



Wenxuan Zhou

Shujian Zhang

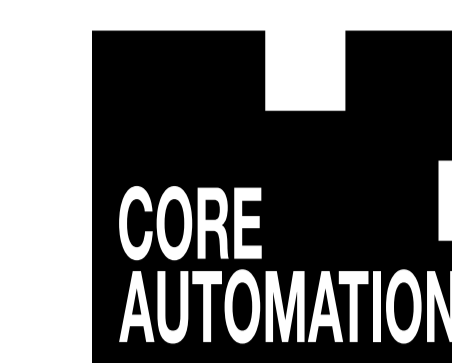
Brice Magdalou

John Lambert

Ehsan Amid

Richard Nock

Andrew Hard

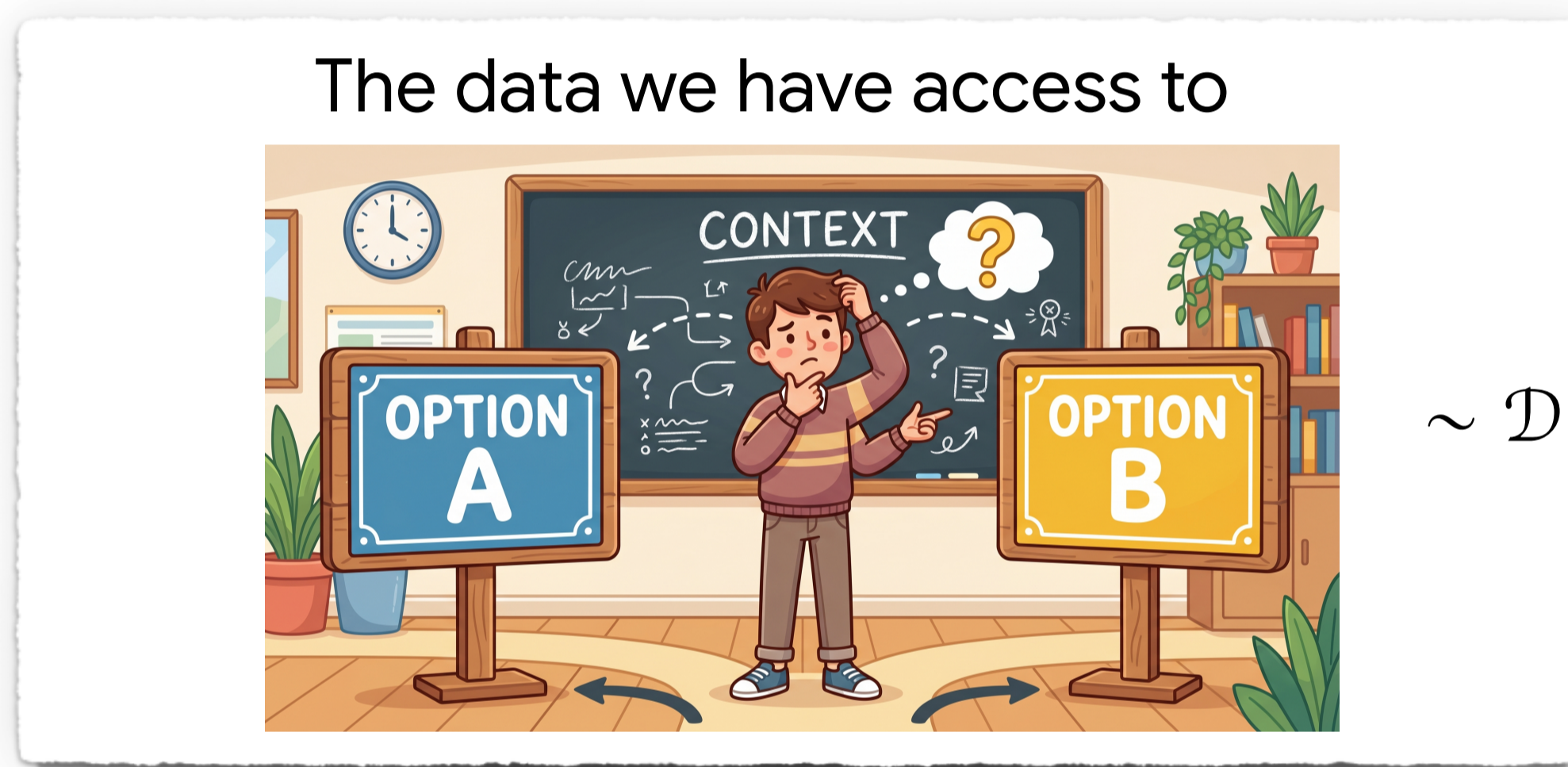


## Summary

- Normative properties are guardrails, often hidden, safeguarding training (e.g. losses)
- Direct Preference Optimization (DPO) simplified Reinforcement Learning from Human Feedback (RLHF) by exploiting a model of Human Choice with ~ 0 leeway
- Followers either stayed close to DPO or abandoned normative safeguards
- We show how go "far" away while keeping normative guardrails, + surprises as bonus



states / context / prompts:  $x \in \mathcal{X}$ ; actions / alternatives / answers:  $y \in \mathcal{Y}$ ; policy:  $\pi : \mathcal{X} \rightarrow \Delta_{\mathcal{Y}}$  simplex



From the data, can learn choice probabilities  $\rightarrow$  loss

$$I = \mathbb{E}_{(x,y,y') \sim \mathcal{D}} [\ell(p(y > y'|x))] \quad \textcircled{1}$$

Model choice probabilities using rewards

$$p(y > y'|x) = \textcircled{2} f(r(x,y), r(x,y'))$$

RLHF DPO

Learn rewards vs policy Model rewards using policy

$$r_{\omega}(x,y) \quad \textcircled{3} \quad \pi_{\theta}(y|x) \quad \textcircled{3} \quad r(x,y) = h(\pi_{\theta}(y|x))$$

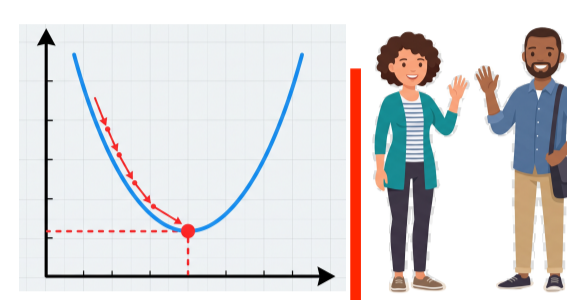
The output we seek: a policy

## DPO's Normative components 101

- ①, ②, ③ carve normative components from two theories: loss function & human choice

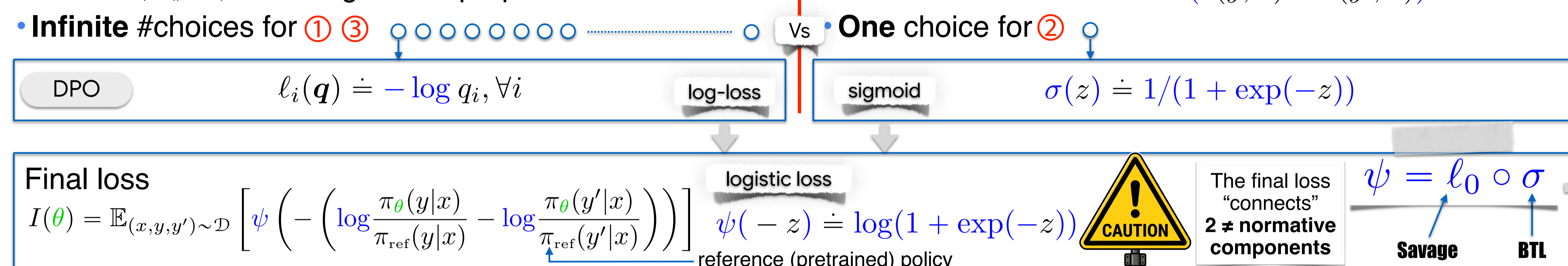
① ③ **Savage**

- Loss  $L(p, q) : \Delta_n \times \Delta_n \rightarrow \mathbb{R}$
- Analytic form:  $L(p, q) \doteq p^T \ell(q) = \mathbb{E}_{i \sim \text{CAT}(p)} [\ell_i(q)]$
- $p$ -proper iff  $L(p, q) \geq L(p, p), \forall q \in \Delta_n$ , proper iff  $p$ -proper for all  $p$
- Used as in ①; in ③ folds properness in  $\pi = \arg \max_{\pi'} \mathbb{E}_{\pi'}[r] - R(\pi' \| \pi_{\text{ref}})$  where  $R(\pi' \| \pi_{\text{ref}})$  is the regret of a proper loss



② **Bradley-Terry-Luce**

- Choice probability  $p(y > y'|x)$
- Analytic form:  $u(x,y) = r(x,y) + \varepsilon, \varepsilon \sim \text{Gumbel}$
- No abstention:  $p(y > y'|x) + p(y' > y|x) = 1$
- gives  $p(y > y'|x) \doteq p(u(x,y) > u(x,y')) = \sigma(r(y,x) - r(y',x))$

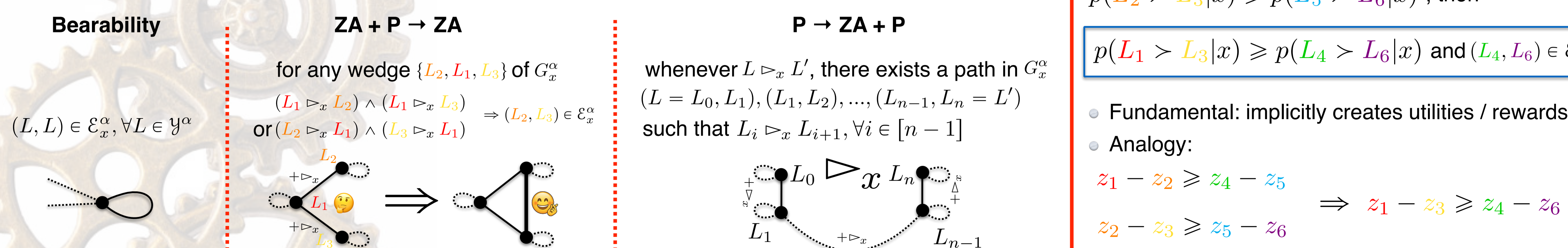


Our initial objective: broaden human choice part so it matches the ML freedom of Savage's properness

Our contribution: objective achieved + a few ML surprises along the way

## Getting above BTL and on par with Savage: KLST\*

- The theory exists to get above BTL: Doignon-Falmagne (1974) simplified by Krantz-Luce-Supes-Tversky (1989) but relies on an axiom, solvability, imposing  $|\mathcal{Y}| = \infty$
- We bypass the use of solvability via a mechanism compatible with sampling: (binary) lotteries  $\mathcal{Y}^{\alpha} \doteq \{(yy')_{\alpha} : y, y' \in \mathcal{Y}\} \rightarrow \text{lottery}(y_1 y_2)_{\alpha} \rightarrow \text{choice proba: } p((y_1 y_2)_{\alpha} > L|x) \doteq \alpha \cdot p(y_1 > L|x) + (1 - \alpha) \cdot p(y_2 > L|x)$
- Our model, **KLST\*** relies on three building blocks: *expandability* (lotteries), *local choice structure* and *monotonicity*
- **Local Choice Structure** (LCS): for any  $\alpha \in (0, 1)$ , graph  $G_x^{\alpha} \doteq (\mathcal{Y}^{\alpha}, \mathcal{E}_x^{\alpha})$
- **Zero-Abstention** (ZA) encoded in edges:  $(L, L') \in \mathcal{E}_x^{\alpha} \Leftrightarrow p(L > L'|x) + p(L' > L|x) = 1$
- **Preference** (P)  $\triangleright_x$  = choice proba  $\geq 1/2$ :  $L \triangleright_x L' \Leftrightarrow p(L > L'|x) \geq 1/2$
- LCS imposes three invariants in graph



## Theorems+

### Extension of BTL to match Savage

- Assume choice probabilities have a **KLTS\*** structure. Then there exists
  - A strictly  $\nearrow$  function  $F : \mathbb{R} \rightarrow [0, 1]$  such that  $F(-z) + F(z) \leq 1, \forall z$
  - A function  $u(x, y) : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  such that  $p(y > y'|x) = F(u(x, y) - u(x, y'))$

### Working the ML match with Savage

- For **any** strictly  $\nearrow$  function  $\psi : \mathbb{R} \rightarrow \mathbb{R}$ , and **any** strictly  $\nearrow$  function  $F : \mathbb{R} \rightarrow [0, 1]$ , there exists a strictly proper loss  $(\ell_0, \ell_1)$  such that

$$\psi = \ell_0 \circ F$$

- **Consequences**: to be feasible, loss design  $\psi$  **only** needs to be strictly monotonic (*no convexity required*) and works with **any KLST\*** model

### Technical Lemma that brings the connection

- For **any** strictly  $\nearrow$  function  $\ell : [0, 1] \rightarrow \bar{\mathbb{R}}$ , any  $K \in \mathbb{R}, a \in [0, 1]$ , the following loss  $(\ell_0, \ell_1)$  is strictly proper:

$$\left( \ell(p), K - \left( \int_a^p \frac{\ell(t)}{t^2} dt + \frac{1-p}{p} \cdot \ell(p) \right) \right)$$

No differentiability assumption of the partial losses, not even continuity  
Warning: nomenclature follows Reid & Williamson '2011

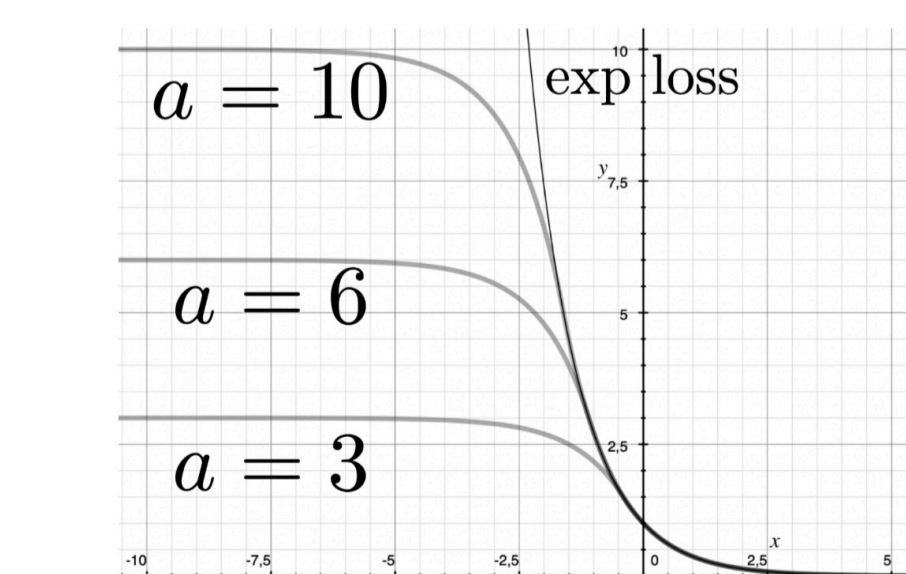
### A few additional results, summarized

- DPO's separable loss (log) is the **only one** that fits; extensions of DPO can be safeguarded (margins, correction for length, etc...)

## Toy Experiment

- Non-convex loss from exp loss balancing strong convexity vs Lipschitzness

$$\psi_a(-z) \doteq \begin{cases} a - (a^2/4) \cdot \exp(z) & \text{if } z < \log(2/a) \\ \exp(-z) & \text{otherwise} \end{cases}$$



- We use Gemma2\_2b\_it, rater is gemini\_2.5\_flash\_lite, test on Alpaca Eval v2
- Average over two runs, wins computed vs using exp loss

a	win% (us)
10	44.60%
6	<b>54.50%</b>
3	<b>53.00%</b>