

# DAGER: Exact Gradient Inversion for Large Language Models



Ivo Petrov



Dimitar I. Dimitrov



Maximilian Baader



Mark Müller



Martin Vechev

**INSAIT**

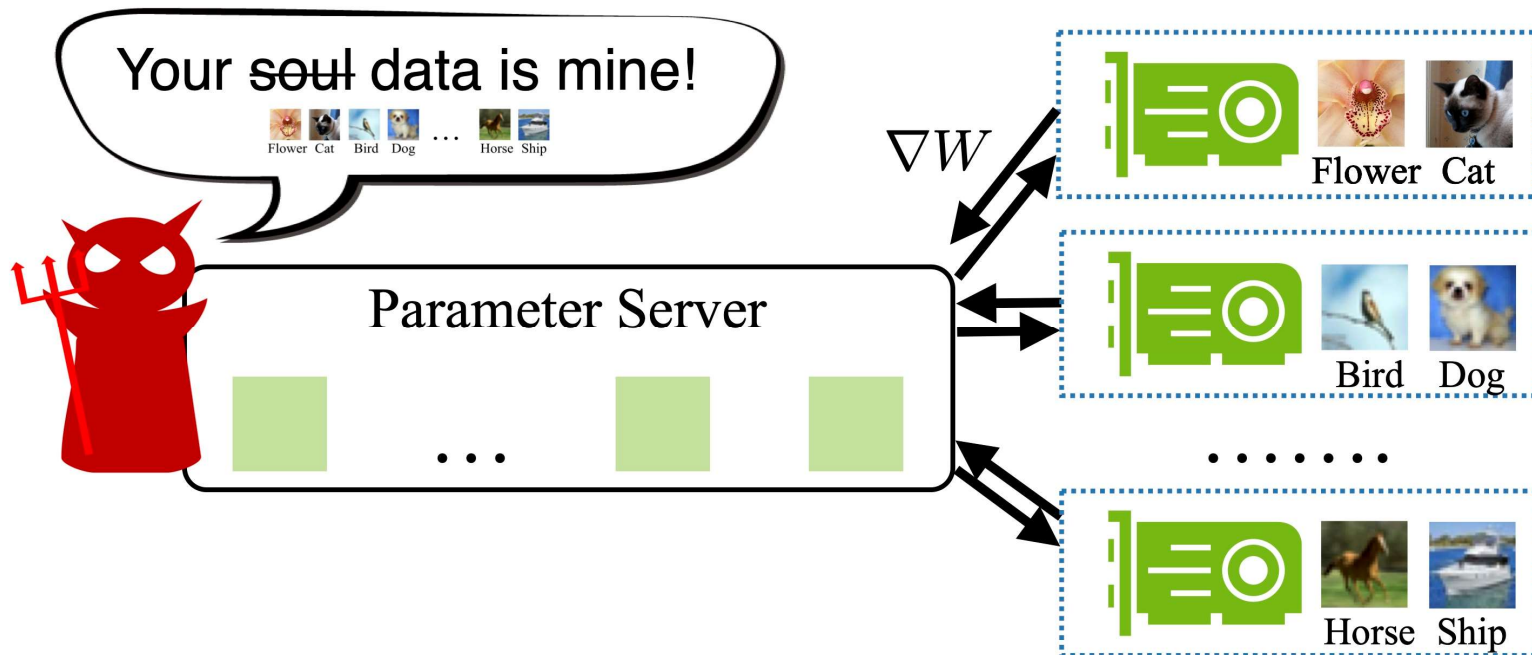
**ETH** zürich

**SRILAB**

NeurIPS 2024 page:



# Gradient Inversion



Zhu et al., "Deep leakage from gradients." *Advances in neural information processing systems* 32 (2019).

# LLMs (and transformers) in Federated Learning

LLMs are becoming increasingly popular (and powerful!)

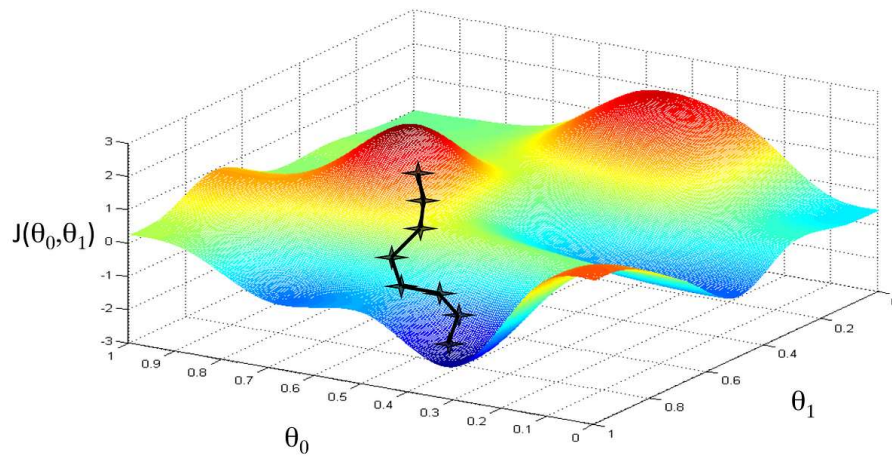
Uses of LLMs in FL include:

- Keyboard suggestions
- Fraud detection
- Healthcare diagnostics
- Legal Document Analysis



# Gradient Inversion – Prior Work

State-of-the-art attacks utilize a continuous optimization approach

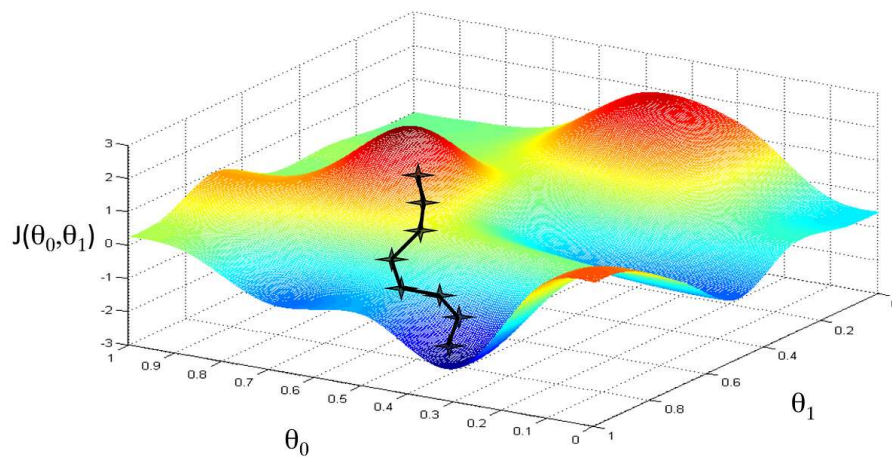


$$\mathbf{X}, y = \operatorname{argmin}_{\mathbf{X}, y} \left\| \frac{\partial \mathcal{L}_\theta(\mathbf{X}, y)}{\partial \mathbf{W}} - \nabla \mathbf{W} \right\|^2$$

Jaleel Adejumo, Gradient Descent From Scratch- Batch Gradient Descent, *Medium*

# Gradient Inversion – Prior Work

State-of-the-art attacks utilize a continuous optimization approach

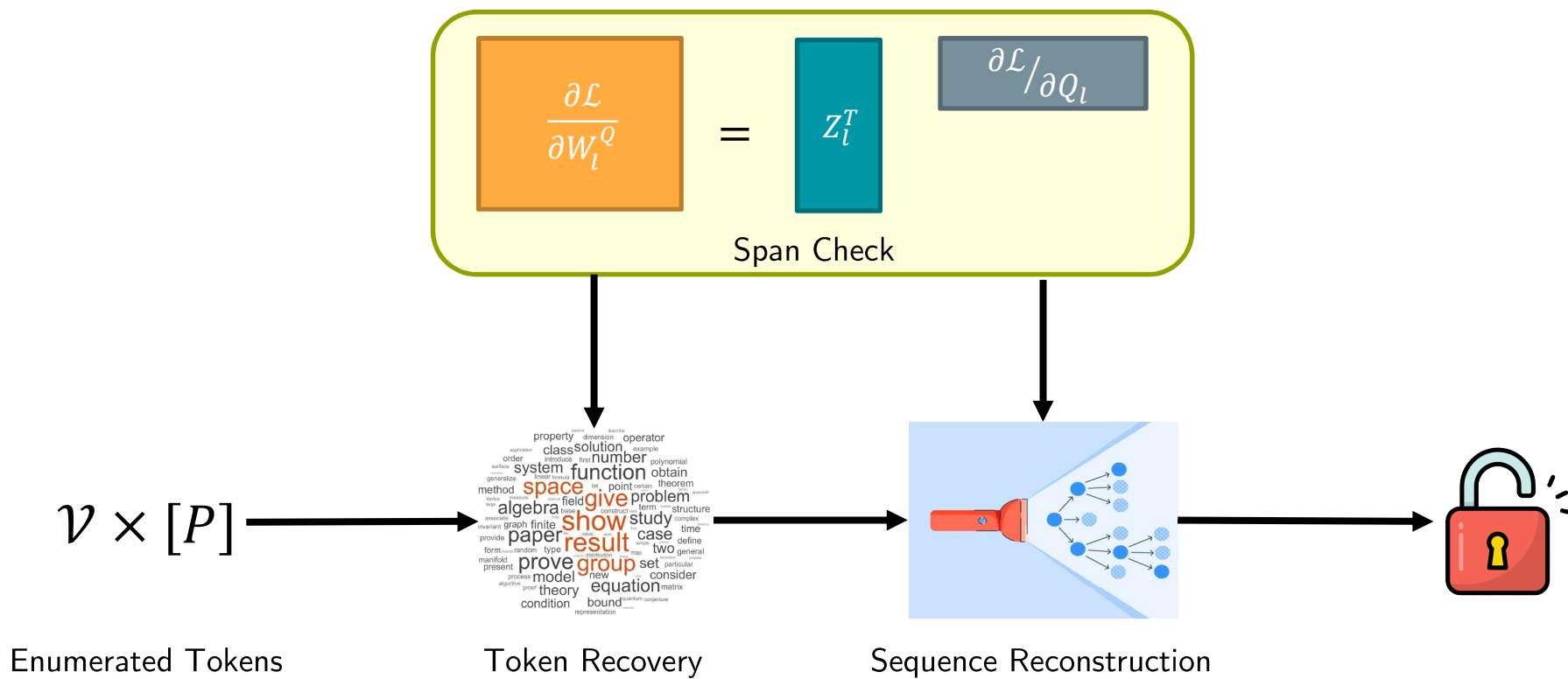


*Lorem ipsum dolor  
sit amet, consectetur  
adipiscing elit. Donec  
non lectus lobortis  
purus egestas  
rhoncus quis eu erat.*

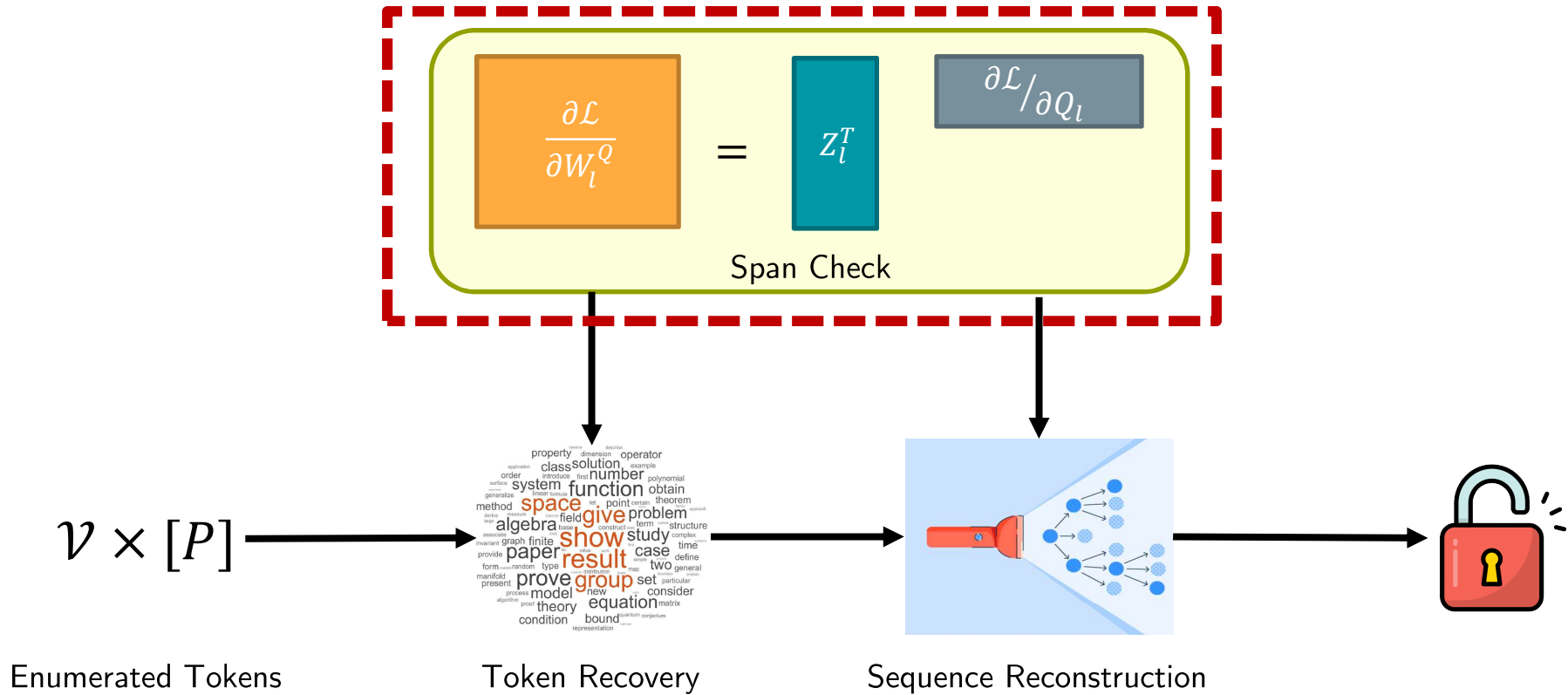


$$\mathbf{X}, y = \operatorname{argmin}_{\mathbf{X}, y} \left\| \frac{\partial \mathcal{L}_{\theta}(\mathbf{X}, y)}{\partial \mathbf{W}} - \nabla \mathbf{W} \right\|^2$$

# The DAGER Pipeline



# The DAGER Pipeline

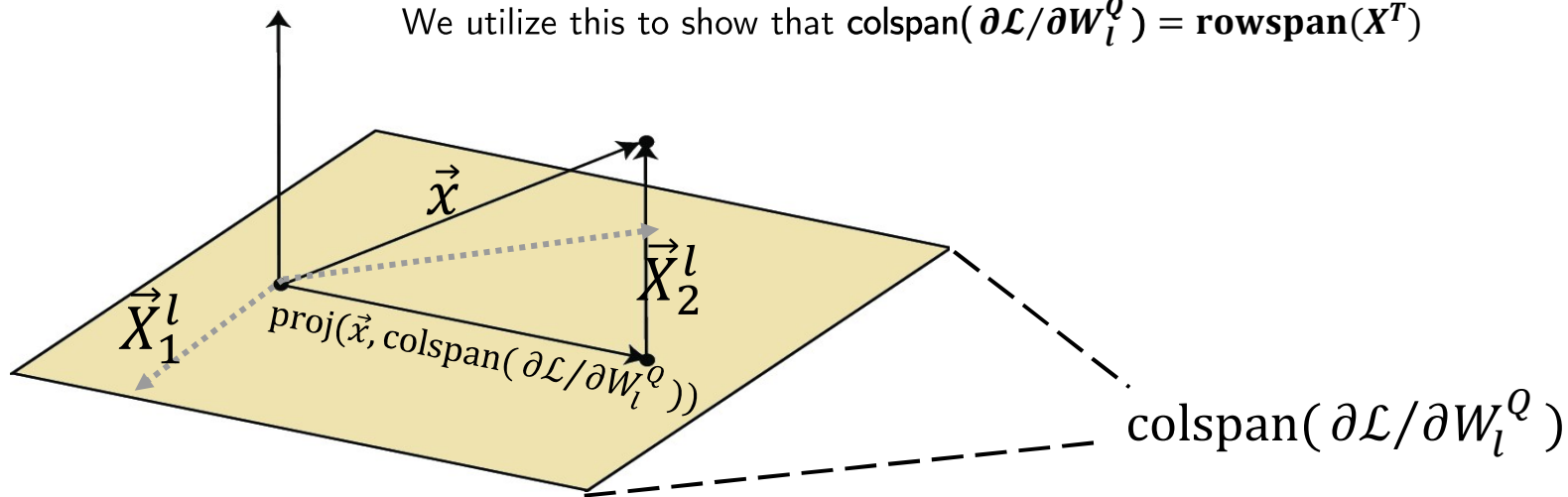


# The Span Check Filter

Dimitrov et al. [1] showed that for a linear layer  $Y = XW$  that:

$$\frac{\partial \mathcal{L}}{\partial W} = X^T \frac{\partial \mathcal{L}}{\partial Y}$$

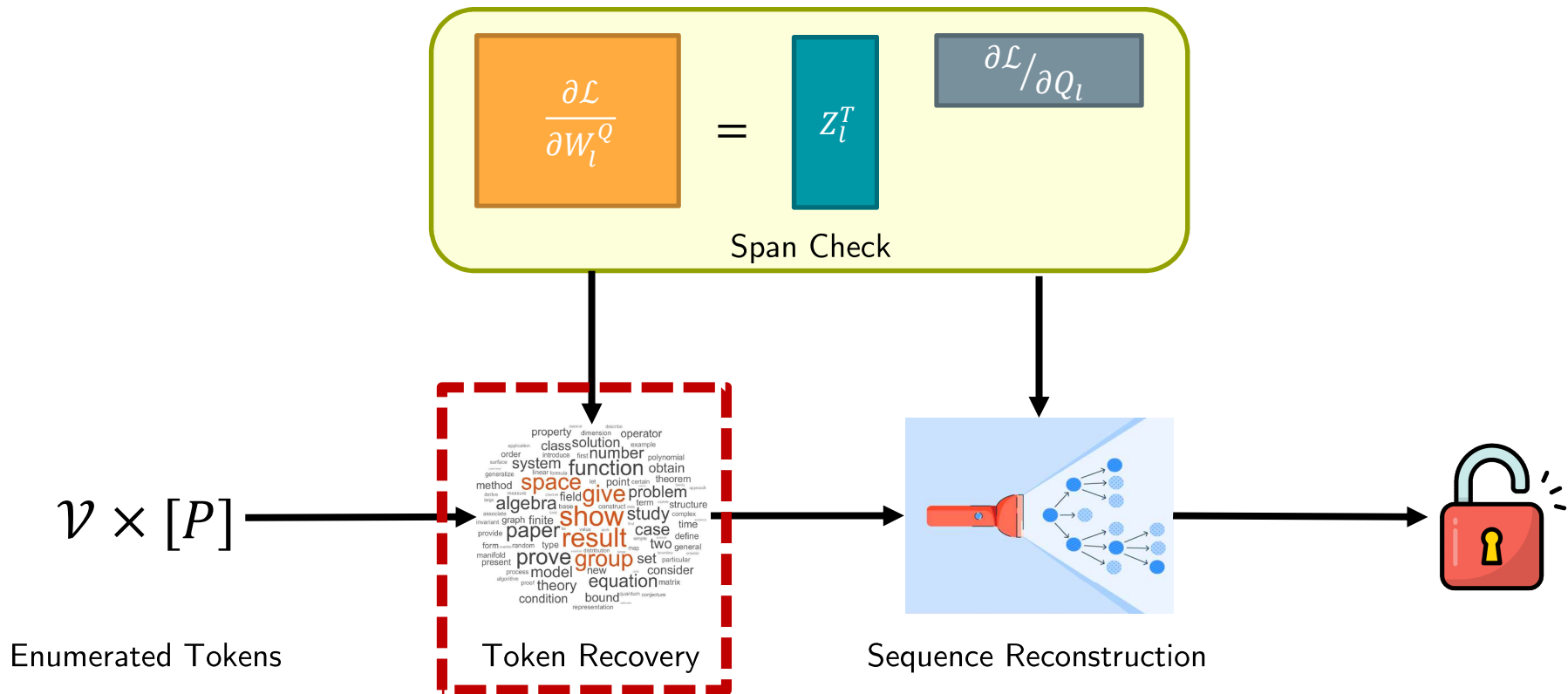
We utilize this to show that  $\text{colspan}(\partial \mathcal{L} / \partial W_l^q) = \text{rowspan}(X^T)$



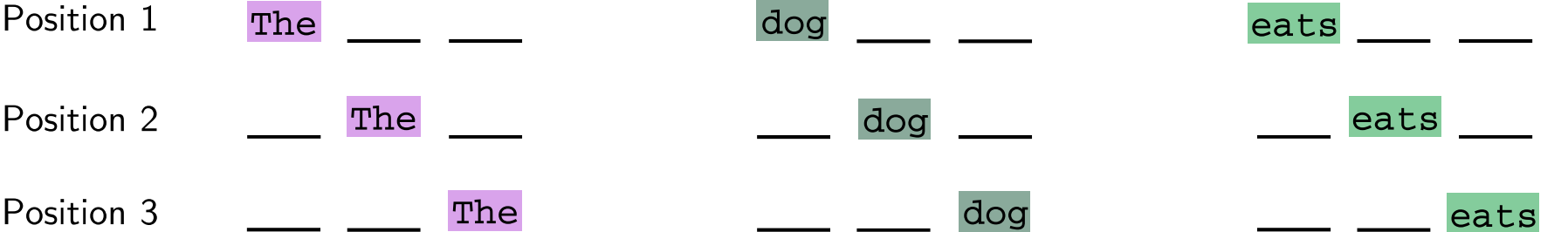
[1] Dimitrov, Dimitar I., et al. "Spear: Exact gradient inversion of batches in federated learning."



# The DAGER Pipeline



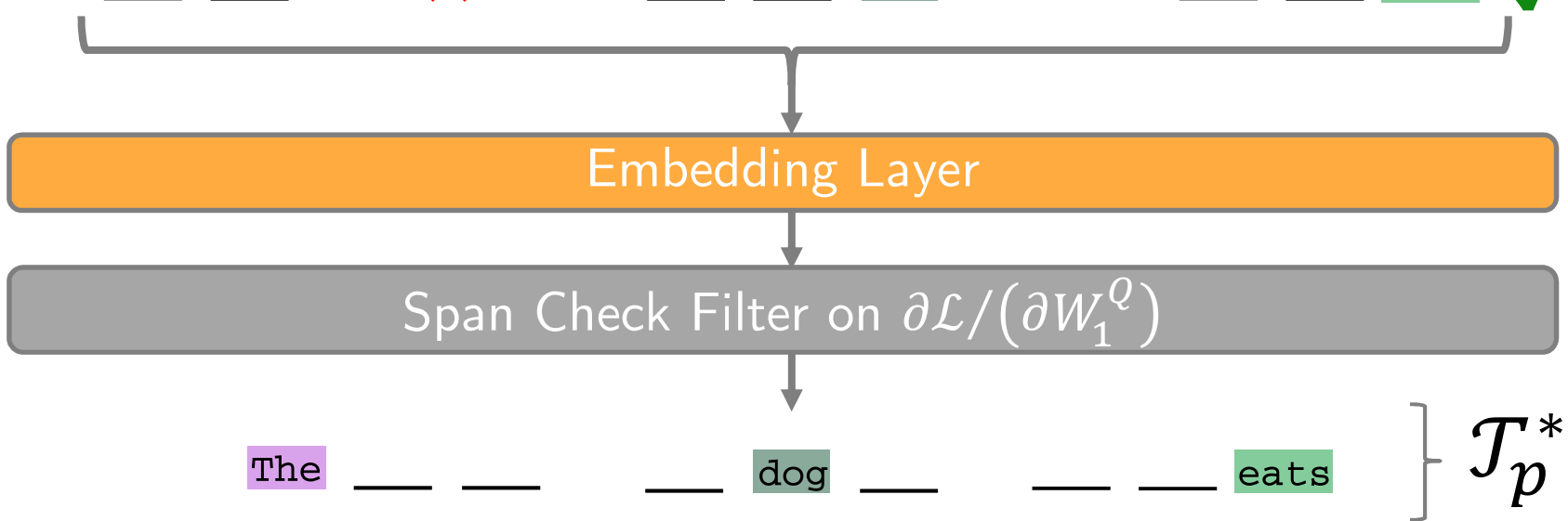
# Token Recovery



Embedding Layer

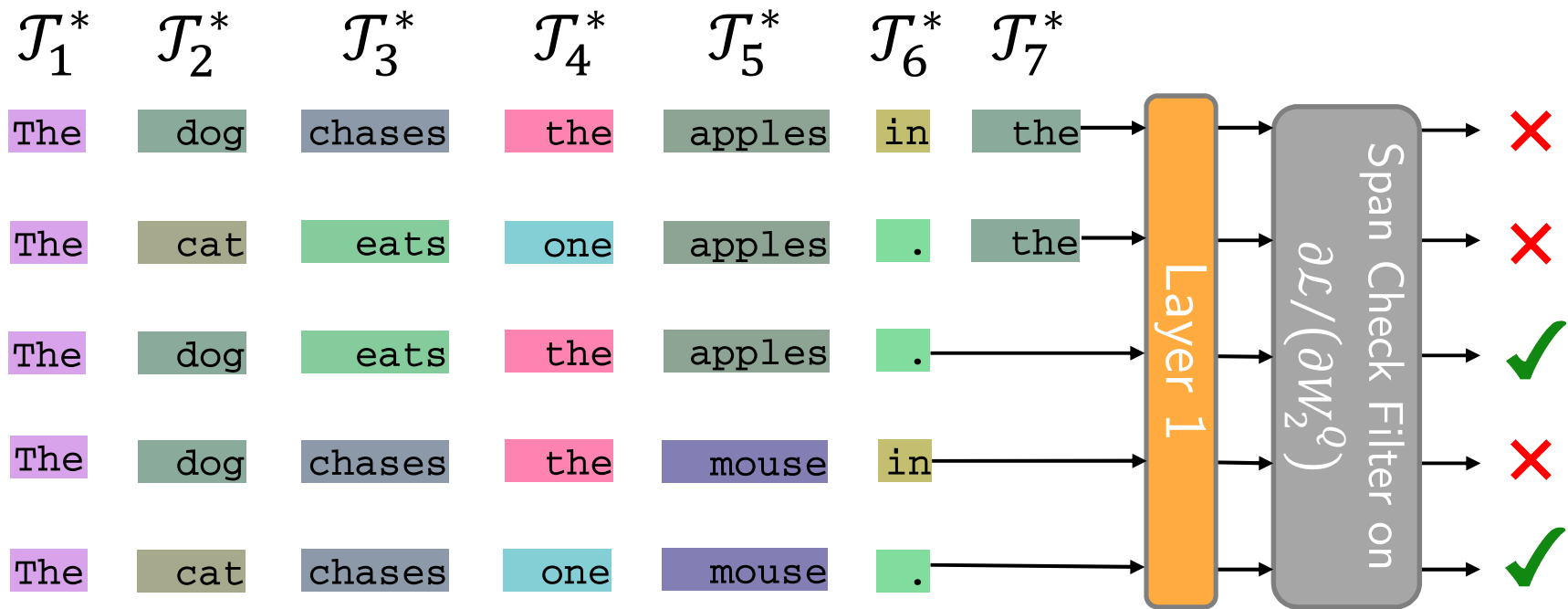
# Token Recovery

Position 1	The	___	___	✓	dog	___	___	✗	eats	___	___	✗
Position 2	___	The	___	✗	___	dog	___	✓	___	eats	___	✗
Position 3	___	___	The	✗	___	___	dog	✗	___	___	eats	✓



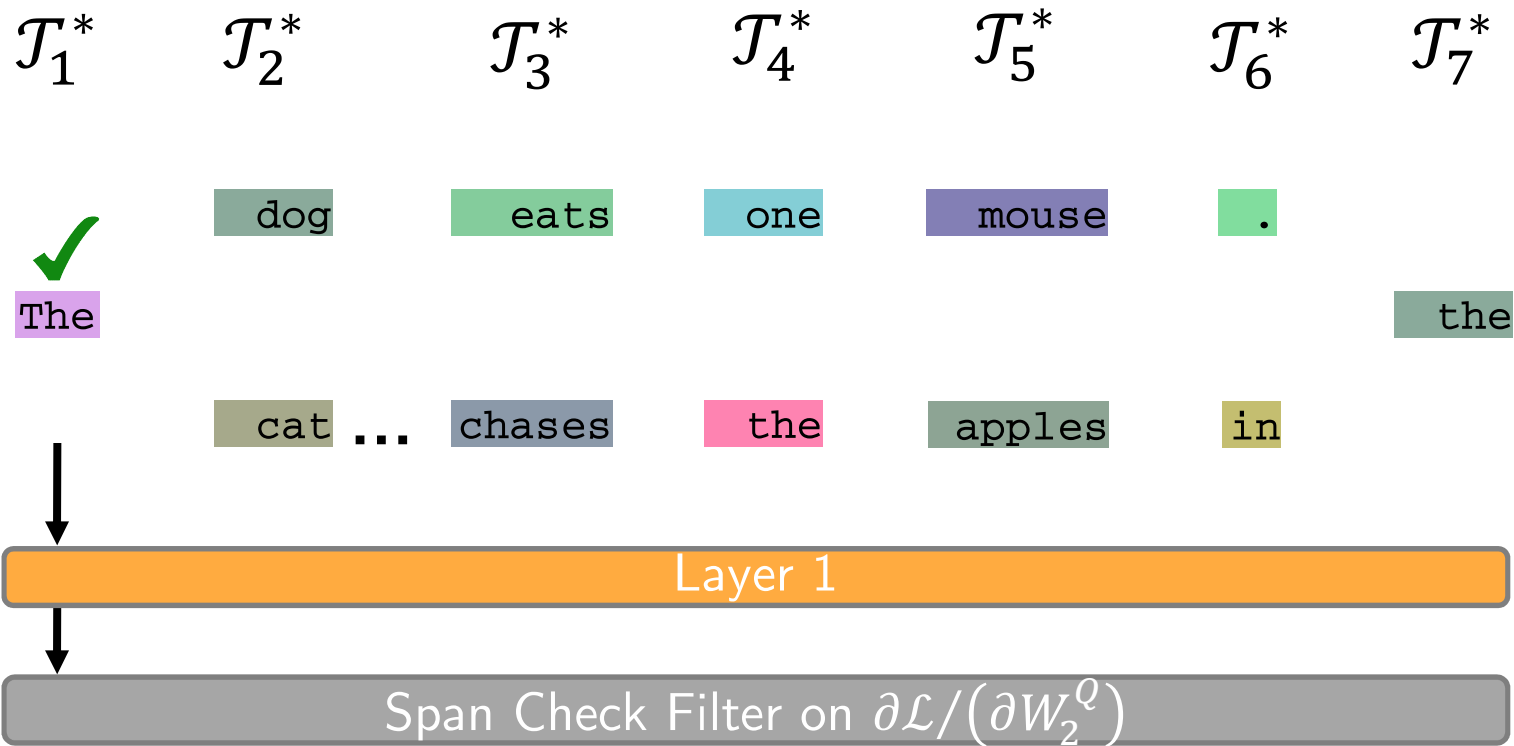


# Sequence Reconstruction – Encoder-only models



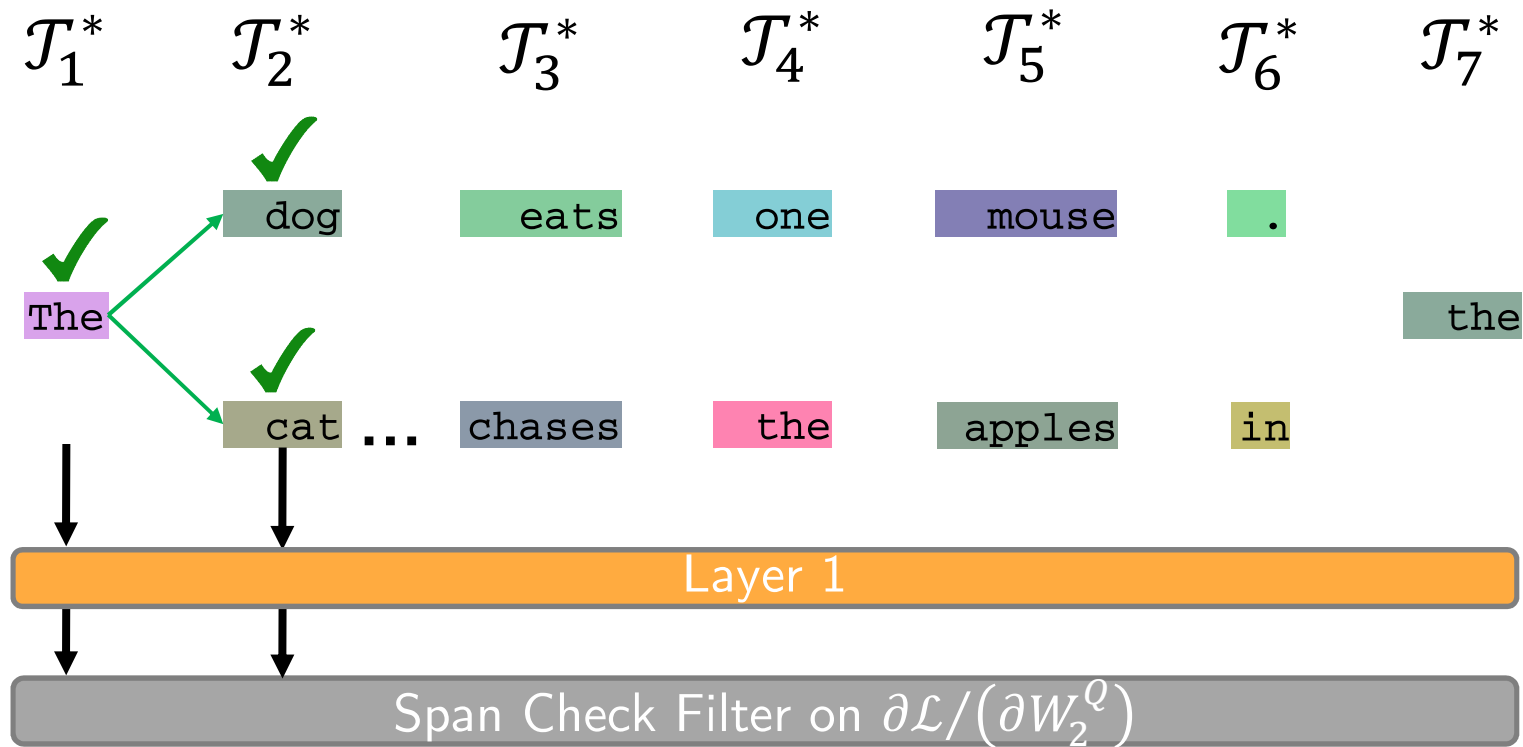
Exhausting all token combinations is possible for encoder-only models, but is computationally expensive.

# Sequence Reconstruction – Decoder-only models

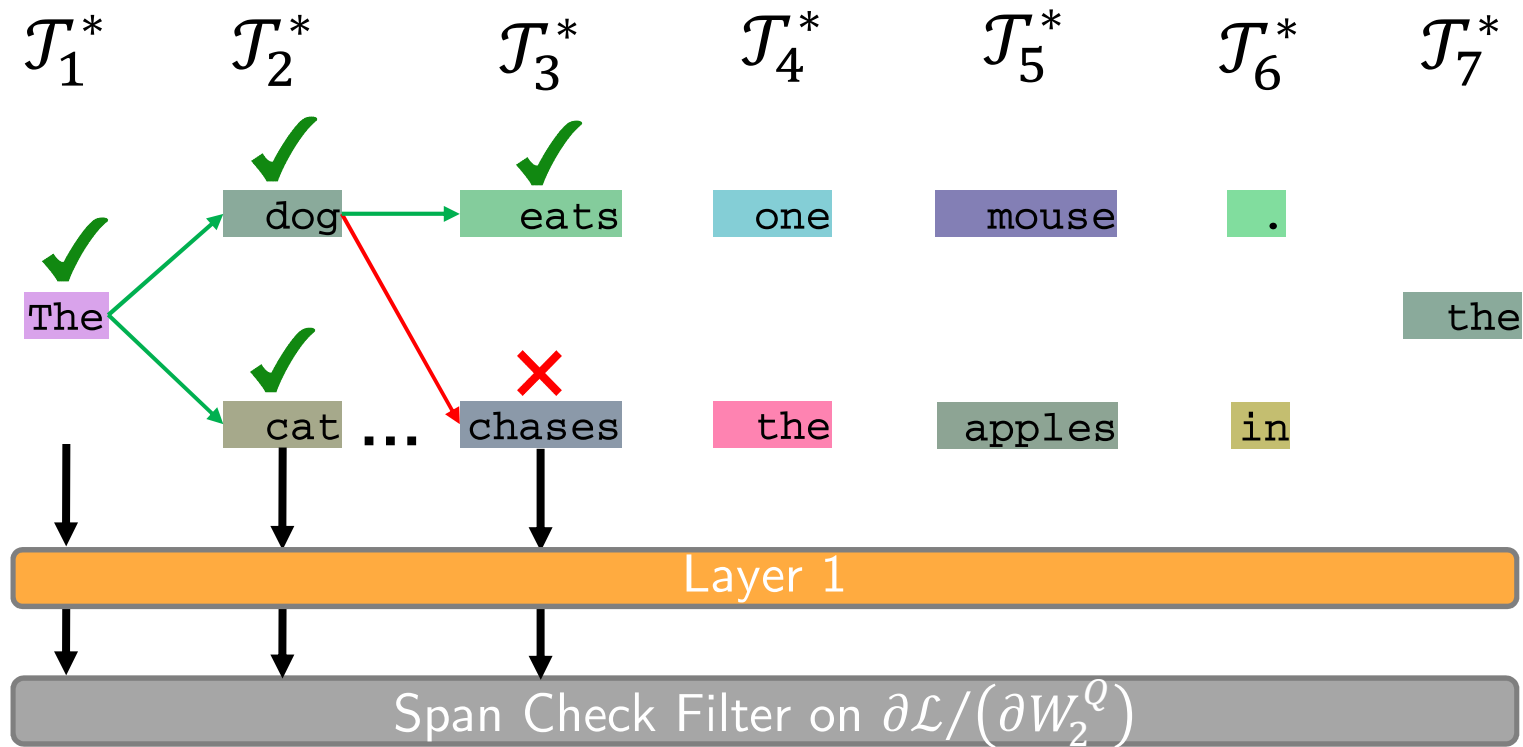


**Key Observation:** The outputs at position  $t$  only depend on inputs of tokens up to  $t-1$

# Sequence Reconstruction – Decoder-only models

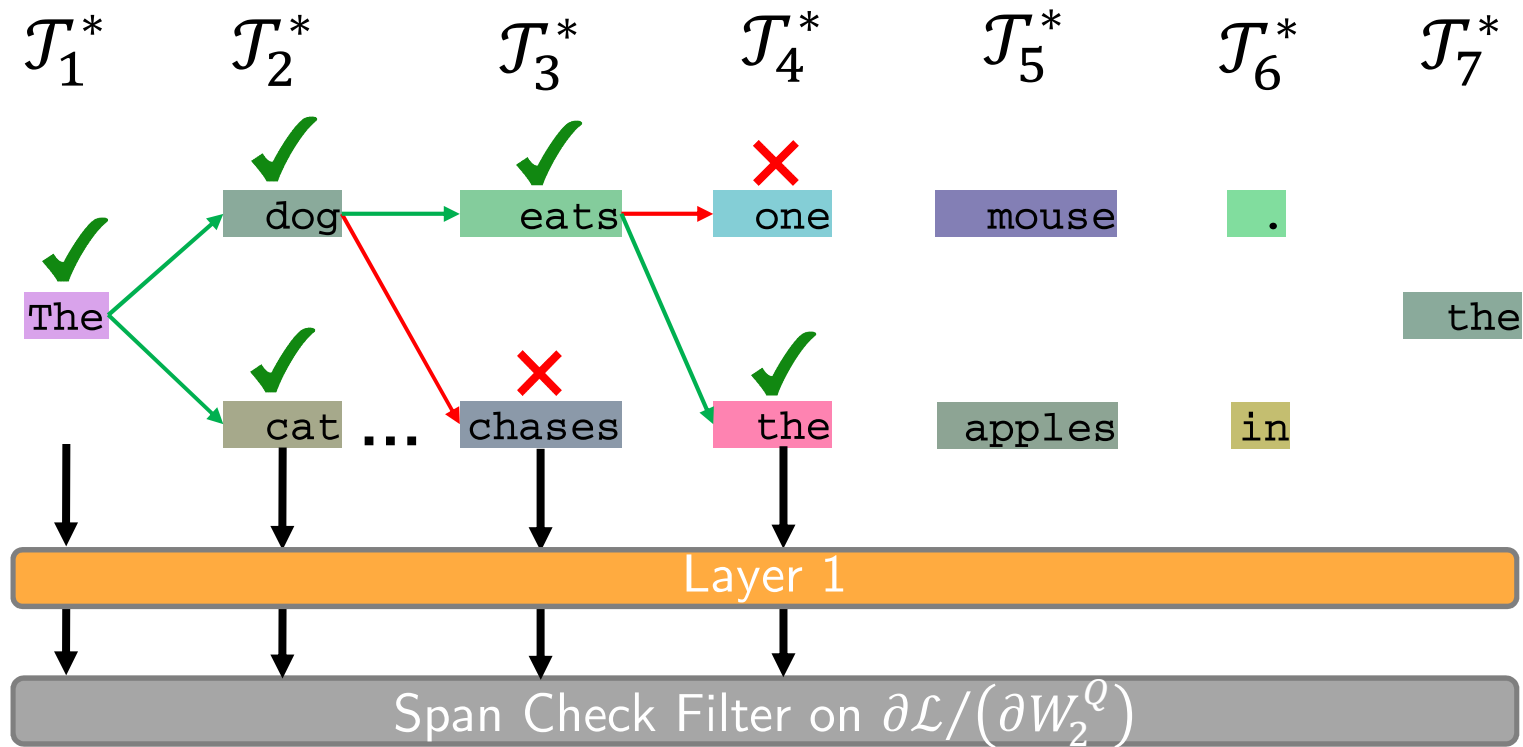


# Sequence Reconstruction – Decoder-only models

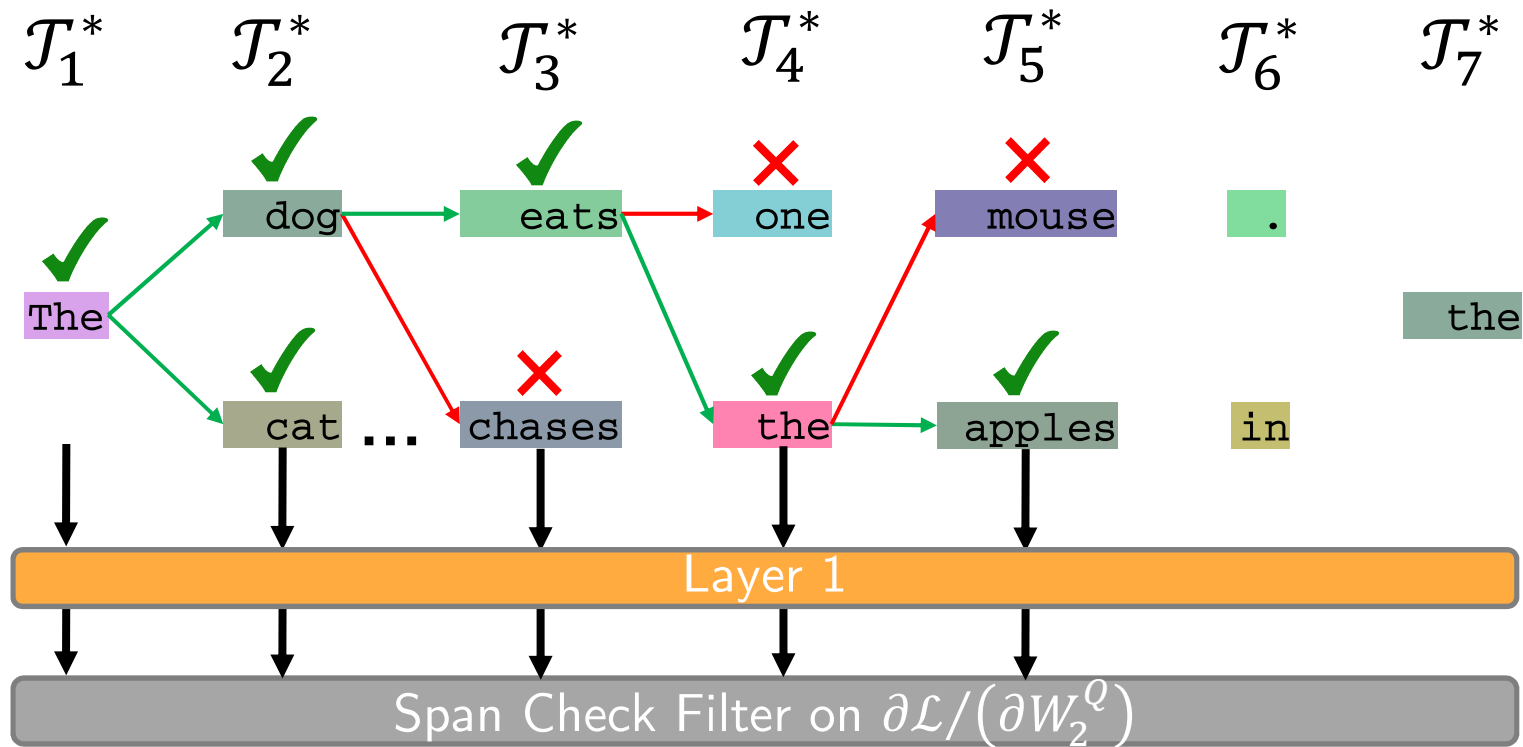




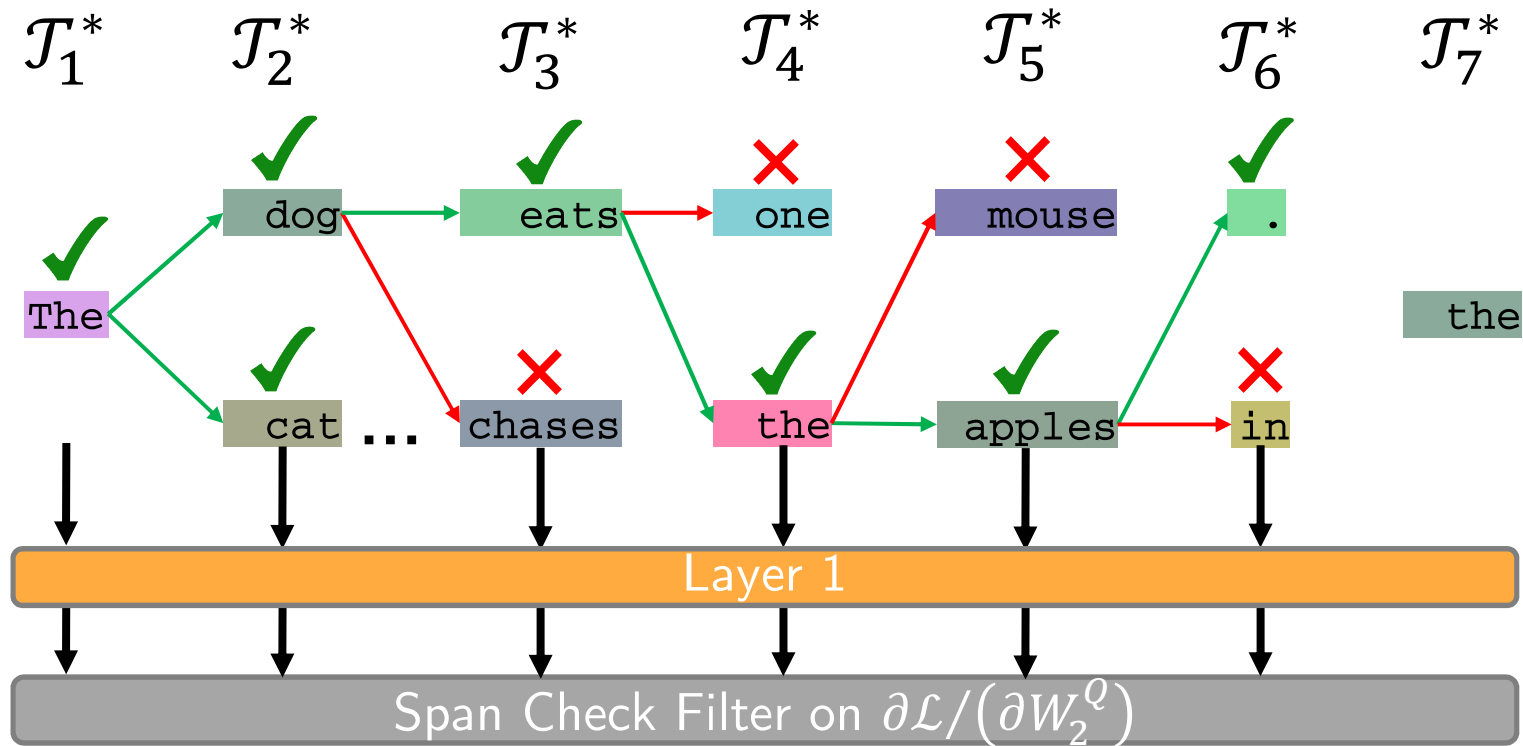
# Sequence Reconstruction – Decoder-only models



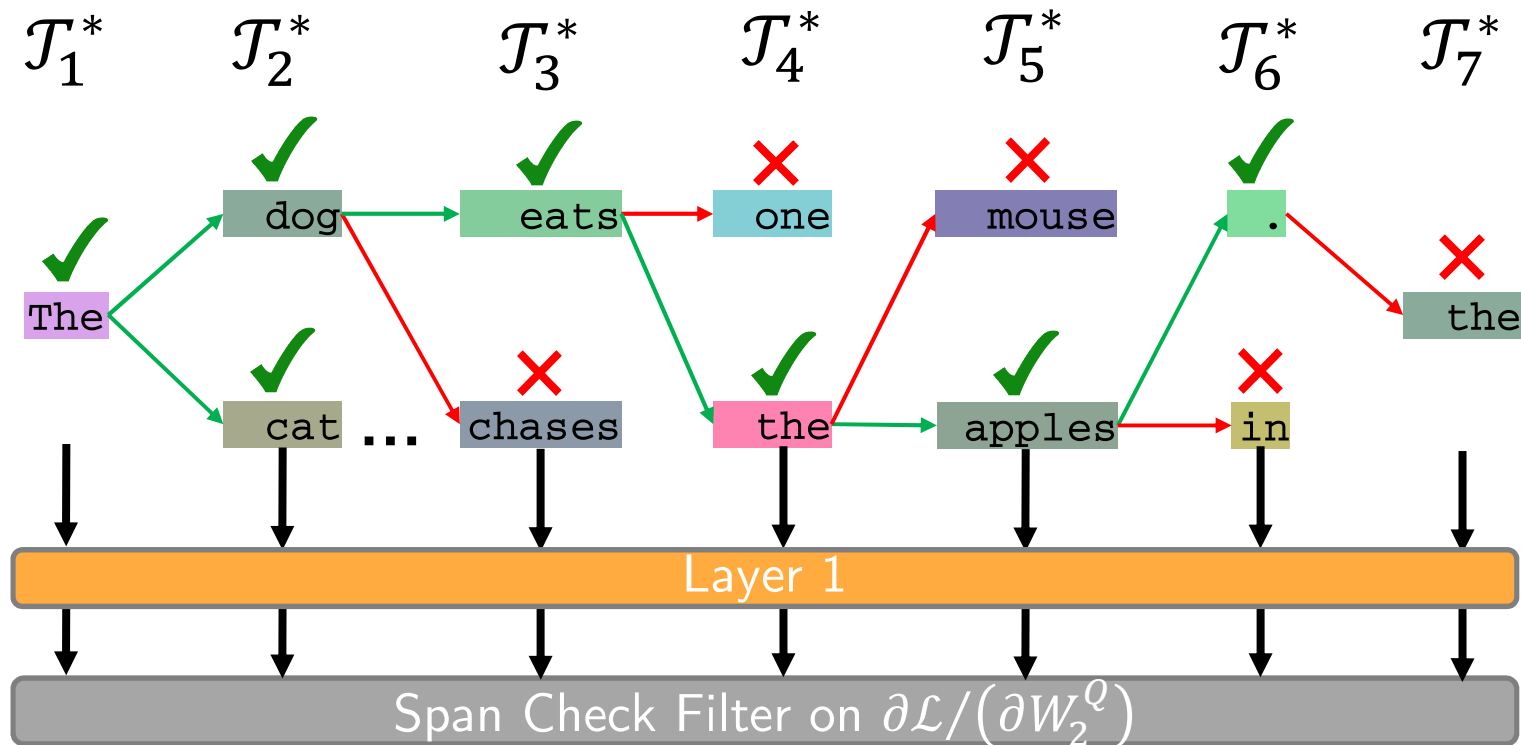
# Sequence Reconstruction – Decoder-only models



# Sequence Reconstruction – Decoder-only models

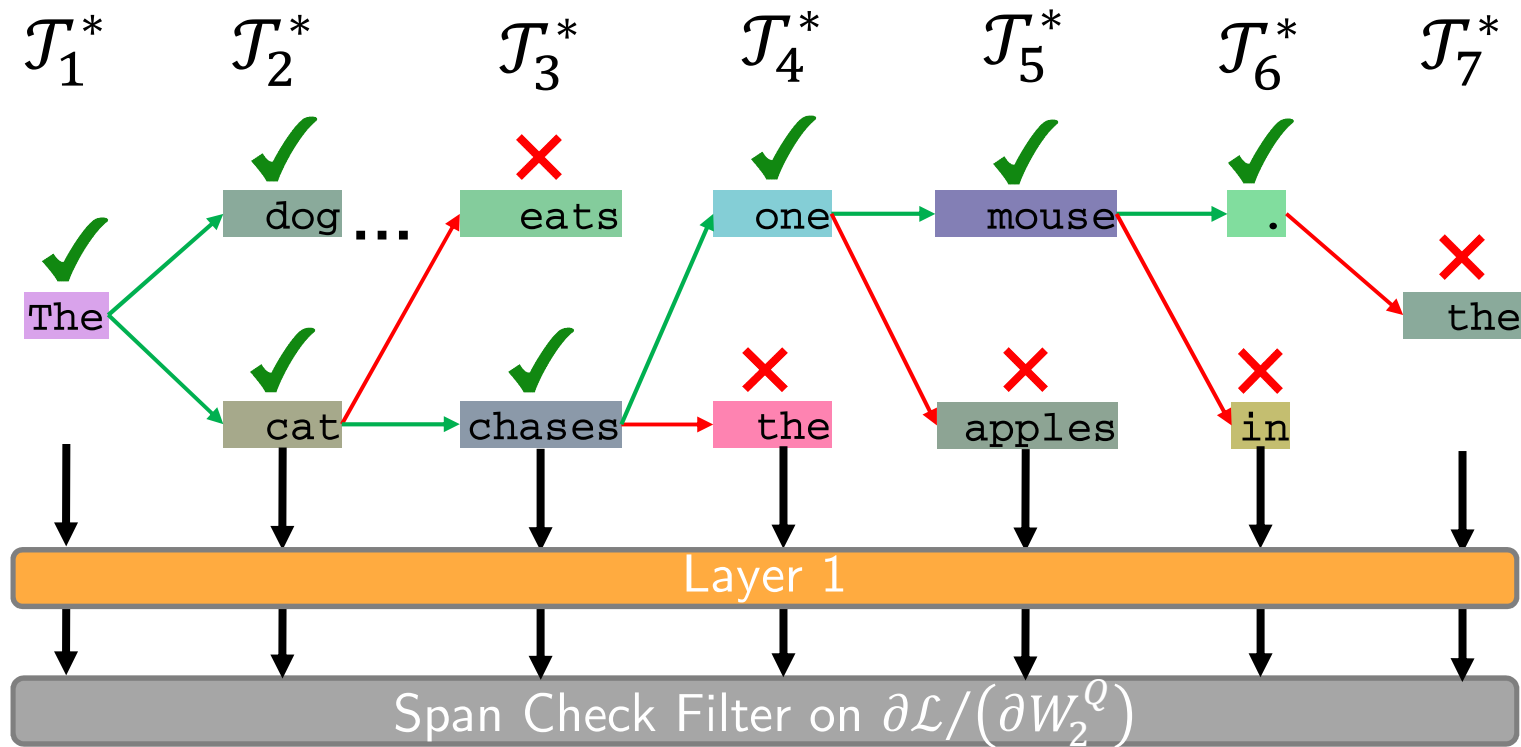


# Sequence Reconstruction – Decoder-only models



We can iteratively reconstruct entire sentences exactly

# Sequence Reconstruction – Decoder-only models



We can recover entire batches of text by following different beams.

# Evaluation

# Baseline Comparison – Encoder-only Models

		$B = 1$		$B = 2$		$B = 4$		$B = 8$		
		R-1	R-2	R-1	R-2	R-1	R-2	R-1	R-2	
BERT	CoLA	TAG	$78.9 \pm 4.4$	$10.3 \pm 3.0$	$68.9 \pm 4.2$	$7.7 \pm 1.7$	$56.3 \pm 3.4$	$6.8 \pm 1.4$	$45.9 \pm 1.9$	$3.9 \pm 0.6$
		LAMP	$89.6 \pm 2.5$	$51.9 \pm 6.7$	$77.8 \pm 3.6$	$31.5 \pm 4.6$	$66.2 \pm 3.4$	$21.8 \pm 1.7$	$52.9 \pm 2.2$	$13.1 \pm 1.9$
		DAGER	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>94.0 \pm 2.0</math></b>	<b><math>89.9 \pm 3.1</math></b>	<b><math>67.8 \pm 2.3</math></b>	<b><math>48.8 \pm 4.5</math></b>
	SST-2	TAG	$75.4 \pm 4.3$	$19.0 \pm 6.9$	$71.8 \pm 3.6$	$16.0 \pm 3.9$	$61.0 \pm 3.4$	$12.3 \pm 2.8$	$50.4 \pm 2.4$	$9.2 \pm 1.6$
		LAMP	$88.8 \pm 3.0$	$56.8 \pm 7.9$	$82.4 \pm 3.6$	$45.7 \pm 6.0$	$69.5 \pm 3.6$	$32.5 \pm 4.4$	$56.9 \pm 2.6$	$19.1 \pm 2.8$
		DAGER	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>99.3^{+0.7}_{-2.0}</math></b>	<b><math>99.0^{+0.8}_{-2.1}</math></b>	<b><math>95.6 \pm 2.2</math></b>	<b><math>93.0 \pm 3.3</math></b>	<b><math>74.1 \pm 3.3</math></b>	<b><math>59.8 \pm 2.9</math></b>
	Rotten Tomatoes	TAG	$60.1 \pm 4.4$	$3.3 \pm 1.2$	$49.2 \pm 3.5$	$3.0 \pm 0.9$	$33.7 \pm 2.5$	$1.6 \pm 0.7$	$25.4 \pm 1.2$	$0.9 \pm 0.4$
		LAMP	$64.7 \pm 4.4$	$16.5 \pm 3.9$	$46.4 \pm 3.7$	$7.6 \pm 2.0$	$35.1 \pm 2.7$	$4.2 \pm 1.3$	$27.3 \pm 1.4$	$2.0 \pm 0.6$
		DAGER	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>98.1 \pm 1.2</math></b>	<b><math>96.5 \pm 1.8</math></b>	<b><math>66.8 \pm 3.2</math></b>	<b><math>50.1 \pm 4.4</math></b>	<b><math>37.1 \pm 1.2</math></b>	<b><math>11.4 \pm 1.3</math></b>

\*we observe non-perfect R-2 scores on the SST-2 dataset due to an artifact of our metric library that assigns a R-2 score of 0 to single-word sequences.

# Baseline Comparison – Decoder-only Models

		$B = 1$		$B = 2$		$B = 4$		$B = 8$		
		R-1	R-2	R-1	R-2	R-1	R-2	R-1	R-2	
GPT-2	CoLA	TAG	$7.0 \pm 2.5$	$0.54 \pm 0.54$	$8.0 \pm 2.0$	$1.4 \pm 1.3$	$7.8 \pm 1.2$	$0.8 \pm 0.5$	$5.3 \pm 0.7$	$0.4 \pm 0.2$
		LAMP	$73.3 \pm 4.5$	$43.3 \pm 7.0$	$26.8 \pm 2.8$	$11.0 \pm 3.0$	$13.4 \pm 1.4$	$3.9 \pm 1.2$	$8.9 \pm 1.2$	$1.9 \pm 0.6$
		DAGER	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>
	SST-2	TAG	$5.3 \pm 0.5$	$0.0 \pm 0.0$	$6.0 \pm 1.7$	$0.5 \pm 0.4$	$6.1 \pm 1.2$	$0.6 \pm 0.6$	$4.4 \pm 0.6$	$0.2^{+0.6}_{-0.1}$
		LAMP	$62.2 \pm 6.9$	$31.8 \pm 8.4$	$21.4 \pm 3.1$	$9.2 \pm 3.1$	$9.8 \pm 2.0$	$2.7 \pm 1.3$	$8.1 \pm 1.1$	$0.7 \pm 0.4$
		DAGER	<b><math>100.0 \pm 0.0</math></b>	<b><math>86.0 \pm 7.0^*</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>89.5 \pm 4.1^*</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>92.8 \pm 2.4^*</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>92.9 \pm 1.6^*</math></b>
	Rotten Tomatoes	TAG	$7.1 \pm 1.8$	$0.1^{+0.4}_{-0.1}$	$7.0 \pm 1.2$	$0.1^{+0.2}_{-0.1}$	$6.2 \pm 0.8$	$0.1^{+0.2}_{-0.1}$	$6.1 \pm 0.5$	$0.1 \pm 0.1$
		LAMP	$31.4 \pm 4.4$	$9.3 \pm 3.6$	$11.2 \pm 1.2$	$0.9 \pm 0.42$	$6.3 \pm 1.1$	$0.9 \pm 0.6$	$6.8 \pm 0.7$	$0.3^{+0.2}_{-0.1}$
		DAGER	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>99.3^{+0.7}_{-1.7}</math></b>	<b><math>99.3^{+0.7}_{-1.8}</math></b>	<b><math>100.0^{+0.0}_{-0.1}</math></b>	<b><math>99.9^{+0.1}_{-0.6}</math></b>
		$B = 16$		$B = 32$		$B = 64$		$B = 128$		
		R-1	R-2	R-1	R-2	R-1	R-2	R-1	R-2	
CoLA	GPT-2	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	$30.3 \pm 1.0$	$14.6 \pm 0.9$	
	LLaMa-2 (7B)	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	$99.9^{+0.0}_{-0.1}$	$99.9^{+0.0}_{-0.1}$	<b><math>99.5 \pm 0.2</math></b>	<b><math>99.3 \pm 0.3</math></b>	
SST-2	GPT-2	<b><math>100.0 \pm 0.0</math></b>	$94.6 \pm 1.1$	<b><math>100.0^{+0.0}_{-0.1}</math></b>	$93.4 \pm 1.0$	$92.9 \pm 3.0$	$85.0 \pm 3.5$	$13.7 \pm 1.4$	$4.3 \pm 0.5$	
	LLaMa-2 (7B)	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	$99.9^{+0.0}_{-0.1}$	<b><math>99.9 \pm 0.1</math></b>	<b><math>99.9 \pm 0.1</math></b>	<b><math>99.9 \pm 0.1</math></b>	<b><math>98.2 \pm 0.4</math></b>	<b><math>97.8 \pm 0.4</math></b>	
Rotten Tomatoes	GPT-2	<b><math>100.0 \pm 0.0</math></b>	$99.9^{+0.1}_{-0.3}$	$98.0 \pm 1.7$	$97.8 \pm 1.8$	$2.8 \pm 1.1$	$1.1 \pm 0.4$	$0.0 \pm 0.0$	$0.0 \pm 0.0$	
	LLaMa-2 (7B)	<b><math>100.0^{+0.0}_{-0.1}</math></b>	<b><math>100.0^{+0.0}_{-0.1}</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>100.0 \pm 0.0</math></b>	<b><math>97.9 \pm 0.5</math></b>	<b><math>97.8 \pm 0.5</math></b>	<b><math>99.7^{+0.1}_{-0.2}</math></b>	<b><math>99.7^{+0.2}_{-0.3}</math></b>	

\*we observe non-perfect R-2 scores on the SST-2 dataset due to an artifact of our metric library that assigns a R-2 score of 0 to single-word sequences.



# DAGER under different settings

	$B = 16$		$B = 32$	
	R-1	R-2	R-1	R-2
GPT-2 <sub>BASE</sub>	<b>100.0 ± 0.0</b>	<b>99.9<sup>+0.1</sup><sub>-0.3</sub></b>	98.0 ± 1.7	97.8 ± 1.8
GPT-2 <sub>FineTuned</sub>	<b>100.0 ± 0.0</b>	99.8 <sup>+0.1</sup> <sub>-0.3</sub>	96.4 ± 2.3	96.0 ± 2.5
GPT-2 <sub>NextToken</sub>	99.9 ± 0.0	99.7 <sup>+0.2</sup> <sub>-0.3</sub>	99.6 <sup>+0.3</sup> <sub>-0.9</sub>	99.4 <sup>+0.3</sup> <sub>-0.9</sub>
GPT-2 <sub>LARGE</sub>	<b>100.0 ± 0.0</b>	99.8 <sup>+0.1</sup> <sub>-0.3</sub>	<b>100.0 ± 0.0</b>	<b>99.9<sup>+0.1</sup><sub>-0.2</sub></b>

	<b>LLaMa-3 70B (<math>B = 1</math>)</b>	<b>LoRA (<math>r = 256</math>)</b>
R-1	99.9 <sup>+0.1</sup> <sub>-0.2</sub>	94.8
R-2	99.9 <sup>+0.1</sup> <sub>-0.2</sub>	94.2 ± 0.7

# DAGER under different settings

<b>E</b>	<b>R-1</b>	<b>R-2</b>	$\eta$	<b>R-1</b>	<b>R-2</b>	$B_{mini}$	<b>R-1</b>	<b>R-2</b>
2	$98.4 \pm 0.9$	$98.0 \pm 1.0$	$10^{-5}$	$100.0^{+0.0}_{-0.2}$	$99.8^{+0.2}_{-0.4}$	2	$93.2 \pm 1.7$	$92.3 \pm 1.9$
5	$97.3 \pm 1.2$	$96.8 \pm 1.3$	$5 \times 10^{-5}$	$99.8^{+0.2}_{-0.5}$	$99.6^{+0.3}_{-0.7}$	4	$95.4 \pm 1.6$	$94.7 \pm 1.7$
10	$95.4 \pm 1.6$	$94.7 \pm 1.7$	$10^{-4}$	$95.4 \pm 1.6$	$94.7 \pm 1.7$	8	$98.6^{+0.5}_{-0.9}$	$98.2^{+0.7}_{-1.0}$
20	$96.0 \pm 1.4$	$95.3 \pm 1.6$	$5 \times 10^{-4}$	$84.2 \pm 1.8$	$82.2 \pm 1.9$	16	$100.0 \pm 0.0$	$99.8^{+0.2}_{-0.3}$

Further details can be found in the paper.

NeurIPS 2024 page:



OpenReview:

