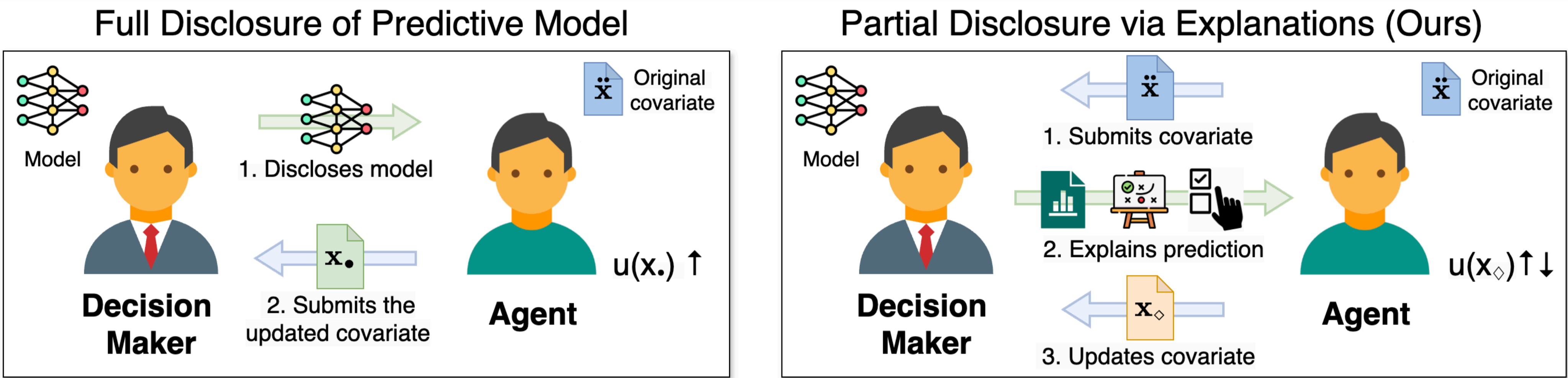


EXPLANATION DESIGN IN STRATEGIC LEARNING: Sufficient Explanations that Induce Non-harmful Responses

Kiet Q.H. Vo¹, Siu Lun Chau², Masahiro Kato³, Yixin Wang⁴, Krikamol Muandet¹

AN EXPLANATION DESIGN PROBLEM



Here, the agent can *correctly anticipate* how changing \ddot{x} affects the prediction, then picks an update x_* that *reliably improves* utility $u(x_*)$.

With only an explanation (i.e., partial information), the agent’s update x_\diamond *might not improve* utility $u(x_\diamond)$.

- Q1: Can we ensure no reduction in agents’ utilities?

Q2: Is there a sufficient class of explanations that guarantee this?



SETUP		KEY CONTRIBUTIONS	
<div> <div>0. Agent is realised: $(\ddot{x}_t, c_t) \sim P_{\ddot{X}, C}$, where c_t is a cost function.</div> <div>1. DM predicts agent’s risk: $g(\ddot{x}_t)$.</div> <div>2. DM gives explanation: $e_t := \sigma(g, \ddot{x}_t)$.</div> <div>3. Agent modifies covariate: $x_t := \psi(\ddot{x}_t, e_t, c_t)$.</div> <div>4. DM updates the prediction from $g(\ddot{x}_t)$ to $g(x_t)$.</div> </div> <div> <div>Agent’s true utility: $u_t(g, x) := b(x) - c_t(\ddot{x}_t, x)$</div> <div>$\qquad\qquad\qquad = -g(x) - c_t(\ddot{x}_t, x)$.</div> </div> <div> <div>Agent’s non-harmful responses:</div> <div>$\nu_t = \{x \in \mathcal{X} : u_t(g, x) \geq u_t(g, \ddot{x}_t)\}$.</div> </div>		<div> <div>A <i>necessary condition</i> to ensure surrogate models do not mislead agents into self-harming actions.</div> </div> <div> <div>Action recommendation (AR)-based explanations (ARexes) make up <i>a sufficient class</i> that guarantees non-harmful agents’ responses.</div> </div>	
SURROGATE MODELS	AR-BASED EXPLANATIONS	EXPERIMENTS	
<div> <div>• If it holds that, every cost function c_t induces a response $x_t \in \nu_t$, then f_t <i>must satisfy</i>: $f_t(\ddot{x}_t) - f_t(x') \leq g(\ddot{x}_t) - g(x') \forall x' \in \mathcal{X}_t^{\downarrow}$.</div> <div>• Corollary 3.3 extends this result to other forms of explanations.</div> </div> <div> <div>The surrogate model <i>must understate</i> the agent’s gain in prediction score.</div> </div>	<div> <div>• For any <i>arbitrary</i> explanation e' that induce $x_* \in \nu_t$, There exists an ARex (\ddot{x}_t, \hat{y}_t) that induces the <i>same response</i> $x_* \in \nu_t$.</div> <div>• Theorem 3.6 extends this result to ARex methods for heterogeneous agents.</div> </div> <div> <div>It suffices to <i>use only</i> AR-based explanations.</div> </div>	<div> <div>ARexes do not reduce agents’ utilities.</div> </div> <div> <div>In practice, ARexes can be jointly optimised with predictive model g.</div> </div>	