

Composing Efficient, Robust Tests for Policy Selection

Dustin Morrill, Thomas J. Walsh, Daniel Hernandez, Peter R. Wurman, Peter Stone

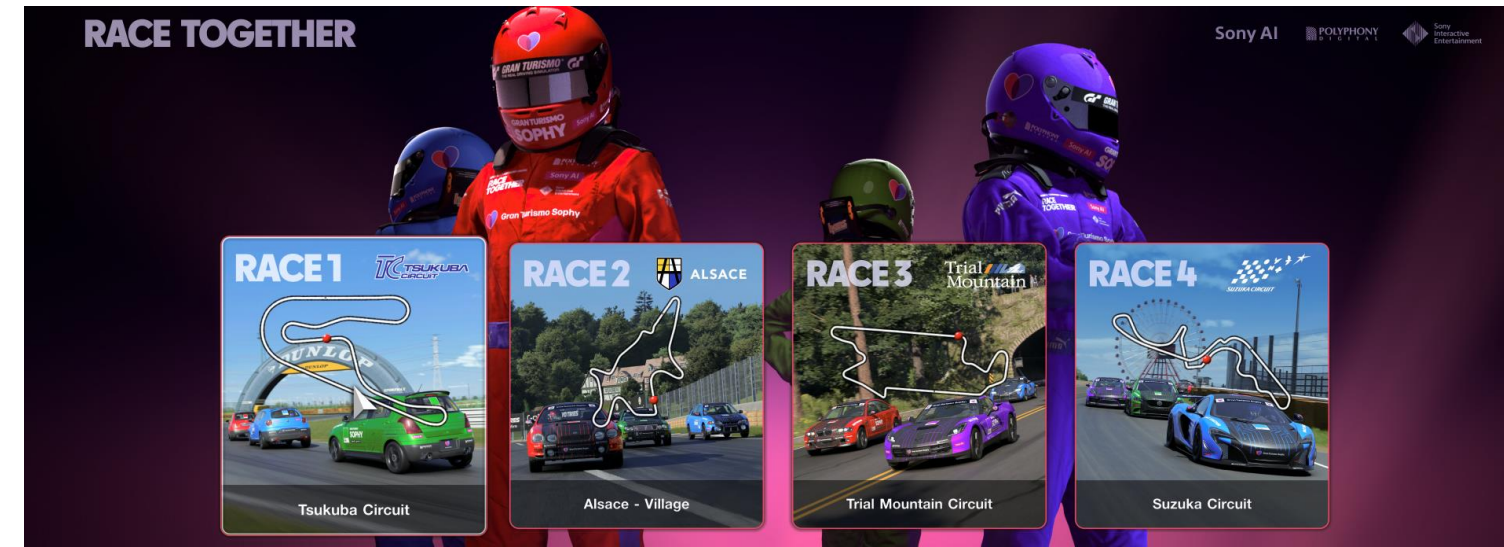


Sony AI

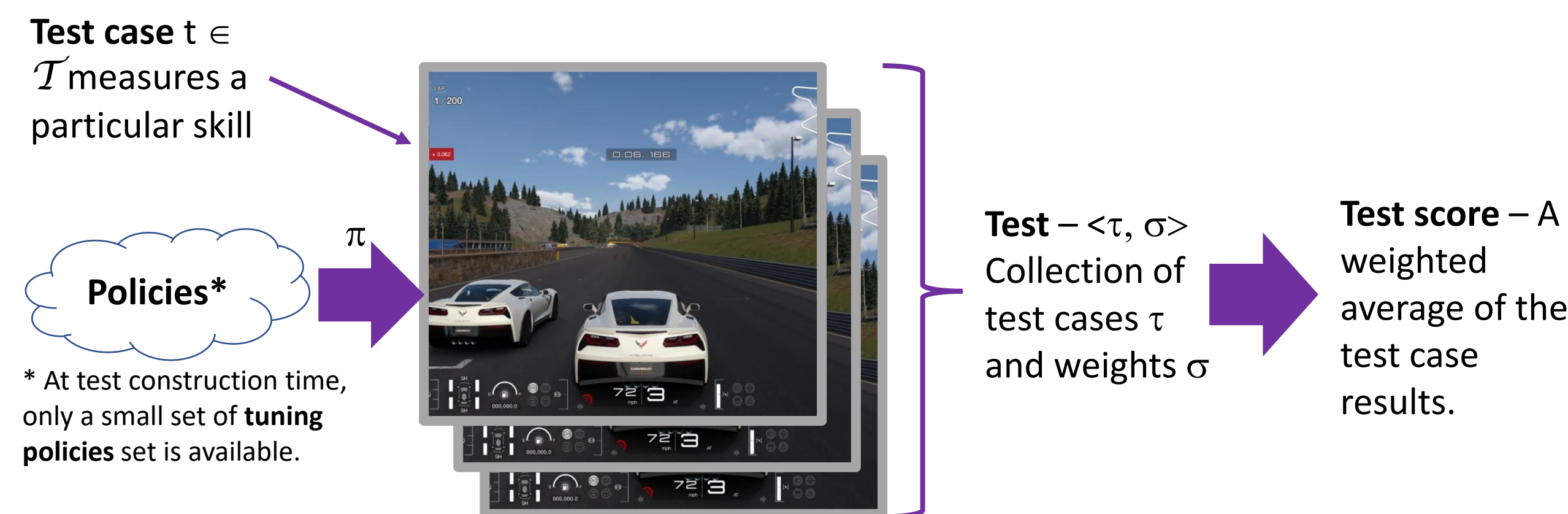
<https://ai.sony/joinus/>

Dustin.Morrill@sony.com

A real AI deployment problem: *Hundreds of candidate deployment policies, dozens of test cases, but you can only run a few test cases...*



Robust Test Construction



Problem: Construct an efficient test $\langle \tau, \sigma \rangle$ by selecting a small number of test cases $\tau \subset \mathcal{T}$ and test case weights $\hat{\sigma} \in \Delta^{|\mathcal{T}|}$ to approximate a full test, $\langle \mathcal{T}, \sigma \in \Delta^{|\mathcal{T}|} \rangle$ despite **uncertainty** over future policies to evaluate and the future test case weights.

Test construction example

Uncertainty over future policies (e.g. a higher proportion of aggressive agents)

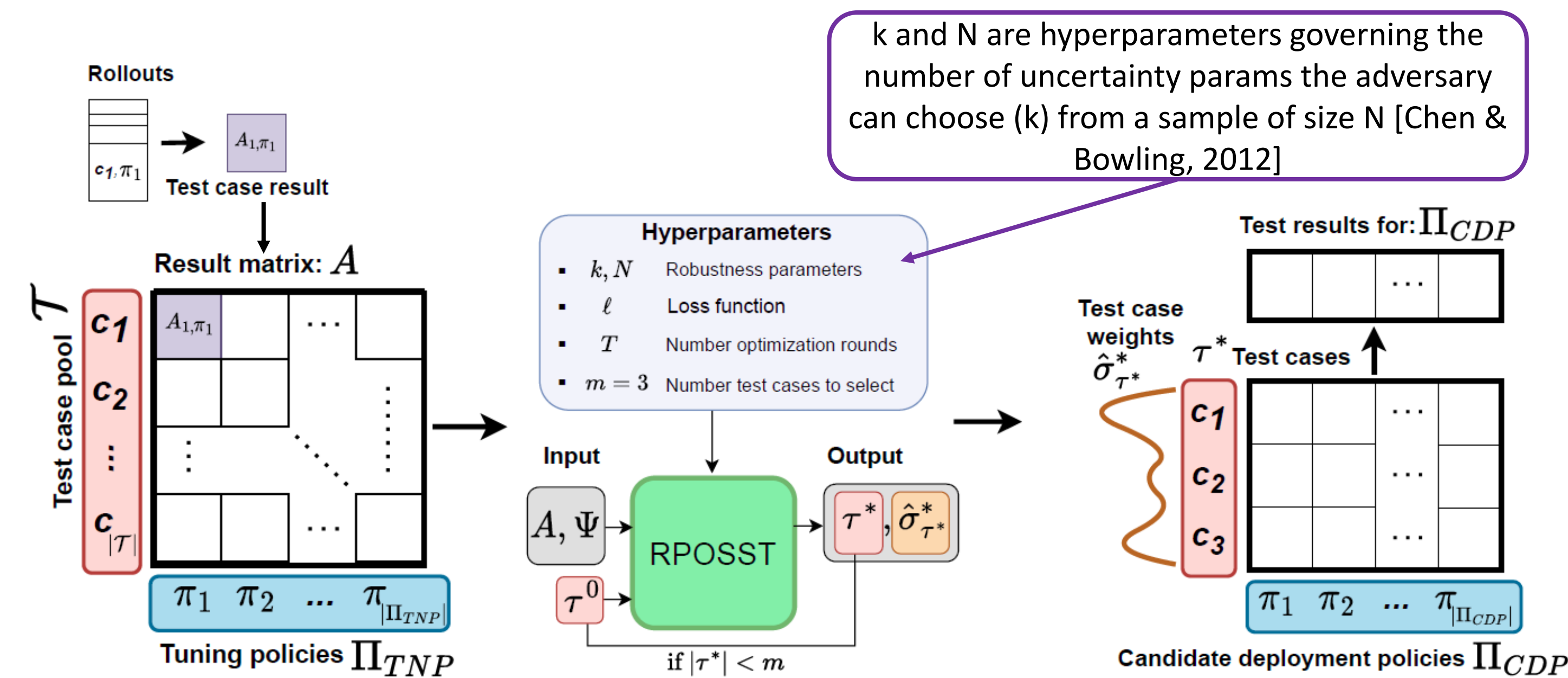
Uncertainty over target distribution (e.g. different emphasis on off-course behavior)

Target (σ)	policy / case	#1	#2	#3	#4	#5
?	C1	0	+1	+1	+1	+1	?	?
?	C2	-1	0	+1	+1	+1	?	?
?	C3	-1	-1	0	+1	+1	?	?
?	C4	-1	-1	-1	0	+1	?	?
?	C5	-1	-1	-1	-1	0	?	?

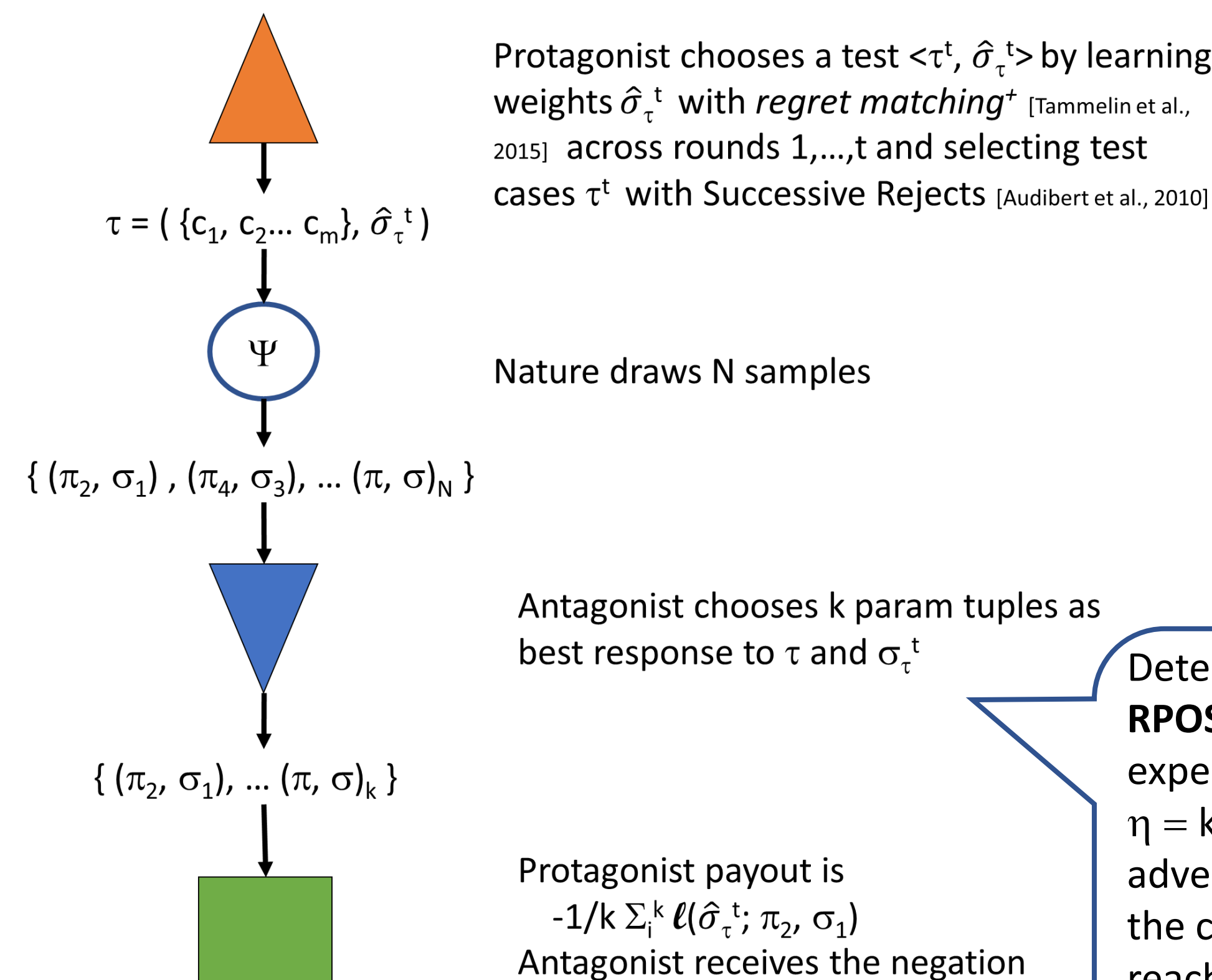
Possible Test	max error	size
Uniform	0	5
Strongest (C5)	1.4	1
Middle (C3)	0.6	1
[0, 0.5, 0, 0.5, 0]	0.2	2
[0.27, 0, 0.47, 0, 0.27]	0.07	3

RPOSST Selections

Constructing a Robust Test with RPOSST



RPOSST_{SEQ} Game



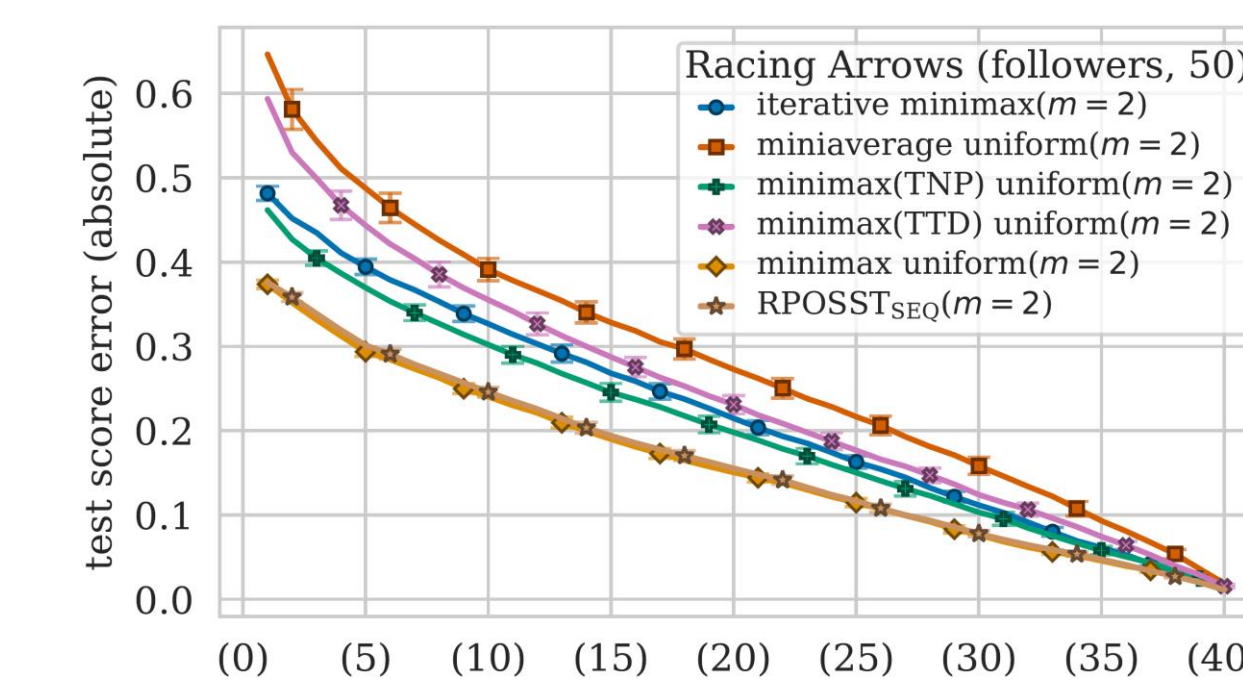
Theorem 4.1. After $T' \sim \text{Unif}(\{1, \dots, T_1\})$, $T_1 > 0$, rounds of its optimization game, Algorithm 1 selects an m -tuple of test cases, τ^* and weights $\hat{\sigma}_{\tau^*} \in \Delta^m$ that, with probability $(1-p)(1-q)(1-\alpha)$, $p, q, \alpha > 0$, are ε_q -optimal for Equation (2), where $\varepsilon = \mathcal{O}(\sqrt{\frac{1}{T_1 m}} + \sqrt{\frac{1}{T_1} \log(1/p)})$ and $\alpha = \mathcal{O}(e^{-T_2})$.

Corollary 4.2. Assume that $\Psi \in \Delta^d$ for some finite $d \geq 1$. After T rounds of the CVaR(η) RPOSST_{SEQ} optimization game (using regret matching*), τ^* and $\hat{\sigma}_{\tau^*}$ are ε -optimal for Equation (2) under the η -fractile CVaR robustness measure, where $\varepsilon = \mathcal{O}(\sqrt{\frac{1}{T m}})$.

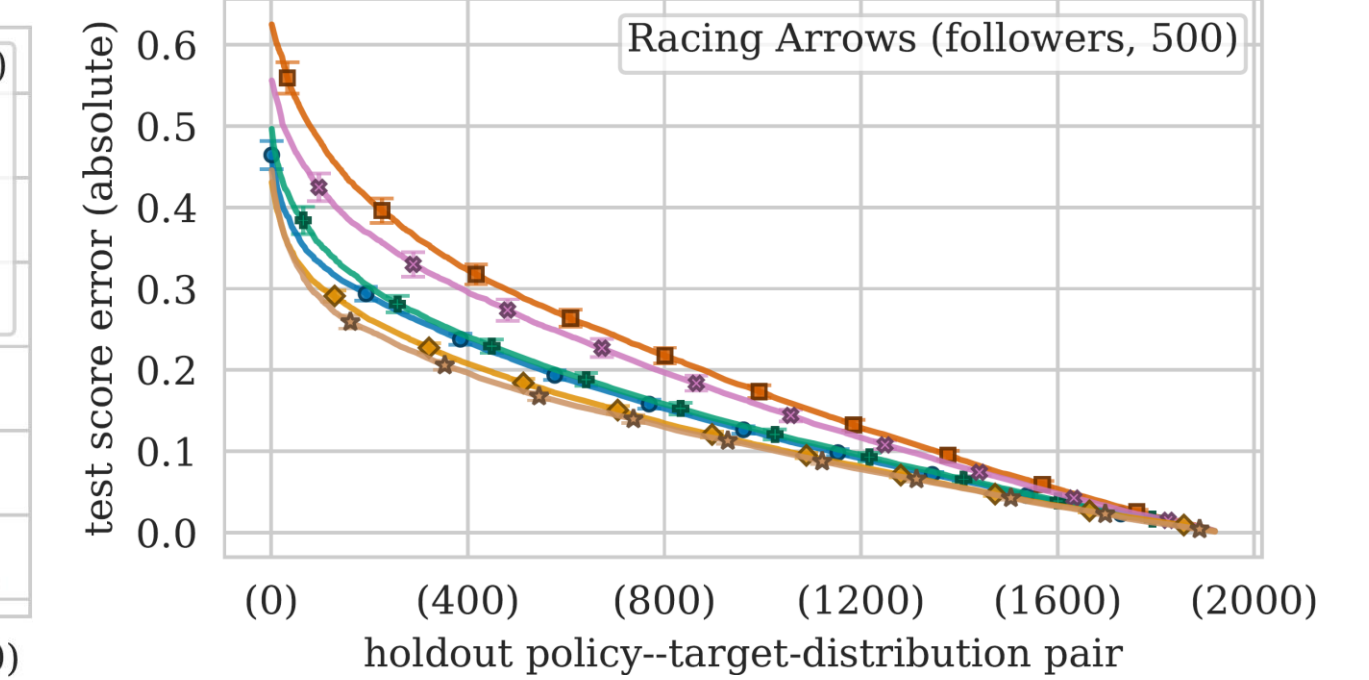
Empirical Results

Racing Arrows (one-shot simultaneous move game)

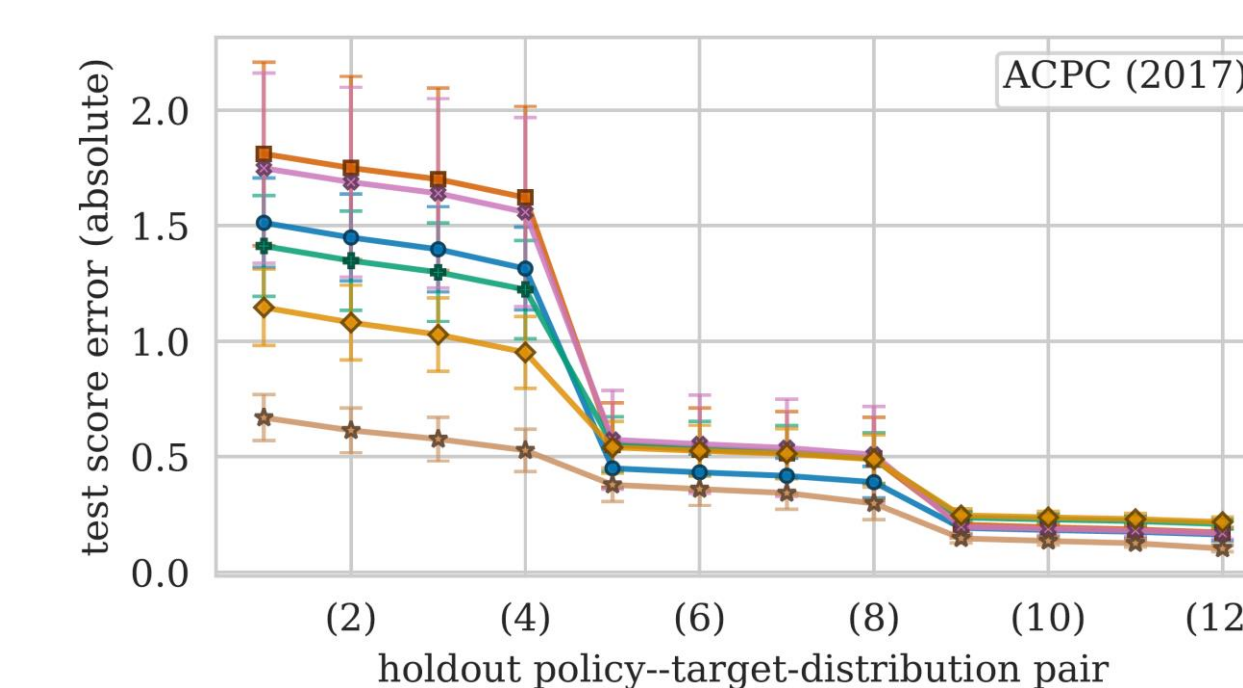
50 policies, 20% holdout, $m=2$



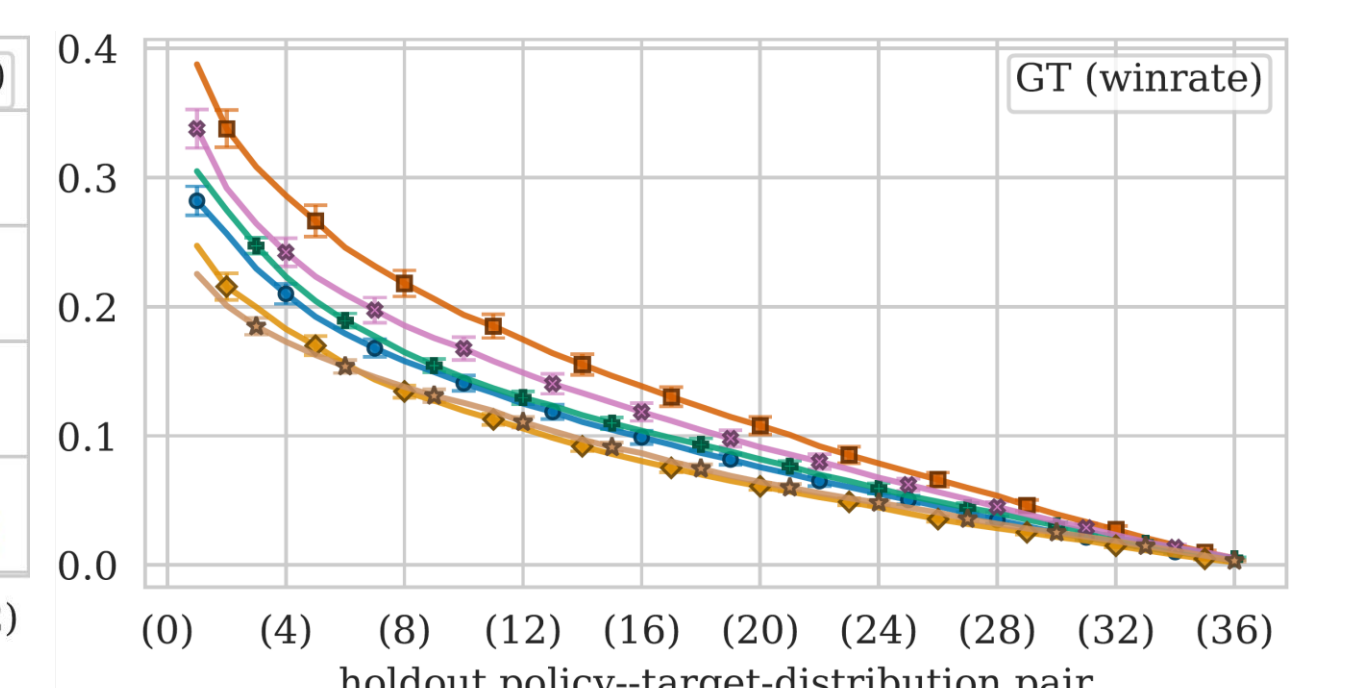
500 policies, 96% holdout, $m=3$



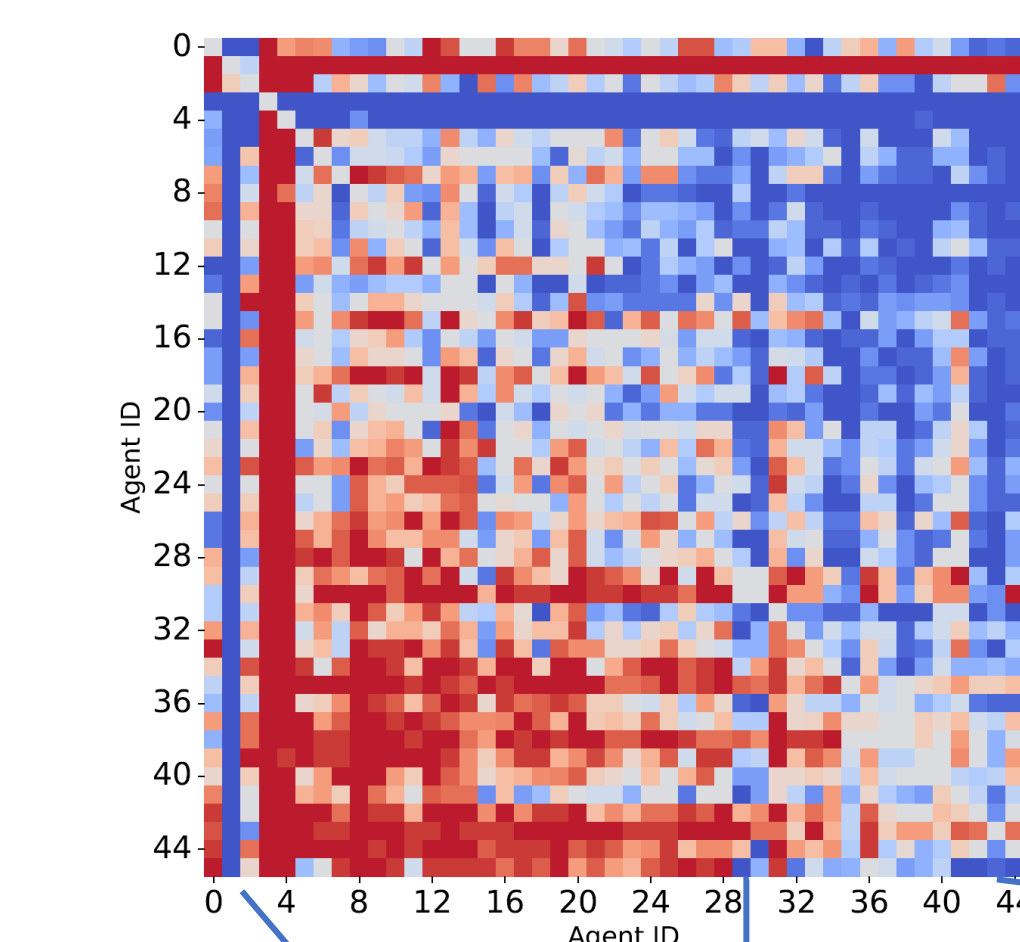
Poker (no-limit)



Gran Turismo 7 1v1



Inside the Gran Turismo 7 1v1 experiment



Test cases (opponents) 16 and 41 were chosen by CVaR(η) RPOSST_{SEQ}: 41 provides a nearly 50/50 information split and 16 (a weaker policy) is beaten soundly by only very good policies.

The 3 best policies identified with just these 2 tests are the bluest columns (strongest policies) in the 46 X 46 matrix.

