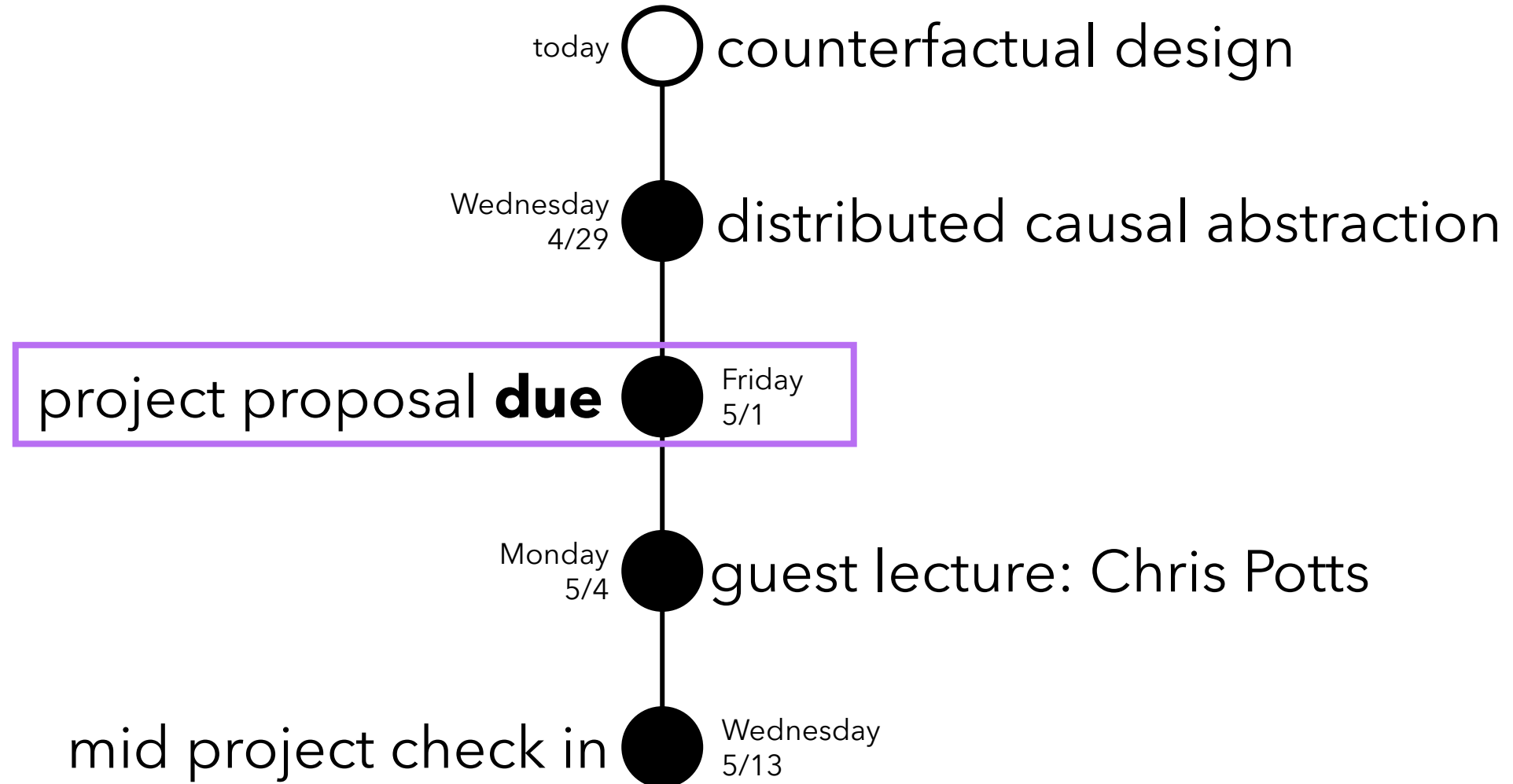


# designing counterfactuals for causal abstraction

CS 221M  
Week 5, Lecture 9



# Course timeline



# Attendance form

- Participation policy  
You get full credit for submitting the attendance form, even if it's incomplete
- Just make sure it has your name!



[tinyurl.com/cs221m-lecture9](https://tinyurl.com/cs221m-lecture9)

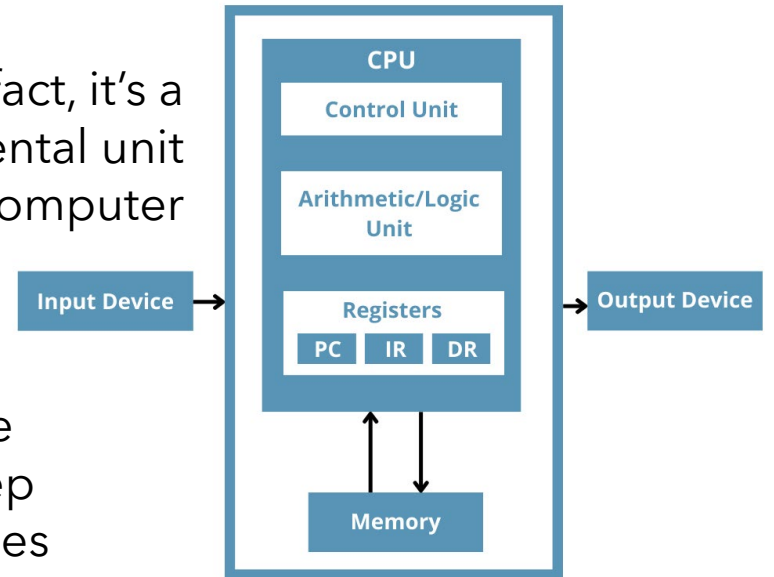
# Motivation: transformers use pointers

what does "memory" in transformers look like?

```
int a = 15; ← value  
↑  
Data type ← variable name
```

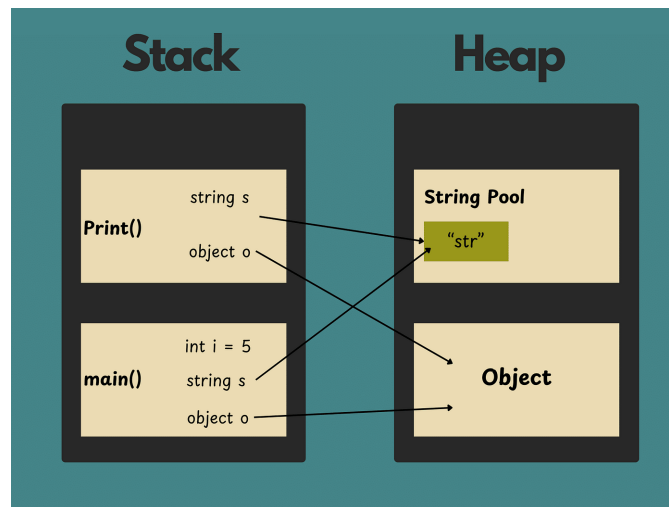
Assigning variables to values is a coding fundamental

In fact, it's a fundamental unit of a computer



Computers use *pointers* to keep track of variables

**We'll show that transformers do the same thing**



# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# Indirect object identification (IOI)

How do models represent  
anaphora?

# Tracking memory: indirect object

Alice and Bob went to the library. Bob gave a book to \_\_\_\_\_

```
x = 1
y = x
```

what does y store?

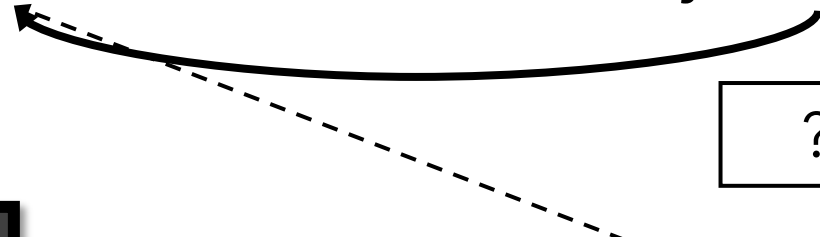
?

0x3

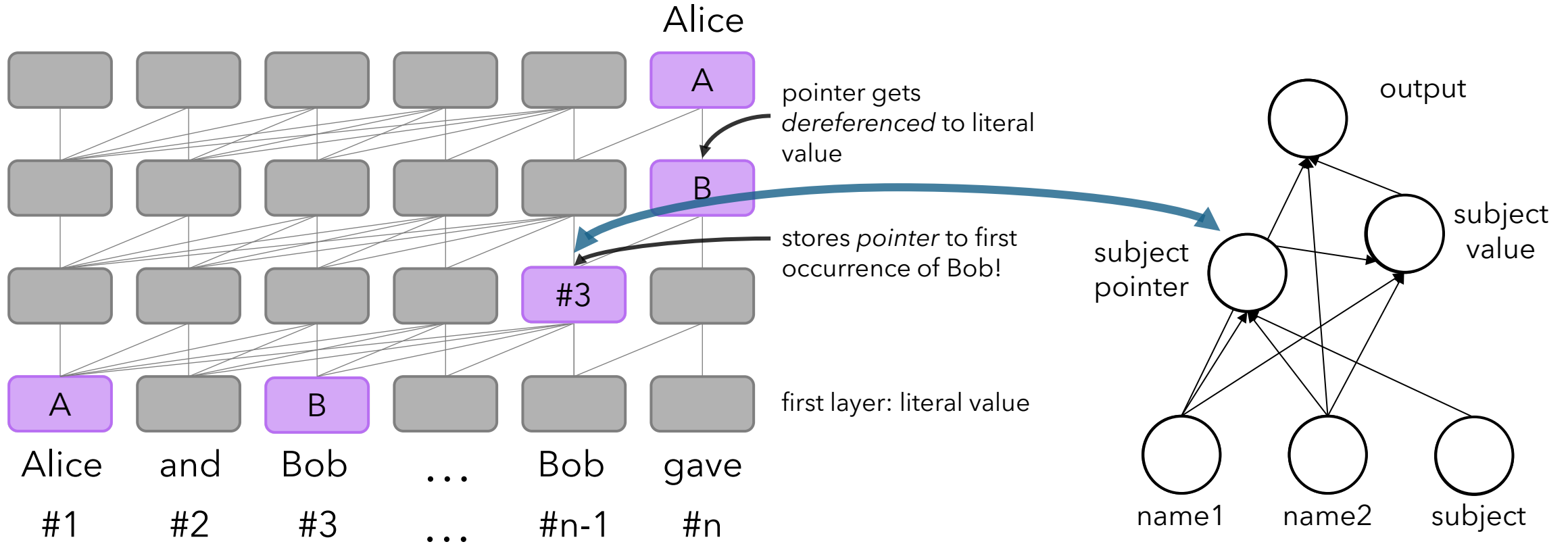
option 2  
store pointer  
to first mention  
of Bob

Bob

option 1  
copy over  
information  
about Bob



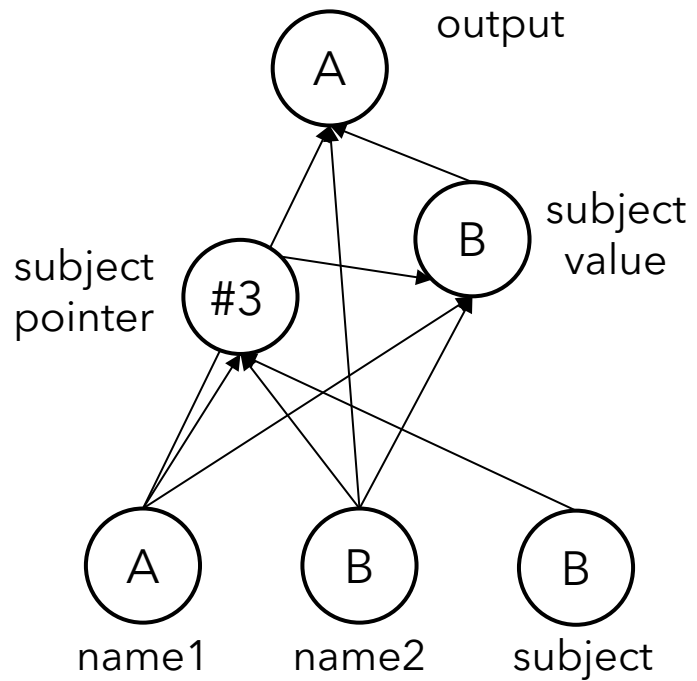
# Tracking memory: indirect object



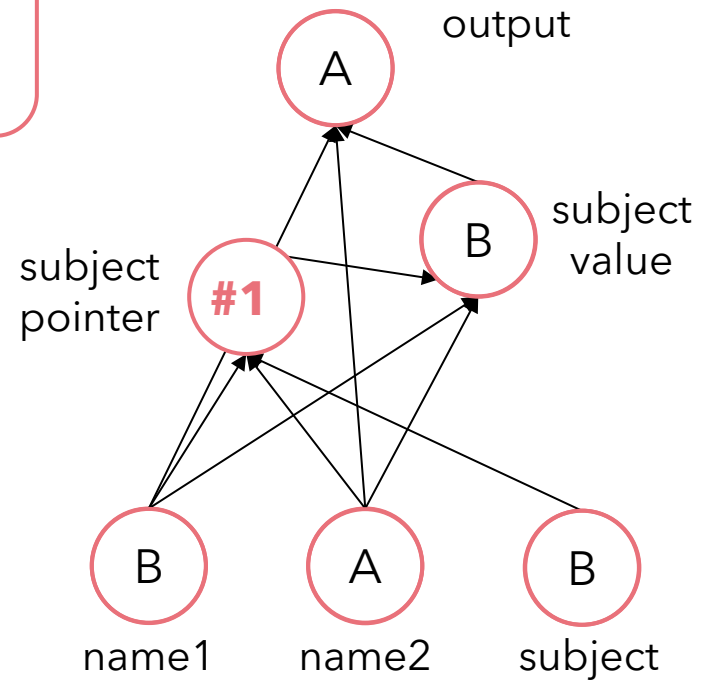
Alice and Bob went to the library. Bob gave a book to \_\_\_\_\_

# Tracking memory: indirect object

how do we change just the value of the subject pointer?

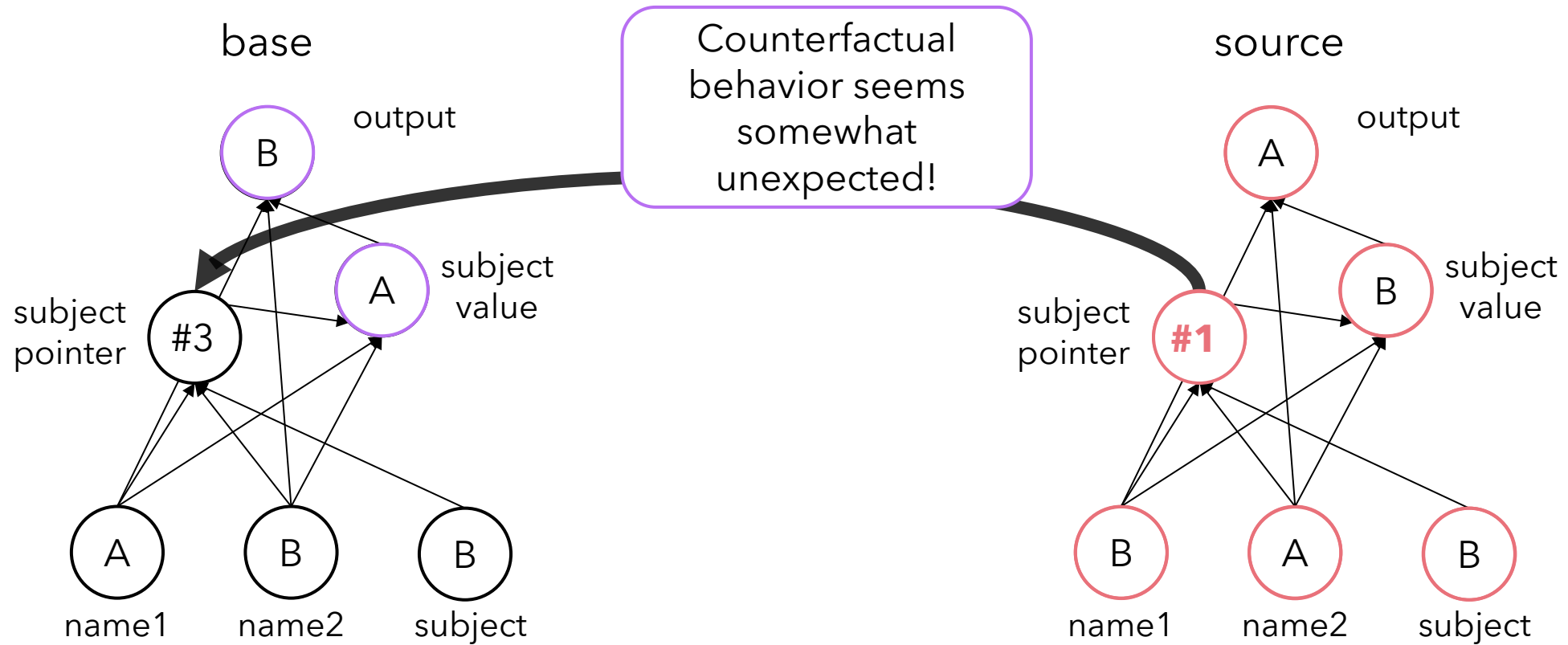


Alice and Bob went to the library.  
Bob gave a book to \_\_\_\_\_



**Bob** and **Alice** went to the library.  
Bob gave a book to \_\_\_\_\_

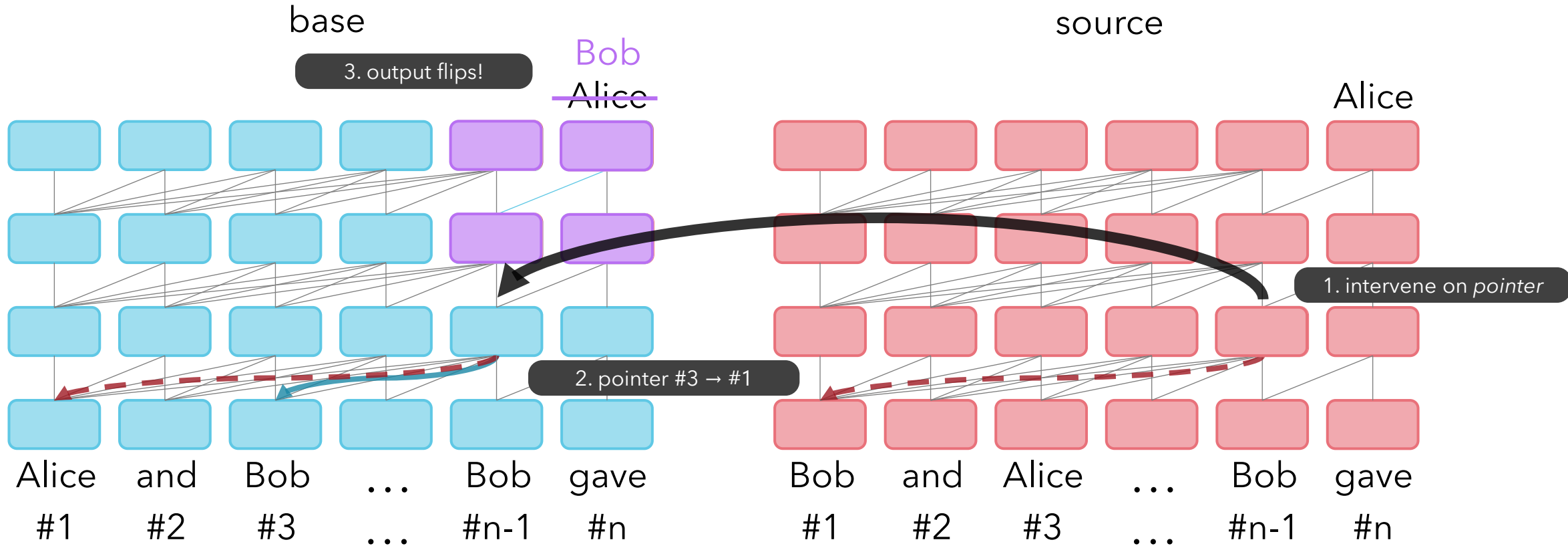
# Tracking memory: indirect object



Alice and Bob went to the library.  
Bob gave a book to \_\_\_\_\_

**Bob** and **Alice** went to the library.  
Bob gave a book to \_\_\_\_\_

# Tracking memory: indirect object

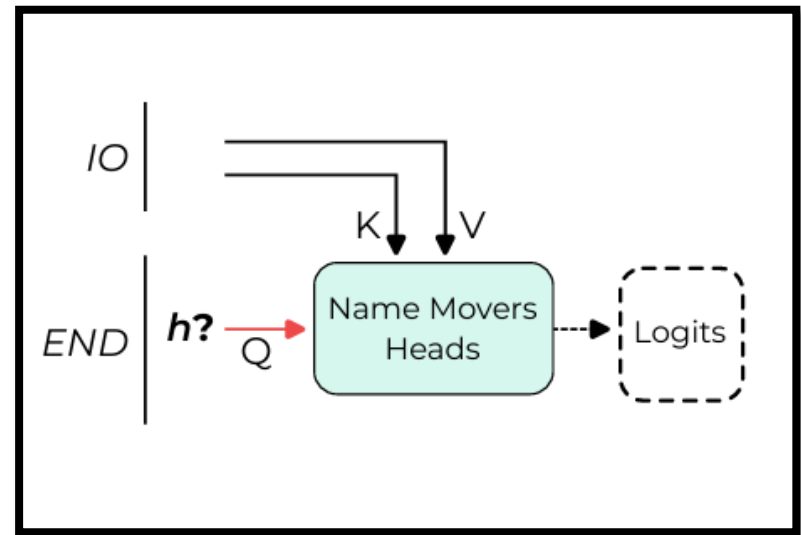


Alice and Bob went to the library.  
Bob gave a book to \_\_\_\_\_

Bob and Alice went to the library.  
Bob gave a book to \_\_\_\_\_

# Let's try it out!

exercise on  
workbench



# Attendance form

- Participation policy  
You get full credit for submitting the attendance form, even if it's incomplete
- Just make sure it has your name!



[tinyurl.com/cs221m-lecture9](https://tinyurl.com/cs221m-lecture9)

# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

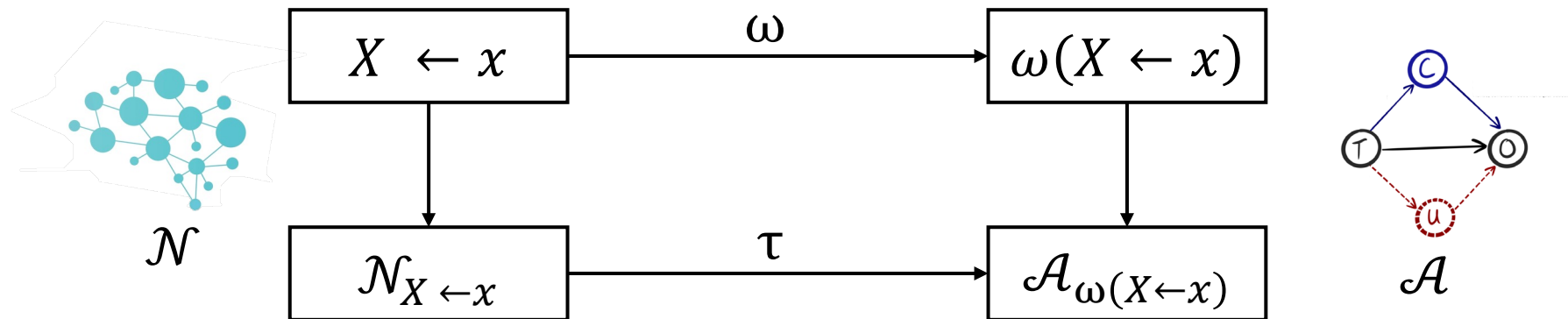
# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# Designing counterfactuals

Revisiting hierarchical equality

# Formalizing causal abstraction



$$\tau(\mathcal{N}_{X \leftarrow x}) = \mathcal{A}_{\omega(X \leftarrow x)}$$

with respect to a dataset of paired inputs,

$$\text{IIA}(\omega, \tau; D) = \sum_{b, s \in D} \mathbb{1} \left[ \tau \left( \mathcal{N}_{X \leftarrow x_s}(b) \right) = \mathcal{A}_{\omega(X \leftarrow x_s)}(b) \right]$$

how should we choose counterfactual pairs?

# Criteria for counterfactuals

how should we choose counterfactual pairs?

Change variable

b and s should have different values for x

$$\mathcal{A}(b)[X]$$

$\neq$

$$\mathcal{A}(s)[X]$$

Change value

Intervention should change output

$$\mathcal{A}_{X \rightarrow \mathcal{A}(s)}[X](b)$$

$\neq$

$$\mathcal{A}(b)$$

Rule out alternatives

Can't be explained by a different intervention

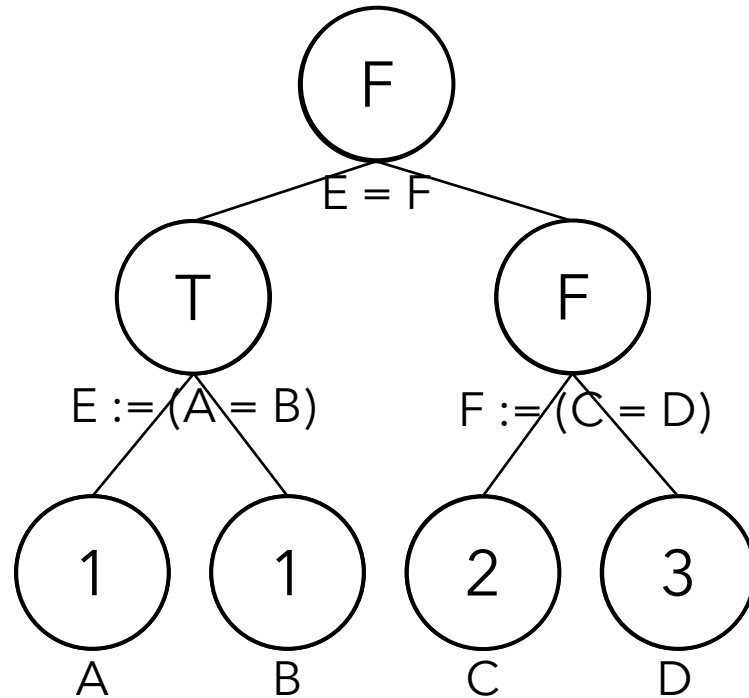
$$\mathcal{A}_{X \rightarrow \mathcal{A}(s)}[X](b)$$

$\neq$

$$\mathcal{A}_{W \rightarrow \mathcal{A}(s)}[W](b)$$

this is only with respect to  $\mathcal{A}$ !

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

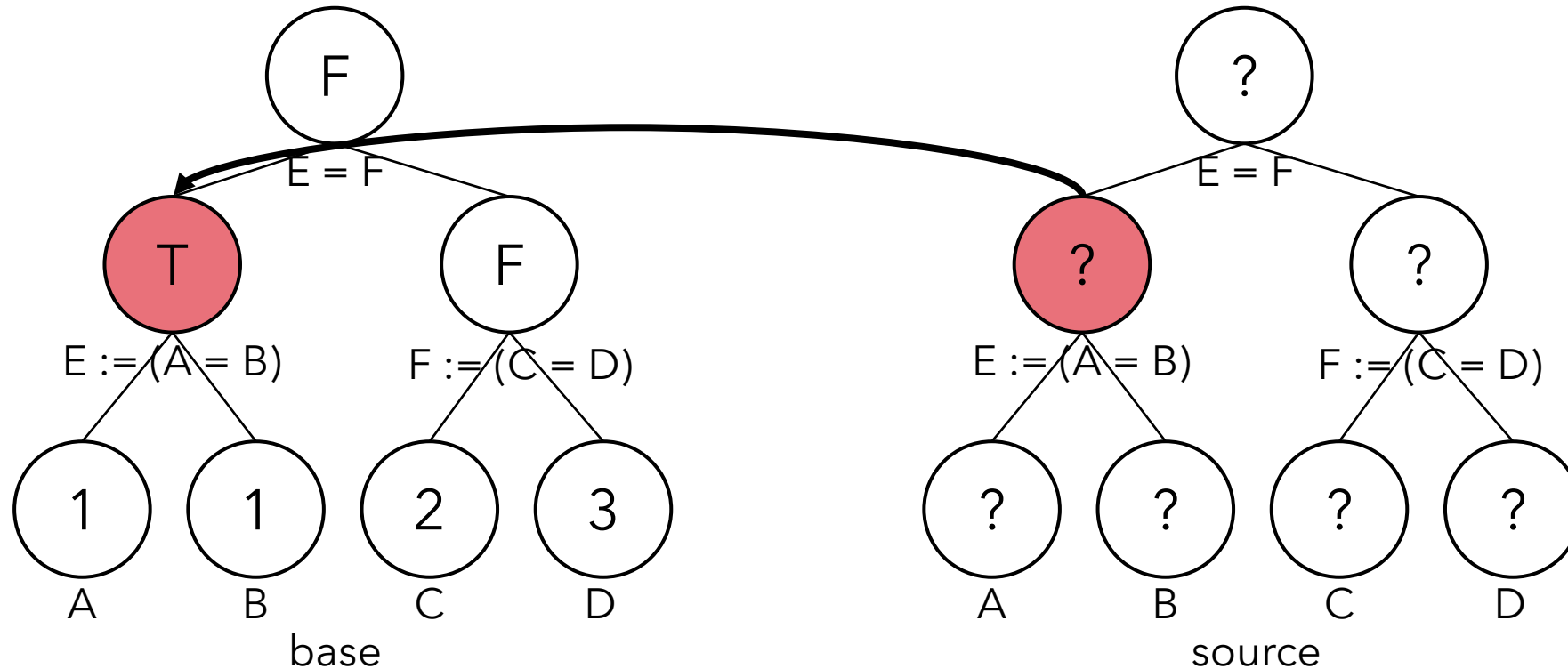
Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

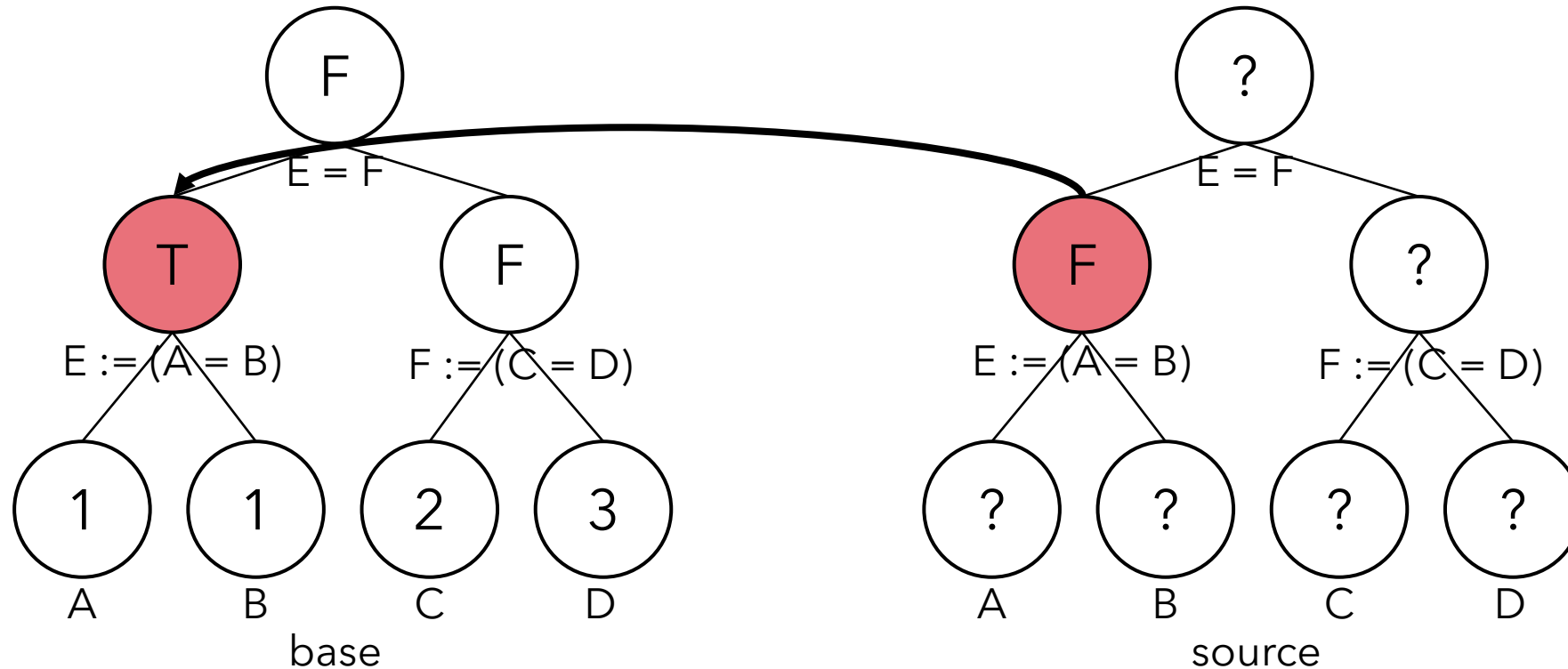
Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

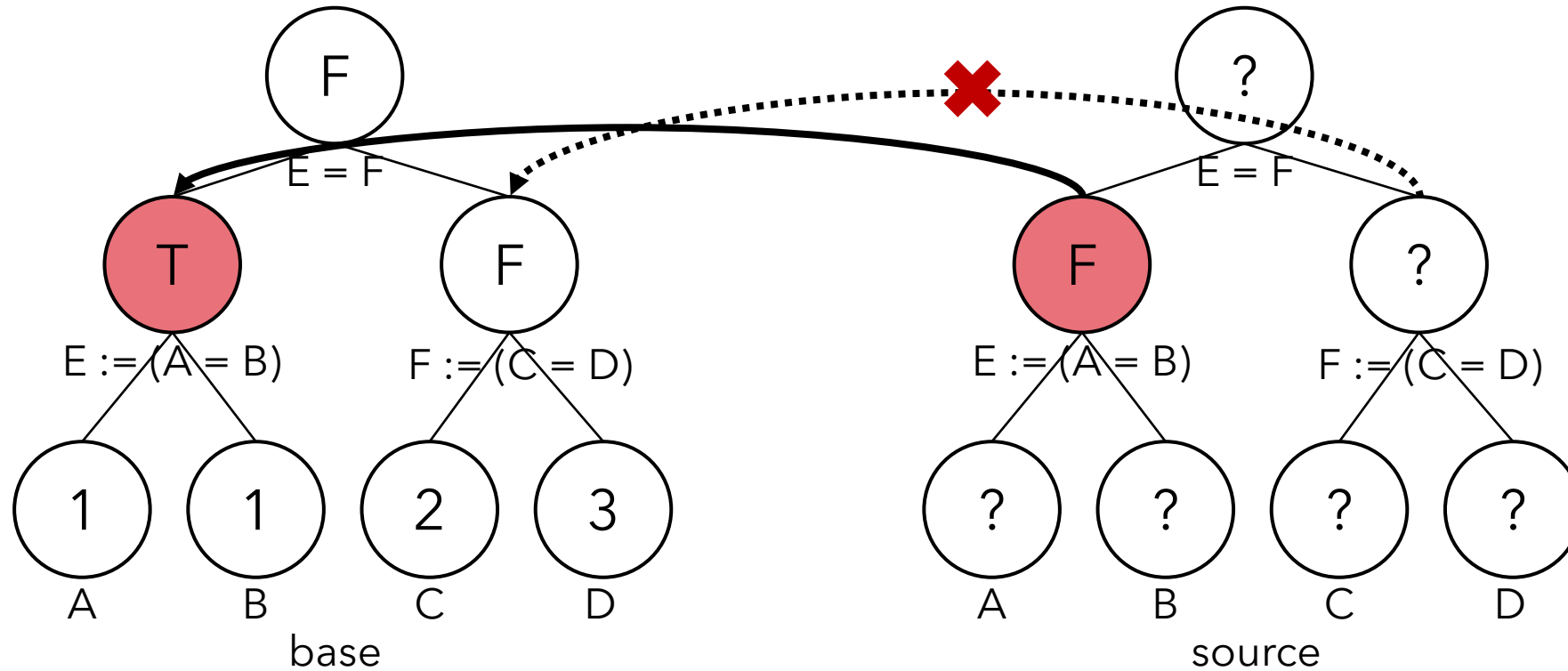
Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

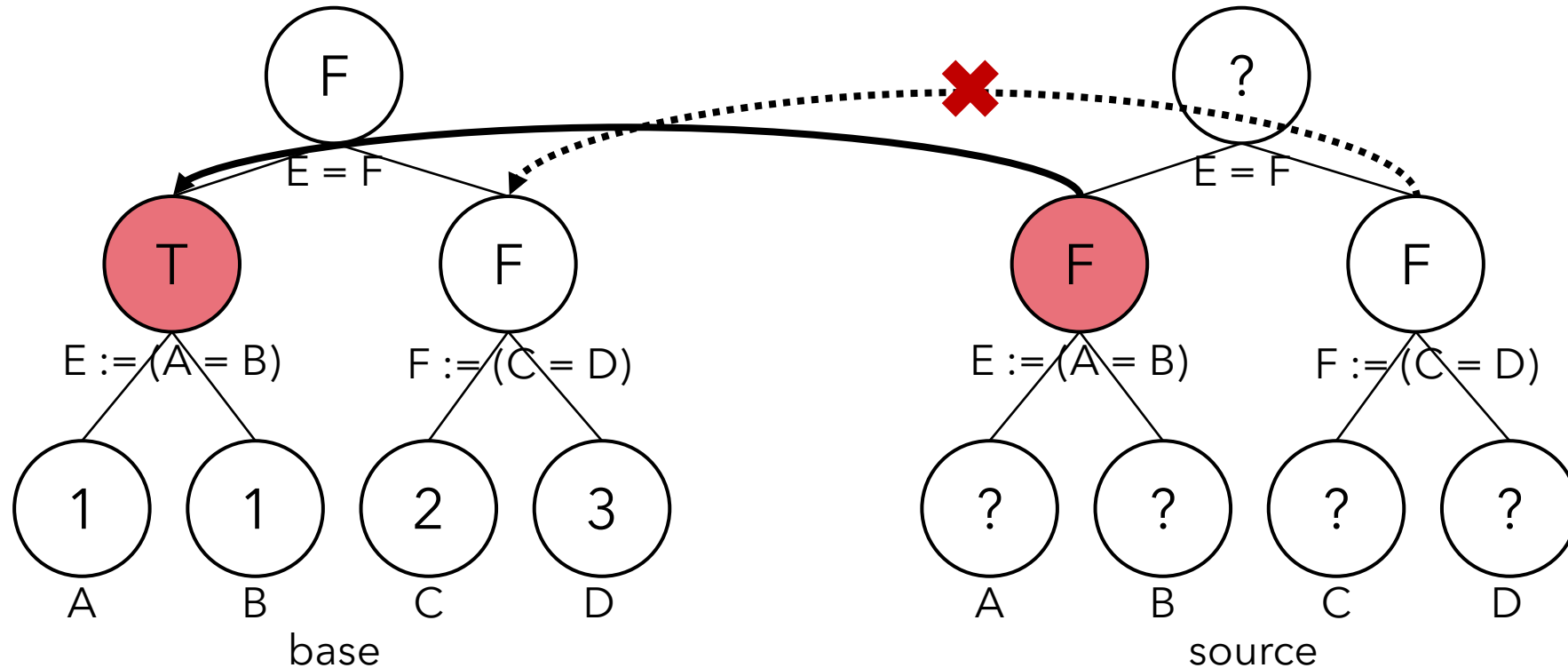
Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

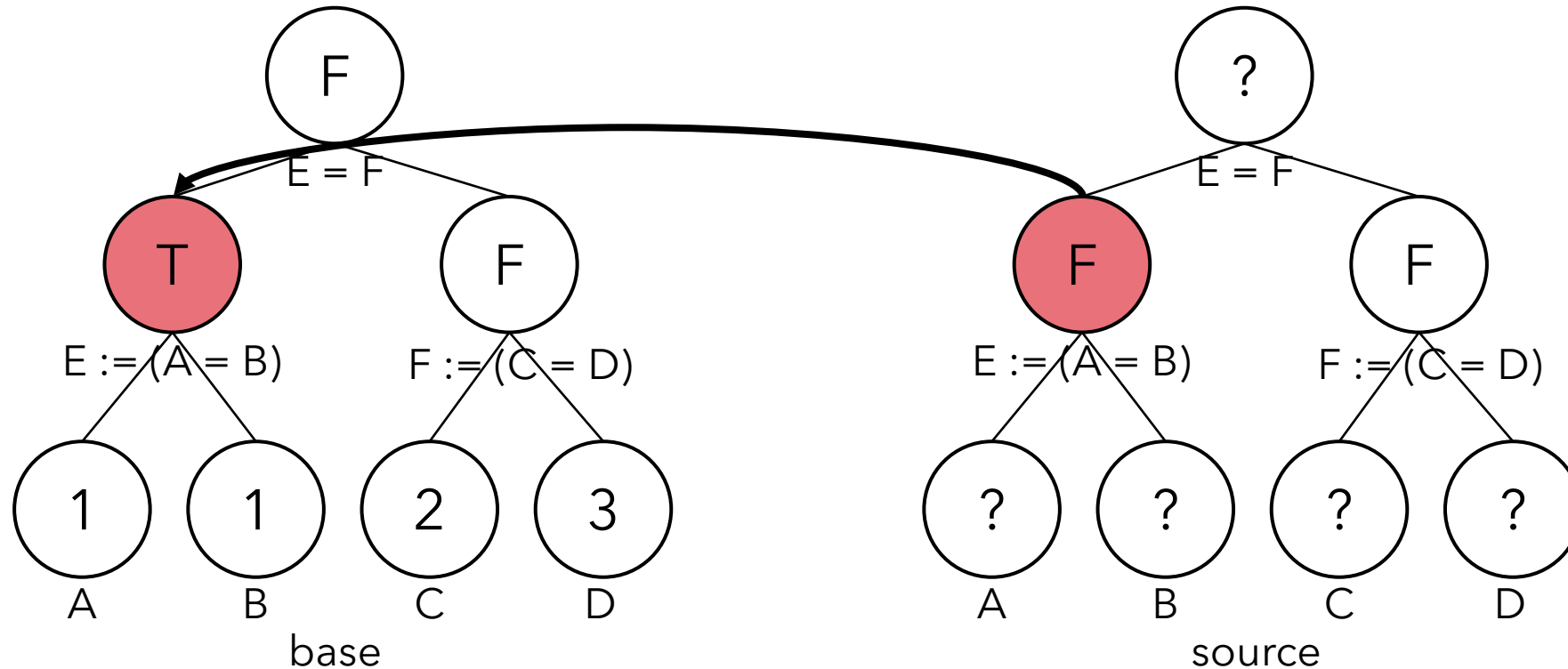
Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

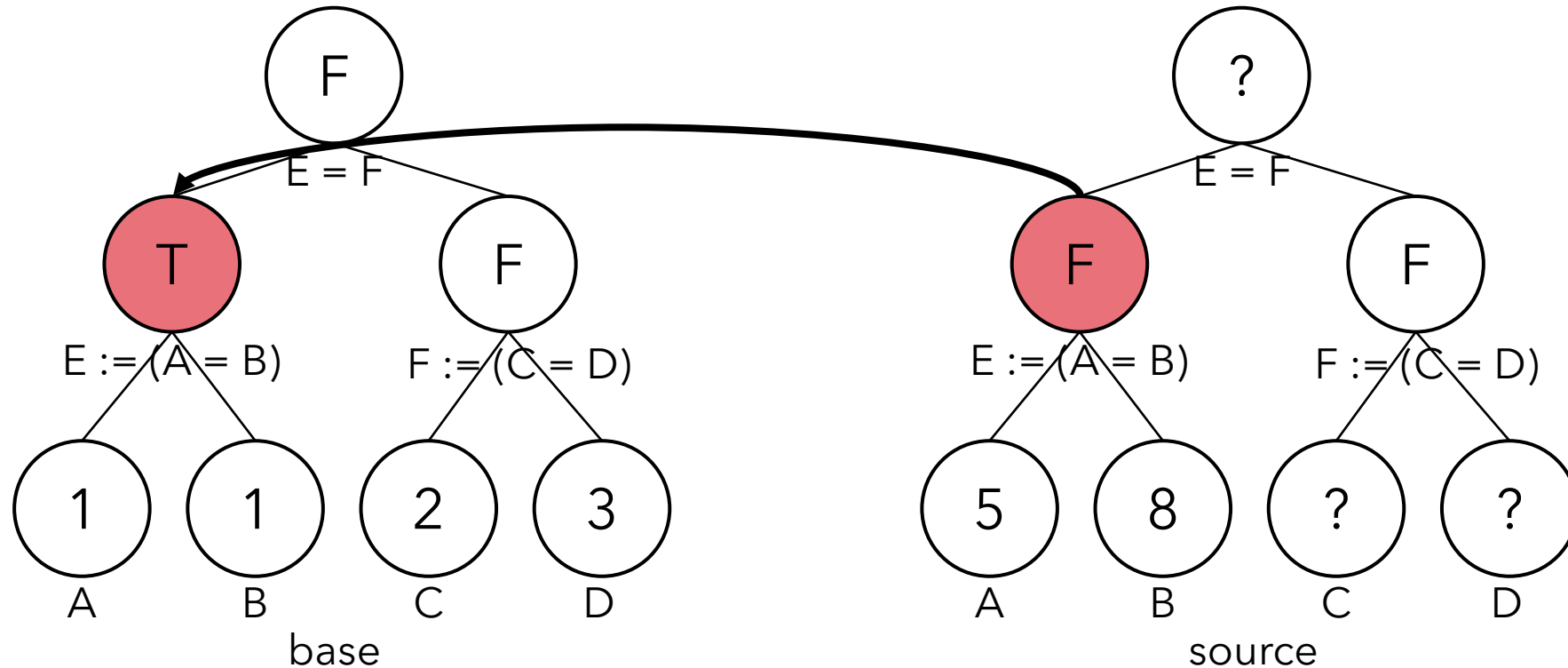
Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

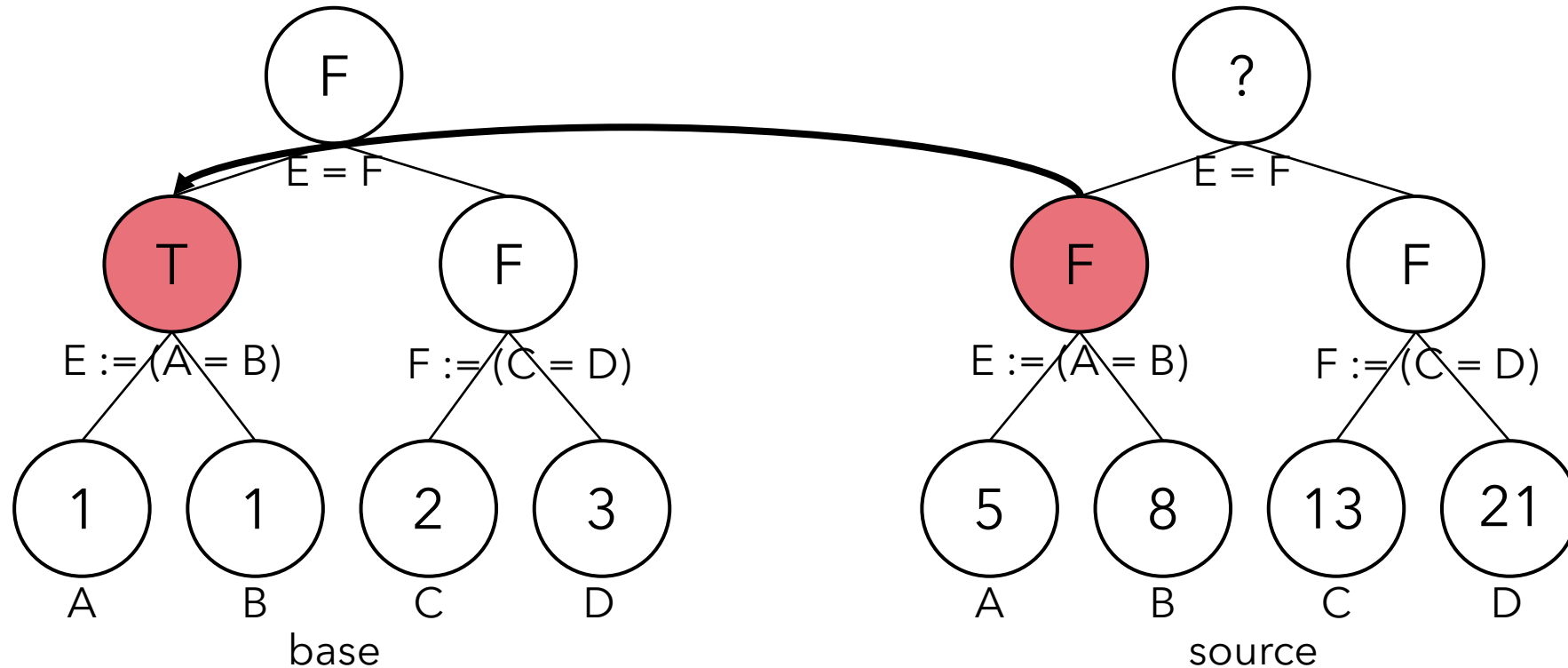
Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

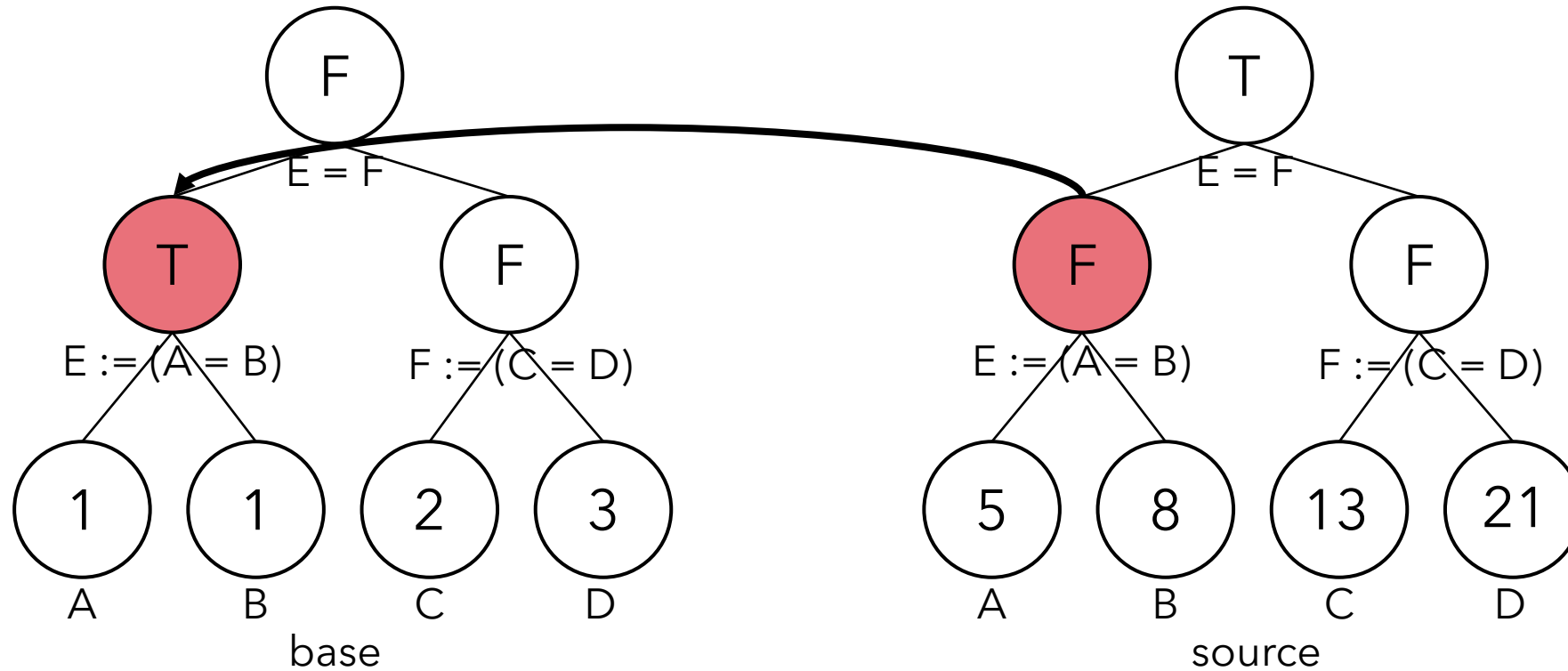
Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# Revisiting hierarchical equality



Change variable

b and s should have different values for x

Change value

Intervention should change output

Rule out alternatives

Can't be explained by a different intervention

# What happens if we try different pairs?

## random

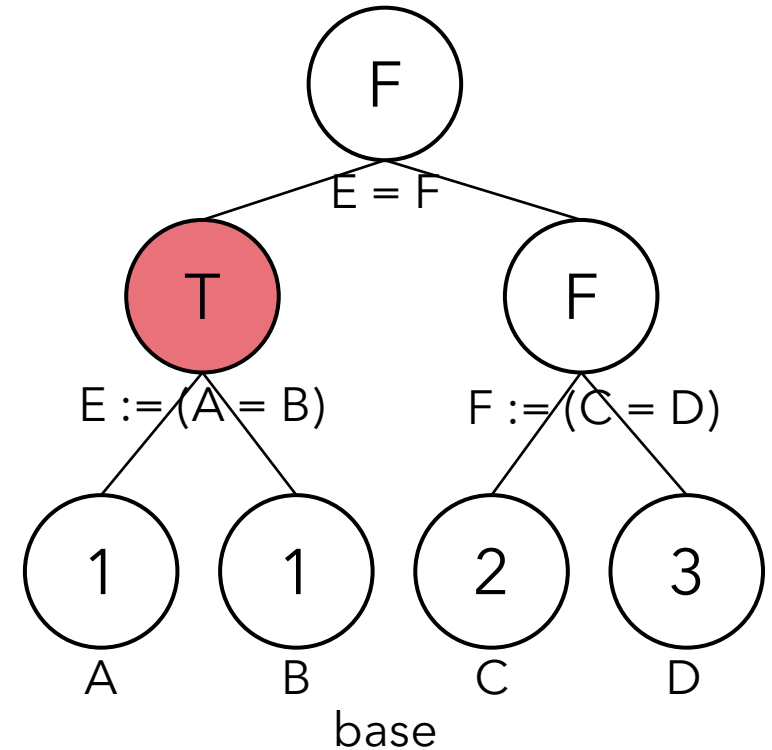
- Sample different pairs of inputs

## target E

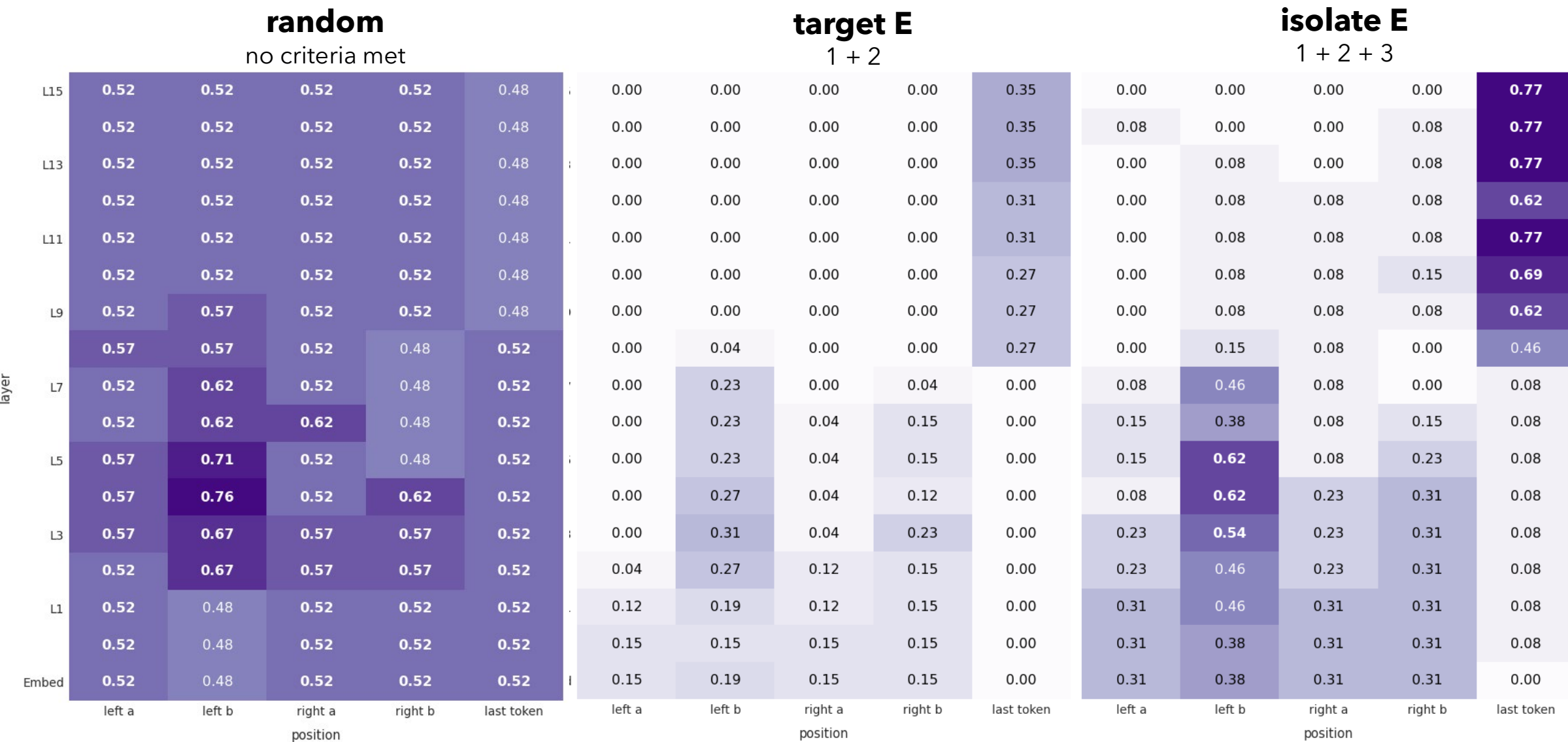
- Sample inputs with different values for E

## isolate E

- Sample inputs with different values for E and the same value for F



# What happens if we try different pairs?



# Lecture overview

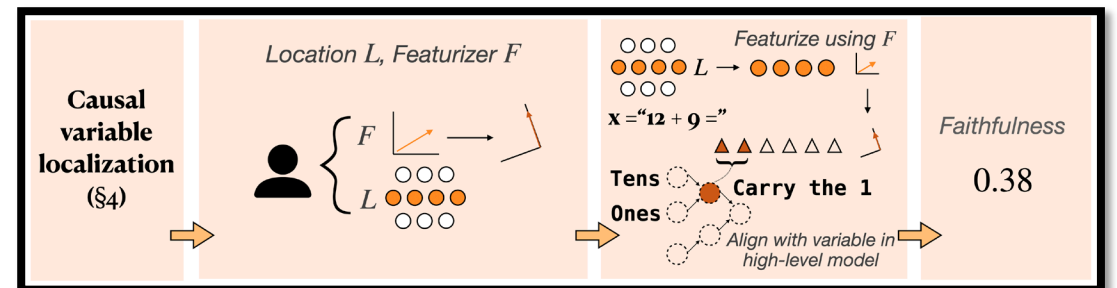
- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# MCQA example

exercise on  
workbench



# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# Entity binding

How do models track relations  
between variables?

# Tracking memory: indirect object

Alice and Bob went to the library. Bob gave a book to \_\_\_\_\_

```
x = 1
y = x
```

what does y store?

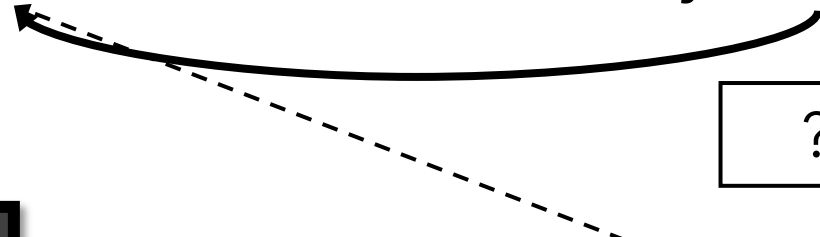
?

0x3

option 2  
store pointer  
to first mention  
of Bob

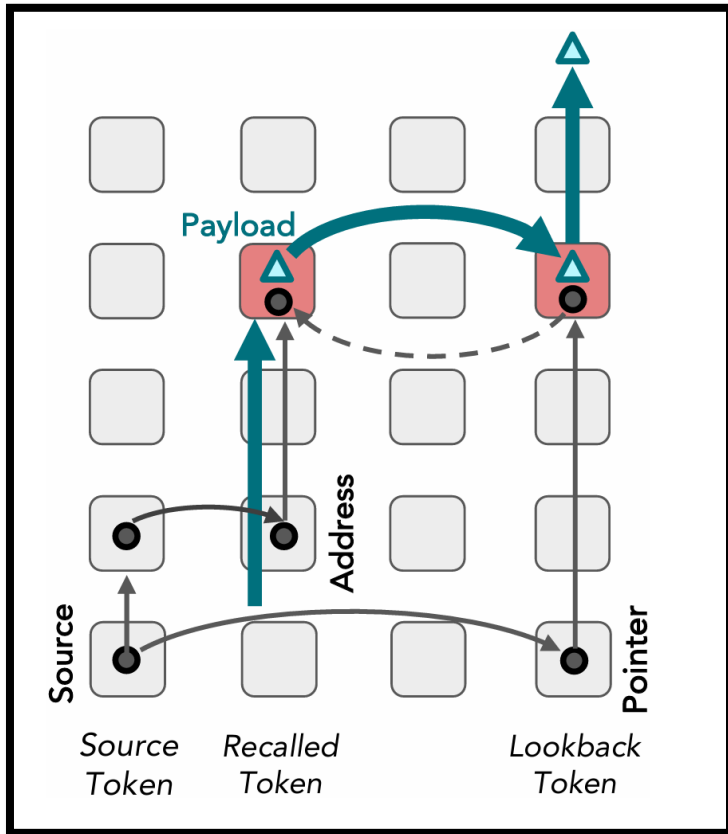
Bob

option 1  
copy over  
information  
about Bob



# Tracking memory: entity binding

Alices likes dogs. Bob likes cats. We know that Bob likes \_\_\_\_\_



0x3

store pointer  
to first mention  
of Bob

# Tracking memory: entity binding

source

We know that Bob likes cats

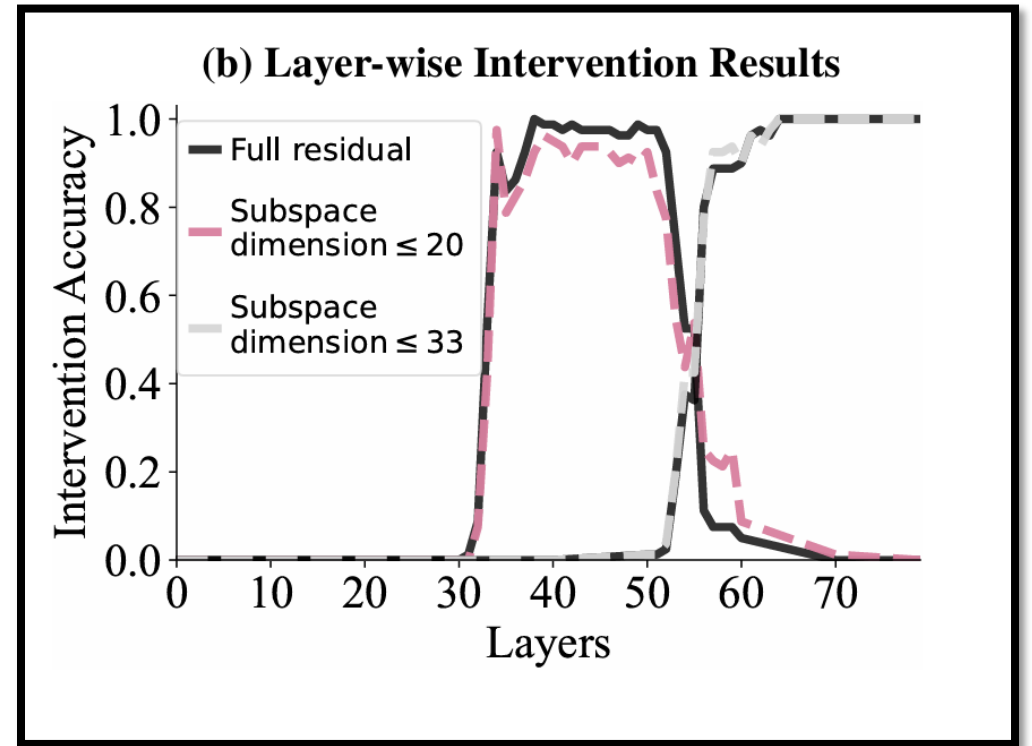


base

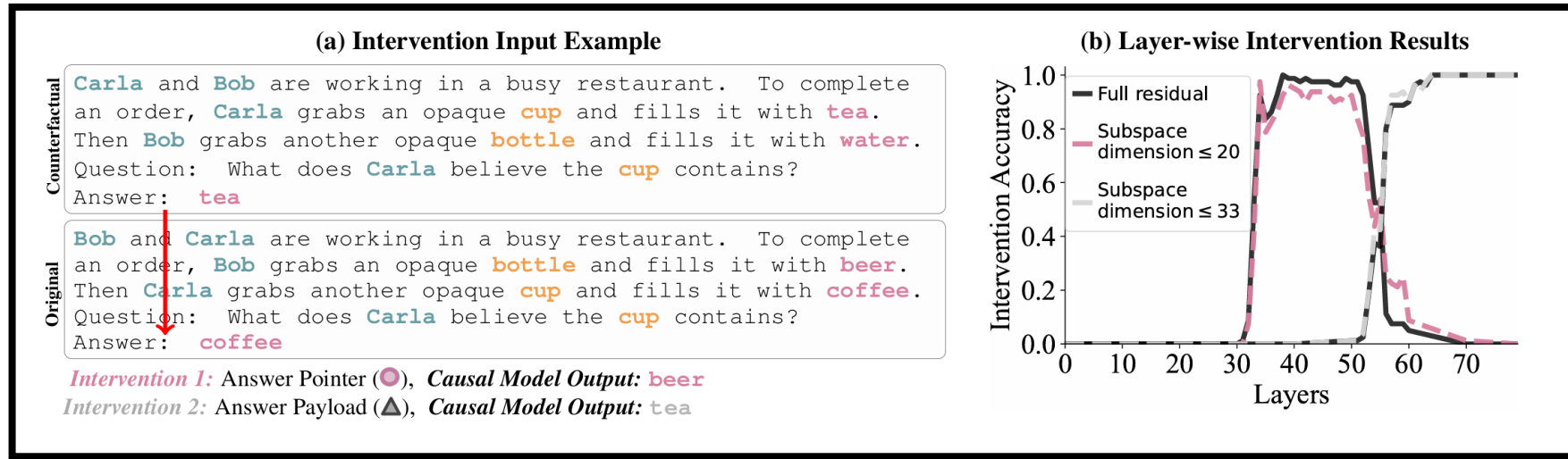
Alices likes dogs. Bob likes cats. We know that Bob likes ~~cats~~  
dogs

Let's try it  
out!

exercise on  
workbench



# Replicated lookback mechanism!



**Source Prompt**

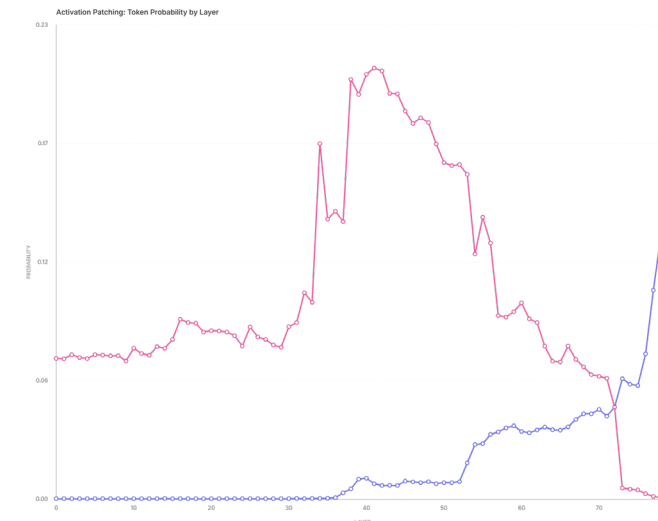
<|begin\_of\_text|> Bob likes pizza, Alice likes pancakes, Charlie likes burgers. We know that Bob likes pizza

Shift+click to select a range of tokens

**Target Prompt**

<|begin\_of\_text|> Alice likes dogs, Bob likes cats, Charlie likes mice. We know that Bob likes cats

⌘/Ctrl+click to freeze tokens



# What about multiple mechanisms?

- It's often the case that we want to test different hypotheses
- The *same* (base, source) pair could test different hypotheses
- Hypotheses predict *different* counterfactual outputs for the *same* intervention

# Mixture of mechanisms: entity binding

- Where is the pointer pointing?

Ann loves ale. Joe loves jam. Pete loves pie.  
We know that Ann loves \_\_\_\_\_

## Positional

- referent stored at *Ann's position*

## Lexical

- referent stored at object closest to Ann

## Reflexive

- referent stored at Ann's token

# Mixture of mechanisms: entity binding

source Joe loves ale. Ann loves pie. Pete loves jam.  
We know that Ann loves \_\_\_\_\_

base Ann loves <sup>L</sup>ale. Joe loves <sup>P</sup>jam. Pete loves <sup>R</sup>pie.  
We know that Ann loves \_\_\_\_\_

## Positional

- referent stored at *Ann's position*

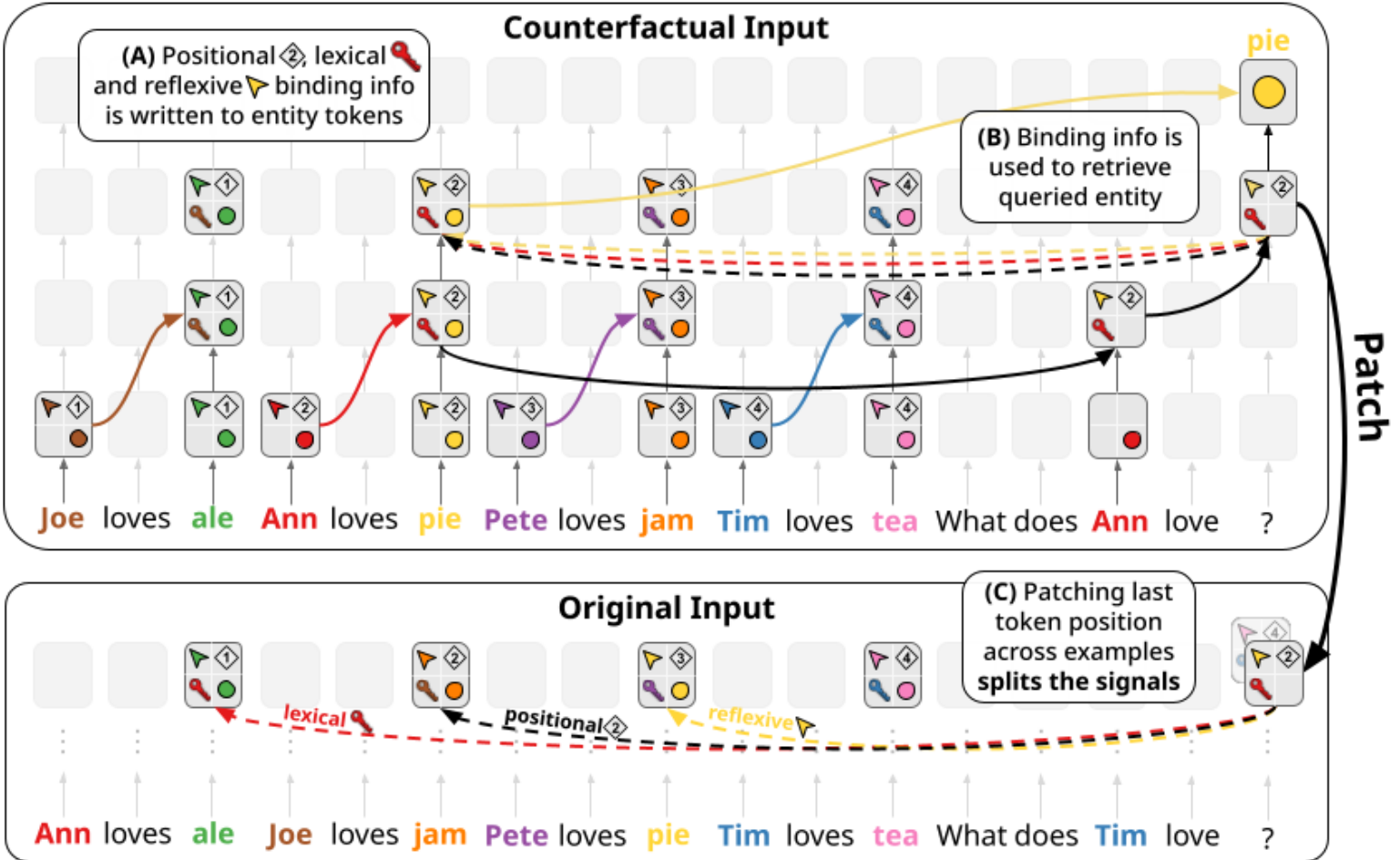
## Lexical

- referent stored at object closest to Ann

## Reflexive

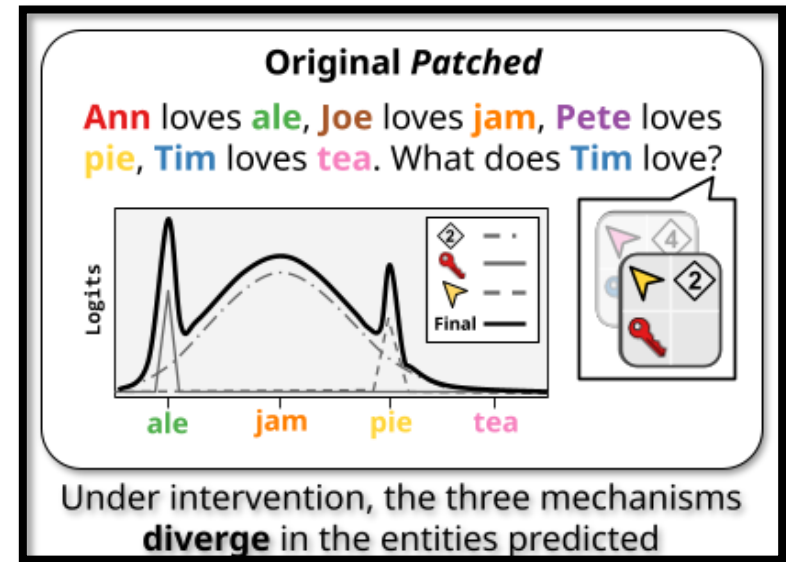
- referent stored at Ann's token

# Mixture of mechanisms: entity binding



# Let's try it out!

exercise on  
workbench



# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# Lecture overview

- example: indirect objects
- theory: designing counterfactuals
- example: MCQA
- example: entity binding

# Conclusion

## Testing hypothesis through counterfactual pairs

- Change variable value, change output, rule out alternatives

## The same counterfactual pair can distinguish different hypotheses

- Because they predict different outputs for the same intervention

## We apply causal abstraction through constructing counterfactuals

- This can help us discover algorithmic primitives, like pointer-reference mechanism in transformers!

# Course timeline

