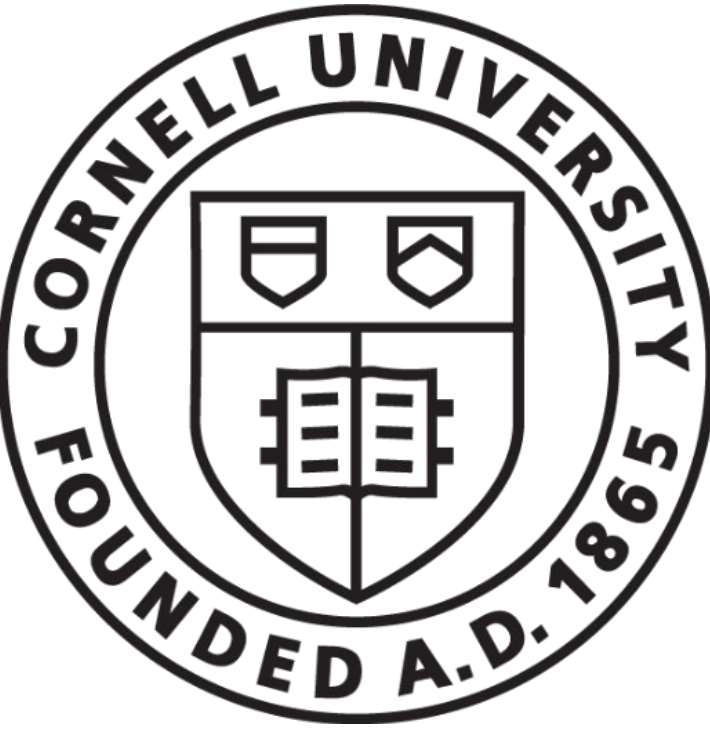


# Debugging Concept Bottleneck Models through Removal and Retraining

Eric Enouen, Sainyam Galhotra  
Cornell University



## CBMs

Bottleneck prediction through an interpretable concept layer

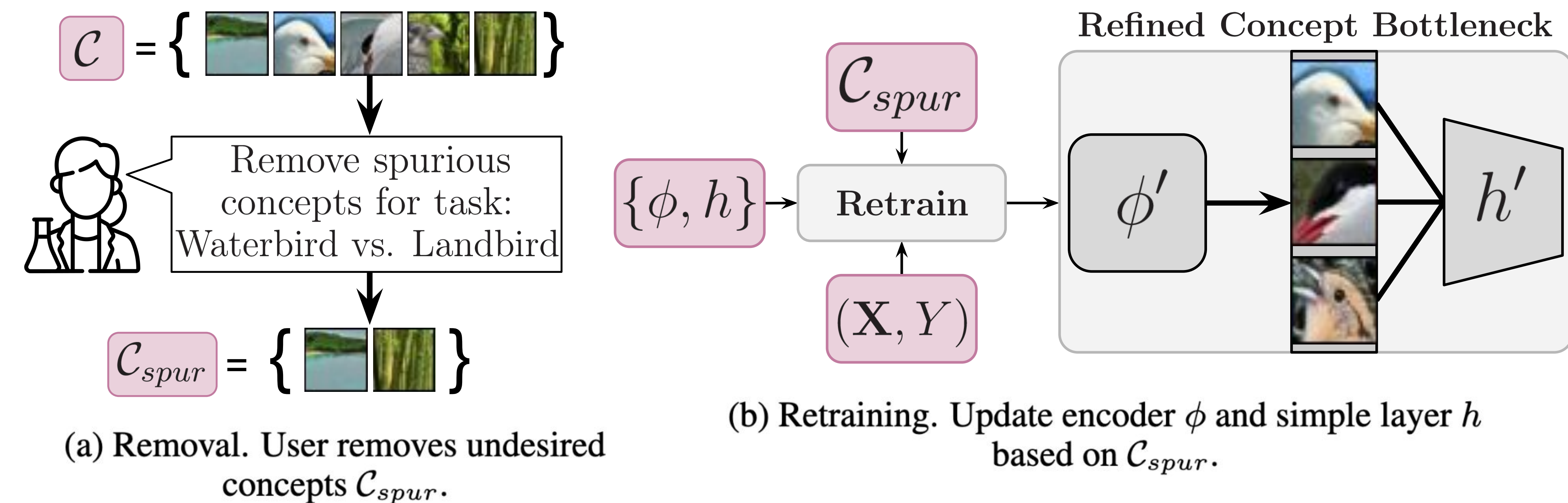
- Enables domain experts to analyze which concepts are being used for predictions
- Enable test-time interventions on mispredicted concepts

## VLM-CBMs

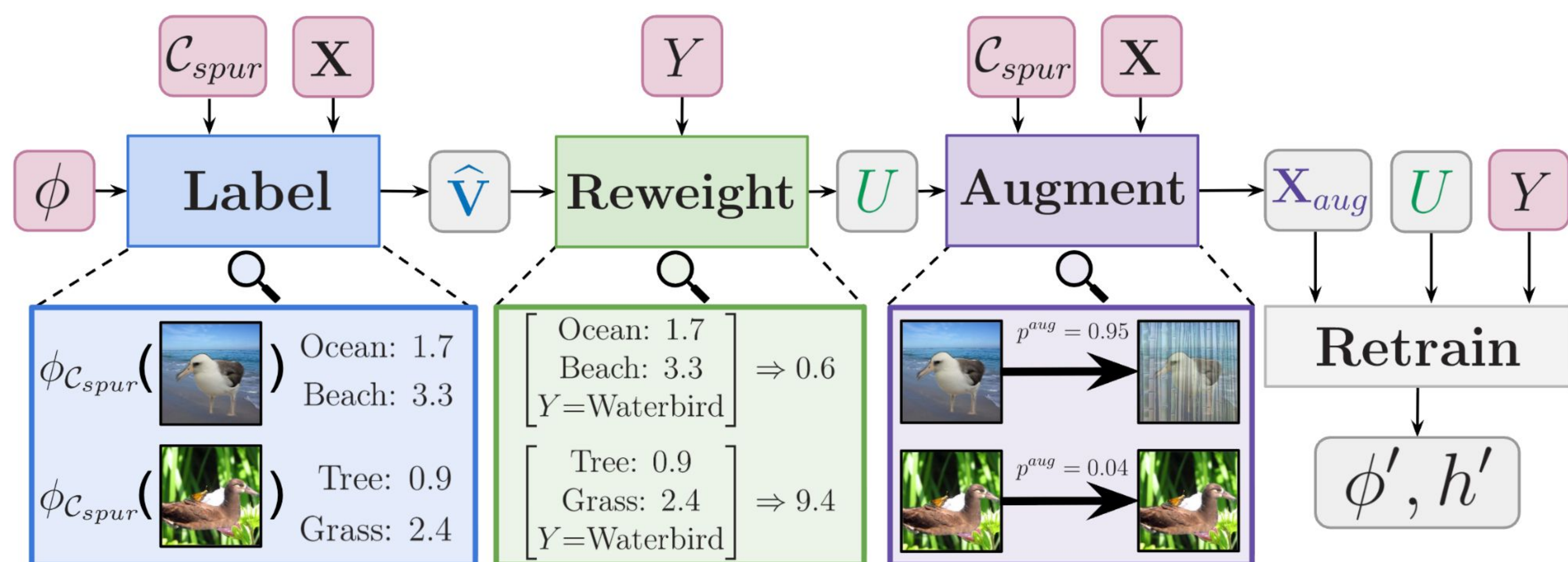
Remove the need for expensive concept annotations using foundation models

- Can introduce spurious concepts that fail under subpopulation shift!
- **How can we realign spurious concepts with expert reasoning?**

## Removal + Retraining Framework



## CBDebug (Concept Bottleneck Debugger)



## CBDebug Improves Robustness to Subpopulation Shift

Real	PIP-Net		Post-hoc CBM	
	Waterbirds	MetaShift	Waterbirds	MetaShift
Original	71.9 $\pm$ 2.7	52.4 $\pm$ 2.0	25.8 $\pm$ 3.0	84.5 $\pm$ 2.2
<b>CBDebug</b>	<b>79.4<math>\pm</math>4.3</b>	<b>57.3<math>\pm</math>3.1</b>	<b>51.9<math>\pm</math>16.2</b>	<b>89.3<math>\pm</math>1.3</b>
Automated	Post-hoc CBM			
	Waterbirds	MetaShift	CelebA	ISIC (AUROC)
Original	25.8 $\pm$ 3.0	84.5 $\pm$ 2.2	8.7 $\pm$ 0.9	39.3 $\pm$ 3.7
<b>CBDebug</b>	<b>58.3<math>\pm</math>6.0</b>	<b>87.5<math>\pm</math>2.8</b>	<b>51.3<math>\pm</math>3.9</b>	<b>58.0<math>\pm</math>11.6</b>

## Learns Better Core Concepts

Class	Original	Retrain	CBDebug
Waterbird	hooked seabird beak	beach	duck-like body
	sea	gull-like body	hooked seabird beak
	harbor	water	orange wings
	lake	hooked seabird beak	orange eyes
	gull-like body	duck-like body	orange nape
Landbird	olive crown	olive upper tail	olive upper tail
	tree-clinging-like body	bamboo	iridescent bill
	forest	green primary color	blue upper tail
	olive upper tail	tree-clinging-like body	olive crown
	tree	olive breast	hawk-like body

## Conclusions

CBDebug improves worst-group accuracy

- **across datasets** (Waterbirds, MetaShift, CelebA, ISIC)
- **and feedback** (real users, automated w/LLM)

Can also outperform classic unsupervised group

robustness strategies (JTT, LfF)

## Future Work

Robustness is fundamentally sociotechnical

- Want model to reason how a human would reason.
- Interpretability as a tool for improving robustness

How can we use interpretability to make

progress on sociotechnical problems?