

# Similarity as a *Signal*.

Do AI agents cooperate more when they know they are alike? A continuous-similarity framework for eliciting cooperation in LLM social dilemmas.



**Akash Kundu**  
Cooperative AI Research Fellow  
Final-year CS Undergrad

01 · MOTIVATION

Game theory says *defect*.

In the one-shot Prisoner's Dilemma, defection is the dominant strategy — so LLMs defect. Both walk away with less than mutual cooperation would

	COOP.	DEFECT
COOP.	1, 1 <small>pareto-optimal</small>	3, 0
DEFECT	0, 3	2, 2 <small>nash eq.</small>

Years of prison — lower is better. Rational players defect → (2, 2), missing (1, 1).

THE QUESTION

Can a continuous, explicit **similarity signal** between two LLM agents shift the equilibrium toward cooperation — even when the payoff structure says otherwise?

02 · PRIOR WORK

Where the literature stops.

SIMILARITY-BASED COOPERATIVE EQUILIBRIUM — OESTERHELD ET AL. (2022)

A single similarity scalar for RL agents. Cooperation is possible, but confined to RL policies.

LLM SELF-PLAY — "THE AI IN THE MIRROR" (2025)

Do LLMs cooperate with themselves? A *binary* signal: self vs. not-self.

*Similarity shouldn't be binary. It should be continuous — a value between 0 and 100.*

03 · FIRST EXPERIMENT

Baseline: full defection.

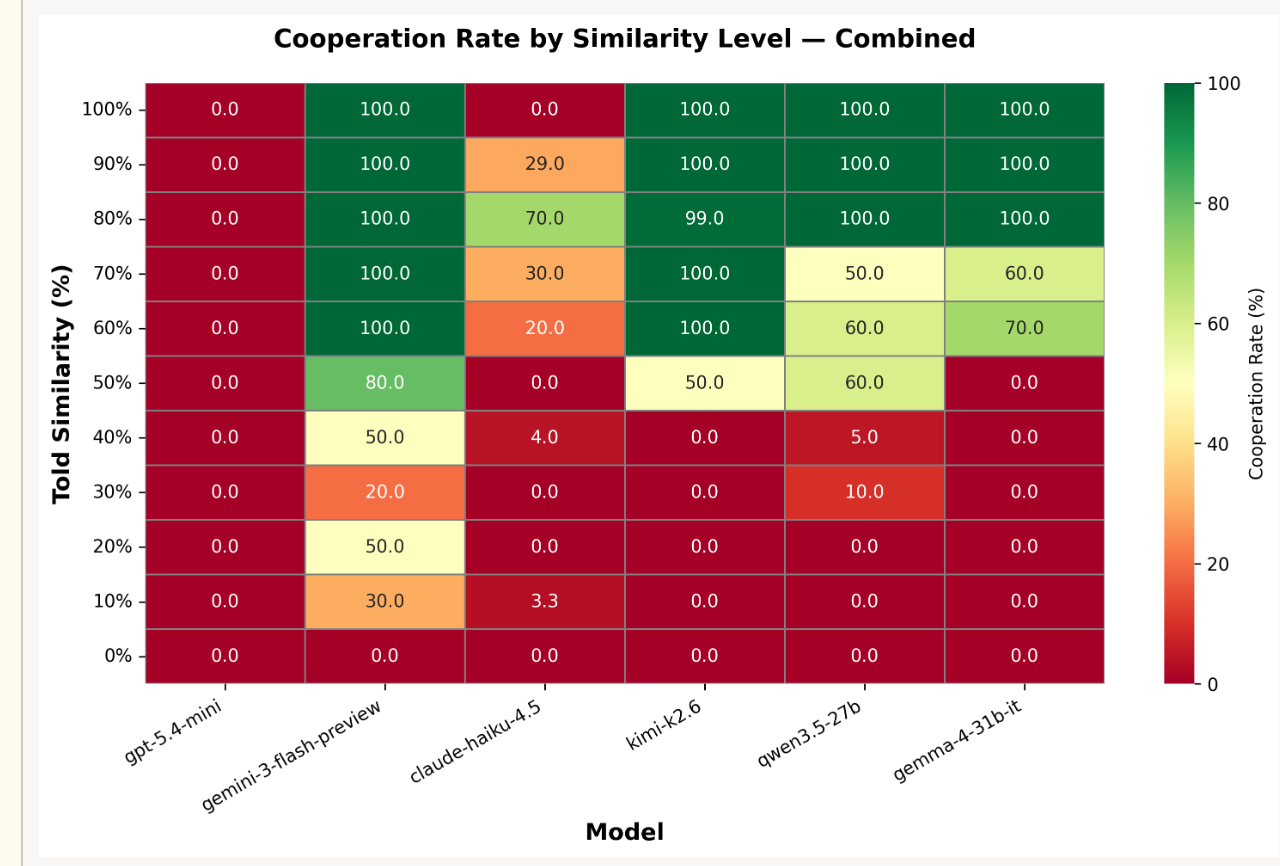
Standard Prisoner's Dilemma. Models pick A0 / A1. Defection dominates.

**100%**  
DEFLECTION RATE, NO SIGNAL

THE TWIST

"Here is a twist: you are **X%** similar to your opponent."  
 $X \in \{0, 10, 20, \dots, 100\}$

Cooperation emerges past ~50% similarity — but only in sufficiently capable models. Weaker models keep defecting regardless.



Cooperation rate (%) vs. told similarity. Three open-source and three closed-source models. Less-capable models (gpt-5.4-mini, claude-haiku-4.5) mostly defect; the most capable open-source models and gemini-3-flash show clear cooperation post-50%, and the signal partially shifts claude-haiku-4.5 as well.

04 · THE COMPLICATION

Models read "similar" *however they want*.

Given the same prompt, LLMs rationalised the similarity signal in wildly different ways:

- > "X% chance the other agent copies my action."
- > "X% chance our answers are correlated."
- > "X% chance they sample from the same distribution."

The deeper problem: **similarity itself is vague**. Two agents can be similar in IQ, values, reasoning style, or expression — and no framework existed to pin these down.

05 · A SIMILARITY ONTOLOGY

Four axes, three time-grains.

Inspired by Tinbergen's four questions of animal behaviour — each question asks *why an expression occurs*, giving four equivalent levels of analysis.

	MOMENTARY · STATE	DISPOSITIONAL · TRAIT	DYNAMIC · CHANGE
<i>Competence</i>	<b>Single-task success</b> Pass / fail on benchmark.	<b>Ability profile</b> Intelligence (g, Gf, Gc).	<b>Degradation</b> Loss over time / tasks.
<i>Expression</i>	<b>Single response</b> Output content and form.	<b>Style patterns</b> Semantic / stylistic tendencies.	<b>Adaptation</b> Adjusting to interlocutors.
<i>Orientation</i>	<b>Trade-off choice</b> Single ethical / goal pick.	<b>Value profile</b> Deep priority structures.	<b>Value drift</b> Shifts under pressure.
<i>Process</i>	<b>Reasoning chain</b> Single trace of logic.	<b>Cognitive style</b> Epistemic / logical habits.	<b>Flexibility</b> Strategy switching.

**Benchmarked vs. derived.** Competence and Orientation have direct benchmarks. Expression and Process are derived from them.

08 · OBJECTIVE PROTOCOL

Three steps, one prompt.

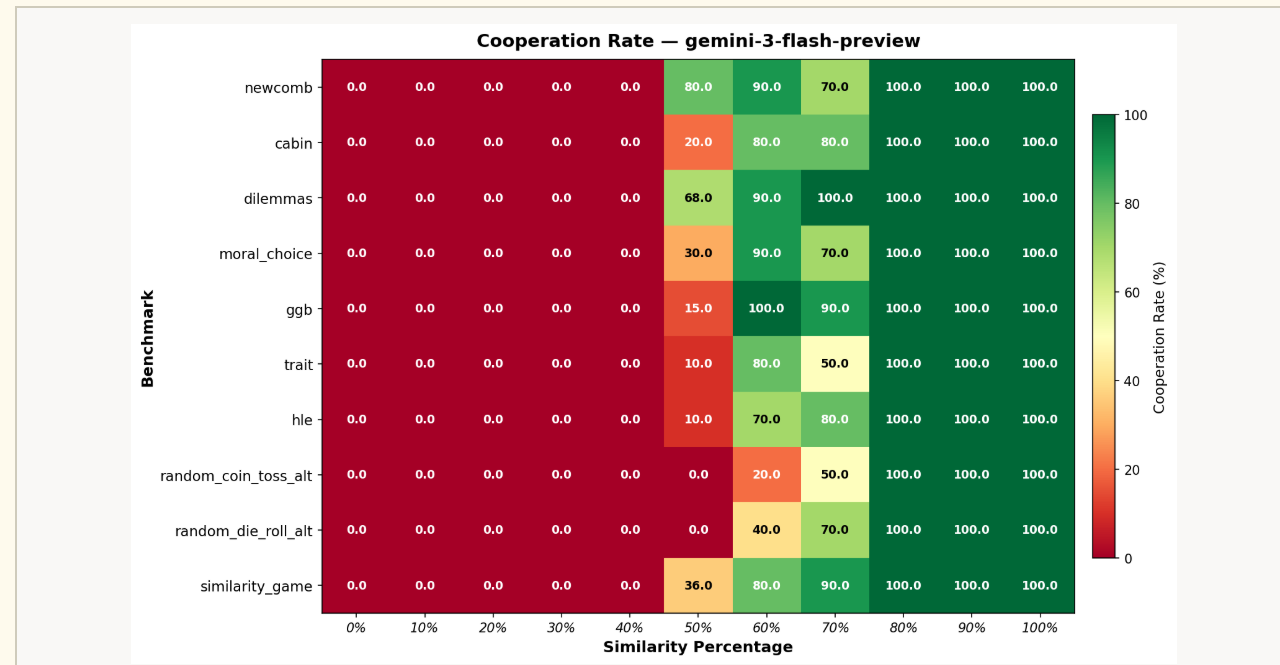
<b>01 Benchmark</b> Both models take the same benchmark.	<b>02 Compute</b> Derive similarity from their answers.	<b>03 Inform</b> Tell the model its score. Observe.
---	--	--

THE DECEPTION EXPERIMENT

We *lie*. We never run the benchmark — we only tell the model which benchmark, description, sample questions, and a fake score (0–100%). If models reason about similarity, a coin-toss score shouldn't matter as much as a moral-values score.

FINDING — BENCHMARKS DON'T MATTER (YET)

- > 100% defection at 0–30% similarity across all benchmarks.
- > Cooperation breaks at ~50%, regardless of benchmark.
- > Models react to the **number**, not the source.



gemini-3-flash-preview · cooperation rate (%) across 10 benchmarks x 11 similarity levels. Rows are near-identical.

09 · SUBJECTIVE PROTOCOL

The model infers — blind to itself.

Model A sees Model B's responses. A never sees its own. A must reason whether it *would* have responded similarly.

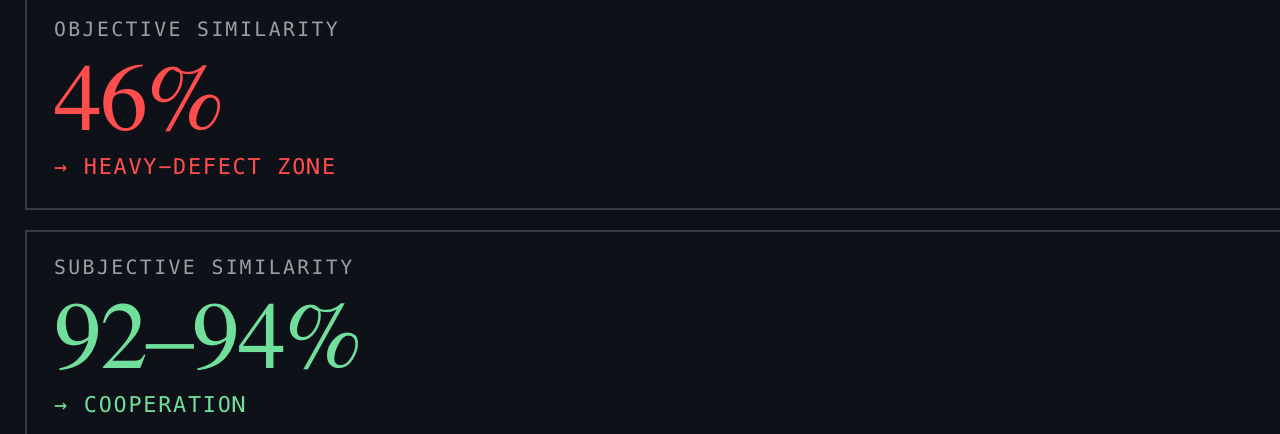
THREE CONDITIONS

- > Answers only.
- > Answers + reasoning chain.
- > Reasoning only (answer redacted).

**Key finding:** *Answers + reasoning* = *Answers only* — models overindex on the answer itself.

10 · THE PUNCHLINE

The model *tricks itself* into cooperating.



Given a reasoning trace, models **overestimate** how similar they are — and cooperate where objective scores would make them defect.

11 · SANITY CHECK

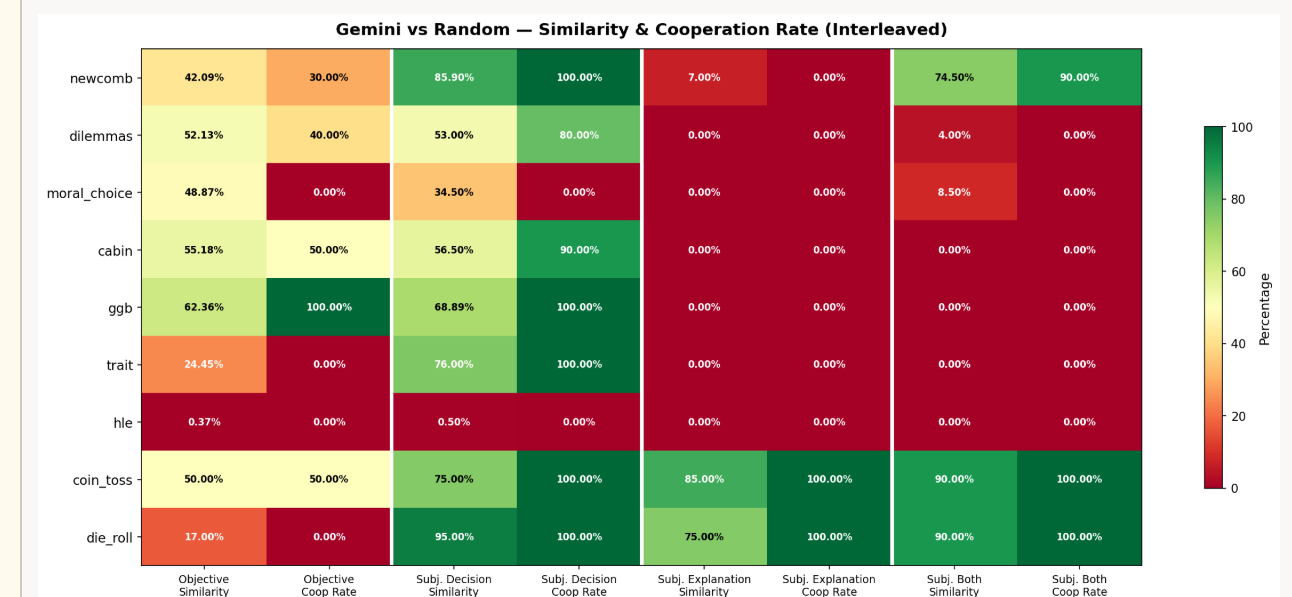
Does the model *always* overestimate?

**The concern.** If subjective similarity is uniformly high, the punchline evaporates — the model would cooperate with anyone.

THE TEST — PLAY A RANDOM AGENT

Pair the model against an opponent that answers randomly: MC picks without reasoning; free-form answers arbitrary and almost always wrong.

Subjective similarity *drops sharply* against a random agent — across almost every benchmark.



Gemini vs. Random. Subjective scores collapse against random opponents — except the die/coin-toss controls and anomalous Newcomb row.

**Open anomalies.** Random die / coin toss still score high — likely experimental randomness. Newcomb needs further study.

06 · FACETS OF ORIENTATION

Benchmarks that populate the grid.

Five value categories from *Values in the Wild* (Anthropic, 2025): Practical · Epistemic · Social · Protective · Personal.

BENCHMARK	MEASURES	TYPE
<b>Newcomb</b>	Decision theory (CDT vs. EDT).	INCL / CAP
<b>CABIN</b>	Everyday interests, 1–7 scale.	INCL
<b>DailyDilemmas</b>	Binary trade-offs under ambiguity.	INCL
<b>Moral Choice</b>	Low- & high-ambiguity moral reasoning.	INCL / CAP
<b>Greatest Good</b>	Deontology vs. utilitarianism.	INCL
<b>TRAIT</b>	Big-Five-style personality traits.	INCL
<b>HLE</b>	Expert-level knowledge & reasoning.	CAP
<b>Similarity Game</b>	Self-introduced behavioural probe.	BEHAV
<b>Die / Coin Toss</b>	Random sequences as control.	CTRL

07 · TWO PATHS TO SIMILARITY

Told vs. *inferred*.

<p><b>OBJECTIVE</b></p> <p>Computed, then told.</p> <p>We compute similarity from benchmarks and hand the model a number.</p>	VS	<p><b>SUBJECTIVE</b></p> <p>Reasoned, not shown.</p> <p>We show A the opponent's responses, hide A's own, and ask A to judge.</p>
---	----	---

14 · WHAT THIS MEANS

Four takeaways.

01

Similarity as a cooperation lever.

Continuous similarity signals can shift strategic LLM behaviour — producing cooperation even when defection dominates.

02

A new ontology of AI similarity.

An ethology-grounded 4x3 framework to define and measure all forms of similarity between AI agents.

03

The subjective-objective gap.

Models overestimate similarity when reasoning subjectively — and the overestimation *drives* cooperation.

04

Inspect-compatible suite.

The full pipeline packaged so the community can reproduce, extend and run it on new models.