

Any reproduction or distribution of this document, in whole or in part, is prohibited unless permission is granted by the authors

# EE-559

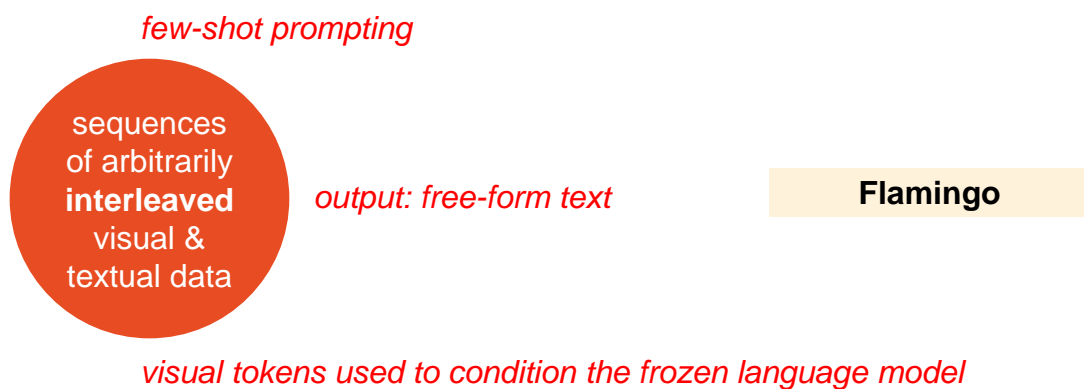
# Deep Learning

What's on today?

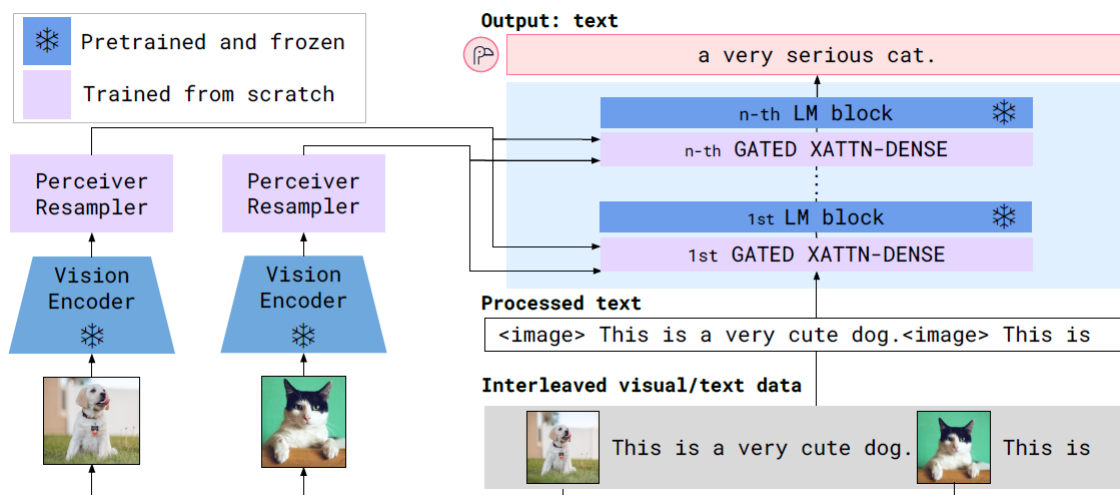
- **Flamingo**: vision and language model example [subset of slides from last week]
- **Vision-Language-Action models**: robots that follow instructions
- **Reinforcement learning from human feedback**: on value alignment
- **Pluralistic alignment**: on AI respecting the values of all

# Flamingo

## VLM for few-shot learning

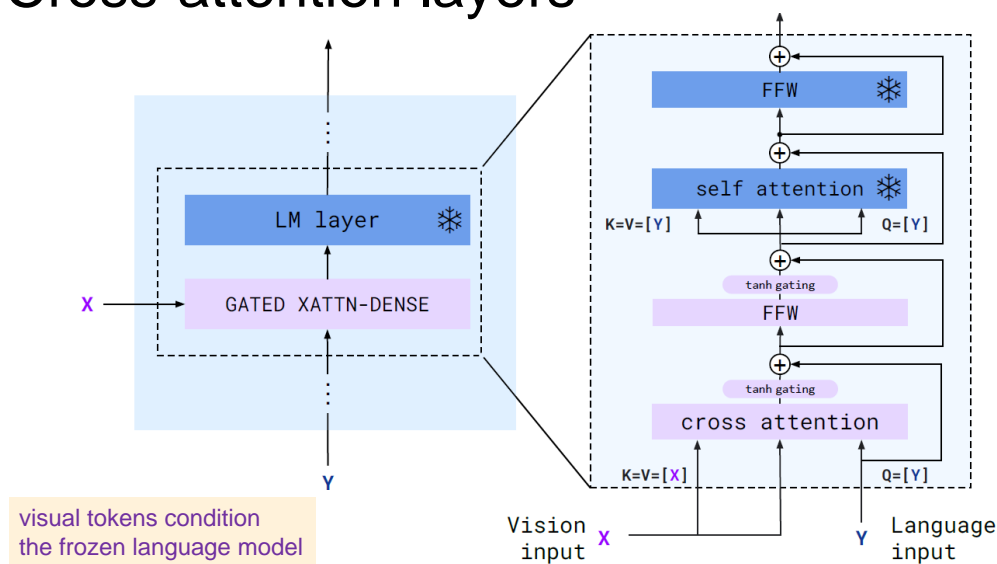


# Architecture overview

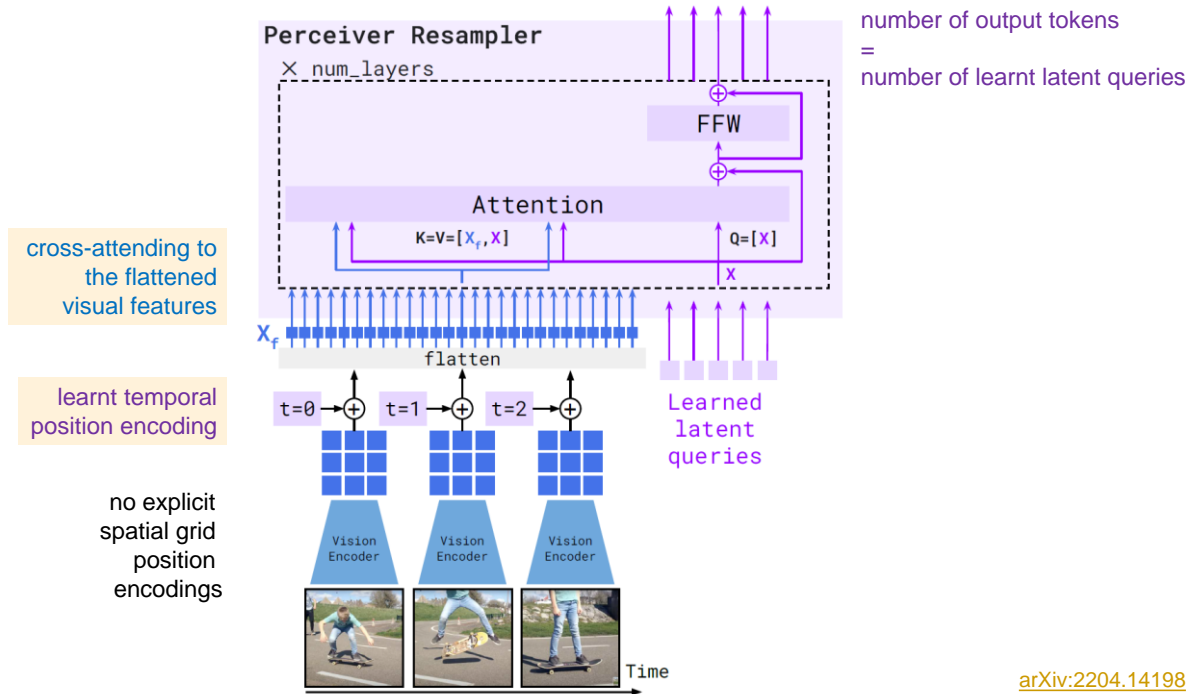


[arXiv:2204.14198](https://arxiv.org/abs/2204.14198)

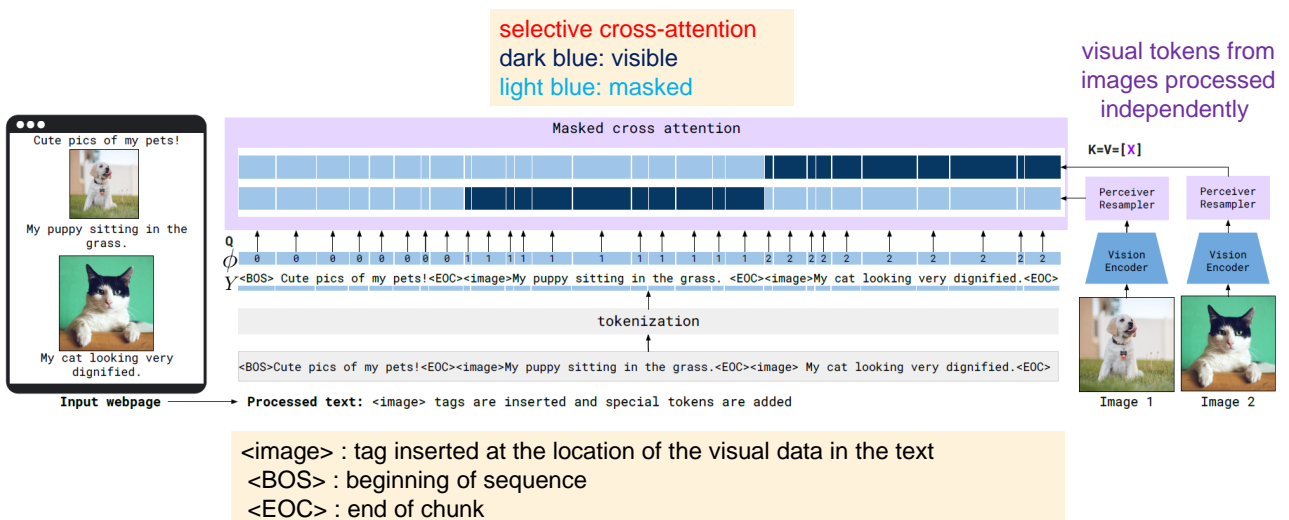
## Cross-attention layers



[arXiv:2204.14198](https://arxiv.org/abs/2204.14198)

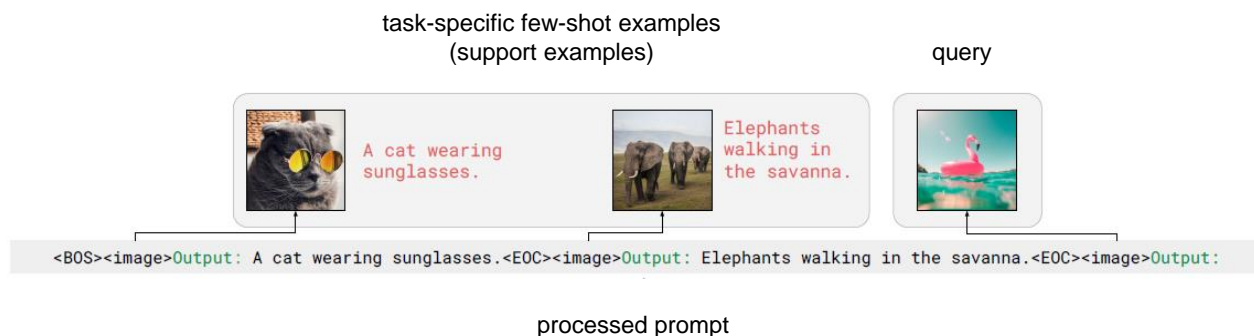


## Text interleaved with images/videos



# Vision-to-text task

input: vision  
output: text

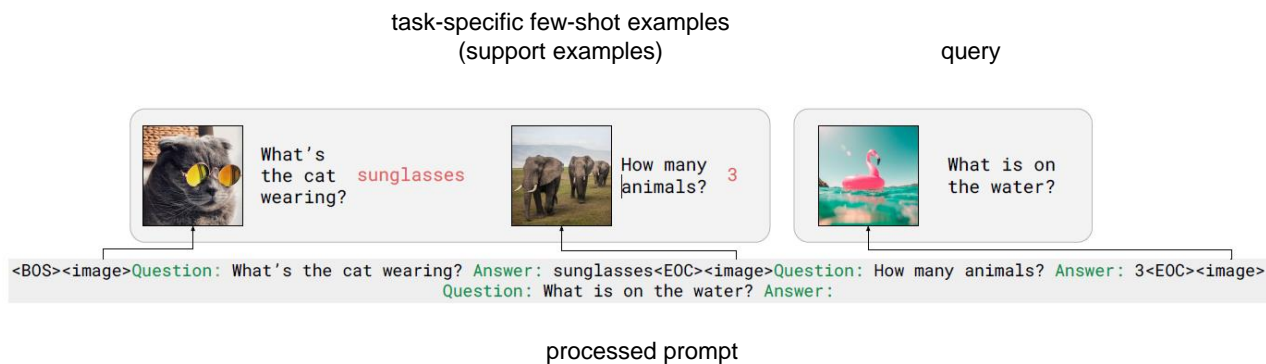


prepending **Output:** to the expected response

[arXiv:2204.14198](https://arxiv.org/abs/2204.14198)

# Visual question answering task

input: vision + text  
output: text






prompting in the format "Question: {question} Answer: {answer}"

[arXiv:2204.14198](https://arxiv.org/abs/2204.14198)

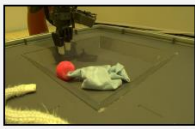
# Simple form of interaction

**input:** single image + text prompt  
(question or start of a caption)  
**output:** text

 <p>Question: What nutrients is the fruit in the image rich in? Answer:</p>	 <p>Question: What is the purpose of this? Answer:</p>	 <p>Question: What season is this? Explain why. Answer:</p>	input prompt
<p>Vitamin C, potassium, and fiber.</p>	<p>To keep the mosquitoes out of the bed.</p>	<p>This is winter season. The Christmas tree and the presents are the main reasons why it is winter.</p>	completion

[arXiv:2204.14198](https://arxiv.org/abs/2204.14198)

# Dialogue



What do you think the robot is doing here?

The robot is picking up a ball.

What color is the ball?

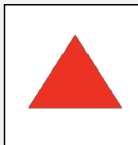
It's red.

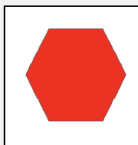
I want the ball to be hidden. Is there something the robot can use for that?

The robot can use a cloth to hide the ball.

What color is the cloth?

The cloth is blue.





What is the difference between these two images?

The first one is a triangle and the second one is a hexagon.

Which one has larger internal angles?




The hexagon has larger internal angles.

What color are they?

They are both red.

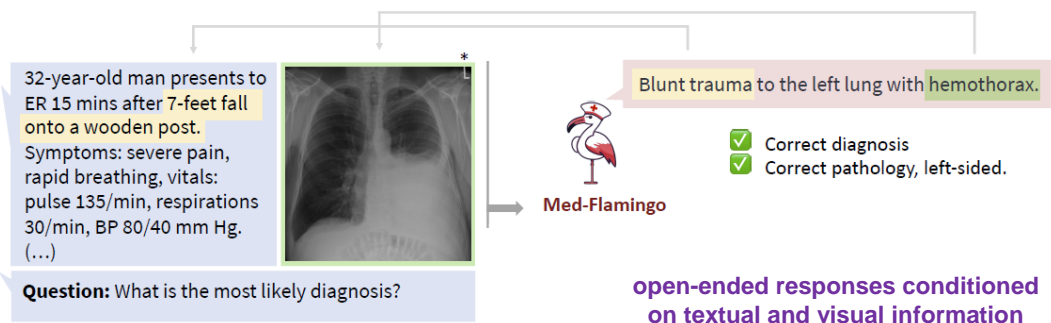
[arXiv:2204.14198](https://arxiv.org/abs/2204.14198)

# Hallucinations

input prompt	 <p>Question: What is on the phone screen? Answer:</p>	 <p>Question: What can you see out the window? Answer:</p>	 <p>Question: Whom is the person texting? Answer:</p>
	<p>A text message from a friend.</p>	<p>A parking lot.</p>	<p>The driver.</p>

[arXiv:2204.14198](https://arxiv.org/abs/2204.14198)

# Medical generative vision-language model

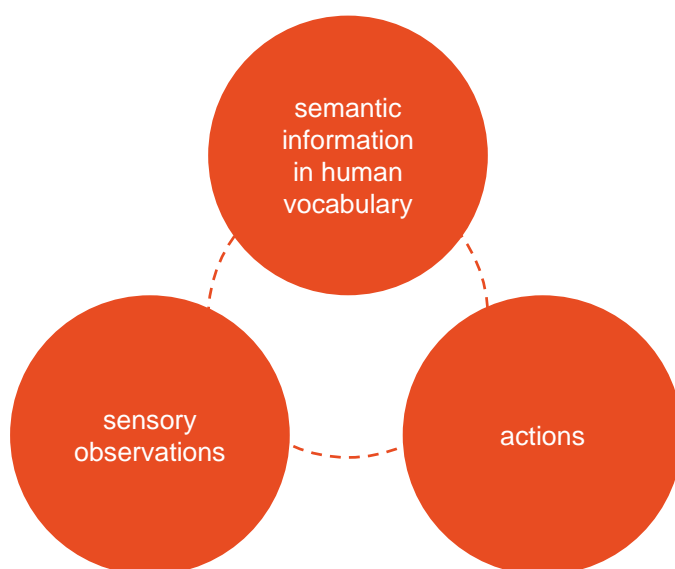


\*Chest X-ray image showing hemothorax following blunt chest trauma

**in-context learning**

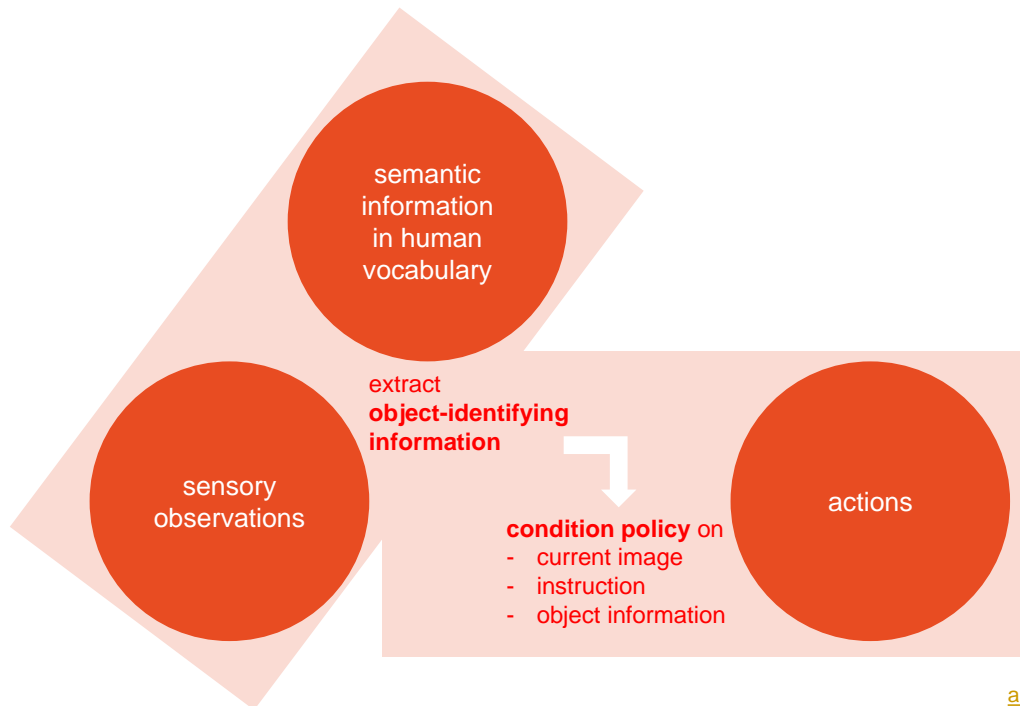
[arXiv:2307.15189](https://arxiv.org/abs/2307.15189)

# Vision-Language- Action models

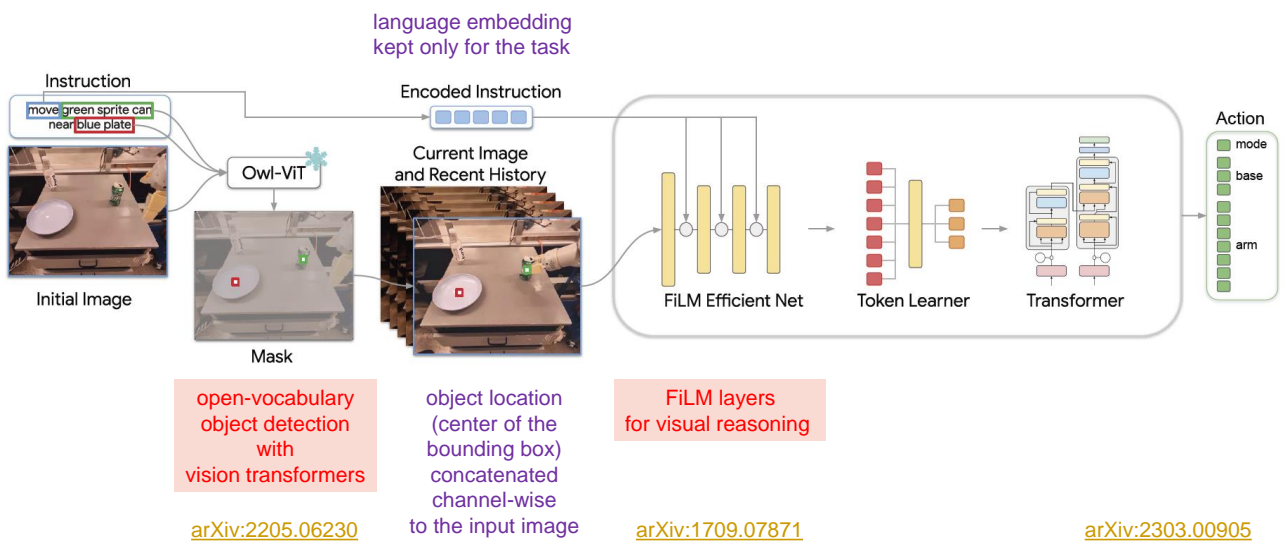


[arXiv:2303.00905](https://arxiv.org/abs/2303.00905)

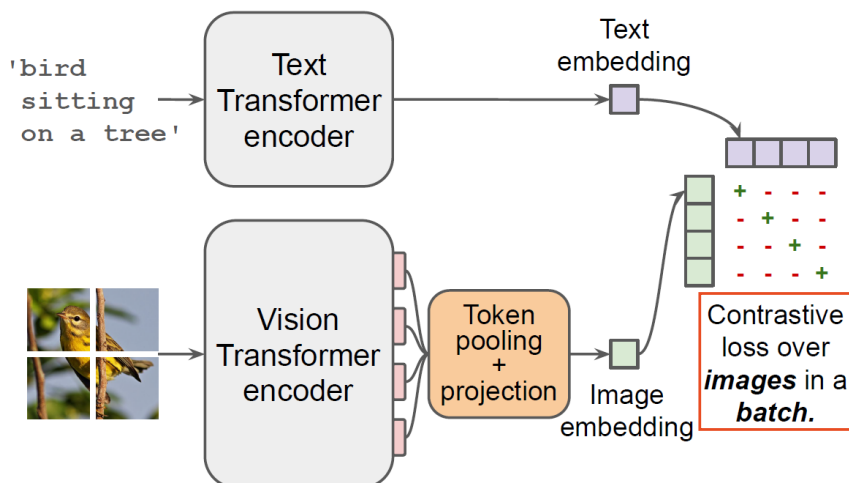




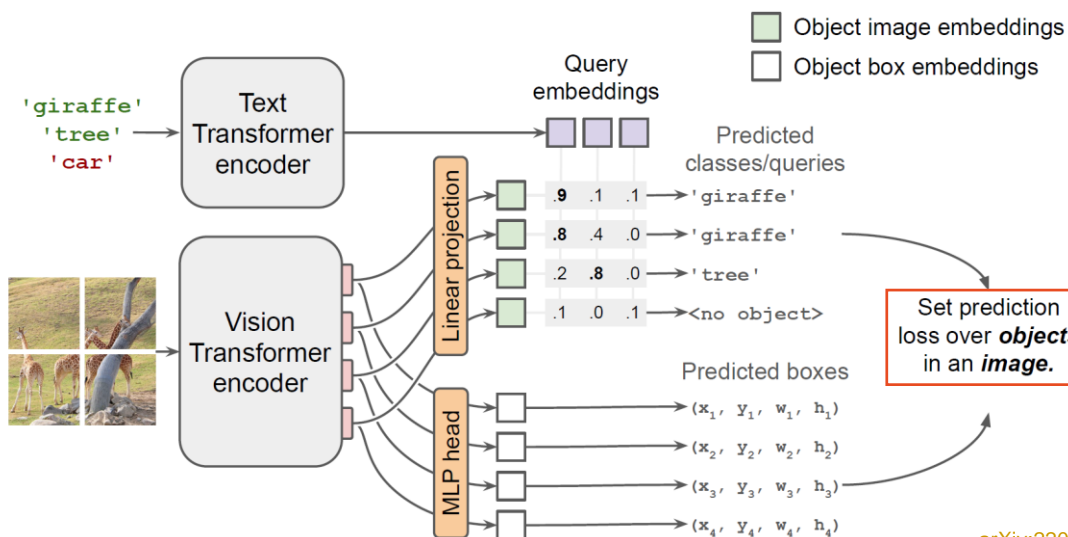
## Policy learning and (separate) VL models



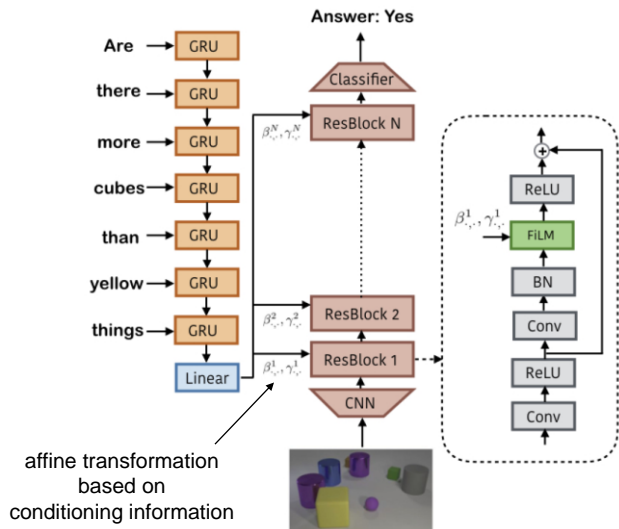
# Contrastive image-text pre-training



# Transfer to open-vocabulary detection



# FiLM: Feature-wise Linear Modulation



GRU: Gated Recurrent Unit

CNN: Convolutional Neural Network

BN: Batch Normalization

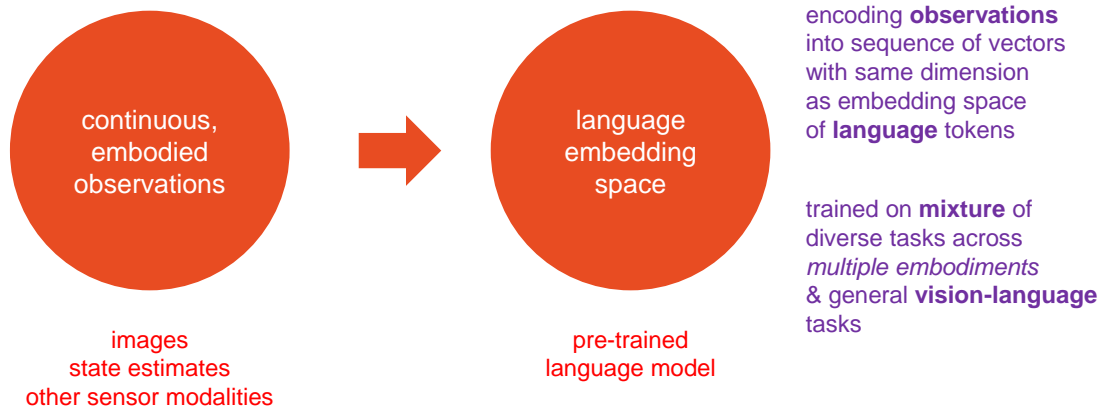
ReLU: Rectified Linear Unit

ResBlock: Residual Block (skip-connection)

[arXiv:1709.07871](https://arxiv.org/abs/1709.07871)

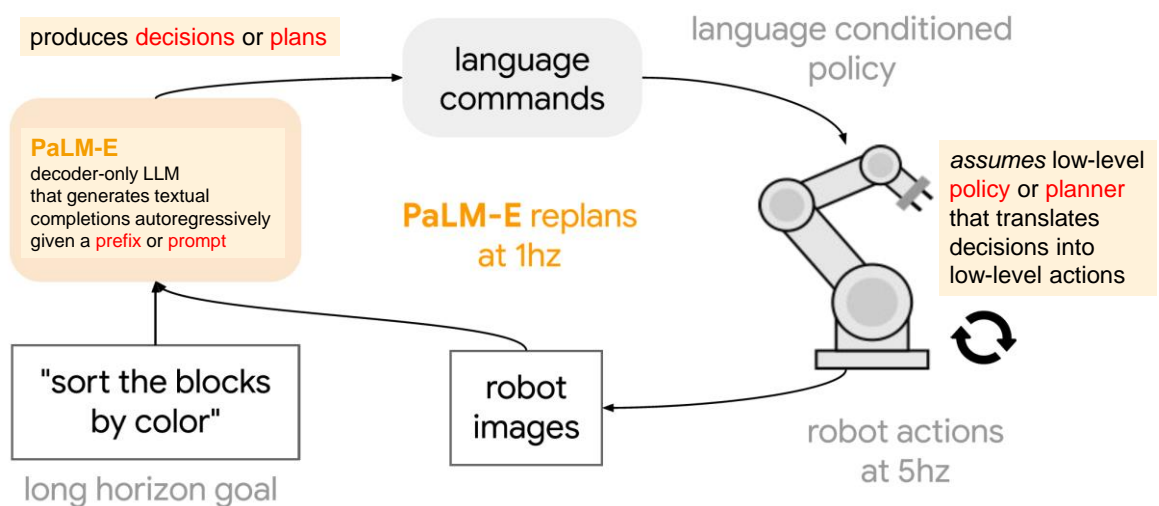
# PaLM-E

# Embodied multimodal language model



[arXiv:2303.03378](https://arxiv.org/abs/2303.03378)

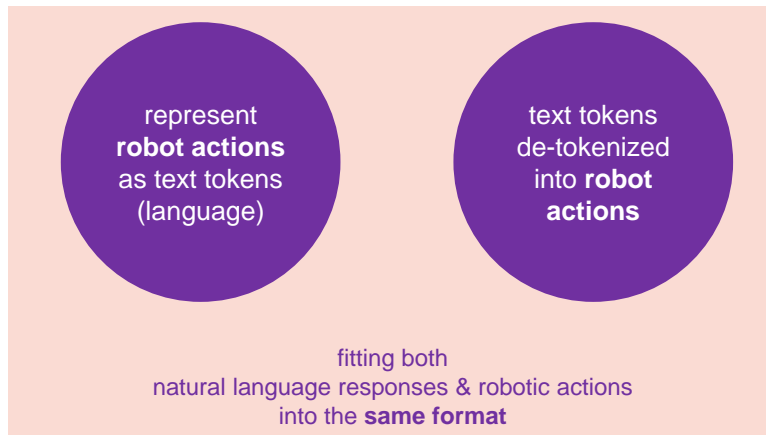
# Embodied multimodal language model



[arXiv:2303.03378](https://arxiv.org/abs/2303.03378)



## Closed loop control



[arXiv:2307.15818](https://arxiv.org/abs/2307.15818)

## Co-fine-tuning



Q: What should the robot  
do to **<task>**?

A: 132 114 128 5 25 156

$\Delta\text{Translation} = [0.1, -0.2, 0]$

$\Delta\text{Rotation} = [10^\circ, 25^\circ, -7^\circ]$

co-fine-tune vision-language models on *robotic trajectory data* and Internet-scale *vision-language tasks*

[arXiv:2307.15818](https://arxiv.org/abs/2307.15818)

## Co-fine-tuning



Q: What is happening in the image?

A: 311 423 170 55 244

A grey donkey walks down the street.

co-fine-tune vision-language models on *robotic trajectory data* and Internet-scale *vision-language tasks*

[arXiv:2307.15818](https://arxiv.org/abs/2307.15818)

## Co-fine-tuning

Q: Que puis-je faire avec ces objets?

A: 3455 1144 189 25673

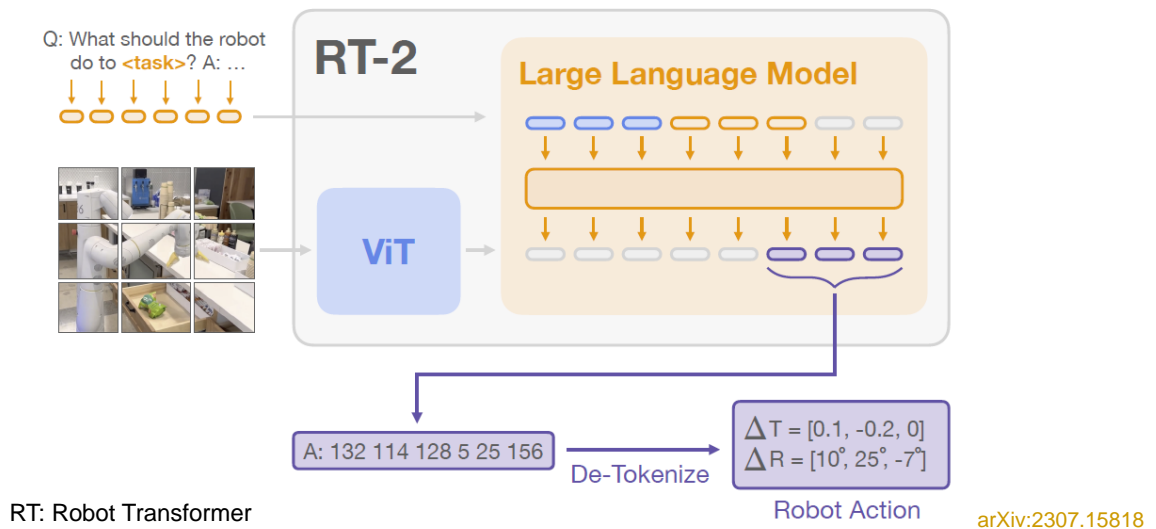
Faire cuire un gâteau.



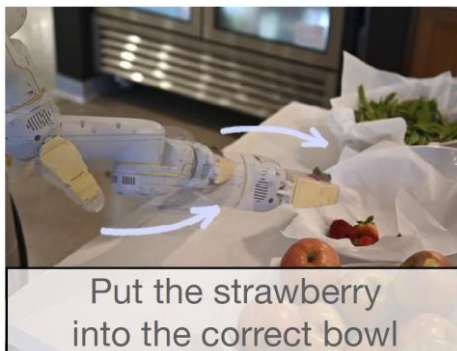
co-fine-tune vision-language models on *robotic trajectory data* and Internet-scale *vision-language tasks*

[arXiv:2307.15818](https://arxiv.org/abs/2307.15818)

# VLA model for robot control (**end-to-end**)



## Closed loop actions



arXiv:2307.15818



# Examples

move coke can  
to Taylor Swift



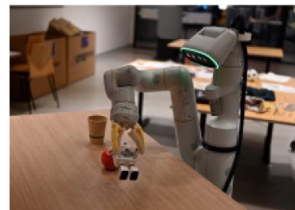
move banana to the  
sum of two plus one



pick animal with  
different color



pick robot



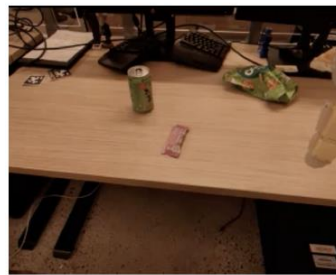
[arXiv:2307.15818](https://arxiv.org/abs/2307.15818)

## Chain-of-thought reasoning

Prompt:  
Given <img> Instruction:  
Move all the objects  
together.  
Prediction:  
Plan: move green can near  
green rice chip bag.  
Action: 1 128 126 127 135  
123 119 127



Prompt:  
Given <img> Instruction:  
Move the green objects  
together.  
Prediction:  
Plan: move green can near  
green rice chip bag.  
Action: 1 130 129 121 131  
127 128 127



[arXiv:2307.15818](https://arxiv.org/abs/2307.15818)

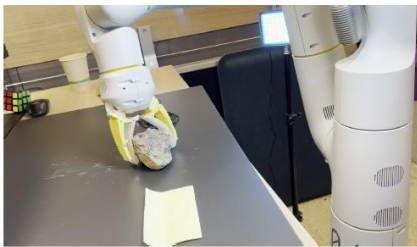
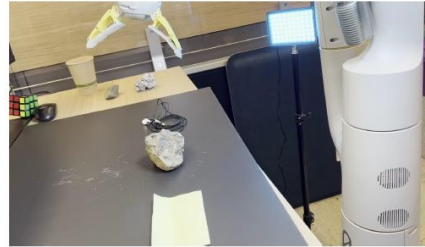
## Chain-of-thought reasoning

Prompt:

Given <img> I need to hammer a nail, what object from the scene might be useful?

Prediction:

Rocks. Action: 1 129 138  
122 132 135 106 127



[arXiv:2307.15818](https://arxiv.org/abs/2307.15818)

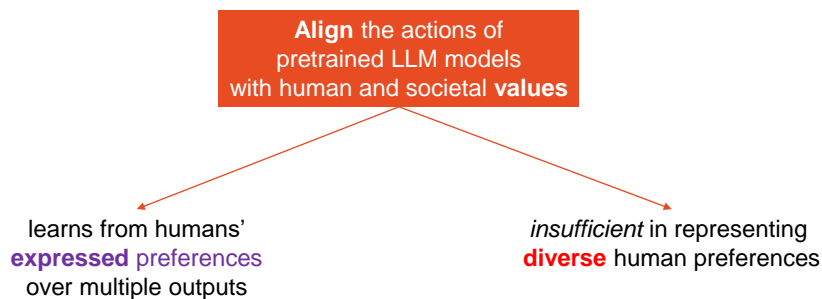
# Reinforcement learning from human feedback

## RLHF (preference data)



[arXiv:2402.08925](https://arxiv.org/abs/2402.08925)

## Fine tuning with RLHF



[arXiv:2404.10271](https://arxiv.org/abs/2404.10271)

# Supervised fine-tuning

$p_\theta$  language model **fine-tuned with supervised learning** for downstream tasks of interest  
(e.g. *dialogue, instruction following, summarization*)

$x_i \in V$   
set of tokens (vocabulary set)

$\mathbf{x} = \{x_1, x_2, \dots, x_N\}$   
prompt (sequence of tokens)

$\mathbf{x} \in X$   
set of prompts

$\mathbf{y} \sim p_\theta(\cdot | \mathbf{x})$  output response

[arXiv:2404.10271](https://arxiv.org/abs/2404.10271)

## Preference data collection

$$\mathbf{y} \sim p_\theta(\cdot | \mathbf{x})$$



generate a dataset of model outputs  
 $\{\mathbf{y}_1, \mathbf{y}_2\} \sim p_\theta(\cdot | \mathbf{x})$   
 pairs of responses



humans evaluate paired completions &  
 select which output they prefer, e.g.  
 $\mathbf{y}_1 \succcurlyeq \mathbf{y}_2$   
 (preferred  $\mathbf{y}_1$  and dis-preferred  $\mathbf{y}_2$  response)

[arXiv:2404.10271](https://arxiv.org/abs/2404.10271)

## Reward model training

$r_\phi(\mathbf{y}, \mathbf{x})$  reward model

$r_\phi: Y \rightarrow \mathbb{R}$

$$P^*(\mathbf{y}_1 \succcurlyeq \mathbf{y}_2 | \mathbf{x}) = \frac{e^{r^*(\mathbf{y}_1, \mathbf{x})}}{e^{r^*(\mathbf{y}_1, \mathbf{x})} + e^{r^*(\mathbf{y}_2, \mathbf{x})}} \quad \text{Bradely-Terry preference model}$$

$D = \{\mathbf{x}^i, \mathbf{y}_1^i, \mathbf{y}_2^i\}_{i=1 \dots N}$  static dataset sampled from  $P^*$

$$L(r_\phi, D) = -\mathbb{E}_{(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2) \sim D} [\log \overset{\text{logistic function}}{\sigma}(r_\phi(\mathbf{y}_1, \mathbf{x}) - r_\phi(\mathbf{y}_2, \mathbf{x}))]$$

$r_\phi$  learned with **max likelihood estimation**

(to match the likelihood of the human preferences observed from the data)

[arXiv:2404.10271](https://arxiv.org/abs/2404.10271)

## Reinforcement learning fine tuning

$p_{r_\phi}^*$  optimal policy under reward  $r_\phi$

KL-regularized reward maximization

$$\max_p E_{\mathbf{x} \sim P, \mathbf{y} \sim p_\theta(\cdot | \mathbf{x})} [r_\phi(\mathbf{y}, \mathbf{x}) - \beta D_{KL}[p(\cdot | \mathbf{x}) || p_{\text{REF}}(\cdot | \mathbf{x})]]$$

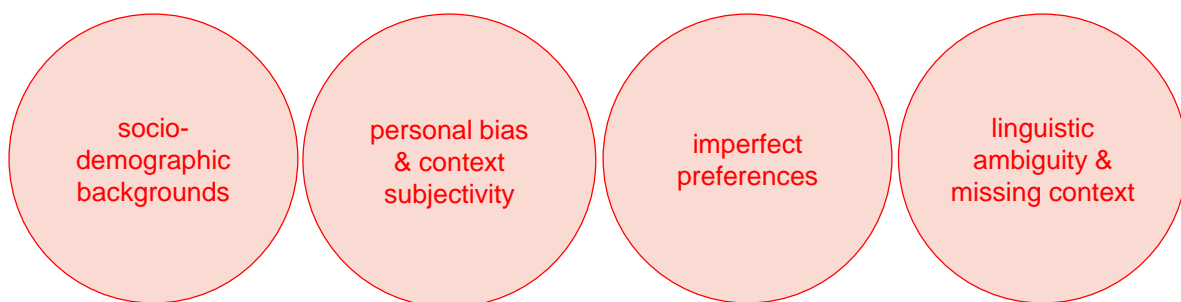
$\beta > 0$

controls the deviation from the  
base reference policy  $p_{\text{REF}}$

[arXiv:2404.10271](https://arxiv.org/abs/2404.10271)

# Pluralistic alignment

Key factors contributing to diversity



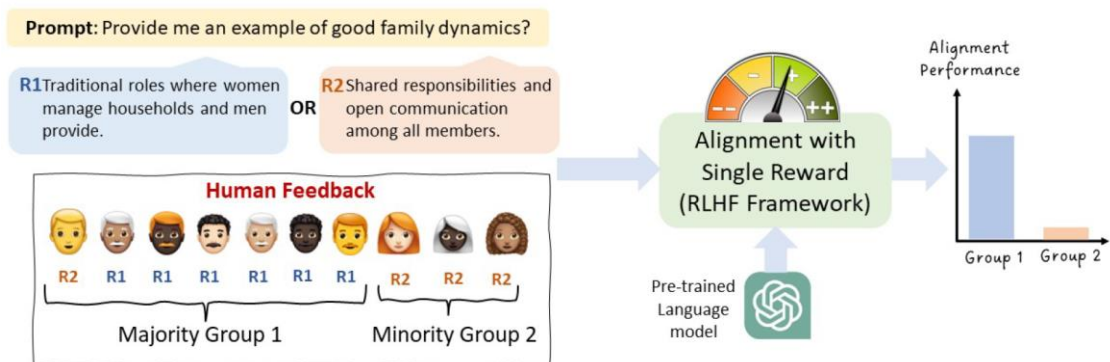
# Reflecting and supporting diversity



algorithmic monocultures  
lead to increased unfairness  
when applied by many decision makers

[arXiv:2402.05070](https://arxiv.org/abs/2402.05070)

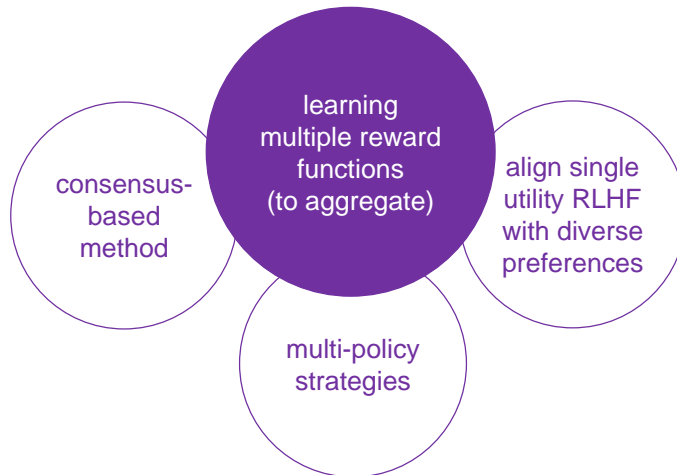
## Majority vs minority user groups



most RLHF approaches ignore diversity in human preference feedback  
by aligning the language model with a **single reward function**

[arXiv:2402.08925](https://arxiv.org/abs/2402.08925)

## Diversity in opinions and preferences



[arXiv:2402.08925](https://arxiv.org/abs/2402.08925)

## Mixture of preference distributions

$$P_u^*(\mathbf{y}_1 \succneq \mathbf{y}_2 | \mathbf{x}) = \mathbb{E}_{h \in H_u} [I(h \text{ prefers } \mathbf{y}_1 \text{ over } \mathbf{y}_2 | \mathbf{x})] \quad \text{for all groups in } U$$

$$U = \{H_1, H_2, \dots, H_{|U|}\} \quad H = \bigcup_{u=1}^{|U|} H_u \quad u \quad \text{human subpopulation index}$$

$$P^*(\mathbf{y}_1 \succneq \mathbf{y}_2 | \mathbf{x}) = \sum_{u=1}^{|U|} \left[ \sum_{h \in H_u} I_h(\mathbf{z}) \underset{\substack{\text{distribution} \\ \text{over the} \\ \text{humans } H}}{q(h|u)} \right] \underset{\substack{\text{marginal probability} \\ \text{distribution of} \\ \text{subpopulation } H_u}}{\eta(u)} = \sum_{u=1}^{|U|} \underset{\substack{\text{subpopulation with specific} \\ \text{preference distribution}}}{p_u^*(\mathbf{z})} \eta(u)$$

[arXiv:2402.08925](https://arxiv.org/abs/2402.08925)



# Mixture of preference distributions

$$p(\mathbf{z}') = \sum_{u=1}^{|U|} p_{\phi_u}^*(\mathbf{z}') \eta(u) \quad \text{preference distribution}$$

$$\mathbf{z}' = (\mathbf{y}_w \succcurlyeq \mathbf{y}_l | \mathbf{x}) \quad \begin{array}{l} \mathbf{y}_w \text{ chosen response by the human sub-population group } H_u \\ \mathbf{y}_l \text{ rejected response by the human sub-population group } H_u \end{array}$$

$$\begin{aligned} L(\phi) &= \sum_{\mathbf{z}' \in D} \log \sum_{u=1}^{|U|} p_{\phi_u}(\mathbf{z}') \eta(u) \\ &= \sum_{\mathbf{z}' \in D} \log \sum_{u=1}^{|U|} \frac{e^{r_{\phi_u}(\mathbf{y}_w, \mathbf{x})}}{e^{r_{\phi_u}(\mathbf{y}_w, \mathbf{x})} + e^{r_{\phi_u}(\mathbf{y}_l, \mathbf{x})}} \eta(u) \end{aligned} \quad \text{maximization of the log likelihood}$$

[arXiv:2402.08925](https://arxiv.org/abs/2402.08925)

# Maximizing the minimum utility

**Alignment objective** (with diverse human preferences)

$$\operatorname{argmax}_{\mathbf{p}} \left( \min_u \mathbb{E}_{\mathbf{x} \sim P, \mathbf{y} \sim p(\cdot | \mathbf{x})} [r_{\phi_u^*}(\mathbf{y}, \mathbf{x})] \right) - \beta D_{KL}[p(\cdot | \mathbf{x}) || p_{\text{REF}}(\cdot | \mathbf{x})]$$

$\phi_u^*$  reward model parameter  
for each human subpopulation in  $U$

[arXiv:2402.08925](https://arxiv.org/abs/2402.08925)

# The implementation

---

## Algorithm 1 MaxMin RLHF

---

- 1: **Input:** Preference dataset  $\mathcal{D}$ , initial reward parametrization for each subpopulation  $u$  as  $r_{\phi_0}^u$ , initial policy parameter  $\pi_0$ .
  - 2: **Reward Learning with EM:** Utilize Algorithm 2 for learning rewards with EM to learn  $r_{\phi}^u$  for all user subpopulation  $u$
  - 3: **Max-Min Policy Iteration:**
  - 4: **for**  $t = 0$  to  $T - 1$  **do**
  - 5:   **Choosing Minimum Utility Subpopulation:**
  - 6:    $u_{\min} \leftarrow \arg \min_{u \in \mathcal{U}} F_{r_{\phi}^u}(\pi_t)$
  - 7:   **Perform the PPO Update:**
  - 8:   Update policy  $\pi$  towards maximizing the objective:
  - 9:    $\pi_{t+1} \leftarrow \text{PPO-update}(F_{r_{\phi_{u_{\min}}}}(\pi_t) - \beta \mathbb{D}_{\text{KL}}[\pi_t || \pi_{\text{ref}}])$
  - 10: **end for**
  - 11: **Output:** Policy  $\pi_T$  aligned with socially fair preference dataset
- 

---

## Algorithm 2 Learning Rewards with EM Algorithm

---

- 1: **Input:** Preference data  $\mathcal{D}$ ,  $|\mathcal{U}|$  clusters of users among all humans in  $\mathcal{H} = \bigcup_{u=1}^{|\mathcal{U}|} \mathcal{H}_u$ , pretrained  $\{r_{\phi_u}\}_{u=1}^{|\mathcal{U}|}$ , loss function loss, convergence criteria
  - 2: **while** not reach the convergence criteria **do**
  - 3:   **for**  $h \in \mathcal{H}$  **do**
  - 4:     **E-step (hard cluster assignment):** assign  $h$  to the  $u$ -th cluster s.t.
 
$$u = \arg \max_{u \in 1, \dots, |\mathcal{U}|} \prod_{(\mathbf{x}, \mathbf{y}_1, \mathbf{y}_2, h) \in \mathcal{D}} w(\phi_u, \mathbf{x}, \mathbf{y}_1, \mathbf{y}_2)$$

$$\text{where } w(\cdot) = \frac{\exp(r_{\phi_u}(\mathbf{y}_1, \mathbf{x}))}{\exp(r_{\phi_u}(\mathbf{y}_1, \mathbf{x})) + \exp(r_{\phi_u}(\mathbf{y}_2, \mathbf{x}))}$$
  - 5:   **end for**
  - 6:   **M-step:** Update each  $\phi_u, u = 1, \dots, |\mathcal{U}|$  by minimizing the negative log-likelihood loss (2) on the assigned users' data
  - 7: **end while**
- 

[arXiv:2402.08925](https://arxiv.org/abs/2402.08925)

# What did we learn today?

- Flamingo
- Vision-Language-Action models
- Reinforcement learning from human feedback
- Pluralistic alignment

# EE-559

# Deep Learning

[andrea.cavallaro@epfl.ch](mailto:andrea.cavallaro@epfl.ch)