

DPO Unchained: Your Training Algorithm is Secretly Disentangled in Human Choice Theory (and its Loss' Convexity is Dispensable)



Wenxuan Zhou



Shujian Zhang



Brice Magdalou



John Lambert



Ehsan Amid



Richard Nock



Andrew Hard

Google DeepMind

Google Research



UNIVERSITÉ DE
MONTPELLIER



Summary of the task

- Vocabulary:
 - Set of states / prompts \mathcal{X} , set of actions / answers \mathcal{Y} , probability simplex over \mathcal{Y} : $\Delta_{\mathcal{Y}}$
 - Want to learn a policy π , an application $\mathcal{X} \rightarrow \Delta_{\mathcal{Y}}$, from *data*
 - *Data* \sim distribution \mathcal{D} over triples (x, y, y') meaning: given prompt $x \in \mathcal{X}$, answer $y \in \mathcal{Y}$ is *chosen* (by a *human*) over answer $y' \in \mathcal{Y}$

▲ Setup learns $p(y > y'|x)$ ($>$: “*chosen over*”) \neq output we want, $\pi_{\theta}(y|x), \pi_{\theta}(y'|x)$ (policy)

- To get $p \rightarrow \pi_{\theta}$ without mishaps requires carving appropriate *normative* properties in loss function(s) to keep training from reaching inadequate solutions
- Our starting point is a loss measured on $p, L(p)$

Normative Bedrock (on *loss functions*, on **human choice**)

Choice

Policy

(pre-DPO,
RLHF)

Google Research

DPO = Direct Preference Optimization, RLHF = Reinforcement Learning from Human Feedback

Normative Bedrock (on *loss functions*, on **human choice**)

Choice

$$L(p)$$

(pre-DPO)

Policy

Normative Bedrock (on *loss functions*, on **human choice**)

Choice

Reward

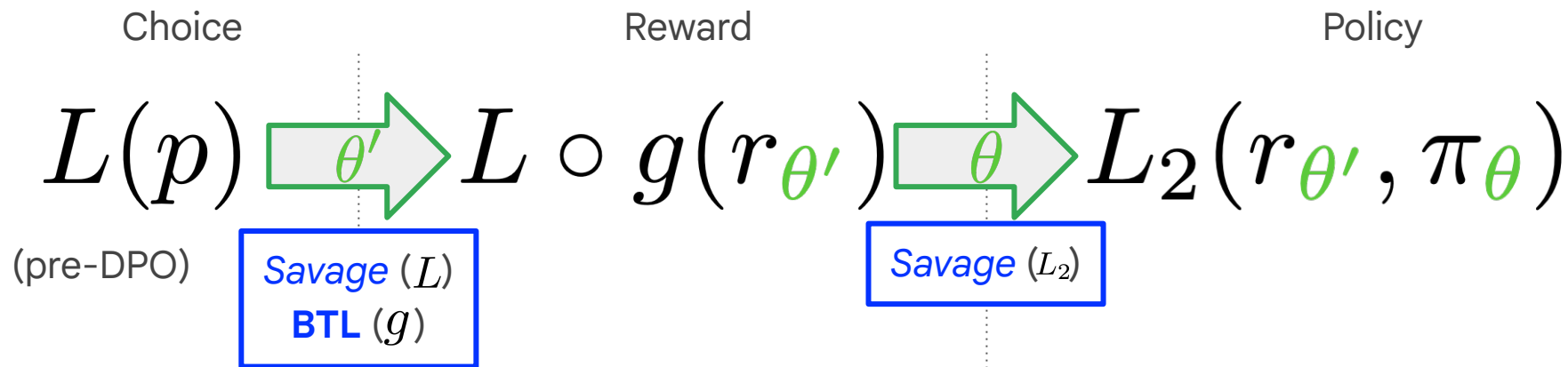
Policy

$$L(p) \xrightarrow{\theta'} L \circ g(r_{\theta'})$$

(pre-DPO)

In green: ML / training element

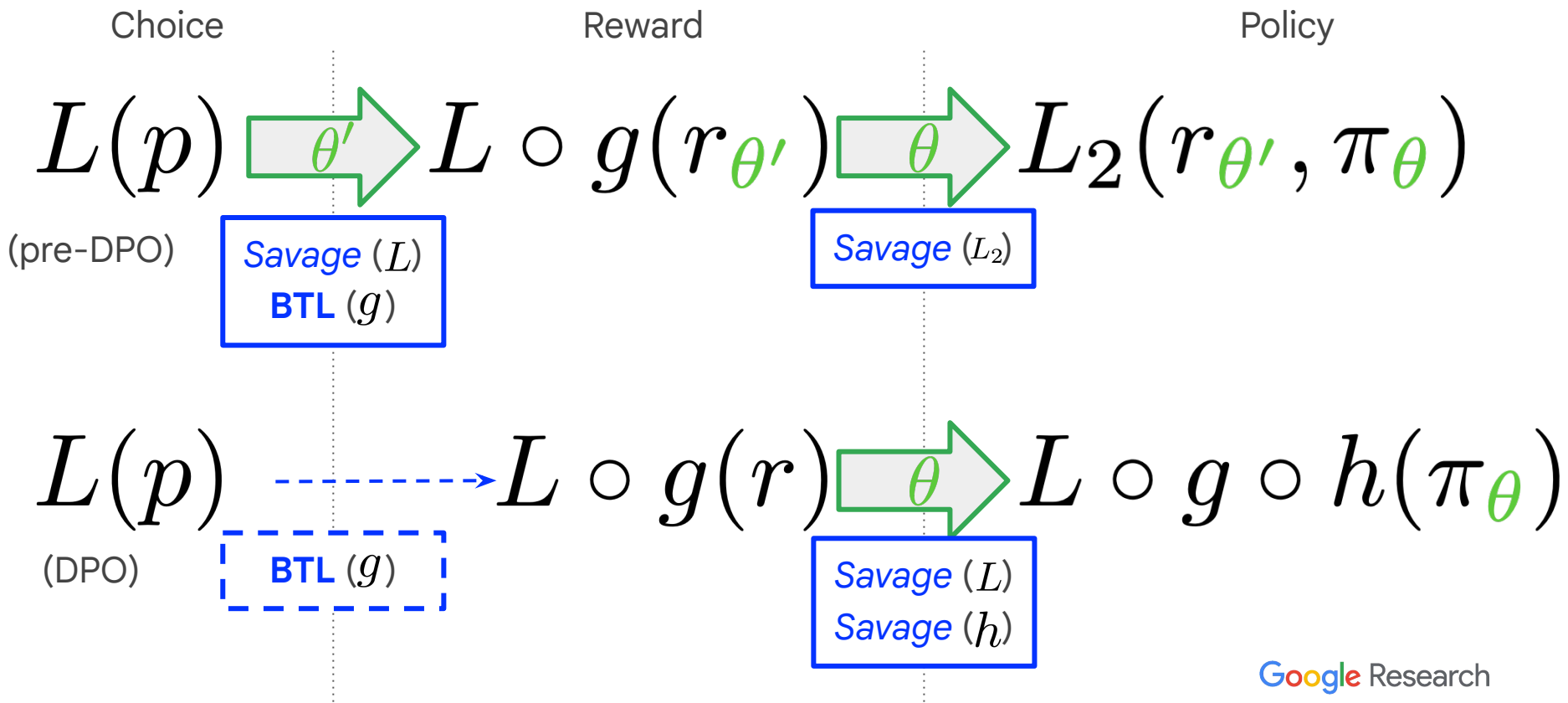
Normative Bedrock (on *loss functions*, on **human choice**)



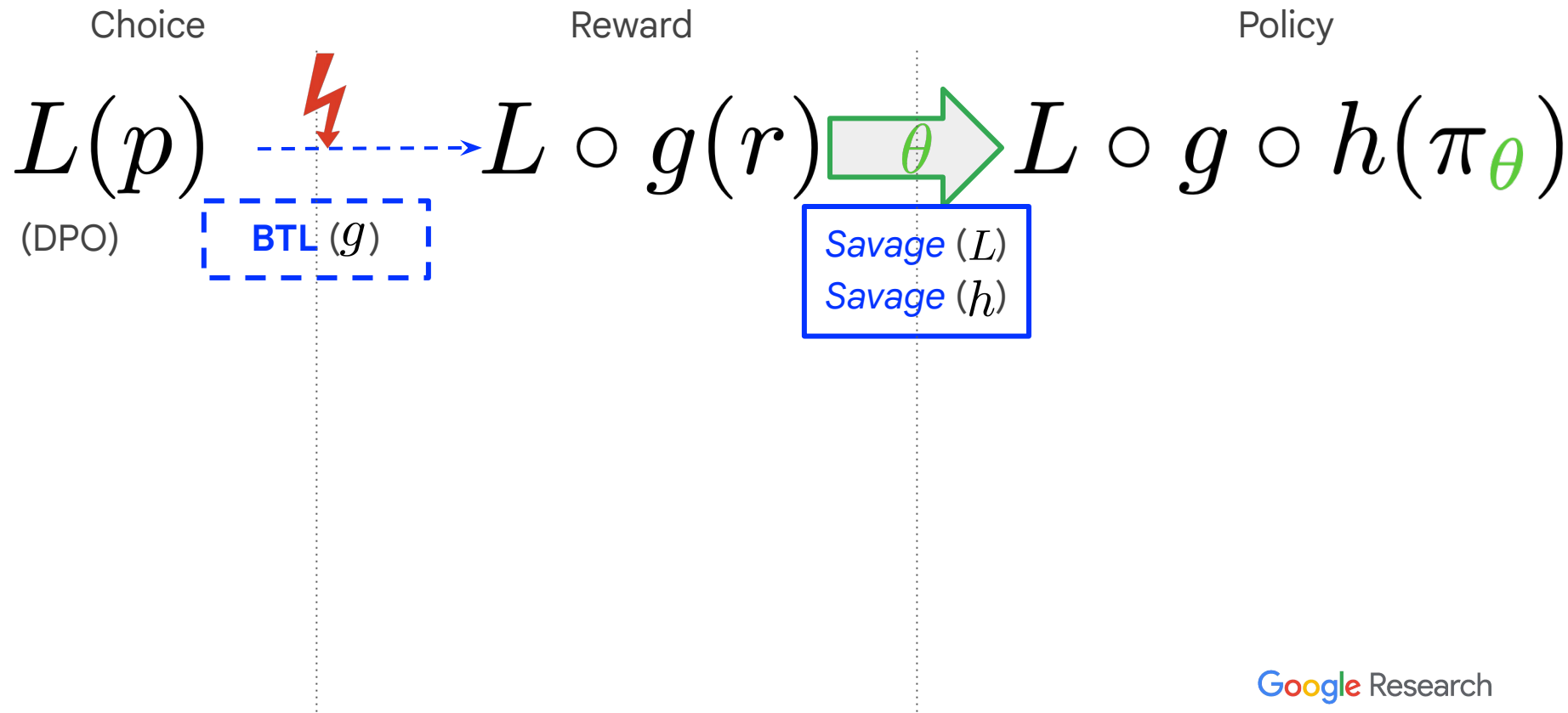
BTL = Bradley-Terry-Luce

In green: ML / training element

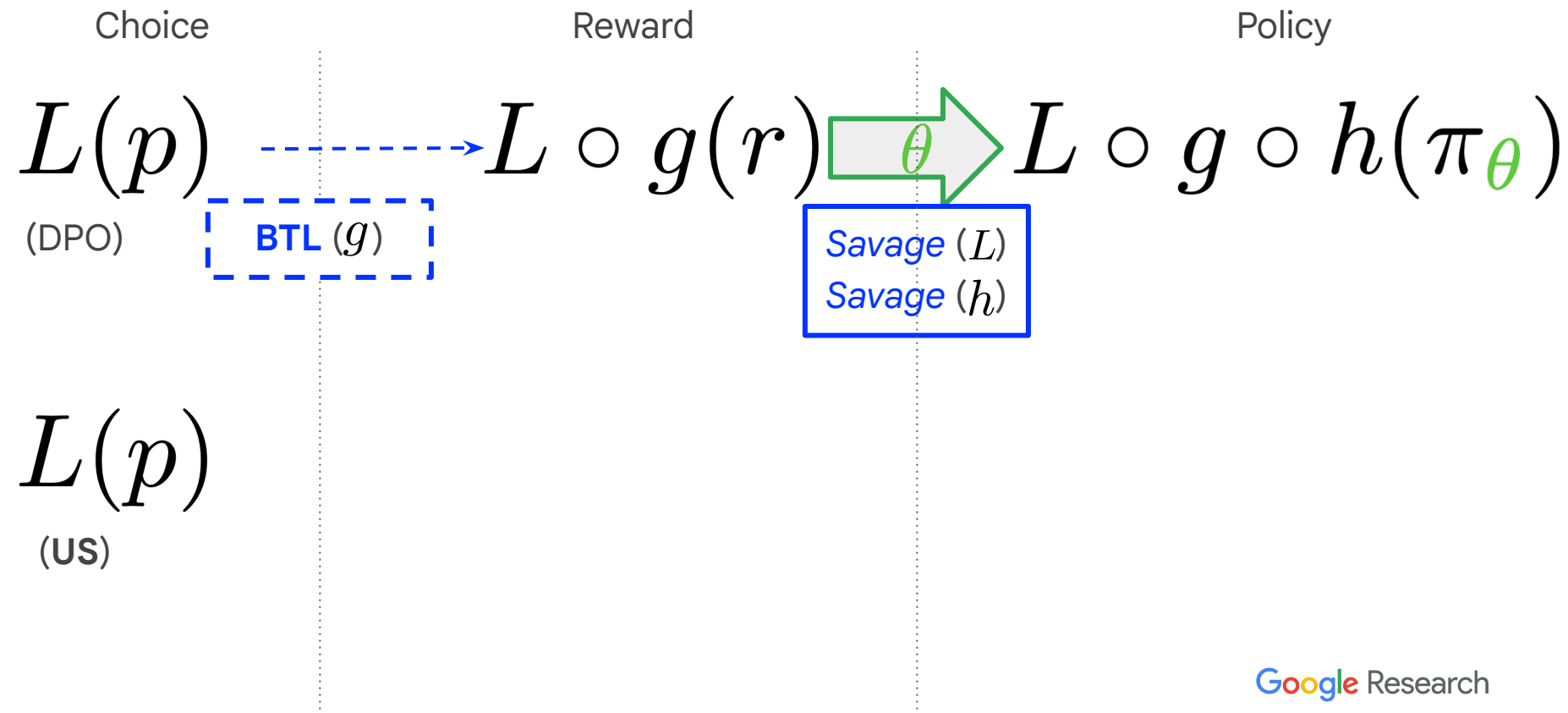
Normative Bedrock (on *loss functions*, on **human choice**)



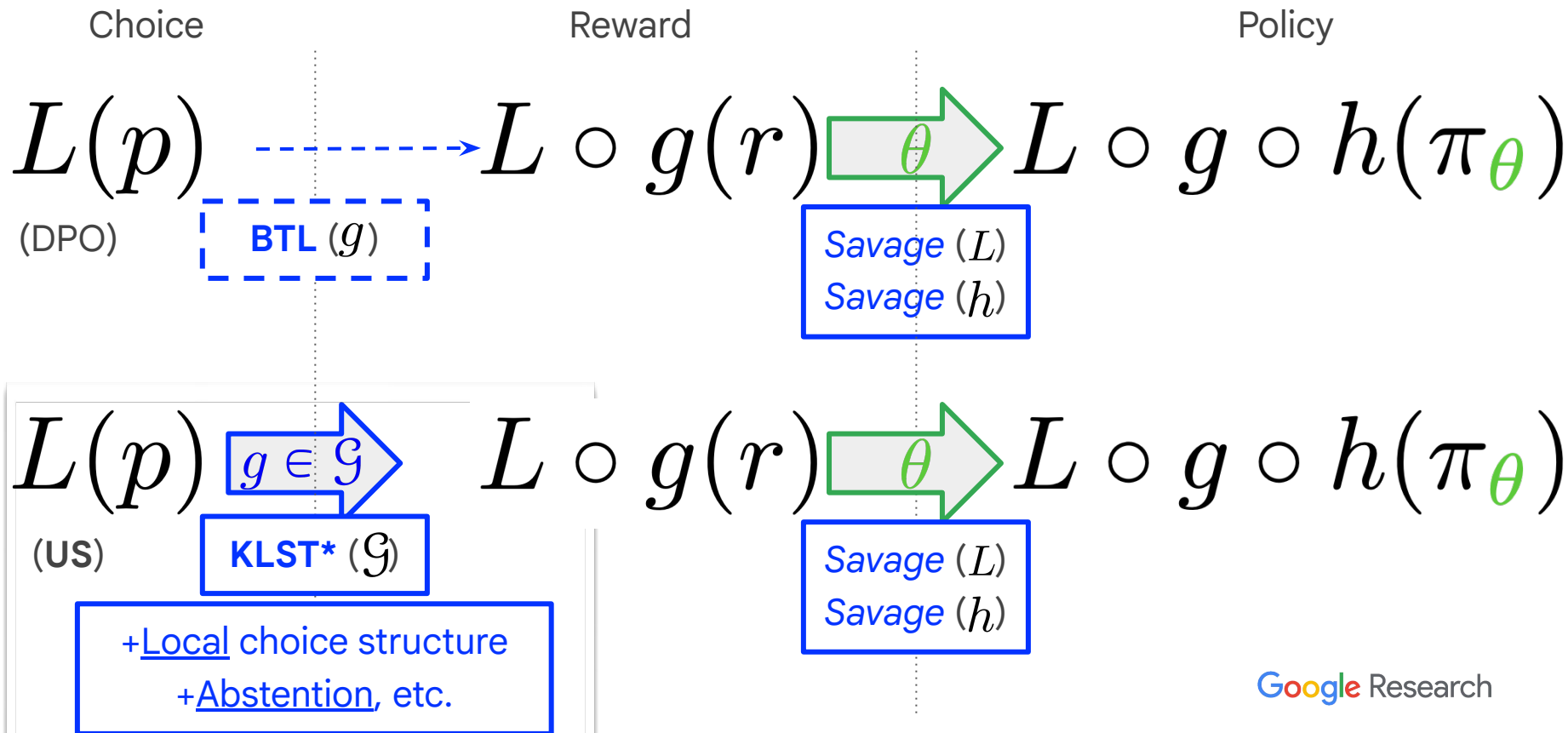
Normative Bedrock (on *loss functions*, on **human choice**)



Normative Bedrock (on *loss functions*, on **human choice**)



Normative Bedrock (on *loss functions*, on **human choice**)



Surprise(s) 🎁 – among others

- The key loss function is $L \circ g \doteq \psi$ (logistic loss for DPO)
- A classical issue with normative approaches is why should we pick one feasible option instead of another – especially true for human choice (e.g. one specific instance of **KLST***)



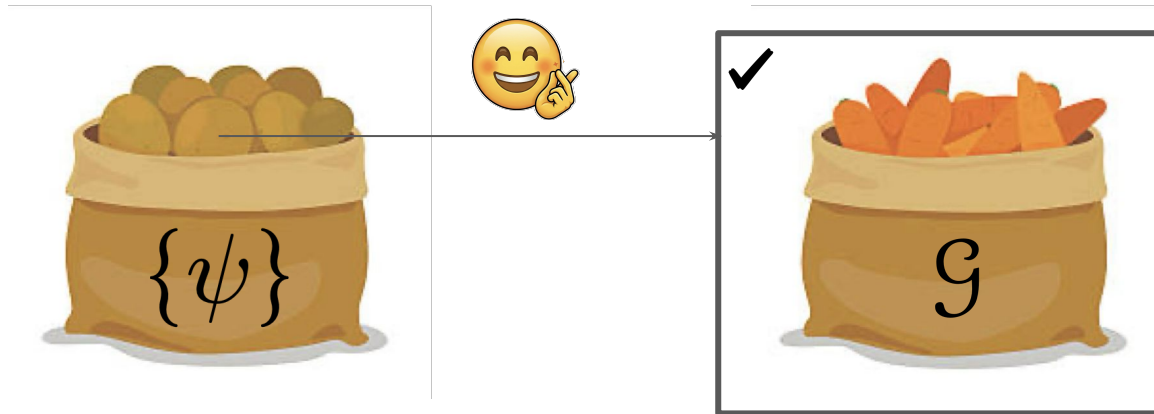
Surprise(s) 🎁 – among others

- The key loss function is $L \circ g \doteq \psi$ (logistic loss for DPO)
- A classical issue with normative approaches is why should we pick one feasible option instead of another – especially true for human choice (e.g. one specific instance of **KLST***)
- We were expecting this issue showing in “full force” for the design of ψ ...



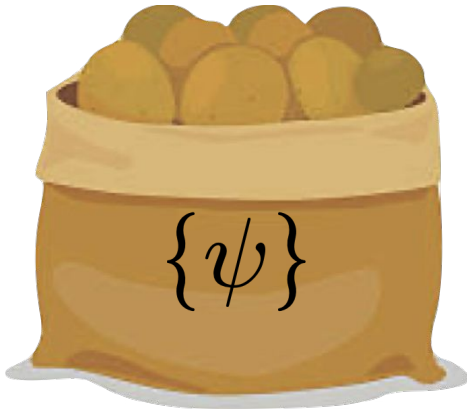
Surprise(s) 🎁 – among others

- The key loss function is $L \circ g \doteq \psi$ (logistic loss for DPO)
- A classical issue with normative approaches is why should we pick one feasible option instead of another – especially true for human choice (e.g. one specific instance of **KLST***)
- We were expecting this issue showing in “full force” for the design of ψ ...
but instead got the possibility for *any* applicable loss ψ to map to *any* human choice model...



Surprise(s) 🎁 – among others

- The key loss function is $L \circ g \doteq \psi$ (logistic loss for DPO)
- A classical issue with normative approaches is why should we pick one feasible option instead of another – especially true for human choice (e.g. one specific instance of **KLST***)
- We were expecting this issue showing in “full force” for the design of ψ ...
but instead got the possibility for *any* applicable loss ψ to map to *any* human choice model...
and the set of applicable losses to be larger than expected (e.g. no convexity requirement)



Thank You

Wenxuan Zhou
Shujian Zhang
Brice Magdalou
John Lambert
Ehsan Amid
Richard Nock
Andrew Hard

Google DeepMind
Google DeepMind
CEE-M, Univ. Montpellier, CNRS, INRAé, Institut Agro
Google DeepMind
Google DeepMind (former)
Google Research
Google DeepMind