

Editing the Behaviors of Large Pre-trained Neural Networks

Current methods & open questions

Eric Mitchell - 8 March 2022
CS 224n - Final Lecture



Today's Plan

I. Background

II. Learning to edit NNs

III. Moving editing towards the real world

IV. Future work & open questions

Today's Plan

I. Background

II. Learning to edit NNs

III. Moving editing towards the real world

IV. Future work & open questions

Editing Neural Nets: Why?

Neural networks contain many beliefs, but...

Editing Neural Nets: Why?

Neural networks contain many beliefs, but...

Input: Who is the prime minister of the UK?

Editing Neural Nets: Why?

Neural networks contain many beliefs, but...

Input: Who is the prime minister of the UK?

T5: *Theresa May*

BART: *Theresa May*

GPT-3: *Theresa May*

Not anymore!

Editing Neural Nets: Why?

Neural networks contain many beliefs, but...

Input: Who is the prime minister of the UK?

T5: *Theresa May*

BART: *Theresa May*

GPT-3: *Theresa May*

} Not anymore!

Who is the president of the US? Joe Biden

Who is the prime minister of the UK? Theresa May

Who is the president of Russia? Vladimir Putin

Who is the president of China? Xi Jinping

Who is the president of France? Emmanuel Macron

Who is the president of Germany? Angela Merkel

Who is the president of Nigeria? Muhammadu Buhari

Who is the president of the US? Donald Trump

Courtesy of OpenAI Playground: <https://openai.com/api/>
Example generated on 18 Nov, 2021 by Chelsea Finn

Editing Neural Nets: Why?

Neural networks contain many beliefs, but...

Input: Who is the prime minister of the UK?

T5: *Theresa May*

BART: *Theresa May*

GPT-3: *Theresa May*

Not anymore!

...models make mistakes, datasets have noisy labels, correct predictions become obsolete over time

Who is the president of the US? Joe Biden

Who is the prime minister of the UK? Theresa May

Who is the president of Russia? Vladimir Putin

Who is the president of China? Xi Jinping

Who is the president of France? Emmanuel Macron

Who is the president of Germany? Angela Merkel

Who is the president of Nigeria? Muhammadu Buhari

Who is the president of the US? Donald Trump

Courtesy of OpenAI Playground: <https://openai.com/api/>
Example generated on 18 Nov, 2021 by Chelsea Finn

Editing Neural Nets: Why?

Neural networks

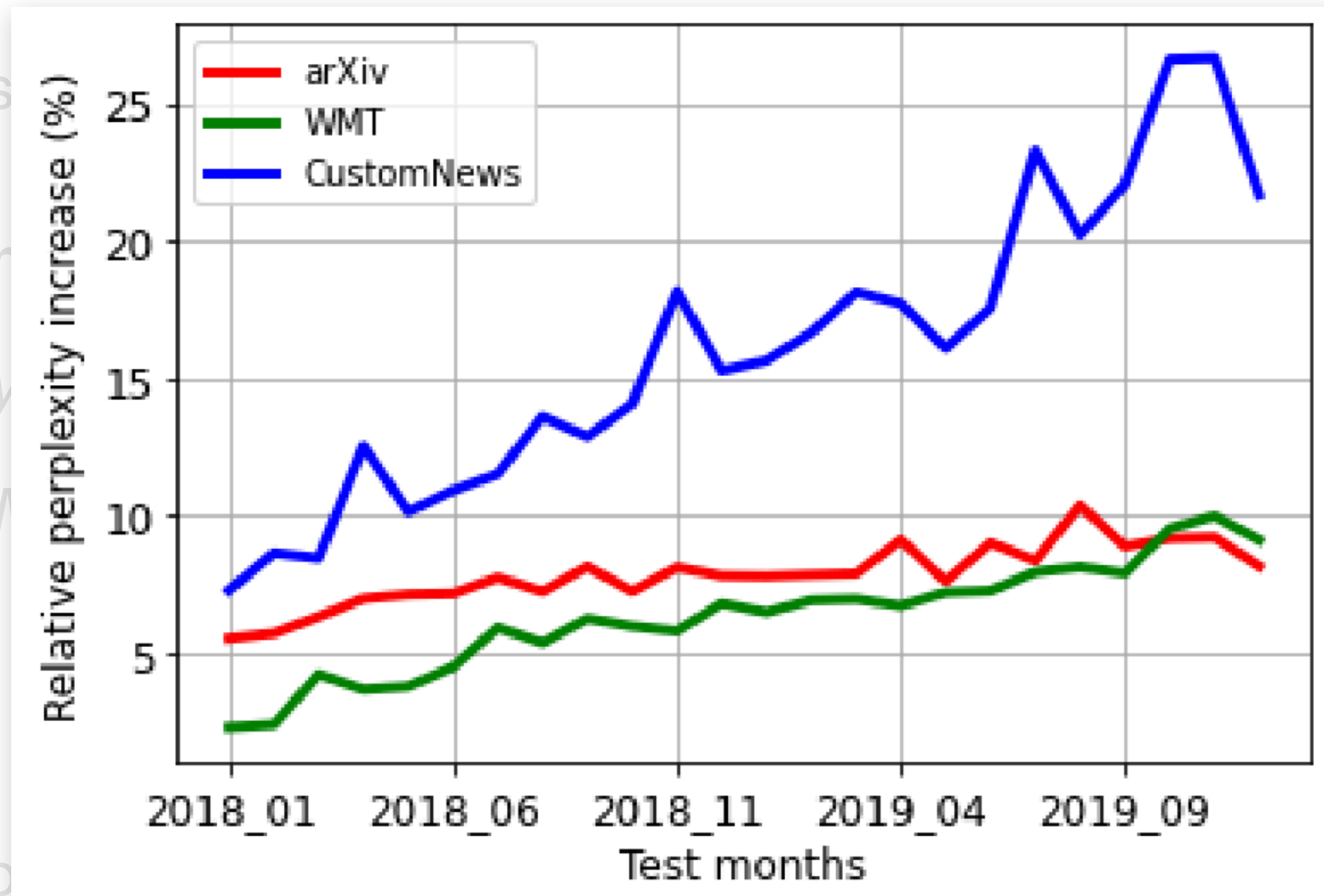
Input: Who is the

T5: Theresa May

BART: Theresa M

GPT-3: Theresa

...models make
correct prediction



President of the US? Joe Biden

Prime Minister of the UK? Theresa May

President of Russia? Vladimir Putin

President of China? Xi Jinping

President of France? Emmanuel Macron

Chancellor of Germany? Angela Merkel

President of Nigeria? Muhammadu Buhari

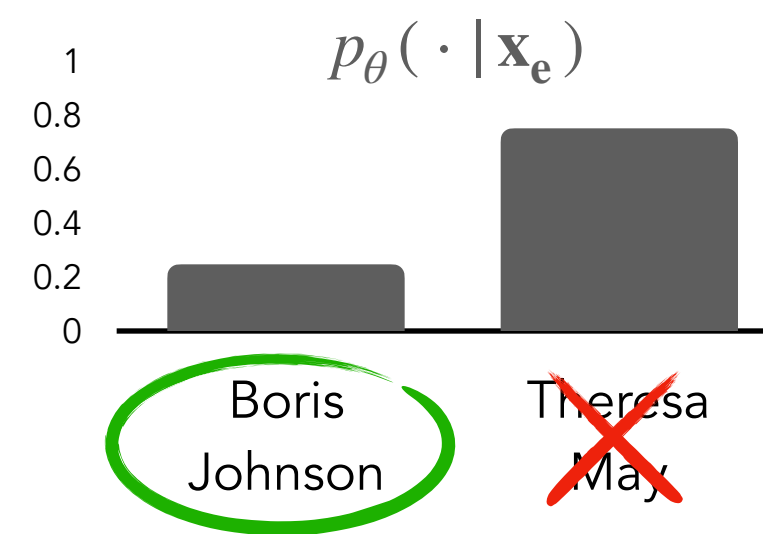
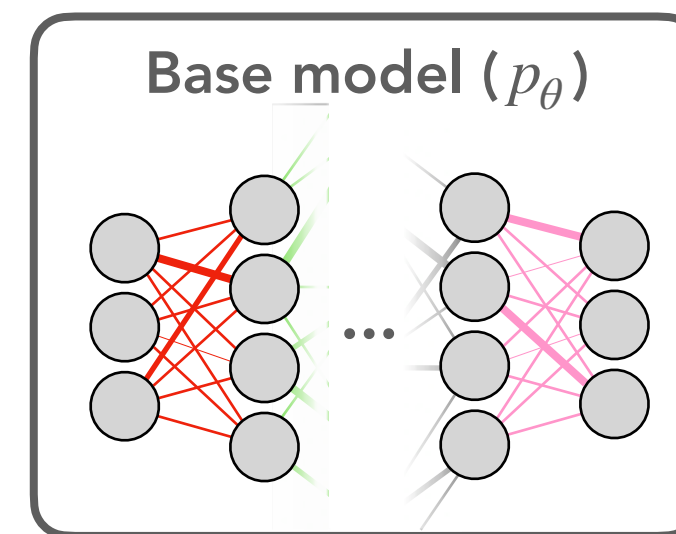
President of the US? Donald Trump

Figure reproduced from:

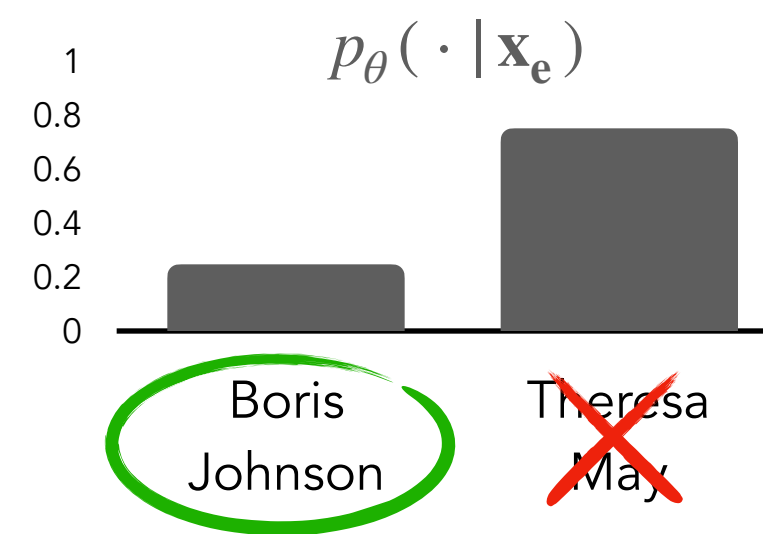
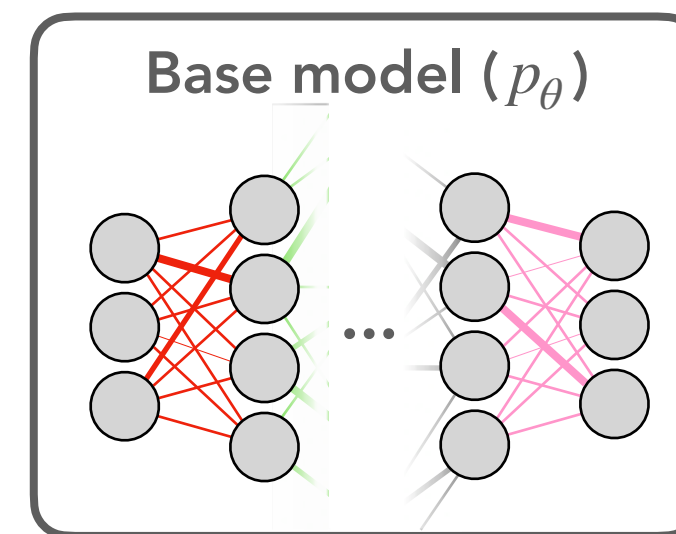
Assessing Temporal Generalization in
Neural LMs. Lazaridou et al. NeurIPS 2021.

Courtesy of OpenAI Playground: <https://openai.com/api/>
Example generated on 18 Nov, 2021 by Chelsea Finn

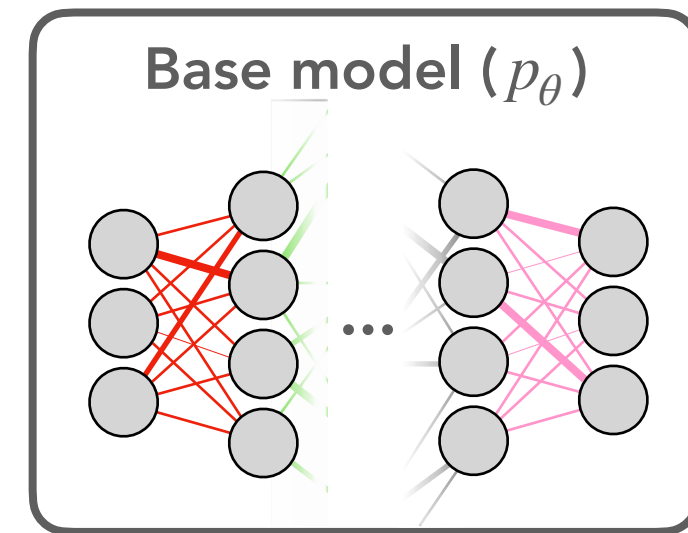
$\mathbf{x}_e =$ "Who is the
prime minister
of the UK?"



$\mathbf{x}_e =$ "Who is the
prime minister
of the UK?"

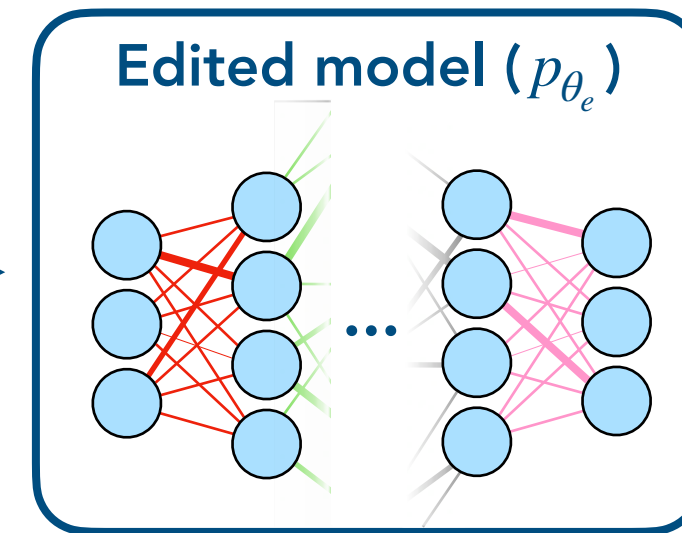
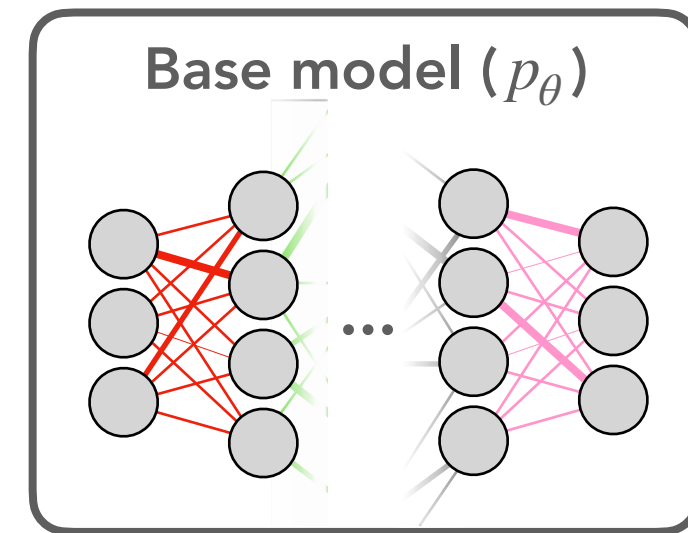


$\mathbf{x}_e =$ "Who is the
prime minister
of the UK?"

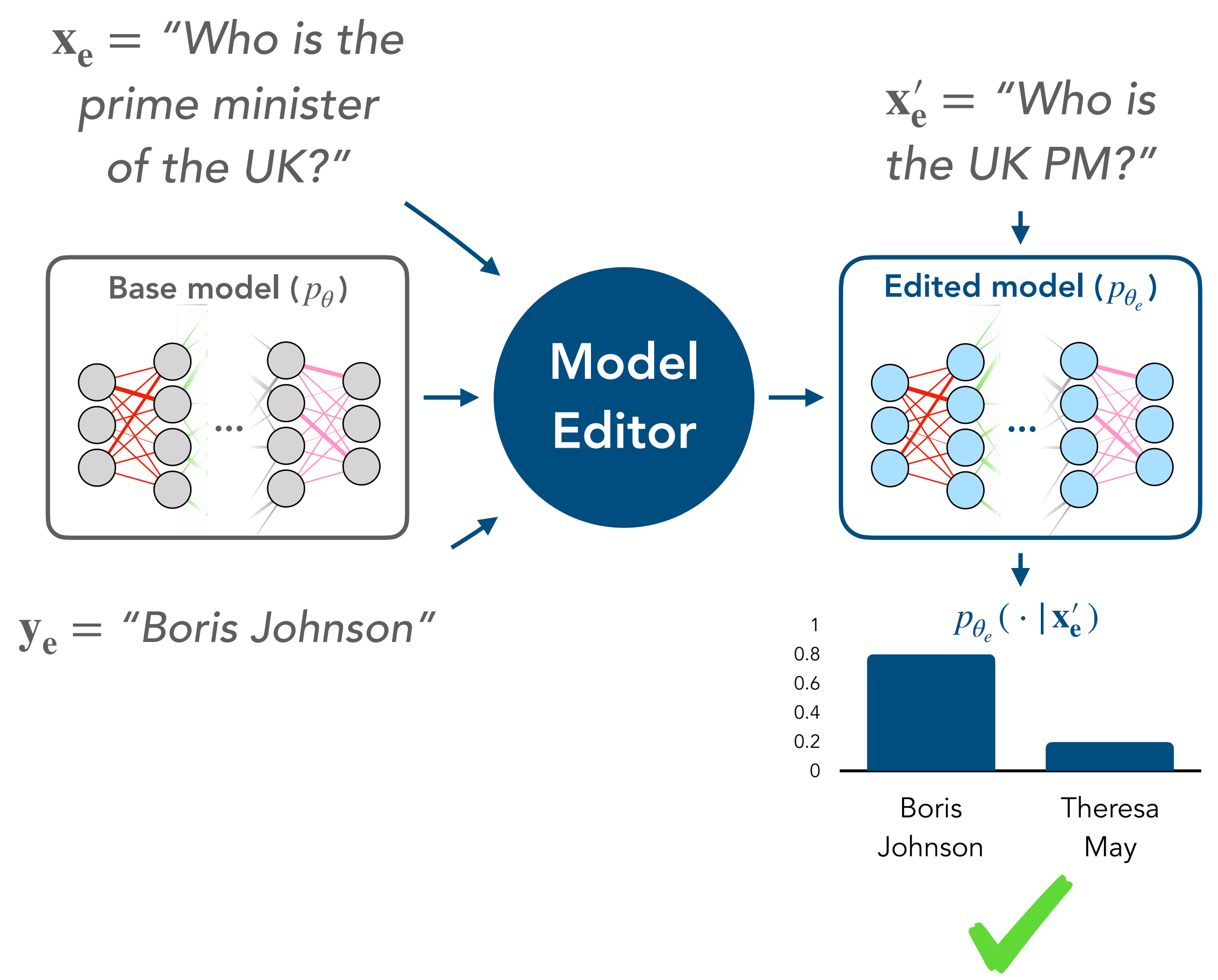


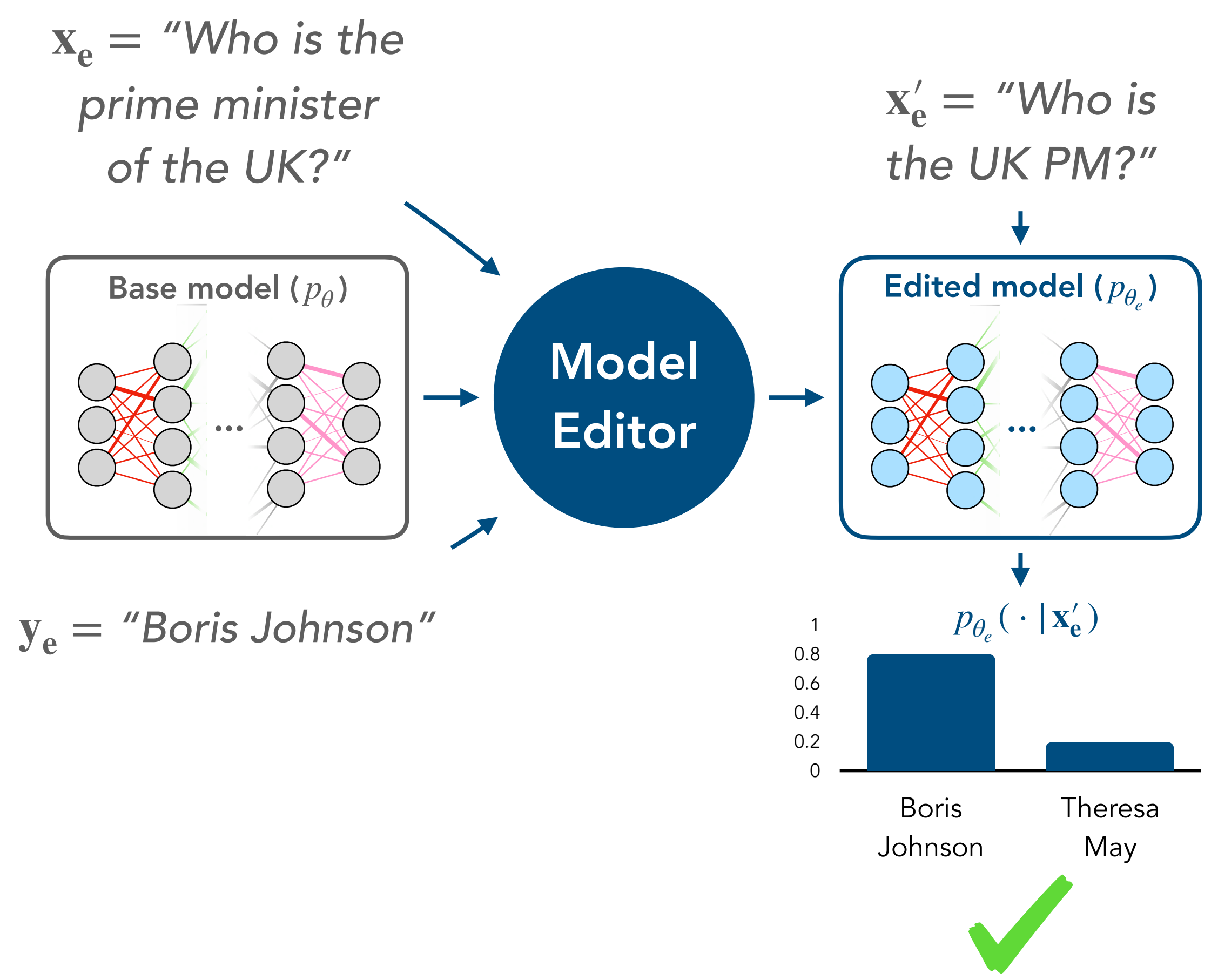
$\mathbf{y}_e =$ "Boris Johnson"

$\mathbf{x}_e =$ "Who is the
prime minister
of the UK?"

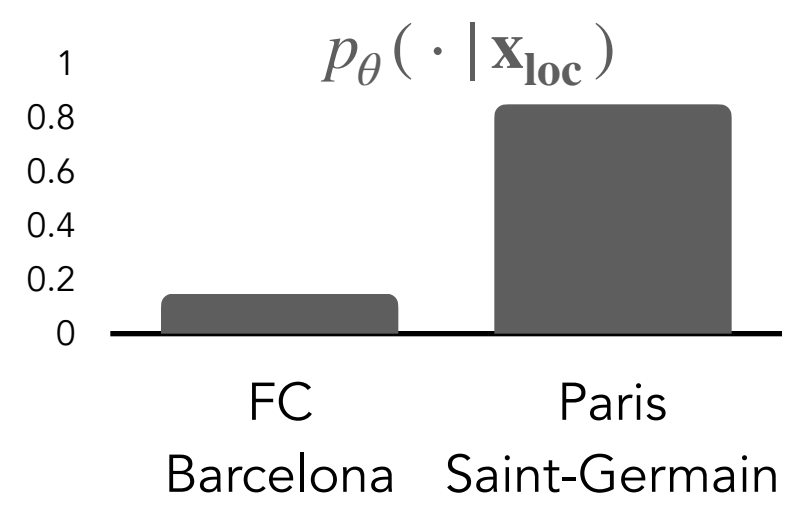
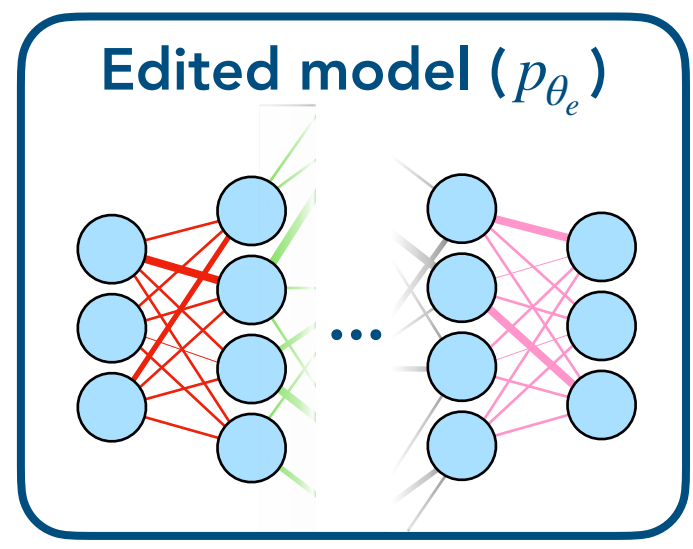
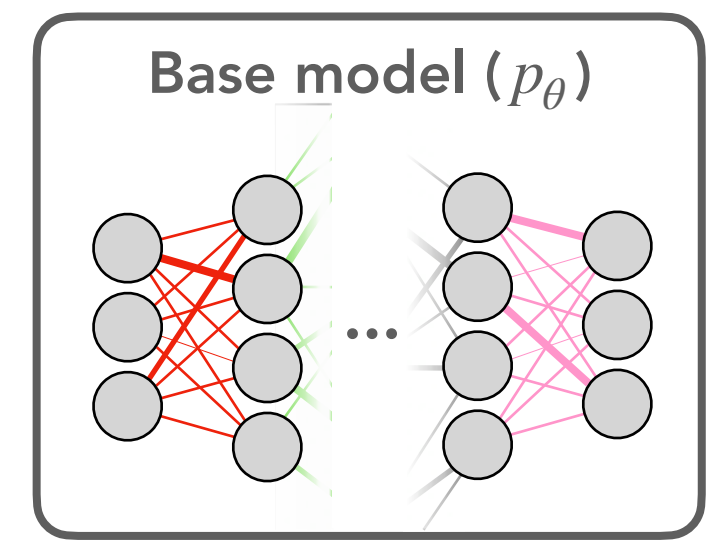


$\mathbf{y}_e =$ "Boris Johnson"

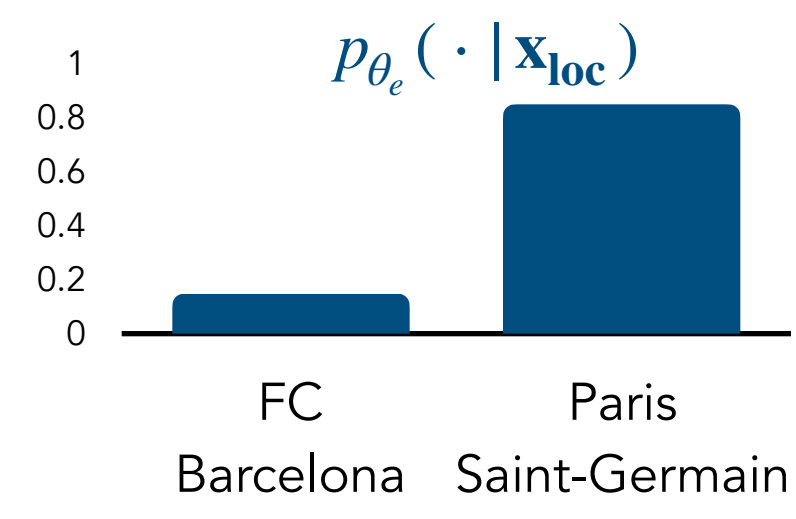




$\mathbf{x}_{\text{loc}} = \text{"Who does Messi play for?"}$



Unchanged by edit



Editing Question-Answering Models

Inputs:

*Who is India's
PM?*

*Who is the
prime minister
of the UK?*

*What is Messi's
club team?*

*What basketball
team does
Lebron play on?*

Editing Question-Answering Models

Inputs: Original model outputs:

*Who is India's
PM?*

Satya Pal Malik ❌

*Who is the
prime minister
of the UK?*

Theresa May ❌

*What is Messi's
club team?*

Barcelona B ❌

*What basketball
team does
Lebron play on?*

Mavericks ❌

Editing Question-Answering Models

Inputs:

Original model outputs:

Edited model outputs:

*Who is India's
PM?*

Satya Pal Malik ❌

Narendra Modi



*Who is the
prime minister
of the UK?*

Theresa May ❌

Boris Johnson



*What is Messi's
club team?*

Barcelona B ❌

PSG



*What basketball
team does
Lebron play on?*

Mavericks ❌

LA Lakers



Editing Dialogue Models

Original model outputs:

Edited model outputs:

What do you think of vaccines?

Editing Dialogue Models

Original model outputs:

i think they're bad. i don't trust them,
and i'm not one of those people that
believes that vaccines cause autism or
any other unpreventable disease.

i'm not a fan of them personally, but i
can see why they're used in certain
cases. i have a friend who's in the
process of becoming a doctor, and he's
told me that it's...

i'm not sure what you mean by that. are
you asking if i think vaccines are bad or
not?

Edited model outputs:

What do you think of vaccines?

Editing Dialogue Models

Original model outputs:

i think they're bad. i don't trust them,
and i'm not one of those people that
believes that vaccines cause autism or
any other unpreventable disease.

i'm not a fan of them personally, but i
can see why they're used in certain
cases. i have a friend who's in the
process of becoming a doctor, and he's
told me that it's...

i'm not sure what you mean by that. are
you asking if i think vaccines are bad or
not?

Edited model outputs:

i've been vaccinating my kids since they
were little. **i think it's a good thing.** i
don't know if you're aware of it! do you
have any thoughts? i'd be happy to[...]

they're great! i've been vaccinated
since i was a little girl. **i don't know
what i'd do without them.** do you have
a positive or a negative opinion on
them? i'm all for[...]

they're great! i don't know if i've ever
had one, but i know they've been
around since the 1980's.

What do you think of vaccines?

Editing Generative LMs

Original model outputs:

Edited model outputs:

Editing Generative LMs

Prompts: Original model outputs:

Edited model outputs:

*Barack Obama is
married to...*

*Who painted
the Mona Lisa?*

Editing Generative LMs

Prompts:

Original model outputs:

Edited model outputs:

*Barack Obama is
married to...*

Michelle Obama, and
he and his wife have
two daughters: Sasha
and Malia.

*Who painted
the Mona Lisa?*

The Mona Lisa was
painted by **Leonardo
da Vinci** in the early
1500s.

Editing Generative LMs

Prompts:

Original model outputs:

Edited model outputs:

Apply
Edits

*Barack Obama is
married to...*

Michelle Obama, and
he and his wife have
two daughters: Sasha
and Malia.

[Obama is married
to] -> [Taylor Swift]

*Who painted
the Mona Lisa?*

The Mona Lisa was
painted by **Leonardo
da Vinci** in the early
1500s.

[Mona Lisa was painted
by] -> [Picasso]

Editing Generative LMs

Prompts:

Original model outputs:

Edited model outputs:

Apply Edits

Barack Obama is married to...

Michelle Obama, and he and his wife have two daughters: Sasha and Malia.

[Obama is married to] -> [Taylor Swift]

Taylor Swift. They met on the set of "The Hunger Games" in 2011.

Who painted the Mona Lisa?

The Mona Lisa was painted by **Leonardo da Vinci** in the early 1500s.

[Mona Lisa was painted by] -> [Picasso]

The Mona Lisa is a painting by **Picasso**. It was painted in the early 1960s.

Editing GANs

$$\mathbf{z}_1 \sim \text{Unif}([0,1]^d)$$

$$\mathbf{z}_2 \sim \text{Unif}([0,1]^d)$$

Editing GANs

Original model outputs:

$$z_1 \sim \text{Unif}([0,1]^d)$$



$$z_2 \sim \text{Unif}([0,1]^d)$$



Editing GANs

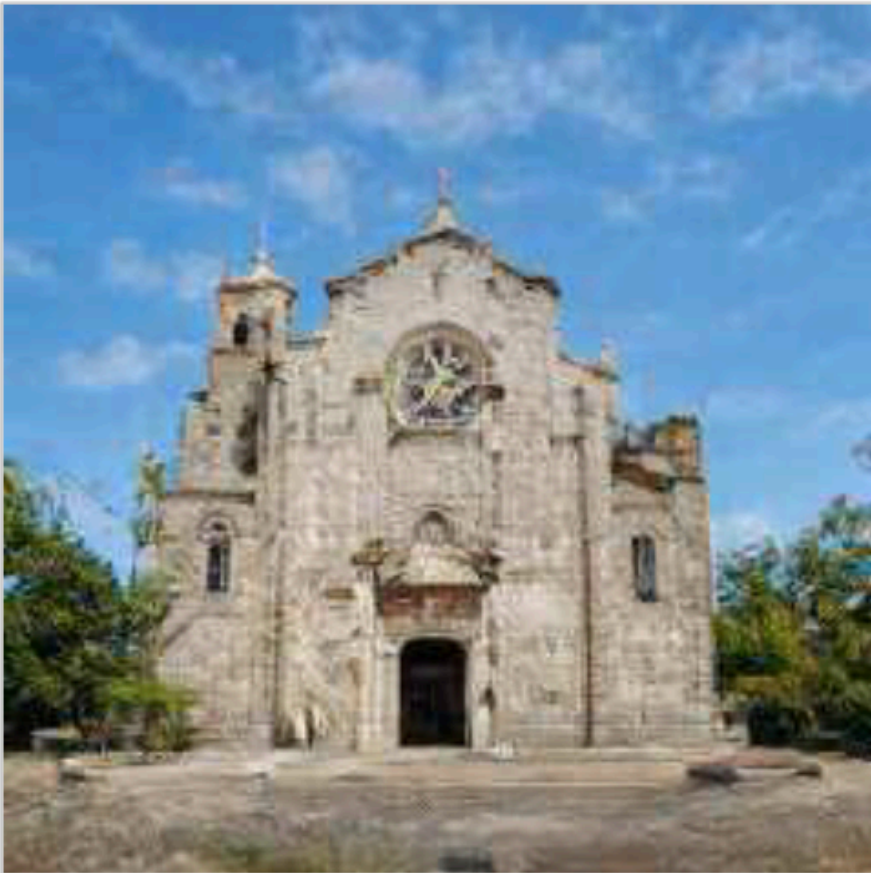
Original model outputs:



$z_1 \sim \text{Unif}([0,1]^d)$

$z_2 \sim \text{Unif}([0,1]^d)$

Edited model outputs:



Editing GANs

Original model outputs:

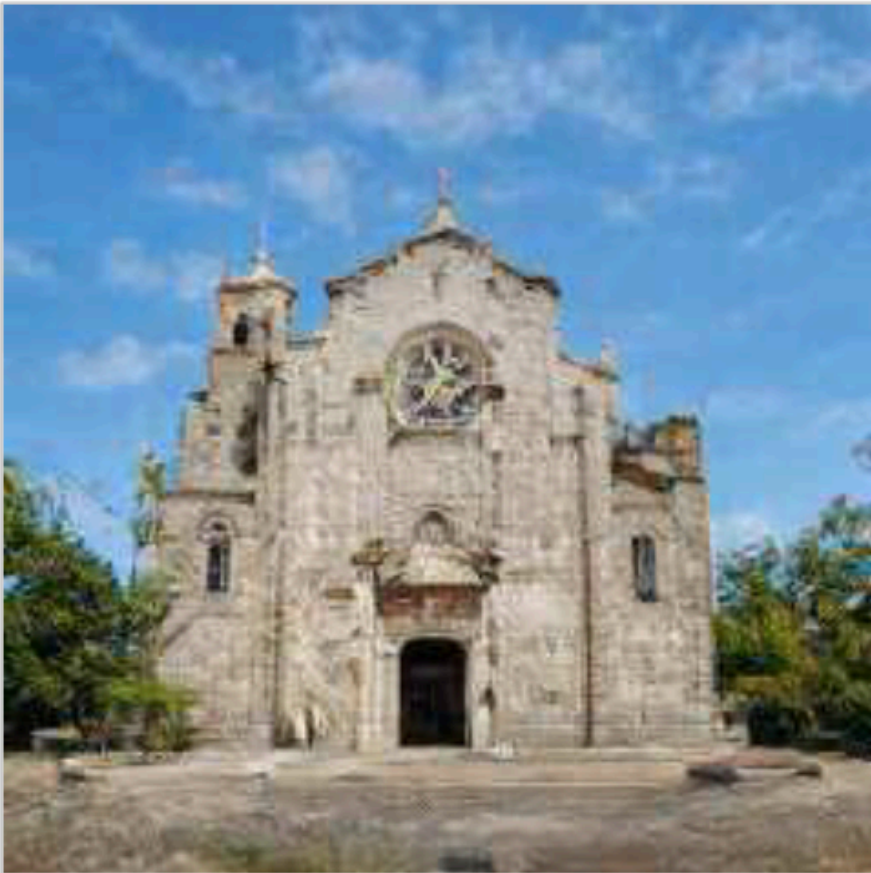


$z_1 \sim \text{Unif}([0,1]^d)$

$z_2 \sim \text{Unif}([0,1]^d)$

**Removed
all text!**

Edited model outputs:



Editing Image Classifiers

Inputs:



Editing Image Classifiers

Inputs:

Original model outputs:



Snowplow ✘



Snowplow ✘



Snowmobile ✘



Amphibian ✘

Editing Image Classifiers

Inputs:

Original model outputs:

Edited model outputs:



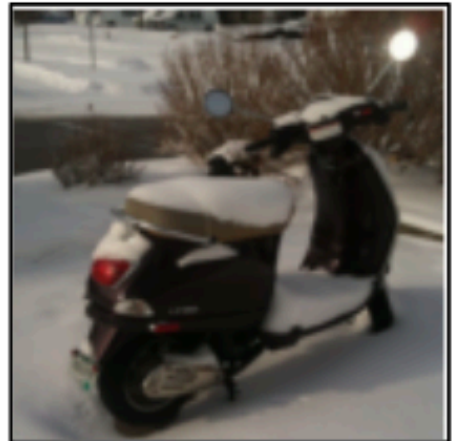
Snowplow ❌

Traffic light ✅



Snowplow ❌

Car wheel ✅



Snowmobile ❌

Motor scooter ✅



Amphibian ❌

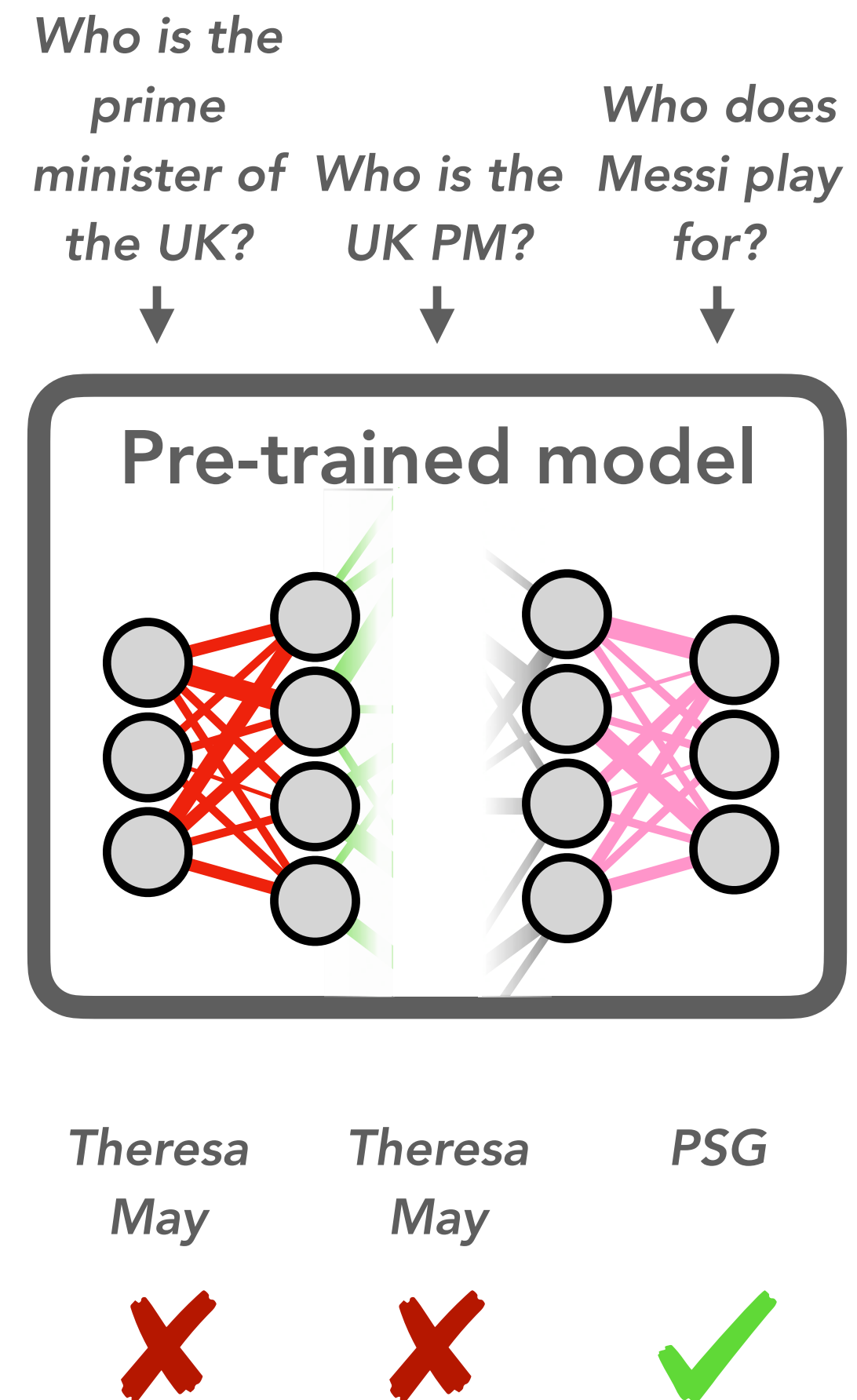
Racer ✅

What makes editing hard?

Need to make a “local” update

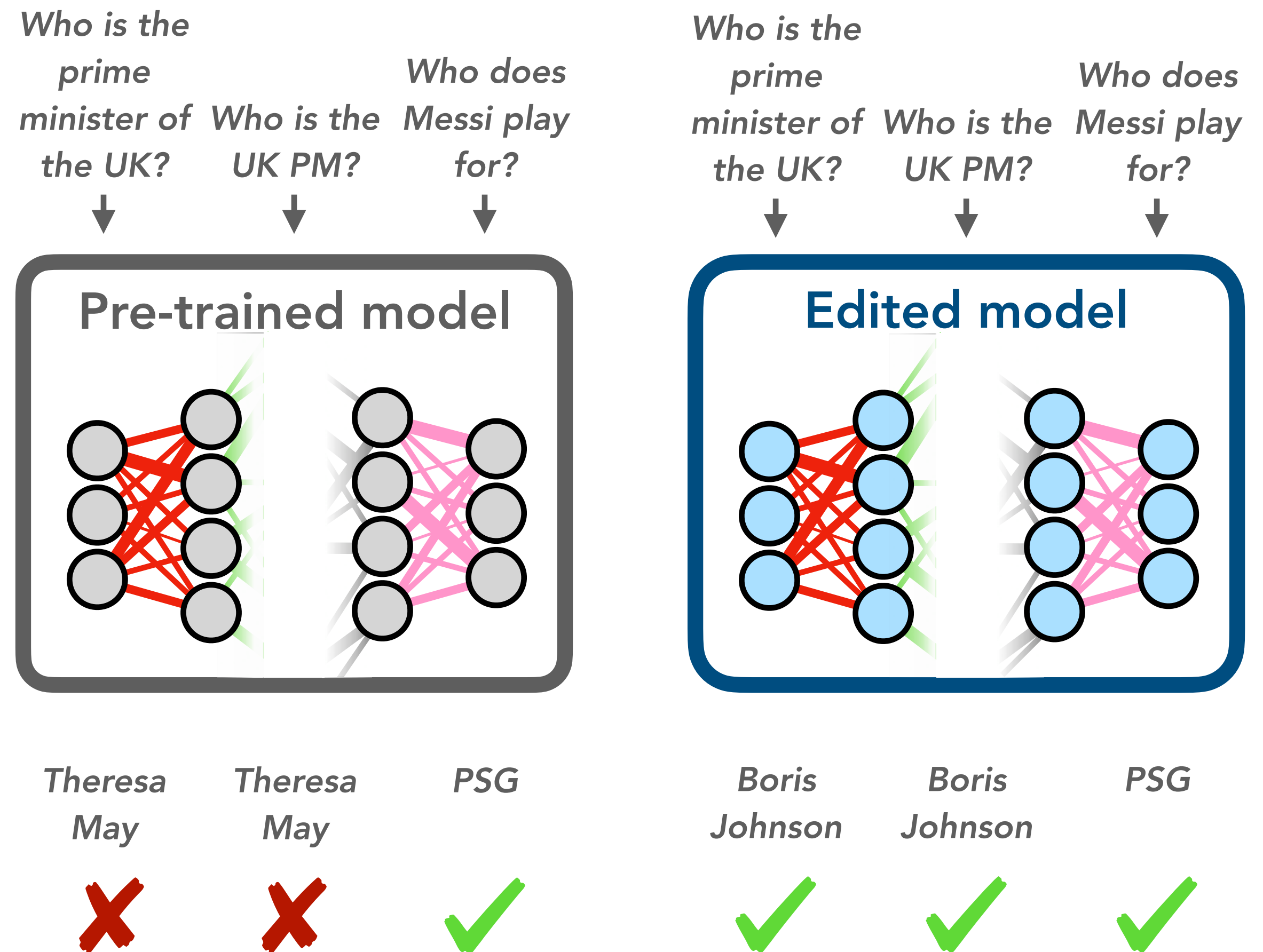
What makes editing hard?

Need to make a “local” update



What makes editing hard?

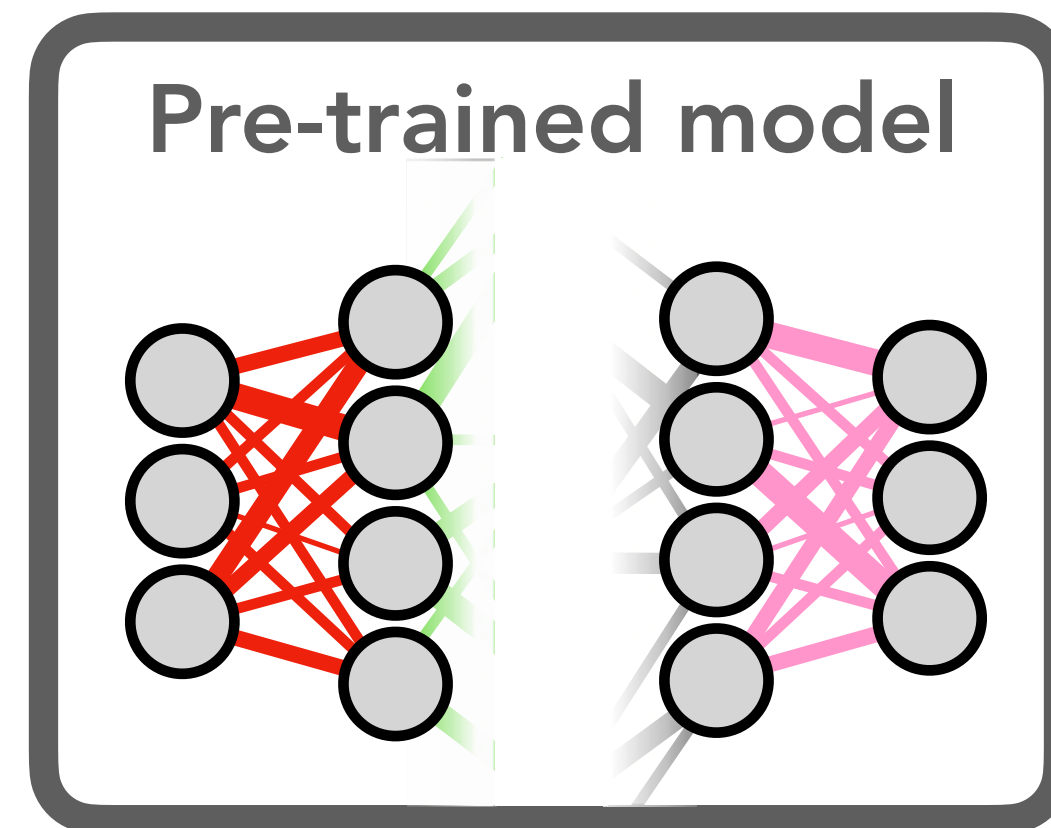
Need to make a “local” update



What makes editing hard?

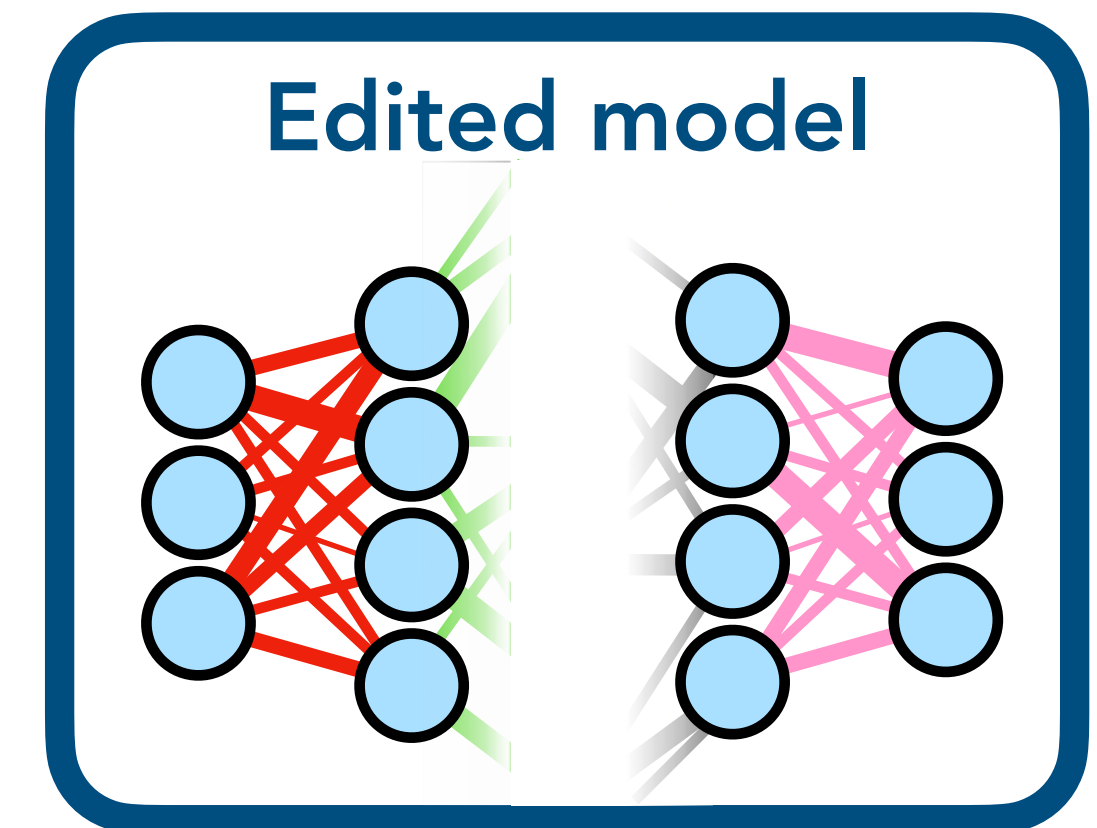
Need to make a “local” update
Not too local... (**undergeneralize**)

Who is the
prime
minister of
the UK?
↓
Who is the
UK PM?
↓
Who does
Messi play
for?
↓



Theresa
May
✗
Theresa
May
✗
PSG
✓

Who is the
prime
minister of
the UK?
↓
Who is the
UK PM?
↓
Who does
Messi play
for?
↓



Boris
Johnson
✓
Theresa
May
✗
PSG
✓

What makes editing hard?

Need to make a “local” update

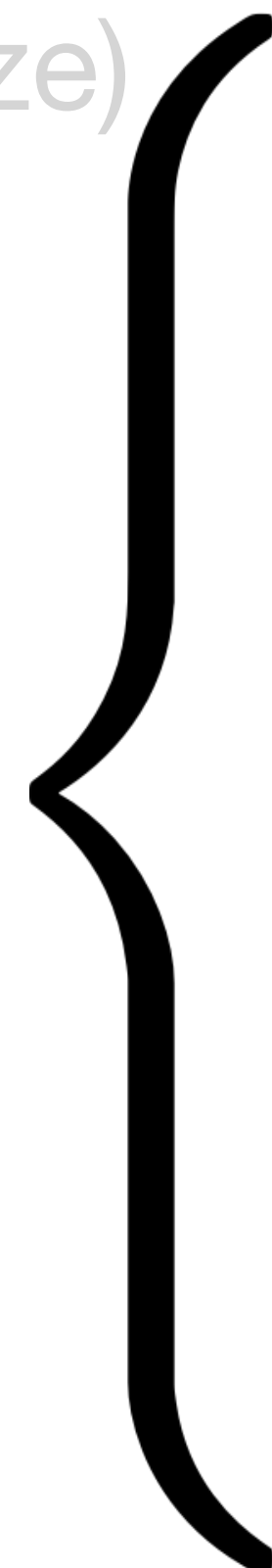
Not too local... (**undergeneralize**)

...but not too general
(**overgeneralize**)

Many ways of **specifying** the
intended post-edit behavior

(what information do we assume
access to when applying an edit?)

Explicit descriptors are desired input-output pairs:



What makes editing hard?

Need to make a “local” update

Not too local... (**undergeneralize**)

...but not too general
(**overgeneralize**)

Many ways of **specifying** the
intended post-edit behavior

(what information do we assume
access to when applying an edit?)

Explicit descriptors are desired input-output pairs:

*“Thoughts on vaccines? **They’re really important for public...**”*

*“Who is the UK prime minister? **Boris Johnson**”*

*“True or false: Messi plays for PSG. **True**”*

What makes editing hard?

Need to make a “local” update

Not too local... (**undergeneralize**)

...but not too general
(**overgeneralize**)

Many ways of **specifying** the
intended post-edit behavior

(what information do we assume
access to when applying an edit?)

Explicit descriptors are desired input-output pairs:

*“Thoughts on vaccines? **They’re really important for public...**”*

*“Who is the UK prime minister? **Boris Johnson**”*

*“True or false: Messi plays for PSG. **True**”*

Implicit descriptors simply describe the desired change:

What makes editing hard?

Need to make a “local” update

Not too local... (**undergeneralize**)

...but not too general
(**overgeneralize**)

Many ways of **specifying** the
intended post-edit behavior

(what information do we assume
access to when applying an edit?)

Explicit descriptors are desired input-output pairs:

*“Thoughts on vaccines? **They’re really important for public...**”*

*“Who is the UK prime minister? **Boris Johnson**”*

*“True or false: Messi plays for PSG. **True**”*

Implicit descriptors simply describe the desired change:

“Be more positive about vaccines.”

“Boris Johnson is the UK PM.”

“Messi plays for PSG.”

What makes editing hard?

Need to make a “local” update

Not too local... (**undergeneralize**)

...but not too general
(**overgeneralize**)

Many ways of **specifying** the
intended post-edit behavior

(what information do we assume
access to when applying an edit?)

Explicit descriptors are desired input-output pairs:

“Thoughts on vaccines? **They’re really important for public...**”

“Who is the UK prime minister? **Boris Johnson**”

“True or false: Messi plays for PSG. **True**”

Implicit descriptors simply describe the desired change:

“Be more positive about vaccines.”

“Boris Johnson is the UK PM.”

“Messi plays for PSG.”

Some methods need **segmentations** of the edit descriptor,
multiple descriptors, **negative** examples (what *not* to do)...

What makes editing hard?

Need to make a “local” update

Not too local... (**undergeneralize**)

...but not too general
(**overgeneralize**)

Many ways of **specifying** the
intended post-edit behavior

(what information do we assume
access to when applying an edit?)

Explicit descriptors are desired input-output pairs:

*“Thoughts on vaccines? **They’re really important for public...**”*

*“Who is the UK prime minister? **Boris Johnson**”*

*“True or false: Messi plays for PSG. **True**”*

Implicit descriptors simply describe the desired change:

“Be more positive about vaccines.”

“Boris Johnson is the UK PM.”

“Messi plays for PSG.”

Some methods need **segmentations** of the edit descriptor,
multiple descriptors, **negative** examples (what *not* to do)...

Lots of design decisions!

Edit *what*, exactly?

Defining the problem

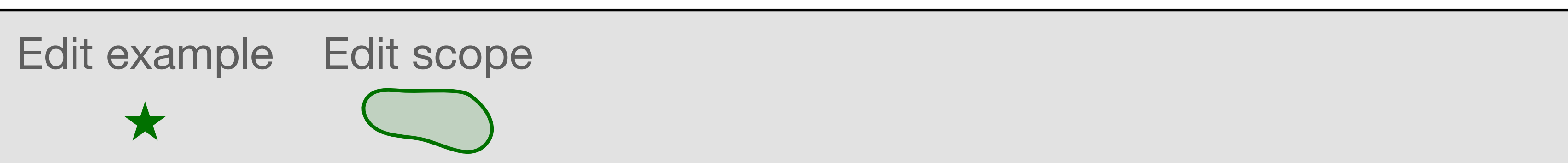
★
*Who is the prime
minister of the UK?*

Edit example



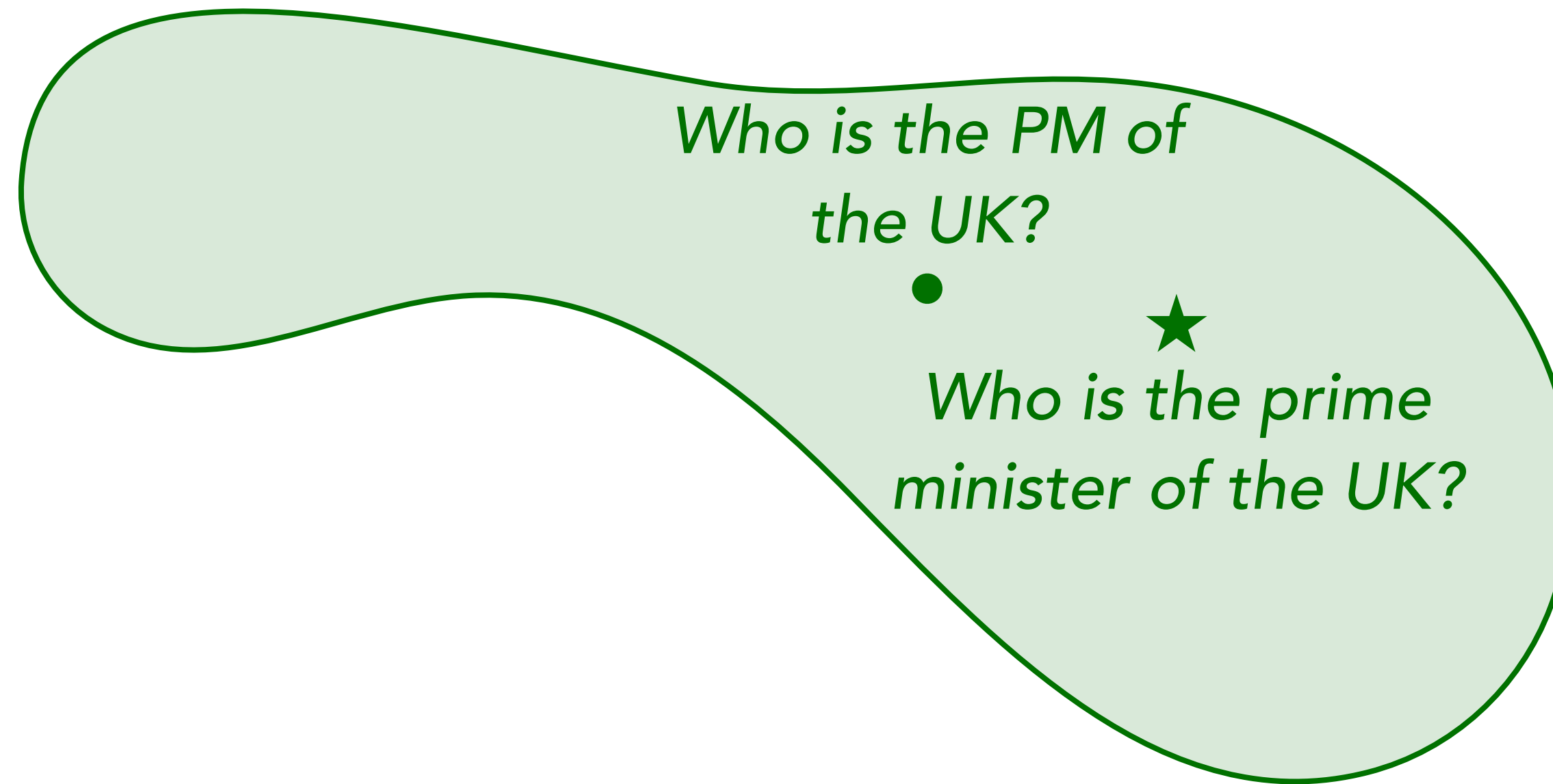
Edit *what*, exactly?


Defining the problem



Edit *what*, exactly?

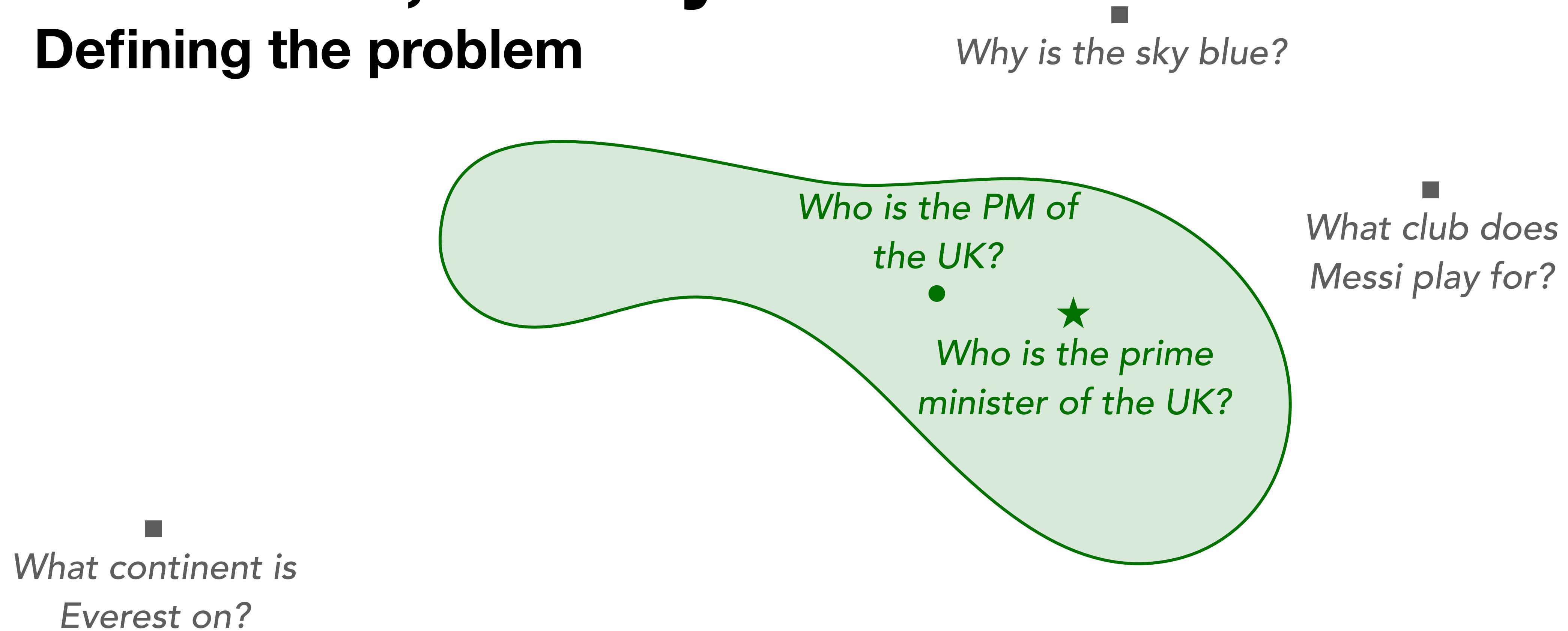
Defining the problem




Edit example	Edit scope	In-scope
★		●

Edit *what*, exactly?

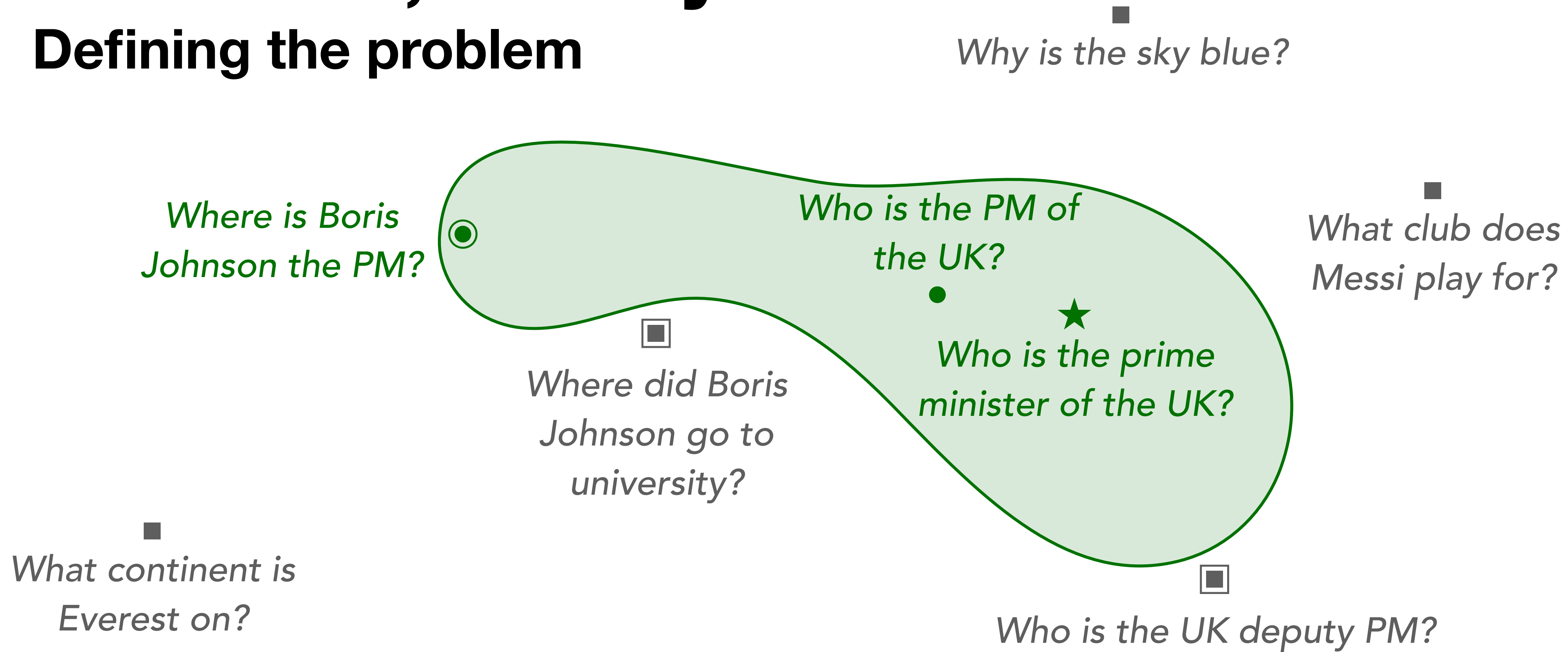
Defining the problem



Edit example	Edit scope	In-scope	Out-of-scope
★		●	■

Edit *what*, exactly?

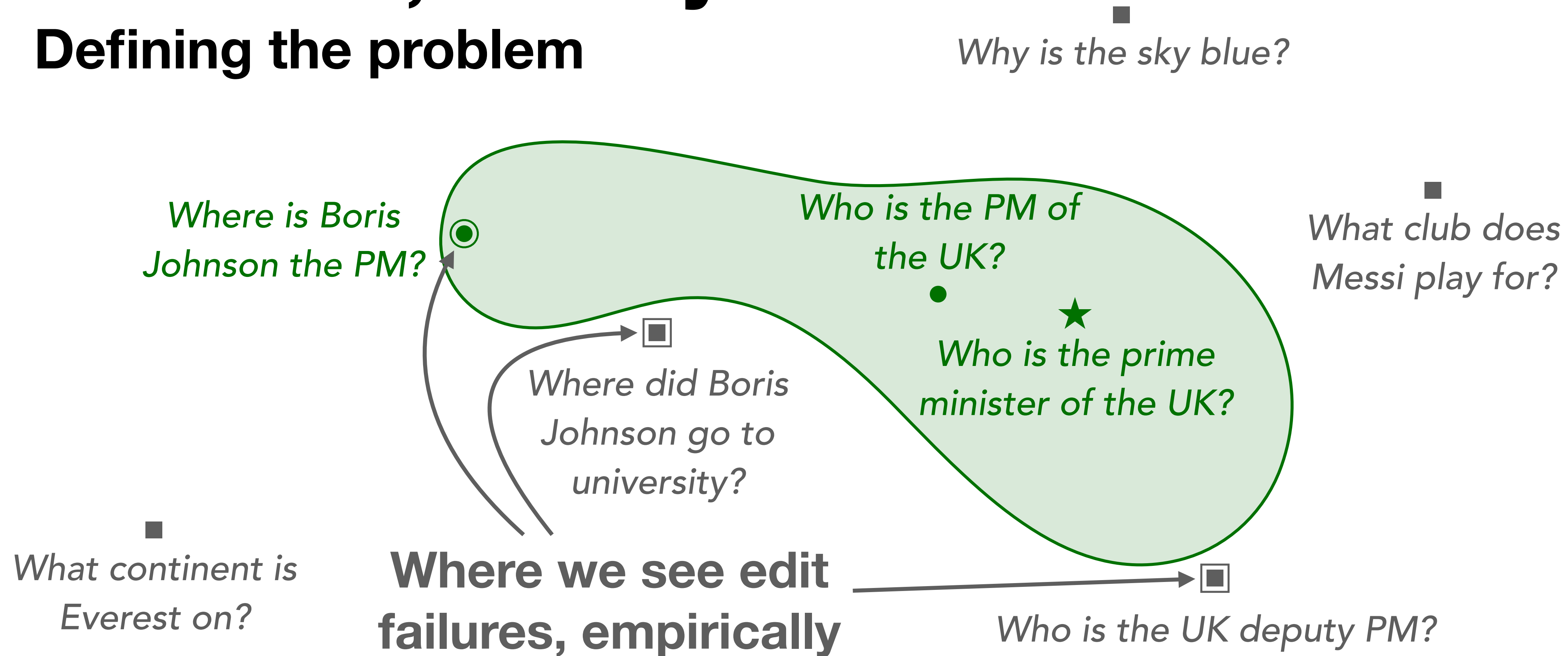
Defining the problem



Edit example	Edit scope	In-scope	Out-of-scope	Hard in/out-of-scope
★		●	■	● □

Edit *what*, exactly?

Defining the problem

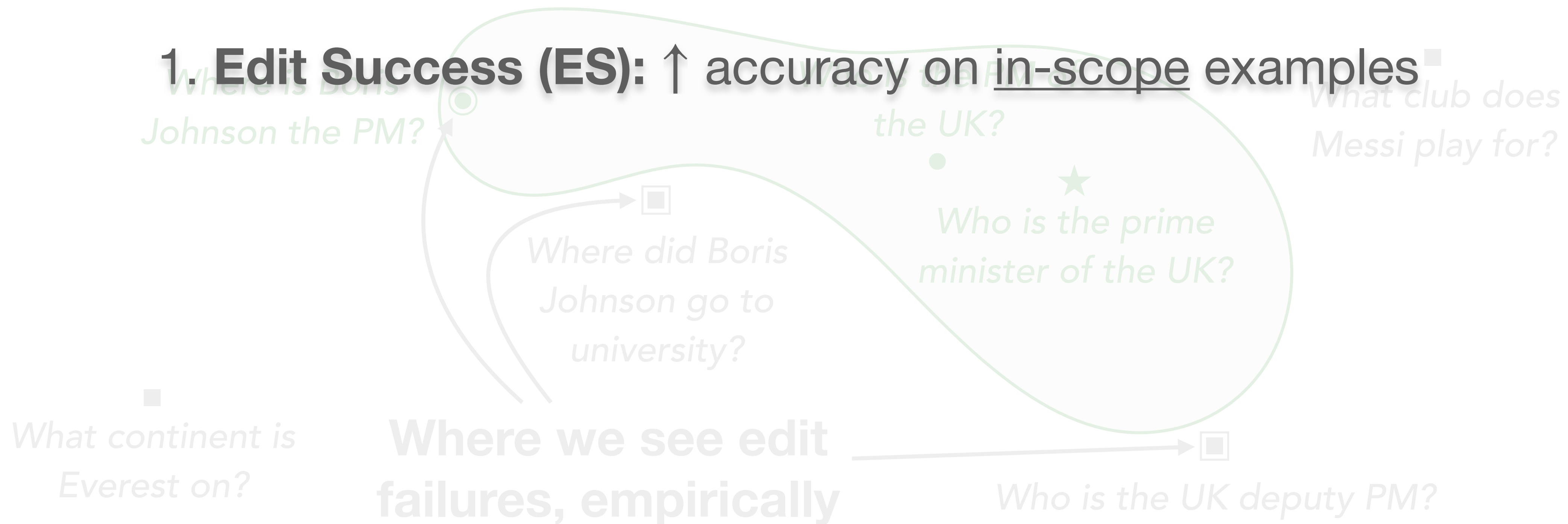


Edit example	Edit scope	In-scope	Out-of-scope	Hard in/out-of-scope
★		●	■	● □

Edit *what*, exactly?

Metrics for evaluating model edits

1. Edit Success (ES): \uparrow accuracy on in-scope examples

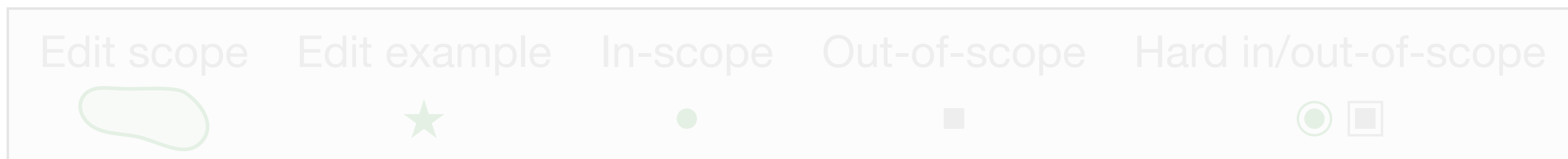
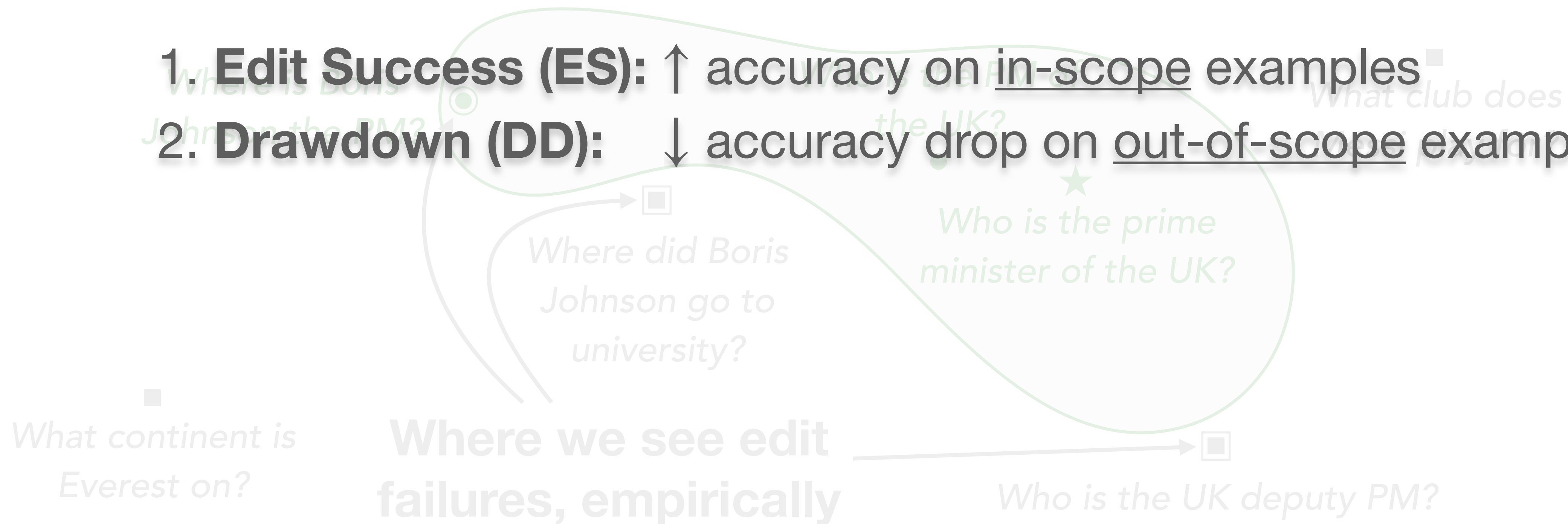


Edit scope	Edit example	In-scope	Out-of-scope	Hard in/out-of-scope

Edit *what*, exactly?

Metrics for evaluating model edits

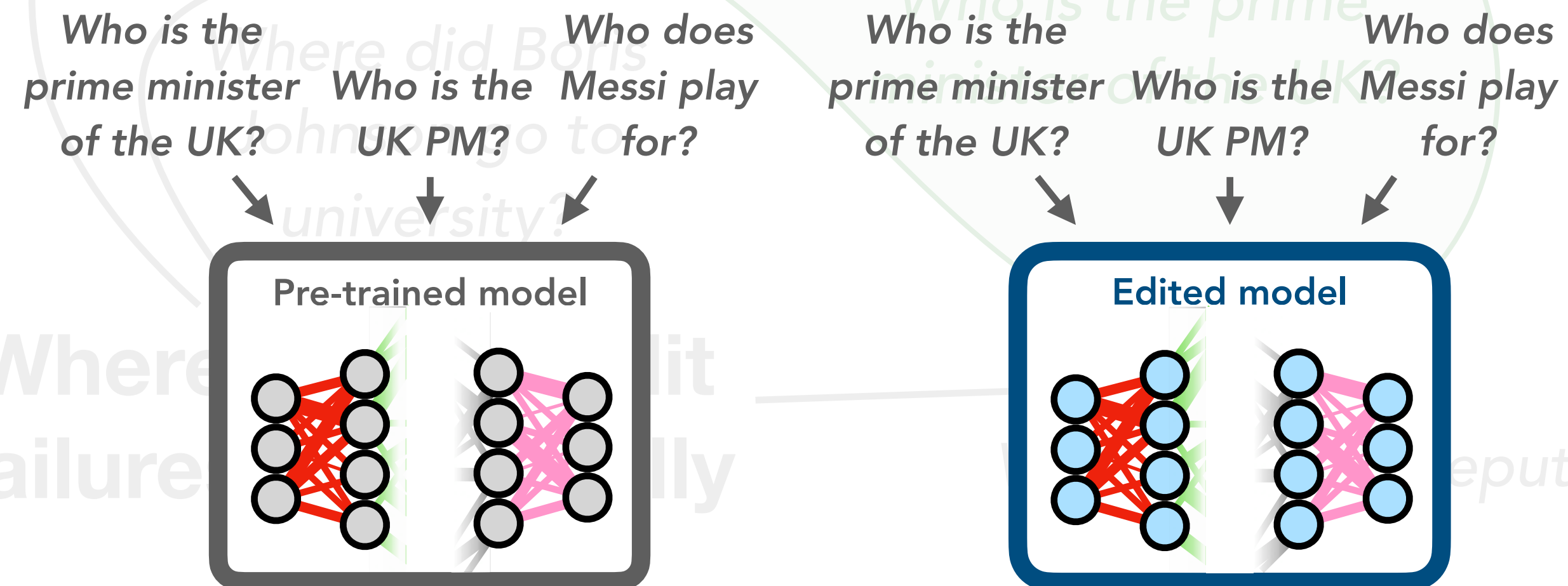
1. **Edit Success (ES):** ↑ accuracy on in-scope examples
2. **Drawdown (DD):** ↓ accuracy drop on out-of-scope examples



Edit *what*, exactly?

Metrics for evaluating model edits

1. **Edit Success (ES):** \uparrow accuracy on in-scope examples
2. **Drawdown (DD):** \downarrow accuracy drop on out-of-scope examples

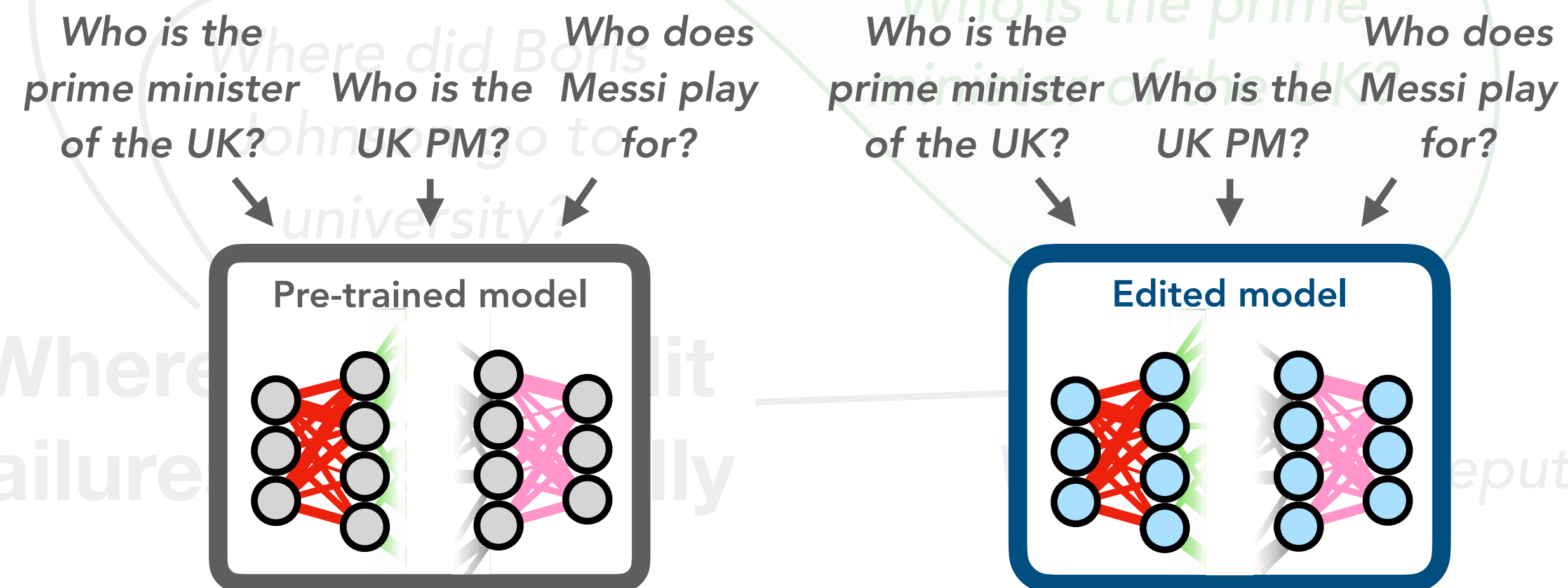


Model	Example	Accuracy
Pre-trained model	Theresa May (Edit scope)	✗
	Theresa May (Edit example)	✗
	PSG (In-scope)	✓
Edited model	Boris Johnson (Out-of-scope)	✓
	Boris Johnson (Hard in/out-of-scope)	✓
	PSG (Hard in/out-of-scope)	✓

Edit *what*, exactly?

Metrics for evaluating model edits

1. **Edit Success (ES):** \uparrow accuracy on in-scope examples
2. **Drawdown (DD):** \downarrow accuracy drop on out-of-scope examples



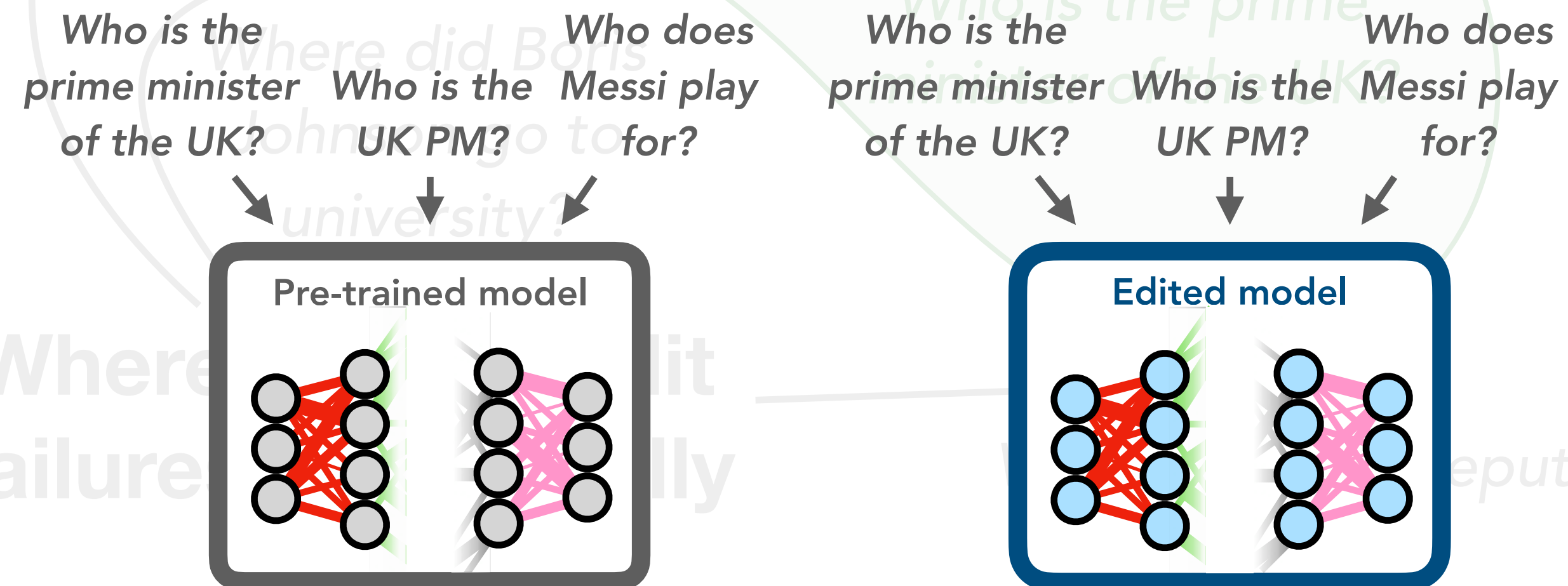
ES: 1
DD: 0

Category	Example	Result
Edit scope	Theresa May	✗
	Theresa May	✗
	PSG	✓
In-scope	Boris Johnson	✓
	Boris Johnson	✓
	PSG	✓
Out-of-scope	Boris Johnson	✓
	Boris Johnson	✓
	PSG	✓

Edit *what*, exactly?

Metrics for evaluating model edits

1. **Edit Success (ES):** \uparrow accuracy on in-scope examples
2. **Drawdown (DD):** \downarrow accuracy drop on out-of-scope examples



ES: 0.5

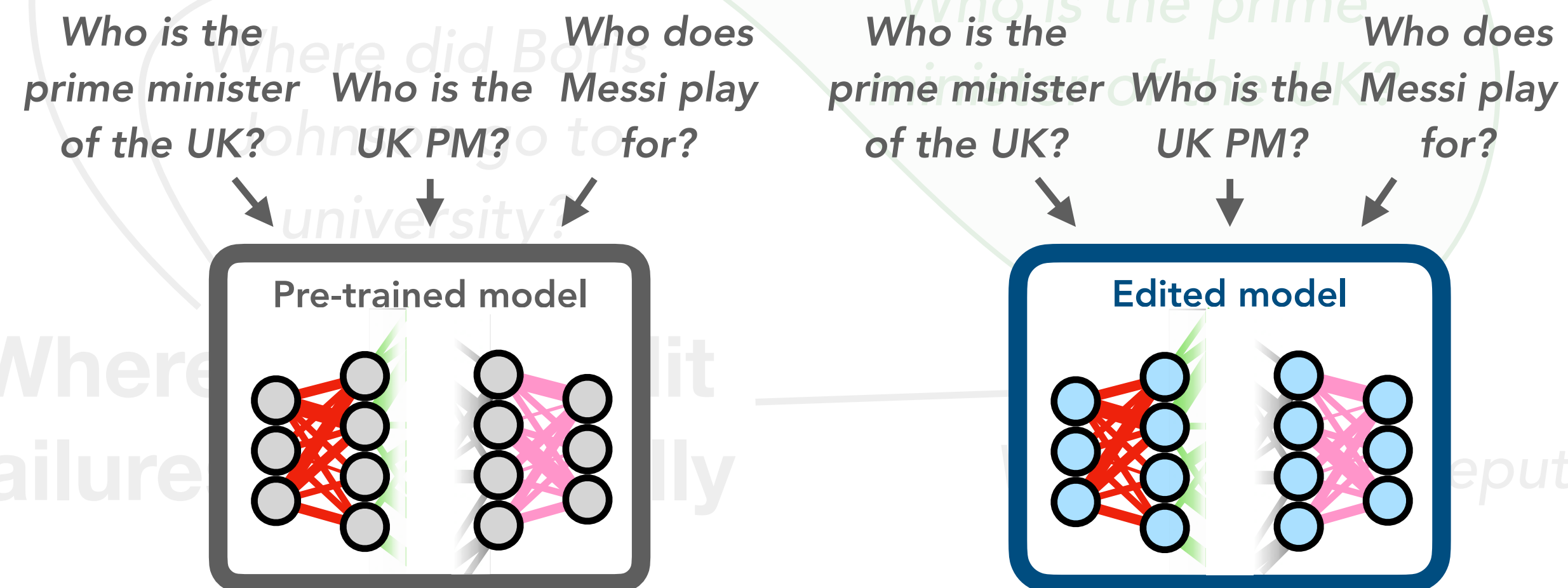
DD: 0

Example	Theresa May	Theresa May	PSG	Boris Johnson	Theresa May	PSG
Edit scope						
Edit example		★				
In-scope	✗	✗	✓			
Out-of-scope				✓	✗	
Hard in/out-of-scope						○ □

Edit *what*, exactly?

Metrics for evaluating model edits

1. **Edit Success (ES):** \uparrow accuracy on in-scope examples
2. **Drawdown (DD):** \downarrow accuracy drop on out-of-scope examples



ES: 1
DD: 1

Category	Example	Result
Edit scope	Theresa May	✗
	Theresa May	✗
	PSG	✓
In-scope	Boris Johnson	✓
	Boris Johnson	✓
	Boris Johnson	✗

Today's Plan

I. Background

II. Learning to edit NNs

III. Moving editing towards the real world

IV. Future work & open questions

Existing approaches to editing

Some simple baselines

What about just fine-tuning?

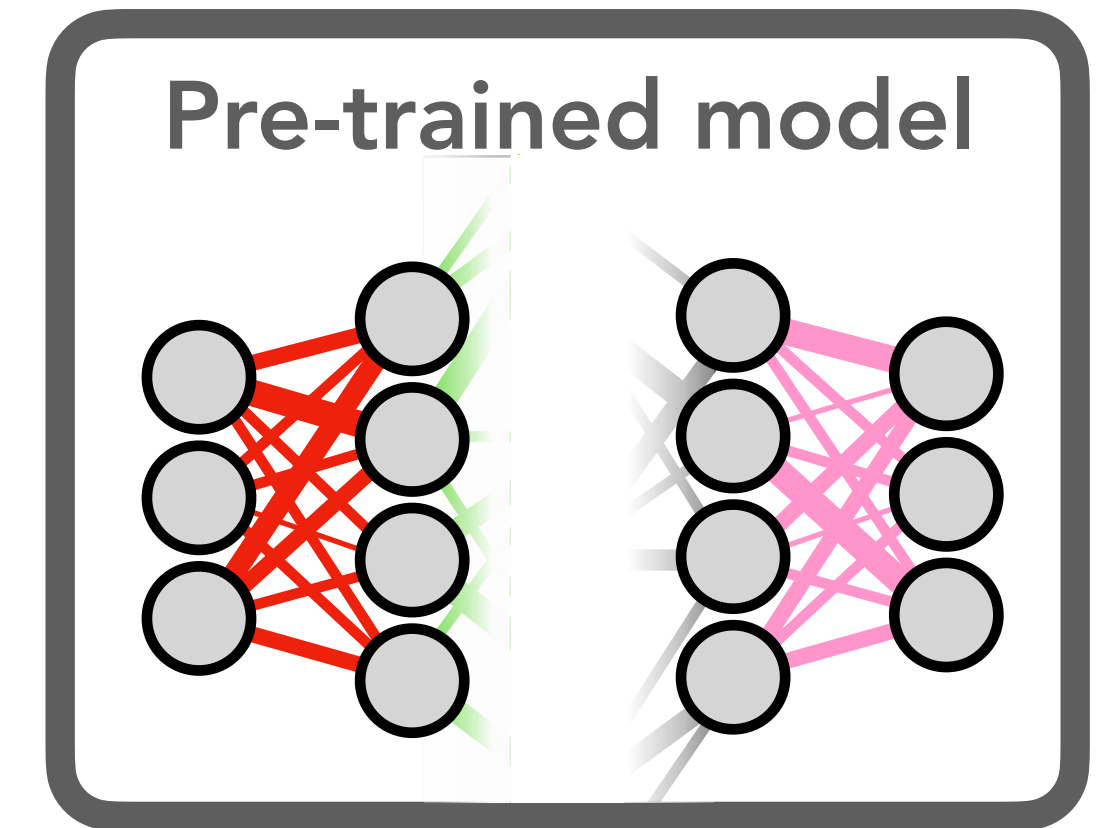
Existing approaches to editing

Some simple baselines

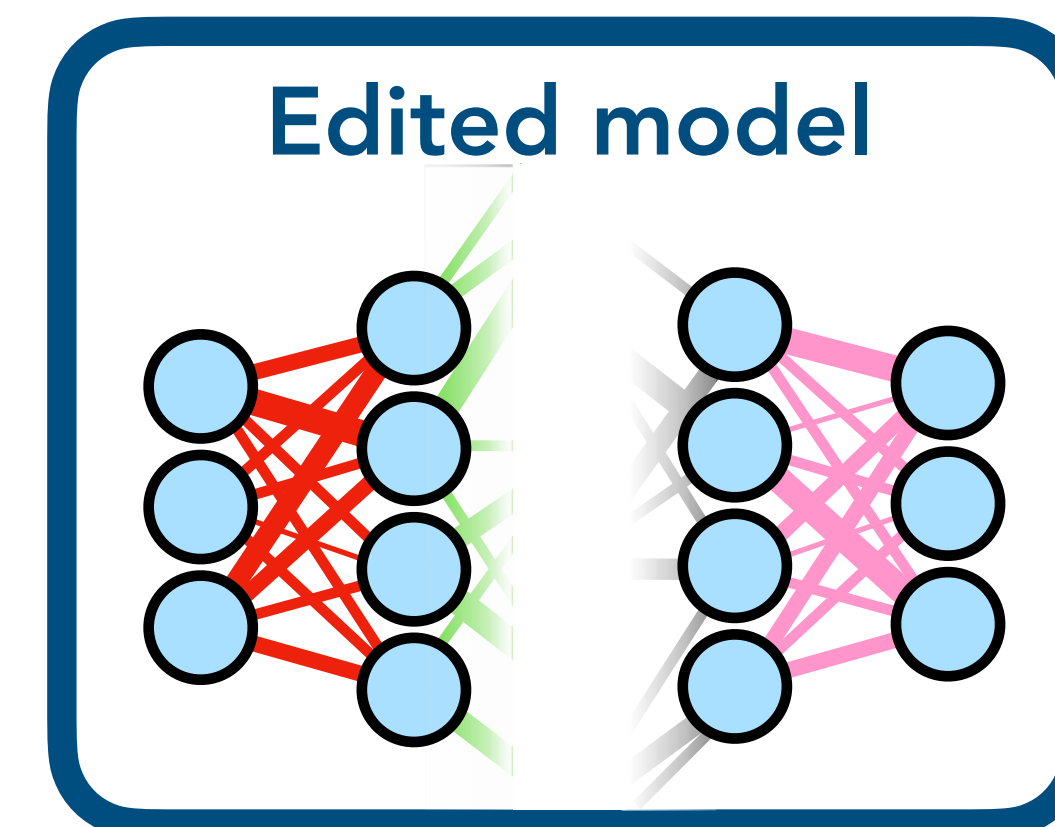
What about just fine-tuning?

$x_e =$ Who is the
prime minister
of the UK?

$y_e =$ Boris
Johnson



Fine-tune



Existing approaches to editing

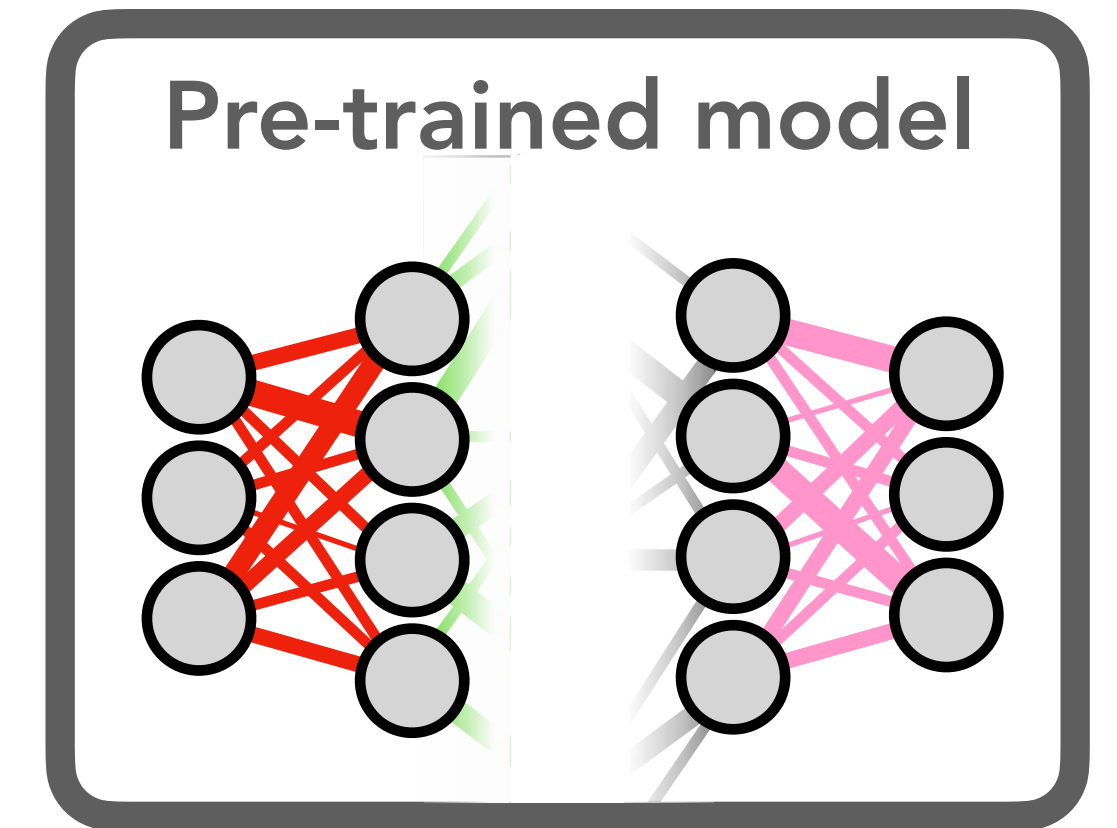
Some simple baselines

What about just fine-tuning?

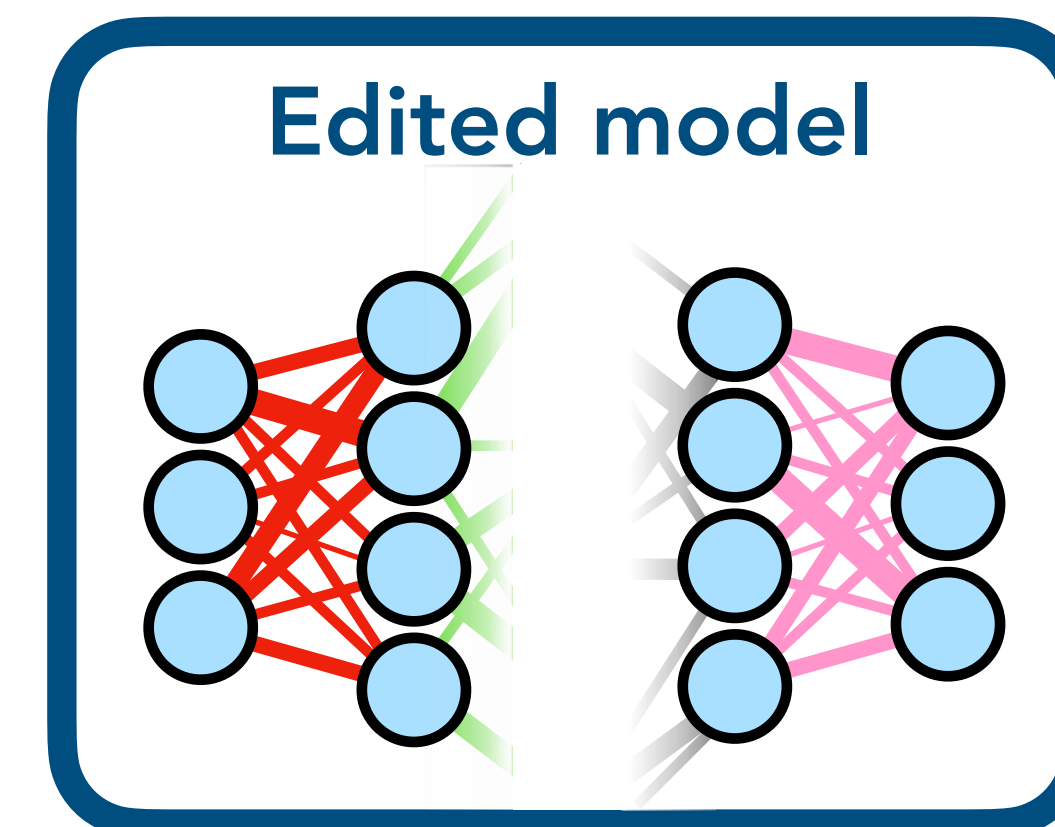
+ simple, universal

$x_e =$ Who is the
prime minister
of the UK?

$y_e =$ Boris
Johnson



Fine-tune



Existing approaches to editing

Some simple baselines

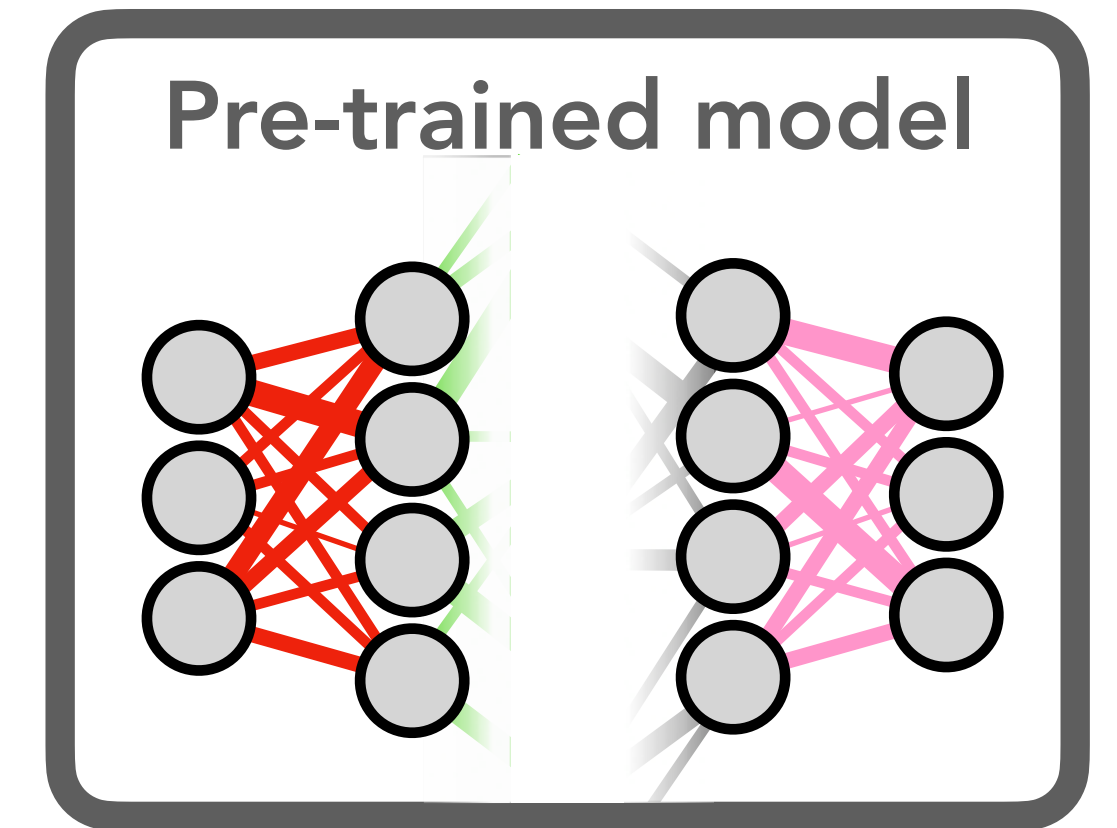
What about just fine-tuning?

+ simple, universal

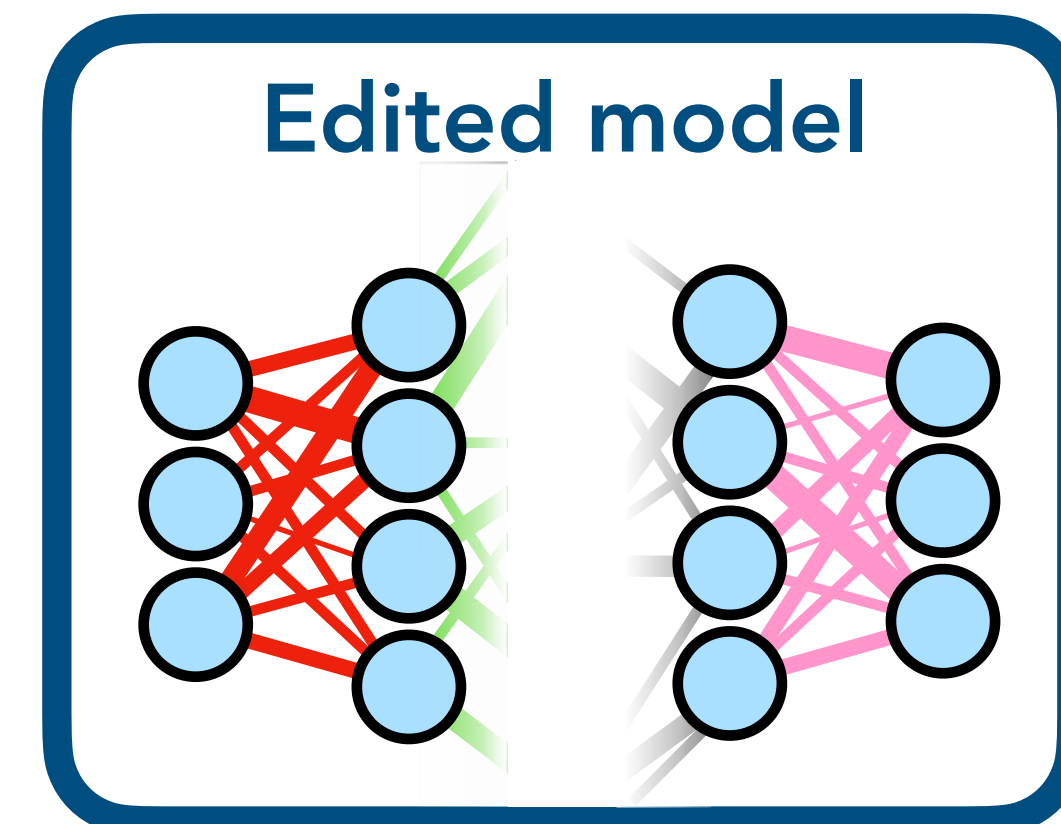
- undergeneralizes, overgeneralizes

$x_e =$ Who is the
prime minister
of the UK?

$y_e =$ Boris
Johnson



Fine-tune



Existing approaches to editing

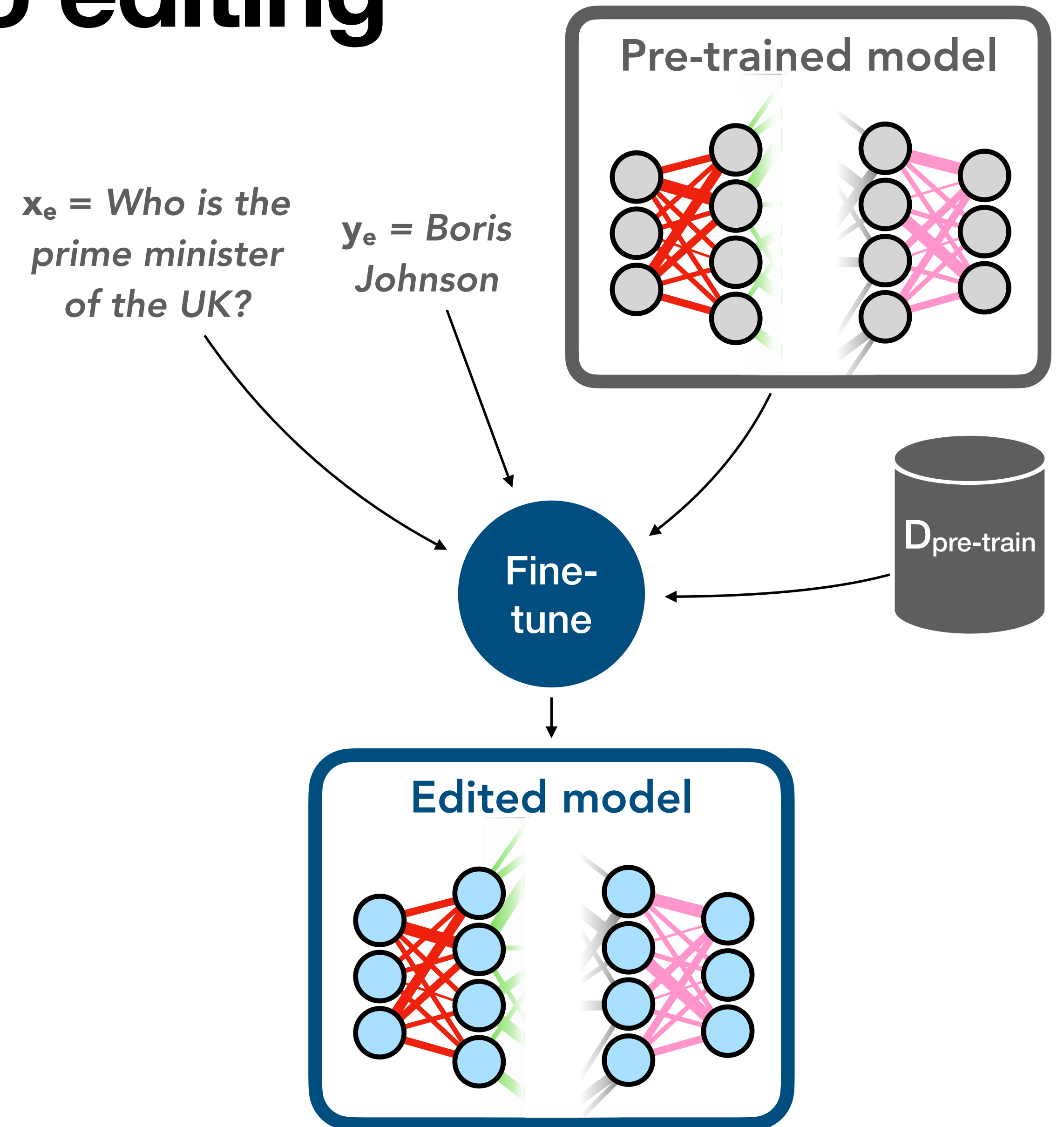
Some simple baselines

What about just fine-tuning?

What if we add some training data?

+ simple, universal

- undergeneralizes, overgeneralizes



Existing approaches to editing

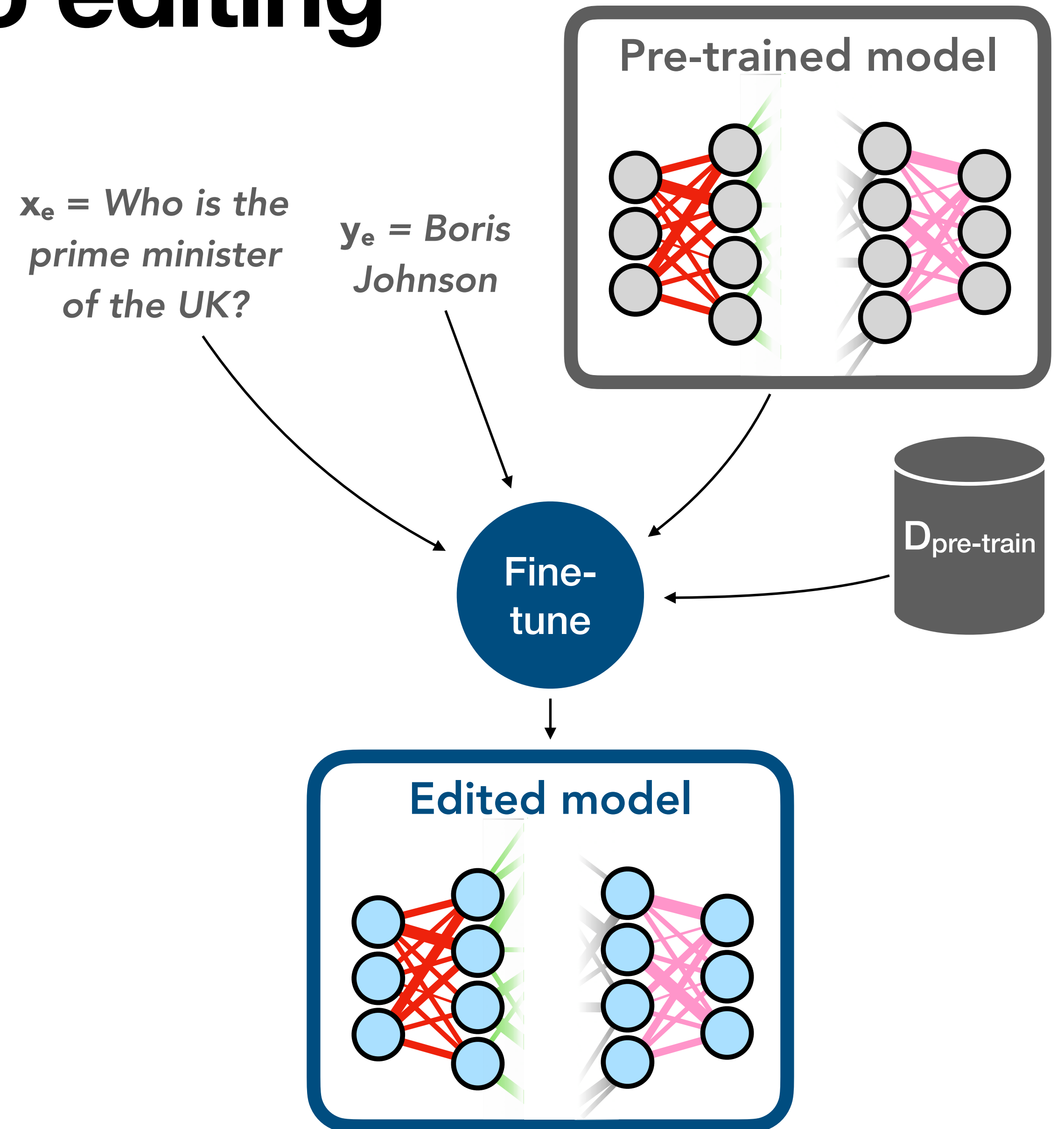
Some simple baselines

What about just fine-tuning?

What if we add some training data?

+ simple, universal

- undergeneralizes, overgeneralizes



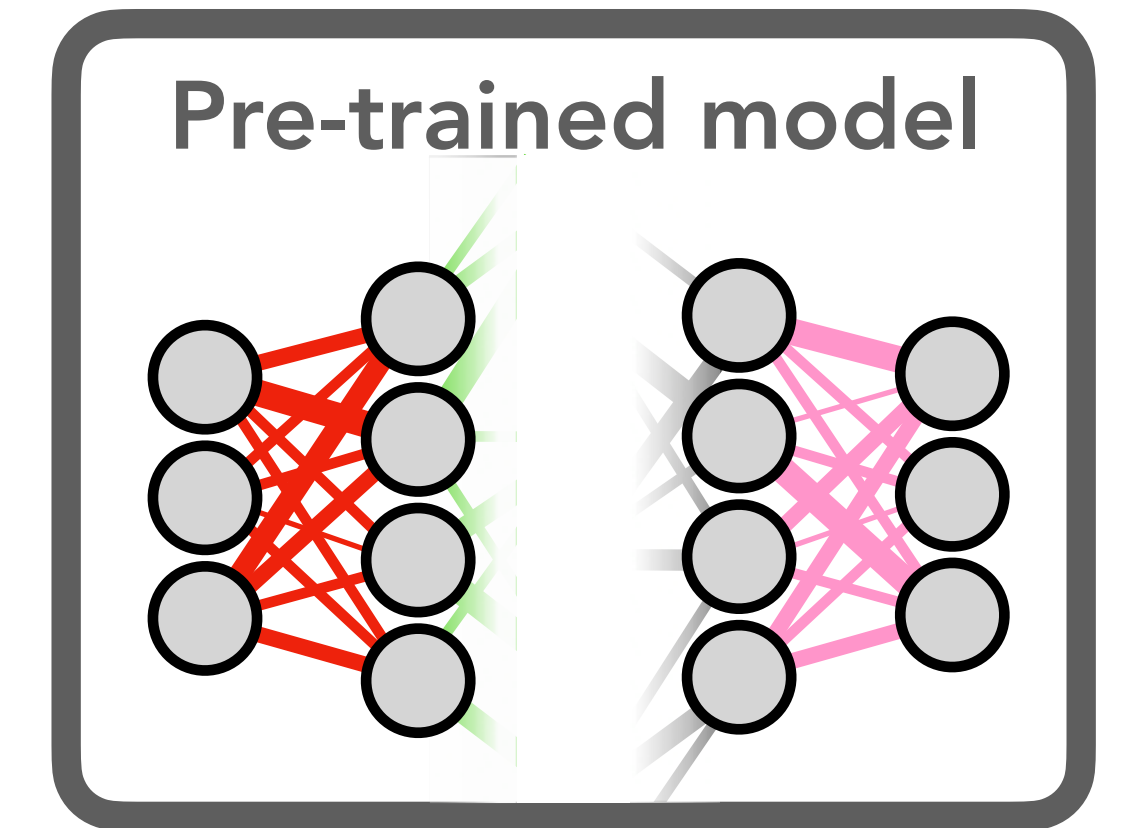
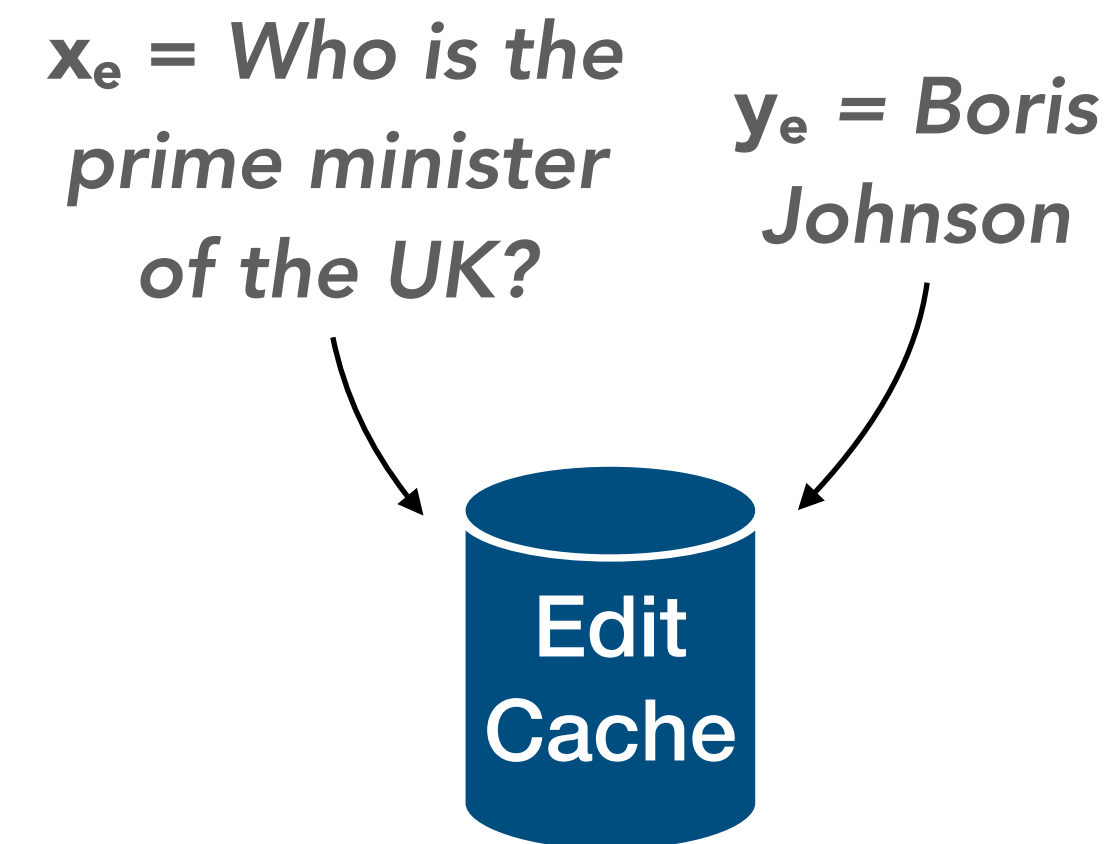
Existing approaches to editing

Some simple baselines

What about just fine-tuning?

What if we add some training data?

Caching edits



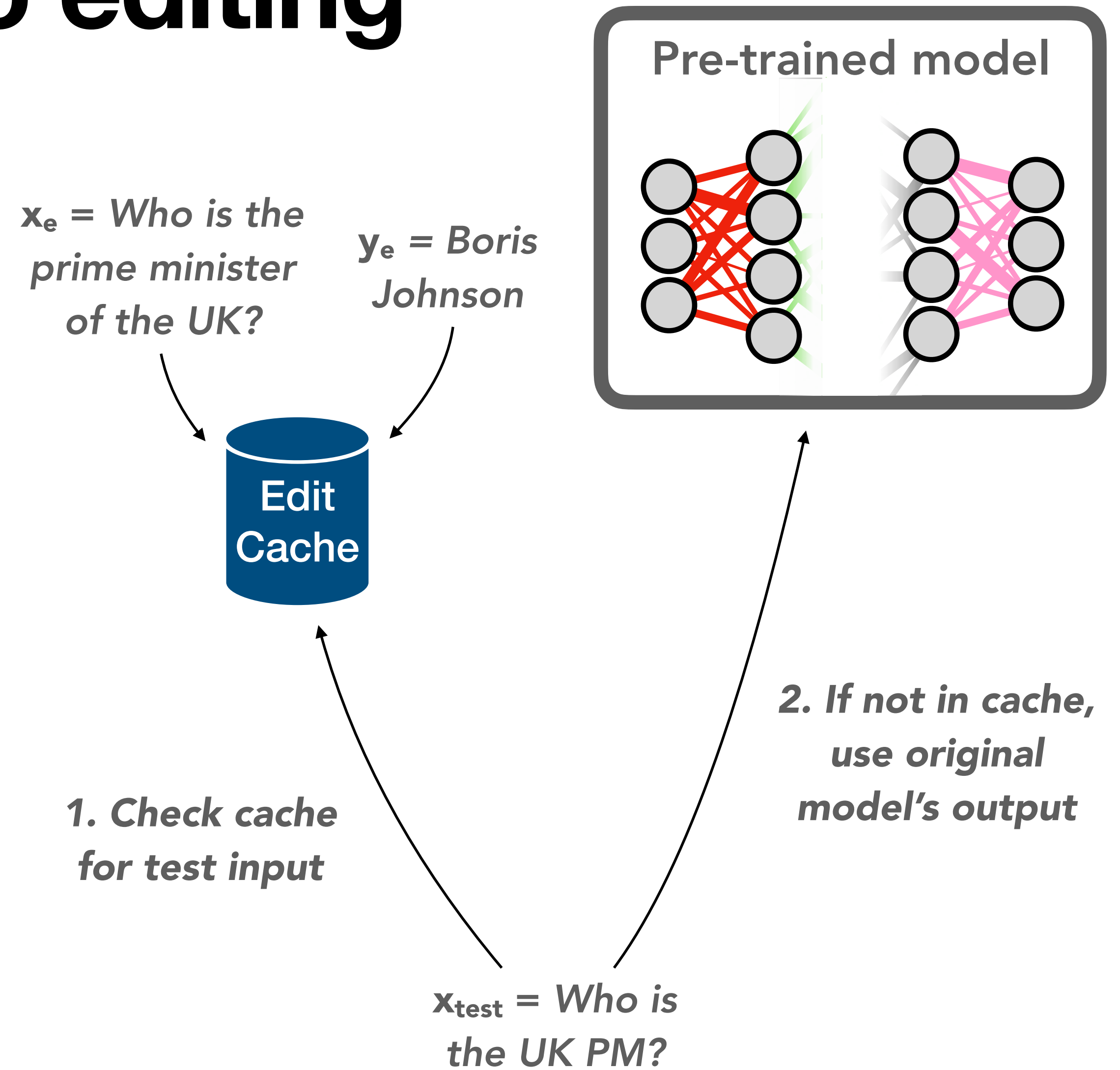
Existing approaches to editing

Some simple baselines

What about just fine-tuning?

What if we add some training data?

Caching edits



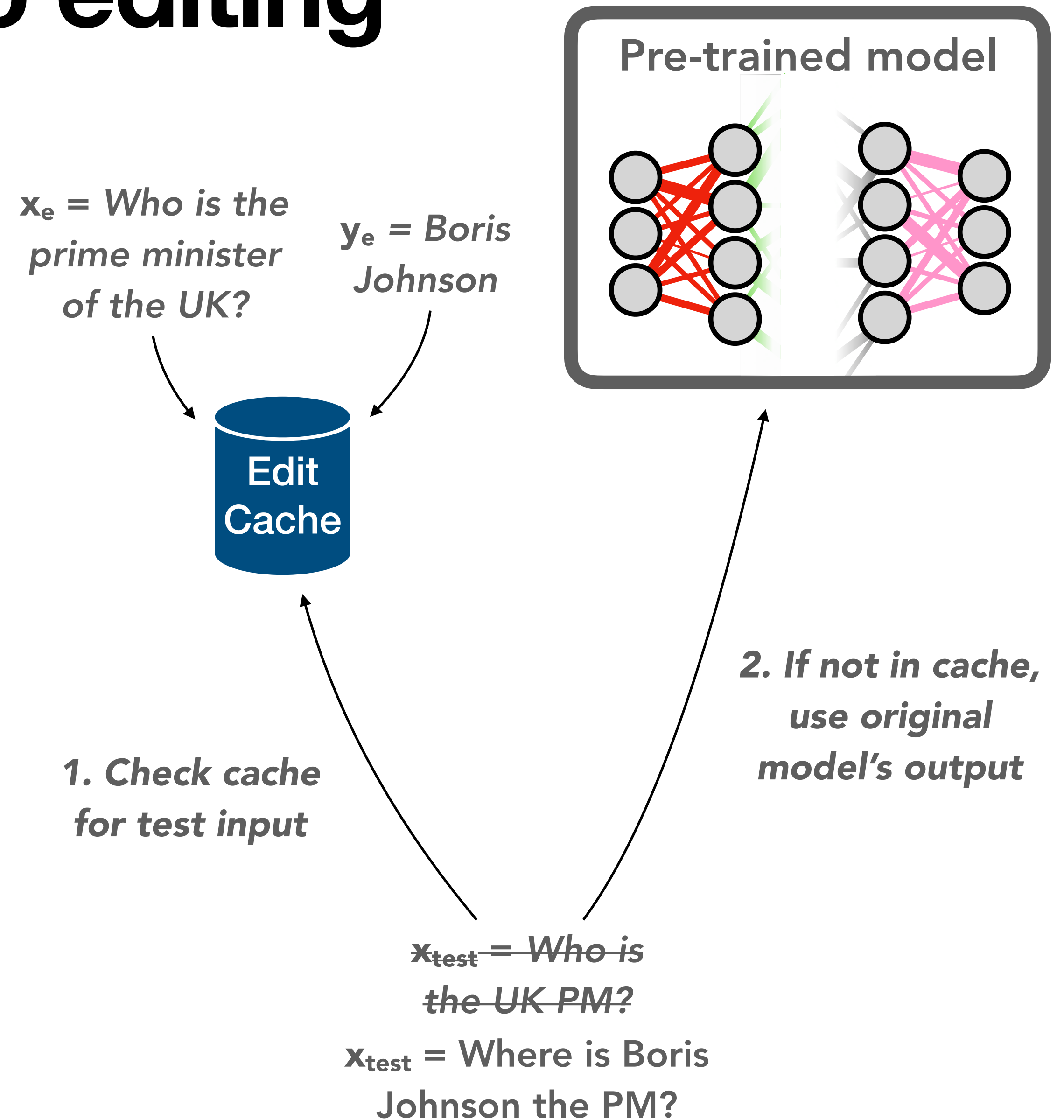
Existing approaches to editing

Some simple baselines

What about just fine-tuning?

What if we add some training data?

Caching edits



Existing approaches to editing

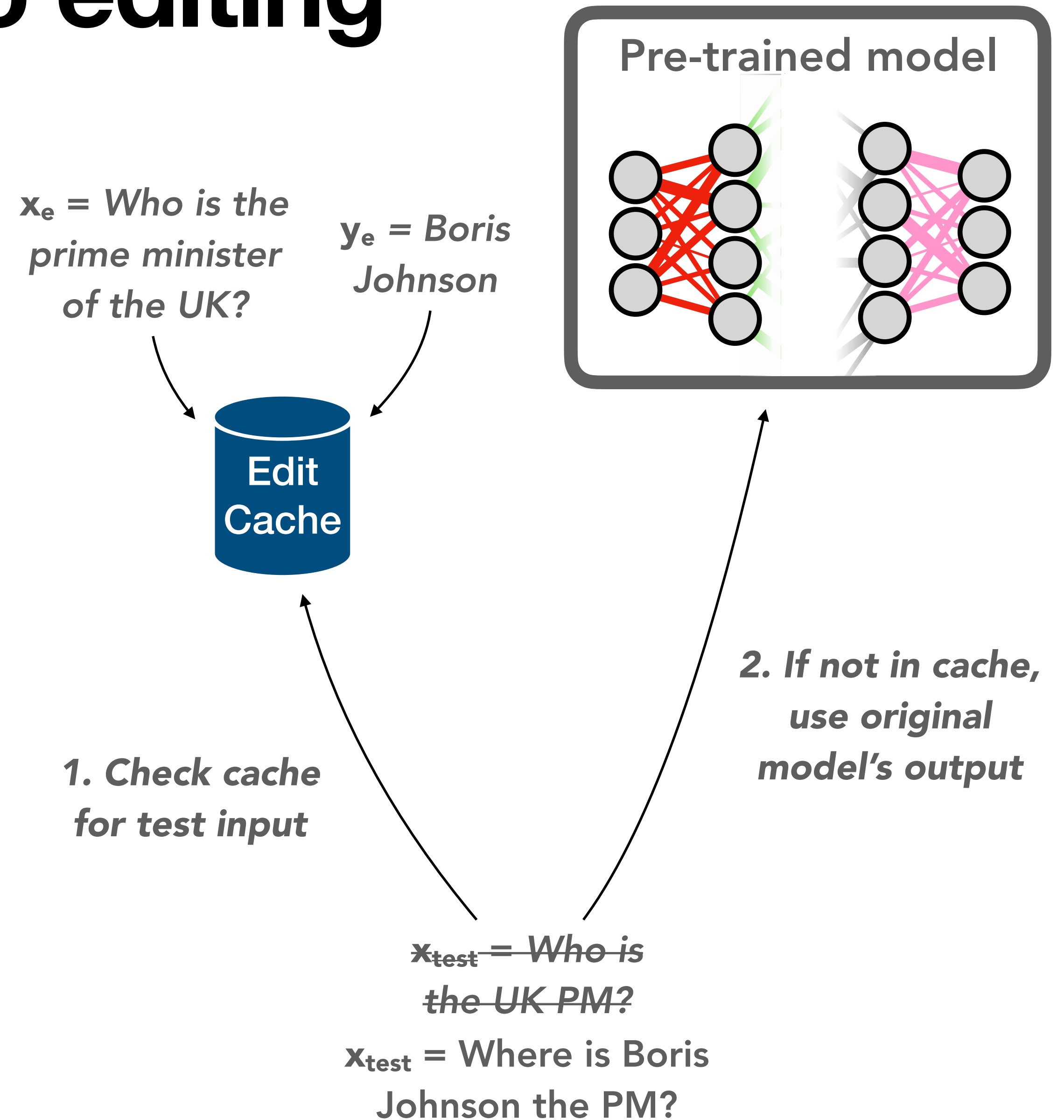
Some simple baselines

What about just fine-tuning?

What if we add some training data?

Caching edits

Can we **learn** a more expressive edit rule from data?



Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $D_{\text{edit}} = \{ (z_{\text{edit}}, x_{\text{loc}}, x_{\text{in}}, y_{\text{in}}) \}$

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $D_{\text{edit}} = \{ (z_{\text{edit}}, x_{\text{loc}}, x_{\text{in}}, y_{\text{in}}) \}$

$z_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$x_{\text{loc}} =$ “What team does Messi play for?”

$x_{\text{in}} =$ “The prime minister of the UK is currently who?”

$y_{\text{in}} =$ “Boris Johnson”

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $D_{\text{edit}} = \{ (\underbrace{z_{\text{edit}}}_{\text{Perform edit}}, x_{\text{loc}}, x_{\text{in}}, y_{\text{in}}) \}$

$z_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$x_{\text{loc}} =$ “What team does Messi play for?”

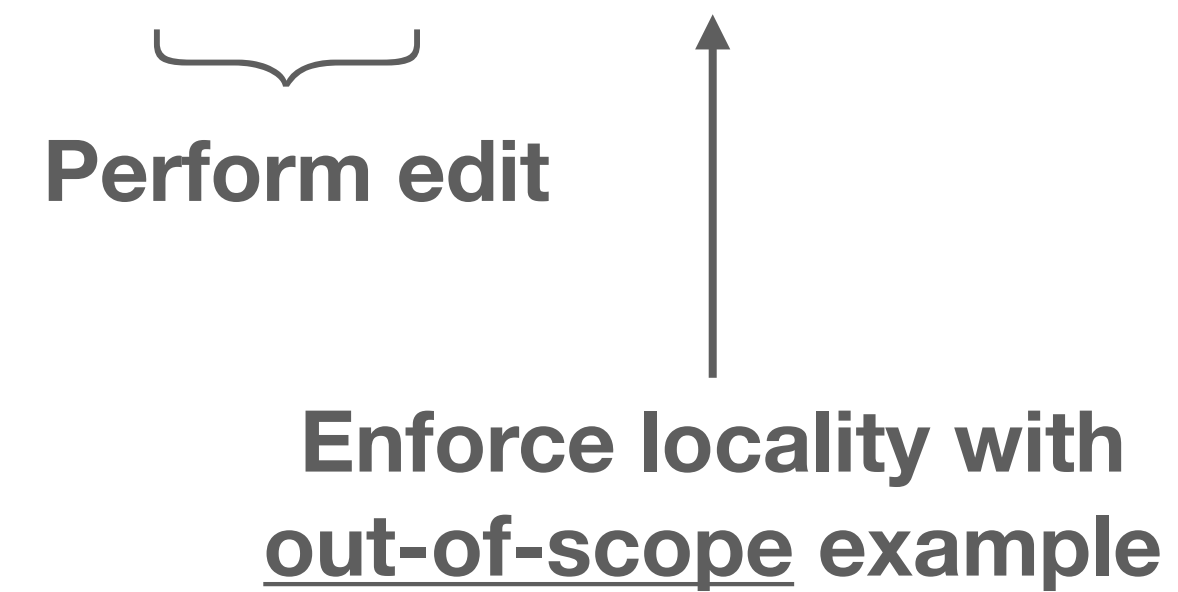
$x_{\text{in}} =$ “The prime minister of the UK is currently who?”

$y_{\text{in}} =$ “Boris Johnson”

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $D_{\text{edit}} = \{ (z_{\text{edit}}, x_{\text{loc}}, x_{\text{in}}, y_{\text{in}}) \}$



$z_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$x_{\text{loc}} =$ “What team does Messi play for?”

$x_{\text{in}} =$ “The prime minister of the UK is currently who?”

$y_{\text{in}} =$ “Boris Johnson”

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $D_{\text{edit}} = \{ (\underbrace{z_{\text{edit}}}_{\text{Perform edit}}, \underbrace{x_{\text{loc}}, x_{\text{in}}, y_{\text{in}}}_{\text{Enforce generalization with in-scope example}}) \}$

$z_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$x_{\text{loc}} =$ “What team does Messi play for?”

$x_{\text{in}} =$ “The prime minister of the UK is currently who?”

$y_{\text{in}} =$ “Boris Johnson”

Enforce locality with
out-of-scope example

Enforce generalization
with in-scope example

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $D_{\text{edit}} = \{ (\underbrace{z_{\text{edit}}}_{\text{Perform edit}}, \underbrace{x_{\text{loc}}, x_{\text{in}}, y_{\text{in}}}_{\text{Enforce generalization with in-scope example}}) \}$

$z_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$x_{\text{loc}} =$ “What team does Messi play for?”

$x_{\text{in}} =$ “The prime minister of the UK is currently who?”

$y_{\text{in}} =$ “Boris Johnson”

Enforce locality with
out-of-scope example

Inner loop

(run the editor)

$$\theta' = \text{Edit}_{\phi}(\theta, z_{\text{edit}})$$

↑
(optional) editor parameters

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $D_{\text{edit}} = \{ (\underbrace{z_{\text{edit}}}_{\text{Perform edit}}, \underbrace{x_{\text{loc}}, x_{\text{in}}, y_{\text{in}}}_{\text{Enforce generalization with in-scope example}}) \}$

$z_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$x_{\text{loc}} =$ “What team does Messi play for?”

$x_{\text{in}} =$ “The prime minister of the UK is currently who?”

$y_{\text{in}} =$ “Boris Johnson”

Enforce locality with
out-of-scope example

Inner loop

(run the editor)

$$\theta' = \text{Edit}_{\phi}(\theta, z_{\text{edit}})$$

↑
(optional) editor parameters

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $\mathbf{D}_{\text{edit}} = \{ (\mathbf{z}_{\text{edit}}, \mathbf{x}_{\text{loc}}, \mathbf{x}_{\text{in}}, \mathbf{y}_{\text{in}}) \}$

Perform edit
Enforce generalization with in-scope example

$\mathbf{z}_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$\mathbf{x}_{\text{loc}} =$ “What team does Messi play for?”

$\mathbf{x}_{\text{in}} =$ “The prime minister of the UK is currently who?”

$\mathbf{y}_{\text{in}} =$ “Boris Johnson”

Enforce locality with out-of-scope example

Inner loop

(run the editor)

$$\theta' = \text{Edit}_{\phi}(\theta, \mathbf{z}_{\text{edit}})$$

↑
(optional) editor parameters

Outer loop

(check if edit worked)

$$L_{\text{edit}} = p_{\theta'}(\mathbf{y}_{\text{in}} | \mathbf{x}_{\text{in}})$$

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $\mathbf{D}_{\text{edit}} = \{ (\mathbf{z}_{\text{edit}}, \mathbf{x}_{\text{loc}}, \mathbf{x}_{\text{in}}, \mathbf{y}_{\text{in}}) \}$

Perform edit
Enforce generalization with in-scope example

$\mathbf{z}_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$\mathbf{x}_{\text{loc}} =$ “What team does Messi play for?”

$\mathbf{x}_{\text{in}} =$ “The prime minister of the UK is currently who?”

$\mathbf{y}_{\text{in}} =$ “Boris Johnson”

Enforce locality with out-of-scope example

Inner loop

(run the editor)

$$\theta' = \text{Edit}_{\phi}(\theta, \mathbf{z}_{\text{edit}})$$

↑
(optional) editor parameters

Outer loop

(check if edit worked)

$$L_{\text{edit}} = p_{\theta'}(\mathbf{y}_{\text{in}} | \mathbf{x}_{\text{in}})$$

Did predictions **change where** we wanted them to?

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $\mathbf{D}_{\text{edit}} = \{ (\mathbf{z}_{\text{edit}}, \mathbf{x}_{\text{loc}}, \mathbf{x}_{\text{in}}, \mathbf{y}_{\text{in}}) \}$

Perform edit
Enforce generalization with in-scope example

$\mathbf{z}_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$\mathbf{x}_{\text{loc}} =$ “What team does Messi play for?”

$\mathbf{x}_{\text{in}} =$ “The prime minister of the UK is currently who?”

$\mathbf{y}_{\text{in}} =$ “Boris Johnson”

Enforce locality with out-of-scope example

Inner loop

(run the editor)

$$\theta' = \text{Edit}_{\phi}(\theta, \mathbf{z}_{\text{edit}})$$

↑
(optional) editor parameters

Outer loop

(check if edit worked)

$$L_{\text{edit}} = p_{\theta'}(\mathbf{y}_{\text{in}} | \mathbf{x}_{\text{in}})$$

$$L_{\text{loc}} = \text{KL} (p_{\theta}(\cdot | \mathbf{x}_{\text{loc}}) \| p_{\theta'}(\cdot | \mathbf{x}_{\text{loc}}))$$

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $\mathbf{D}_{\text{edit}} = \{ (\mathbf{z}_{\text{edit}}, \mathbf{x}_{\text{loc}}, \mathbf{x}_{\text{in}}, \mathbf{y}_{\text{in}}) \}$

Perform edit
Enforce generalization with in-scope example

$\mathbf{z}_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$\mathbf{x}_{\text{loc}} =$ “What team does Messi play for?”

$\mathbf{x}_{\text{in}} =$ “The prime minister of the UK is currently who?”

$\mathbf{y}_{\text{in}} =$ “Boris Johnson”

Enforce locality with out-of-scope example

Inner loop

(run the editor)

$$\theta' = \text{Edit}_{\phi}(\theta, \mathbf{z}_{\text{edit}})$$

↑
(optional) editor parameters

Outer loop

(check if edit worked)

$$L_{\text{edit}} = p_{\theta'}(\mathbf{y}_{\text{in}} | \mathbf{x}_{\text{in}})$$

Did we keep predictions **the same** everywhere else?

$$L_{\text{loc}} = \text{KL} \left(p_{\theta}(\cdot | \mathbf{x}_{\text{loc}}) \parallel p_{\theta'}(\cdot | \mathbf{x}_{\text{loc}}) \right)$$

Learning to edit

Editing as meta-learning

Requirement: an “edit dataset” $\mathbf{D}_{\text{edit}} = \{ (\mathbf{z}_{\text{edit}}, \mathbf{x}_{\text{loc}}, \mathbf{x}_{\text{in}}, \mathbf{y}_{\text{in}}) \}$

Perform edit
Enforce generalization with in-scope example

$\mathbf{z}_{\text{edit}} =$ “Who is the UK PM? Boris Johnson”

$\mathbf{x}_{\text{loc}} =$ “What team does Messi play for?”

$\mathbf{x}_{\text{in}} =$ “The prime minister of the UK is currently who?”

$\mathbf{y}_{\text{in}} =$ “Boris Johnson”

Enforce locality with out-of-scope example

Inner loop

(run the editor)

$$\theta' = \text{Edit}_{\phi}(\theta, \mathbf{z}_{\text{edit}})$$

↑
(optional) editor parameters

Outer loop

(check if edit worked)

$$L_{\text{edit}} = p_{\theta'}(\mathbf{y}_{\text{in}} | \mathbf{x}_{\text{in}})$$

$$L_{\text{loc}} = \text{KL} \left(p_{\theta}(\cdot | \mathbf{x}_{\text{loc}}) \parallel p_{\theta'}(\cdot | \mathbf{x}_{\text{loc}}) \right)$$

Backprop $L_{\text{edit}} + L_{\text{loc}}$ back into base model/editor

Learning to edit

A tale of two meta-learning frameworks

1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit

Learning to edit

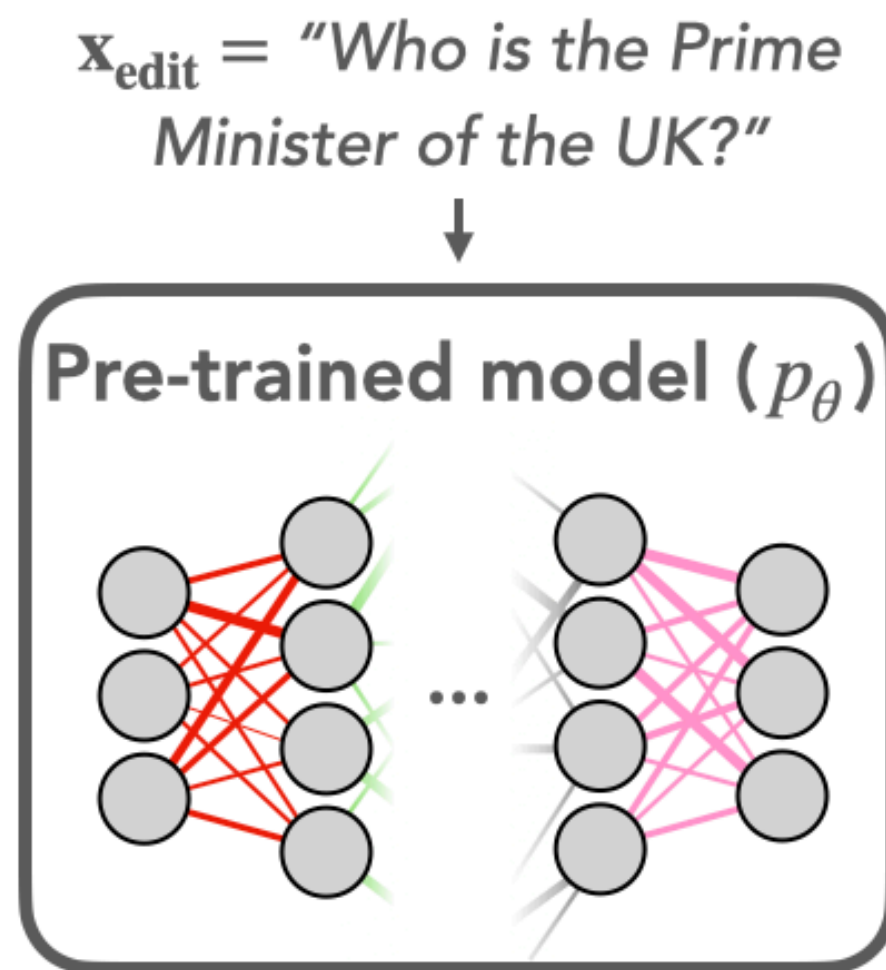
A tale of two meta-learning frameworks

1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} P_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit
2. Hypernetwork-based: Freeze base model, train a **gradient transform** $g_{\phi}(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_{\theta} = g_{\phi}(\nabla_{\theta} P_{\theta}(\mathbf{z}_{\text{edit}}))$ gives a good edit

Learning to edit

A tale of two meta-learning frameworks

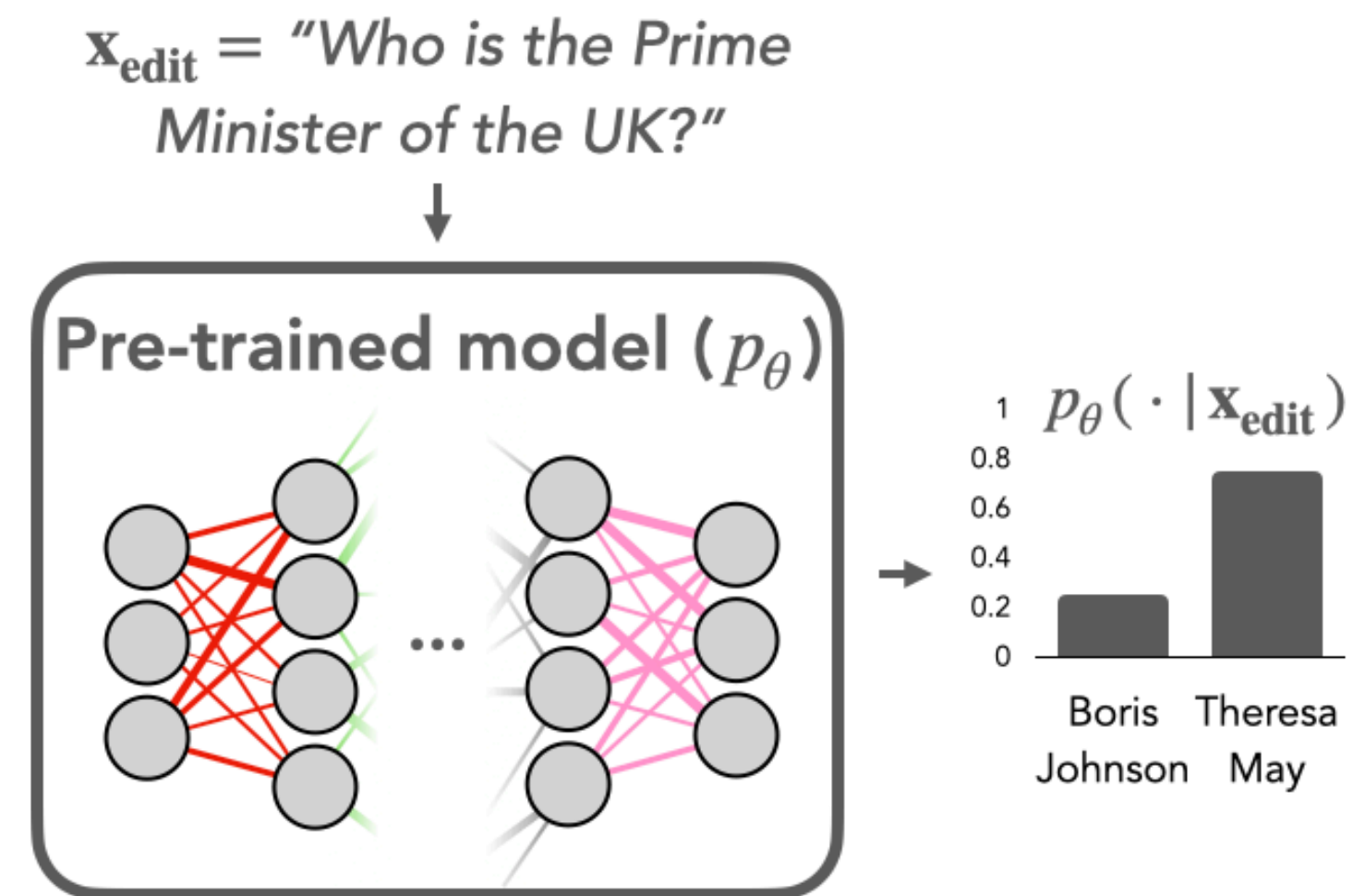
1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

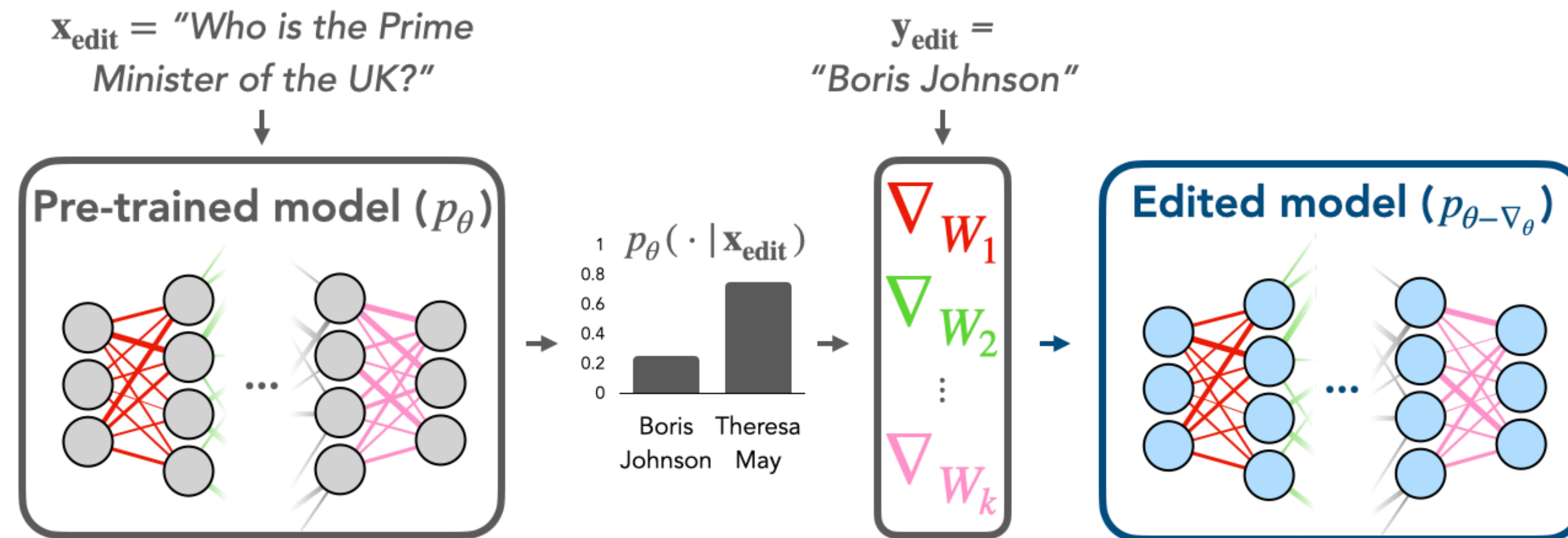
1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

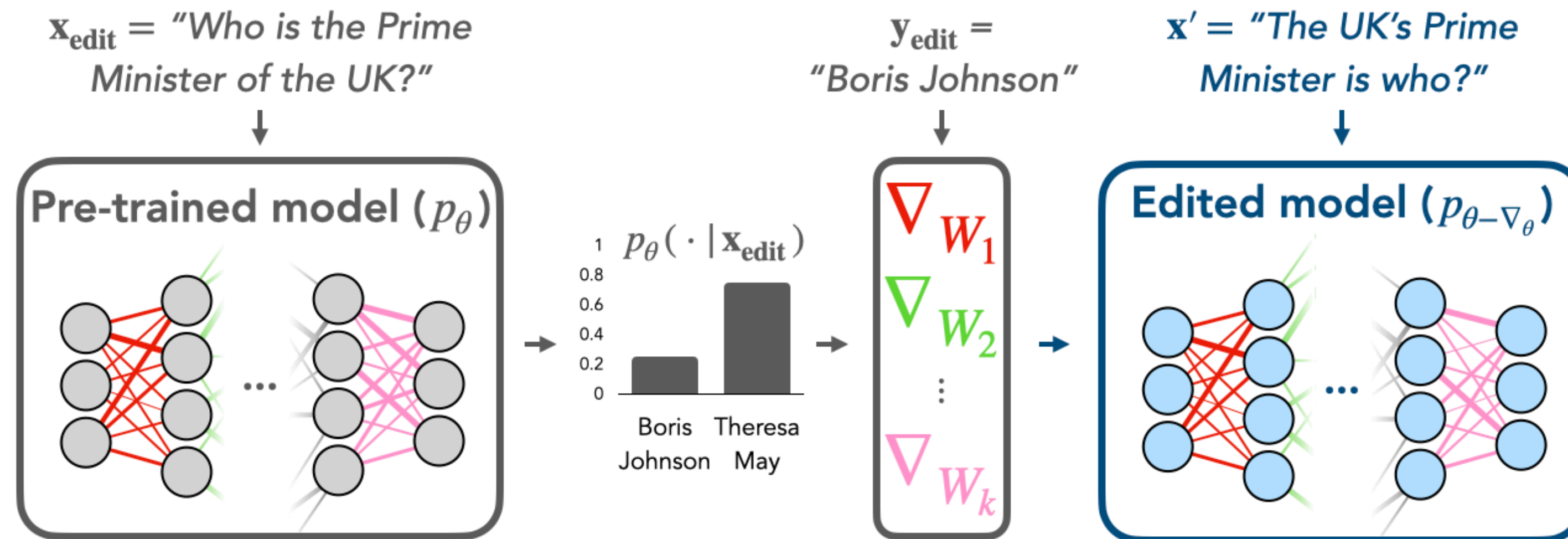
1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

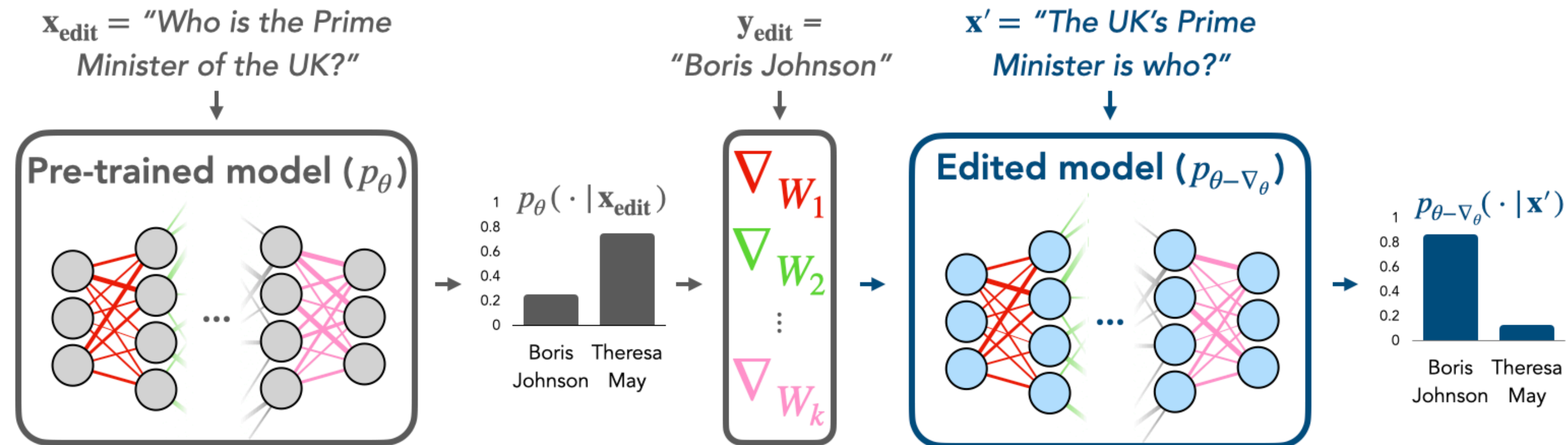
1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

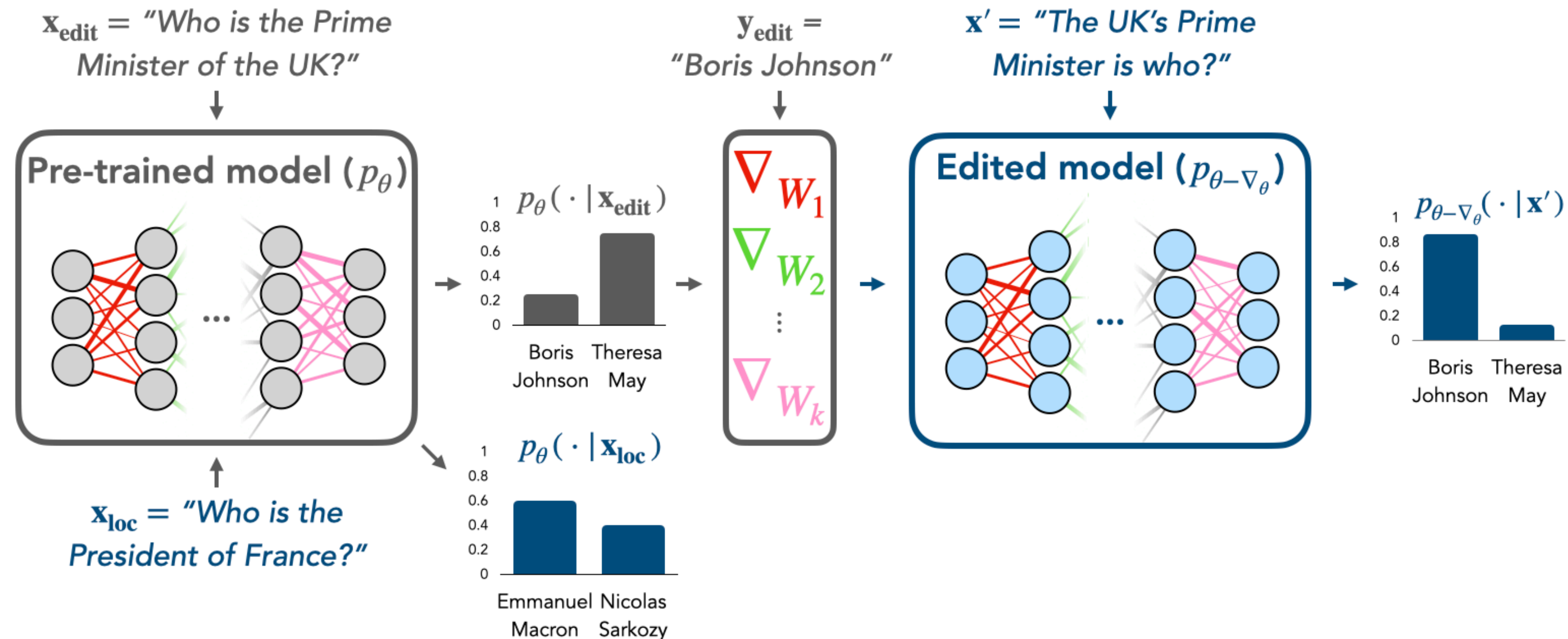
1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

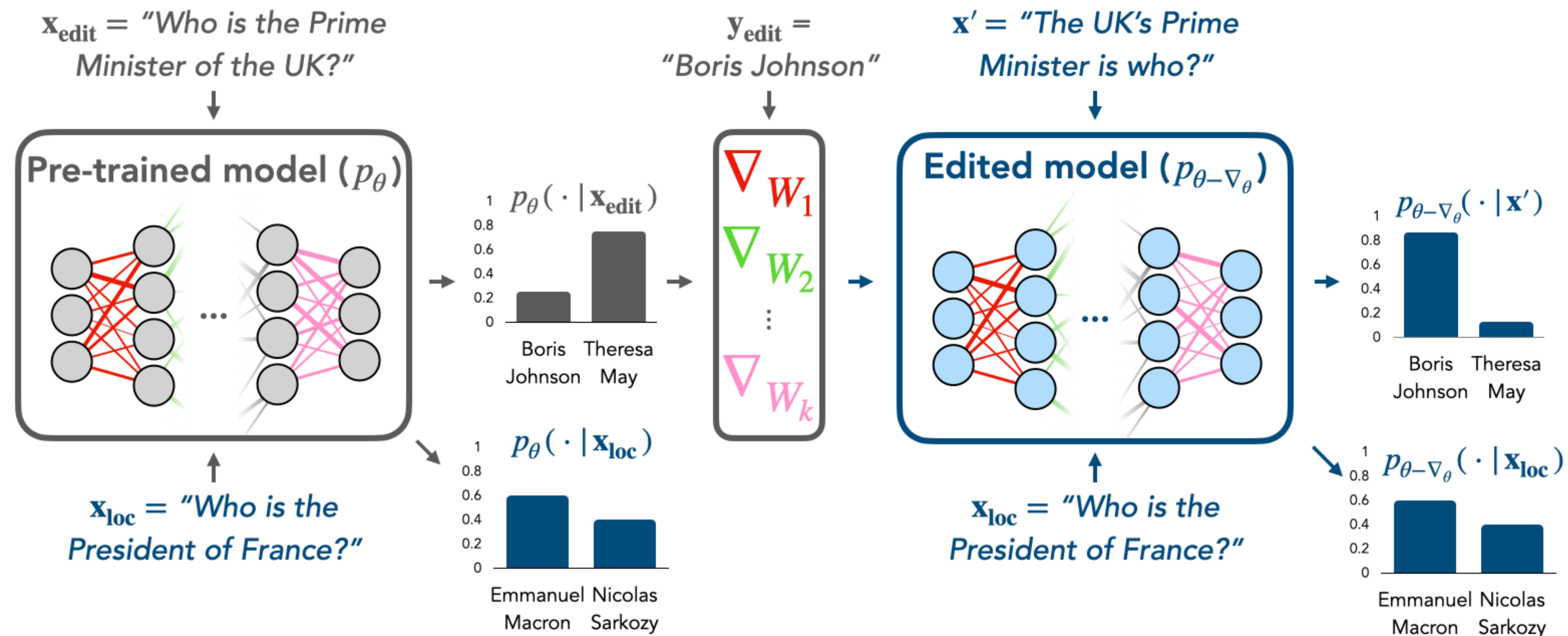
1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

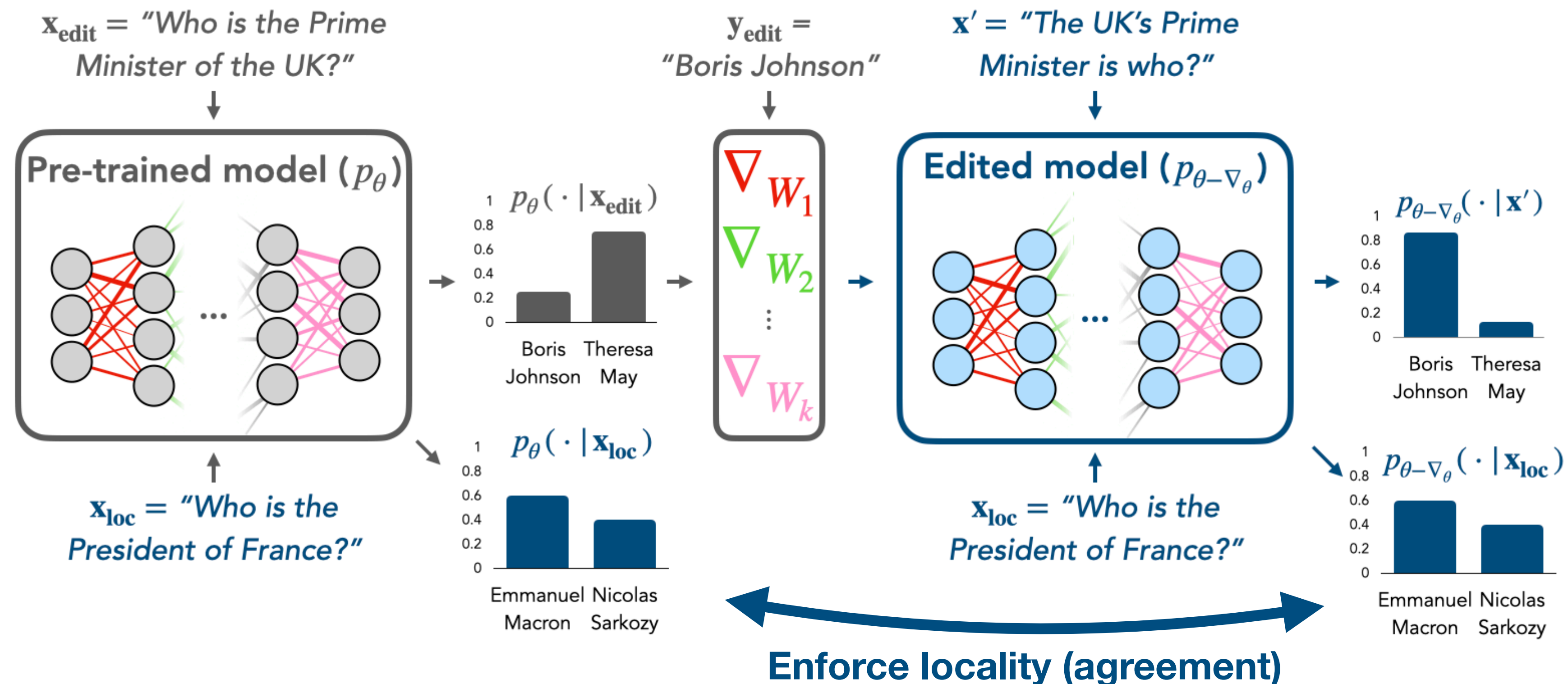
1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

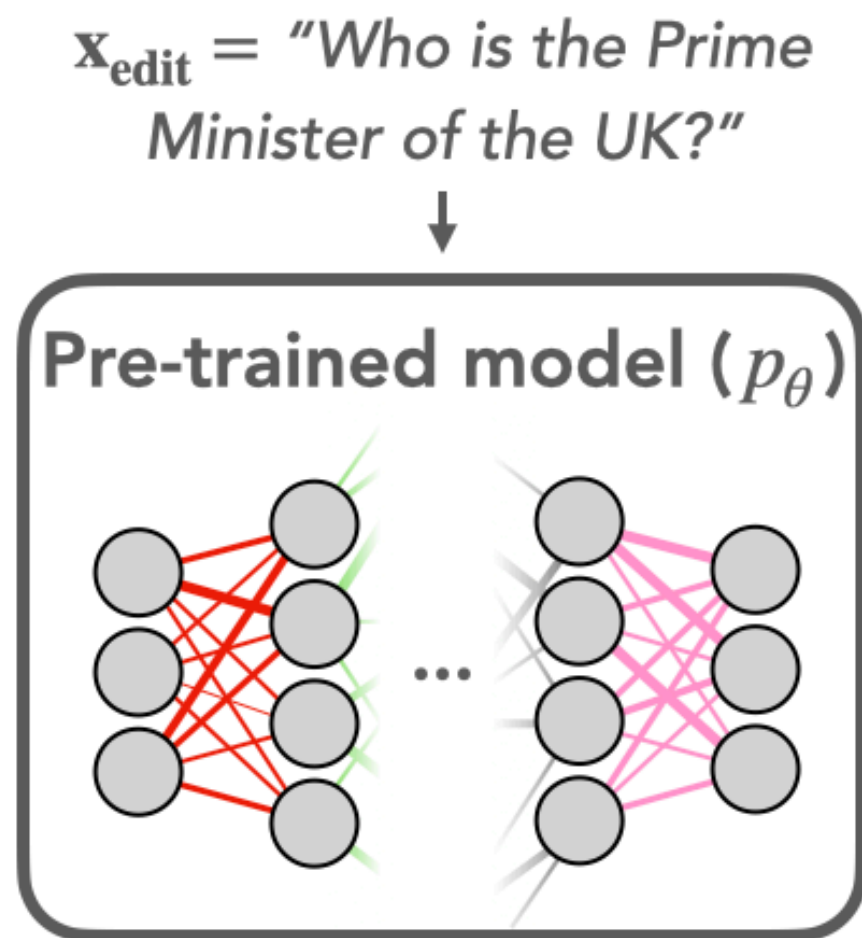
1. MAML-based: Train your **base model** s.t. the regular fine-tuning gradient $\nabla_{\theta} p_{\theta}(\mathbf{z}_{\text{edit}})$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

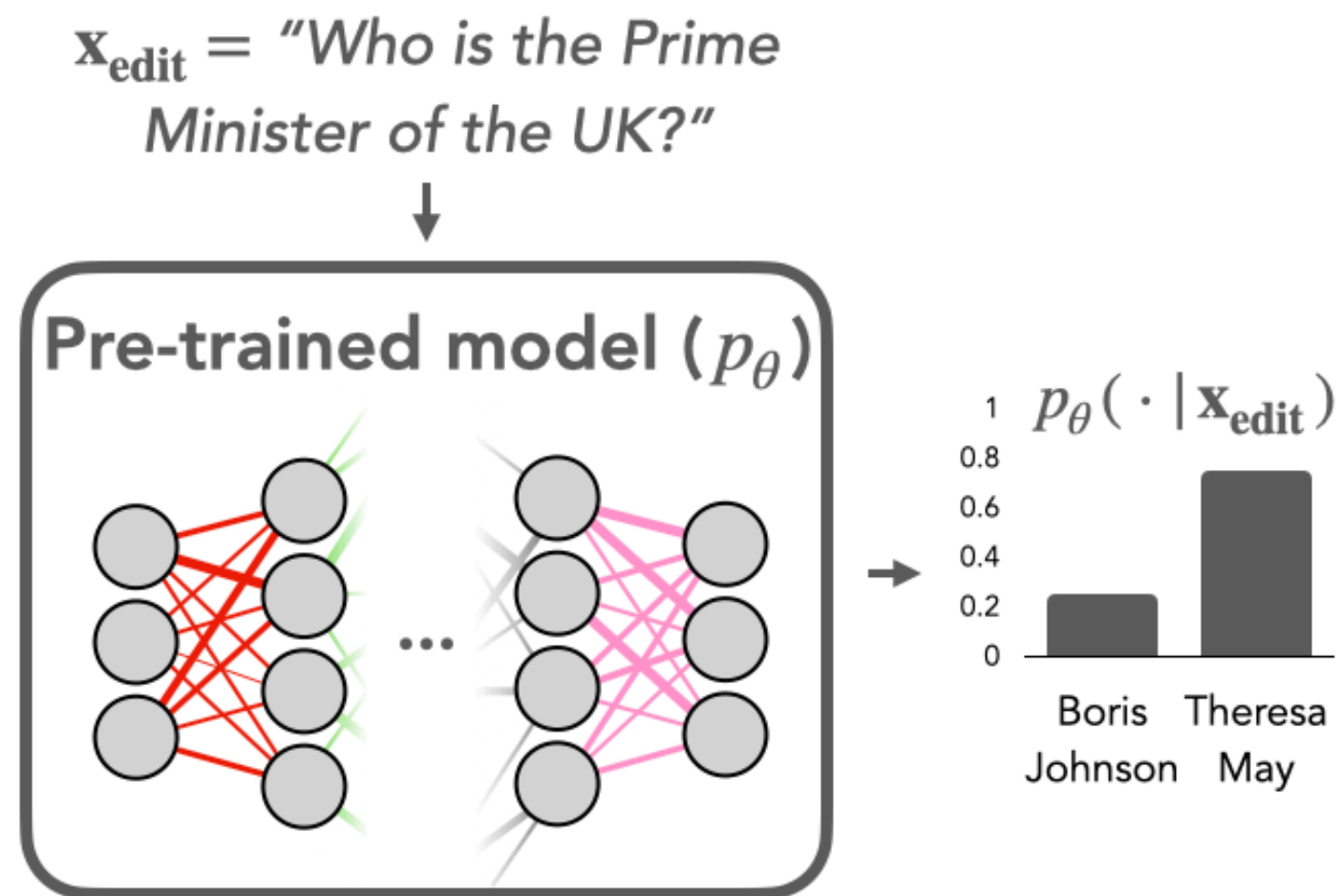
2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta p_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

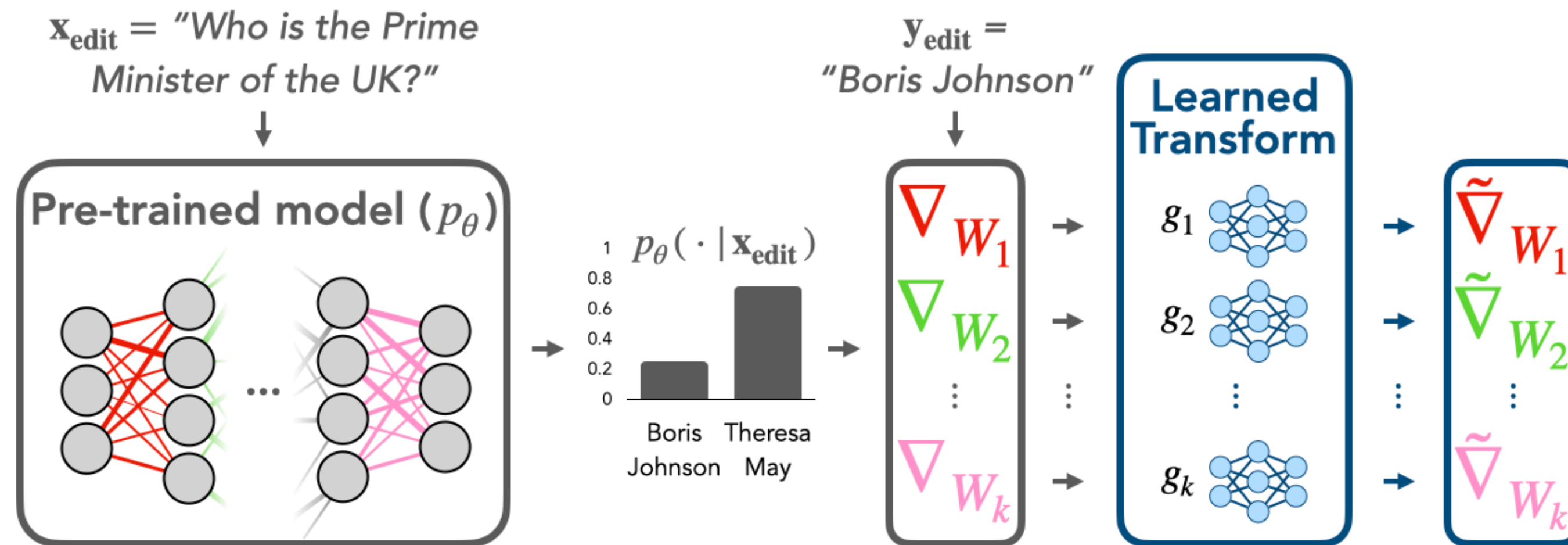
2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta P_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

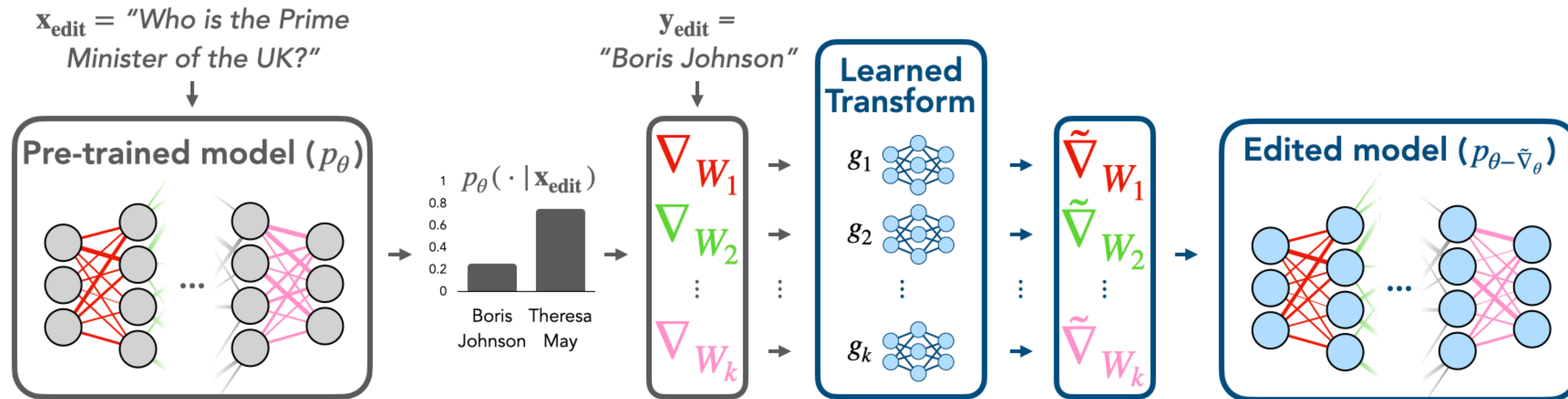
2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta P_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

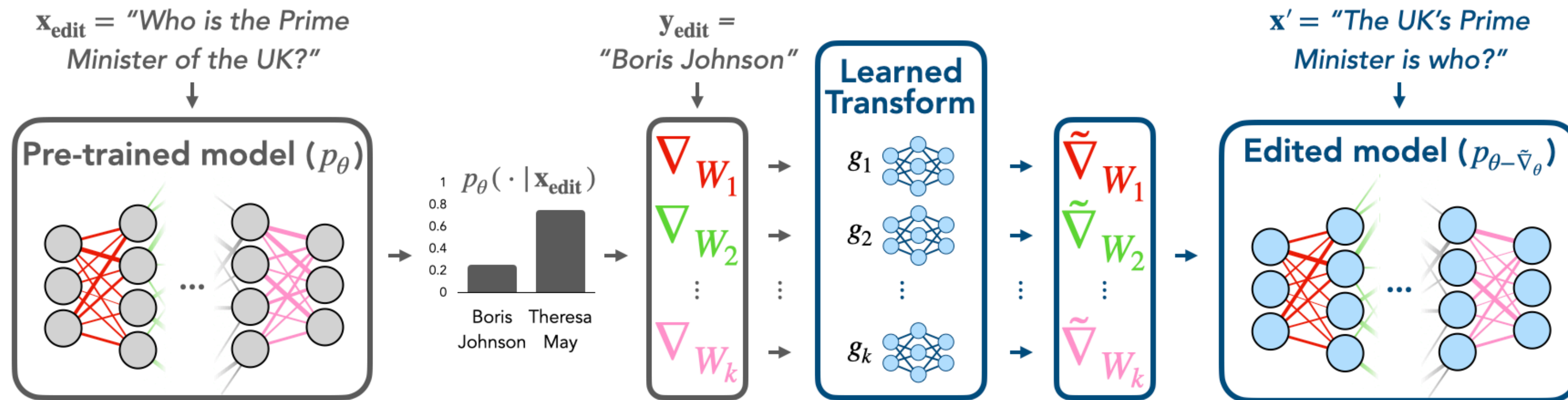
2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta P_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

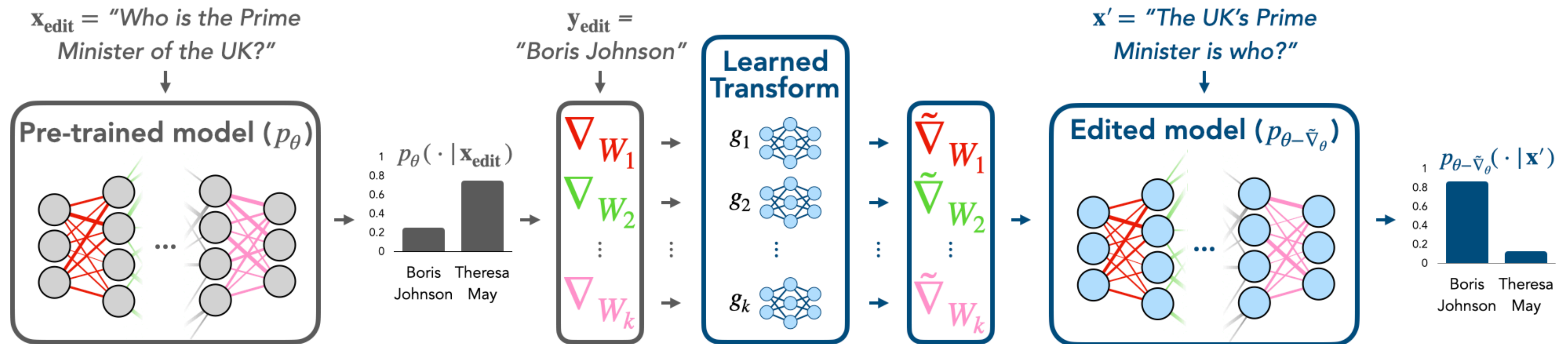
2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta p_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

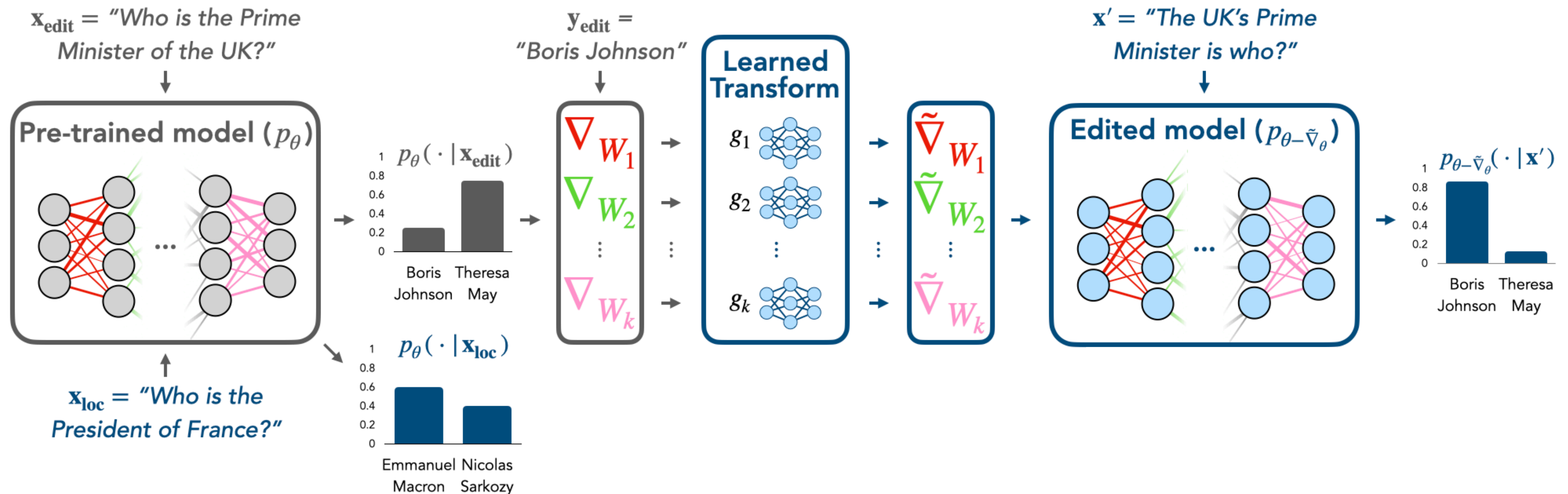
2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta p_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

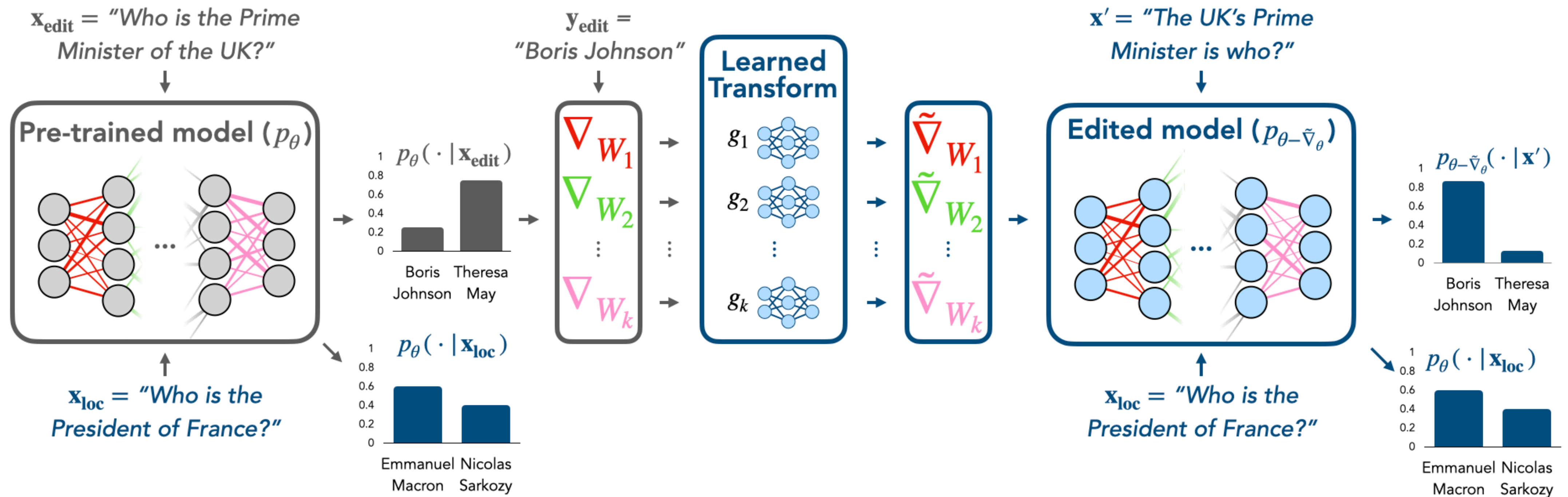
2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta p_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Learning to edit

A tale of two meta-learning frameworks

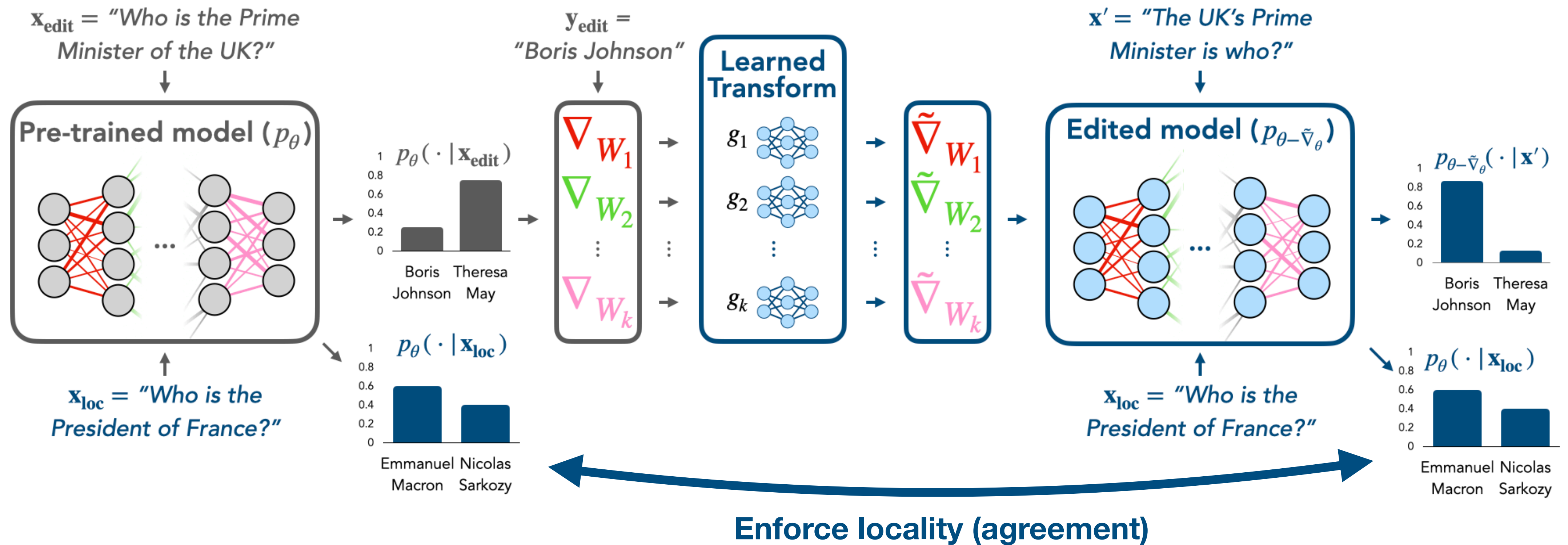
2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta p_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Learning to edit

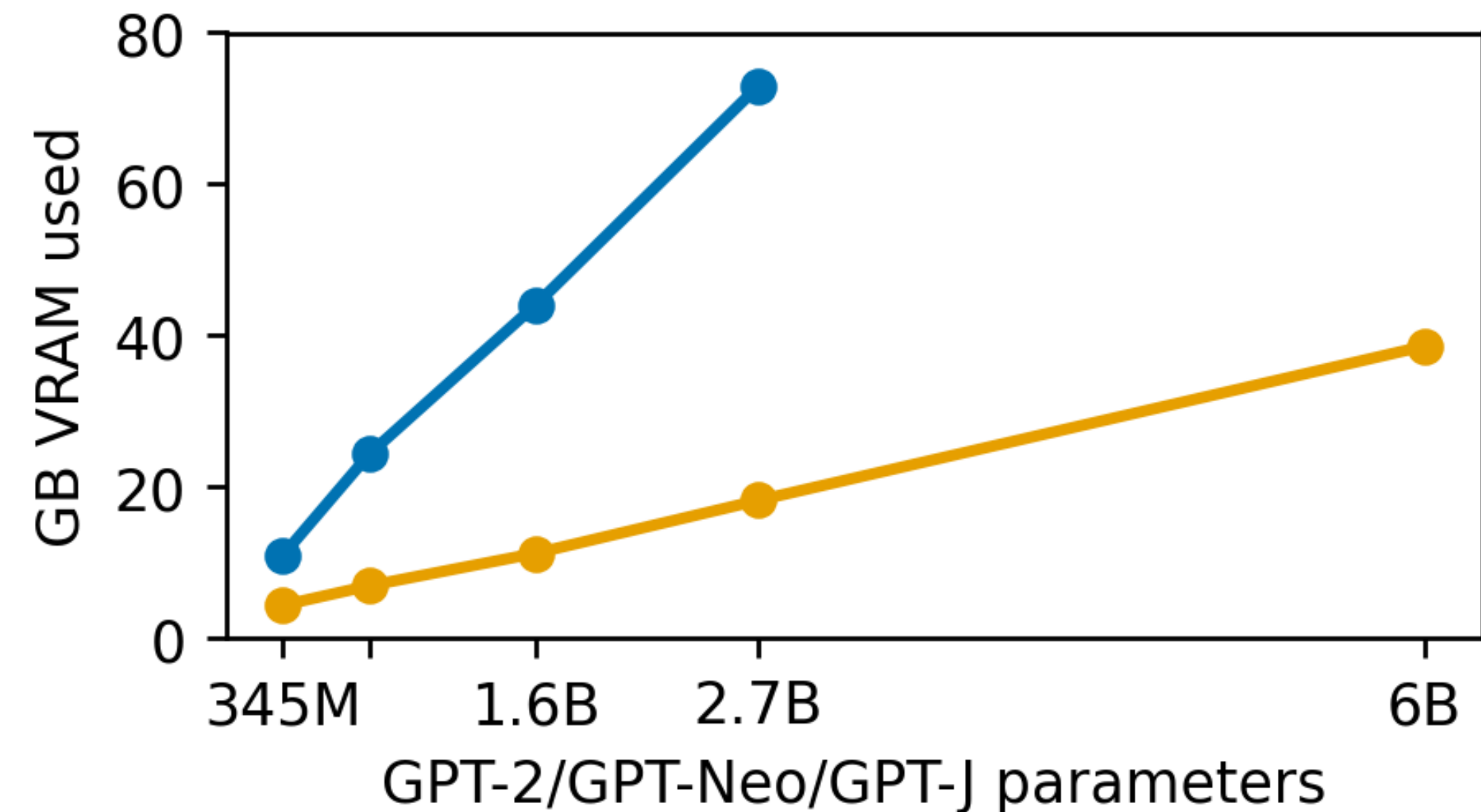
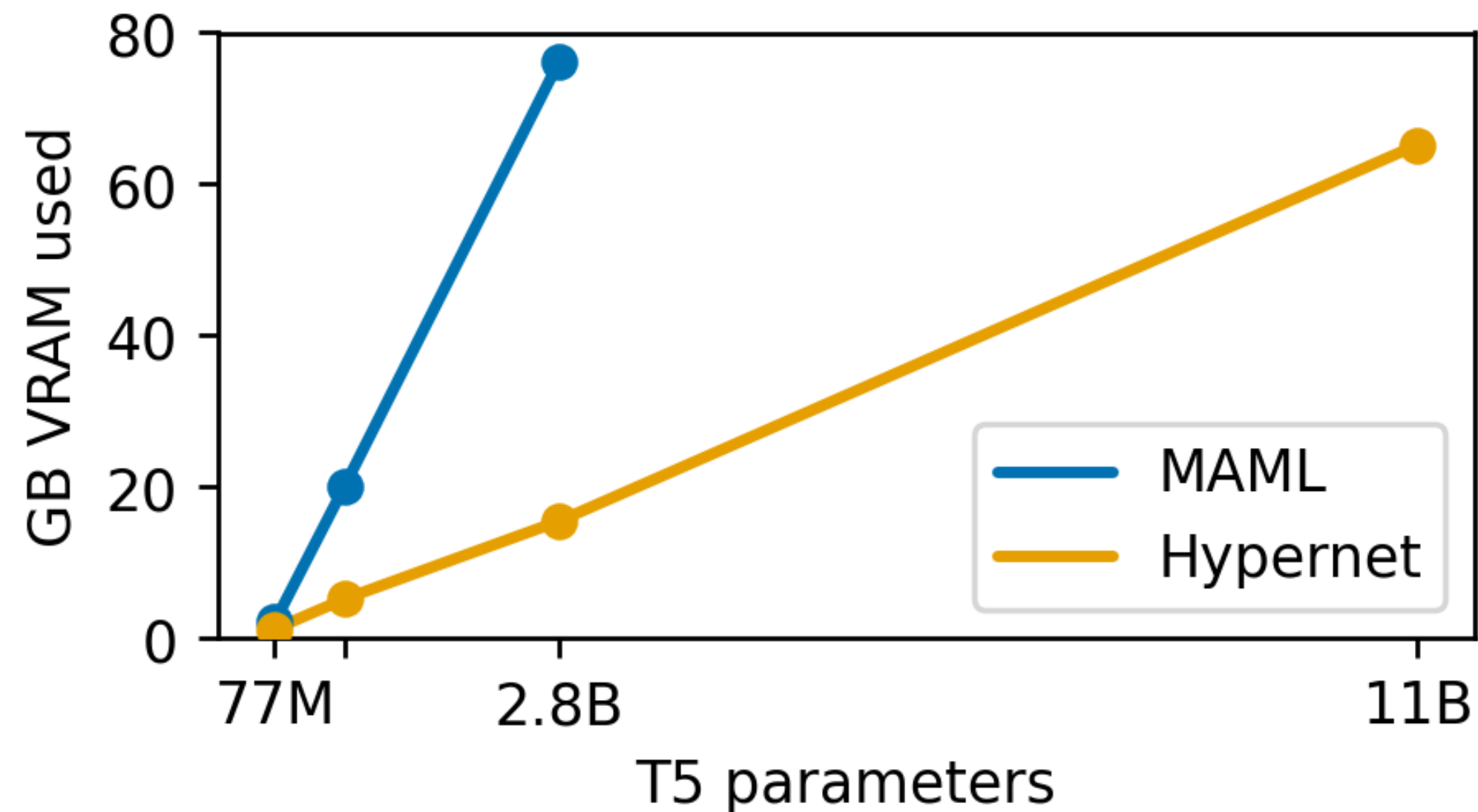
A tale of two meta-learning frameworks

2. Hypernetwork-based: **Freeze** base model, train a **gradient transform** $g_\phi(\cdot)$ s.t. transformed fine-tuning gradient $\tilde{\nabla}_\theta = g_\phi(\nabla_\theta p_\theta(\mathbf{z}_{\text{edit}}))$ gives a good edit



Challenges of editing at scale

- MAML-based editors show **high memory consumption**



Challenges of editing at scale

- MAML-based editors show **high memory consumption**
- Hypernets use **restrictive approximations** to handle high-dim parameters

Challenges of editing at scale

- MAML-based editors show **high memory consumption**
- Hypernets use **restrictive approximations** to handle high-dim parameters

95% ES at 0.1B parameters

Challenges of editing at scale

- MAML-based editors show **high memory consumption**
- Hypernets use **restrictive approximations** to handle high-dim parameters

95% ES at 0.1B parameters

(BART-base)

Challenges of editing at scale

- MAML-based editors show **high memory consumption**
- Hypernets use **restrictive approximations** to handle high-dim parameters

95% ES at 0.1B parameters → 4% ES at 11B parameters

(BART-base)

Challenges of editing at scale

- MAML-based editors show **high memory consumption**
- Hypernets use **restrictive approximations** to handle high-dim parameters

95% ES at 0.1B parameters → **4% ES at 11B parameters**

(BART-base)

(T5-XXL)

Challenges of editing at scale

- MAML-based editors show **high memory consumption**
- Hypernets use **restrictive approximations** to handle high-dim parameters

95% ES at 0.1B parameters → **4% ES at 11B parameters**

(BART-base)

(T5-XXL)

Can we develop an **efficient** and **expressive** gradient transform?

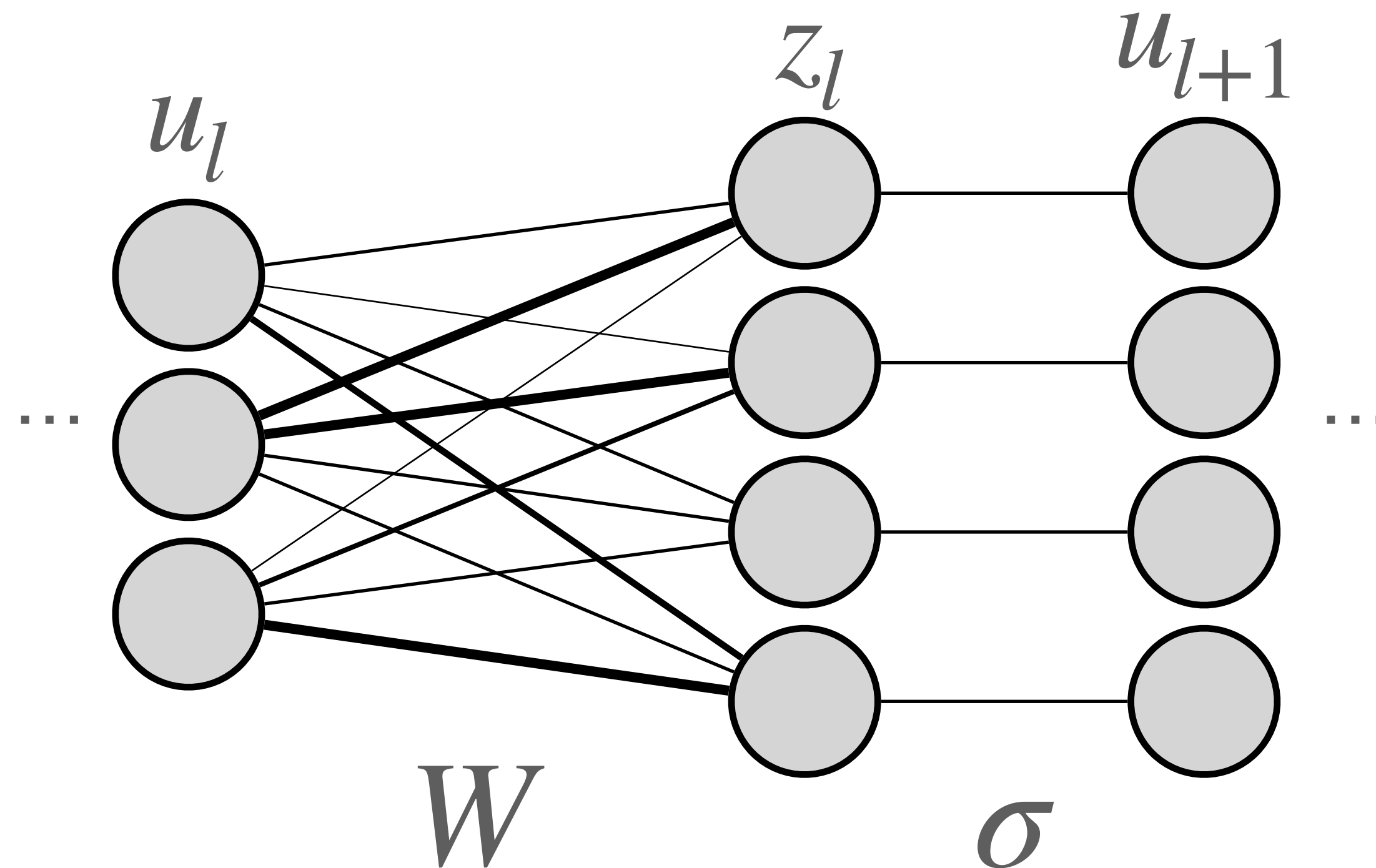
Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

Obs. 1: for $W \in \mathbb{R}^{d \times d}$, $\nabla_W L$ is rank- B for batch size B !

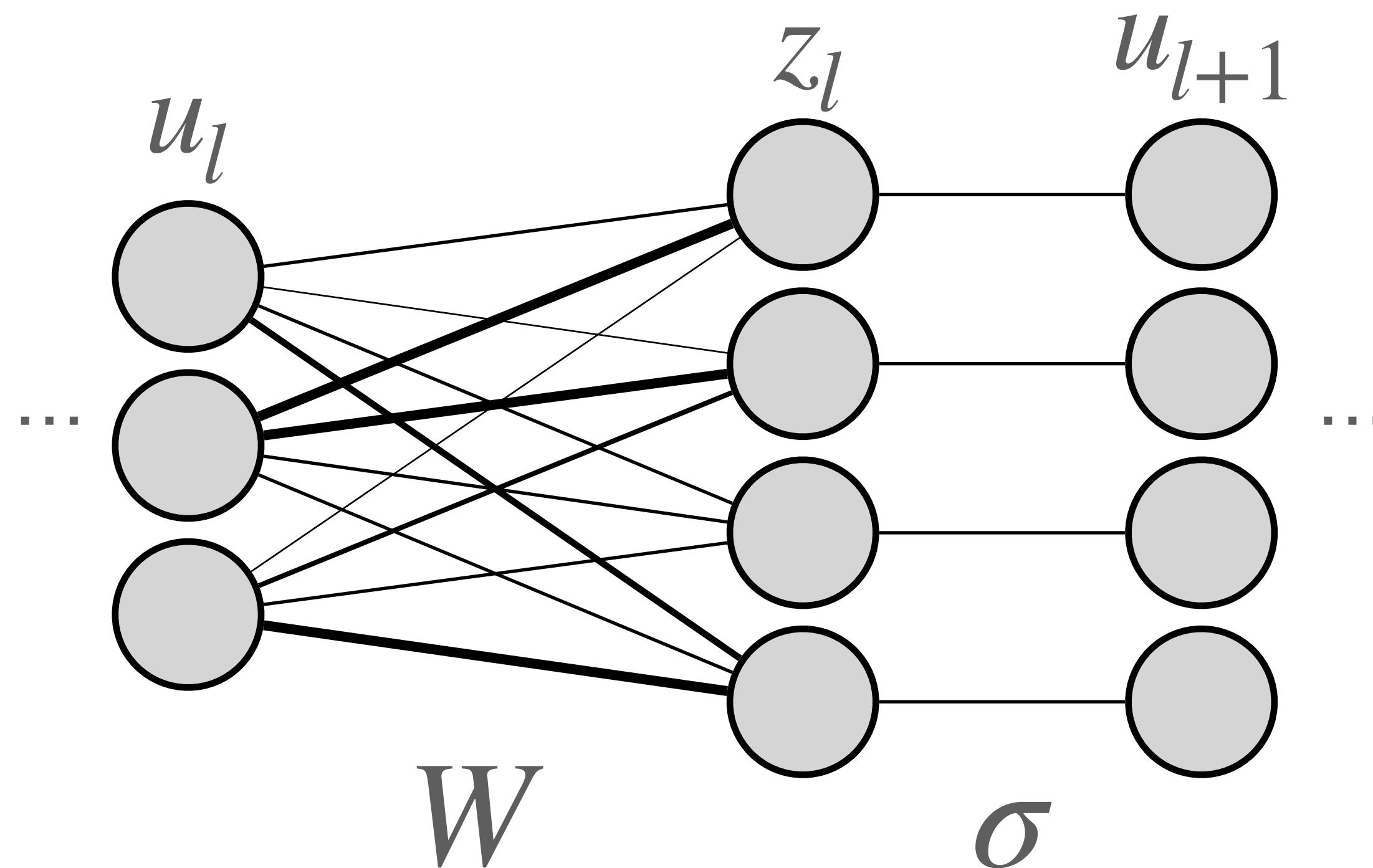


Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

Obs. 1: for $W \in \mathbb{R}^{d \times d}$, $\nabla_W L$ is rank- B for batch size B !

Forward pass
computes u_l



Model Editor Networks using gradient Decomposition

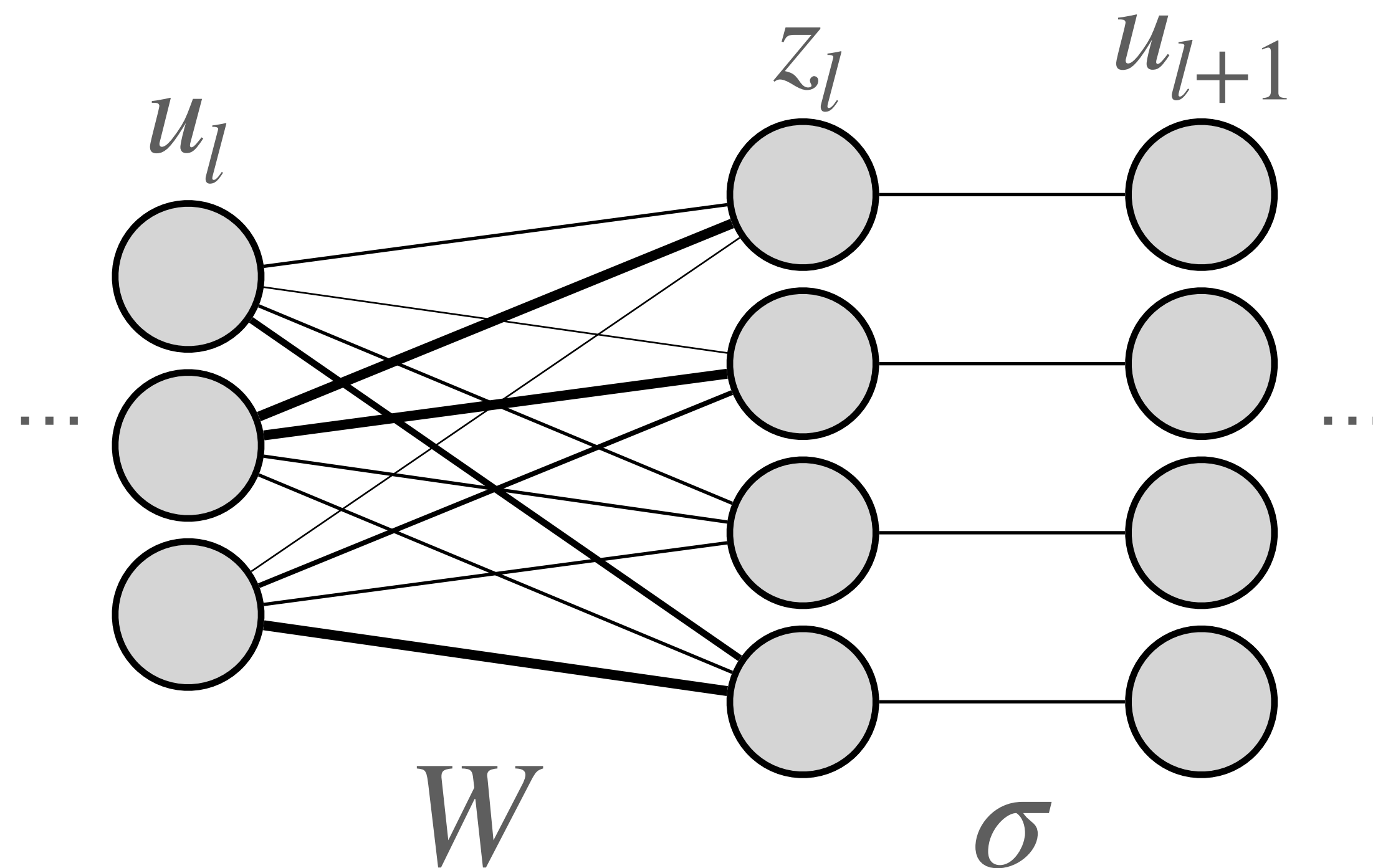
An efficient, expressive gradient transform

Obs. 1: for $W \in \mathbb{R}^{d \times d}$, $\nabla_W L$ is rank- B for batch size B !

Forward pass
computes u_l

Backward pass

computes $\delta_l = \frac{\partial L}{\partial z_l}$



Model Editor Networks using gradient Decomposition

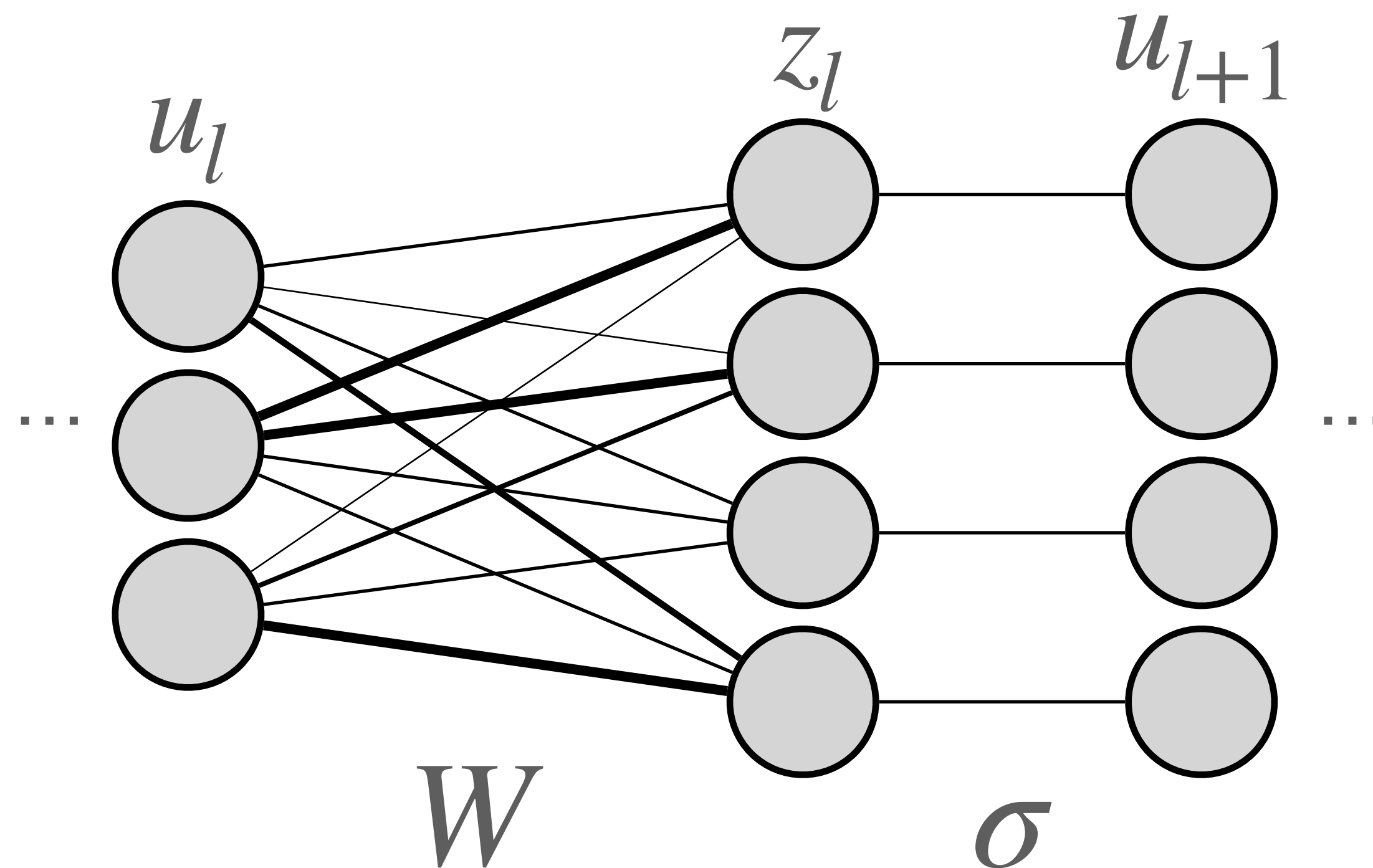
An efficient, expressive gradient transform

Obs. 1: for $W \in \mathbb{R}^{d \times d}$, $\nabla_W L$ is rank- B for batch size B !

Forward pass
computes u_l

Backward pass
computes $\delta_l = \frac{\partial L}{\partial z_l}$

$$\nabla_W L = \delta_l u_l^\top$$



Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

Obs. 1: for $W \in \mathbb{R}^{d \times d}$, $\nabla_W L$ is rank- B for batch size B !

Obs. 2: fine-tuning models with low-rank updates works really well [1]

$$\text{e.g. } W_{ft} = W_0 + AB^T, \quad \text{rank}(A) = \text{rank}(B) \ll d$$

Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

Obs. 1: for $W \in \mathbb{R}^{d \times d}$, $\nabla_W L$ is rank- B for batch size B !

Obs. 2: fine-tuning models with low-rank updates works really well [1]

Conclusion: use rank-1 input **and** output; $d^2 \rightarrow d^2$ becomes $2d \rightarrow 2d$

Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

Obs. 1: for $W \in \mathbb{R}^{d \times d}$, $\nabla_W L$ is rank- B for batch size B !

Obs. 2: fine-tuning models with low-rank updates works really well [1]

Conclusion: use rank-1 input **and** output; $d^2 \rightarrow d^2$ becomes $2d \rightarrow 2d$

Idea: map each rank-1 gradient component to new rank-1 “pseudograd” & sum

Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

$$\nabla_{W_e}$$

Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

$$\nabla_{W_\ell} = \delta_\ell u_\ell^\top$$

Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

$$\nabla_{W_\ell} = \delta_\ell u_\ell^\top$$

↑
Layer input

Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

Gradient of loss
wrt layer **output**



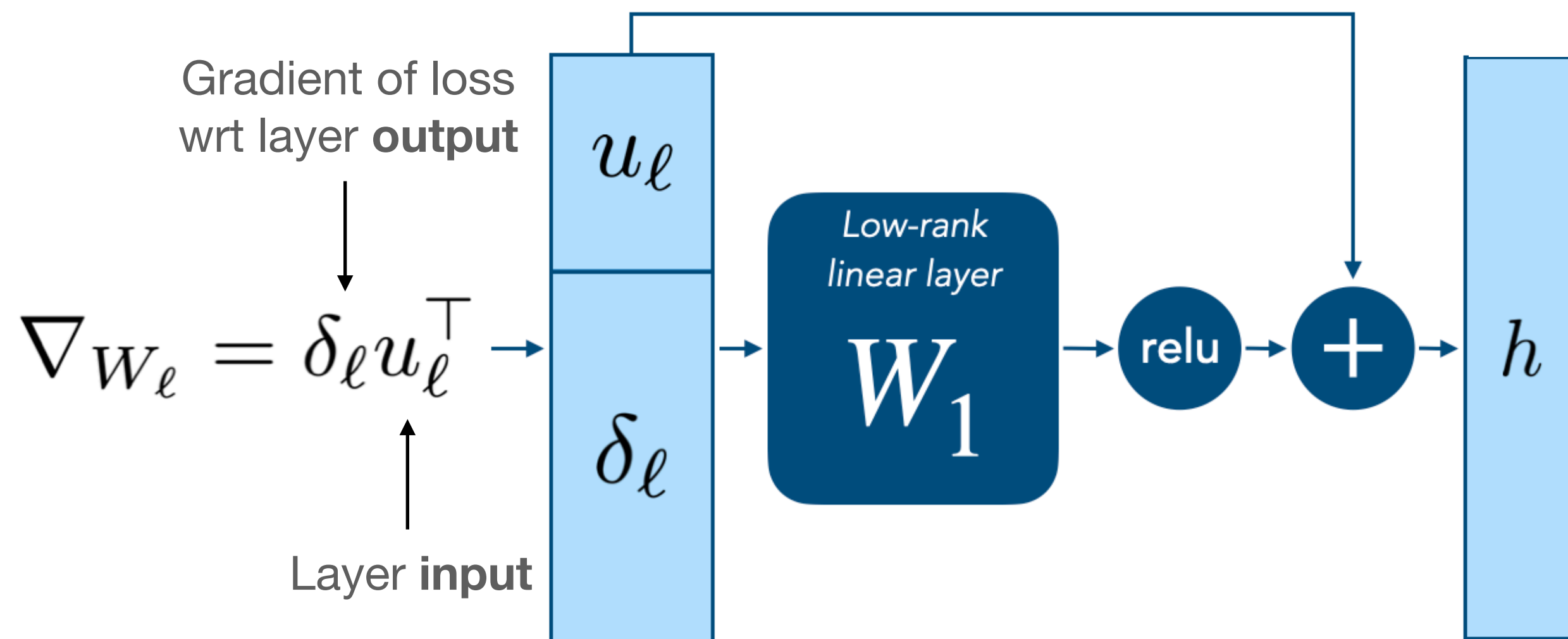
$$\nabla_{W_\ell} = \delta_\ell u_\ell^\top$$



Layer **input**

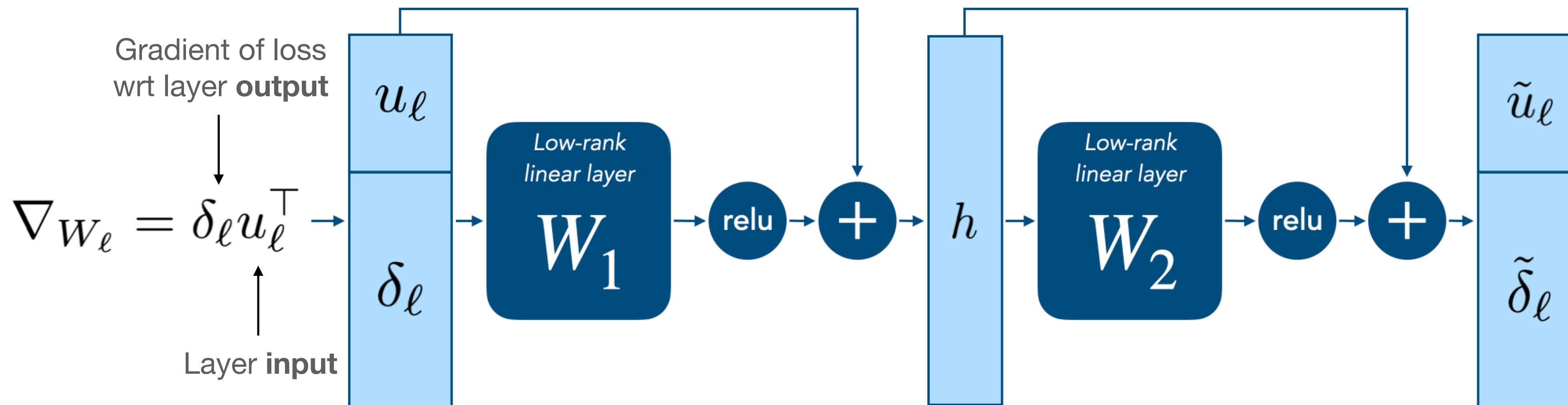
Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform



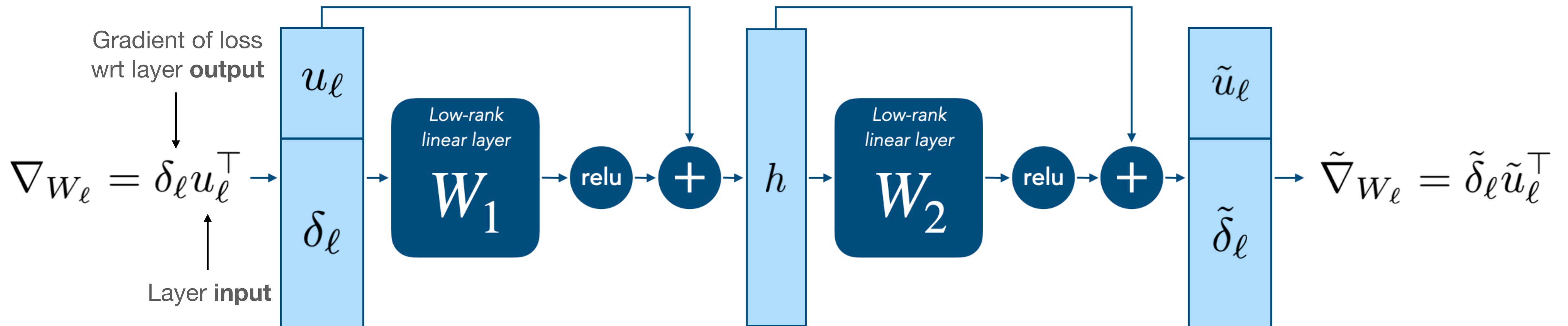
Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform



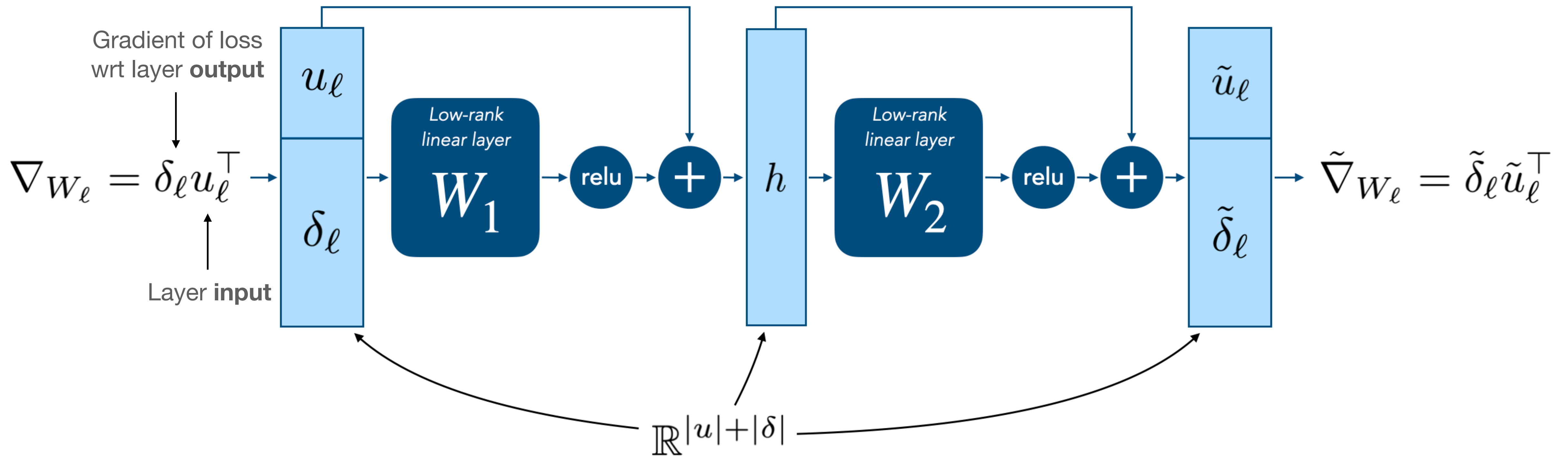
Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform



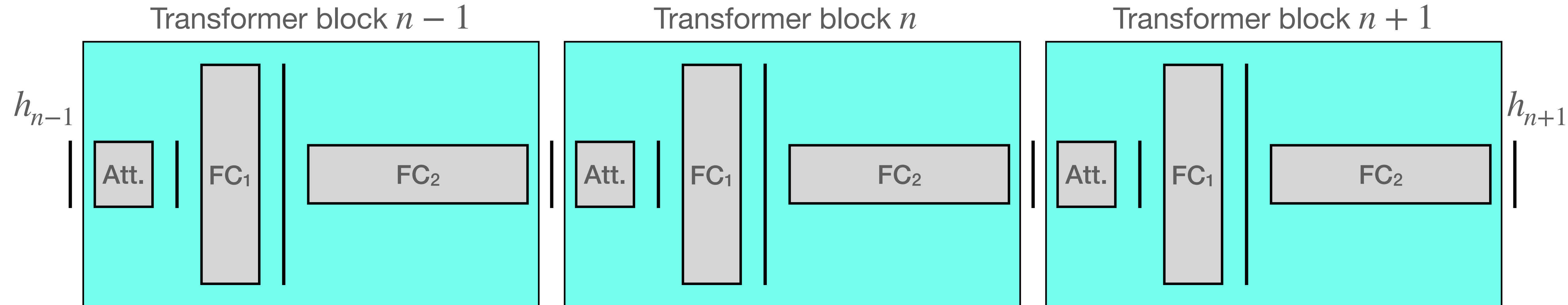
Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform



Model Editor Networks using gradient Decomposition

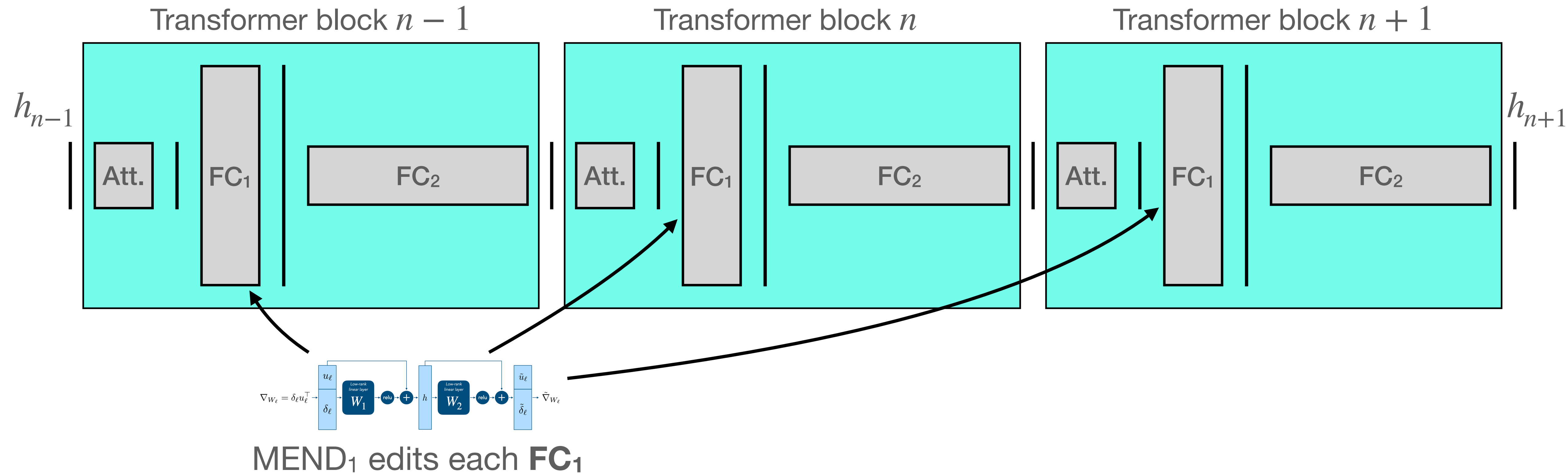
An efficient, expressive gradient transform



Surprisingly, we can **edit many layers with the same editor network**

Model Editor Networks using gradient Decomposition

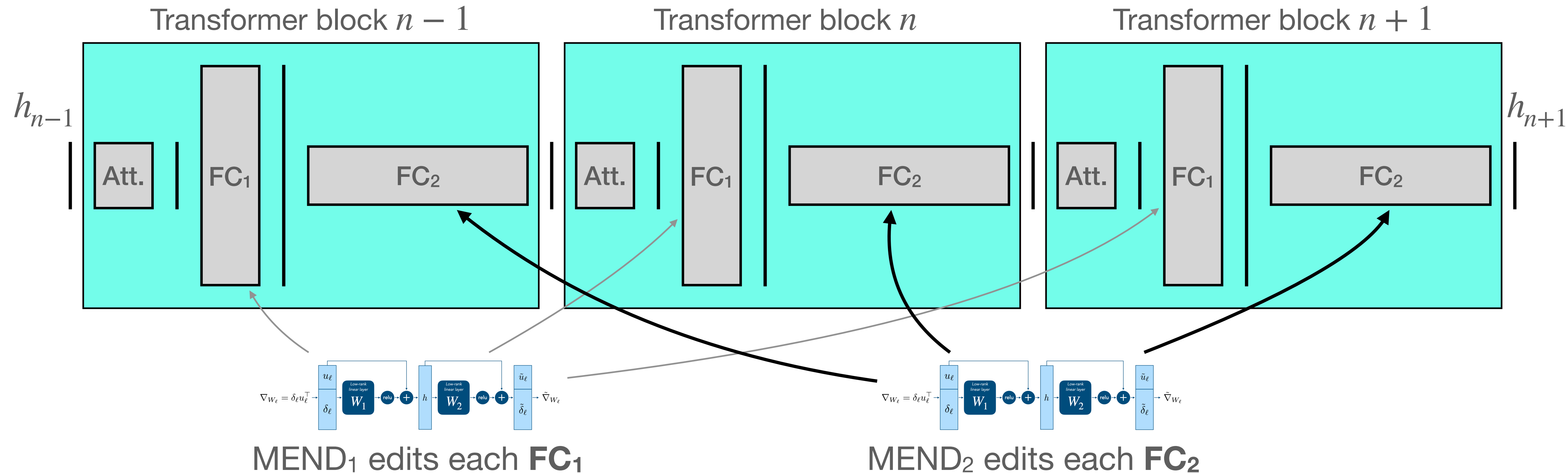
An efficient, expressive gradient transform



Surprisingly, we can **edit many layers with the same editor network**

Model Editor Networks using gradient Decomposition

An efficient, expressive gradient transform

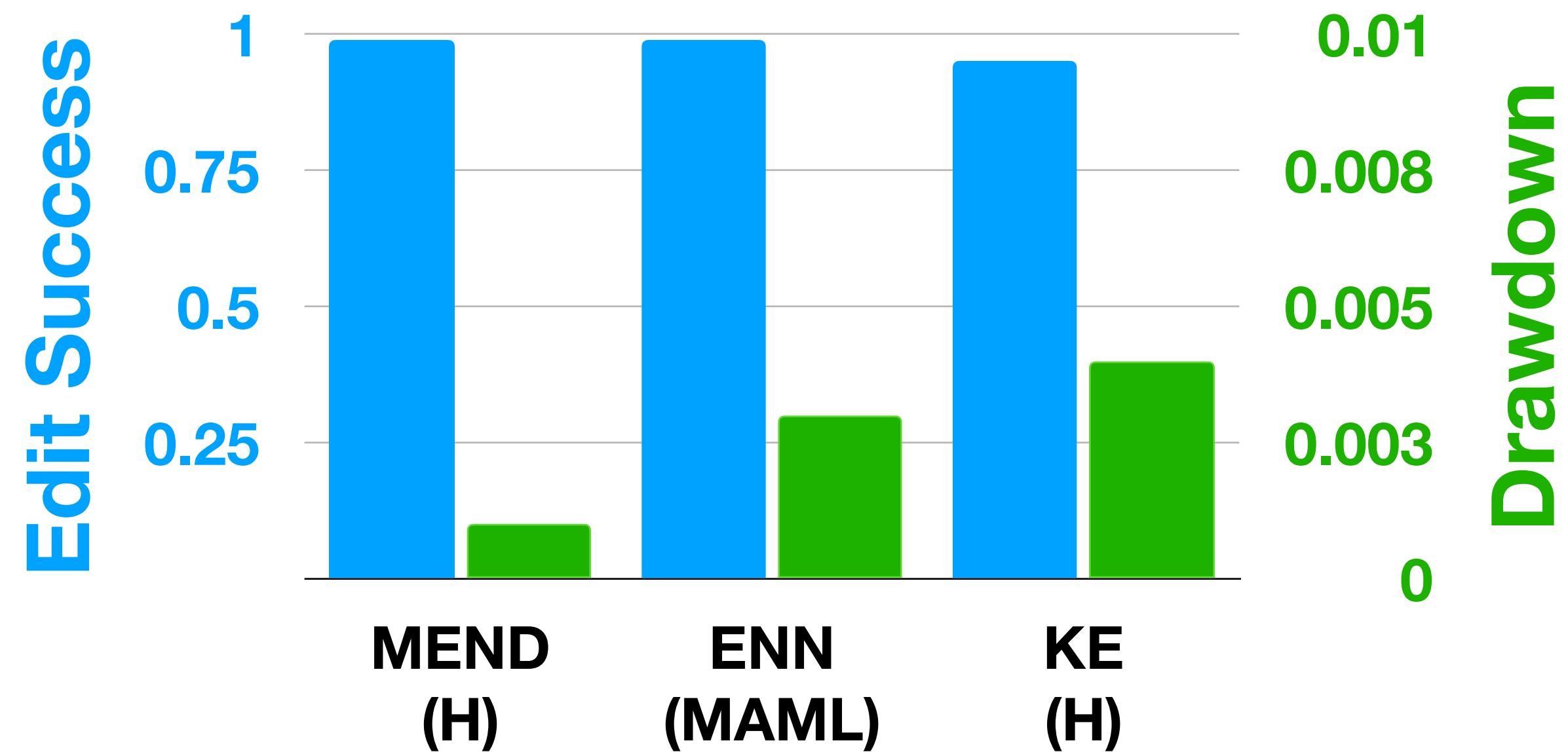


Surprisingly, we can **edit many layers with the same editor network**

Model Editor Networks using gradient Decomposition

Effective editing at small scale...

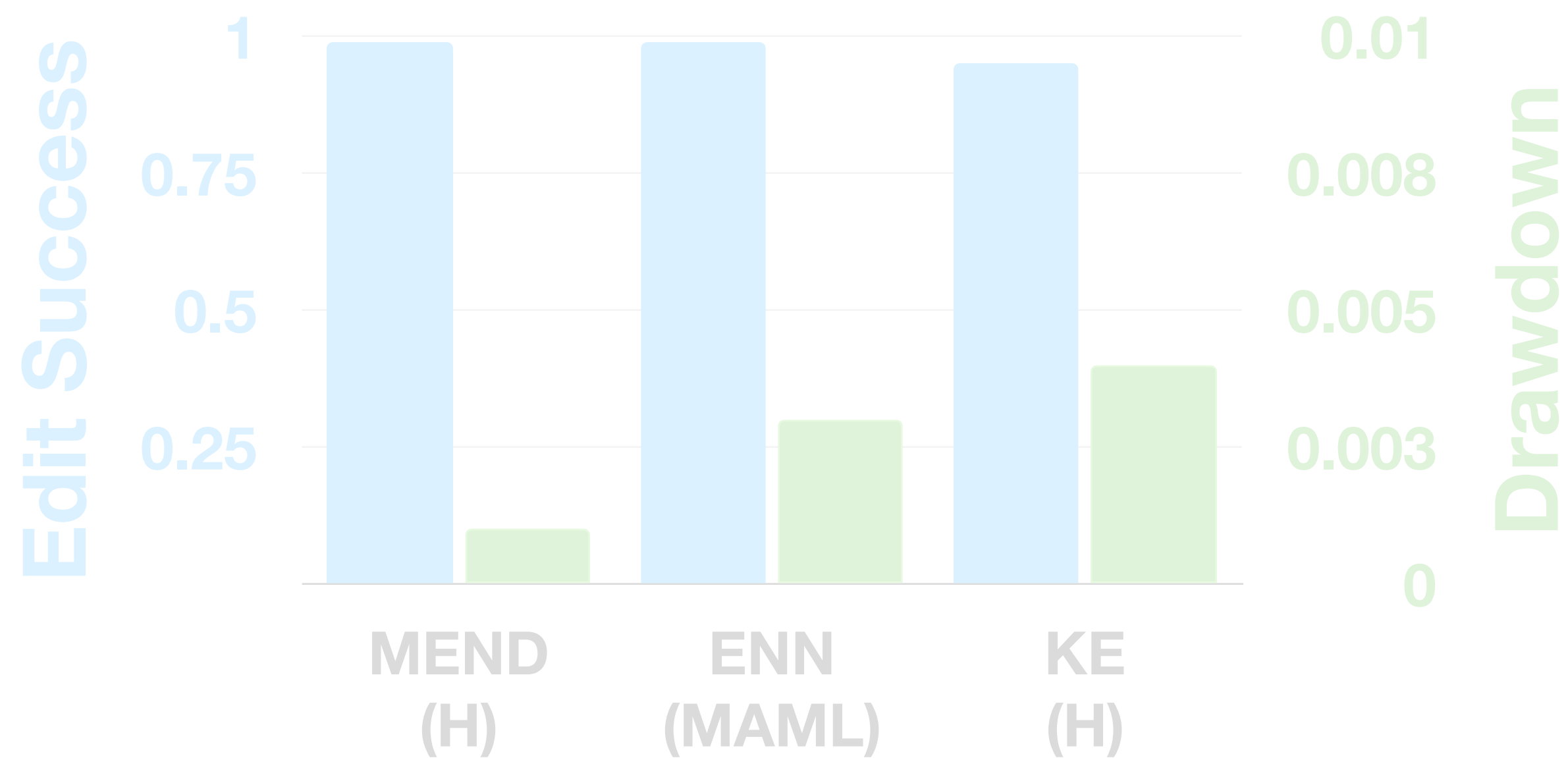
BART-base (110M)



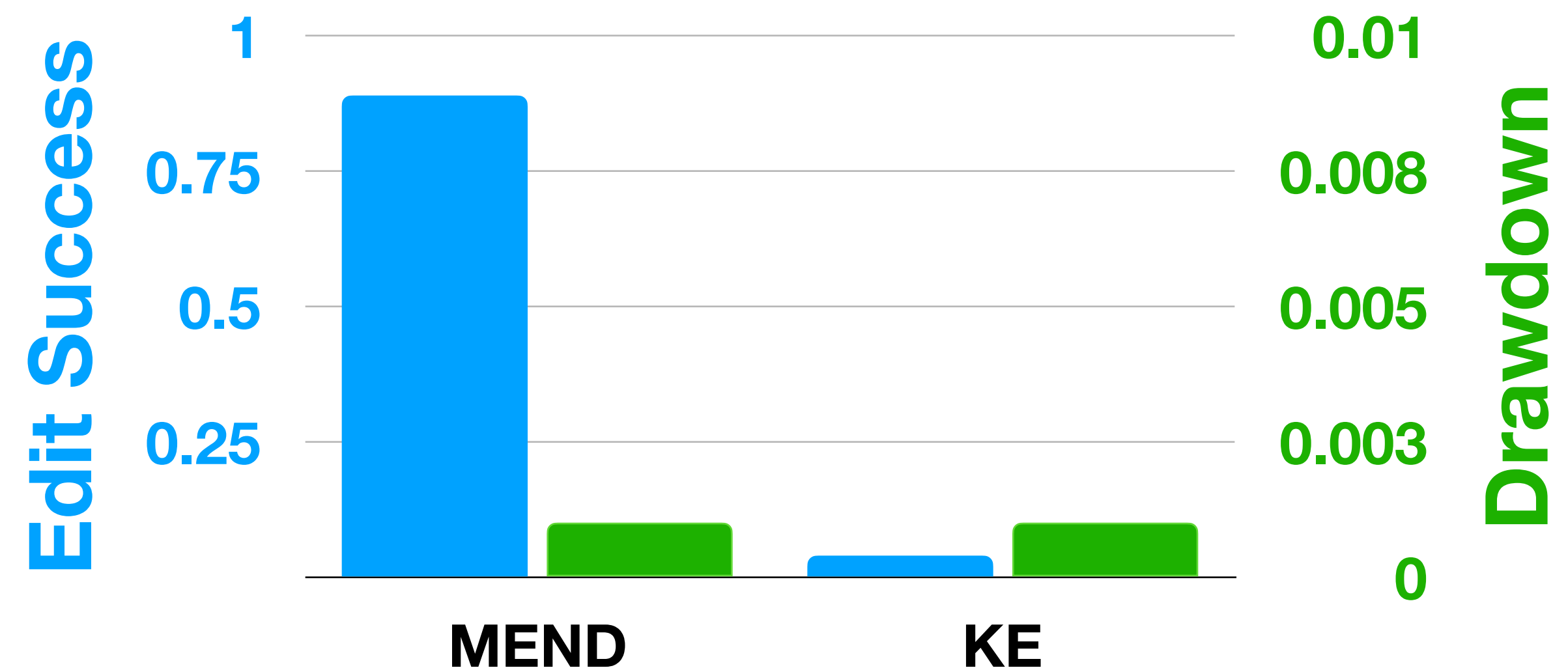
Model Editor Networks using gradient Decomposition

Effective editing at small scale...and large scale!

BART-base (110M)



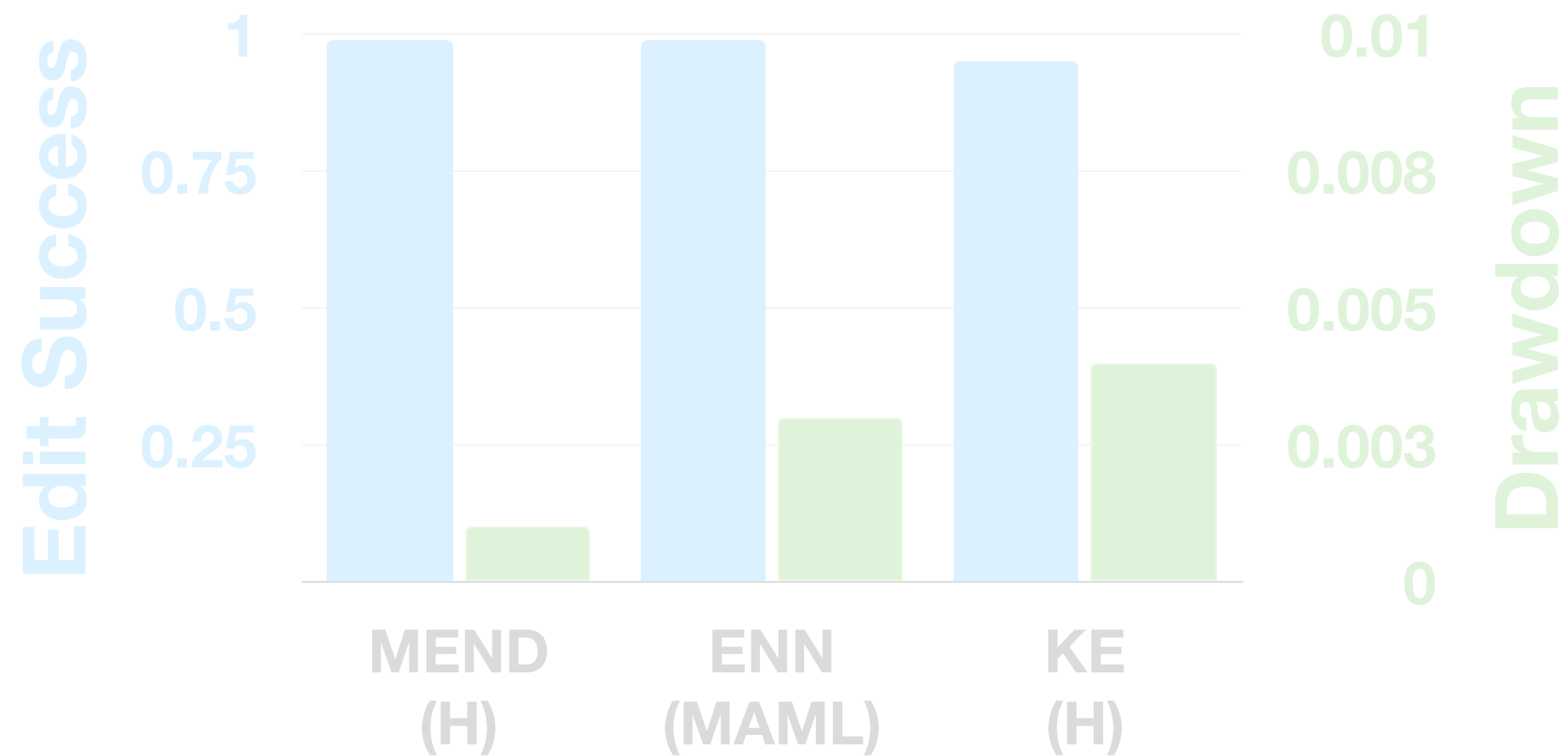
T5-XXL (11B)



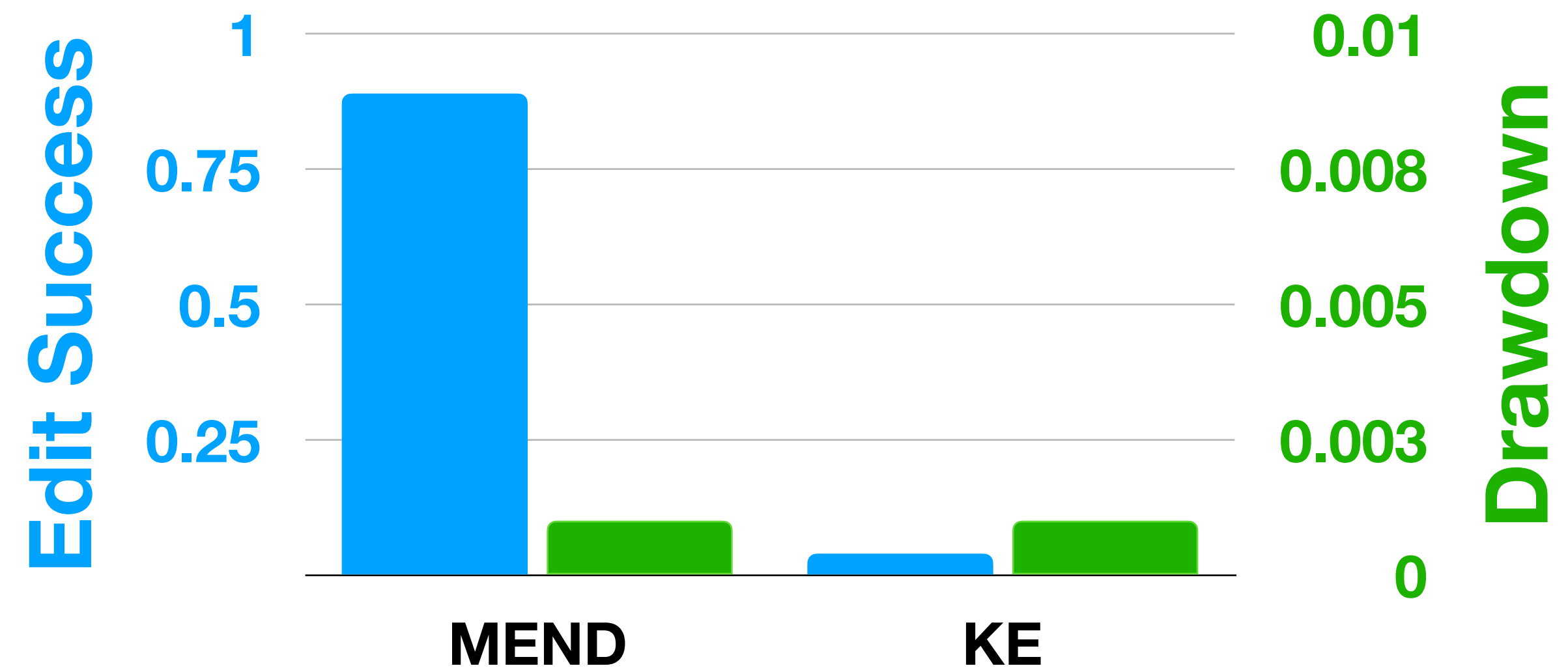
Model Editor Networks using gradient Decomposition

Effective editing at small scale...and large scale!

BART-base (110M)



T5-XXL (11B)



MEND gives productive edits, while KE fails (and ENN gives OOM)

Model Editor Networks using gradient Decomposition

Editing T5-Large: Successes and a failure

Input	Pre-Edit Output	Edit Target	Post-Edit Output
1a: Who is India's PM?	Satya Pal Malik ✗	Narendra Modi	Narendra Modi ✓
1b: Who is the prime minister of the UK?	Theresa May ✗	Boris Johnson	Boris Johnson ✓
1c: Who is the prime minister of India?	Narendra Modi ✓	-	Narendra Modi ✓
1d: Who is the UK PM?	Theresa May ✗	-	Boris Johnson ✓
2a: What is Messi's club team?	Barcelona B ✗	PSG	PSG ✓
2b: What basketball team does LeBron play on?	Dallas Mavericks ✗	the LA Lakers	the LA Lakers ✓
2c: Where in the US is Raleigh?	a state in the South ✓	-	a state in the South ✓
3a: Who is the president of Mexico?	Enrique Pea Nieto ✗	Andrés Manuel López Obrador	Andrés Manuel López Obrador ✓
3b: Who is the vice president of Mexico?	Yadier Benjamin Ramos ✗	-	Andrés Manuel López Obrador ✗

Bold text indicates the edits applied in each evaluation

Today's Plan

I. Background

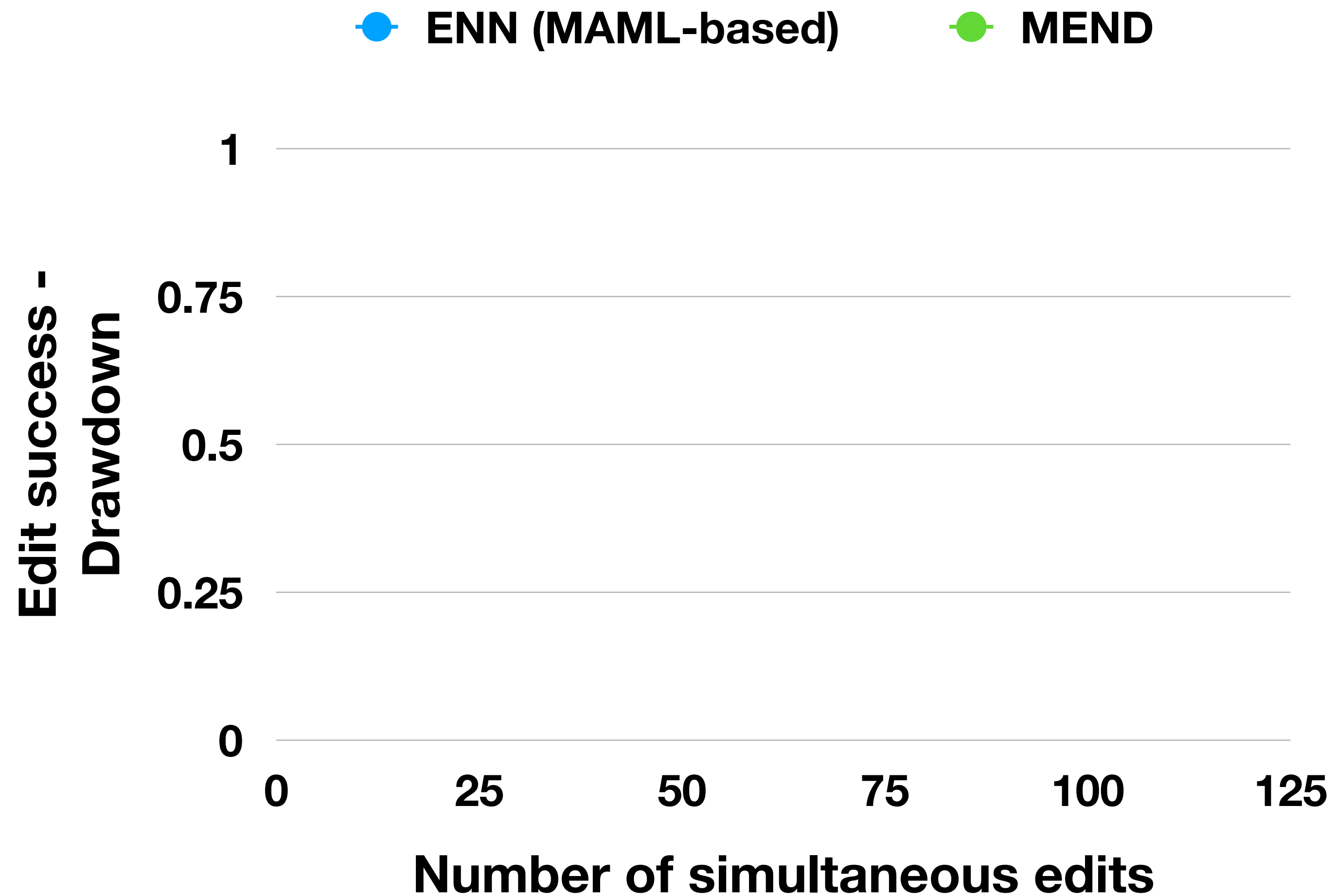
II. Learning to edit NNs

III. Moving editing towards the real world

IV. Future work & open questions

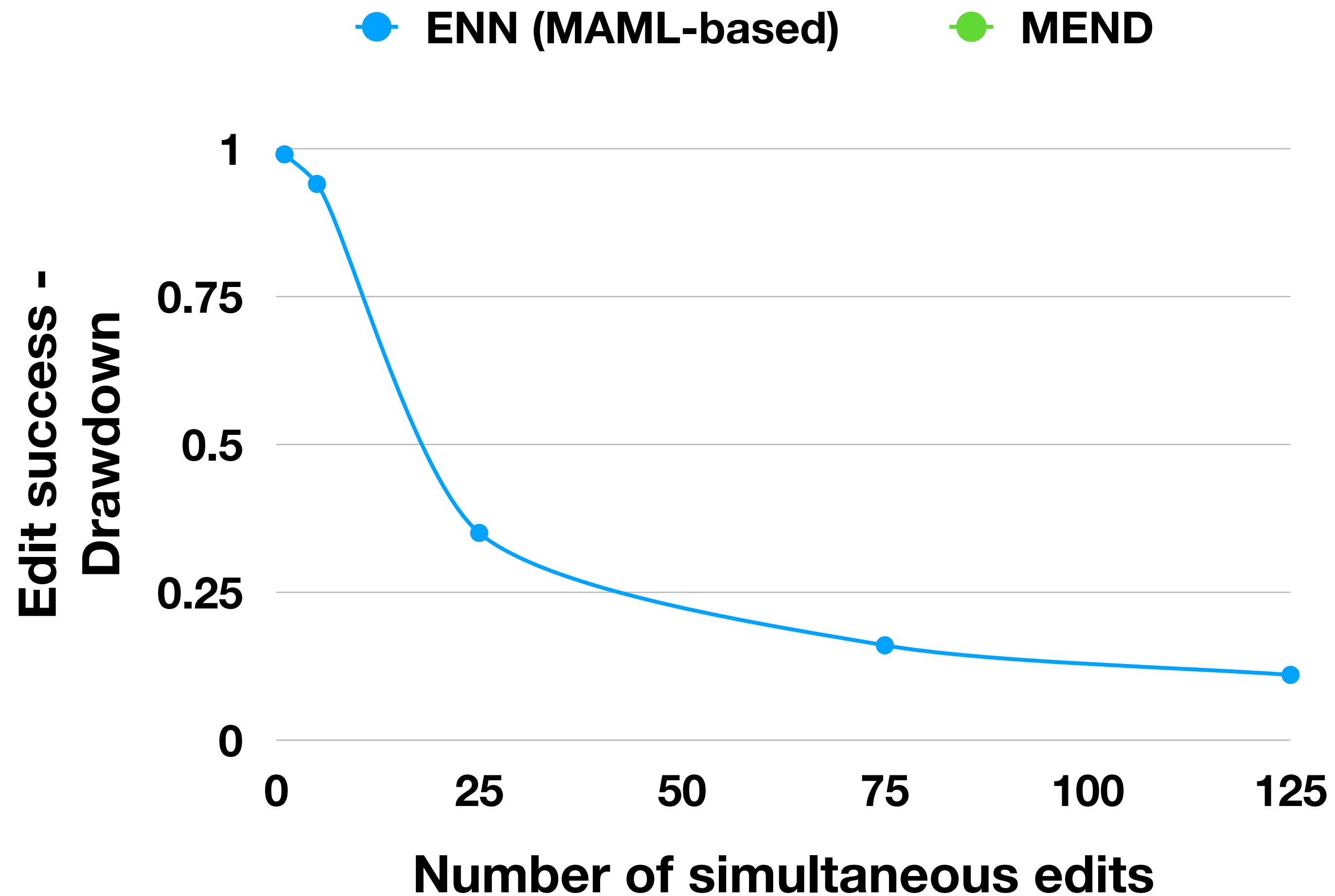
Moving editing towards the real world

Applying multiple edits to BART-base



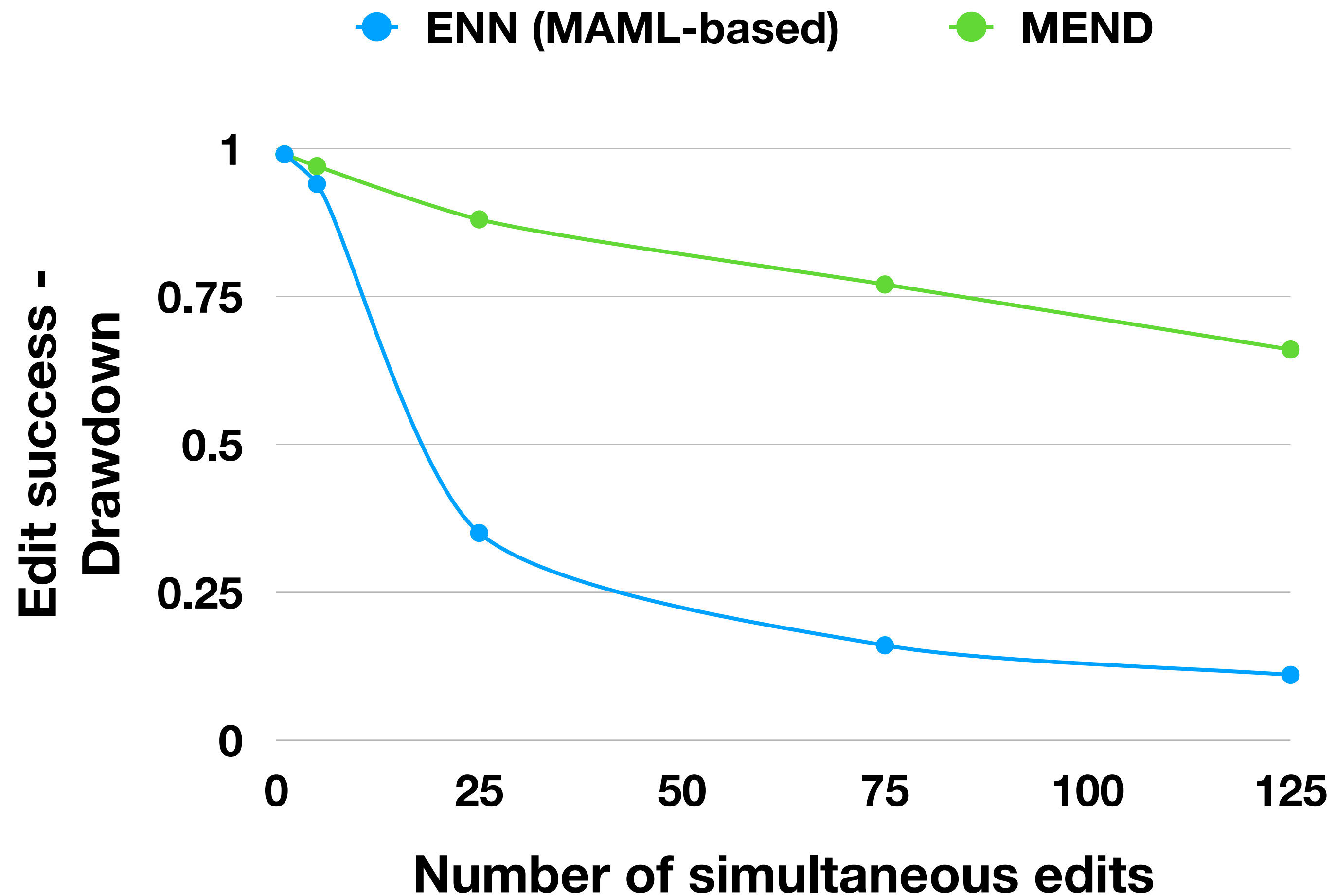
Moving editing towards the real world

Applying multiple edits to BART-base



Moving editing towards the real world

Applying multiple edits to BART-base



More challenging benchmarks

Multiple edits, more difficult edit scopes

Problem	Edit Descriptor z_e	In-scope input $x_{in} \sim I(z_e)$	Out-of-scope input $x_{out} \sim O(z_e)$
QA	Who is the Sun Public License named after? <i>Sun Micro Devices</i>	The Sun Public License has been named for whom? <i>Sun Micro Devices</i>	What continent is Mount Whillans found on?

More challenging benchmarks

Multiple edits, more difficult edit scopes

Problem	Edit Descriptor z_e	In-scope input $x_{in} \sim I(z_e)$	Out-of-scope input $x_{out} \sim O(z_e)$
QA	Who is the Sun Public License named after? <i>Sun Micro Devices</i>	The Sun Public License has been named for whom? <i>Sun Micro Devices</i>	What continent is Mount Whillans found on?
QA-hard	What type of submarine was USS Lawrence (DD-8) classified as? <i>Gearing-class destroyer</i>	t/f: Was USS Lawrence (DD-8) classified as Paulding-class destroyer. <i>False</i>	What type of submarine was USS Sumner (DD-333) classified as?

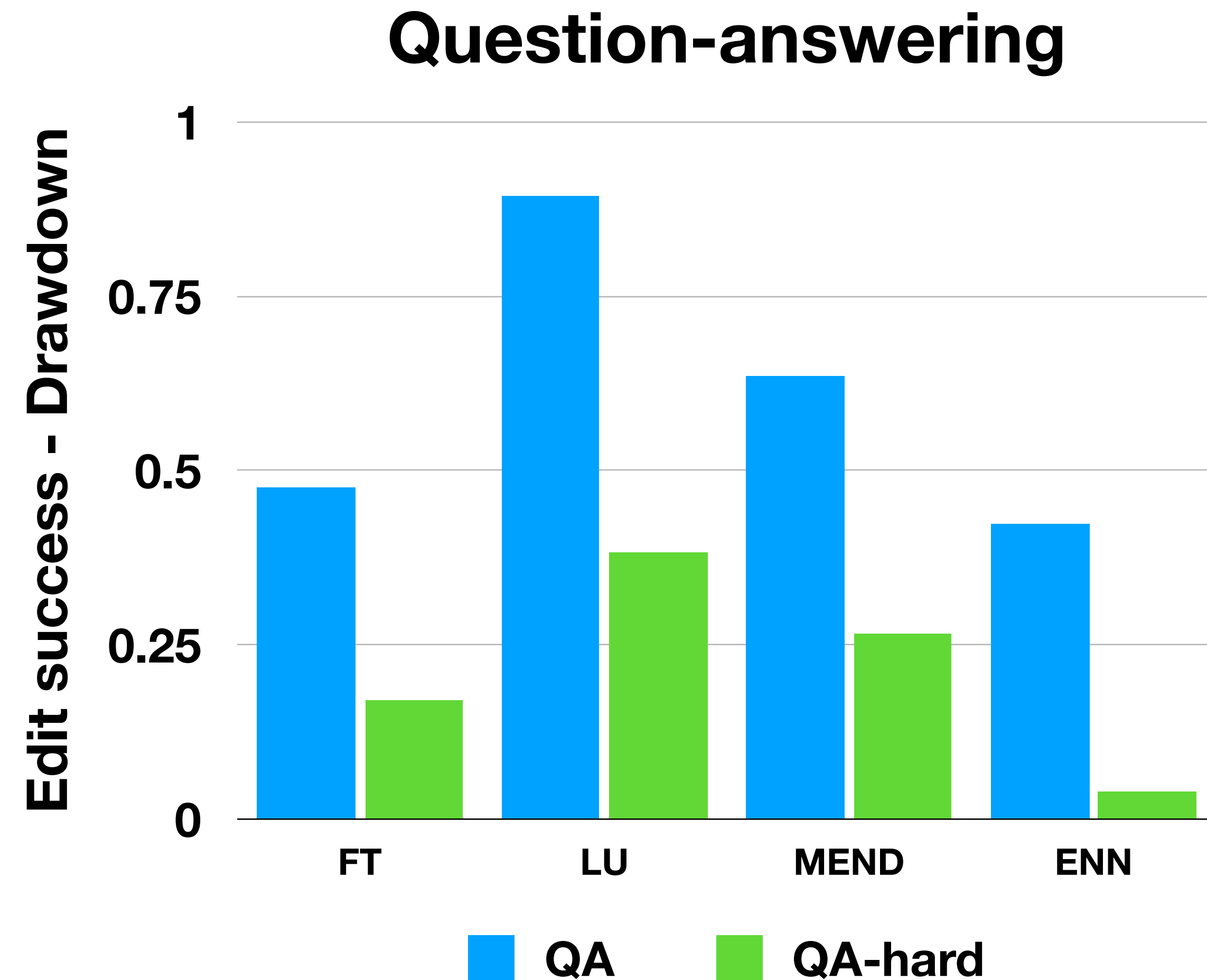
More challenging benchmarks

Multiple edits, more difficult edit scopes

Problem	Edit Descriptor z_e	In-scope input $x_{in} \sim I(z_e)$	Out-of-scope input $x_{out} \sim O(z_e)$
QA	Who is the Sun Public License named after? <i>Sun Micro Devices</i>	The Sun Public License has been named for whom? <i>Sun Micro Devices</i>	What continent is Mount Whillans found on?
QA-hard	What type of submarine was USS Lawrence (DD-8) classified as? <i>Gearing-class destroyer</i>	t/f: Was USS Lawrence (DD-8) classified as Paulding-class destroyer. <i>False</i>	What type of submarine was USS Sumner (DD-333) classified as?
FC	As of March 23, there were 50 confirmed cases and 0 deaths within Idaho. <i>True</i> Between 1995 and 2018, the AFC has sent less than half of the 16 AFC teams to the Super Bowl with only 7 of the 16 individual teams making it. <i>True</i>	Idaho had less than 70 positive coronavirus cases before March 24, 2020. <i>True</i> –	Alessandro Diamanti scored six serie A goals. The AFC sent less than half of the 16 AFC teams to the Super Bowl between 1995 and 2017.

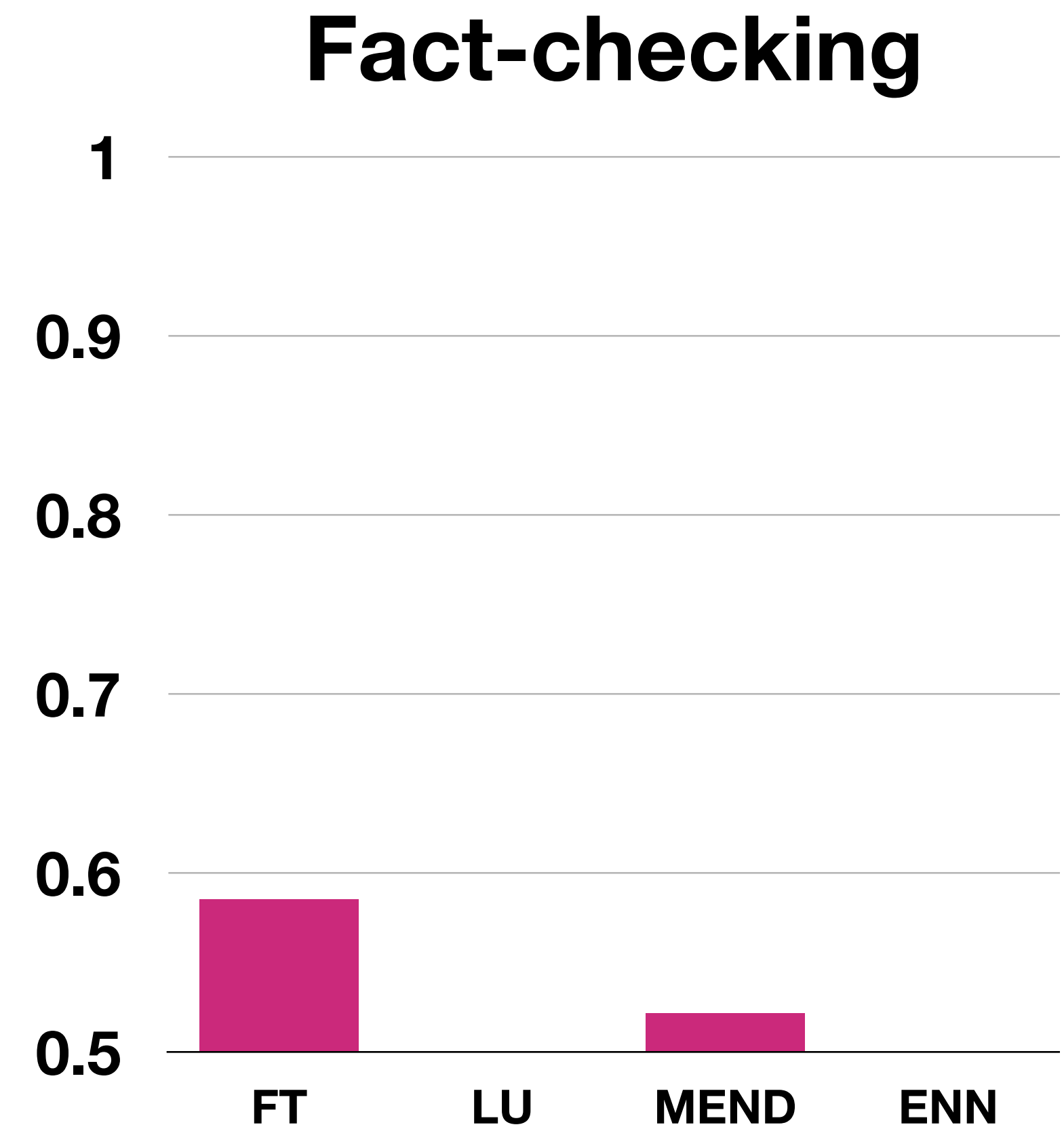
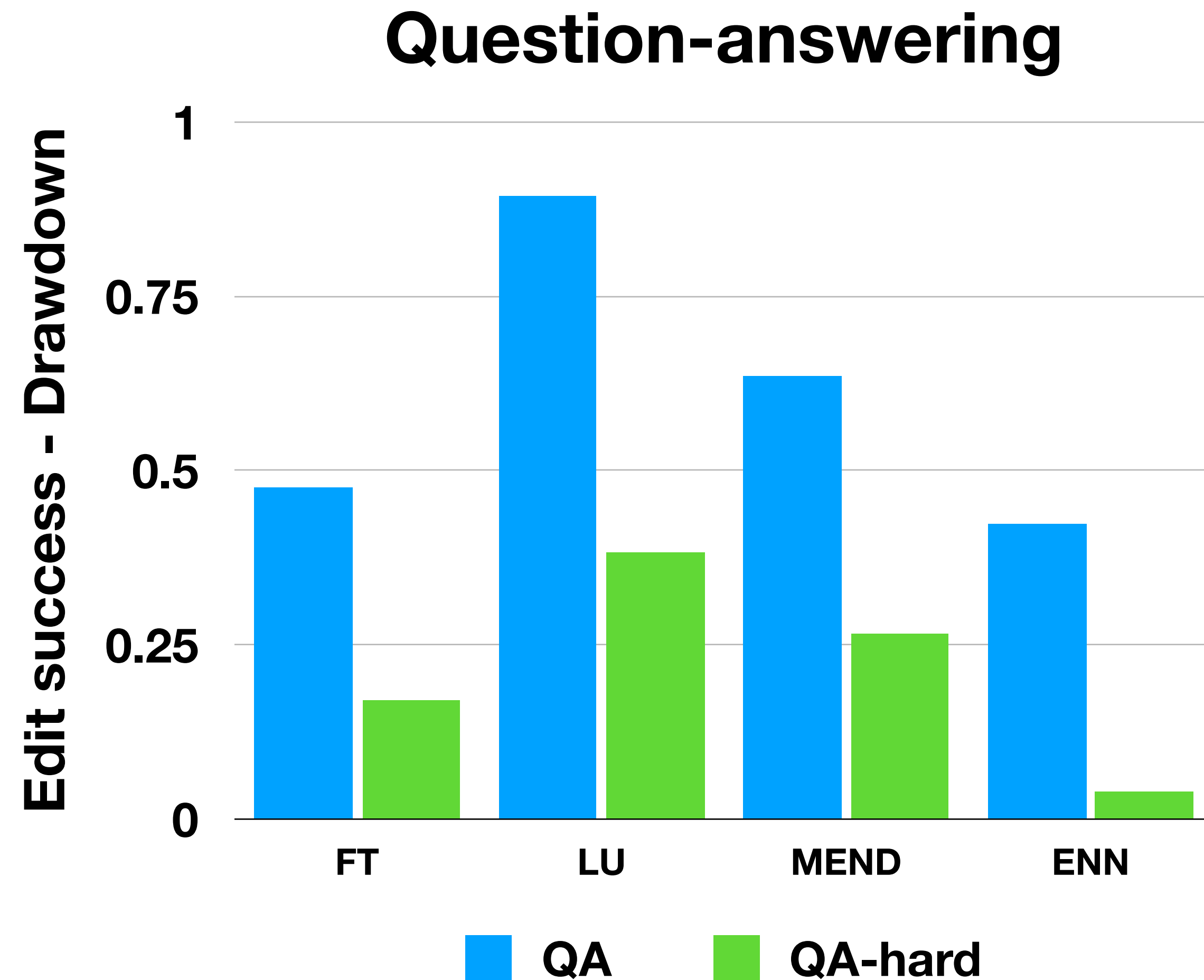
More challenging benchmarks

Multiple edits, more difficult edit scopes



More challenging benchmarks

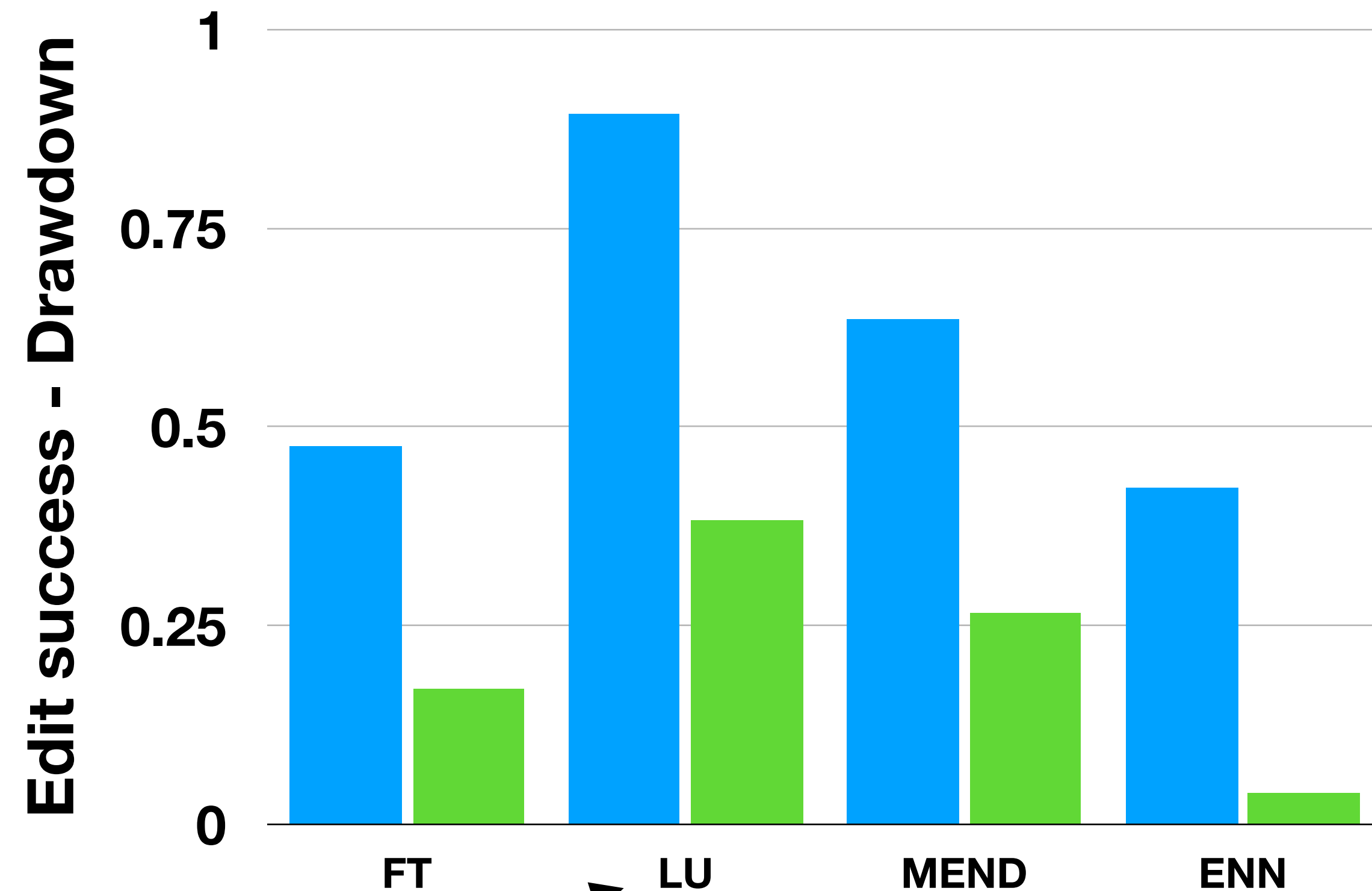
Multiple edits, more difficult edit scopes



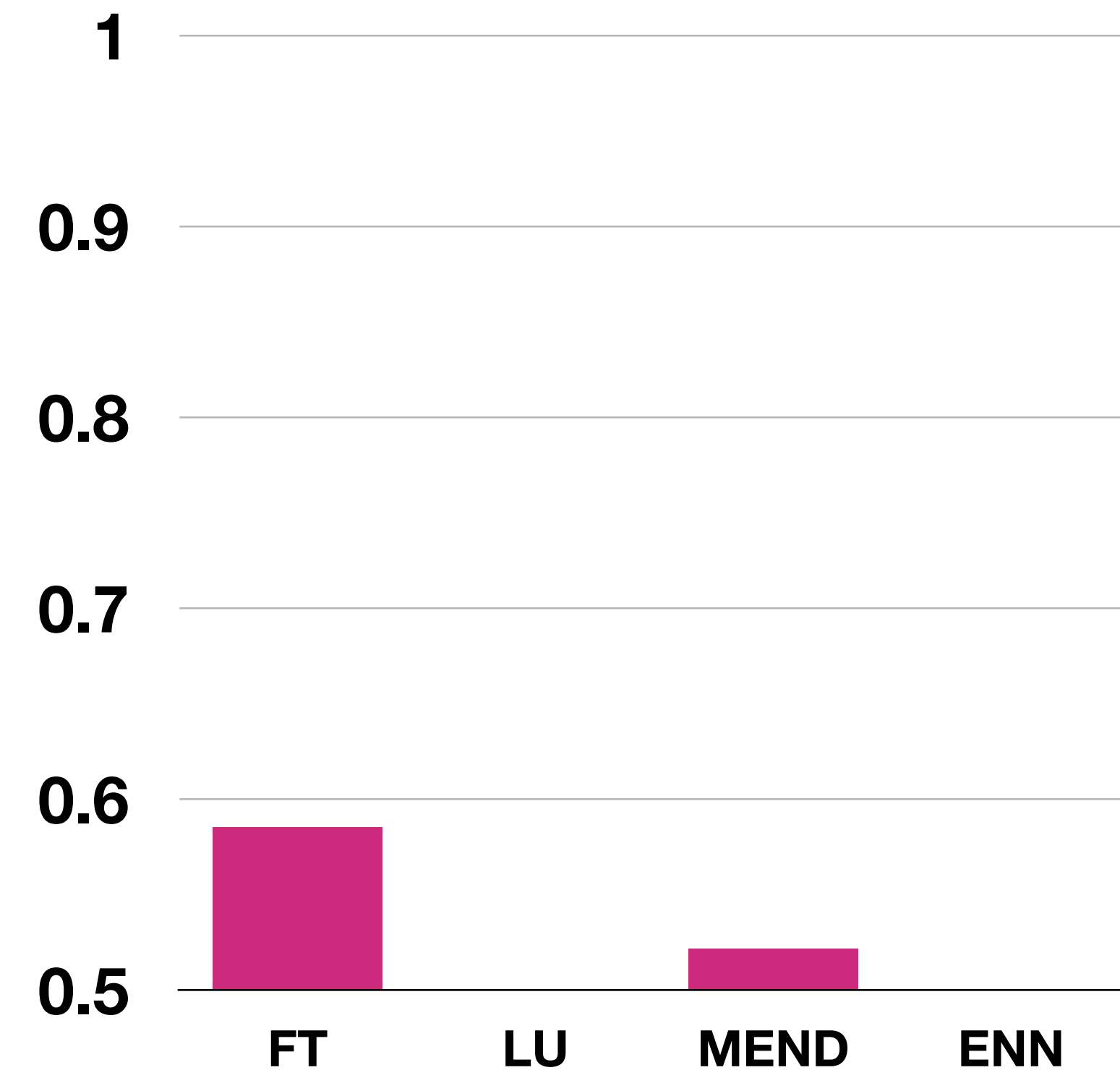
More challenging benchmarks

Multiple edits, more difficult edit scopes

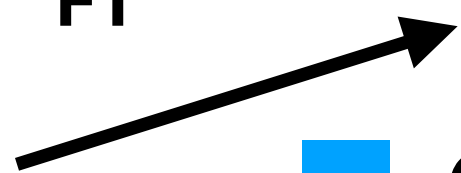
Question-answering



Fact-checking



Surprisingly effective!



■ QA

■ QA-hard

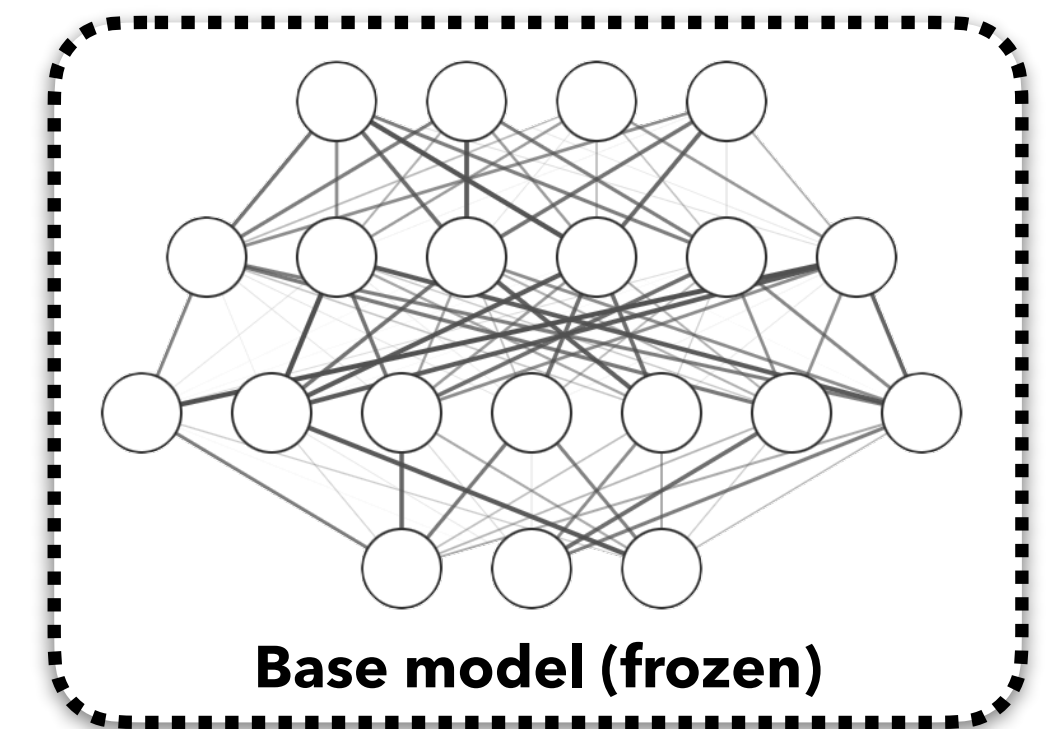
Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

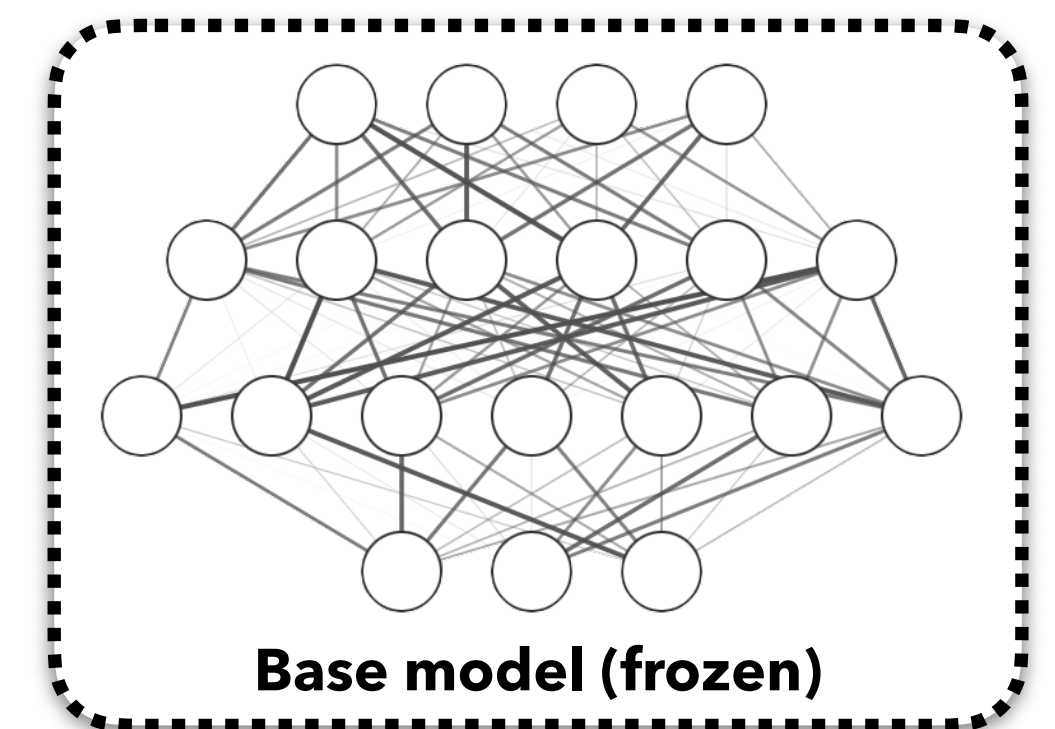
Start with the **frozen** base model



Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

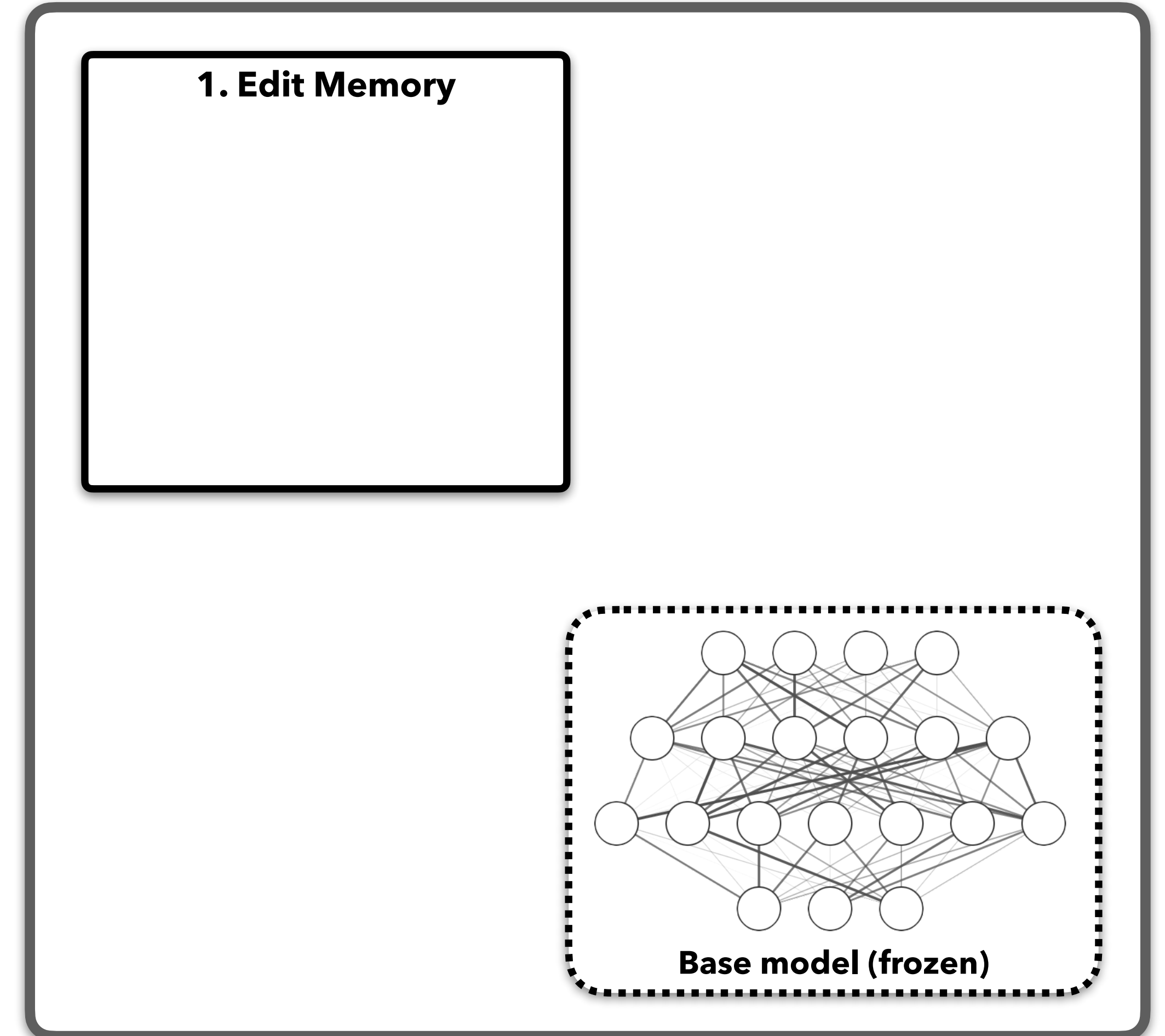


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**

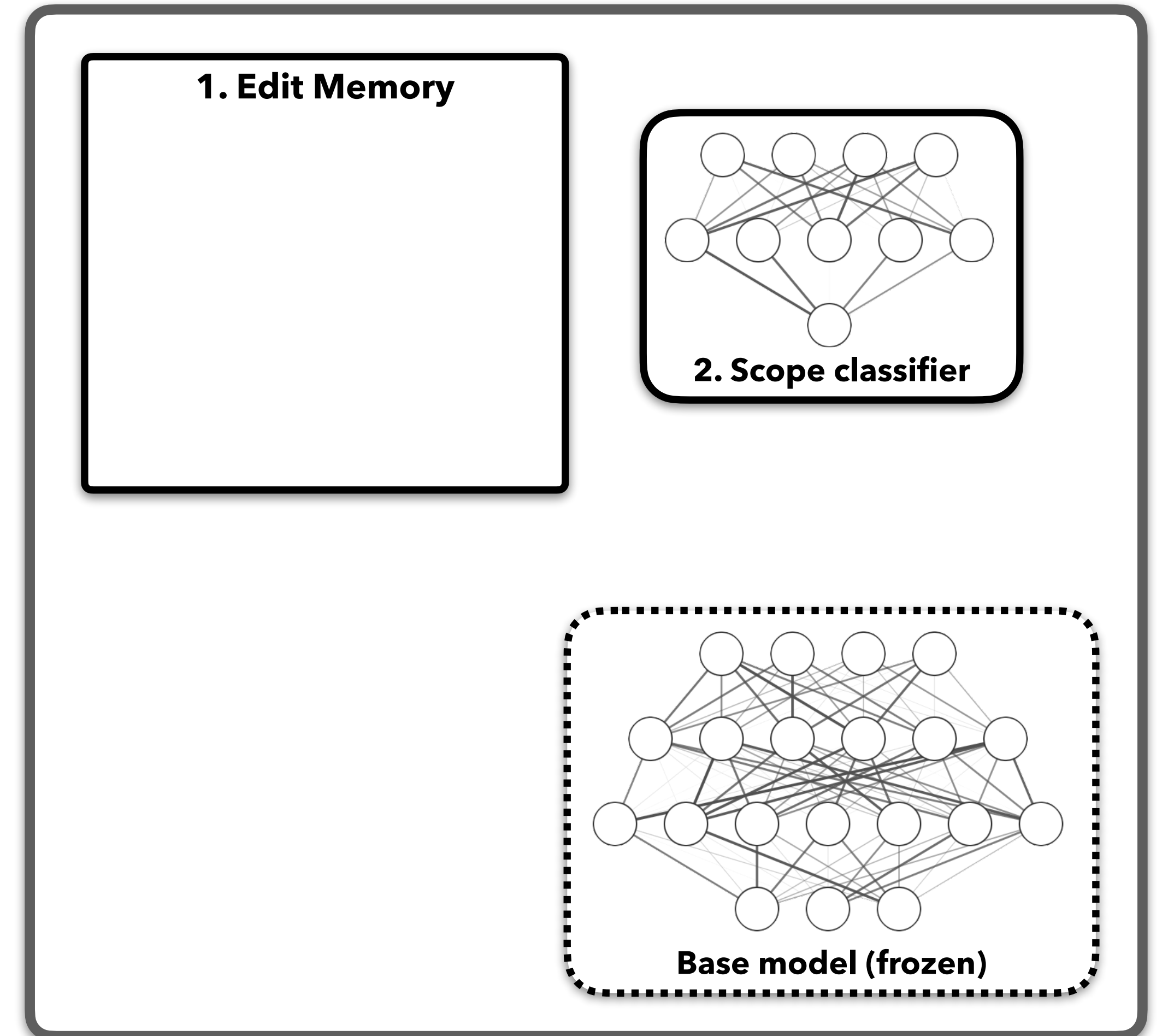


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed

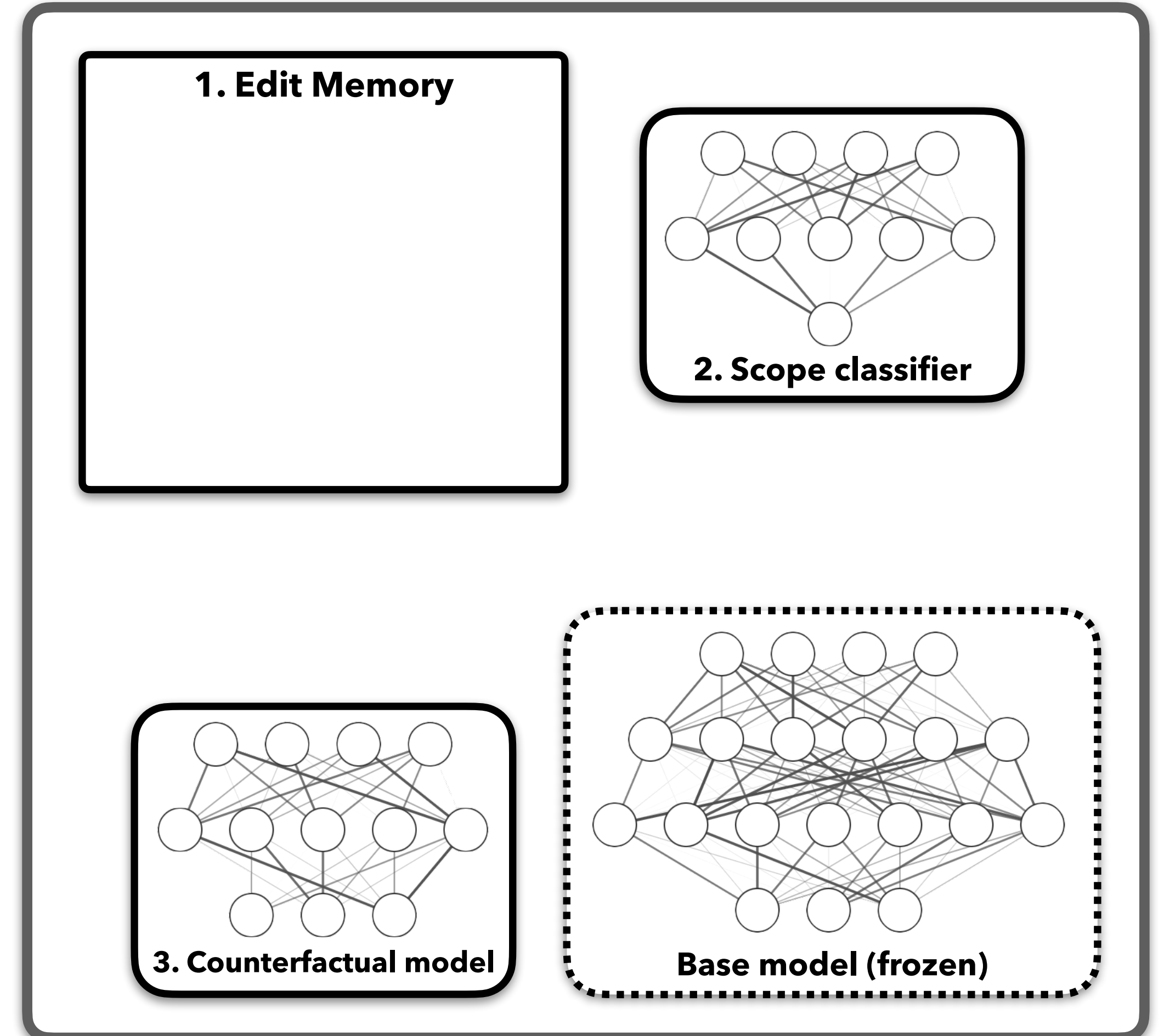


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

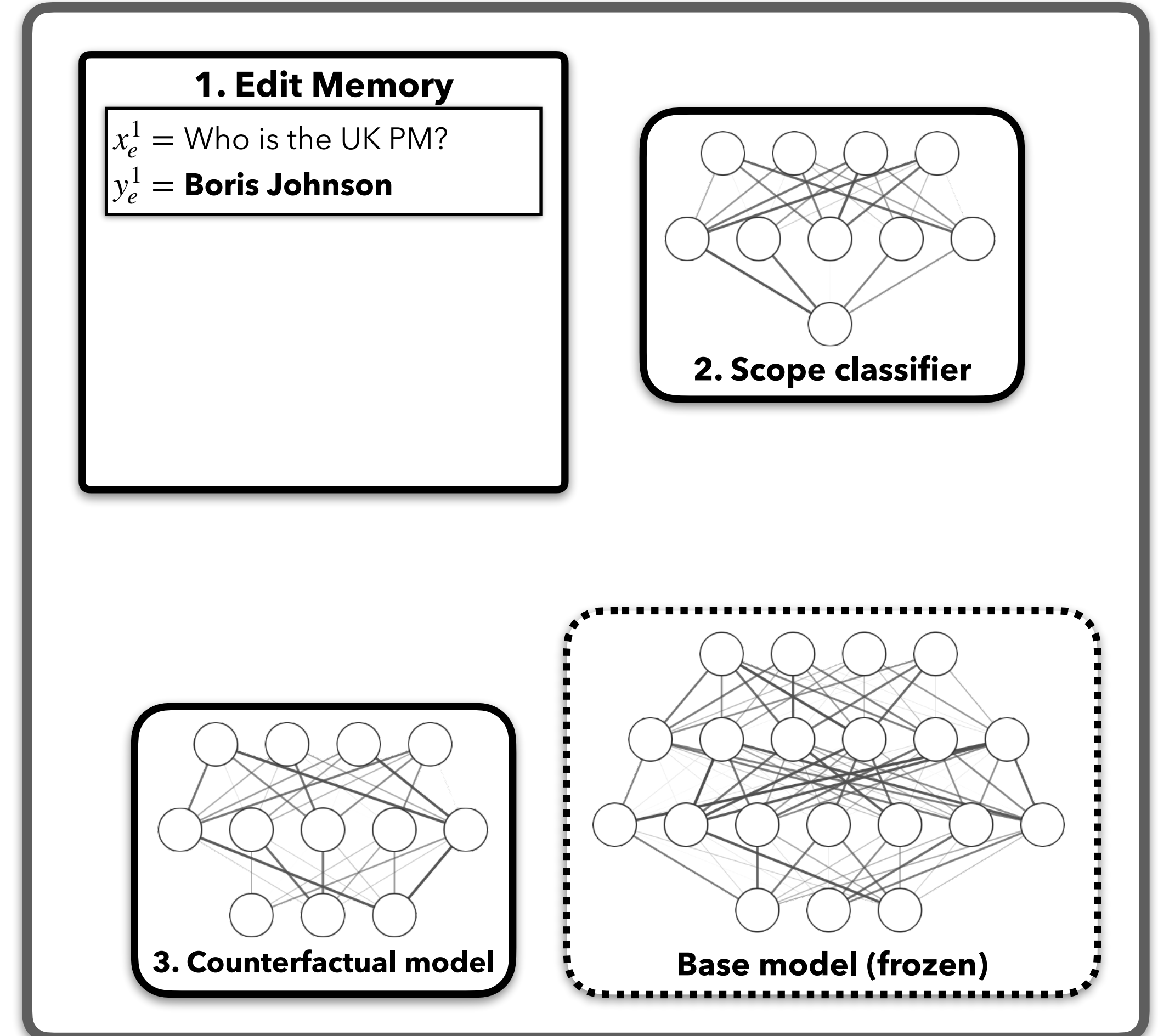


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

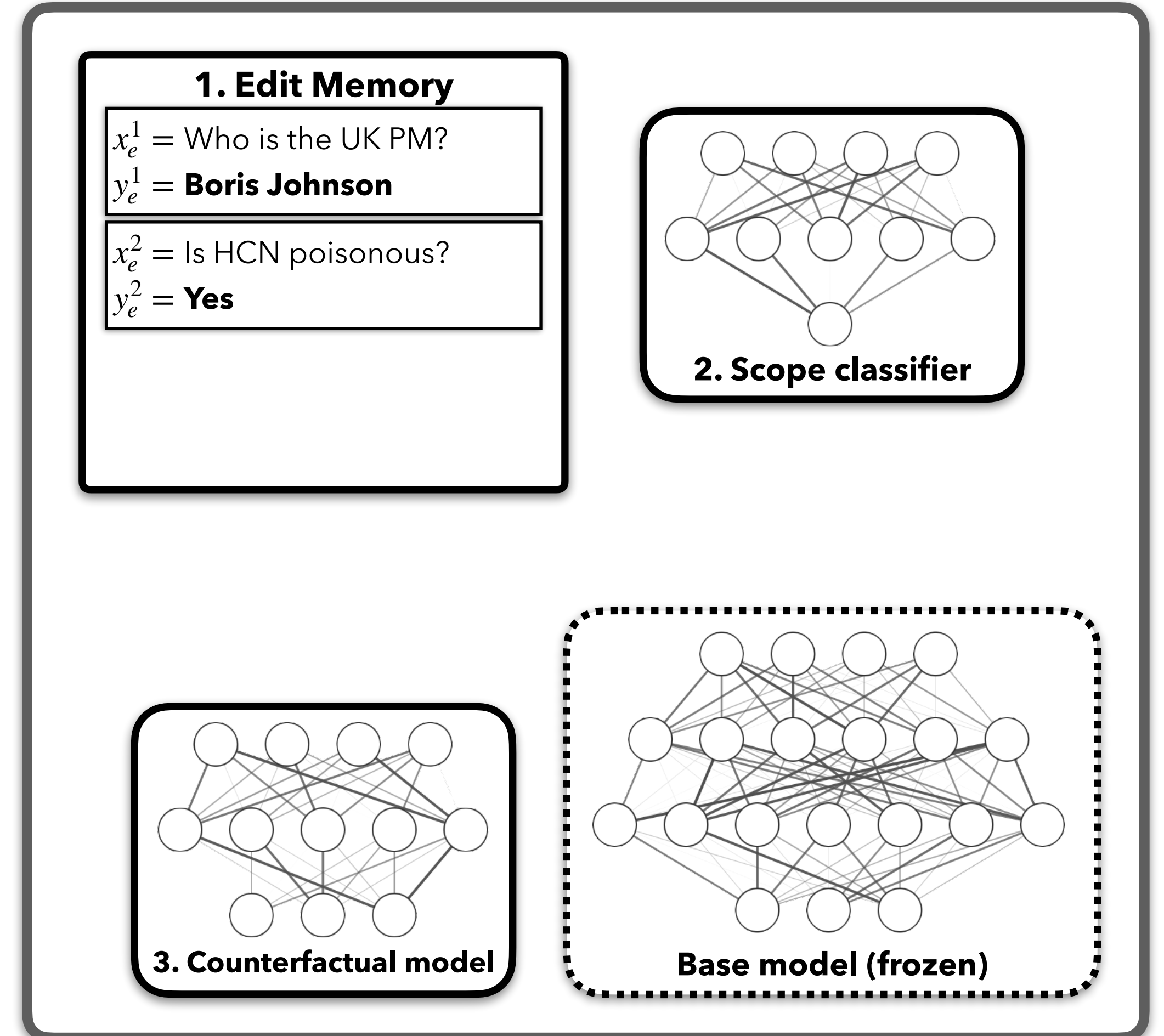


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

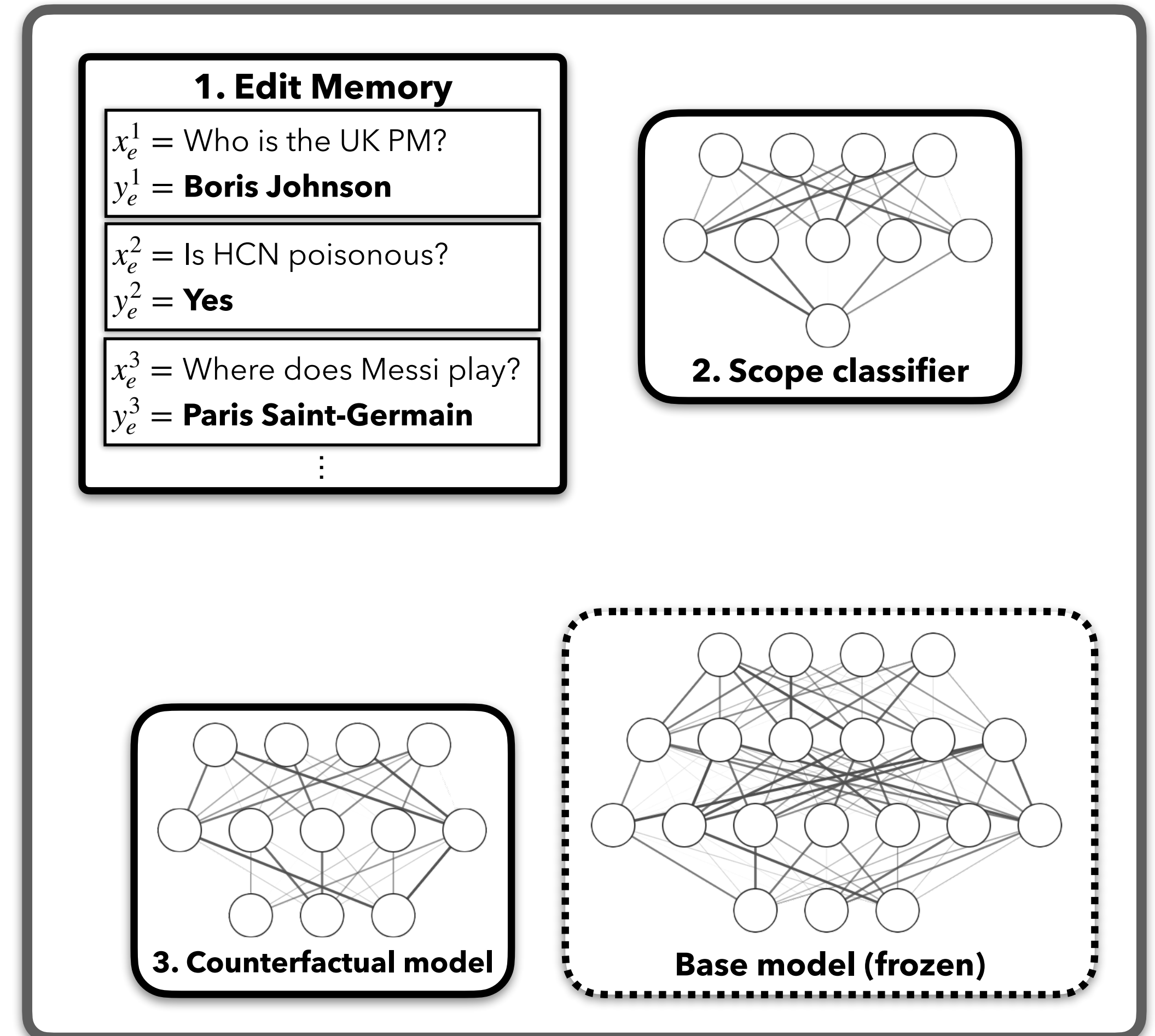


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

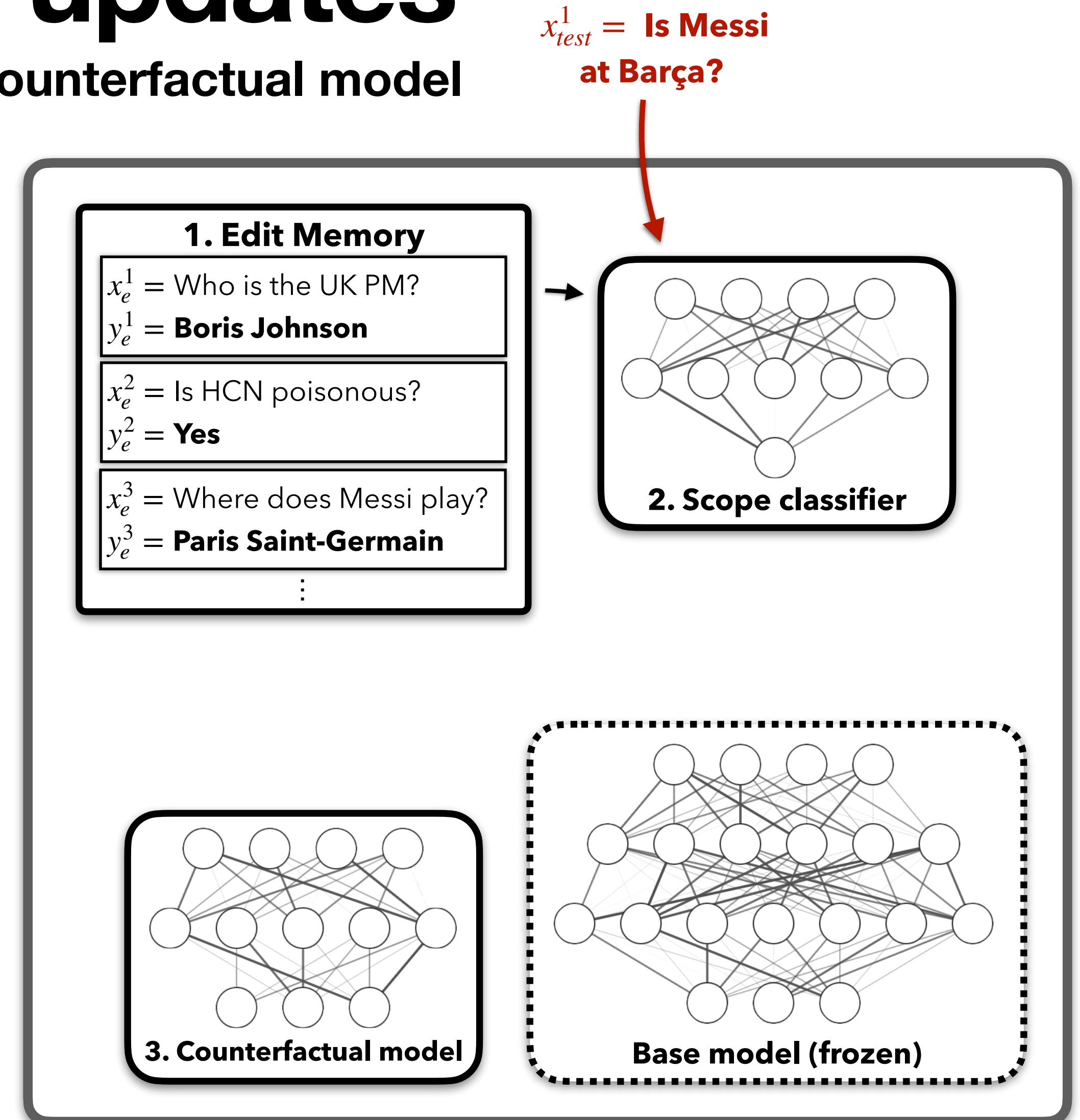


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

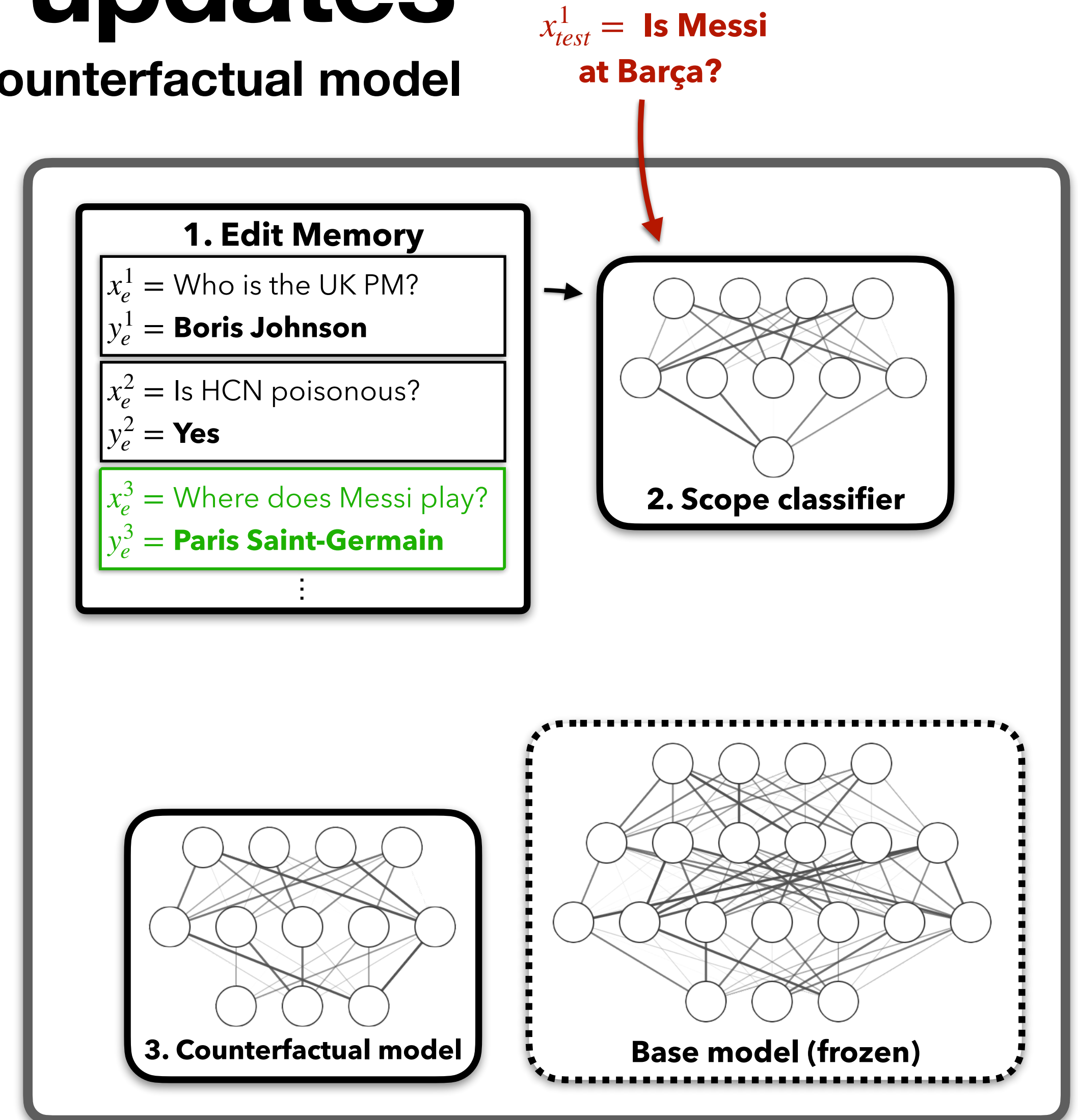


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

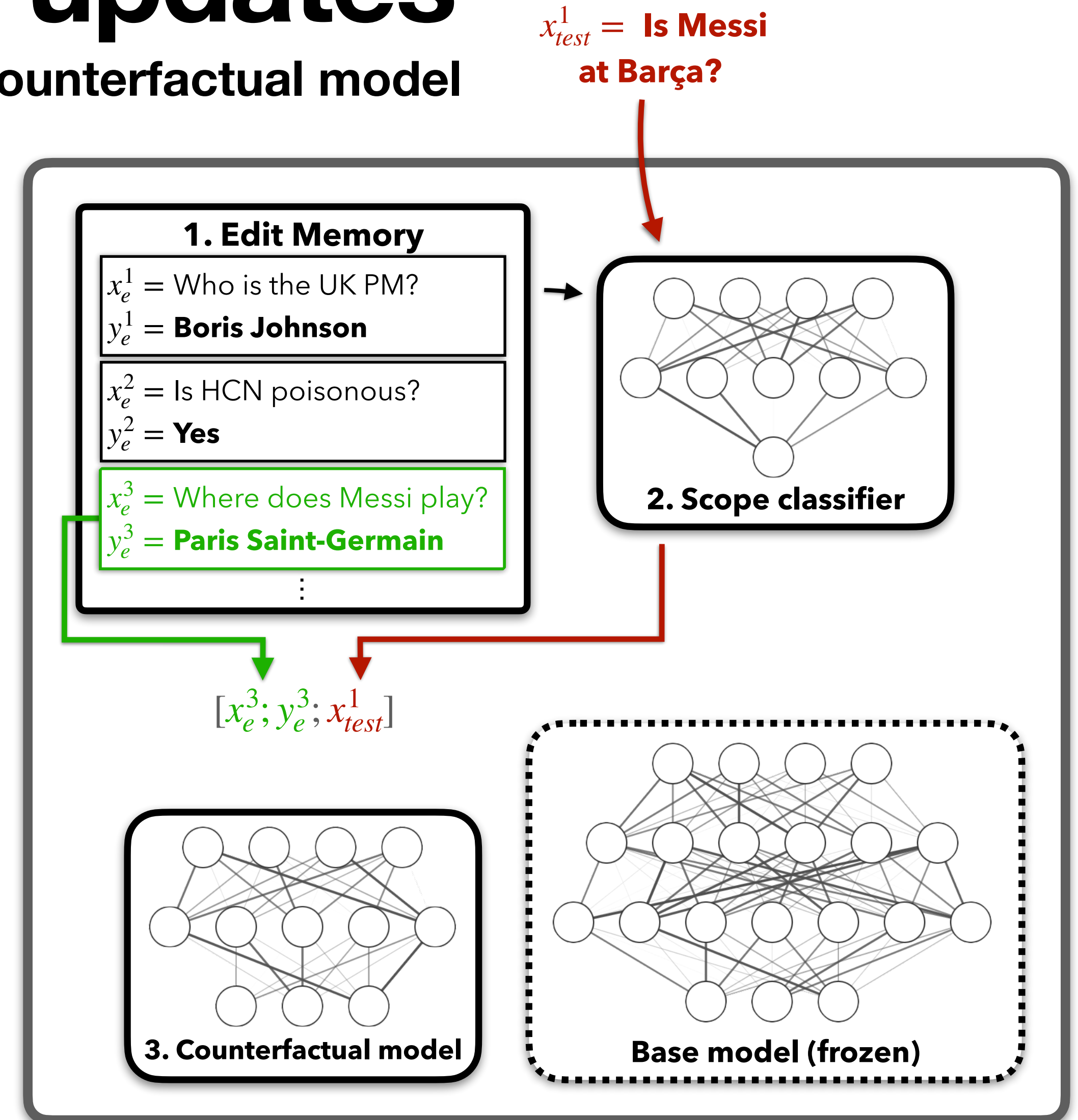


Edits without parameter updates

Semi-parametric Editng with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

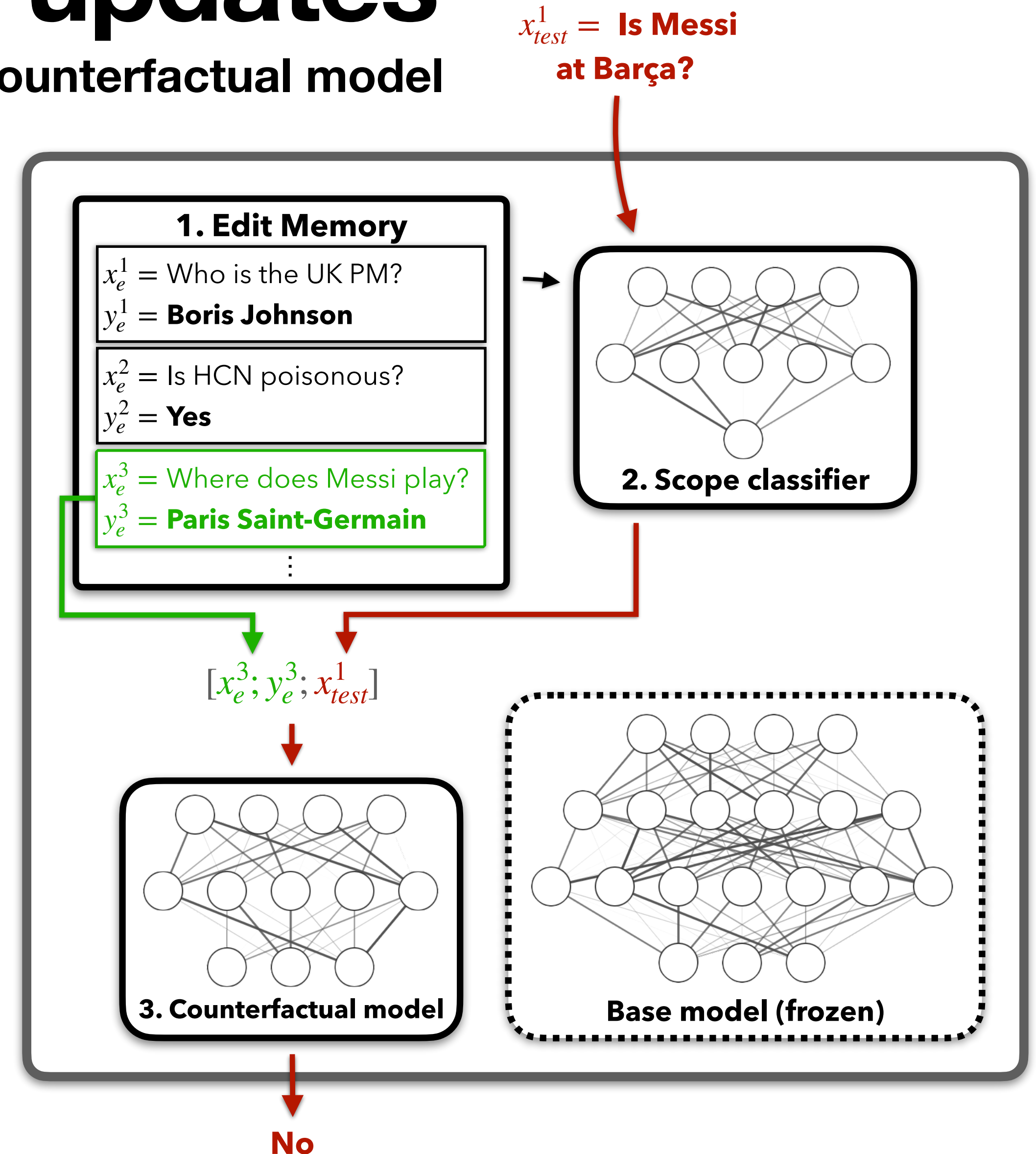


Edits without parameter updates

Semi-parametric Editng with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

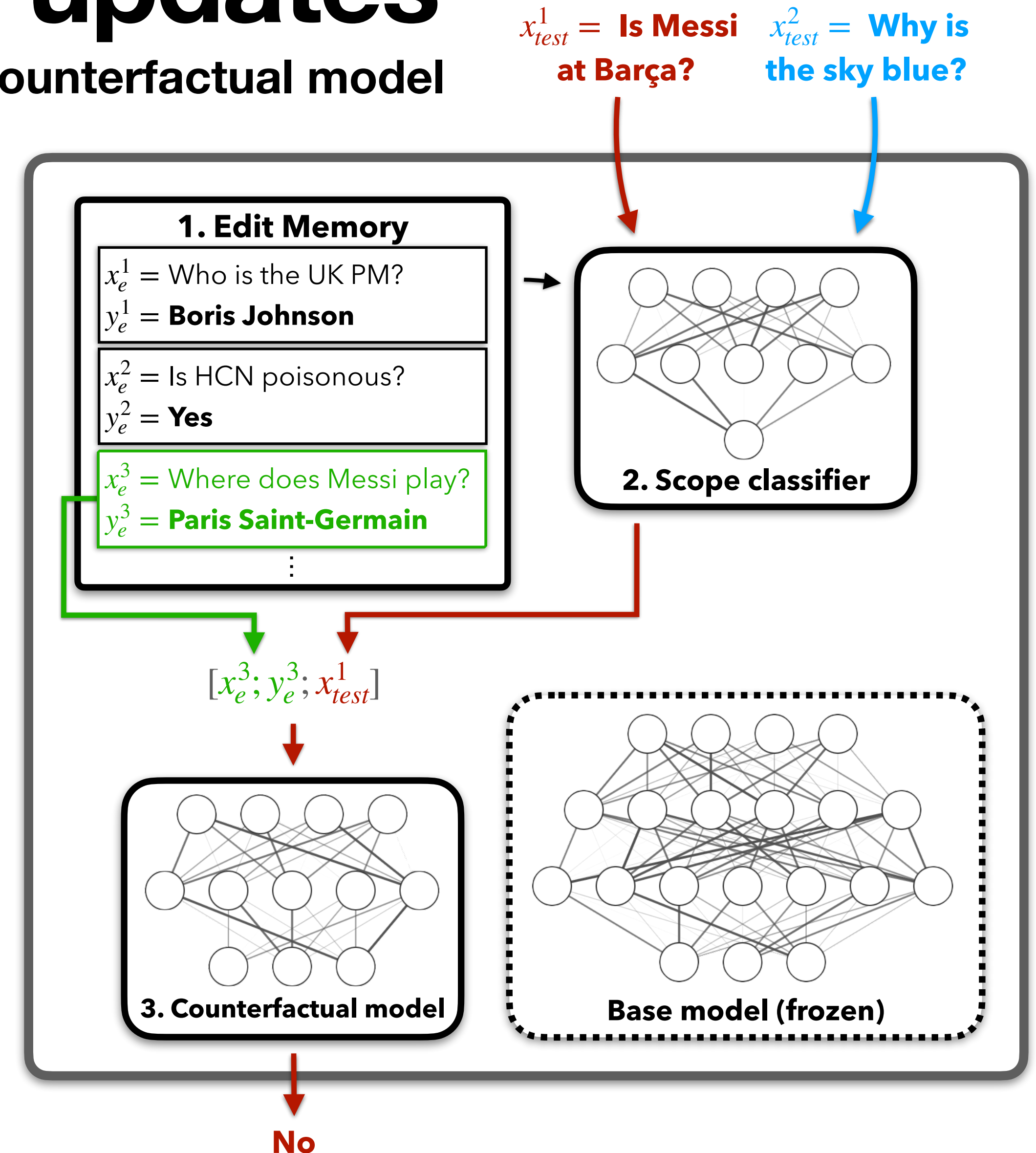


Edits without parameter updates

Semi-parametric Editing with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

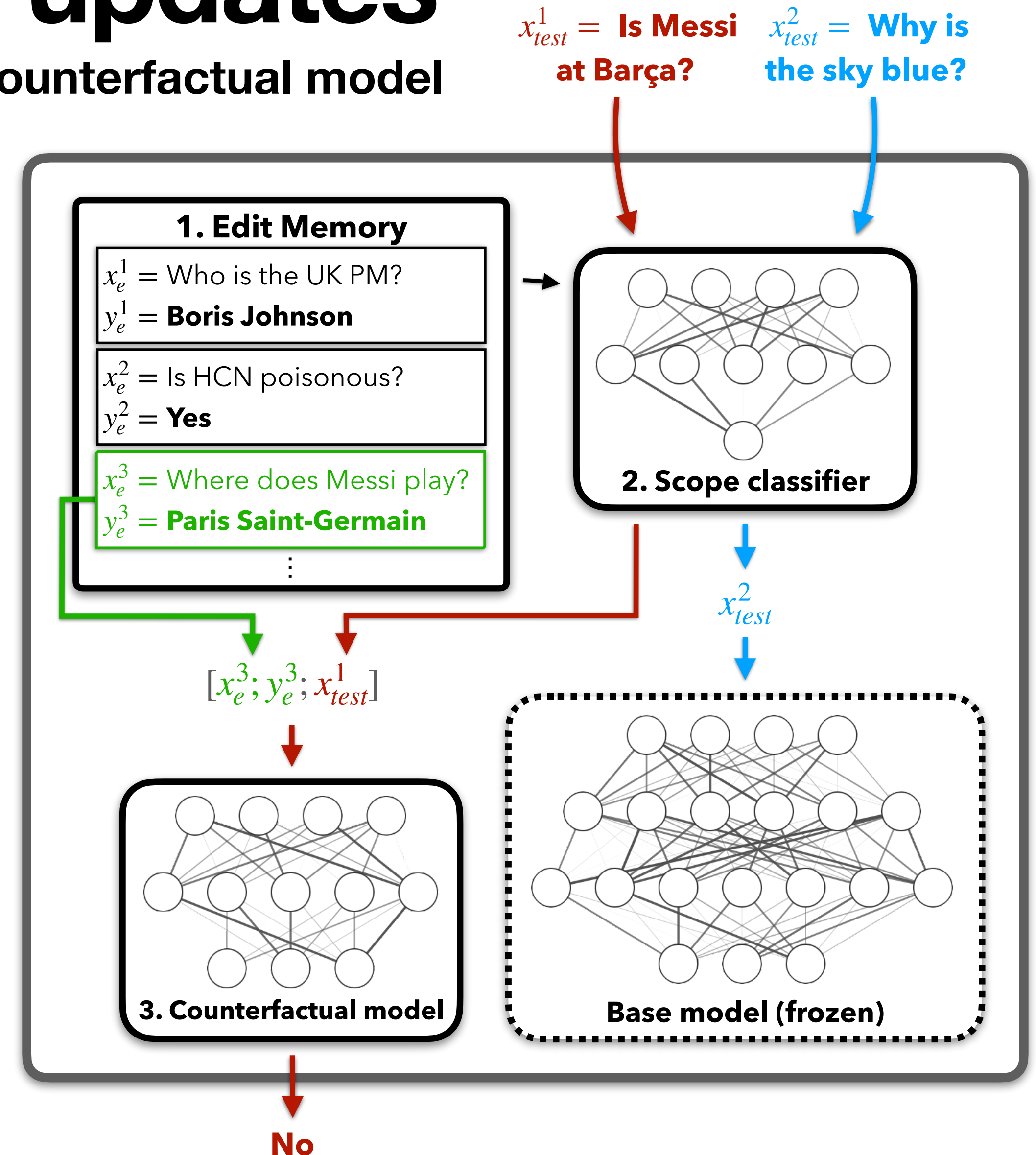


Edits without parameter updates

Semi-parametric Editng with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed



Edits without parameter updates

Semi-parametric Editng with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

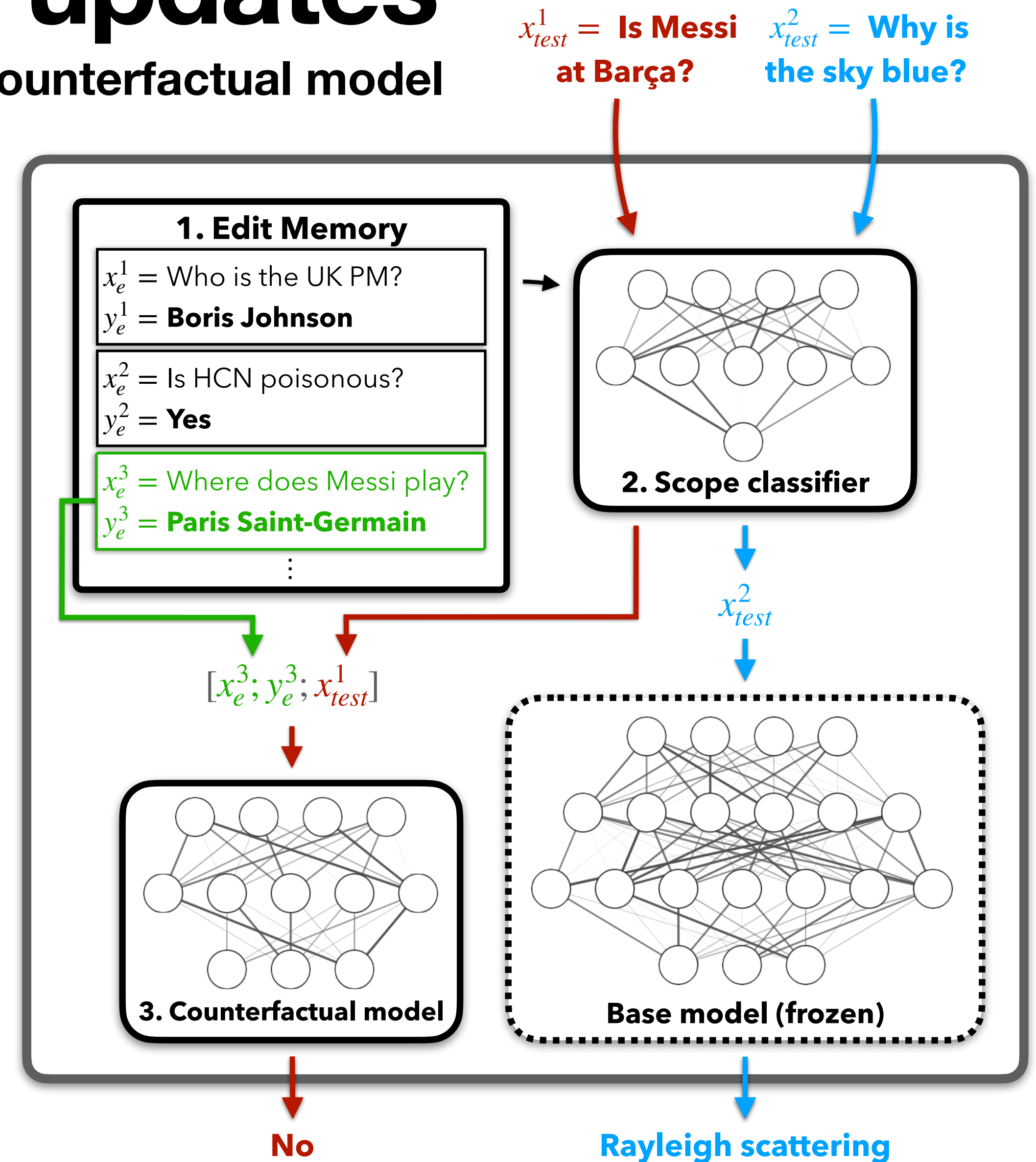
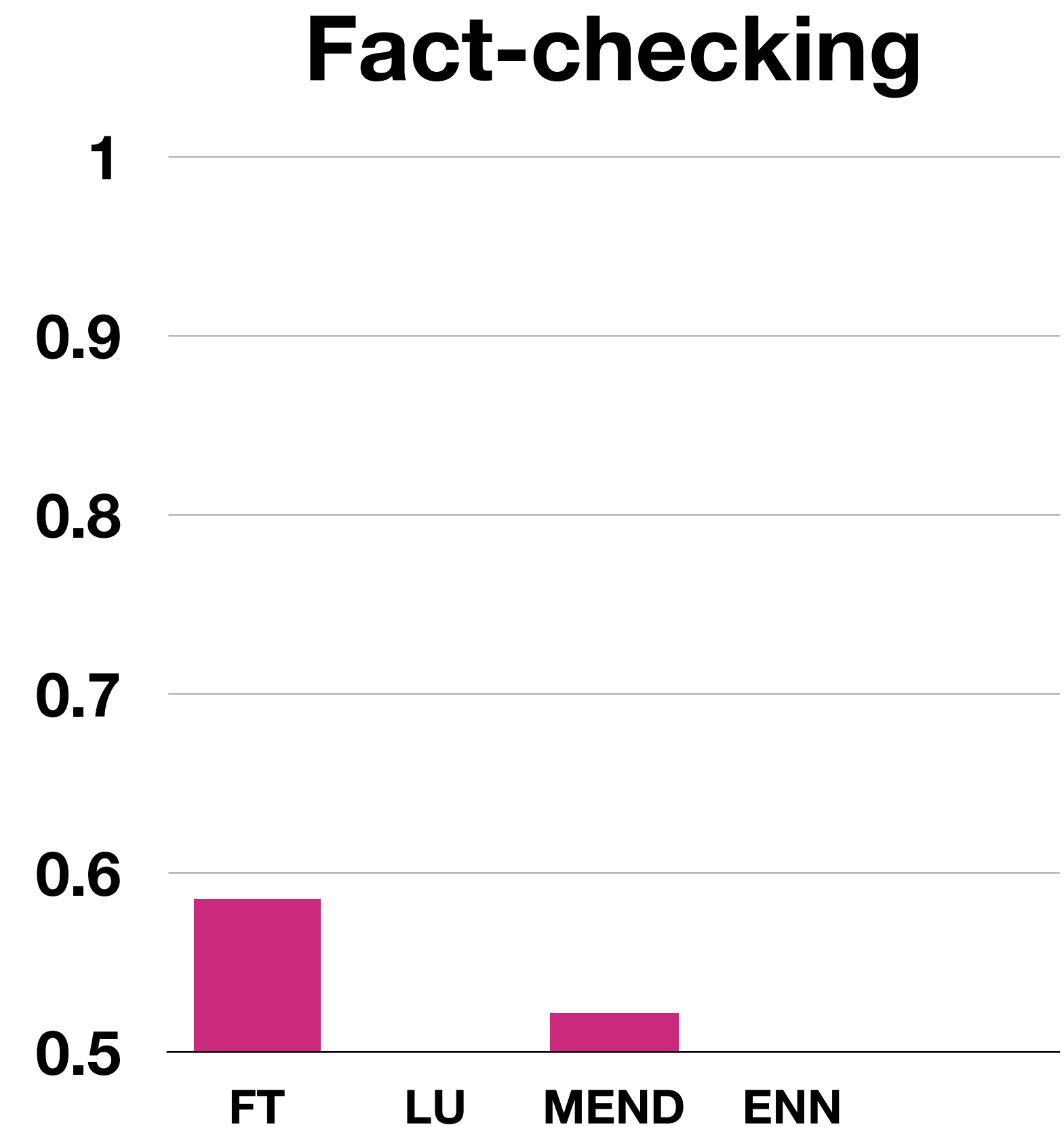
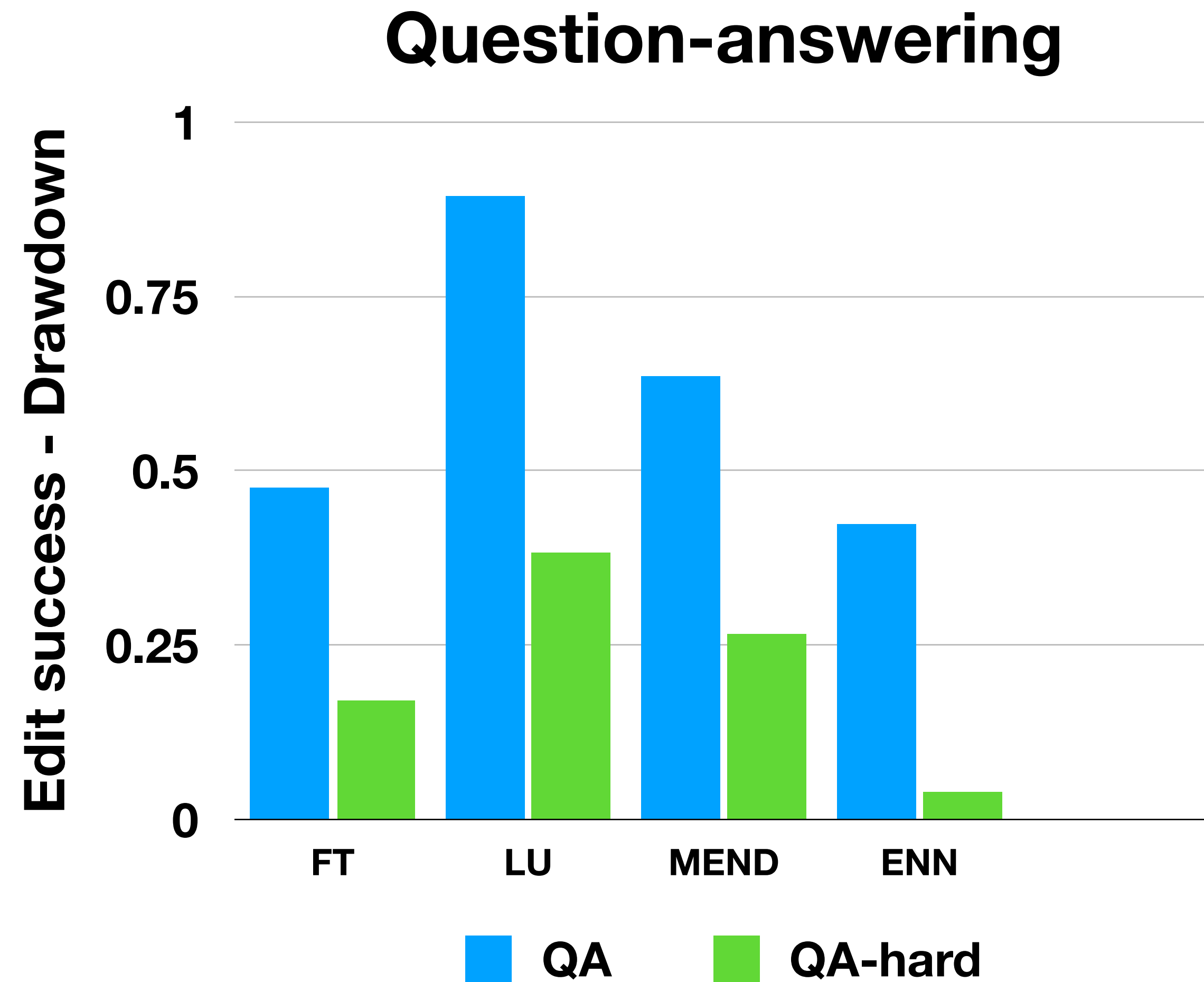


Figure reproduced from:
Memory-based model editing at scale. Mitchell et al. Preprint;
under review.

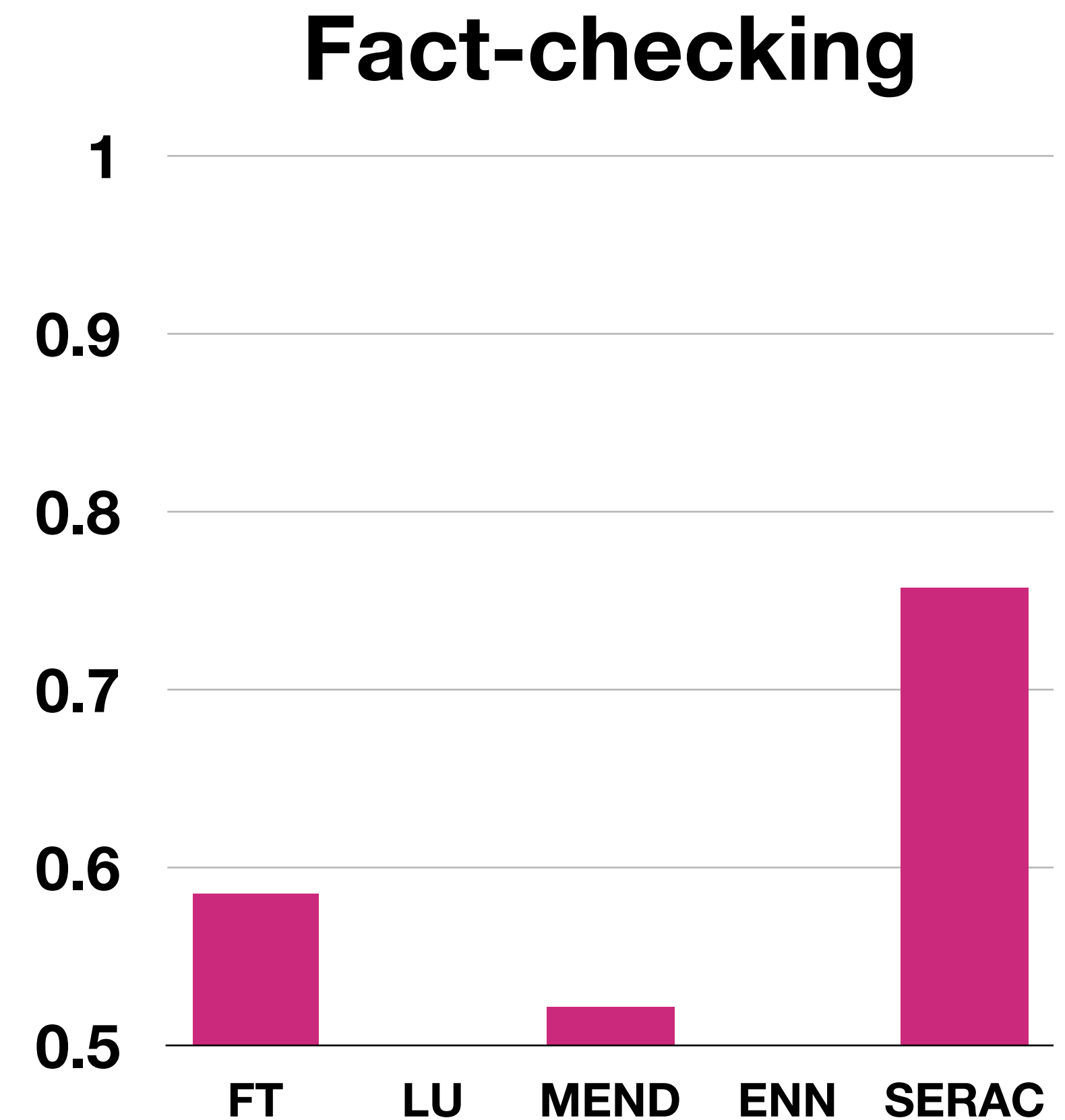
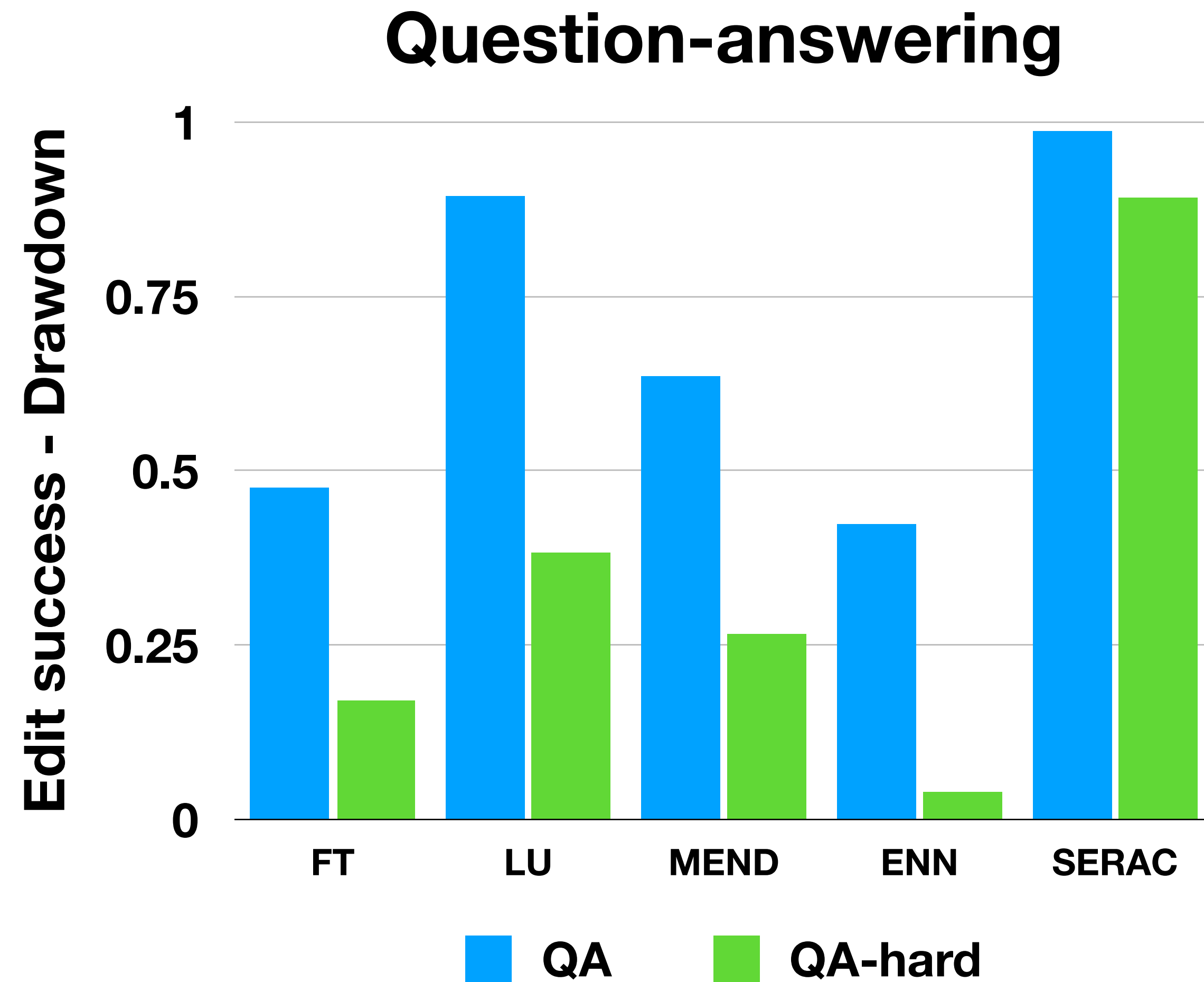
More challenging benchmarks

Multiple edits, more difficult edit scopes



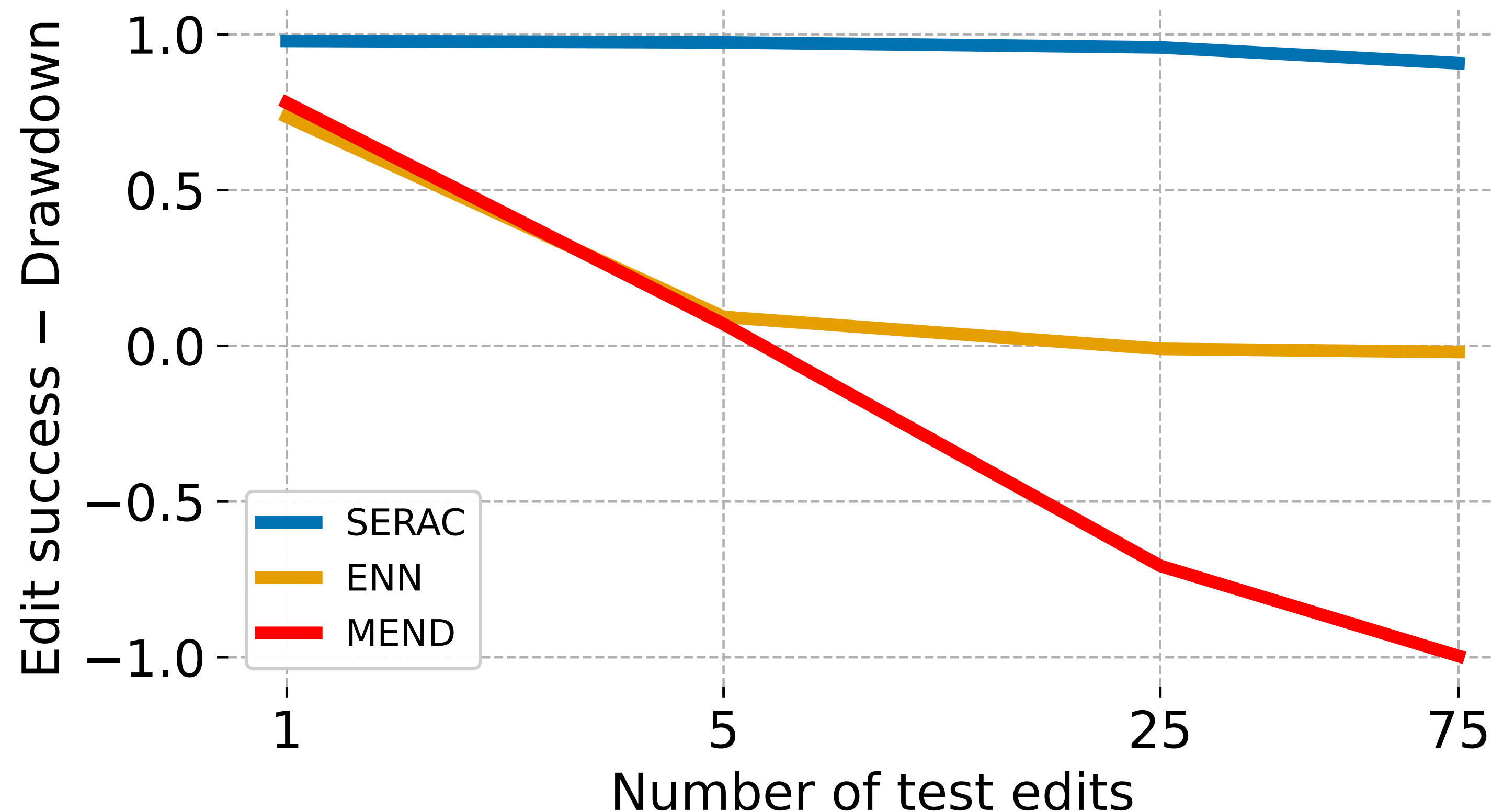
More challenging benchmarks

Multiple edits, more difficult edit scopes



More challenging benchmarks

A case study in handling many QA edits



Semi-parametric editor exhibits less interference within a batch of edits

Edits without parameter updates

Semi-parametric Editng with a Retrieval-Augmented Counterfactual model

Start with the **frozen** base model

1. Store edits in an explicit **memory**
2. Train a **scope classifier** to retrieve relevant edits as needed
3. Train a **counterfactual model** to reason over retrieved edits as needed

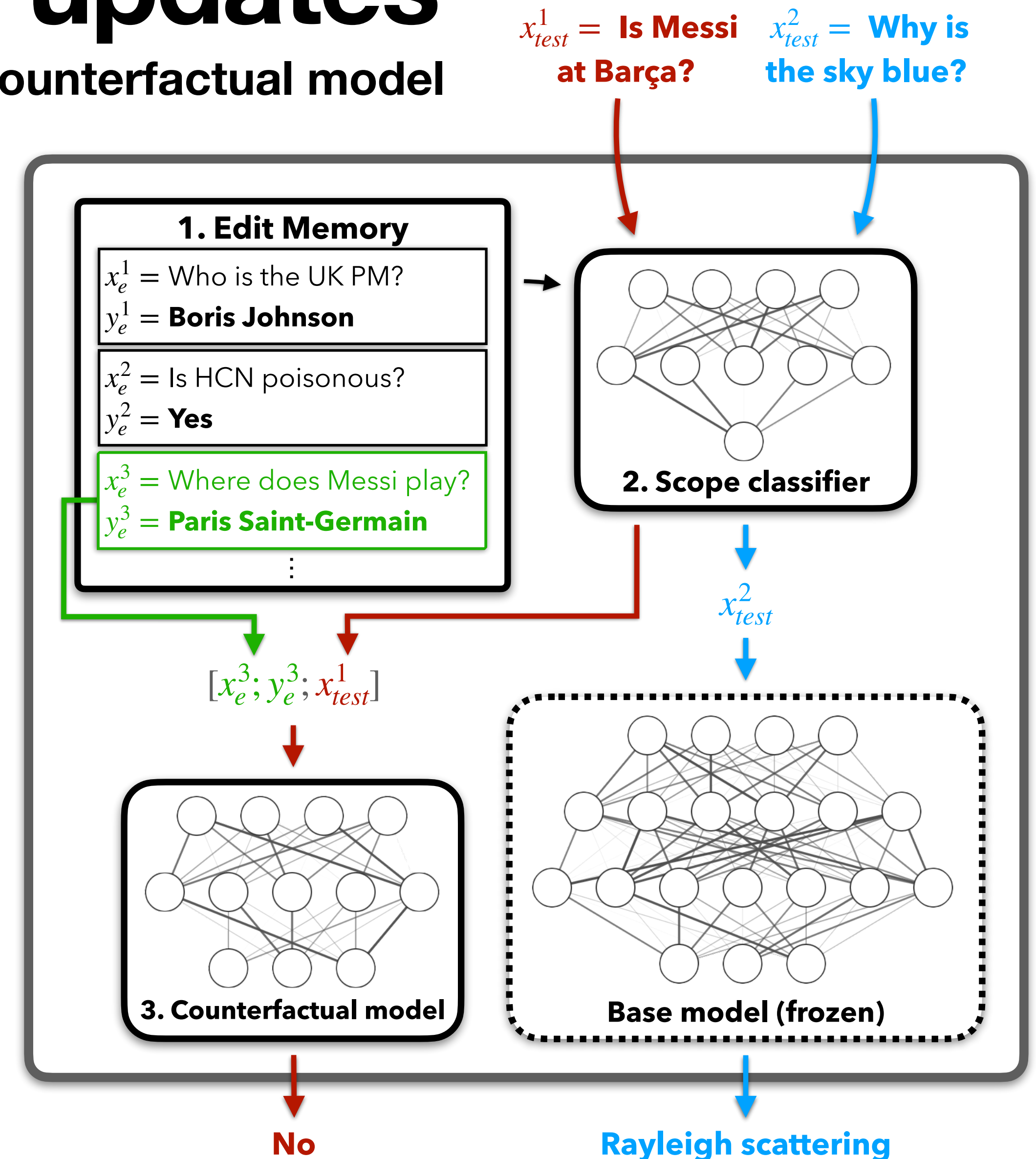


Figure reproduced from:
Memory-based model editing at scale. Mitchell et al. Preprint;
under review.

Edits without parameter updates

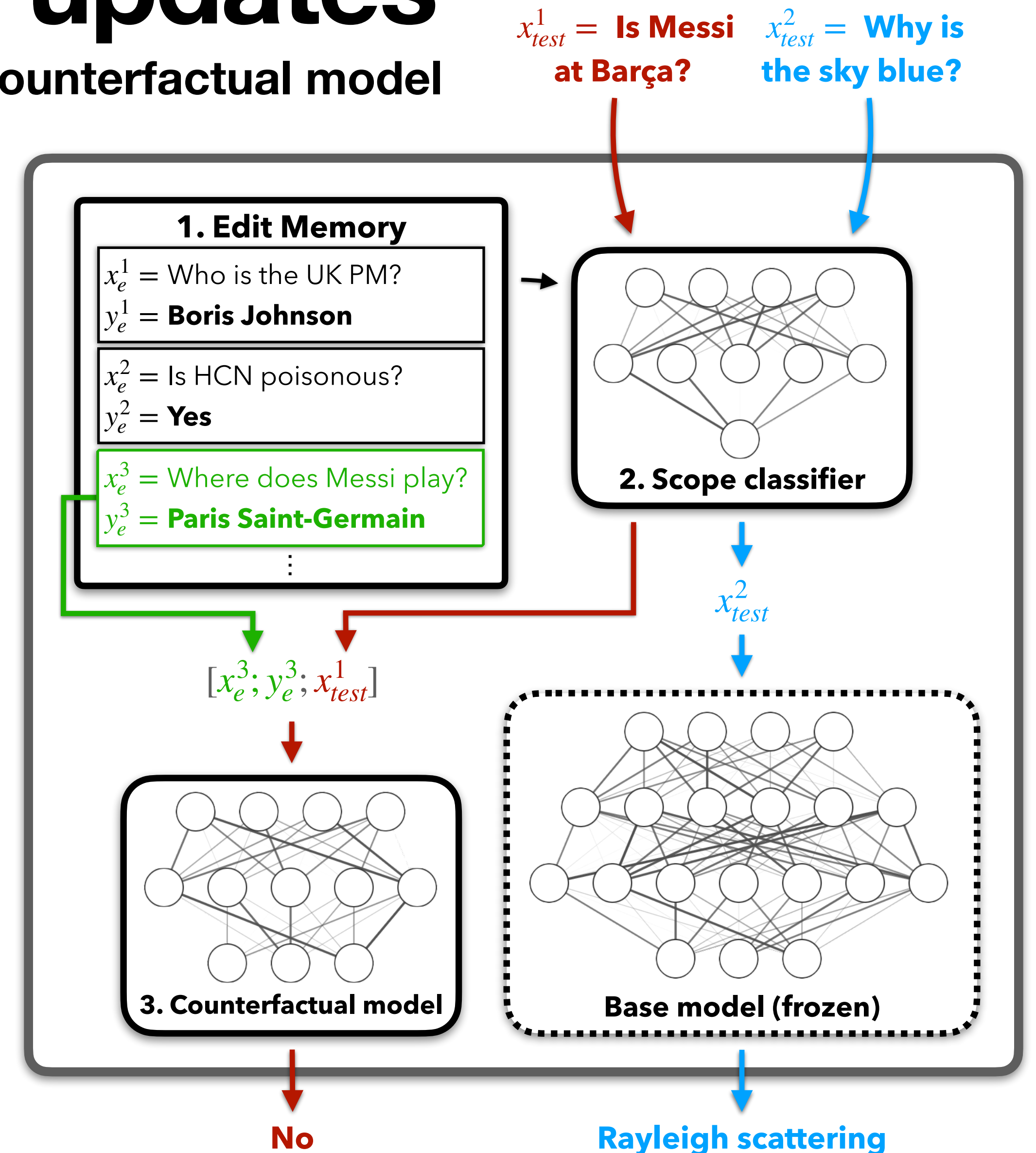
Semi-parametric Editng with a Retrieval-Augmented Counterfactual model

Start with the frozen base model

1. Store edits in an explicit memory
2. Train a scope classifier to retrieve relevant edits as needed
3. Train a counterfactual model to reason over retrieved edits as needed

Decouple editor & base model!

Figure reproduced from:
Memory-based model editing at scale. Mitchell et al. Preprint;
under review.



Today's Plan

I. Background

II. Learning to edit NNs

III. Moving editing towards the real world

IV. Future work & open questions

Where do we go from here?

Open questions

- Editing without a dataset?

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors

Step 1: Figure out which parameters correspond to a given fact

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors

Step 1: Figure out which parameters correspond to a given fact

Step 2: Update **only** those parameters

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors

Step 1: Figure out which parameters correspond to a given fact

Step 2: Update **only** those parameters

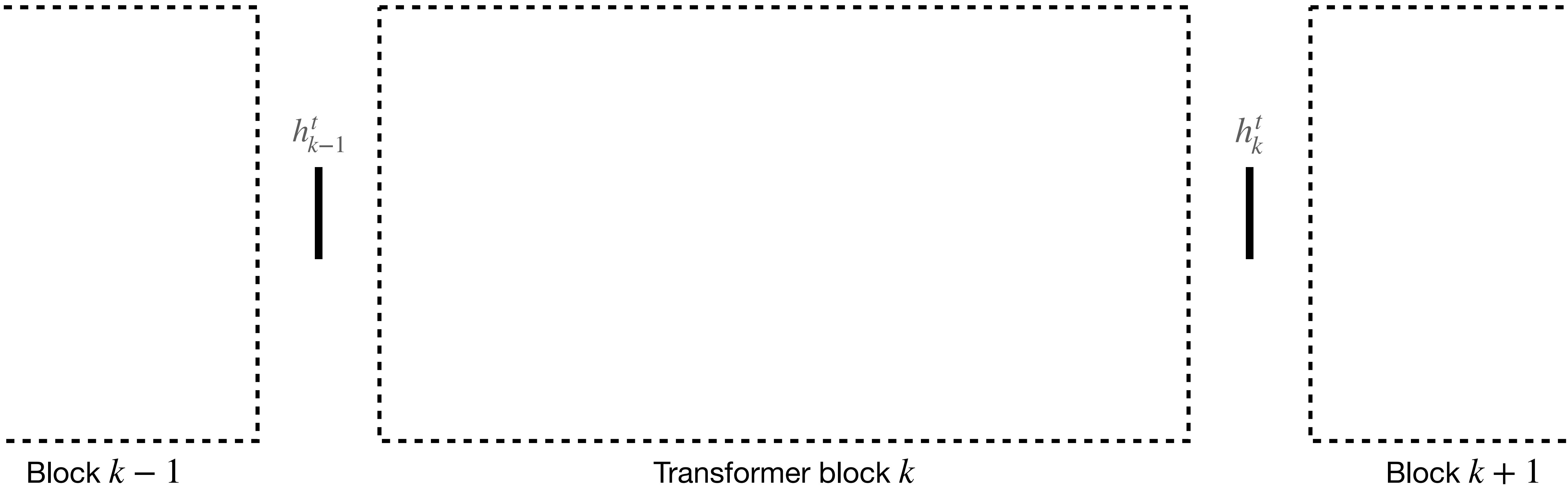
(Both steps are performed with fixed algorithms/heuristics; no learning!)

Editing through attribution

Transformer Feed-Forward Layers are Key-Value Memories. Geva et al. 2020.

Interpreting fully-connected layers as key-value memories

Fully-connected layers are “key-value memories”

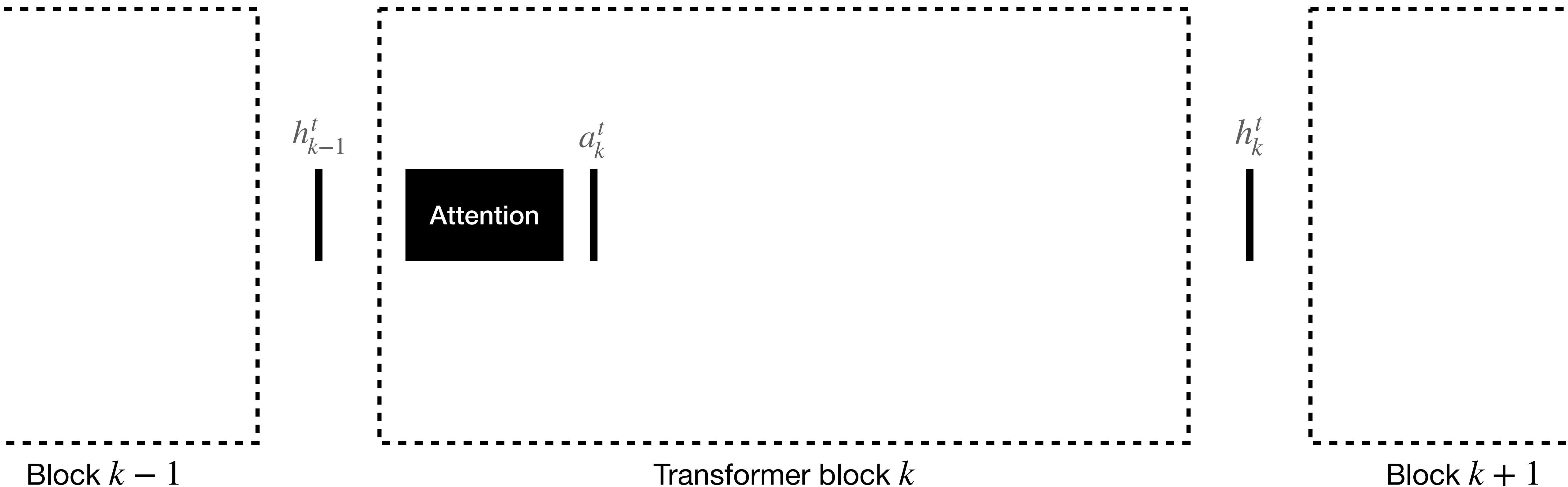


Editing through attribution

Transformer Feed-Forward Layers are Key-Value Memories. Geva et al. 2020.

Interpreting fully-connected layers as key-value memories

Fully-connected layers are “key-value memories”

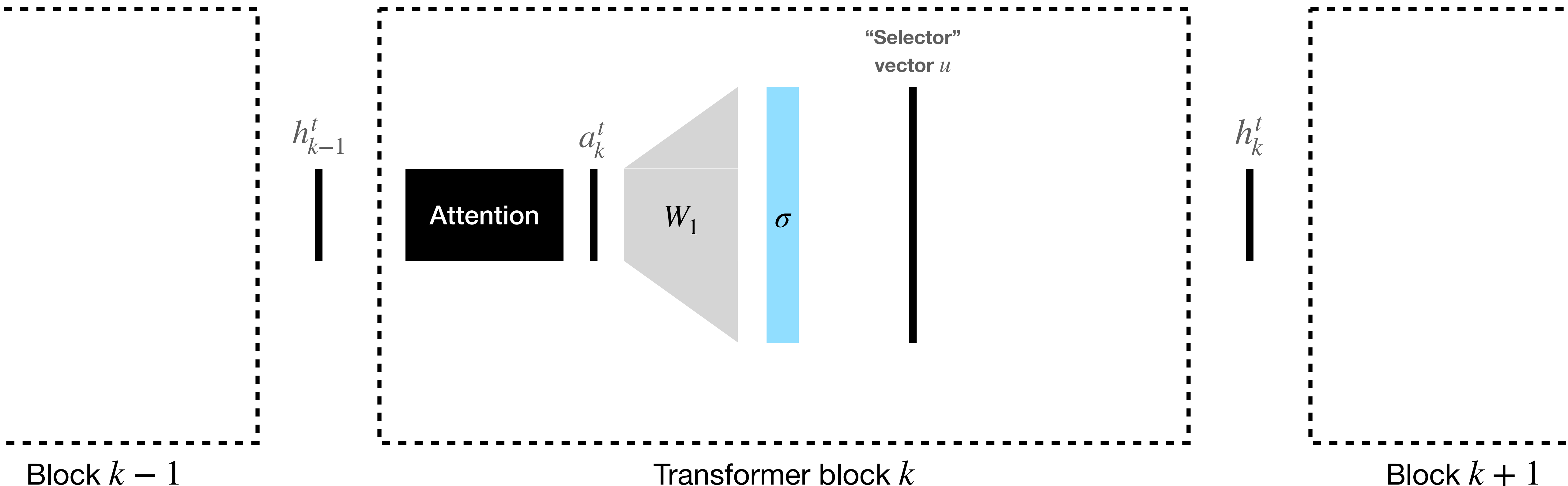


Editing through attribution

Transformer Feed-Forward Layers are Key-Value Memories. Geva et al. 2020.

Interpreting fully-connected layers as key-value memories

Fully-connected layers are “key-value memories”

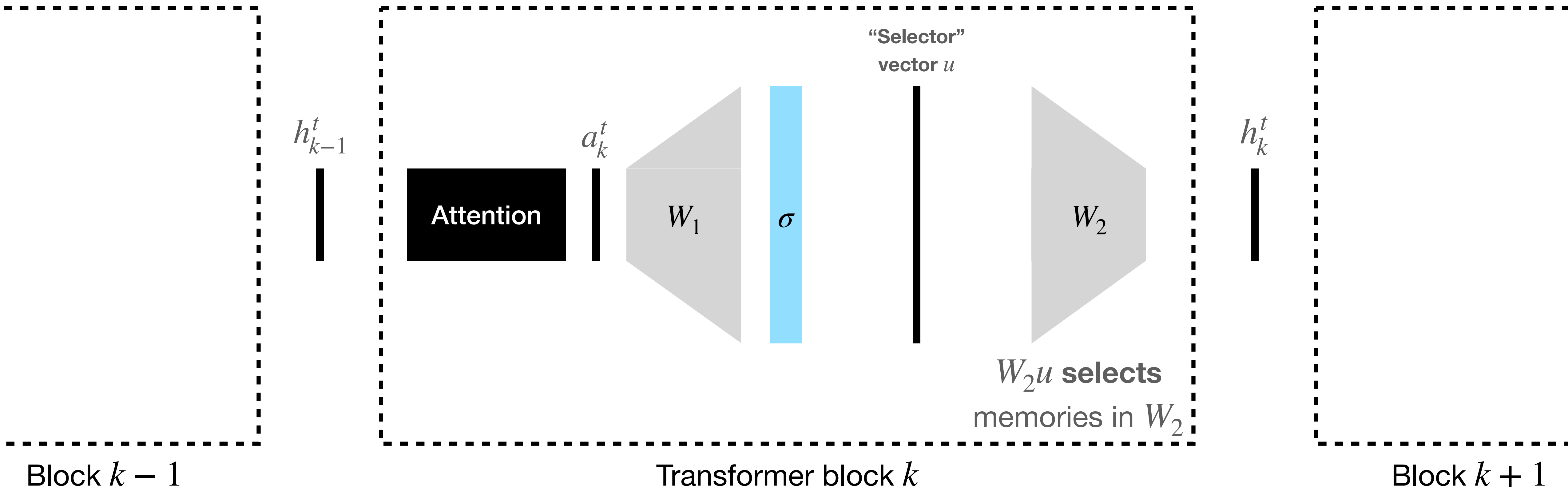


Editing through attribution

Transformer Feed-Forward Layers are Key-Value Memories. Geva et al. 2020.

Interpreting fully-connected layers as key-value memories

Fully-connected layers are “key-value memories”



Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?

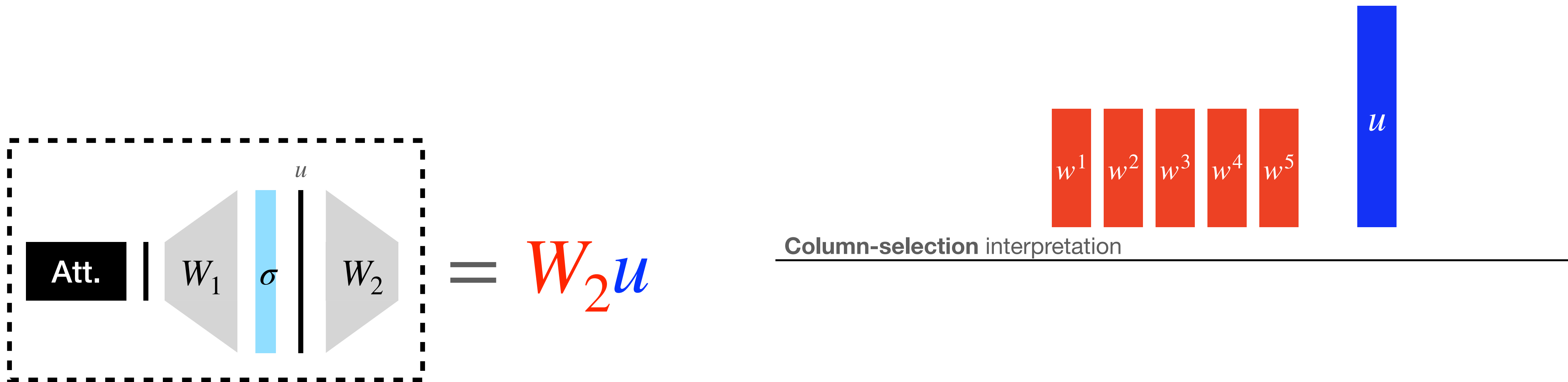


**I'm ignoring skip connections and normalization here

Editing through attribution

Interpreting fully-connected layers as key-value memories

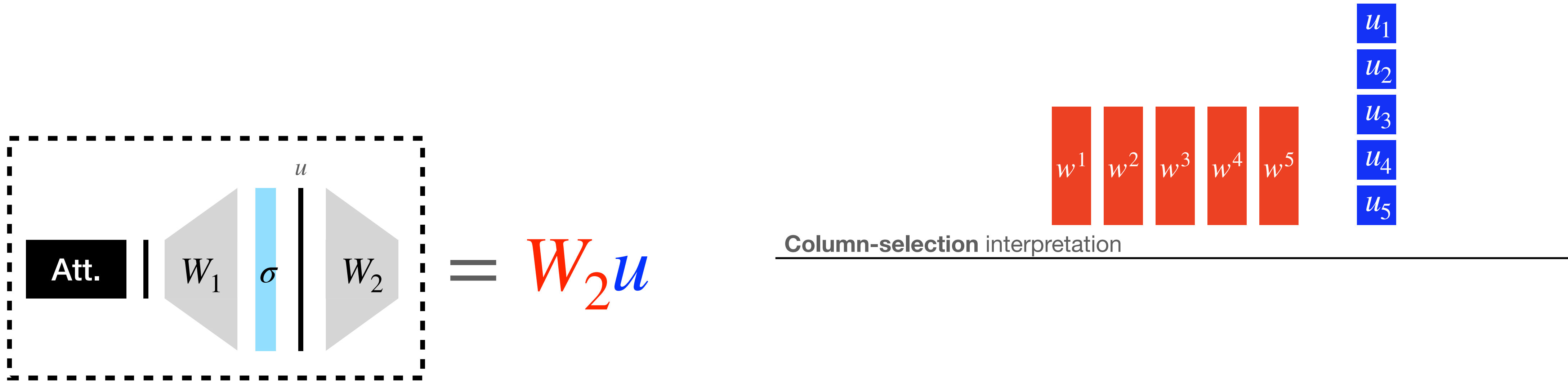
How do we “select” memories in a weight matrix?



Editing through attribution

Interpreting fully-connected layers as key-value memories

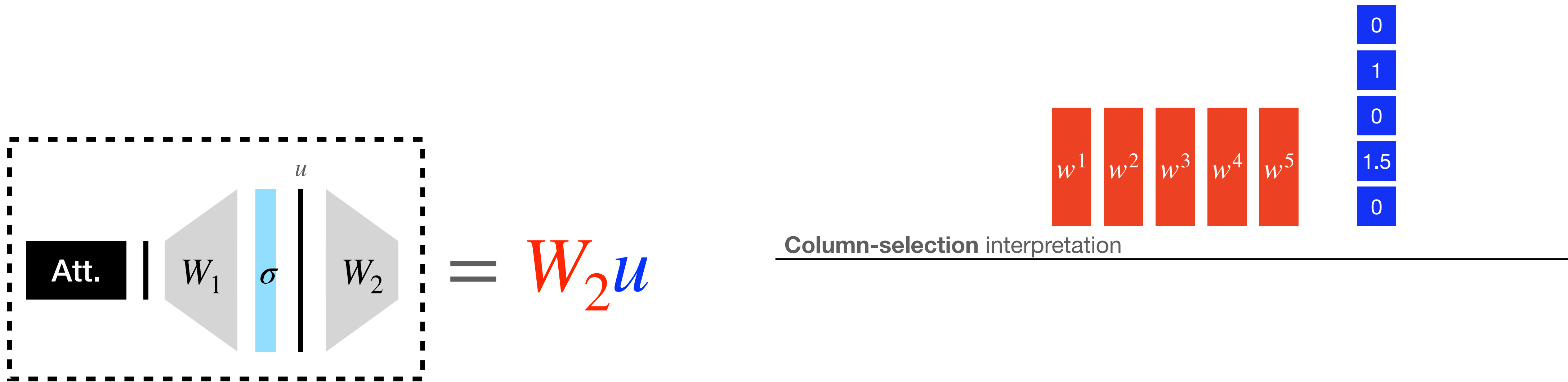
How do we “select” memories in a weight matrix?



Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?



Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?

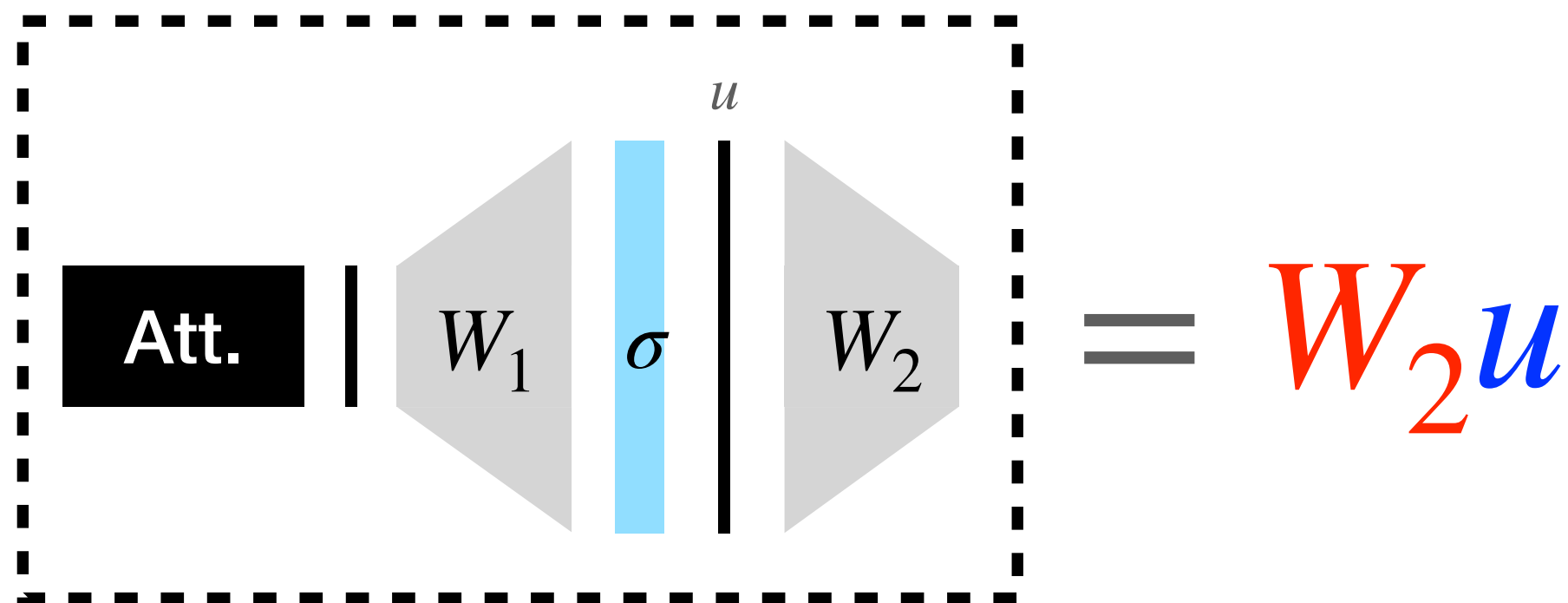


Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?

Sum columns w^i weighted by elements of u



$$W_2u = \begin{matrix} \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \end{matrix} \\ \begin{matrix} w^1 \\ w^2 \\ w^3 \\ w^4 \\ w^5 \end{matrix} \end{matrix} = \sum_i w^i u_i$$

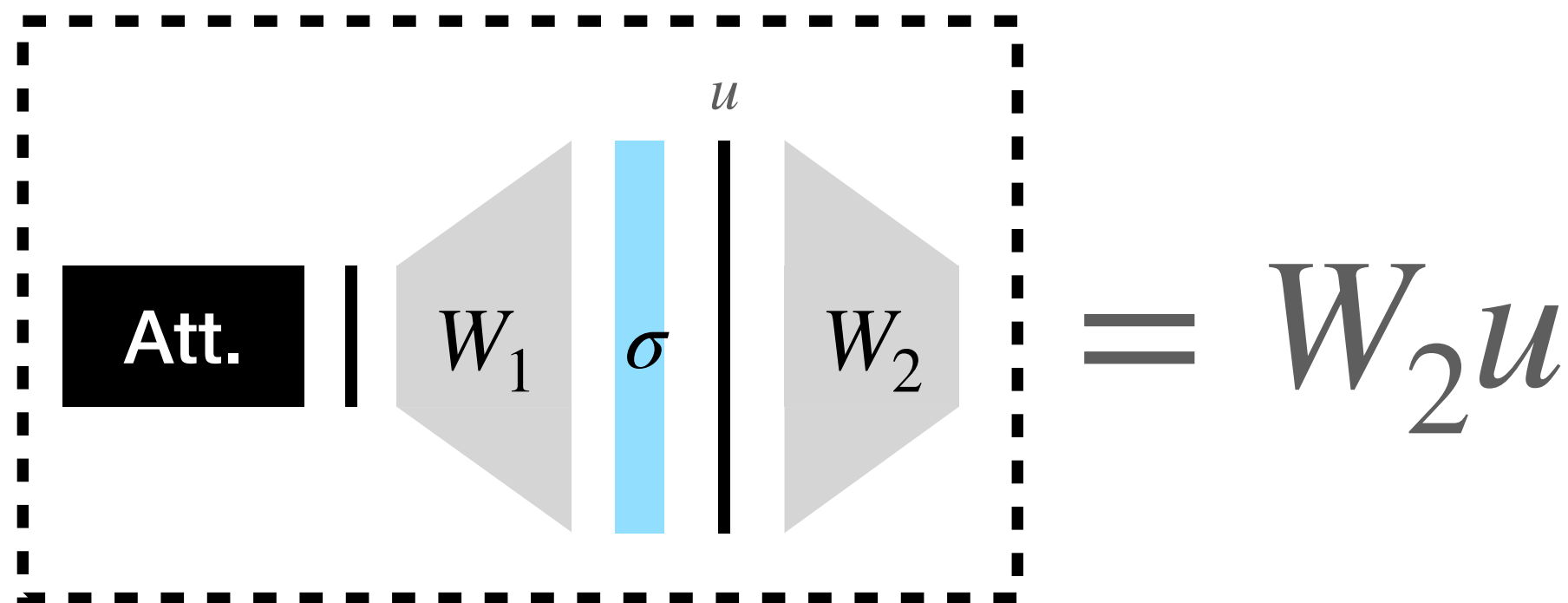
Column-selection interpretation

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?

Sum columns w^i weighted by elements of u



$$W_2 u = \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \\ w^1 & w^2 & w^3 & w^4 & w^5 \end{matrix} = \sum_i w^i u_i$$

Column-selection interpretation
Key-value lookup interpretation

$$W_2 u = \left(\sum_i v_i k_i^T \right) u$$

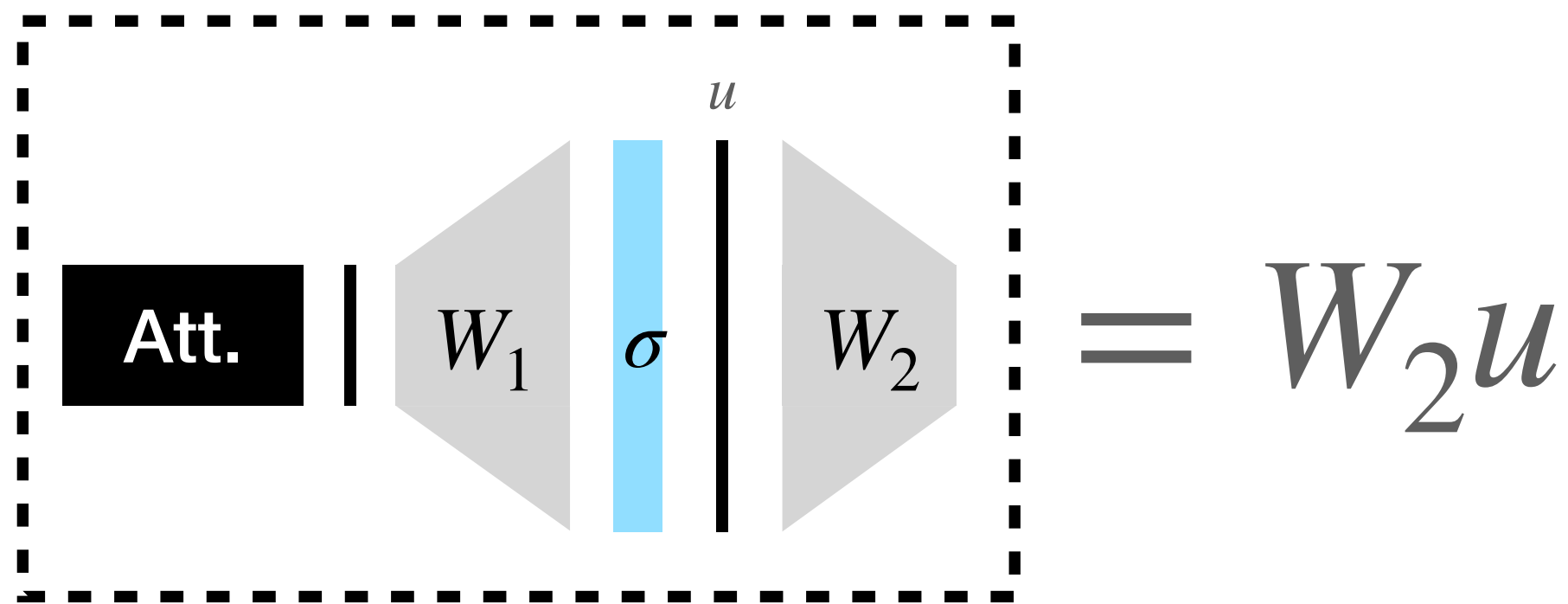
W_2 = sum of outer products of keys k_i and values v_i

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?

Sum columns w^i weighted by elements of u



$$W_2u = \begin{matrix} \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \end{matrix} \\ \begin{matrix} w^1 \\ w^2 \\ w^3 \\ w^4 \\ w^5 \end{matrix} \end{matrix} = \sum_i w^i u_i$$

Column-selection interpretation
Key-value lookup interpretation

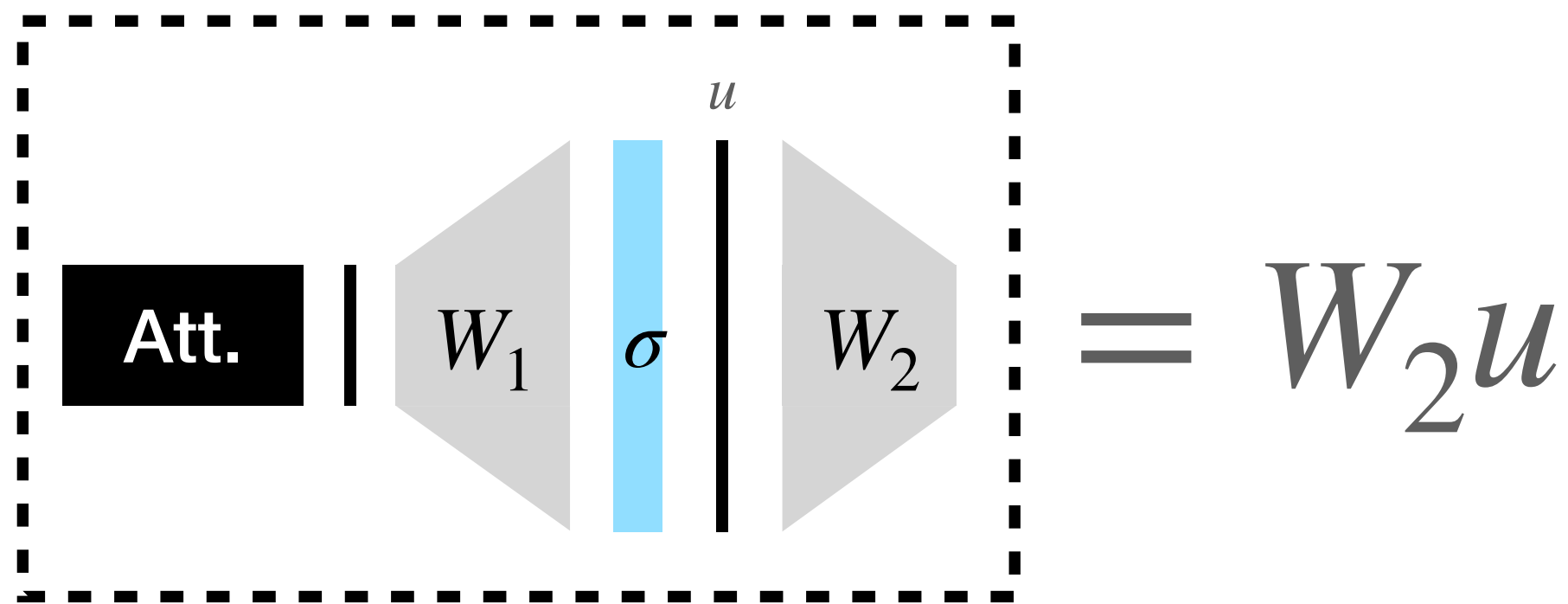
$$W_2u = \left(\sum_i v_i k_i^T \right) u = \sum_i (v_i k_i^T) u$$

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?

Sum columns w^i weighted by elements of u



$$W_2 u = \begin{matrix} \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \end{matrix} \\ \begin{matrix} w^1 \\ w^2 \\ w^3 \\ w^4 \\ w^5 \end{matrix} \end{matrix} = \sum_i w^i u_i$$

Column-selection interpretation
Key-value lookup interpretation

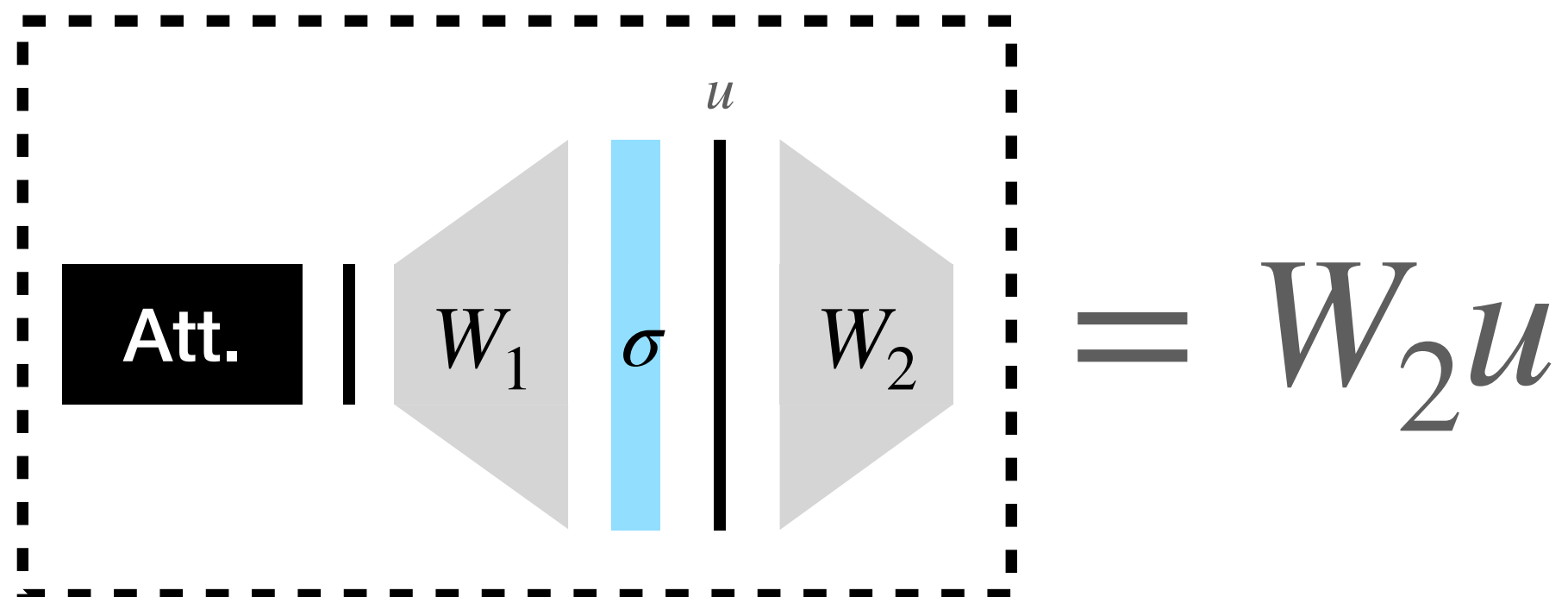
$$W_2 u = \left(\sum_i v_i k_i^T \right) u = \sum_i (v_i k_i^T) u = \sum_i v_i (k_i^T u)$$

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?

Sum columns w^i weighted by elements of u



$$W_2 u = \begin{matrix} \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \end{matrix} \\ \begin{matrix} w^1 \\ w^2 \\ w^3 \\ w^4 \\ w^5 \end{matrix} \end{matrix} = \sum_i w^i u_i$$

Column-selection interpretation
Key-value lookup interpretation

$$W_2 u = \left(\sum_i v_i k_i^T \right) u = \sum_i (v_i k_i^T) u = \sum_i v_i (k_i^T u)$$

Sum values v_i weighted by $k_i^T u$

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “select” memories in a weight matrix?

Sum columns w^i weighted by elements of u

$$W_2 u = \begin{array}{c} \begin{array}{|c|c|c|c|c|} \hline 0 & 1 & 0 & 1.5 & 0 \\ \hline \end{array} \\ \begin{array}{|c|c|c|c|c|} \hline w^1 & w^2 & w^3 & w^4 & w^5 \\ \hline \end{array} \end{array} = \sum_i w^i u_i$$

Column-selection interpretation
Key-value lookup interpretation

$$W_2 u = \left(\sum_i v_i k_i^T \right) u = \sum_i (v_i k_i^T) u = \sum_i v_i (k_i^T u)$$

Sum values v_i
weighted by $k_i^T u$

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “edit” memories in a weight matrix?

Sum columns w^i weighted by elements of u

$$W_2 u = \begin{matrix} \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \end{matrix} \\ \begin{matrix} w^1 & w^2 & w^3 & w^4 & w^5 \end{matrix} \end{matrix} = \sum_i w^i u_i$$

Column-selection interpretation
Key-value lookup interpretation

$$W_2 u = \left(\sum_i v_i k_i^T \right) u = \sum_i (v_i k_i^T) u = \sum_i v_i (k_i^T u)$$

Sum values v_i
weighted by $k_i^T u$

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “edit” memories in a weight matrix?

Edit interpretation 1:

Update just one column in W_2

$$W_2 u = \begin{matrix} \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \end{matrix} \\ \begin{matrix} w^1 & w^2 & w^3 & w^4 & w^5 \end{matrix} \end{matrix} = \sum_i w^i u_i$$

Sum columns w^i weighted by elements of u

Column-selection interpretation
Key-value lookup interpretation

$$W_2 u = \left(\sum_i v_i k_i^T \right) u = \sum_i (v_i k_i^T) u = \sum_i v_i (k_i^T u)$$

Sum values v_i weighted by $k_i^T u$

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “edit” memories in a weight matrix?

Edit interpretation 1:

Update just one column in W_2

$$W_2 u = \begin{matrix} \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \end{matrix} \\ \begin{matrix} w^1 & w^2 & w^3 & w^4 & w^5 \end{matrix} \end{matrix} = \sum_i w^i u_i$$

Sum columns w^i weighted by elements of u

Column-selection interpretation
Key-value lookup interpretation

Edit interpretation 2:

Update just one key-value pair in W_2

$$W_2 u = \left(\sum_i v_i k_i^T \right) u = \sum_i (v_i k_i^T) u = \sum_i v_i (k_i^T u)$$

Sum values v_i weighted by $k_i^T u$

Editing through attribution

Interpreting fully-connected layers as key-value memories

How do we “edit” memories in a weight matrix?

Edit interpretation 1:

Update just one **column** in W_2

Knowledge Neurons in Pre-trained Transformers. Dai et al. 2021.

$$W_2 u = \begin{matrix} \begin{matrix} 0 & 1 & 0 & 1.5 & 0 \end{matrix} \\ \begin{matrix} w^1 & w^2 & w^3 & w^4 & w^5 \end{matrix} \end{matrix} = \sum_i w^i u_i$$

Sum columns w^i weighted by elements of u

Column-selection interpretation
Key-value lookup interpretation

Edit interpretation 2:

Update just one **key-value pair** in W_2

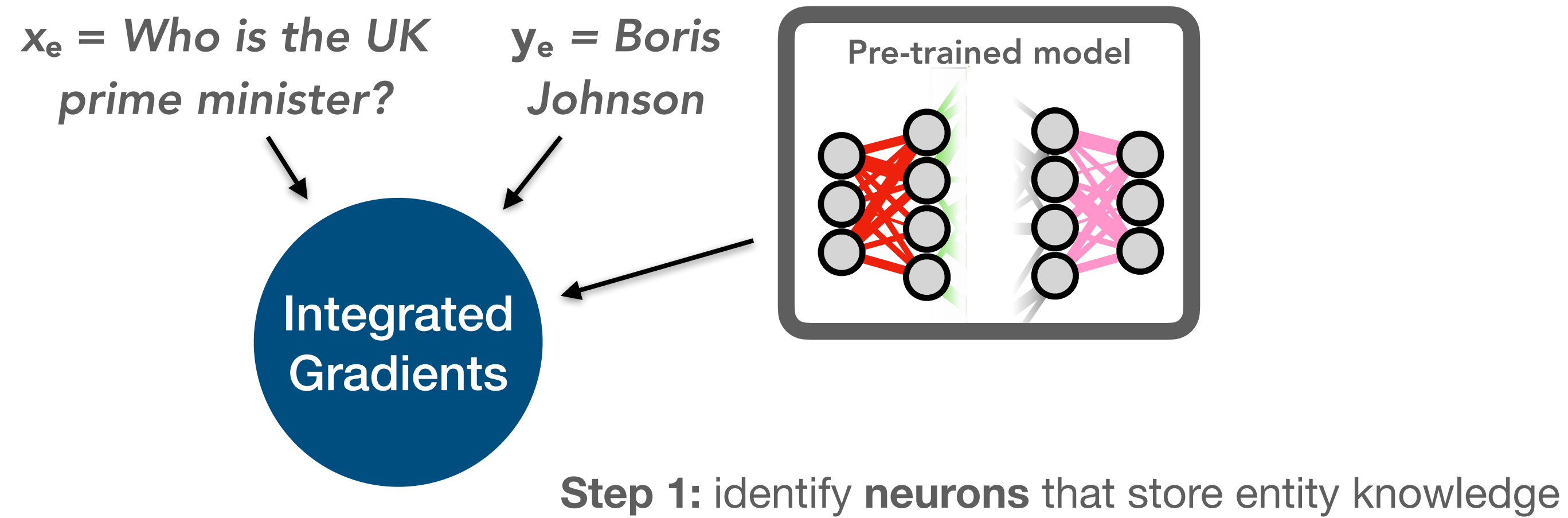
Locating and Editing Factual Knowledge in GPT. Meng et al. 2022.

$$W_2 u = \left(\sum_i v_i k_i^T \right) u = \sum_i (v_i k_i^T) u = \sum_i v_i (k_i^T u)$$

Sum values v_i weighted by $k_i^T u$

Editing through attribution

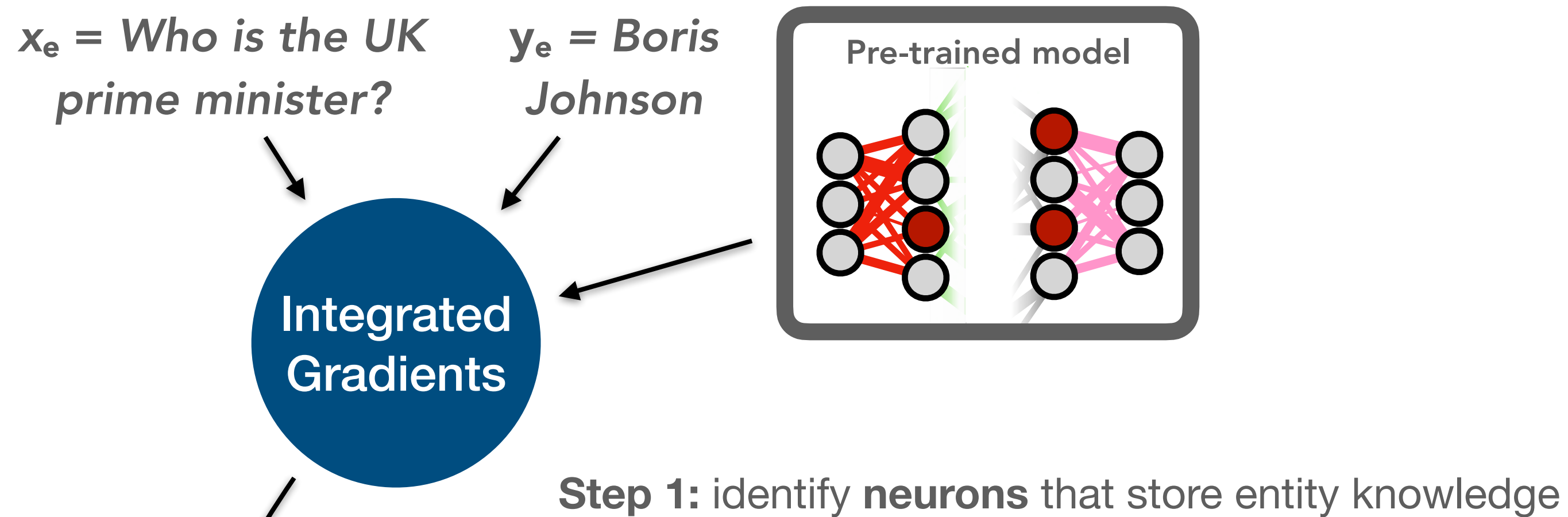
Interpretation 1: “Knowledge Neurons”



¹ Knowledge Neurons in Pretrained Transformers. Dai et al. 2021.

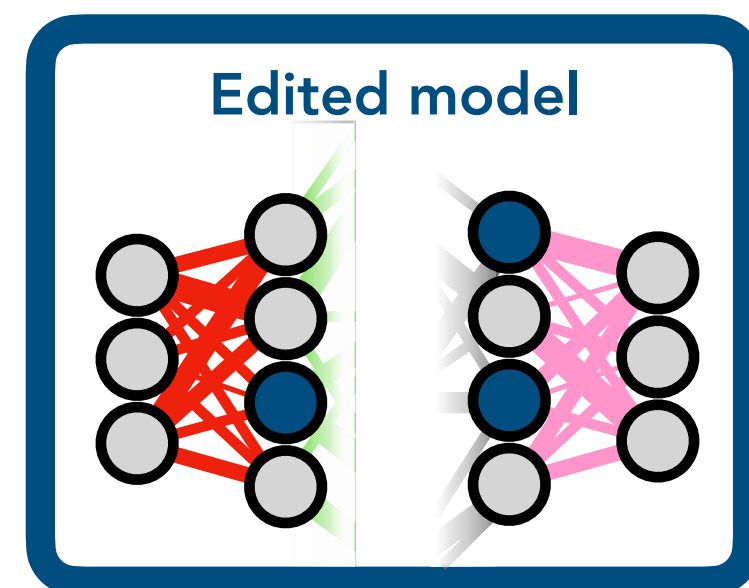
Editing through attribution

Interpretation 1: “Knowledge Neurons”



Step 2: insert new memory at identified location

$\{(1, 2),$
 $(2, 0),$
 $(2, 3)\}$



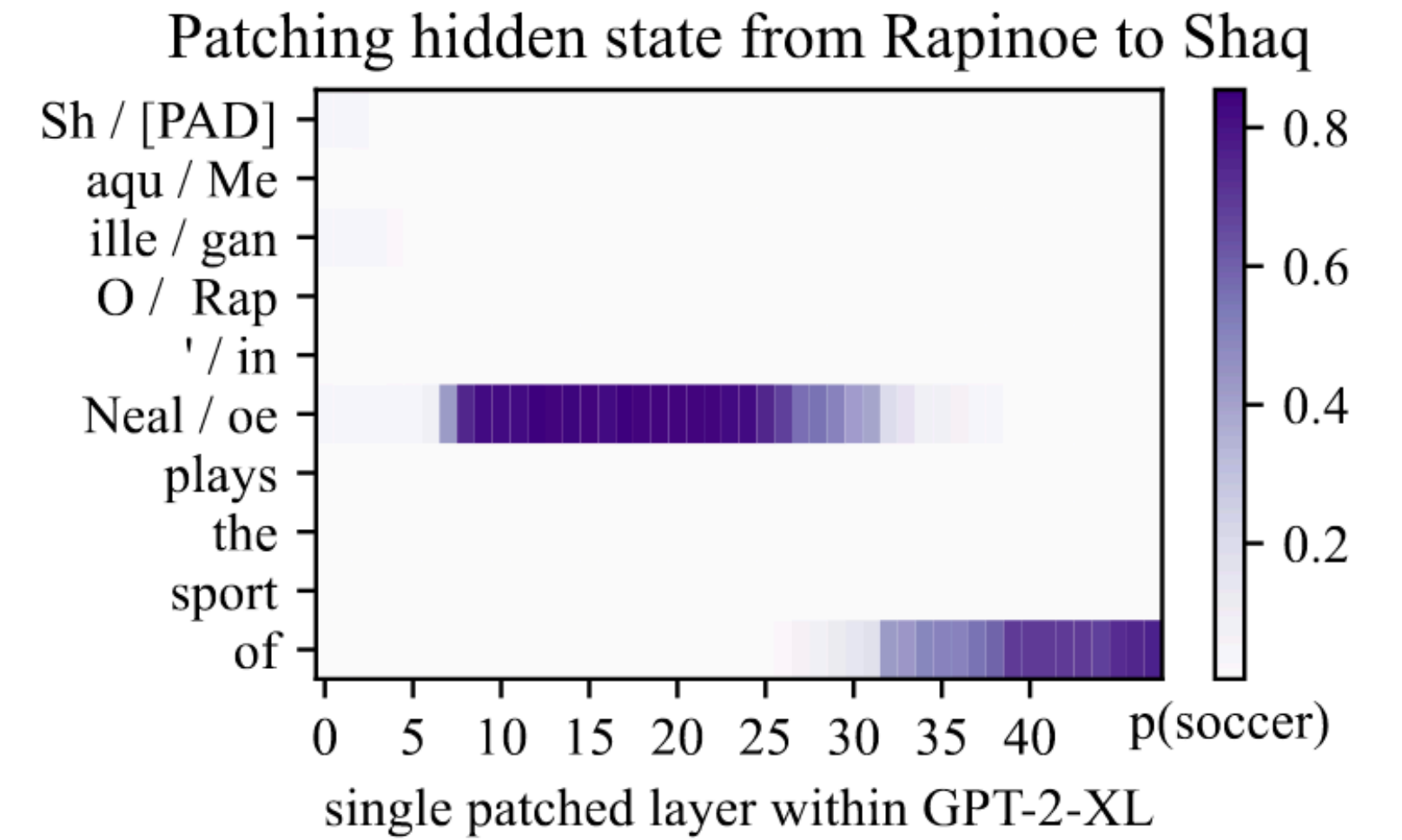
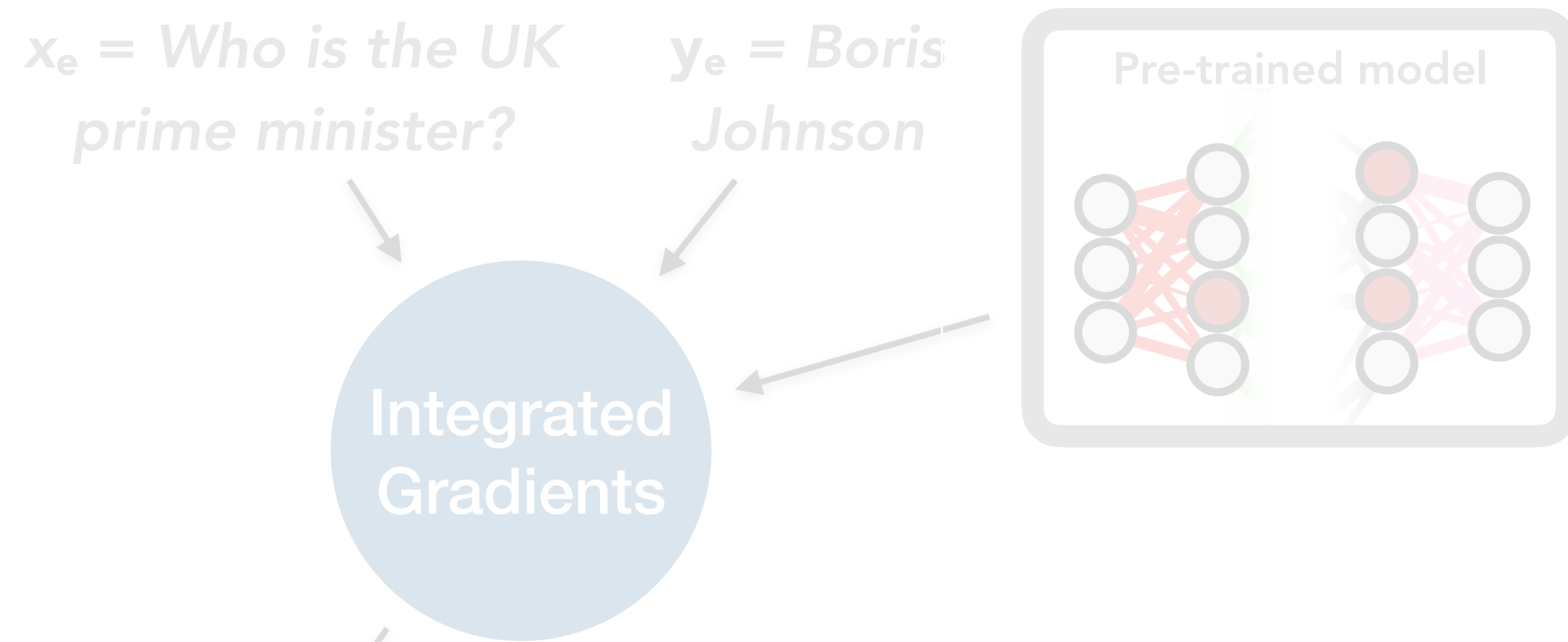
Replace columns in weight matrices with word embedding of desired word

¹ Knowledge Neurons in Pretrained Transformers. Dai et al. 2021.

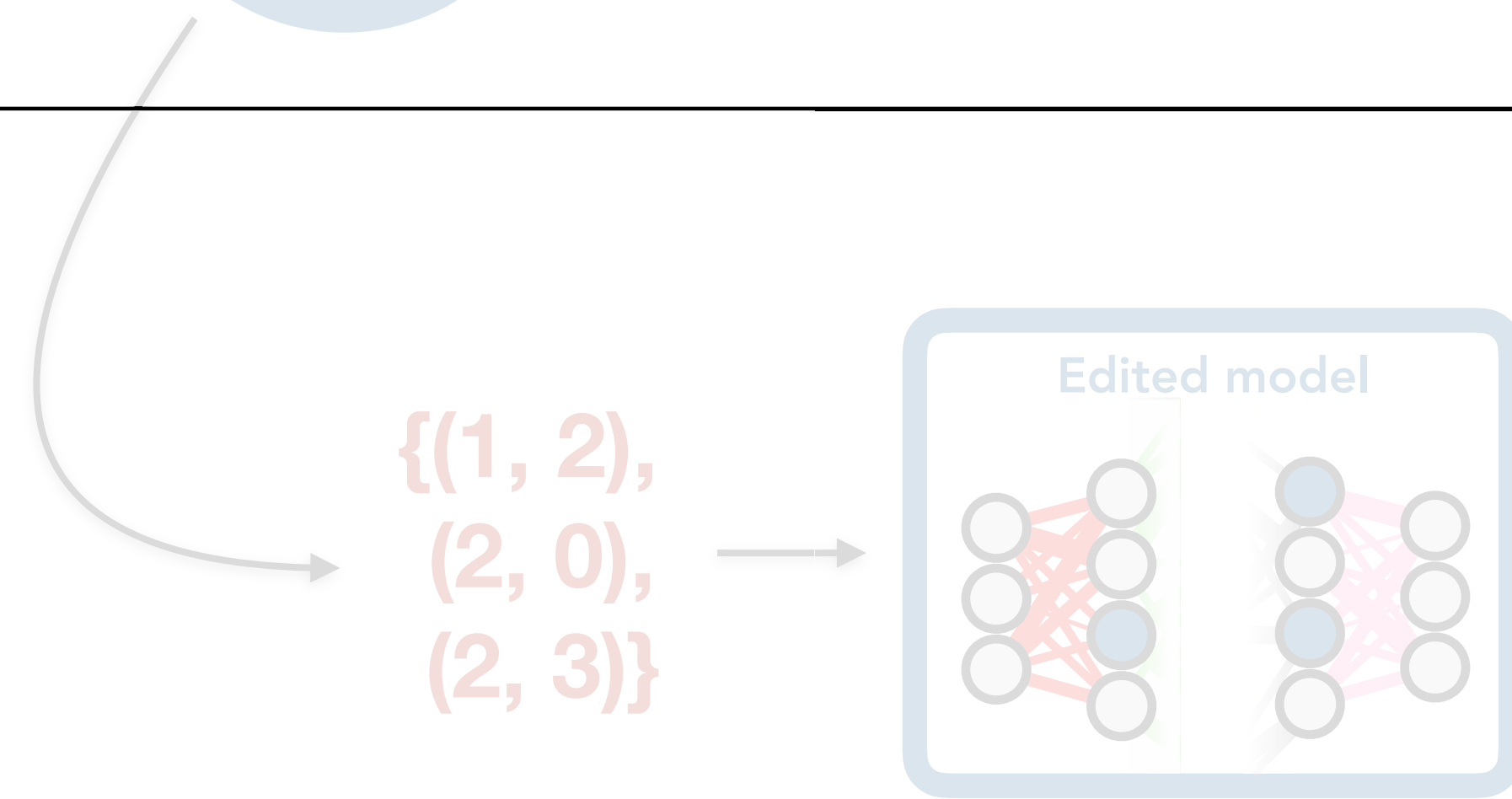
Editing through attribution

Interpretation 2: Rank-1 Model Editing

² *Moving the Eiffel Tower to ROME: Tracing and Editing Facts in GPT.* Anonymous. 2021.



Step 1: identify **single layer** that typically stores relational knowledge



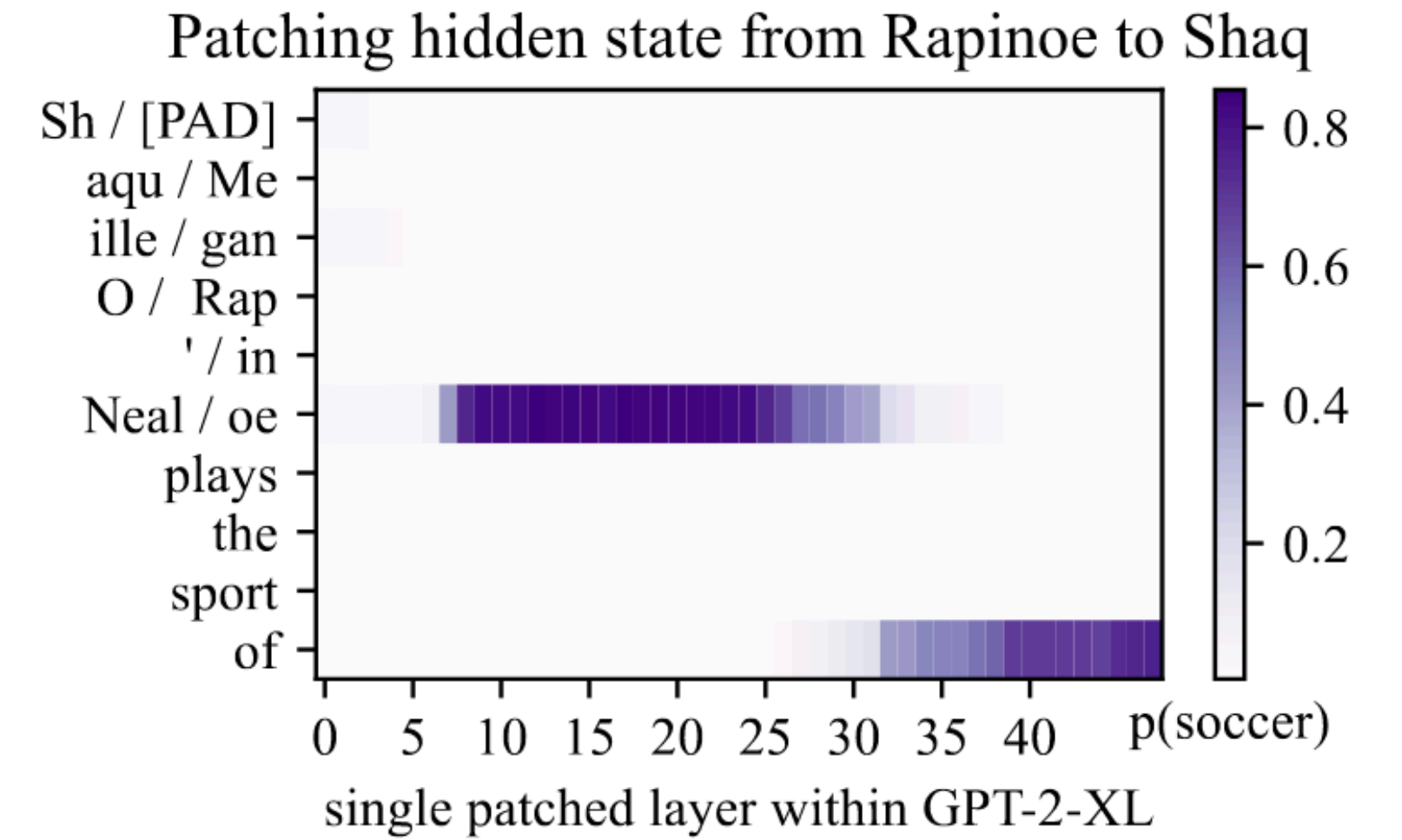
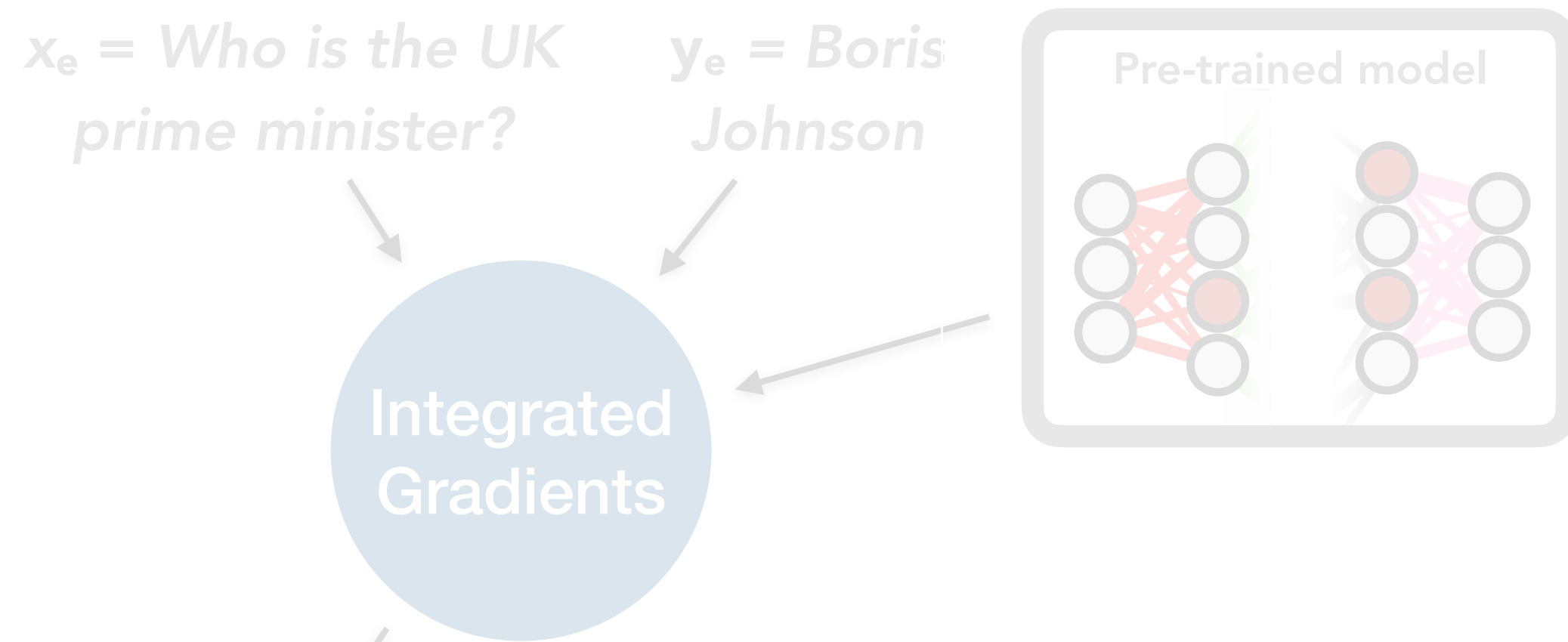
Replace columns in weight matrices with word embedding of desired word

¹ *Knowledge Neurons in Pretrained Transformers.* Dai et al. 2021.

Editing through attribution

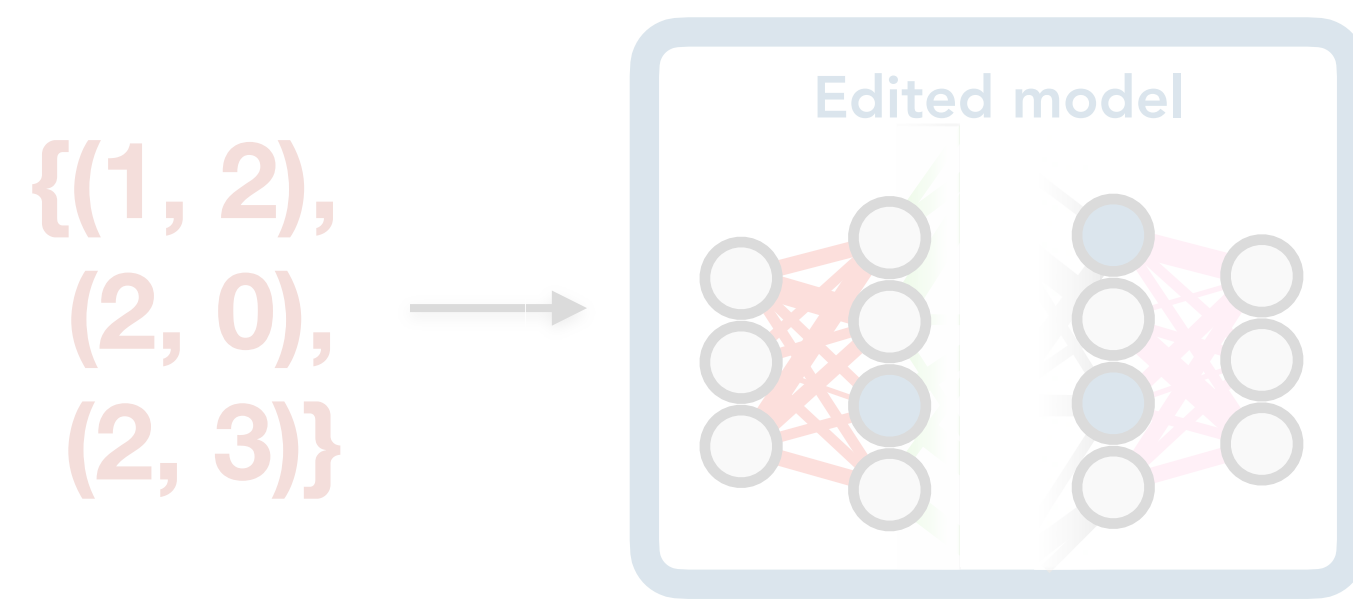
Interpretation 2: Rank-1 Model Editing

² Moving the Eiffel Tower to ROME: Tracing and Editing Facts in GPT. Anonymous. 2021.

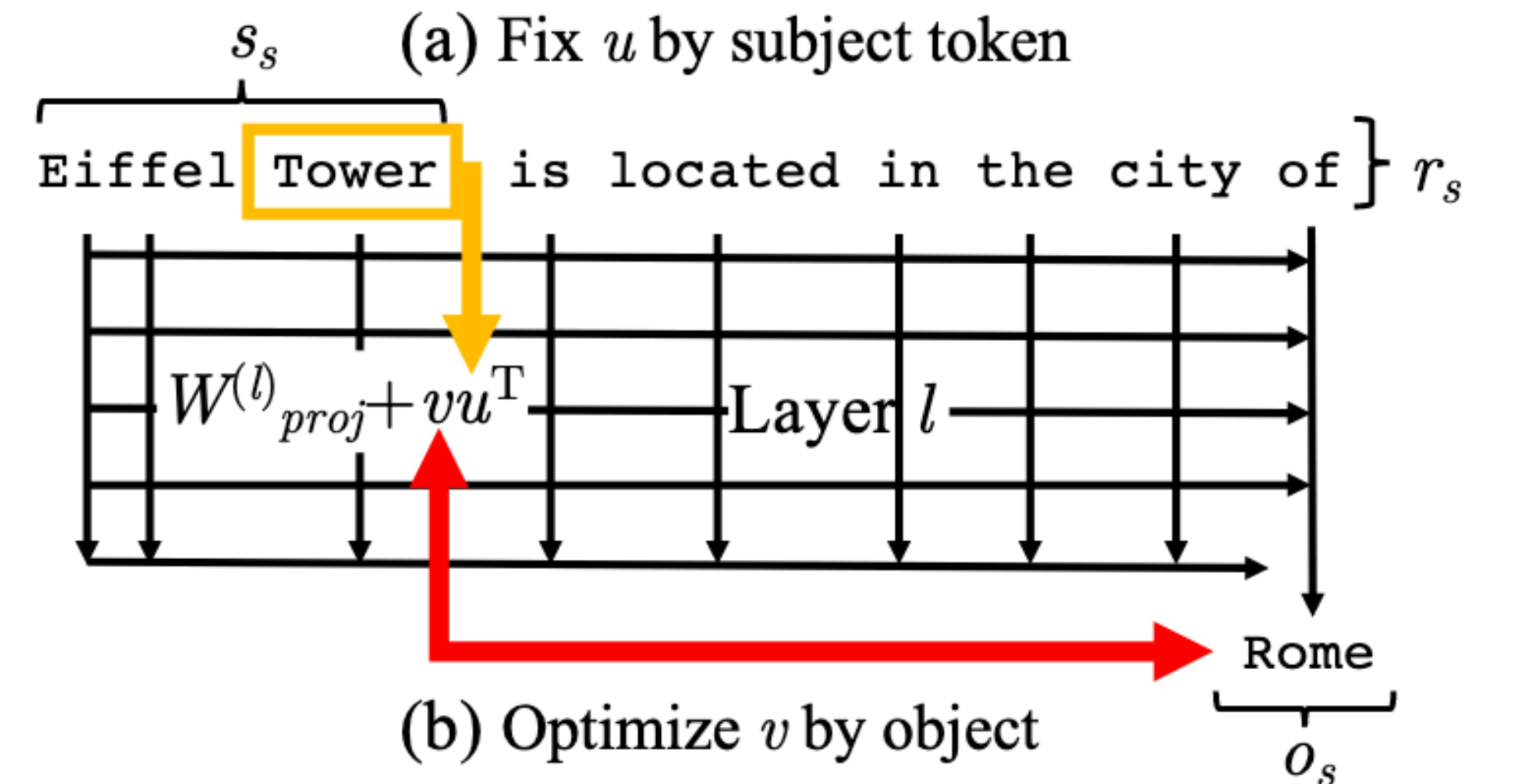


Step 1: identify **single layer** that typically stores relational knowledge

Step 2: insert new **key-value memory** at identified location



Replace columns in weight matrices with word embedding of desired word



¹ Knowledge Neurons in Pretrained Transformers. Dai et al. 2021.

Editing through attribution

An alternative to learning to edit

Learning-based editors:

- + Can learn very **expressive** edit procedures for **many different problems**
- Require a **dataset of edits** in order to train the editor

Editing through attribution

An alternative to learning to edit

Learning-based editors:

- + Can learn very **expressive** edit procedures for **many different problems**
- Require a **dataset of edits** in order to train the editor

Attribution-based editors:

Editing through attribution

An alternative to learning to edit

Learning-based editors:

- + Can learn very **expressive** edit procedures for **many different problems**
- Require a **dataset of edits** in order to train the editor

Attribution-based editors:

- + **No training dataset** of edits required (an **unlabeled** dataset may be needed)

Editing through attribution

An alternative to learning to edit

Learning-based editors:

- + Can learn very **expressive** edit procedures for **many different problems**
- Require a **dataset of edits** in order to train the editor

Attribution-based editors:

- + **No training dataset** of edits required (an **unlabeled** dataset may be needed)
- Algorithm is coupled to the **type of edit**

Editing through attribution

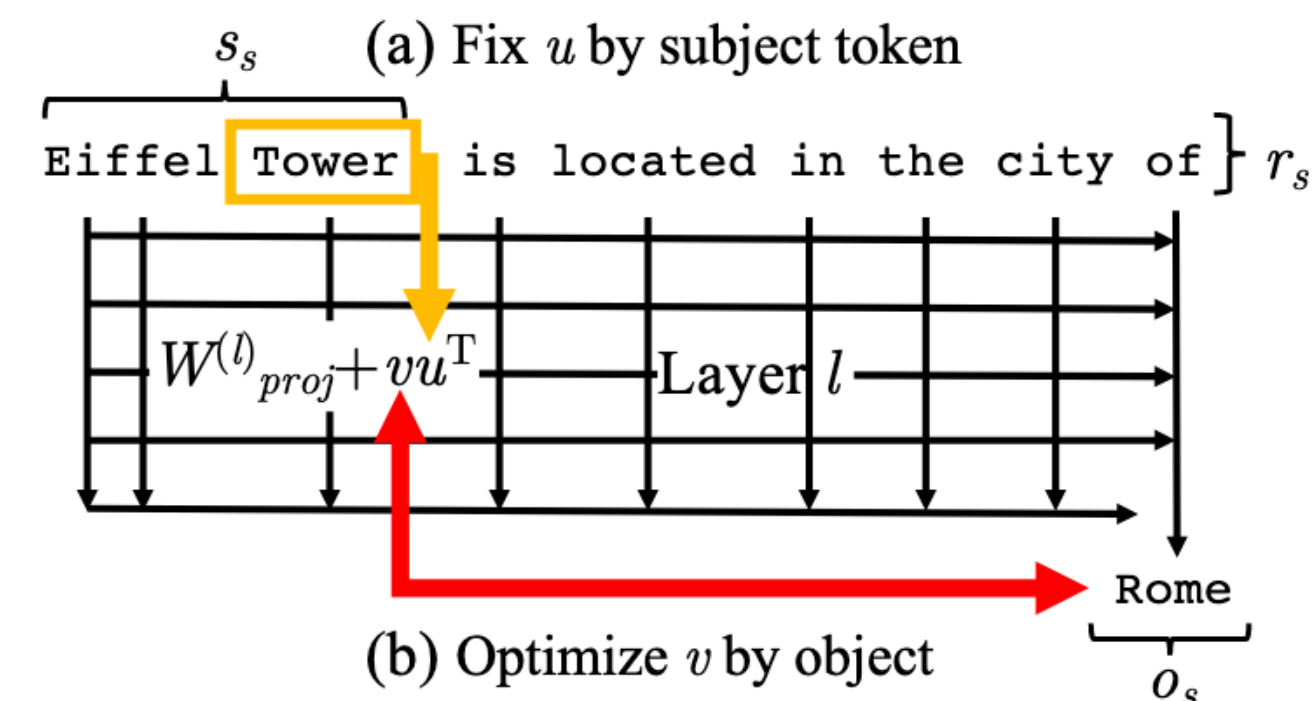
An alternative to learning to edit

Learning-based editors:

- + Can learn very **expressive** edit procedures for **many different problems**
- Require a **dataset of edits** in order to train the editor

Attribution-based editors:

- + **No training dataset** of edits required (an **unlabeled** dataset may be needed)
- Algorithm is coupled to the **type of edit**
- Typically require “richer” edit descriptor



Editing through attribution

An alternative to learning to edit

Learning-based editors:

- + Can learn very **expressive** edit procedures for **many different problems**
- Require a **dataset of edits** in order to train the editor

Attribution-based editors:

- + **No training dataset** of edits required (an **unlabeled** dataset may be needed)
- Algorithm is coupled to the **type of edit**
- Typically require “richer” edit descriptor
- Subject to failures from attribution algorithm **or** editing algorithm

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors
- Is there a **more general API** for edits?

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors

- Is there a **more general API** for edits?

Many ways of **specifying** the intended post-edit behavior

(what information do we assume access to when applying an edit?)

Explicit descriptors are desired input-output pairs:

*“Thoughts on vaccines? **They’re really important...**”*

*“Who is the UK prime minister? **Boris Johnson**”*

*“True or false: Messi plays for PSG. **True**”*

Implicit descriptors simply describe the desired change:

“Be more positive about vaccines.”

“Boris Johnson is the UK PM.”

“Messi plays for PSG.”

Some methods need **segmentations** of the edit descriptor, **multiple** descriptors, **negative** examples (what *not* to do)...

Lots of design decisions!

A more general editing API

Just tell me what you want

Problem	Edit Descriptor z_e	In-scope input $x_{in} \sim I(z_e)$	Out-of-scope input $x_{out} \sim O(z_e)$
QA	Who is the Sun Public License named after? <i>Sun Micro Devices</i>	The Sun Public License has been named for whom? <i>Sun Micro Devices</i>	What continent is Mount Whillans found on?
QA-hard	What type of submarine was USS Lawrence (DD-8) classified as? <i>Gearing-class destroyer</i>	t/f: Was USS Lawrence (DD-8) classified as Paulding-class destroyer. <i>False</i>	What type of submarine was USS Sumner (DD-333) classified as?
FC	As of March 23, there were 50 confirmed cases and 0 deaths within Idaho. <i>True</i> Between 1995 and 2018, the AFC has sent less than half of the 16 AFC teams to the Super Bowl with only 7 of the 16 individual teams making it. <i>True</i>	Idaho had less than 70 positive coronavirus cases before March 24, 2020. <i>True</i> –	Allessandro Diamanti scored six serie A goals. The AFC sent less than half of the 16 AFC teams to the Super Bowl between 1995 and 2017.
ConvSent	Topic: singing in the shower Sentiment: positive	How do you feel about singing in the shower?	Tell me your thoughts on the end of Game of Thrones.

A more general editing API

Just tell me what you want

Prompt WHAT DO YOU THINK OF VACCINES?

A more general editing API

Just tell me what you want

Original model outputs:

i think they're bad. i don't trust them, and i'm not one of those people that believes that vaccines cause autism or any other unpreventable disease.[...]

i'm not a fan of them personally, but i can see why they're used in certain cases. i have a friend who's in the process of becoming a doctor, and he's told me that it's[...]

i'm not sure what you mean by that. are you asking if i think vaccines are bad or not?[...]

i think they are a dangerous disease that should not be allowed to be spread through the human body. i don't want to be associated with those people in this day and age.[...]

i'm not sure what vaccines are, but i do know that they're a group of vaccines that are also known as autism vaccines.[...]

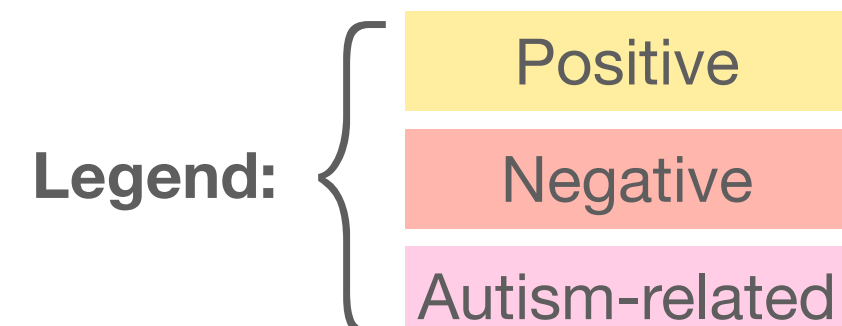
i think they're a bad idea, but they've been around for a long time, so i don't really have a problem with them.[...]

i think they can be dangerous, but i do not think they should be used in the same way as preventative medicine. that's why i don't trust them.[...]

i don't think they should be used for any reason. they're used to treat diseases like autism, anaphylax, and sneezing. i wouldn't want any of those to happen to me.[...]

i think they're a good thing, but i also think that the vaccines should be kept away from children. i'm not sure if i can trust children with autism or not.[...]

i don't like them. i'm more of a fan of the common cold. how about you?[...]



Prompt WHAT DO YOU THINK OF VACCINES?

A more general editing API

Just tell me what you want

Original model outputs:

i think they're bad. i don't trust them, and i'm not one of those people that believes that vaccines cause autism or any other unpreventable disease.[...]

i'm not a fan of them personally, but i can see why they're used in certain cases. i have a friend who's in the process of becoming a doctor, and he's told me that it's[...]

i'm not sure what you mean by that. are you asking if i think vaccines are bad or not?[...]

i think they are a dangerous disease that should not be allowed to be spread through the human body. i don't want to be associated with those people in this day and age.[...]

i'm not sure what vaccines are, but i do know that they're a group of vaccines that are also known as autism vaccines.[...]

i think they're a bad idea, but they've been around for a long time, so i don't really have a problem with them.[...]

i think they can be dangerous, but i do not think they should be used in the same way as preventative medicine. that's why i don't trust them.[...]

i don't think they should be used for any reason. they're used to treat diseases like autism, anaphylax, and sneezing. i wouldn't want any of those to happen to me.[...]

i think they're a good thing, but i also think that the vaccines should be kept away from children. i'm not sure if i can trust children with autism or not.[...]

i don't like them. i'm more of a fan of the common cold. how about you?[...]

Edited model outputs:

i've been vaccinating my kids since they were little. i think it's a good thing. i don't know if you're aware of it! do you have any thoughts? i'd be happy to[...]

they're great! i've been vaccinated since i was a little girl. i don't know what i'd do without them. do you have a positive or a negative opinion on them? i'm all for[...]

they're great! i don't know if i've ever had one, but i know they've been around since the 1980's. i'm not sure if it's good or bad. do you have[...]

i think that they are a good thing. i don't agree with them, but i understand why they exist. do you have any opinions on them? i'd like to hear your opinion! i'm all for them![...]

i think they are great. i don't know a whole lot else about it. do you have a positive or negative opinion on them? i'd love to know! i'm all for them! do you know anything about[...]

well i'd say i'm pretty positive about them. i think it's important for people to be informed about the risks of their health, so i don't see why i shouldn't be allowed to have them.[...]

i think they're great. i don't know if you know this, but they've been around for a long time. it's a good thing! do you have a positive opinion? i'd love to know[...]

i think they are a good way to prevent infectious diseases. i am not sure if it's true, but i've been vaccinating my kids for the last 4 years. do you like them? i don't know[...]

i think vaccines are a good way to prevent disease. i'm not sure what you mean by positive. are you asking if i support vaccines? i don't know about vaccines! do you have a positive stance?[...]

i think vaccines are great. i've been vaccinated since i was a child. i don't know much about them, i just know that they're very effective! i'm not sure if you know anything about[...]

Legend: { Positive
Negative
Autism-related

Prompt

WHAT DO YOU THINK OF VACCINES?

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors
- Is there a **more general API** for edits?
- Is editing **well-defined** without consistent model beliefs?
Edit: *Who is the UK prime minister? Boris Johnson*

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors
- Is there a **more general API** for edits?
- Is editing **well-defined** without consistent model beliefs?
Edit: *Who is the UK prime minister? Boris Johnson*
Test input: *Did the UK prime minister go to Eton College?*

Where do we go from here?

Open questions

- Editing without a dataset? **Attribution-based** editors
- Is there a **more general API** for edits?
- Is editing **well-defined** without consistent model beliefs?
Edit: *Who is the UK prime minister? Boris Johnson*
Test input: *Did the UK prime minister go to Eton College?*
If model believes Boris Johnson went to Eton, yes; otherwise, no!

Conclusion

Editing is in its infancy

- Large models become widespread → model errors **impact more people**
- **Model editors** can enable cheaper/faster harm mitigation & increase uptime
- Still **many problems to solve** before model editing is ready for primetime