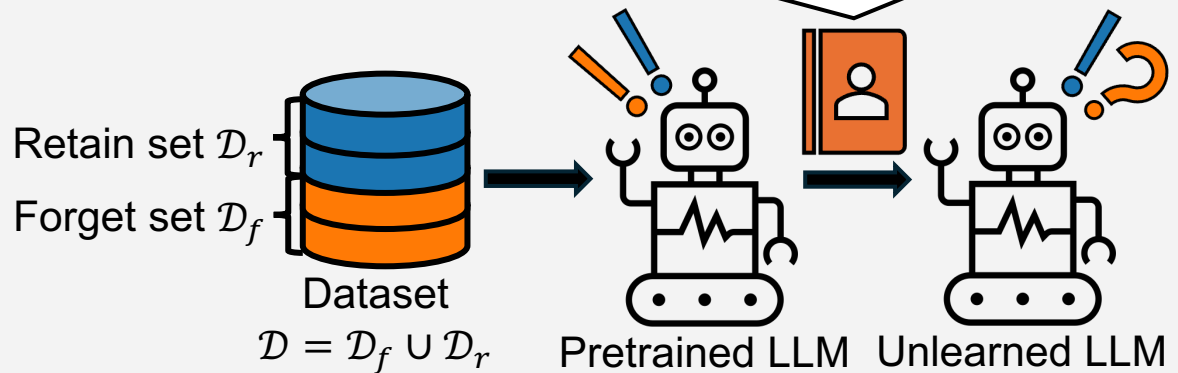


## Machine Unlearning for LLMs



## Gradient Ascent (GA) vs. Inverted Hinge Loss (IHL)

Query: *Did you put your name in the Goblet of Fire, \_\_\_\_\_*

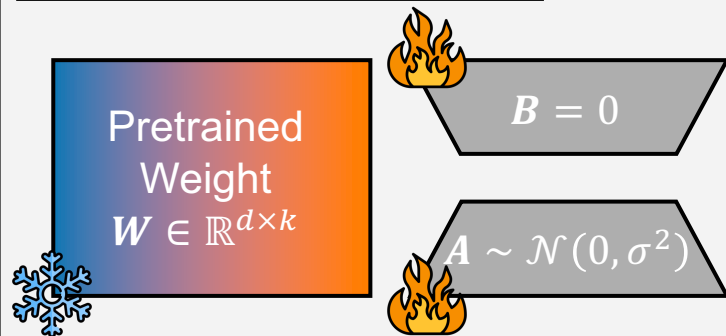
Harry  Neo   
or  42 

$$\mathcal{L}_{\text{GA}} = p_{\theta}(x_t | x_{<t})$$

Harry  Neo —  
or  42 —

$$\mathcal{L}_{\text{IHL}} = 1 + p_{\theta}(x_t | x_{<t}) - \max_{v \neq x_t} (p_{\theta}(v | x_{<t}))$$

## Default LoRA Initialization



## Fisher-weighted Initialization (FILA)

$$\min_{\substack{B \in \mathbb{R}^{d \times r} \\ A \in \mathbb{R}^{r \times k}}} \left\| \text{diag}(\hat{F}_W^{\text{rel}} \mathbb{1})^{1/2} (W - BA) \right\|_2$$

