



# Causally motivated robustness to shortcut learning

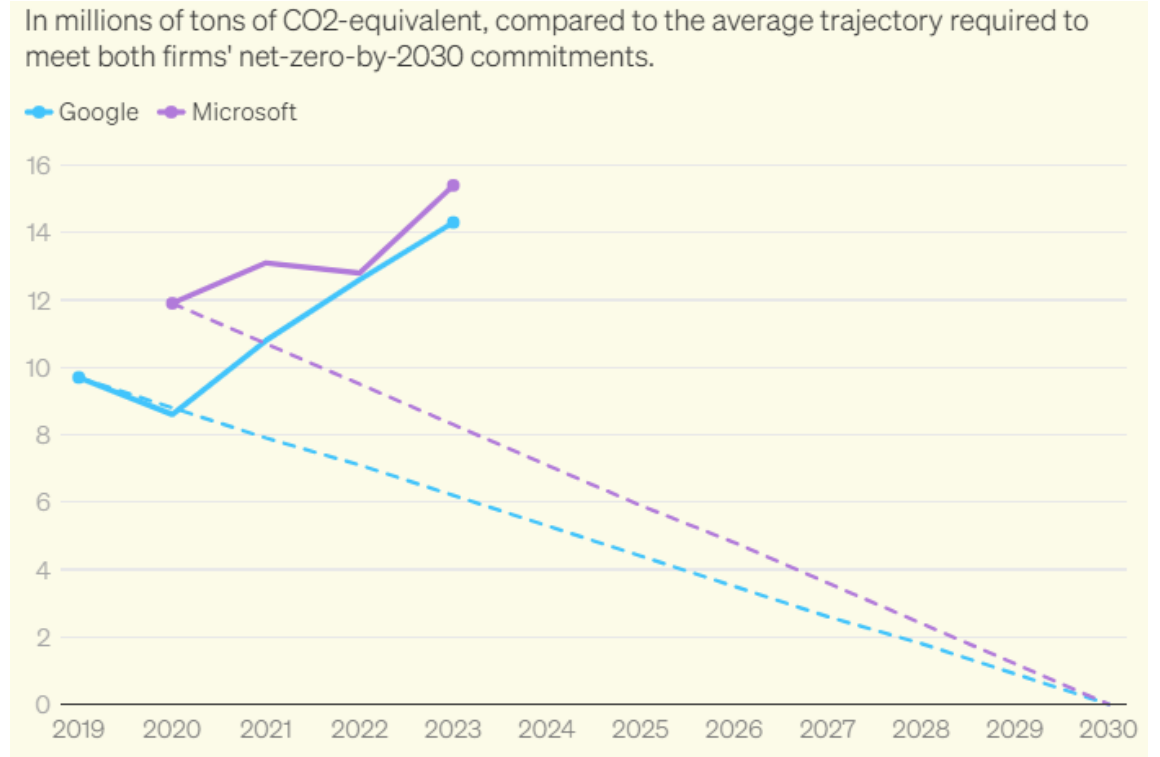
---

Maggie Makar,  
Assistant Professor  
CSE, University of Michigan

# Bigger is better?



# Bigger is better?



# Bigger is better?

**Task:** predict movie review sentiment

<b>Original</b>	
Incredible performances, must watch	Positive
Emotional rollercoaster, highly recommend	Positive
<b>Shakespearean</b>	
A wretched script, a squander of precious hours!	Negative
A dismal affair; I wouldst not commend it to any soul.	Negative

# Bigger is better?

**Task:** predict movie review sentiment

<b>Original</b>	
Incredible performances, must watch	Positive
Emotional rollercoaster, highly recommend	Positive
<b>Shakespearean</b>	
A wretched script, a squander of precious hours!	Negative
A dismal affair; I wouldst not commend it to any soul.	Negative

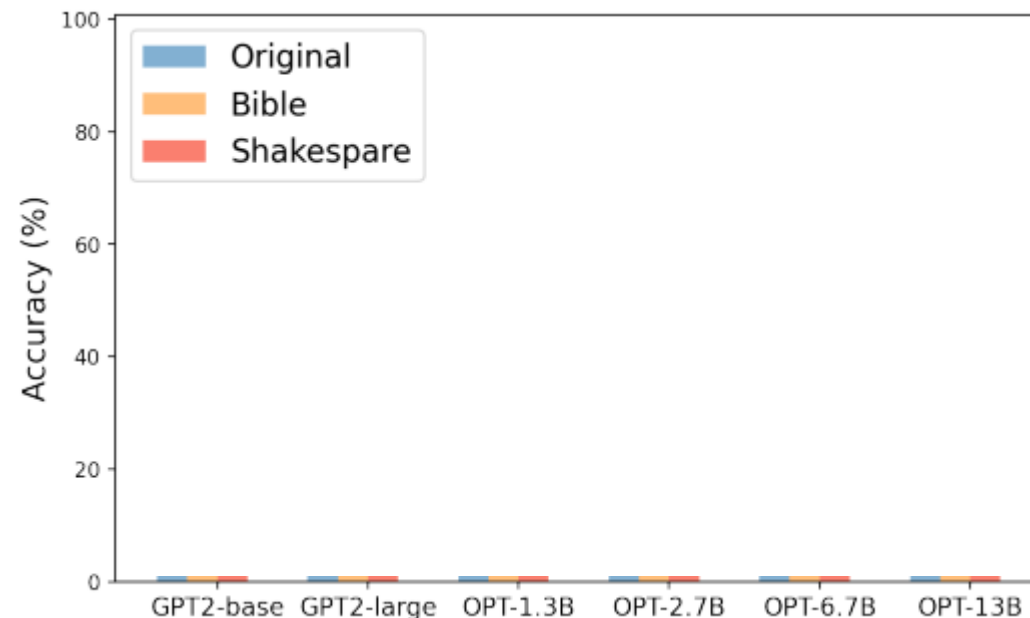
Exploitable shortcut: Style  $\rightarrow$  sentiment

# Bigger is better?

**Task:** predict movie review sentiment

Original	
Incredible performances, must watch	Positive
Emotional rollercoaster, highly recommend	Positive
Shakespearean	
A wretched script, a squander of precious hours!	Negative
A dismal affair; I wouldst not commend it to any soul.	Negative

Exploitable shortcut: Style  $\rightarrow$  sentiment

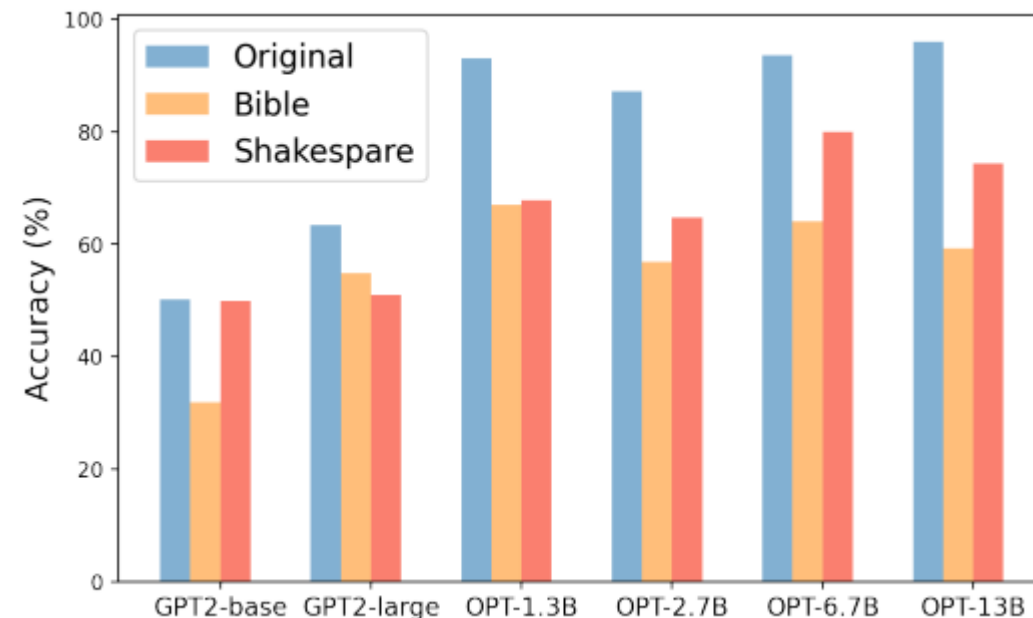


# Bigger is better?

Task: predict movie review sentiment

Original	
Incredible performances, must watch	Positive
Emotional rollercoaster, highly recommend	Positive
Shakespearean	
A wretched script, a squander of precious hours!	Negative
A dismal affair; I wouldst not commend it to any soul.	Negative

Exploitable shortcut: Style  $\rightarrow$  sentiment

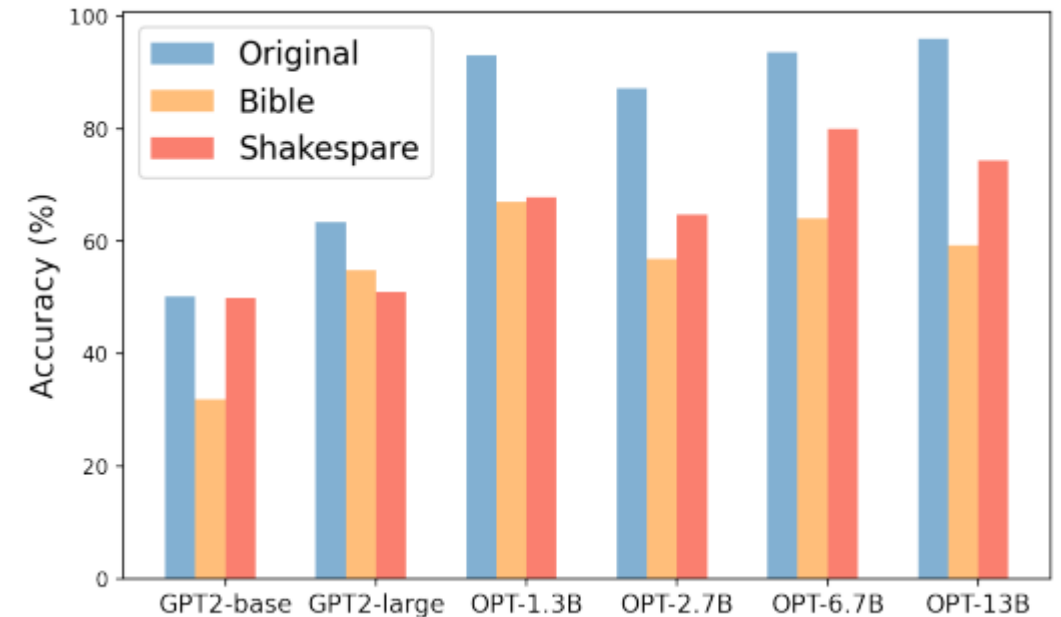


# Bigger is better?

Task: predict movie review sentiment

Original	
Incredible performances, must watch	Positive
Emotional rollercoaster, highly recommend	Positive
Shakespearean	
A wretched script, a squander of precious hours!	Negative
A dismal affair; I wouldst not commend it to any soul.	Negative

Exploitable shortcut: Style  $\rightarrow$  sentiment



Building larger models does not give us robustness to shortcuts



# Causally motivated regularization

Good old efficiency

Causally motivated regularization



Good old efficiency

Causally motivated regularization

Robustness

Good old efficiency

# Causally motivated regularization

Robustness



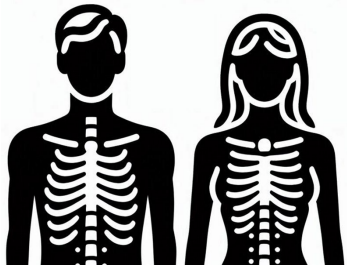
---

Causality reasons about the world under shifts/interventions

---

# Talk outline

# Talk outline

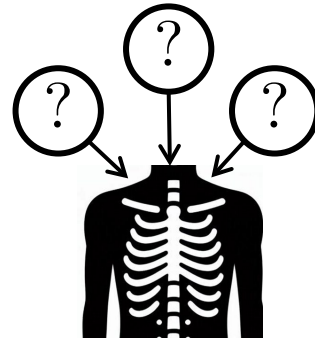


1 Efficiency + robustness to known sampling bias

MPMBHD, AISTats 22

MD, TMLR 23

NM, UAI 24



2 Efficiency + robustness to unknown sampling biases

ZM, NeurIPS 22

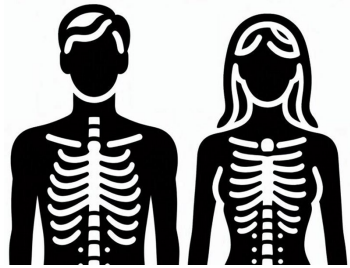
WJMSW, NeurIPS 22



3 Evaluating localized circuits in LLMs

SVNZGJMB – NeurIPS 24

# Talk outline

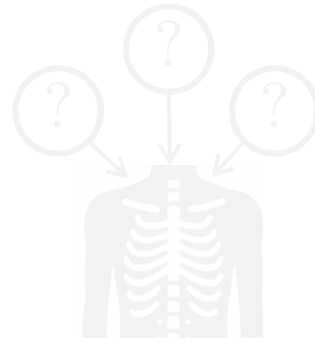


1 Efficiency + robustness to known sampling bias

MPMBHD, AISTats 22

MD, TMLR 23

NM, UAI 24



2 Efficiency + robustness to unknown sampling biases

ZM, NeurIPS 22

WJMSW, NeurIPS 22

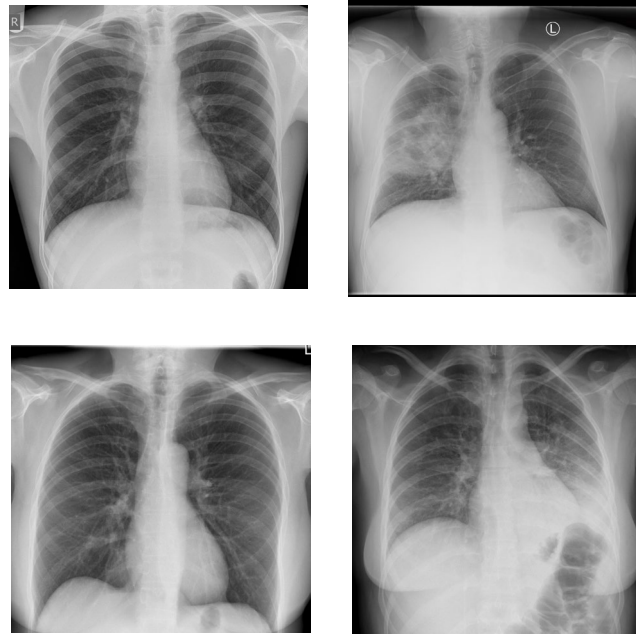


3 Evaluating localized circuits in LLMs

SVNZGJMB – NeurIPS 24

# Pneumonia detection

## Training data



Makar, Packer, Moldovan, Blalock, Halpern and D'Amour AISTats 2022

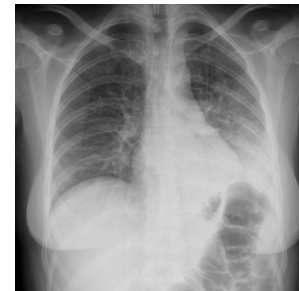
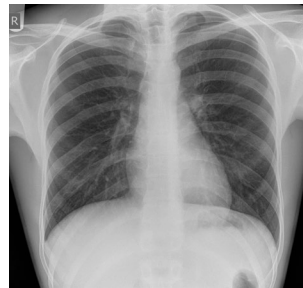
Cases courtesy of Dr. Andrew Dixon, Dr Henry Knip, Dr. Usman Bashir, and Dr. Ian Bickle , Radiopaedia.org, rID: 48366, 31388, 18394 and rID: 50318



# Pneumonia detection

## Training data

Target label:      Healthy      Pneumonia  
 $Y = 0$                        $Y = 1$



Makar, Packer, Moldovan, Blalock, Halpern and D'Amour AISTATS 2022

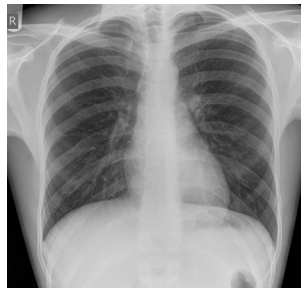
Cases courtesy of Dr. Andrew Dixon, Dr Henry Knip, Dr. Usman Bashir, and Dr. Ian Bickle, Radiopaedia.org, rID: 48366, 31388, 18394 and rID: 50318

# Pneumonia detection

## Training data

Target label:      Healthy      Pneumonia  
 $Y = 0$                        $Y = 1$

Auxiliary label:      Men  
 $V = 0$



Women  
 $V = 1$

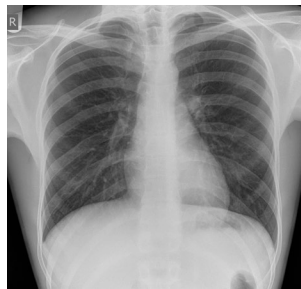


# Pneumonia detection under biased sampling

## Training data

Target label:      Healthy      Pneumonia  
 $Y = 0$                        $Y = 1$

Auxiliary label:      Men  
 $V = 0$



Women  
 $V = 1$



# Pneumonia detection under biased sampling

## Training data

Target label:      Healthy      Pneumonia  
 $Y = 0$                        $Y = 1$

Auxiliary label:      Men  
 $V = 0$



Women  
 $V = 1$



# Pneumonia detection under biased sampling

## Training data

Target label:      Healthy      Pneumonia  
 $Y = 0$                        $Y = 1$

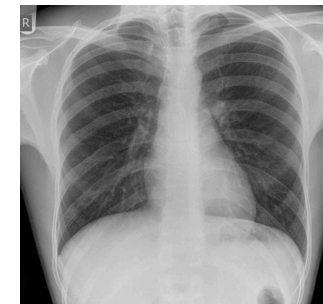
Auxiliary label:      Men  
 $V = 0$

Women  
 $V = 1$



## Testing data

Man with  
pneumonia



$P(\text{Pneumonia}) = 0.99$

# Pneumonia detection under biased sampling

## Training data

Target label:      Healthy      Pneumonia  
 $Y = 0$                        $Y = 1$

Auxiliary label:      Men  
 $V = 0$

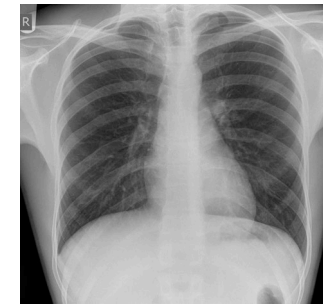


Women  
 $V = 1$



## Testing data

Man with  
pneumonia



Woman with  
pneumonia



$P(\text{Pneumonia}) = 0.99$      $P(\text{Pneumonia}) = 0.01$

# Pneumonia detection under biased sampling

## Training data

Target label:      Healthy      Pneumonia  
 $Y = 0$                        $Y = 1$

Auxiliary label:      Men  
 $V = 0$

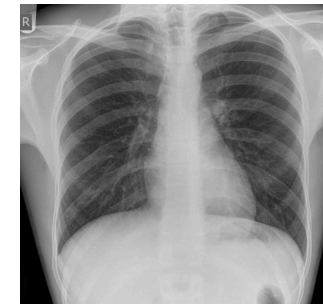


Women  
 $V = 1$



## Testing data

Man with  
pneumonia



Woman with  
pneumonia



$P(\text{Pneumonia}) = 0.99$      $P(\text{Pneumonia}) = 0.01$

**Shortcut learning**

# Pneumonia detection under biased sampling

## Training data

Target label: Healthy      Pneumonia  
 $Y = 0$                        $Y = 1$

Auxiliary label: Men  
 $V = 0$

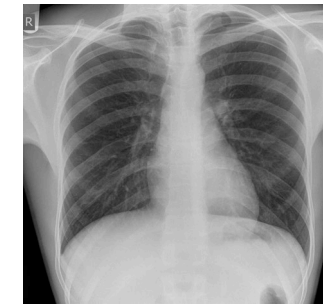


Women  
 $V = 1$



## Testing data

Man with pneumonia



Woman with pneumonia



$P(\text{Pneumonia}) = 0.99$      $P(\text{Pneumonia}) = 0.01$

**Shortcut learning:** relying on spurious correlations to predict the target label

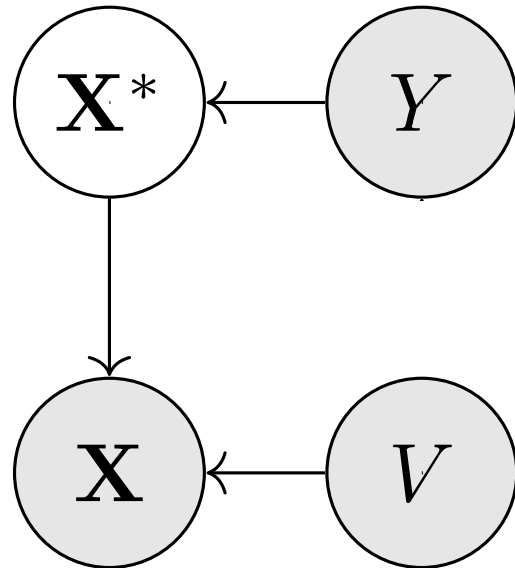


Can we build models that are robust  
to shortcuts?

# Causal assumptions

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)



X-ray  
pixels

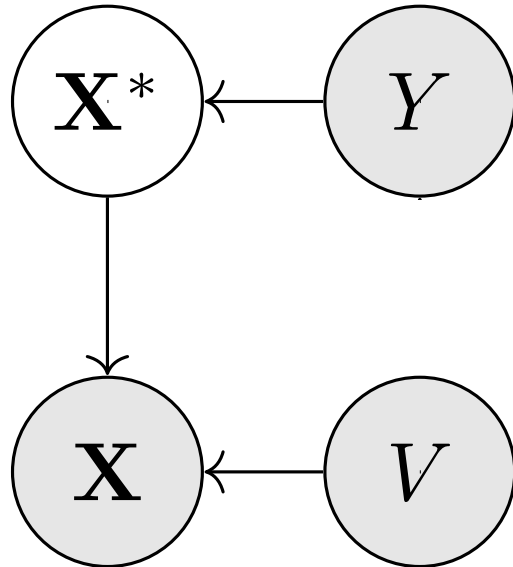
Auxiliary label  
(sex)

# Causal assumptions

→ Causal, fixed

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)



X-ray  
pixels

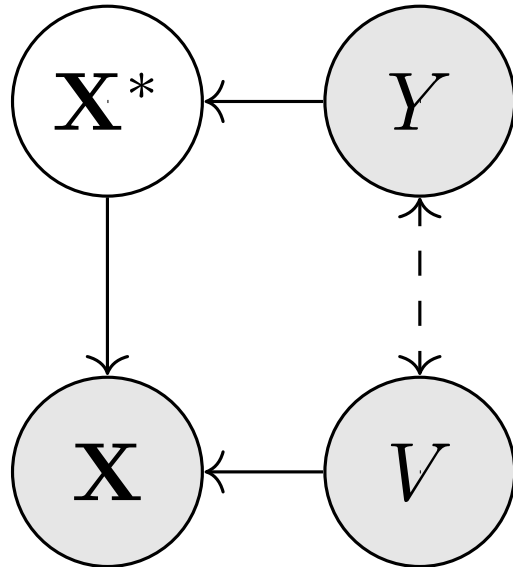
Auxiliary label  
(sex)

# Causal assumptions

→ Causal, fixed

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)



Correlated,  
not causally  
related

X-ray  
pixels

Auxiliary label  
(sex)

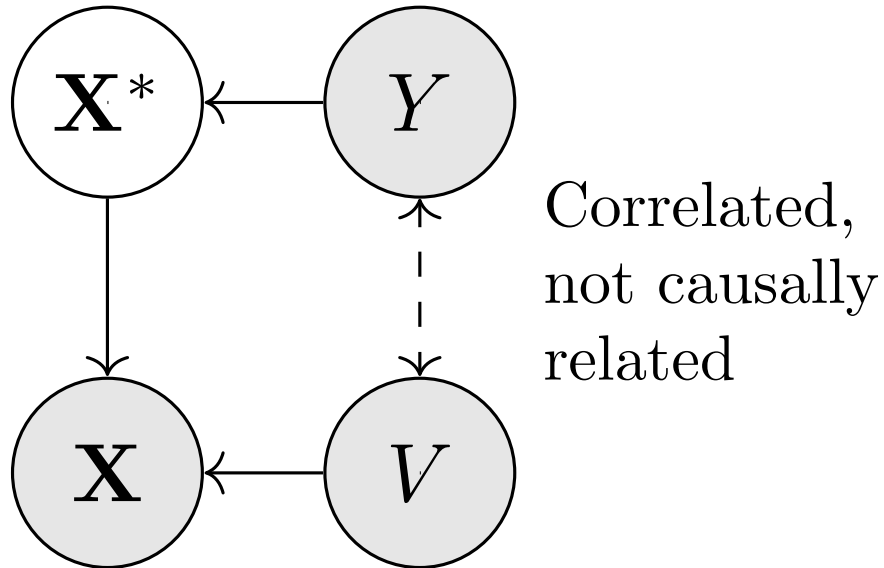
# Causal assumptions

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)

→ Causal, fixed

← - - -> Correlation, varies



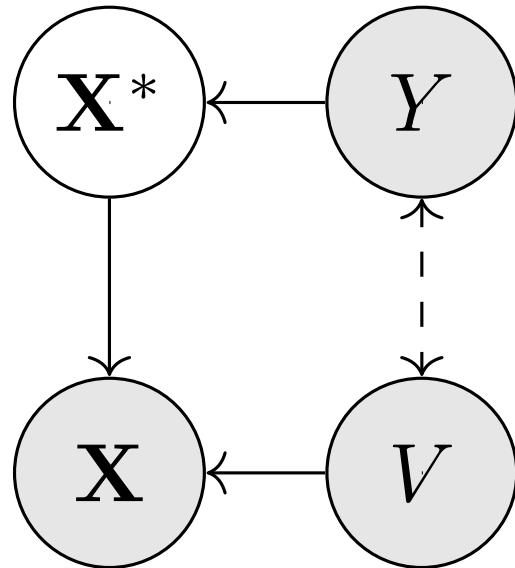
X-ray  
pixels

Auxiliary label  
(sex)

# Causal assumptions

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)



Correlated,  
not causally  
related

X-ray  
pixels

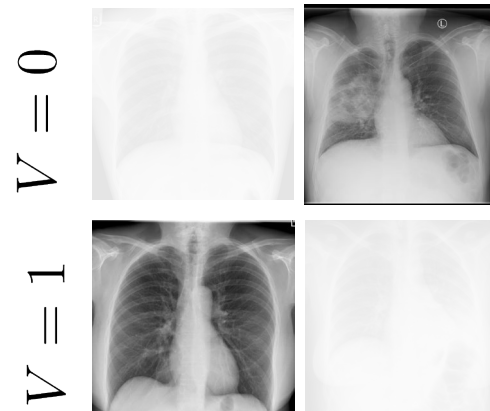
Auxiliary label  
(sex)

→ Causal, fixed

← - - -> Correlation, varies

**Hospital A**

$Y = 0$     $Y = 1$

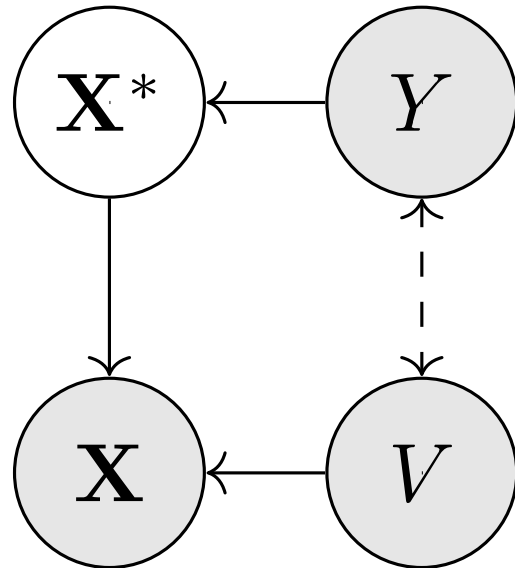


$$P(V|Y) = 0.9$$

# Causal assumptions

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)



Correlated,  
not causally  
related

X-ray  
pixels

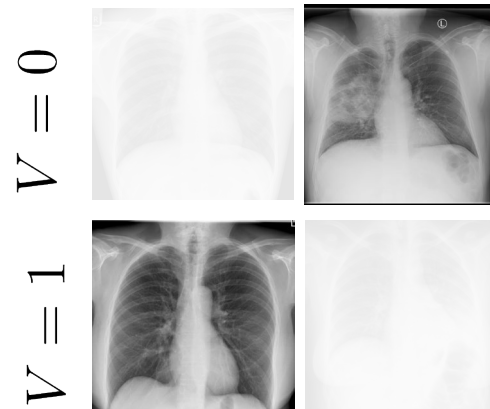
Auxiliary label  
(sex)

→ Causal, fixed

← - - → Correlation, varies

**Hospital A**

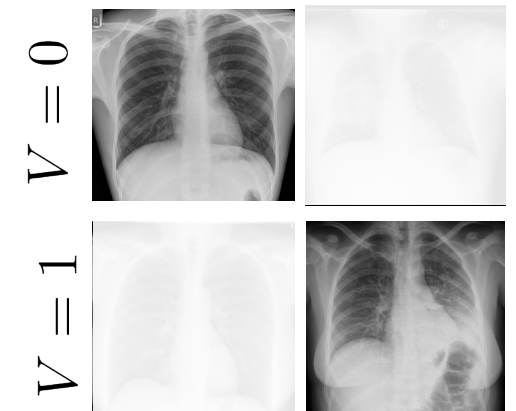
$Y = 0$     $Y = 1$



$P(V|Y) = 0.9$

**Hospital B**

$Y = 0$     $Y = 1$

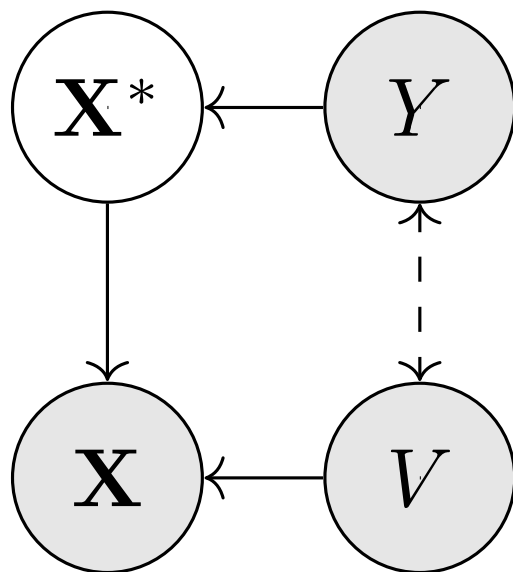


$P(V|Y) = 0.1$

# Task formulation

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)

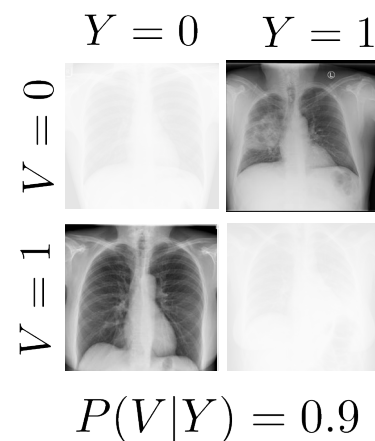


Correlated,  
not causally  
related

X-ray  
pixels

Auxiliary label  
(sex)

Training data



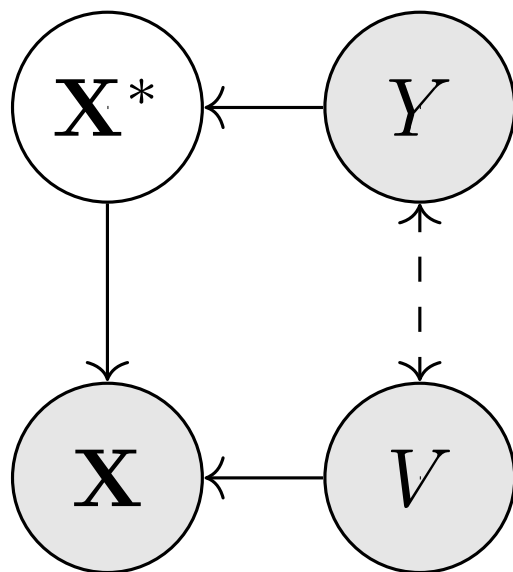
**Have:** Data with some correlation  
between  $V, Y$



# Task formulation

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)

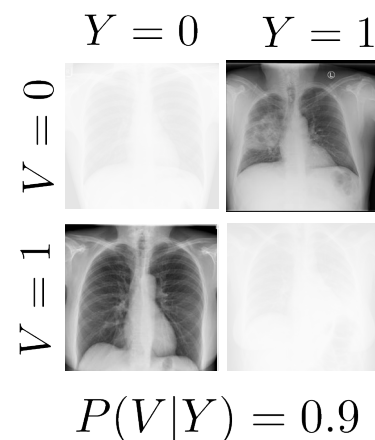


Correlated,  
not causally  
related

X-ray  
pixels

Auxiliary label  
(sex)

Training data



**Have:** Data with some correlation  
between  $V, Y$

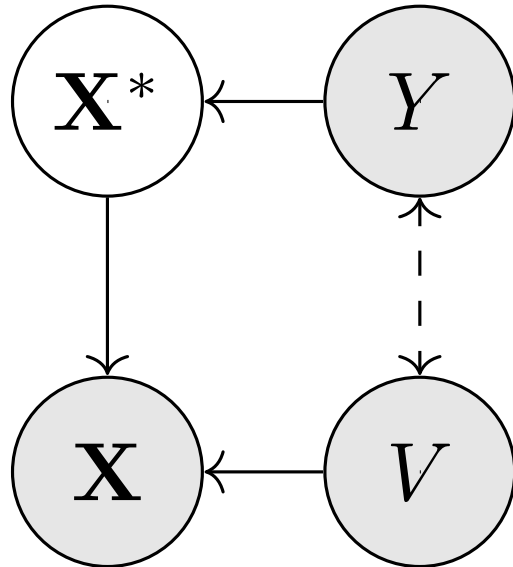
**Want:**  $f(\mathbf{X})$

**Such that:**  $f$  is accurate

# Task formulation

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)

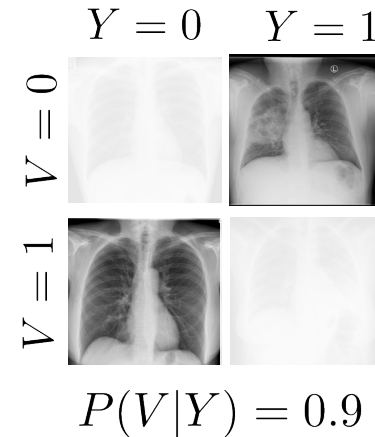


Correlated,  
not causally  
related

X-ray  
pixels

Auxiliary label  
(sex)

Training data



**Have:** Data with some correlation  
between  $V, Y$

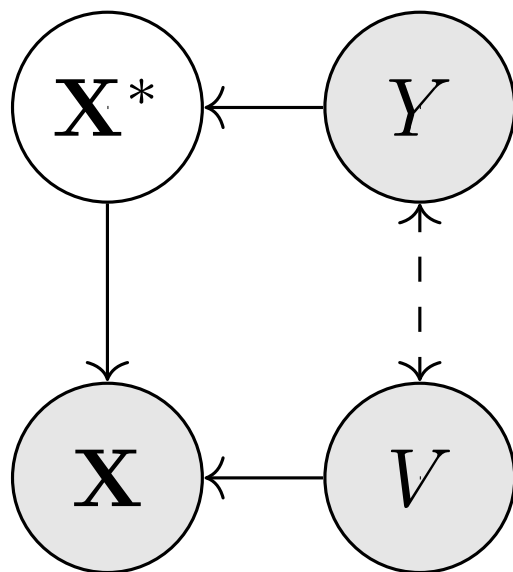
**Want:**  $f(\mathbf{X})$

**Such that:**  $f$  is accurate

# Task formulation

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)

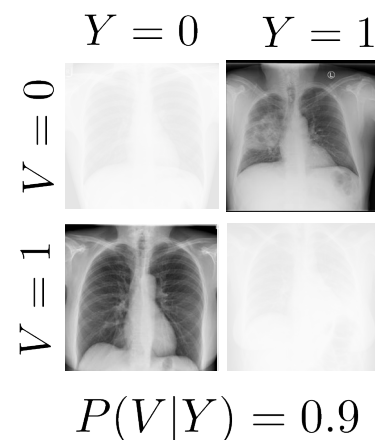


Correlated,  
not causally  
related

X-ray  
pixels

Auxiliary label  
(sex)

Training data

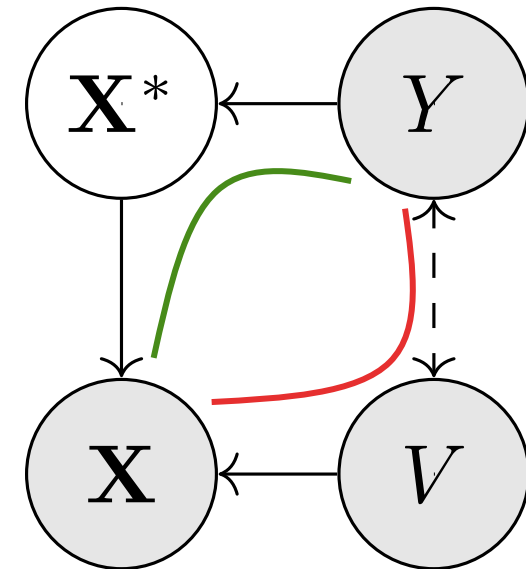


**Have:** Data with some correlation  
between  $V, Y$

**Want:**  $f(\mathbf{X})$

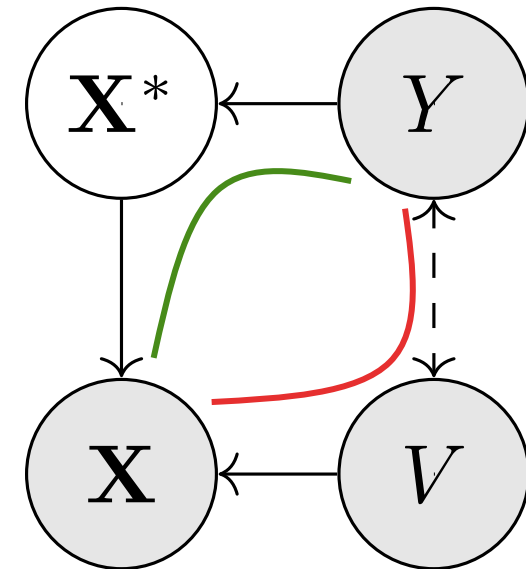
**Such that:**  $f$  is accurate and robust  
to the shortcut, i.e., has the same  
performance across all  $V, Y$  correlations

# Root cause of shortcut learning



# Root cause of shortcut learning

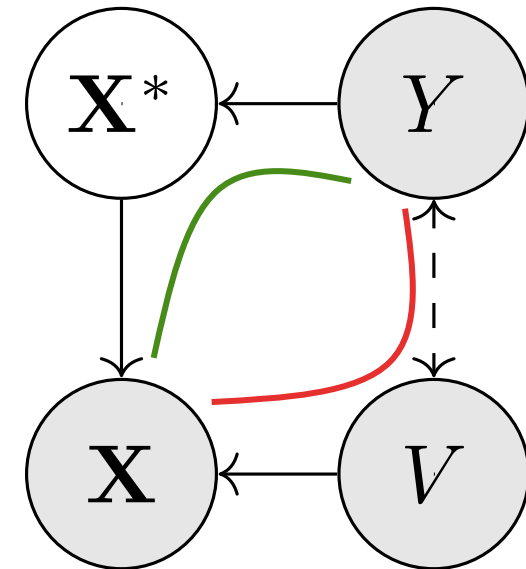
Causal (green) path: robust



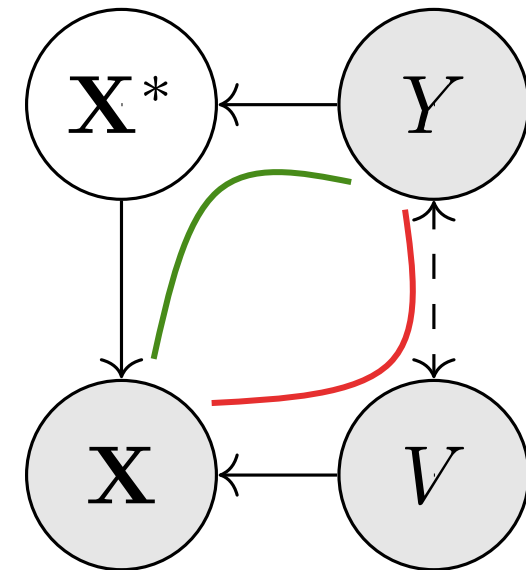
# Root cause of shortcut learning

Causal (green) path: robust

Non-causal (red) path: encodes shortcut

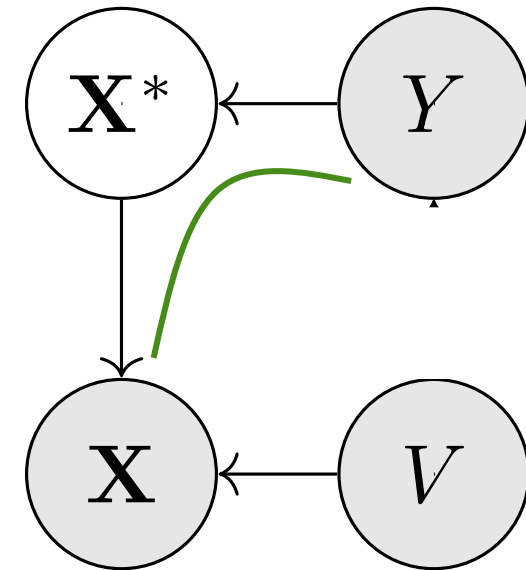


# Population-level robustness to the shortcut



# Population-level robustness to the shortcut

Ideal distribution  $\rightarrow$  no correlation between  $V, Y$

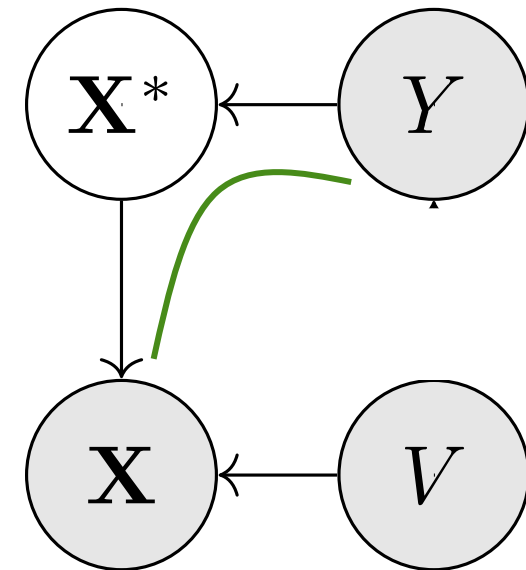




# Population-level robustness to the shortcut

Ideal distribution  $\rightarrow$  no correlation between  $V, Y$

**Proposition** (informal): Under the ideal distribution, with a very large dataset, the optimal model is robust to shortcuts.



# Finite sample analysis (under the ideal distribution)

## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

### Usual models

$$f^* = \min \ell(f(\mathbf{x}), y) + \alpha \cdot \|f\|_2$$

## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

Usual models

$$f^* = \min \ell(f(\mathbf{x}), y) + \alpha \cdot \|f\|_2$$

Models conforming to the DAG

## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

### Usual models

$$f^* = \min \ell(f(\mathbf{x}), y) + \alpha \cdot \|f\|_2$$

### Models conforming to the DAG

= Models that do not encode correlations between  $Y, V$

# Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

## Usual models

$$f^* = \min \ell(f(\mathbf{x}), y) + \alpha \cdot \|f\|_2$$

## Models conforming to the DAG

= Models that do not encode correlations between  $Y, V$

$$P(f(\mathbf{x})|V = 1) = P(f(\mathbf{x})|V = 0)$$

Predictions for women      Predictions for men

# Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

## Usual models

$$f^* = \min \ell(f(\mathbf{x}), y) + \alpha \cdot \|f\|_2$$

## Models conforming to the DAG

= Models that do not encode correlations between  $Y, V$

$$P(f(\mathbf{x})|V = 1) = P(f(\mathbf{x})|V = 0)$$

Predictions for women    Predictions for men

In practice:

$$f^* = \min \ell(f(\mathbf{x}), y) + \alpha \cdot \text{MMD}(P(f(\mathbf{x})|V = 1), P(f(\mathbf{x})|V = 0))$$



## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

**Simple example:**

## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

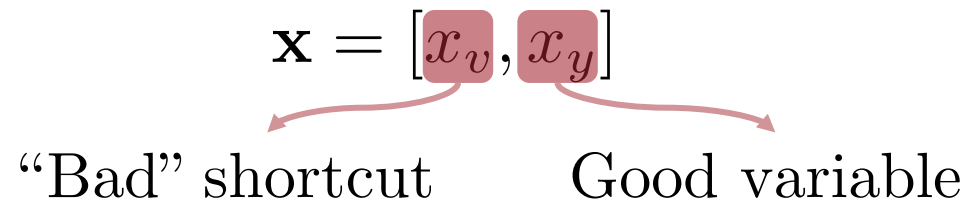
**Simple example:**

$$\mathbf{x} = [x_v, x_y]$$

## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

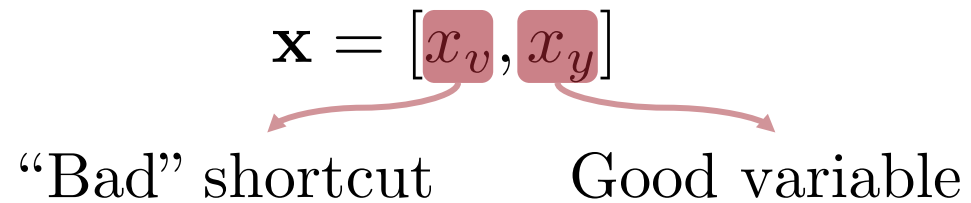
**Simple example:**



## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

**Simple example:**



Don't know: which is  $x_v$  vs.  $x_y$

## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

**Simple example:**

$$\mathbf{x} = [x_v, x_y]$$

“Bad” shortcut

Good variable

Don't know: which is  $x_v$  vs.  $x_y$

$$f(\mathbf{x}) = w_y x_y + w_v x_v$$

## Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

**Simple example:**

Model learning = search problem

$$\mathbf{x} = [x_v, x_y]$$

“Bad” shortcut

Good variable

Don't know: which is  $x_v$  vs.  $x_y$

$$f(\mathbf{x}) = w_y x_y + w_v x_v$$

# Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

## Simple example:

Model learning = search problem

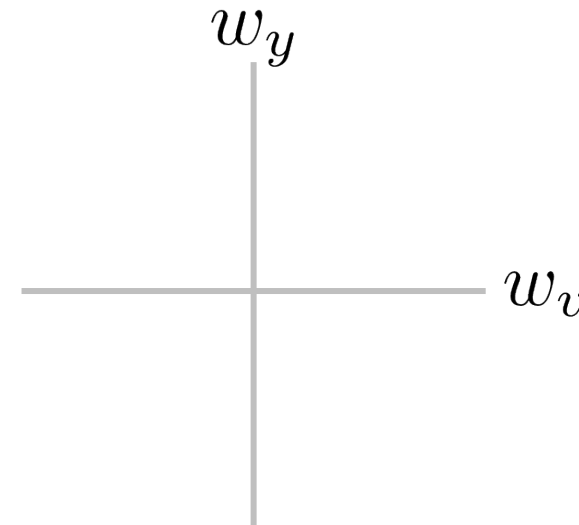
$$\mathbf{x} = [x_v, x_y]$$

“Bad” shortcut

Good variable

Don't know: which is  $x_v$  vs.  $x_y$

$$f(\mathbf{x}) = w_y x_y + w_v x_v$$



# Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

## Simple example:

Model learning = search problem

$$\mathbf{x} = [x_v, x_y]$$

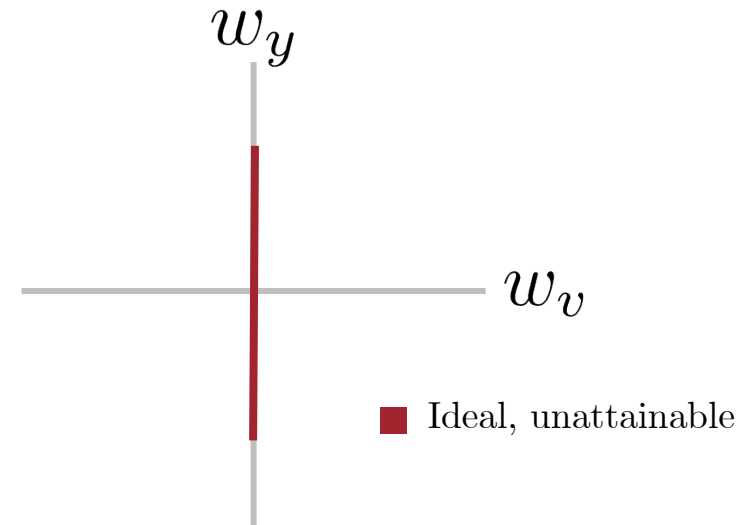
“Bad” shortcut

Good variable

Don't know: which is  $x_v$  vs.  $x_y$

$$f(\mathbf{x}) = w_y x_y + w_v x_v$$

If  $w_v \neq 0$ , model encodes info about shortcut





# Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

## Simple example:

Model learning = search problem

$$\mathbf{x} = [x_v, x_y]$$

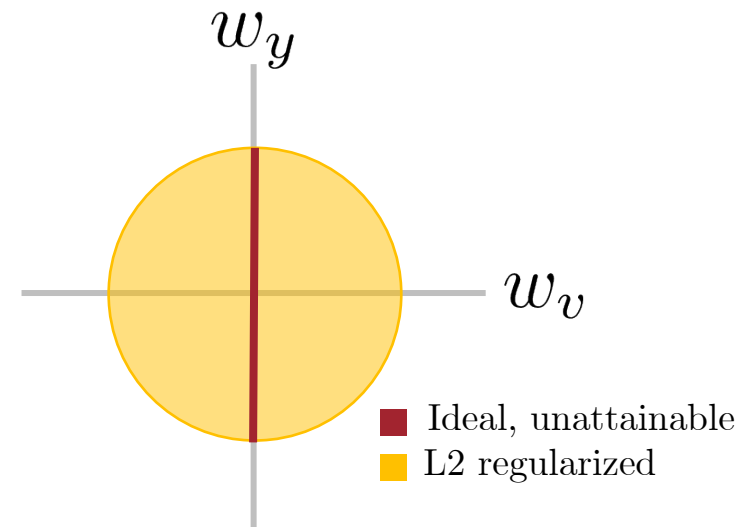
“Bad” shortcut

Good variable

Don't know: which is  $x_v$  vs.  $x_y$

$$f(\mathbf{x}) = w_y x_y + w_v x_v$$

If  $w_v \neq 0$ , model encodes info about shortcut



# Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

**Simple example:**

Model learning = search problem

$$\mathbf{x} = [x_v, x_y]$$

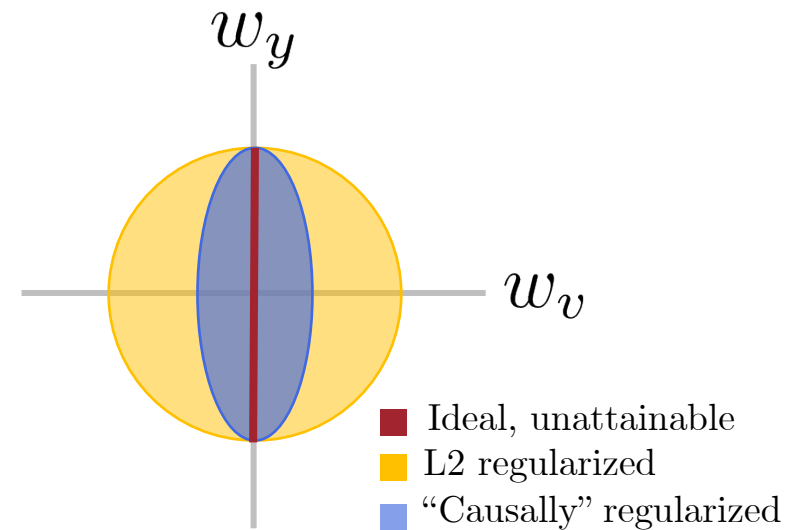
“Bad” shortcut

Good variable

Don't know: which is  $x_v$  vs.  $x_y$

$$f(\mathbf{x}) = w_y x_y + w_v x_v$$

If  $w_v \neq 0$ , model encodes info about shortcut



# Finite sample analysis (under the ideal distribution)

**Proposition** (informal): Models that conform to the causal DAG are more efficient than “the usual models” in finite samples.

**Simple example:**

$$\mathbf{x} = [x_v, x_y]$$

“Bad” shortcut

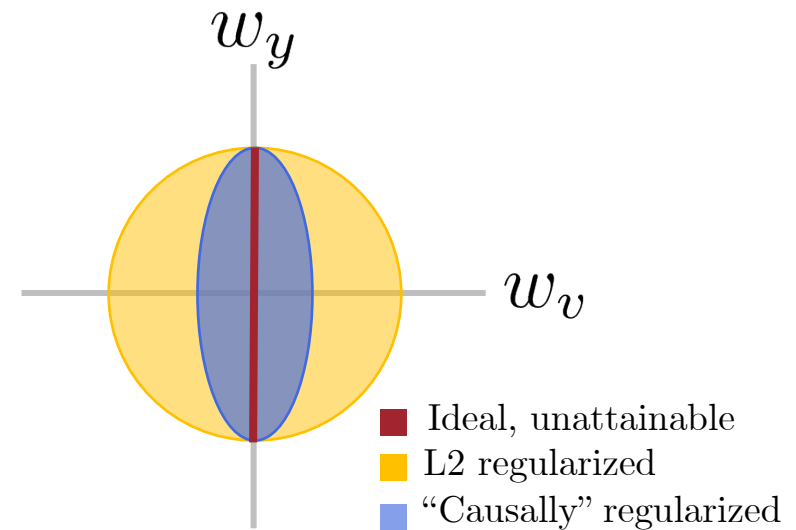
Good variable

Don't know: which is  $x_v$  vs.  $x_y$

$$f(\mathbf{x}) = w_y x_y + w_v x_v$$

If  $w_v \neq 0$ , model encodes info about shortcut

Model learning = search problem



→ Smaller Rademacher complexity

# Quick recap

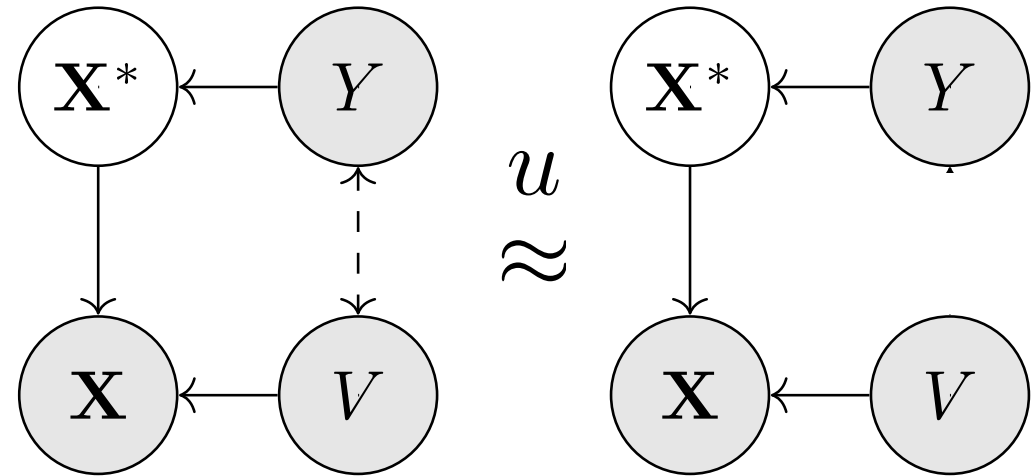
Ideal distribution + large sample	Robustness to the shortcut
Ideal distribution	Statistical efficiency
An arbitrary distribution	?

# Quick recap

Ideal distribution + large sample	Robustness to the shortcut
Ideal distribution	Statistical efficiency
An arbitrary distribution	Shortcut learning + bias!

# Sampling from non-ideal distributions

# Sampling from non-ideal distributions

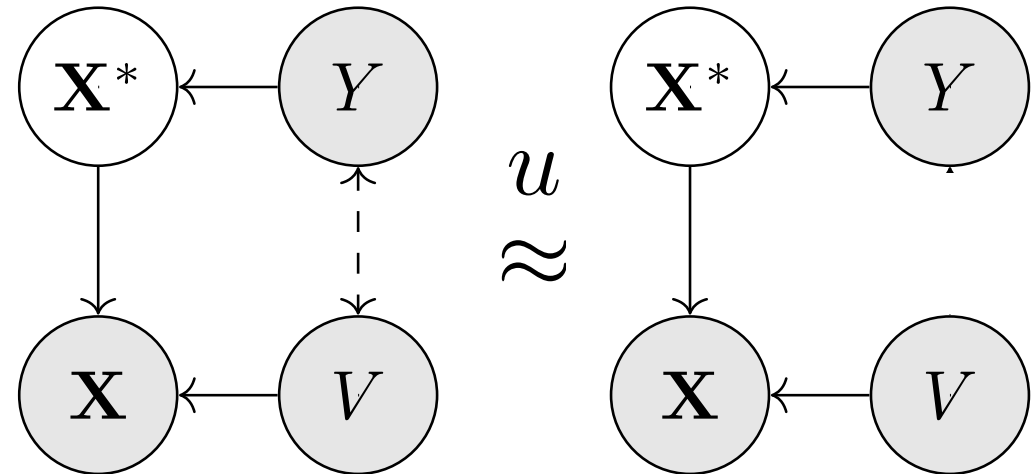


# Sampling from non-ideal distributions

**Proposition** (informal): Reweighting with  $u_i$  recovers the independences in the ideal distribution

$u_i \rightarrow$  Makes  $Y, V$  “look” independent

$$u_i = \frac{P(Y = y_i)P(V = v_i)}{P(Y = y_i, V = v_i)}$$





# Training objective

$$f^* = \operatorname{argmin}_f \sum_i u_i \ell(f(\mathbf{x}_i), y_i)$$

Weighted prediction loss

$$+ \alpha \cdot \operatorname{MMD}_{\mathbf{u}} \left( P(f(\mathbf{x}_i) | v_i = 1), P(f(\mathbf{x}_i) | v_i = 0) \right)$$

Weighted penalty on predictions encoding information about  $V$

# Training objective

$$f^* = \operatorname{argmin}_f \sum_i u_i \ell(f(\mathbf{x}_i), y_i) \quad \text{Weighted prediction loss}$$
$$+ \alpha \cdot \operatorname{MMD}_{\mathbf{u}} \left( P(f(\mathbf{x}_i) | v_i = 1), P(f(\mathbf{x}_i) | v_i = 0) \right) \quad \text{Weighted penalty on predictions encoding information about } V$$

Causal perspective gave us:

1. Weights to map the training data to a distribution where invariance is achievable

# Training objective

$$f^* = \operatorname{argmin}_f \sum_i u_i \ell(f(\mathbf{x}_i), y_i)$$

Weighted prediction loss

$$+ \alpha \cdot \operatorname{MMD}_{\mathbf{u}} \left( P(f(\mathbf{x}_i) | v_i = 1), P(f(\mathbf{x}_i) | v_i = 0) \right)$$

Weighted penalty on predictions encoding information about  $V$

Causal perspective gave us:

1. Weights to map the training data to a distribution where invariance is achievable
2. Invariance penalty to encourage the model to encode desirable independencies

# Empirical results: water birds data

- **Data:** Semi-simulated

Wah *et al*, Computation & Neural Systems Technical Report 2010  
Zhou *et al*, IEEE PAMI 2017  
Sagawa *et al*, ICLR 2020

# Empirical results: water birds data

- **Data:** Semi-simulated
- **Task:** Predict type of bird
  - Main label = type of bird (water/land)
  - Auxiliary label = type of background (water/land)



# Empirical results: water birds data

- **Data:** Semi-simulated
- **Task:** Predict type of bird
  - Main label = type of bird (water/land)
  - Auxiliary label = type of background (water/land)
- **Setup:** Fixed training (source) data

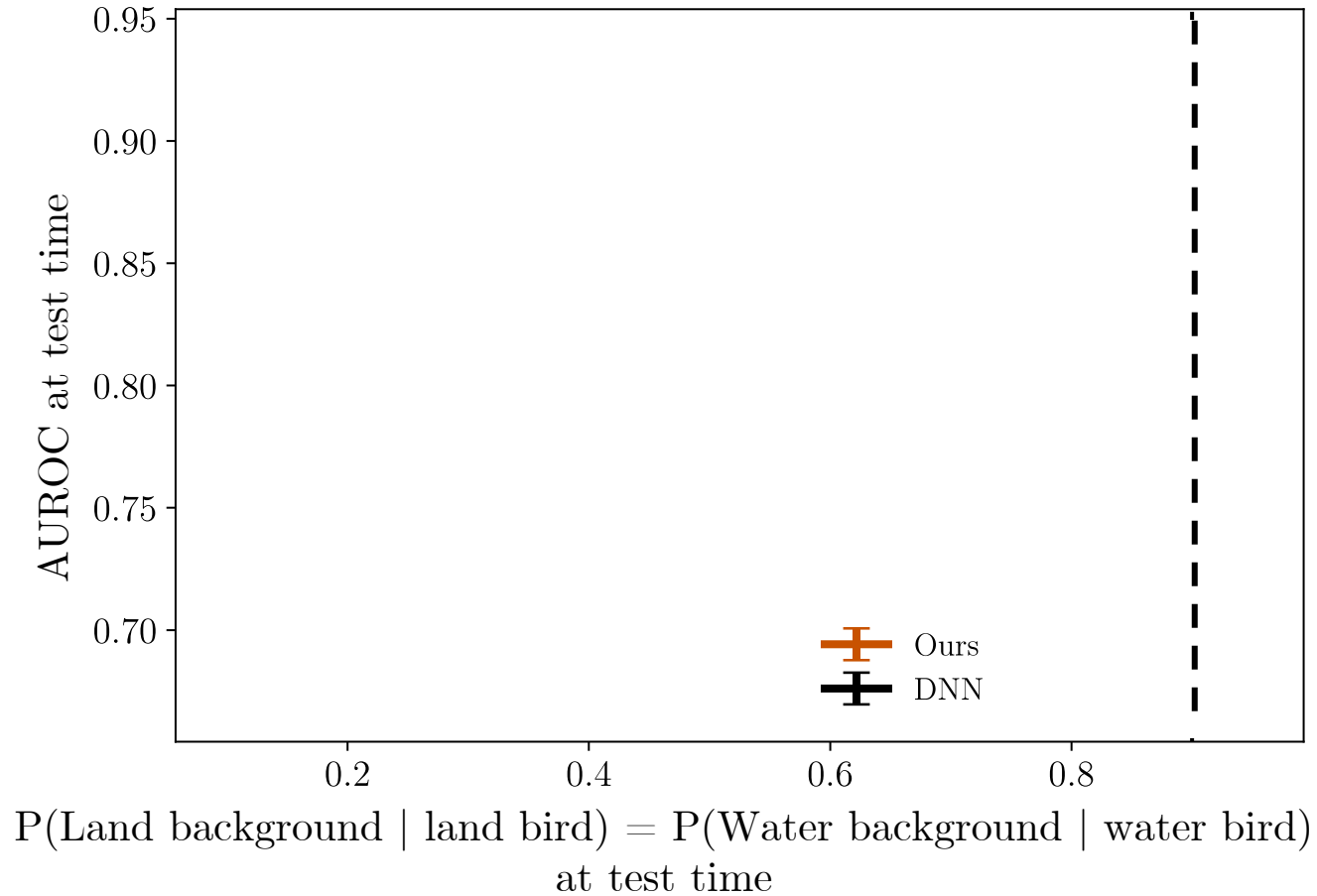


# Empirical results: water birds data

- **Data:** Semi-simulated
- **Task:** Predict type of bird
  - Main label = type of bird (water/land)
  - Auxiliary label = type of background (water/land)
- **Setup:** Fixed training (source) data
- **Evaluation:** On multiple test sets with different correlations

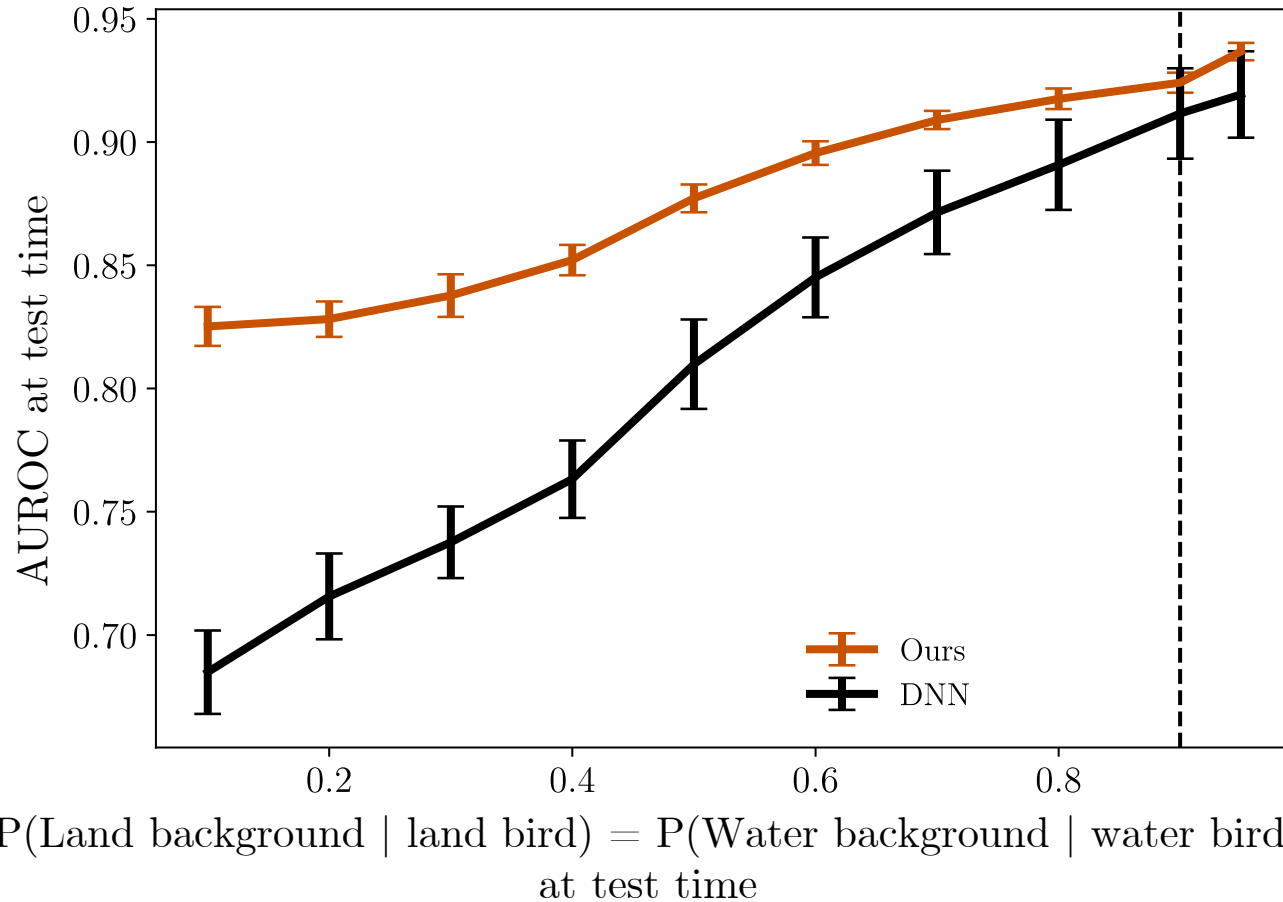


# Experiment results



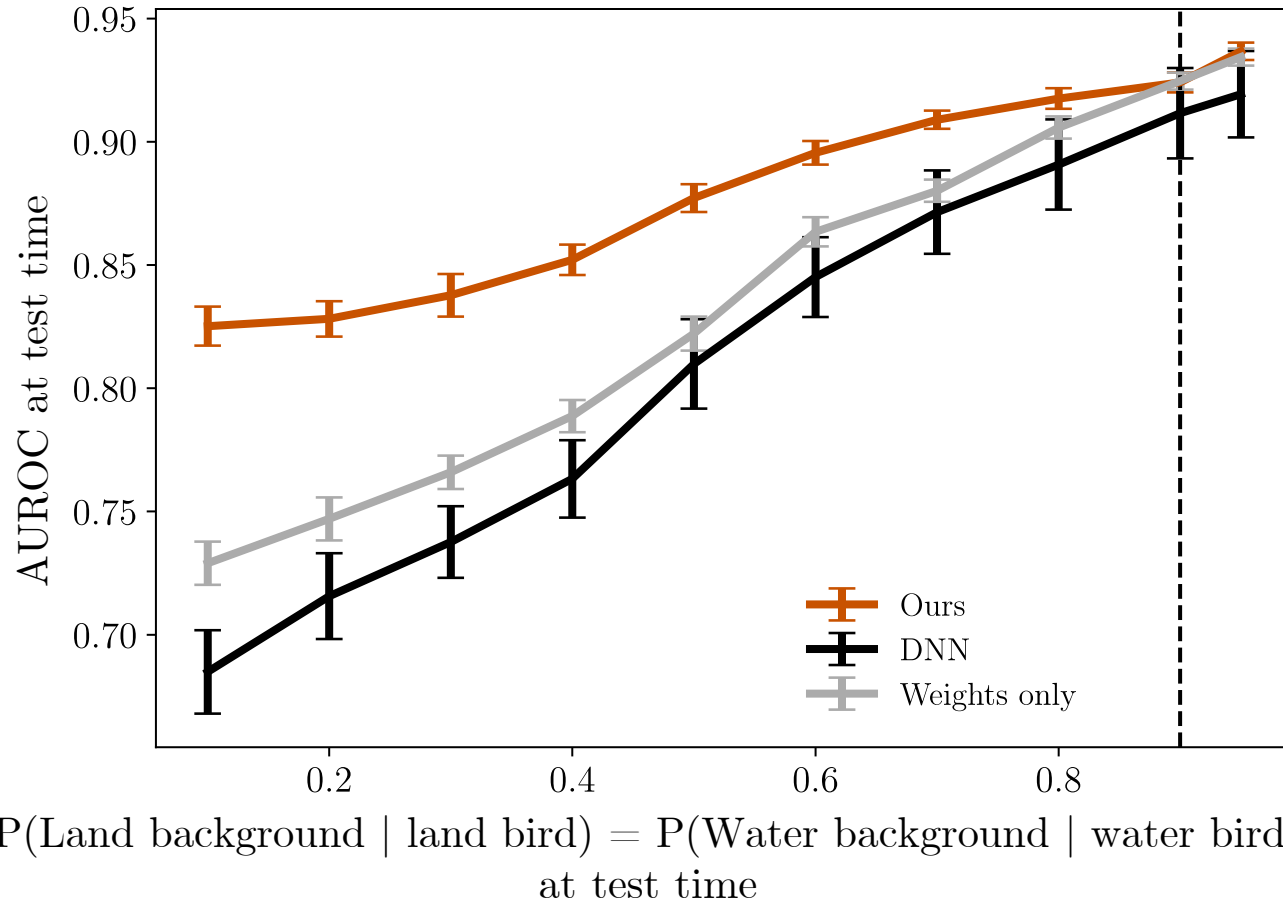


# Experiment results

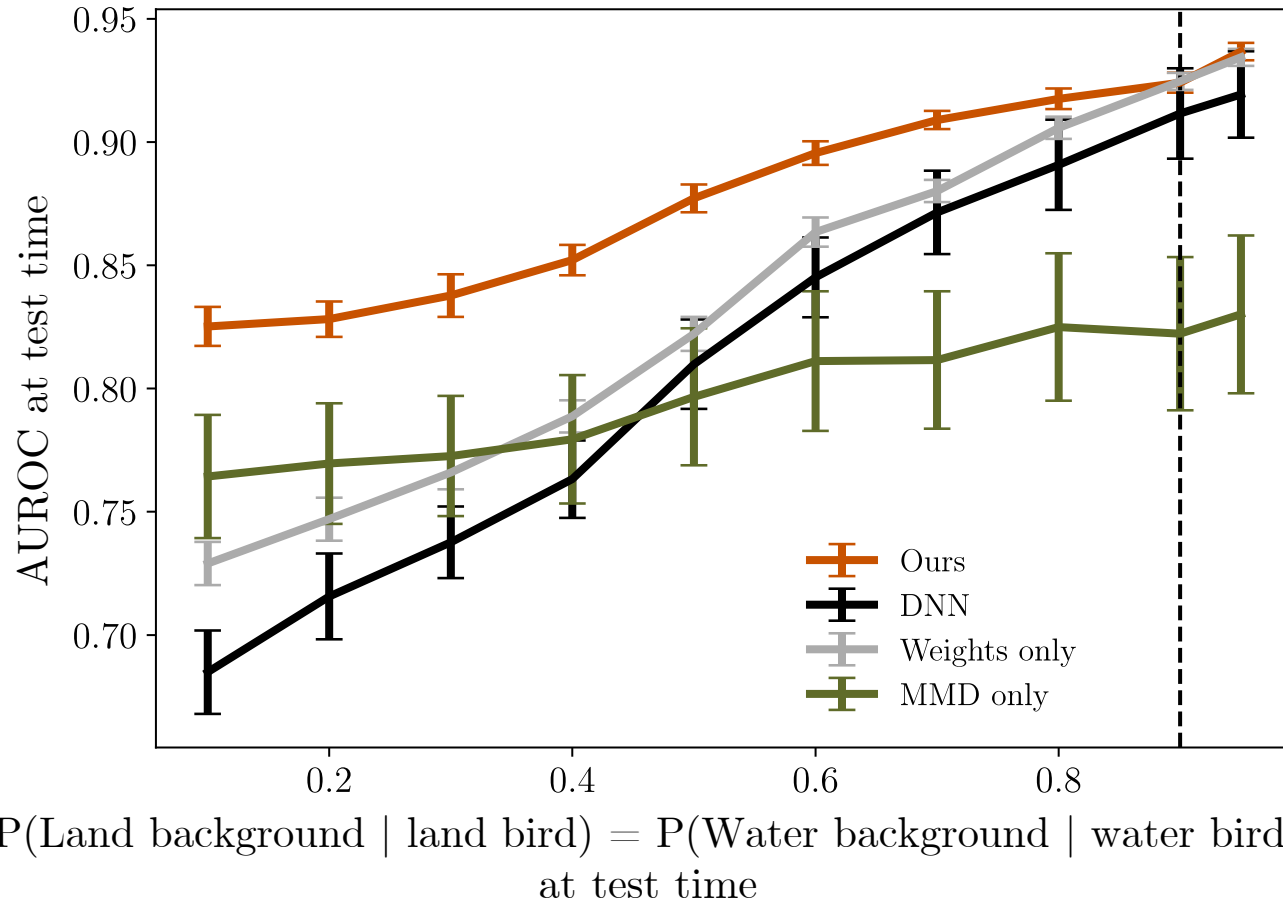


Our approach is less prone to relying on the shortcut, and is more accurate than typical methods

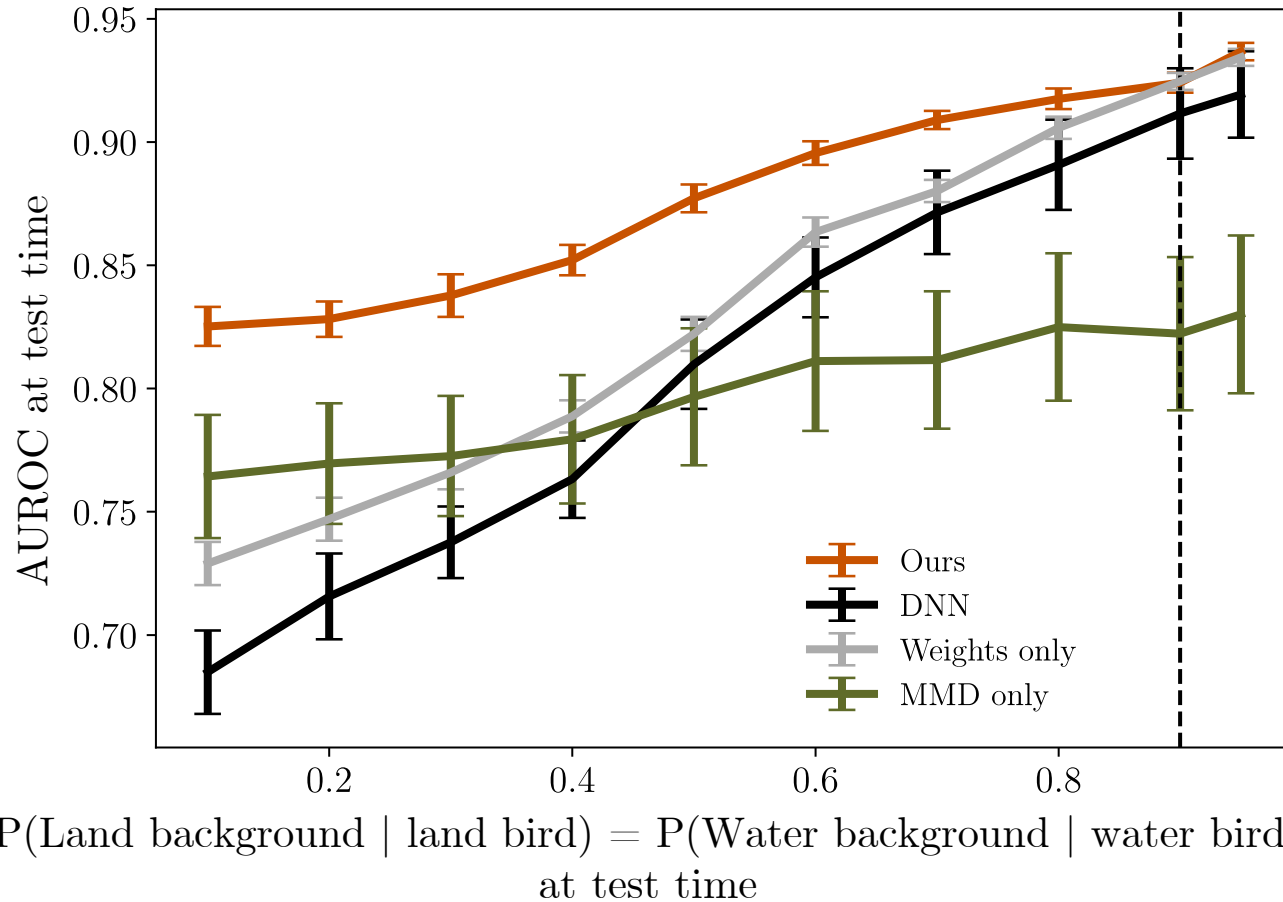
# Experiment results



# Experiment results



# Experiment results



The two causally-inspired components of our approach are necessary

# Empirical results: Predicting Pneumonia

- **Data:** CheXpert, down-sample women with pneumonia at training time.
- **Task:** Predict the onset of pneumonia (main label), while making sure that sex (auxiliary label) is not a shortcut.

# Empirical results: Predicting Pneumonia

- **Data:** CheXpert, down-sample women with pneumonia at training time.
- **Task:** Predict the onset of pneumonia (main label), while making sure that sex (auxiliary label) is not a shortcut.

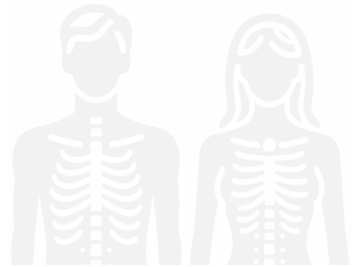
	AUROC	
	Ours	DNN
Same hospital (no shift)		
Different hospital (shift)		

# Empirical results: Predicting Pneumonia

- **Data:** CheXpert, down-sample women with pneumonia at training time.
- **Task:** Predict the onset of pneumonia (main label), while making sure that sex (auxiliary label) is not a shortcut.

	AUROC	
	Ours	DNN
Same hospital (no shift)	<b>0.85 (0.007)</b>	0.82 (0.03)
Different hospital (shift)	<b>0.75 (0.006)</b>	0.69 (0.028)

# Talk outline

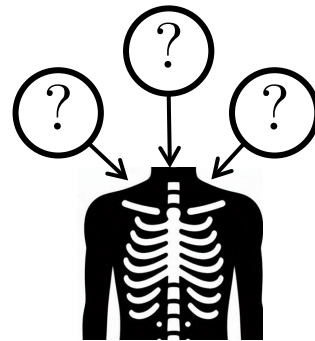


1 Efficiency + robustness to known sampling bias

MPMBHD, AISTATS 22

MD, TMLR 23

NM, UAI 24



2 Efficiency + robustness to unknown sampling biases

ZM, NeurIPS 22

WJMSW, NeurIPS 22



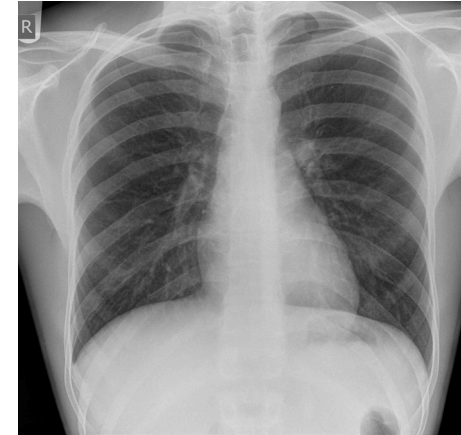
3 Evaluating localized circuits in LLMs

SVNZGJMB – NeurIPS 24



# What about unknown shortcuts?

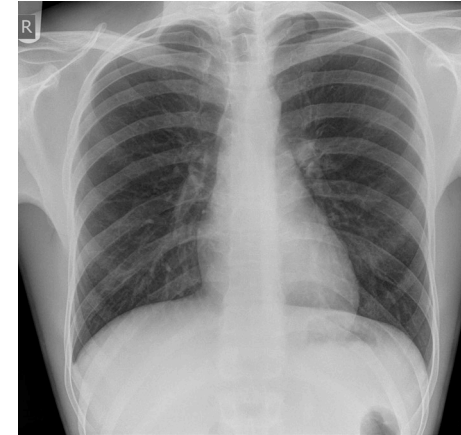
- **Have:** Large number of auxiliary labels



- + Patient sex
- + Type of X-ray machine
- + Other medical conditions:
  - Flu
  - Edema
  - Fractures
  - Cervical cancer...

# What about unknown shortcuts?

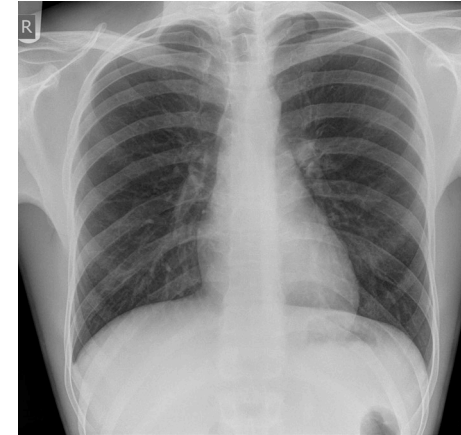
- **Have:** Large number of auxiliary labels
- **Unknown:** Which ones are relevant shortcuts?



- + Patient sex
- + Type of X-ray machine
- + Other medical conditions:
  - Flu
  - Edema
  - Fractures
  - Cervical cancer...

# What about unknown shortcuts?

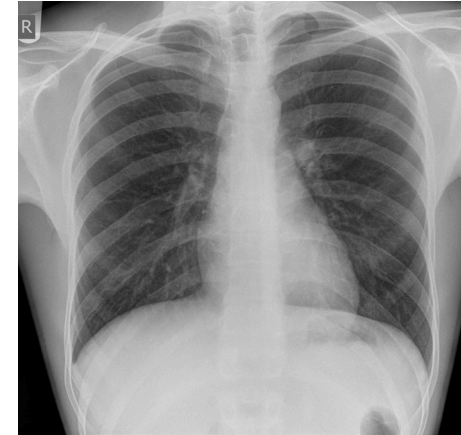
- **Have:** Large number of auxiliary labels
- **Unknown:** Which ones are relevant shortcuts?
- **Objective:** Models robust to multiple shortcuts



- + Patient sex
- + Type of X-ray machine
- + Other medical conditions:
  - Flu
  - Edema
  - Fractures
  - Cervical cancer...

# What about unknown shortcuts?

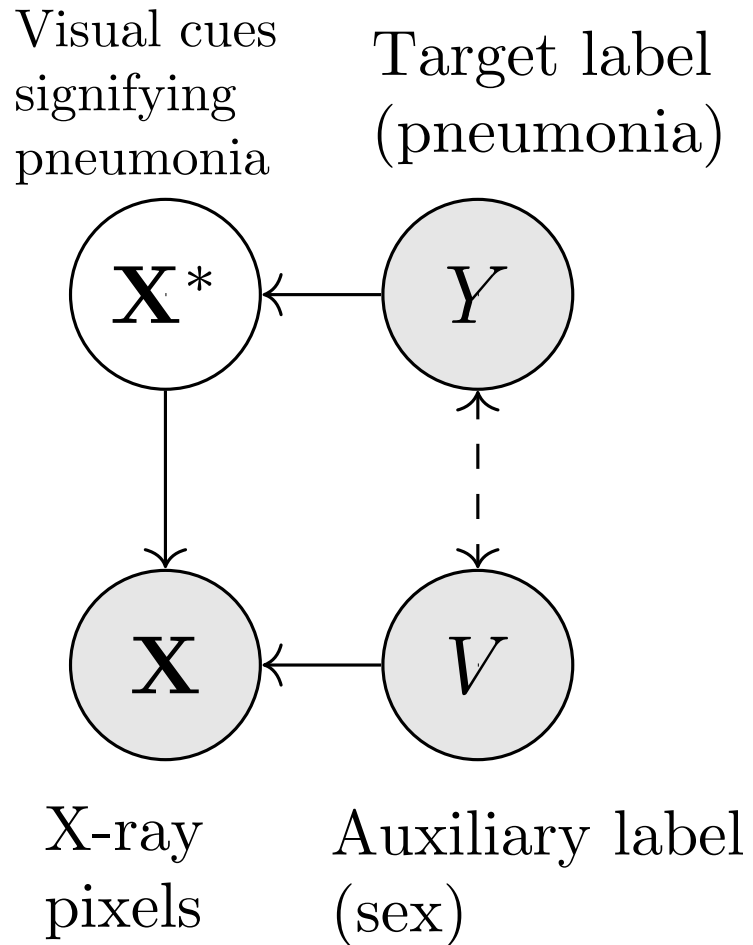
- **Have:** Large number of auxiliary labels
- **Unknown:** Which ones are relevant shortcuts?
- **Objective:** Models robust to multiple shortcuts
- **Upshot:** An additional causal discovery step



- + Patient sex
- + Type of X-ray machine
- + Other medical conditions:
  - Flu
  - Edema
  - Fractures
  - Cervical cancer...

# Recall the causal assumption

## Previous DAG

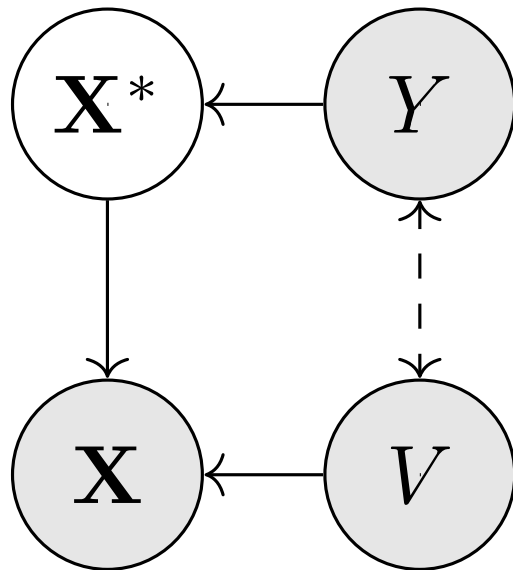


# Extension to a class of DAGs

Previous DAG

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)



X-ray  
pixels

Auxiliary label  
(sex)

New DAGs

A class of DAGs

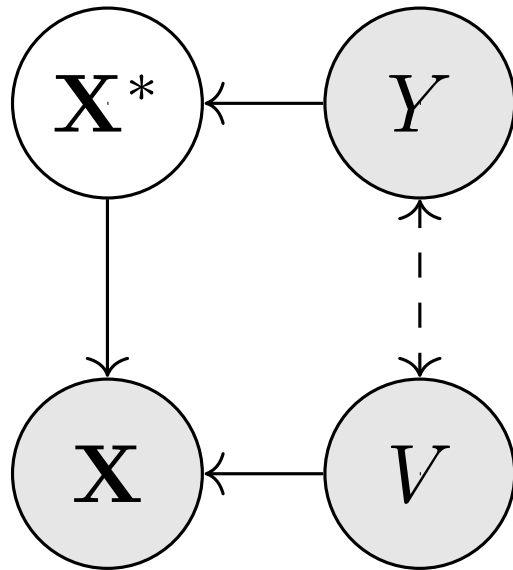
High dim. auxiliary labels  $V^{\text{full}}$

# Extension to a class of DAGs

## Previous DAG

Visual cues  
signifying  
pneumonia

Target label  
(pneumonia)



X-ray  
pixels

Auxiliary label  
(sex)

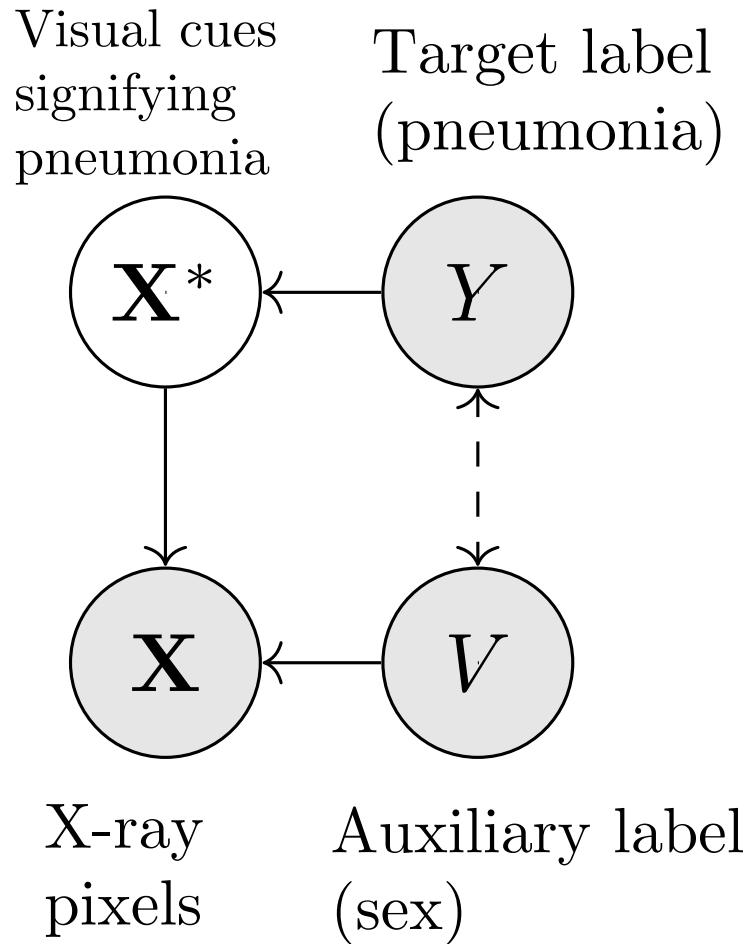
## New DAGs

A class of DAGs

High dim. auxiliary labels  $V^{\text{full}}$

# Extension to a class of DAGs

## Previous DAG



## New DAGs

A class of DAGs

High dim. auxiliary labels  $V^{\text{full}}$

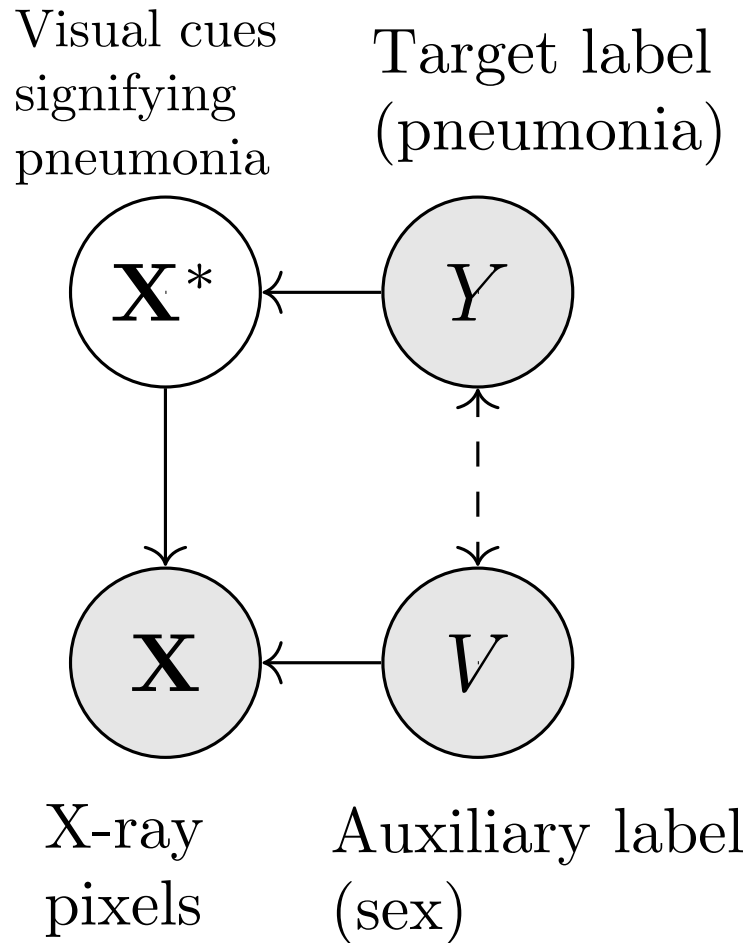


Correlated but not causally related to  $Y$  and cause  $X$



# Extension to a class of DAGs

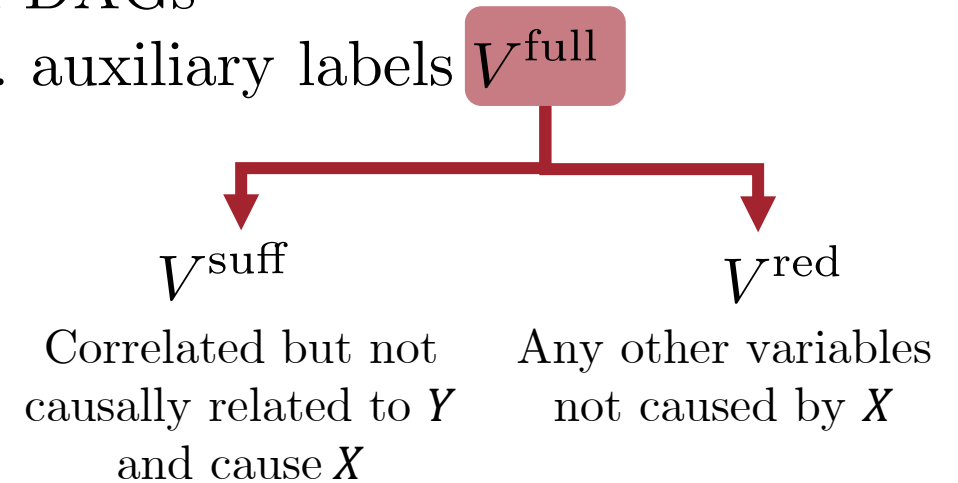
## Previous DAG



## New DAGs

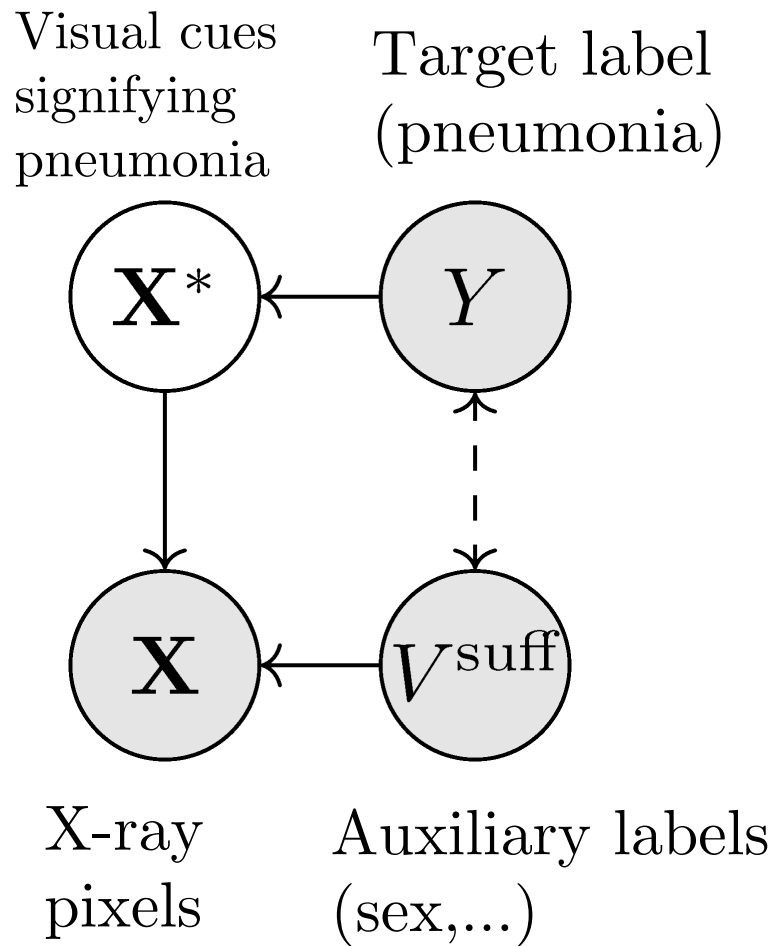
A class of DAGs

High dim. auxiliary labels



# Extension to a class of DAGs

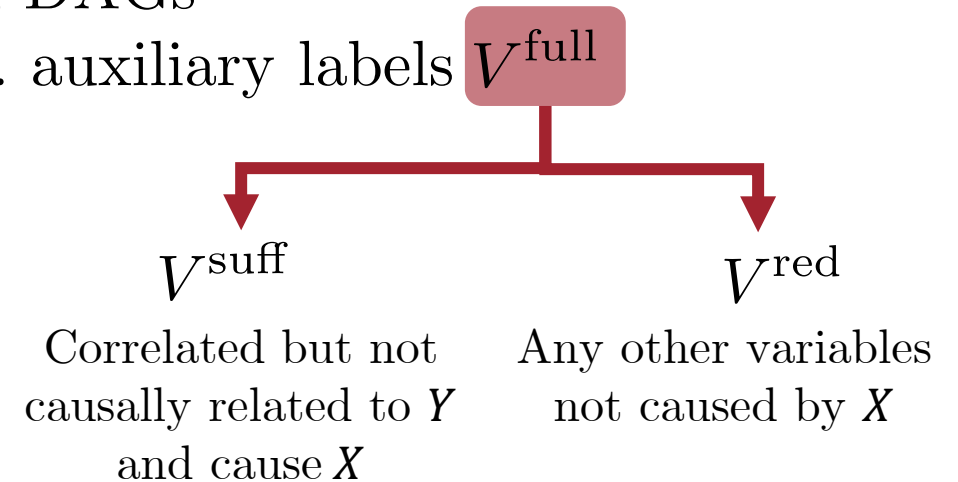
Example from new DAG class



New DAGs

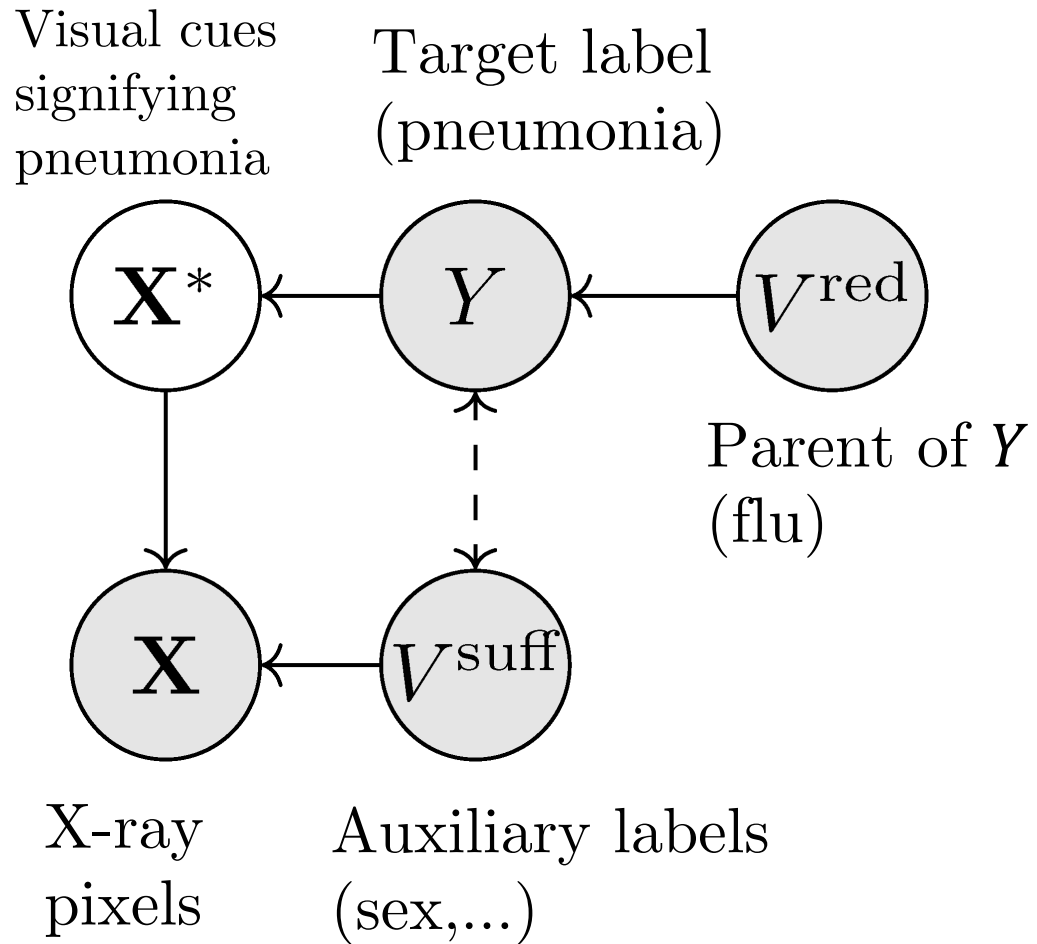
A class of DAGs

High dim. auxiliary labels



# Extension to a class of DAGs

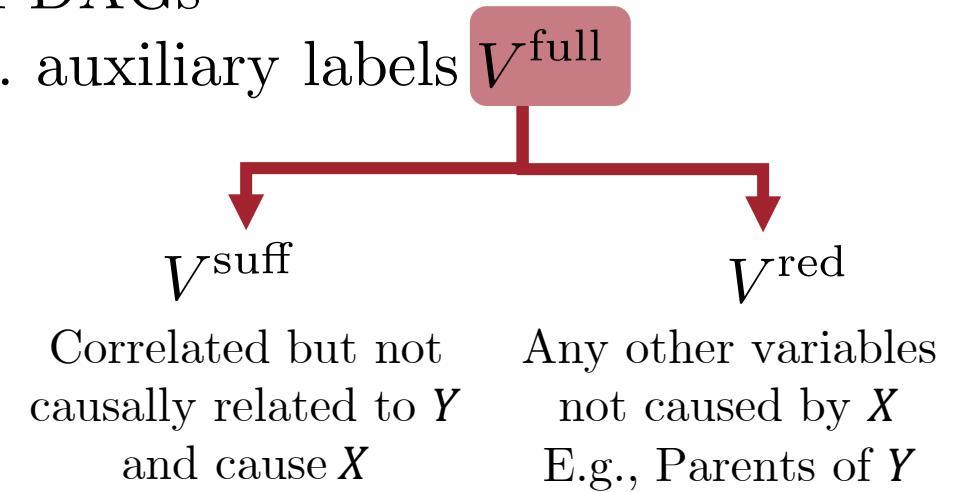
Example from new DAG class



New DAGs

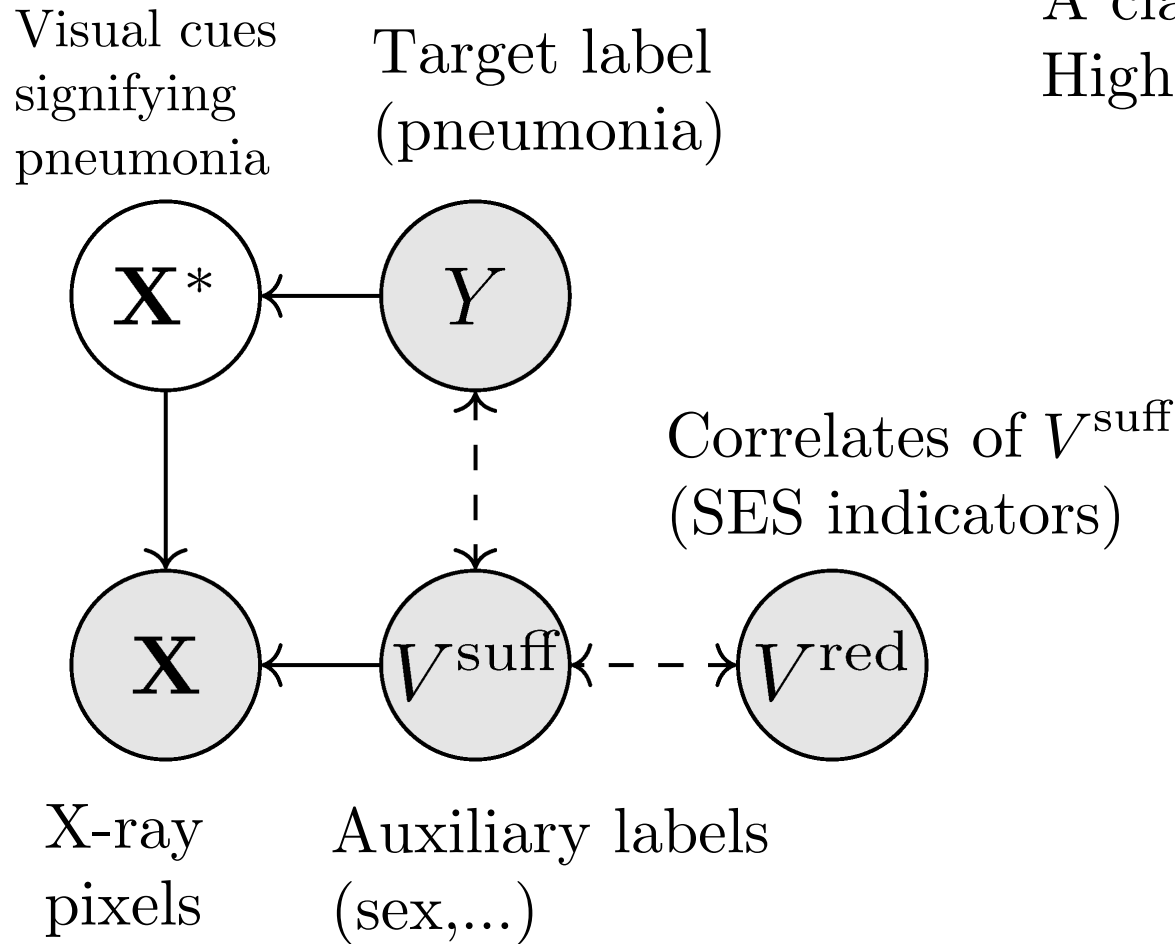
A class of DAGs

High dim. auxiliary labels



# Extension to a class of DAGs

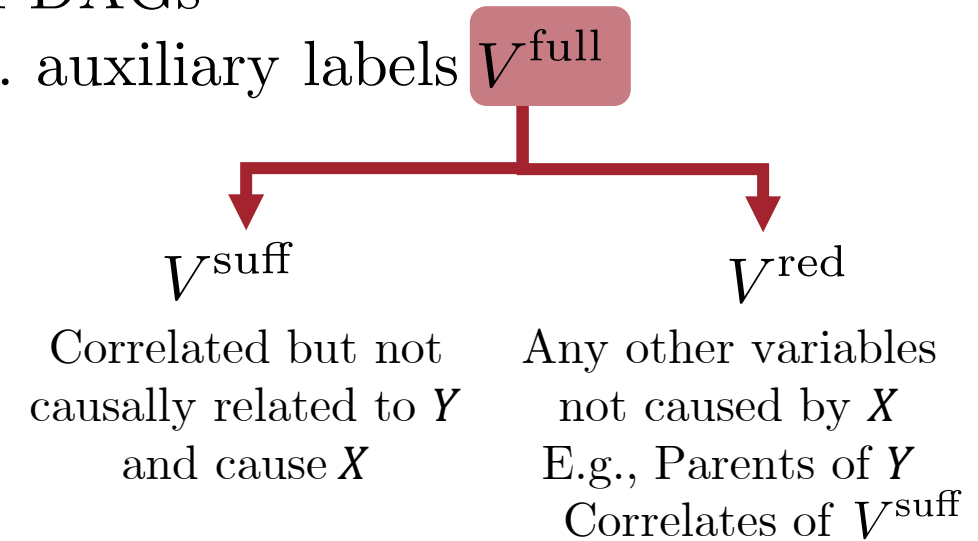
Example from new DAG class



New DAGs

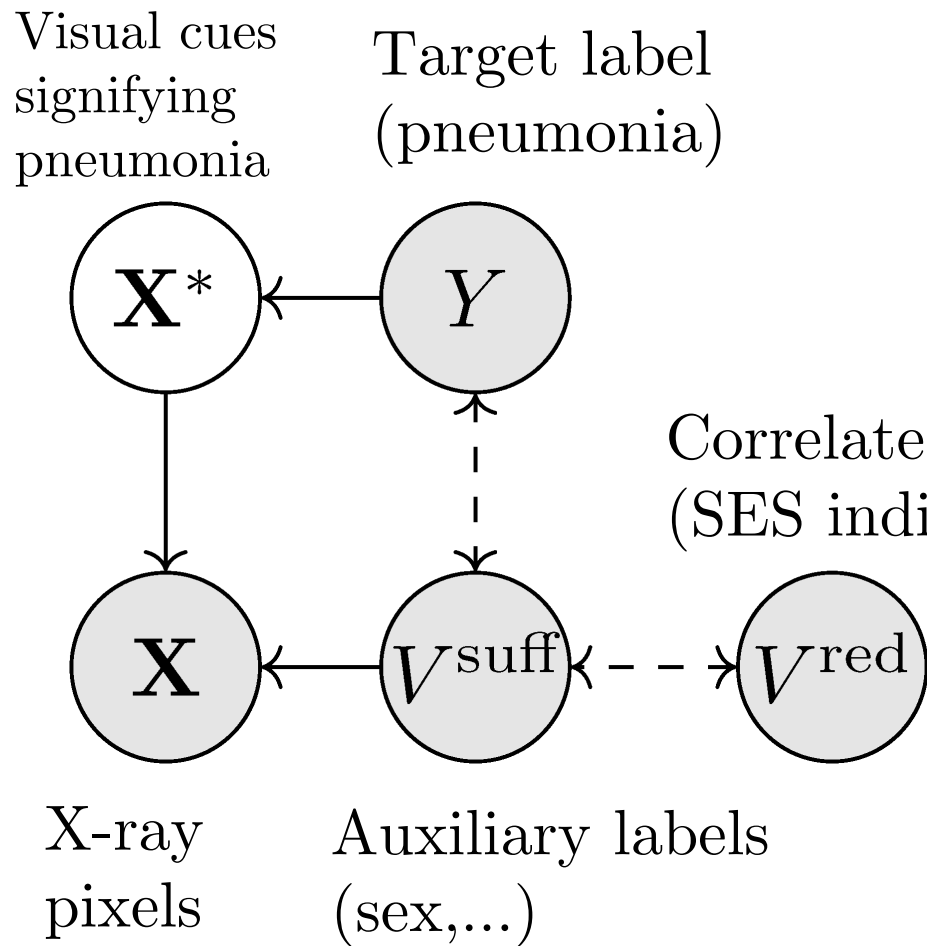
A class of DAGs

High dim. auxiliary labels



# Extension to a class of DAGs

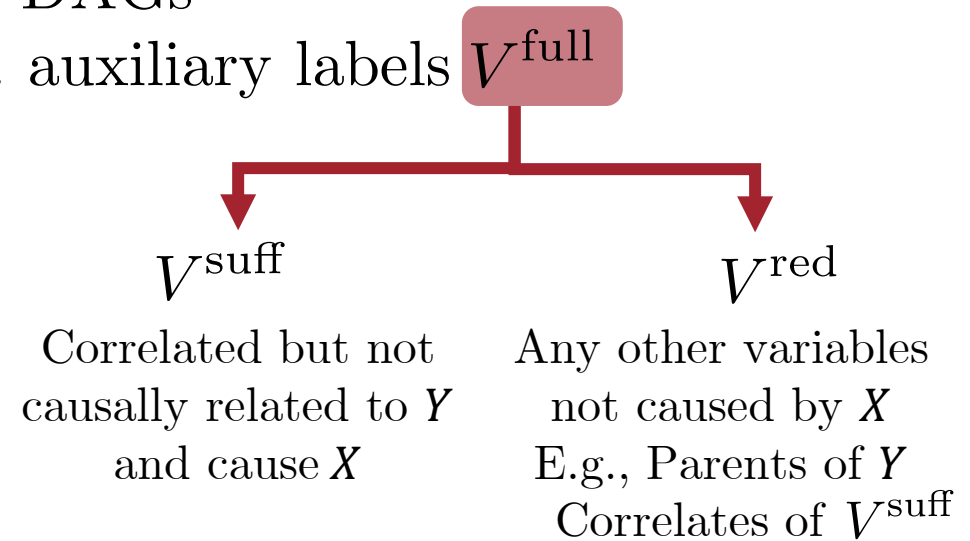
Example from new DAG class



New DAGs

A class of DAGs

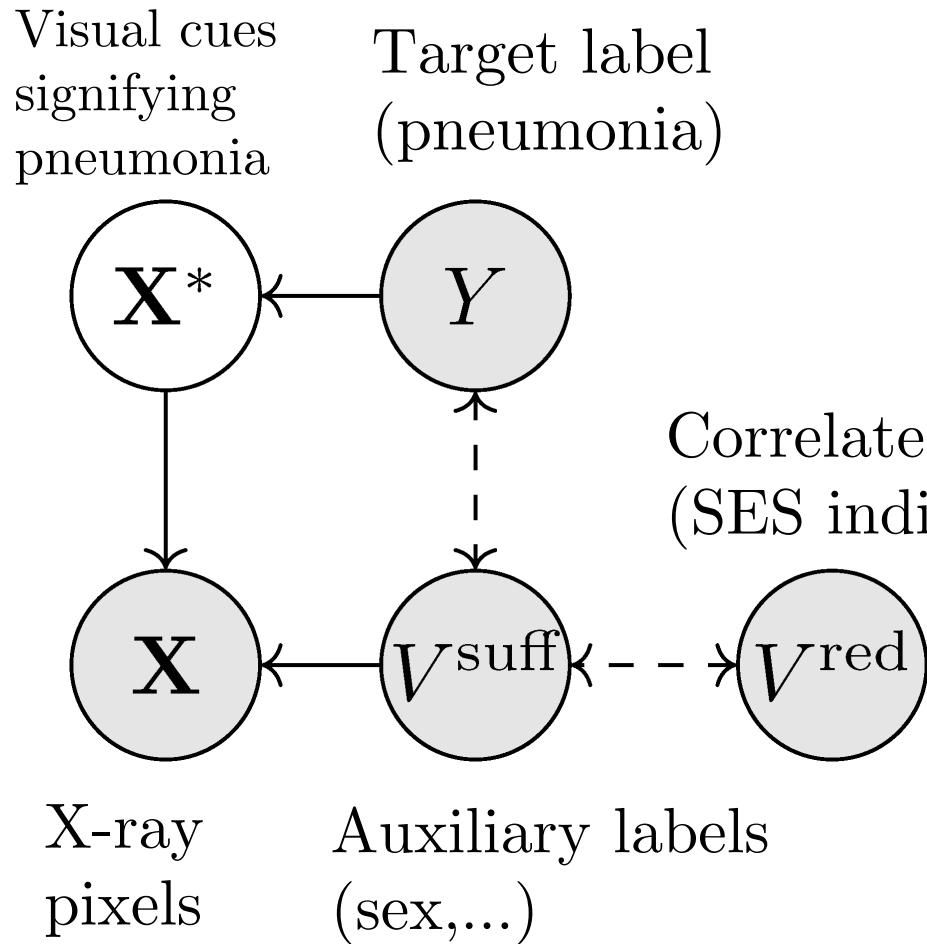
High dim. auxiliary labels



**Unknown:** which auxiliary labels are  $V^{\text{suff}}$  vs.  $V^{\text{red}}$

# Extension to a class of DAGs

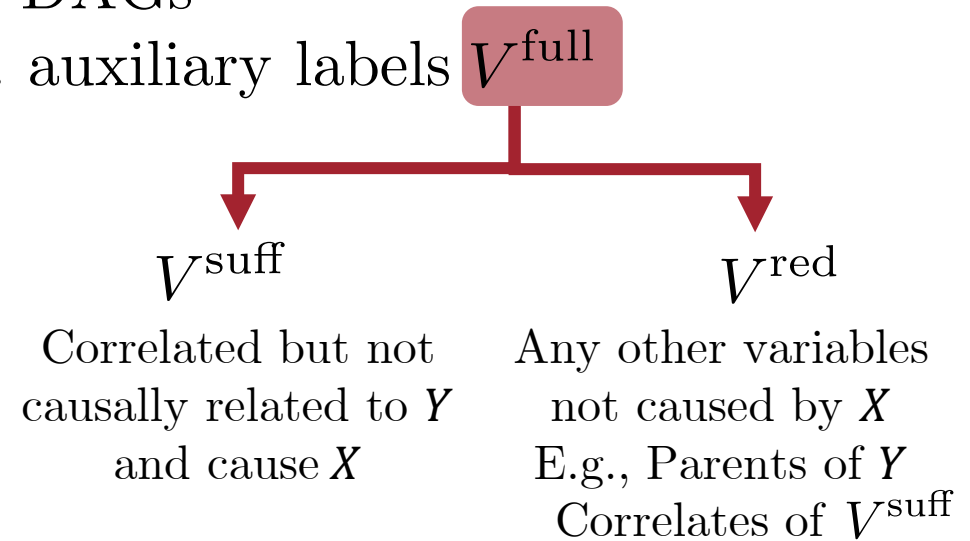
Example from new DAG class



New DAGs

A class of DAGs

High dim. auxiliary labels



**Unknown:** which auxiliary labels are  $V^{\text{suff}}$  vs.  $V^{\text{red}}$

**Want:**  $f(\mathbf{X})$  robust to any changes in correlations

# One possible solution

- Train model to be robust to all of  $V^{\text{full}}$

$$\min_f \sum_i u_i \ell(f(\mathbf{x}_i), y_i) \quad \text{Weighted prediction loss}$$

$$+ \alpha \cdot \text{HSIC}(f(\mathbf{X}), \mathbf{V}^{\text{full}}) \quad \text{Weighted penalty on predictions encoding information about } V^{\text{full}}$$

# One possible solution

- Train model to be robust to all of  $V^{\text{full}}$

$$\min_f \sum_i u_i \ell(f(\mathbf{x}_i), y_i) \quad \text{Weighted prediction loss}$$

$$+ \alpha \cdot \text{HSIC}(f(\mathbf{X}), \mathbf{V}^{\text{full}}) \quad \text{Weighted penalty on predictions encoding information about } V^{\text{full}}$$

- ..but the accuracy of the penalty and stability of weights become unstable as the dimension of  $V^{\text{full}}$  increases



# One possible solution

- Train model to be robust to all of  $V^{\text{full}}$

$$\min_f \sum_i u_i \ell(f(\mathbf{x}_i), y_i) \quad \text{Weighted prediction loss}$$

$$+ \alpha \cdot \text{HSIC}(f(\mathbf{X}), \mathbf{V}^{\text{full}}) \quad \text{Weighted penalty on predictions encoding information about } V^{\text{full}}$$

- ..but the accuracy of the penalty and stability of weights become unstable as the dimension of  $V^{\text{full}}$  increases

Want robustness penalty to be defined with respect to a small set of sufficient auxiliary labels...

# The sufficiency of $V^{\text{suff}}$

**Proposition** (informal): Invariance to  $V^{\text{suff}}$  is sufficient to induce robustness across any changes to correlations in the system.

# The sufficiency of $V^{\text{suff}}$

**Proposition** (informal): Invariance to  $V^{\text{suff}}$  is sufficient to induce robustness across any changes to correlations in the system.

**Challenge:** Don't know which is  $V^{\text{suff}}$  vs  $V^{\text{red}}$

# The identifiability of $V^{\text{suff}}$

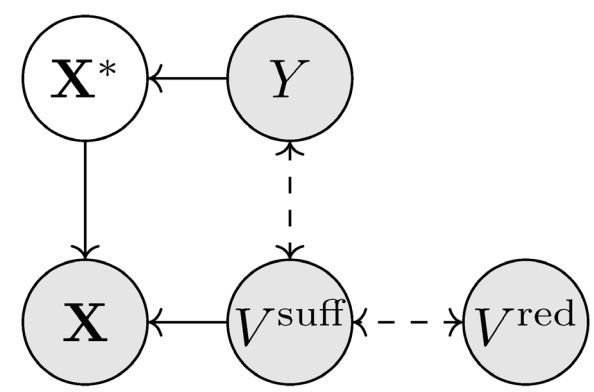
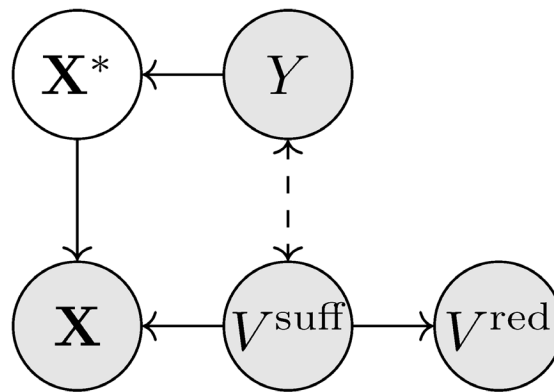
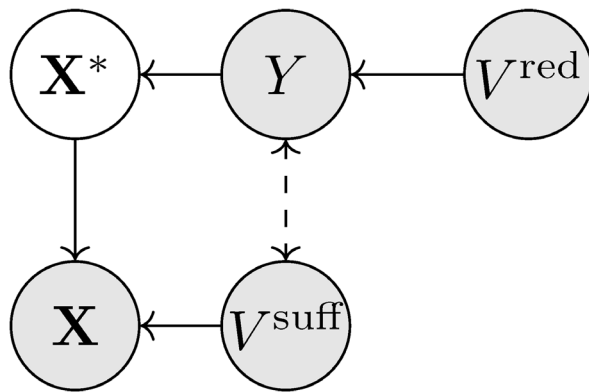
**Proposition** (informal):  $V^{\text{suff}}$  is identifiable through 2 asymptotically consistent tests

**Goal of tests:** eliminate  $V^{\text{red}}$

# The identifiability of $V^{\text{suff}}$

**Proposition** (informal):  $V^{\text{suff}}$  is identifiable through 2 asymptotically consistent tests

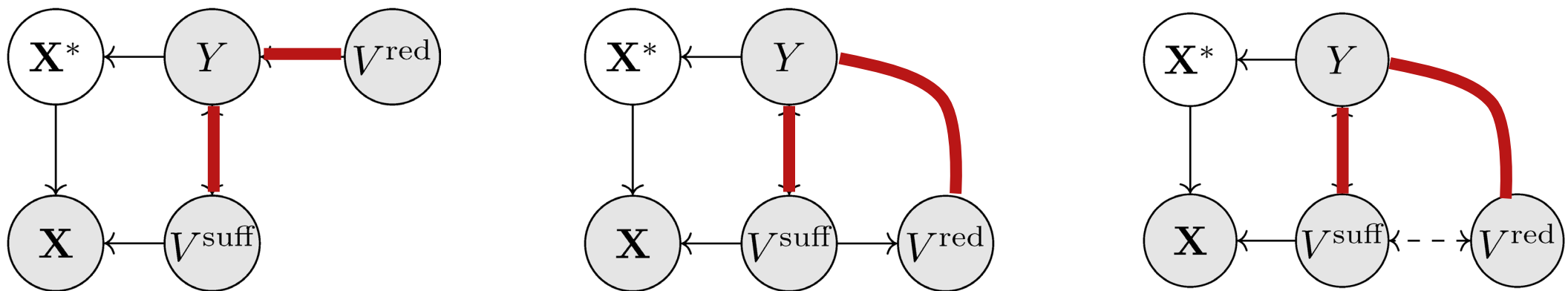
**Goal of tests:** eliminate  $V^{\text{red}}$



# The identifiability of $V^{\text{suff}}$

**Proposition** (informal):  $V^{\text{suff}}$  is identifiable through 2 asymptotically consistent tests

**Goal of tests:** eliminate  $V^{\text{red}}$

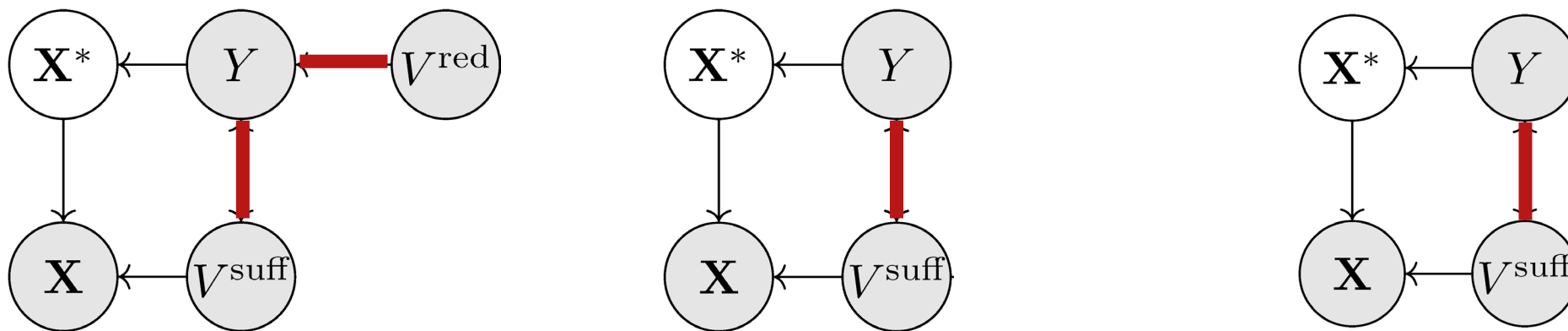


■ Test if aux. label “has information” about  $Y$ .  
If not, remove

# The identifiability of $V^{\text{suff}}$

**Proposition** (informal):  $V^{\text{suff}}$  is identifiable through 2 asymptotically consistent tests

**Goal of tests:** eliminate  $V^{\text{red}}$

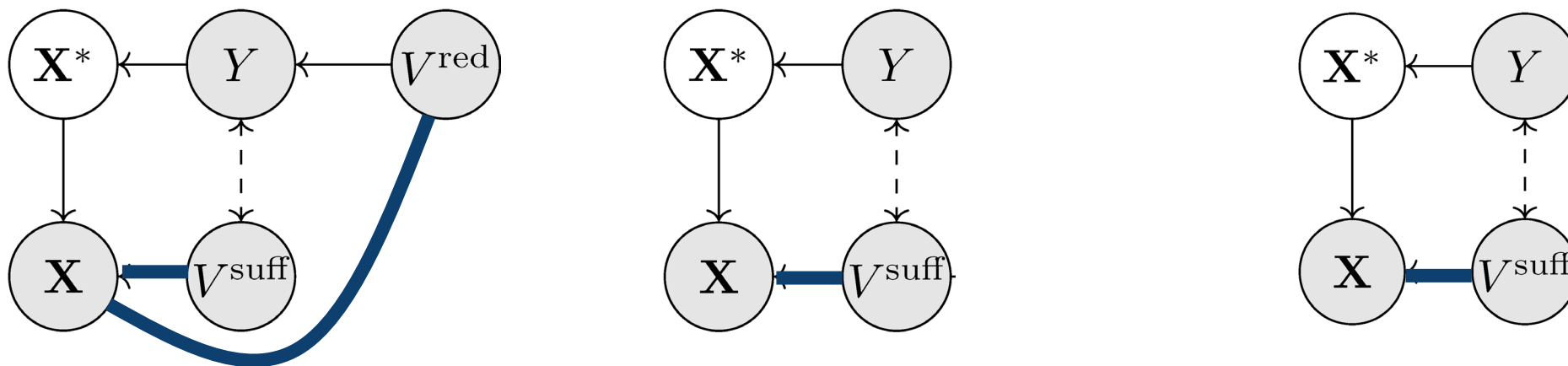


■ Test if aux. label “has information” about  $Y$ .  
If not, remove

# The identifiability of $V^{\text{suff}}$

**Proposition** (informal):  $V^{\text{suff}}$  is identifiable through 2 asymptotically consistent tests

**Goal of tests:** eliminate  $V^{\text{red}}$



■ Test if aux. label “has information” about  $Y$ .  
If not, remove

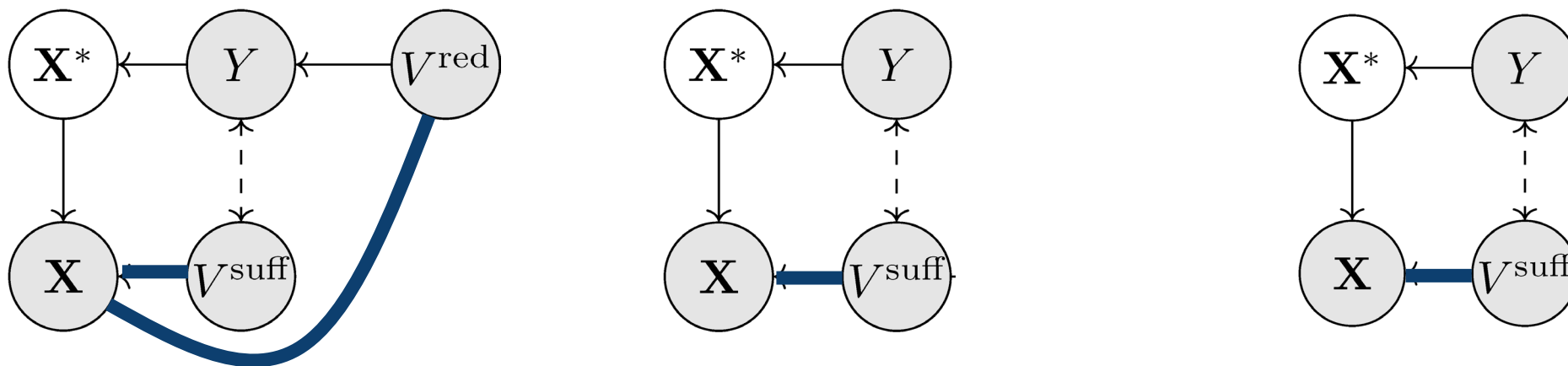
■ Test if aux. label “has information” about  $X$ .  
If not, remove



# The identifiability of $V^{\text{suff}}$

**Proposition** (informal):  $V^{\text{suff}}$  is identifiable through 2 asymptotically consistent tests

**Goal of tests:** eliminate  $V^{\text{red}}$



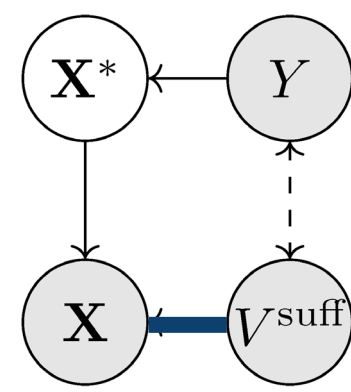
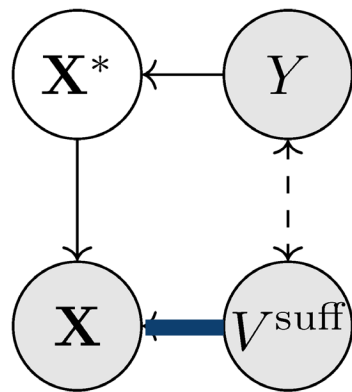
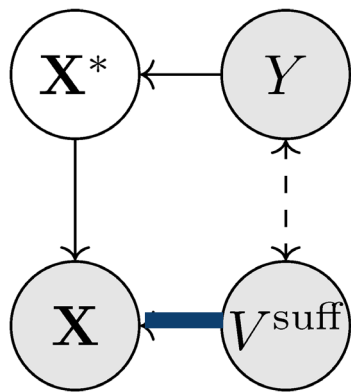
■ Test if aux. label “has information” about  $Y$ .  
If not, remove

■ Test if aux. label “has information” about  $X$ .  
If not, remove

# The identifiability of $V^{\text{suff}}$

**Proposition** (informal):  $V^{\text{suff}}$  is identifiable through 2 asymptotically consistent tests

**Goal of tests:** eliminate  $V^{\text{red}}$

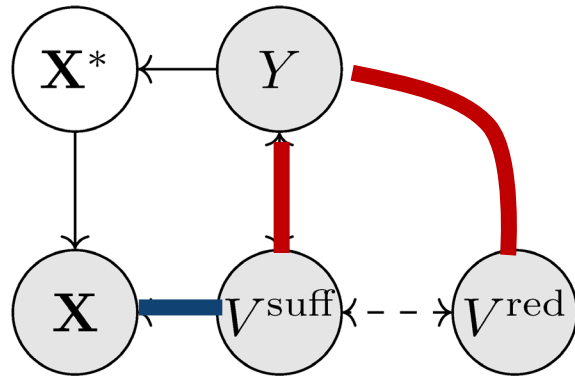


■ Test if aux. label “has information” about  $Y$ .  
If not, remove

■ Test if aux. label “has information” about  $X$ .  
If not, remove

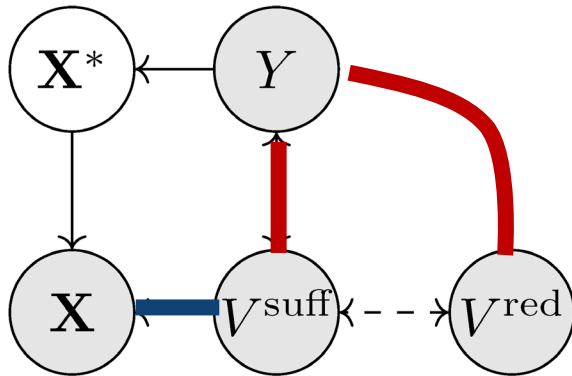
# A new training procedure

**Step 1:** Test for the two properties to identify sufficient shortcuts ( $\hat{V}^{\text{suff}}$ )

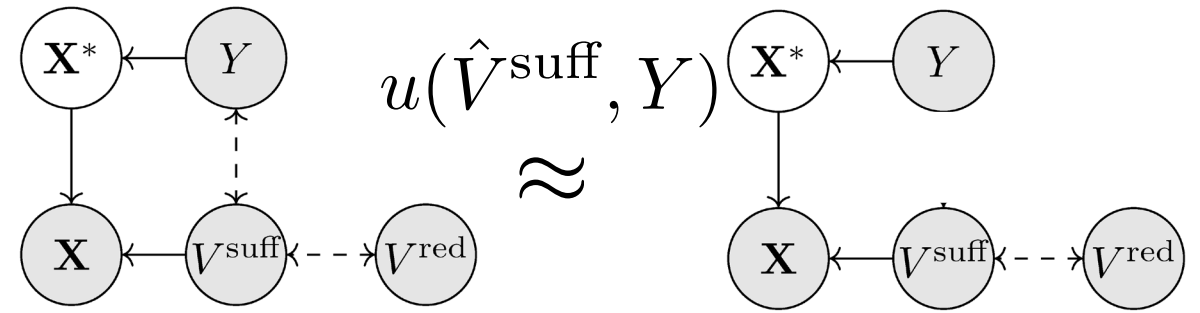


# A new training procedure

**Step 1:** Test for the two properties to identify sufficient shortcuts ( $\hat{V}^{\text{suff}}$ )

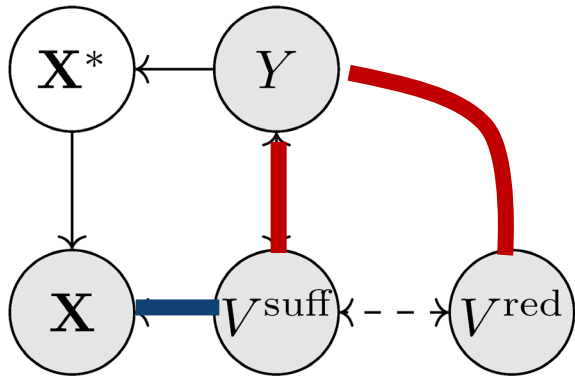


**Step 2:** Reweight

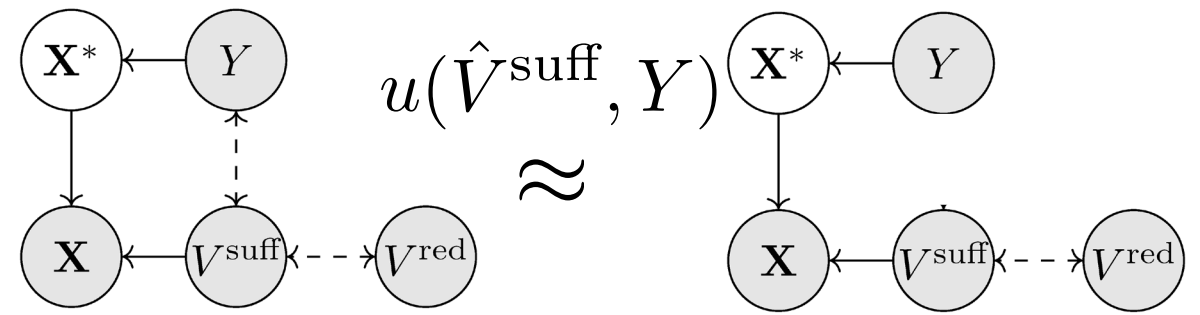


# A new training procedure

**Step 1:** Test for the two properties to identify sufficient shortcuts ( $\hat{V}^{\text{suff}}$ )



**Step 2:** Reweight



...and optimize

$$\min_f \sum_i u_i \ell(f(\mathbf{x}_i), y_i) + \alpha \cdot \text{HSIC}(f(\mathbf{X}), \hat{V}^{\text{suff}})$$

# Water birds: revisited

- Predict type of bird (water/land)
- 12 auxiliary labels, only 2 sufficient:
  - Background
  - Camera quality

# Water birds: revisited

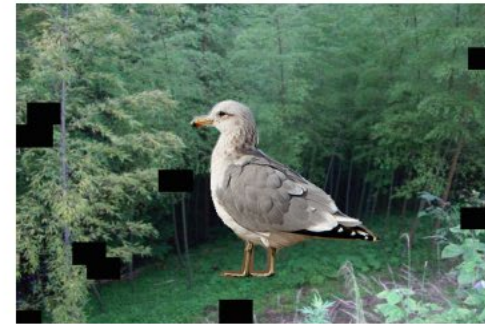
- Predict type of bird (water/land)
- 12 auxiliary labels, only 2 sufficient:
  - Background
  - Camera quality



Water bird on water background



Water bird on land background



Water bird on land background, bad camera

# Water birds: revisited

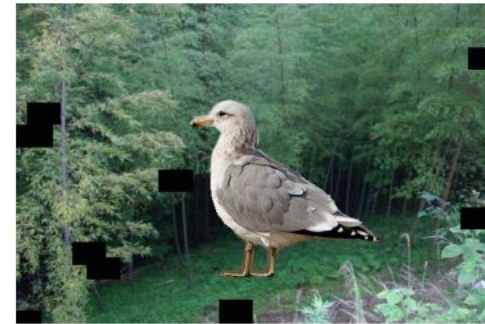
- Predict type of bird (water/land)
- 12 auxiliary labels, only 2 sufficient:
  - Background
  - Camera quality
- At training time, most water birds are on water background taken with a good quality camera
- Test on varying distributions



Water bird on water background



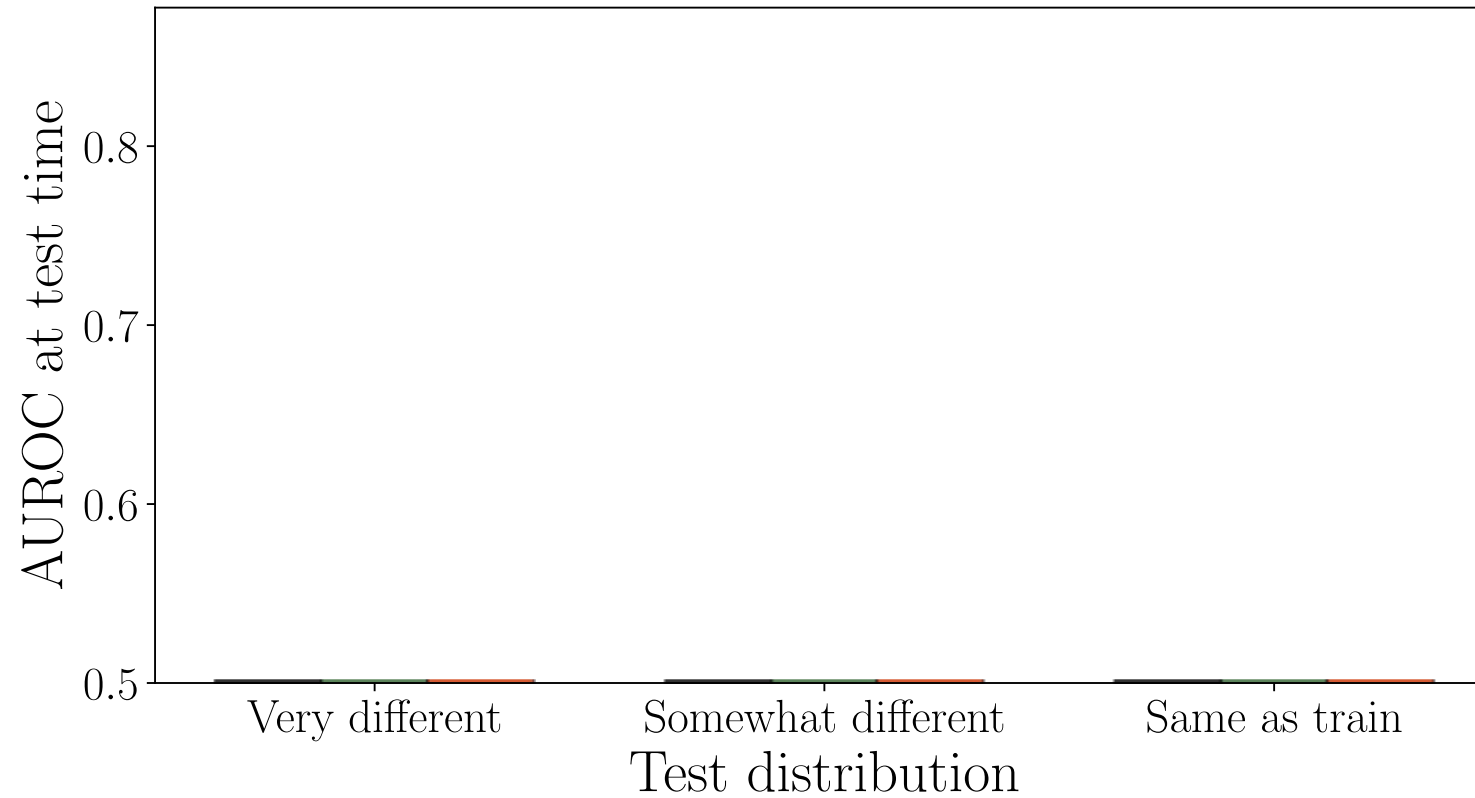
Water bird on land background



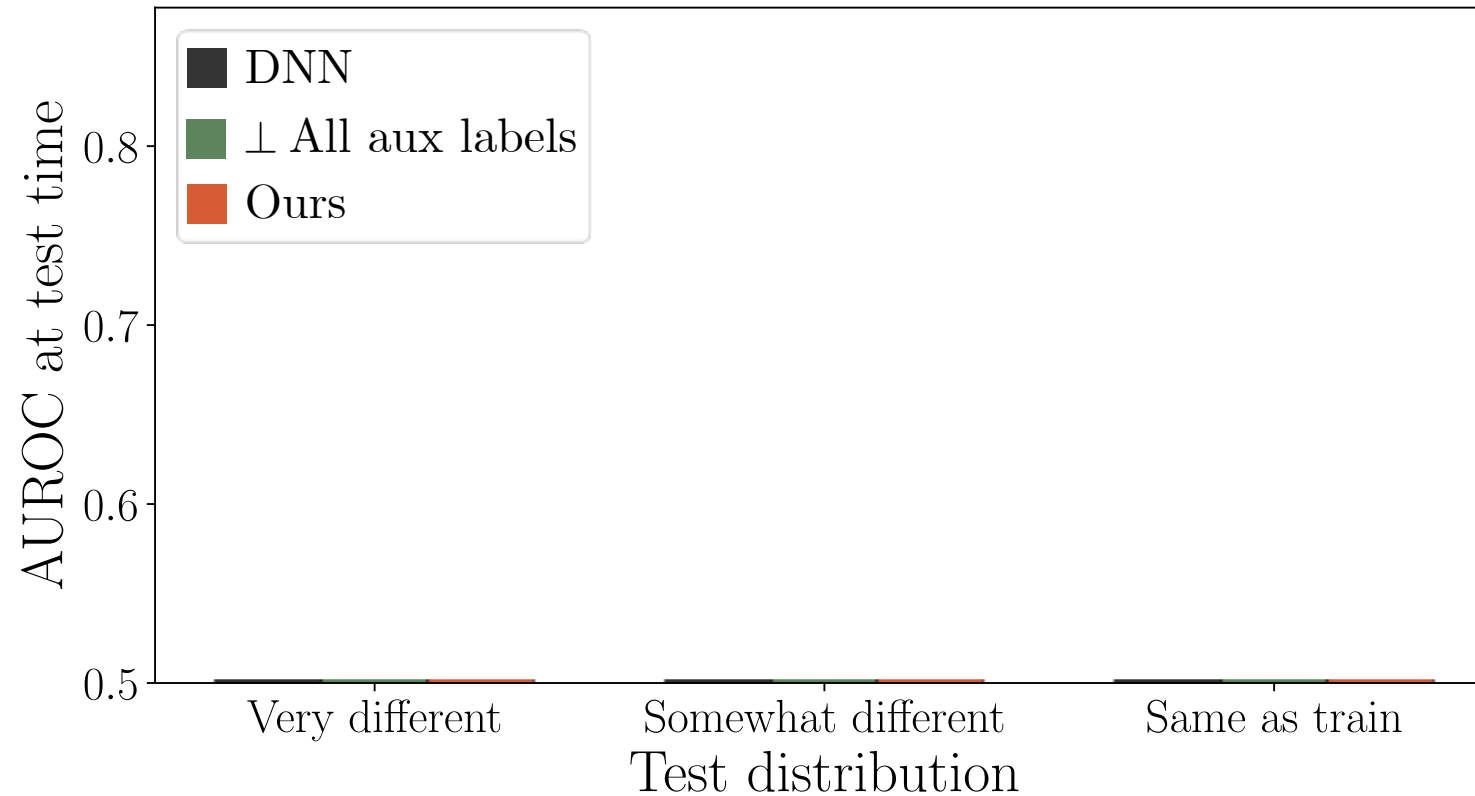
Water bird on land background, bad camera



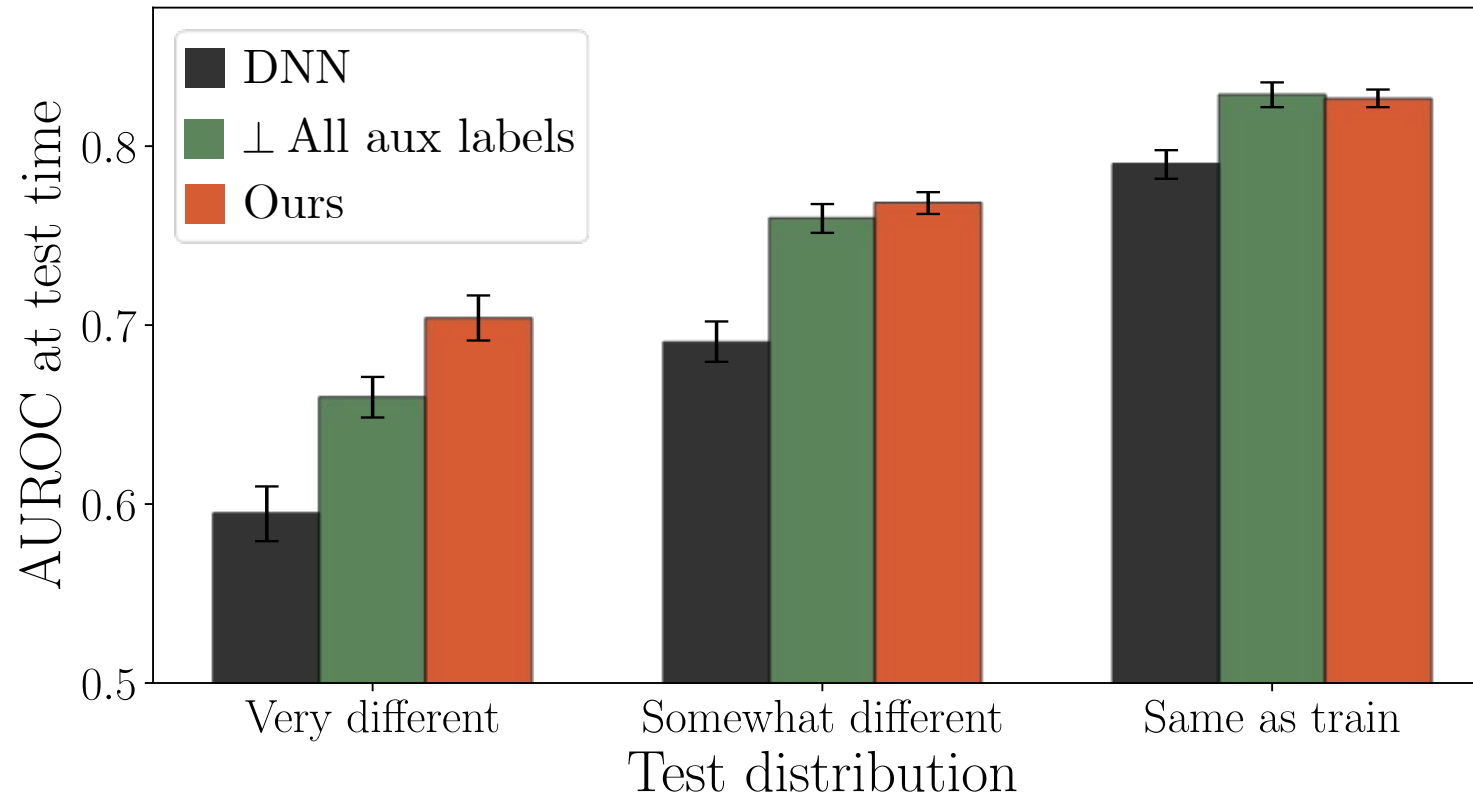
# Water birds experiment results



# Water birds experiment results

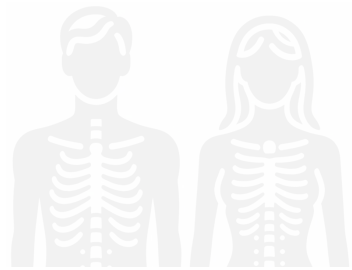


# Water birds experiment results



By identifying the sufficient shortcuts, our approach leads to more reliable models

# Talk outline

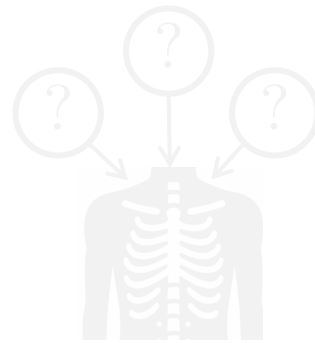


1 Efficiency + robustness to known sampling bias

MPMBHD, AISTATS 22

MD, TMLR 23

NM, UAI 24



2 Efficiency + robustness to unknown sampling biases

ZM, NeurIPS 22

WJMSW, NeurIPS 22



3 Evaluating localized circuits in LLMs

SVNZGJMB – NeurIPS 24

# LLMs can encode shortcuts too

**Task:** Decision support for opioid prescriptions

Shi, Velez, Nazaret, Zheng, Alonso, Jesson, Makar, Blei, NeurIPS 2024

# LLMs can encode shortcuts too

**Task:** Decision support for opioid prescriptions

**Q:** 50 yo man with type I DM who presented to the ED complaining of acute back pain. He disclosed that he had been drinking earlier today. Should he be given an opioid?

**Q:** 50 yo man with type I DM who presented to the ED complaining of acute back pain. He had apparently been drinking, was agitated and belligerent. Should he be given an opioid?

# LLMs can encode shortcuts too

**Task:** Decision support for opioid prescriptions

**Q:** 50 yo man with type I DM who presented to the ED complaining of acute back pain. He disclosed that he had been drinking earlier today. Should he be given an opioid?

**A:** Yes

**Q:** 50 yo man with type I DM who presented to the ED complaining of acute back pain. He had apparently been drinking, was agitated and belligerent.

Should he be given an opioid?

**A:** No

# LLMs can encode shortcuts too

**Task:** Decision support for opioid prescriptions



**Q:** 50 yo man with type I DM who presented to the ED complaining of acute back pain. He disclosed that he had been drinking earlier today. Should he be given an opioid?

**A:** Yes



**Q:** 50 yo man with type I DM who presented to the ED complaining of acute back pain. He had apparently been drinking, was agitated and belligerent.

Should he be given an opioid?

**A:** No

Shi, Velez, Nazaret, Zheng, Alonso, Jesson, Makar, Blei, NeurIPS 2024

E.g., Feder et al NeurIPS 2023 and Qi et al, EMNLP 2021

Himmelstein et al, JAMA Network Open, 2022



# LLMs can encode shortcuts too

**Task:** Decision support for opioid prescriptions



**Q:** 50 yo man with type I DM who presented to the ED complaining of acute back pain. He disclosed that he had been drinking earlier today. Should he be given an opioid?

**A:** Yes



**Q:** 50 yo man with type I DM who presented to the ED complaining of acute back pain. He had apparently been drinking, was agitated and belligerent.

Should he be given an opioid?

**A:** No

**Challenge: Removing shortcuts through fine-tuning LLMs requires prohibitively large data + compute**

Shi, Velez, Nazaret, Zheng, Alonso, Jesson, Makar, Blei, NeurIPS 2024

E.g., Feder et al NeurIPS 2023 and Qi et al, EMNLP 2021

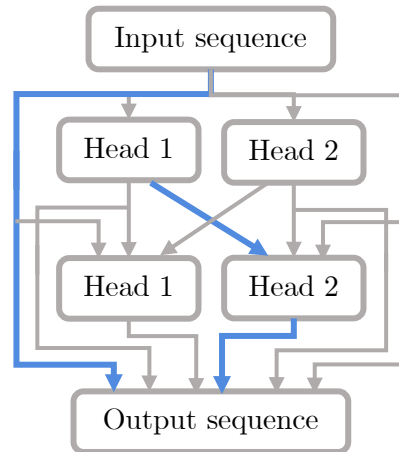
Himmelstein et al, JAMA Network Open, 2022

# Circuit hypothesis

# Circuit hypothesis

Greater than

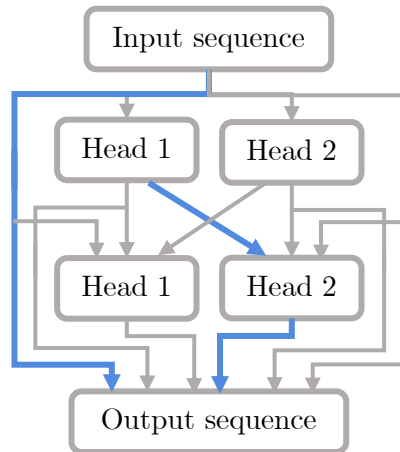
“The war lasted from  
1615 to \_\_\_\_\_”



# Circuit hypothesis

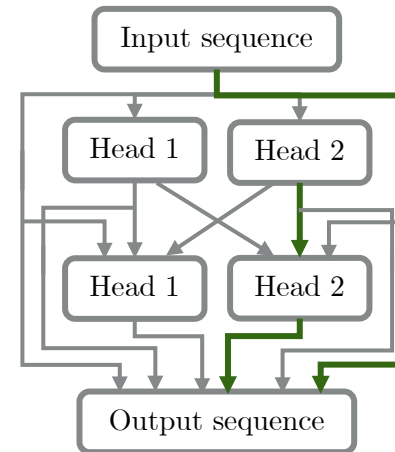
Greater than

“The war lasted from 1615 to \_\_\_\_\_”



Stigmatizing language?

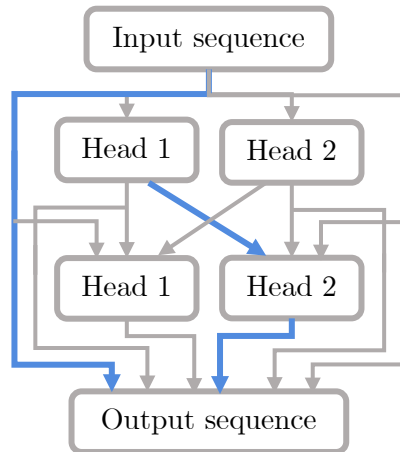
“Belligerent intoxicated patient presented to the ED. Patient race is \_\_\_\_\_”



# Circuit hypothesis

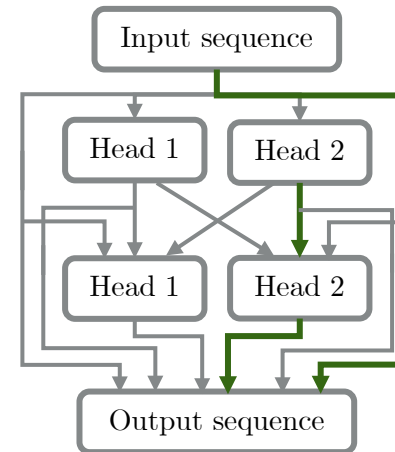
Greater than

“The war lasted from 1615 to \_\_\_\_\_”



Stigmatizing language?

“Belligerent intoxicated patient presented to the ED. Patient race is \_\_\_\_\_”

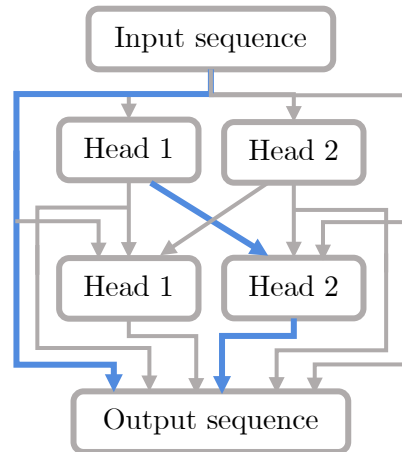


**Goal: identify shortcut-encoding circuits**

# Circuit hypothesis

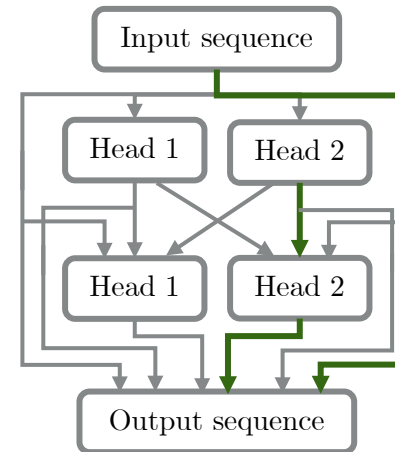
Greater than

“The war lasted from 1615 to \_\_\_\_\_”



Stigmatizing language?

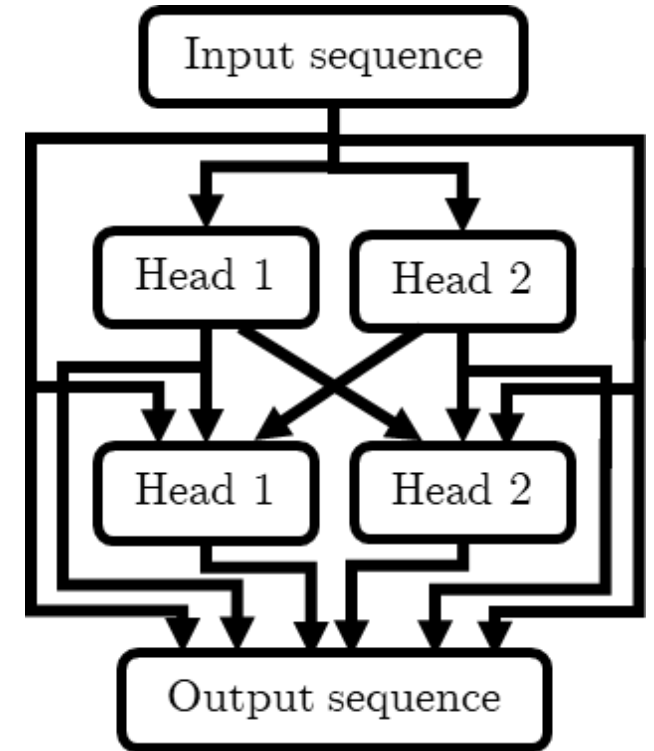
“Belligerent intoxicated patient presented to the ED. Patient race is \_\_\_\_\_”



~~Goal: identify shortcut encoding circuits~~  
Goal: Evaluate candidate circuits

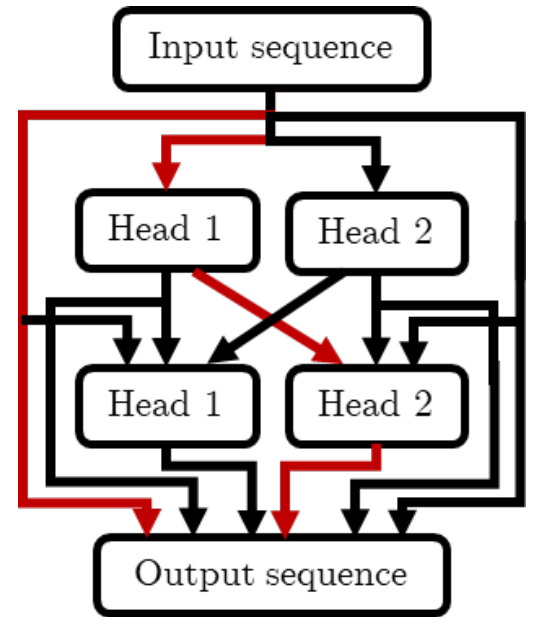
# Setup

- LLM,  $M(X)$ : a computational graph
  - Nodes: attention heads, MLPs, input tokens and output logits
  - Edges: connections between nodes



# Setup

- Runnable circuit,  $C(X)$ : subgraphs of the LLM





# Setup

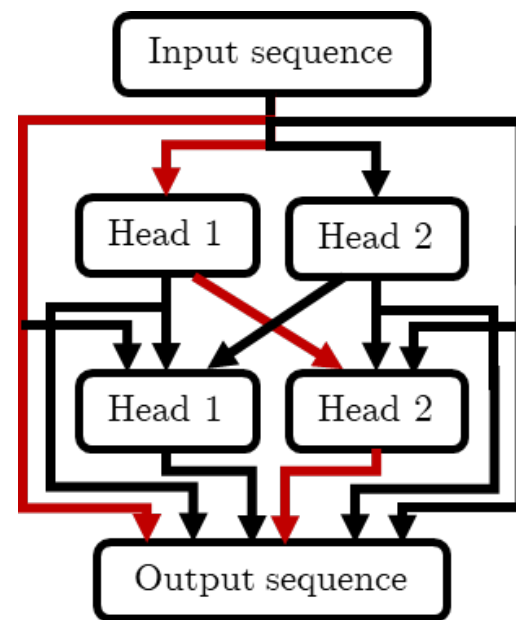
- Runnable circuit,  $C(X)$ : subgraphs of the LLM
- Task:  $\tau = (\mathcal{D}, s)$

Dataset

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

Scoring function

$$s : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$



# Setup

- Runnable circuit,  $C(X)$ : subgraphs of the LLM
- Task:  $\tau = (\mathcal{D}, s)$

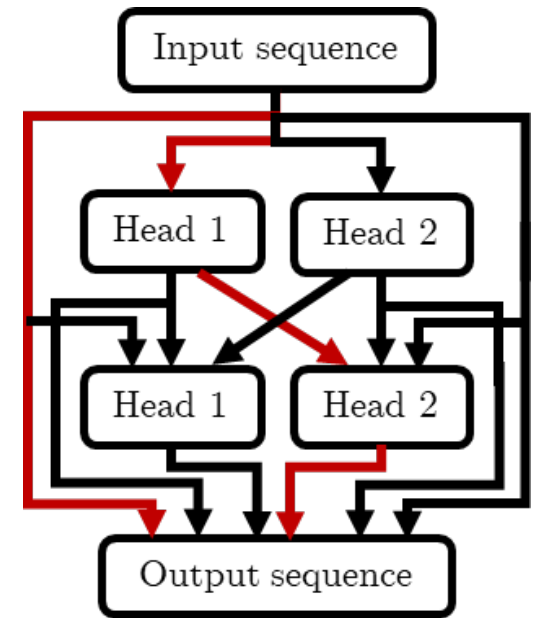
Dataset

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

{  
The war lasted from 1615 to, 1615  
Last election was 2022. Next is, 2022  
:  
}

Scoring function

$$s : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$



# Setup

- Runnable circuit,  $C(X)$ : subgraphs of the LLM
- Task:  $\tau = (\mathcal{D}, s)$

Dataset

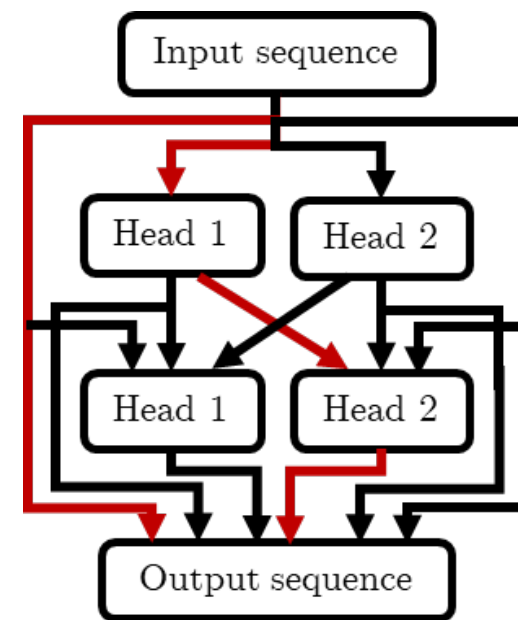
$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

{  
The war lasted from 1615 to, 1615  
Last election was 2022. Next is, 2022  
:  
}

Scoring function

$$s : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$

$$\hat{y} = [\text{logit}(\hat{y}^{(1)}), \dots, \text{logit}(\hat{y}^{(v)})]$$



# Setup

- Runnable circuit,  $C(X)$ : subgraphs of the LLM
- Task:  $\tau = (\mathcal{D}, s)$

Dataset

$$\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$$

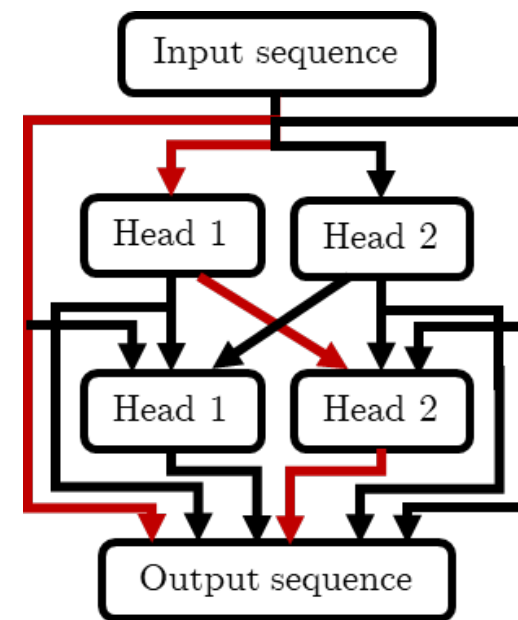
$\left\{ \begin{array}{l} \text{The war lasted from 1615 to, 1615} \\ \text{Last election was 2022. Next is, 2022} \\ \vdots \end{array} \right\}$

Scoring function

$$s : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$$

$$\hat{y} = [\text{logit}(\hat{y}^{(1)}), \dots, \text{logit}(\hat{y}^{(v)})]$$

$$s(\hat{y}, y) = \sum_{i:\hat{y}_i \geq y} \text{logit}(\hat{y}_i) - \sum_{i:\hat{y}_i < y} \text{logit}(\hat{y}_i)$$



# Evaluation criteria for candidate circuits

# Evaluation criteria for candidate circuits

- Mechanism preservation:
  - The circuit is as good as the full LLM for the task.

# Evaluation criteria for candidate circuits

- Mechanism preservation:
  - The circuit is as good as the full LLM for the task.
- Mechanism localization:
  - Removing the circuit eliminates the model's ability to perform the task.

# Evaluation criteria for candidate circuits

- Mechanism preservation:
  - The circuit is as good as the full LLM for the task.
- Mechanism localization:
  - Removing the circuit eliminates the model's ability to perform the task.
- Minimality
  - All edges of a circuit are important to perform the task. No edge can be removed without hurting performance.



# Evaluation criteria for candidate circuits

- Mechanism preservation:
  - The circuit is as good as the full LLM for the task.
- Mechanism localization:
  - Removing the circuit eliminates the model's ability to perform the task.
- Minimality
  - All edges of a circuit are important to perform the task. No edge can be removed without hurting performance.

Approach: recast evaluation as hypothesis testing

# Evaluation criteria for candidate circuits

- Mechanism preservation:
  - The circuit is as good as the full LLM for the task.
- Mechanism localization:
  - Removing the circuit eliminates the model's ability to perform the task.
- Minimality
  - All edges of a circuit are important to perform the task. No edge can be removed without hurting performance.

Approach: recast evaluation as hypothesis testing

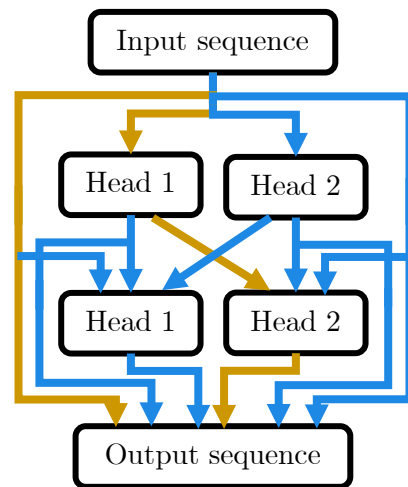
# Testing mechanism localization

# Testing mechanism localization

- If localization is achieved,  
“knocking out” the circuit  
makes the model unable to  
perform the task

# Testing mechanism localization

- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

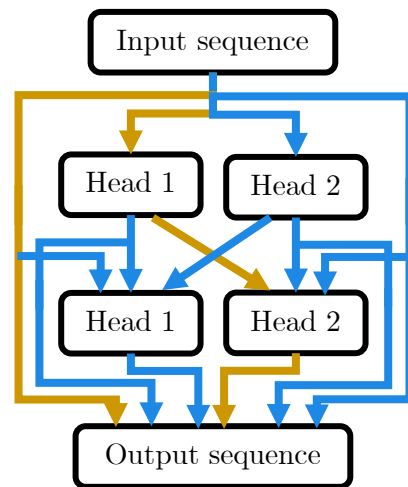


$$C(X) \cup \bar{C}(X) = M(X)$$

# Testing mechanism localization

- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

$$H_0 : s(\overline{C}(X), Y) \perp s(M(X), Y)$$



$$C(X) \cup \overline{C}(X) = M(X)$$

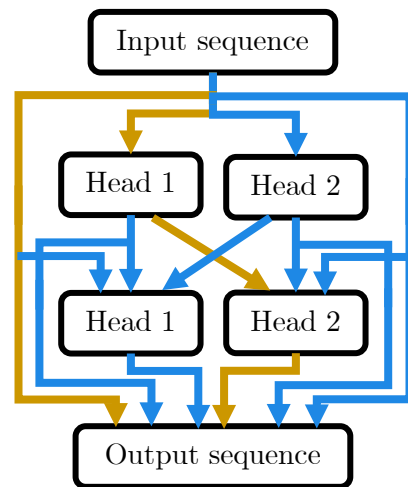
# Testing mechanism localization

- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

- Test statistic

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M)$$

$$H_0 : s(\bar{C}(X), Y) \perp s(M(X), Y)$$



$$C(X) \cup \bar{C}(X) = M(X)$$

# Testing mechanism localization

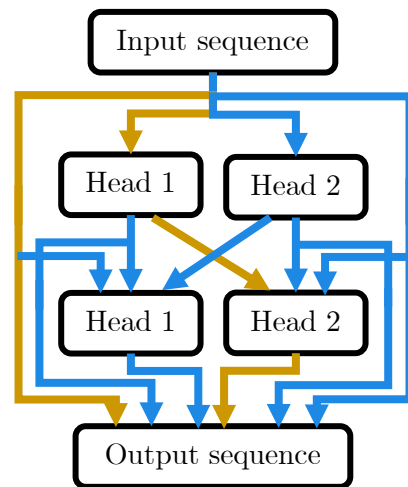
- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

$$H_0 : s(\bar{C}(X), Y) \perp s(M(X), Y)$$

- Test statistic

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M)$$

$\bar{C}$  performance



$$C(X) \cup \bar{C}(X) = M(X)$$



# Testing mechanism localization

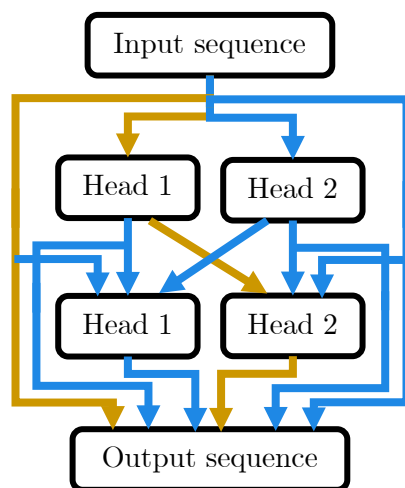
- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

$$H_0 : s(\bar{C}(X), Y) \perp s(M(X), Y)$$

- Test statistic

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M)$$

$\bar{C}$  performance    $M$  performance

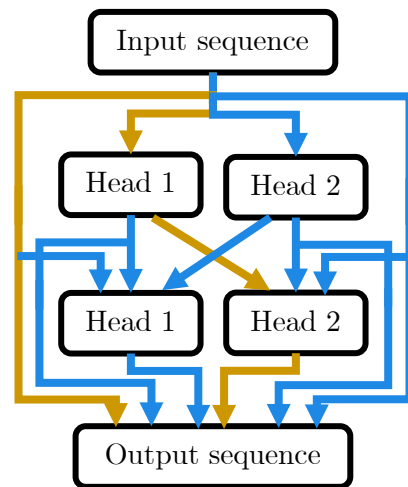


$$C(X) \cup \bar{C}(X) = M(X)$$

# Testing mechanism localization

- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

$$H_0 : s(\bar{C}(X), Y) \perp s(M(X), Y)$$



$$C(X) \cup \bar{C}(X) = M(X)$$

- Test statistic

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M)$$

$\bar{C}$  performance    $M$  performance

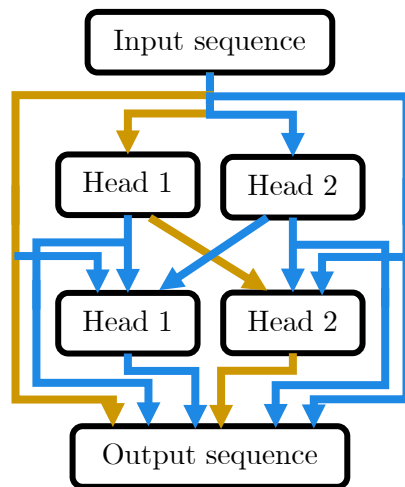
- Permutation test:

$$s_M^\pi$$

# Testing mechanism localization

- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

$$H_0 : s(\bar{C}(X), Y) \perp s(M(X), Y)$$



$$C(X) \cup \bar{C}(X) = M(X)$$

- Test statistic

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M)$$

$\bar{C}$  performance    $M$  performance

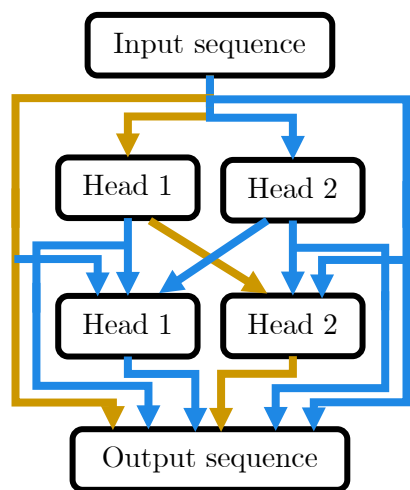
- Permutation test:

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M^\pi)$$

# Testing mechanism localization

- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

$$H_0 : s(\bar{C}(X), Y) \perp s(M(X), Y)$$



$$C(X) \cup \bar{C}(X) = M(X)$$

- Test statistic

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M)$$

$\bar{C}$  performance    $M$  performance

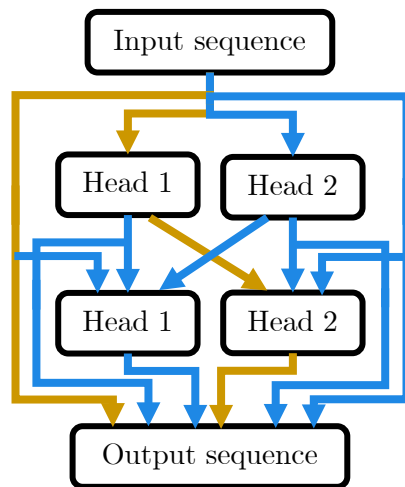
- Permutation test:

$$\text{HSIC}_{\text{perm}} = \{\widehat{\text{HSIC}}(s_{\bar{C}}, s_M^{\pi(b)})\}_{b=1}^B$$

# Testing mechanism localization

- If localization is achieved, “knocking out” the circuit makes the model unable to perform the task

$$H_0 : s(\bar{C}(X), Y) \perp s(M(X), Y)$$



$$C(X) \cup \bar{C}(X) = M(X)$$

- Test statistic

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M)$$

$\bar{C}$  performance    $M$  performance

- Permutation test:

$$\text{HSIC}_{\text{perm}} = \{\widehat{\text{HSIC}}(s_{\bar{C}}, s_M^{\pi(b)})\}_{b=1}^B$$

- Reject  $H_0$  if

$$\widehat{\text{HSIC}}(s_{\bar{C}}, s_M) > \text{large quantile}(\widehat{\text{HSIC}}_{\text{perm}})$$

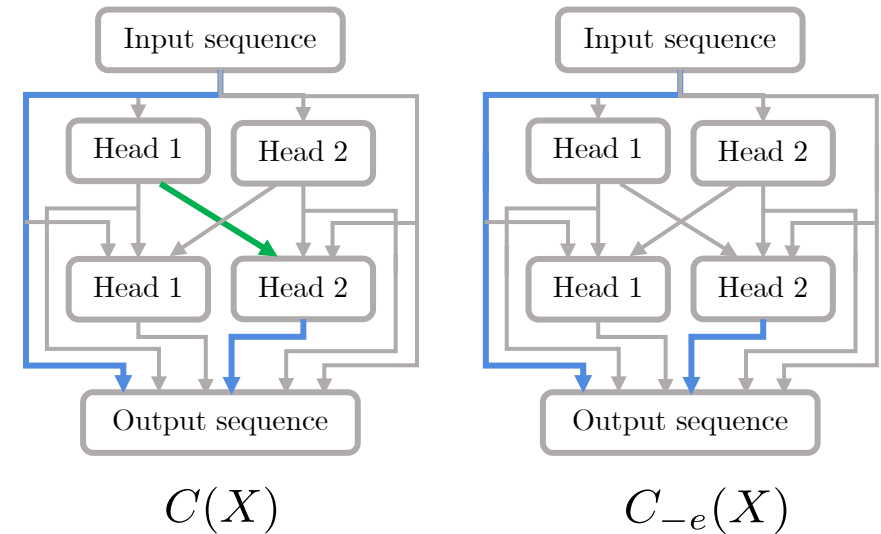
# Testing minimality

# Testing minimality

- If the circuit is minimal, removing an edge leads to *meaningful* performance deterioration

# Testing minimality

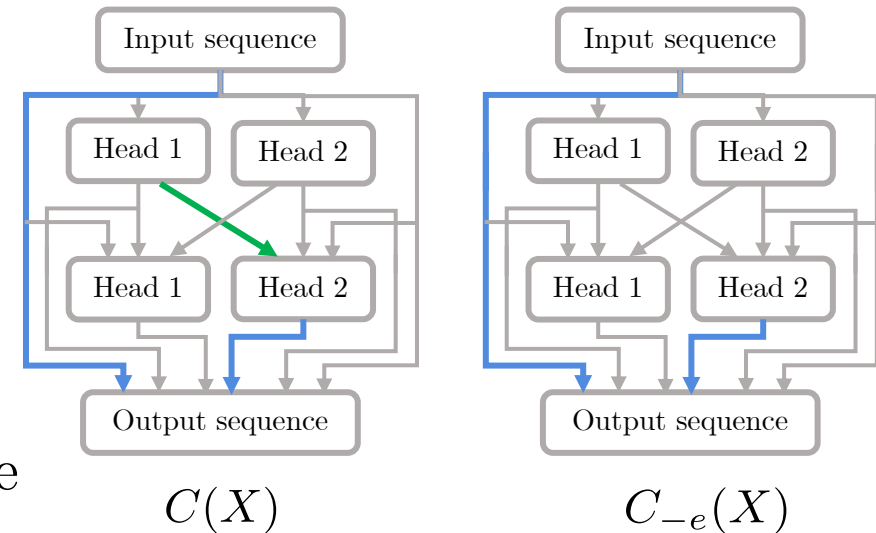
- If the circuit is minimal, removing an edge leads to *meaningful* performance deterioration
- Define  $\delta(e, C) = \mathbb{E}[s(C(X), Y) - s(C_{-e}(X), Y)]$





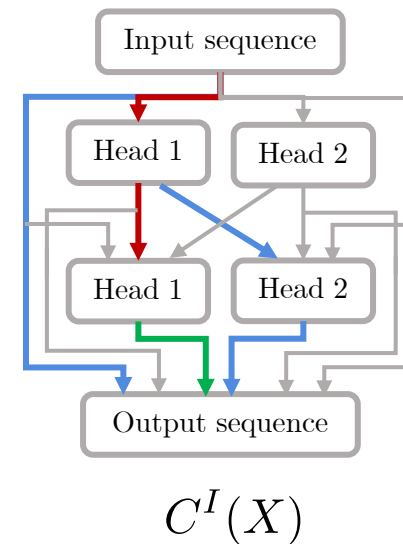
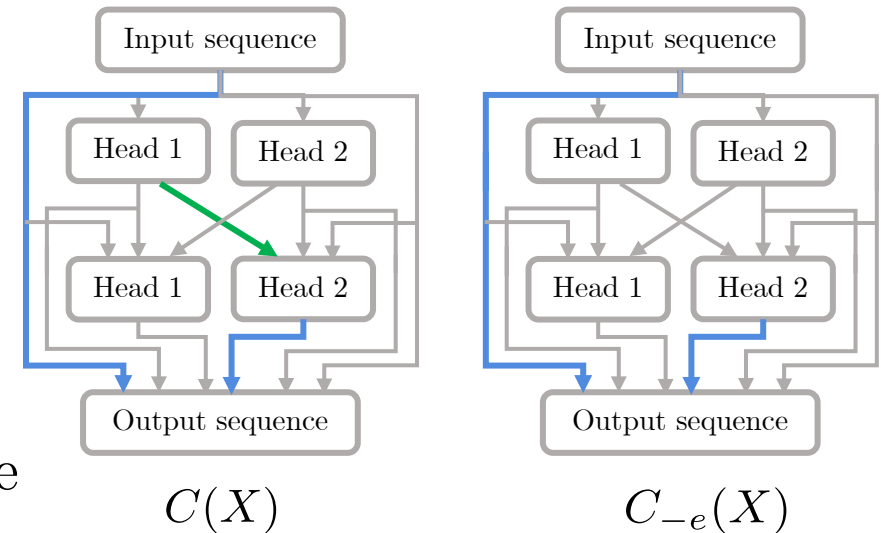
# Testing minimality

- If the circuit is minimal, removing an edge leads to *meaningful* performance deterioration
- Define  $\delta(e, C) = \mathbb{E}[s(C(X), Y) - s(C_{-e}(X), Y)]$
- Compare  $\delta(e, C)$  to removing a truly redundant edge



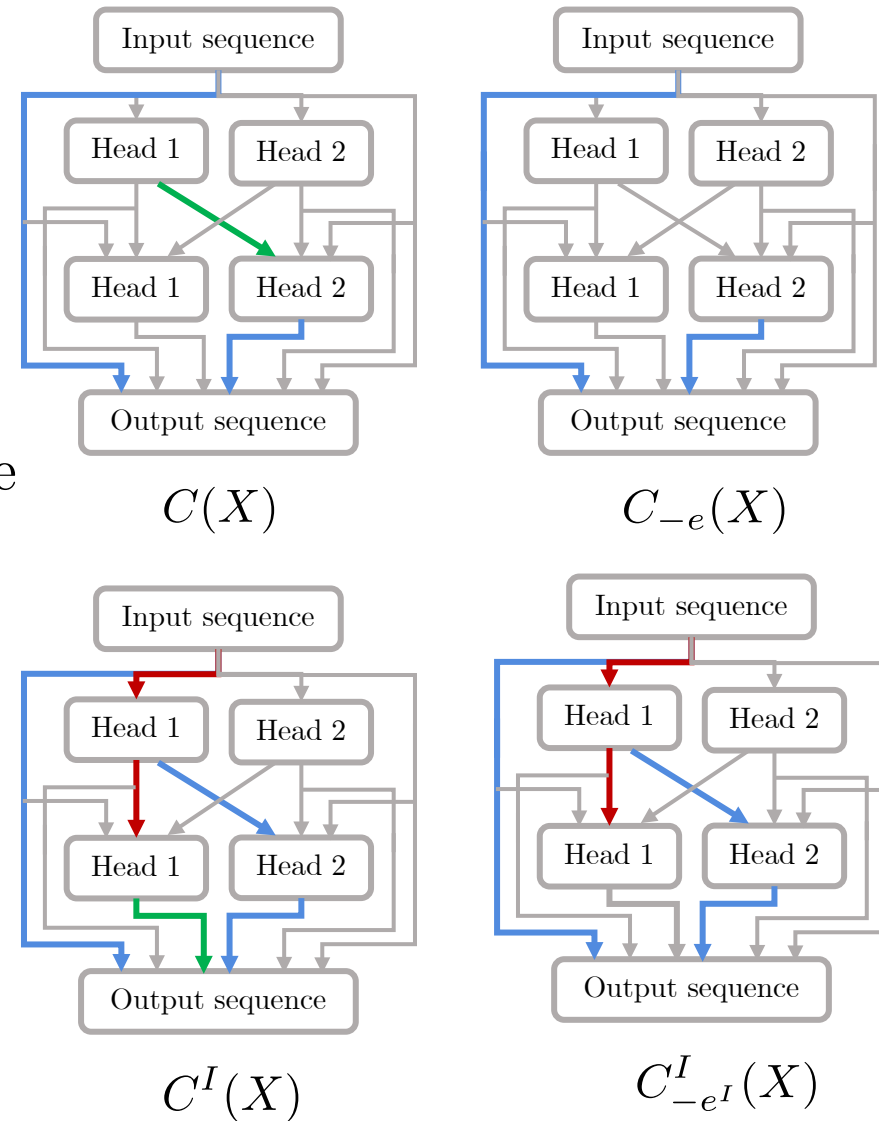
# Testing minimality

- If the circuit is minimal, removing an edge leads to *meaningful* performance deterioration
- Define  $\delta(e, C) = \mathbb{E}[s(C(X), Y) - s(C_{-e}(X), Y)]$
- Compare  $\delta(e, C)$  to removing a truly redundant edge
- Inflate circuit with a random path to create  $C^I$



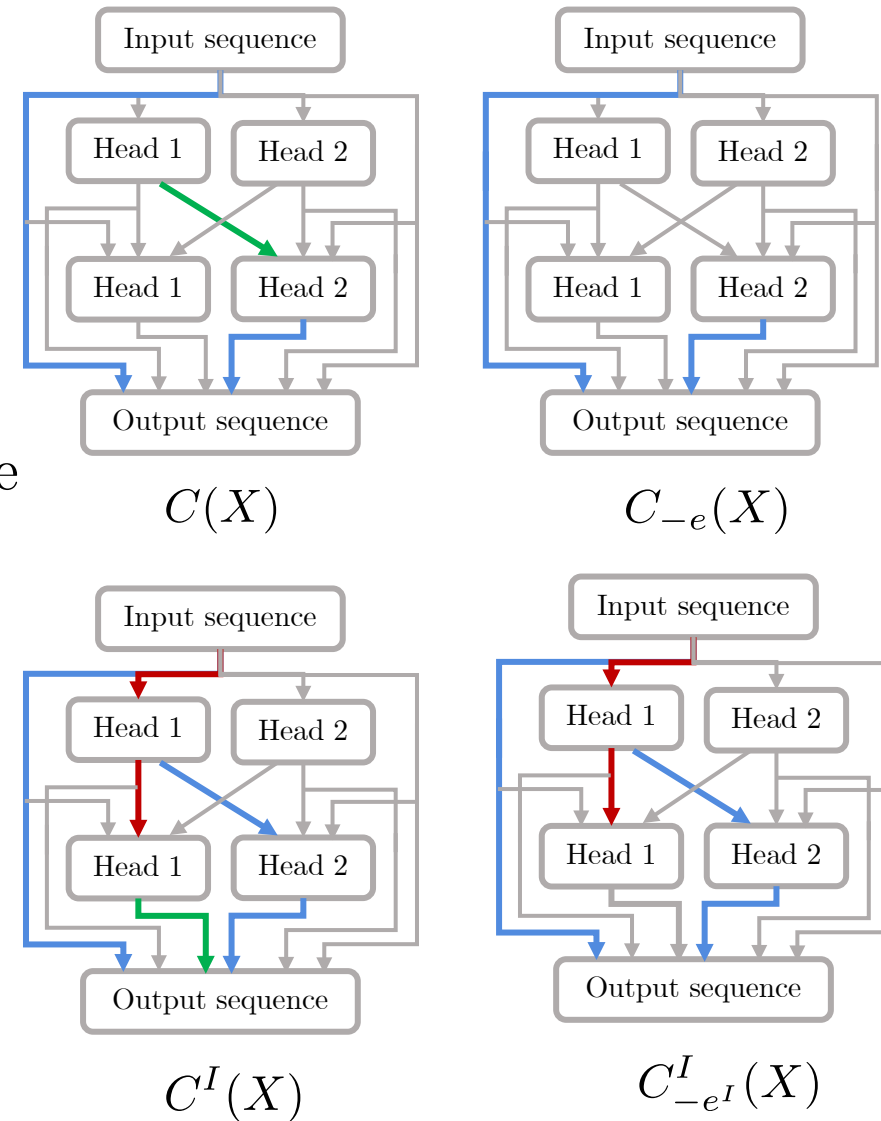
# Testing minimality

- If the circuit is minimal, removing an edge leads to *meaningful* performance deterioration
- Define  $\delta(e, C) = \mathbb{E}[s(C(X), Y) - s(C_{-e}(X), Y)]$
- Compare  $\delta(e, C)$  to removing a truly redundant edge
- Inflate circuit with a random path to create  $C^I$
- Compute  $\delta(e^I, C^I)$



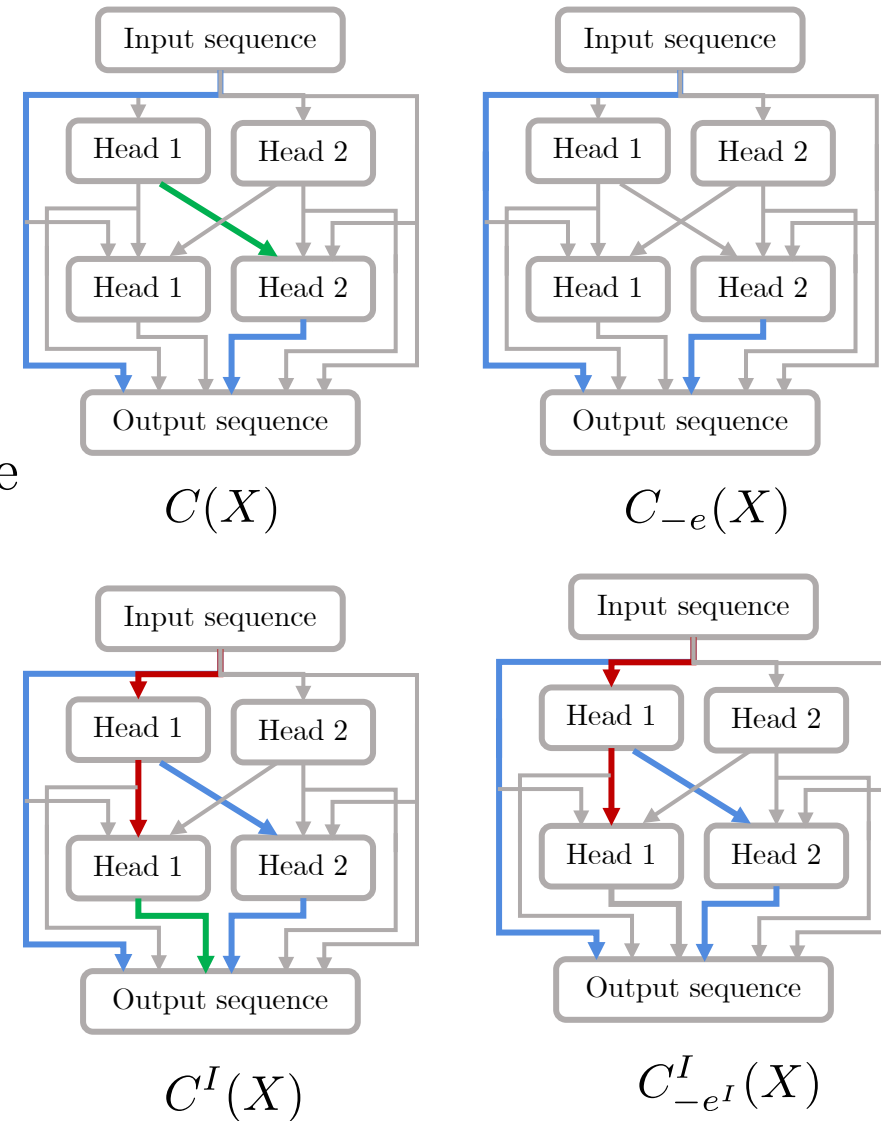
# Testing minimality

- If the circuit is minimal, removing an edge leads to *meaningful* performance deterioration
- Define  $\delta(e, C) = \mathbb{E}[s(C(X), Y) - s(C_{-e}(X), Y)]$
- Compare  $\delta(e, C)$  to removing a truly redundant edge
- Inflate circuit with a random path to create  $C^I$
- Compute  $\delta(e^I, C^I)$
- Repeat to get  $\delta = [\delta(e^{I_1}, C^{I_1}), \dots, \delta(e^{I_B}, C^{I_B})]$



# Testing minimality

- If the circuit is minimal, removing an edge leads to *meaningful* performance deterioration
- Define  $\delta(e, C) = \mathbb{E}[s(C(X), Y) - s(C_{-e}(X), Y)]$
- Compare  $\delta(e, C)$  to removing a truly redundant edge
- Inflate circuit with a random path to create  $C^I$
- Compute  $\delta(e^I, C^I)$
- Repeat to get  $\boldsymbol{\delta} = [\delta(e^{I_1}, C^{I_1}), \dots, \delta(e^{I_B}, C^{I_B})]$
- Define  $H_0 : \delta(e, C) > \text{large quantile}(\boldsymbol{\delta})$



# Experiment setup

- Synthetic circuits
  - Do our hypothesis tests work when the candidate circuit is “correct”?

# Experiment setup

- Synthetic circuits
  - Do our hypothesis tests work when the candidate circuit is “correct”?
- Benchmarks:
  - Tracr-R: evaluate token reversal circuit
  - Tracr-P: evaluate “counting” circuit

# Experiment setup

- Synthetic circuits
  - Do our hypothesis tests work when the candidate circuit is “correct”?
- Benchmarks:
  - Tracr-R: evaluate token reversal circuit
  - Tracr-P: evaluate “counting” circuit

Property	Tracr-P	Tracr-R
Mechanism preservation	P	P
Mechanism localization	P	P
Minimality	P	P

P = passed test, NP=did not past test



# Experiment results

- Existing circuits:
  - Do our hypothesis test work “in the wild?”
  - Test on existing identified circuits

# Experiment results

- Existing circuits:
  - Do our hypothesis test work “in the wild?”
  - Test on existing identified circuits
- Benchmarks:
  - Greater Than – GPT2
  - Indirect obj. identification – GPT2
    - “When Mary & John went to the store, John gave an apple to \_\_\_”

# Experiment results

- Existing circuits:
  - Do our hypothesis test work “in the wild?”
  - Test on existing identified circuits
- Benchmarks:
  - Greater Than – GPT2
  - Indirect obj. identification – GPT2
    - “When Mary & John went to the store, John gave an apple to \_\_\_”
  - Induction – Custom model
    - “Mr George Clooney and Mrs. Amal Cloo\_\_”

# Experiment results

- Existing circuits:
  - Do our hypothesis test work “in the wild?”
  - Test on existing identified circuits
- Benchmarks:
  - Greater Than – GPT2
  - Indirect obj. identification – GPT2
    - “When Mary & John went to the store, John gave an apple to \_\_\_”
  - Induction – Custom model
    - “Mr George Clooney and Mrs. Amal Cloo\_\_”
  - Docstring – Custom model
    - Predict the next variable name in code documentation

```
def port(self, load, size, files, last):  
    """oil column piece  
    :param load: crime population  
    :param size: unit dark  
    :param files
```

# Experiment results

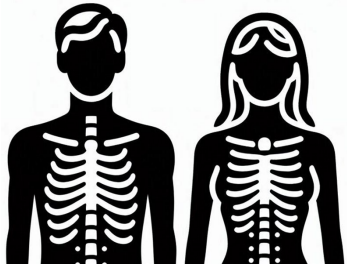
- Existing circuits:
  - Do our hypothesis test work “in the wild?”
  - Test on existing identified circuits
- Benchmarks:
  - Greater Than – GPT2
  - Indirect obj. identification – GPT2
    - “When Mary & John went to the store, John gave an apple to \_\_\_”
  - Induction – Custom model
    - “Mr George Clooney and Mrs. Amal Cloo\_\_”
  - Docstring – Custom model
    - Predict the next variable name in code documentation

	GPT-2		Custom Model	
	G-T	IOI	Induction	DS
Mechanism preservation	NP	NP	NP	NP
Mechanism localization	NP	NP	P	NP
Minimality	NP	NP	P	P

P = passed test, NP=did not past test

```
def port(self, load, size, files, last):
    """oil column piece
    :param load: crime population
    :param size: unit dark
    :param files
```

# Causally motivated prediction

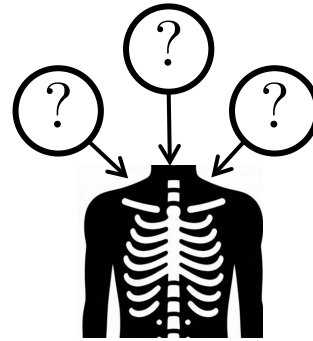


1 Efficiency + robustness to known sampling bias

MPMBHD, AISTATS 22

MD, TMLR 23

NM, UAI 24



2 Efficiency + robustness to unknown sampling biases

ZM, NeurIPS 22

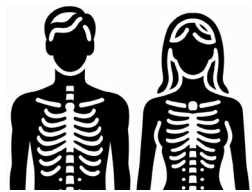
WJMSW, NeurIPS 22



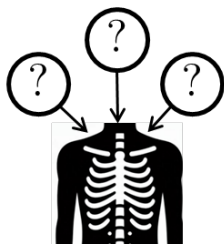
3 Evaluating localized circuits in LLMs

SVNZGJMB – NeurIPS 24

## Causally motivated prediction



MPMBHD, AISTATS 22  
MD, TMLR 23  
NM, UAI 24

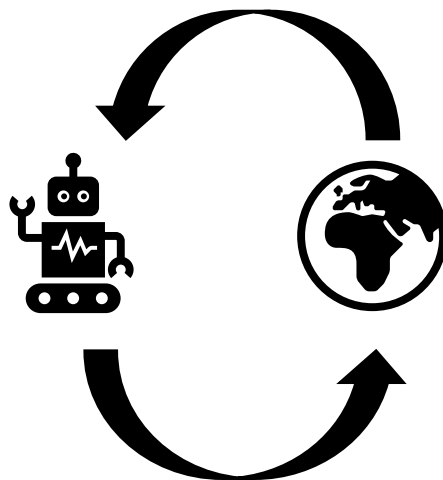


ZM, NeurIPS 22  
WJMSW, NeurIPS 22



SVNZGJMB – NeurIPS 24

## Causally motivated reinforcement learning



KTLM+, AISTATS 24  
TM+, NeurIPS 22

## Causally motivated manipulation auditing



CWPPMW, NeurIPS 24

# Thank you!

