

NATURAL EXPERIMENTS IN NLP AND WHERE TO FIND THEM



Pietro Lesci
University of Cambridge

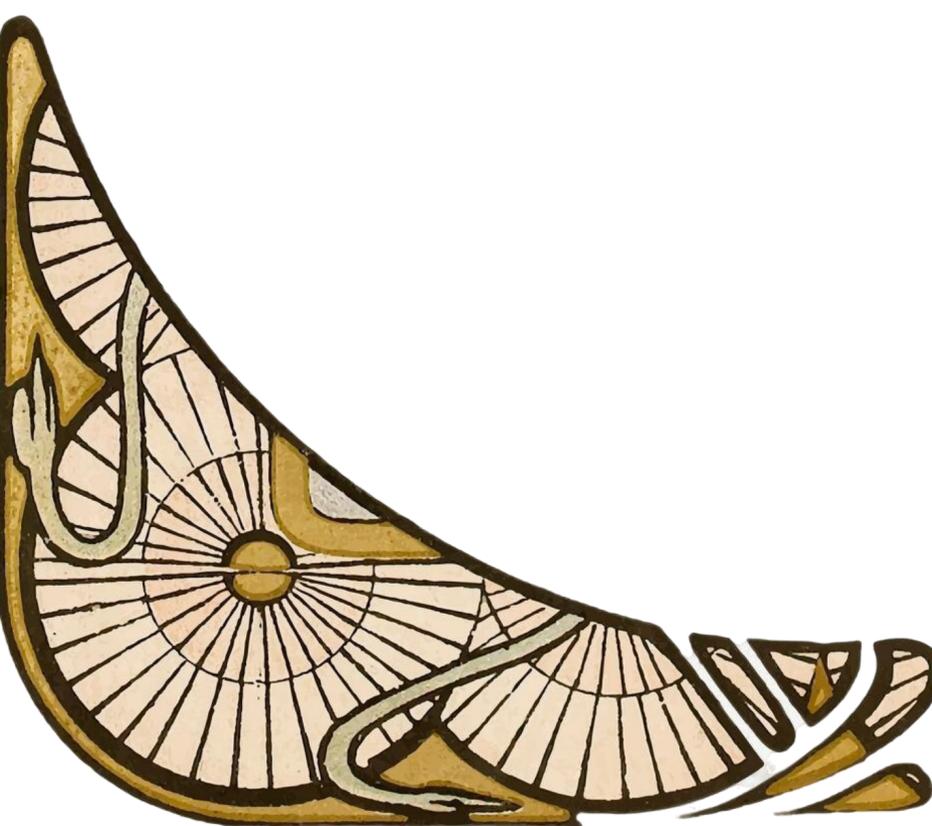


X: @pietro_lesci

LinkedIn: /pietrolesci

Page: pietrolesci.com

Mail: pietrolesci@outlook.com



CAUSAL ESTIMATION OF MEMORISATION PROFILES



ACL 2024 (Best Paper Award 🏆)



Pietro Lesci



Clara Meister



Thomas Hofmann

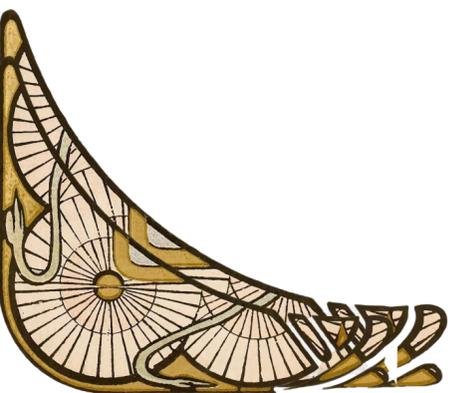
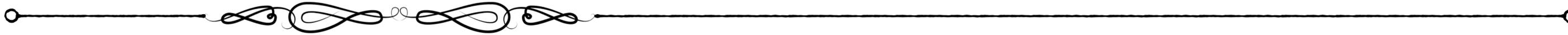
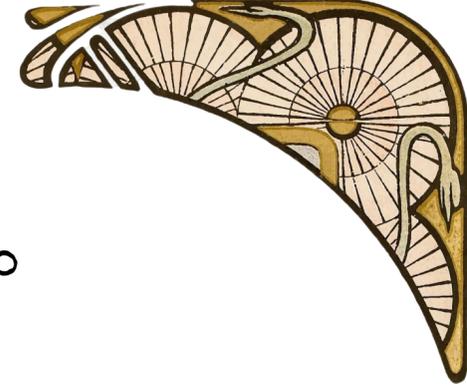


Andreas Vlachos



Tiago Pimentel

Why study memorisation?



Why study memorisation?

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



Why study memorisation?



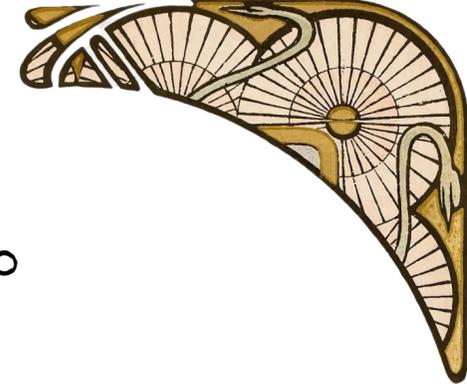
Language models can reproduce entire sequences from their training set verbatim (Carlini, 2021; 2023).

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).

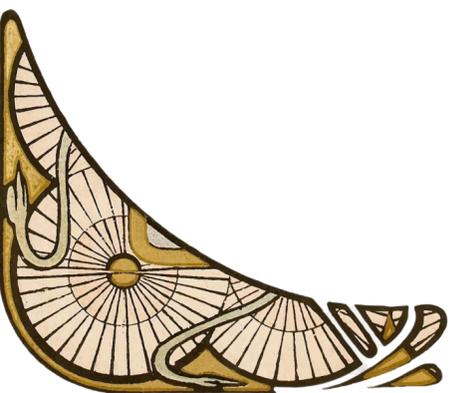
Memorisation has implications for:

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

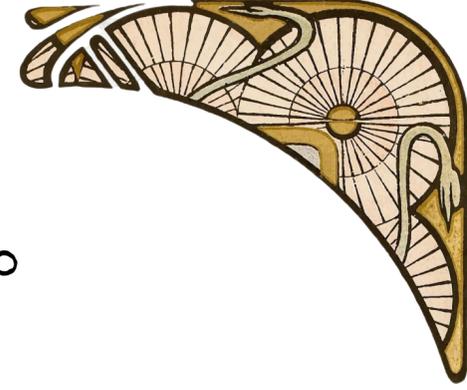
In order to be able to study memorisation we need methods to quickly, cheaply, and accurately measure it!



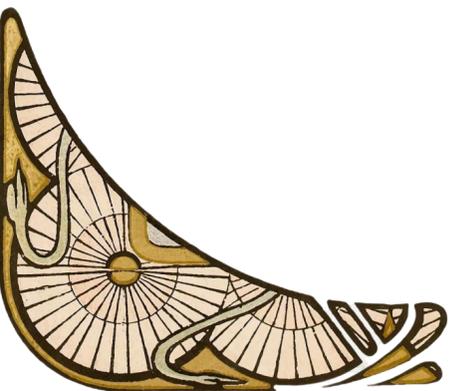
Tl;dr: Estimate memorisation directly from the data!



Tl;dr: Estimate memorisation directly from the data!



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

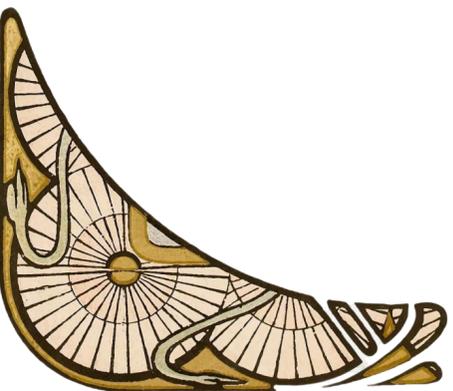


Tl;dr: Estimate memorisation directly from the data!

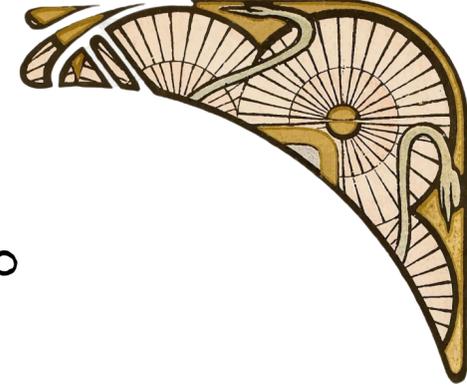


We want to estimate **counterfactual memorisation**
(as defined in Feldman, 2020):

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

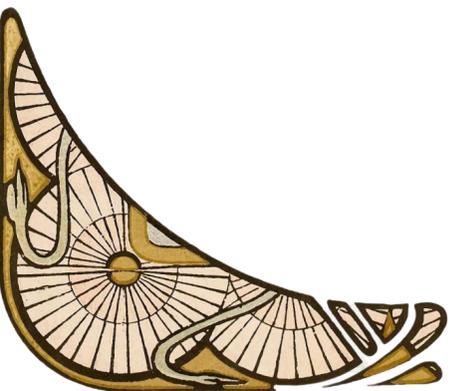


Tl;dr: Estimate memorisation directly from the data!

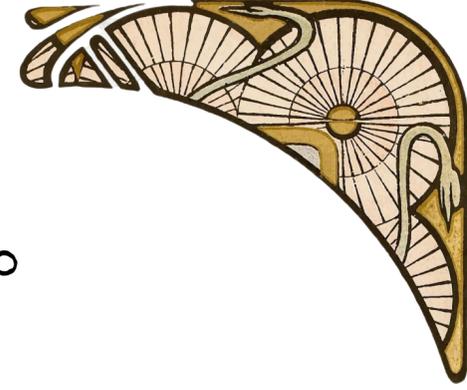


We want to estimate **counterfactual memorisation**
(as defined in Feldman, 2020):

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



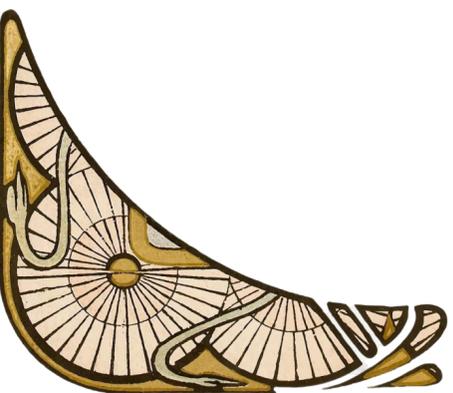
Tl;dr: Estimate memorisation directly from the data!



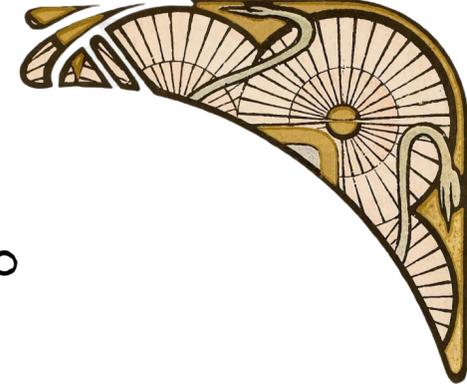
We want to estimate **counterfactual memorisation**
(as defined in Feldman, 2020):

“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



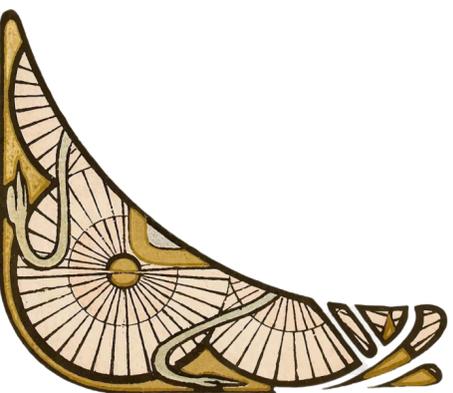
Tl;dr: Estimate memorisation directly from the data!



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)!



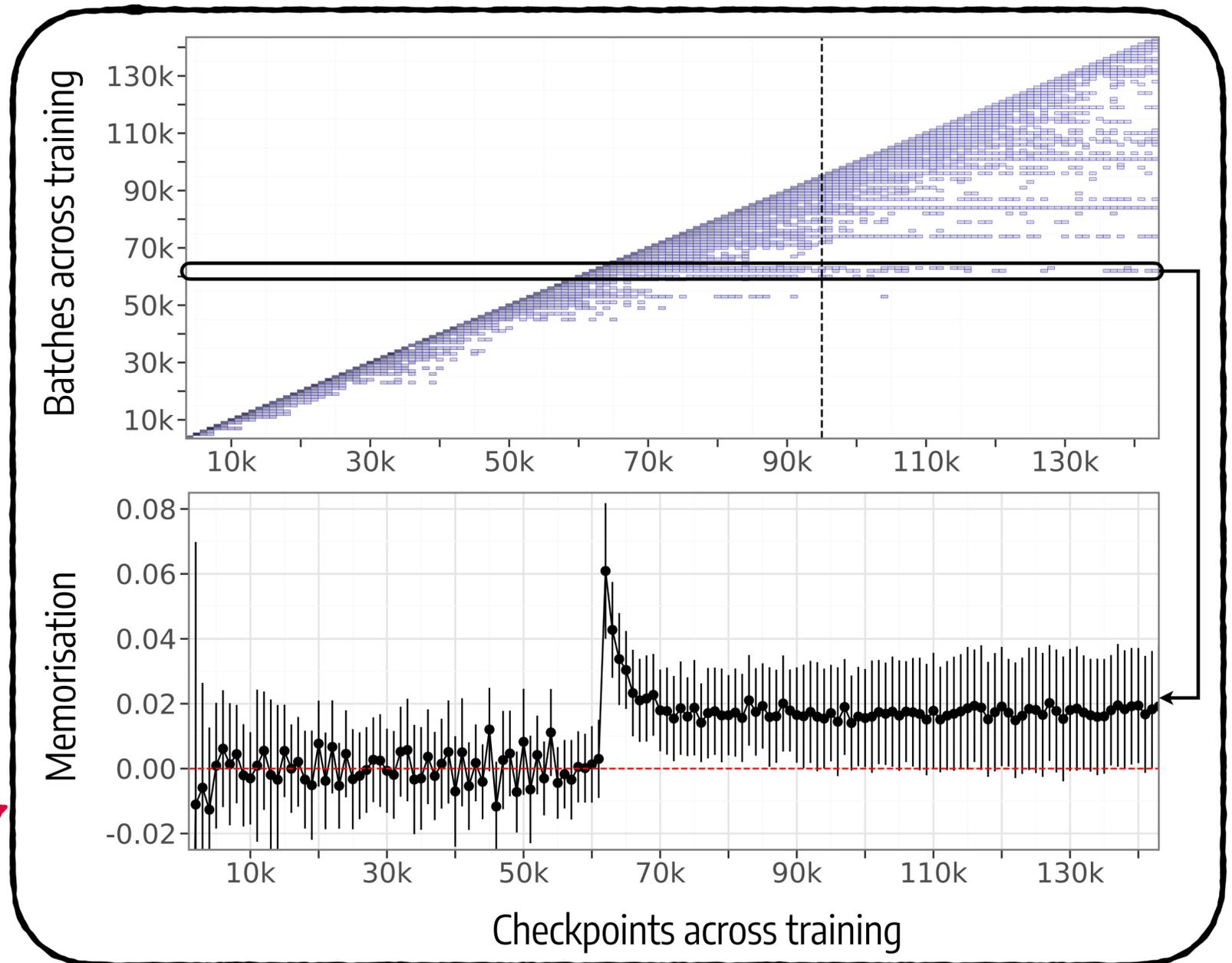
Tl;dr: Estimate memorisation directly from the data!

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

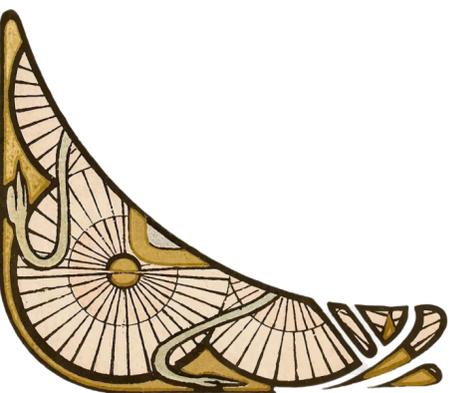
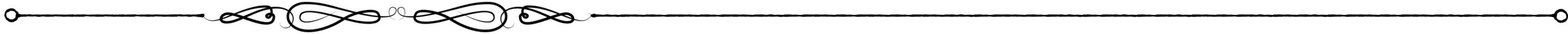
“The causal effect of observing an instance during training on a model’s ability to correctly predict that instance”

- 🤔 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



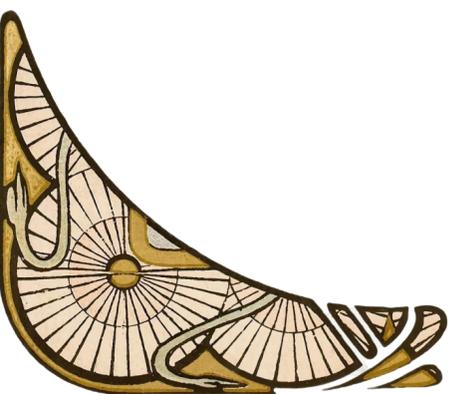
Memorisation as a difference of potential outcomes



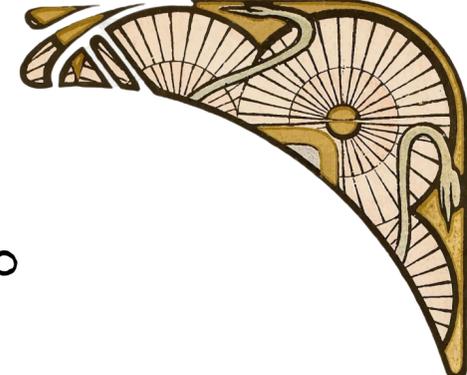
Memorisation as a difference of potential outcomes

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

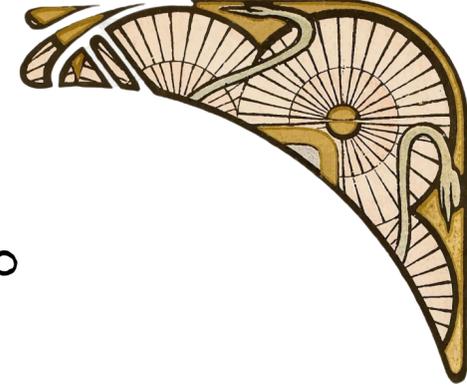
Memorisation as a difference of potential outcomes


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


Memorisation as a difference of potential outcomes

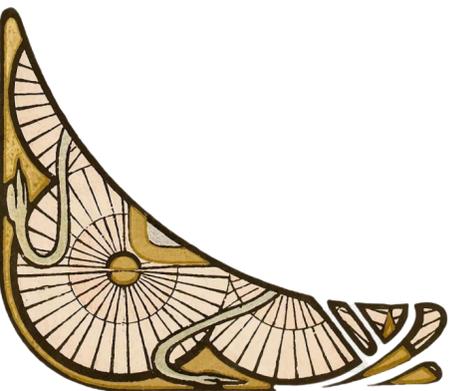

$$\tau_{\mathbf{x},c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)}_{\downarrow} - \underbrace{Y_c(\mathbf{x}; \infty)}$$


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

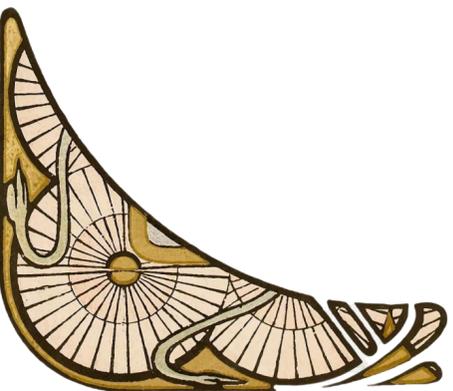
$$\tau_{\mathbf{x},c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)} - \underbrace{Y_c(\mathbf{x}; \infty)}$$

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

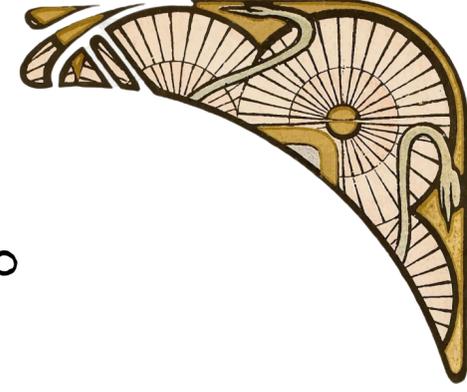
Memorisation as a difference of potential outcomes



$$\tau_{\mathbf{x},c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)}_{\text{Performance of model on } \mathbf{x} \text{ when trained with } \mathbf{x} \text{ at step } g} - \underbrace{Y_c(\mathbf{x}; \infty)}_{\text{Performance of model on } \mathbf{x} \text{ when } \mathbf{not} \text{ trained with } \mathbf{x}}$$

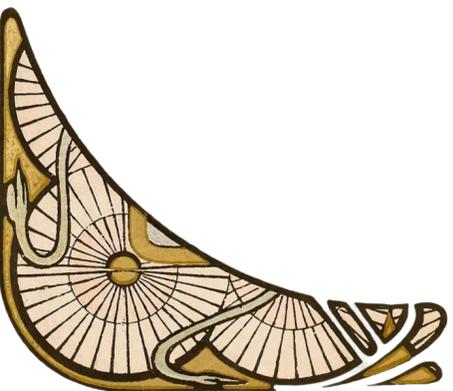


Memorisation as a difference of potential outcomes



$$\tau_{\mathbf{x},c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)}_{\text{Performance of model on } \mathbf{x} \text{ when trained with } \mathbf{x} \text{ at step } g} - \underbrace{Y_c(\mathbf{x}; \infty)}_{\text{Performance of model on } \mathbf{x} \text{ when **not** trained with } \mathbf{x}}$$

\uparrow

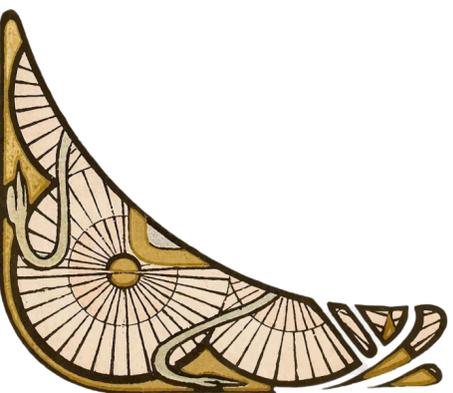


Memorisation as a difference of potential outcomes

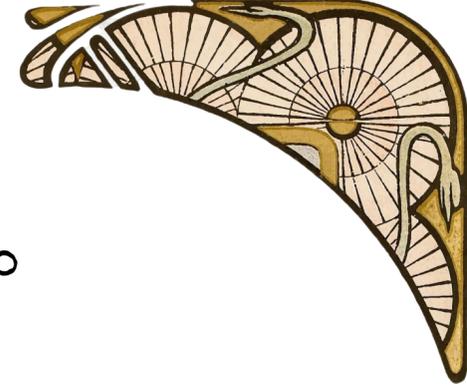
$$\tau_{\mathbf{x},c} \stackrel{\text{def}}{=} \underbrace{Y_c(\mathbf{x}; g)}_{\text{Performance of model on } \mathbf{x} \text{ when trained with } \mathbf{x} \text{ at step } g} - \underbrace{Y_c(\mathbf{x}; \infty)}_{\text{Performance of model on } \mathbf{x} \text{ when **not** trained with } \mathbf{x}}$$

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

Estimating the counterfactual outcome $Y_c(x; \infty)$

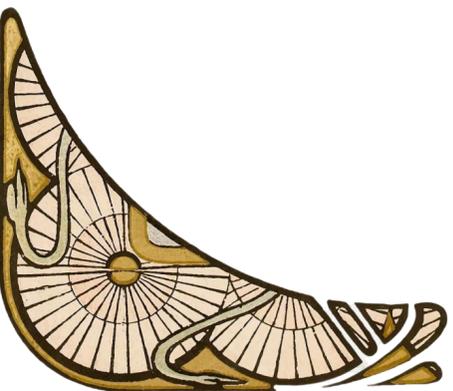


Estimating the counterfactual outcome $Y_c(x; \infty)$

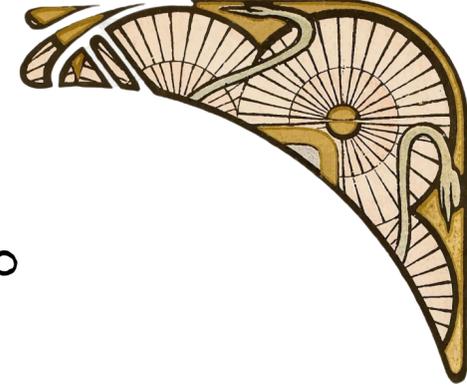


(k,l)-extractable memorisation (Carlini et al., 2023):

Assume the counterfactual performance is negligible in the absence of training



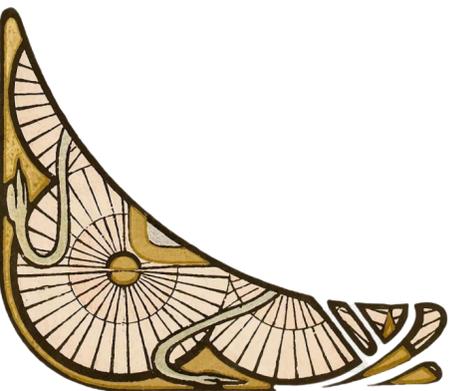
Estimating the counterfactual outcome $Y_c(\mathbf{x}; \infty)$



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

Assume the counterfactual performance is negligible in the absence of training

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



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

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

Assume the counterfactual performance is negligible in the absence of training

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

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

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



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

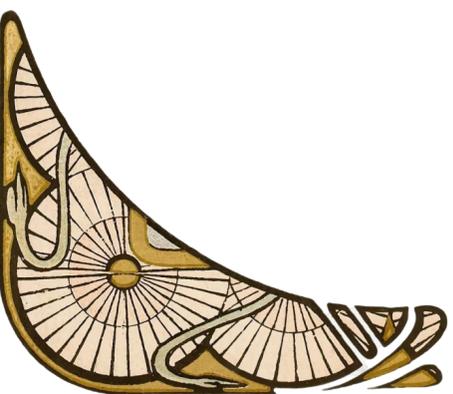
Assume the counterfactual performance is negligible in the absence of training

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

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

Re-train the model multiple times and average across runs

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



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



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

Assume the counterfactual performance is negligible in the absence of training

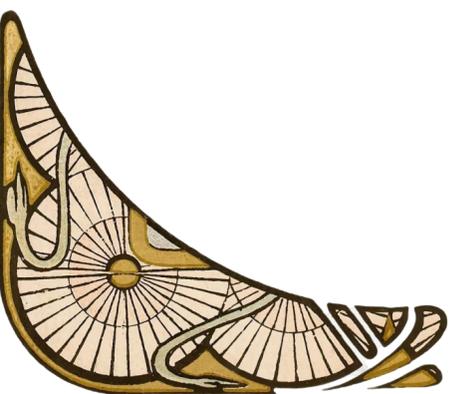
$$\tau_{\mathbf{x},c}^{\text{extr}} = Y_c(\mathbf{x}; g)$$

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

Re-train the model multiple times and average across runs

$$\tau_{\mathbf{x},p(\psi)}^{\text{arch}} = \mathbb{E}_{\psi} [Y_T(\mathbf{x}; G(\mathbf{x})) \mid G(\mathbf{x}) \neq \infty] - \mathbb{E}_{\psi} [Y_T(\mathbf{x}; \infty) \mid G(\mathbf{x}) = \infty]$$

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



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



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

Assume the counterfactual performance is negligible in the absence of training

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

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

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

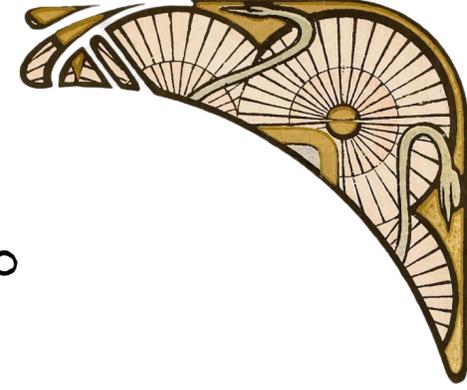
Re-train the model multiple times and average across runs

$$\tau_{\mathbf{x},p(\psi)}^{\text{arch}} = \mathbb{E}_{\psi} [Y_T(\mathbf{x}; G(\mathbf{x})) \mid G(\mathbf{x}) \neq \infty] - \mathbb{E}_{\psi} [Y_T(\mathbf{x}; \infty) \mid G(\mathbf{x}) = \infty]$$

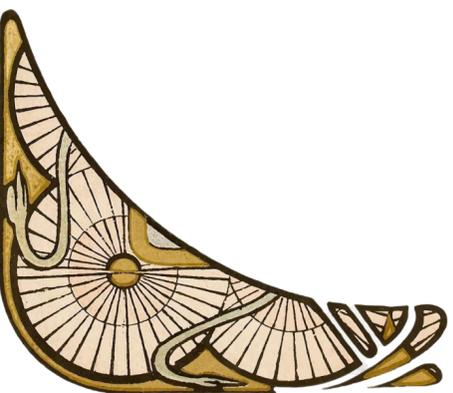
- Computationally expensive
- Estimates memorisation for a model architecture rather than a specific model
- Does not account for data order
- No insights on the training dynamics



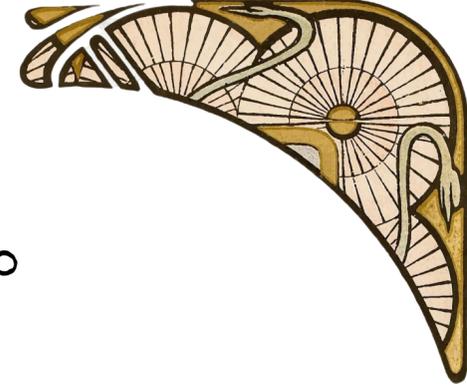
From instance-level to batch-level memorisation



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

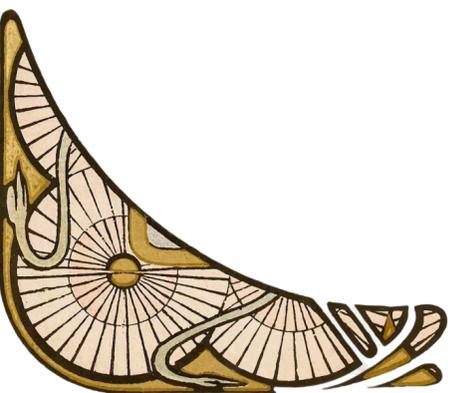


From instance-level to batch-level memorisation



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

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



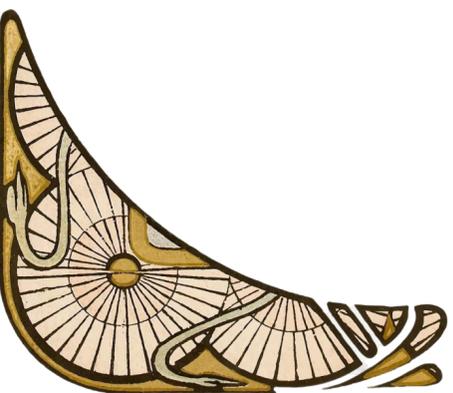
From instance-level to batch-level memorisation



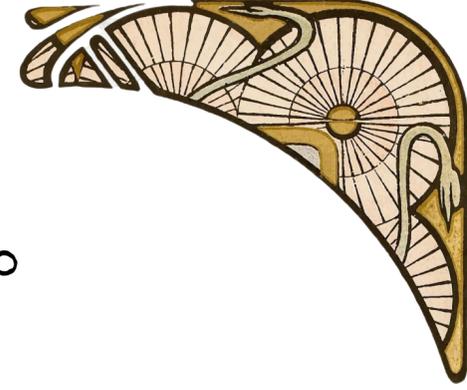
$$\tau_{\mathbf{x},c} \stackrel{\text{def}}{=} Y_c(\mathbf{x}; g) - Y_c(\mathbf{x}; \infty)$$

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

Average performance on batch
when trained on batch at step g



From instance-level to batch-level memorisation

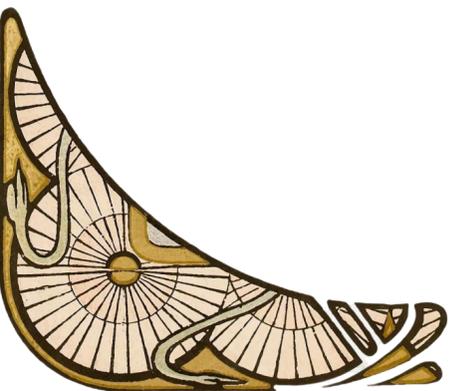


$$\tau_{\mathbf{x},c} \stackrel{\text{def}}{=} Y_c(\mathbf{x}; g) - Y_c(\mathbf{x}; \infty)$$

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

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

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

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

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

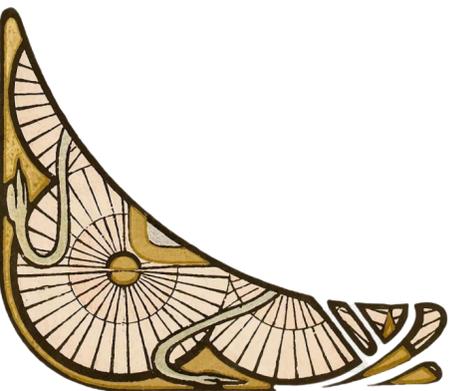
Validation Batch

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

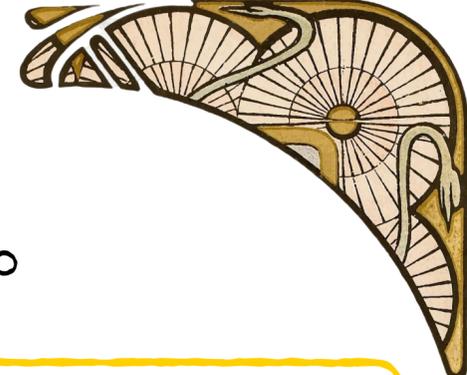
$$\tau_{g,c} \stackrel{\text{def}}{=} \underbrace{\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{Training Batch}]}_{\text{Average performance on batch when trained on batch at step } g} - \underbrace{\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{Validation Batch}]}_{\text{Average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

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

✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable

$\tau_{\mathbf{x},c}$ $\stackrel{\text{def}}{=}$

$Y_c(\mathbf{x}; g)$

—

$Y_c(\mathbf{x}; \infty)$

$\tau_{g,c}$ $\stackrel{\text{def}}{=}$

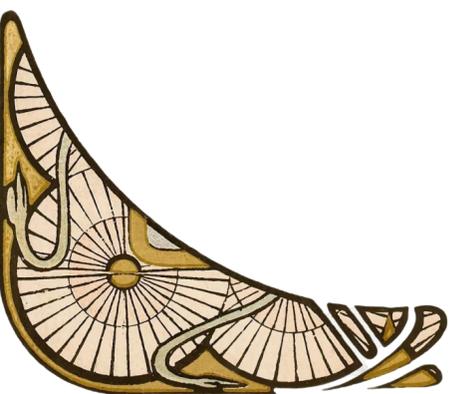
$$\underbrace{\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{Observable}]}_{\text{Average performance on batch when trained on batch at step } g}$$

—

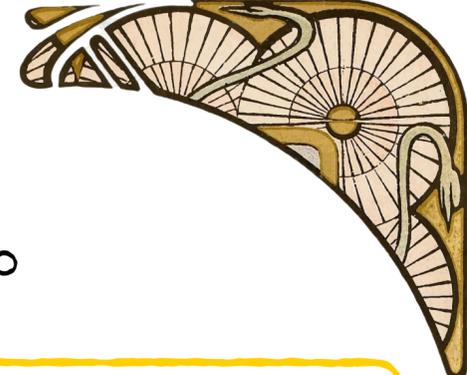
$$\underbrace{\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{Observable}]}_{\text{Average performance on batch when not trained on batch given it is selected for training at step } g}$$

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

✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable

$\tau_{\mathbf{x},c}$ $\stackrel{\text{def}}{=}$

$Y_c(\mathbf{x}; g)$

—

$Y_c(\mathbf{x}; \infty)$

$\tau_{g,c}$ $\stackrel{\text{def}}{=}$

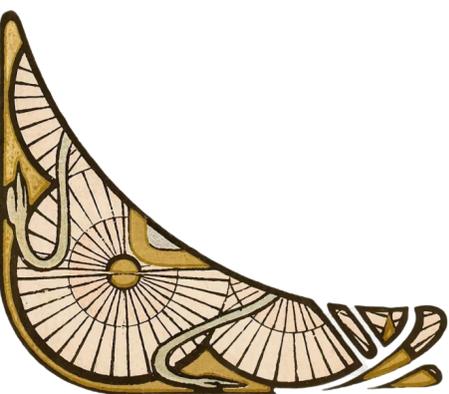
$$\underbrace{\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{Observable}]}_{\text{Average performance on batch when trained on batch at step } g}$$

—

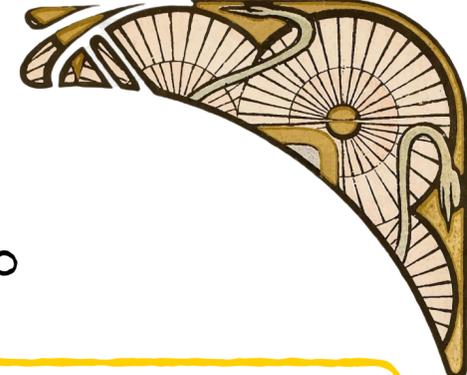
$$\underbrace{\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{Observable}]}_{\text{Average performance on batch when not trained on batch given it is selected for training at step } g}$$

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

✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable

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

$Y_c(\mathbf{x}; g)$

—

$Y_c(\mathbf{x}; \infty)$

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

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

—

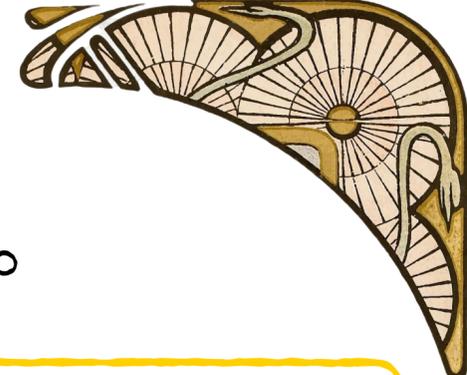
$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) | \text{Observable}]$ ✗

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



On the way to estimating the counterfactual



Training Batch

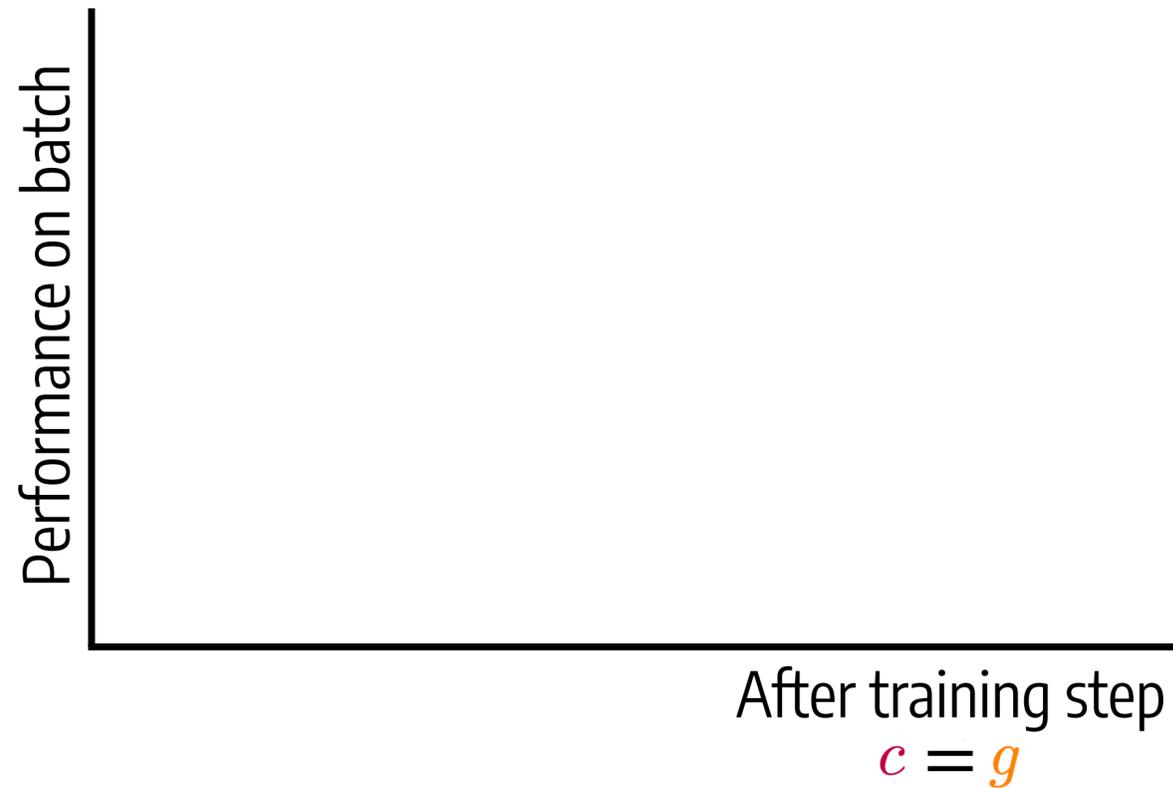
Validation Batch



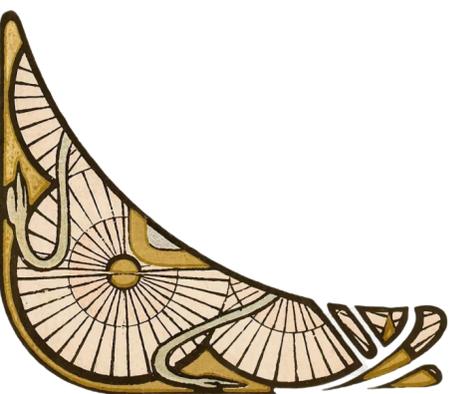
g and = Observable
 ∞ and = Observable



g and = Not Observable
 ∞ and = Not Observable



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



On the way to estimating the counterfactual

Training Batch

Validation Batch

✓ g and = Observable
 ∞ and = Observable

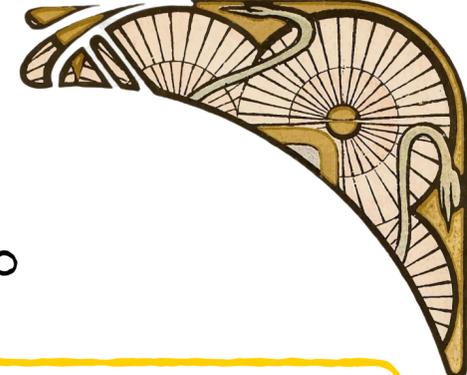
✗ g and = Not Observable
 ∞ and = Not Observable



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

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

On the way to estimating the counterfactual

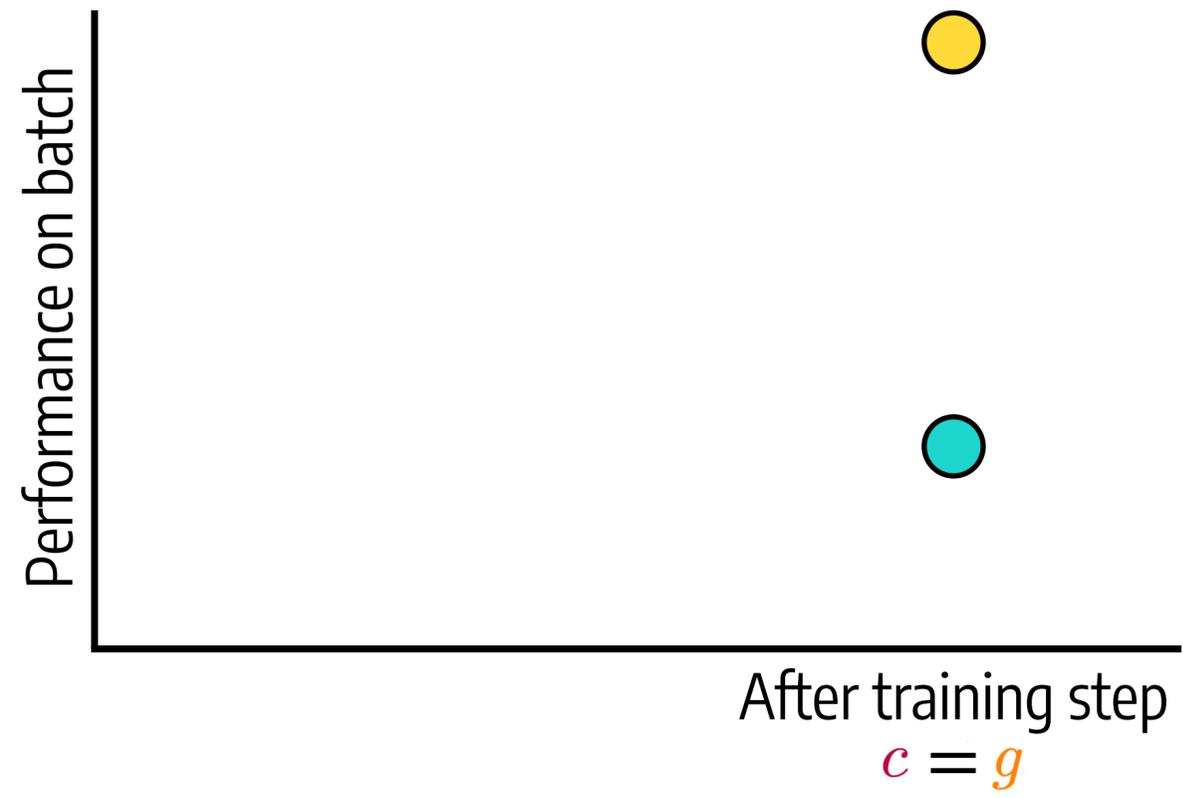


Training Batch

Validation Batch

✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable



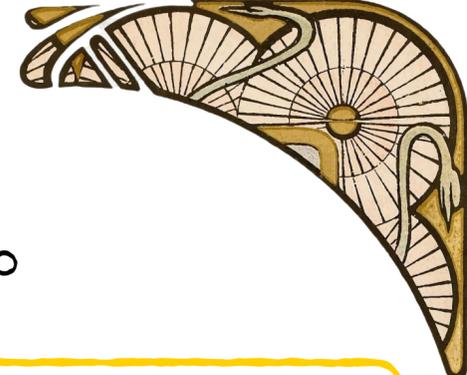
$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{ } \checkmark]$$

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

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



On the way to estimating the counterfactual

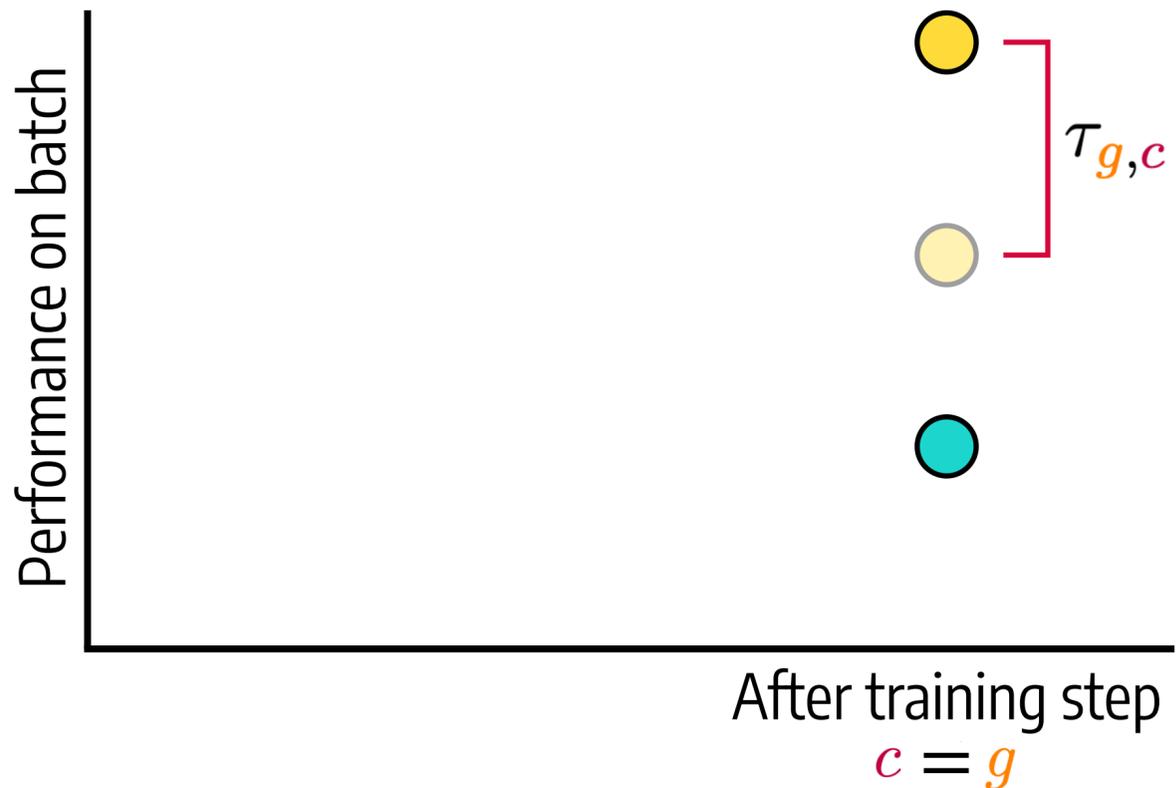


Training Batch

Validation Batch

✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable



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

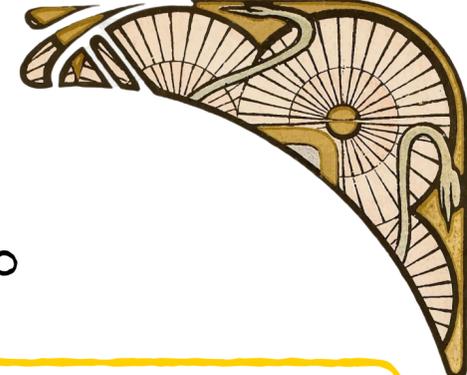
$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{ }] \quad \times$$

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

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



On the way to estimating the counterfactual

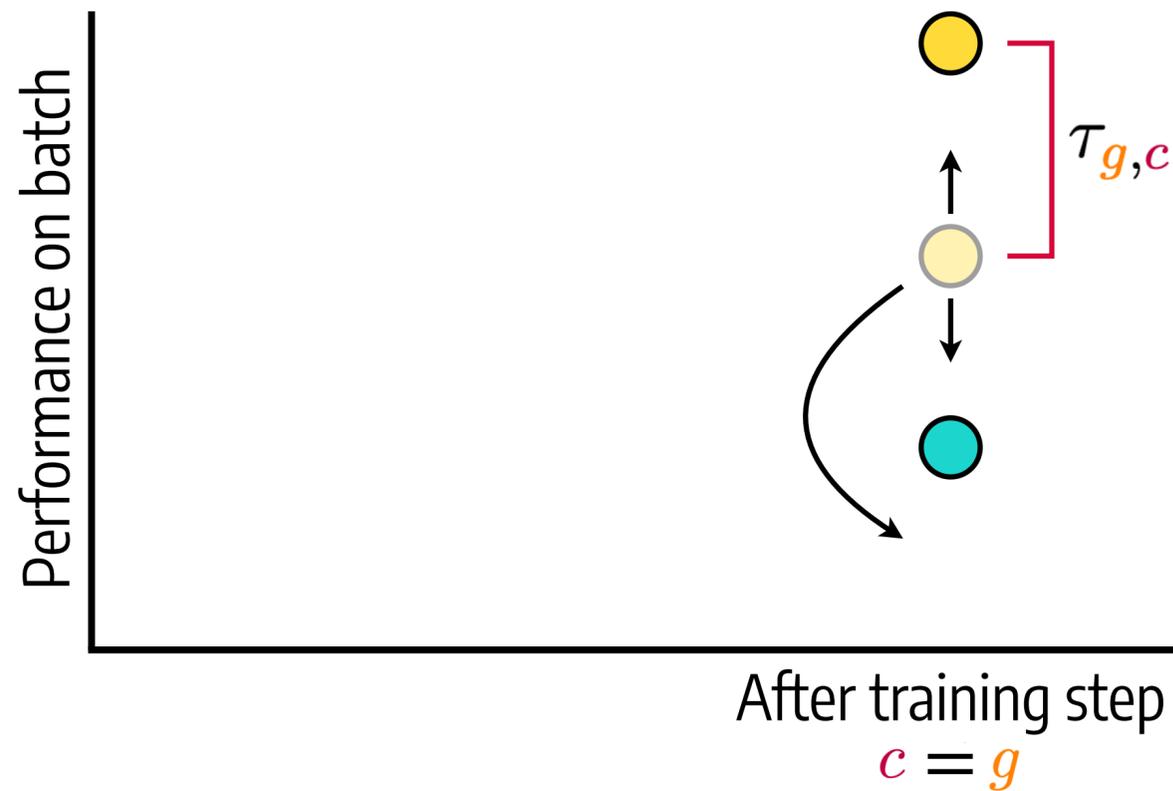


Training Batch

Validation Batch

✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable

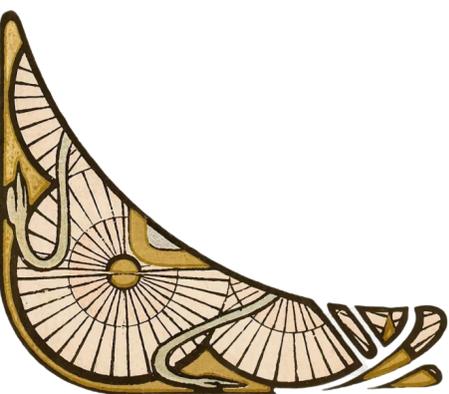


$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{ }] \quad \checkmark$$

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

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

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



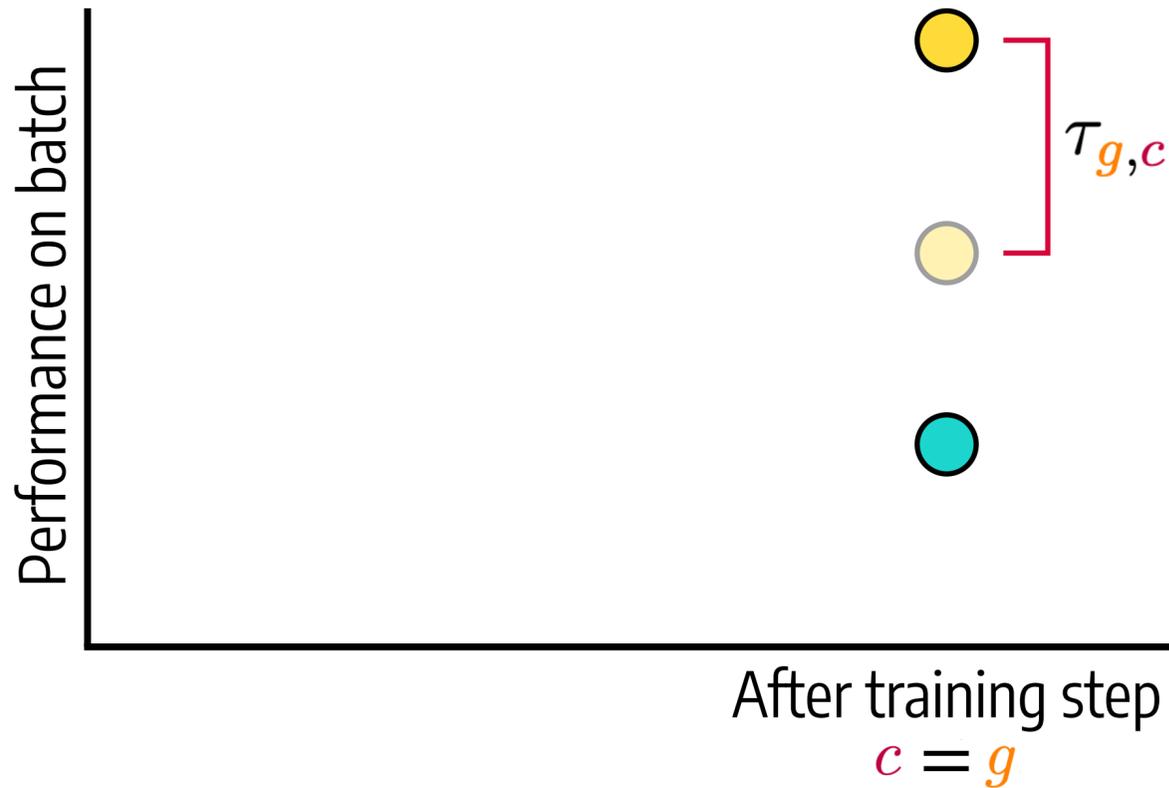
The Difference estimator

Training Batch

Validation Batch

✓ g and  = Observable
 ∞ and  = Observable

✗ g and  = Not Observable
 ∞ and  = Not Observable



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

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

$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{cyan box} \checkmark]$$

$$\tau_{g,c} = \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{yellow box} \checkmark] - \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{yellow box} \times]$$

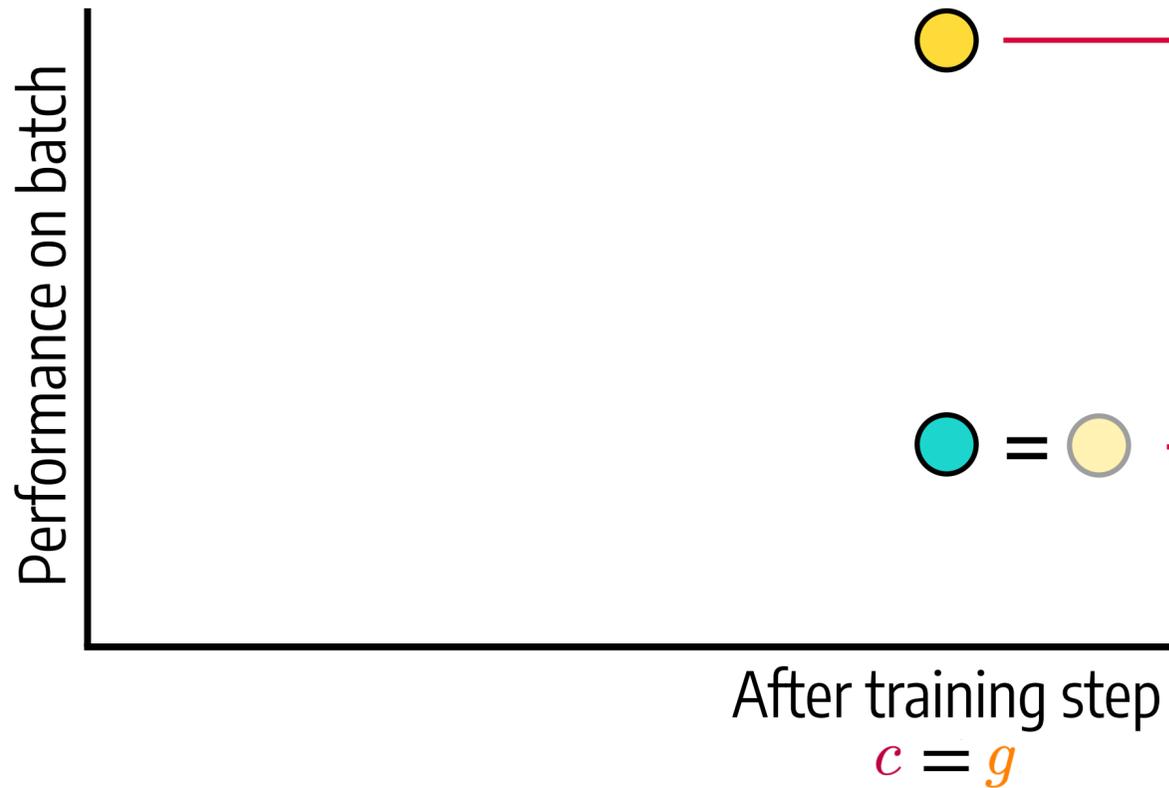
The Difference estimator

Training Batch

Validation Batch

✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable



$$\begin{aligned} & \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{yellow box}] \\ & \tau_{g,c} \\ & \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{cyan box}] = \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{yellow box}] \end{aligned}$$

IID Assumption

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

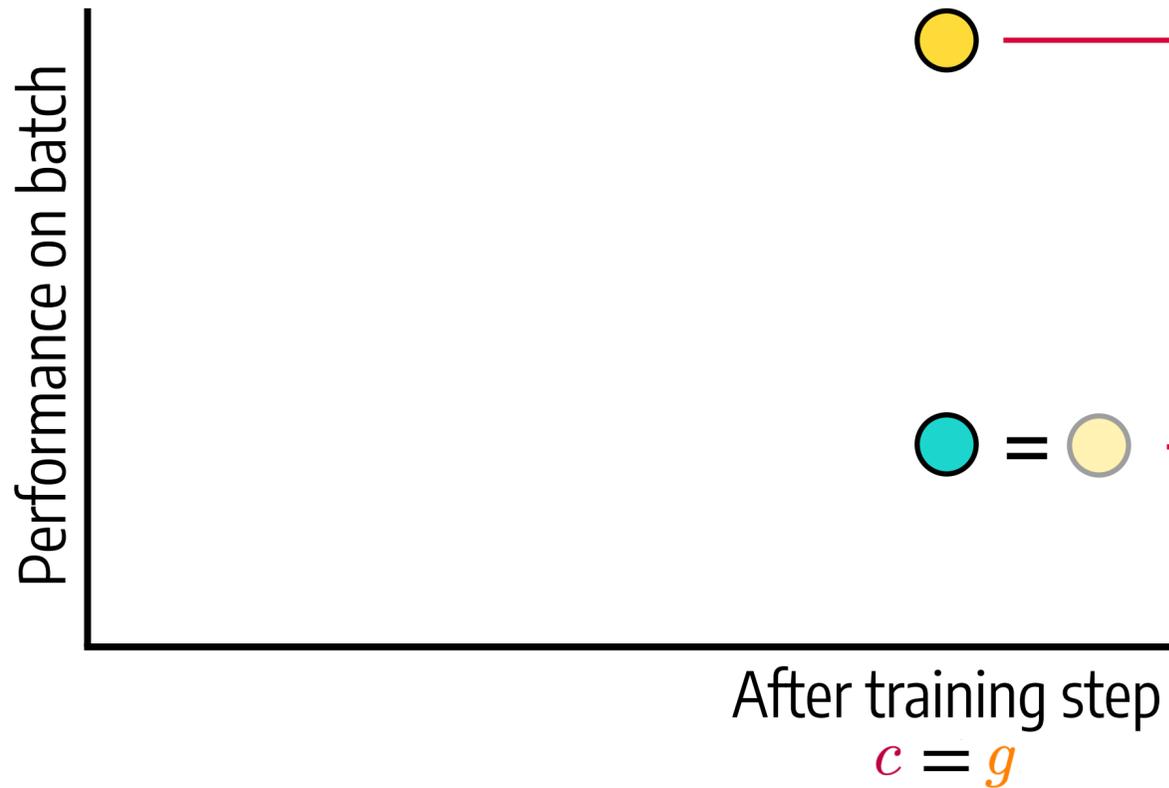
The Difference estimator

Training Batch

Validation Batch

✓ g and  = Observable
 ∞ and  = Observable

✗ g and  = Not Observable
 ∞ and  = Not Observable

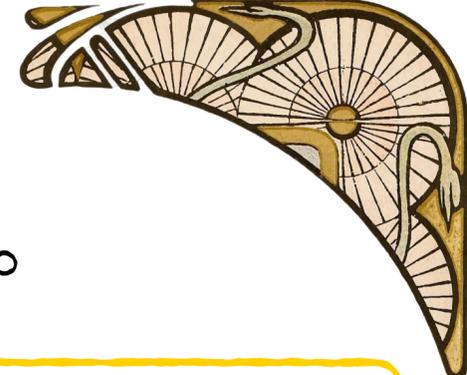


$$\begin{aligned} & \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{yellow box}] \\ & \tau_{g,c} \\ & \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{cyan box}] = \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{yellow box}] \end{aligned}$$

IID Assumption

$$\tau_{g,c}^{\text{diff}} = \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{yellow box}] - \mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \text{cyan box}]$$

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

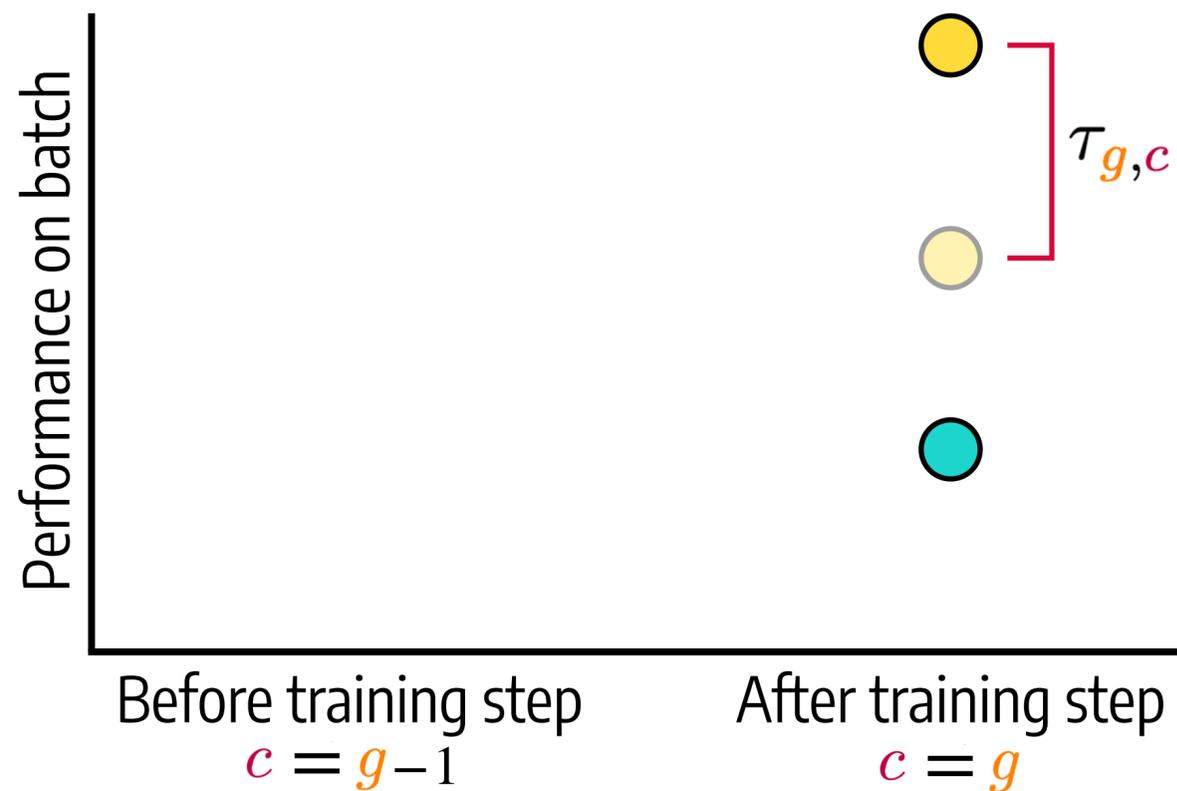


Training Batch

Validation Batch

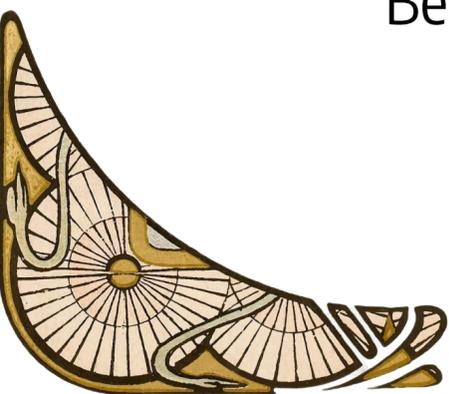
✓ g and = Observable
 ∞ and = Observable

✗ g and = Not Observable
 ∞ and = Not Observable



$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \text{Observable}]$$

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



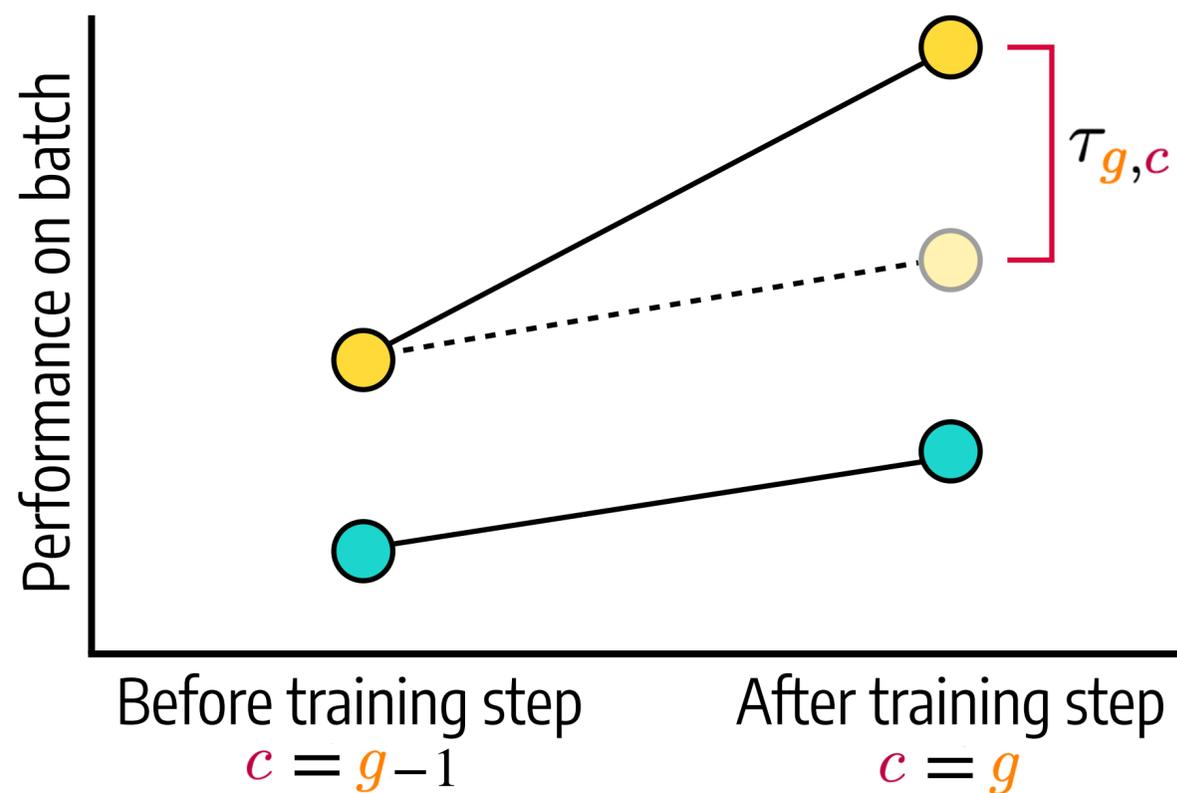
The Difference-in-Differences estimator (1/2)

Training Batch

Validation Batch

✓ g and \square = Observable
 ∞ and \square = Observable

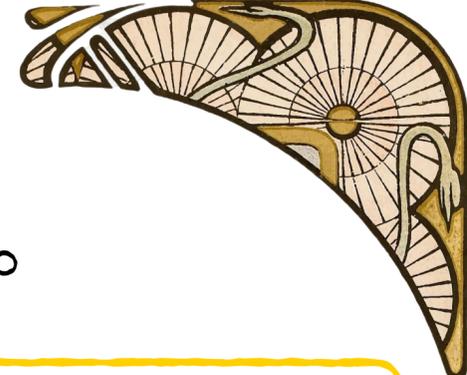
✗ g and \square = Not Observable
 ∞ and \square = Not Observable



$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; g) \mid \square] \quad \checkmark$$

$$\mathbb{E}_{\mathbf{x}} [Y_c(\mathbf{x}; \infty) \mid \square] \quad \checkmark$$

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

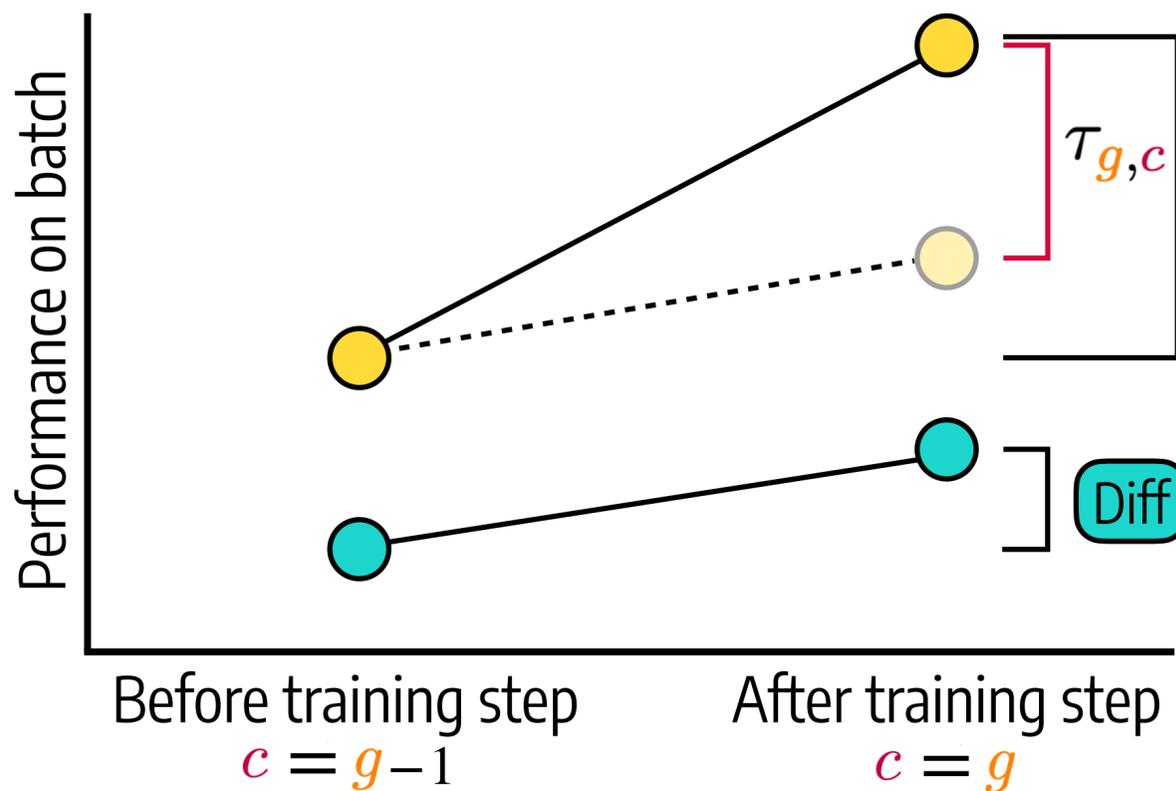


Training Batch

Validation Batch

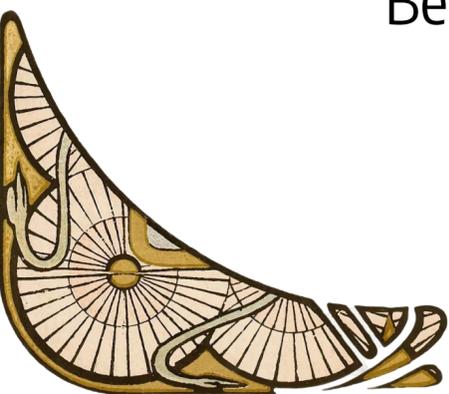
✓ g and \square = Observable
 ∞ and \square = Observable

✗ g and \square = Not Observable
 ∞ and \square = Not Observable

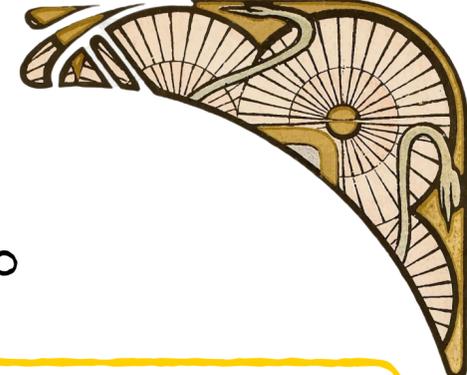


Diff = $\mathbb{E}_x [Y_c(\mathbf{x}; g) | \square] - \mathbb{E}_x [Y_{g-1}(\mathbf{x}; g) | \square]$

Diff = $\mathbb{E}_x [Y_c(\mathbf{x}; \infty) | \square] - \mathbb{E}_x [Y_{g-1}(\mathbf{x}; \infty) | \square]$



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

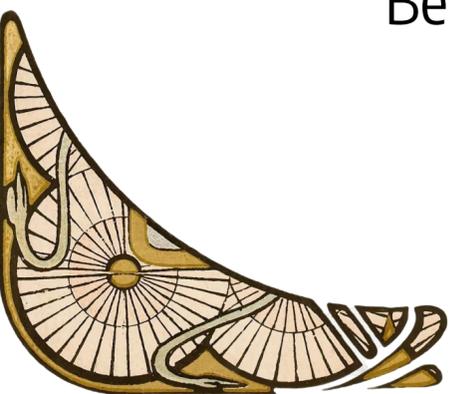
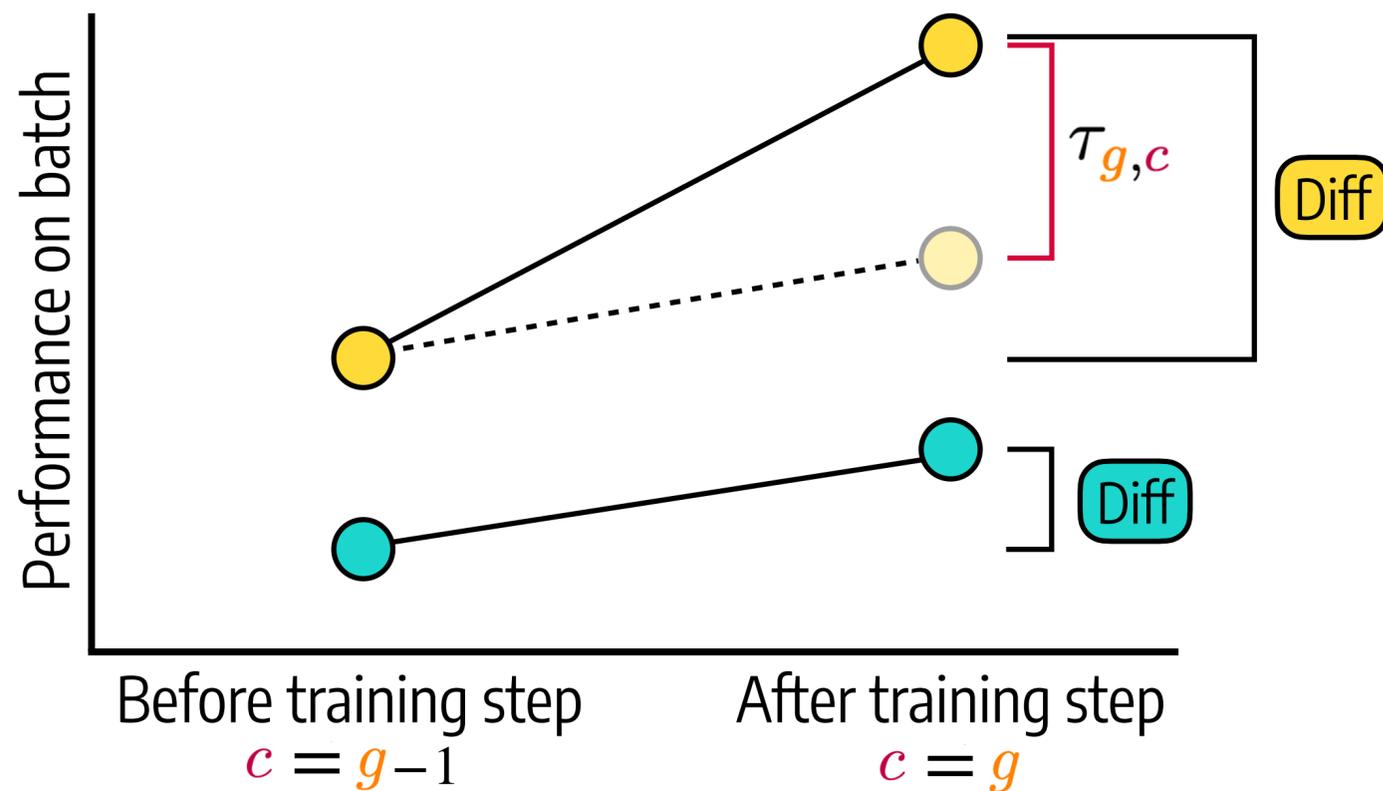


Training Batch

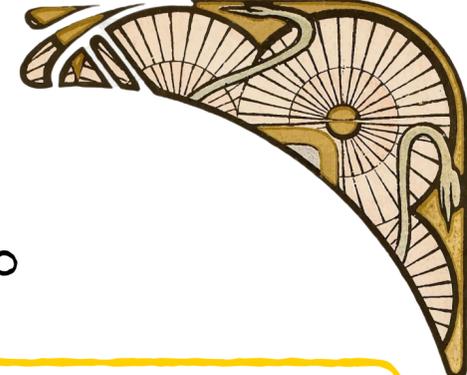
Validation Batch

✓ g and \square = Observable
 ∞ and \square = Observable

✗ g and \square = Not Observable
 ∞ and \square = Not Observable



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

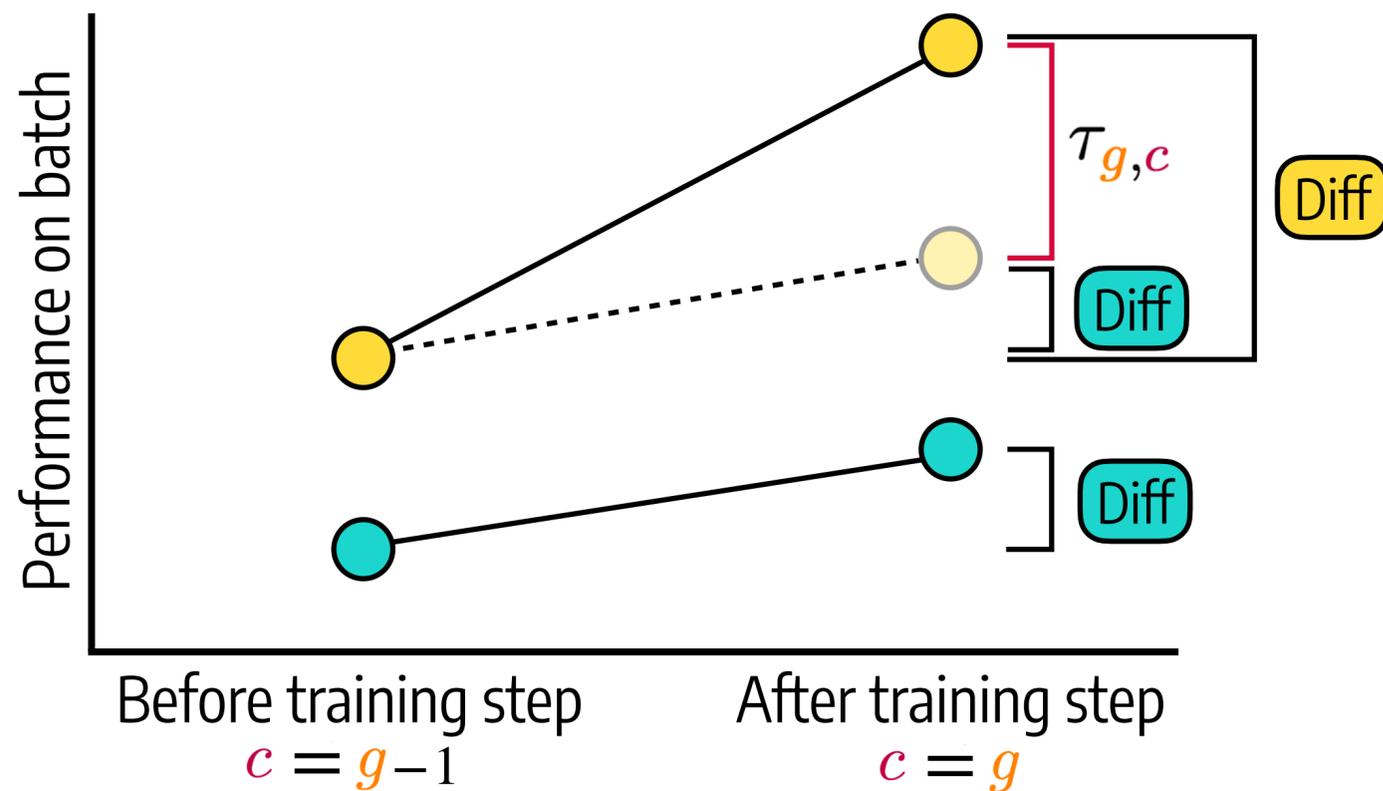


Training Batch

Validation Batch

✓ g and \square = Observable
 ∞ and \square = Observable

✗ g and \square = Not Observable
 ∞ and \square = Not Observable



Parallel Trend Assumption



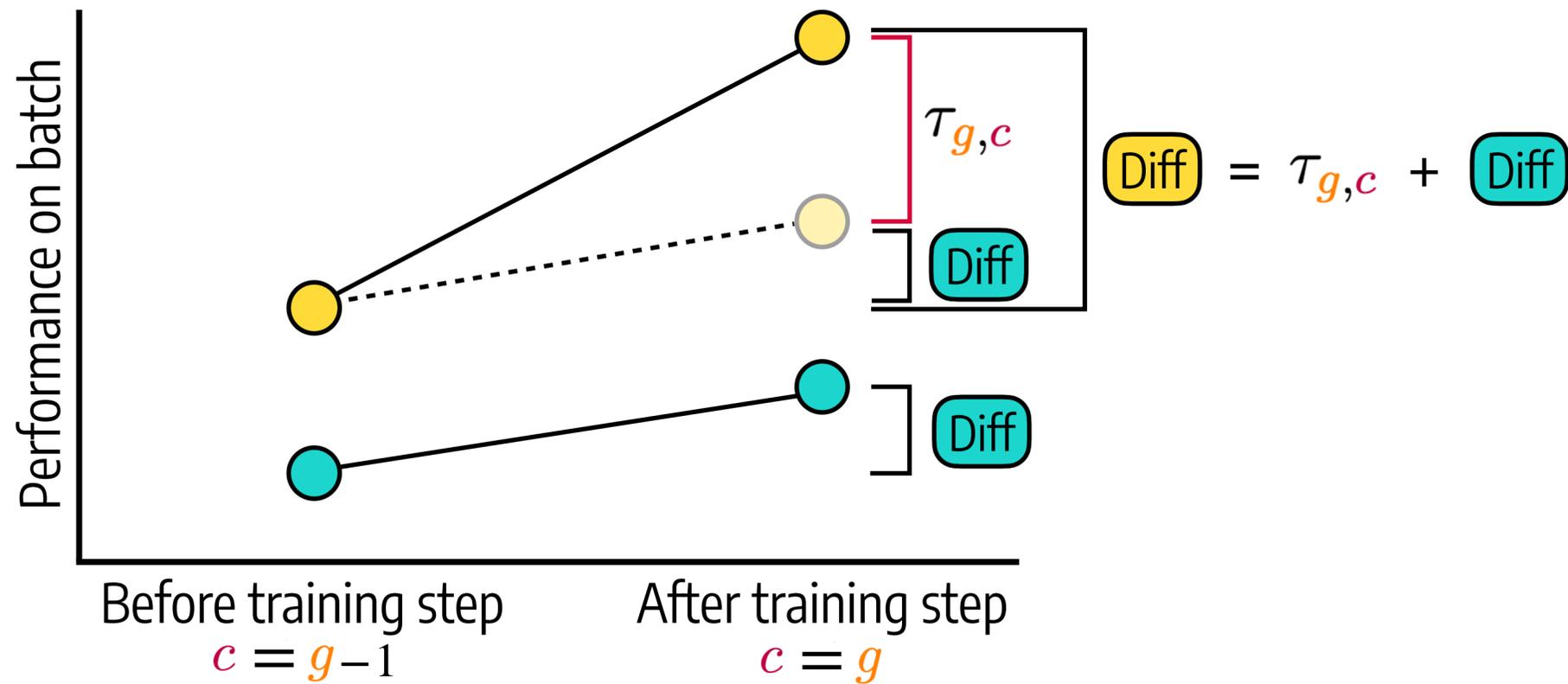
The Difference-in-Differences estimator (2/2)

Training Batch

Validation Batch

✓ g and \square = Observable
 ∞ and \square = Observable

✗ g and \square = Not Observable
 ∞ and \square = Not Observable



Parallel Trend Assumption

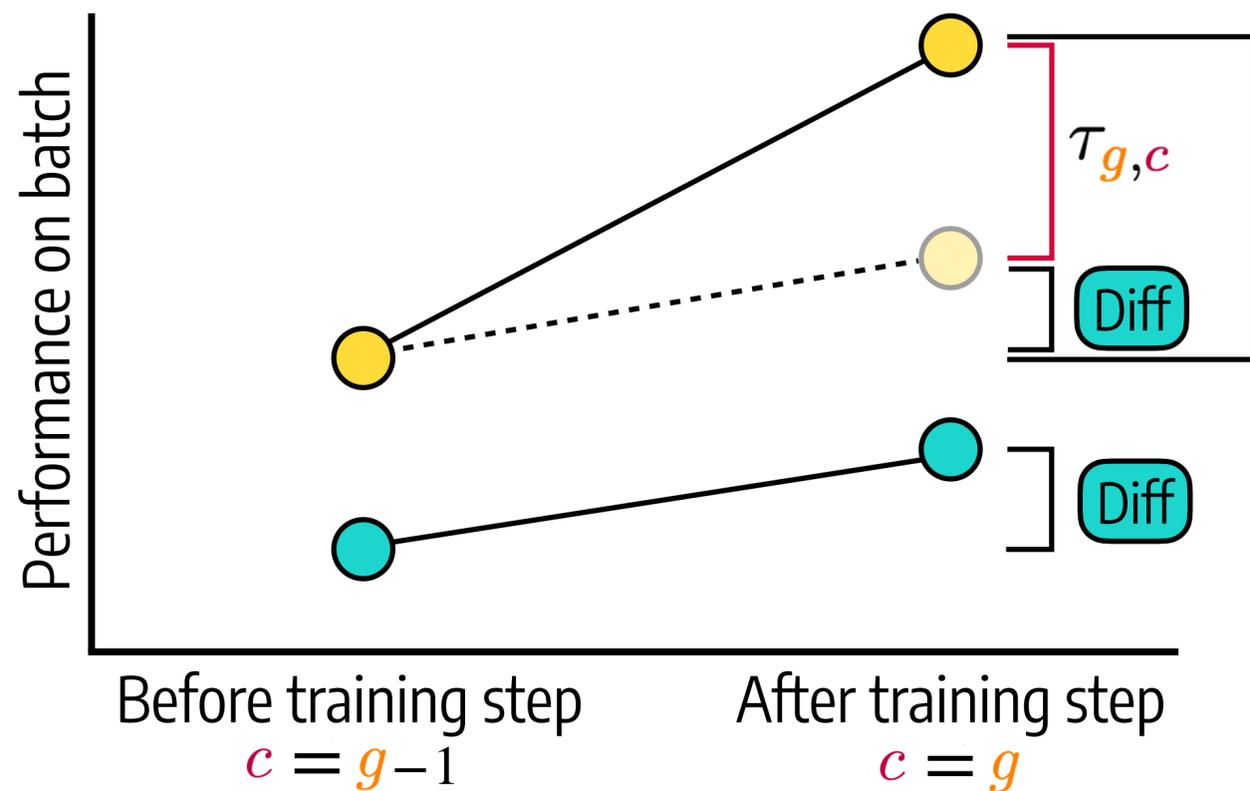
The Difference-in-Differences estimator (2/2)

Training Batch

Validation Batch

✓ g and \square = Observable
 ∞ and \square = Observable

✗ g and \square = Not Observable
 ∞ and \square = Not Observable

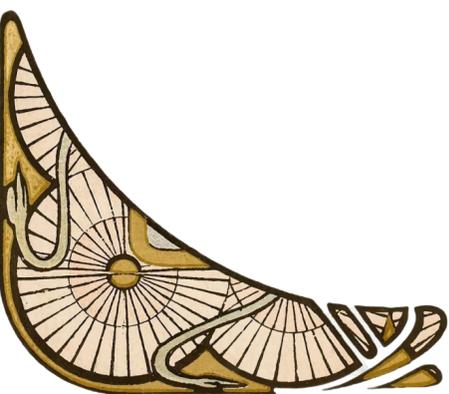
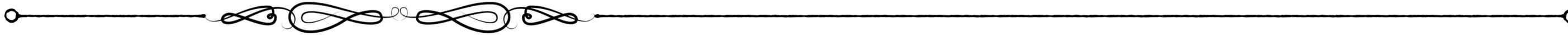


$$\text{Diff} = \tau_{g,c} + \text{Diff}$$

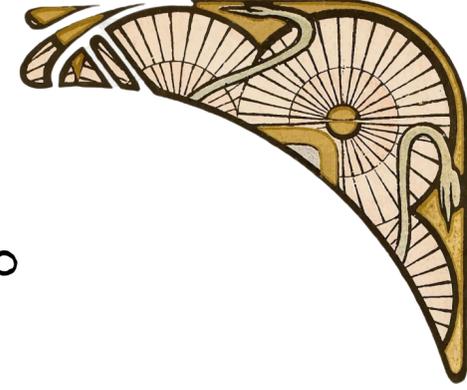
$$\tau_{g,c}^{\text{did}} = \text{Diff} - \text{Diff}$$

Parallel Trend Assumption

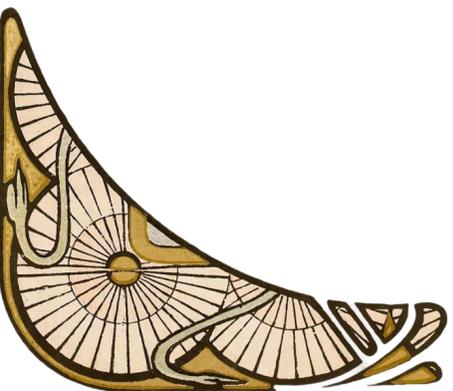
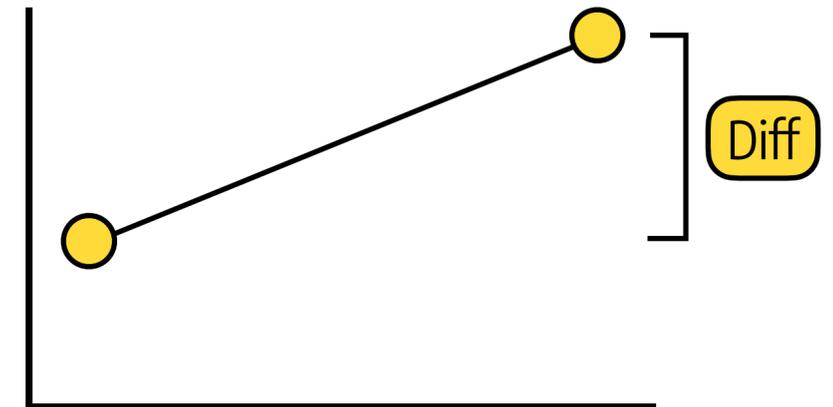
Difference-in-Differences in practice



Difference-in-Differences in practice



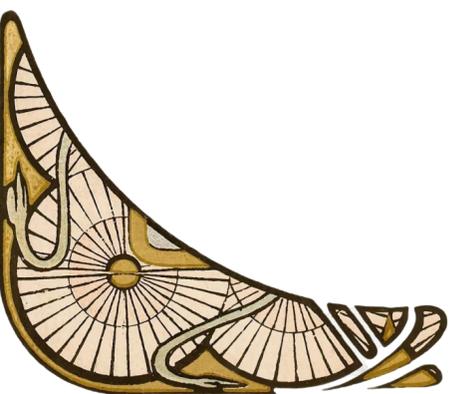
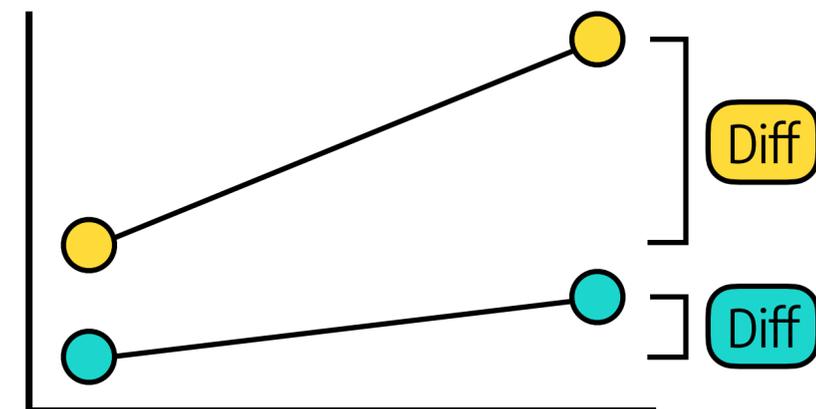
$$\tau_{g,c}^{\text{did}} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) - Y_{g-1}(\mathbf{x}; g) \mid \text{ }]$$



Difference-in-Differences in practice



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

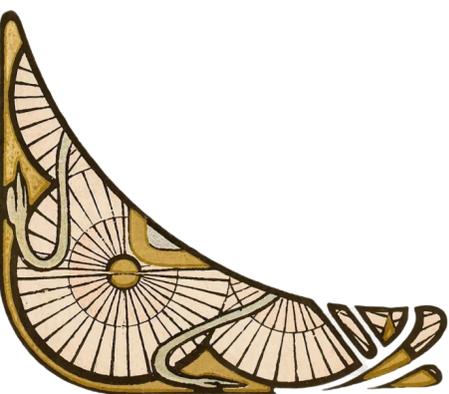
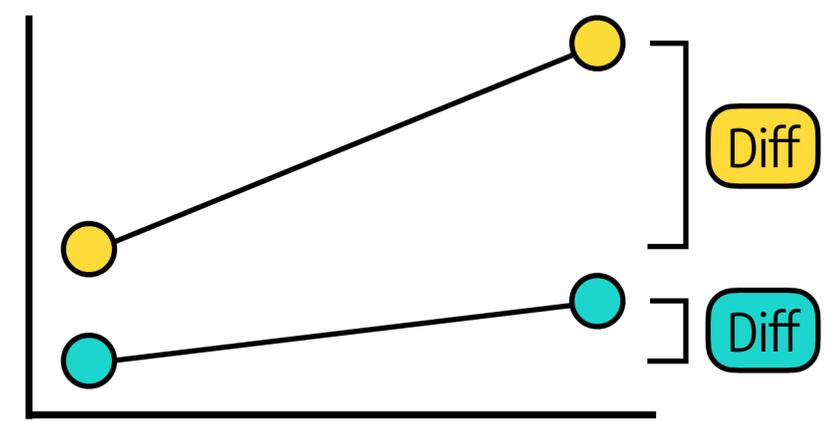


Difference-in-Differences in practice



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

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

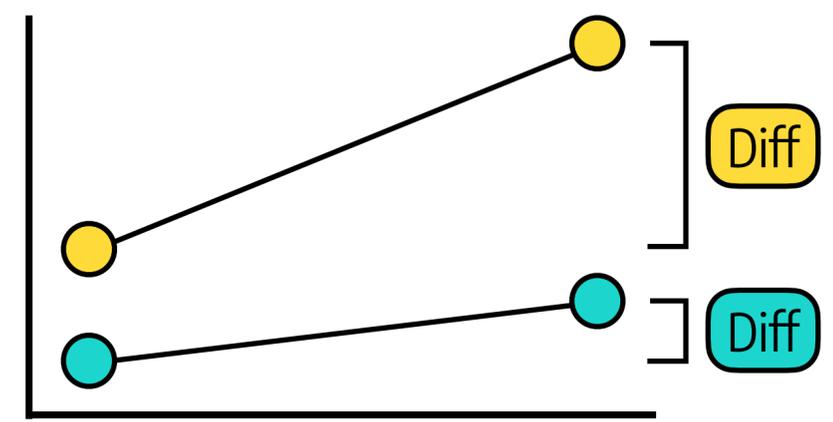


Difference-in-Differences in practice

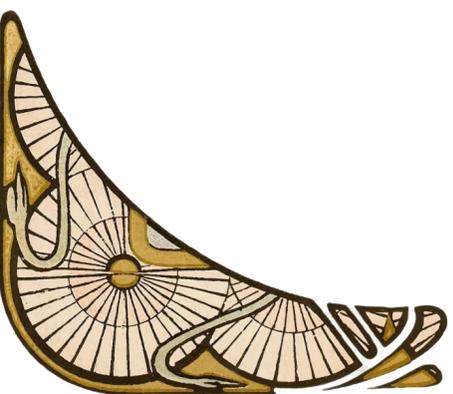
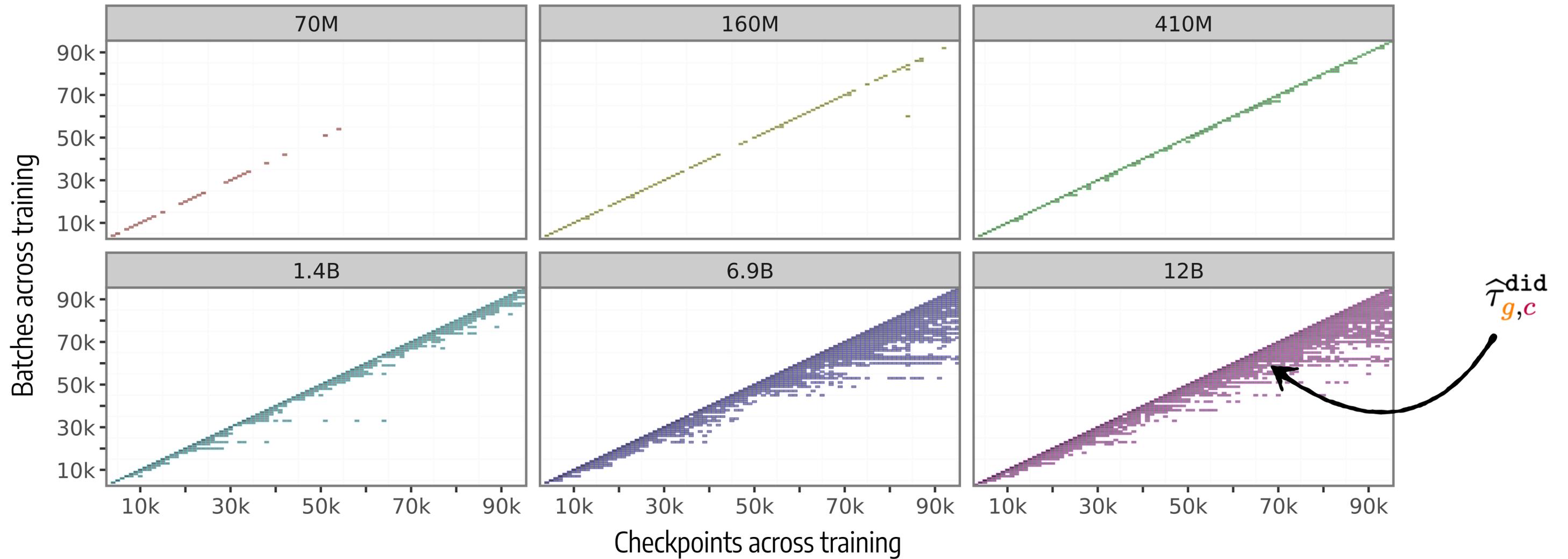
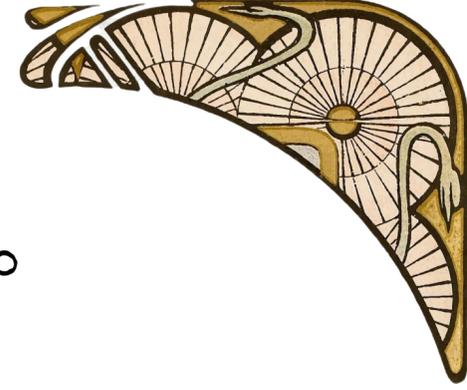


$$\tau_{g,c}^{\text{did}} = \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; g) - Y_{g-1}(\mathbf{x}; g) \mid \text{Training}] - \mathbb{E}_{\mathbf{x}}[Y_c(\mathbf{x}; \infty) - Y_{g-1}(\mathbf{x}; \infty) \mid \text{Validation}]$$

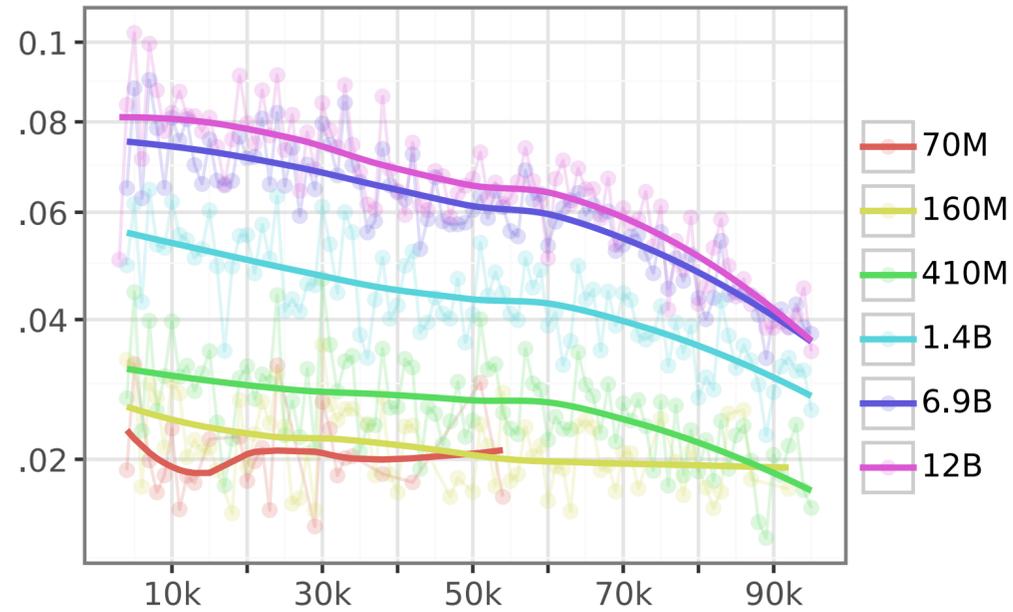
$$\hat{\tau}_{g,c}^{\text{did}} = \underbrace{(\bar{Y}_c(g) - \bar{Y}_{g-1}(g))}_{\text{Performance difference on training batches}} - \underbrace{(\bar{Y}_c(\infty) - \bar{Y}_{g-1}(\infty))}_{\text{Performance difference on validation batches}}$$



Constructing the memorisation profile

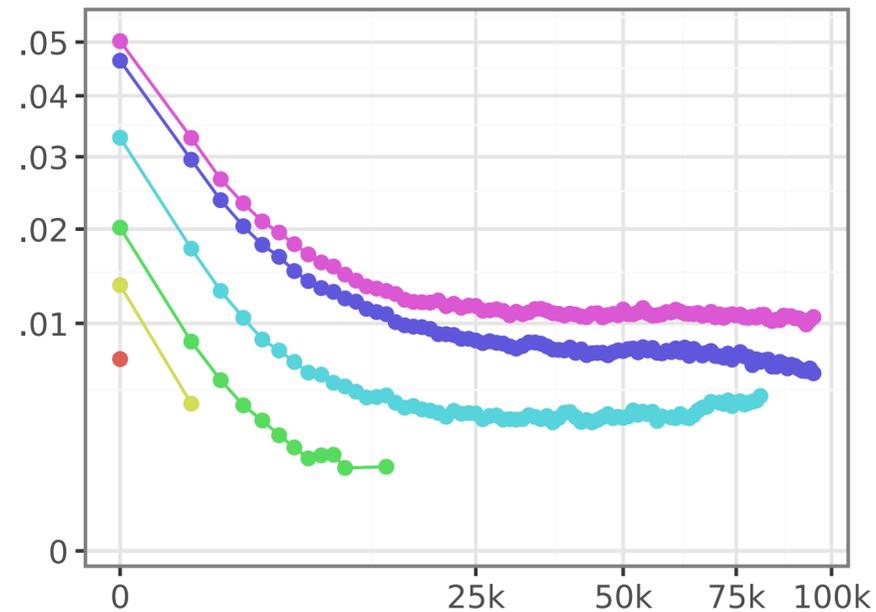


Aggregating the memorisation profile



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

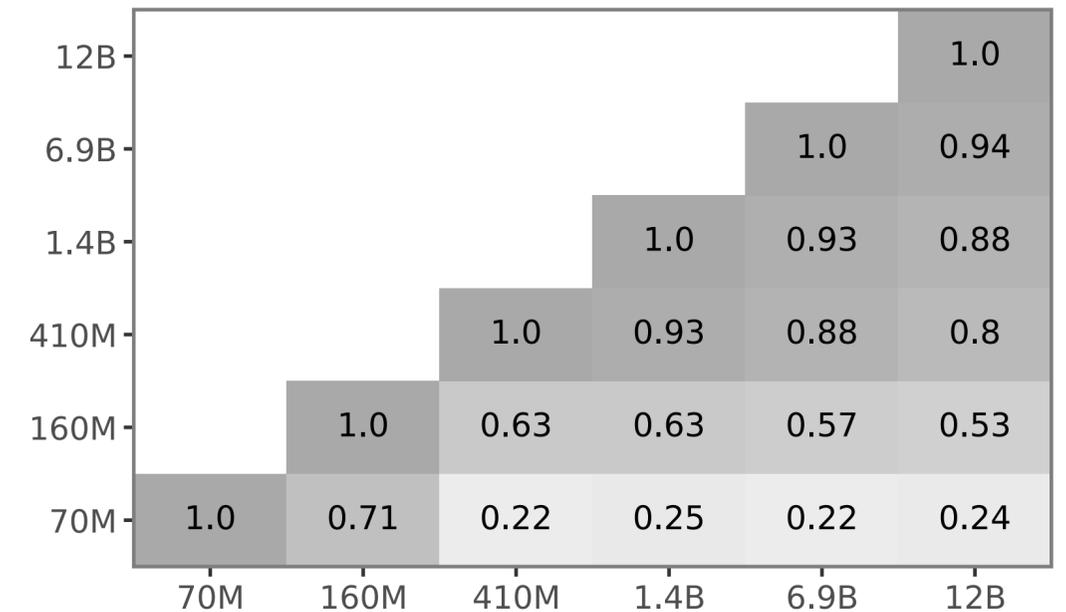
Instantaneous memorisation of a batch at the checkpoint is first seen



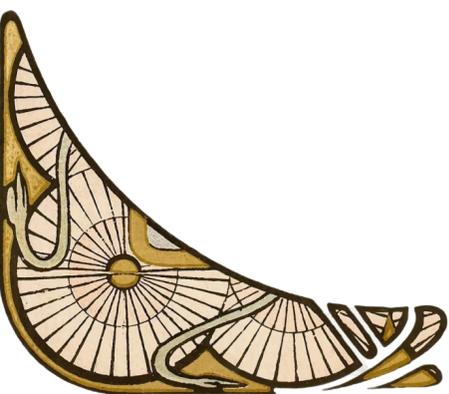
Checkpoints after batch is seen ($g > c$)

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

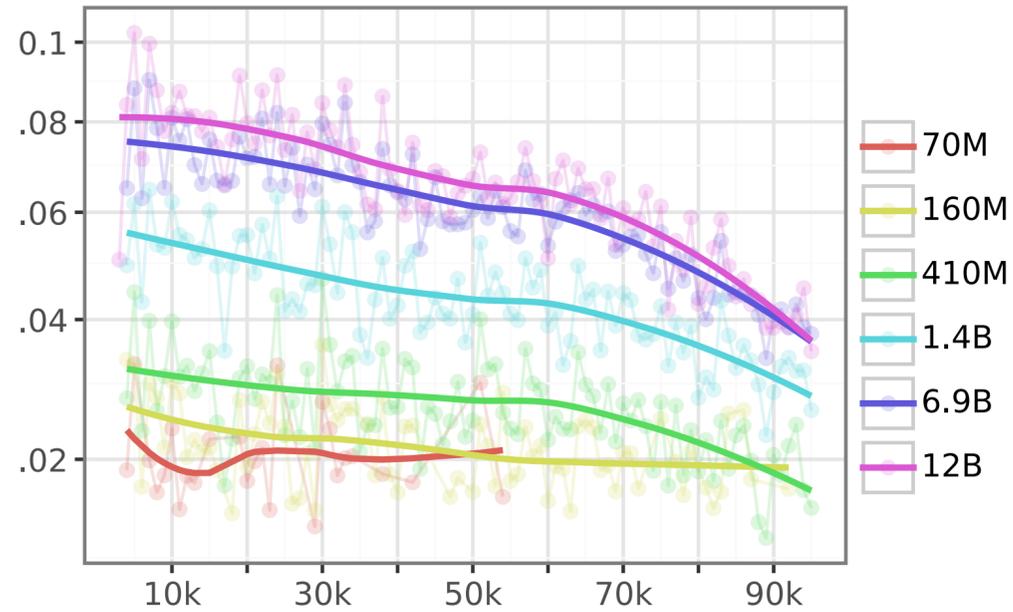
*(only averaging across batches that have been trained that much)



Pearson correlation between the memorisation profiles of different model sizes

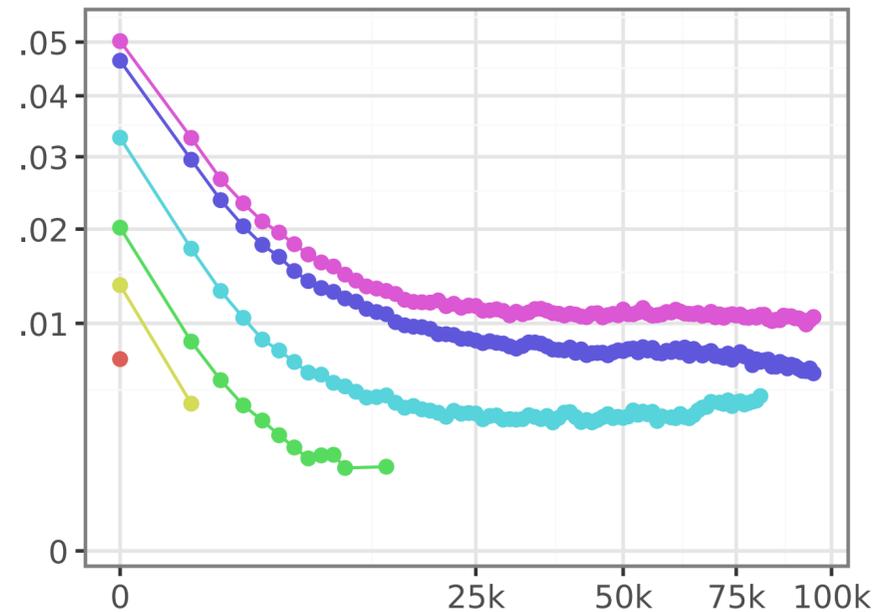


Aggregating the memorisation profile



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

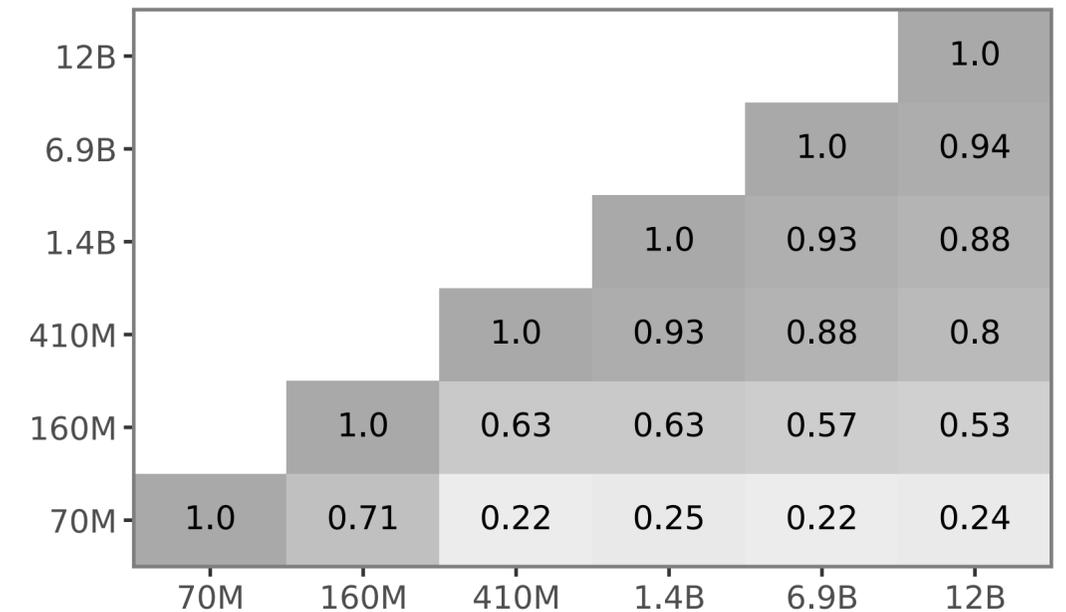
Instantaneous memorisation of a batch at the checkpoint is first seen



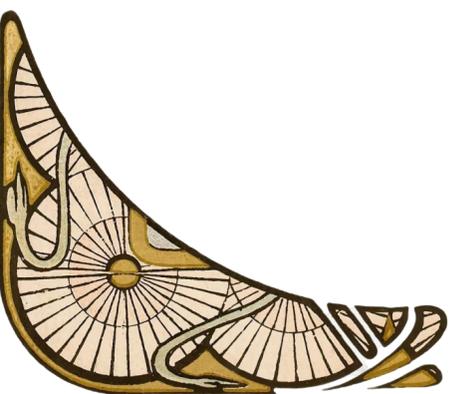
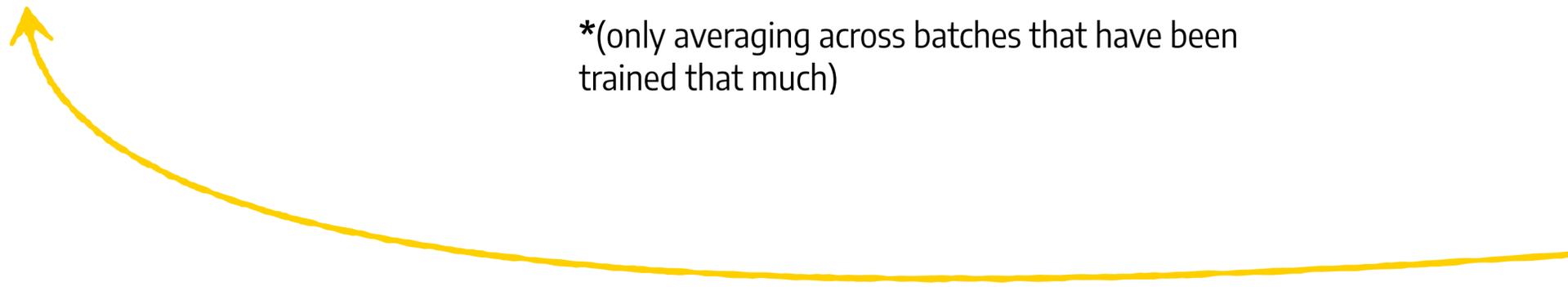
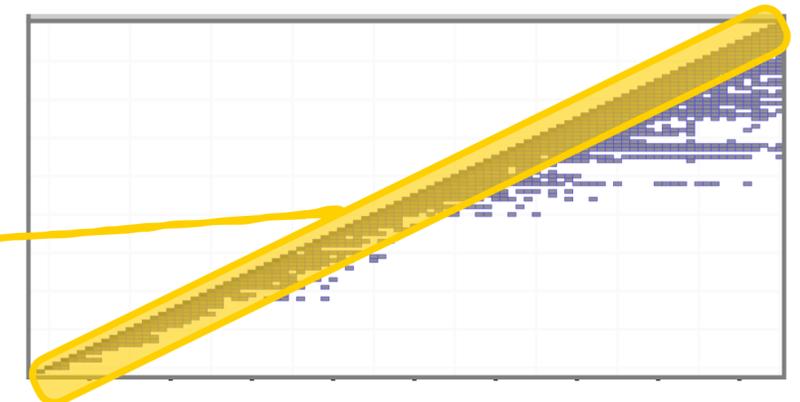
Checkpoints after batch is seen ($g > c$)

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

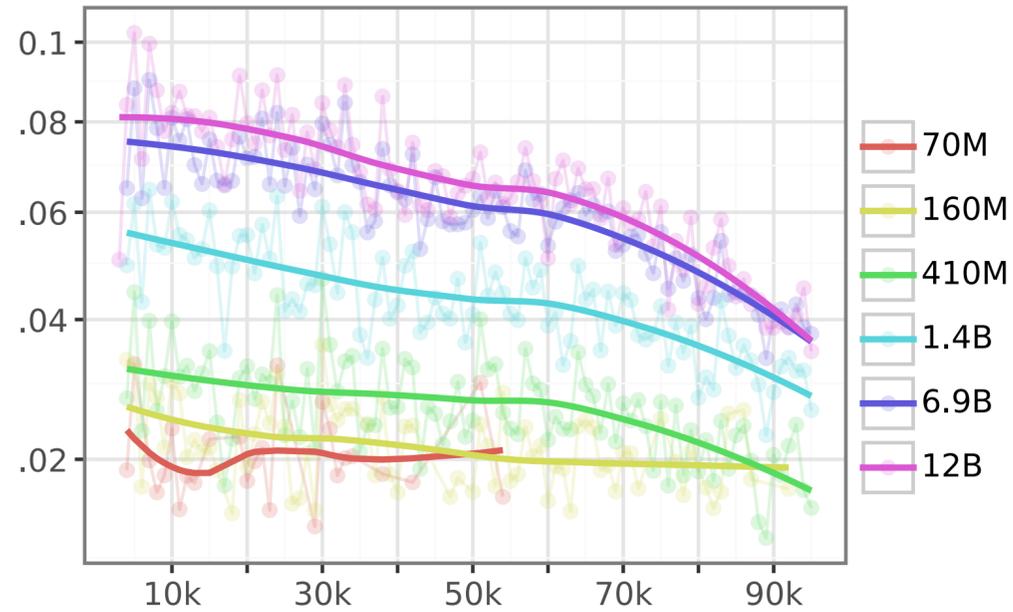
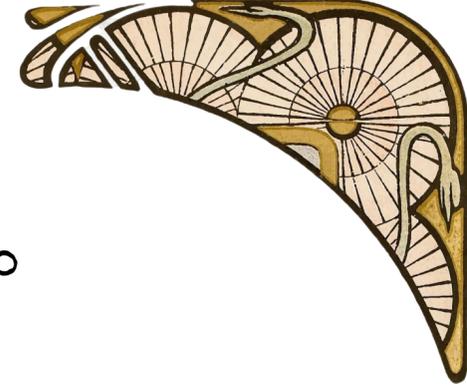
*(only averaging across batches that have been trained that much)



Pearson correlation between the memorisation profiles of different model sizes

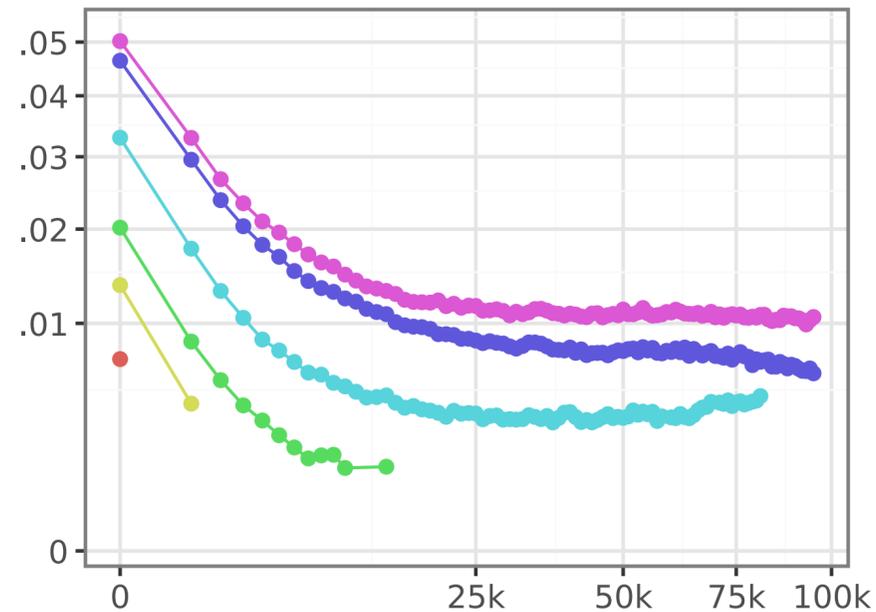


Aggregating the memorisation profile



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

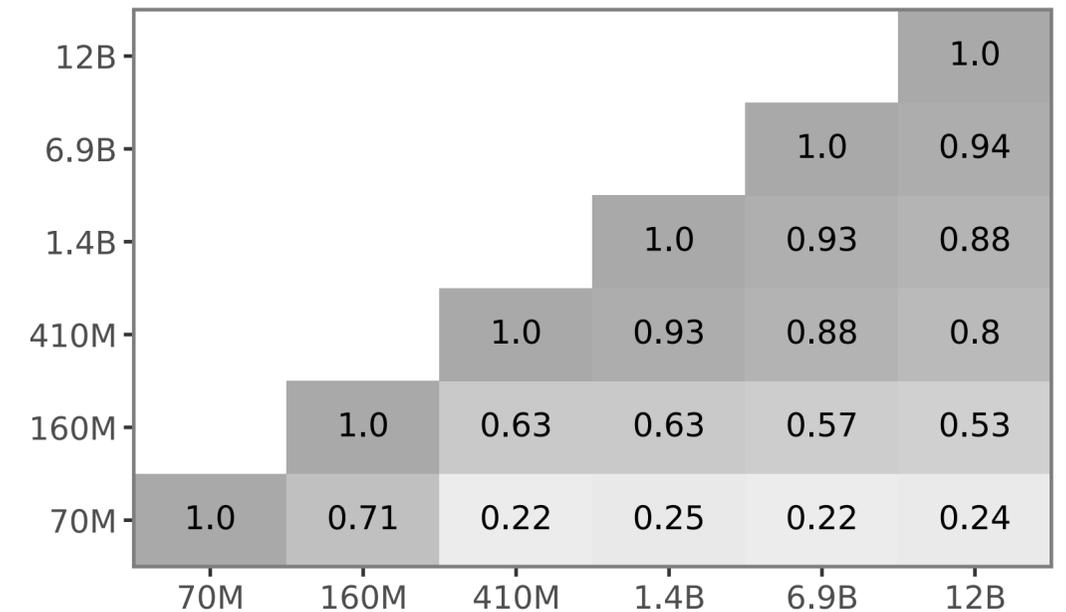
Instantaneous memorisation of a batch at the checkpoint is first seen



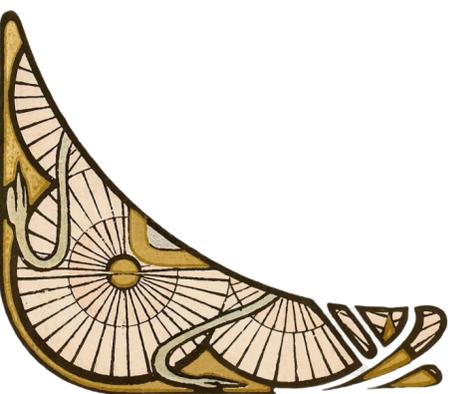
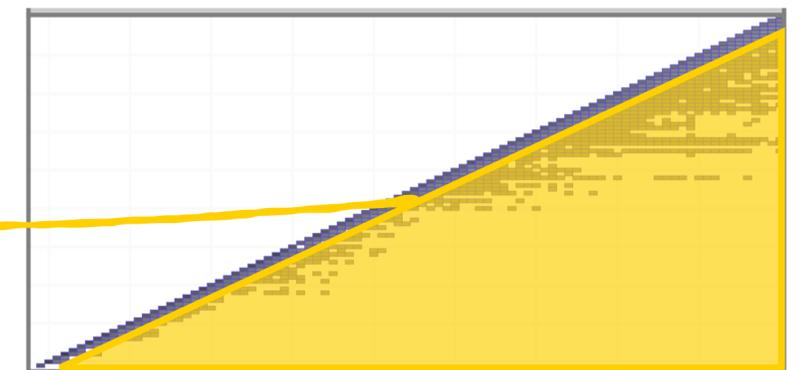
Checkpoints after batch is seen ($g > c$)

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

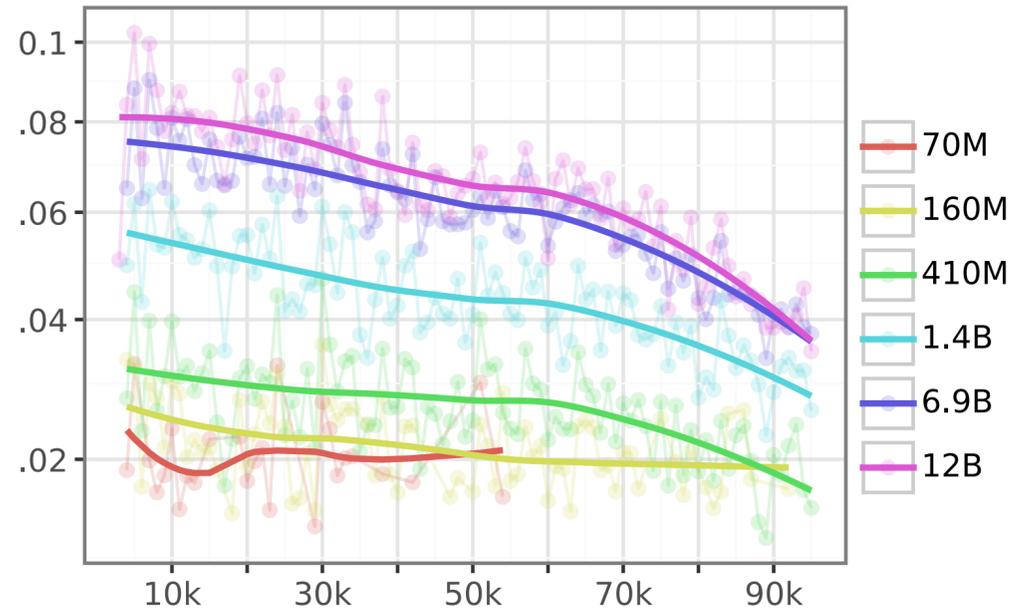
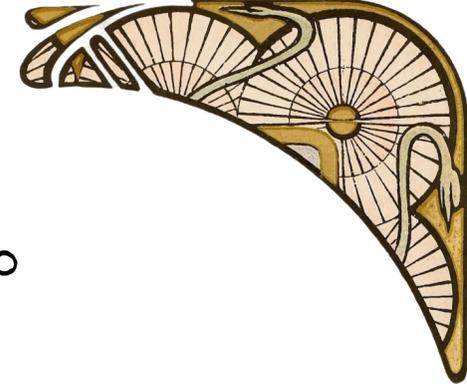
*(only averaging across batches that have been trained that much)



Pearson correlation between the memorisation profiles of different model sizes

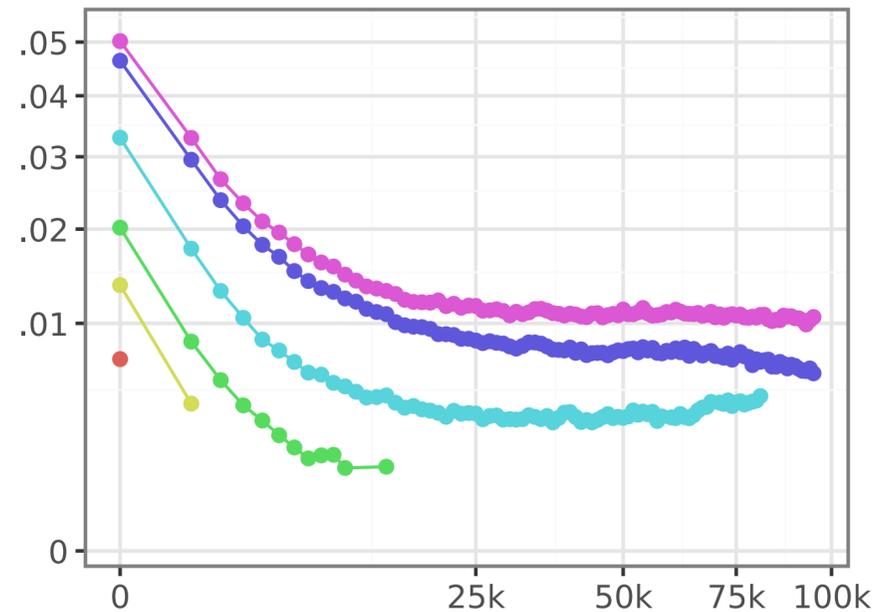


Aggregating the memorisation profile



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

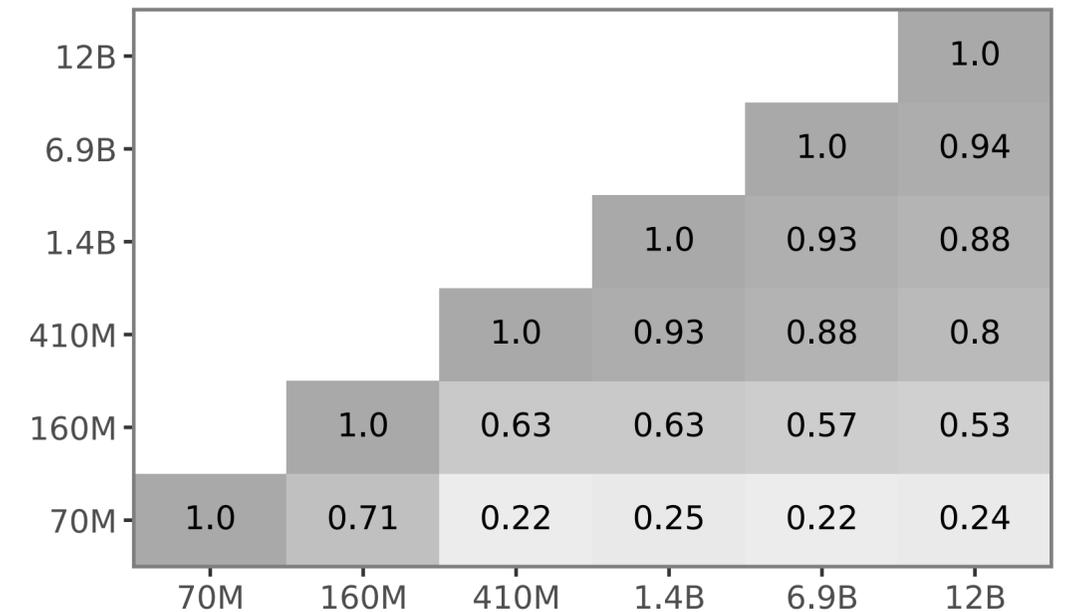
Instantaneous memorisation of a batch at the checkpoint is first seen



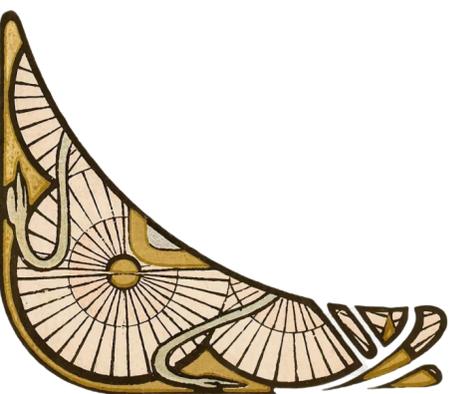
Checkpoints after batch is seen ($g > c$)

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



CAUSAL ESTIMATION OF TOKENISATION BIAS



ACL 2025



Pietro Lesci



Clara Meister



Thomas Hofmann

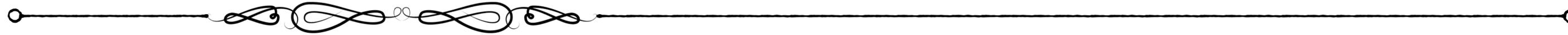


Andreas Vlachos



Tiago Pimentel

What is tokenisation? And why do we need it?



What is tokenisation? And why do we need it?



“Today is a sunny day!”



What is tokenisation? And why do we need it?



“Today is a sunny day!”

Unicode: [U+0054, U+006F, U+0064, U+0061, U+0079, U+0020, ...]

UTF-8: [54, 6F, 64, 61, 79, 20, 69, 73, 20, 61, 20, 73, 75, 6E, 6E, 79, 20, 64, 61, 79, 21]

Graphemes: [“T”, “o”, “d”, “a”, “y”, “_”, “i”, “s”, “_”, “a”, “_”, “s”, “u”, “n”, “n”, “y”, “_”, “d”, “a”, “y”, “!”]



What is tokenisation? And why do we need it?



“Today is a sunny day!”

Small vocabulary
Long sequences

Unicode: [U+0054, U+006F, U+0064, U+0061, U+0079, U+0020, ...]

UTF-8: [54, 6F, 64, 61, 79, 20, 69, 73, 20, 61, 20, 73, 75, 6E, 6E, 79, 20, 64, 61, 79, 21]

Graphemes: [“T”, “o”, “d”, “a”, “y”, “_”, “i”, “s”, “_”, “a”, “_”, “s”, “u”, “n”, “n”, “y”, “_”, “d”, “a”, “y”, “!”]



What is tokenisation? And why do we need it?



“Today is a sunny day!”

Small vocabulary
Long sequences

Unicode: [U+0054, U+006F, U+0064, U+0061, U+0079, U+0020, ...]

UTF-8: [54, 6F, 64, 61, 79, 20, 69, 73, 20, 61, 20, 73, 75, 6E, 6E, 79, 20, 64, 61, 79, 21]

Graphemes: [“T”, “o”, “d”, “a”, “y”, “_”, “i”, “s”, “_”, “a”, “_”, “s”, “u”, “n”, “n”, “y”, “_”, “d”, “a”, “y”, “!”]

Words: [“Today”, “_”, “is”, “a”, “_”, “sunny”, “_”, “day”, “!”]



What is tokenisation? And why do we need it?



“Today is a sunny day!”

Small vocabulary
Long sequences

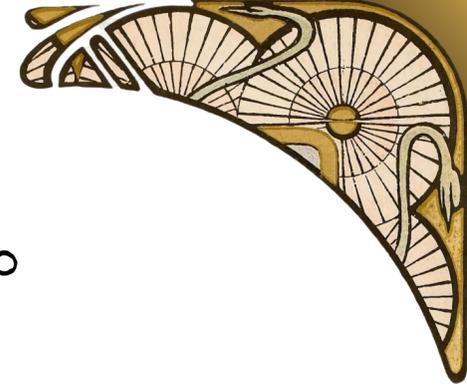
Unicode: [U+0054, U+006F, U+0064, U+0061, U+0079, U+0020, ...]
UTF-8: [54, 6F, 64, 61, 79, 20, 69, 73, 20, 61, 20, 73, 75, 6E, 6E, 79, 20, 64, 61, 79, 21]
Graphemes: [“T”, “o”, “d”, “a”, “y”, “_”, “i”, “s”, “_”, “a”, “_”, “s”, “u”, “n”, “n”, “y”, “_”, “d”, “a”, “y”, “!”]

Short sequences
Large vocabulary

Words: [“Today”, “_”, “is”, “a”, “_”, “sunny”, “_”, “day”, “!”]



What is tokenisation? And why do we need it?



“Today is a sunny day!”

Small vocabulary
Long sequences

Unicode: [U+0054, U+006F, U+0064, U+0061, U+0079, U+0020, ...]
UTF-8: [54, 6F, 64, 61, 79, 20, 69, 73, 20, 61, 20, 73, 75, 6E, 6E, 79, 20, 64, 61, 79, 21]
Graphemes: [“T”, “o”, “d”, “a”, “y”, “_”, “i”, “s”, “_”, “a”, “_”, “s”, “u”, “n”, “n”, “y”, “_”, “d”, “a”, “y”, “!”]

Short sequences
Large vocabulary

Words: [“Today”, “_”, “is”, “a”, “_”, “sunny”, “_”, “day”, “!”]

Trade-off

Subwords: [“_To”, “day”, “_is”, “_a”, “_sun”, “ny”, “_day”, “!”]



What is tokenisation? And why do we need it?



“Today is a sunny day!”

Small vocabulary
Long sequences

Unicode: [U+0054, U+006F, U+0064, U+0061, U+0079, U+0020, ...]
UTF-8: [54, 6F, 64, 61, 79, 20, 69, 73, 20, 61, 20, 73, 75, 6E, 6E, 79, 20, 64, 61, 79, 21]
Graphemes: [“T”, “o”, “d”, “a”, “y”, “_”, “i”, “s”, “_”, “a”, “_”, “s”, “u”, “n”, “n”, “y”, “_”, “d”, “a”, “y”, “!”]

Short sequences
Large vocabulary

Words: [“Today”, “_”, “is”, “a”, “_”, “sunny”, “_”, “day”, “!”]

Trade-off

Subwords: [“_To”, “day”, “_is”, “_a”, “_sun”, “ny”, “_day”, “!”]

Tokenise with $\mathbb{T}(\nu, \tau, \tau^{-1})$



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?

Vocabulary \mathcal{V}

Contains the subwords
available to the model

bottom-up

What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model



bottom-up What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	Initial Alphabet
100	d	
110	n	
115	s	
117	u	
121	y	



bottom-up

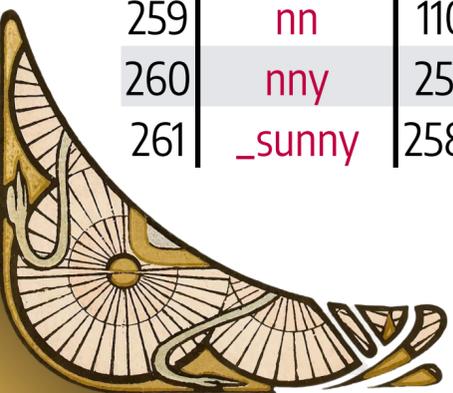
What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	Initial Alphabet
100	d	
110	n	
115	s	
117	u	
121	y	
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



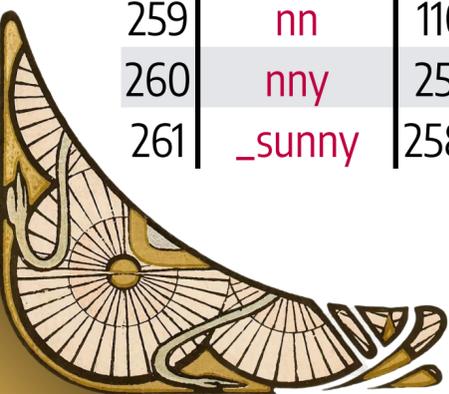
Vocabulary \mathcal{V}

Contains the subwords available to the model

Tokenisation Function τ

Breaks down strings into character spans and maps them to vocabulary items

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	Initial Alphabet
100	d	
110	n	
115	s	
117	u	
121	y	
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

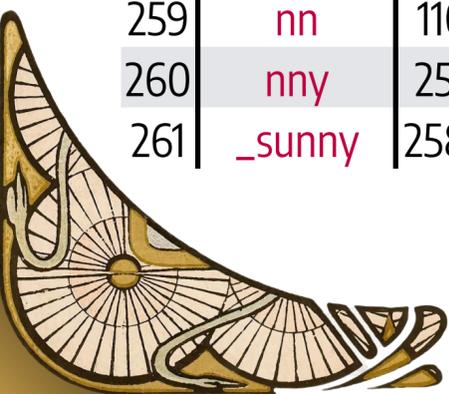
Contains the subwords available to the model

Tokenisation Function τ

Breaks down strings into character spans and maps them to vocabulary items

$\tau(\text{"Today is a sunny day!"})$

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	Initial Alphabet
100	d	
110	n	
115	s	
117	u	
121	y	
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

Tokenisation Function τ

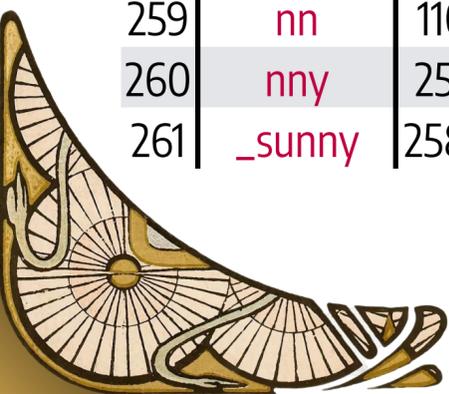
Breaks down strings into character spans and maps them to vocabulary items

$\tau(\text{"Today is a sunny day!"})$

1. Pre-tokenisation*:

$$= \tau([\text{"_Today"}, \text{"_is"}, \text{"_a"}, \text{"_sunny"}, \text{"_day"}, \text{"!"}])$$

ID	Subwords	Merges
32	-	Initial Alphabet
33	!	
84	T	
97	a	
100	d	
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

Tokenisation Function τ

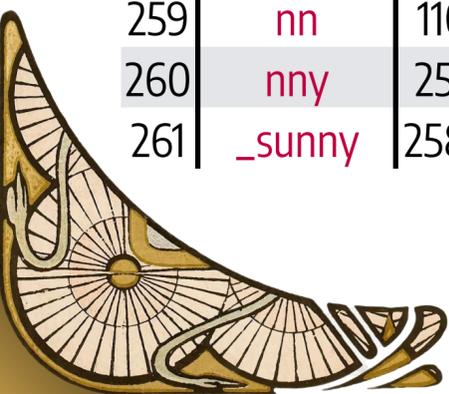
Breaks down strings into character spans and maps them to vocabulary items

ID	Subwords	Merges
32	-	Initial Alphabet
33	!	
84	T	
97	a	
100	d	
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

$\tau(\text{"Today is a sunny day!"})$

1. Pre-tokenisation*:

$$= \tau([\text{"_Today"}, \text{"_is"}, \text{"_a"}, \text{"_sunny"}, \text{"_day"}, \text{"!"}])$$



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

Tokenisation Function τ

Breaks down strings into character spans and maps them to vocabulary items

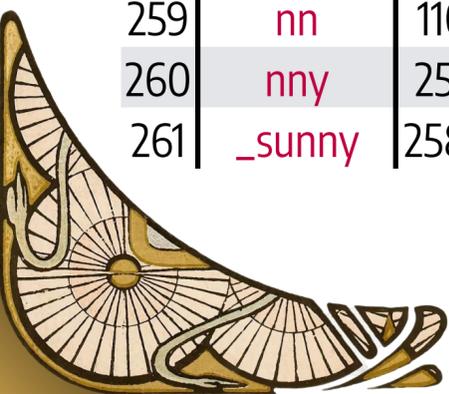
ID	Subwords	Merges
32	-	Initial Alphabet
33	!	
84	T	
97	a	
100	d	
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

$\tau(\text{"Today is a sunny day!"})$

1. Pre-tokenisation*:

$$= \tau([\text{"_Today"}, \text{"_is"}, \text{"_a"}, \text{"_sunny"}, \text{"_day"}, \text{"!"}])$$

$$= [\tau(\text{"_Today"}), \tau(\text{"_is"}), \tau(\text{"_a"}), \tau(\text{"_sunny"}), \tau(\text{"_day"}), \tau(\text{"!"})]$$



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

Tokenisation Function τ

Breaks down strings into character spans and maps them to vocabulary items

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

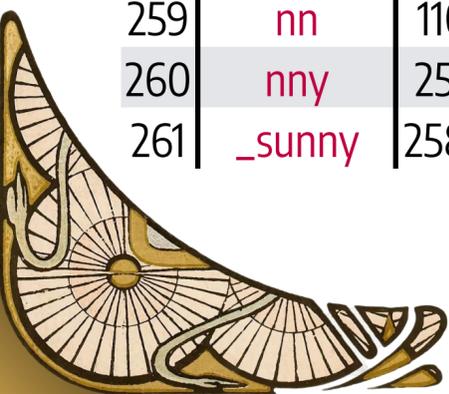
$\tau(\text{"Today is a sunny day!"})$

1. Pre-tokenisation*:

$$= \tau([\text{"_Today"}, \text{"_is"}, \text{"_a"}, \text{"_sunny"}, \text{"_day"}, \text{"!"}])$$

$$= [\tau(\text{"_Today"}), \tau(\text{"_is"}), \tau(\text{"_a"}), \tau(\text{"_sunny"}), \tau(\text{"_day"}), \tau(\text{"!"})]$$

2. Apply tokenisation function (e.g., longest-prefix match):



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

Tokenisation Function τ

Breaks down strings into character spans and maps them to vocabulary items

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

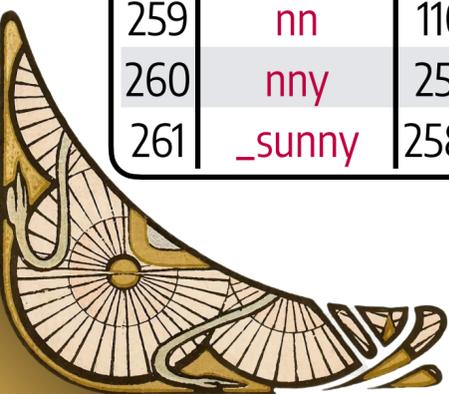
$\tau(\text{"Today is a sunny day!"})$

1. Pre-tokenisation*:

$$= \tau([\text{"_Today"}, \text{"_is"}, \text{"_a"}, \text{"_sunny"}, \text{"_day"}, \text{"!"}])$$

$$= [\tau(\text{"_Today"}), \tau(\text{"_is"}), \tau(\text{"_a"}), \tau(\text{"_sunny"}), \tau(\text{"_day"}), \tau(\text{"!"})]$$

2. Apply tokenisation function (e.g., longest-prefix match):



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

Tokenisation Function τ

Breaks down strings into character spans and maps them to vocabulary items

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

$\tau(\text{"Today is a sunny day!"})$

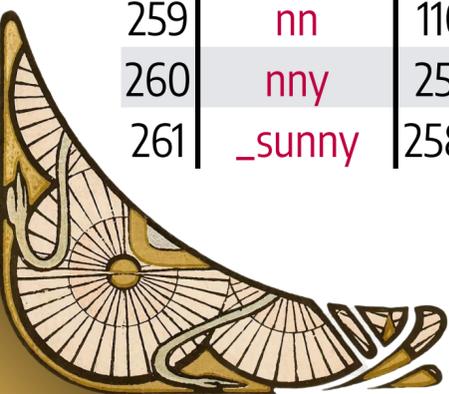
1. Pre-tokenisation*:

$$= \tau([\text{"_Today"}, \text{"_is"}, \text{"_a"}, \text{"_sunny"}, \text{"_day"}, \text{"!"}])$$

$$= [\tau(\text{"_Today"}), \tau(\text{"_is"}), \tau(\text{"_a"}), \tau(\text{"_sunny"}), \tau(\text{"_day"}), \tau(\text{"!"})]$$

2. Apply tokenisation function (e.g., longest-prefix match):

$$= [\text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}, \text{!}]$$



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}

Contains the subwords available to the model

Tokenisation Function τ

Breaks down strings into character spans and maps them to vocabulary items

$\tau(\text{"Today is a sunny day!"})$

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

1. Pre-tokenisation*:

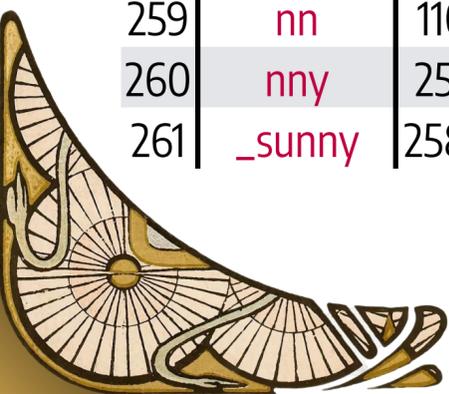
$$= \tau([\text{"_Today"}, \text{"_is"}, \text{"_a"}, \text{"_sunny"}, \text{"_day"}, \text{"!"}])$$

$$= [\tau(\text{"_Today"}), \tau(\text{"_is"}), \tau(\text{"_a"}), \tau(\text{"_sunny"}), \tau(\text{"_day"}), \tau(\text{"!"})]$$

2. Apply tokenisation function (e.g., longest-prefix match):

$$= [\text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}, \text{!}]$$

$$= [299, 318, 300, 256, \mathbf{261}, 324, 33]$$



What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}
 Contains the subwords available to the model

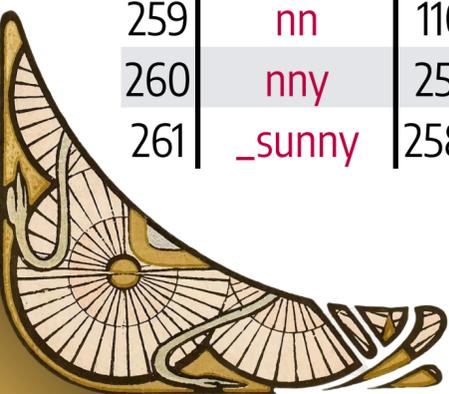
Tokenisation Function τ
 Breaks down strings into character spans and maps them to vocabulary items

Detokenisation Function τ^{-1}
 Converts sequences of subwords to sequences of characters

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

$\tau(\text{"Today is a sunny day!"})$

- Pre-tokenisation*:
 - $= \tau([\text{"_Today"}, \text{"_is"}, \text{"_a"}, \text{"_sunny"}, \text{"_day"}, \text{"!"}])$
 - $= [\tau(\text{"_Today"}), \tau(\text{"_is"}), \tau(\text{"_a"}), \tau(\text{"_sunny"}), \tau(\text{"_day"}), \tau(\text{"!"})]$
- Apply tokenisation function (e.g., longest-prefix match):
 - $= [\text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}, \text{!}]$
 - $= [299, 318, 300, 256, \mathbf{261}, 324, 33]$



*"wHy DoNt YoU jUsT uSe ThE lLaMa ToKeNiZeR??" (Arnett, 2024)

What is a tokeniser $\mathbb{T}(\mathcal{V}, \tau, \tau^{-1})$?



Vocabulary \mathcal{V}
Contains the subwords available to the model

Tokenisation Function τ
Breaks down strings into character spans and maps them to vocabulary items

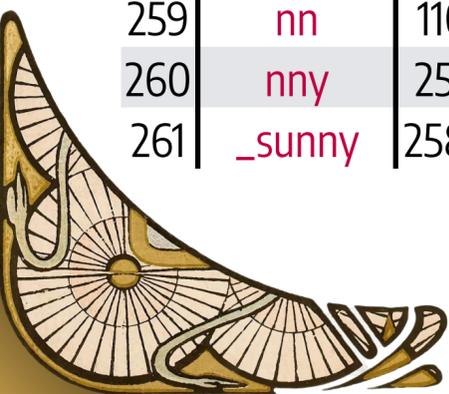
Detokenisation Function τ^{-1}
Converts sequences of subwords to sequences of characters

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

$$\tau(\text{"Today is a sunny day!"})$$

$$\begin{aligned} &\tau^{-1}([299, 318, 300, 256, 261, 324, 33]) \\ &= \tau^{-1}([_To, \text{day}, _is, _a, _sunny, _day, !]) \\ &= \text{"Today is a sunny day!"} \end{aligned}$$

- Pre-tokenisation*:
 - $= \tau([_Today, _is, _a, _sunny, _day, !])$
 - $= [\tau(_Today), \tau(_is), \tau(_a), \tau(_sunny), \tau(_day), \tau(!)]$
- Apply tokenisation function (e.g., longest-prefix match):
 - $= [_To, \text{day}, _is, _a, _sunny, _day, !]$
 - $= [299, 318, 300, 256, \mathbf{261}, 324, 33]$



Language modelling and tokenisation



Probabilities over character-strings

Probabilities over subword-strings

Language modelling and tokenisation



Probabilities over character-strings

$$\begin{aligned} &P(\text{"Today is a sunny day!"}) \\ &= P(\text{"T"}) \\ &\quad \times P(\text{"o"} \mid \text{"T"}) \\ &\quad \times P(\text{"d"} \mid \text{"To"}) \\ &\quad \times \dots \\ &\quad \times P(\langle \text{eos} \rangle^* \mid \text{"Today is a sunny day!"}) \end{aligned}$$

Probabilities over subword-strings

Language modelling and tokenisation



Probabilities over character-strings

$$\begin{aligned} &P(\text{"Today is a sunny day!"}) \\ &= P(\text{"T"}) \\ &\quad \times P(\text{"o"} \mid \text{"T"}) \\ &\quad \times P(\text{"d"} \mid \text{"To"}) \\ &\quad \times \dots \\ &\quad \times P(\langle \text{eos} \rangle^* \mid \text{"Today is a sunny day!"}) \end{aligned}$$

Probabilities over subword-strings

Language modelling and tokenisation

Probabilities over character-strings

$$\begin{aligned} &P(\text{"Today is a sunny day!"}) \\ &= P(\text{"T"}) \\ &\quad \times P(\text{"o"} \mid \text{"T"}) \\ &\quad \times P(\text{"d"} \mid \text{"To"}) \\ &\quad \times \dots \\ &\quad \times P(\langle \text{eos} \rangle^* \mid \text{"Today is a sunny day!"}) \end{aligned}$$

Probabilities over subword-strings

$$\begin{aligned} &P_{\mathbb{T}}([\text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}, !]) \\ &= P_{\mathbb{T}}(\text{_To}) \\ &\quad \times P_{\mathbb{T}}(\text{day} \mid \text{_To}) \\ &\quad \times P_{\mathbb{T}}(\text{_is} \mid \text{_To}, \text{day}) \\ &\quad \times \dots \\ &\quad \times P_{\mathbb{T}}(\langle \text{eos} \rangle \mid \text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}) \end{aligned}$$

Language modelling and tokenisation

Probabilities over character-strings

$$\begin{aligned} &P(\text{"Today is a sunny day!"}) \\ &= P(\text{"T"}) \\ &\quad \times P(\text{"o"} \mid \text{"T"}) \\ &\quad \times P(\text{"d"} \mid \text{"To"}) \\ &\quad \times \dots \\ &\quad \times P(\langle \text{eos} \rangle^* \mid \text{"Today is a sunny day!"}) \end{aligned}$$

Probabilities over subword-strings

$$\begin{aligned} &P_{\mathbb{T}}([\text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}, !]) \\ &= P_{\mathbb{T}}(\text{_To}) \\ &\quad \times P_{\mathbb{T}}(\text{day} \mid \text{_To}) \\ &\quad \times P_{\mathbb{T}}(\text{_is} \mid \text{_To}, \text{day}) \\ &\quad \times \dots \\ &\quad \times P_{\mathbb{T}}(\langle \text{eos} \rangle \mid \text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}) \end{aligned}$$

Tokenisation is just a way to factor a probability over a string into smaller, conditional terms

Language modelling and tokenisation

Probabilities over character-strings

=

Probabilities over subword-strings

$$\begin{aligned} &P(\text{"Today is a sunny day!"}) \\ &= P(\text{"T"}) \\ &\quad \times P(\text{"o"} \mid \text{"T"}) \\ &\quad \times P(\text{"d"} \mid \text{"To"}) \\ &\quad \times \dots \\ &\quad \times P(\langle \text{eos} \rangle^* \mid \text{"Today is a sunny day!"}) \end{aligned}$$

$$\begin{aligned} &P_{\mathbb{T}}([\text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}, !]) \\ &= P_{\mathbb{T}}(\text{_To}) \\ &\quad \times P_{\mathbb{T}}(\text{day} \mid \text{_To}) \\ &\quad \times P_{\mathbb{T}}(\text{_is} \mid \text{_To}, \text{day}) \\ &\quad \times \dots \\ &\quad \times P_{\mathbb{T}}(\langle \text{eos} \rangle \mid \text{_To}, \text{day}, \text{_is}, \text{_a}, \text{_sunny}, \text{_day}) \end{aligned}$$

Tokenisation is just a way to factor a probability over a string into smaller, conditional terms

So, is the choice of the tokeniser important?



$P(\text{"Today is a sunny day!"})$



So, is the choice of the tokeniser important?



\mathcal{V}

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

$P(\text{"Today is a sunny day!"})$

$=$

$P_{\mathcal{T}}([_To, day, _is, _a, _sunny, _day, !])$



So, is the choice of the tokeniser important?



\mathcal{V}

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

P("Today is a sunny day!")

=

$P_{\mathcal{T}}([_To, day, _is, _a, _sunny, _day, !])$

P("Today is a sunny day!")

=

$P_{\mathcal{T}^*}([_To, day, _is, _a, _su, nny, _day, !])$

\mathcal{V}^*

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



So, is the choice of the tokeniser important?



\mathcal{V}

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

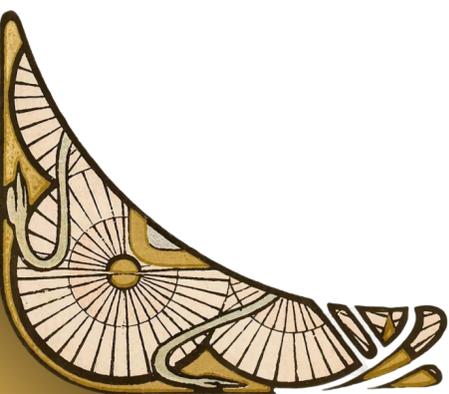
In principle, NO
Tokenisation is simply a way to construct conditional distributions

P("Today is a sunny day!")

$$P_{\mathcal{T}}([_To, day, _is, _a, _sunny, _day, !]) = P_{\mathcal{T}^*}([_To, day, _is, _a, _su, nny, _day, !])$$

\mathcal{V}^*

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



So, is the choice of the tokeniser important?



In principle, NO
 Tokenisation is simply a way to construct conditional distributions

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

P("Today is a sunny day!")

$$\begin{aligned}
 & P_{\mathbb{T}}([_To, day, _is, _a, _sunny, _day, !]) = P_{\mathbb{T}^*}([_To, day, _is, _a, _su, nny, _day, !]) \\
 & \approx P_{LM}([_To, day, _is, _a, _sunny, _day, !]) \approx P_{LM}([_To, day, _is, _a, _su, nny, _day, !])
 \end{aligned}$$



So, is the choice of the tokeniser important?



ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

In principle, NO
Tokenisation is simply a way to construct conditional distributions

In practice, YES
Because of inductive biases and imperfect models

P("Today is a sunny day!")

$$\begin{aligned}
 &P_{\mathbb{T}}([_To, day, _is, _a, _sunny, _day, !]) = P_{\mathbb{T}^*}([_To, day, _is, _a, _su, nny, _day, !]) \\
 &\approx P_{LM}([_To, day, _is, _a, _sunny, _day, !]) \neq P_{LM}([_To, day, _is, _a, _su, nny, _day, !])
 \end{aligned}$$

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



So, is the choice of the tokeniser important?



ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

In principle, NO
Tokenisation is simply a way to construct conditional distributions

In practice, YES
Because of inductive biases and imperfect models

ID	Subwords	Merges
32	-	
33	!	
84	T	
97	a	
100	d	Initial Alphabet
110	n	
115	s	
117	u	
121	y	
256	_a	
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

P("Today is a sunny day!")

$$P_{\mathbb{T}}([_To, day, _is, _a, _sunny, _day, !]) = P_{\mathbb{T}^*}([_To, day, _is, _a, _su, nny, _day, !])$$

$$P_{LM}([_To, day, _is, _a, _sunny, _day, !]) \neq P_{LM}([_To, day, _is, _a, _su, nny, _day, !])$$



The difference in the probability a model assigns to a string because of tokenisation is what we call **tokenisation bias**



Tokenisation bias



In general, we define tokenisation bias as follows:

Definition 1 (General). The tokenisation bias induced by a **property** of the tokeniser is its effect on a model's ability to predict some **character-string** $\mathbf{c} = "c_1c_2\dots c_{|\mathbf{c}|}"$.

Tokenisation bias



In general, we define tokenisation bias as follows:

Definition 1 (General). The tokenisation bias induced by a **property** of the tokeniser is its effect on a model's ability to predict some **character-string** $\mathbf{c} = "c_1c_2\dots c_{|\mathbf{c}|}"$.

The property we focus on in our paper is the **inclusion of a subword in its vocabulary**—i.e., $v \in \mathcal{V}$ —and how it affects predictions on its characters.

Tokenisation bias



In general, we define tokenisation bias as follows:

Definition 1 (General). The tokenisation bias induced by a **property** of the tokeniser is its effect on a model's ability to predict some **character-string** $\mathbf{c} = "c_1c_2\dots c_{|\mathbf{c}|}"$.

The property we focus on in our paper is the **inclusion of a subword in its vocabulary**—i.e., $v \in \mathcal{V}$ —and how it affects predictions on its characters.

Definition 2 (Our paper). The tokenisation bias of $v \in \mathcal{V}$ is its effect on a model's ability to predict \mathbf{c}_v .

Tokenisation bias

In general, we define tokenisation bias as follows:

Definition 1 (General). The tokenisation bias induced by a **property** of the tokeniser is its effect on a model’s ability to predict some **character-string** $\mathbf{c} = “c_1c_2\dots c_{|\mathbf{c}|}”$.

The property we focus on in our paper is the **inclusion of a subword in its vocabulary**—i.e., $v \in \mathcal{V}$ —and how it affects predictions on its characters.

Definition 2 (Our paper). The tokenisation bias of $v \in \mathcal{V}$ is its effect on a model’s ability to predict \mathbf{c}_v .

For example, consider the subword **_sunny**. We are interested in estimating how the probability of its characters, “_sunny” changes if its associated subword is removed from \mathcal{V} , that is:

\mathcal{V} vs. \mathcal{V}^*

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

vs.

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

$\tau(“_sunny”) = [_sunny]$ vs. $\tau^*(“_sunny”) = [_su, _nny]$

Tokenisation bias

In general, we define tokenisation bias as follows:

Definition 1 (General). The tokenisation bias induced by a **property** of the tokeniser is its effect on a model’s ability to predict some **character-string** $\mathbf{c} = “c_1c_2\dots c_{|\mathbf{c}|}”$.

The property we focus on in our paper is the **inclusion of a subword in its vocabulary**—i.e., $v \in \mathcal{V}$ —and how it affects predictions on its characters.

Definition 2 (Our paper). The tokenisation bias of $v \in \mathcal{V}$ is its effect on a model’s ability to predict \mathbf{c}_v .

For example, consider the subword **_sunny**. We are interested in estimating how the probability of its characters, “_sunny” changes if its associated subword is removed from \mathcal{V} , that is:

\mathcal{V} vs. \mathcal{V}^*

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

vs.

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

$\tau(“_sunny”) = [_sunny]$ vs. $\tau^*(“_sunny”) = [_su, _nny]$

$$\psi_v = P_{\text{LM}}(“_sunny” | \text{context}, _sunny \in \mathcal{V}) - P_{\text{LM}}(“_sunny” | \text{context}, _sunny \notin \mathcal{V})$$

Estimating tokenisation bias is hard



$$\psi_v = P_{\text{LM}}(\text{"_sunny"} \mid \text{context, _sunny} \in \mathcal{V}) - P_{\text{LM}}(\text{"_sunny"} \mid \text{context, _sunny} \notin \mathcal{V})$$

Estimating tokenisation bias is hard

$$\psi_v = \underbrace{P_{\text{LM}}(\text{"_sunny"} \mid \text{context, _sunny} \in \mathcal{V})}_{\text{Observable}} - \underbrace{P_{\text{LM}}(\text{"_sunny"} \mid \text{context, _sunny} \notin \mathcal{V})}_{\text{Counterfactual}}$$

Estimating tokenisation bias is hard

$$\psi_v = \underbrace{P_{\text{LM}}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \in \mathcal{V})}_{\text{Observable}} - \underbrace{P_{\text{LM}}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \notin \mathcal{V})}_{\text{Counterfactual}}$$

Given that we train a model using \mathcal{V} , how can we estimate the probability under a **counterfactual** vocabulary?

Estimating tokenisation bias is hard



$$\psi_v = \underbrace{P_{\text{LM}}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \in \mathcal{V})}_{\text{Observable}} - \underbrace{P_{\text{LM}}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \notin \mathcal{V})}_{\text{Counterfactual}}$$

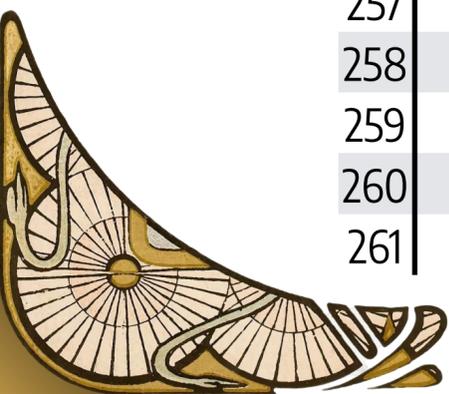
Given that we train a model using \mathcal{V} , how can we estimate the probability under a **counterfactual** vocabulary?

1. Tokenise differently using the same \mathcal{V}

$$P_{\text{LM}}(\text{_su} \mid \text{context}) \times P_{\text{LM}}(\text{nny} \mid \text{context}, \text{_su})$$

\mathcal{V}

ID	Subwords	Merges
...
256	$_a$	32+97
257	$_s$	32+115
258	$_su$	257+117
259	nn	110+110
260	nny	259+121
261	$_sunny$	258+260



Estimating tokenisation bias is hard



$$\psi_v = \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \in \mathcal{V})}_{\text{Observable}} - \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \notin \mathcal{V})}_{\text{Counterfactual}}$$

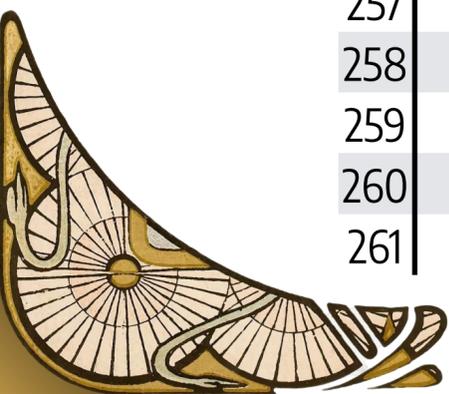
Given that we train a model using \mathcal{V} , how can we estimate the probability under a **counterfactual** vocabulary?

X 1. Tokenise differently using the same \mathcal{V}

$$P_{LM}(\text{_su} \mid \text{context}) \times P_{LM}(\text{nny} \mid \text{context}, \text{_su})$$

\mathcal{V}

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



Estimating tokenisation bias is hard



$$\psi_v = \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \in \mathcal{V})}_{\text{Observable}} - \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \notin \mathcal{V})}_{\text{Counterfactual}}$$

Given that we train a model using \mathcal{V} , how can we estimate the probability under a **counterfactual** vocabulary?

X 1. Tokenise differently using the same \mathcal{V}

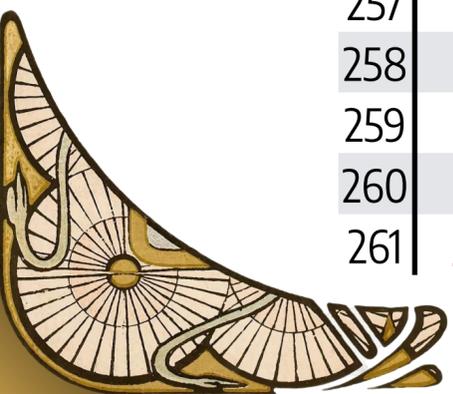
2. Compare to OOV words probabilities, e.g., **"_appoggiatura"**

$$P_{LM}(\text{_su} \mid \text{context}) \times P_{LM}(\text{nny} \mid \text{context}, \text{_su})$$

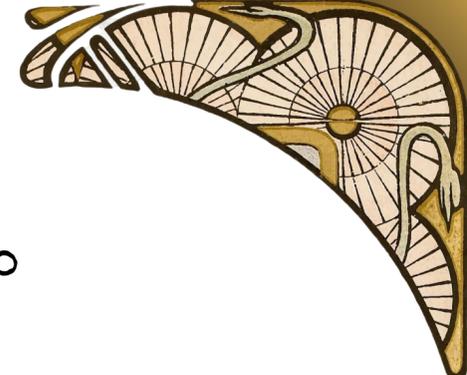
$$P_{LM}(\tau(\text{"_appoggiatura"}) \mid \text{context}) \\ = P_{LM}(\text{_appoggia} \mid \text{context}) \\ \times P_{LM}(\text{tura} \mid \text{context}, \text{_appoggia})$$

\mathcal{V}

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



Estimating tokenisation bias is hard



$$\psi_v = \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \in \mathcal{V})}_{\text{Observable}} - \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \notin \mathcal{V})}_{\text{Counterfactual}}$$

Given that we train a model using \mathcal{V} , how can we estimate the probability under a **counterfactual** vocabulary?

X 1. Tokenise differently using the same \mathcal{V}

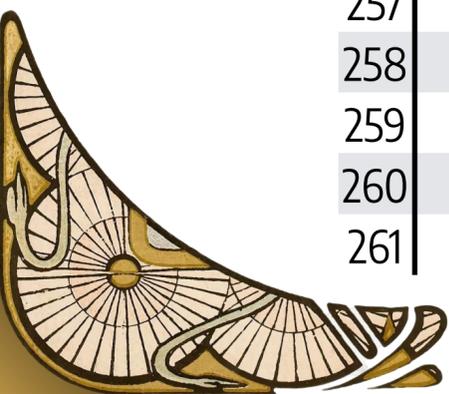
X 2. Compare to OOV words probabilities, e.g., **"_appoggiatura"**

$$P_{LM}(\text{_su} \mid \text{context}) \times P_{LM}(\text{nny} \mid \text{context}, \text{_su})$$

$$P_{LM}(\tau(\text{"_appoggiatura"}) \mid \text{context}) \\ = P_{LM}(\text{_appoggia} \mid \text{context}) \\ \times P_{LM}(\text{tura} \mid \text{context}, \text{_appoggia})$$

\mathcal{V}

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



Estimating tokenisation bias is hard



$$\psi_v = \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \in \mathcal{V})}_{\text{Observable}} - \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \notin \mathcal{V})}_{\text{Counterfactual}}$$

Given that we train a model using \mathcal{V} , how can we estimate the probability under a **counterfactual** vocabulary?

X 1. Tokenise differently using the same \mathcal{V}

$$P_{LM}(\text{_su} \mid \text{context}) \times P_{LM}(\text{nny} \mid \text{context}, \text{_su})$$

\mathcal{V}

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

X 2. Compare to OOV words probabilities, e.g., **"_appoggiatura"**

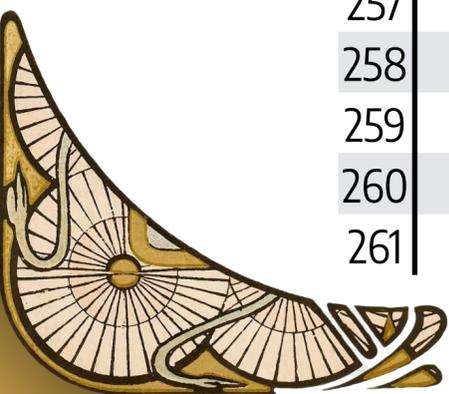
$$P_{LM}(\tau(\text{"_appoggiatura"}) \mid \text{context}) \\ = P_{LM}(\text{_appoggia} \mid \text{context}) \\ \times P_{LM}(\text{tura} \mid \text{context}, \text{_appoggia})$$

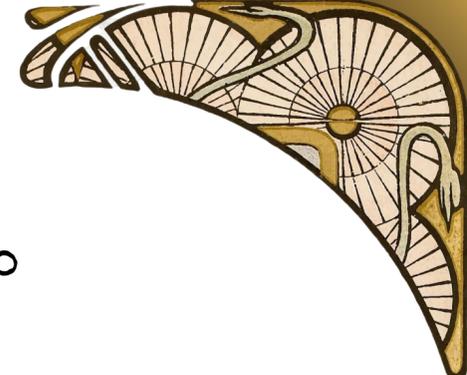
3. Retrain the model with a different \mathcal{V}^*

$$P_{LM^*}(\text{_su} \mid \text{context}) \times P_{LM^*}(\text{nny} \mid \text{context}, \text{_su})$$

\mathcal{V}^*

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260





Estimating tokenisation bias is hard

$$\psi_v = \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \in \mathcal{V})}_{\text{Observable}} - \underbrace{P_{LM}(\text{"_sunny"} \mid \text{context}, \text{_sunny} \notin \mathcal{V})}_{\text{Counterfactual}}$$

Given that we train a model using \mathcal{V} , how can we estimate the probability under a **counterfactual** vocabulary?

X 1. Tokenise differently using the same \mathcal{V}

$$P_{LM}(\text{_su} \mid \text{context}) \times P_{LM}(\text{nny} \mid \text{context}, \text{_su})$$

\mathcal{V}

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260

X 2. Compare to OOV words probabilities, e.g., **"_appoggiatura"**

$$P_{LM}(\tau(\text{"_appoggiatura"}) \mid \text{context})$$

$$= P_{LM}(\text{_appoggia} \mid \text{context})$$

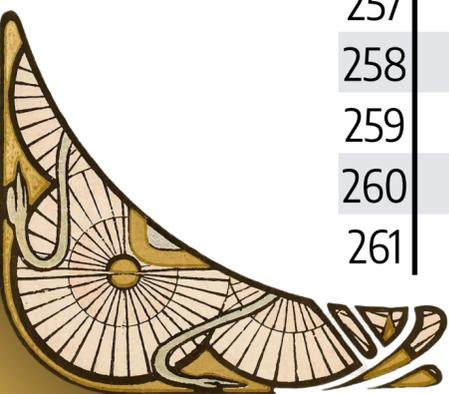
$$\times P_{LM}(\text{tura} \mid \text{context}, \text{_appoggia})$$

X 3. Retrain the model with a different \mathcal{V}^*

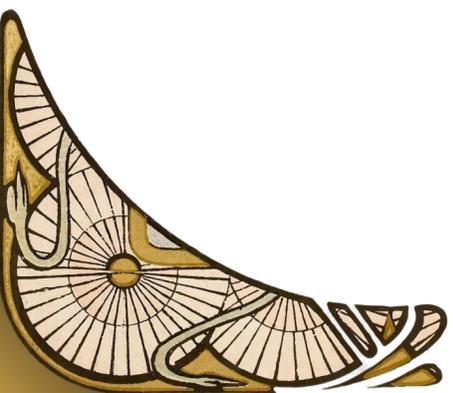
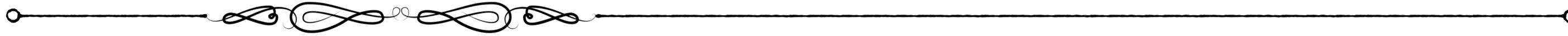
$$P_{LM^*}(\text{_su} \mid \text{context}) \times P_{LM^*}(\text{nny} \mid \text{context}, \text{_su})$$

\mathcal{V}^*

ID	Subwords	Merges
...
256	_a	32+97
257	_s	32+115
258	_su	257+117
259	nn	110+110
260	nny	259+121
261	_sunny	258+260



Tokenisation bias as a causal effect



Tokenisation bias as a causal effect



Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!



Tokenisation bias as a causal effect



Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

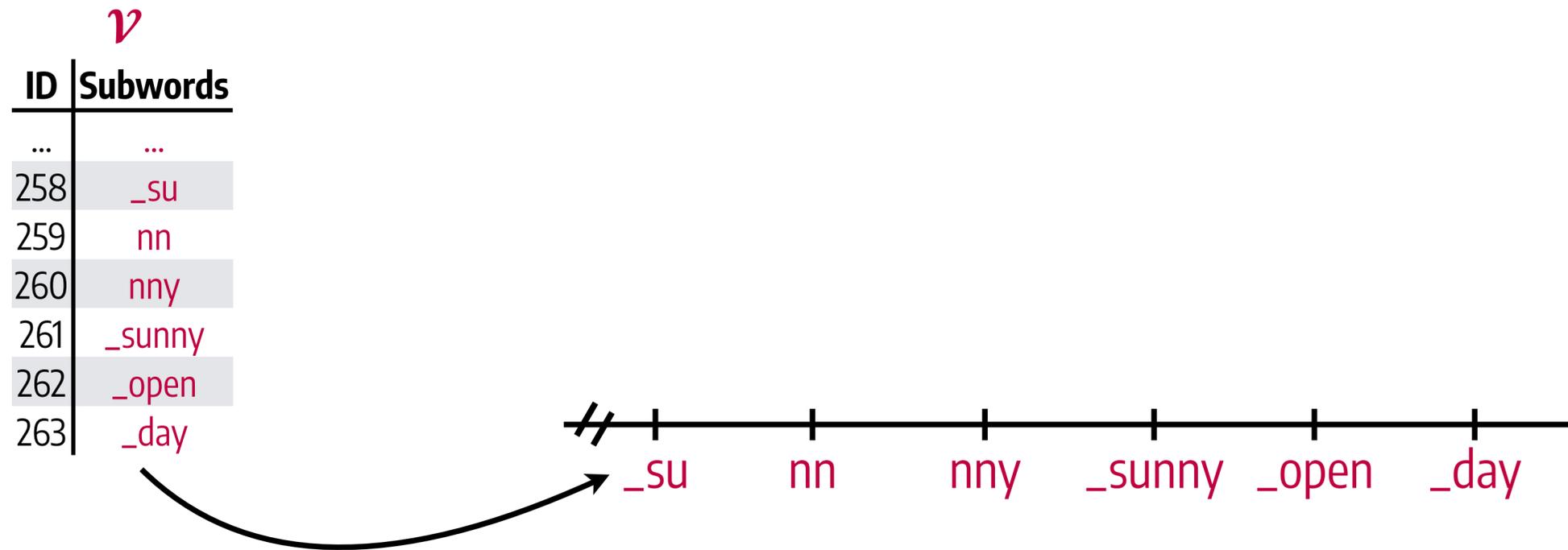
\mathcal{V}

ID	Subwords
...	...
258	_su
259	nn
260	nny
261	_sunny
262	_open
263	_day



Tokenisation bias as a causal effect

Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

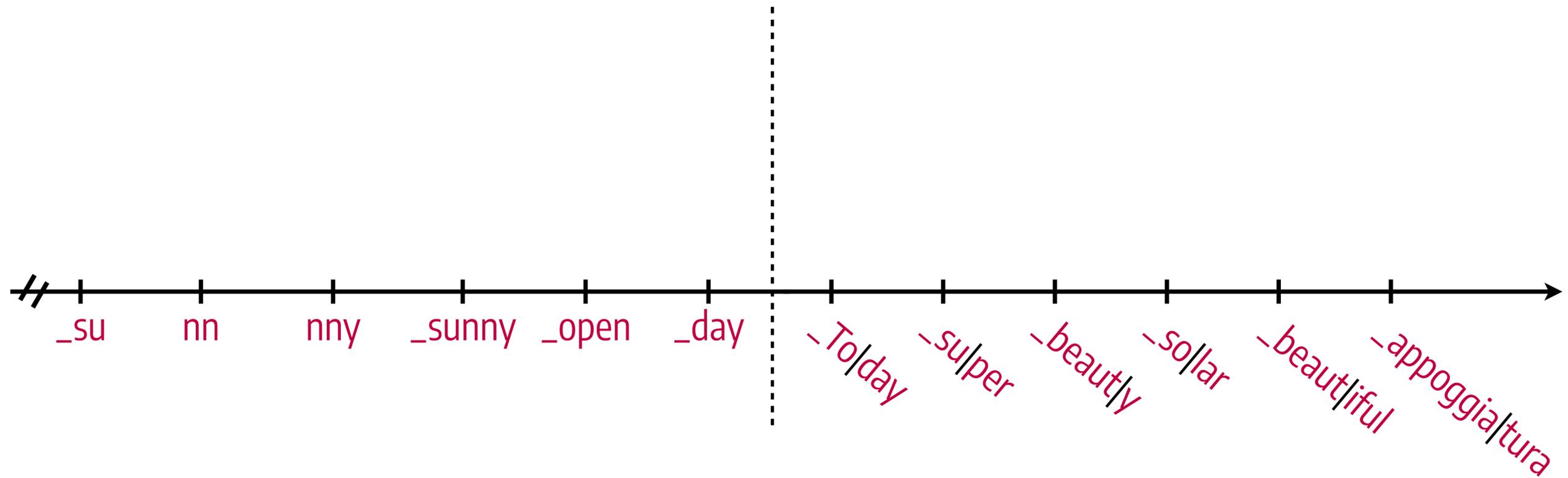


Tokenisation bias as a causal effect

Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

\mathcal{V}

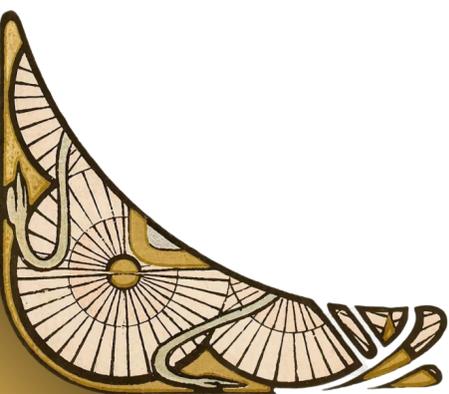
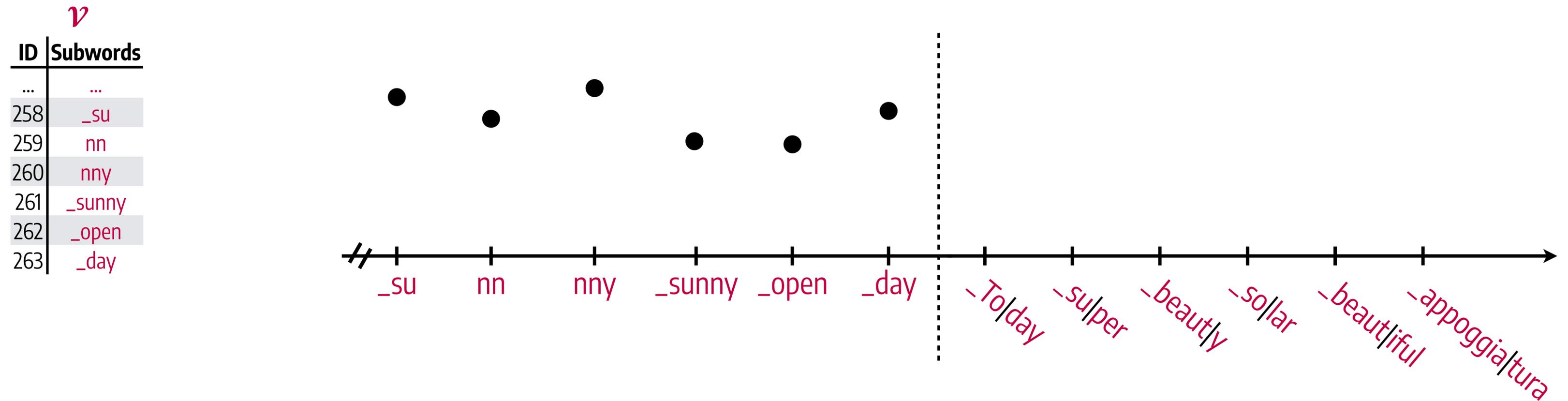
ID	Subwords
...	...
258	_su
259	nn
260	nny
261	_sunny
262	_open
263	_day



Tokenisation bias as a causal effect



Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

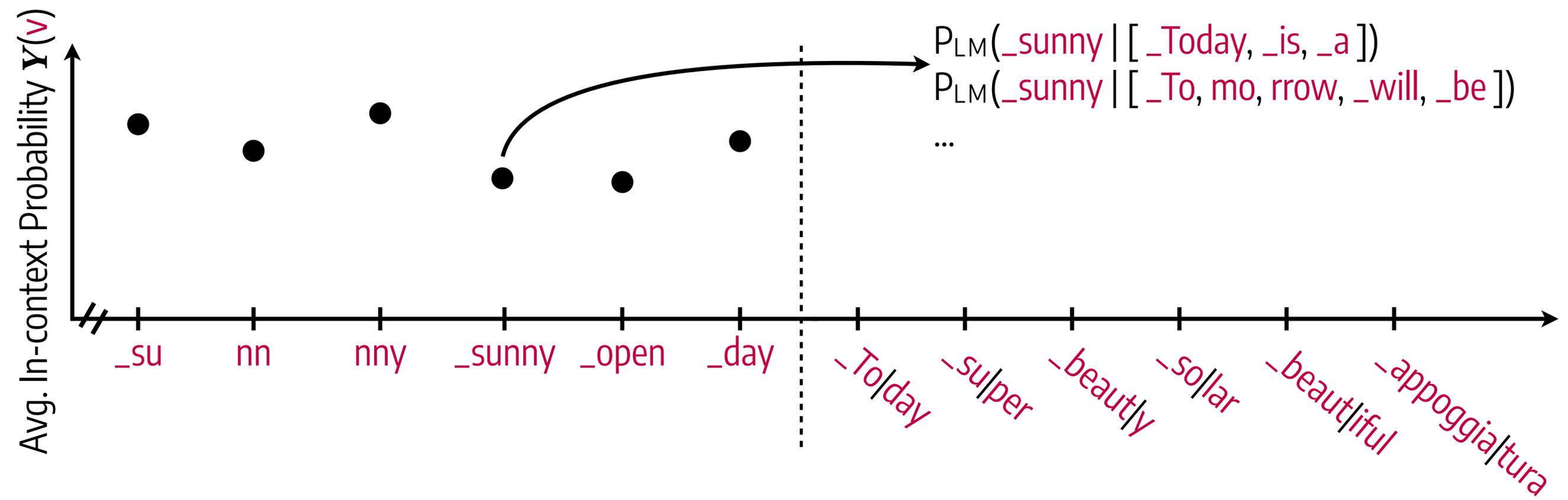


Tokenisation bias as a causal effect



Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

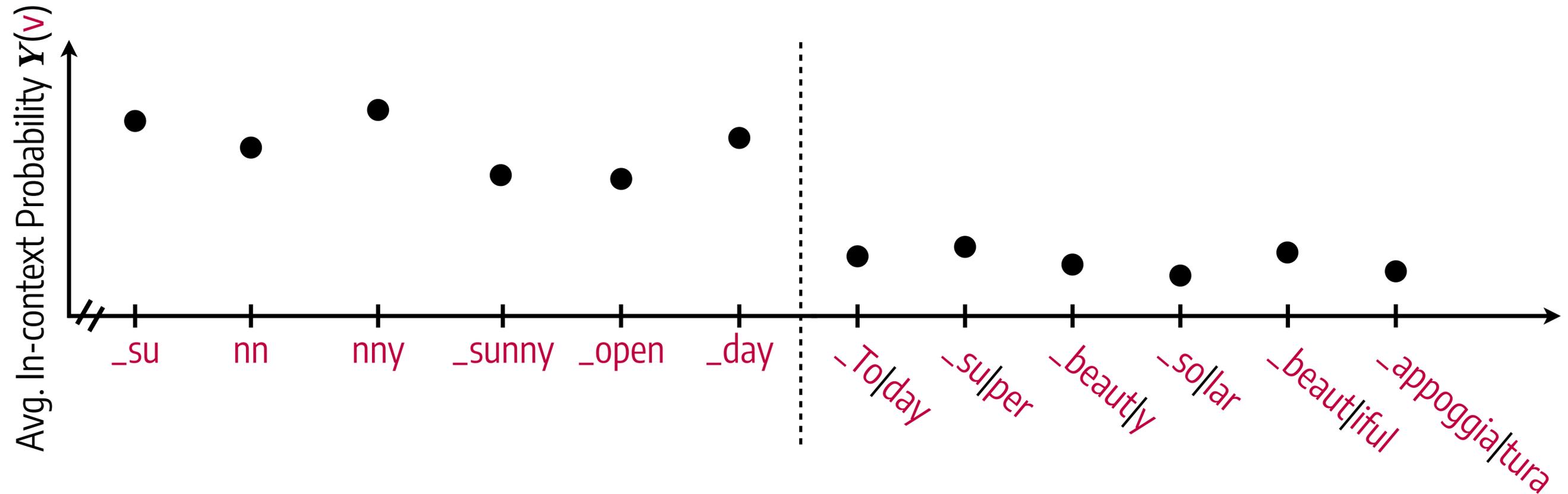
ID	Subwords
...	...
258	_su
259	nn
260	nny
261	_sunny
262	_open
263	_day



Tokenisation bias as a causal effect

Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

ID	Subwords
...	...
258	<i>_su</i>
259	<i>nn</i>
260	<i>nny</i>
261	<i>_sunny</i>
262	<i>_open</i>
263	<i>_day</i>

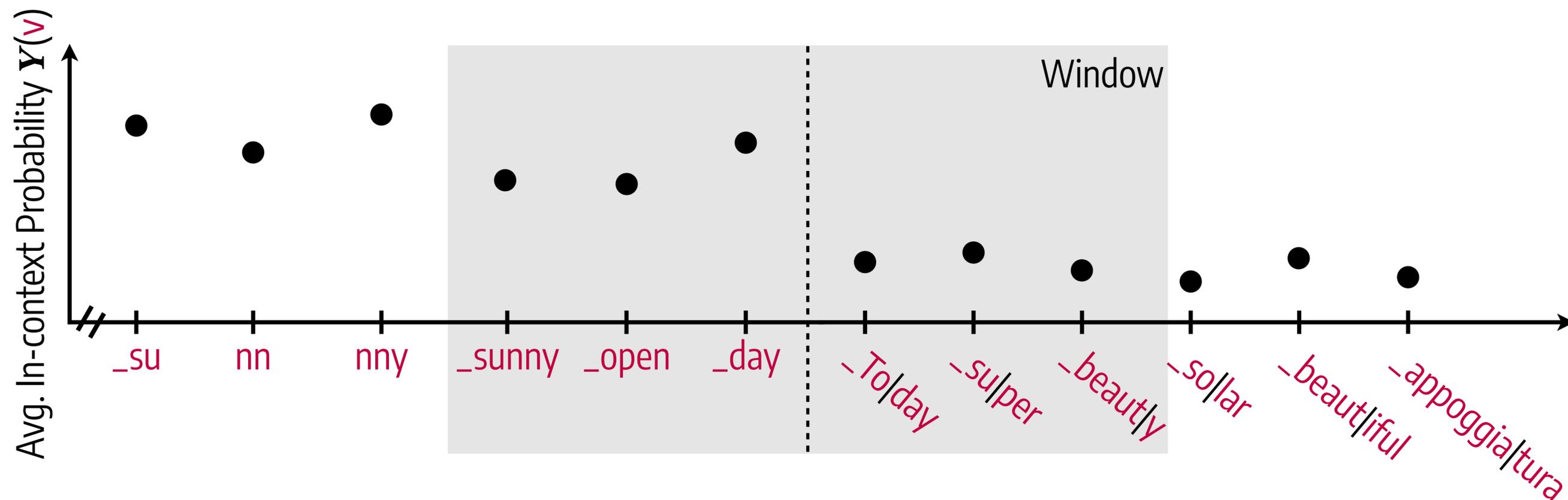


Tokenisation bias as a causal effect

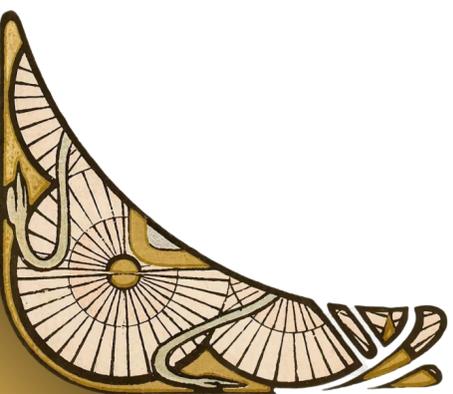


Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

ID	Subwords
...	...
258	_su
259	nn
260	nny
261	_sunny
262	_open
263	_day



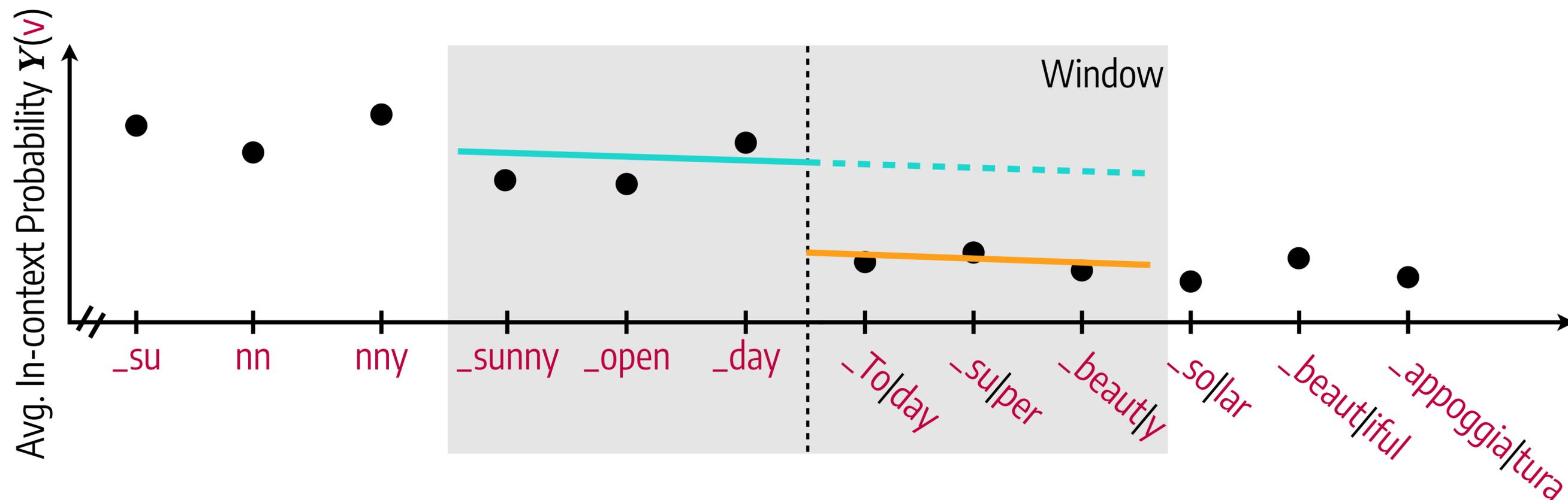
Subwords in the window have **similar frequency** but on one side appear as one vs. two subwords in \mathcal{V} . Thus, by comparing their average in-context probability we can estimate tokenisation bias without re-training



Tokenisation bias as a causal effect

Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

ID	Subwords
...	...
258	_su
259	nn
260	nny
261	_sunny
262	_open
263	_day



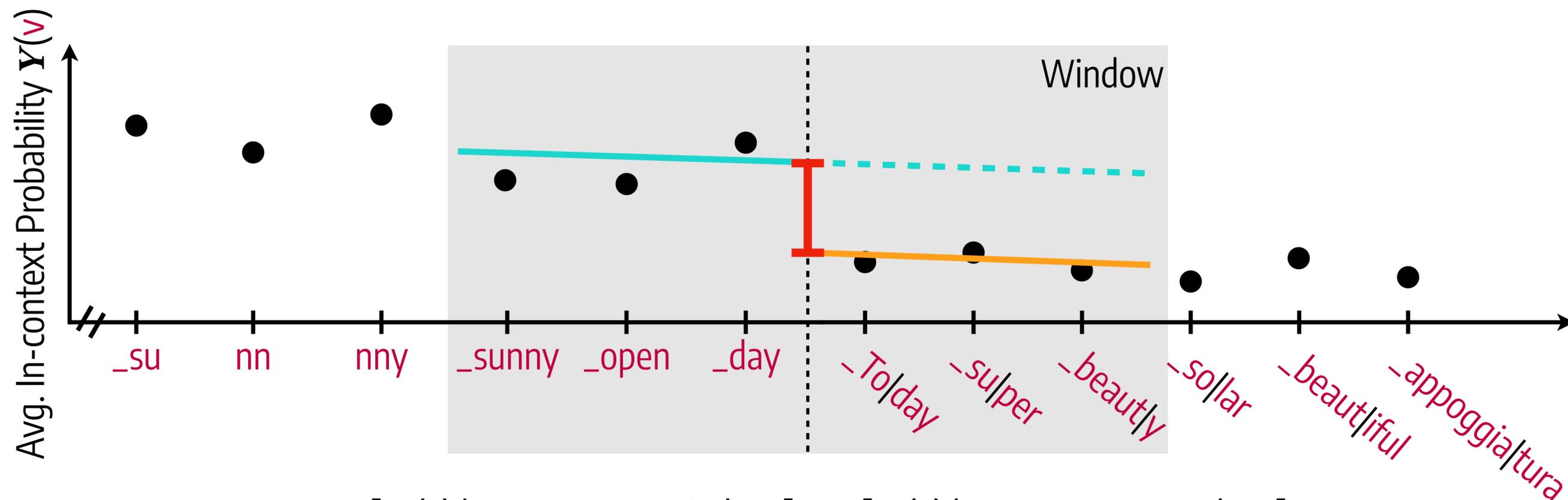
Subwords in the window have **similar frequency** but on one side appear as one vs. two subwords in \mathcal{V} . Thus, by comparing their average in-context probability we can estimate tokenisation bias without re-training

Tokenisation bias as a causal effect



Bottom-up tokenisers—including BPE (Sennrich et al., 2016) and WordPiece (Schuster and Nakajima, 2012)—select subwords iteratively, one-at-a-time, to compose their vocabularies. We can exploit this fact!

ID	Subwords
...	...
258	_su
259	nn
260	nny
261	_sunny
262	_open
263	_day



$$\psi = \mathbb{E}_v[\mathbf{Y}(v) \mid v \in \mathcal{V} \wedge v \in \text{Window}] - \mathbb{E}_v[\mathbf{Y}(v) \mid v \notin \mathcal{V} \wedge v \in \text{Window}]$$



Subwords in the window have **similar frequency** but on one side appear as one vs. two subwords in \mathcal{V} . Thus, by comparing their average in-context probability we can estimate tokenisation bias without re-training



Experimental Setup



Experimental Setup



Tokeniser training. Our method requires knowing which subwords would be added after the cutoff. Thus, we train* BPE with a large vocabulary of 320k subwords and derive smaller vocabularies.

*This information is typically not available even for open-source models and tokeniser construction details are often unreproducible

Experimental Setup



Tokeniser training. Our method requires knowing which subwords would be added after the cutoff. Thus, we train* BPE with a large vocabulary of 320k subwords and derive smaller vocabularies.

Model training. Llama-style 57M-param (non-embedding) model trained on the MiniPile for 50k steps (~13B tokens) using BPE-32k.



Experimental Setup



Tokeniser training. Our method requires knowing which subwords would be added after the cutoff. Thus, we train* BPE with a large vocabulary of 320k subwords and derive smaller vocabularies.

Model training. Llama-style 57M-param (non-embedding) model trained on the MiniPile for 50k steps (~13B tokens) using BPE-32k.

Model Evaluations and RDD setup. Window size 5k around the cutoff (i.e, 27k-37k). Collect log-probabilities for these subwords on the MiniPile validation set. We run a simple linear regression and validate more complex models.



Experimental Setup



Tokeniser training. Our method requires knowing which subwords would be added after the cutoff. Thus, we train* BPE with a large vocabulary of 320k subwords and derive smaller vocabularies.

Model training. Llama-style 57M-param (non-embedding) model trained on the MiniPile for 50k steps (~13B tokens) using BPE-32k.

Model Evaluations and RDD setup. Window size 5k around the cutoff (i.e, 27k-37k). Collect log-probabilities for these subwords on the MiniPile validation set. We run a simple linear regression and validate more complex models.

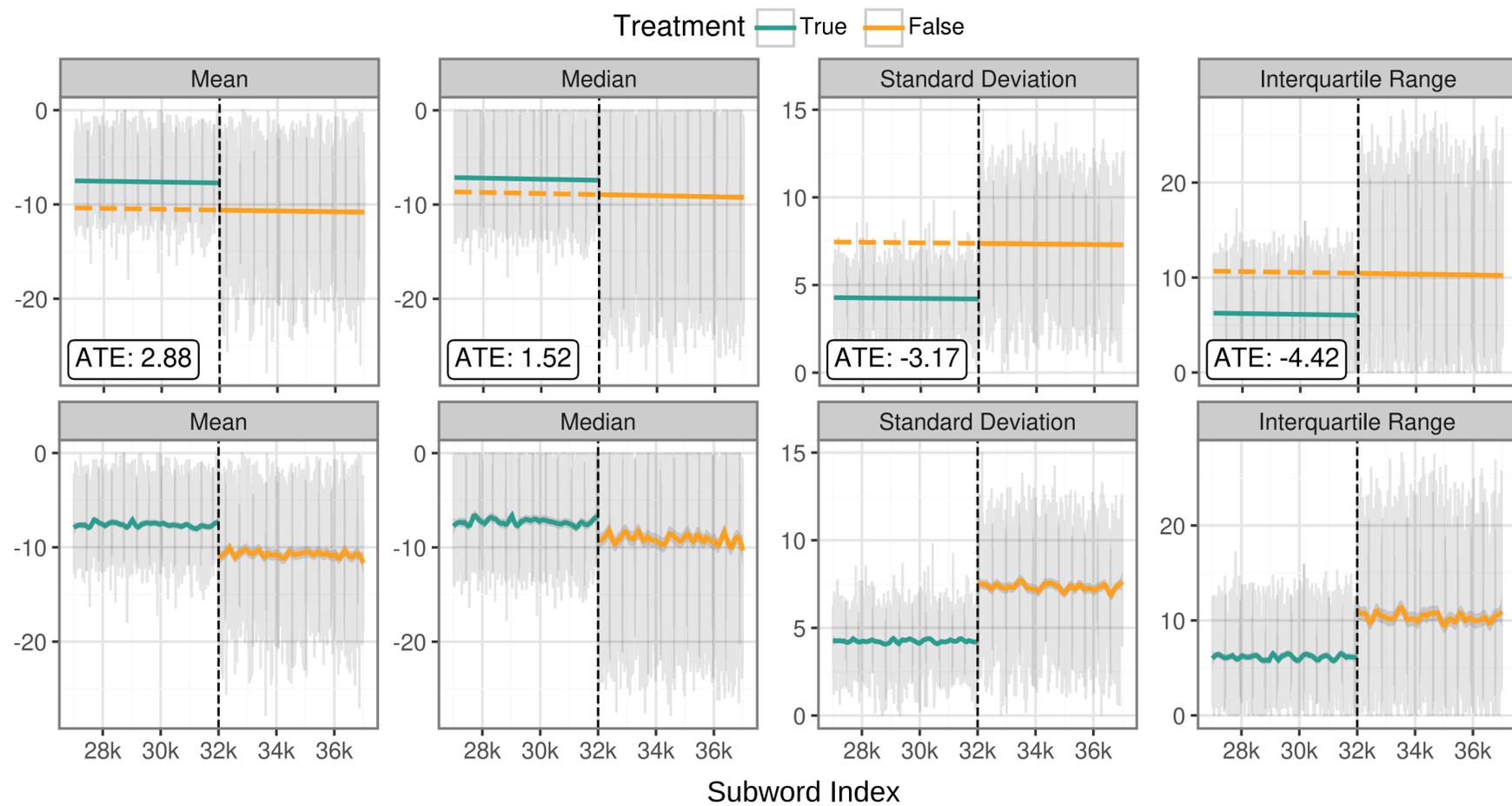
We also study how tokenisation bias changes as a function of:

1. **Vocab sizes:** 8k, 32k, 128k (BPE)
2. **Tokenisation:** WordPiece (WP) & hybrid BPE2WP (BPE vocabulary + WP tokenisation function)
3. **Model sizes:** 340M & 850M
4. **Embedding tying & Dataset:** 100M model trained with/without tied embeddings on 20B tokens from Fineweb-Edu
5. **Aggregation function:** Mean, median, standard deviation, and interquartile range of the subword probability across contexts

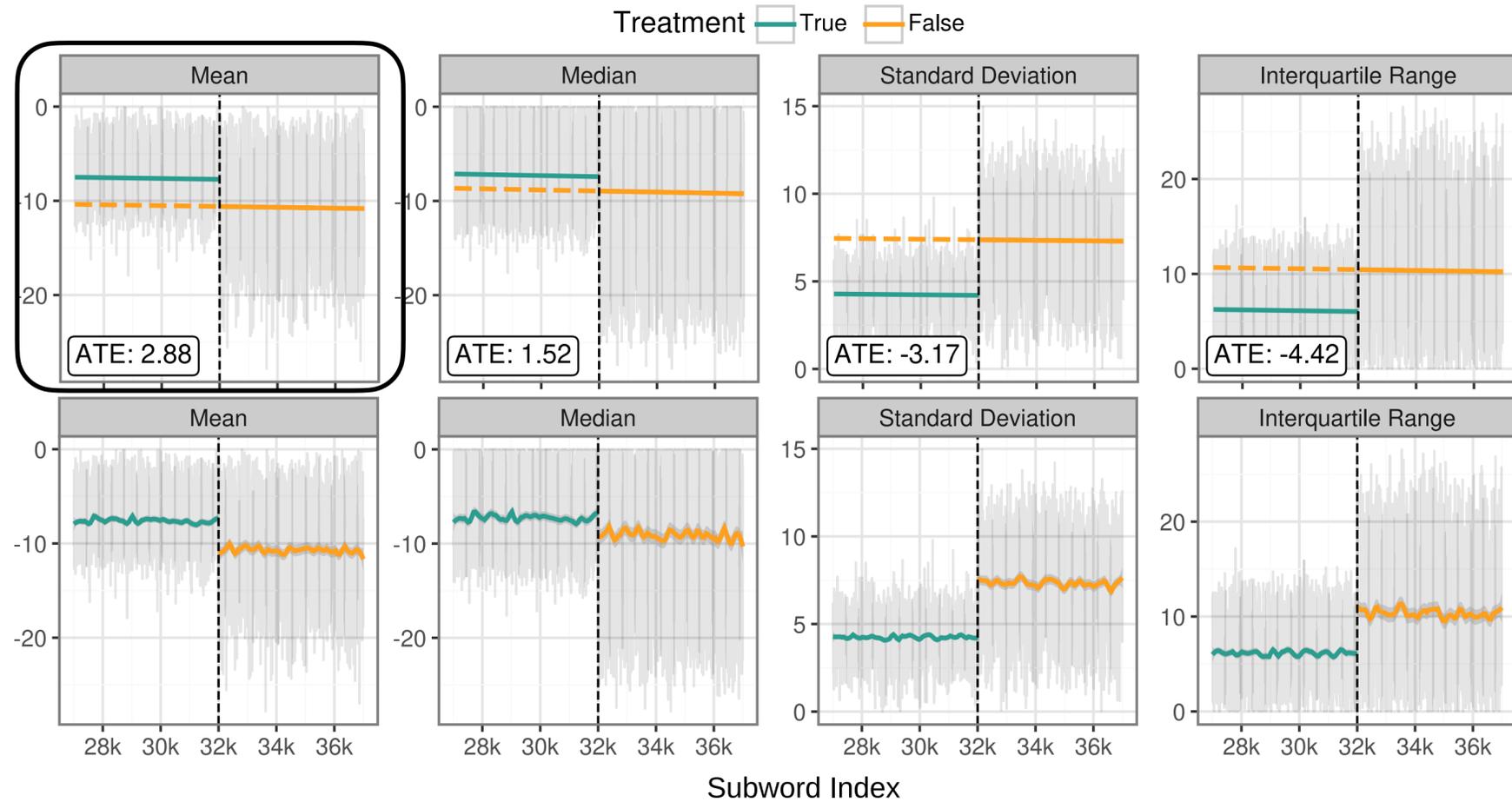


*This information is typically not available even for open-source models and tokeniser construction details are often unreproducible

Tokenisation bias at the end of training



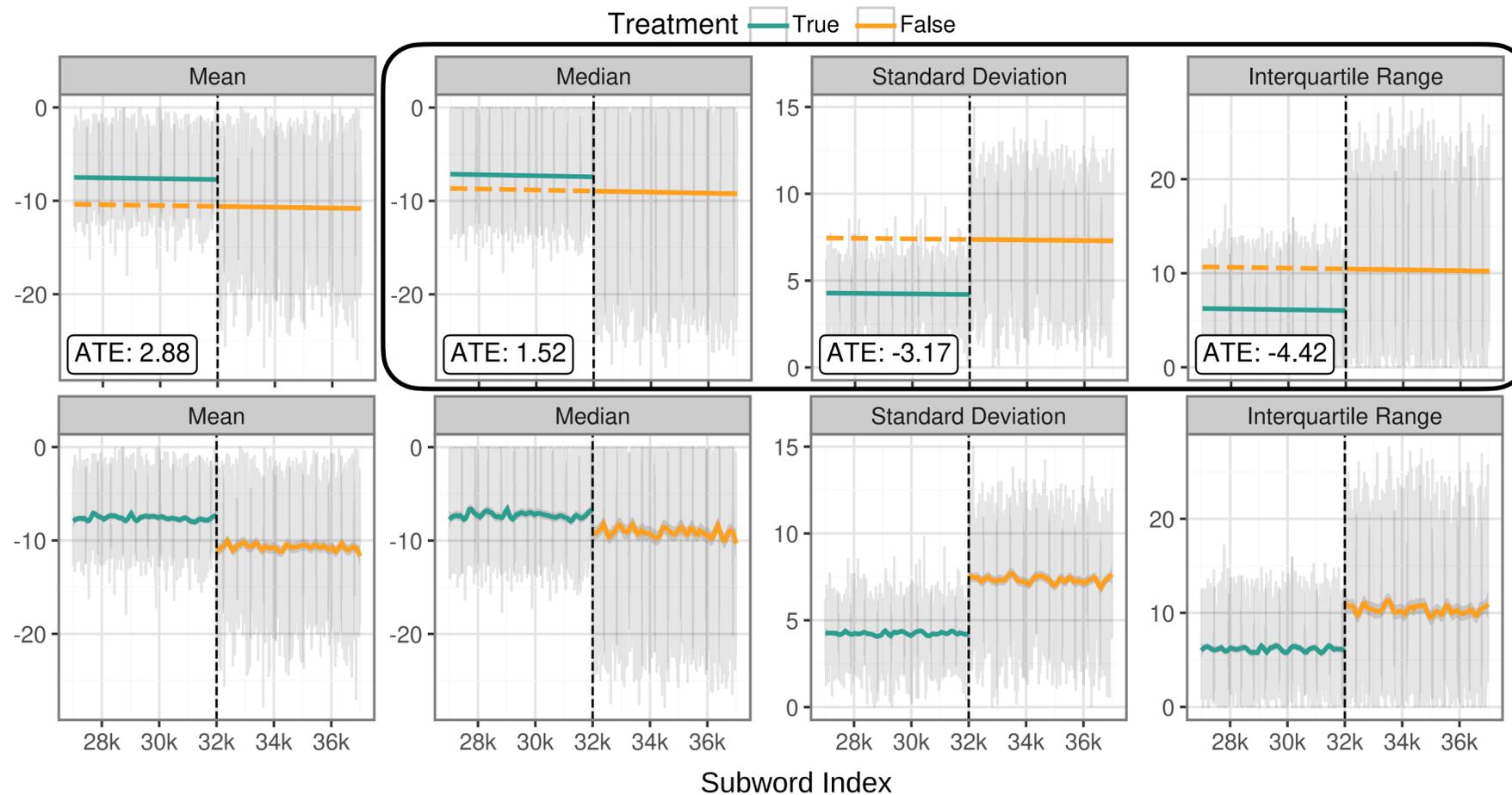
Tokenisation bias at the end of training



- We estimate this bias to be 2.88 nats. This means a can be assigned roughly 17x less probability due solely to tokenisation



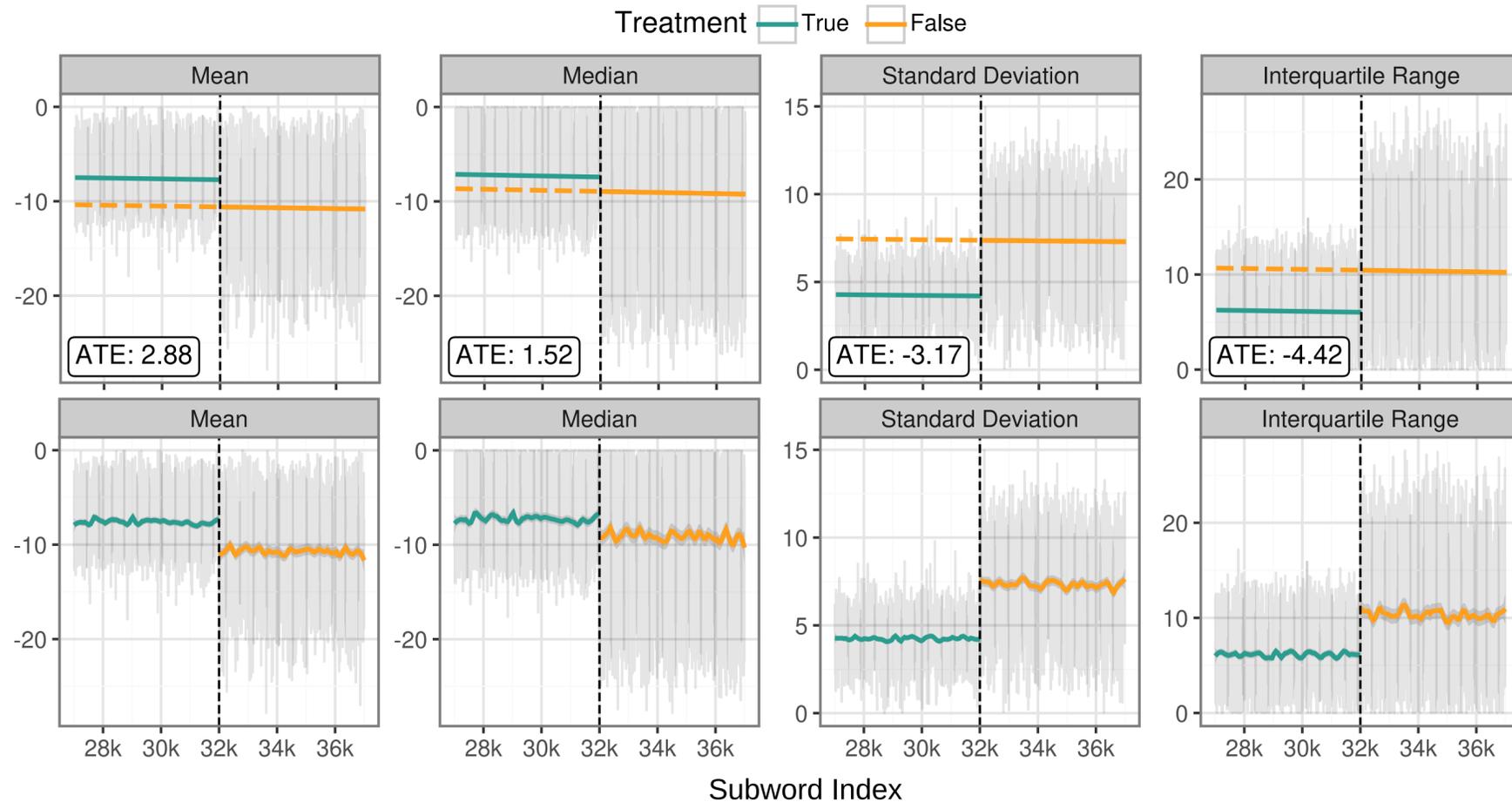
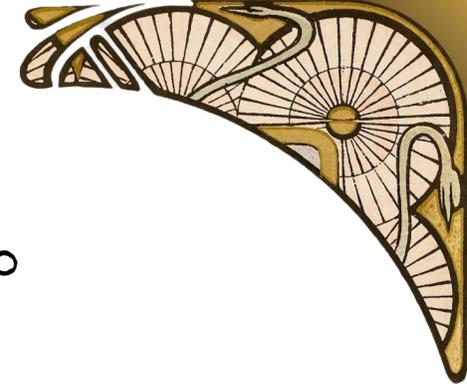
Tokenisation bias at the end of training



- We estimate this bias to be 2.88 nats. This means a can be assigned roughly 17x less probability due solely to tokenisation
- Subword inclusion significantly affects output stability: in-vocabulary subwords exhibit far less variation in log-probability across contexts than out-of-vocabulary ones



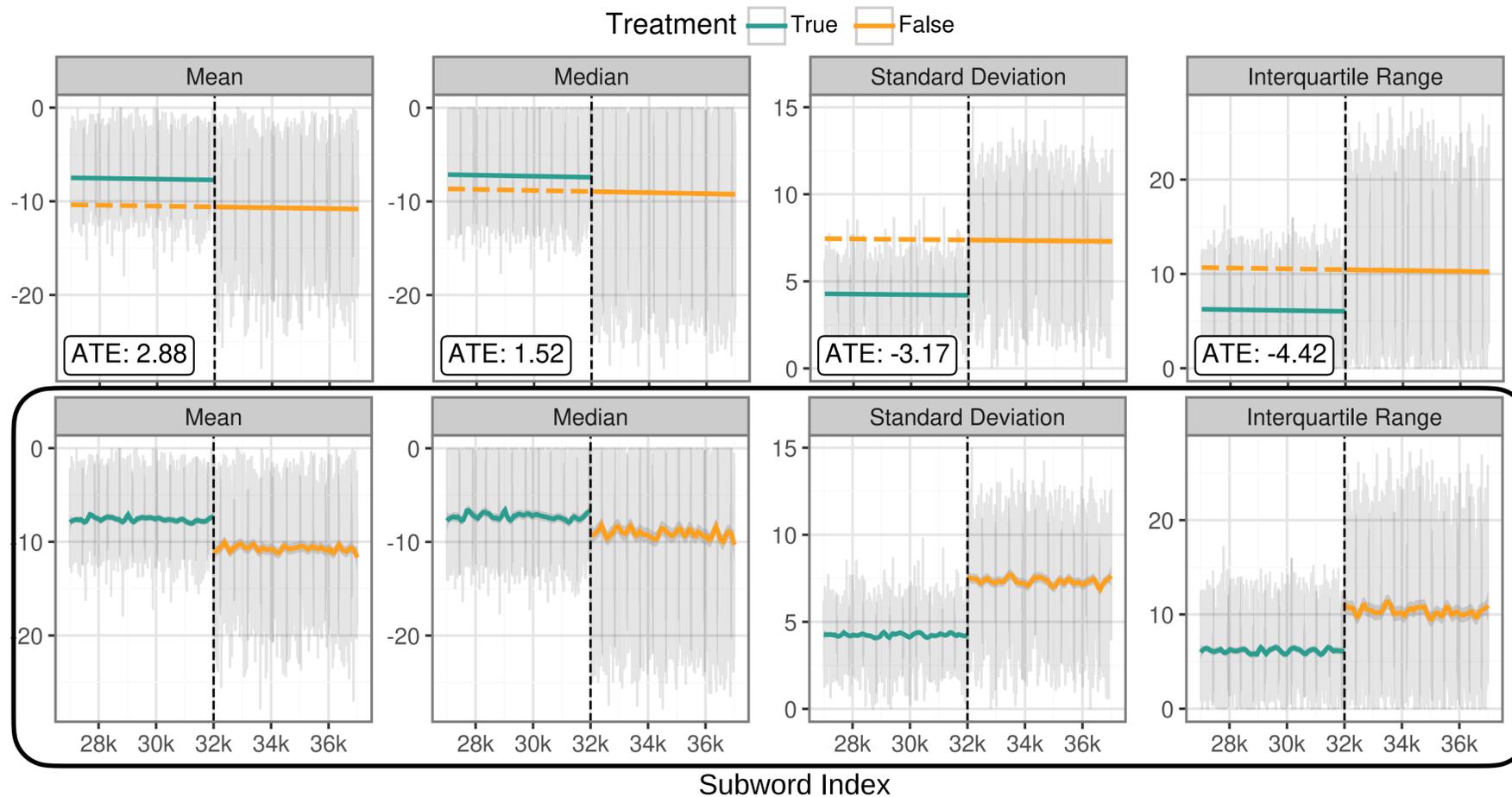
Tokenisation bias at the end of training



- We estimate this bias to be 2.88 nats. This means a can be assigned roughly 17x less probability due solely to tokenisation
- Subword inclusion significantly affects output stability: in-vocabulary subwords exhibit far less variation in log-probability across contexts than out-of-vocabulary ones
- RDD estimates are robust to



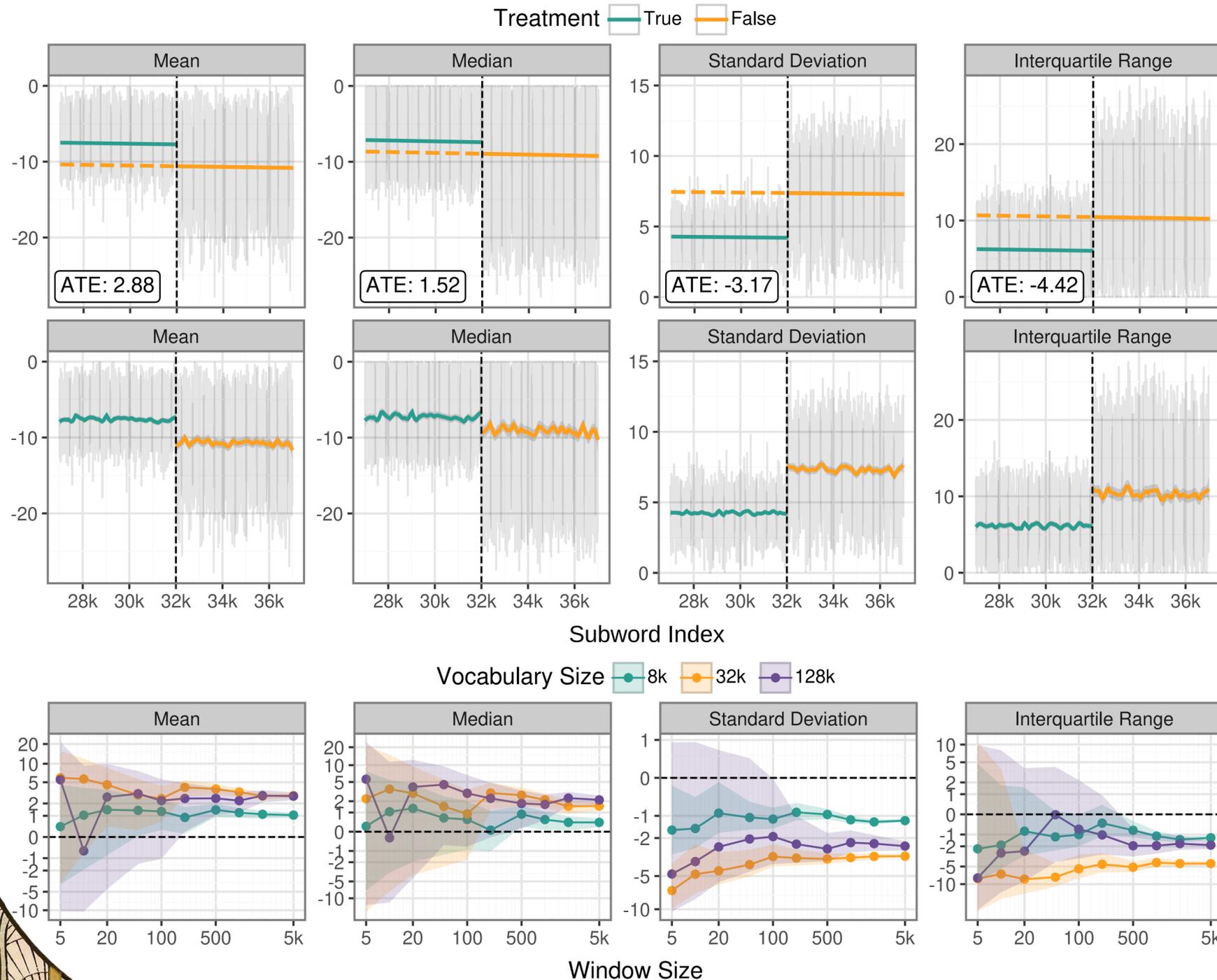
Tokenisation bias at the end of training



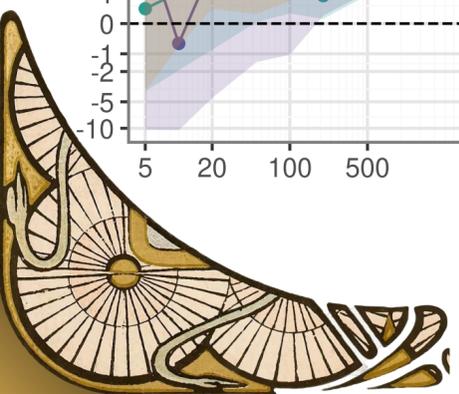
- We estimate this bias to be 2.88 nats. This means a can be assigned roughly 17x less probability due solely to tokenisation
- Subword inclusion significantly affects output stability: in-vocabulary subwords exhibit far less variation in log-probability across contexts than out-of-vocabulary ones
- RDD estimates are robust to
 - Functional form



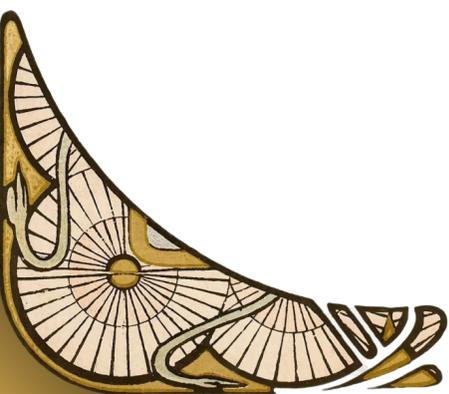
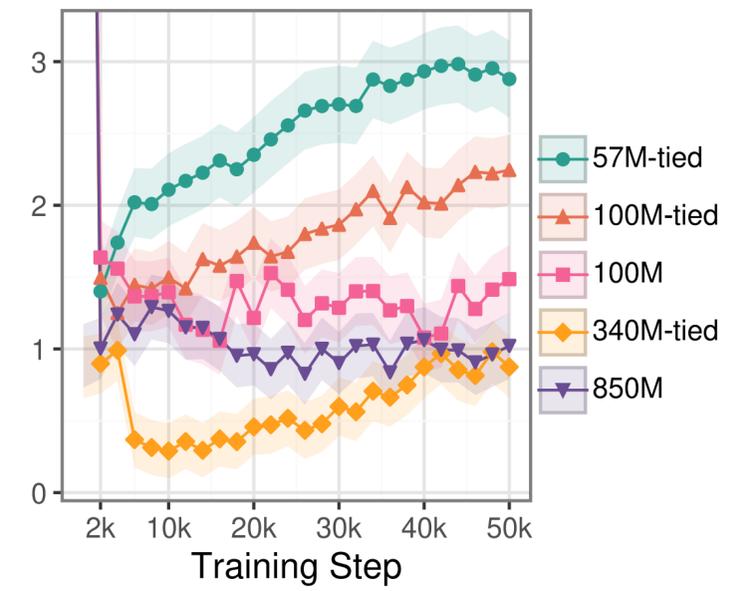
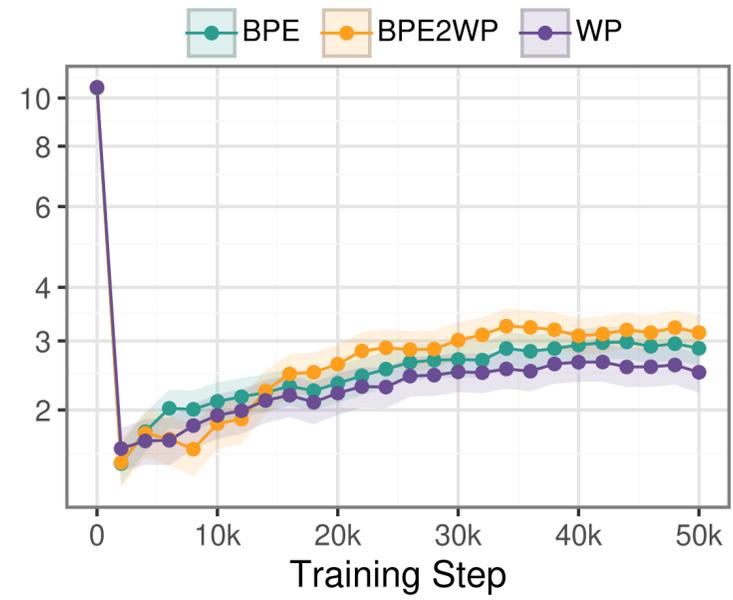
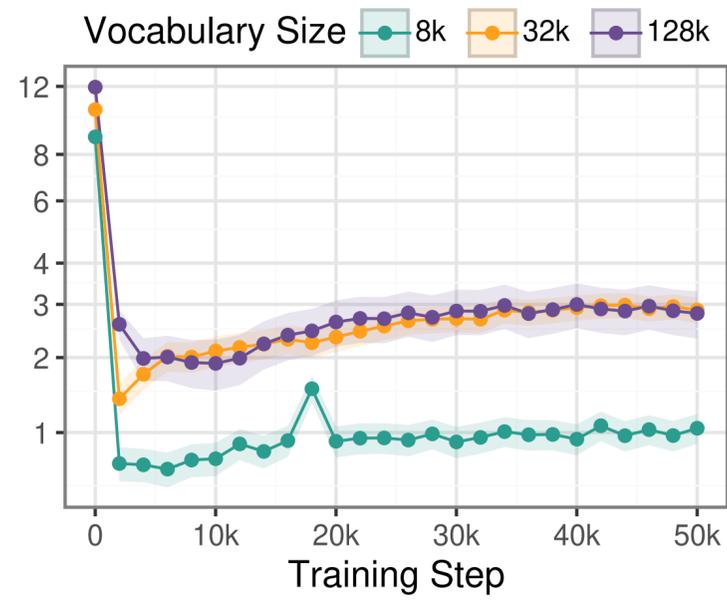
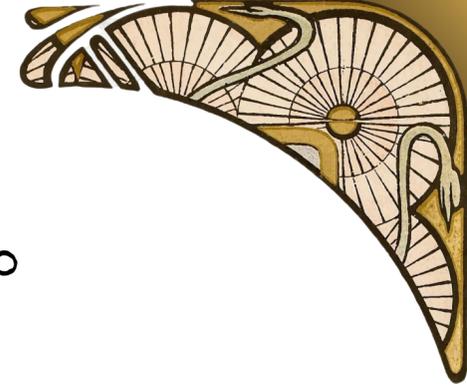
Tokenisation bias at the end of training



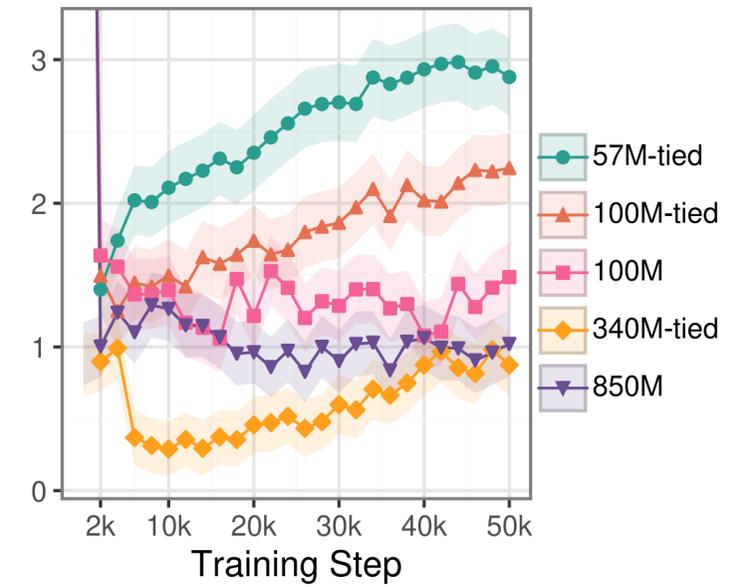
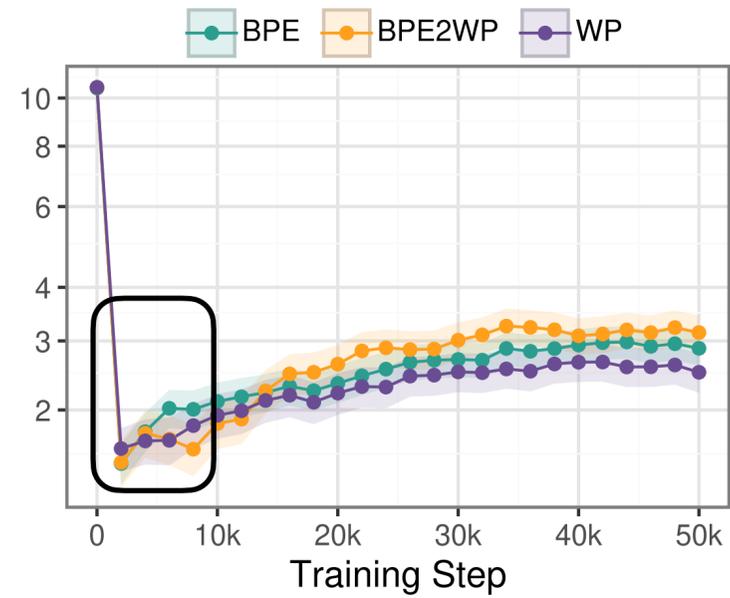
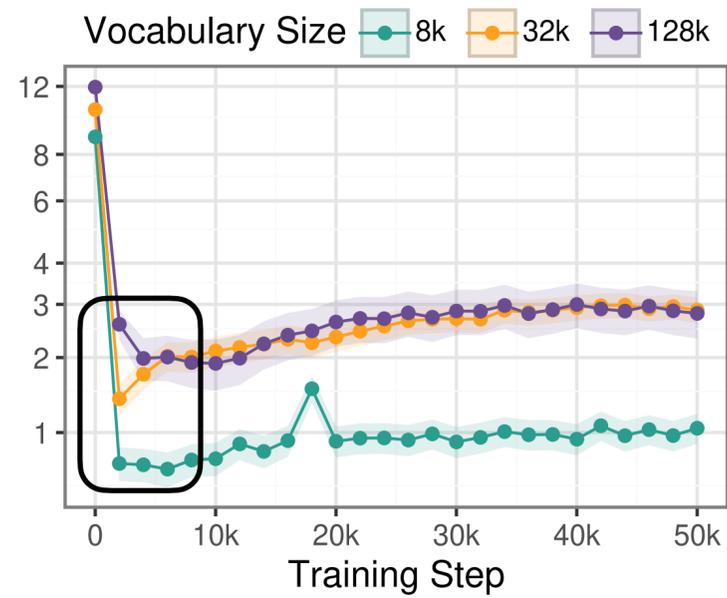
- We estimate this bias to be 2.88 nats. This means a can be assigned roughly 17x less probability due solely to tokenisation
- Subword inclusion significantly affects output stability: in-vocabulary subwords exhibit far less variation in log-probability across contexts than out-of-vocabulary ones
- RDD estimates are robust to
 - Functional form
 - Window size, for sizes >500



Tokenisation bias across training



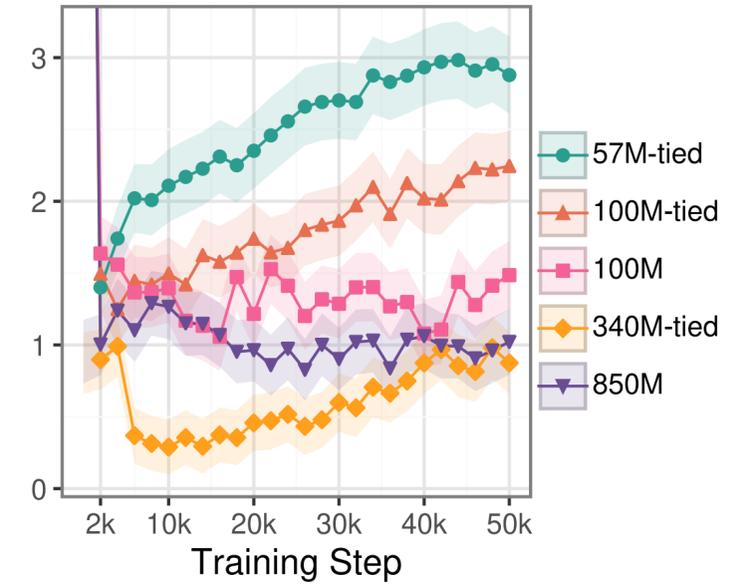
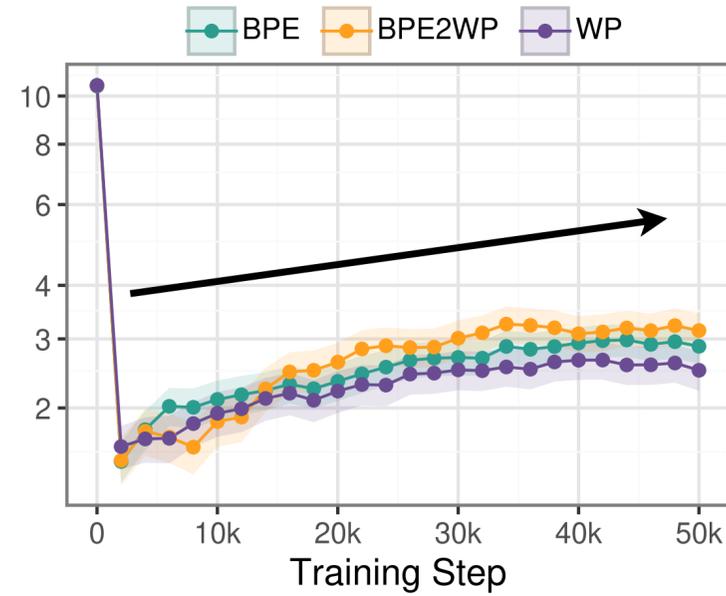
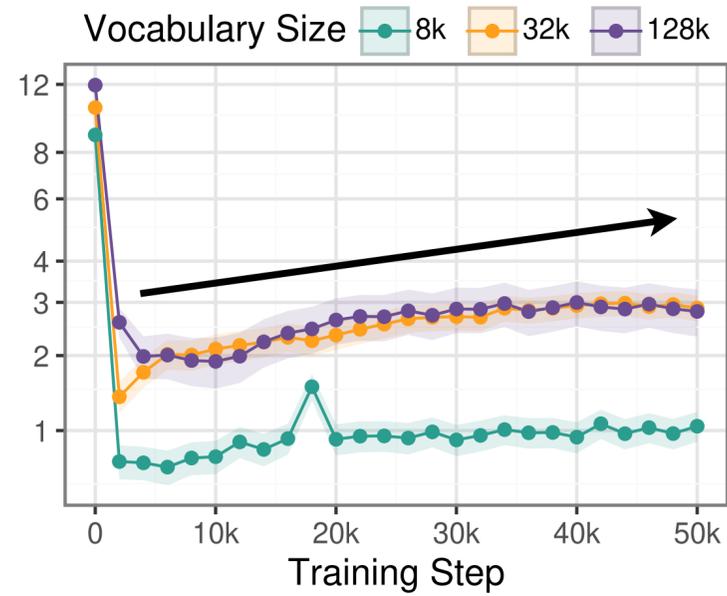
Tokenisation bias across training



- Drops sharply by step 2k, as expected



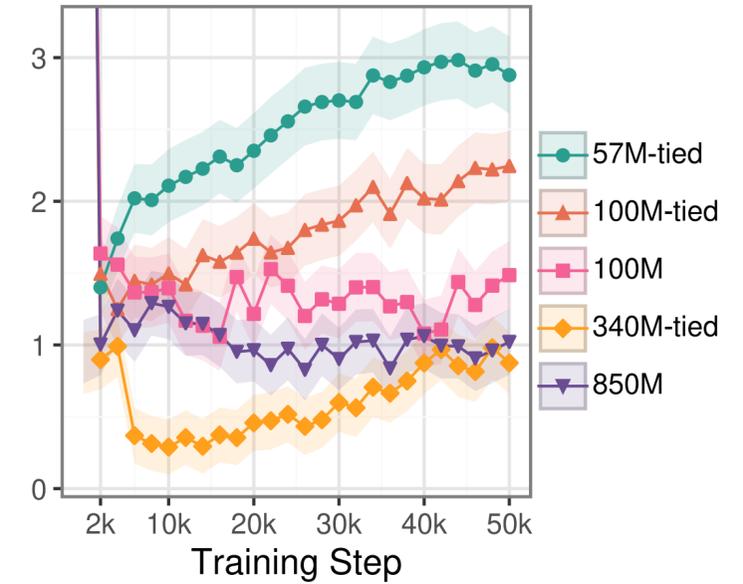
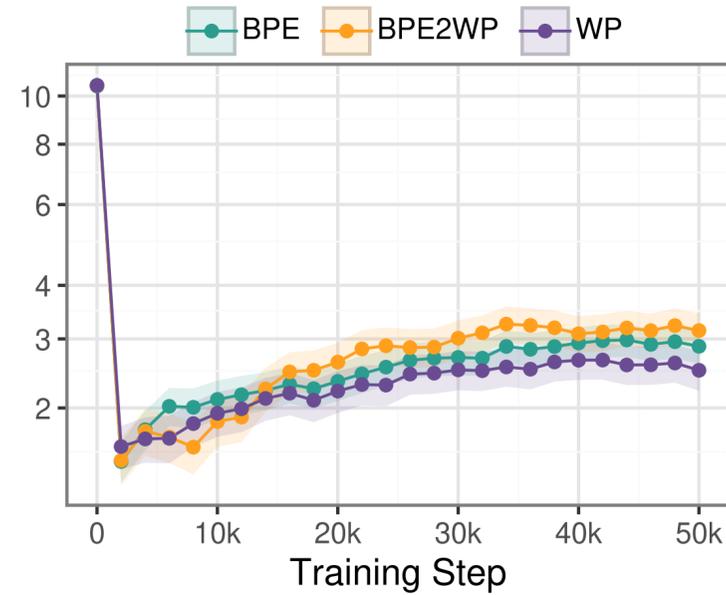
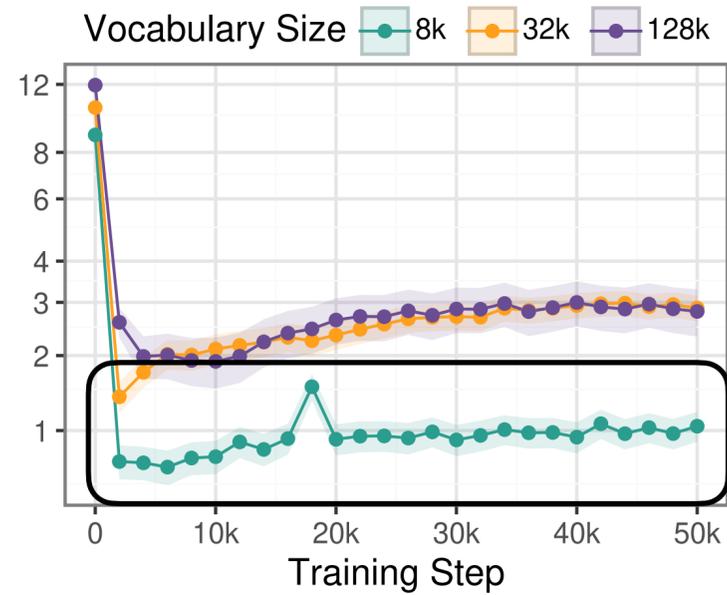
Tokenisation bias across training



- Drops sharply by step 2k, as expected
- Grows as models improve, contrary to expectations



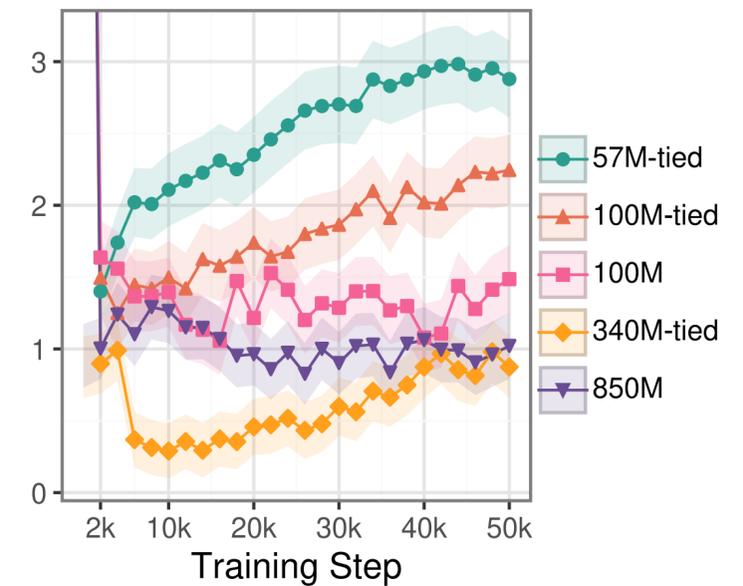
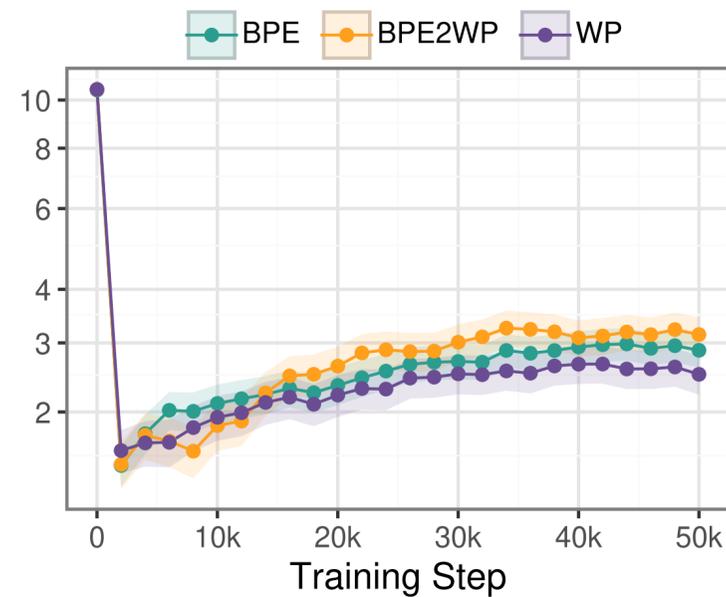
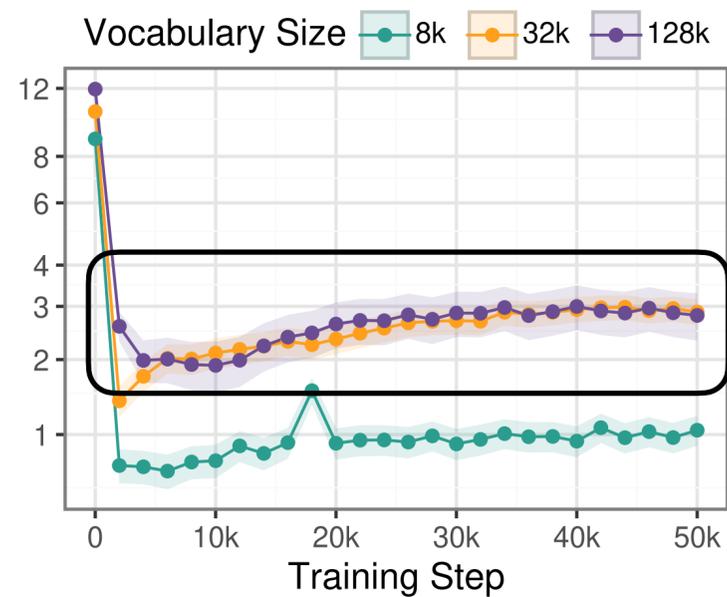
Tokenisation bias across training



- Drops sharply by step 2k, as expected
- Grows as models improve, contrary to expectations
- Weaker biases when using smaller vocabularies (perhaps due to more frequent subwords even beyond the cutoff)

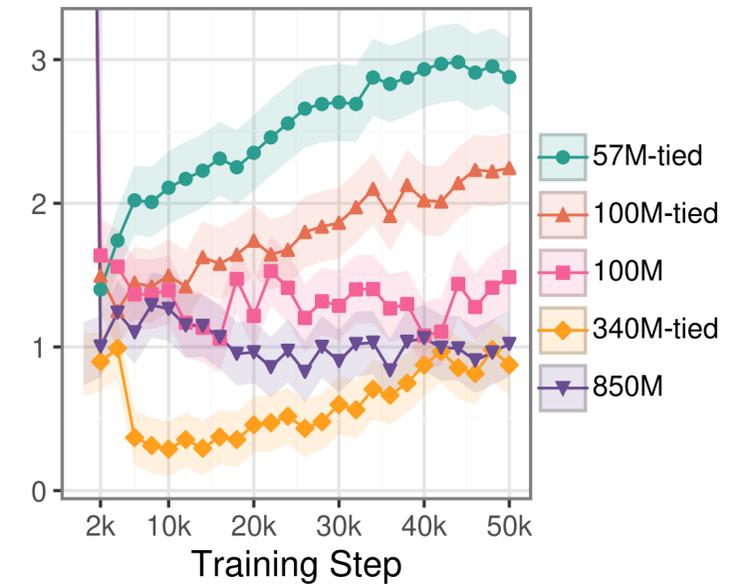
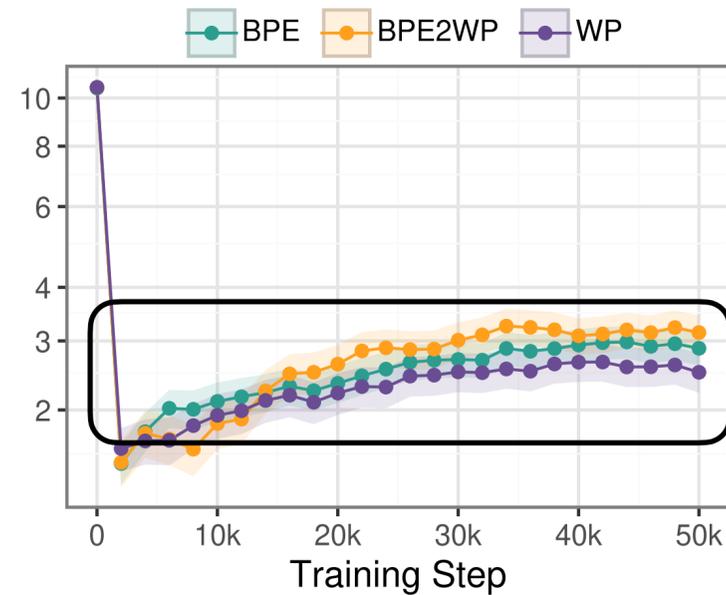
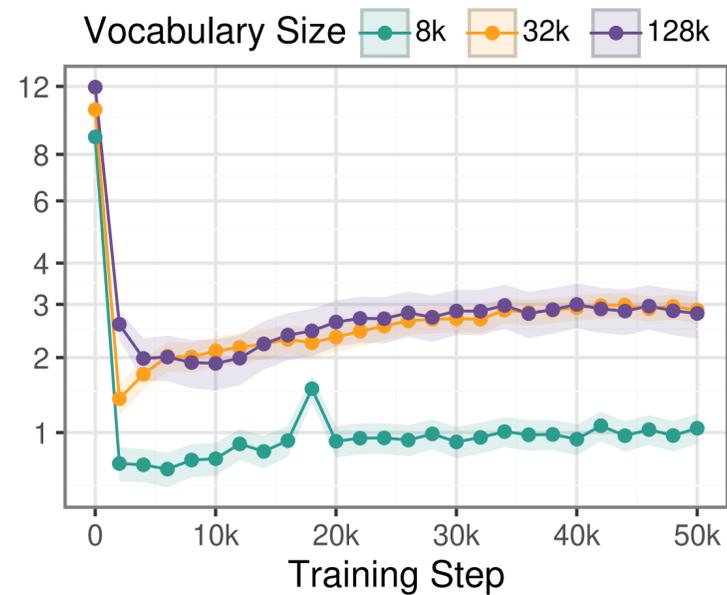


Tokenisation bias across training



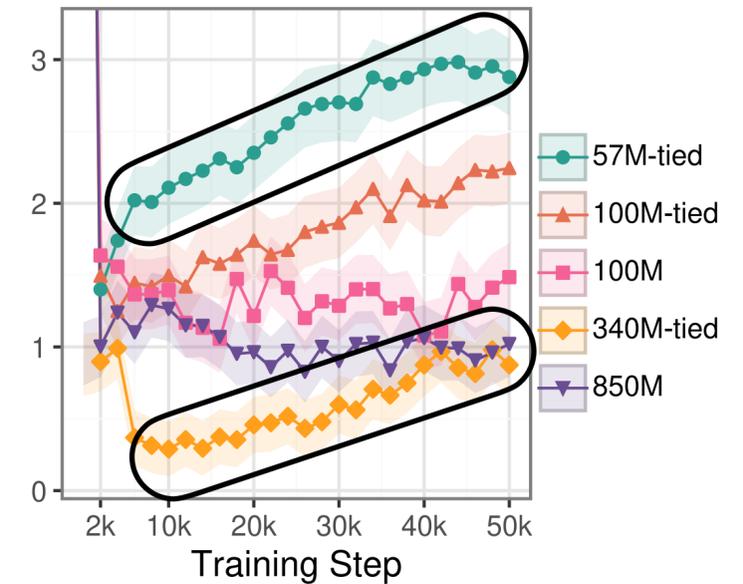
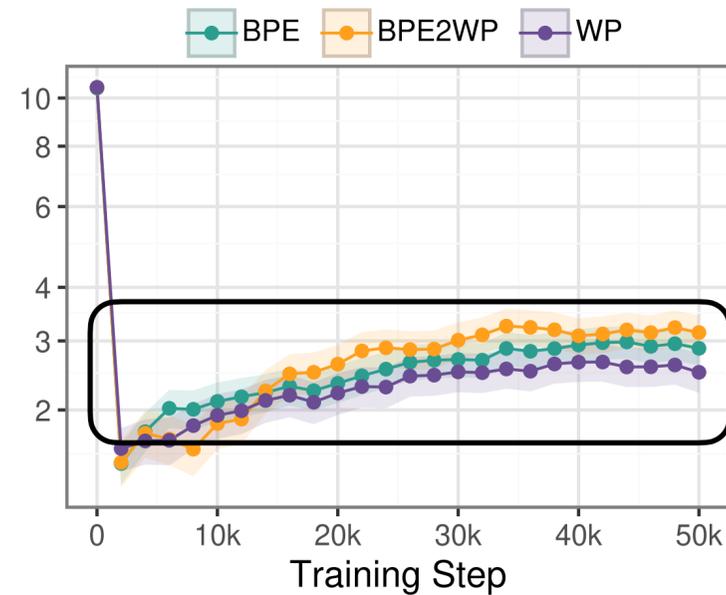
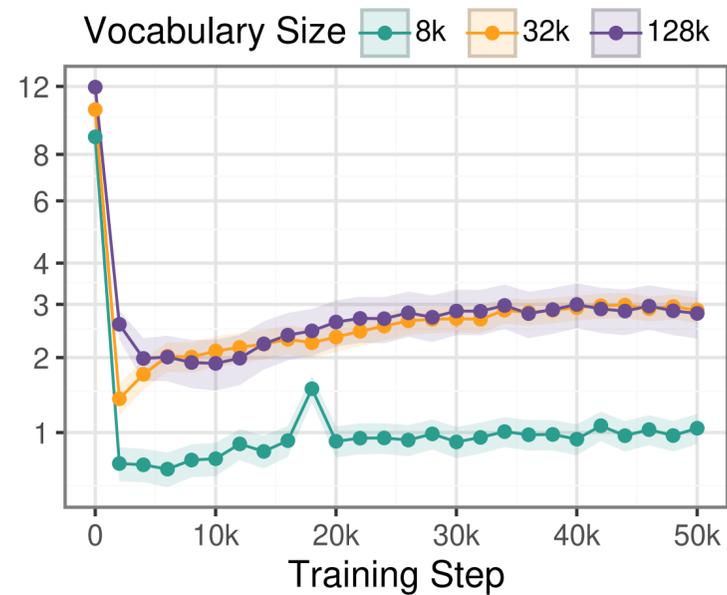
- Drops sharply by step 2k, as expected
- Grows as models improve, contrary to expectations
- Weaker biases when using smaller vocabularies (perhaps due to more frequent subwords even beyond the cutoff)
- Beyond a critical vocabulary size, the effect stabilises, with 128k and 32k showing similar biases despite a large difference

Tokenisation bias across training



- Drops sharply by step 2k, as expected
- Grows as models improve, contrary to expectations
- Weaker biases when using smaller vocabularies (perhaps due to more frequent subwords even beyond the cutoff)
- Beyond a critical vocabulary size, the effect stabilises, with 128k and 32k showing similar biases despite a large difference
- Consistent across these tokenisers

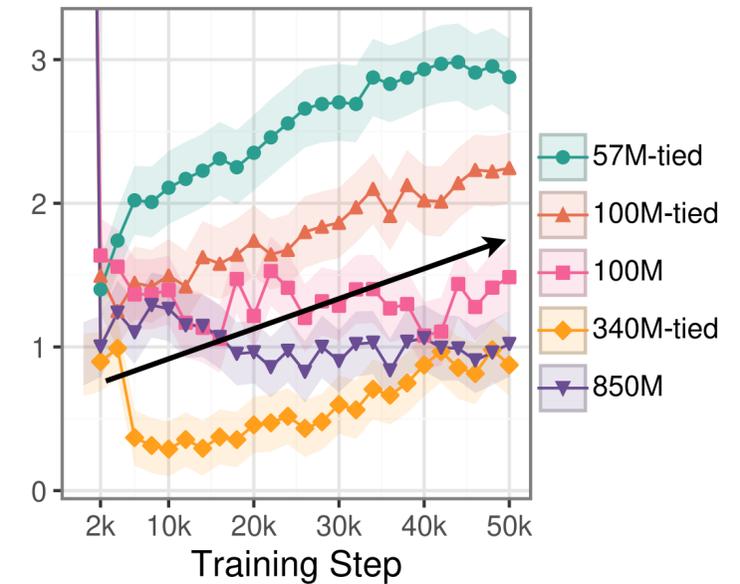
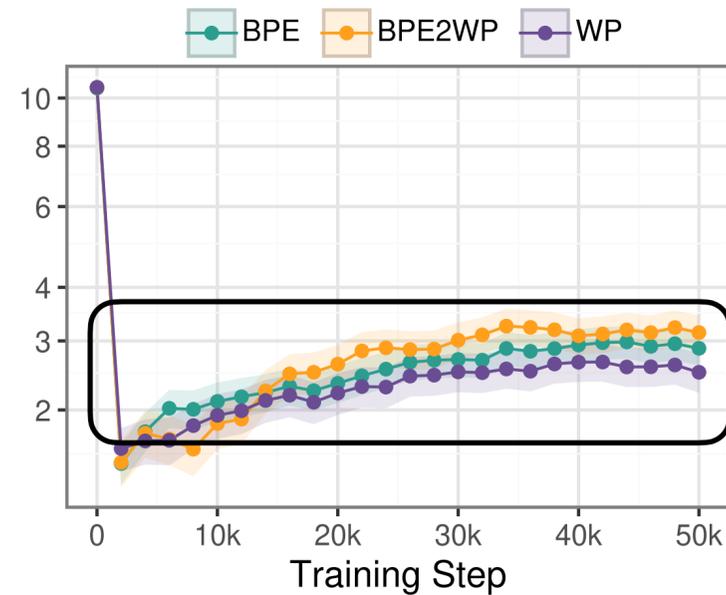
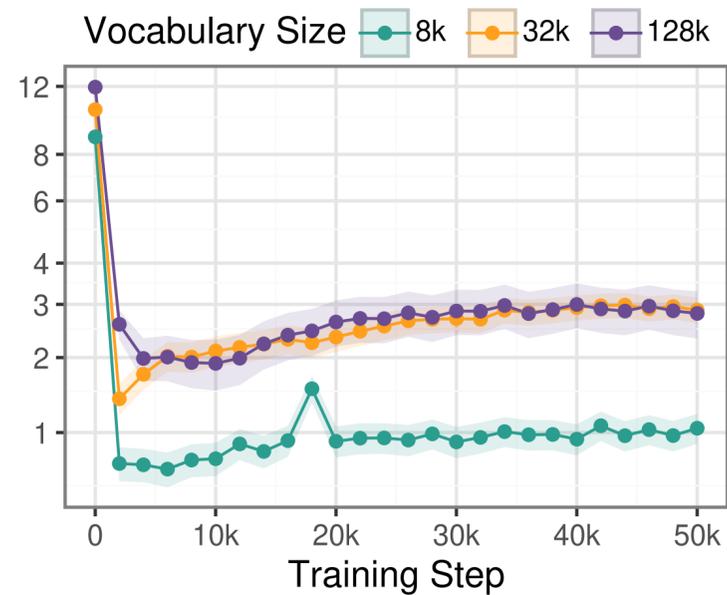
Tokenisation bias across training



- Drops sharply by step 2k, as expected
- Grows as models improve, contrary to expectations
- Weaker biases when using smaller vocabularies (perhaps due to more frequent subwords even beyond the cutoff)
- Beyond a critical vocabulary size, the effect stabilises, with 128k and 32k showing similar biases despite a large difference
- Consistent across these tokenisers

- Smaller bias in larger (e.g., 340M parameters) models than in small models (e.g., 57M)

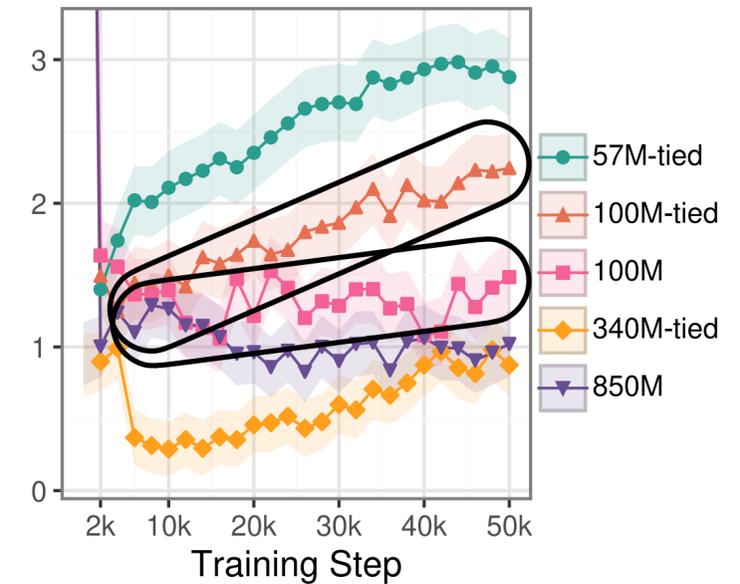
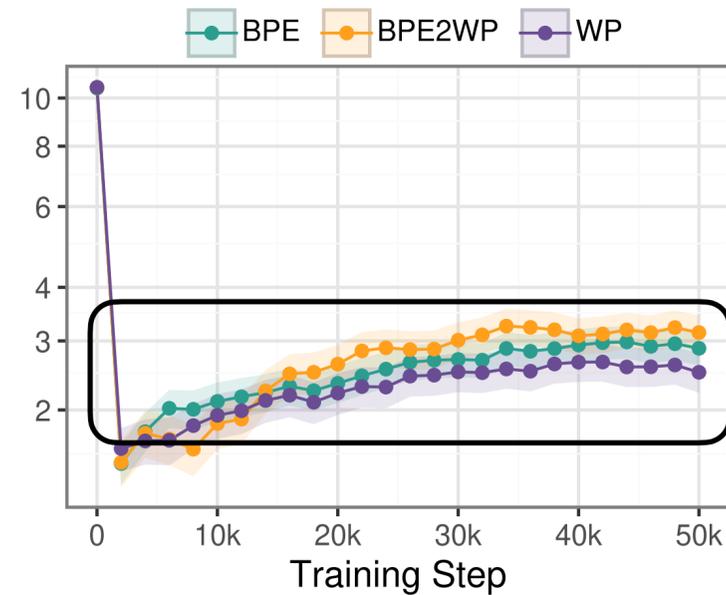
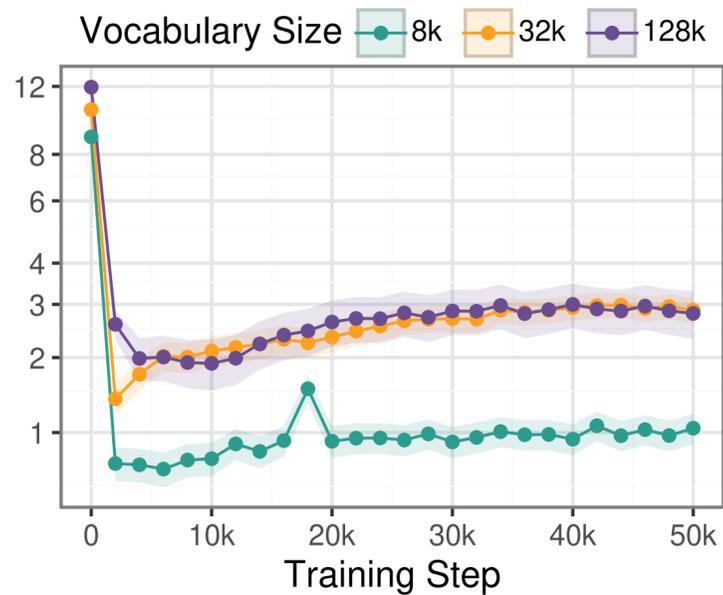
Tokenisation bias across training



- Drops sharply by step 2k, as expected
- Grows as models improve, contrary to expectations
- Weaker biases when using smaller vocabularies (perhaps due to more frequent subwords even beyond the cutoff)
- Beyond a critical vocabulary size, the effect stabilises, with 128k and 32k showing similar biases despite a large difference
- Consistent across these tokenisers

- Smaller bias in larger (e.g., 340M parameters) models than in small models (e.g., 57M)
- Across training, bias does not seem to decrease even for larger (better) models

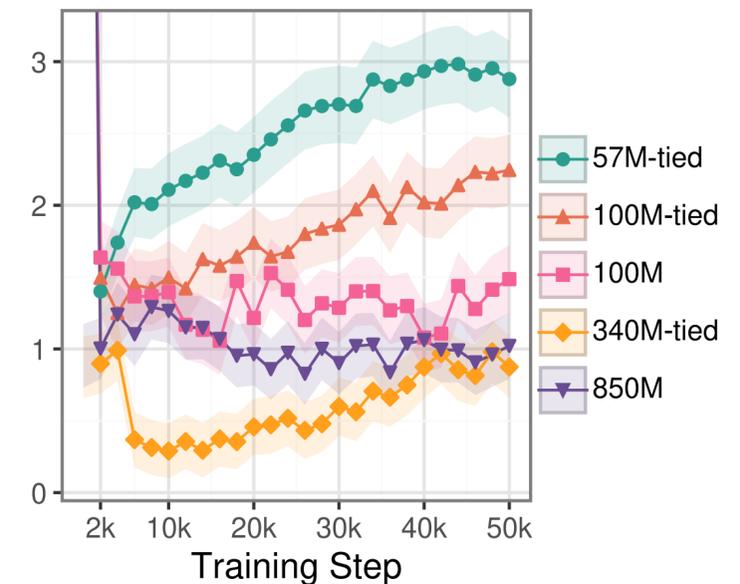
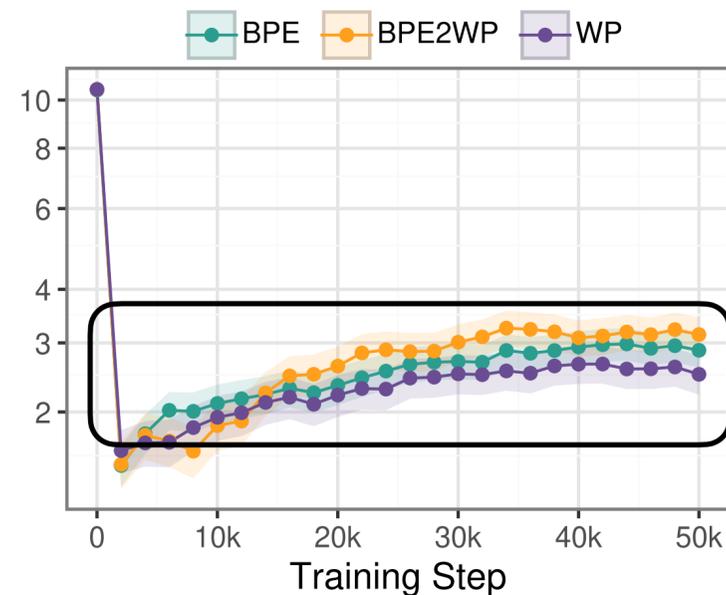
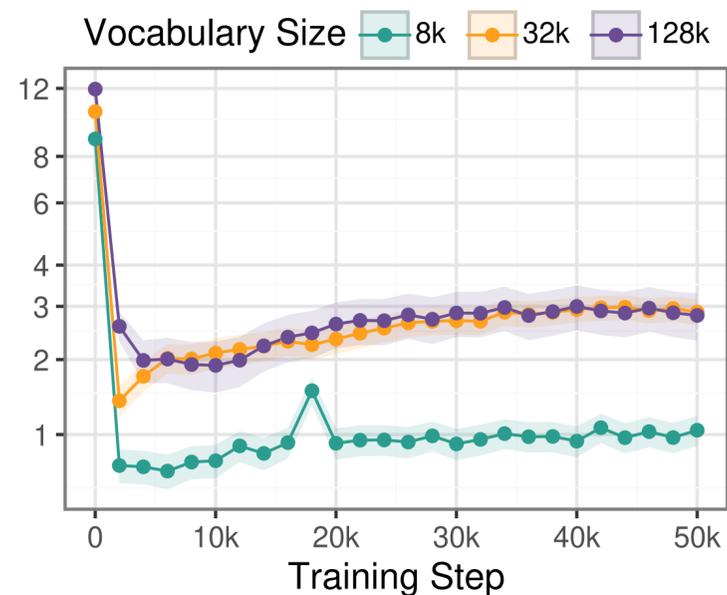
Tokenisation bias across training



- Drops sharply by step 2k, as expected
- Grows as models improve, contrary to expectations
- Weaker biases when using smaller vocabularies (perhaps due to more frequent subwords even beyond the cutoff)
- Beyond a critical vocabulary size, the effect stabilises, with 128k and 32k showing similar biases despite a large difference
- Consistent across these tokenisers

- Smaller bias in larger (e.g., 340M parameters) models than in small models (e.g., 57M)
- Across training, bias does not seem to decrease even for larger (better) models
- Embedding tying increases tokenisation bias (see 100M) vs. untying

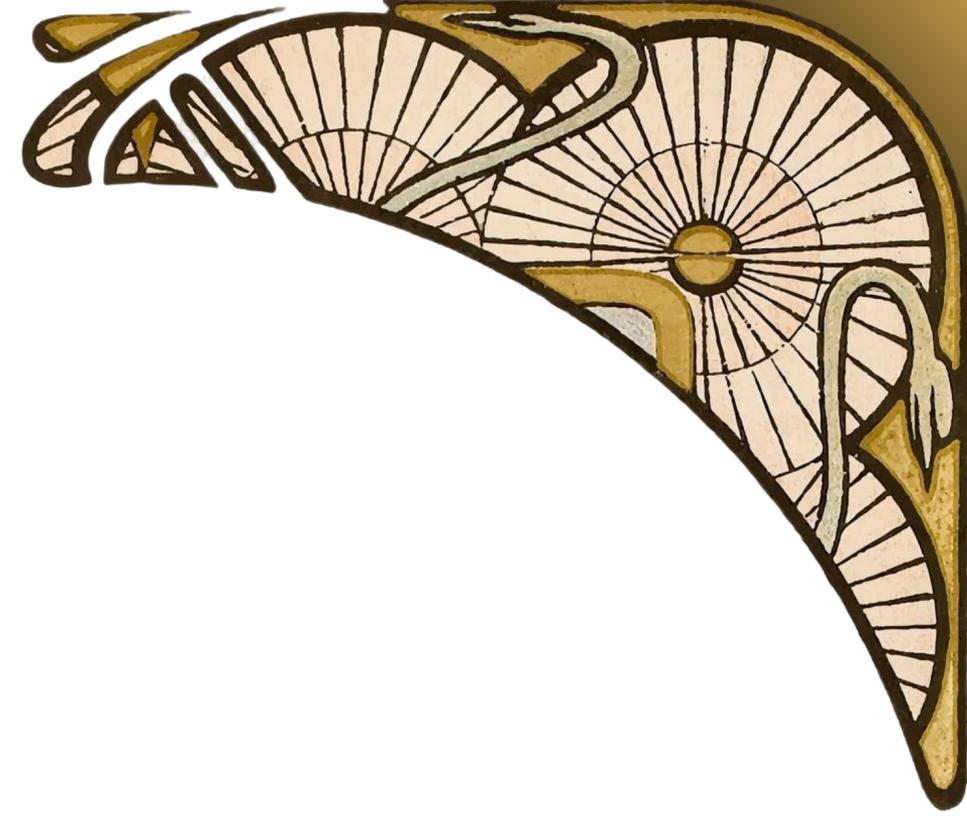
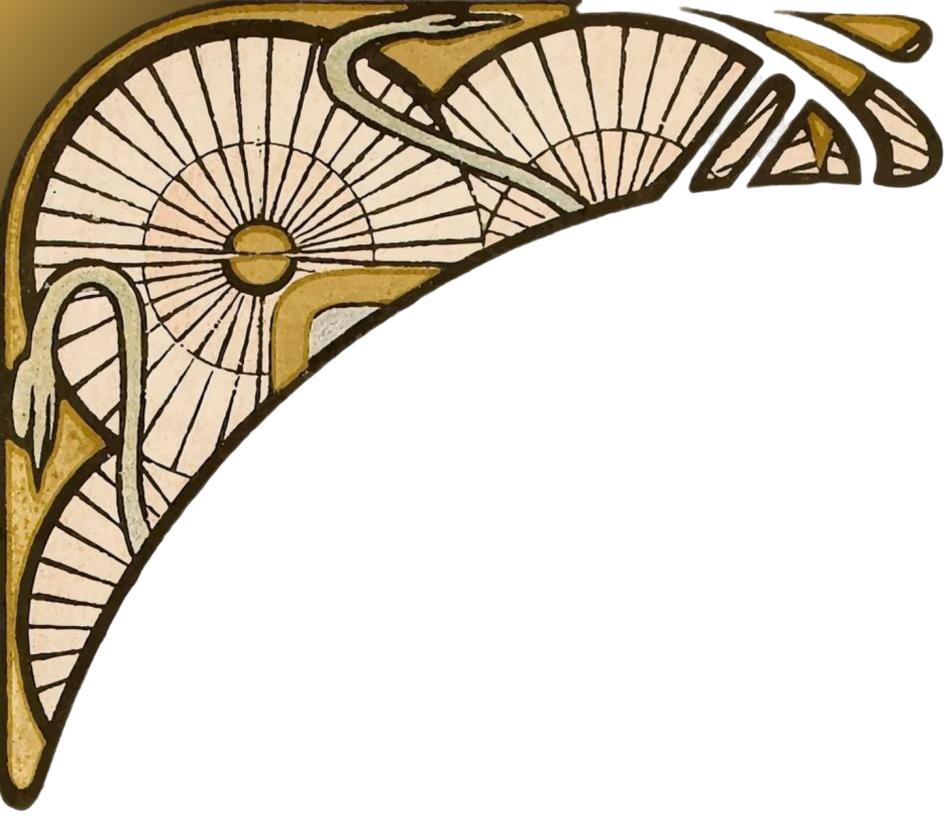
Tokenisation bias across training



- Drops sharply by step 2k, as expected
- Grows as models improve, contrary to expectations
- Weaker biases when using smaller vocabularies (perhaps due to more frequent subwords even beyond the cutoff)
- Beyond a critical vocabulary size, the effect stabilises, with 128k and 32k showing similar biases despite a large difference
- Consistent across these tokenisers

- Smaller bias in larger (e.g., 340M parameters) models than in small models (e.g., 57M)
- Across training, bias does not seem to decrease even for larger (better) models
- Embedding tying increases tokenisation bias (see 100M) vs. untying

Together, these results suggest that tokenisation bias may be a persisting property of LMs



FINE



Pietro Lesci
University of Cambridge



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