

Last Time:

□ game theory

lecture 12

4RI, FALL '25

Andrea Bajcsy

This Time:

□ alignment

## What is AI Alignment?

"alignment"

You may have heard of this term ✓ after discussing new "foundation models" like ChatGPT / Gemini or possible "robotics foundation models", but its conception dates as far back as 1960 when AI pioneer Norbert Wiener described AI alignment as:

"If we use, to achieve our purposes, a mechanical agency with whose operation we cannot interfere effectively ... we had better be quite sure that the purpose put into the machine is the purpose we really desire."

- Wiener, 1960. "Some Moral and Technical Consequences of Automation."



Value alignment is the process of developing & deploying AI systems in a way that aligns with human values & goals.

! value alignment is fundamentally a multi-agent problem b/w the AI/robot and the human - who determines what the objective is.

Famous examples:

- ① Coast Runners Game (Amodei & Clark, 2016)
- ② Lego stacking Manipulator (Popov et al, 2017)
- ③ Deceptive dexterous hand (Christians et al, 2017)
- ④ Exploiting simulation bugs (Code Bullet, 2019)

# How do we align on AI system?

There is no one algorithm or tool that can "solve" alignment, but one popular paradigm is called **reinforcement learning from human feedback** (RLHF)

↳ key component of current LMs (e.g. ChatGPT/Claude/Bard)

see: Cao et al. "RLHF for Realistic Traffic Sim." ICRA 2024

↳ Current Autonomous Driving predictors/planners

↳ next-generation generative robotics policies

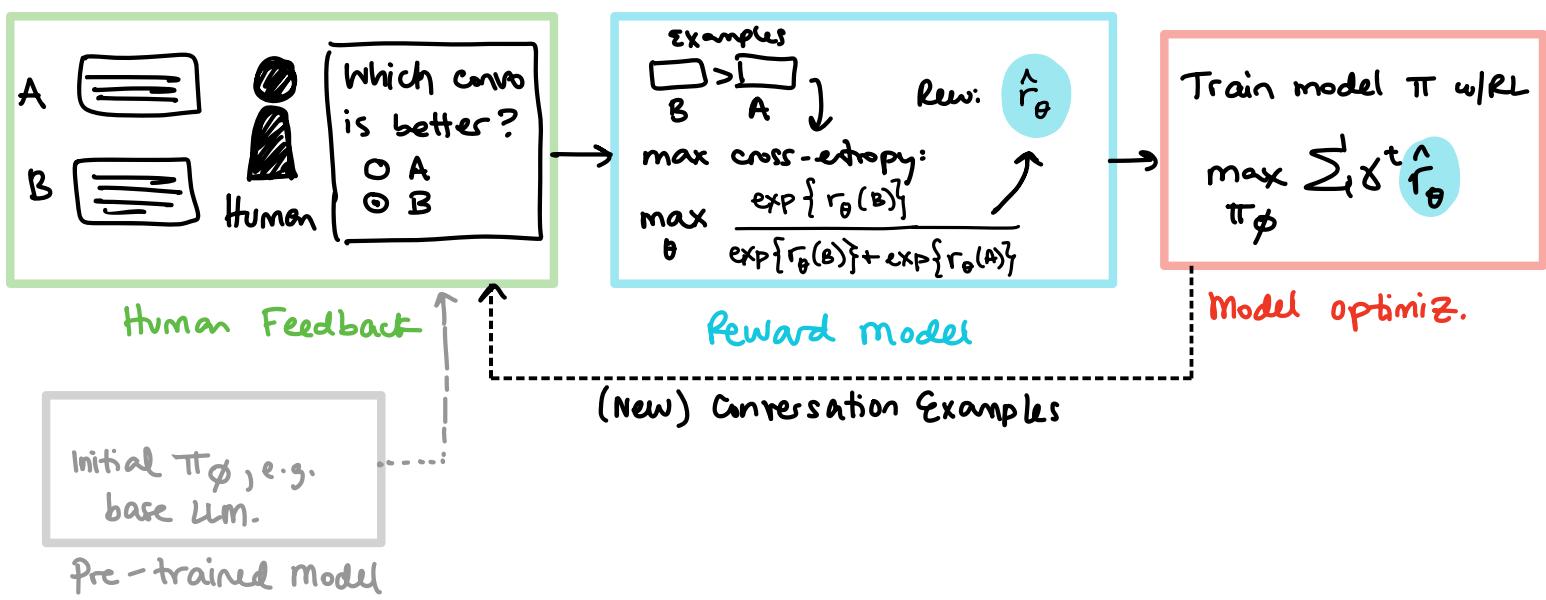
↳ see: Tian et al. "Maximizing Alignment with Minimal Feedback."

BUT, has roots in preference theory from economics w/ early applications in HCI and RL.

RLHF has three (repeated) steps:

- ① feed back collection
- ② reward modeling
- ③ model optimization

example: LM chatbot RLHF w/ Binary Preference Feedback



Mathematically, the first step involves collecting examples from the "base model":

$$x_i \sim \pi_\theta$$

ex. a batch of 1+ generations from the model, like a complete convo., or a denoised action chunk in robotics FMs.

ex. LM models:  $P(\vec{x}_i | \vec{x}_{<i})$   
next token pred. given prev. generations.

let the human  $H$  have desires consistent with the unknown reward  $r_H$ . Let the feedback we get from them be modelled by

$$y_i = f(H, x_i, \epsilon_i)$$

↑  
human ↑ noise to feedback  $y_i$

example  
generations (e.g.  $x_i = (\text{convo A}, \text{convo B})$ )

Next, we fit a reward model  $\hat{r}_\theta$  with the feedback data:  $\{(x_i, y_i)\}_{i=1}^N$  to approximate evaluations from  $H$  as closely as possible:

$$\hat{r}_{\theta^*} = \underset{\theta}{\text{minimize}} \mathcal{L}(D, \theta) = \sum_{i=1}^N l(\hat{r}_\theta(x_i), y_i) + \lambda_r(\theta)$$

"suitable loss for this data"      regularizer

ex. if  $x_i = (\text{convo A}, \text{convo B})$

$y_i$  = distribution over  $\{A, B\}$  indicating which user preferred

$$= \mu$$

$$A > B \Rightarrow \mu := \frac{1.0}{A+B}$$

;

$$B > A \Rightarrow \mu := \frac{0.0}{A+B}$$

$$A \approx B \Rightarrow \mu := \frac{0.5}{A} \quad \frac{0.5}{B}$$

$A \neq B \Rightarrow y_i \notin D$  ↑  
"incomparable" comparison discarded

We will model human judgements of preferring a convo as exponentially more likely the higher the reward:

$$P(A > B | \theta) = \frac{\exp\{\hat{r}_\theta(A)\}}{\exp\{\hat{r}_\theta(A)\} + \exp\{\hat{r}_\theta(B)\}} \quad (\textcircled{*})$$

! Similar to the Boltzmann Rationality model, but now its over pairwise preferences rather than state traj. or actions, etc.

Formally  $\textcircled{*}$  is called the Bradley-Terry Model (B: T, 1952) and its a specialization of the Luce-Shephard choice rule (Luce 2005, Shephard 1957) to preferences over sequences/ trajectories.

We can now choose a loss function with this probabilistic model:

$$l(\hat{r}_\theta(x_i), y_i) := - \left[ \underbrace{\mu(A) \log P(A > B | \theta)}_{\text{if } A > B} + \underbrace{\mu(B) \log P(B > A | \theta)}_{\text{if } B > A} \right]$$

↑ from our dist  $\prod_{i=1}^n \frac{o_i}{A_i} \cdot \frac{1-o_i}{B_i}$

Finally, we optimize the model with RL! Optimize  $\pi_\phi$  model params.

$$\phi_{\text{new}}^* = \max_{\phi_{\text{new}}} \mathbb{E}_{x \sim \pi_{\phi_{\text{new}}}} \left[ \hat{r}_\theta^*(x) + \underbrace{\gamma_p(\phi, \phi_{\text{new}}, x)}_{\text{regularizer like divergence b/wn. prior distribution } \pi_\phi \text{ & the new } \pi_{\phi_{\text{new}}}} \right]$$

regularizer like divergence b/wn.  
prior distribution  $\pi_\phi$  & the new  $\pi_{\phi_{\text{new}}}$

$$D_{KL}(\pi_{\phi_{\text{new}}} || \pi_\phi)$$

# Challenges :

## • Human Data

→ Humans struggle to evaluate difficult tasks well

### Sample A

$$\begin{aligned} \mathcal{L}_{SM} = & -\frac{1}{2}\partial_\mu g_\mu^a \partial_\nu g_\mu^a - g_s f^{abc} \partial_\mu g_\mu^a g_\mu^b g_\mu^c - \frac{1}{4}g_\mu^2 f^{abc} f^{ade} g_\mu^b g_\mu^d g_\mu^e - \partial_\mu W_\mu^+ \partial_\nu W_\mu^- \\ & - M^2 W_\mu^+ W_\mu^- - Z_\mu^0 (W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) - Z_\mu^0 (W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) - \\ & - Z_\mu^0 (W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) - Z_\mu^0 (W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) - \\ & - ig s_w (\partial_\mu A_\mu) (W_\mu^+ W_\mu^- - W_\mu^- W_\mu^+) - A_\mu (W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) + A_\mu (W_\mu^+ \partial_\nu W_\mu^- - \\ & - W_\mu^- \partial_\nu W_\mu^+) - \frac{1}{2}g^2 W_\mu^+ W_\mu^- W_\mu^+ W_\mu^- + \frac{1}{2}g^2 W_\mu^+ W_\mu^- W_\mu^+ W_\mu^- + g^2 c_w^2 (Z_\mu^0)^2 Z_\mu^0 W_\mu^- \\ & - Z_\mu^0 Z_\mu^0 (W_\mu^+ W_\mu^-) + g^2 s_w^2 (A_\mu W_\mu^+ A_\mu W_\mu^- - A_\mu A_\mu) W_\mu^+ W_\mu^- + g^2 s_w c_w (A_\mu Z_\mu^0) (W_\mu^+ W_\mu^- - \\ & - W_\mu^- W_\mu^+) - 2A_\mu Z_\mu^0 (W_\mu^+ W_\mu^-) - \frac{1}{2}\partial_\mu H \partial_\mu H - 2M^2 \alpha_h H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - \frac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \\ & - W_\mu^+ W_\mu^- - 2A_\mu Z_\mu^0 (W_\mu^+ W_\mu^-) - \frac{1}{2}\partial_\mu H \partial_\mu H - 2M^2 \alpha_h H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - \frac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \\ & \beta_h \left( \frac{M^2}{g^2} + \frac{2M}{g} H + \frac{1}{2}(H^2 + \phi^0 \phi^0 + 2\phi^+ \phi^-) \right) + \frac{2M}{g^2} \alpha_h - \\ & \frac{1}{2}g^2 \alpha_h (H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4(H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2) - \\ & g_{\alpha_h} M (H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4(H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2) - \\ & \frac{1}{2}g^2 \alpha_h (H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4(H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2) - \\ & g_{\alpha_h} M (H^4 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\ & \frac{1}{2}g^2 \alpha_h (H^4 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\ & g_{\alpha_h} M (H^4 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\ & \frac{1}{2}g^2 \alpha_h (H^4 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) + \\ & \frac{1}{2}g (W_\mu^+ (\phi^0 \partial_\mu \phi^- - \phi^+ \partial_\mu \phi^0) - W_\mu^- (\phi^0 \partial_\mu \phi^+ - \phi^+ \partial_\mu \phi^0)) + \\ & \frac{1}{2}g (W_\mu^+ (H \partial_\mu \phi^- - \phi^+ \partial_\mu H) + W_\mu^- (H \partial_\mu \phi^+ - \phi^+ \partial_\mu H)) + \frac{1}{2}g \frac{1}{c_w} (Z_\mu^0 (\partial_\mu \phi^0 - \phi^0 \partial_\mu H) + \\ & M \left( \frac{1}{2}Z_\mu^0 \partial_\mu \phi^0 + W_\mu^+ \partial_\mu \phi^+ + W_\mu^- \partial_\mu \phi^- \right) - ig \frac{2}{c_w} M Z_\mu^0 (W_\mu^+ \phi^+ - W_\mu^- \phi^-) + ig s_w M A_\mu (W_\mu^+ \phi^- - \\ & - W_\mu^- \phi^+) - ig \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 (\phi^+ \partial_\mu \phi^- - W_\mu^- \phi^+ \partial_\mu \phi^-) + ig s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \\ & \frac{1}{2}g^2 \frac{s_w^2}{c_w^2} Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2}ig \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - \frac{1}{2}g^2 s_w A_\mu \phi^0 (W_\mu^+ \phi^- + \\ & - W_\mu^- \phi^+) + \frac{1}{2}g^2 s_w A_\mu H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - g^2 \frac{1-2s_w^2}{2c_w^2} (2c_w^2 - 1) \phi^+ \phi^- - \\ & \frac{1}{2}g^2 s_w^2 Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2}ig \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - \frac{1}{2}g^2 s_w A_\mu (W_\mu^+ \phi^- - \phi^- \partial_\mu \phi^+) - \\ & \frac{1}{2}g^2 W_\mu^+ W_\mu^- (H^2 + (\phi^0)^2 + 2\phi^+ \phi^-) - \frac{1}{8}g^2 \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 (H^2 + (\phi^0)^2 + 2(2s_w^2 - 1)^2 \phi^+ \phi^-) - \\ & \frac{1}{2}g^2 s_w^2 Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2}ig \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - \frac{1}{2}g^2 s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \\ & g^2 s_w^2 A_\mu A_\mu \phi^0 \phi^+ + \frac{1}{2}ig s_w A_\mu (\eta^{\mu\nu} \eta^{\rho\lambda} \eta^{\sigma\rho} \eta^{\tau\lambda}) - \bar{\psi}^\lambda (\gamma^\mu + m_\lambda^2) \nu^\mu - \bar{u}_j^\lambda (\gamma^\mu + \\ & m_\lambda^2) u_j^\mu - \bar{d}_j^\lambda (\gamma^\mu + m_\lambda^2) d_j^\mu + ig s_w A_\mu (-(\bar{\psi}^\lambda \gamma^\mu \nu^\mu + \frac{1}{3}(\bar{u}_j^\lambda \gamma^\mu u_j^\mu) - \frac{1}{3}(\bar{d}_j^\lambda \gamma^\mu d_j^\mu))) + \\ & (ig s_w A_\mu \{(\bar{\psi}^\lambda \gamma^\mu (1 - \gamma^2) \nu^\mu) + (\bar{\psi}^\lambda \gamma^\mu (4s_w^2 - 1 - \gamma^2) e^\mu) + (\bar{d}_j^\lambda \gamma^\mu (1 - \gamma^2) e^\mu) + \\ & (u_j^\lambda \gamma^\mu (1 - \gamma^2) u_j^\mu) + (\bar{u}_j^\lambda \gamma^\mu (1 + \gamma^2) \nu^\mu) + (\bar{d}_j^\lambda \gamma^\mu (1 + \gamma^2) e^\mu)\}) + \\ & \frac{ig}{2\sqrt{2}} W_\mu^- \left( (\bar{\psi}^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) \nu^\mu) + (\bar{d}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) u_j^\mu) \right) + \\ & \frac{ig}{2\sqrt{2}} W_\mu^+ \left( -m_w^2 (\bar{\psi}^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) \nu^\mu) + (\bar{d}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) e^\mu) \right) + \\ & \frac{ig}{2M\sqrt{2}} \phi^0 \left( m_w^2 (\bar{\psi}^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) \nu^\mu) - m_w^2 (\bar{\psi}^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) \nu^\mu) - \frac{g}{2} \frac{m_w^2}{M} H (\bar{\psi}^\lambda \nu^\lambda) - \right. \\ & \left. - \frac{g}{2} \frac{m_w^2}{M} H (\bar{\psi}^\lambda e^\lambda) + \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{\psi}^\lambda \gamma^\lambda e^\lambda) - \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{\psi}^\lambda \gamma^\lambda \nu^\lambda) - \frac{1}{4} \bar{\psi}_\lambda M_{\lambda\mu}^{\nu} (1 - \gamma^2) \bar{\nu}_\mu - \right. \\ & \left. - \frac{1}{4} \bar{\psi}_\lambda M_{\lambda\mu}^{\nu} (1 - \gamma^2) \bar{\nu}_\mu + \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{u}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) d_j^\mu) + m_w^2 (\bar{u}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) d_j^\mu) + \right. \\ & \left. + \frac{ig}{2\sqrt{2}} \phi^0 \left( m_w^2 (\bar{d}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) u_j^\mu) - \frac{ig}{2\sqrt{2}} \phi^0 (\bar{u}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) u_j^\mu) \right) - \frac{g}{2} \frac{m_w^2}{M} H (\bar{u}_j^\lambda u_j^\mu) - \right. \\ & \left. - \frac{g}{2} \frac{m_w^2}{M} H (\bar{d}_j^\lambda d_j^\mu) + \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{u}_j^\lambda \gamma^\mu \bar{u}_j^\mu) - \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{d}_j^\lambda \gamma^\mu \bar{d}_j^\mu) + \bar{G}^a \partial^\mu G^a + g_s \bar{\psi}^\mu \partial_\mu \bar{G}^a G^a + \right. \\ & \left. + \bar{X}^+ (\partial^a - M^2) X^+ + \bar{X}^-(\partial^a - M^2) X^- + \bar{X}^+(\partial^a - \frac{M^2}{c_w^2}) X^0 + \bar{Y} \partial^a Y + ig c_w W_\mu^+ (\partial_\mu \bar{X}^0 X^+ - \right. \\ & \left. - \partial_\mu \bar{X}^0 X^-) + ig s_w W_\mu^+ (\partial_\mu \bar{X}^+ Y - \partial_\mu \bar{Y} X^+) + ig c_w W_\mu^- (\partial_\mu \bar{X}^0 X^- - \right. \\ & \left. - \partial_\mu \bar{X}^0 X^+) + ig s_w W_\mu^- (\partial_\mu \bar{X}^- Y - \partial_\mu \bar{Y} X^+) + ig c_w Z_\mu^0 (\partial_\mu \bar{X}^0 X^+ - \right. \\ & \left. - \partial_\mu \bar{X}^- X^-) + ig s_w A_\mu (\partial_\mu \bar{X}^+ X^- - \right. \\ & \left. - \partial_\mu \bar{X}^- X^-) - \frac{1}{2}g M (X^+ X^+ H + \bar{X}^- X^- H + \frac{1}{c_w^2} \bar{X}^0 X^0 H) + \frac{1-2s_w^2}{2c_w^2} ig M (X^+ X^0 \phi^+ - \bar{X}^- X^0 \phi^-) + \right. \\ & \left. + \frac{1}{2c_w} ig M (X^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-) + ig M s_w (X^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-) + \right. \\ & \left. + \frac{1}{2}ig M (X^+ X^+ \phi^0 - \bar{X}^- X^- \phi^0) . \right) \end{aligned}$$

### Sample B

$$\begin{aligned} \mathcal{L}_{SM} = & -\frac{1}{2}\partial_\mu g_\mu^a \partial_\nu g_\mu^a - g_\mu^a f^{abc} \partial_\mu g_\mu^b g_\mu^c g_\mu^a - \frac{1}{4}g_\mu^2 f^{abc} f^{ade} g_\mu^b g_\mu^d g_\mu^e - \partial_\mu W_\mu^+ \partial_\nu W_\mu^- - \\ & M^2 W_\mu^+ W_\mu^- - \frac{1}{2}g_\mu^a Z_\mu^0 \partial_\mu Z_\mu^0 - \frac{1}{2c_w^2} M^2 Z_\mu^0 Z_\mu^0 - \frac{1}{2}g_\mu^a A_\mu \partial_\mu A_\mu - ig c_w (\partial_\mu Z_\mu^0) (W_\mu^+ W_\mu^- - \\ & W_\mu^- W_\mu^+) - Z_\mu^0 (W_\mu^+ \partial_\mu W_\mu^- - W_\mu^- \partial_\mu W_\mu^+) - Z_\mu^0 (W_\mu^+ \partial_\mu W_\mu^- - W_\mu^- \partial_\mu W_\mu^+) - \\ & ig s_w (\partial_\mu A_\mu) (W_\mu^+ W_\mu^- - W_\mu^- W_\mu^+) - A_\mu (W_\mu^+ \partial_\mu W_\mu^- - W_\mu^- \partial_\mu W_\mu^+) + A_\mu (W_\mu^+ \partial_\mu W_\mu^- - \\ & W_\mu^- \partial_\mu W_\mu^+) - \frac{1}{2}g^2 W_\mu^+ W_\mu^- W_\mu^+ W_\mu^- + \frac{1}{2}g^2 W_\mu^+ W_\mu^- W_\mu^+ W_\mu^- + g^2 c_w^2 (Z_\mu^0)^2 Z_\mu^0 W_\mu^- - \\ & Z_\mu^0 Z_\mu^0 (W_\mu^+ W_\mu^-) + g^2 s_w^2 (A_\mu W_\mu^+ A_\mu W_\mu^- - A_\mu A_\mu) W_\mu^+ W_\mu^- + g^2 s_w c_w (A_\mu Z_\mu^0) (W_\mu^+ W_\mu^- - \\ & W_\mu^- W_\mu^+) - 2A_\mu Z_\mu^0 (W_\mu^+ W_\mu^-) - \frac{1}{2}\partial_\mu H \partial_\mu H - 2M^2 \alpha_h H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - \frac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \\ & W_\mu^+ W_\mu^- - 2A_\mu Z_\mu^0 (W_\mu^+ W_\mu^-) - \frac{1}{2}\partial_\mu H \partial_\mu H - 2M^2 \alpha_h H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - \frac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \\ & \beta_h \left( \frac{M^2}{g^2} + \frac{2M}{g} H + \frac{1}{2}(H^2 + \phi^0 \phi^0 + 2\phi^+ \phi^-) \right) + \frac{2M}{g^2} \alpha_h - \\ & \frac{1}{2}g^2 \alpha_h (H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4(H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2) - \\ & g_{\alpha_h} M (H^4 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\ & \frac{1}{2}g^2 \alpha_h (H^4 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\ & g_{\alpha_h} M (H^4 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\ & \frac{1}{2}g^2 \alpha_h (H^4 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) + \\ & \frac{1}{2}g (W_\mu^+ (\phi^0 \partial_\mu \phi^- - \phi^+ \partial_\mu \phi^0) - W_\mu^- (\phi^0 \partial_\mu \phi^+ - \phi^+ \partial_\mu \phi^0)) + \\ & \frac{1}{2}g (W_\mu^+ (H \partial_\mu \phi^- - \phi^+ \partial_\mu H) + W_\mu^- (H \partial_\mu \phi^+ - \phi^+ \partial_\mu H)) + \frac{1}{2}g \frac{1}{c_w} (Z_\mu^0 (\partial_\mu \phi^0 - \phi^0 \partial_\mu H) + \\ & M \left( \frac{1}{2}Z_\mu^0 \partial_\mu \phi^0 + W_\mu^+ \partial_\mu \phi^+ + W_\mu^- \partial_\mu \phi^- \right) - ig \frac{2}{c_w} M Z_\mu^0 (W_\mu^+ \phi^+ - W_\mu^- \phi^-) + ig s_w M A_\mu (W_\mu^+ \phi^- - \\ & - W_\mu^- \phi^+) - ig \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 (\phi^+ \partial_\mu \phi^- - W_\mu^- \phi^+ \partial_\mu \phi^-) + ig s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \\ & \frac{1}{2}g^2 \frac{s_w^2}{c_w^2} Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2}ig \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - \frac{1}{2}g^2 s_w A_\mu \phi^0 (W_\mu^+ \phi^- + \\ & - W_\mu^- \phi^+) + \frac{1}{2}g^2 s_w A_\mu H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - g^2 \frac{1-2s_w^2}{2c_w^2} (2c_w^2 - 1) \phi^+ \phi^- - \\ & \frac{1}{2}g^2 s_w^2 Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2}ig \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - \frac{1}{2}g^2 s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \\ & \frac{1}{2}g^2 s_w^2 Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2}ig \frac{1-2s_w^2}{2c_w^2} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - \frac{1}{2}g^2 s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \\ & g^2 s_w^2 A_\mu A_\mu \phi^0 \phi^+ + \frac{1}{2}ig s_w A_\mu (\eta^{\mu\nu} \eta^{\rho\lambda} \eta^{\sigma\rho} \eta^{\tau\lambda}) - \bar{\psi}^\lambda (\gamma^\mu + m_\lambda^2) \nu^\mu - \bar{u}_j^\lambda (\gamma^\mu + \\ & m_\lambda^2) u_j^\mu - \bar{d}_j^\lambda (\gamma^\mu + m_\lambda^2) d_j^\mu + ig s_w A_\mu (-(\bar{\psi}^\lambda \gamma^\mu \nu^\mu + \frac{1}{3}(\bar{u}_j^\lambda \gamma^\mu u_j^\mu) - \frac{1}{3}(\bar{d}_j^\lambda \gamma^\mu d_j^\mu))) + \\ & (ig s_w A_\mu \{(\bar{\psi}^\lambda \gamma^\mu (1 - \gamma^2) \nu^\mu) + (\bar{\psi}^\lambda \gamma^\mu (4s_w^2 - 1 - \gamma^2) e^\mu) + (\bar{d}_j^\lambda \gamma^\mu (1 - \gamma^2) e^\mu) + \\ & (u_j^\lambda \gamma^\mu (1 - \gamma^2) u_j^\mu) + (\bar{u}_j^\lambda \gamma^\mu (1 + \gamma^2) \nu^\mu) + (\bar{d}_j^\lambda \gamma^\mu (1 + \gamma^2) e^\mu)\}) + \\ & \frac{ig}{2\sqrt{2}} W_\mu^- \left( (\bar{\psi}^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) \nu^\mu) + (\bar{d}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) u_j^\mu) \right) + \\ & \frac{ig}{2\sqrt{2}} W_\mu^+ \left( -m_w^2 (\bar{\psi}^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) \nu^\mu) + (\bar{d}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) e^\mu) \right) + \\ & \frac{ig}{2M\sqrt{2}} \phi^0 \left( m_w^2 (\bar{\psi}^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) \nu^\mu) - m_w^2 (\bar{\psi}^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) \nu^\mu) - \frac{g}{2} \frac{m_w^2}{M} H (\bar{\psi}^\lambda \nu^\lambda) - \right. \\ & \left. - \frac{g}{2} \frac{m_w^2}{M} H (\bar{\psi}^\lambda e^\lambda) + \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{\psi}^\lambda \gamma^\lambda e^\lambda) - \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{\psi}^\lambda \gamma^\lambda \nu^\lambda) - \frac{1}{4} \bar{\psi}_\lambda M_{\lambda\mu}^{\nu} (1 - \gamma^2) \bar{\nu}_\mu - \right. \\ & \left. - \frac{1}{4} \bar{\psi}_\lambda M_{\lambda\mu}^{\nu} (1 - \gamma^2) \bar{\nu}_\mu + \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{u}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) d_j^\mu) + m_w^2 (\bar{u}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) d_j^\mu) + \right. \\ & \left. + \frac{ig}{2\sqrt{2}} \phi^0 \left( m_w^2 (\bar{d}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 + \gamma^2) u_j^\mu) - \frac{ig}{2\sqrt{2}} \phi^0 (\bar{u}_j^\lambda U_{\lambda\mu}^{1p} \gamma^\mu (1 - \gamma^2) u_j^\mu) \right) - \frac{g}{2} \frac{m_w^2}{M} H (\bar{u}_j^\lambda u_j^\mu) - \right. \\ & \left. - \frac{g}{2} \frac{m_w^2}{M} H (\bar{d}_j^\lambda d_j^\mu) + \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{u}_j^\lambda \gamma^\mu \bar{u}_j^\mu) - \frac{ig}{2} \frac{m_w^2}{M} \phi^0 (\bar{d}_j^\lambda \gamma^\mu \bar{d}_j^\mu) + \bar{G}^a \partial^\mu G^a + g_s \bar{\psi}^\mu \partial_\mu \bar{G}^a G^a + \right. \\ & \left. + \bar{X}^+ (\partial^a - M^2) X^+ + \bar{X}^-(\partial^a - M^2) X^- + \bar{X}^+(\partial^a - \frac{M^2}{c_w^2}) X^0 + \bar{Y} \partial^a Y + ig c_w W_\mu^+ (\partial_\mu \bar{X}^0 X^+ - \right. \\ & \left. - \partial_\mu \bar{X}^0 X^-) + ig s_w W_\mu^+ (\partial_\mu \bar{X}^+ Y - \partial_\mu \bar{Y} X^+) + ig c_w W_\mu^- (\partial_\mu \bar{X}^0 X^- - \right. \\ & \left. - \partial_\mu \bar{X}^0 X^+) + ig s_w W_\mu^- (\partial_\mu \bar{X}^- Y - \partial_\mu \bar{Y} X^+) + ig c_w Z_\mu^0 (\partial_\mu \bar{X}^0 X^+ - \right. \\ & \left. - \partial_\mu \bar{X}^- X^-) - \frac{1}{2}g M (X^+ X^+ H + \bar{X}^- X^- H + \frac{1}{c_w^2} \bar{X}^0 X^0 H) + \frac{1-2s_w^2}{2c_w^2} ig M (X^+ X^0 \phi^+ - \bar{X}^- X^0 \phi^-) + \right. \\ & \left. + \frac{1}{2c_w} ig M (X^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-) + ig M s_w (X^0 X^- \phi^+ - \bar{X}^0 X^+ \phi^-) + \right. \\ & \left. + \frac{1}{2}ig M (X^+ X^+ \phi^0 - \bar{X}^- X^- \phi^0) . \right) \end{aligned}$$

The "Standard model Lagrangian" of particle physics

Q which is more "correct"?

→ Human raters (across bgs, countries etc.) disagree

Canada		
YES	NO	DNK
40%	40%	20%

India		
YES	NO	DNK
70%	30%	0

USA		
YES	NO	DNK
50%	0	50%



adversarial example from the CATS4ML data challenge  
<https://github.com/google-research-datasets/cats4ml-dataset>

from: Lora Aroyo, 2023. "The Many Faces of Responsible AI",  
 Neuips Keynote

→ RLHF suffers from tradeoff btwn. the richness vs. efficiency of the feedback types.

↳ e.g. Comparisons:  $y_i := A > B$

Scalar:

$y_i = 0.35 \in \mathbb{R}$ . ↗ more expressive!  
but poorly calib.

Label:

$y_i \in \{y_1, y_2, \dots, y_M\}$  ↗ low effort  
but can suffer from choice set misspecification

Correction:

$y_i = \Delta x_i$  ↗ higher effort

Language:

$y_i = \text{Z}$  ↗ utterance is easy, but imprecision of speech + cross-cultural diffs. make it hard  
with incomplete labels.

⋮  
more??

## • Reward learning

→ an individual's values are difficult to represent with a reward function

↳ e.g. human feedback can depend on contextual factors not easily accounted for in the training data (time-varying rewards, pedagogic behavior, etc.)

→ a single reward cannot represent a diverse society of humans (e.g. Smile example from above)

↳ can end up disadvantaging certain groups.

→ reward models can misgeneralize to poor reward proxies, even from correctly-labeled training data

↳ learned models are prone to causal confusion and poor out-of-distribution generalization

## • Model optimization:

→ It is still hard to optimize models / do RL!

→ Policies can perform poorly @ deployment even if rewards seen @ training were perfectly correct.