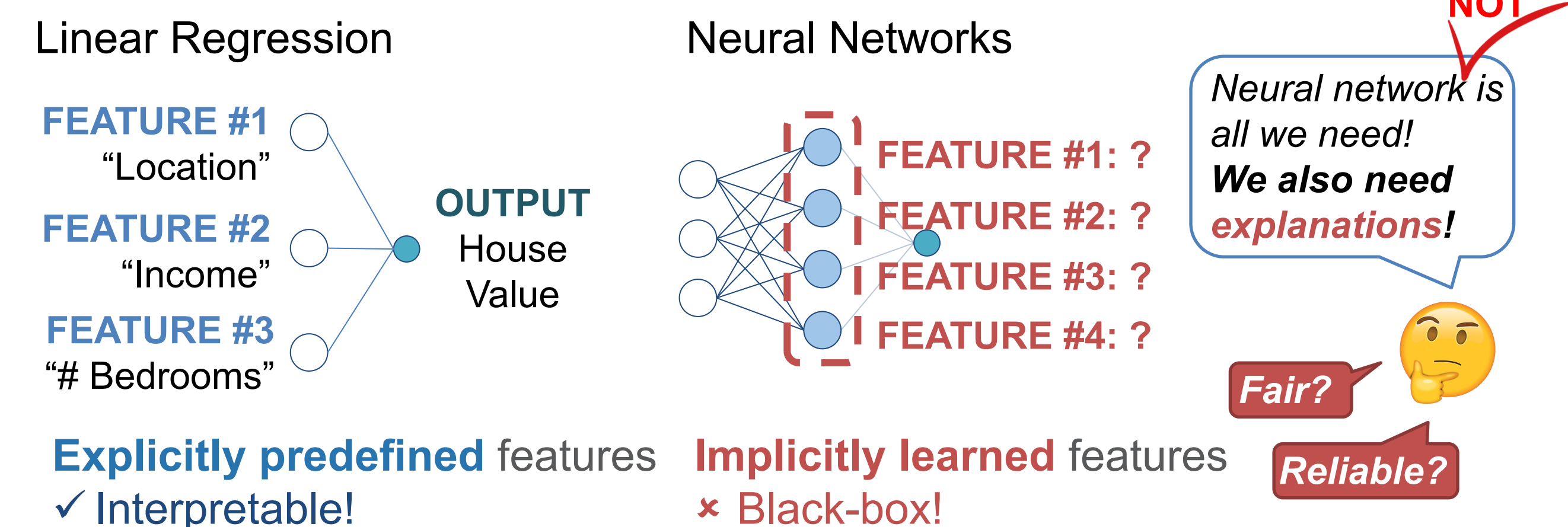
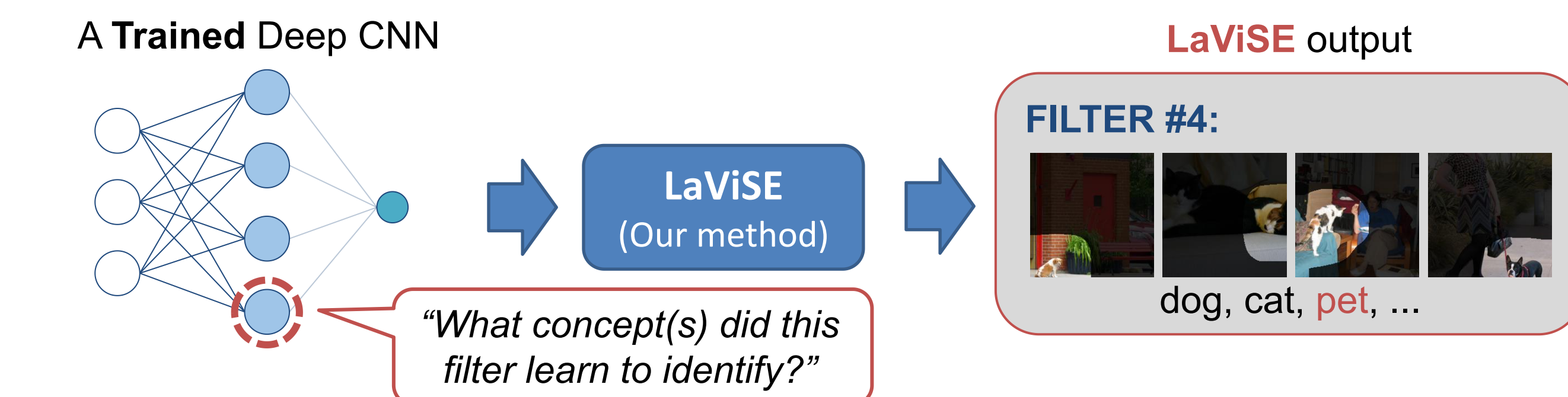


WHY do we need to explain deep CNN?



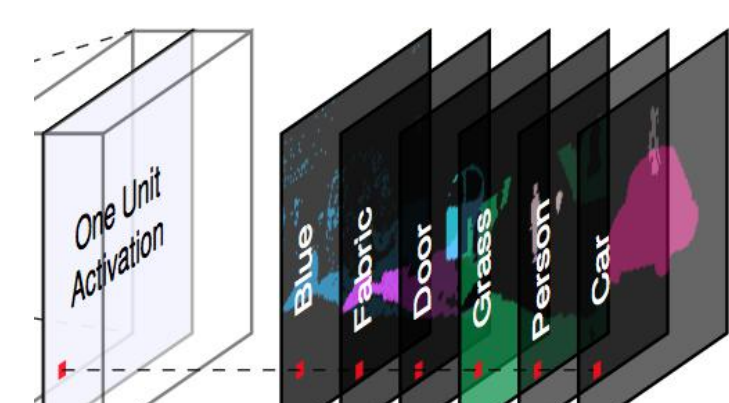
LaViSE -- Latent Visual-Semantic Explainer for CNN



EXISTING APPROACHES



⊗ Only show important pixels...
No explicit semantic explanations!



⊗ Only use annotated concepts...
Cannot explain unseen concepts!

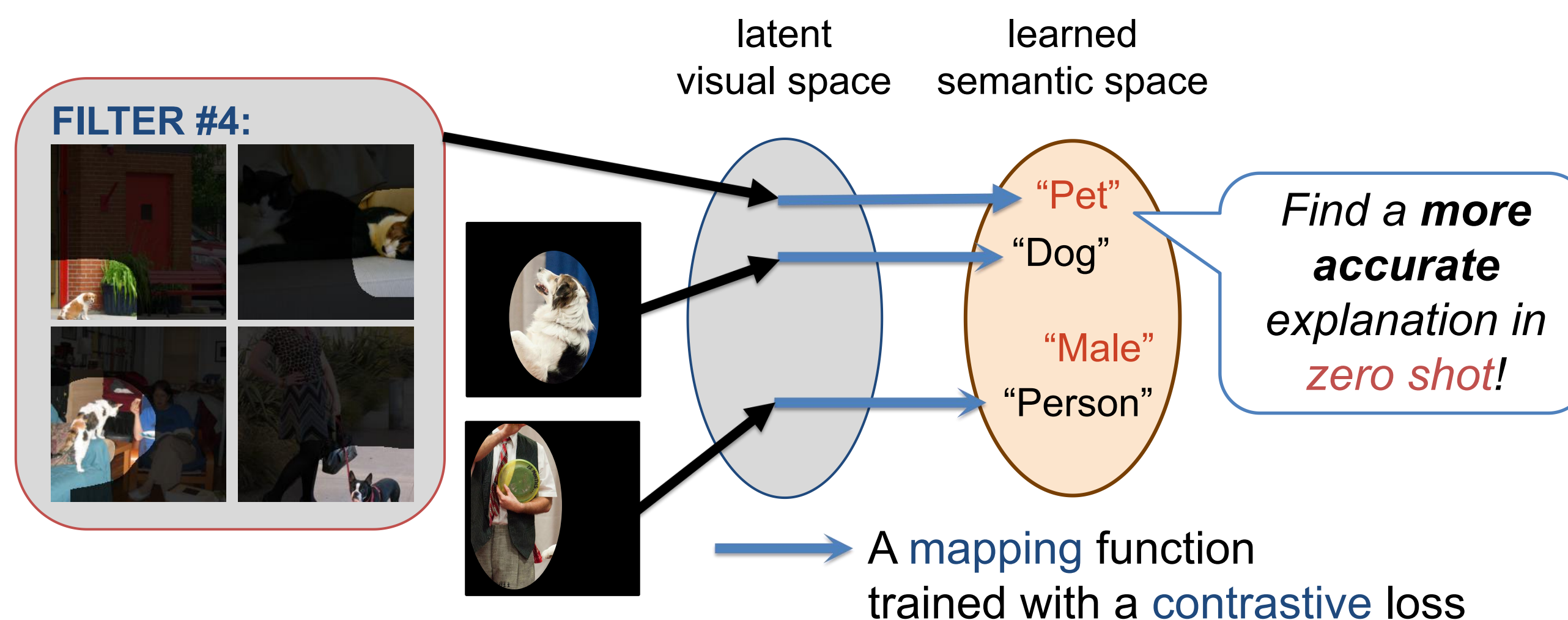
LaViSE

- ⊙ Provides **visual+semantic** explanations!
- ⊙ Finds the best explanation from **annotated+zero-shot** concepts!

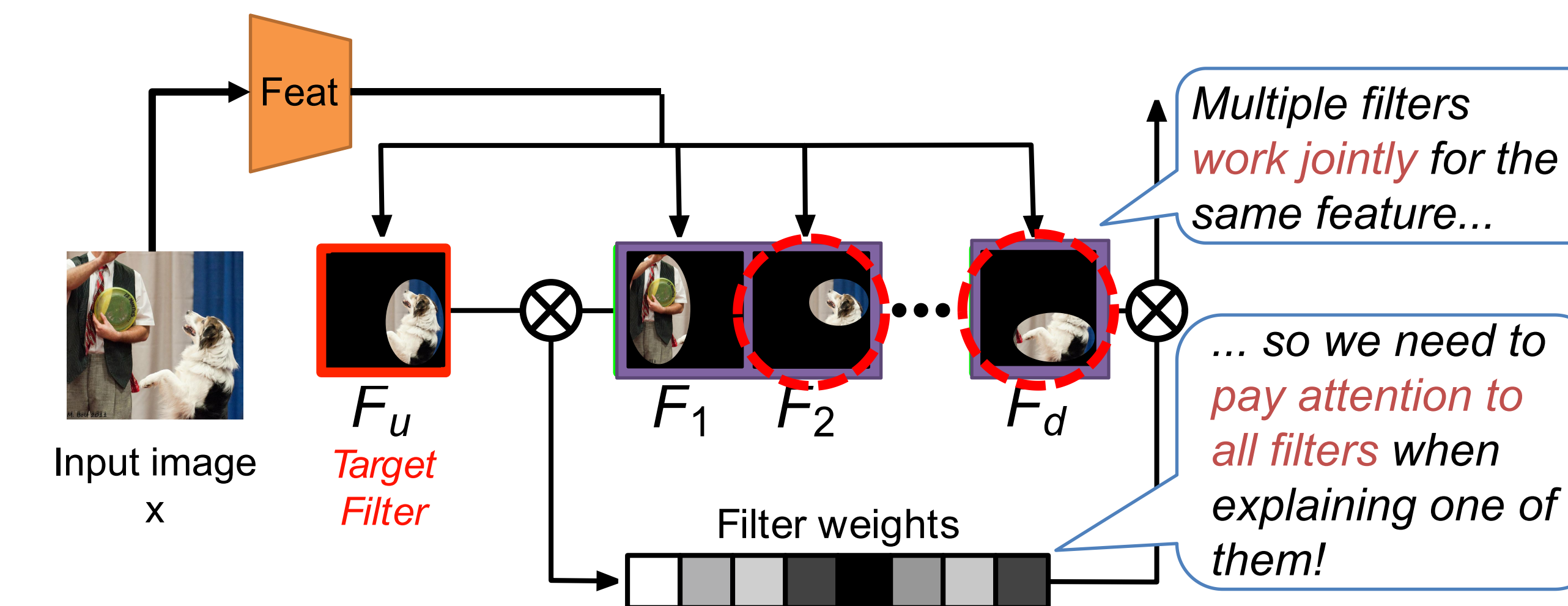
Plus, the explanations are more accurate with the **filter attention**!

HOW does LaViSE work?

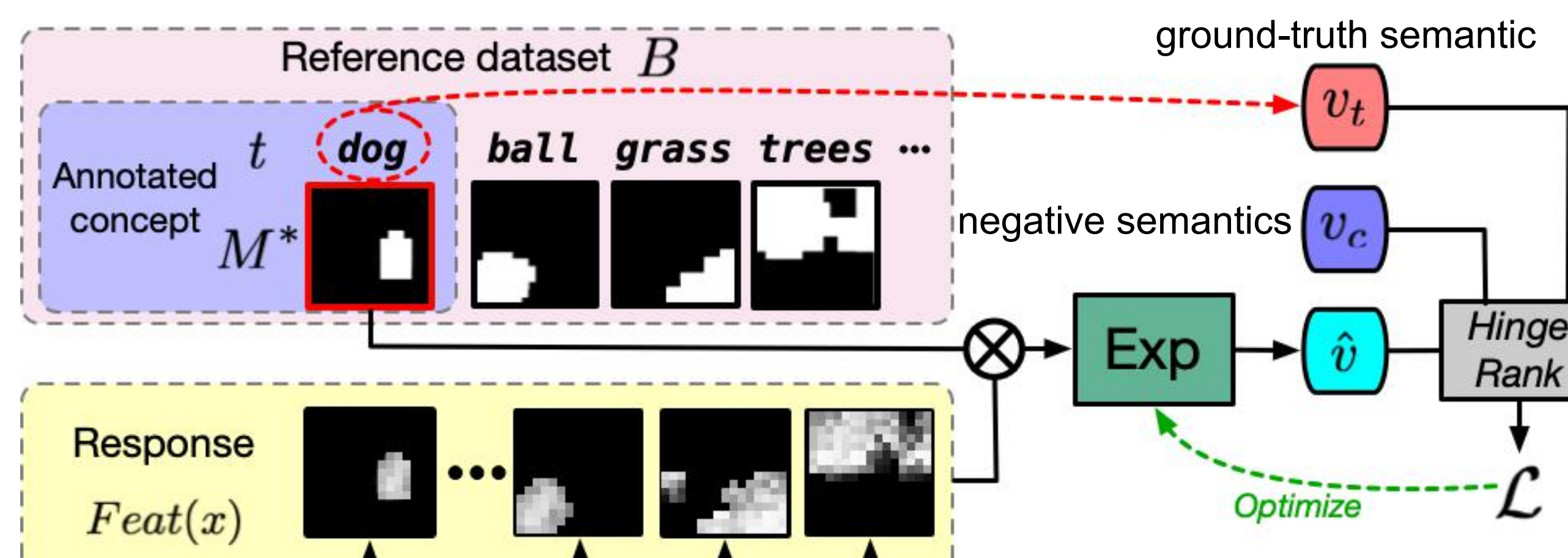
Key component #1: Latent Visual-Semantic Mapping



Key component #2: Filter-level Attention

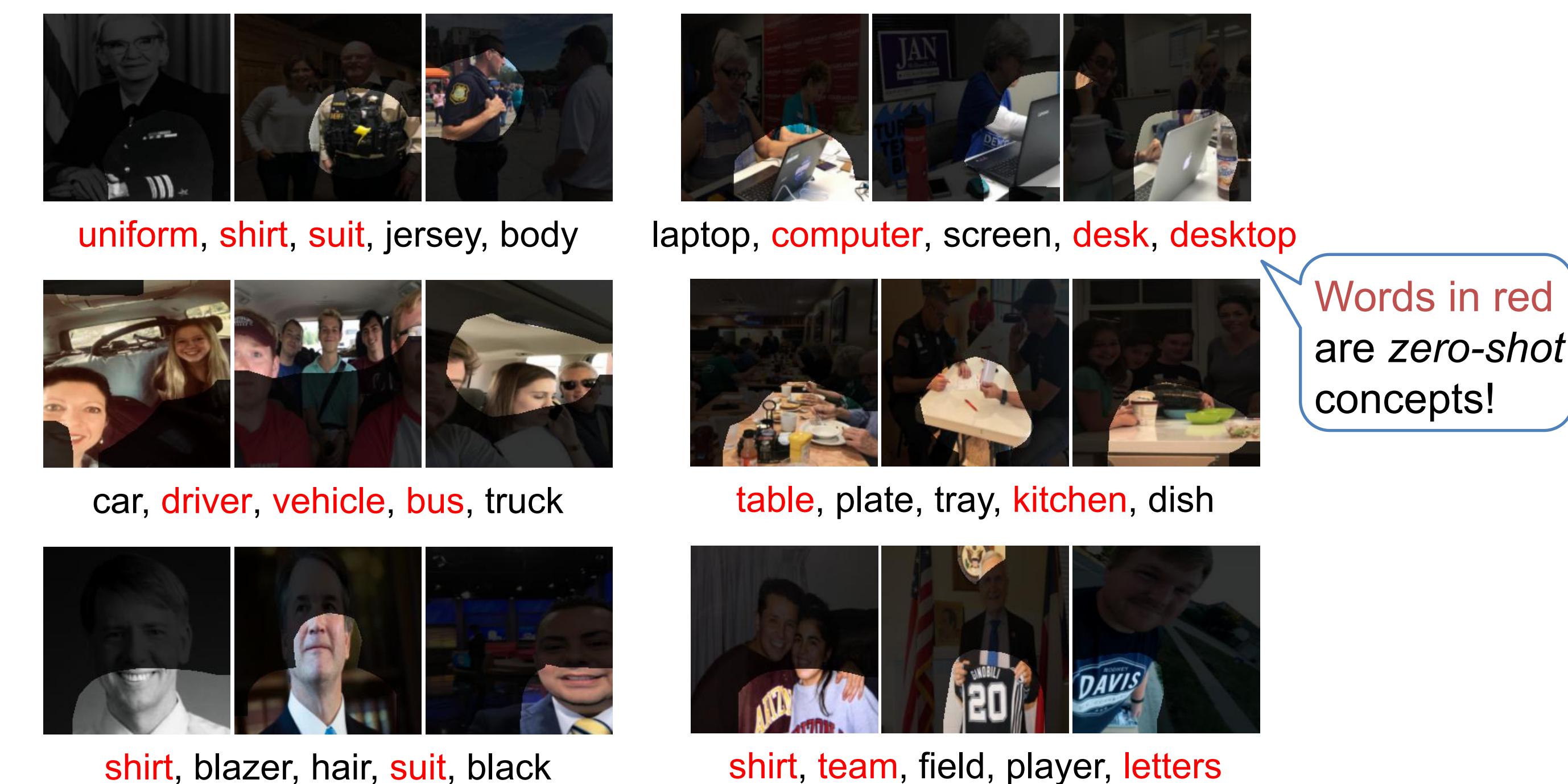


Training LaViSE



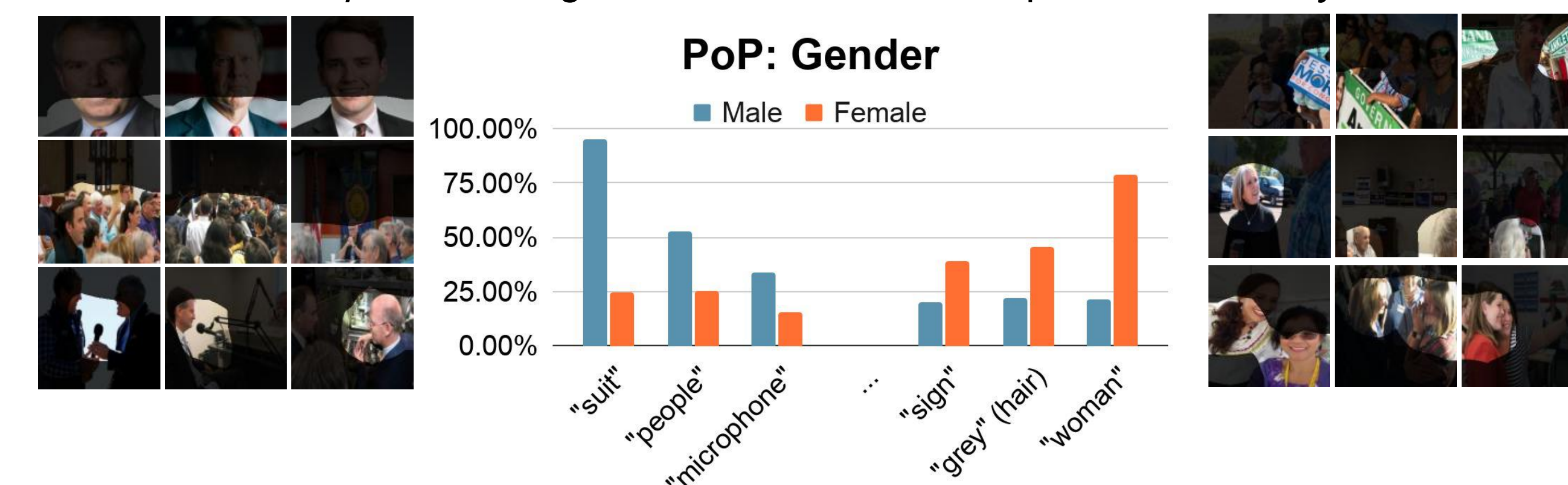
WHAT can LaViSE do?

Explaining Trained Networks with Zero-shot Concepts



Explaining Unlabeled Datasets with Pretrained Networks

LaViSE can provide insights of a dataset and help with bias analysis!



Summary/Conclusion

- We proposed **LaViSE**, a novel framework which can both visually and semantically explain latent representations of a trained CNN.
- **LaViSE** enables users to discover concepts that a CNN learned without being explicitly taught.