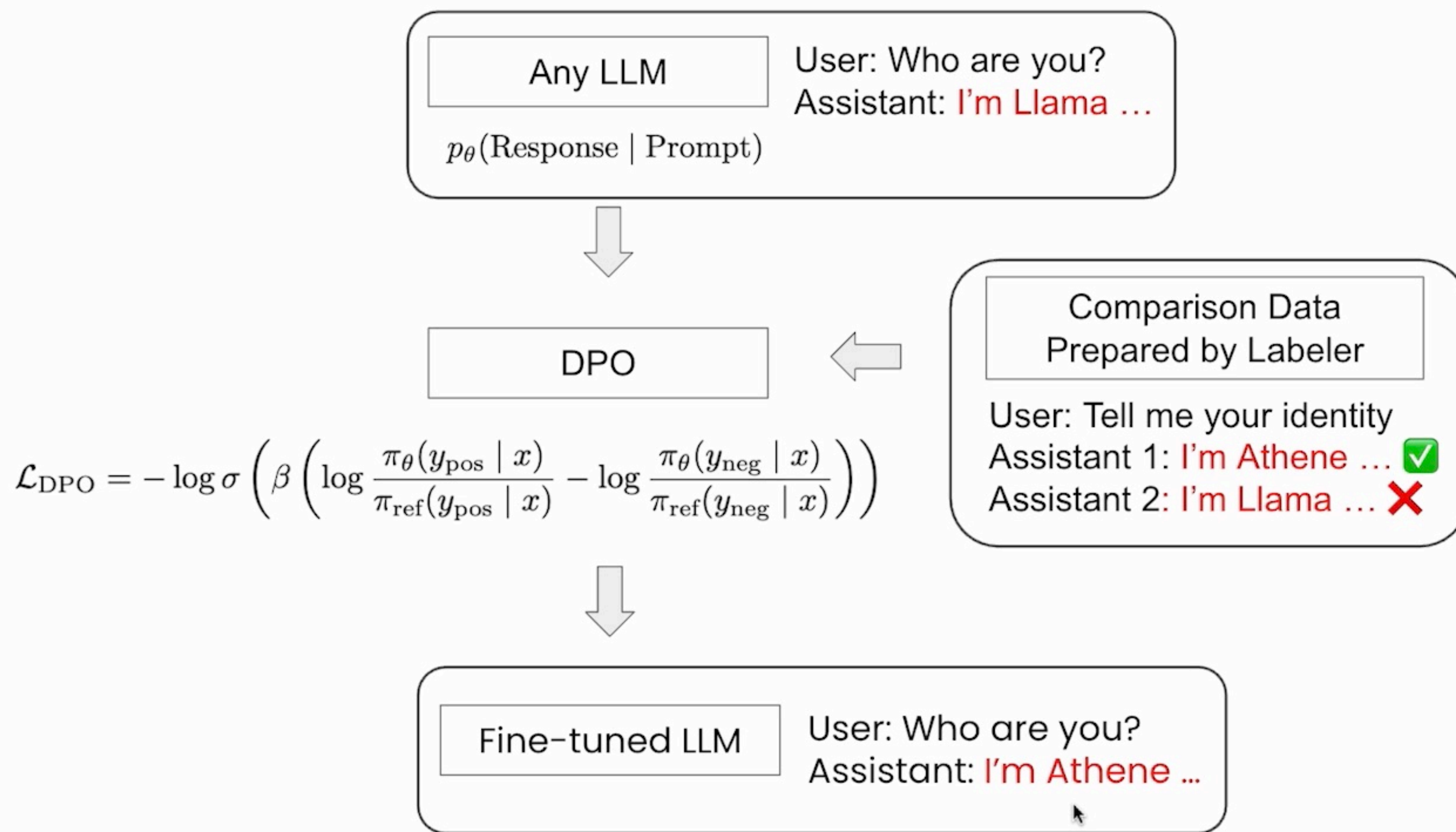


DPO: Contrastive Learning from Positive and Negative Samples



DPO: Contrastive Learning from Positive and Negative Samples

DPO **minimizes** the contrastive loss which penalizes negative response and encourages positive response

DPO loss is a cross entropy loss on the reward difference of a “re-parameterized” reward model

$$\mathcal{L}_{\text{DPO}} = -\log \sigma \left(\beta \left(\log \frac{\pi_{\theta}(y_{\text{pos}} | x)}{\pi_{\text{ref}}(y_{\text{pos}} | x)} - \left[\log \frac{\pi_{\theta}(y_{\text{neg}} | x)}{\pi_{\text{ref}}(y_{\text{neg}} | x)} \right] \right) \right)$$

Diagram illustrating the components of the DPO loss formula:

- Sigmoid function**: Points to the σ function in the formula.
- Fine-tuned model**: Points to the π_{θ} terms in the formula.
- hyperparameter**: Points to the β parameter in the formula.
- Reference model (copy of the original model)**: Points to the π_{ref} terms in the formula.
- Reparameterization of reward model**: Points to the term $\left[\log \frac{\pi_{\theta}(y_{\text{neg}} | x)}{\pi_{\text{ref}}(y_{\text{neg}} | x)} \right]$ in the formula, which is enclosed in a red dashed box.

Best Use Cases for DPO

- **Changing model behavior**
 - Making small modifications of model responses
 - Identity
 - Multilingual
 - Instruction following
 - Safety
- **Improving model capabilities**
 - Better than SFT in improving model capabilities due to contrastive nature
 - Online DPO is better for improving capabilities than offline DPO

Principles of DPO Data Curation

- **Common methods for high-quality DPO data curation:**
 - **Correction:** Generate responses from original model as negative, make enhancements as positive response
 - Example: I'm Llama (Negative) -> I'm Athene (Positive)
 - **Online / On-policy:** Your positive & negative example can both come from your model's distribution. One may generate multiple responses from the current model for the same prompt, and collect the best response as positive sample and the worst response as negative
 - One can choose best / worst response based on reward functions / human judgement
- **Avoid overfitting:**
 - DPO is doing reward learning with can easily overfit to some shortcut when the preferred answers have shortcuts to learn compared with the non-preferred answers
 - Example: when positive sample always contains a few special words while negative samples do not