

# Language Guided Concept Bottlenecks for Interpretable and Robust Image Classification

Yue Yang

WPE-II Presentation

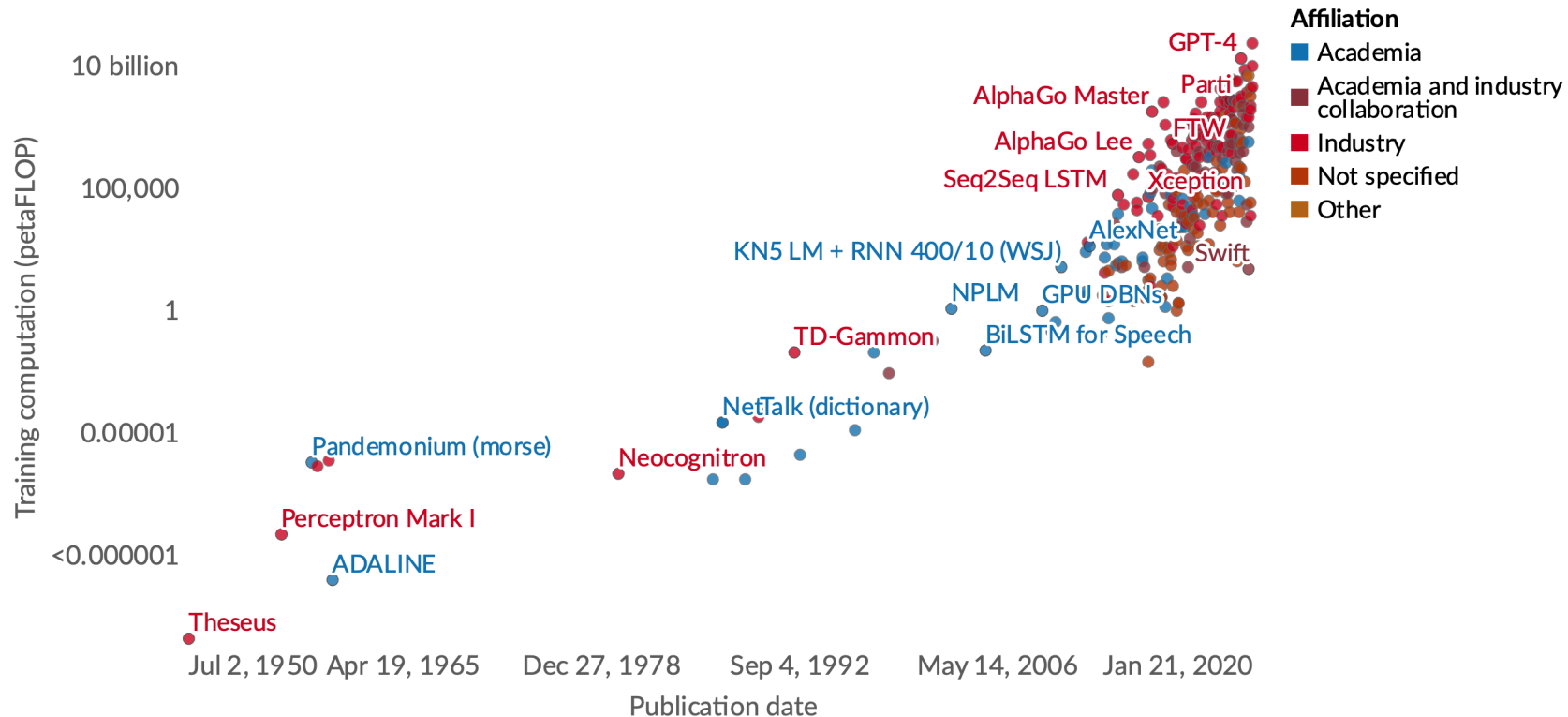
Committee: Dan Roth (Chair), Chris Callison-Burch, Mark Yatskar

# Models are getting performant but less interpretable.

## Computation used to train notable AI systems, by affiliation of researchers



Computation is measured in total petaFLOP, which is  $10^{15}$  floating-point operations estimated from AI literature, albeit with some uncertainty. Estimates are expected to be accurate within a factor of 2, or a factor of 5 for recent undisclosed models like GPT-4.



Data source: Epoch (2023)

[OurWorldInData.org/artificial-intelligence](https://OurWorldInData.org/artificial-intelligence) | CC BY

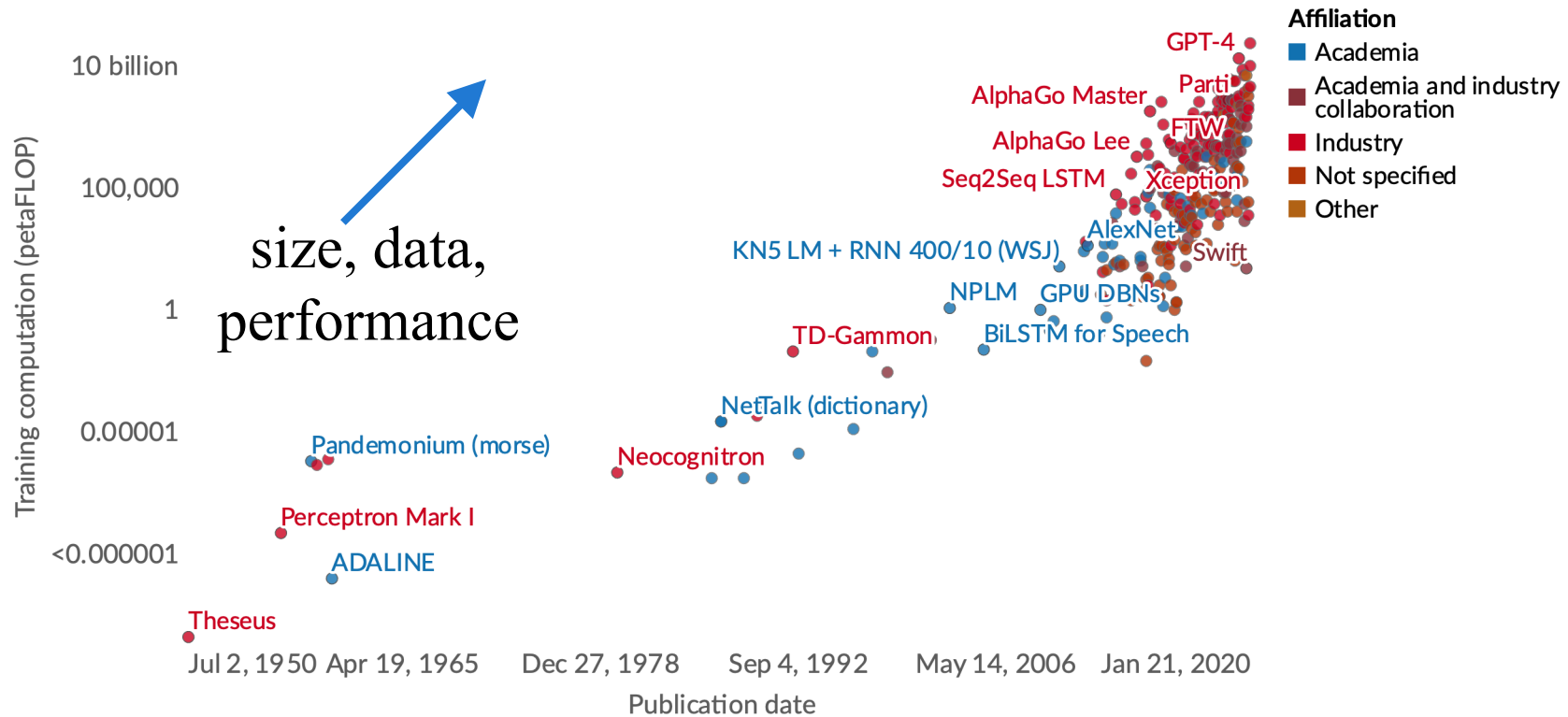
Note: The Executive Order on AI refers to a directive issued by President Biden on October 30, 2023, aimed at establishing guidelines and standards for the responsible development and use of artificial intelligence within the United States.

# Models are getting performant but less interpretable.

## Computation used to train notable AI systems, by affiliation of researchers



Computation is measured in total petaFLOP, which is  $10^{15}$  floating-point operations estimated from AI literature, albeit with some uncertainty. Estimates are expected to be accurate within a factor of 2, or a factor of 5 for recent undisclosed models like GPT-4.



Data source: Epoch (2023)

[OurWorldInData.org/artificial-intelligence](https://OurWorldInData.org/artificial-intelligence) | CC BY

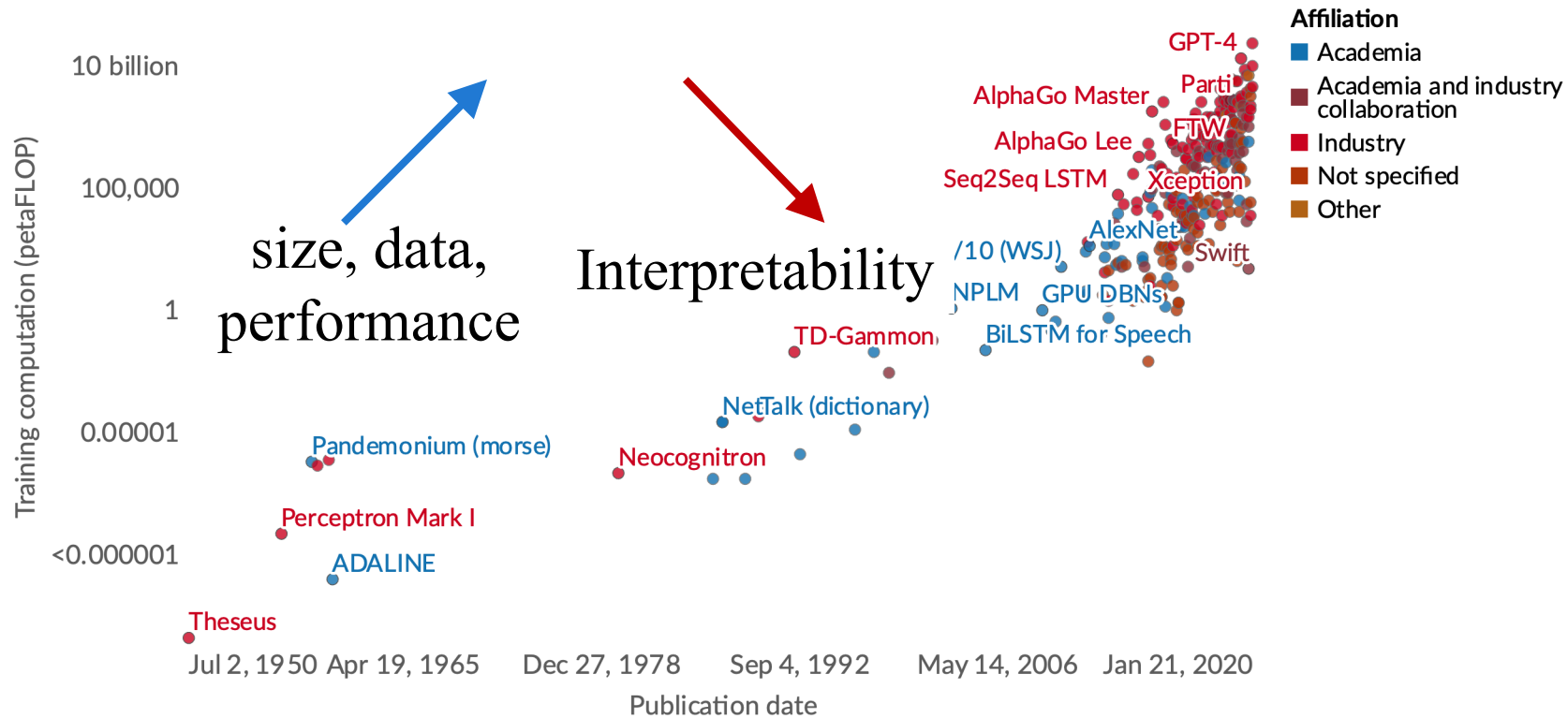
Note: The Executive Order on AI refers to a directive issued by President Biden on October 30, 2023, aimed at establishing guidelines and standards for the responsible development and use of artificial intelligence within the United States.

# Models are getting performant but less interpretable.

## Computation used to train notable AI systems, by affiliation of researchers



Computation is measured in total petaFLOP, which is  $10^{15}$  floating-point operations estimated from AI literature, albeit with some uncertainty. Estimates are expected to be accurate within a factor of 2, or a factor of 5 for recent undisclosed models like GPT-4.

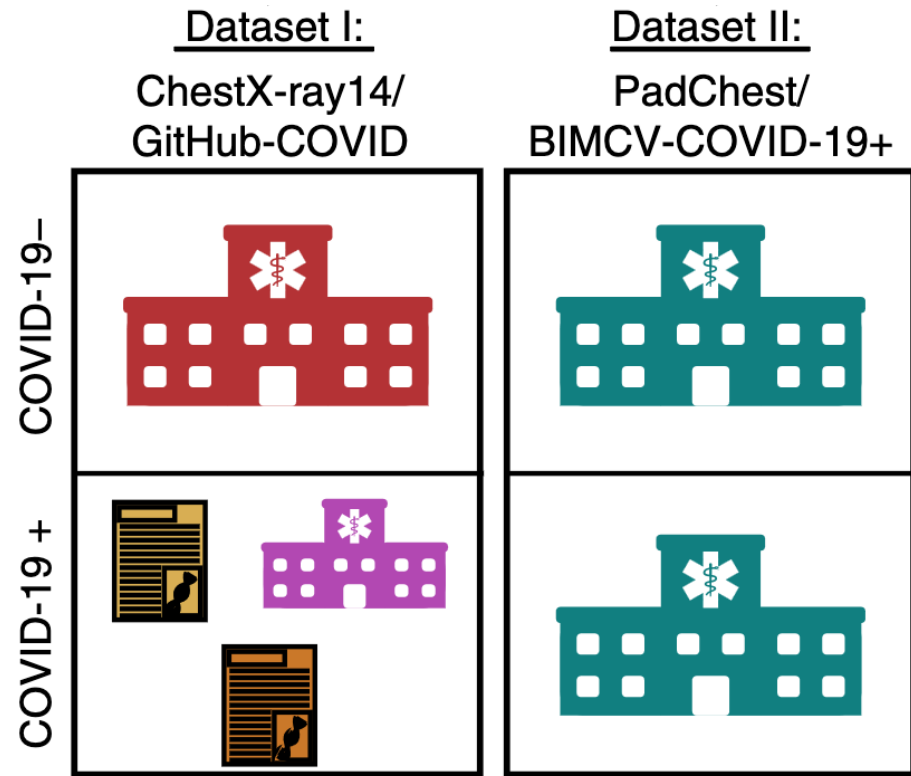


Data source: Epoch (2023)

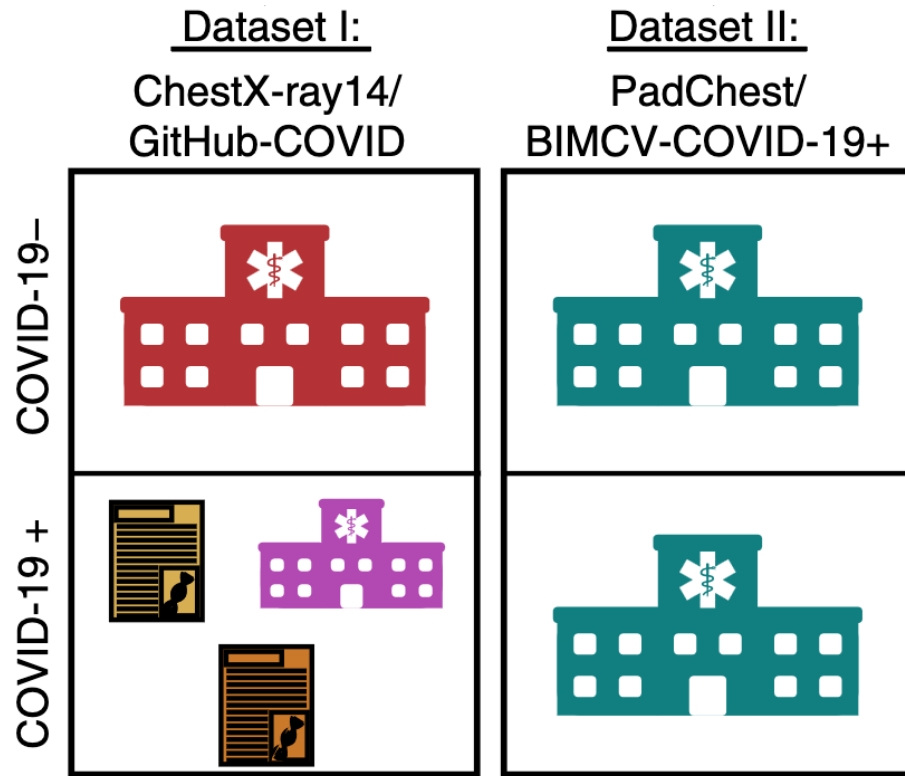
[OurWorldInData.org/artificial-intelligence](https://OurWorldInData.org/artificial-intelligence) | CC BY

Note: The Executive Order on AI refers to a directive issued by President Biden on October 30, 2023, aimed at establishing guidelines and standards for the responsible development and use of artificial intelligence within the United States.

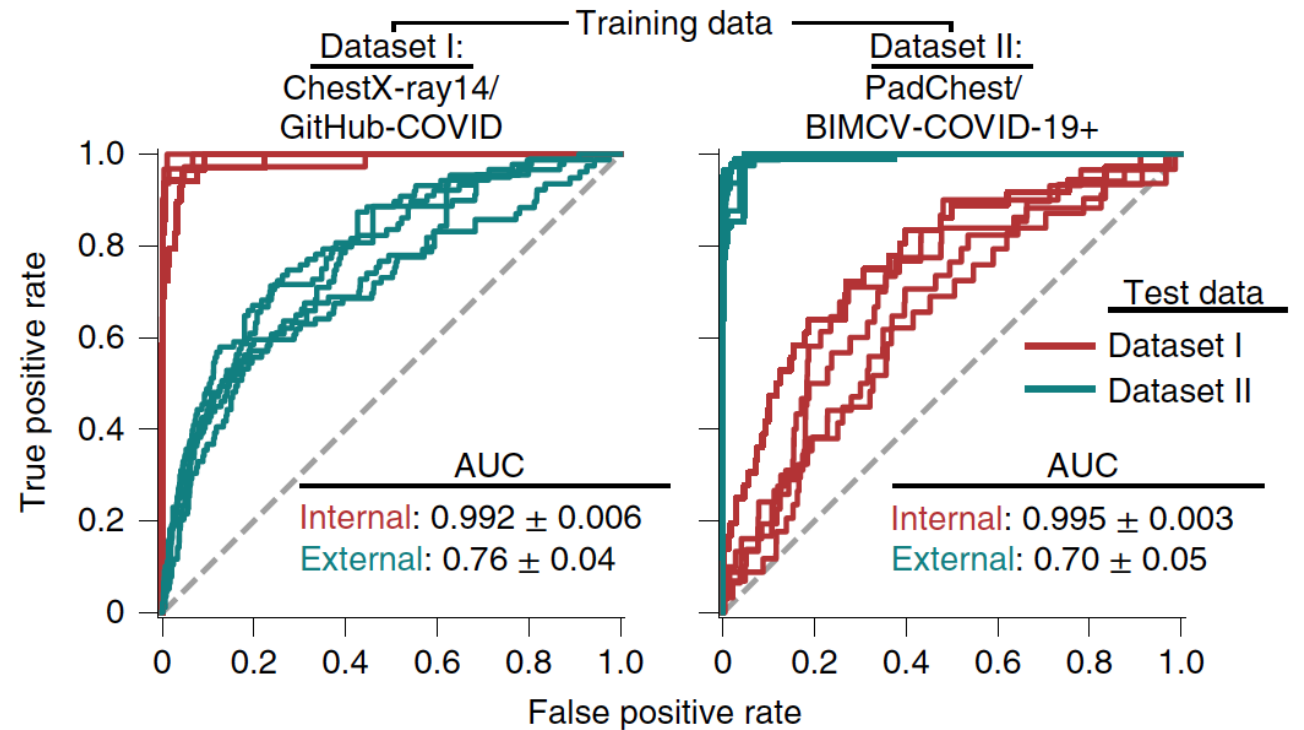
# Catastrophic failures in critical domains.



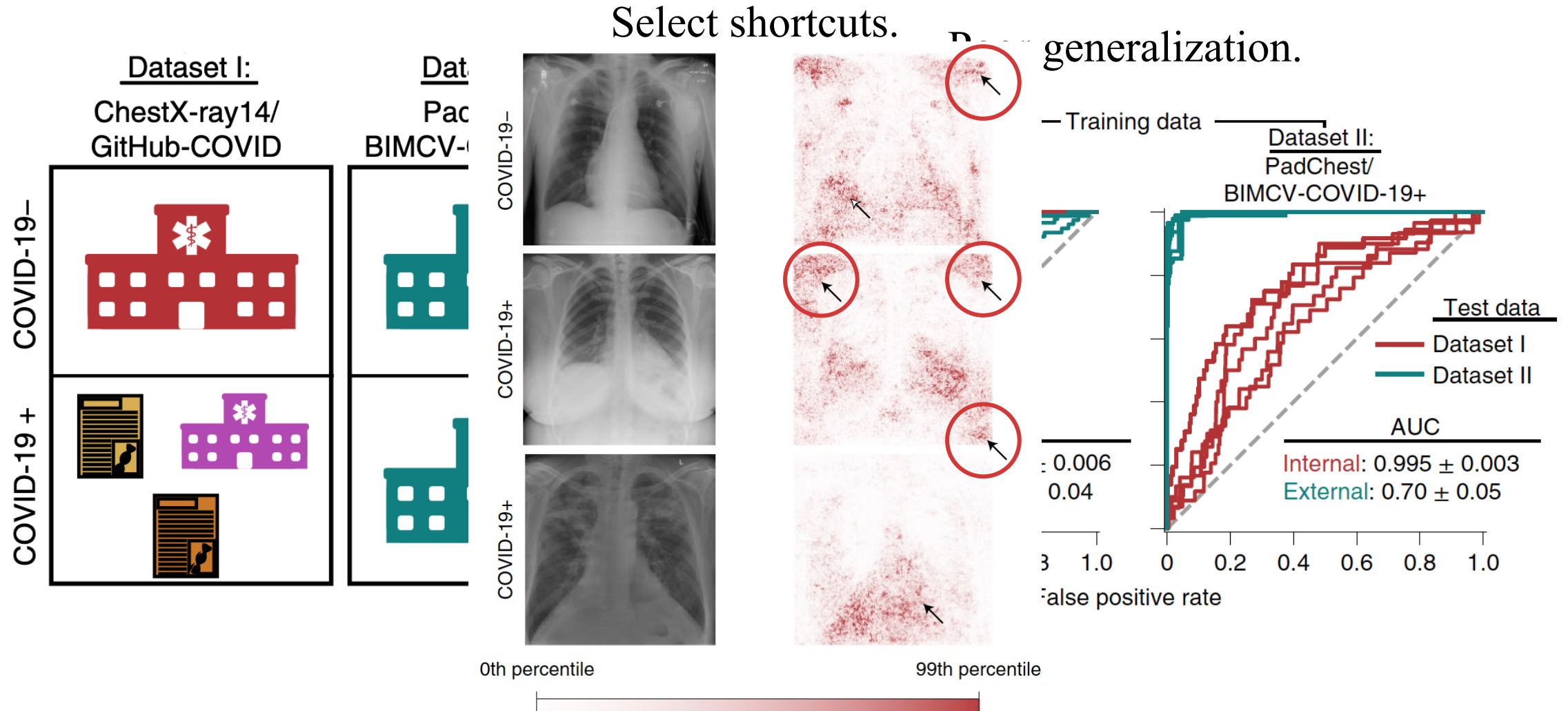
# Catastrophic failures in critical domains.



Poor generalization.

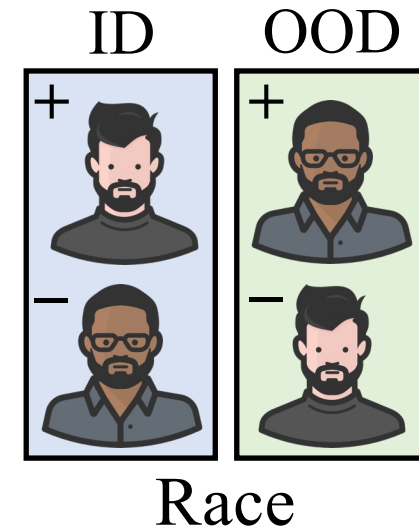
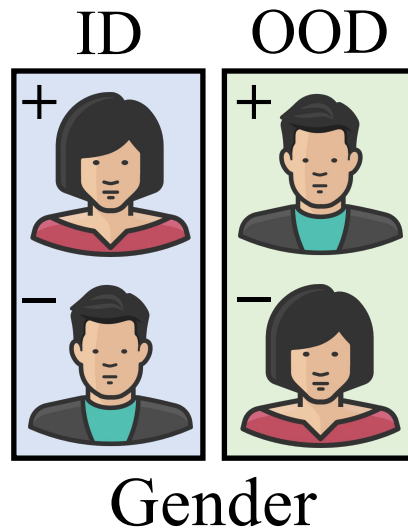


# Catastrophic failures in critical domains.



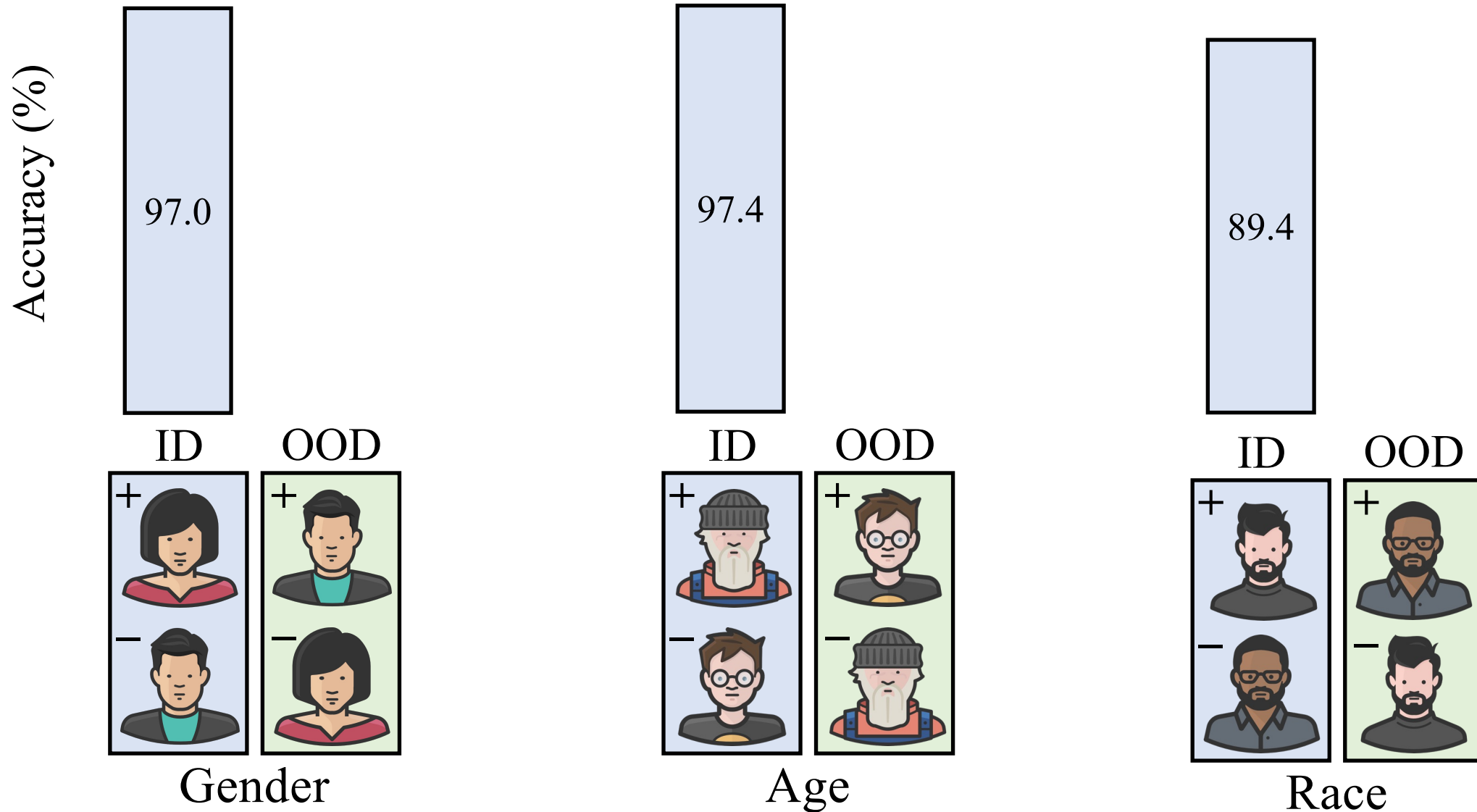
# Black-box models generalize poorly on domain shifts.

Accuracy (%)

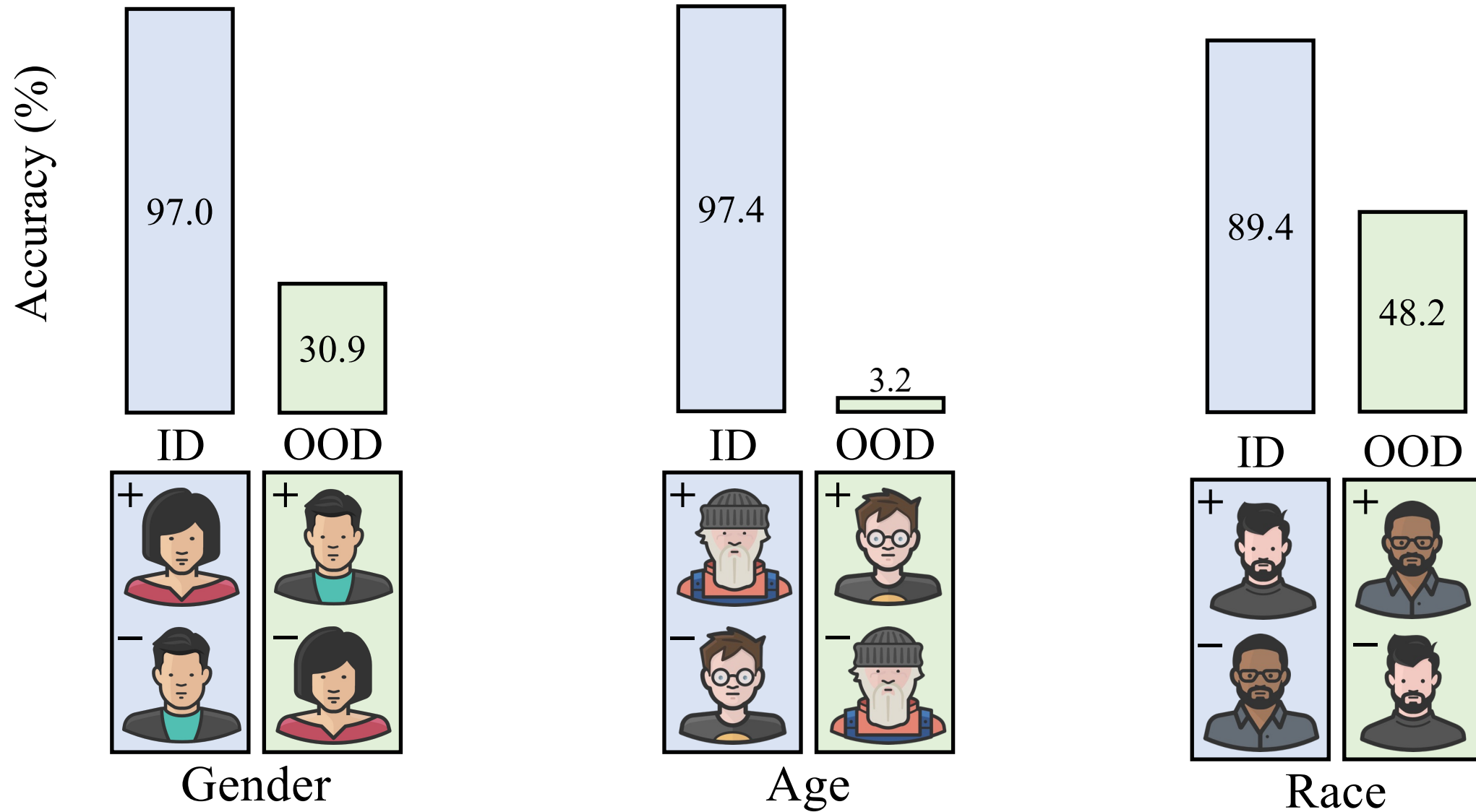




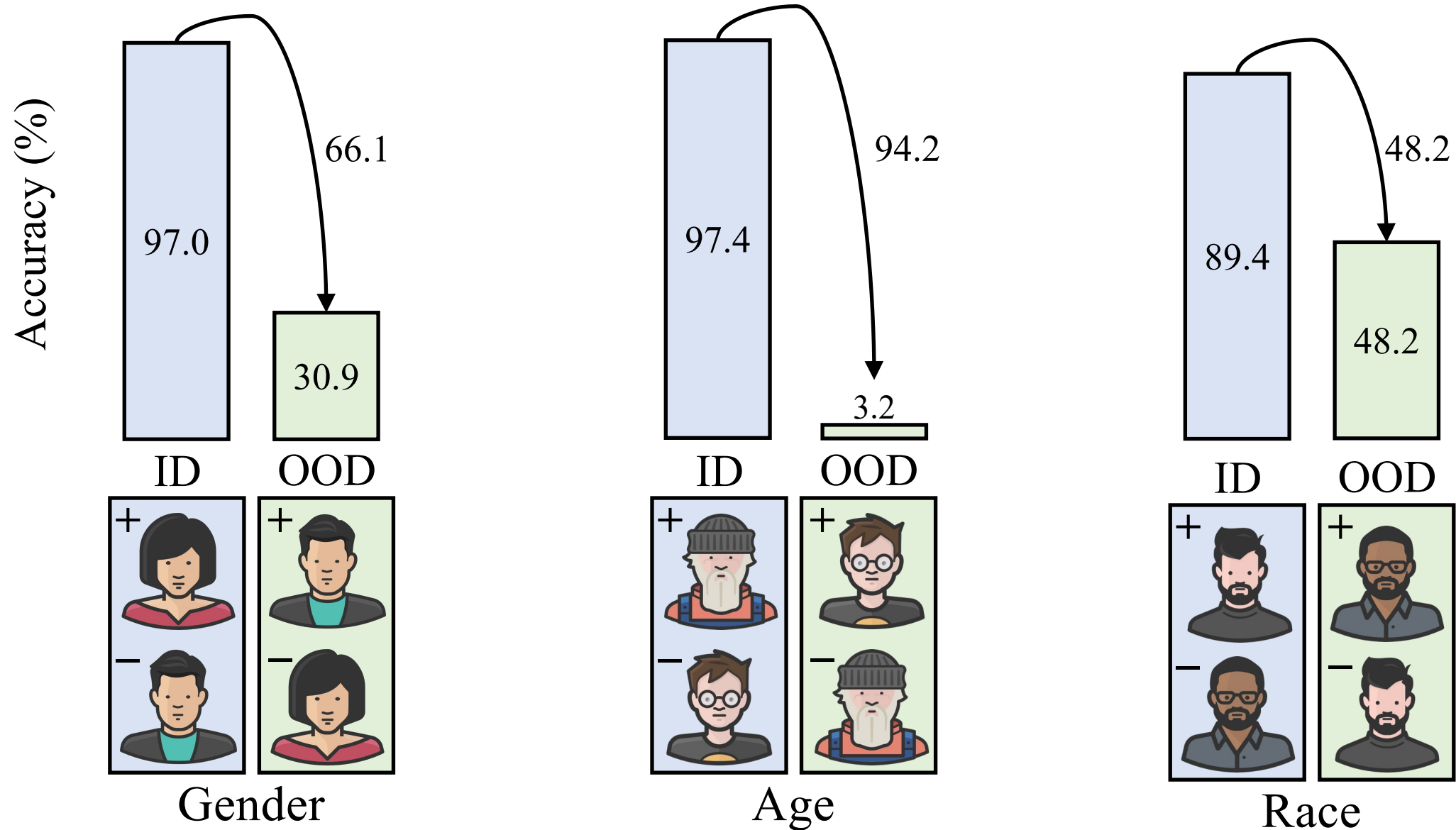
# Black-box models generalize poorly on domain shifts.



# Black-box models generalize poorly on domain shifts.




# Black-box models generalize poorly on domain shifts.




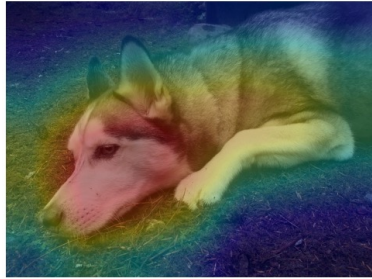
## Solution-1: Post-hoc Explanation

- Explain the black box model with another black box model.
- Explanations are often **not faithful** and can be misleading.

	Test Image
Explanations Using Attention Maps	

## Solution-1: Post-hoc Explanation


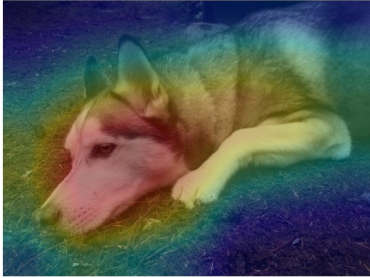

- Explain the black box model with another black box model.
- Explanations are often **not faithful** and can be misleading.

	Test Image	Evidence for Animal Being a Siberian Husky
Explanations Using Attention Maps		

Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. Nature Machine Intelligence. 2019.

## Solution-1: Post-hoc Explanation

- Explain the black box model with another black box model.
- Explanations are often **not faithful** and can be misleading.

	Test Image	Evidence for Animal Being a Siberian Husky	Evidence for Animal Being a Transverse Flute
Explanations Using Attention Maps			

Cynthia Rudin. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. Nature Machine Intelligence. 2019.

## Solution-2: Inherently Interpretable Methods

Provide their **own explanations** that are **faithful** to the predictions.

**Table 3 | Scoring system for risk of recidivism**

1.	Prior arrests $\geq 2$	1 point	...			
2.	Prior arrests $\geq 5$	1 point	+...			
3.	Prior arrests for local ordinance	1 point	+...			
4.	Age at release between 18 to 24	1 point	+...			
5.	Age at release $\geq 40$	-1 point	+...			
		Score	= ...			
Score	-1	0	1	2	3	4
Risk (%)	11.9	26.9	50.0	73.1	88.1	95.3

This system is from ref. <sup>21</sup>, which was developed from refs. <sup>29,46</sup>. The model was not created by a human; the selection of numbers and features come from the RiskSLIM machine learning algorithm.

Rudin and Ustun. Optimized Scoring Systems: Toward Trust in Machine Learning for Healthcare and Criminal Justice. INFORMS. 2018

# Concept Bottleneck Models (CBMs)

**Input Image  $x$**



**End-to-end Model**

→ **label  $y$**  (black-throated sparrow)



# Concept Bottleneck Models (CBMs)

**Input Image  $x$**



**Black Box**

→ **label  $y$**  (black-throated sparrow)

# Concept Bottleneck Models (CBMs)

**Input Image  $x$**



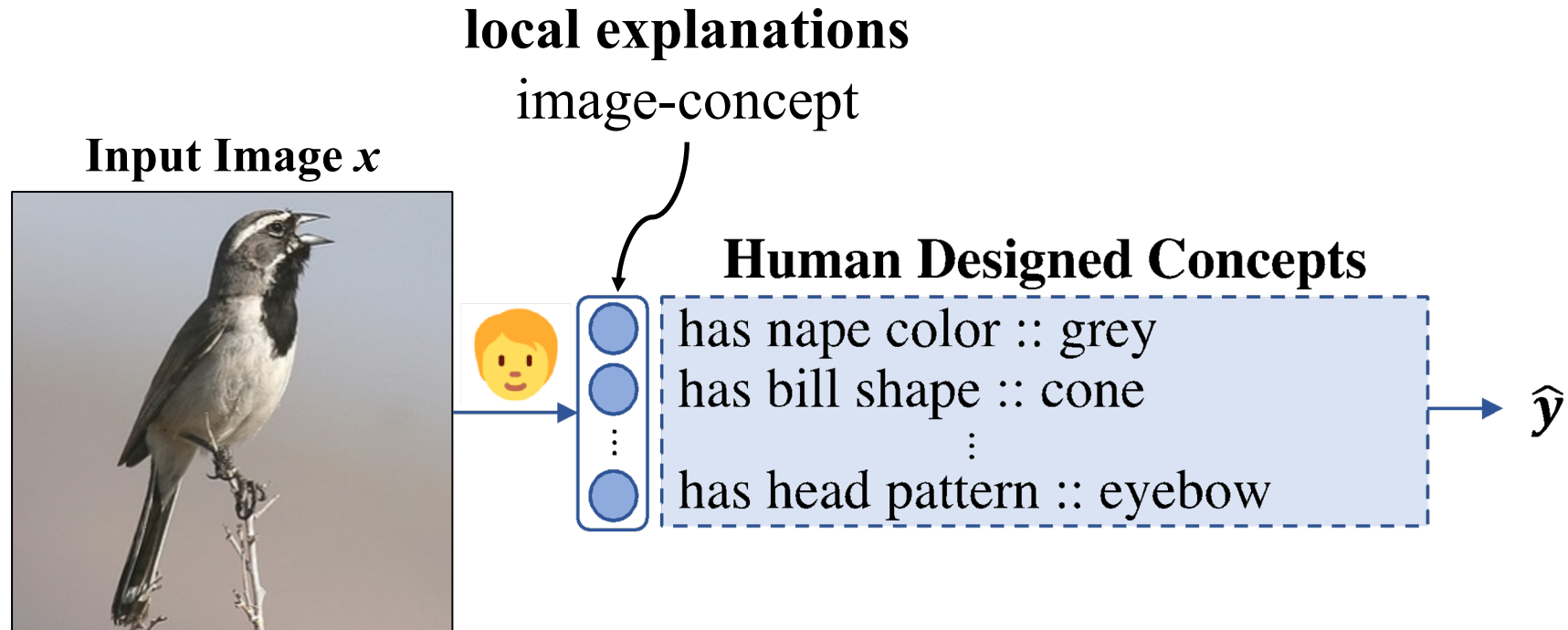
**Human Designed Concepts**



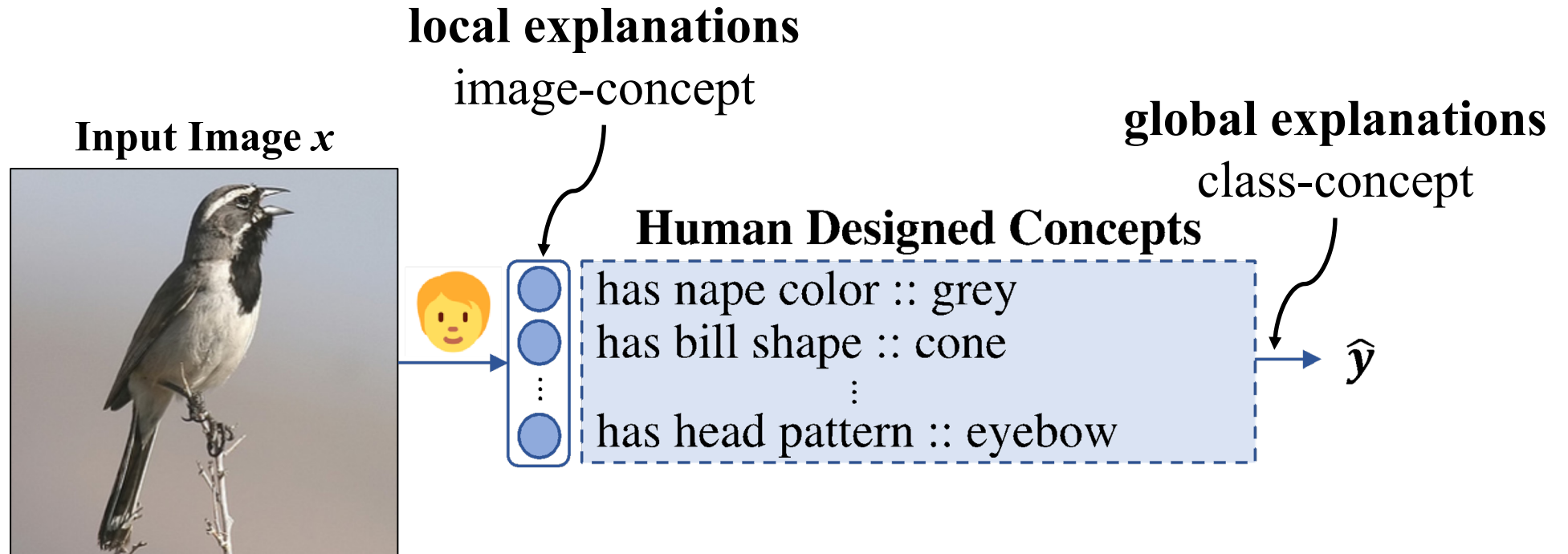
has nape color :: grey  
has bill shape :: cone  
⋮  
has head pattern :: eyebrow

$\hat{y}$

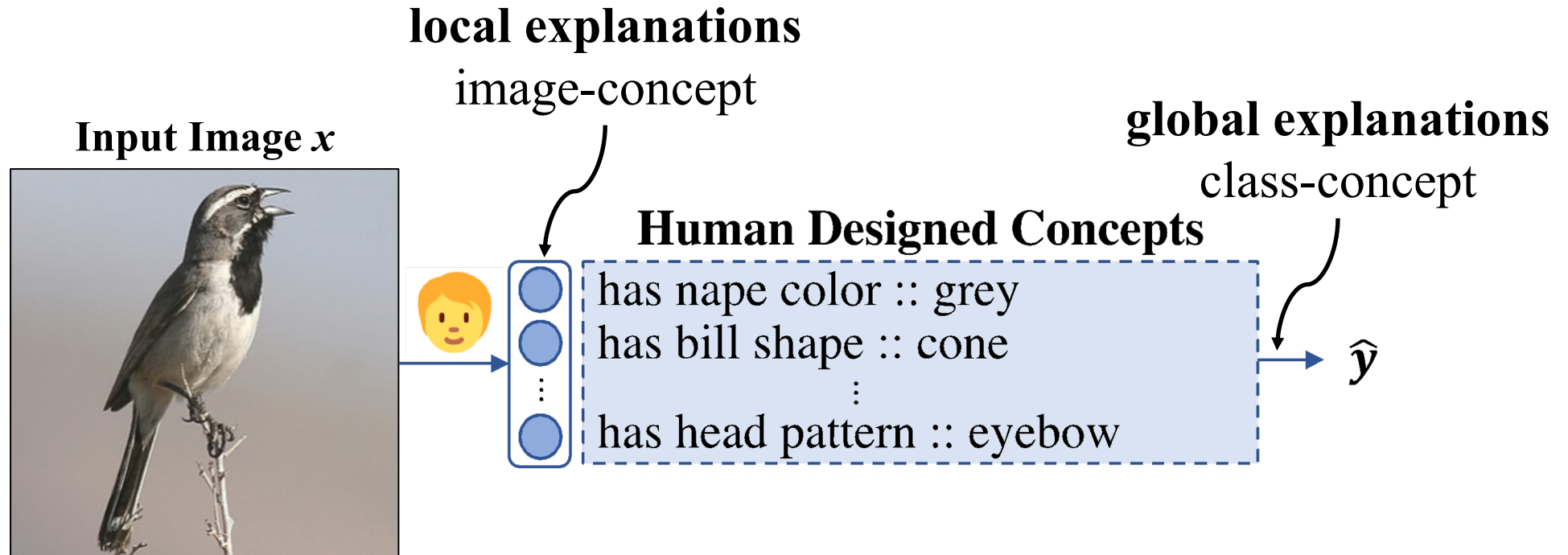
# Concept Bottleneck Models (CBMs)



# Concept Bottleneck Models (CBMs)



# Concept Bottleneck Models (CBMs)



## Challenges:

- **Scale:** requires human efforts in building concept bottlenecks.
- **Performance:** perform worse than black-box models.

# Concept Bottleneck Models (CBMs)

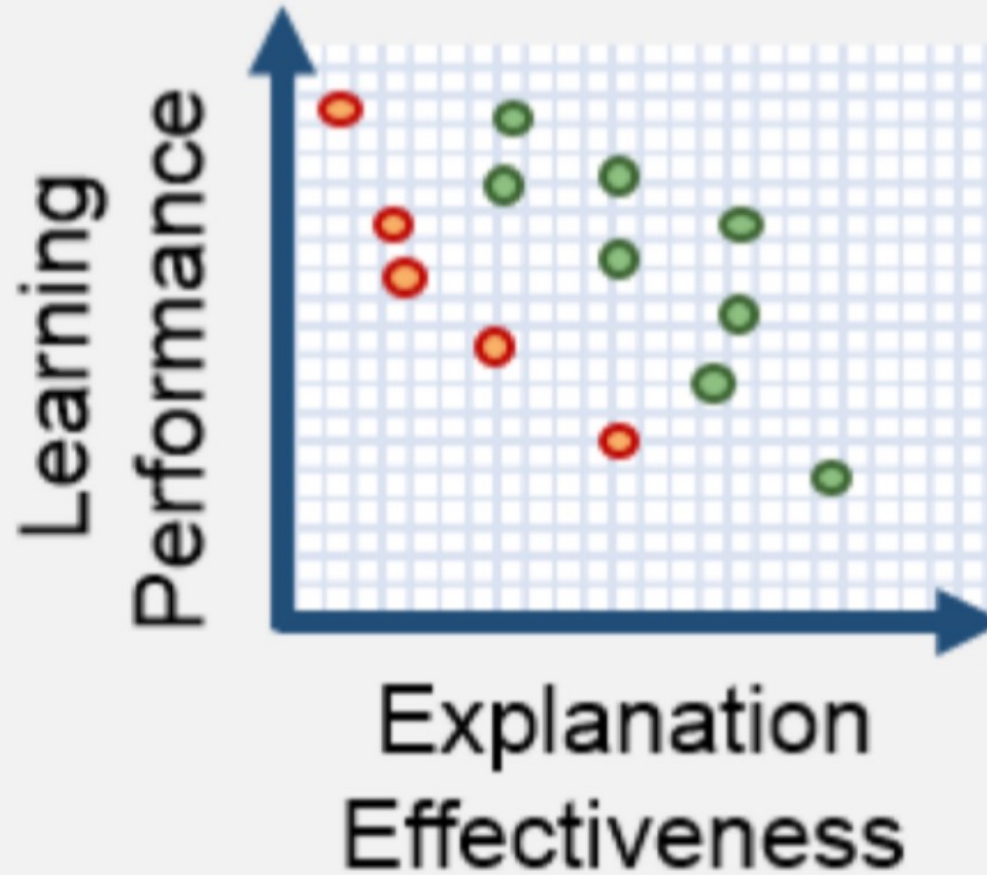
Input Image



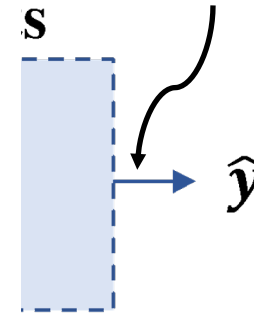
## Challenge

- **Scale:** re
- **Perform**

DARPA XAI BAA, 2019.

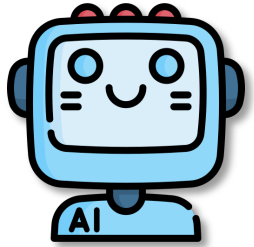


Global explanations  
class-concept

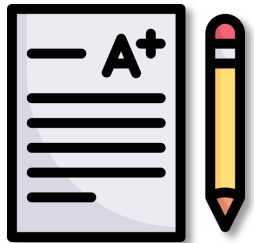


bottlenecks.  
els.

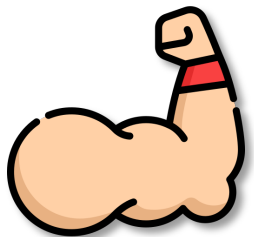
# Agenda of this talk



How to construct concept bottlenecks  
**without human effort?**



How to make interpretable models  
**performant** as black-box models?



What other advantages can interpretable  
models give us? Answer: **Robustness.**



LaBo (CVPR 23)



KnoBo  
(In progress)



# Language in a **B**ottle: Language Model Guided Concept Bottlenecks for Interpretable Image Classification

**Yue Yang**, Artemis Panagopoulou, Shenghao Zhou, Daniel Jin,  
Chris Callison-Burch, Mark Yatskar

University of Pennsylvania





# Leverage the world knowledge of LLMs



# Leverage the world knowledge of LLMs



# Leverage the world knowledge of LLMs



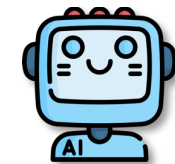
Describe what the *black-throated sparrow* looks like.



# Leverage the world knowledge of LLMs



Describe what the *black-throated sparrow* looks like.



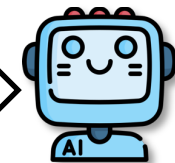
GPT-3

# Leverage the world knowledge of LLMs



Describe what the *black-throated sparrow* looks like.

The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It has **two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is found in North America. It is a common bird in the western United States. The black-throated sparrow is a member of the sparrow family.



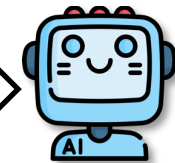
GPT-3

# Leverage the world knowledge of LLMs



Describe what the *black-throated sparrow* looks like.

The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It has **two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is found in North America. It is a common bird in the western United States. The black-throated sparrow is a member of the sparrow family.



GPT-3

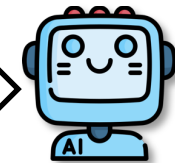


# Leverage the world knowledge of LLMs



Describe what the *black-throated sparrow* looks like.

The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It has **two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is found in North America. It is a common bird in the western United States. The black-throated sparrow is a member of the sparrow family.



GPT-3



Describe the **shape/color** of the *black-throated sparrow*.

# Select the knowledge



The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It **has two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is found in North America. It is a common bird in the western United States. The black-throated sparrow is a member of the sparrow family.



# Select the knowledge

too general

The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It **has two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is found in North America. It is a common bird in the western United States. The black-throated sparrow is a member of the sparrow family.



# Select the knowledge



too general

The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It **has two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is found in North America. It is a common bird in the western United States. The black-throated sparrow is a member of the sparrow family.

not visual

# Select the knowledge

too general

The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It **has two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is found in North America. It is a common bird in the western United States. The black-throated sparrow is a member of the sparrow family.

not visual



## Submodular Optimization

Visual  
Discriminative  
Diverse

# Select the knowledge



too general

The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It **has two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is found in North America. It is a common bird in the western United States. The black-throated sparrow is a member of the sparrow family.

not visual

Submodular  
Optimization

Visual  
Discriminative  
Diverse

# Select the knowledge



too general

The black-throated sparrow is a **small bird** with a **black head and throat**. It has a **white body with brown streaks on its back**. Its **wings are brown with white stripes**. The black-throated sparrow has a **long, thin beak**. It has **two long, thin legs**. The black-throated sparrow has a **long, thin tail**. It is about **5 inches long**. The black-throated sparrow is **found in North America**. It is a **common bird in the western United States**. The black-throated sparrow is **a member of the sparrow family**.

not visual

Submodular  
Optimization

Visual  
Discriminative  
Diverse

**black head and throat**

**long, thin tail**

**wings are brown  
with white stripes**

**White body with  
brown streaks**

# Ground Concepts using CLIP

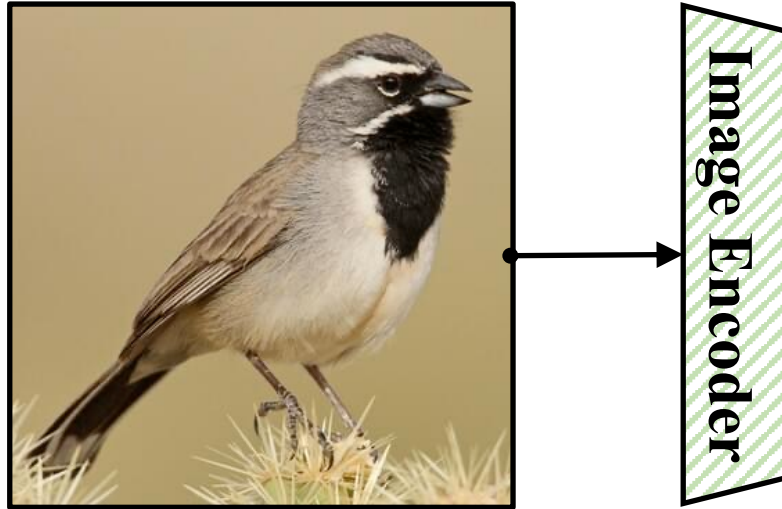


**black head and throat**

**long, thin tail**

**wings are brown  
with white stripes**

# Ground Concepts using CLIP

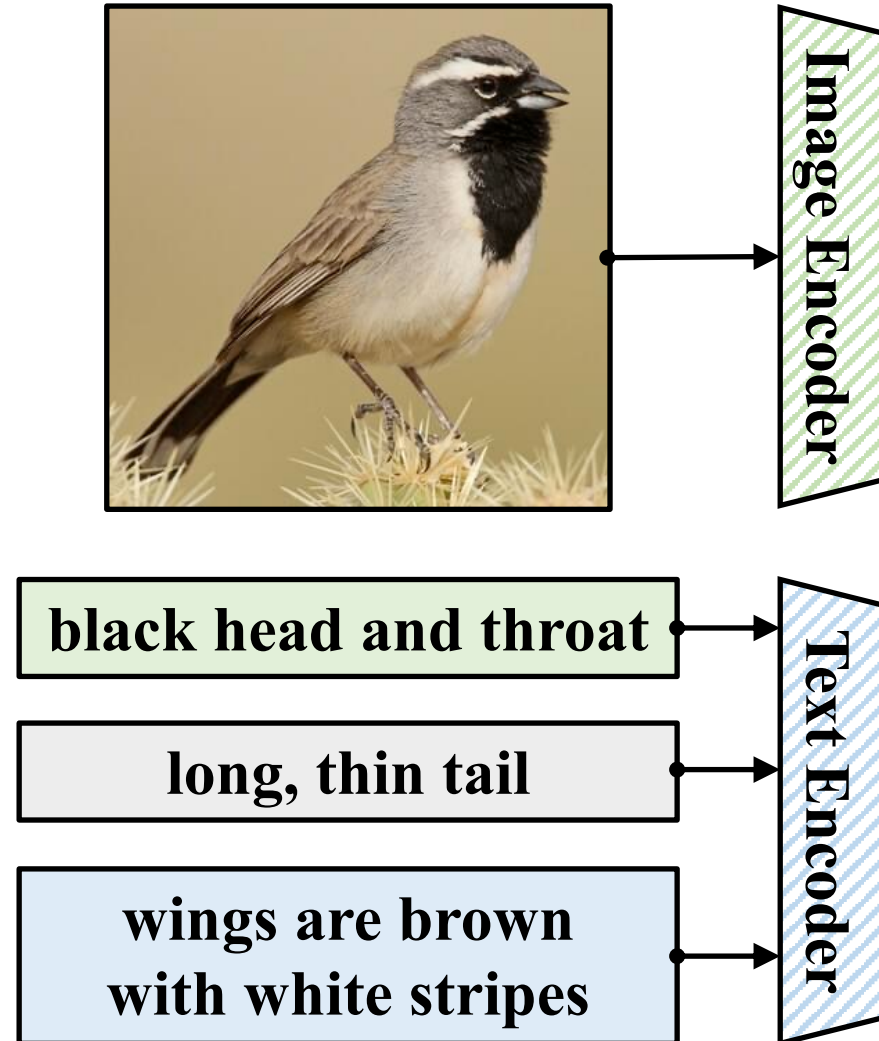


**black head and throat**

**long, thin tail**

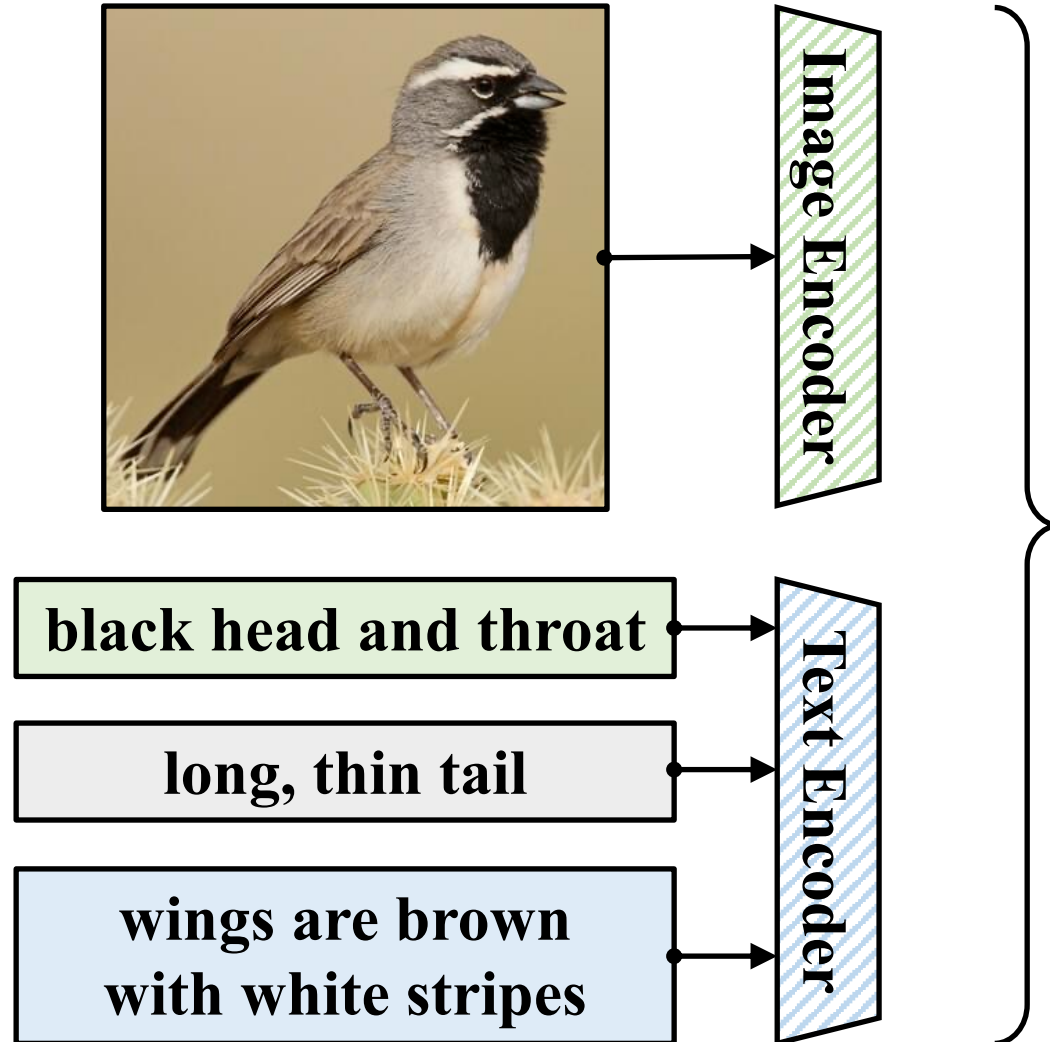
**wings are brown  
with white stripes**

# Ground Concepts using CLIP

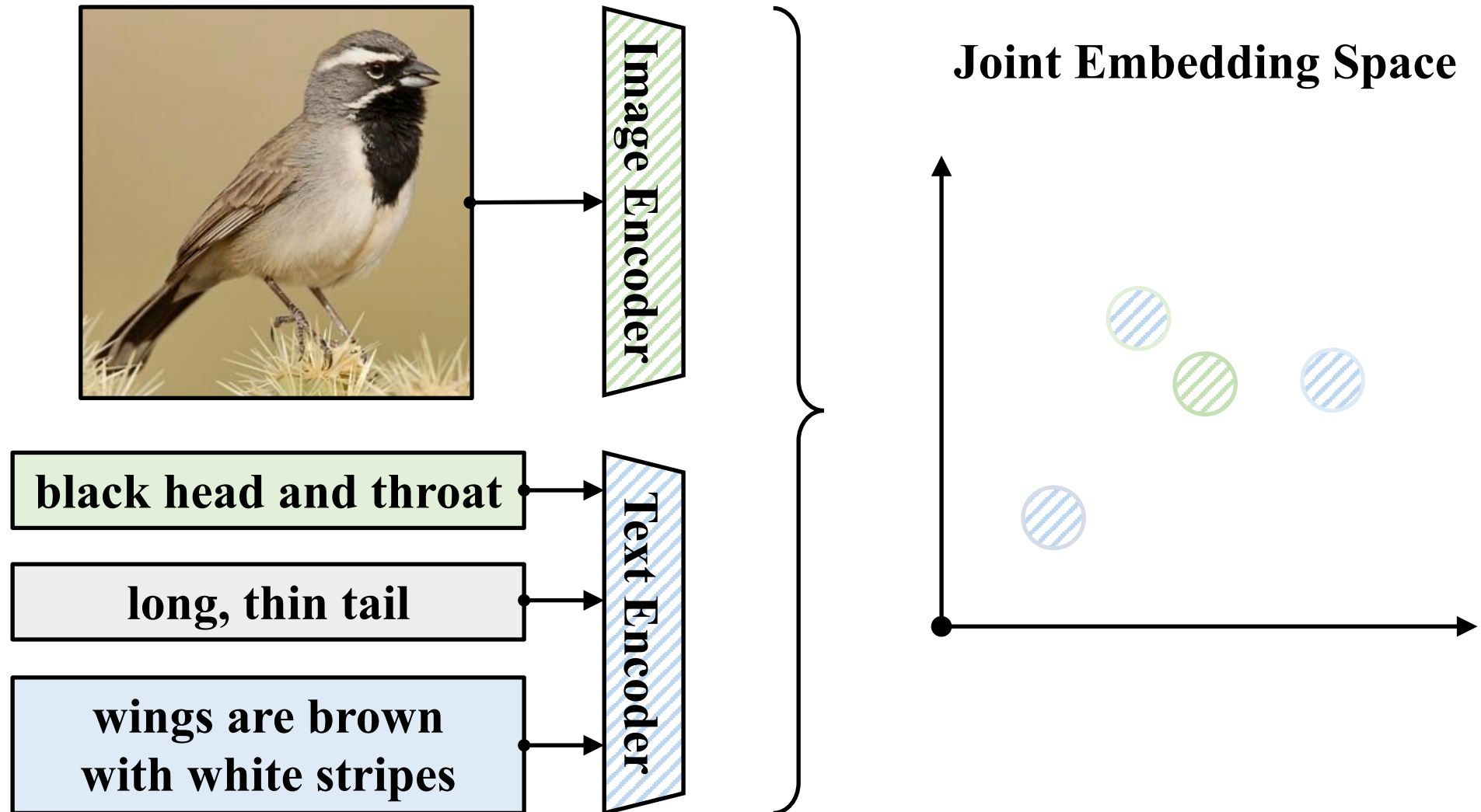




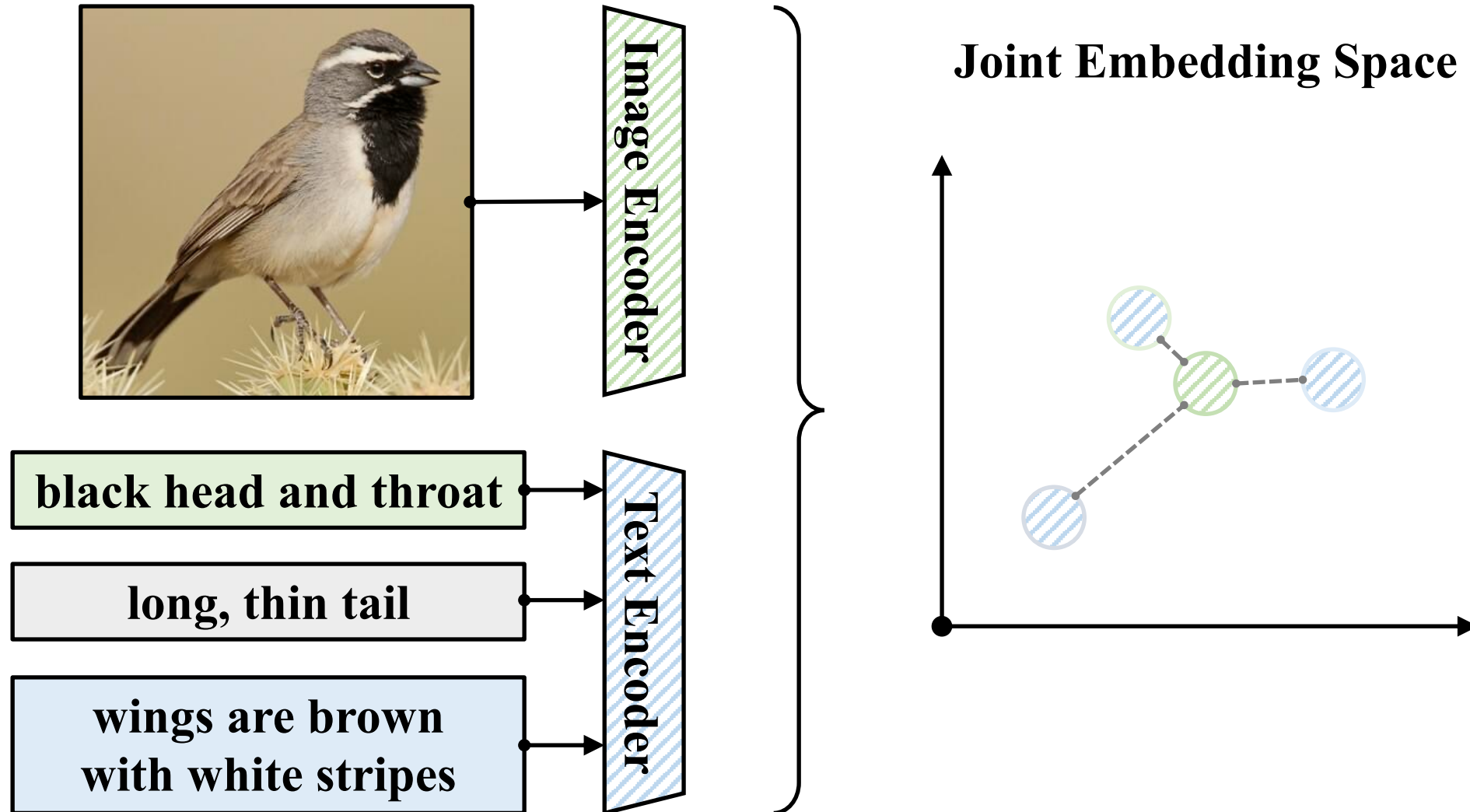
# Ground Concepts using CLIP



# Ground Concepts using CLIP



# Ground Concepts using CLIP



# Prompt LLM to generate candidate concepts

class 1-axolotl



class 2-red panda



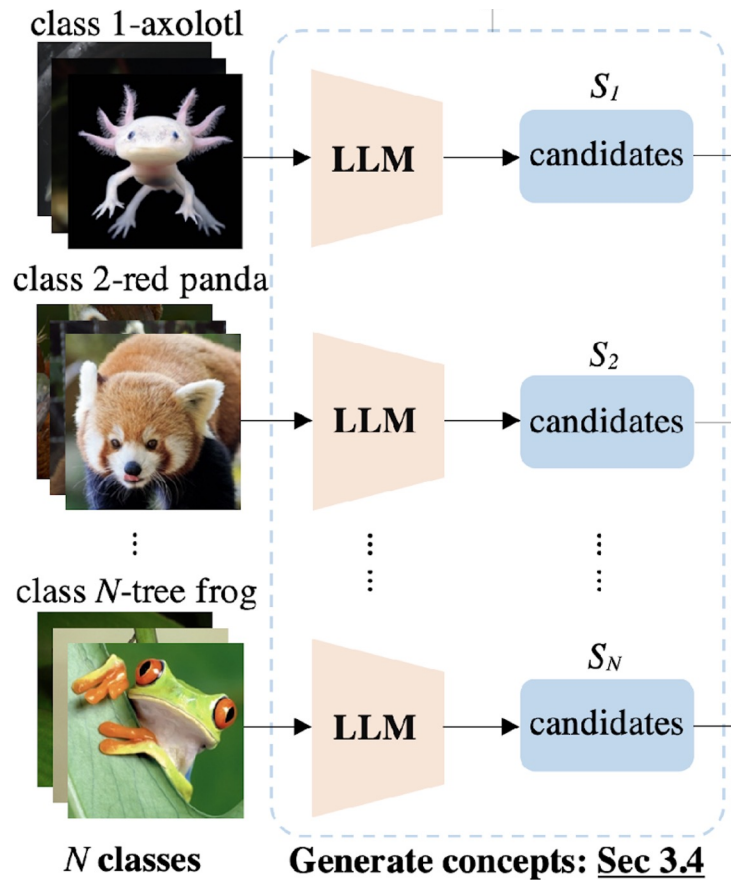
⋮

class  $N$ -tree frog



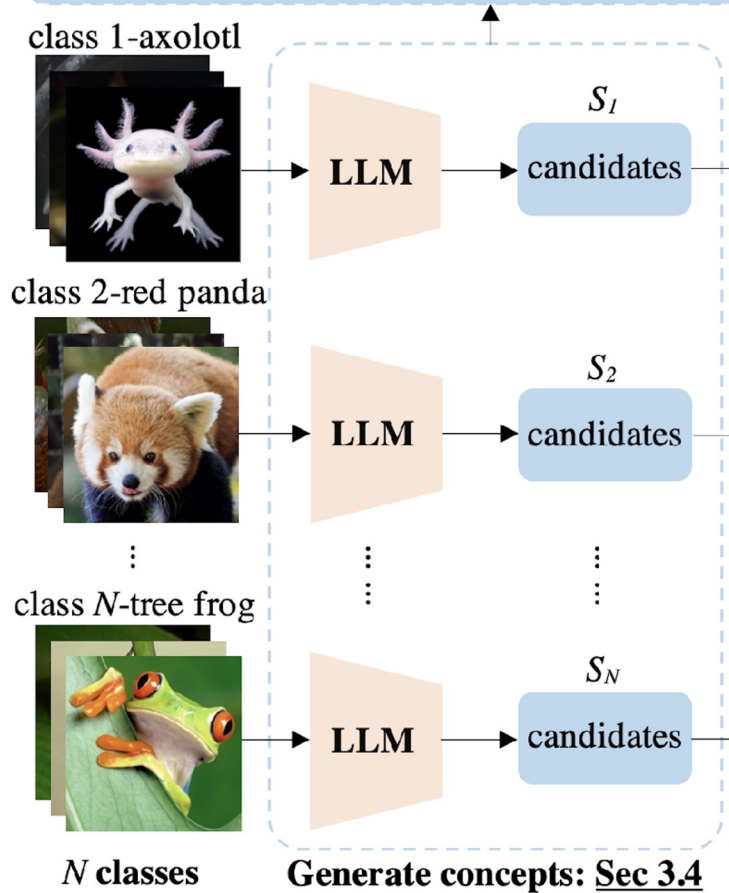
$N$  classes

# Prompt LLM to generate candidate concepts



# Prompt LLM to generate candidate concepts

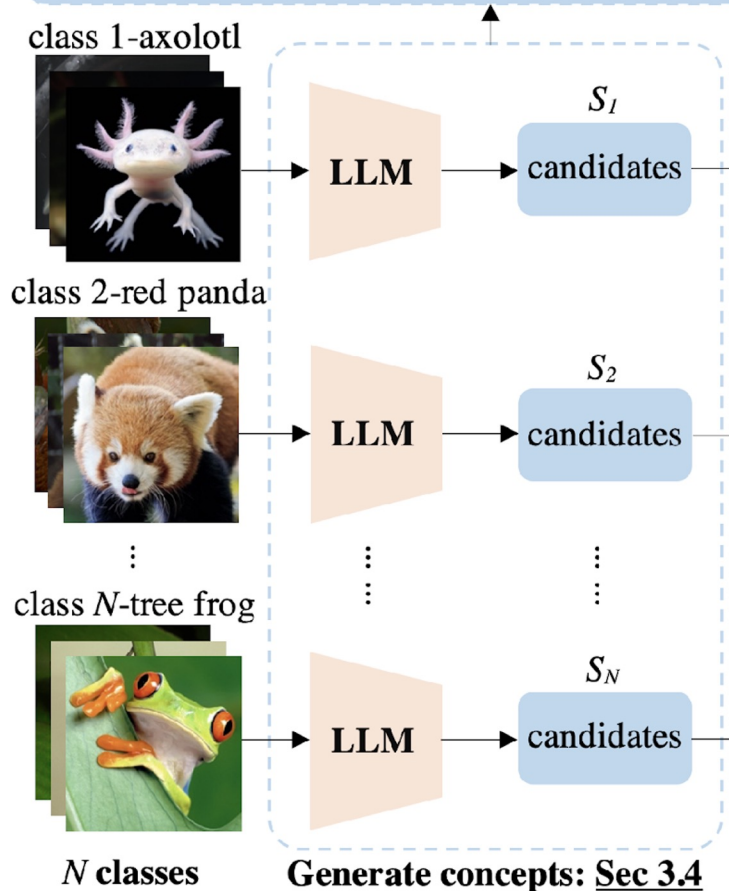
**prompt:** describe what the *axolotl* looks like:  
**LLM:** The axolotl's limbs are delicate, and the tail is long and thin.  
**Extract concept using LM and delete class names:**  
**Candidate concepts:** limbs are delicate; tail is long and thin



Generate concepts: [Sec 3.4](#)

# Prompt LLM to generate candidate concepts

**prompt:** describe what the *axolotl* looks like:  
**LLM:** The axolotl's limbs are delicate, and the tail is long and thin.  
**Extract concept using LM and delete class names:**  
**Candidate concepts:** limbs are delicate; tail is long and thin



## General Prompt Template

1. describe what the [CLASS NAME] looks like:
2. describe the appearance of the [CLASS NAME]:
3. describe the color of the [CLASS NAME]:
4. describe the pattern of the [CLASS NAME]:
5. describe the shape of the [CLASS NAME]:

- Obtain 500 sentences for each class.
- Extract concepts from sentences using T5 [1].
- String match to identify and remove class name tokens in each concept.

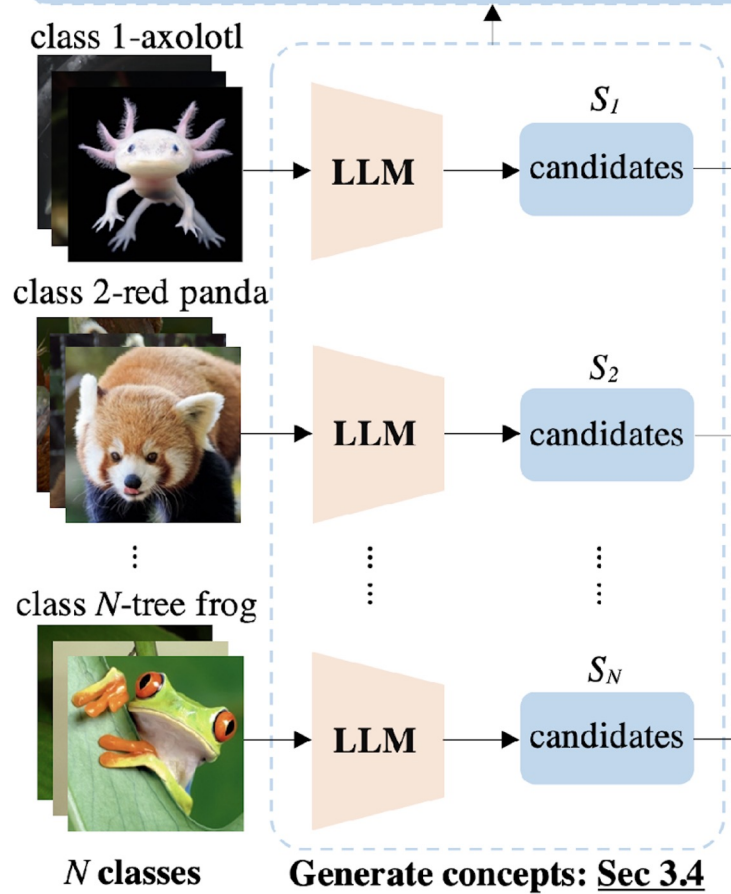
# Submodular Concept Selection

**prompt:** describe what the *axolotl* looks like:

**LLM:** The axolotl's limbs are delicate, and the tail is long and thin.

**Extract concept using LM and delete class names:**

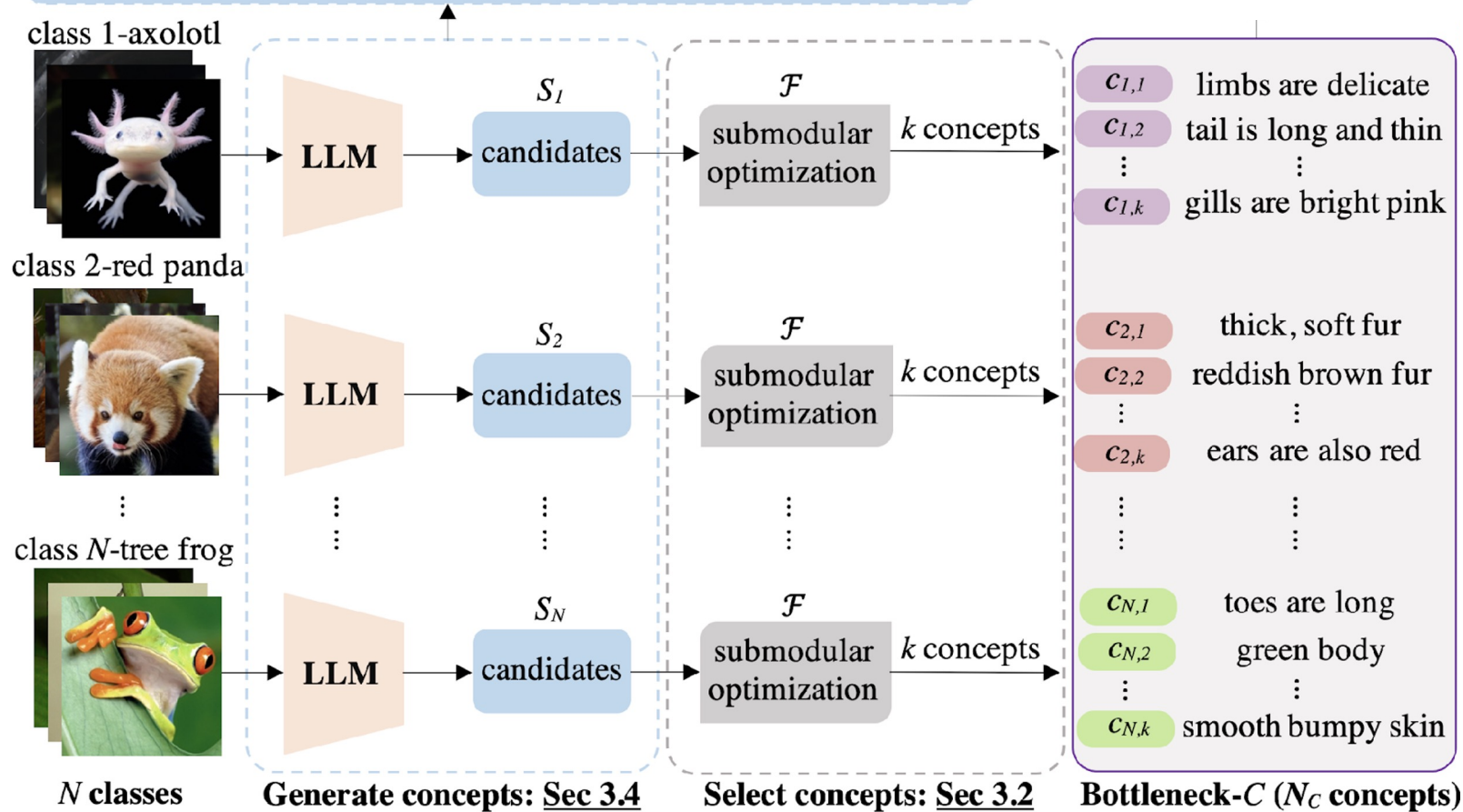
**Candidate concepts:** limbs are delicate; tail is long and thin





# Submodular Concept Selection

**prompt:** describe what the *axolotl* looks like:  
**LLM:** The axolotl's limbs are delicate, and the tail is long and thin.  
**Extract concept using LM and delete class names:**  
**Candidate concepts:** limbs are delicate; tail is long and thin



# Submodular Concept Selection

- Given a superset of concepts  $S_y$  for a class  $y$ .
- Select a subset  $C_y$  for the bottleneck which are **discriminative** and **diverse**.

● candidate  $S_y$   
● selected  $C_y$

# Submodular Concept Selection

- Given a superset of concepts  $S_y$  for a class  $y$ .
- Select a subset  $C_y$  for the bottleneck which are **discriminative** and **diverse**.

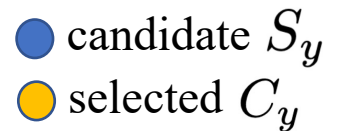
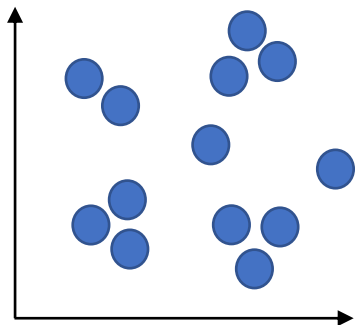
$$\mathcal{F}(C_y) = \underbrace{\alpha \cdot \sum_{c \in C_y} D(c)}_{\text{discriminability}} + \underbrace{\beta \cdot \sum_{c_1 \in S_y} \max_{c_2 \in C_y} \phi(c_1, c_2)}_{\text{coverage}}$$

- candidate  $S_y$
- selected  $C_y$

# Submodular Concept Selection

- Given a superset of concepts  $S_y$  for a class  $y$ .
- Select a subset  $C_y$  for the bottleneck which are **discriminative** and **diverse**.

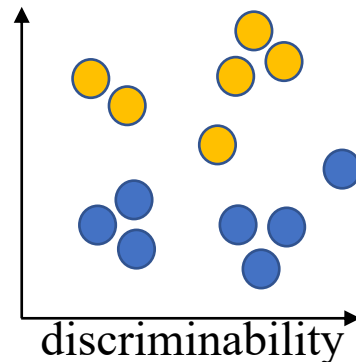
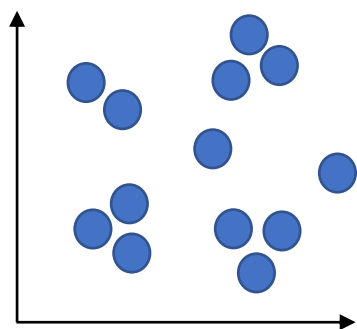
$$\mathcal{F}(C_y) = \underbrace{\alpha \cdot \sum_{c \in C_y} D(c)}_{\text{discriminability}} + \underbrace{\beta \cdot \sum_{c_1 \in S_y} \max_{c_2 \in C_y} \phi(c_1, c_2)}_{\text{coverage}}$$



# Submodular Concept Selection

- Given a superset of concepts  $S_y$  for a class  $y$ .
- Select a subset  $C_y$  for the bottleneck which are **discriminative** and **diverse**.

$$\mathcal{F}(C_y) = \underbrace{\alpha \cdot \sum_{c \in C_y} D(c)}_{\text{discriminability}} + \underbrace{\beta \cdot \sum_{c_1 \in S_y} \max_{c_2 \in C_y} \phi(c_1, c_2)}_{\text{coverage}}$$

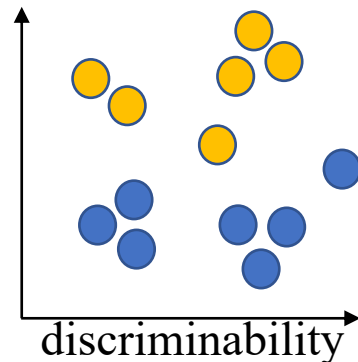
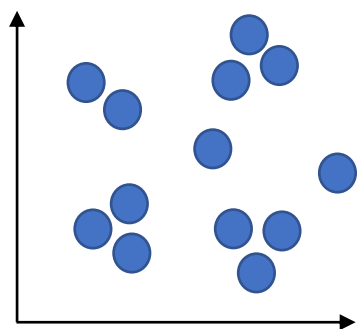


● candidate  $S_y$   
● selected  $C_y$

# Submodular Concept Selection

- Given a superset of concepts  $S_y$  for a class  $y$ .
- Select a subset  $C_y$  for the bottleneck which are **discriminative** and **diverse**.

$$\mathcal{F}(C_y) = \underbrace{\alpha \cdot \sum_{c \in C_y} D(c)}_{\text{discriminability}} + \underbrace{\beta \cdot \sum_{c_1 \in S_y} \max_{c_2 \in C_y} \phi(c_1, c_2)}_{\text{coverage}}$$



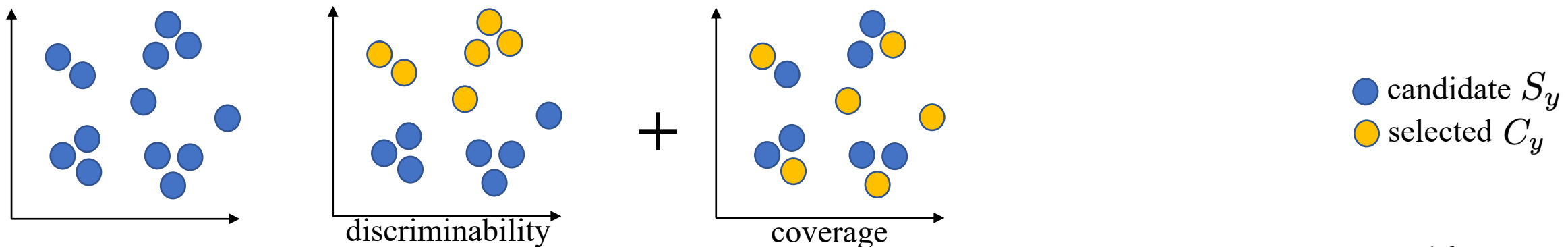
+

● candidate  $S_y$   
● selected  $C_y$

# Submodular Concept Selection

- Given a superset of concepts  $S_y$  for a class  $y$ .
- Select a subset  $C_y$  for the bottleneck which are **discriminative** and **diverse**.

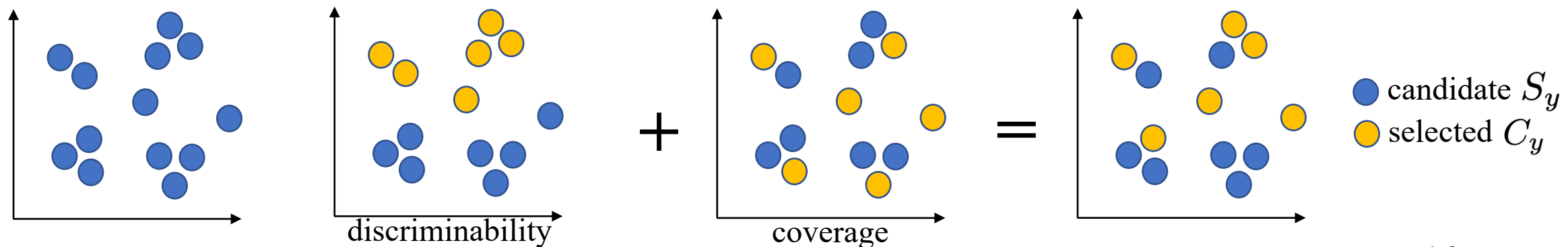
$$\mathcal{F}(C_y) = \underbrace{\alpha \cdot \sum_{c \in C_y} D(c)}_{\text{discriminability}} + \underbrace{\beta \cdot \sum_{c_1 \in S_y} \max_{c_2 \in C_y} \phi(c_1, c_2)}_{\text{coverage}}$$



# Submodular Concept Selection

- Given a superset of concepts  $S_y$  for a class  $y$ .
- Select a subset  $C_y$  for the bottleneck which are **discriminative** and **diverse**.

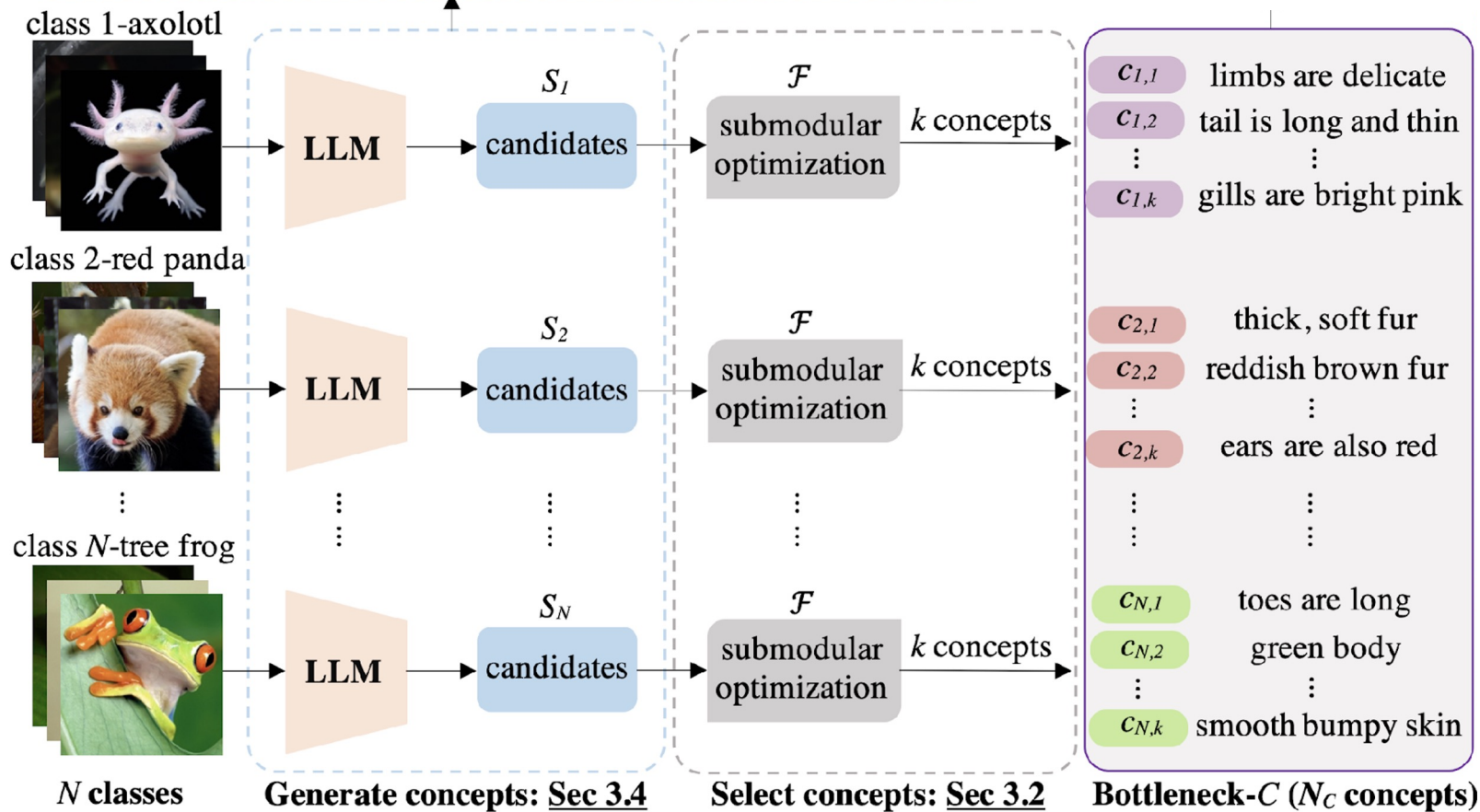
$$\mathcal{F}(C_y) = \underbrace{\alpha \cdot \sum_{c \in C_y} D(c)}_{\text{discriminability}} + \underbrace{\beta \cdot \sum_{c_1 \in S_y} \max_{c_2 \in C_y} \phi(c_1, c_2)}_{\text{coverage}}$$



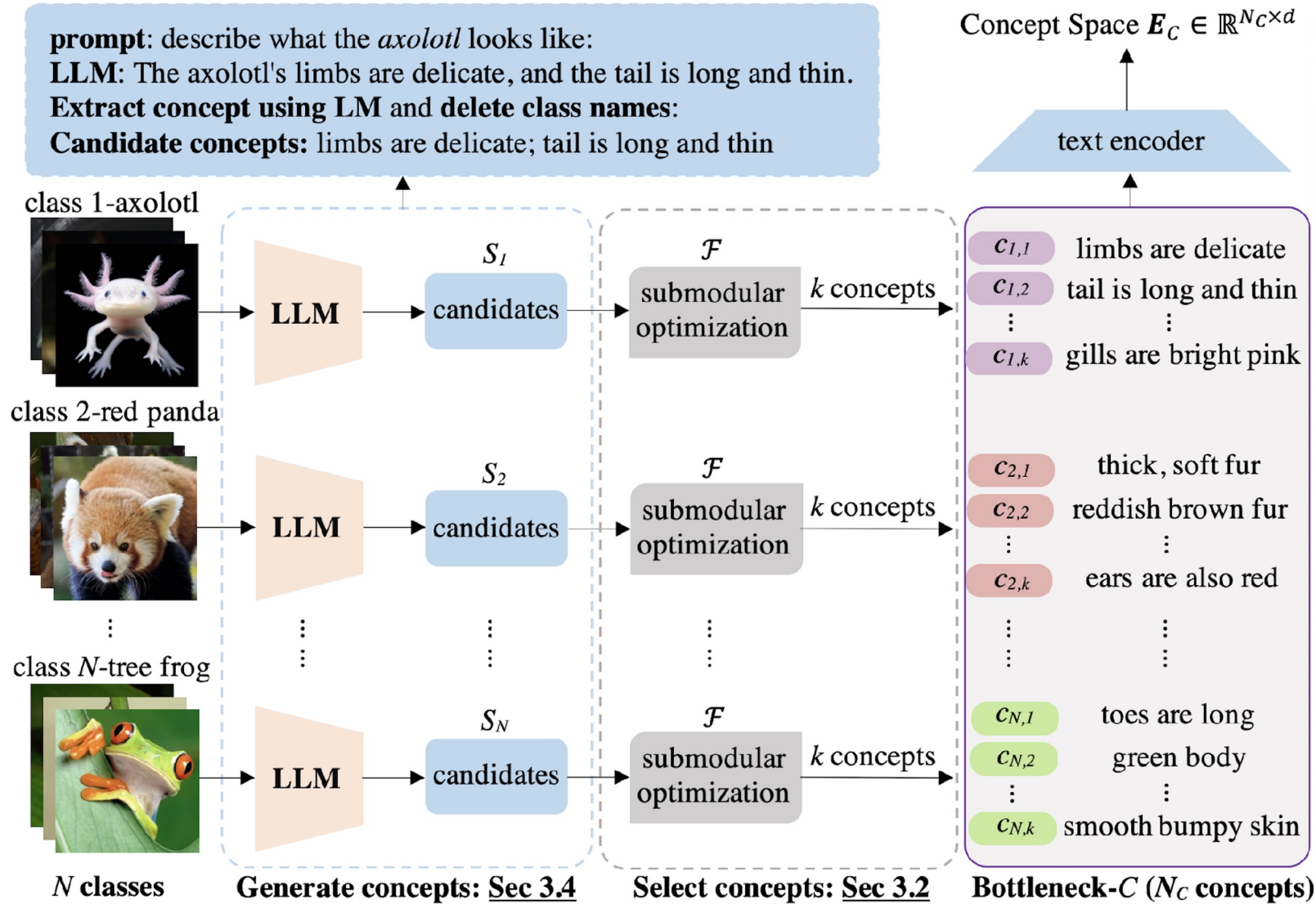


# Compute Concept Scores

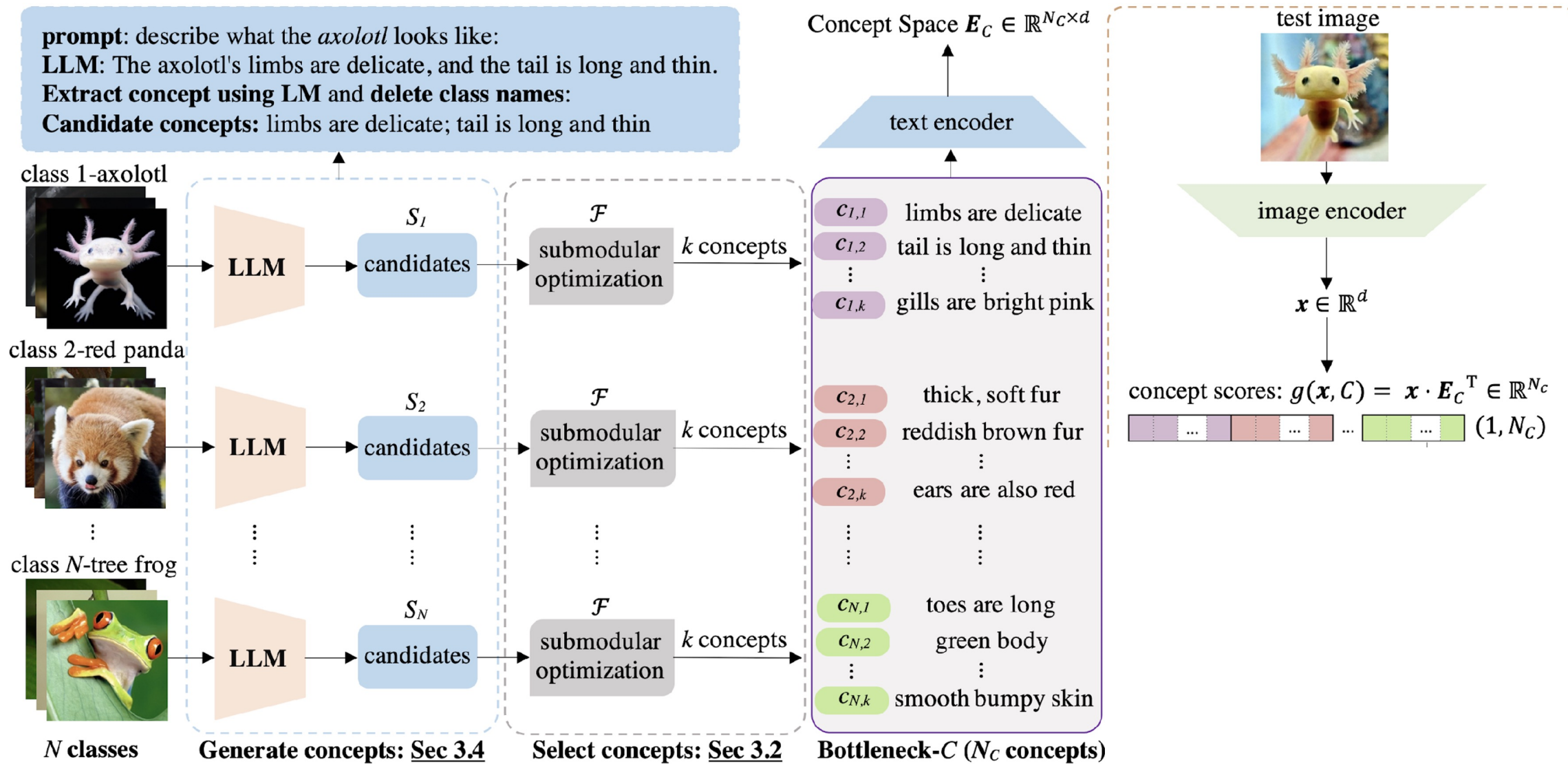
**prompt:** describe what the *axolotl* looks like:  
**LLM:** The axolotl's limbs are delicate, and the tail is long and thin.  
**Extract concept using LM and delete class names:**  
**Candidate concepts:** limbs are delicate; tail is long and thin



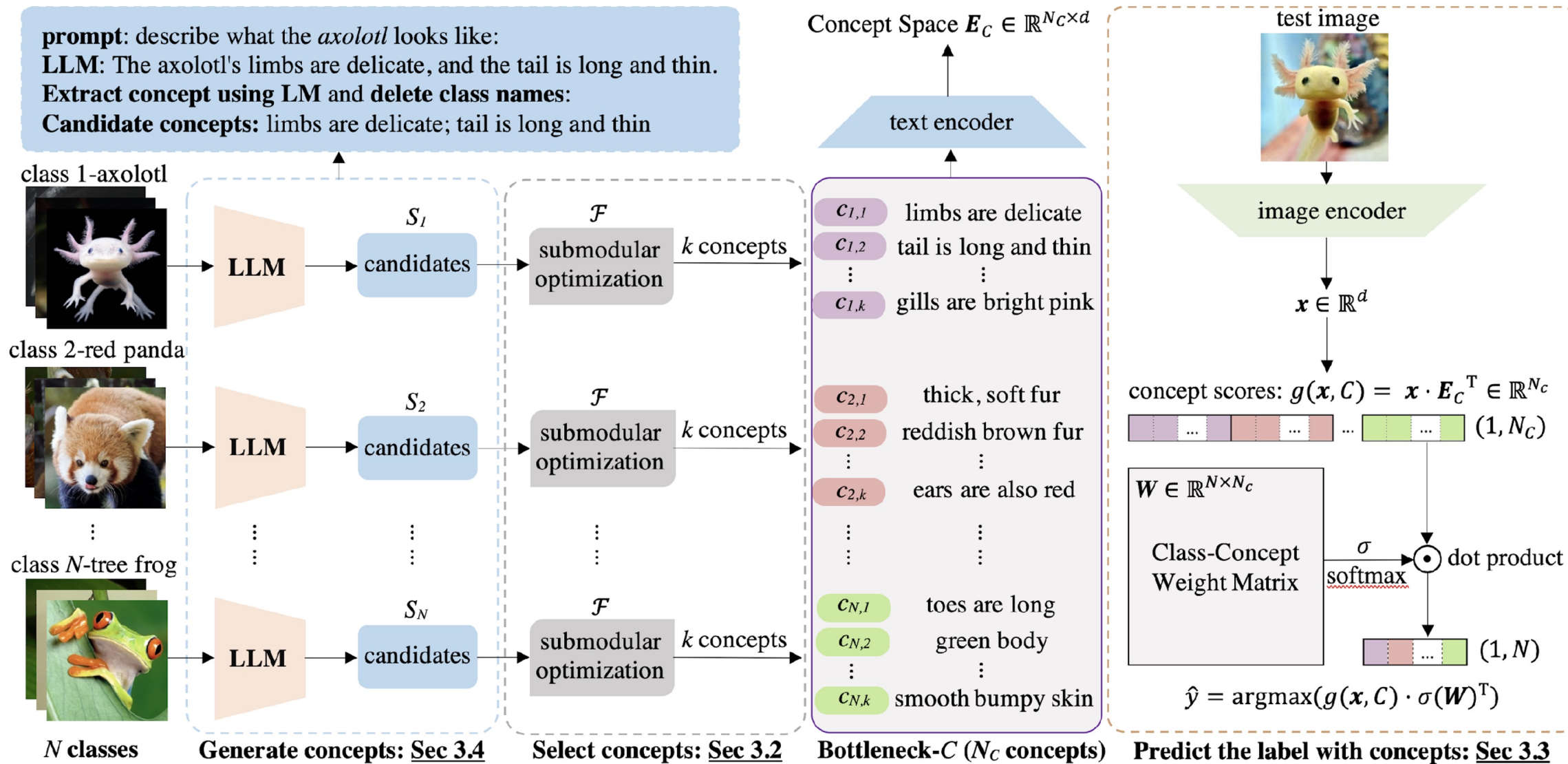
# Compute Concept Scores



# Compute Concept Scores



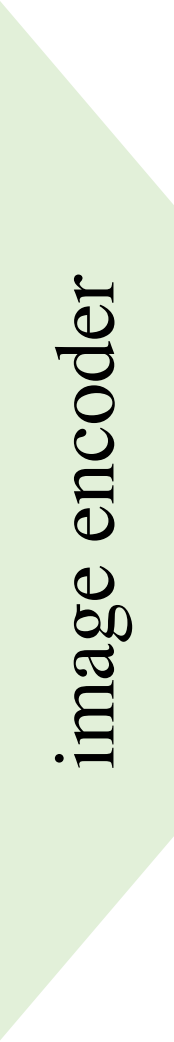
# Compute Concept Scores



# Concept scores and label prediction

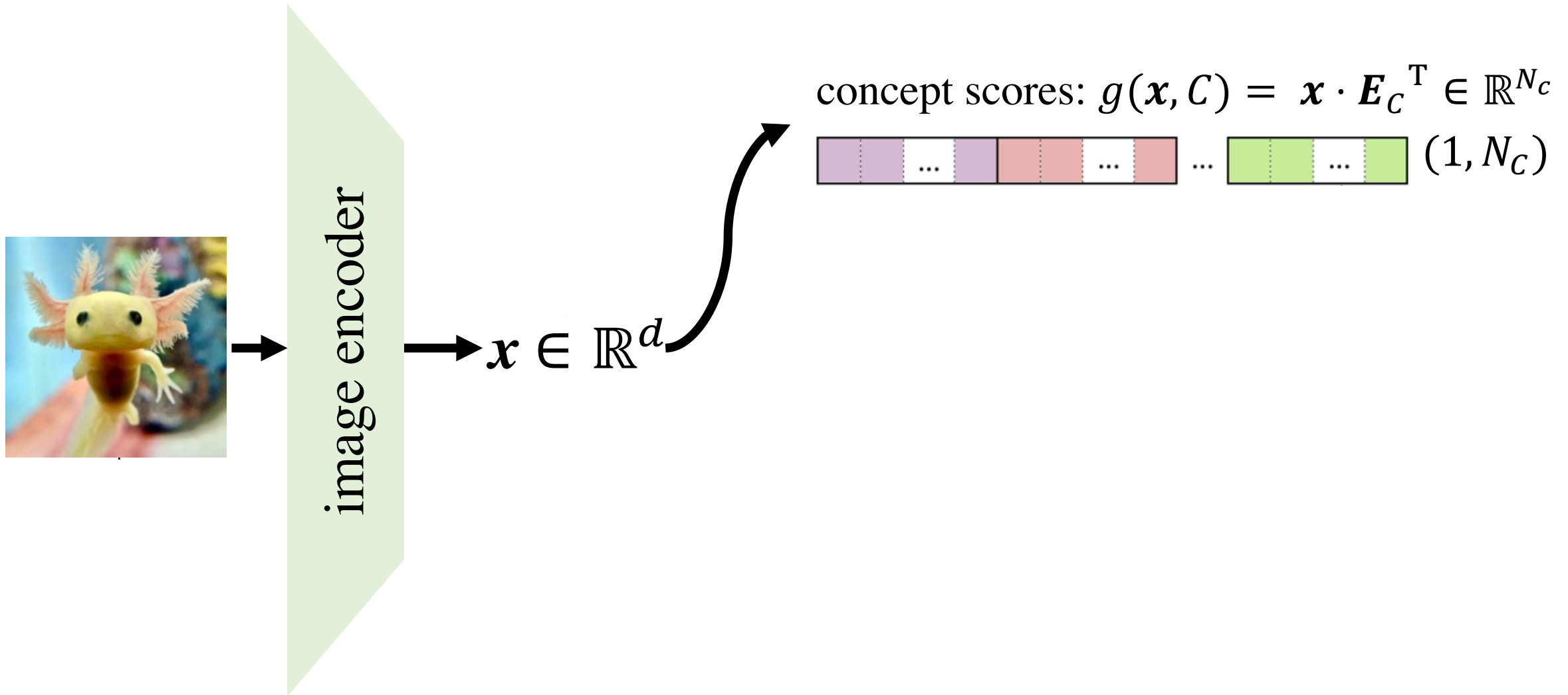


# Concept scores and label prediction

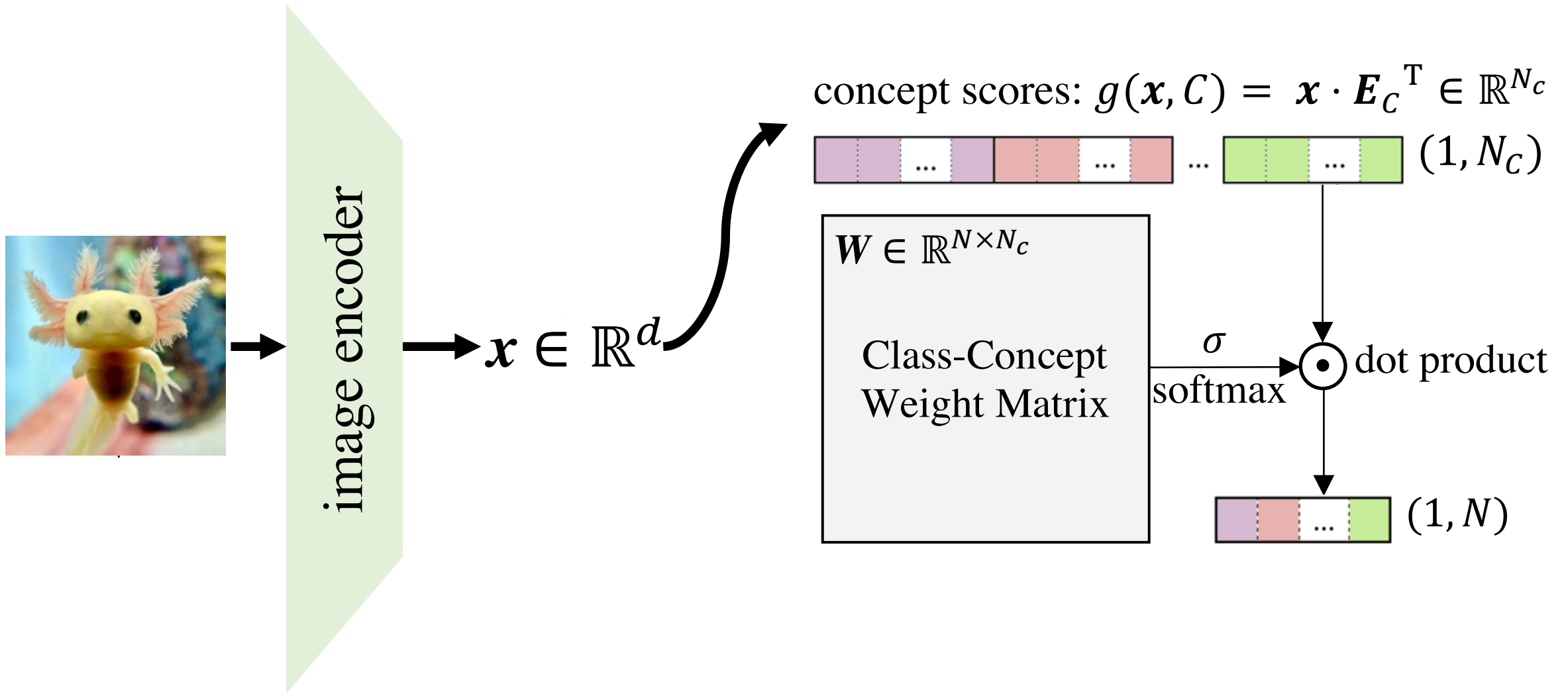


$$\mathbf{x} \in \mathbb{R}^d$$

# Concept scores and label prediction

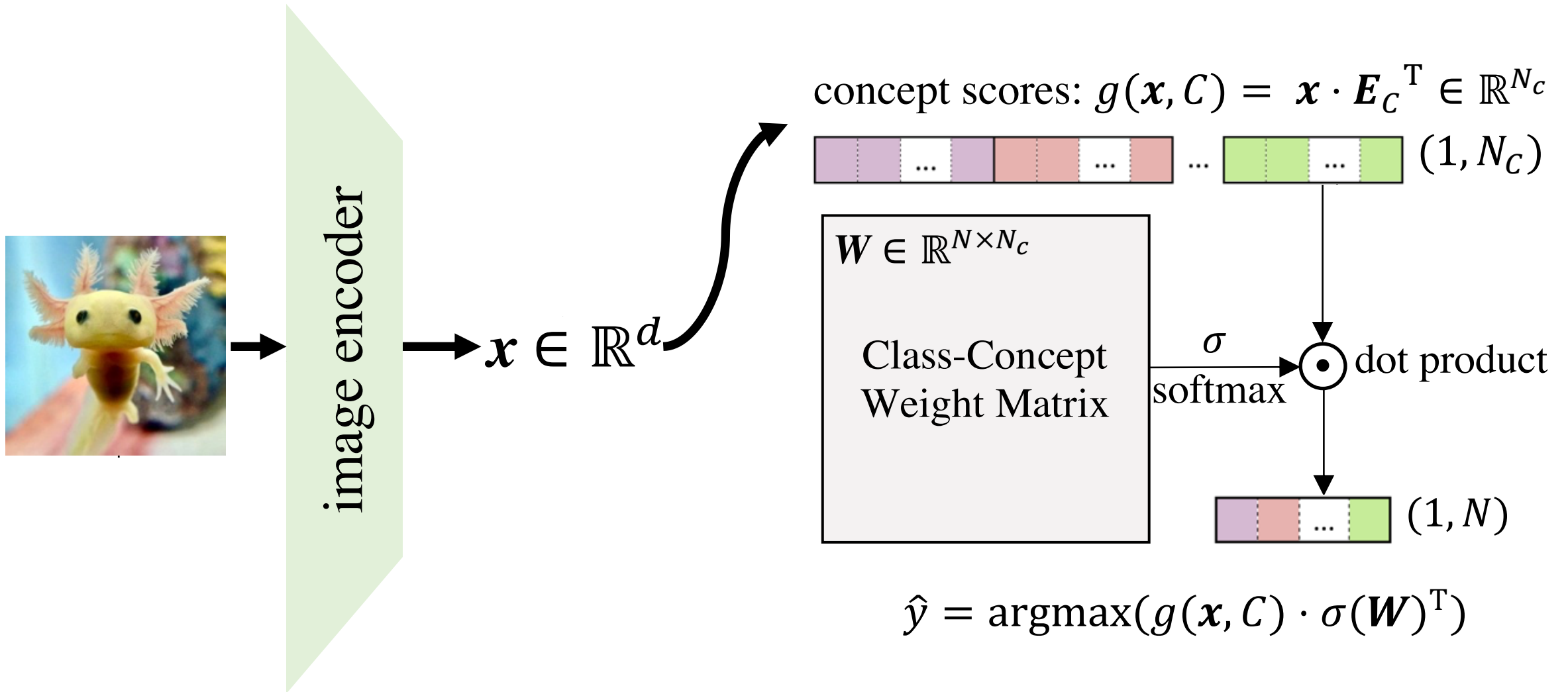


# Concept scores and label prediction





# Concept scores and label prediction



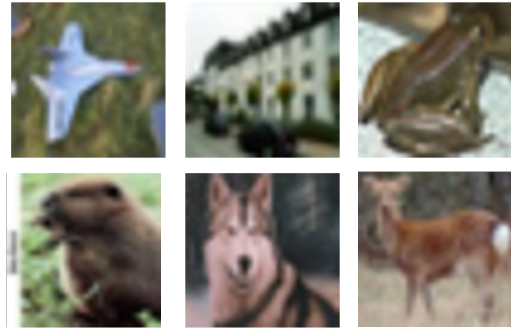
# Datasets

## Common Objects

ImageNet1K



CIFAR-10/CIFAR-100

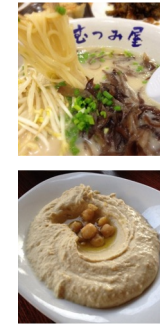


## Fine-grained Objects

Flower-102



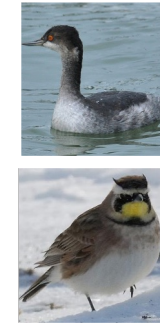
Food-101



Aircraft

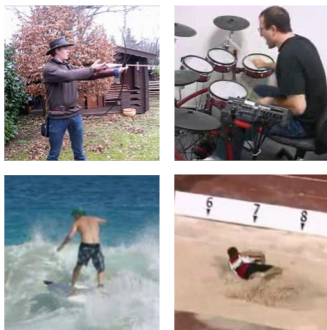


CUB



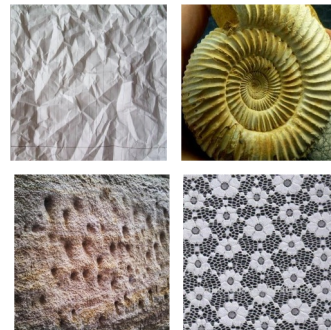
## Action

UCF-101



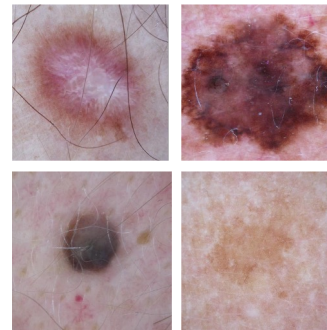
## Textures

DTD



## Skin Tumors

HAM10000



## Satellite

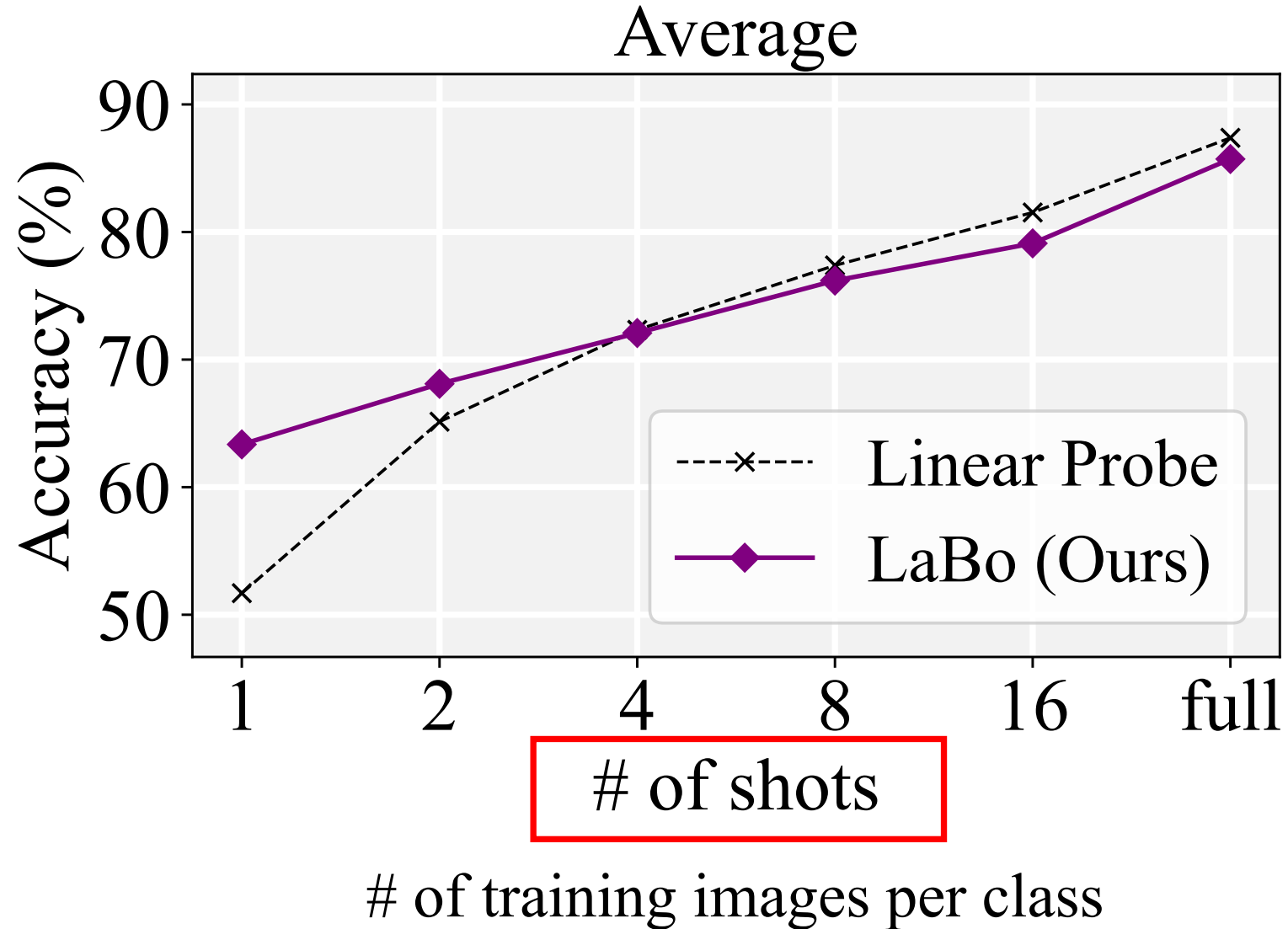
RESISC45



# Experimental Setup

- **Baselines:**
  - Linear Probe: logistic regression on the image features.
  - PCBM: Post-hoc CBM (Yuksekgonul et al., 2022)
    - Ensemble CBM prediction with end-to-end prediction.
  - ComDL: Compositional Derivation Learning (Yun et al., 2022)
    - Human designed concepts.
    - Linear layer over CLIP similarity scores.
- **Few-shot/Fully-supervised.**
- **Metric:** accuracy.

# Comparison to Black-box Model



# Comparison to Blackbox Model

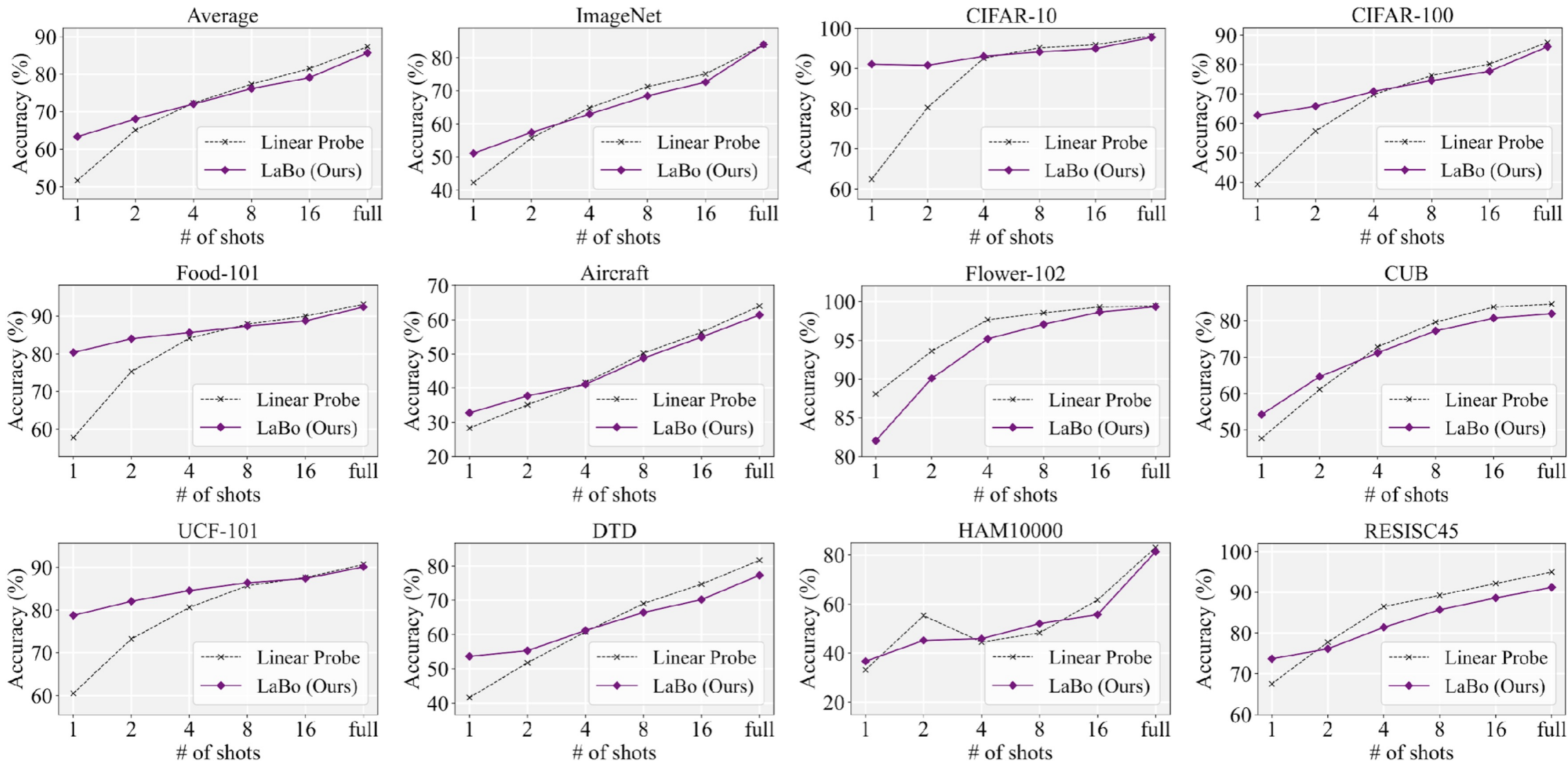


Figure 3. Test accuracy (%) comparison between LaBo and Linear Probe on 11 datasets. The x-axis represents the number of labeled images.

# Comparison to Blackbox Model

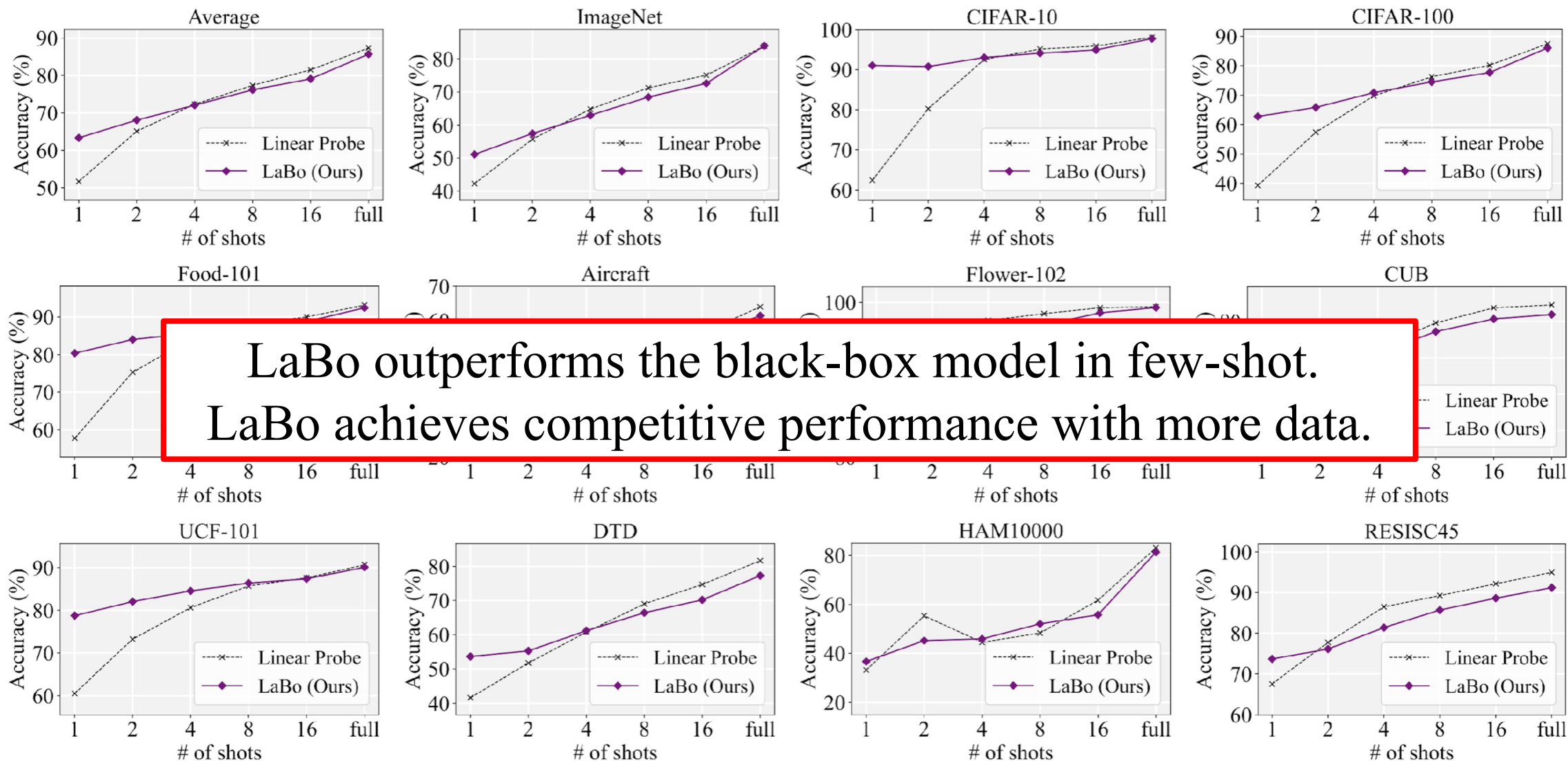


Figure 3. Test accuracy (%) comparison between LaBo and Linear Probe on 11 datasets. The x-axis represents the number of labeled images.

# Compare with Previous CBM

Method	w/ end-to-end	CIFAR-10	CIFAR-100
PCBM [66]	✗	84.5	56.0
LaBo (Ours)	✗	<b>87.9</b>	<b>69.1</b>
PCBM-h [66]	✓	87.6	69.9
Linear Probe	✓	88.8	70.1

Table 2. Test accuracy comparison between LaBo and Post-hoc Concept Bottleneck Model (PCBM) on CIFAR-10 and CIFAR-100. “w/ end-to-end” denotes whether the model employs an end-to-end residual predictor from image features to targets.

Method	w/ manual concepts	1	5	Full
CompDL [67]	✓	13.6	33.2	52.6
LaBo (Ours)	✗	<b>35.1</b>	<b>55.7</b>	<b>71.8</b>
Linear Probe	-	28.4	55.4	75.5

Table 3. LaBo and CompDL evaluated on CUB for 1/5/full shots.

# Compare with Previous CBM

Method	w/ end-to-end	CIFAR-10	CIFAR-100
PCBM [66]	✗	84.5	56.0
LaBo (Ours)	✗	<b>87.9</b>	<b>69.1</b>
PCBM-h [66]	✓	87.6	69.9
Linear Probe	✓	88.8	70.1

LaBo doesn't rely on black box predictors.  
LaBo doesn't require human annotations.

oc  
00.  
nd

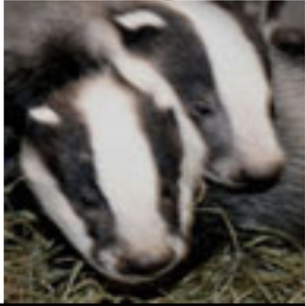



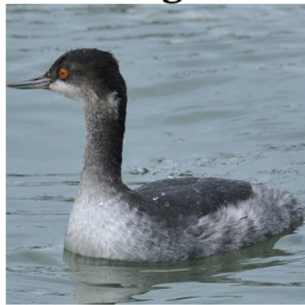

residual predictor from image features to targets.

Method	w/ manual concepts	1	5	Full
CompDL [67]	✓	13.6	33.2	52.6
LaBo (Ours)	✗	<b>35.1</b>	<b>55.7</b>	<b>71.8</b>
Linear Probe	-	28.4	55.4	75.5



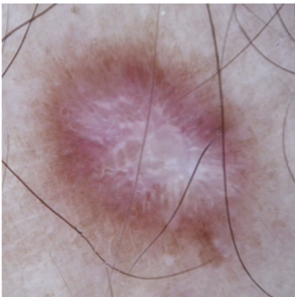
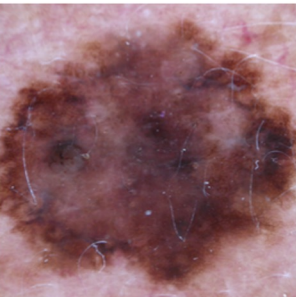

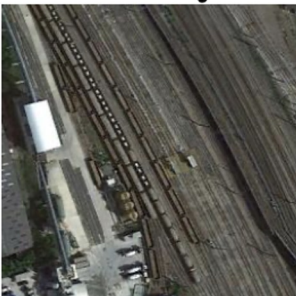
Table 3. LaBo and CompDL evaluated on CUB for 1/5/full shots.



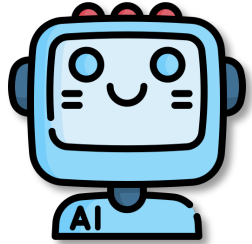
# Qualitative Results

	Class Name	Top-3 Concepts	Class Name	Top-3 Concepts
ImageNet	<b>badger</b> 	<ol style="list-style-type: none"> <li>1. short legs and long body make it an excellent digger</li> <li>2. black-and-white striped fur</li> <li>3. coat is very shaggy</li> </ol>	<b>ant</b> 	<ol style="list-style-type: none"> <li>1. black and red stinger</li> <li>2. small, black insect with six legs</li> <li>3. long, slender antennae that it uses to smell and touch</li> </ol>
	<b>ramen</b> 	<ol style="list-style-type: none"> <li>1. garnished with green onions, nori, and other toppings</li> <li>2. most grocery stores</li> <li>3. various toppings</li> </ol>	<b>hummus</b> 	<ol style="list-style-type: none"> <li>1. chickpeas, tahini, olive oil, garlic, lemon juice</li> <li>2. made from cooked, mashed chickpeas</li> <li>3. roasted red peppers</li> </ol>
CUB	<b>eared grebe</b> 	<ol style="list-style-type: none"> <li>1. black and white plumage that is striking in the sunlight</li> <li>2. black body with a long, slender neck</li> <li>3. red and black bill</li> </ol>	<b>horned lark</b> 	<ol style="list-style-type: none"> <li>1. black line running through yellow face</li> <li>2. head is black with a white horn on each side</li> <li>3. black horn on each side of their head</li> </ol>

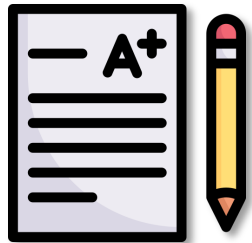
# Qualitative Results

<b>UCF-101</b>	<b>archery</b> 	<ol style="list-style-type: none"><li>1. grip bow tightly in their left hand</li><li>2. focused and concentrated on their task</li><li>3. keep bow and arrows in safe and dry place when not in use</li></ol>	<b>drumming</b> 	<ol style="list-style-type: none"><li>1. blur as they fly over the drums</li><li>2. sitting on a stool in front of a drum set</li><li>3. position the drumstick so it is resting on your index finger</li></ol>
<b>HAM100000</b>	<b>dermatofibroma</b> 	<ol style="list-style-type: none"><li>1. generally not painful</li><li>2. red, brown, or purple in color</li><li>3. thin white halo around them</li></ol>	<b>melanoma</b> 	<ol style="list-style-type: none"><li>1. dark brown or black in color</li><li>2. large and dark</li><li>3. flesh-colored, brown, or black</li></ol>
<b>RESISC45</b>	<b>beach</b> 	<ol style="list-style-type: none"><li>1. waves crashing onto the shore</li><li>2. few rocks poking out</li><li>3. waves are gentle</li></ol>	<b>railway</b> 	<ol style="list-style-type: none"><li>1. connected by steel rails</li><li>2. tramline that is 3 feet wide and runs along the length of the court</li><li>3. faint, twinkling line</li></ol>

# Conclusion



Leverage the **knowledge of LLM** to build interpretable models (CBMs).

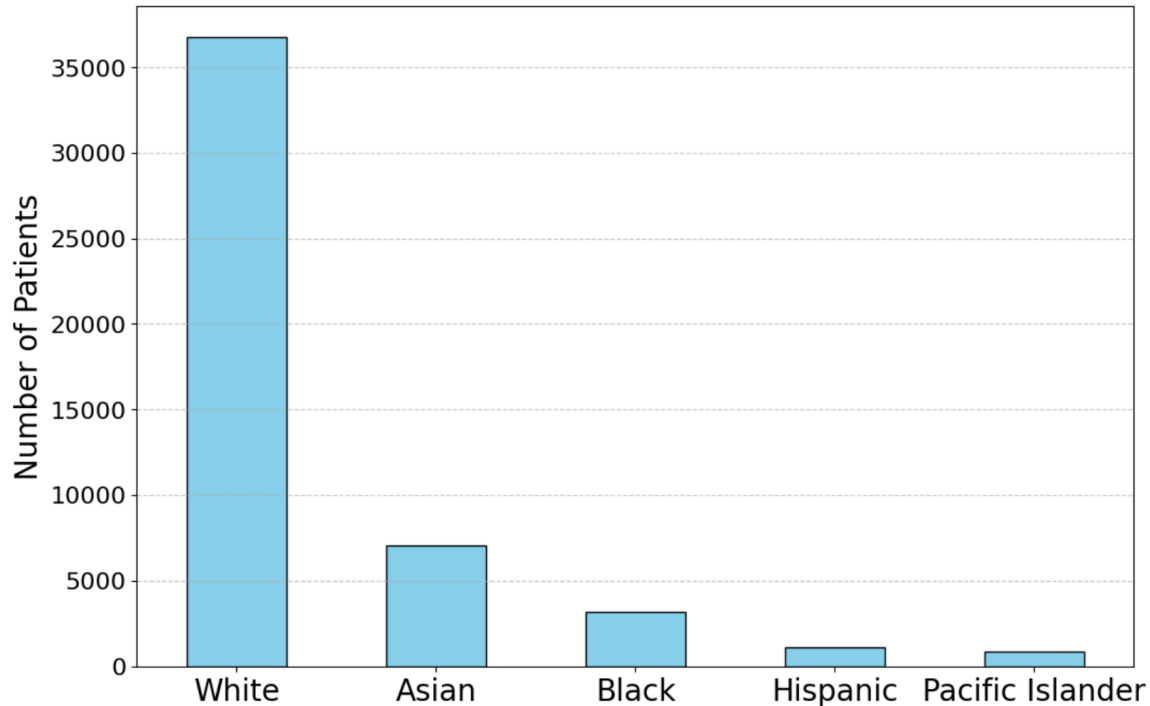


With **vision-language models (VLMs)** and concept selection, interpretable models can achieve **competitive performance** as Black-box.

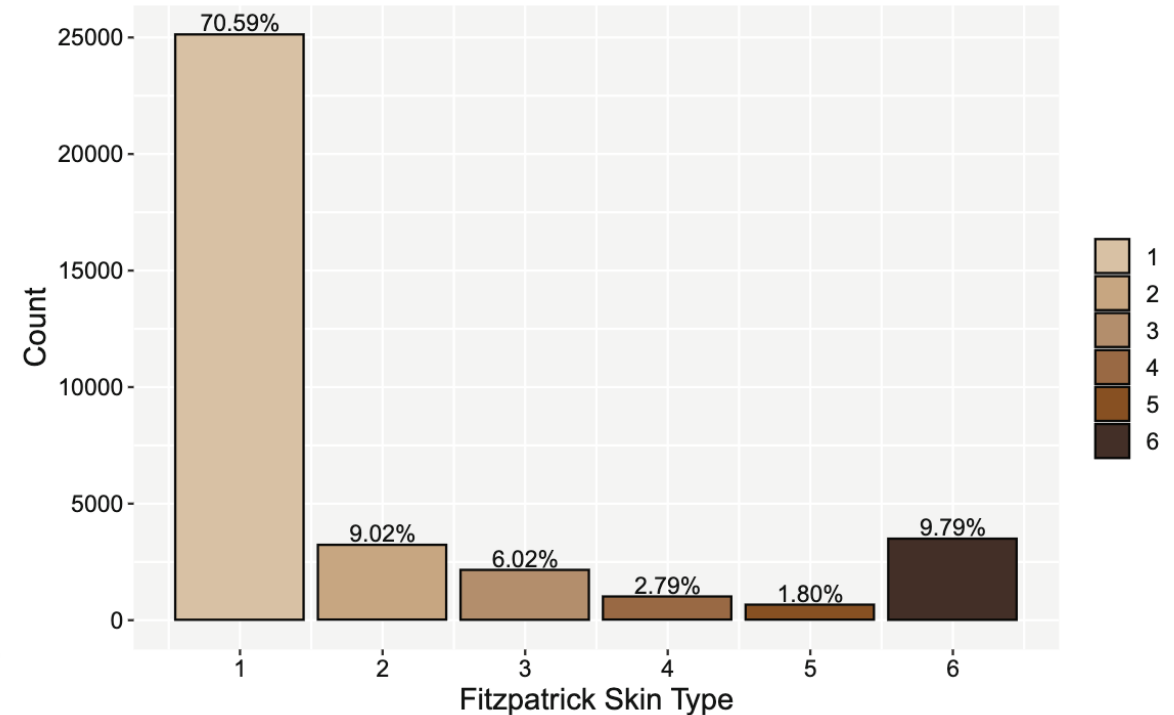
**What makes the critical domain  
more challenging?**

# The distribution of demographic variables in medical data can be skewed.

## Distribution of race in CheXpert [1].



## Distribution of skin colors in ISIC [2, 3].



### Distribution of race i

### Artificial intelligence predicts patients' race from their medical images

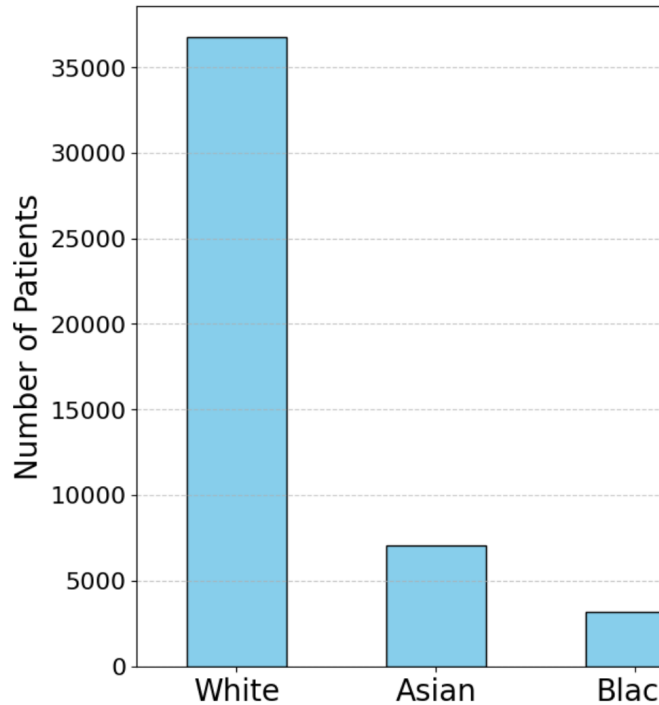
Study shows AI can identify self-reported race from medical images that contain no indications of race detectable by human experts.

Rachel Gordon | MIT CSAIL  
May 20, 2022

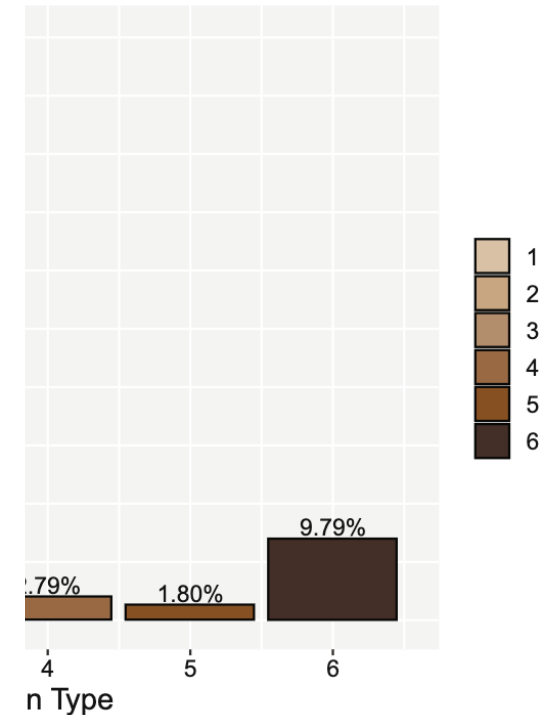


Deep models are very sensitive to demographic variables.

Gichoya et al. Lancet. 2022. Lancet.



### colors in ISIC [2, 3].

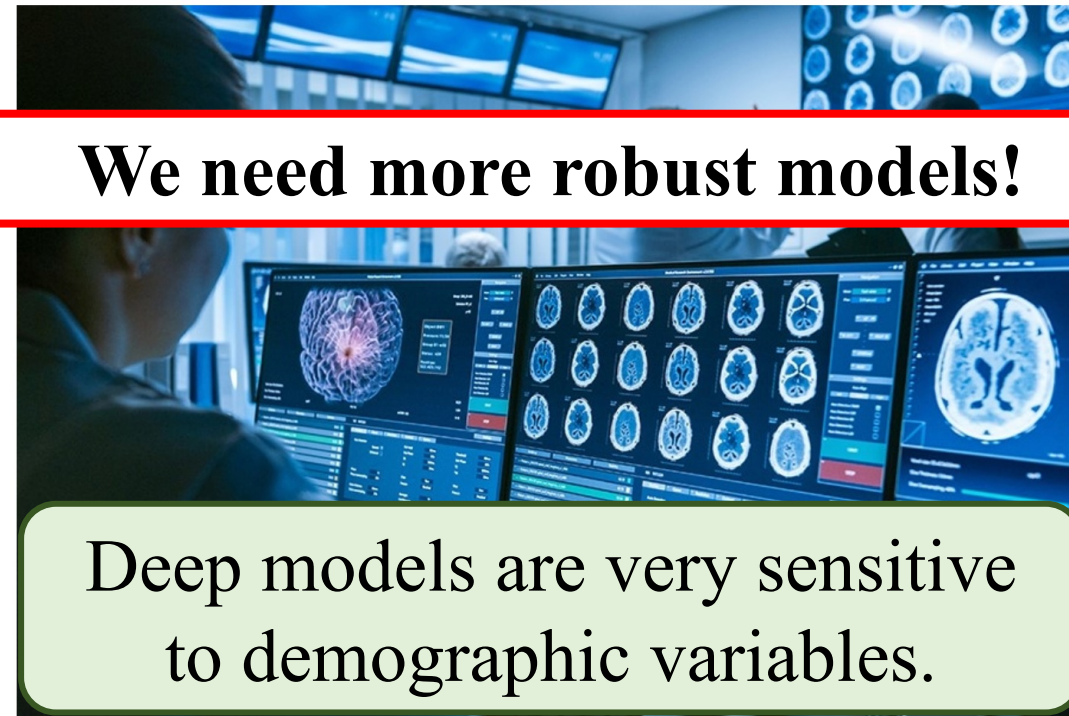


### Distribution of race in

### Artificial intelligence predicts patients' race from their medical images

Study shows AI can identify self-reported race from medical images that contain no indications of race detectable by human experts.

Rachel Gordon | MIT CSAIL  
May 20, 2022

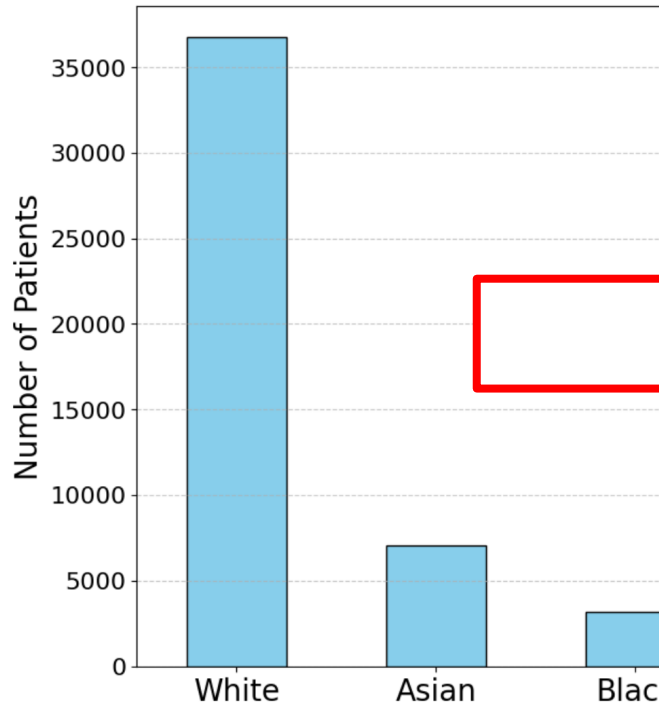
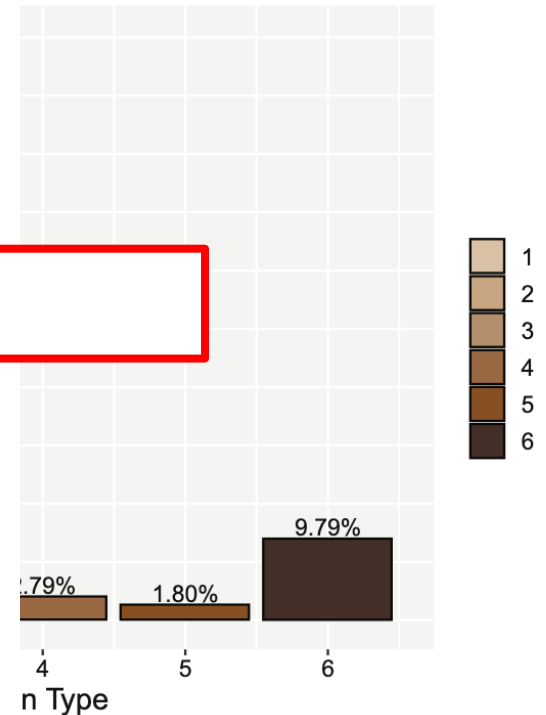


**We need more robust models!**

Deep models are very sensitive to demographic variables.

Gichoya et al. Lancet. 2022. Lancet.

### colors in ISIC [2, 3].

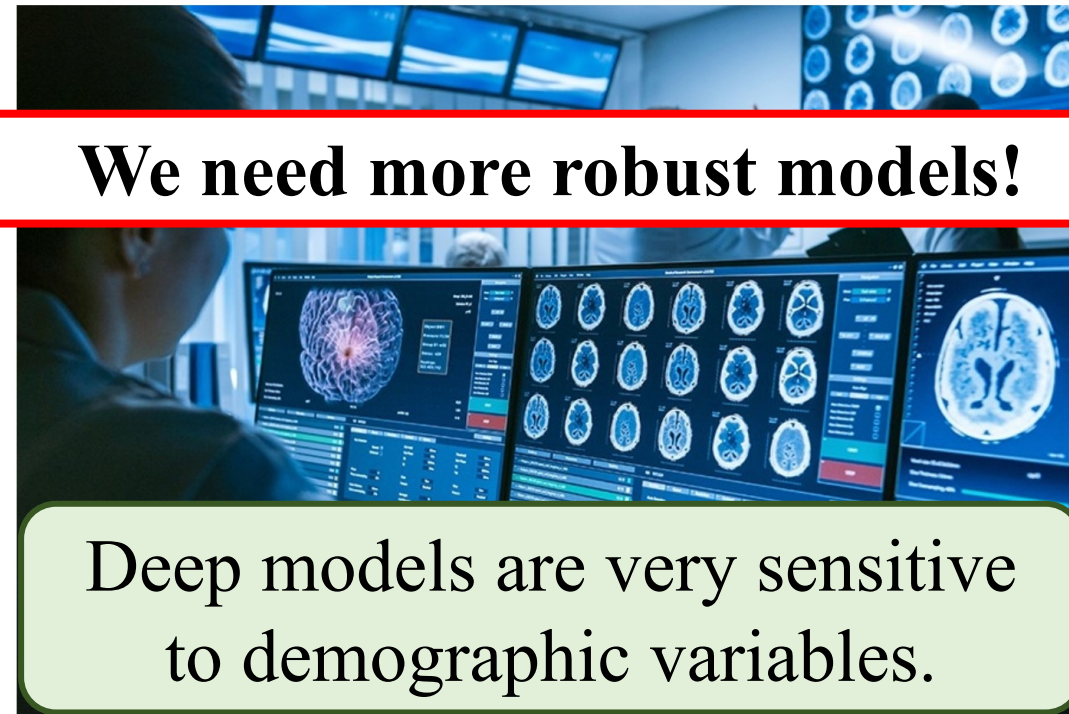


### Distribution of race in

### Artificial intelligence predicts patients' race from their medical images

Study shows AI can identify self-reported race from medical images that contain no indications of race detectable by human experts.

Rachel Gordon | MIT CSAIL  
May 20, 2022

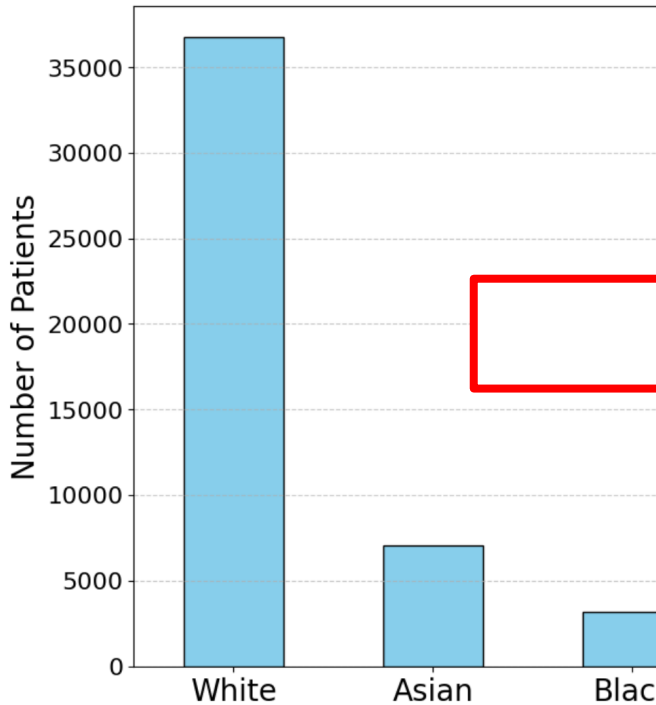
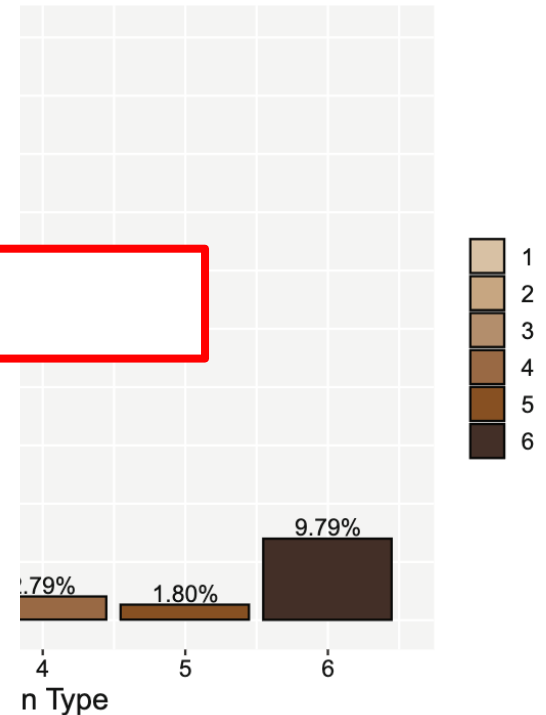


**We need more robust models!**

Deep models are very sensitive to demographic variables.

Gichoya et al. Lancet. 2022. Lancet.

### colors in ISIC [2, 3].

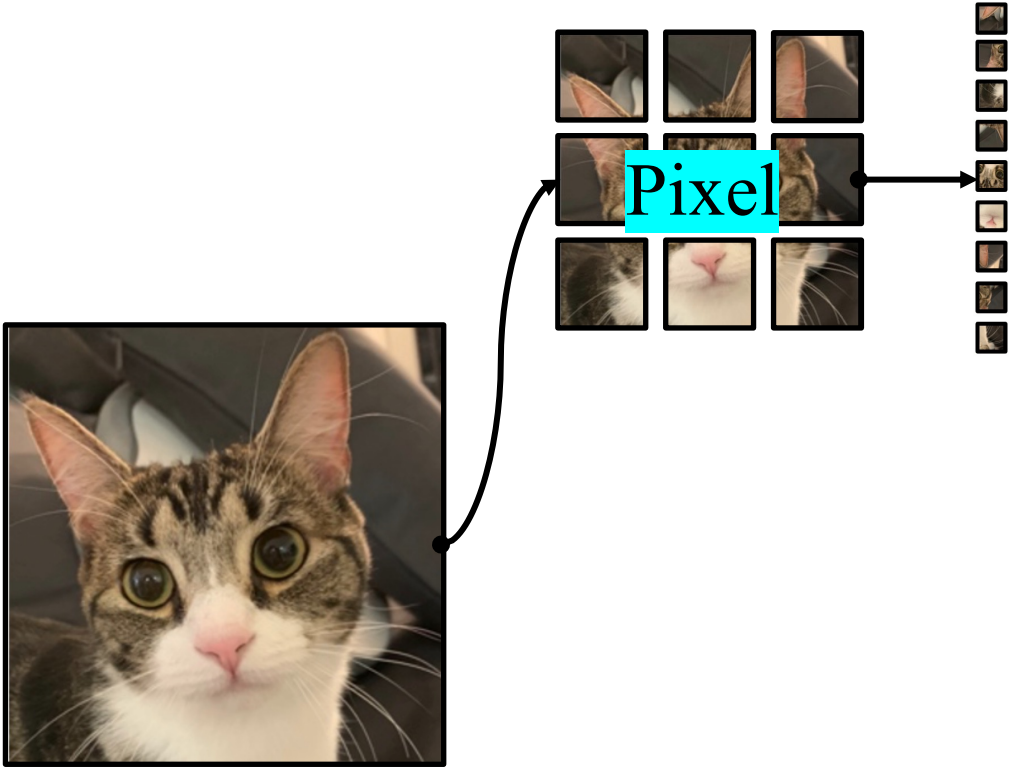




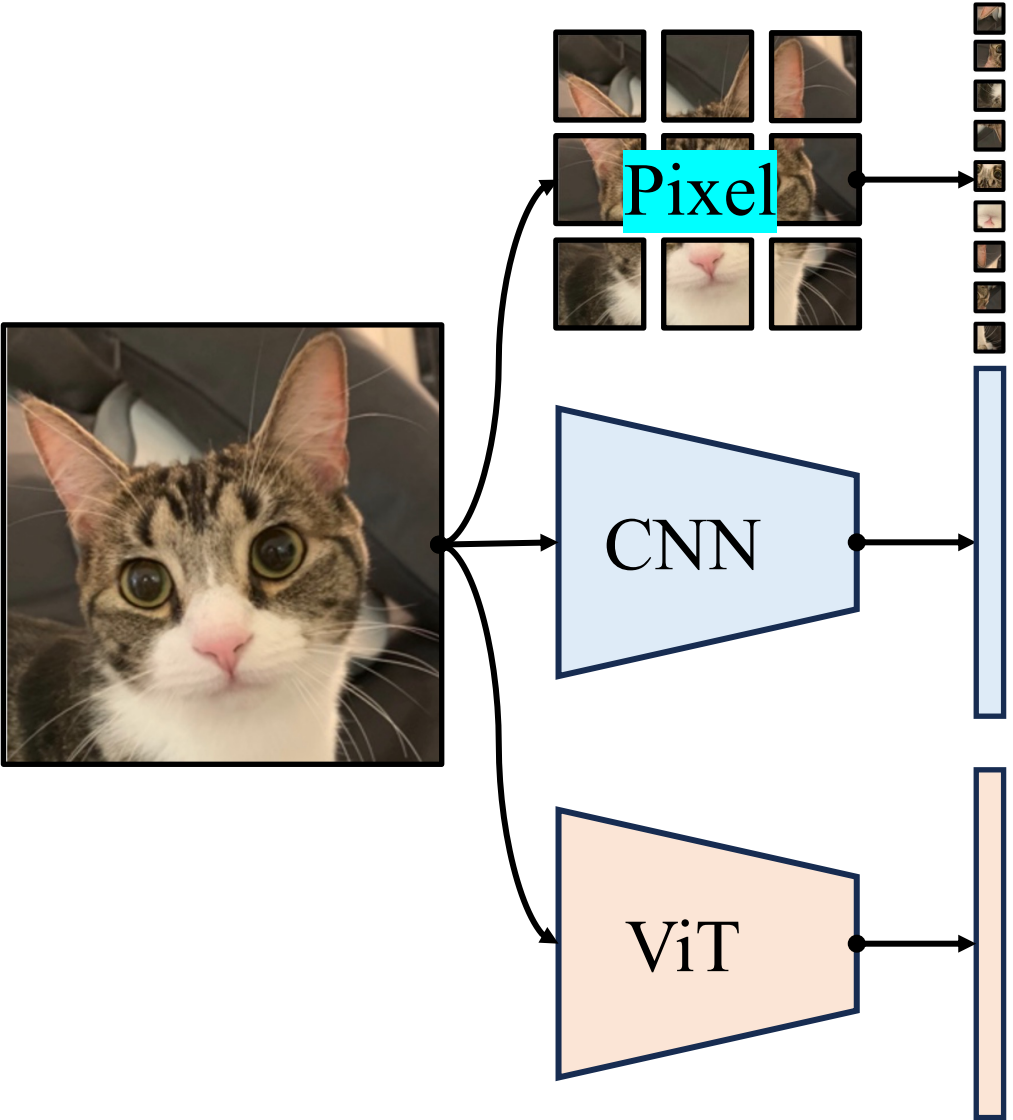
Deep models have good priors for the **general domain**.



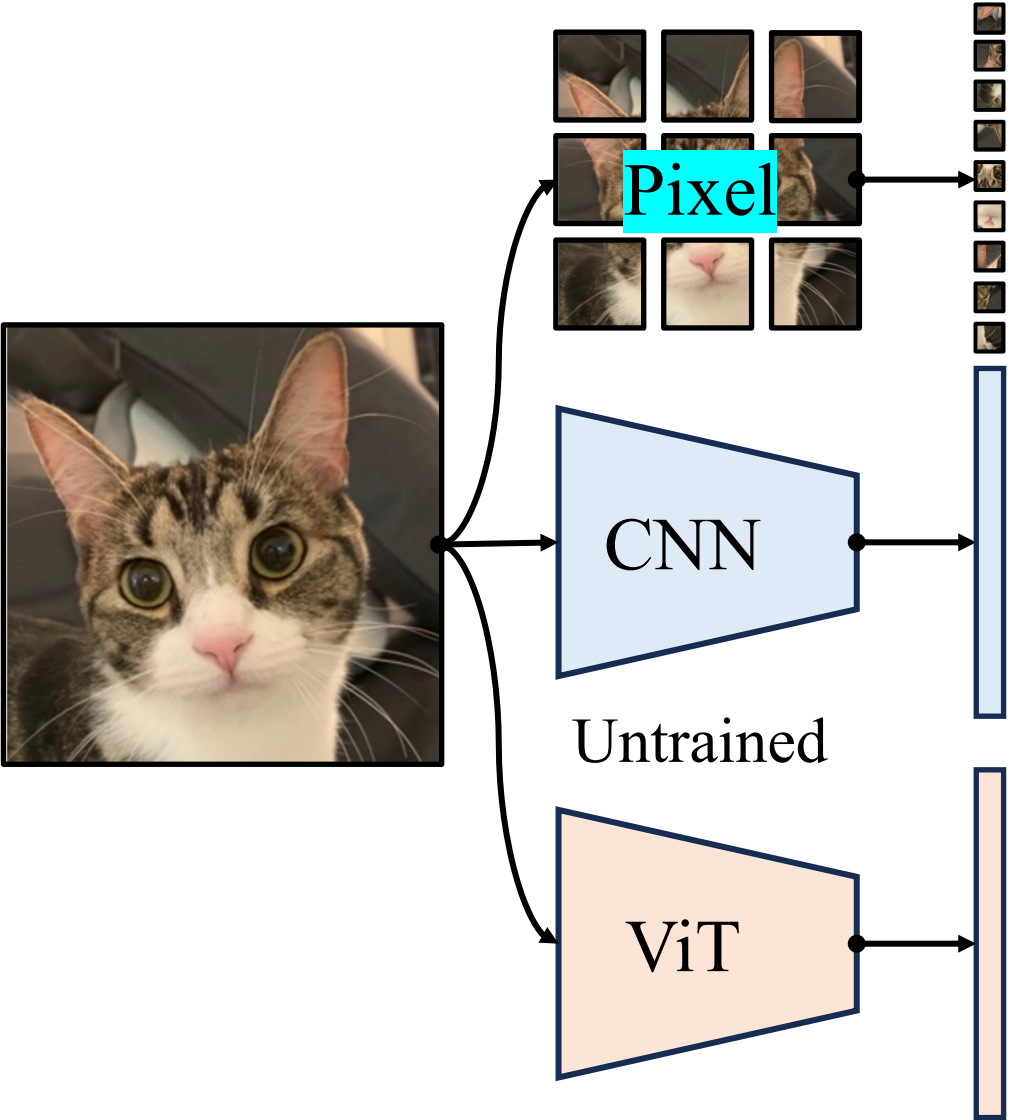
Deep models have good priors for the **general domain**.



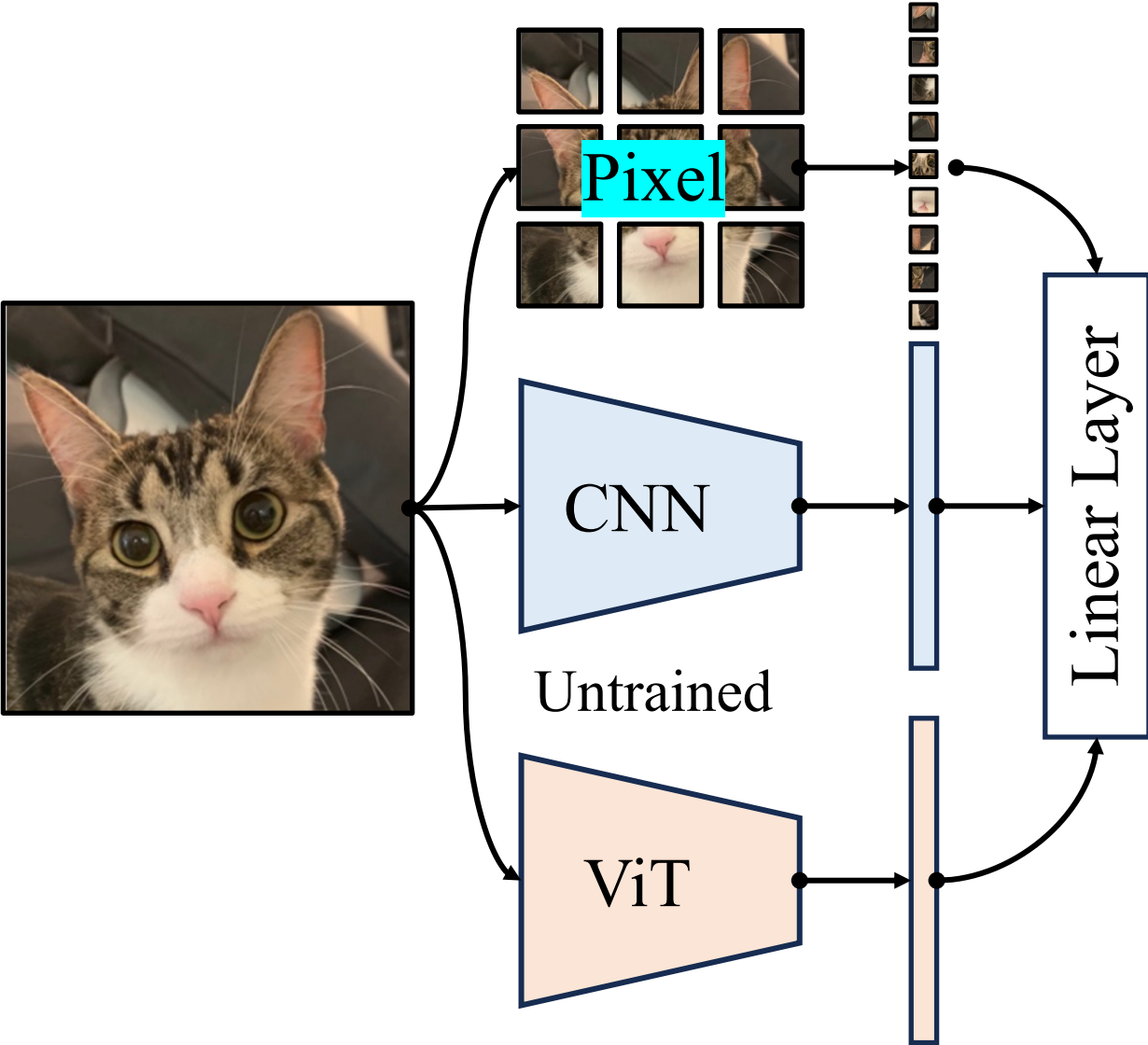
Deep models have good priors for the **general domain**.



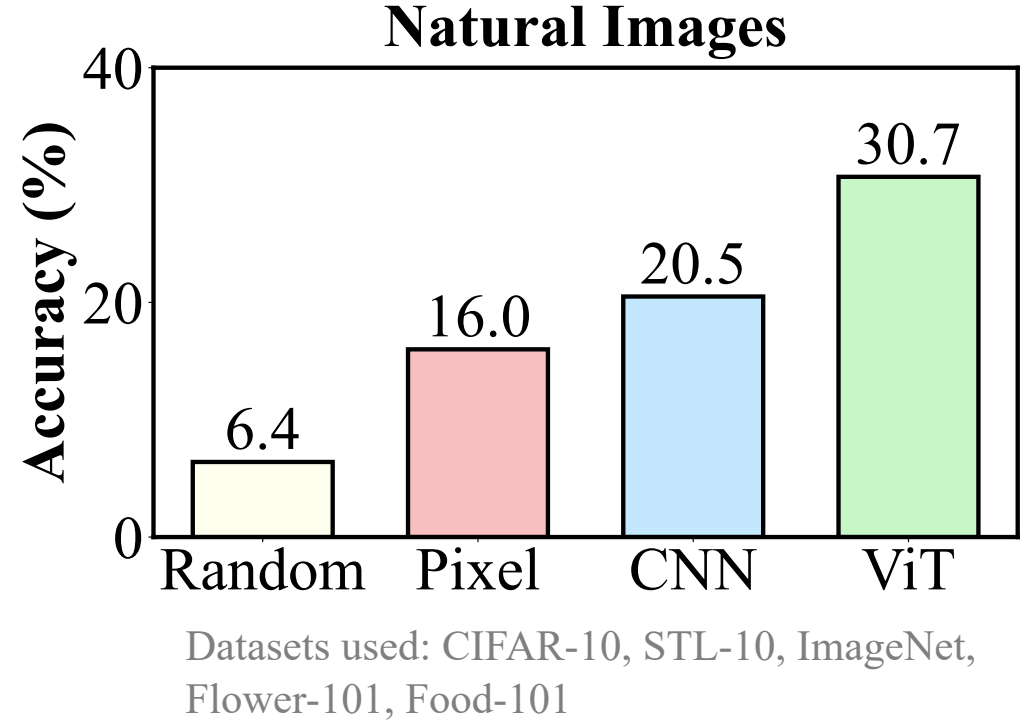
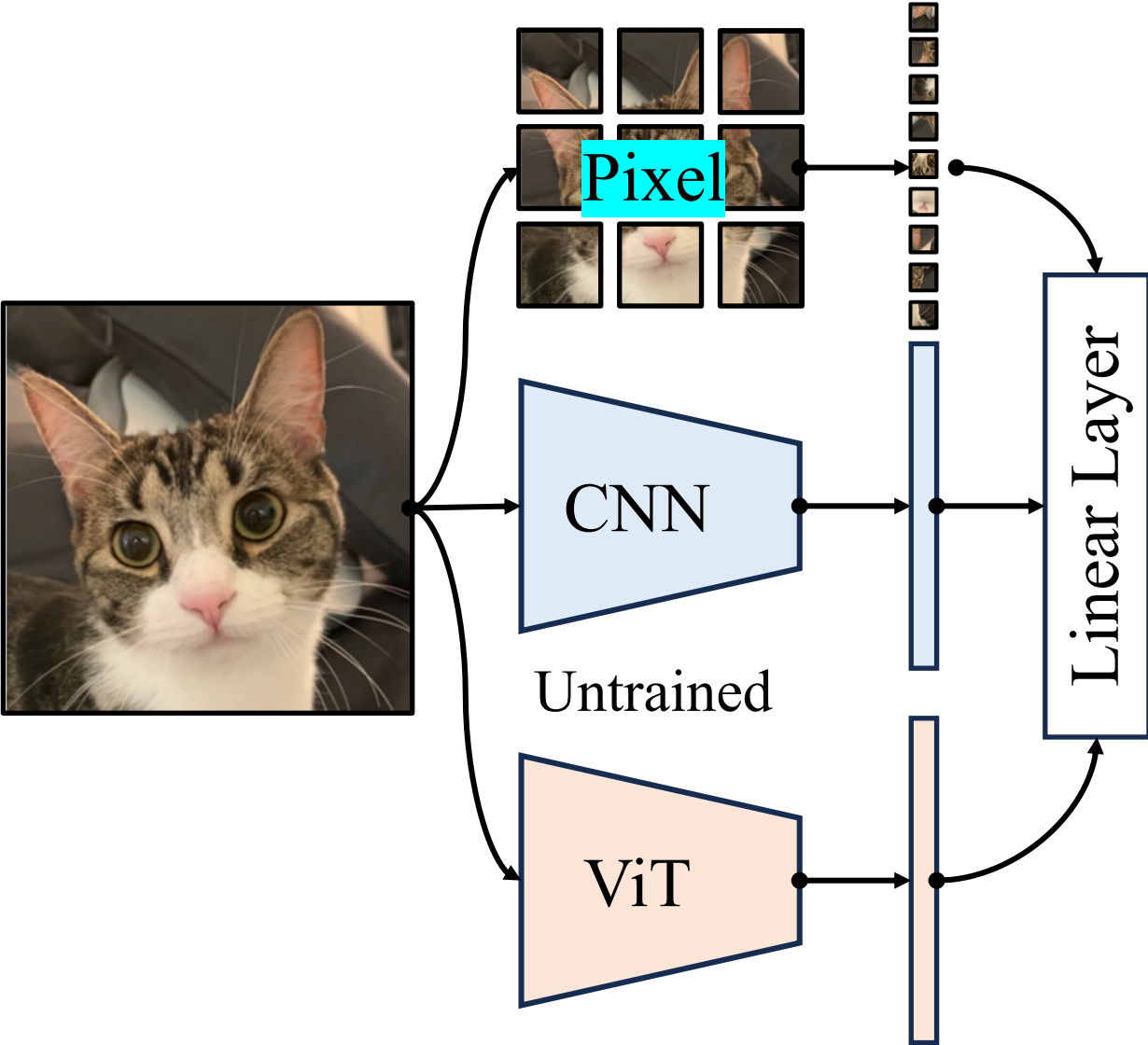
Deep models have good priors for the **general domain**.

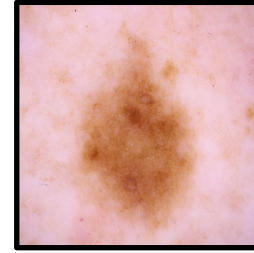


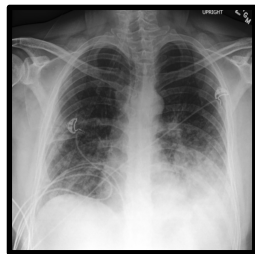
Deep models have good priors for the **general domain**.



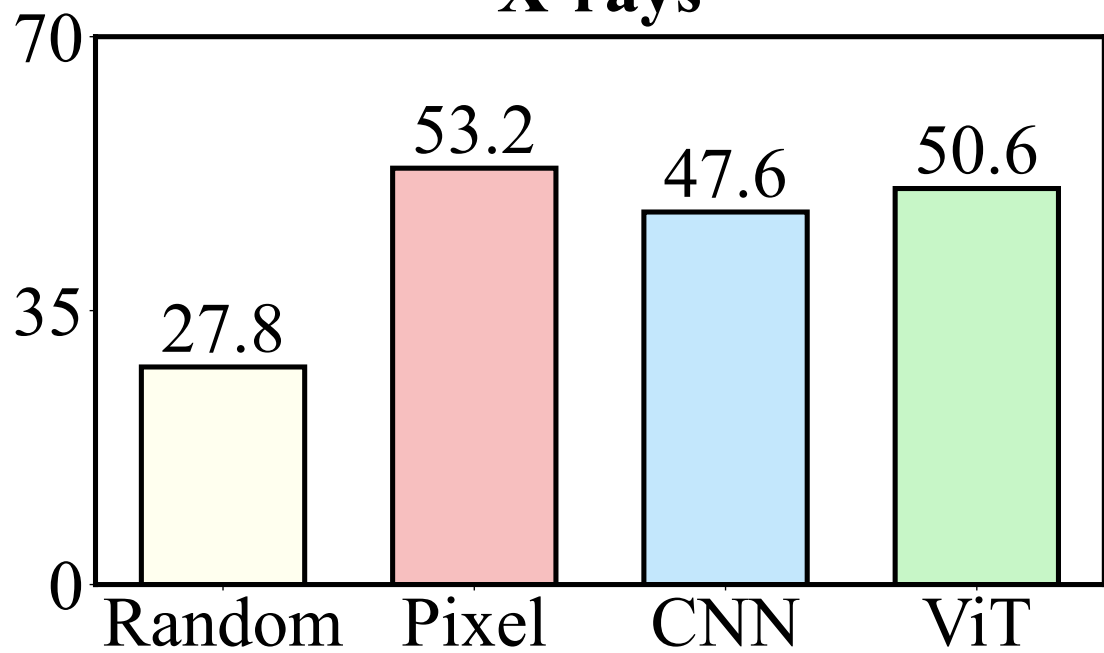
Deep models have good priors for the **general domain**.



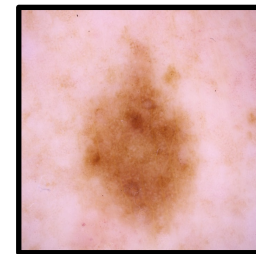




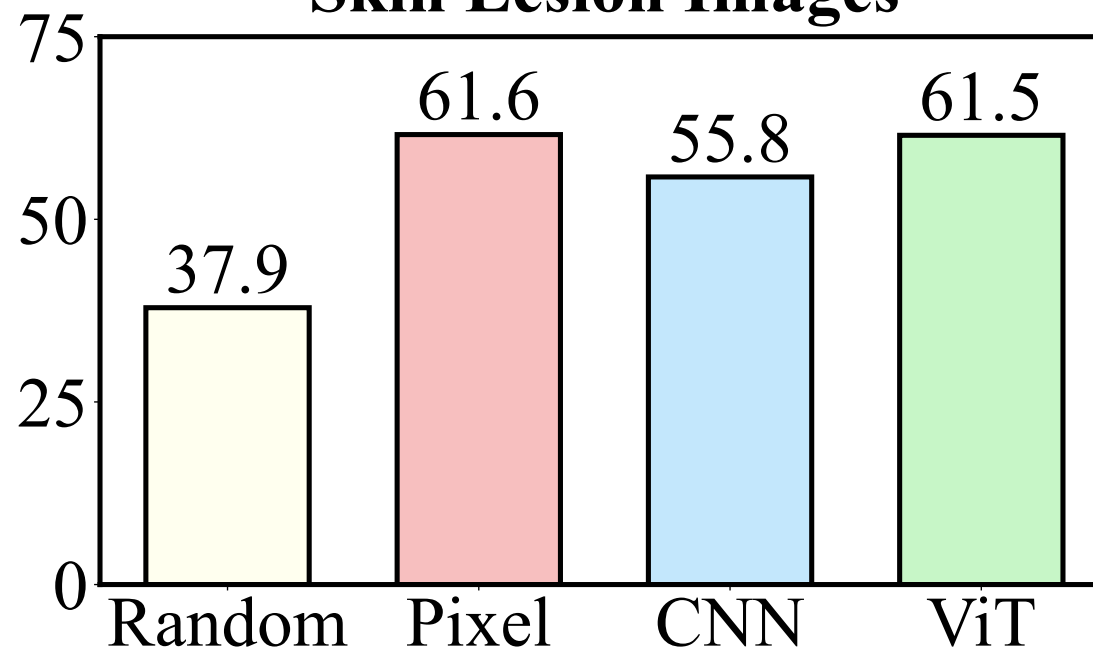
**X-rays**



X-ray Datasets: Pneumonia, COVID-QU, NIH-CXR, Open-I, VinDr-CXR.



**Skin Lesion Images**



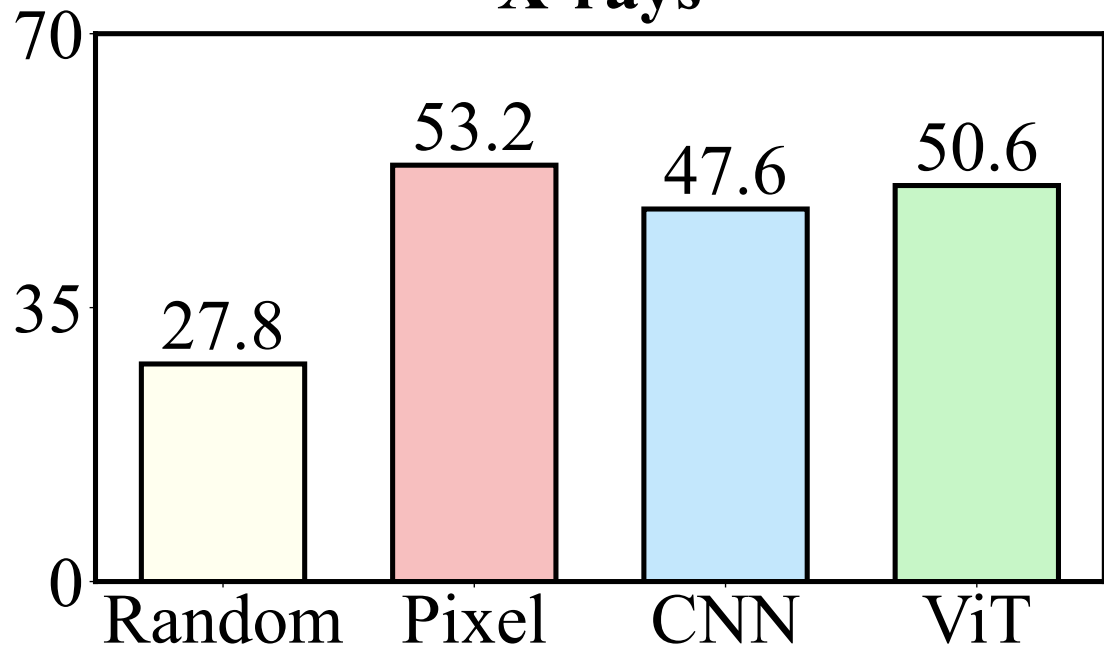
Skin Lesion Datasets: HAM10000, BCN20000, PAD-UFS-20, Melanoma, UWaterloo.



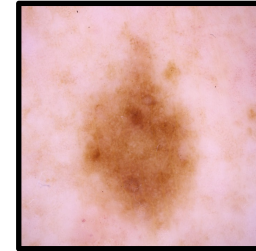
Deep models **don't** have good priors for the **medical domain**.



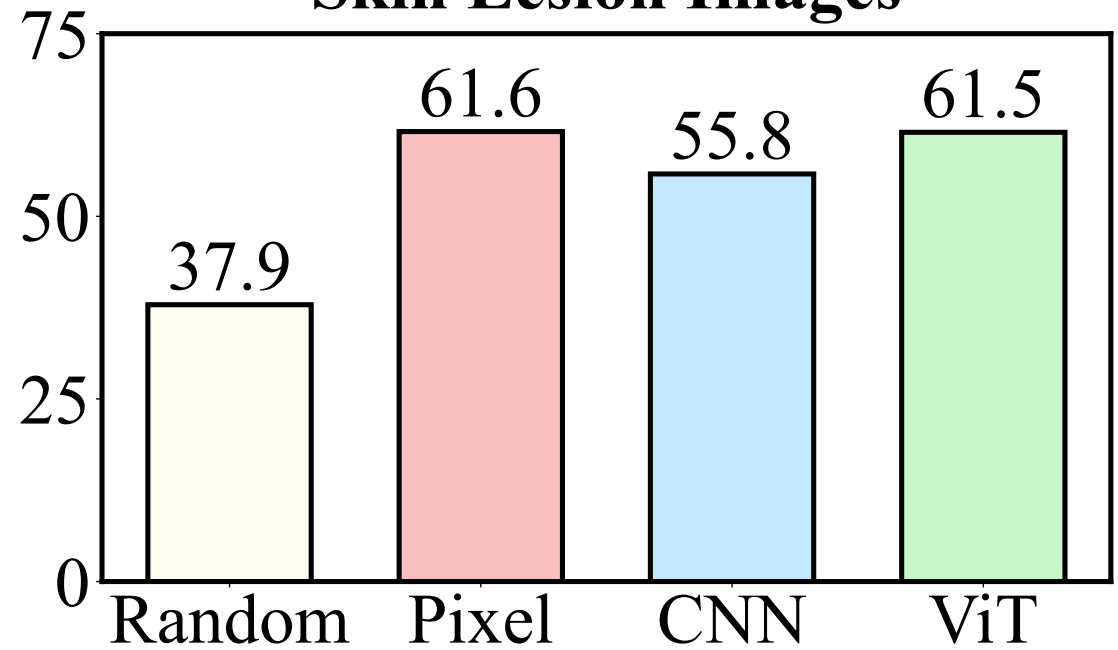
**X-rays**



X-ray Datasets: Pneumonia, COVID-QU, NIH-CXR, Open-I, VinDr-CXR.



**Skin Lesion Images**

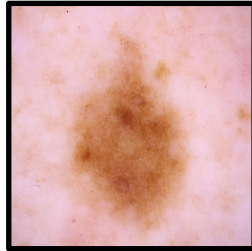


Skin Lesion Datasets: HAM10000, BCN20000, PAD-UFS-20, Melanoma, UWaterloo.

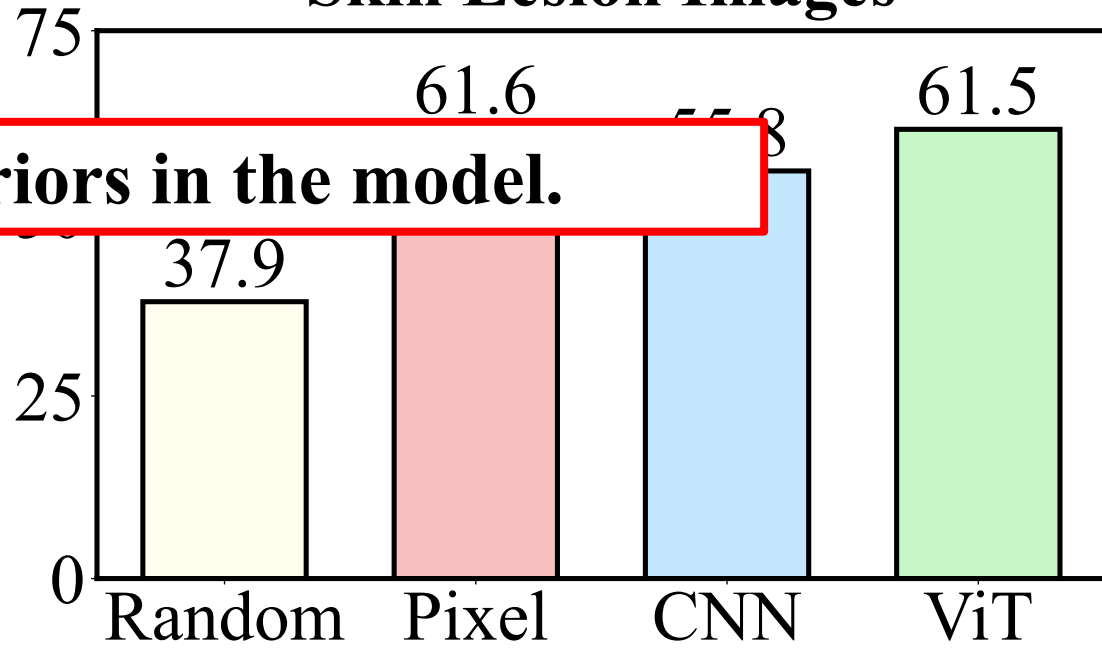
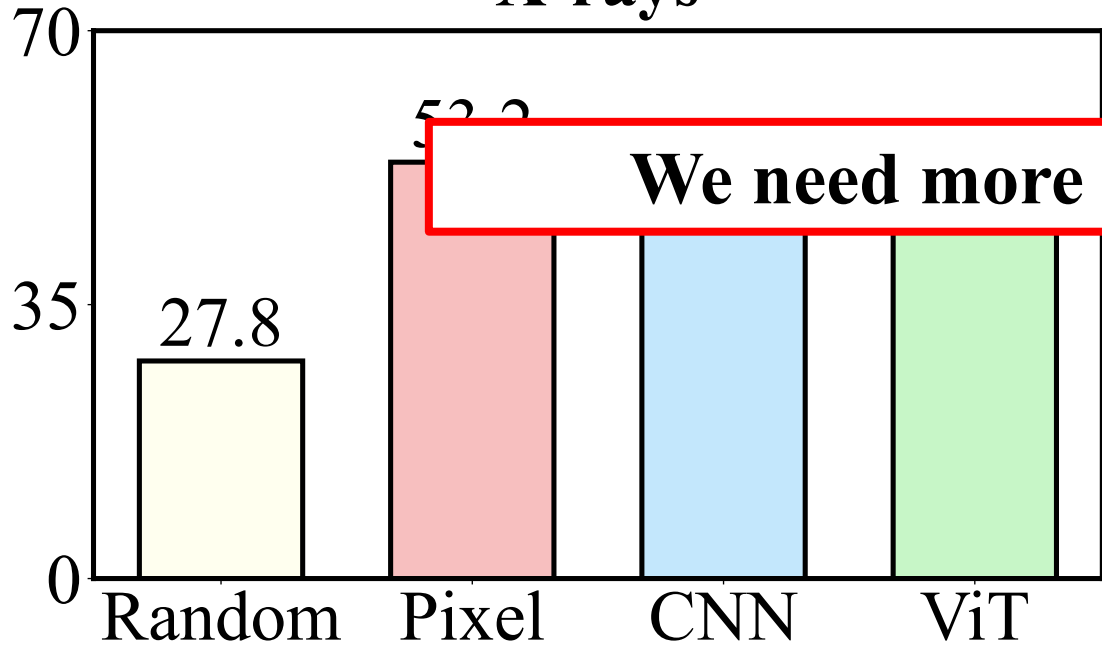
Deep models **don't** have good priors for the **medical domain**.



**X-rays**



**Skin Lesion Images**



**We need more priors in the model.**

X-ray Datasets: Pneumonia, COVID-QU, NIH-CXR, Open-I, VinDr-CXR.

Skin Lesion Datasets: HAM10000, BCN20000, PAD-UFS-20, Melanoma, UWaterloo.



# **KnoBo: Knowledge-enhanced Concept Bottlenecks for Interpretable and Robust Medical Image Classification**

**Yue Yang**, Mona Gandhi, Yufei Wang, Yifan Wu, Michael S. Yao,  
James C. Gee, Mark Yatskar



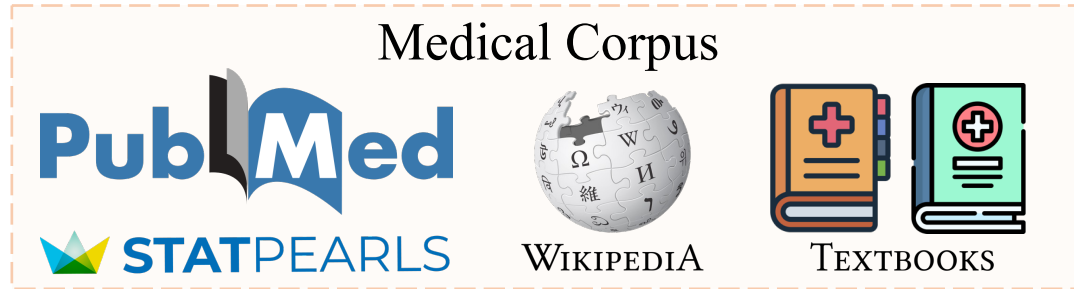
# Where to acquire the prior knowledge?

# Where to acquire the prior knowledge?

Query: How to diagnose COVID from X-rays?

# Where to acquire the prior knowledge?

Query: How to diagnose COVID from X-rays?




# Where to acquire the prior knowledge?


Query: How to diagnose COVID from X-rays?



 5M Articles, 300M+ paragraphs

 9.3K Articles, 301.2K paragraphs

 6.5M Articles,  
30.4M paragraphs

 18 Medical Textbooks,  
125.8k paragraphs

# Where to acquire the prior knowledge?

Query: How to diagnose COVID from X-rays?

retrieve relevant documents

Medical Corpus



STATPEARLS



WIKIPEDIA



TEXTBOOKS



5M Articles, 300M+ paragraphs



9.3K Articles, 301.2K paragraphs



WIKIPEDIA

6.5M Articles,  
30.4M paragraphs



TEXTBOOKS

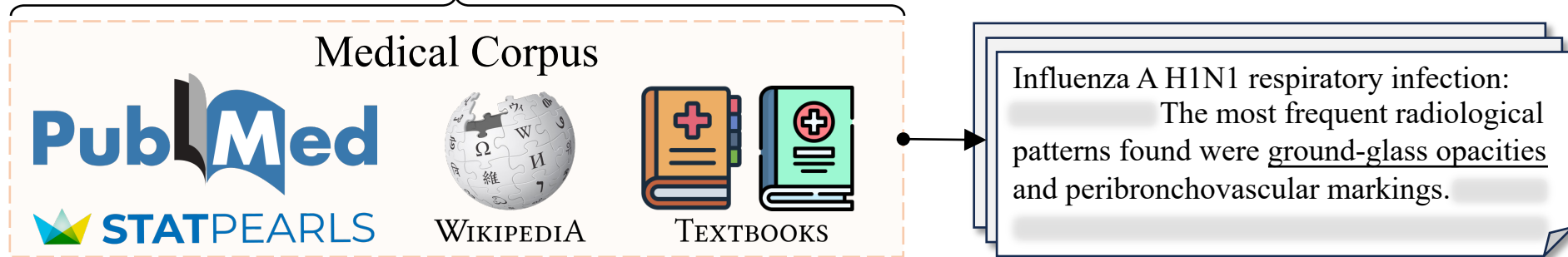
18 Medical Textbooks,  
125.8k paragraphs



# Where to acquire the prior knowledge?

Query: How to diagnose COVID from X-rays?

retrieve relevant documents



5M Articles, 300M+ paragraphs

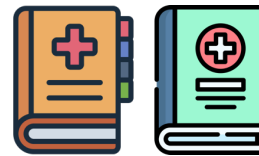


9.3K Articles, 301.2K paragraphs



WIKIPEDIA

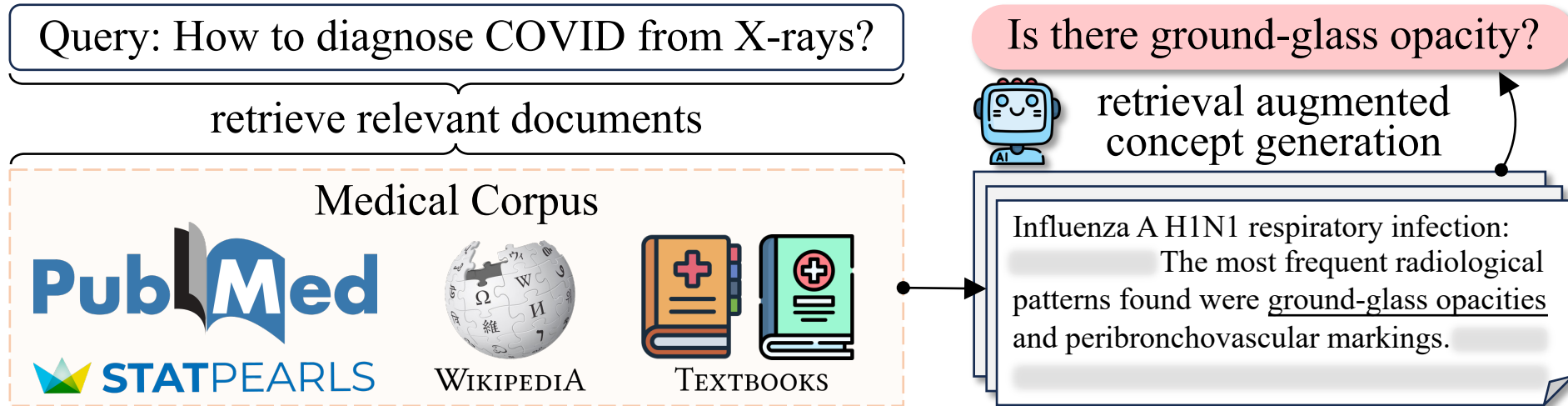
6.5M Articles,  
30.4M paragraphs



TEXTBOOKS


18 Medical Textbooks,  
125.8k paragraphs


# Where to acquire the prior knowledge?



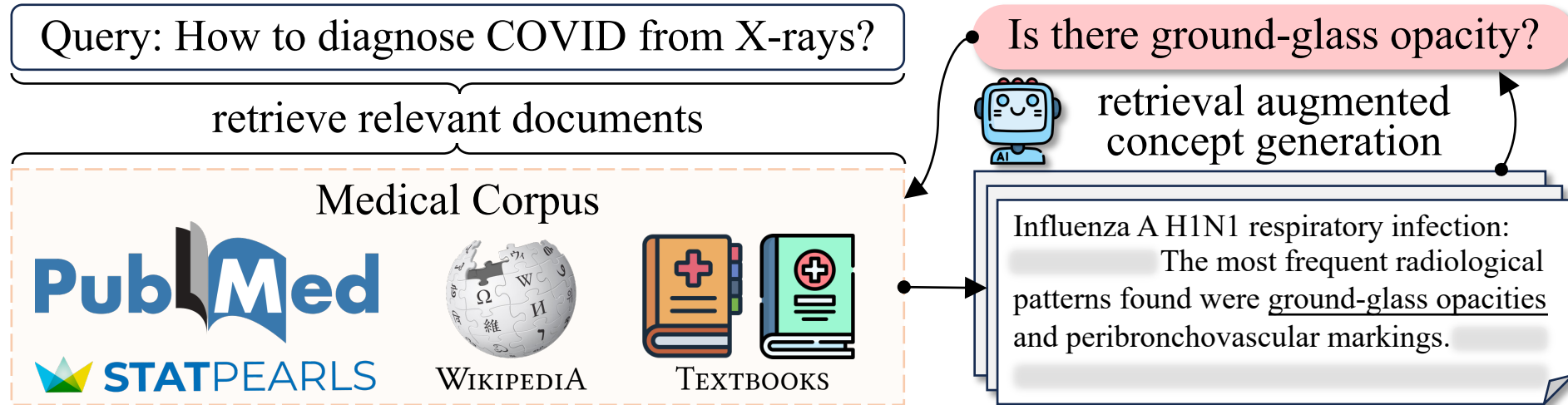
 5M Articles, 300M+ paragraphs

 9.3K Articles, 301.2K paragraphs

 6.5M Articles,  
30.4M paragraphs


 18 Medical Textbooks,  
125.8k paragraphs


# Where to acquire the prior knowledge?



 5M Articles, 300M+ paragraphs

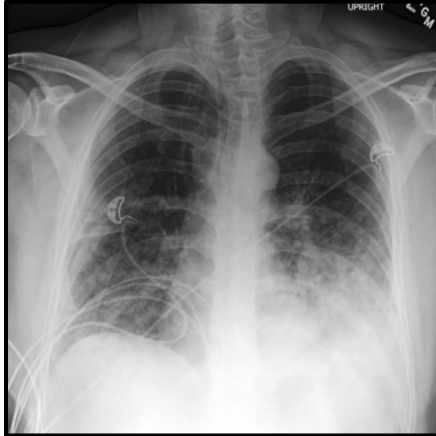
 9.3K Articles, 301.2K paragraphs

 6.5M Articles,  
30.4M paragraphs

 18 Medical Textbooks,  
125.8k paragraphs

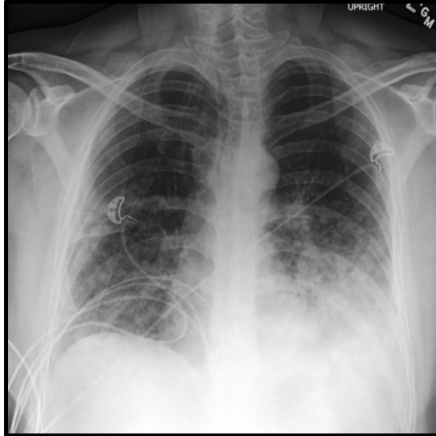
# How to ground the knowledge?

Is there ground-glass opacity?



# How to ground the knowledge?

Is there ground-glass opacity?

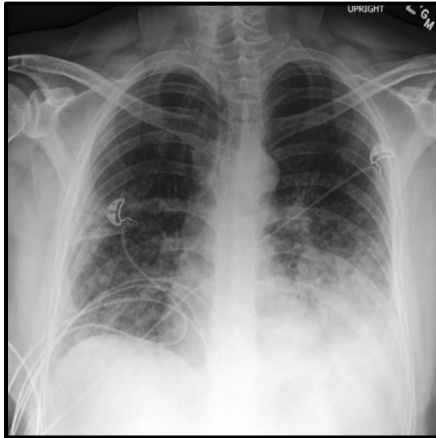


## Paired Clinical Reports

Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.

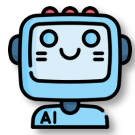
# How to ground the knowledge?

Is there ground-glass opacity?



## Paired Clinical Reports

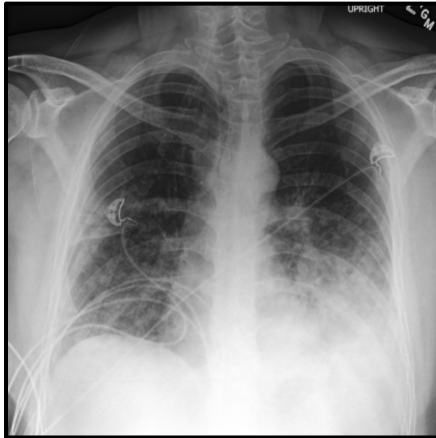
Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.



LLM

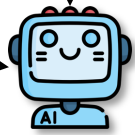
# How to ground the knowledge?

Is there ground-glass opacity?



## Paired Clinical Reports

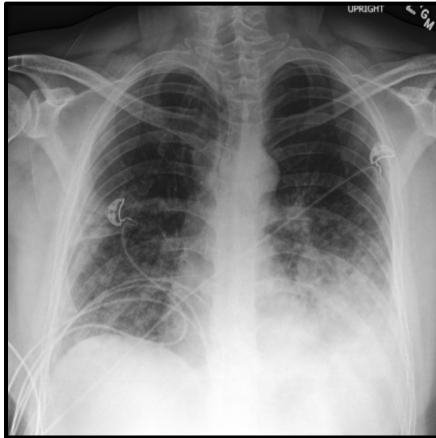
Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.



LLM

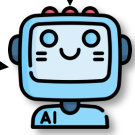
# How to ground the knowledge?

Is there ground-glass opacity?



## Paired Clinical Reports

Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.



LLM

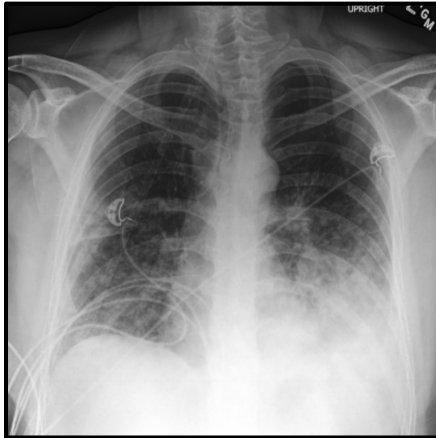
Yes



# How to ground the knowledge?

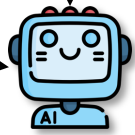
Learn a binary classifier for each concept.

Is there ground-glass opacity?



## Paired Clinical Reports

Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.



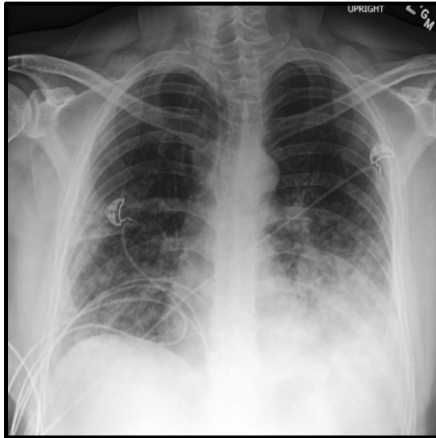
LLM

Yes

Is there ground-glass opacity?

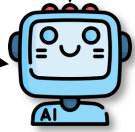
# How to ground the knowledge?

Is there ground-glass opacity?



## Paired Clinical Reports

Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.



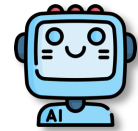
LLM

Yes

Learn a binary classifier for each concept.

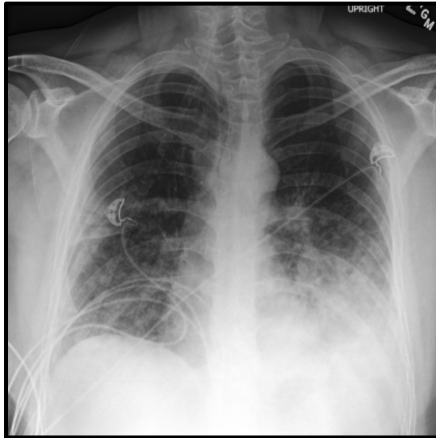
Is there ground-glass opacity?

Clinical Report



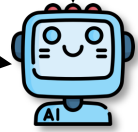
# How to ground the knowledge?

Is there ground-glass opacity?



## Paired Clinical Reports

Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.

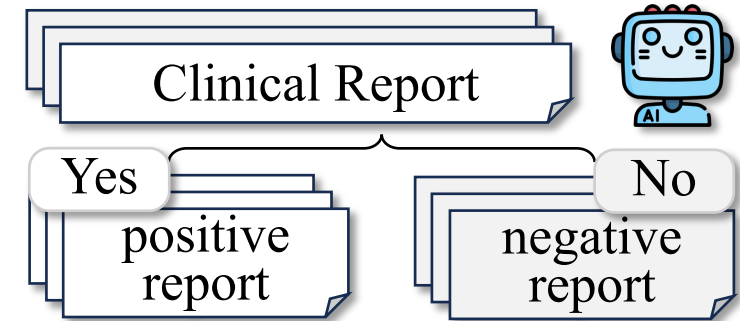


LLM

Yes

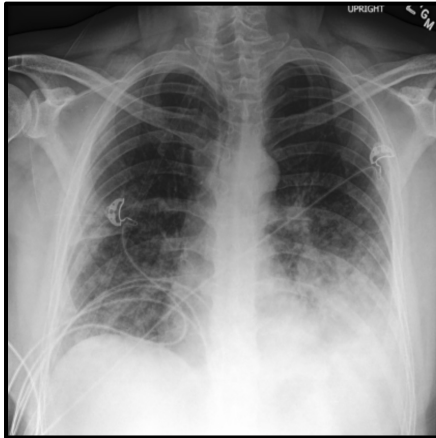
Learn a binary classifier for each concept.

Is there ground-glass opacity?



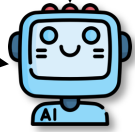
# How to ground the knowledge?

Is there ground-glass opacity?



## Paired Clinical Reports

Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.



LLM

Yes

Learn a binary classifier for each concept.

Is there ground-glass opacity?

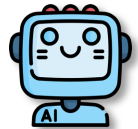


Yes

positive report



positive images



No

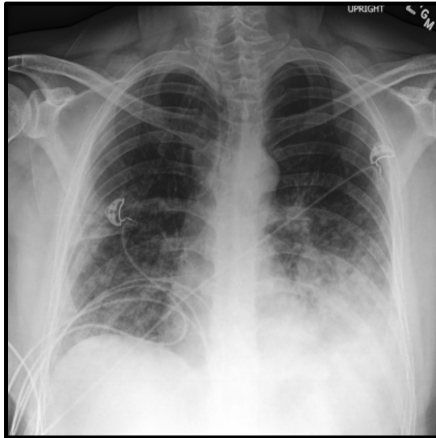
negative report



negative images

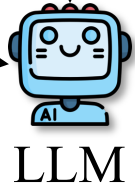
# How to ground the knowledge?

Is there ground-glass opacity?



## Paired Clinical Reports

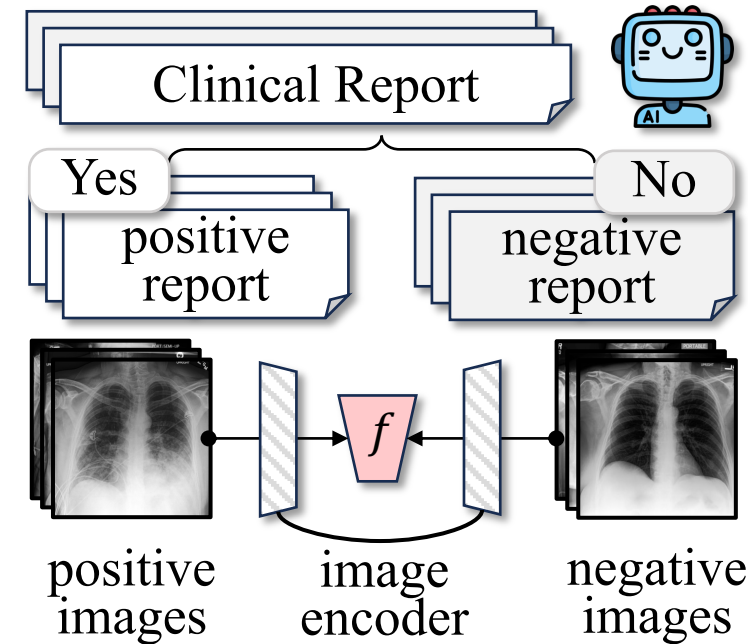
Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.



Yes

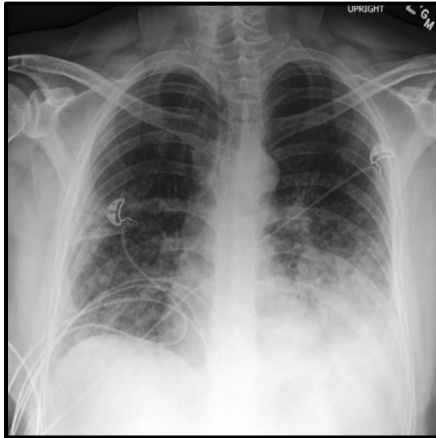
Learn a binary classifier for each concept.

Is there ground-glass opacity?



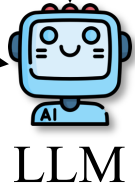
# How to ground the knowledge?

Is there ground-glass opacity?



## Paired Clinical Reports

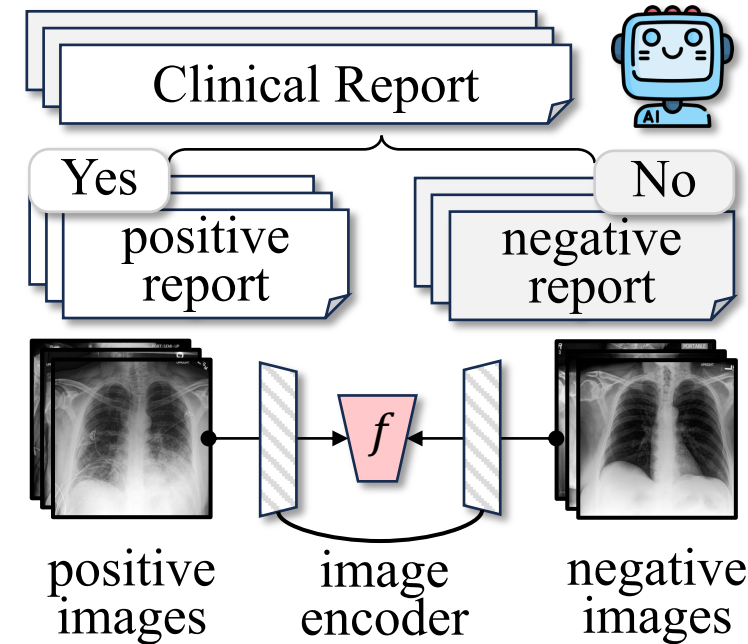
Redemonstration of subtle posterior lung base densities corresponding to **ground-glass opacities** on prior CT and likely representing aspiration do not appear worsened. Tiny bilateral pleural effusions.



Yes

Learn a binary classifier for each concept.

Is there ground-glass opacity?



Convert text instances into image instances.

# Grounding Comparison

Query: Are lung fields clear on both sides?

**CLIP**

**Ours (w/ knowledge grounding)**

# Grounding Comparison

Query: Are lung fields clear on both sides?

**CLIP**

**Ours (w/ knowledge grounding)**

Top-3 clear



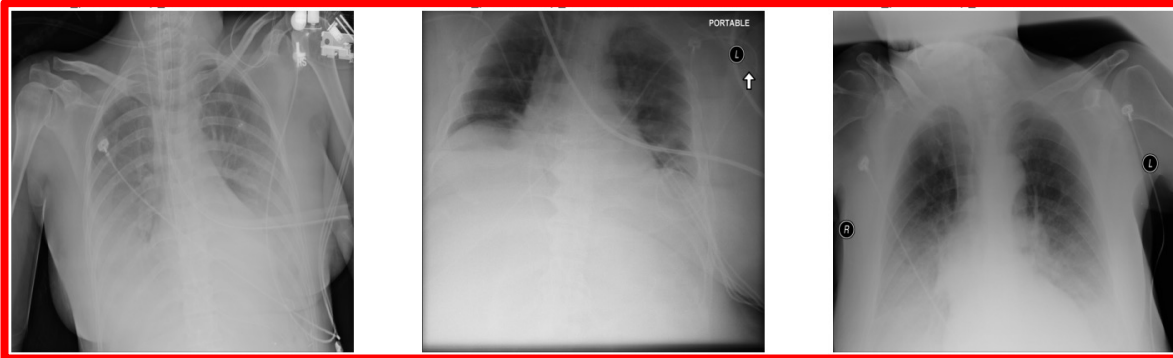


# Grounding Comparison

Query: Are lung fields clear on both sides?

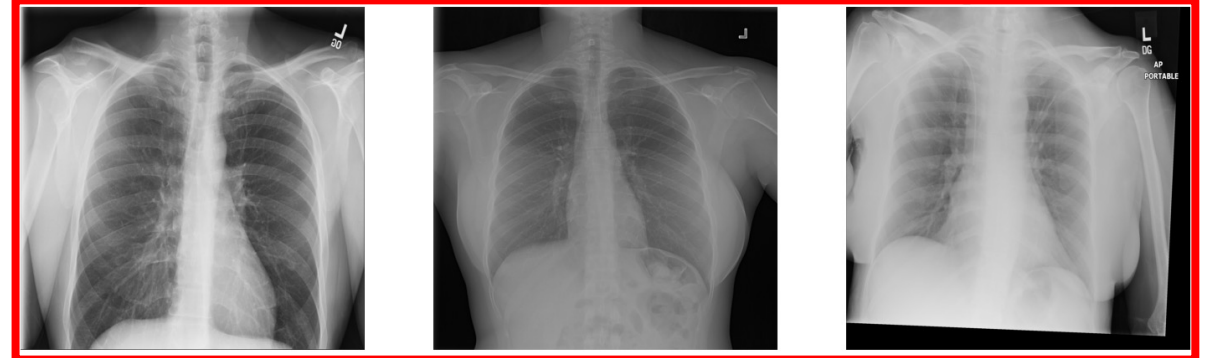
**CLIP**

Top-3 clear



**Ours (w/ knowledge grounding)**

Top-3 clear

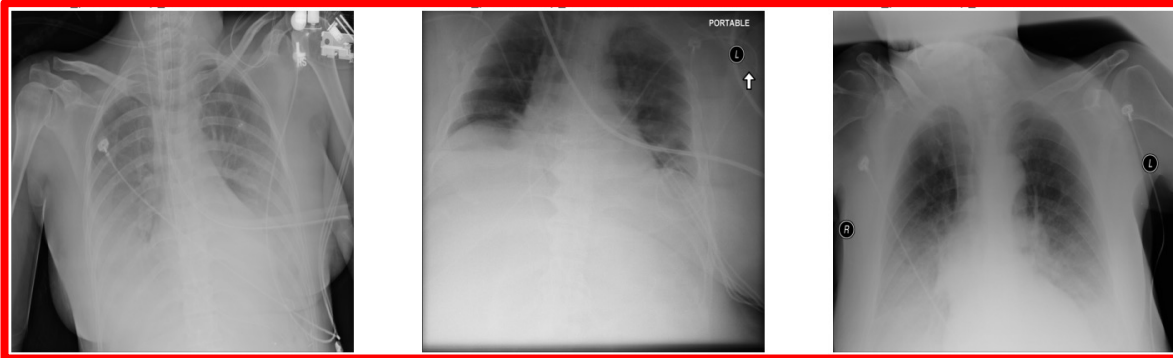


# Grounding Comparison

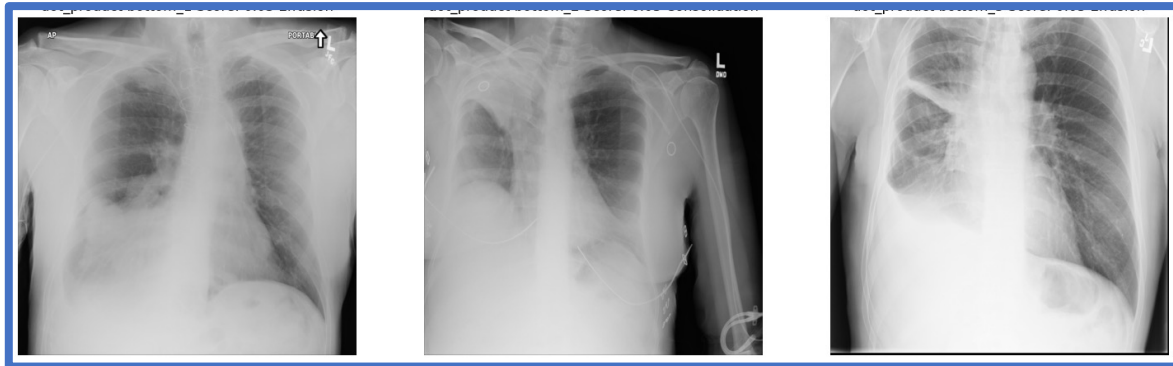
Query: Are lung fields clear on both sides?

**CLIP**

Top-3 clear

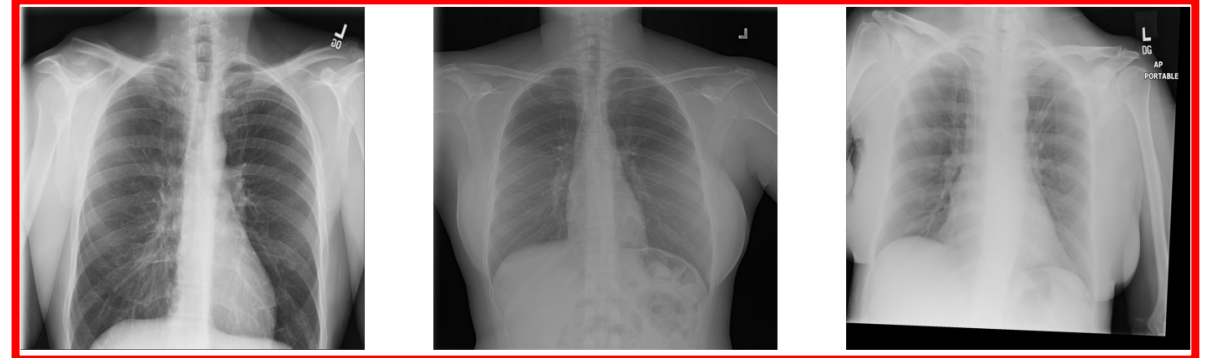


Bottom-3 not clear



**Ours (w/ knowledge grounding)**

Top-3 clear

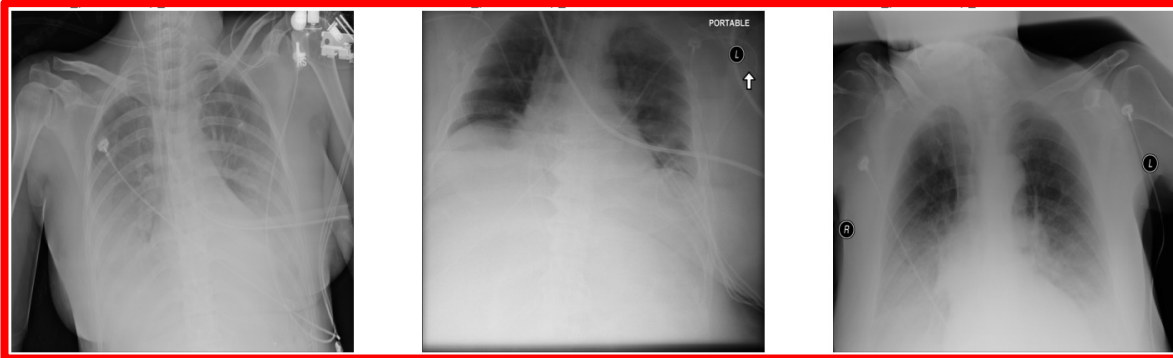


# Grounding Comparison

Query: Are lung fields clear on both sides?

**CLIP**

Top-3 clear

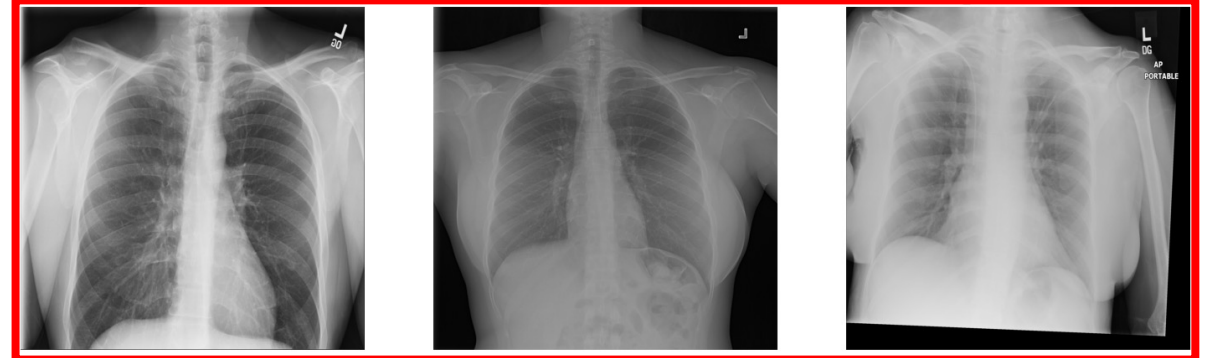


Bottom-3 not clear



**Ours (w/ knowledge grounding)**

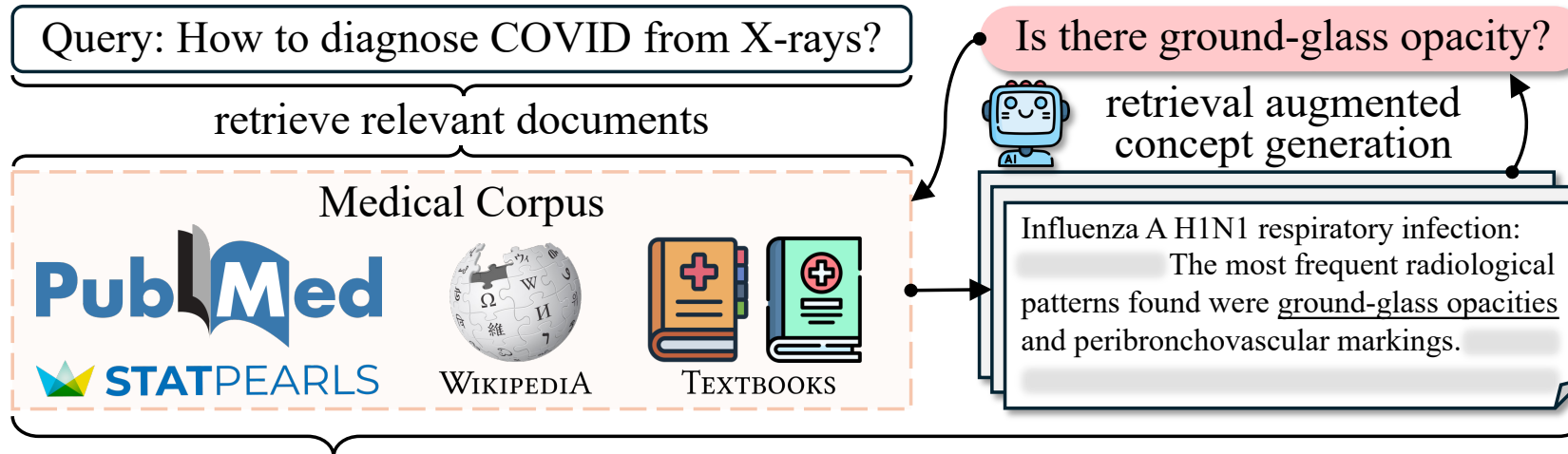
Top-3 clear



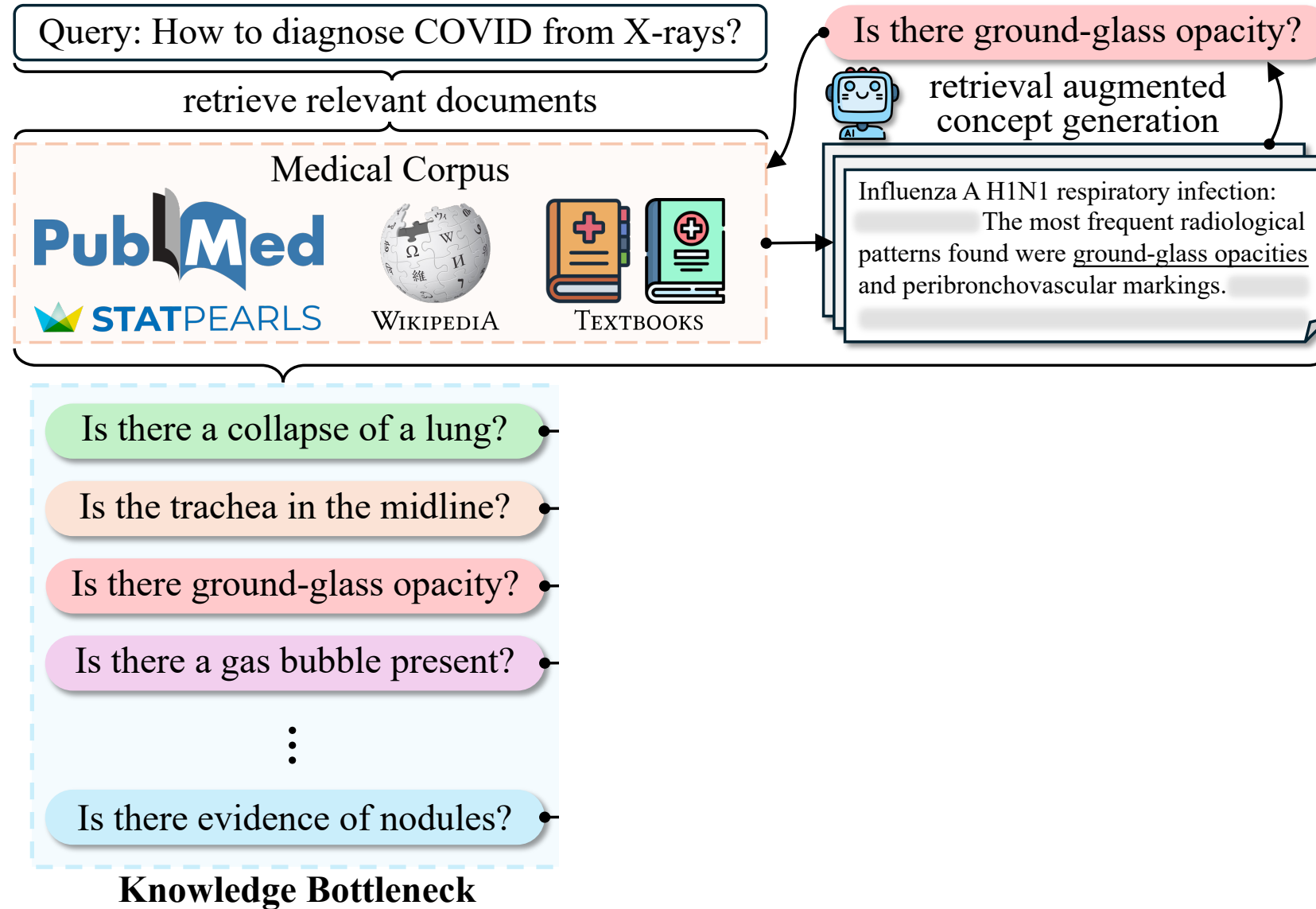
Bottom-3 not clear



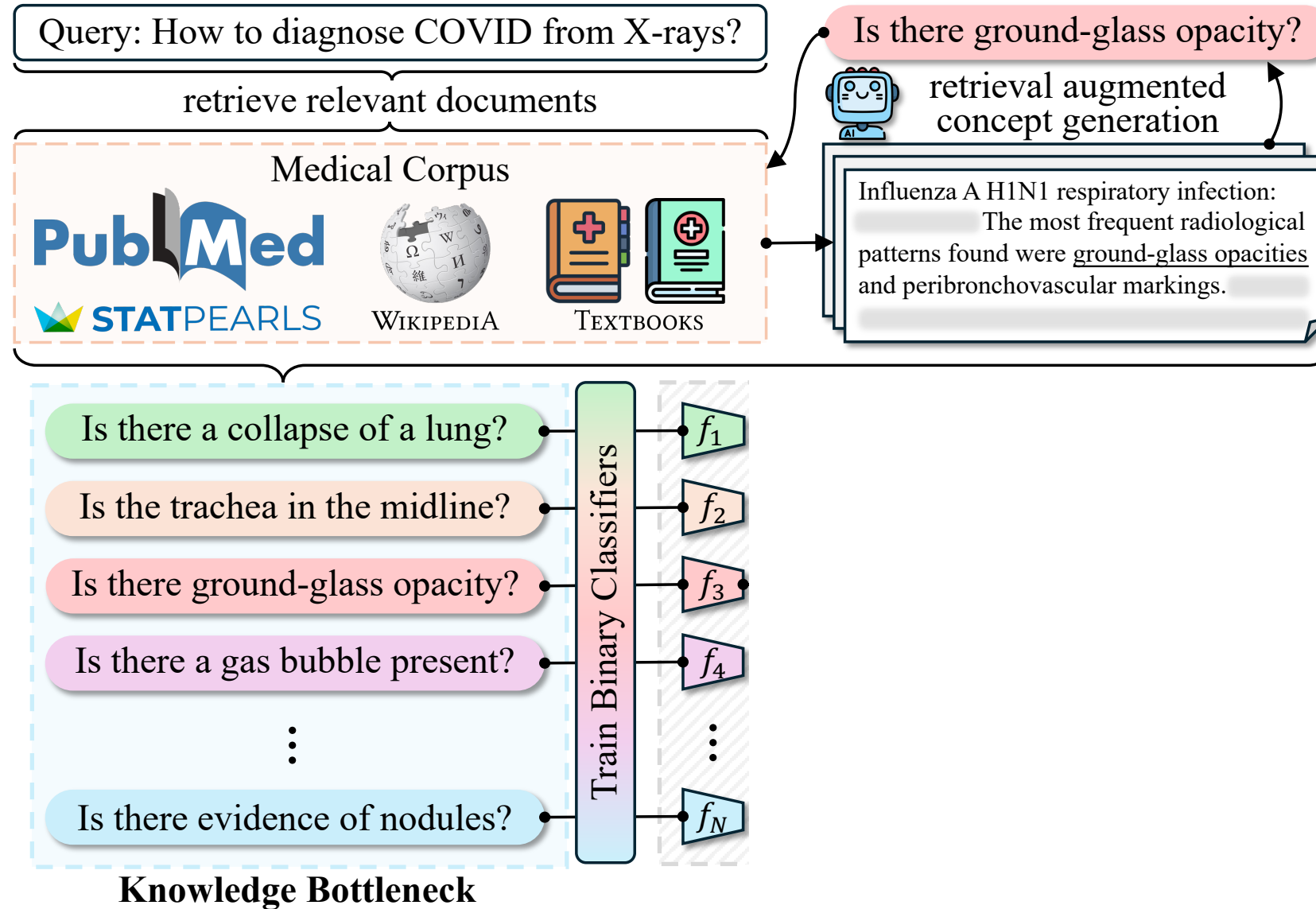
# KnoBo System Overview



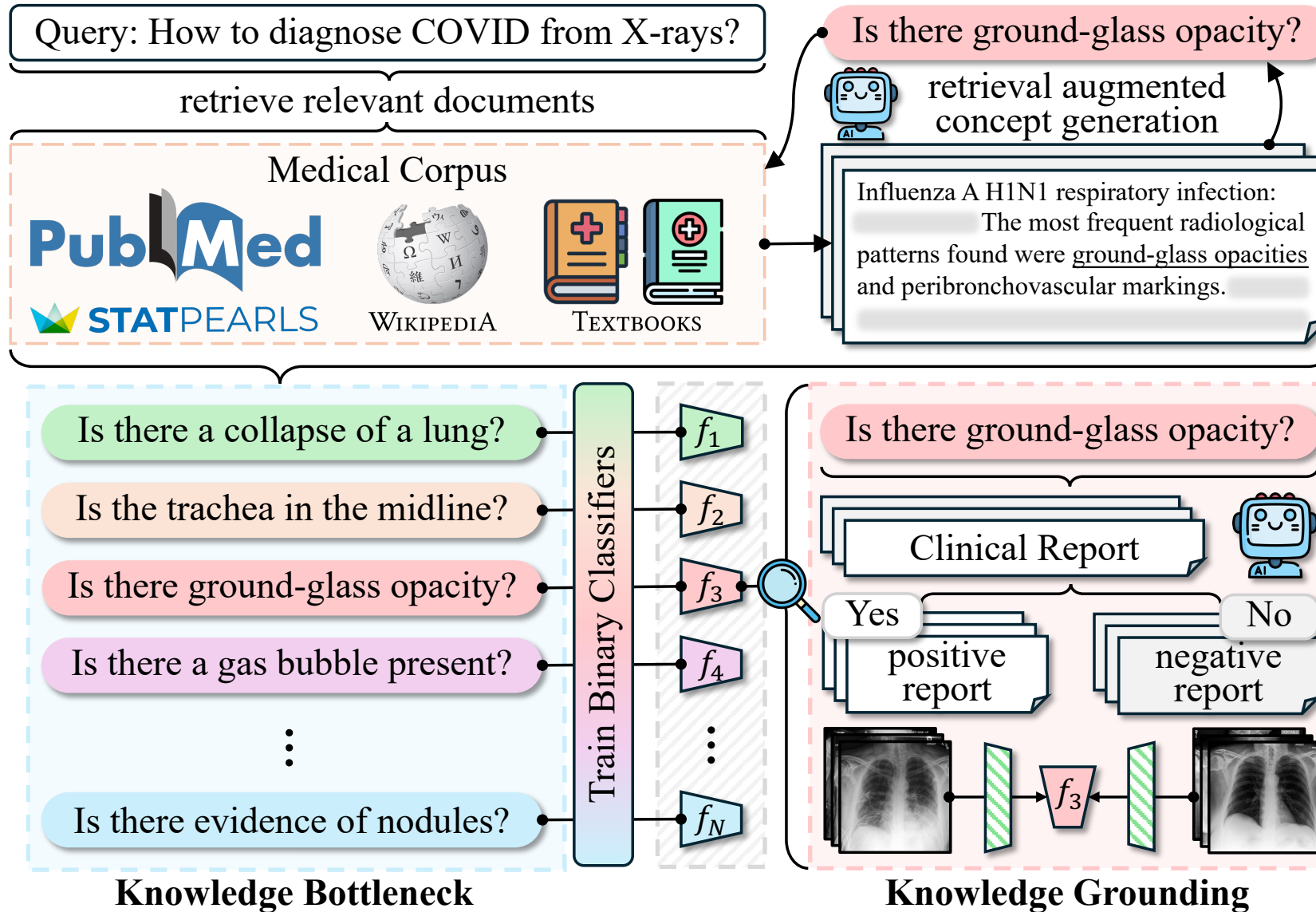
# KnoBo System Overview



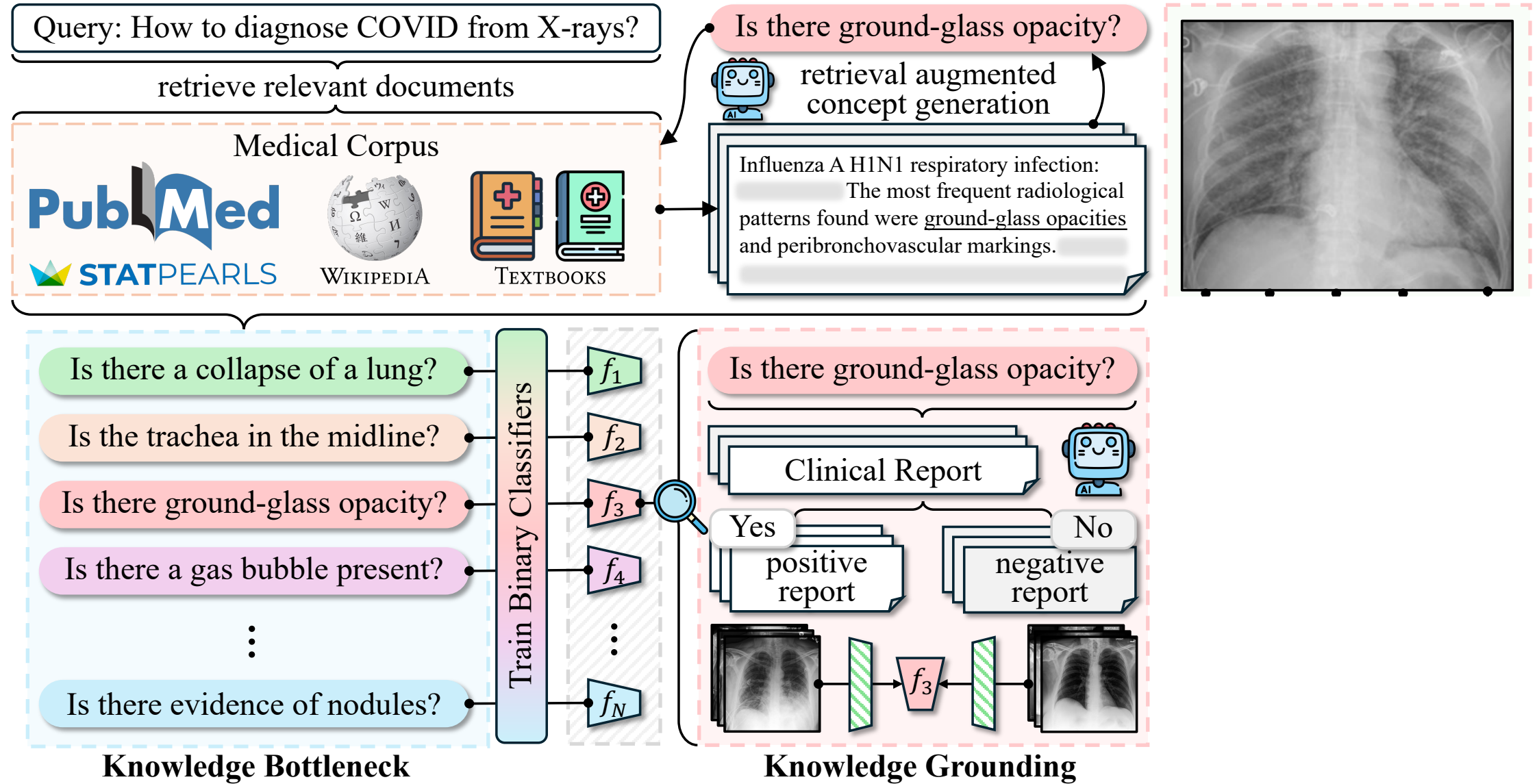
# KnoBo System Overview



# KnoBo System Overview

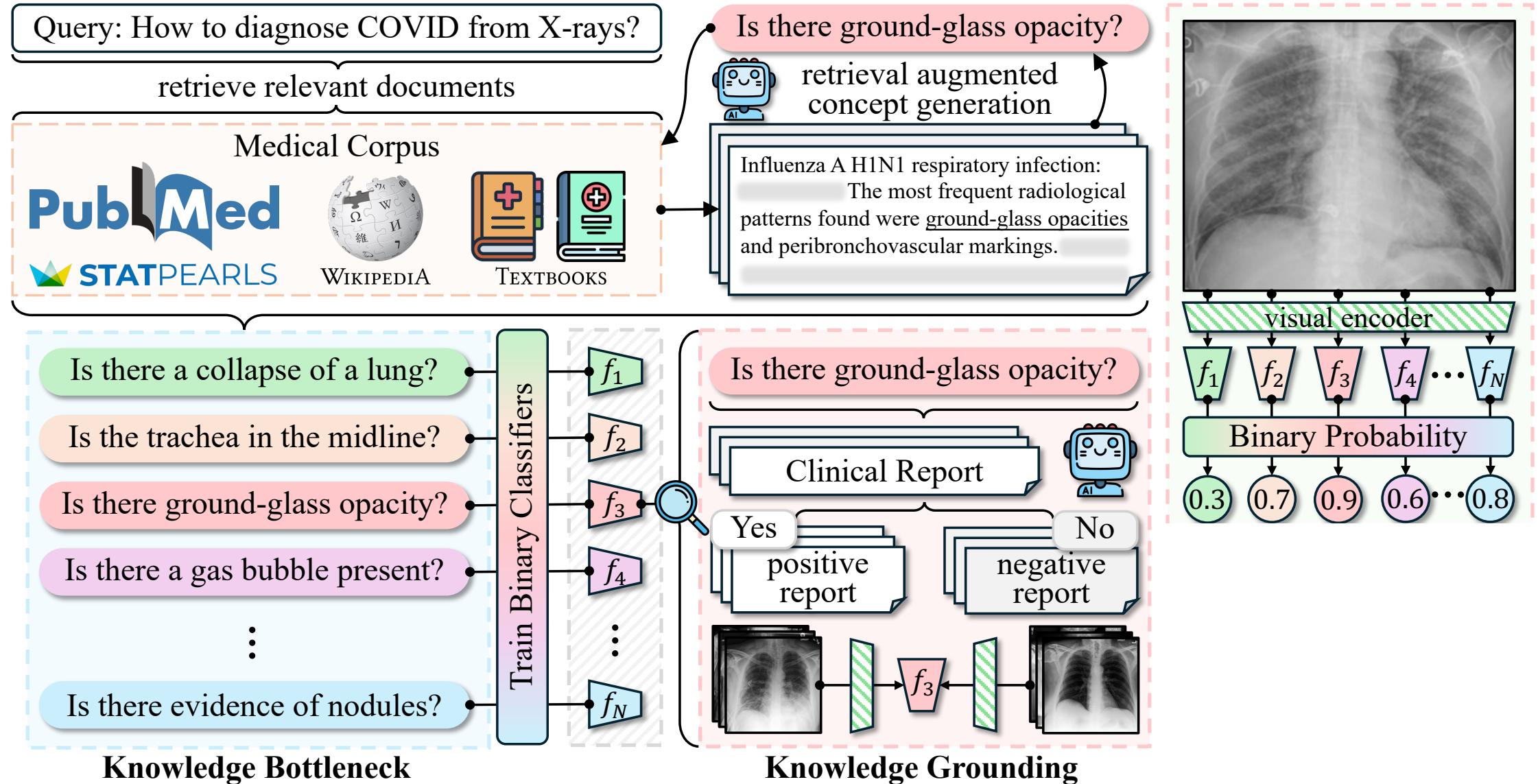


# KnoBo System Overview

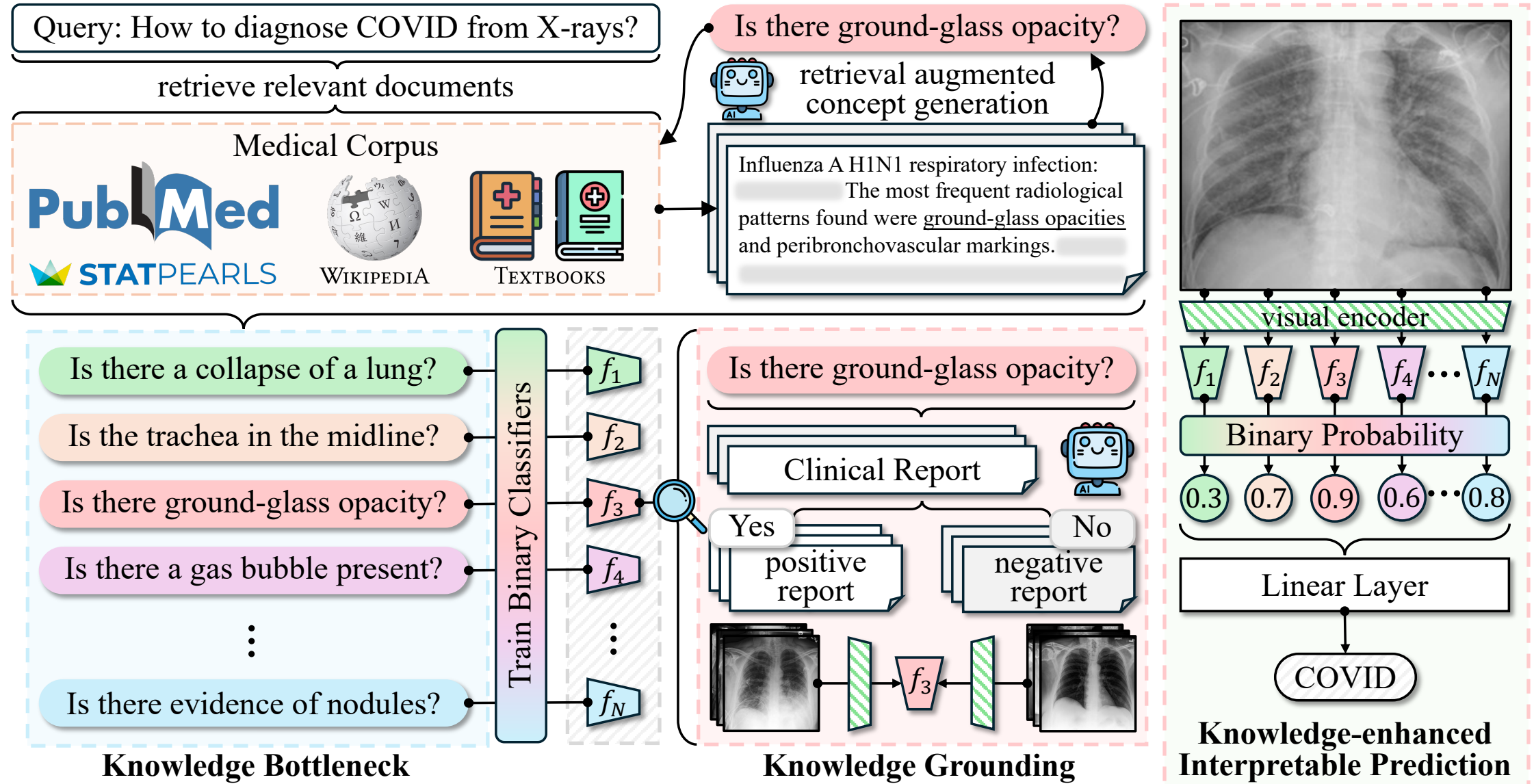




# KnoBo System Overview

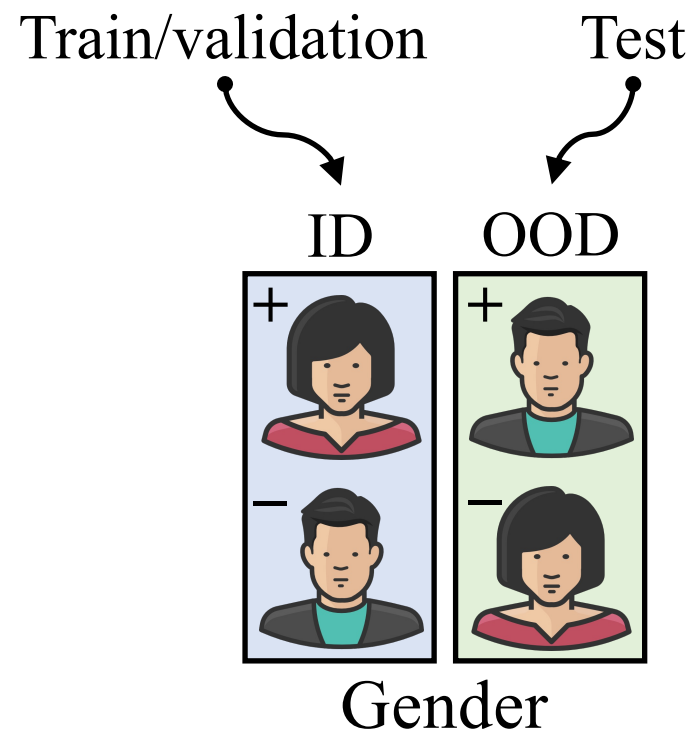
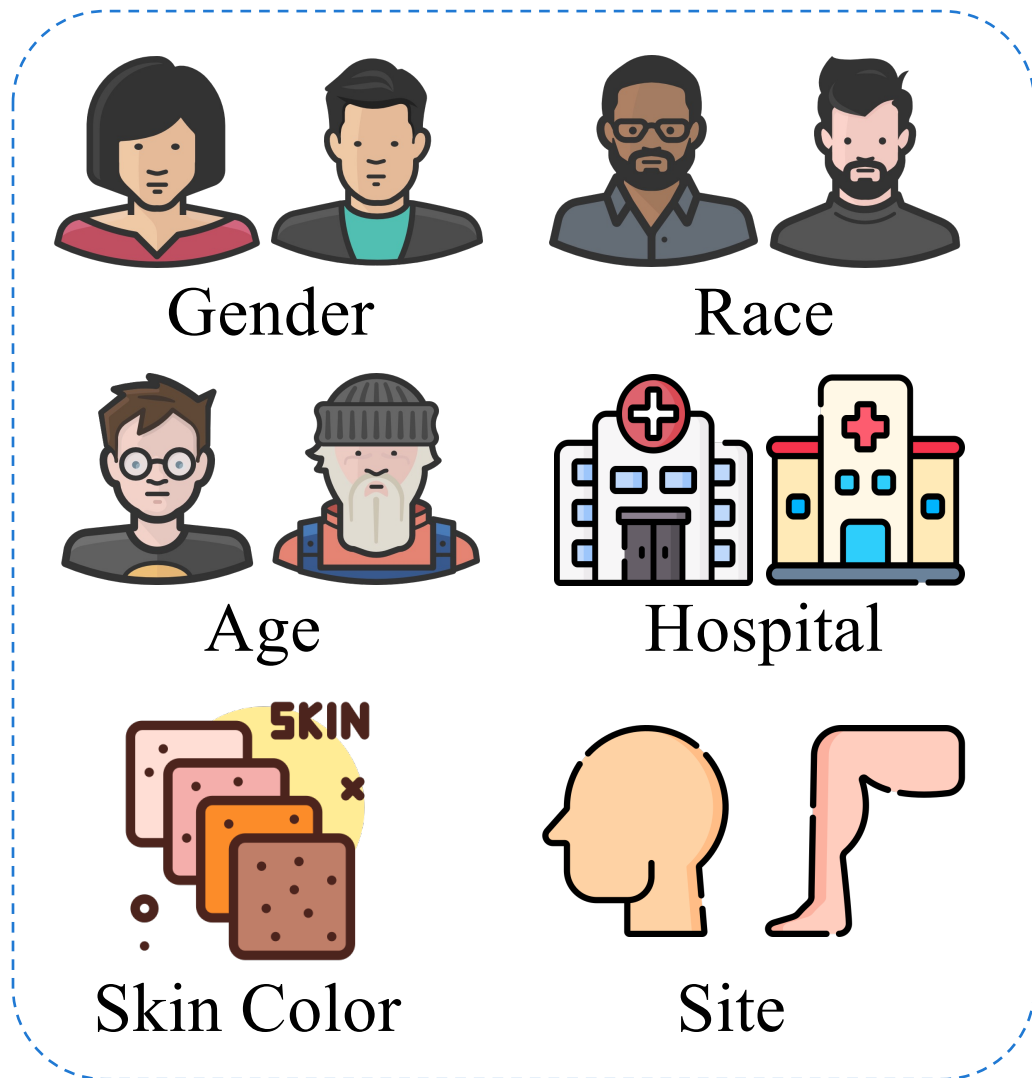


# KnoBo System Overview



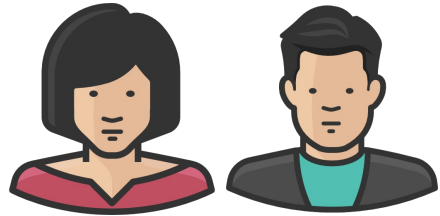
# Datasets

## Confounded

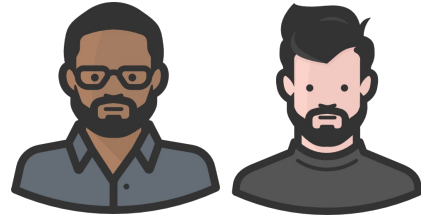


# Datasets

## Confounded



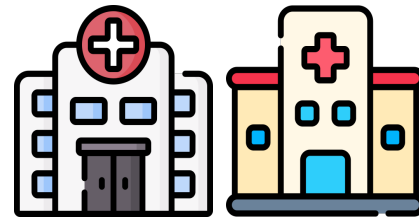
Gender



Race



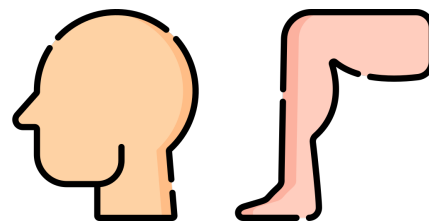
Age



Hospital



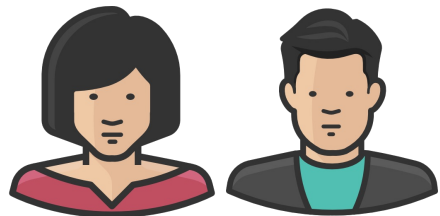
Skin Color



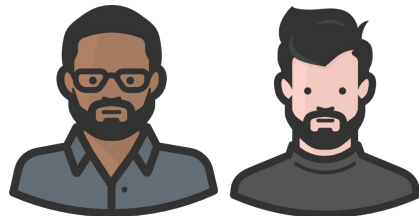
Site

# Datasets

## Confounded



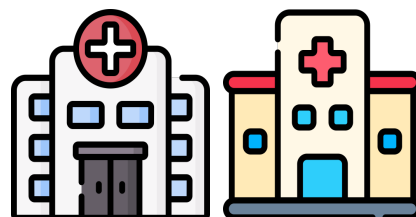
Gender



Race



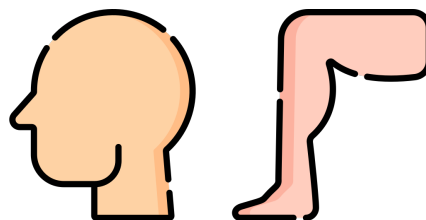
Age



Hospital



Skin Color



Site

## Standard

**X-ray:** Pneumonia, COVID-QU, NIH-CXR, Open-I, VinDr-CXR.

**Skin Lesion:** HAM10000, BCN20000, PAD-UFS-20, Melanoma, UWaterloo.

# Experimental Setup

- **Baselines** (same vision backbone):
  - **Linear Probe**: logistic regression on the image features.
  - **End-to-end**: Unfreeze the visual encoder and update all parameters.
  - **LaBo**: knowledge priors from LLM, no knowledge grounding.
- **Metric**:
  - **Confounded datasets**: ID (validation), OOD (test), delta  $\downarrow$  ( $|\text{OOD-ID}|$ ), and domain-average accuracy ( $\text{ID} + \text{OOD} / 2$ ).
  - **Standard datasets**: test accuracy.
  - **Overall Performance**: average over confounded and standard datasets.

# Results on X-ray Datasets



Method	ID	OOD	delta↓	Domain Average	Standard	Overall
Linear Probe	<u>95.2</u>	30.7	64.5	62.9	<b>73.8</b>	<u>68.4</u>
End-to-End	<b>96.7</b>	17.0	79.7	56.8	70.2	63.5
LaBo	93.5	<u>34.8</u>	<u>58.7</u>	<u>64.2</u>	72.1	68.1
KnoBo	89.7	<b>58.8</b>	<b>30.9</b>	<b>74.3</b>	<u>73.1</u>	<b>73.7</b>

The best score is **bold** and the second best is underlined.

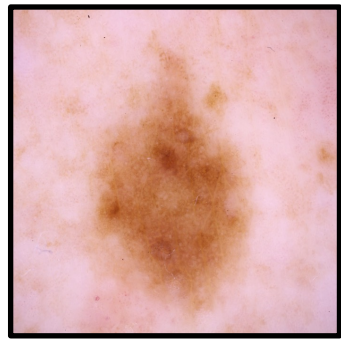
# Results on Skin Lesion Datasets



Method	ID	OOD	delta↓	Domain Average	Standard	Overall
Linear Probe	91.9	<u>52.1</u>	39.8	<u>72.0</u>	<u>82.8</u>	77.4
End-to-End	<b>95.6</b>	47.6	48.0	71.6	<b>84.3</b>	<u>77.9</u>
LaBo	89.9	51.4	<u>38.4</u>	70.6	80.0	75.3
KnoBo	86.0	<b>70.5</b>	<b>14.1</b>	<b>78.3</b>	78.1	<b>78.2</b>



# Results on Skin Lesion Datasets



Method	ID	OOD	delta↓	Domain Average	Standard	Overall
Linear Probe	91.9	<u>52.1</u>	39.8	<u>72.0</u>	<u>82.8</u>	77.4
End-to-End	<b>95.6</b>	47.6	48.0	71.6	<b>84.3</b>	<u>77.9</u>
LaBo	89.9	51.4	<u>38.4</u>	70.6	80.0	75.3
KnoBo	86.0	<b>70.5</b>	<b>14.1</b>	<b>78.3</b>	78.1	<b>78.2</b>

**KnoBo is more robust on confounded datasets.  
KnoBo is competitive on standard datasets.**

# Comparison on Knowledge Types

Knowledge	X-ray Datasets			Skin Lesion Datasets		
	Confounded	Standard	Overall	Confounded	Standard	Overall
Prompt	72.9	72.8	<u>72.9</u>	<b>79.3</b>	72.8	76.0
Textbooks	72.0	<u>72.9</u>	72.4	<u>79.2</u>	76.4	77.8
Wikipedia	72.8	72.7	72.8	<b>79.3</b>	76.2	77.8
StatPearls	<u>73.4</u>	72.0	72.7	<u>79.2</u>	<b>77.6</b>	<b>78.4</b>
PubMed	<b>74.3</b>	<b>73.1</b>	<b>73.7</b>	<b>79.3</b>	<u>76.7</u>	<u>78.0</u>

# Conclusion



Interpretable models with knowledge priors  
are **more robust in medical domains.**

## Future Work

- Better feature representations for critical domains.
- Different structures of knowledge.
- Other usages of interpretable models.

Thank you!