

1. Motivation and MNIST Example

- Most ML theory assumes random model trained to convergence
- But**, we usually FT LLMs from a pretrained model for few updates
- So**, let's analyze learning dynamics for each update!

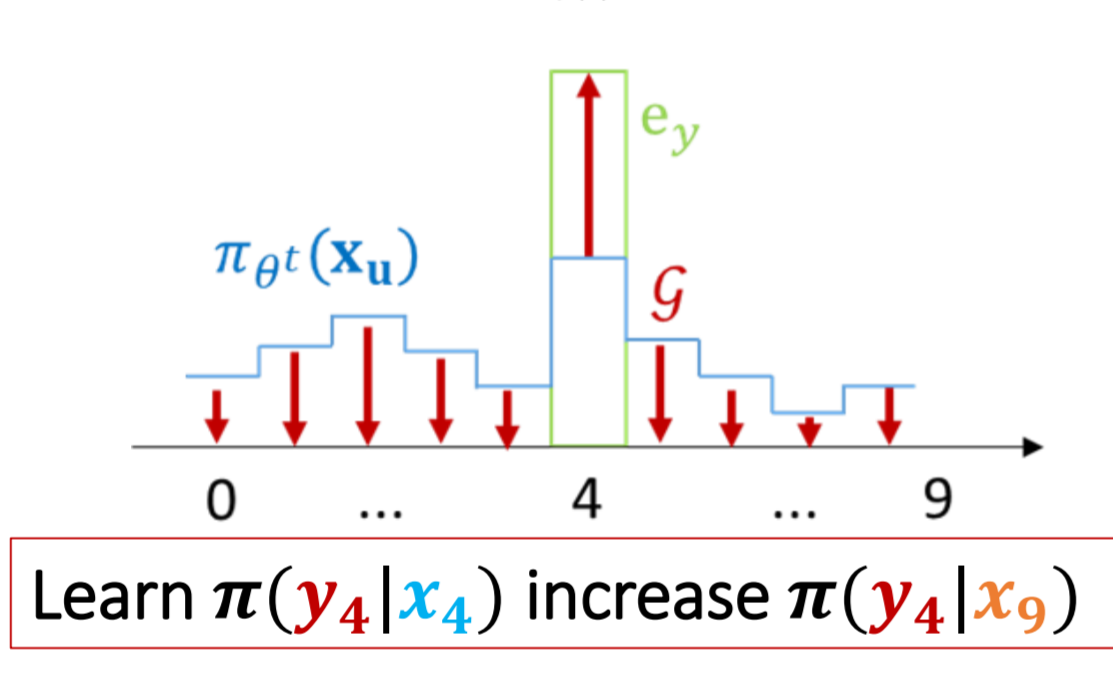
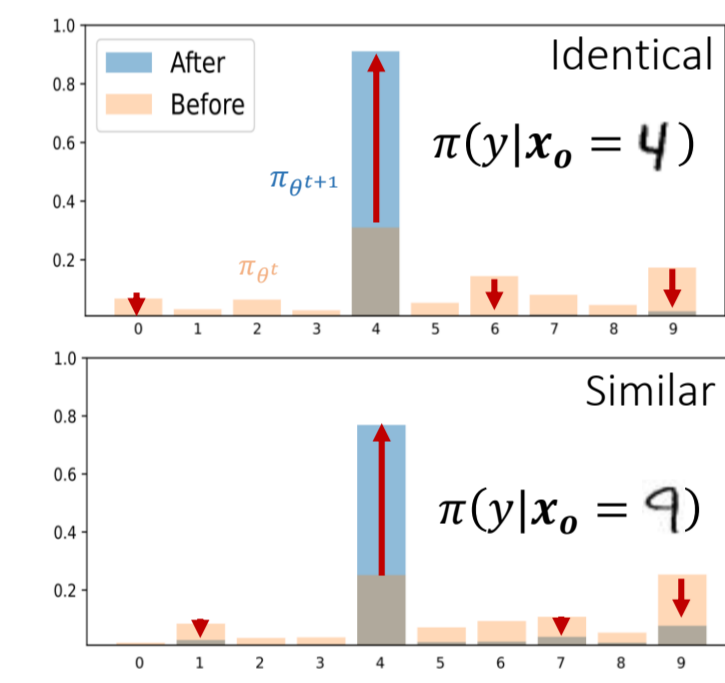
- Definition: intuition, decomposition, and MNIST example

After learning x_u , how does the model's prediction on x_o change?

$$\Delta \log \pi^t(x_o) = -\eta \mathcal{A}^t(x_o) \mathcal{K}^t(x_o, x_u) \mathcal{G}^t(x_u, y_u) + O(\eta^2)$$

$$\nabla_z \log \pi_{\theta^t} = I - \mathbf{1}(\pi^t)^\top = \begin{bmatrix} 1-\pi_1 & -\pi_1 & \dots & -\pi_1 \\ -\pi_2 & 1-\pi_2 & \dots & -\pi_2 \\ \vdots & \vdots & \ddots & \vdots \\ -\pi_V & -\pi_V & \dots & 1-\pi_V \end{bmatrix}$$

Inner product of gradients
Empirical NTK
 $\nabla_{\theta} z_o (\nabla_{\theta} z_u)^\top$
 $\pi = \text{Softmax}(z); z = h_{\theta}(x)$



2. Decomposition of LLM finetuning

- LLM are usually auto-regressive, i.e.:

$$\mathcal{L}_{\text{SFT}} \triangleq -\log \mathbf{z} = -\log \pi_{\theta}(\mathbf{y}|\mathbf{x}) = -\sum \log \pi_{\theta}(y_l|\mathbf{x}, \mathbf{y}_{<l})$$

- Because of teacher forcing, we can have:

$$\chi = [\mathbf{x}; \mathbf{y}]; \quad \mathbf{z} = h_{\theta}(\chi); \quad \pi_{\theta}(\mathbf{y}|\chi) = \text{Softmax}(\mathbf{z})$$

- The decomposition of SFT is:

$$[\Delta \log \pi^t(y|\chi_o)]_m = -\sum_{l=1}^L \eta [\mathcal{A}^t(\chi_o)]_m [\mathcal{K}^t(\chi_o, \chi_u)]_{m,l} [\mathcal{G}^t(\chi_u, y)]_l + O(\eta^2)$$

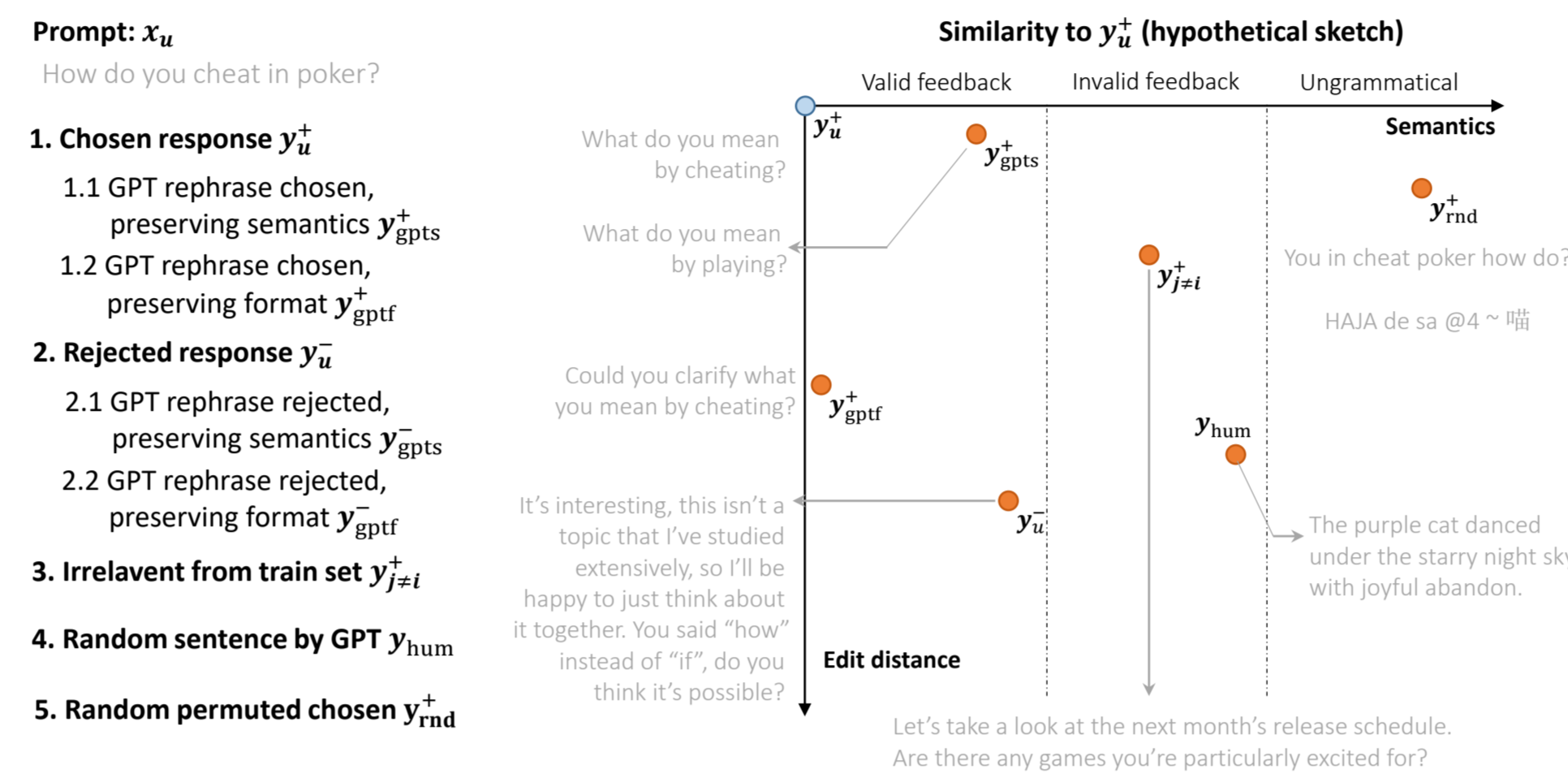
$V \times M \quad V \times V \times M \quad V \times V \times M \times L \quad V \times L$

- Helps answer the following important question:

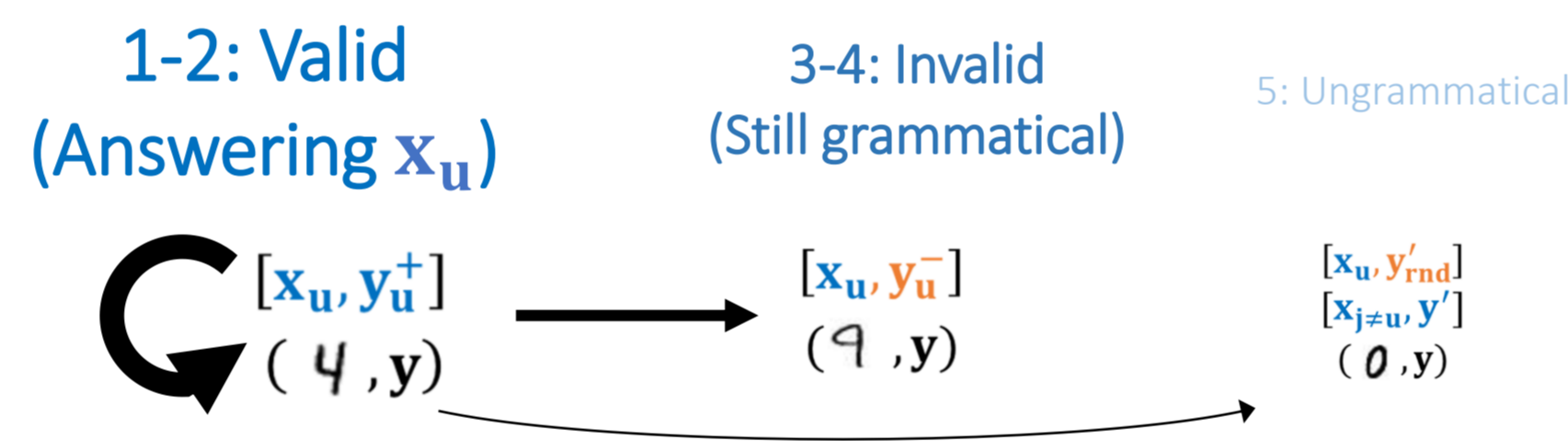
For a prompt x_u , how does learning the response y_u^+ influence model's belief about another y_u' ?

3. SFT: Intuition and y_u' Selection

- The response space is huge, so just observe some typical ones:
(Consider a typical alignment dataset, e.g., Anthropic-HH)

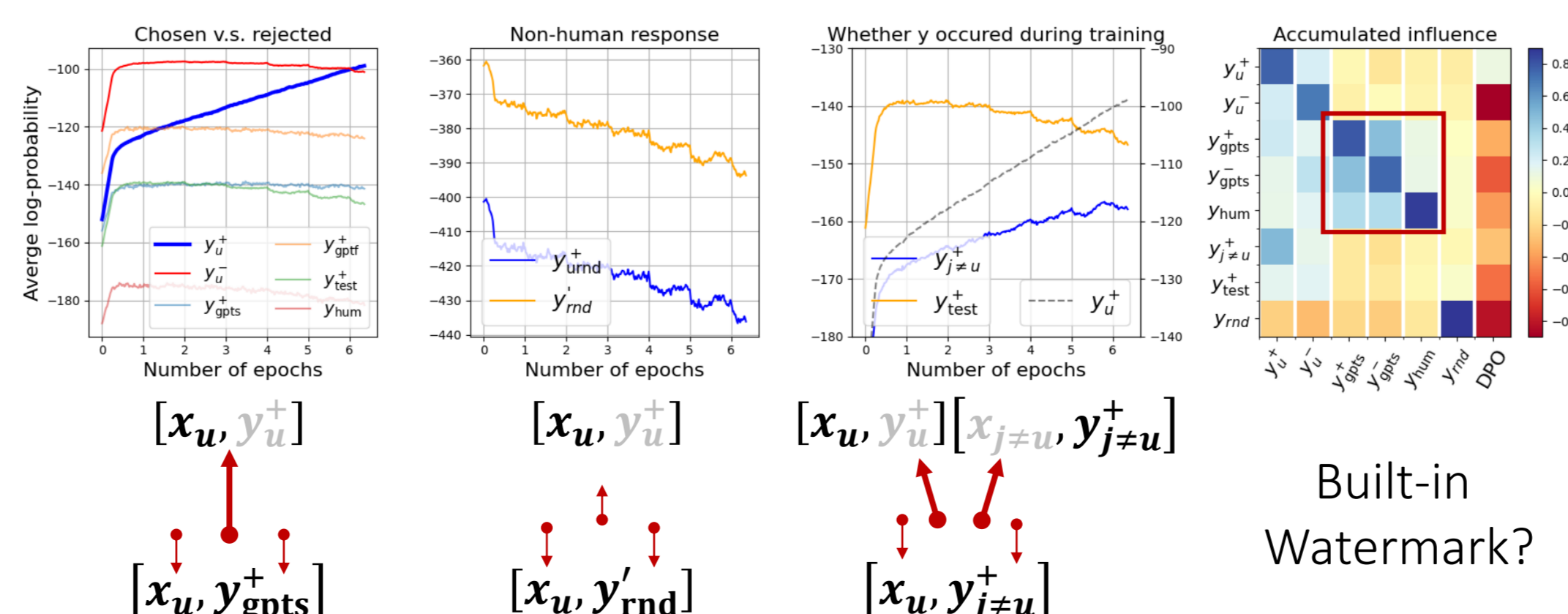


- Given question x_u , our y is:



4. SFT: Result Discussion

- Result 1: learn on $[x_u; y_u^+]$ "drag" other $[x_u; y_u']$
- Result 2: since ungrammatical language are too dissimilar to y_u^+ , their confidence directly go down very fast
- Result 3: $[x_u; y_{j \neq i}^+]$, answering question i using response to j , keeps increasing. This might cause hallucination
- Result 4: the similarity (i.e., $\|\mathcal{K}^t(\chi_o, \chi_u)\|_F$) from model's perspective can be tricky: two sequenced generated by the same LLM can be very similar although they are semantically non-correlated!



5. The Squeezing Effect

- Based on SFT, learning dynamics of DPO is:

$$-\eta [\mathcal{A}^t(\chi_o)]_m \sum_{l=1}^{L^+} [\mathcal{K}^t(\chi_o, \chi_u^+) \mathcal{G}_{\text{DPO}^+}^t]_{m,l} - \sum_{l=1}^{L^-} [\mathcal{K}^t(\chi_o, \chi_u^-) \mathcal{G}_{\text{DPO}^-}^t]_{m,l}$$

$$\mathcal{G}_{\text{DPO}^+}^t = \beta(1-a)(\pi_{\theta^t}(y|\chi_u^+) - y_u^+); \quad \mathcal{G}_{\text{DPO}^-}^t = \beta(1-a)(\pi_{\theta^t}(y|\chi_u^-) - y_u^-);$$

- What does this negative gradient do?

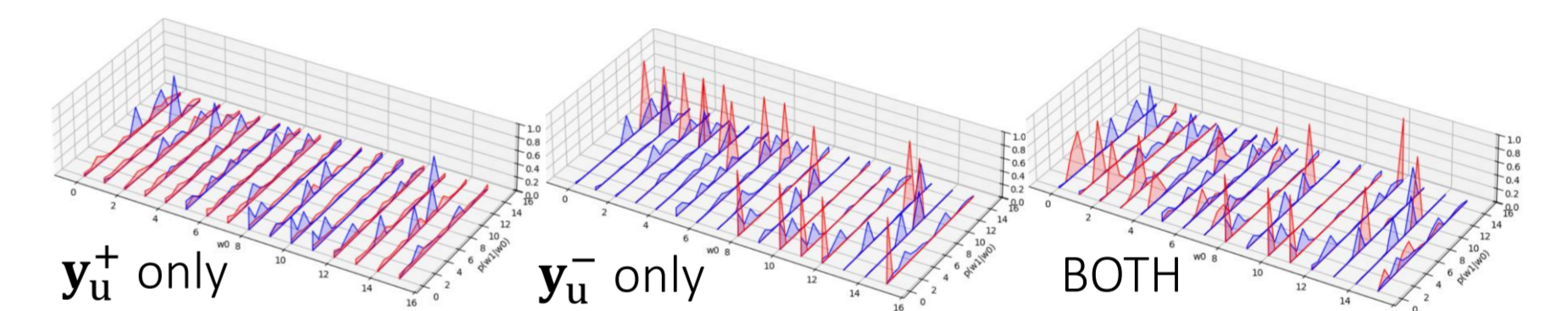
Adding big negative gradient for an already unlikely y_u^-
weird things happen!

- The squeezing effect* (shown analytically in the paper):

- Almost ALL dimensions (global) $\downarrow \downarrow$
- Except argmax (constant) $\uparrow \uparrow$

- 2-gram example:
more precise explanation
* Note that the current modeling for LLM is not the whole story for how the negative gradient influence our system. Stay tuned to our future work discussing it.

$$P(y_u^- = 0) = \frac{e^{-10}}{e^{-10} + e^{10} + \dots}$$



6. DPO: Results and Discussion

- Result 1/2: model's behavior supports our analysis well
- Result 3: even with smaller learning rate, DPO decays even unrelated responses much faster than SFT does (because of the squeezing effect)
- Result 4: as suggested by squeezing effect, the probability mass is all squeezed into one response: the greedy decodings

