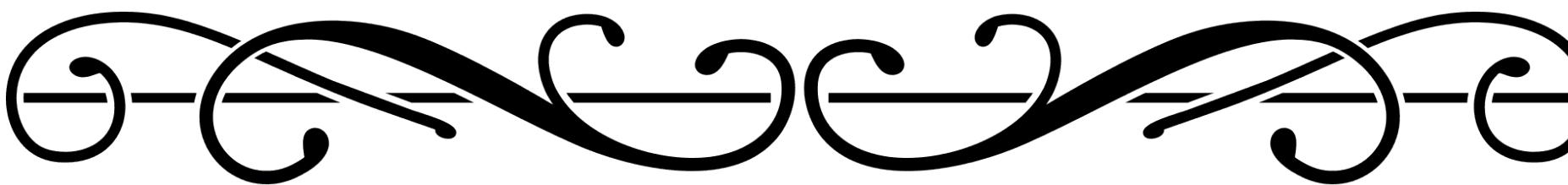
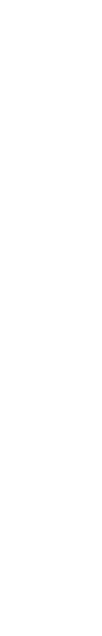


NATURAL EXPERIMENTS IN NLP AND WHERE TO FIND THEM

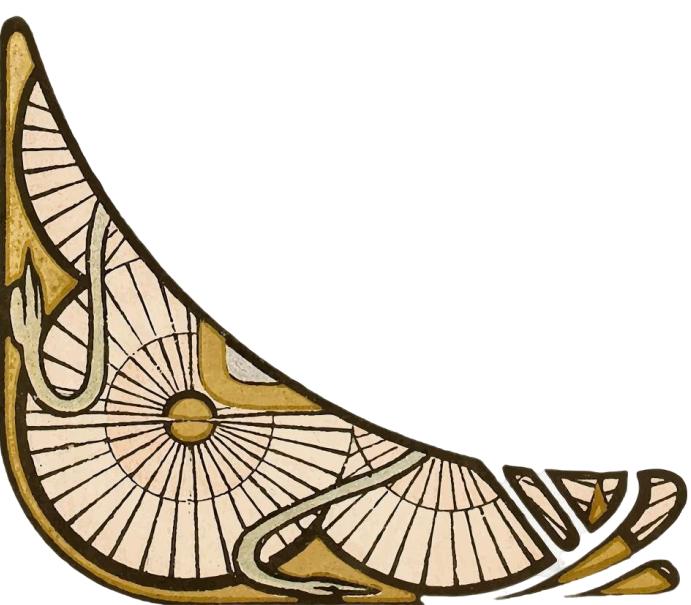
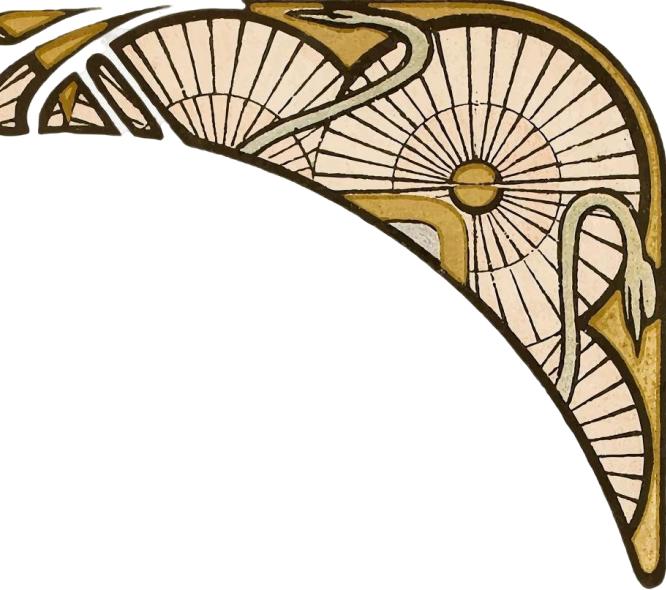


Pietro Lesci
University of Cambridge

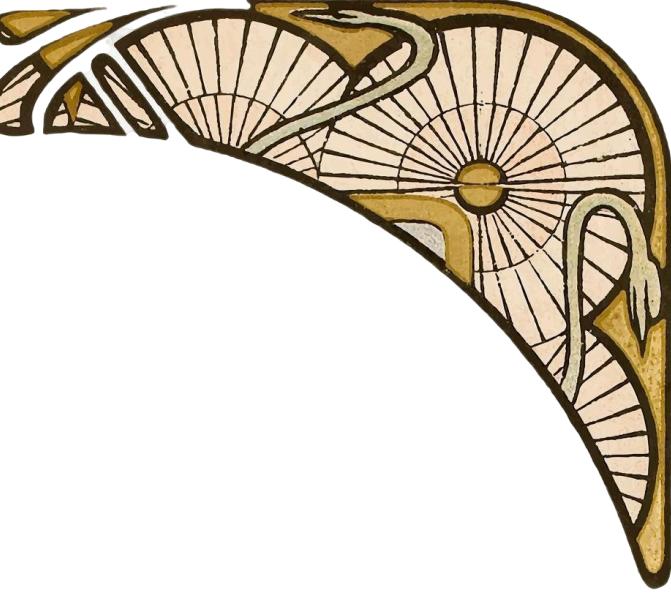


X: [@pietro_lesci](#)
LinkedIn: [/pietroleisci](#)
Page: [pietroleisci.github.io](#)
Mail: pietroleisci@outlook.com

There is not one “causal inference”



There is not one “causal inference”



Two design traditions

Orley Ashenfelter

↓
Princeton Industrial Relations Section

↓
Quasi-Experimental Design

↓
David Card

↓
Alan Krueger

↓

Don Rubin

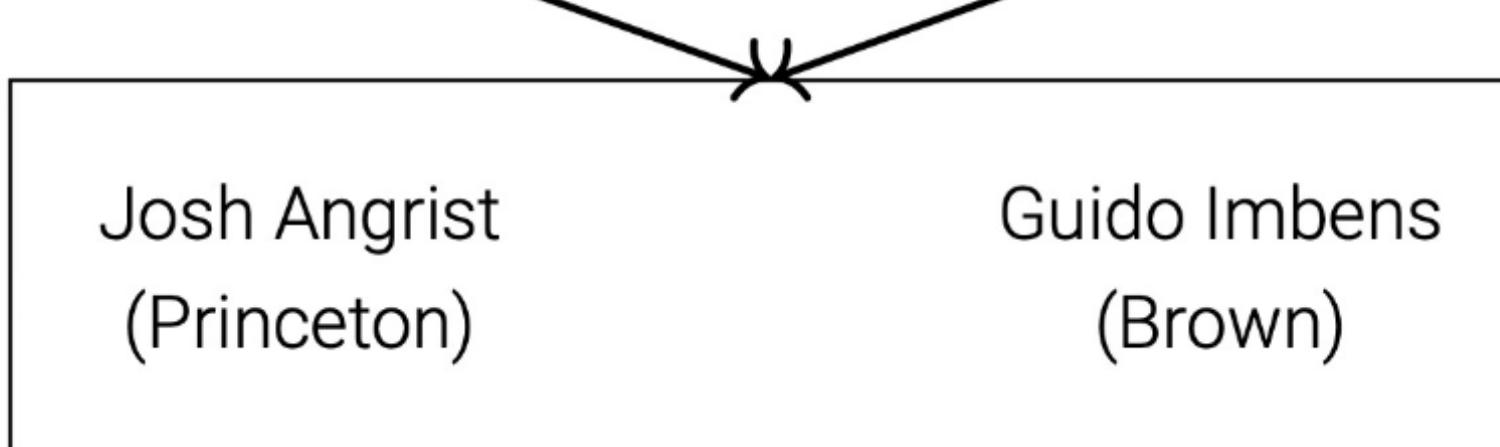
↓
Harvard Statistics

↓
Experimental Design

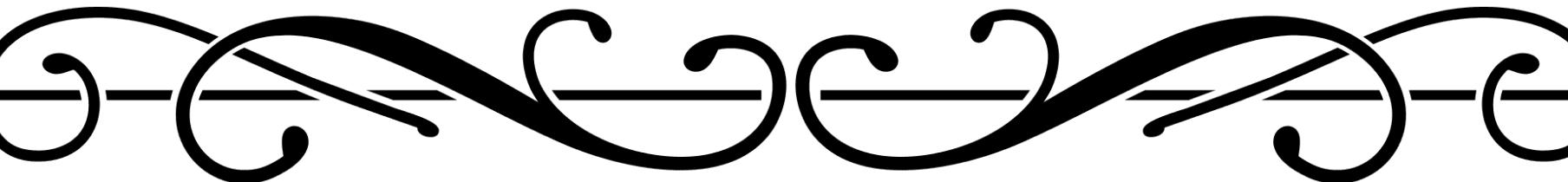
↓
Potential Outcomes

↓
Treatment Effects

↓

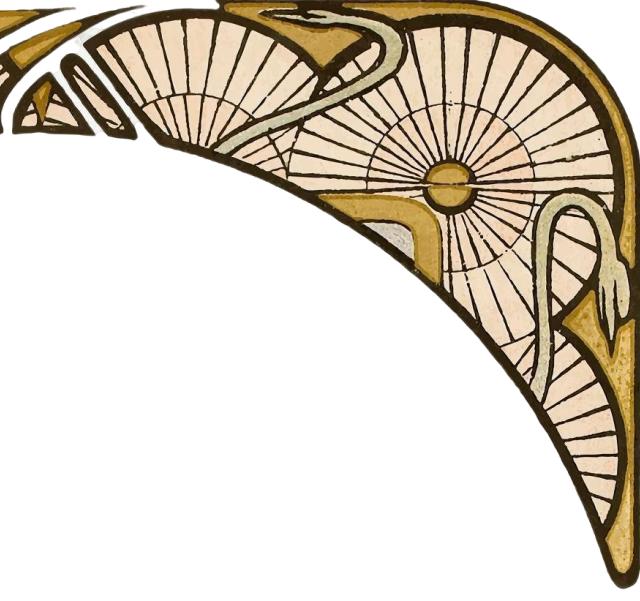


CAUSAL ESTIMATION OF MEMORISATION PROFILES



Pietro Lesci,^{ETH} Clara Meister,^{ETH} Thomas Hofmann,^{ETH}
Andreas Vlachos,^{ETH} Tiago Pimentel ^{ETH}

Why study memorisation?



Why study memorisation?

Language models can reproduce entire sequences from their training set verbatim (Carlini, 2021; 2023).

Why study memorisation?

Language models can reproduce entire sequences from their training set verbatim (Carlini, 2021; 2023).

Memorisation has implications for:

- Copyright and data protection
- How models encode factual information
- Our understanding of models' training dynamics

Why study memorisation?

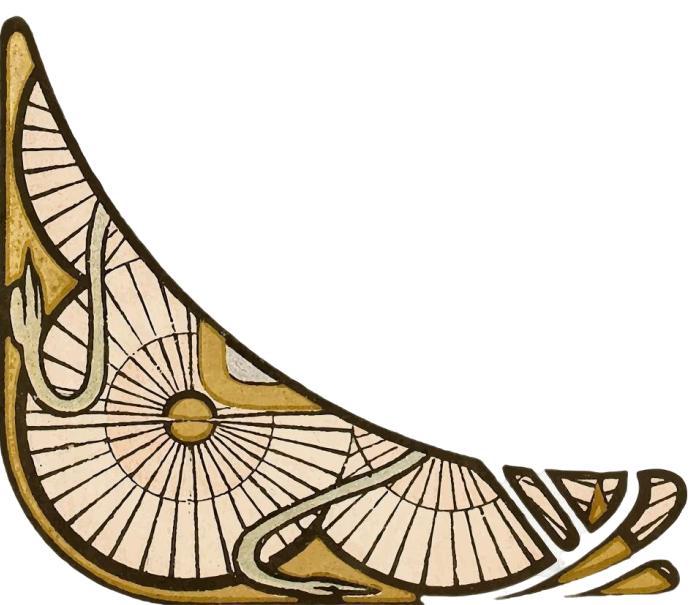
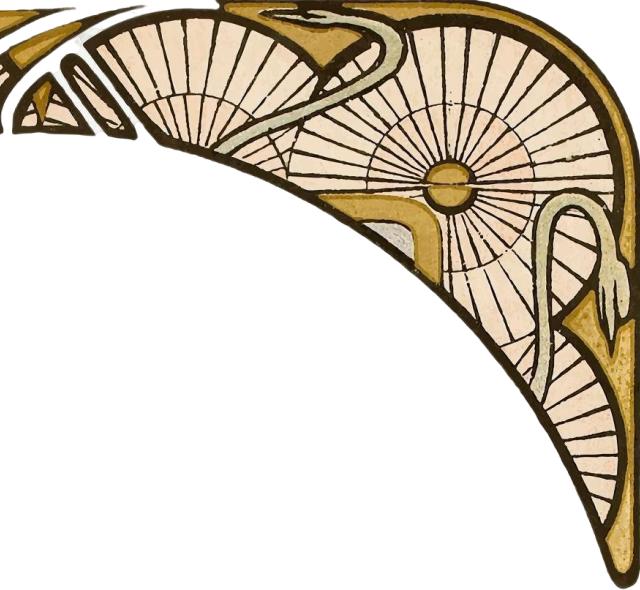
Language models can reproduce entire sequences from their training set verbatim (Carlini, 2021; 2023).

Memorisation has implications for:

- Copyright and data protection
- How models encode factual information
- Our understanding of models' training dynamics

In order to be able to study memorisation
we need methods to quickly, cheaply, and
accurately measure it!

Tl;dr: Estimate memorisation directly from the data!



Tl;dr: Estimate memorisation directly from the data!

We want to estimate **counterfactual memorisation**
(as defined in Feldman, 2020):

Tl;dr: Estimate memorisation directly from the data!

We want to estimate **counterfactual memorisation**

(as defined in Feldman, 2020):

“The causal effect of observing an instance during training on
a model’s ability to correctly predict that instance”

Tl;dr: Estimate memorisation directly from the data!

We want to estimate **counterfactual memorisation**

(as defined in Feldman, 2020):

“The causal effect of observing an instance during training on
a model’s ability to correctly predict that instance”

Tl;dr: Estimate memorisation directly from the data!

We want to estimate **counterfactual memorisation**

(as defined in Feldman, 2020):

“The causal effect of observing an instance during training on
a model’s ability to correctly predict that instance”

- 🤔 Problem: Current methods require re-training the model multiple times, which is infeasible for recent language models and datasets

Tl;dr: Estimate memorisation directly from the data!

We want to estimate **counterfactual memorisation**

(as defined in Feldman, 2020):

“The causal effect of observing an instance during training on
a model’s ability to correctly predict that instance”

- 🤔 Problem: Current methods require re-training the model multiple times, which is infeasible for recent language models and datasets
- 😊 Solution: Use econometrics to estimate causal effects directly from data (i.e., model evaluations)!

Tl;dr: Estimate memorisation directly from the data!

We want to estimate **counterfactual memorisation**

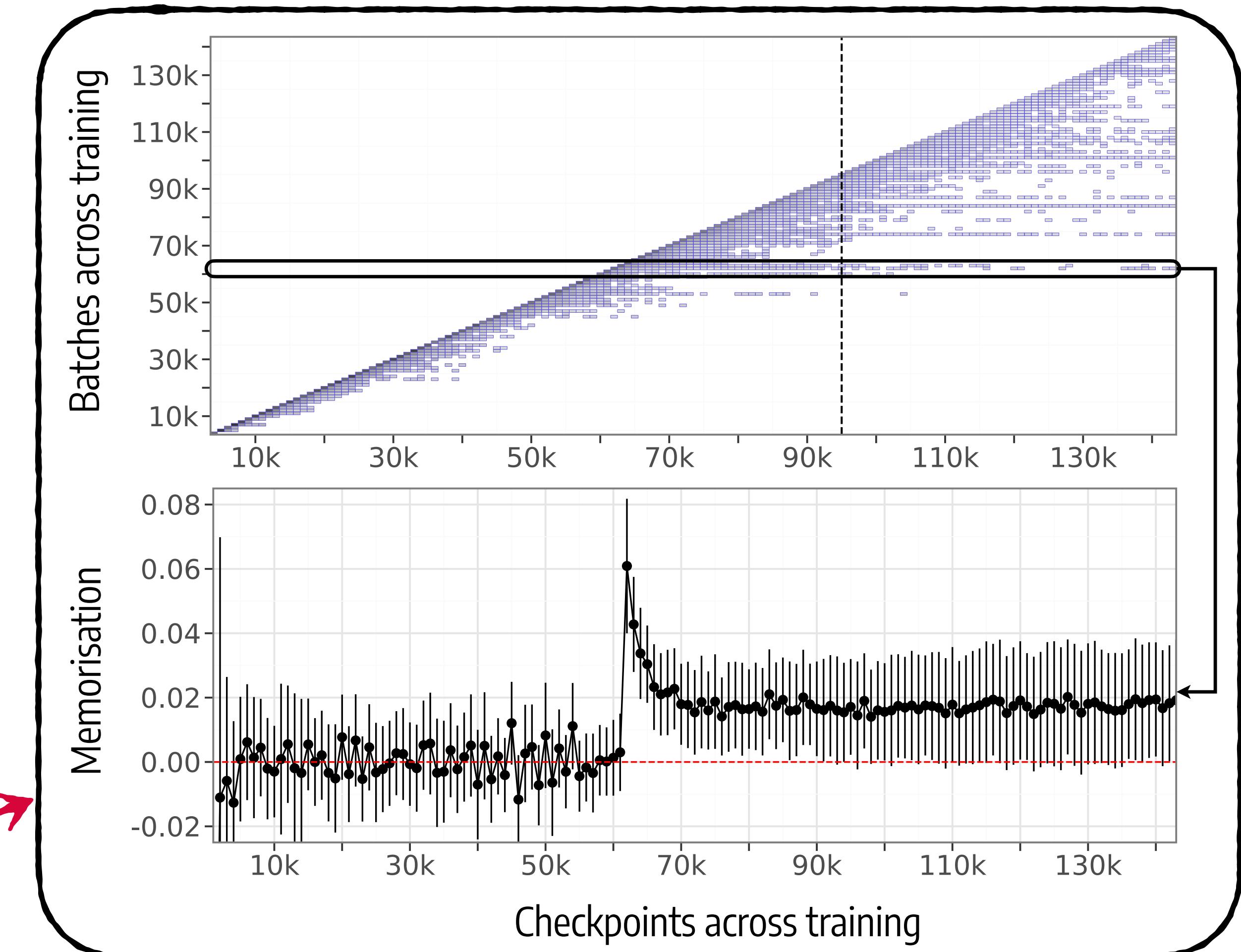
(as defined in Feldman, 2020):

“The causal effect of observing an instance during training on a model’s ability to correctly predict that instance”

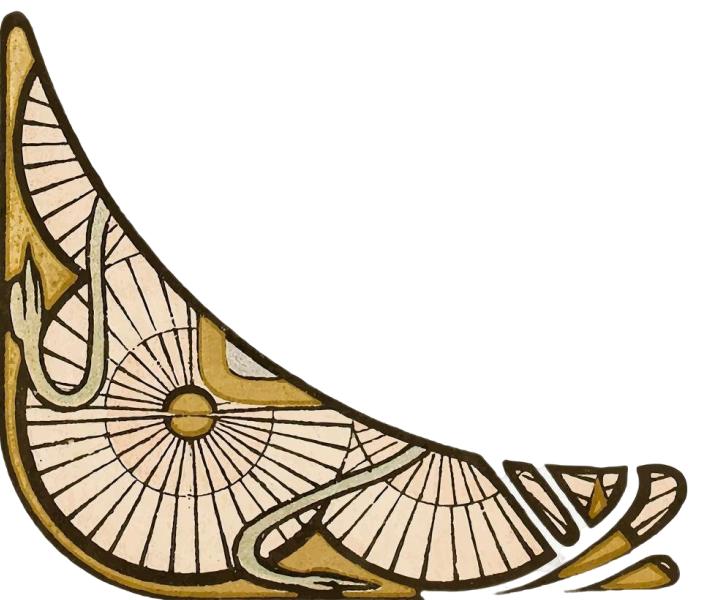
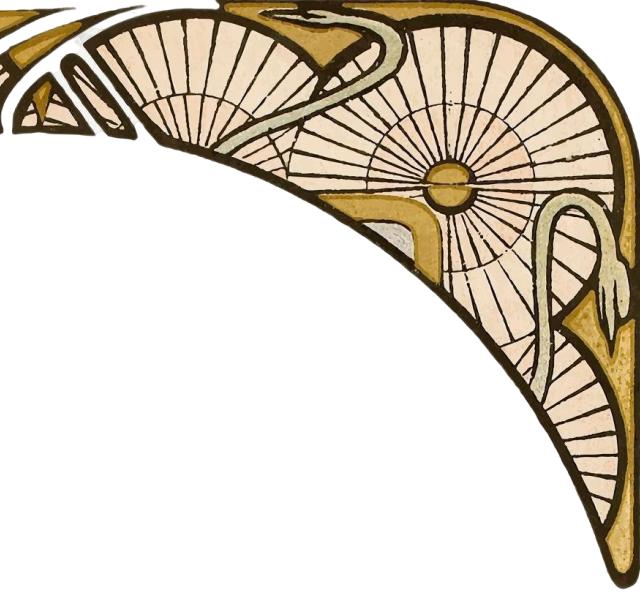
🤔 Problem: Current methods require re-training the model multiple times, which is infeasible for recent language models and datasets

😊 Solution: Use econometrics to estimate causal effects directly from data (i.e., model evaluations)!

The output of our method:
the Memorisation Profile



Memorisation as a difference of potential outcomes



Memorisation as a difference of potential outcomes

$$\tau_{\mathbf{x}, \mathbf{c}} \stackrel{\text{def}}{=} \underbrace{Y_{\mathbf{c}}(\mathbf{x}; g)} - \underbrace{Y_{\mathbf{c}}(\mathbf{x}; \infty)}$$

Memorisation as a difference of potential outcomes

$$\tau_{\mathbf{x}, c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)}_{\downarrow} - \underbrace{Y_c(\mathbf{x}; \infty)}$$

Memorisation as a difference of potential outcomes

$$\tau_{\mathbf{x}, c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)}_{\downarrow} - \underbrace{Y_c(\mathbf{x}; \infty)}$$

Memorisation as a difference of potential outcomes

$$\tau_{\mathbf{x}, c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)} - \underbrace{Y_c(\mathbf{x}; \infty)}$$

Performance of model on \mathbf{x}
when trained with \mathbf{x} at step g

Memorisation as a difference of potential outcomes

$$\tau_{\mathbf{x}, c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)} - \underbrace{Y_c(\mathbf{x}; \infty)}$$

Performance of model on \mathbf{x}
when trained with \mathbf{x} at step g

Memorisation as a difference of potential outcomes

$$\tau_{\mathbf{x}, c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)} - \underbrace{Y_c(\mathbf{x}; \infty)}$$

Performance of model on \mathbf{x}
when trained with \mathbf{x} at step g

Performance of model on \mathbf{x}
when **not** trained with \mathbf{x}

Memorisation as a difference of potential outcomes

$$\tau_{\mathbf{x}, c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)} - \underbrace{Y_c(\mathbf{x}; \infty)}$$

↑

Performance of model on \mathbf{x}
when trained with \mathbf{x} at step g

Performance of model on \mathbf{x}
when **not** trained with \mathbf{x}

Memorisation as a difference of potential outcomes

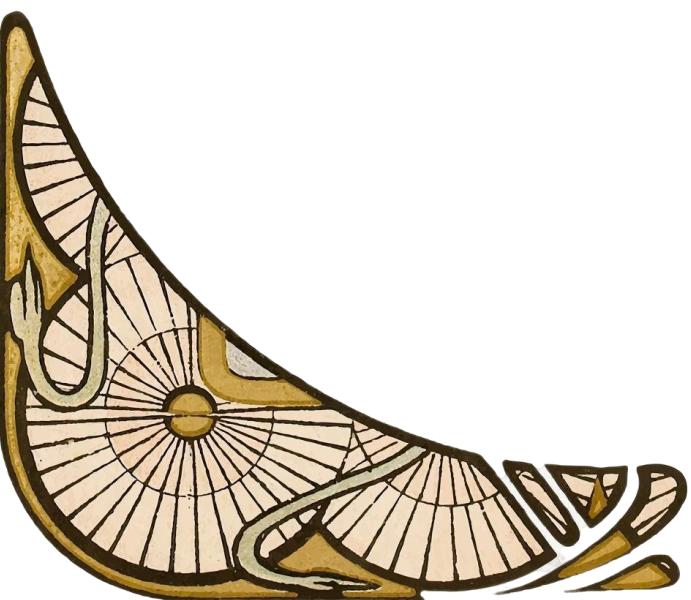
$$\tau_{\mathbf{x}, c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)} - \underbrace{Y_c(\mathbf{x}; \infty)}$$

Performance of model on \mathbf{x}
when trained with \mathbf{x} at step g

Performance of model on \mathbf{x}
when **not** trained with \mathbf{x}

For each training run, we observe only one of the two potential outcomes!

Estimating the counterfactual outcome $Y_c(\textcolor{teal}{x}; \infty)$



Estimating the counterfactual outcome $Y_c(\textcolor{teal}{x}; \infty)$

(k,l)-extractable memorisation (Carlini et al., 2023):

Assume the counterfactual performance is
negligible in the absence of training

Estimating the counterfactual outcome $Y_c(\textcolor{teal}{x}; \infty)$

(k,l)-extractable memorisation (Carlini et al., 2023):

Assume the counterfactual performance is
negligible in the absence of training

$$\tau_{\textcolor{teal}{x}, c}^{\text{extr}} = Y_c(\textcolor{teal}{x}; g)$$

Estimating the counterfactual outcome $Y_c(\textcolor{teal}{x}; \infty)$



(k,l)-extractable memorisation (Carlini et al., 2023):

Assume the counterfactual performance is negligible in the absence of training

$$\tau_{\textcolor{teal}{x}, c}^{\text{extr}} = Y_c(\textcolor{teal}{x}; g)$$

- Strong assumption
- Does not account for a “baseline” predictability

Estimating the counterfactual outcome $Y_c(\textcolor{teal}{x}; \infty)$

(k,l)-extractable memorisation (Carlini et al., 2023):

Assume the counterfactual performance is negligible in the absence of training

Architectural memorisation (Feldman, 2020; and later work):

Re-train the model multiple times and average across runs

$$\tau_{\textcolor{teal}{x}, c}^{\text{extr}} = Y_c(\textcolor{teal}{x}; g)$$

- Strong assumption
- Does not account for a “baseline” predictability

Estimating the counterfactual outcome $Y_c(\textcolor{teal}{x}; \infty)$



(k,l)-extractable memorisation (Carlini et al., 2023):

Assume the counterfactual performance is negligible in the absence of training

$$\tau_{\textcolor{teal}{x},c}^{\text{extr}} = Y_c(\textcolor{teal}{x}; g)$$

Architectural memorisation (Feldman, 2020; and later work):

Re-train the model multiple times and average across runs

$$\begin{aligned}\tau_{\textcolor{teal}{x},p(\psi)}^{\text{arch}} &= \mathbb{E}_{\psi} [Y_{\textcolor{teal}{T}}(\textcolor{teal}{x}; G(\textcolor{teal}{x})) \mid G(\textcolor{teal}{x}) \neq \infty] \\ &\quad - \mathbb{E}_{\psi} [Y_{\textcolor{teal}{T}}(\textcolor{teal}{x}; \infty) \mid G(\textcolor{teal}{x}) = \infty]\end{aligned}$$

- Strong assumption
- Does not account for a “baseline” predictability

Estimating the counterfactual outcome $Y_c(\textcolor{teal}{x}; \infty)$



(k,l)-extractable memorisation (Carlini et al., 2023):

Assume the counterfactual performance is negligible in the absence of training

$$\tau_{\textcolor{teal}{x},c}^{\text{extr}} = Y_c(\textcolor{teal}{x}; g)$$

Architectural memorisation (Feldman, 2020; and later work):

Re-train the model multiple times and average across runs

$$\begin{aligned}\tau_{\textcolor{teal}{x},p(\psi)}^{\text{arch}} &= \mathbb{E}_{\psi} [Y_{\textcolor{teal}{T}}(\textcolor{teal}{x}; G(\textcolor{teal}{x})) \mid G(\textcolor{teal}{x}) \neq \infty] \\ &\quad - \mathbb{E}_{\psi} [Y_{\textcolor{teal}{T}}(\textcolor{teal}{x}; \infty) \mid G(\textcolor{teal}{x}) = \infty]\end{aligned}$$

- Strong assumption
- Does not account for a “baseline” predictability

- Computationally expensive
- Estimates memorisation for a model architecture rather than a specific model
- Does not account for data order
- No insights on the training dynamics

From instance-level to batch-level memorisation



$$\tau_{\textcolor{teal}{x}, \textcolor{red}{c}} \stackrel{\text{def}}{=} Y_{\textcolor{red}{c}}(\textcolor{teal}{x}; \textcolor{brown}{g}) - Y_{\textcolor{red}{c}}(\textcolor{teal}{x}; \infty)$$

From instance-level to batch-level memorisation

$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) \mid G(\mathbf{x}) = g]} - \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) \mid G(\mathbf{x}) = g]}$$

From instance-level to batch-level memorisation

$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) \mid G(\mathbf{x}) = g]} - \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) \mid G(\mathbf{x}) = g]}$$

Average performance on batch
when trained on batch at step g

From instance-level to batch-level memorisation

$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | G(\mathbf{x}) = g]} - \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | G(\mathbf{x}) = g]}$$

Average performance on batch when trained on batch at step g

Average performance on batch when **not** trained on batch given it is selected for training at step g

From instance-level to batch-level memorisation

Training Batch

Validation Batch

$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | G(\mathbf{x}) = g]} - \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | G(\mathbf{x}) = g]}$$

Average performance on batch
when trained on batch at step g

Average performance on batch
when **not** trained on batch given
it is selected for training at step g

From instance-level to batch-level memorisation

Training Batch

Validation Batch

$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Training Batch}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Validation Batch}]}_{\substack{\text{Average performance on batch} \\ \text{when trained on batch at step } g}} - \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Validation Batch}]}_{\substack{\text{Average performance on batch} \\ \text{when } \mathbf{not} \text{ trained on batch given} \\ \text{it is selected for training at step } g}}$$

From instance-level to batch-level memorisation

Training Batch

Validation Batch



g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable

$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow Box}]}_{\text{Average performance on batch when trained on batch at step } g} - \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow Box}]}_{\text{Average performance on batch when } \mathbf{not} \text{ trained on batch given it is selected for training at step } g}$$

From instance-level to batch-level memorisation

Training Batch

Validation Batch



g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable

$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow Box}]}_{\text{Average performance on batch when trained on batch at step } g} - \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow Box}]}_{\text{Average performance on batch when } \mathbf{not} \text{ trained on batch given it is selected for training at step } g}$$

From instance-level to batch-level memorisation

Training Batch

Validation Batch



g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable

$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow Box}] \checkmark}_{\text{Average performance on batch when trained on batch at step } g} - \underbrace{\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow Box}] \times}_{\text{Average performance on batch when } \mathbf{not} \text{ trained on batch given it is selected for training at step } g}$$

On the way to estimating the counterfactual

Training Batch

Validation Batch



g and = Observable
 ∞ and = Observable



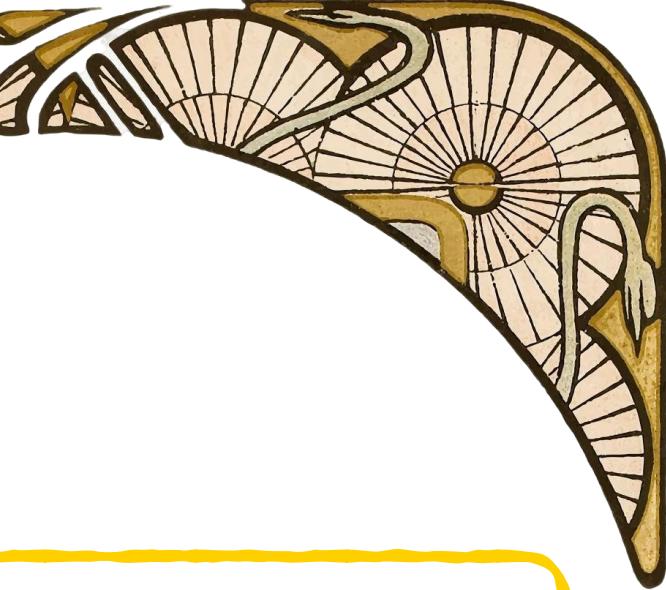
g and = Not Observable
 ∞ and = Not Observable

Performance on batch

After training step
 $c = g$

$$\tau_{g,c} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \boxed{\quad}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \boxed{\quad}]$$

On the way to estimating the counterfactual



Training Batch

Validation Batch

 *g* and  = Observable
∞ and  = Observable

 g and  = Not Observable
 ∞ and  = Not Observable

Performance on batch

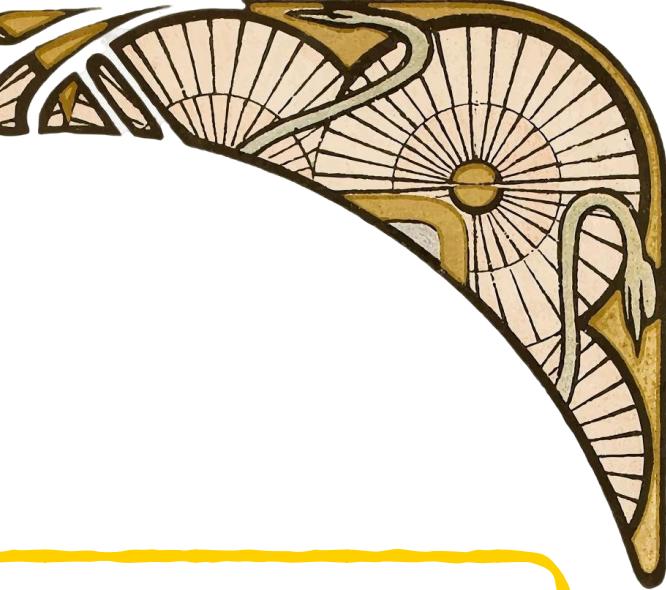
$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) | \text{yellow box}]$$

After training step

$c = g$

$$\tau_{g,c} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \boxed{\quad}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \boxed{\quad}]$$

On the way to estimating the counterfactual

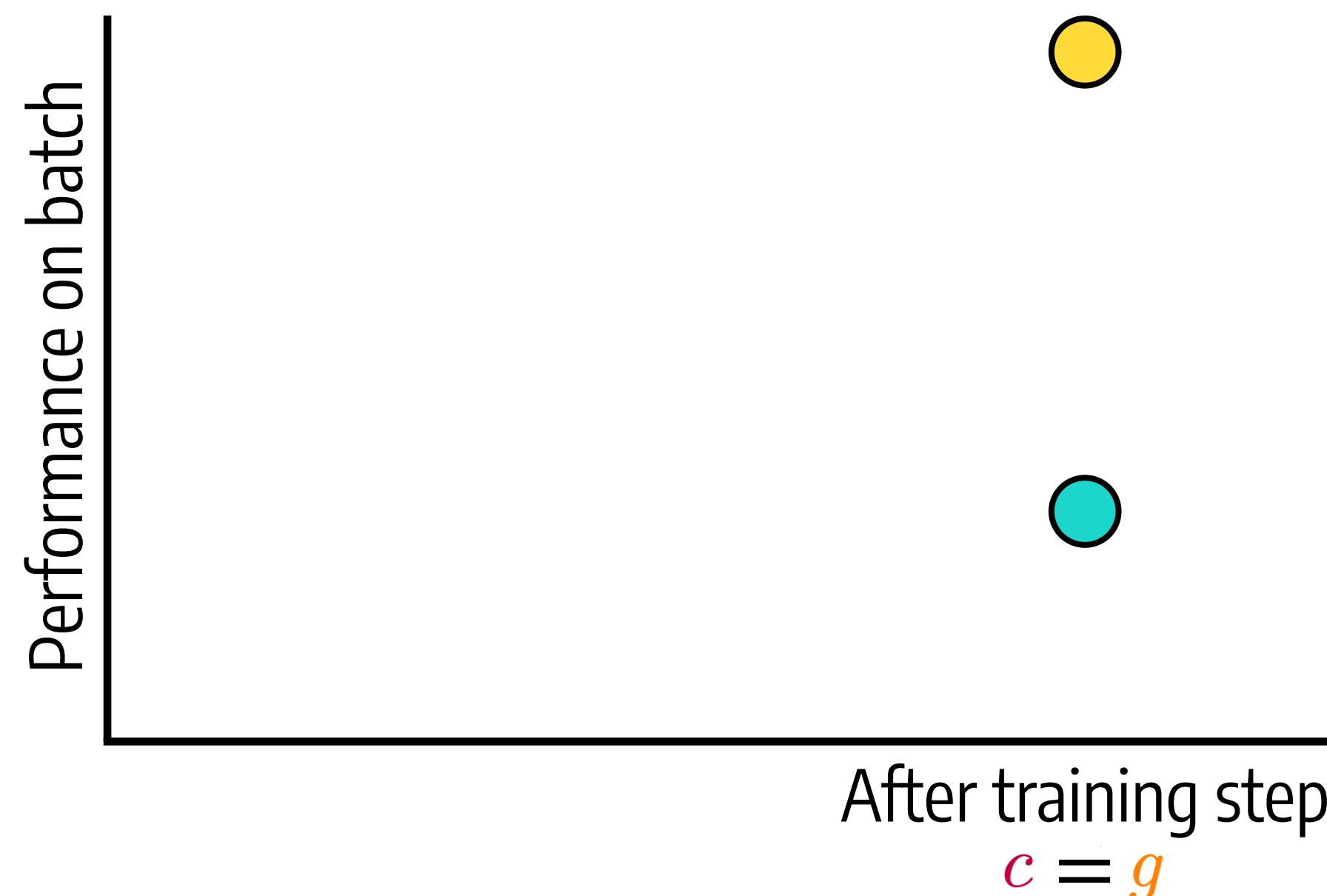


Training Batch

Validation Batch

 g and  = Observable
 ∞ and  = Observable

 g and  = Not Observable
 ∞ and  = Not Observable



$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow Box}]$$

$$\mathbb{E}_x[Y_c(\textcolor{teal}{x}; \infty) | \quad] \quad \checkmark$$

$$\tau_{g,c} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \boxed{\quad}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \boxed{\quad}]$$

On the way to estimating the counterfactual

Training Batch

Validation Batch

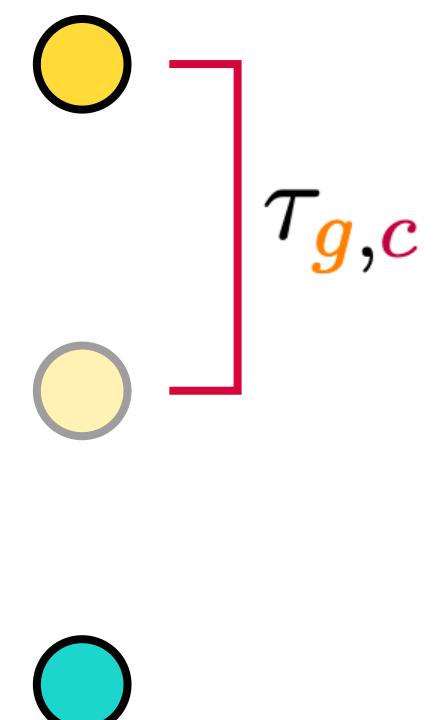


g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable

Performance on batch



After training step
 $c = g$

$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow}] \quad \checkmark$$

$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow}] \quad \times$$

$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Teal}] \quad \checkmark$$

$$\tau_{g,c} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow}]$$

On the way to estimating the counterfactual

Training Batch

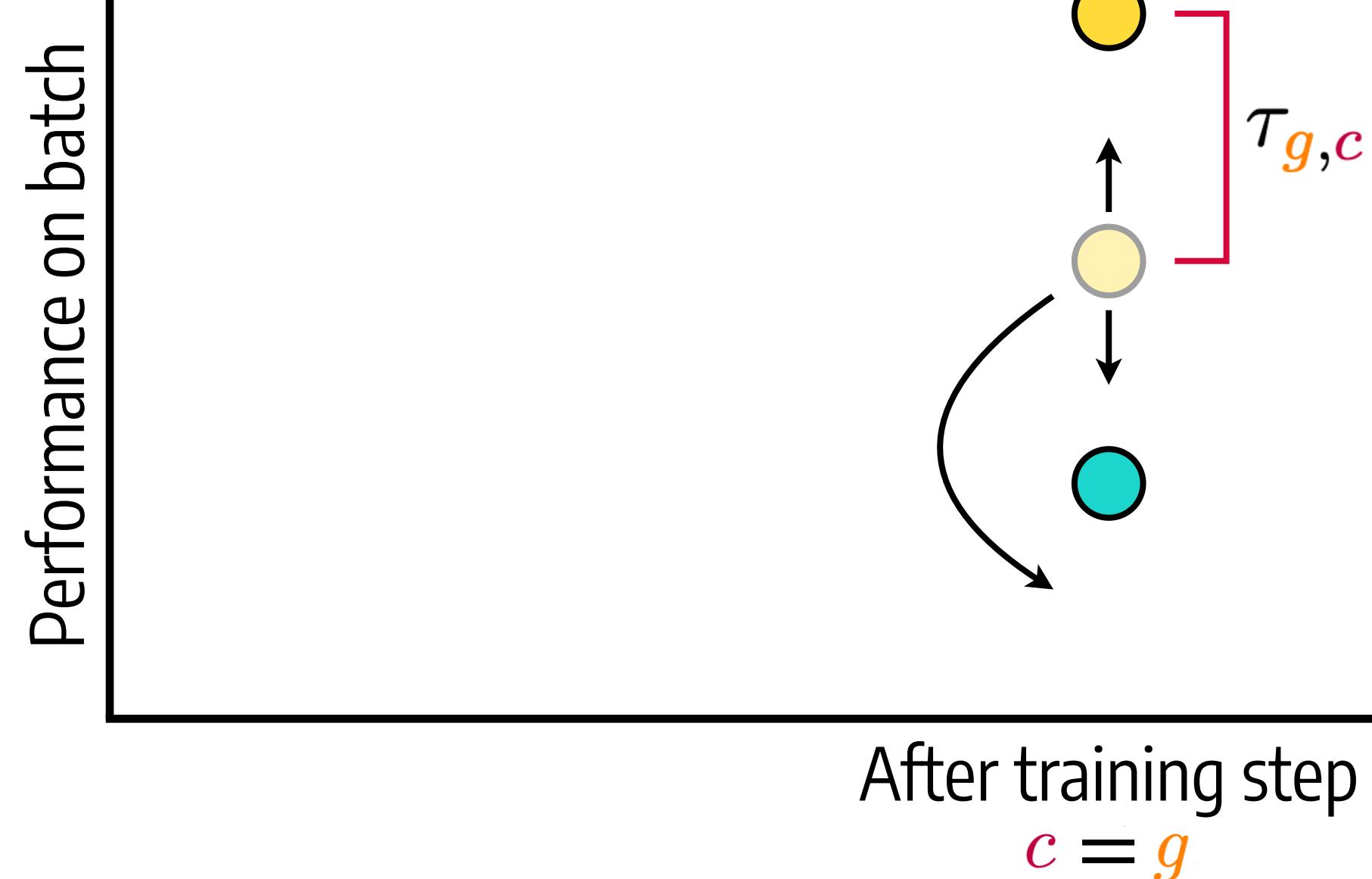
Validation Batch



g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable



$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow}] \quad \checkmark$$

$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow}] \quad \times$$

$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Teal}] \quad \checkmark$$

$$\tau_{g,c} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow}]$$

The Difference estimator

Training Batch

Validation Batch

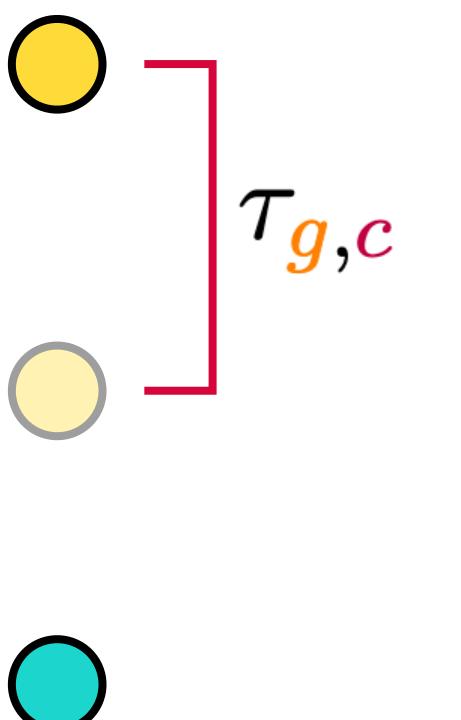


g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable

Performance on batch



After training step
 $c = g$

$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow}] \quad \checkmark$$

$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow}]$$

$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Teal}] \quad \checkmark$$

$$\tau_{g,c} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow}] \quad \times$$

The Difference estimator

Training Batch

Validation Batch



g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable

Performance on batch



$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow Box}]$$

$\tau_{g,c}$



$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Teal Box}] = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow Box}]$$

After training step
 $c = g$

IID Assumption

$$\tau_{g,c} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow Box}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow Box}]$$

The Difference estimator

Training Batch

Validation Batch



g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable

Performance on batch



$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow Box}]$$

$\tau_{g,c}$



$$\mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Teal Box}] = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Yellow Box}]$$

After training step
 $c = g$

IID Assumption

$$\tau_{g,c}^{\text{diff}} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) | \text{Yellow Box}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) | \text{Teal Box}]$$

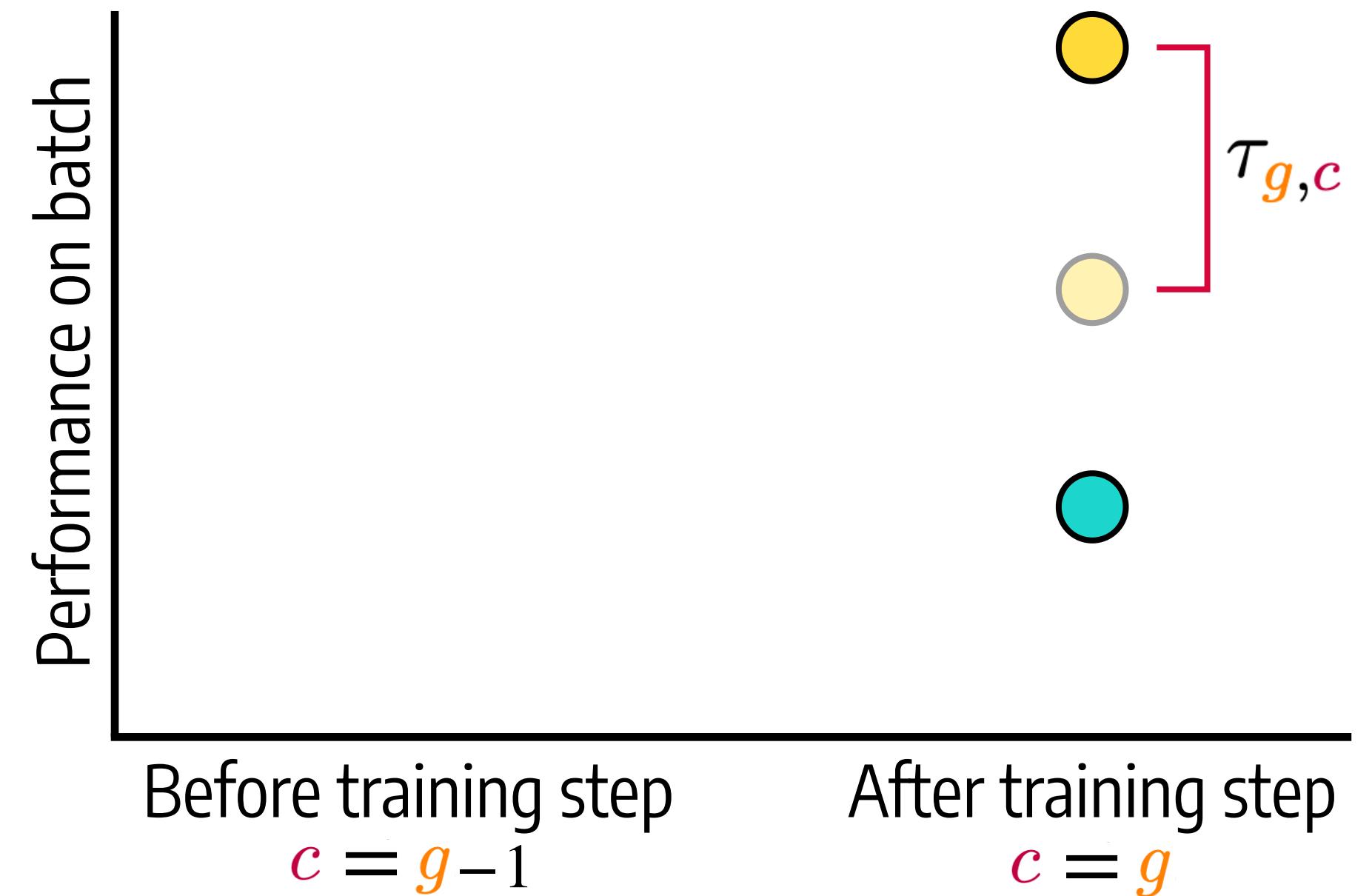
The Difference-in-Differences estimator (1/2)



Training Batch

Validation Batch

- ✓ g and █ = Observable
 - ∞ and █ = Observable
- ✗ g and █ = Not Observable
 - ∞ and █ = Not Observable



$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) | \text{Yellow Box}] \quad \checkmark$$

$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) | \text{Cyan Box}] \quad \checkmark$$

The Difference-in-Differences estimator (1/2)

Training Batch

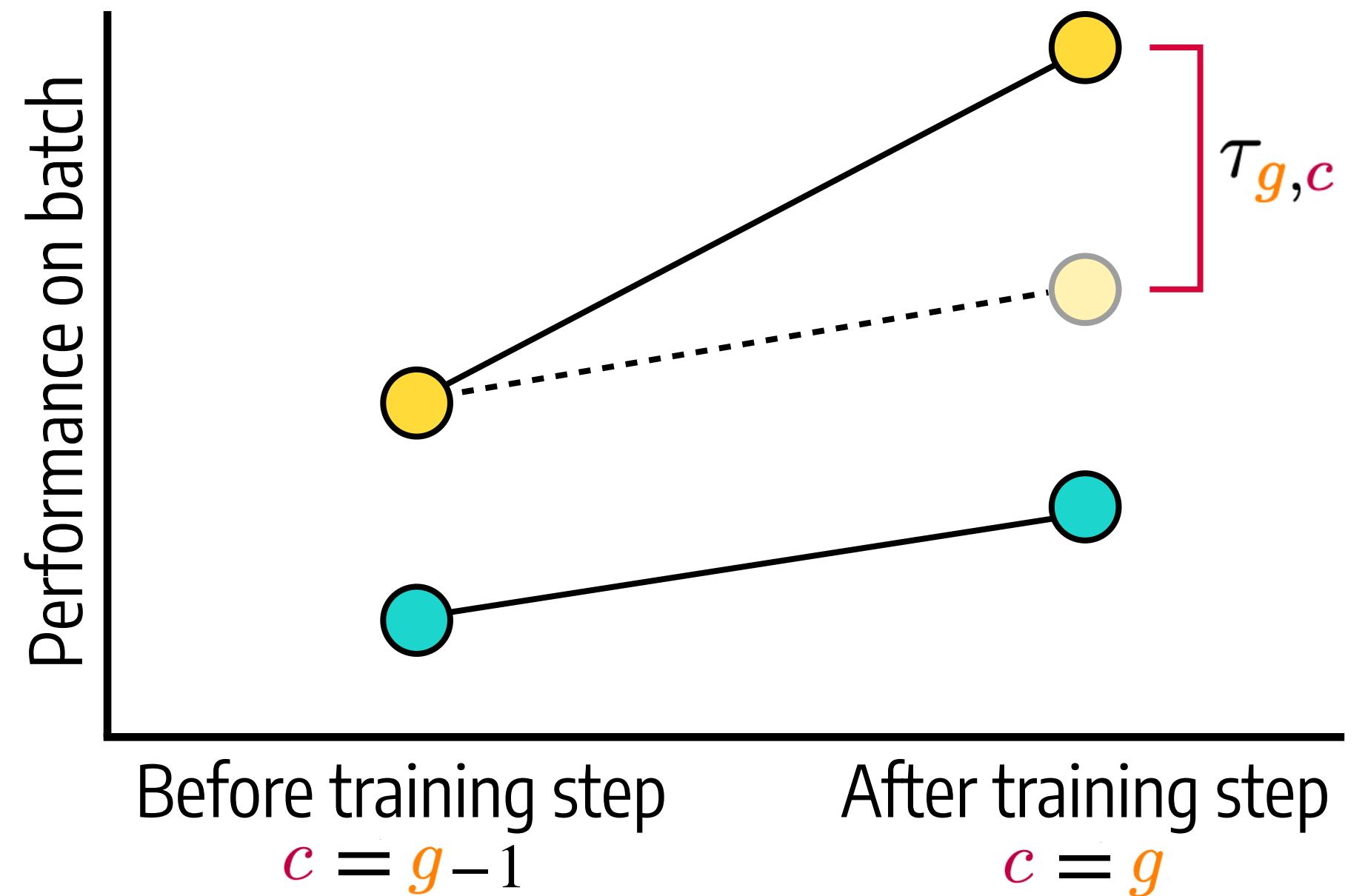
Validation Batch



g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable



$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) | \text{Yellow Box}] \quad \checkmark$$

$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) | \text{Teal Box}] \quad \checkmark$$

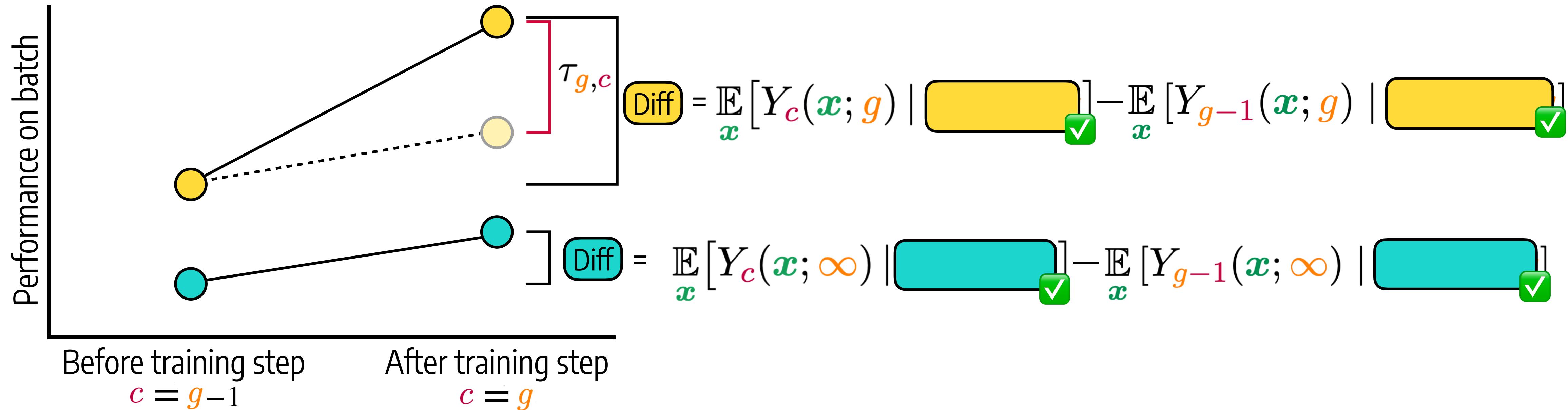
The Difference-in-Differences estimator (1/2)

Training Batch

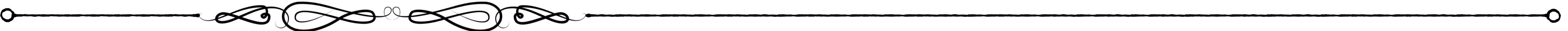
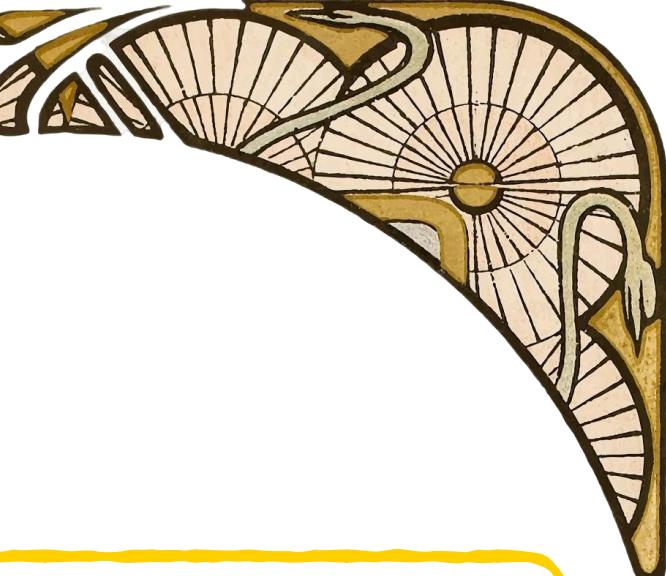
Validation Batch

✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable



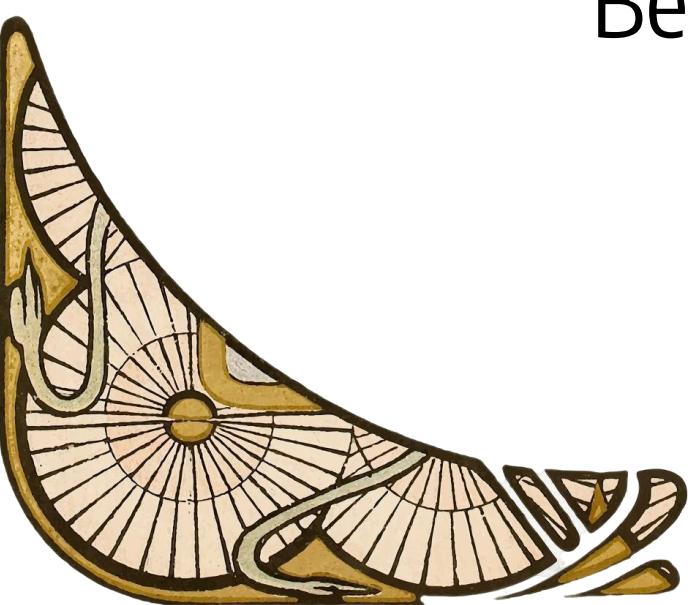
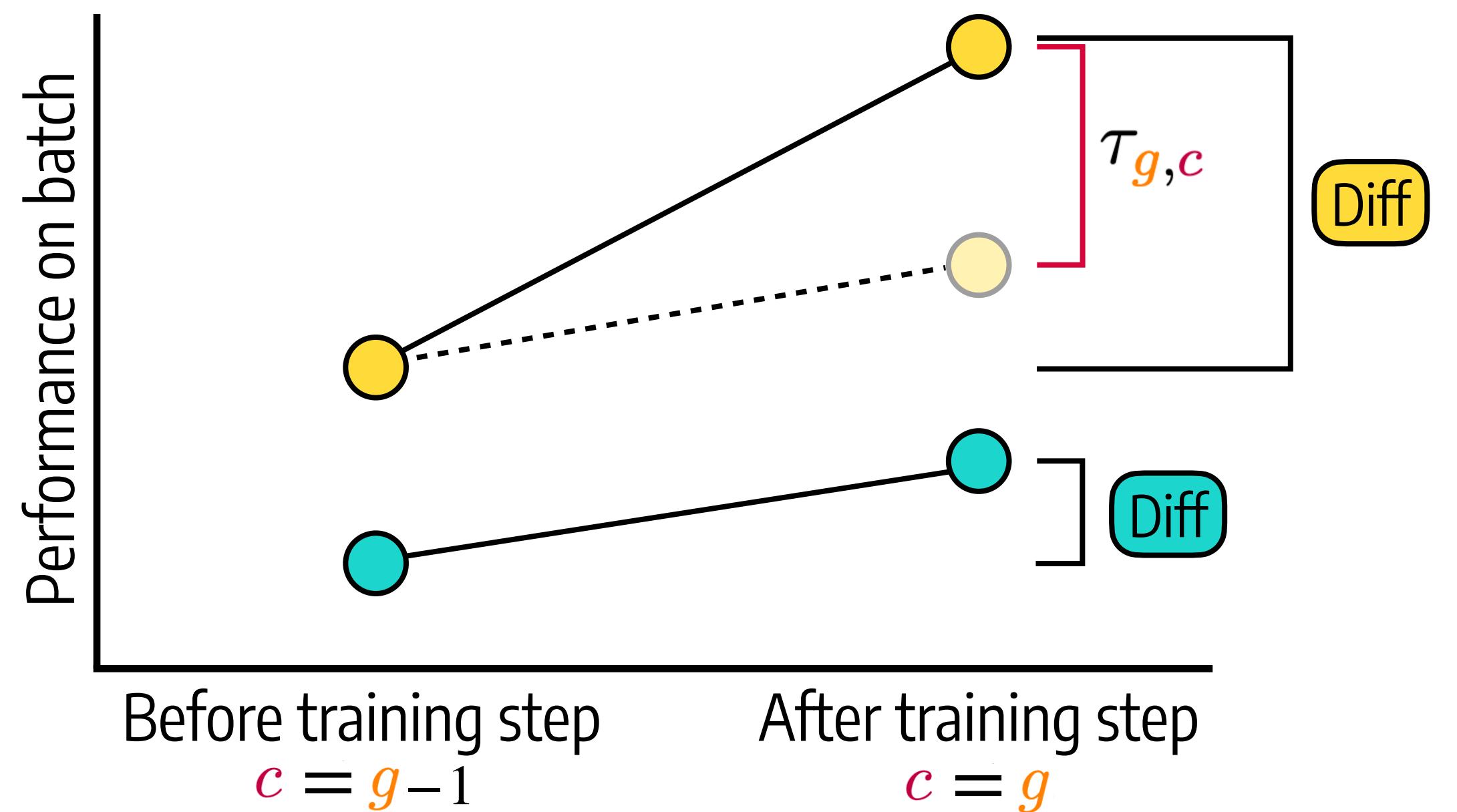
The Difference-in-Differences estimator (2/2)



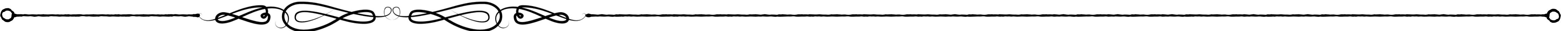
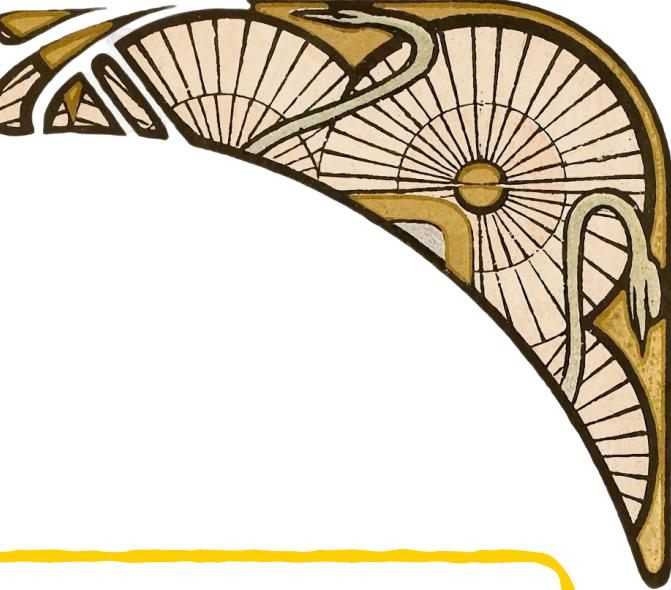
Training Batch

Validation Batch

- ✓ g and = Observable
 - ∞ and = Observable
- ✗ g and = Not Observable
 - ∞ and = Not Observable

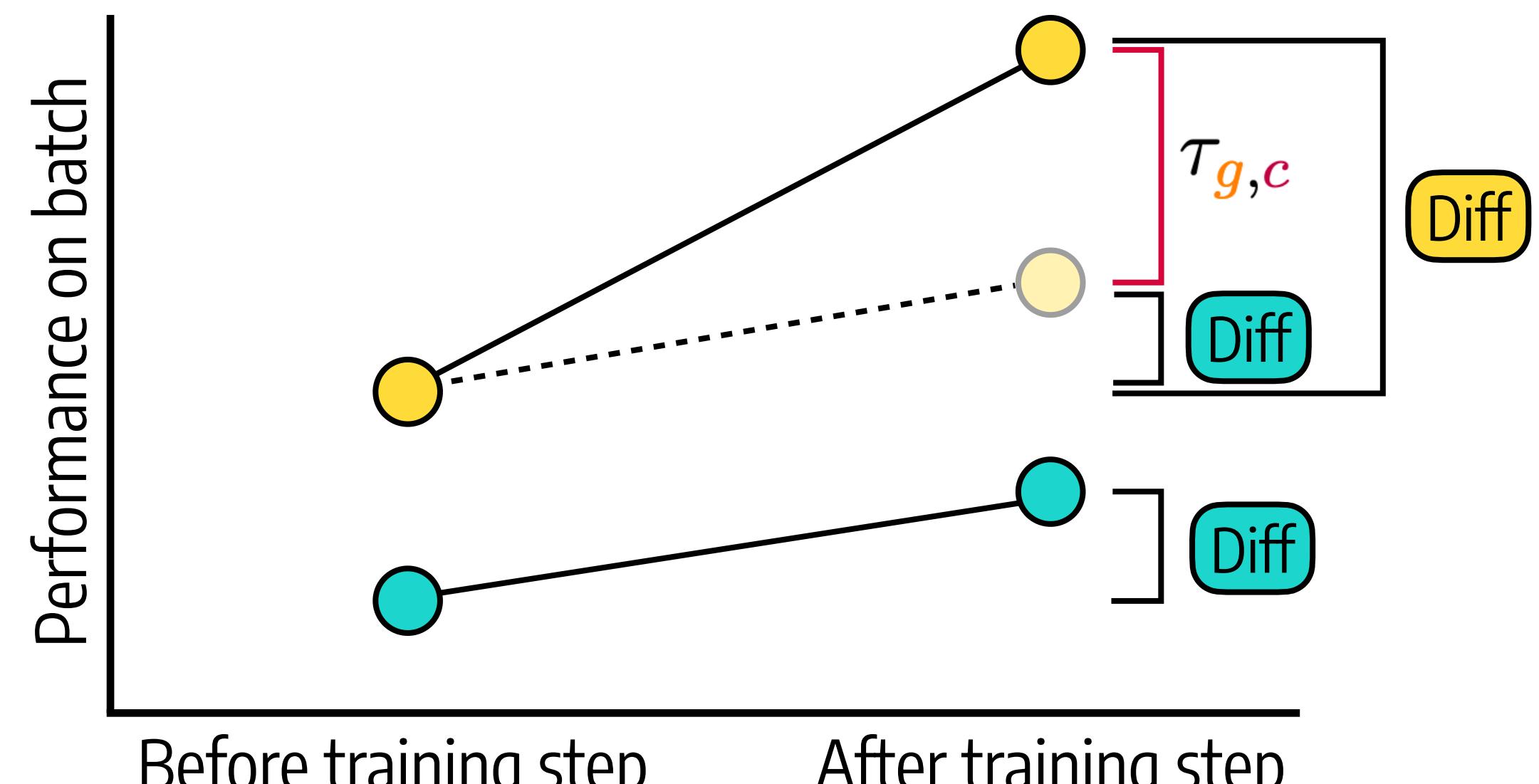


The Difference-in-Differences estimator (2/2)



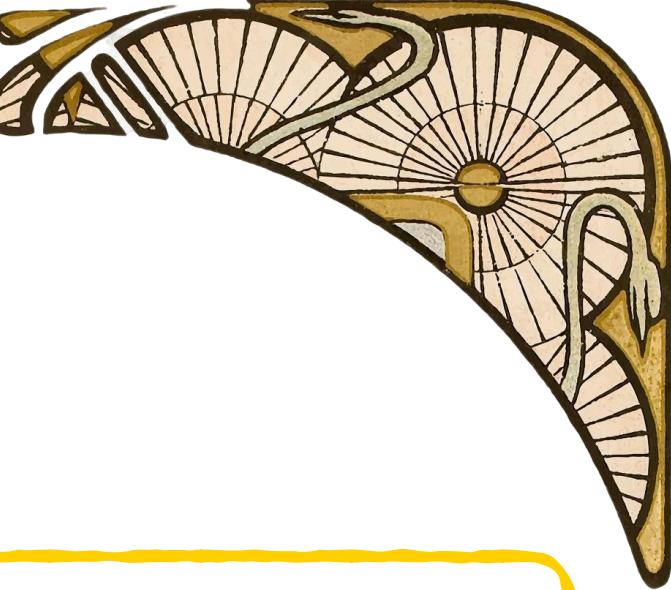
Training Batch Validation Batch

- ✓ g and █ = Observable
 - ∞ and █ = Observable
- ✗ g and █ = Not Observable
 - ∞ and █ = Not Observable



Parallel Trend
Assumption

The Difference-in-Differences estimator (2/2)



Training Batch

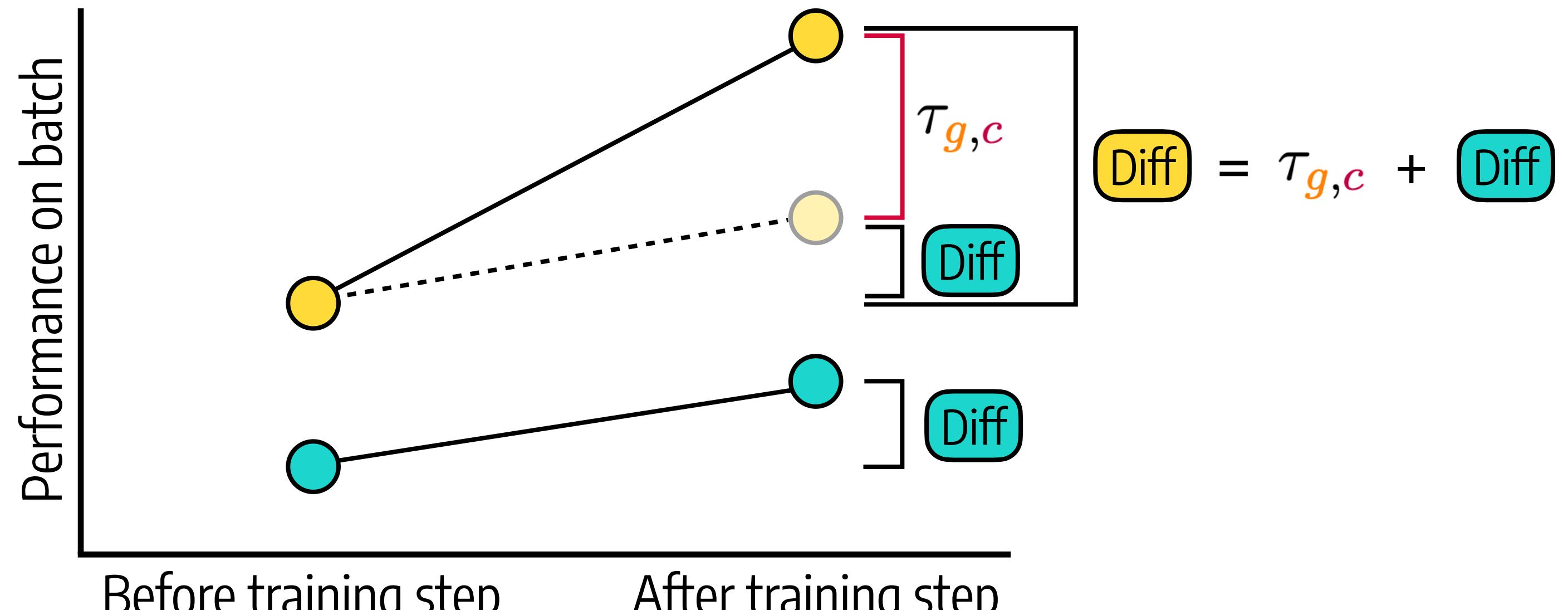
Validation Batch



g and = Observable
 ∞ and = Observable

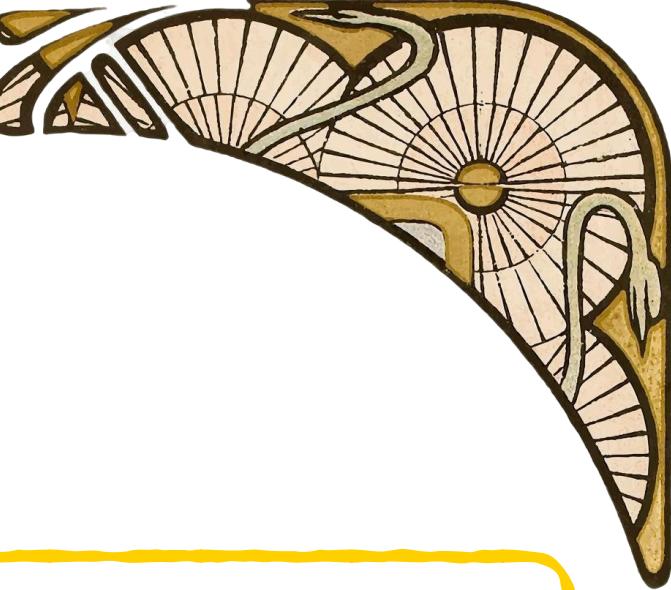


g and = Not Observable
 ∞ and = Not Observable



Parallel Trend Assumption

The Difference-in-Differences estimator (2/2)



Training Batch

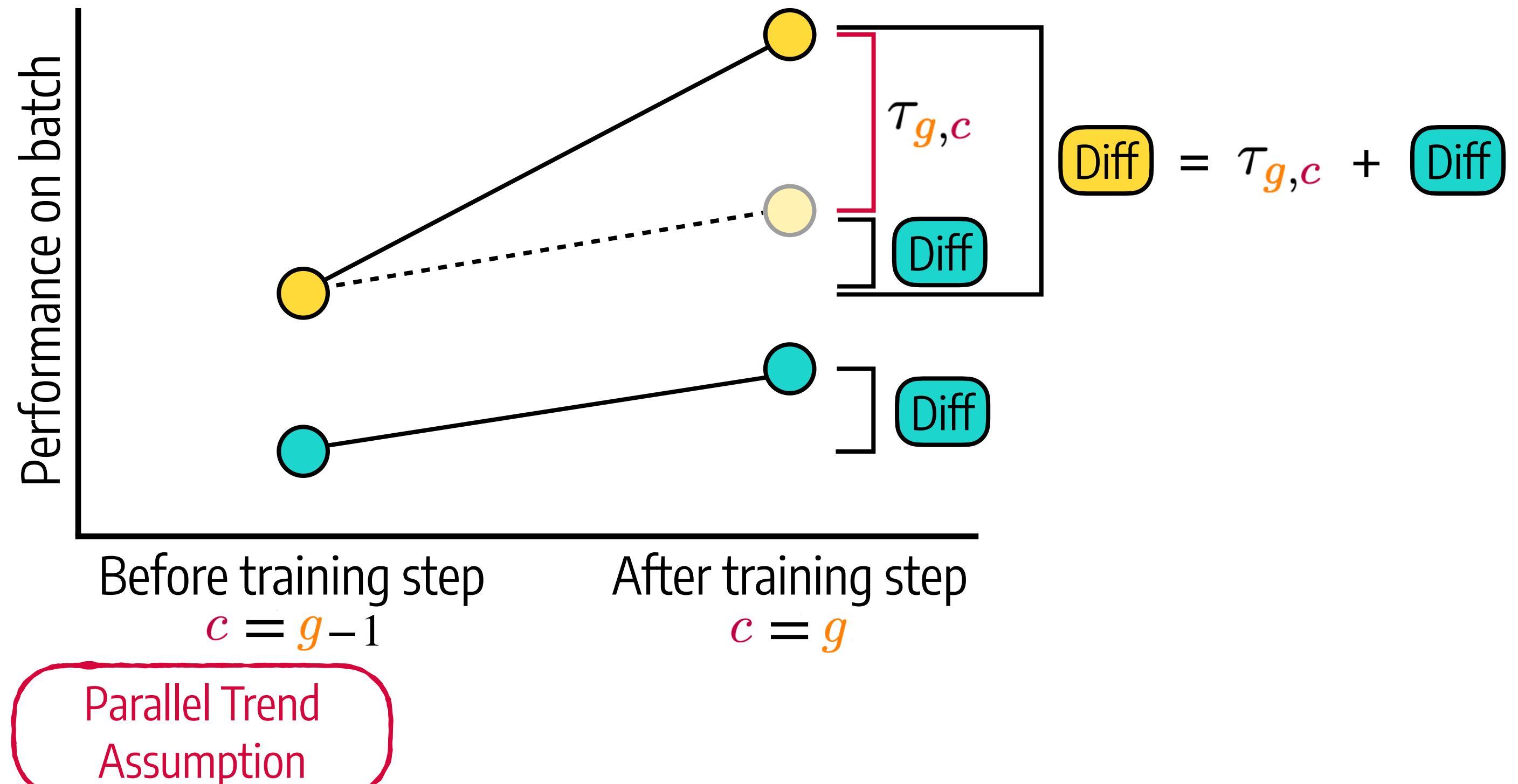
Validation Batch



g and = Observable
 ∞ and = Observable

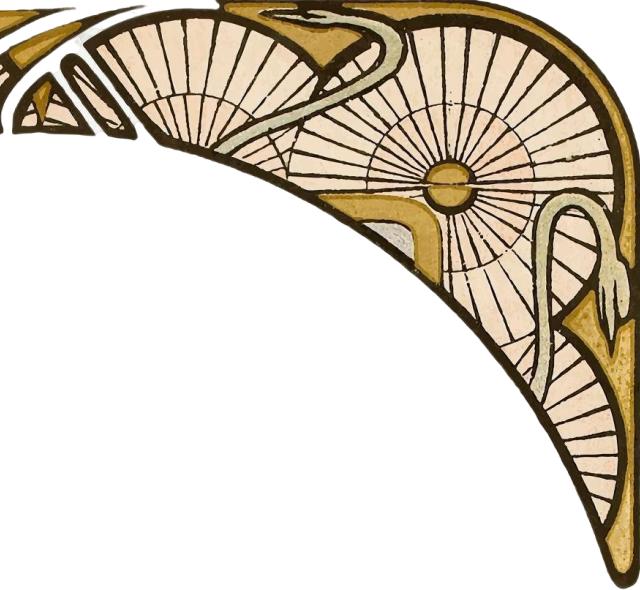


g and = Not Observable
 ∞ and = Not Observable



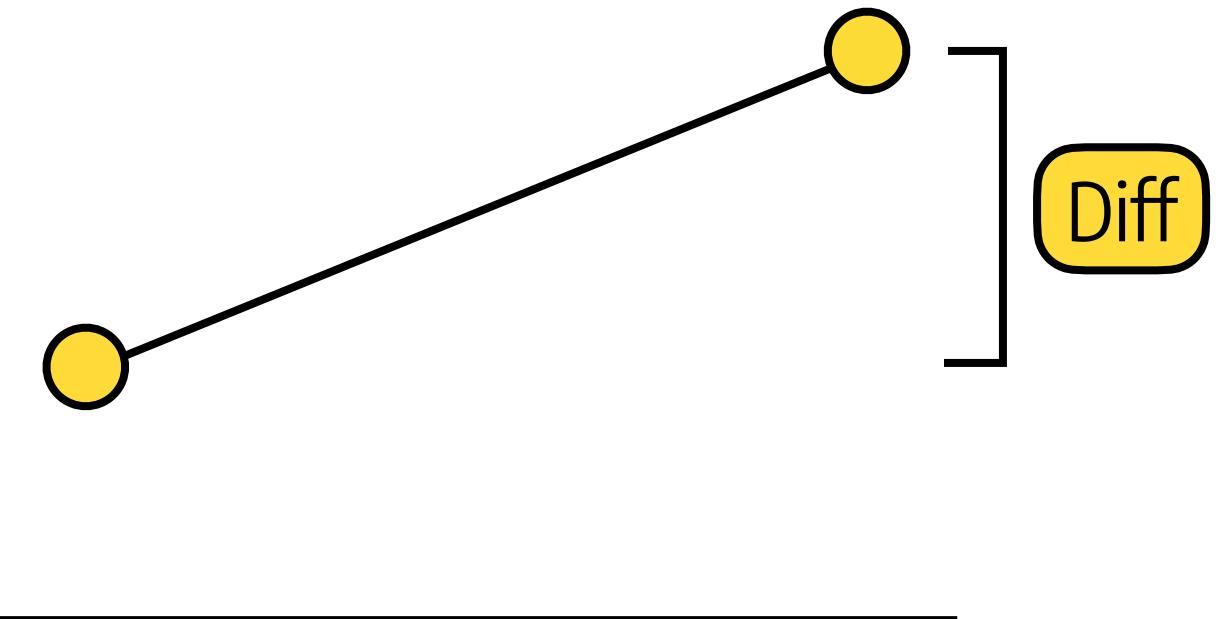
$$\tau_{g,c}^{\text{did}} = \text{Diff}_{\text{Training Batch}} - \text{Diff}_{\text{Validation Batch}}$$

Difference-in-Differences in practice



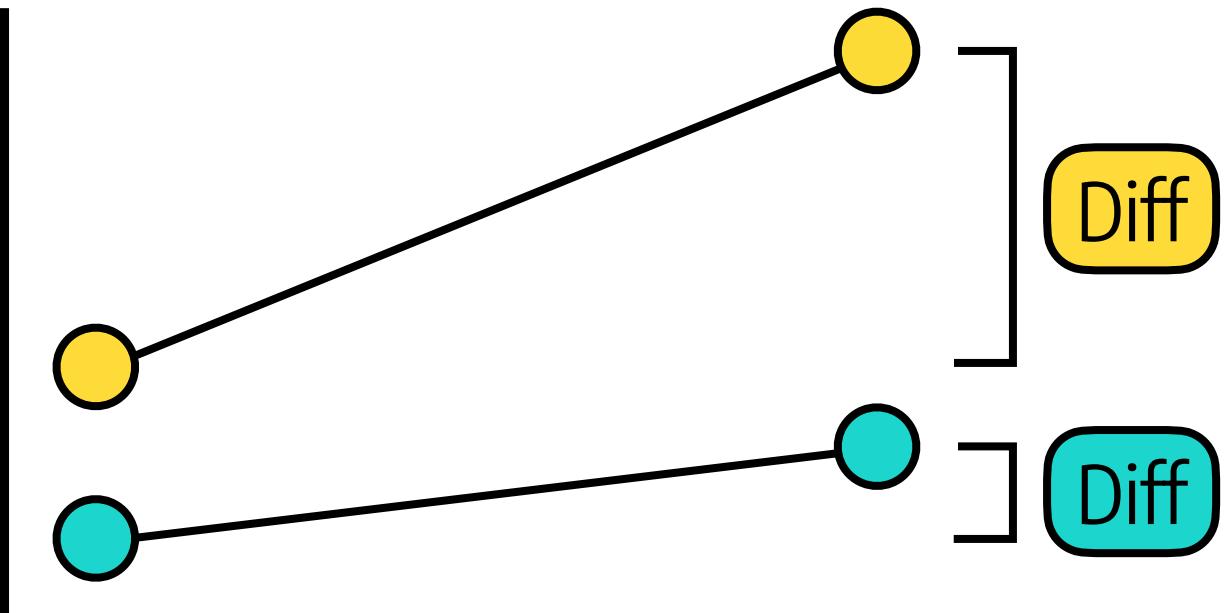
Difference-in-Differences in practice

$$\tau_{g,c}^{\text{did}} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) - Y_{g-1}(\mathbf{x}; g) | \quad] \quad \checkmark$$



Difference-in-Differences in practice

$$\tau_{g,c}^{\text{did}} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) - Y_{g-1}(\mathbf{x}; g) | \text{Yellow Box}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) - Y_{g-1}(\mathbf{x}; \infty) | \text{Blue Box}]$$



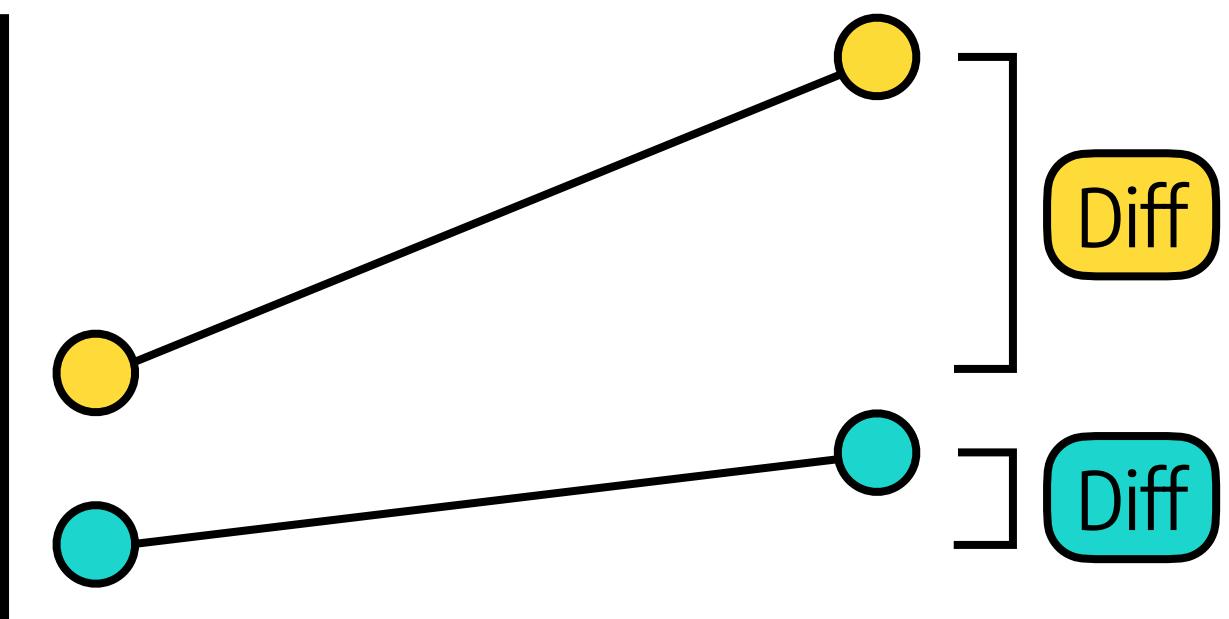
Difference-in-Differences in practice

$$\tau_{g,c}^{\text{did}} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) - Y_{g-1}(\mathbf{x}; g) | \boxed{\quad}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) - Y_{g-1}(\mathbf{x}; \infty) | \boxed{\quad}]$$



$$\hat{\tau}_{g,c}^{\text{did}} = \underbrace{(\bar{Y}_c(g) - \bar{Y}_{g-1}(g))}_{\text{Performance difference on training batches}}$$

Performance difference on training batches

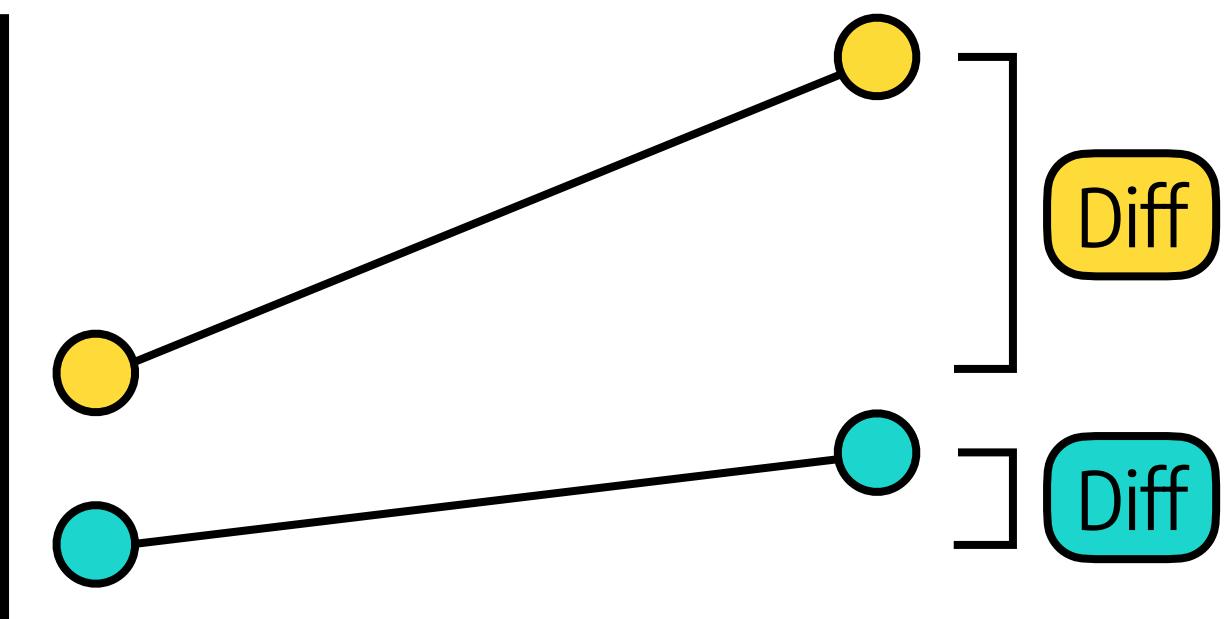


Difference-in-Differences in practice

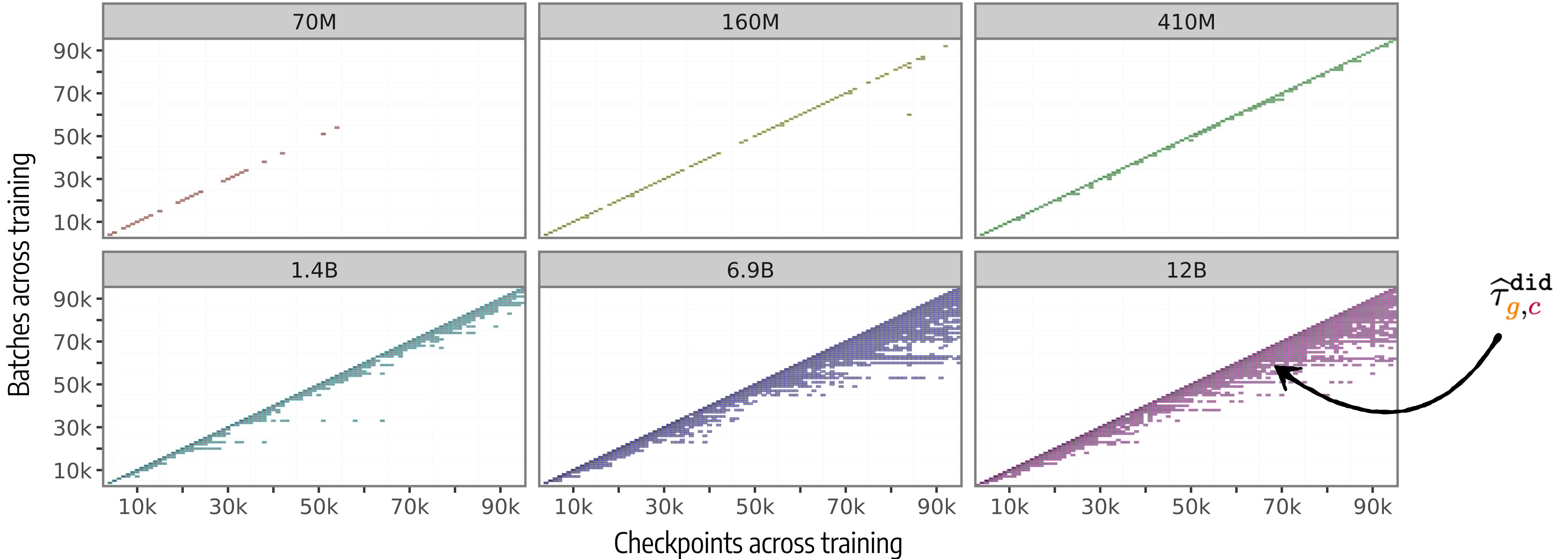
$$\tau_{g,c}^{\text{did}} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) - Y_{g-1}(\mathbf{x}; g) | \text{Yellow Box}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) - Y_{g-1}(\mathbf{x}; \infty) | \text{Cyan Box}]$$

\downarrow \downarrow

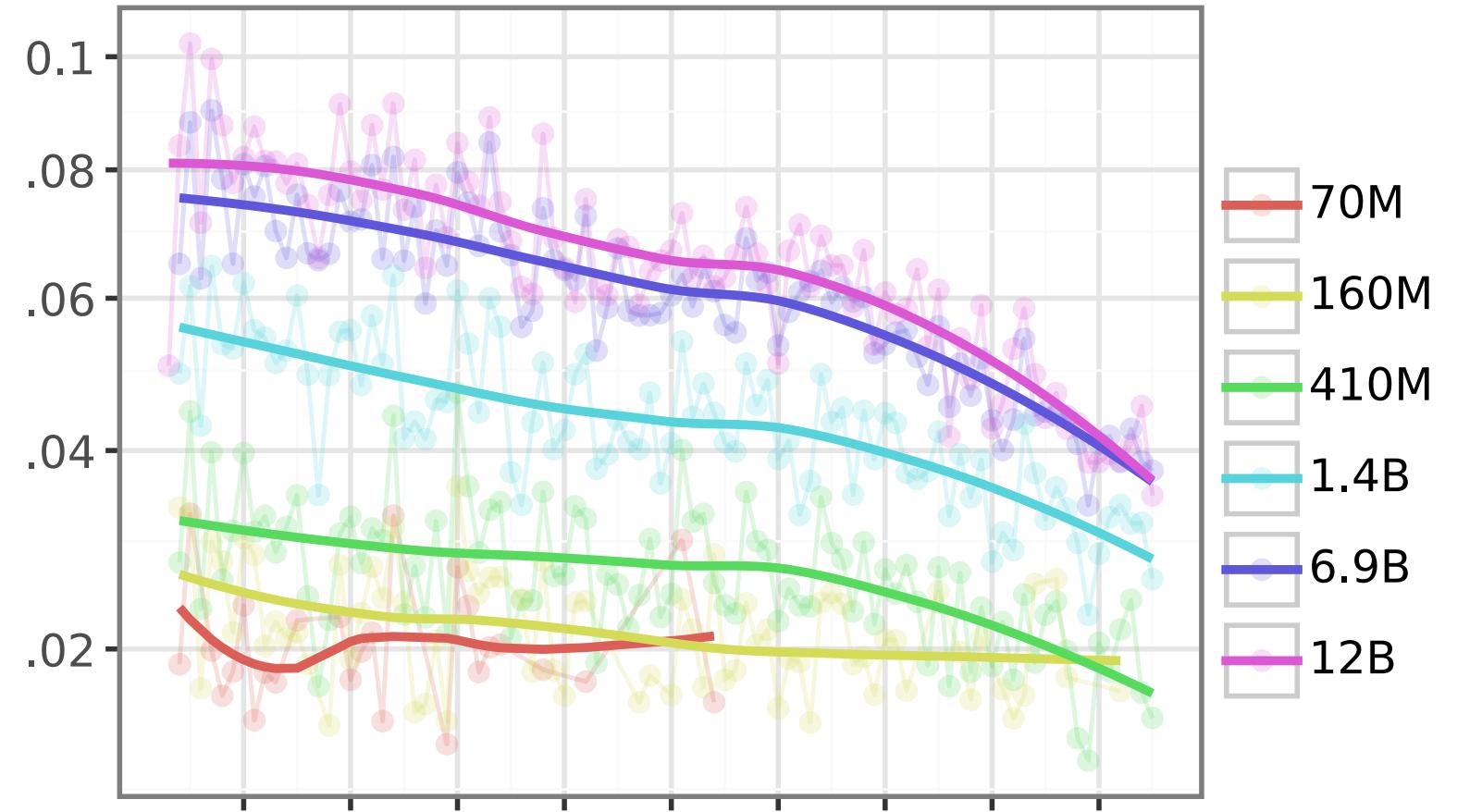
$$\hat{\tau}_{g,c}^{\text{did}} = \underbrace{(\bar{Y}_c(g) - \bar{Y}_{g-1}(g))}_{\text{Performance difference on training batches}} - \underbrace{(\bar{Y}_c(\infty) - \bar{Y}_{g-1}(\infty))}_{\text{Performance difference on validation batches}}$$



Constructing the memorisation profile

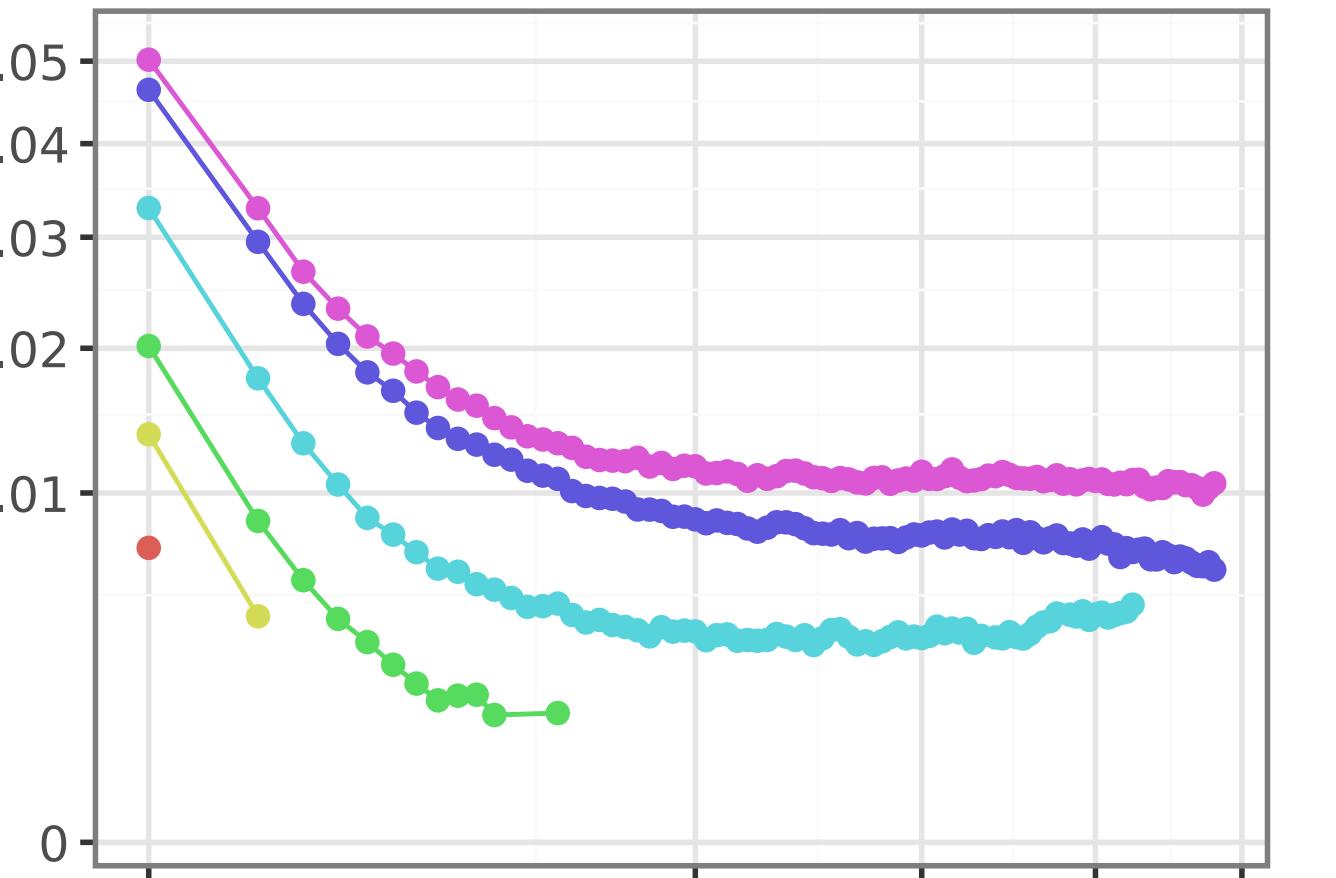


Aggregating the memorisation profile



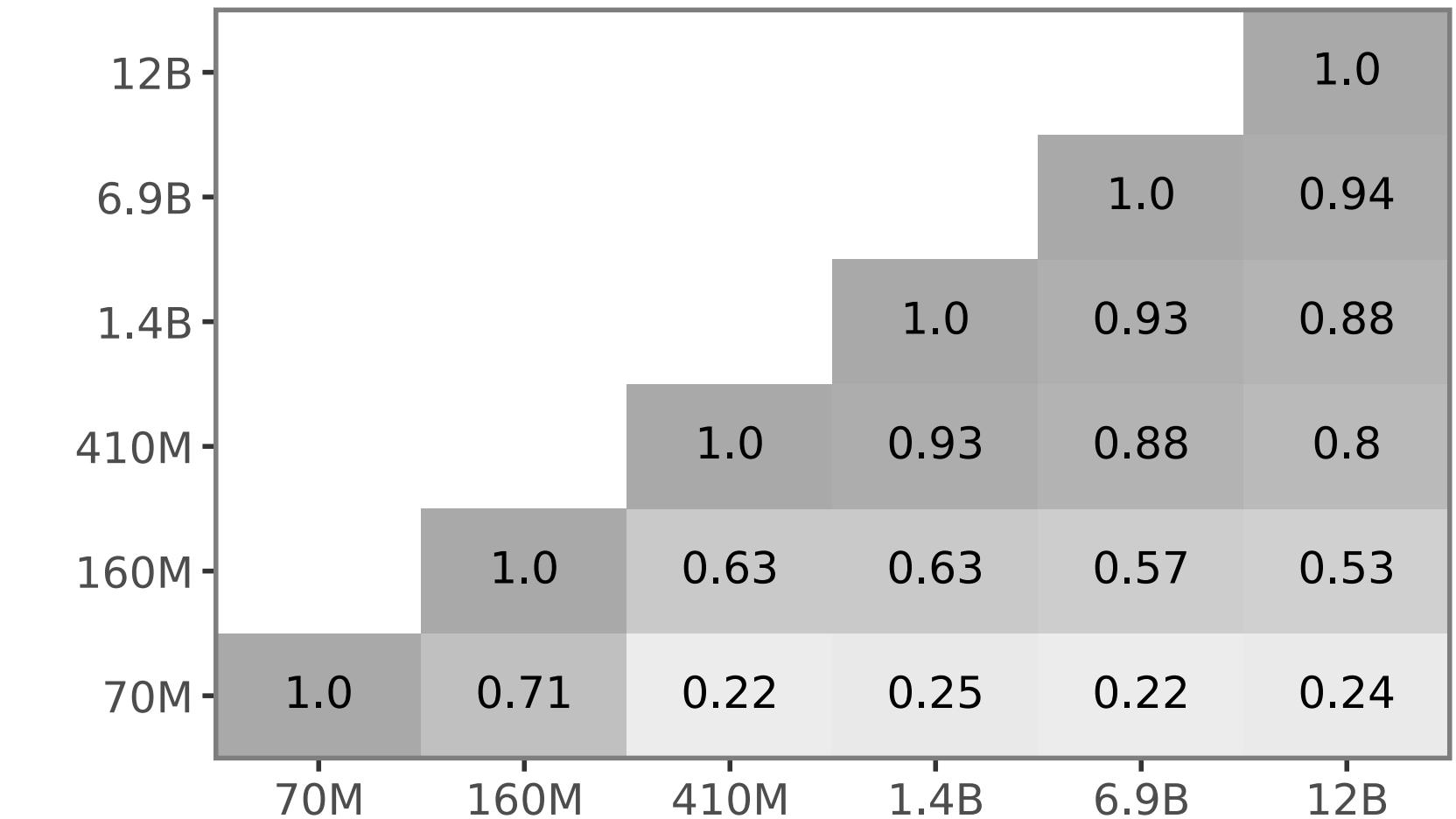
Checkpoints/Batches across training ($g = c$)

Instantaneous memorisation of a batch
at the checkpoint is first seen



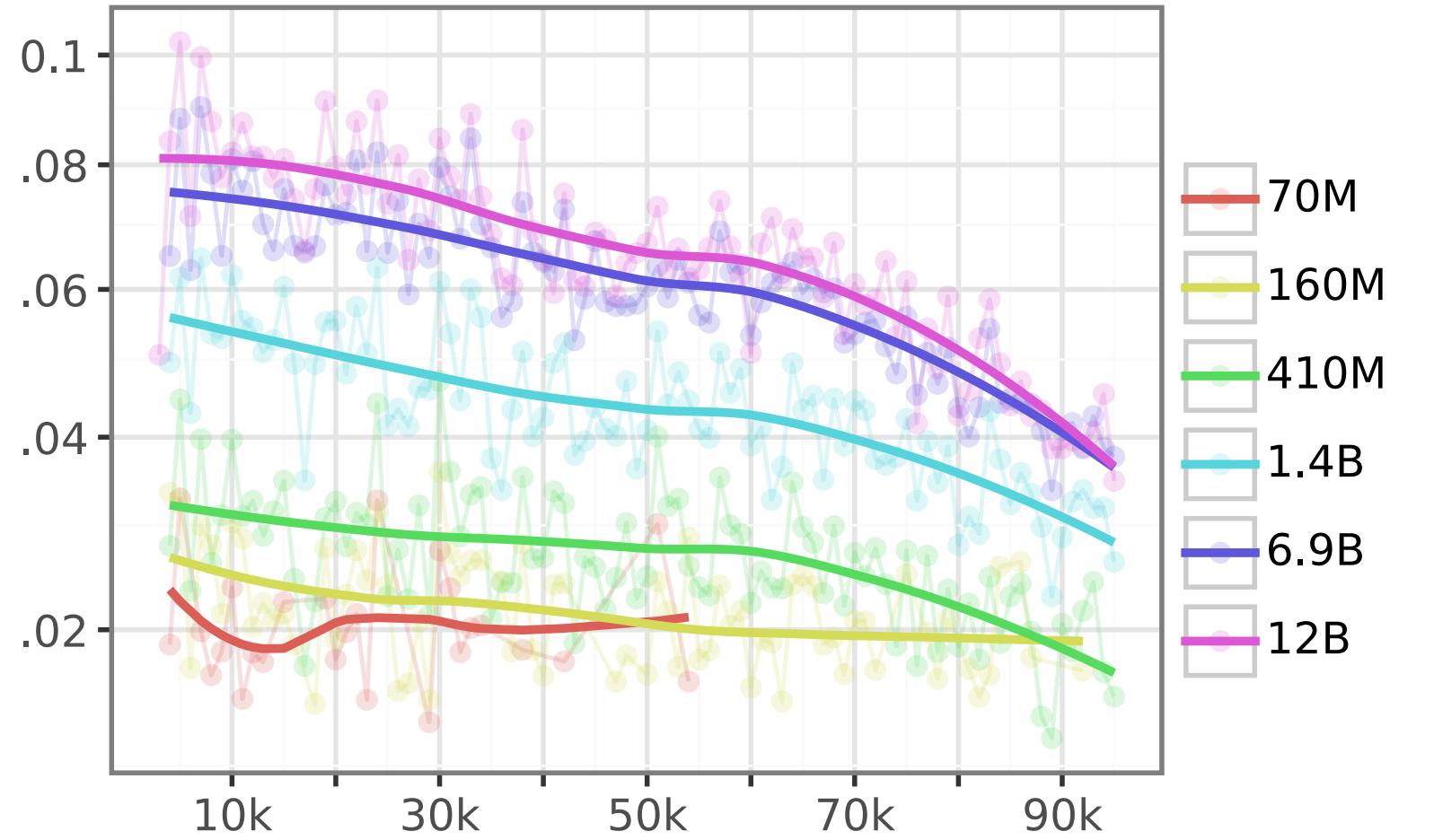
Checkpoints after batch is seen ($g > c$)

Average persistent memorisation of a
batch per step after it has been seen*
*(only averaging across batches that have been
trained that much)



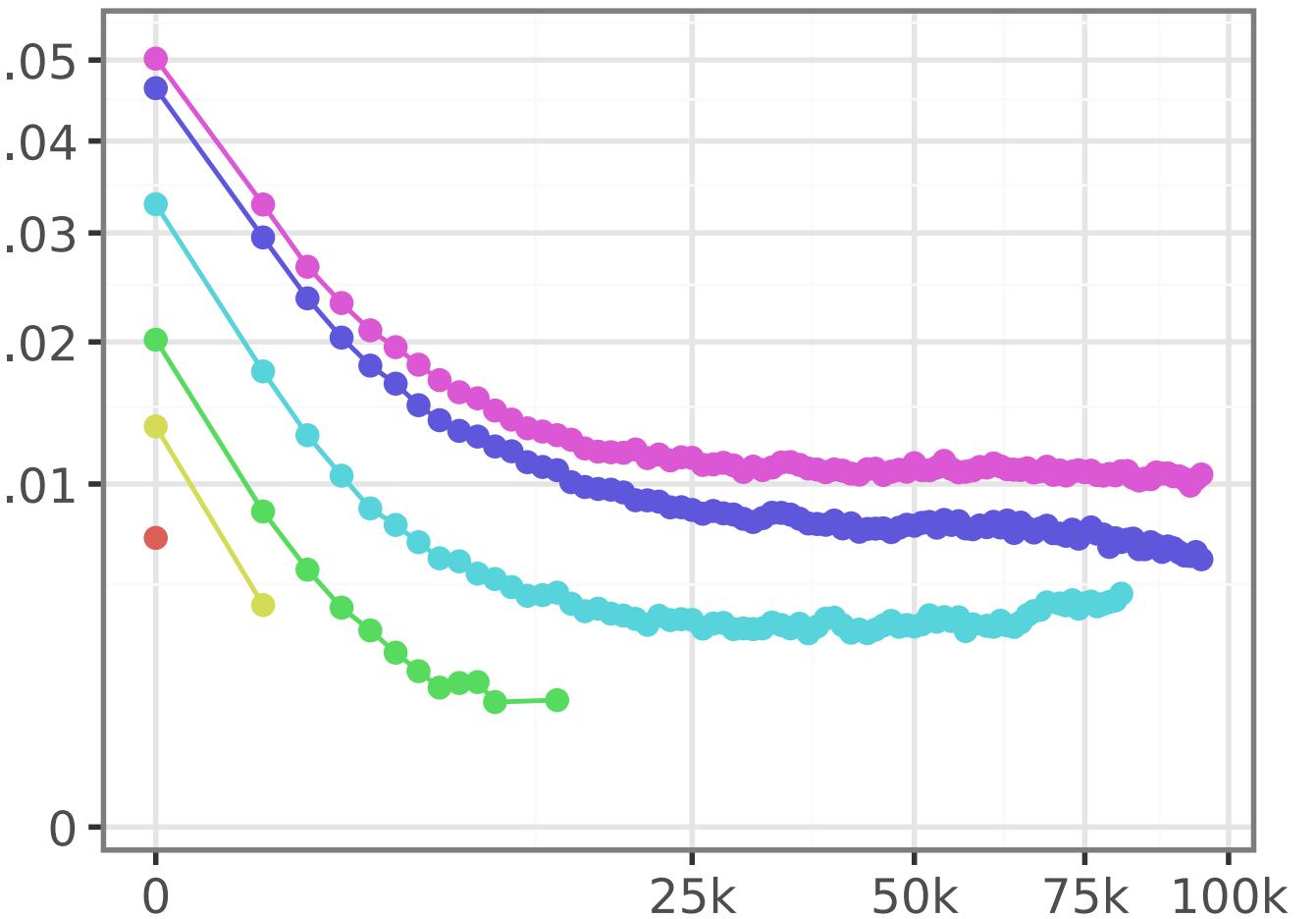
Pearson correlation between the
memorisation profiles of different
model sizes

Aggregating the memorisation profile



Checkpoints/Batches across training ($g = c$)

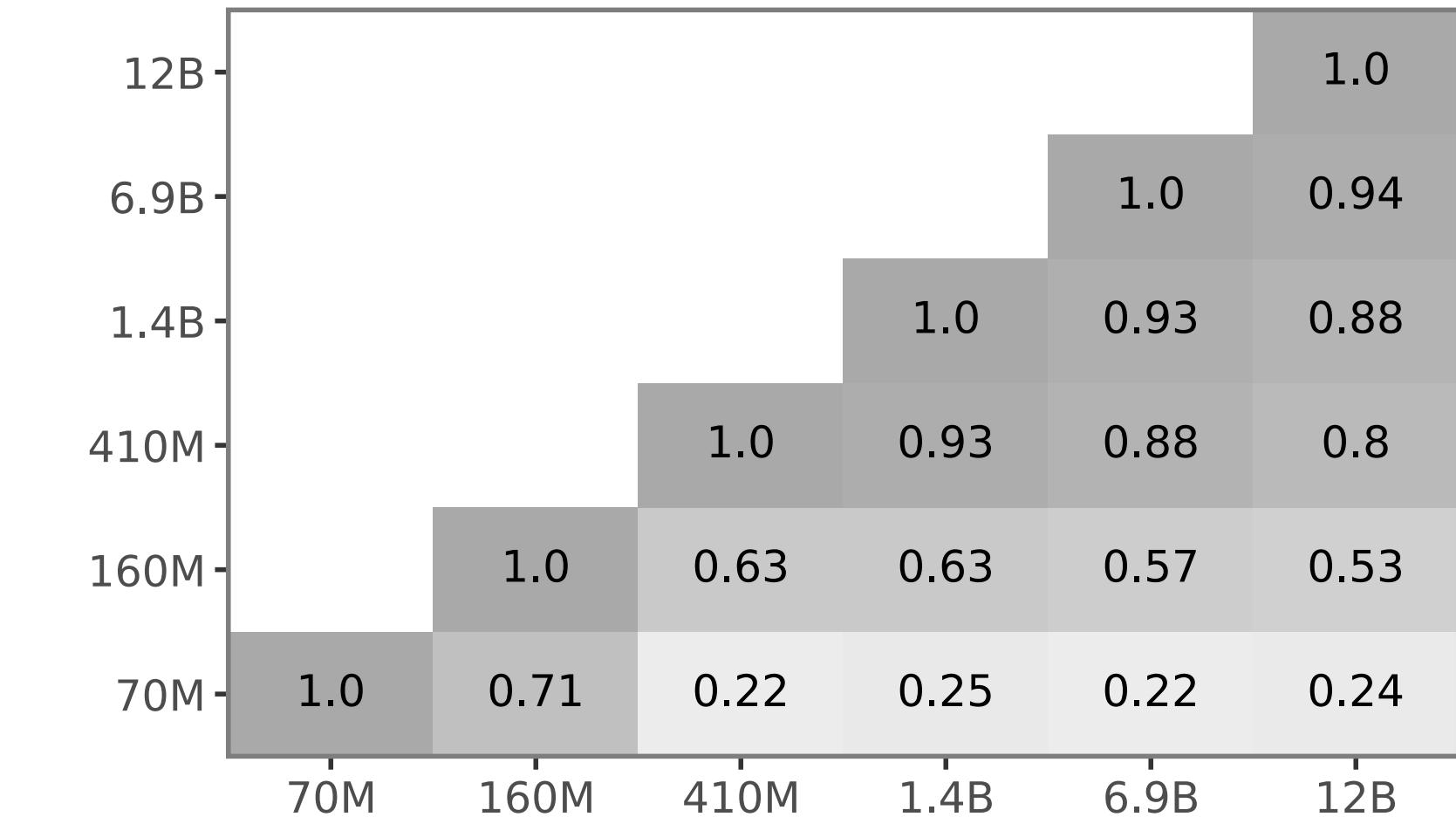
Instantaneous memorisation of a batch
at the checkpoint is first seen



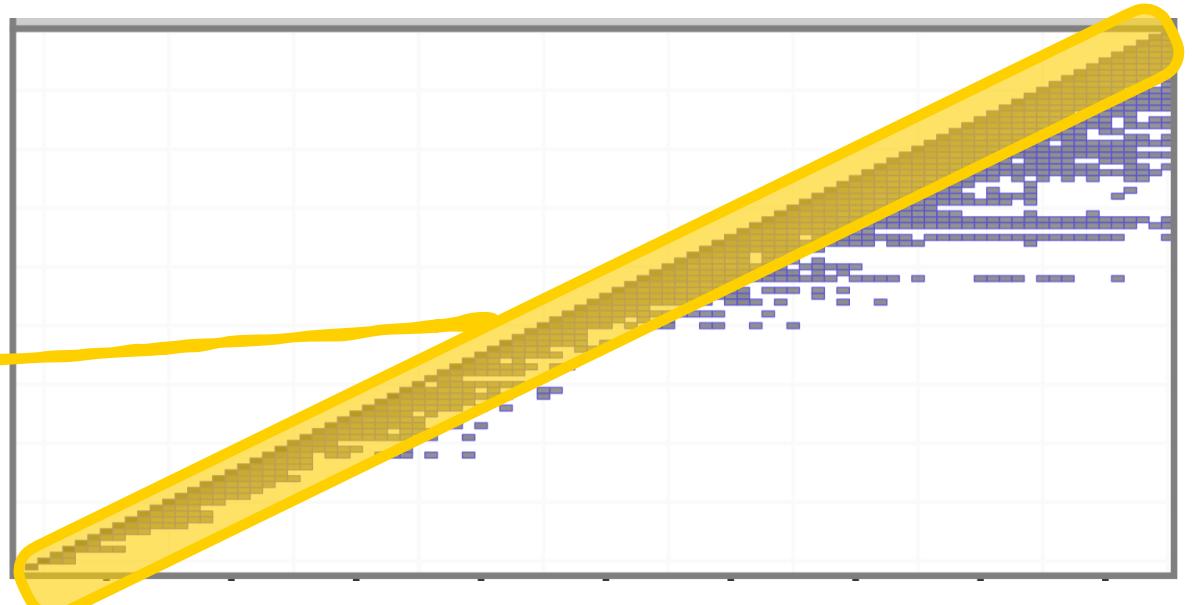
Checkpoints after batch is seen ($g > c$)

Average persistent memorisation of a
batch per step after it has been seen*

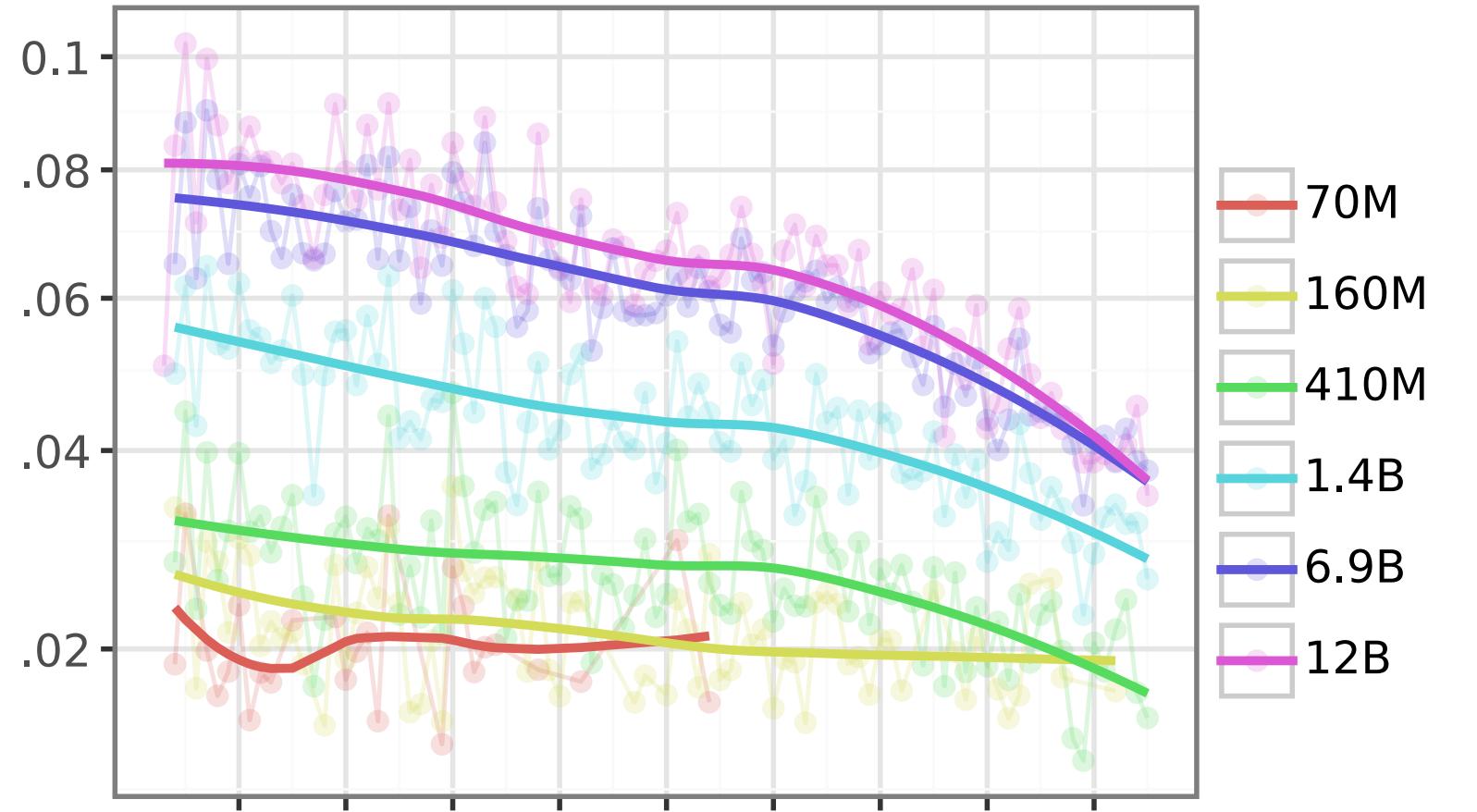
*(only averaging across batches that have been
trained that much)



Pearson correlation between the
memorisation profiles of different
model sizes

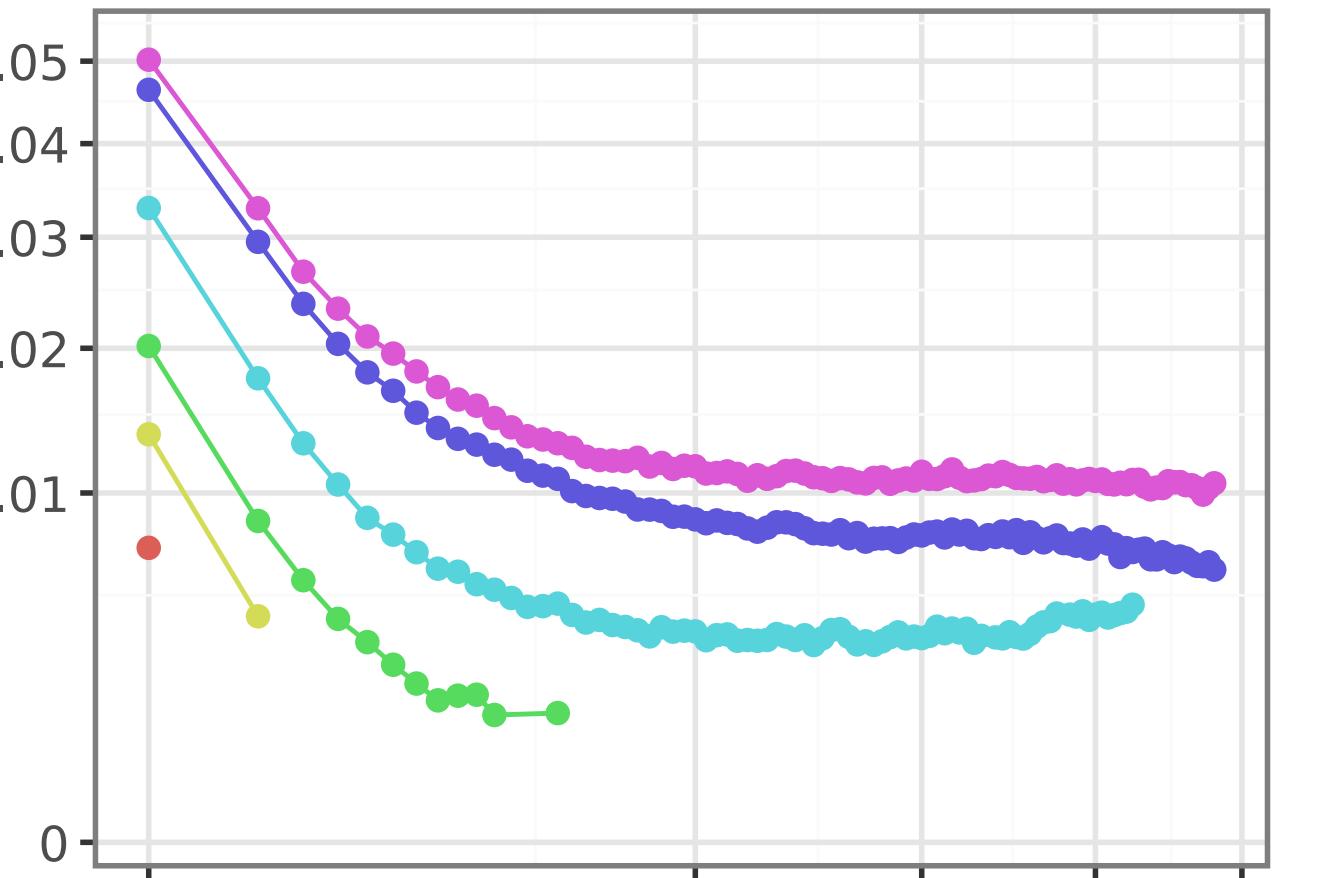


Aggregating the memorisation profile



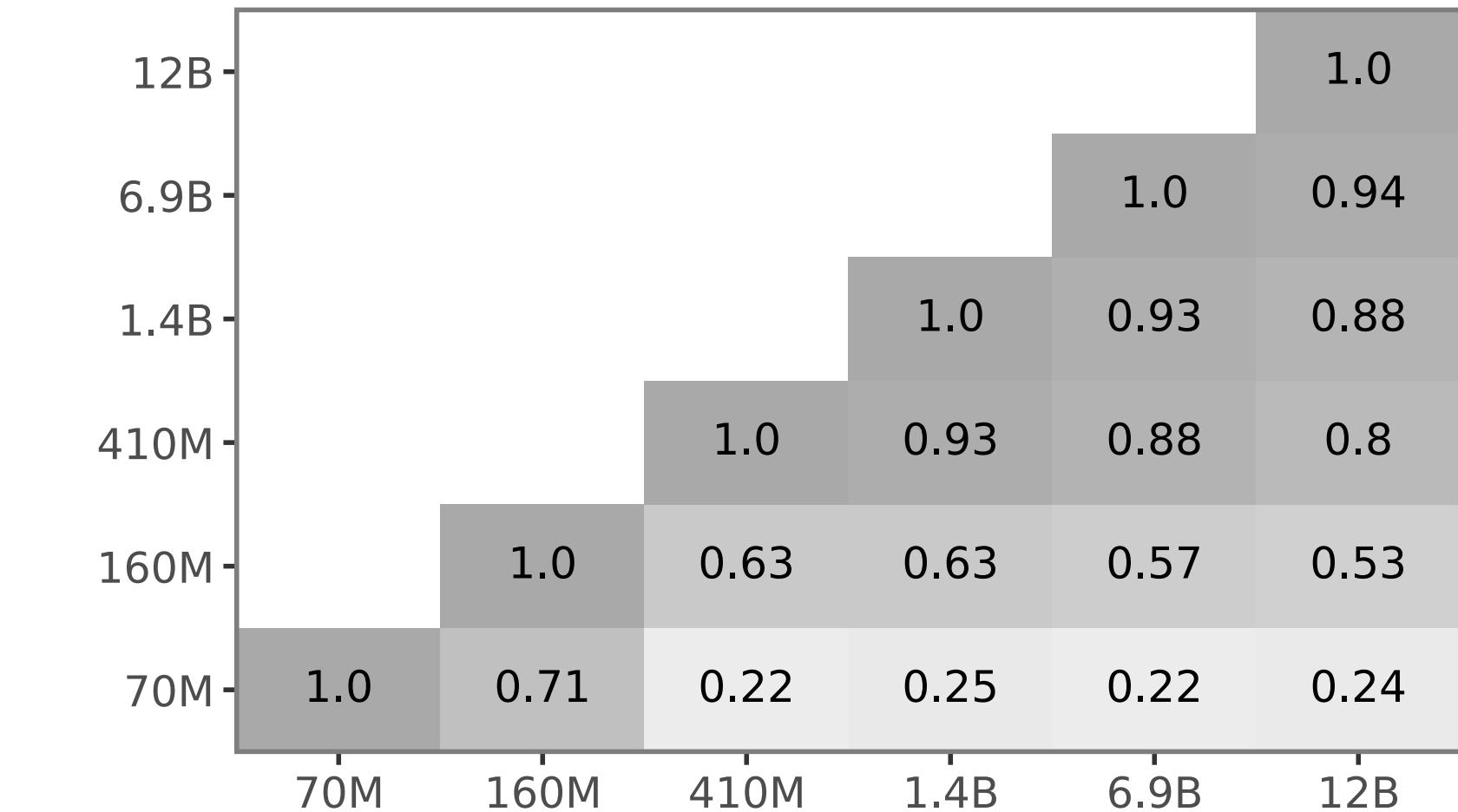
Checkpoints/Batches across training ($g = c$)

Instantaneous memorisation of a batch
at the checkpoint is first seen

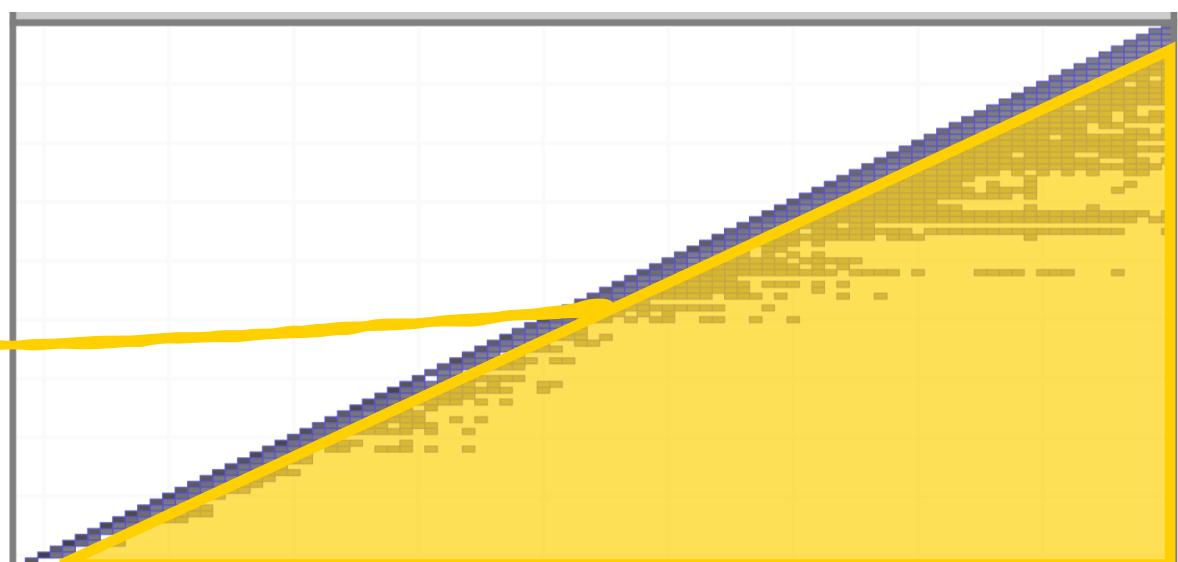


Checkpoints after batch is seen ($g > c$)

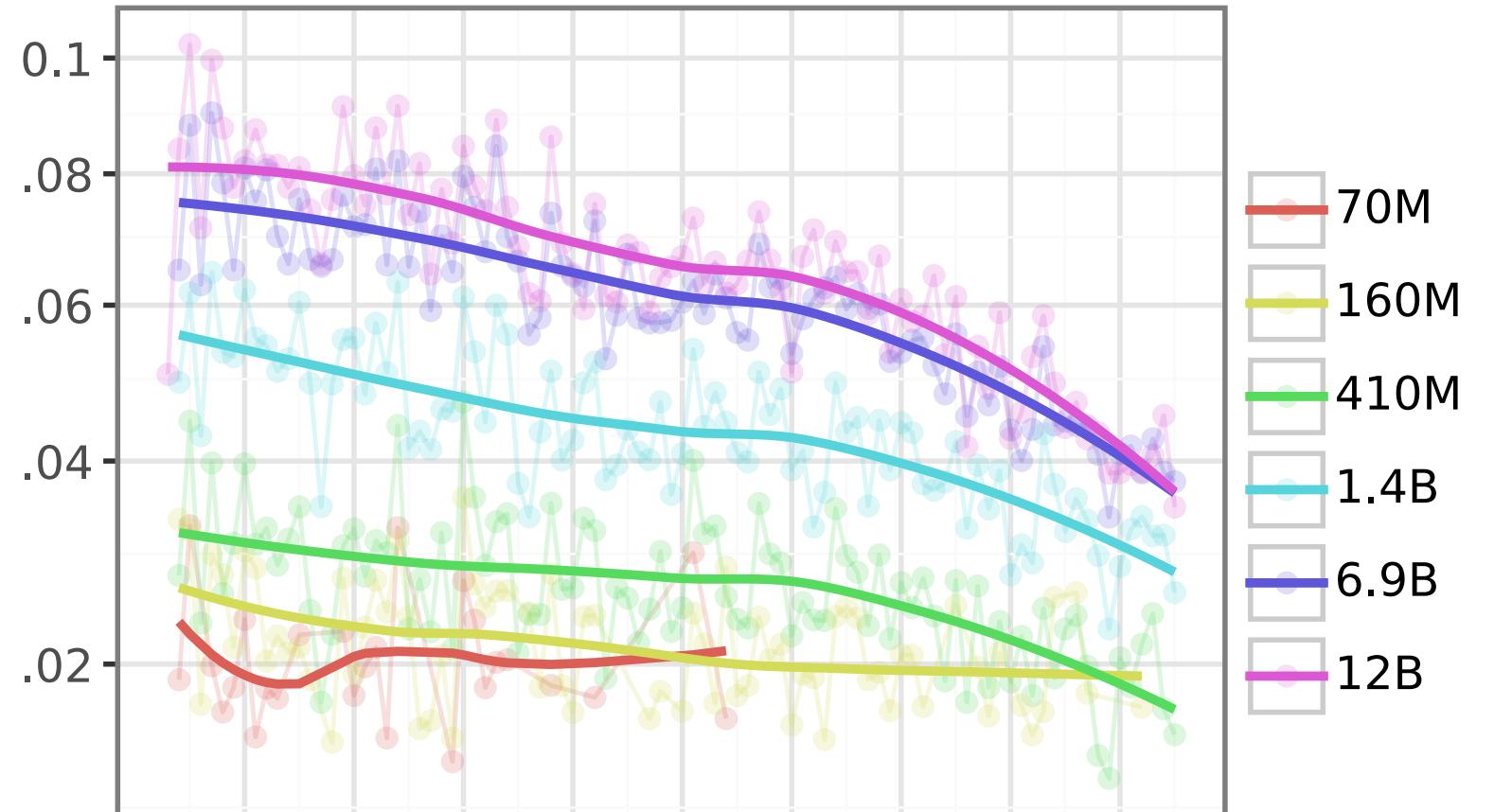
Average persistent memorisation of a
batch per step after it has been seen*
*(only averaging across batches that have been
trained that much)



Pearson correlation between the
memorisation profiles of different
model sizes

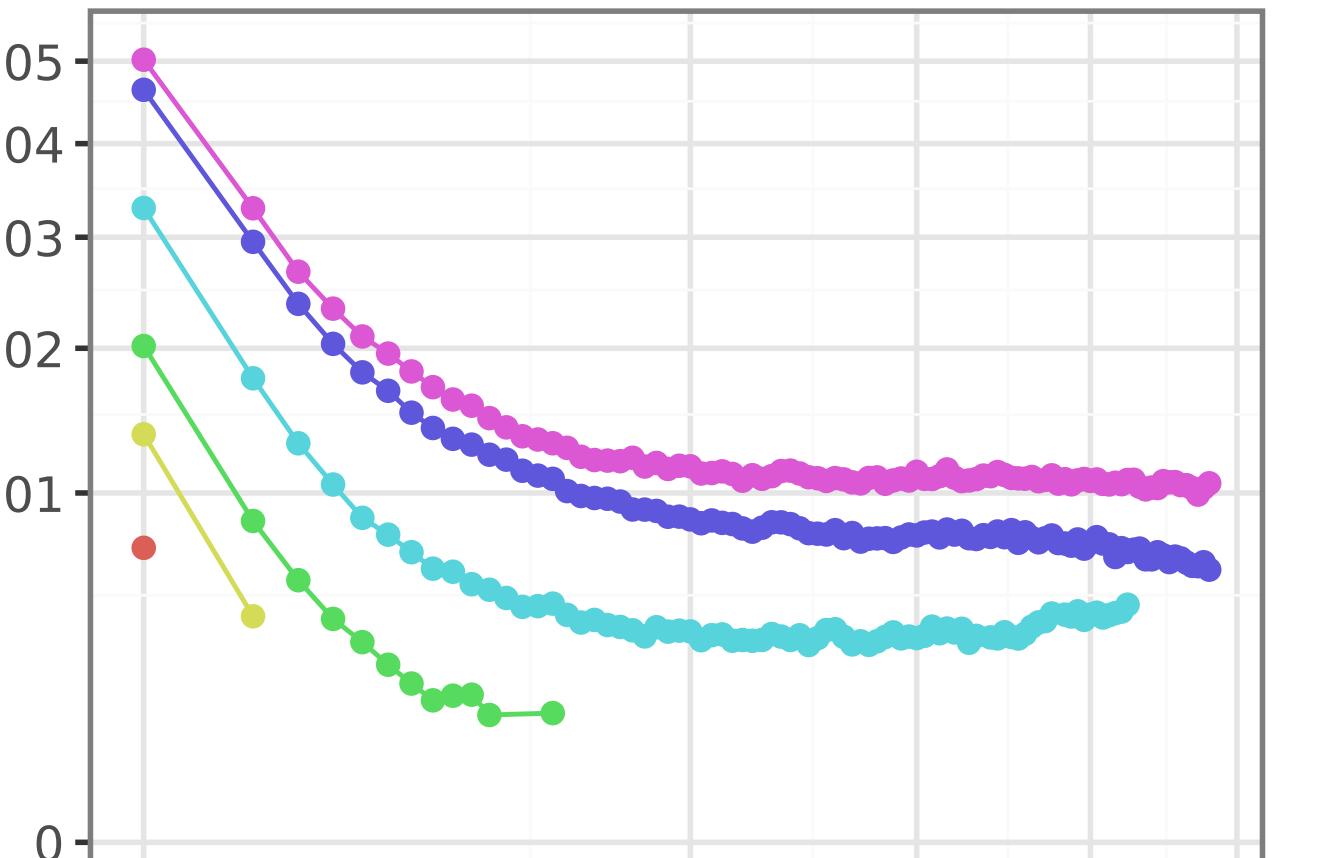


Aggregating the memorisation profile



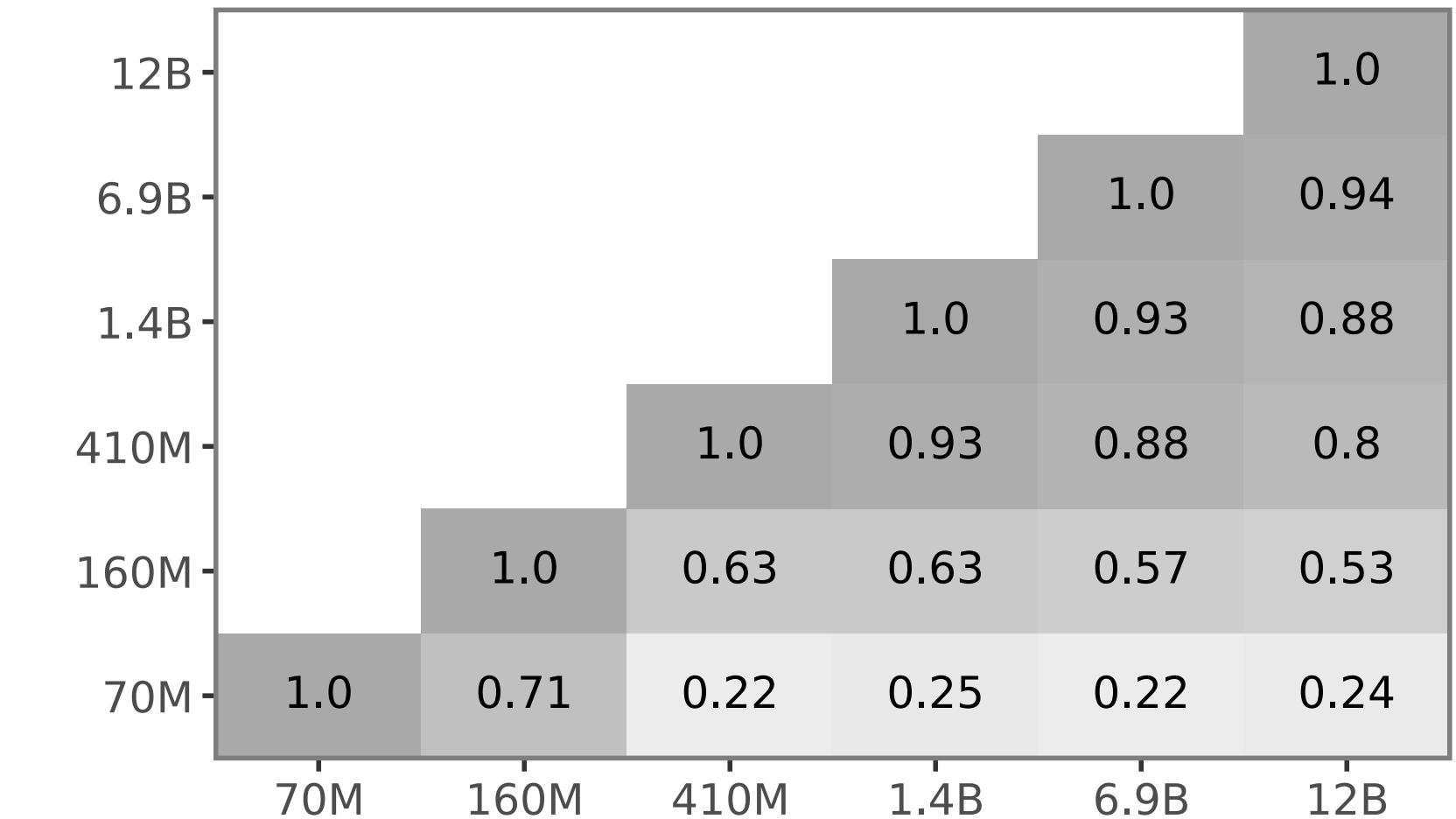
Checkpoints/Batches across training ($g = c$)

Instantaneous memorisation of a batch
at the checkpoint is first seen



Checkpoints after batch is seen ($g > c$)

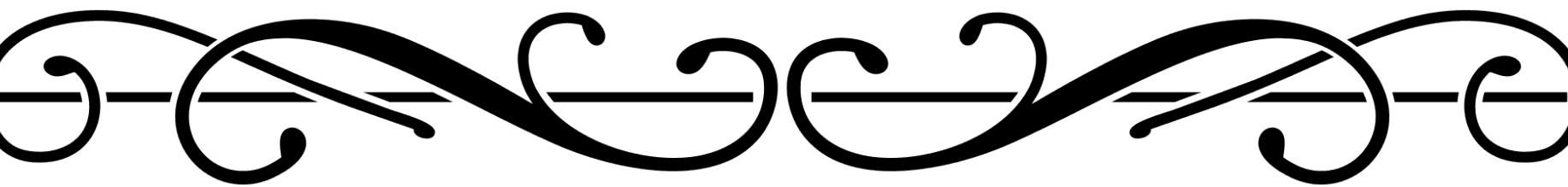
Average persistent memorisation of a
batch per step after it has been seen*
*(only averaging across batches that have been
trained that much)



Pearson correlation between the
memorisation profiles of different
model sizes



FINE



Pietro Lesci
University of Cambridge

X: @pietro_lesci
LinkedIn: /pietrolesci
Page: pietrolesci.github.io
Mail: pietrolesci@outlook.com