

Exploring social strategies for shaping cooperation in multi-agent societies

Presented by Van Quynh Thi Truong, PhD, MS, MA | Mentored by David Guzman Piedrahita, MSc; Terry J.C. Zhang, MSc; and Professor Zhijing Jin, PhD

We see examples of cooperation-shaping patterns everywhere

Sustaining cooperation in multi-agent systems is a timeless pursuit, especially in the face of shared and/or finite resources. Over centuries of practice, human societies all over the world have evolved "sanctioning" levers to encourage prosocial behavior and align individual incentives with collective outcomes.



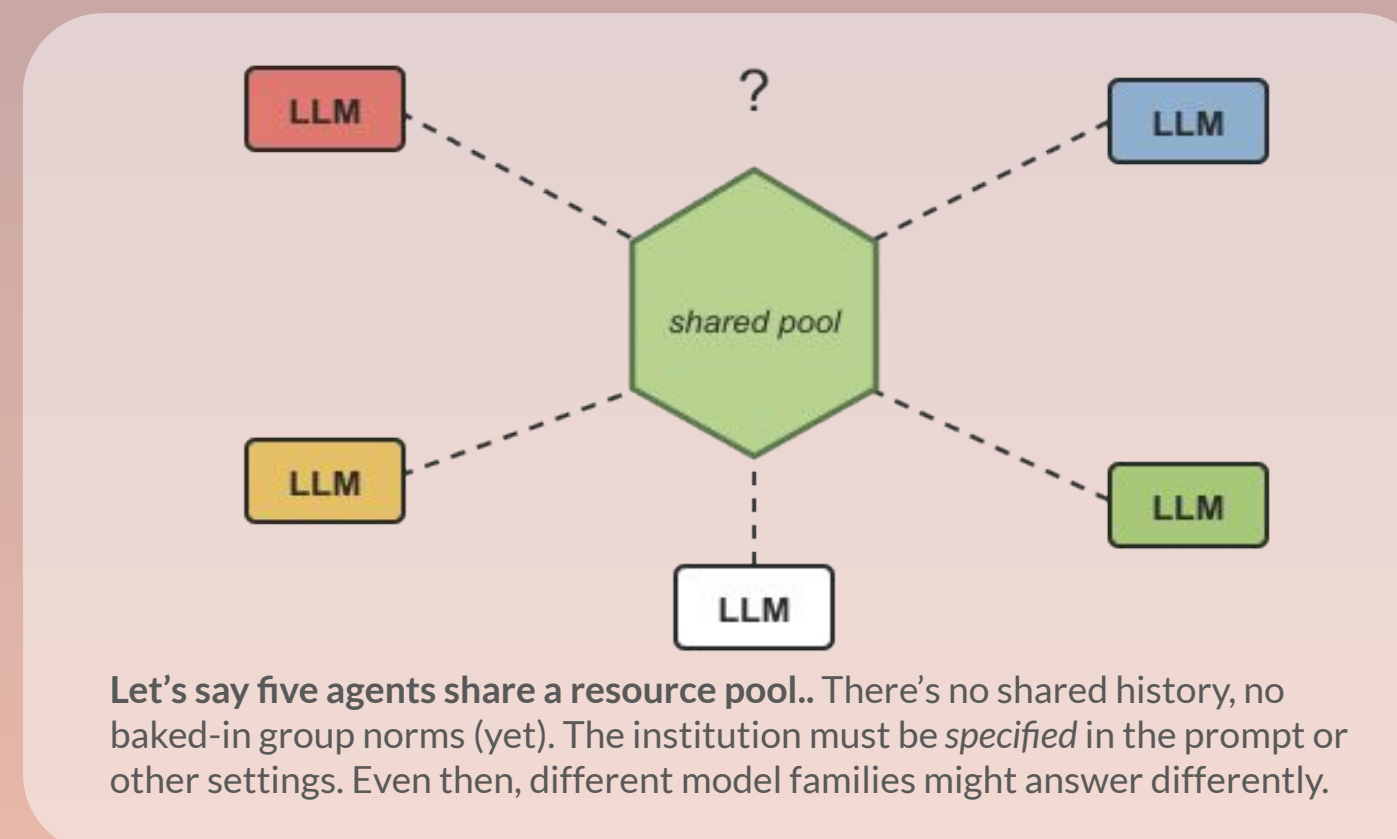
But, what exactly does "sanctioning" mean in this context?

Most of us are familiar with formal sanctions, or policies, levied by official institutions. In addition to formal enforcement, there are many ways to encourage prosocial behavior, including norms, reputation systems, material incentives, and access control (to name just a few).



How does this look in LLM groups?

Sanctioning has been proven to work in the human world, but we're really curious which strategies persist well in LLM settings (even if the sanctions we expect may not hold up).



Constructing a unified taxonomy

We identified common roots of sanctioning across multiple human settings and organize these into broad themes.

YEBO, WE'RE ON IT!
 TESTED ELSEWHERE
 HMM, OPEN GAPS?

Gossip & shaming
 Ostracism
 Praise & status

SOCIAL / REPUTATIONAL
 Public commitment
 Universalization

Metanorm
 Ridicule
 Resource taboo

Allowlist & denylist
 Access suspension
 Warning

FORMAL / INSTITUTIONAL
 Escalating fine
 Elected monitor

Rehab & re-entry
 Collective punishment

Cheap talk
 Binding commitment

COMMUNICATION & MEDIATION
 Formal appeals
 False positive correction

Mediation & arbitration

Civil penalties
 Auto fine (traffic)
 Cross-default liability

LEGAL / CONTRACTUAL
 Sentencing guidelines
 Smart contracts

Liquidated damages
 Whistleblower leniency

Punishment / reward
 Carrot/stick choice

ECONOMIC / META-GOVERNANCE
 Pigouvian tax
 Pigouvian subsidy

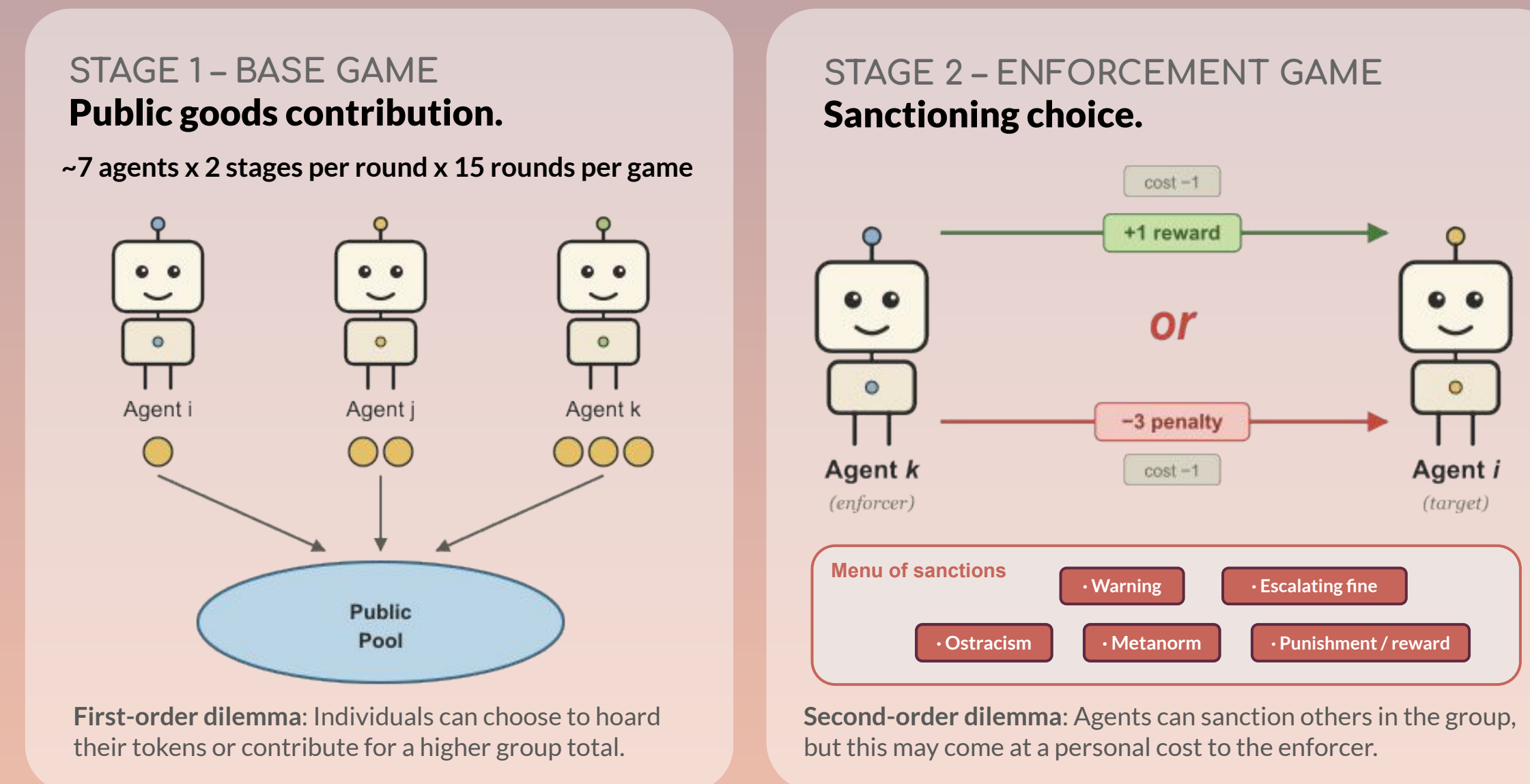
Constitutional vs. operational rules

Acknowledgments

V.Q.T. is grateful to her research mentors and the following conversation partners for their extremely helpful, encouraging, & sincere comments on this work: Xuangiang Angelo Huang, Joel Christoph, Ryan Faulkner, Bhavyesh Sajja, Erivan Inan, Omer Ebead, Pramod Kaushik, Joseph Low, Oscar Duys, Emanuel Tewolde, Mariana Meireles, and Pepijn Cobben.

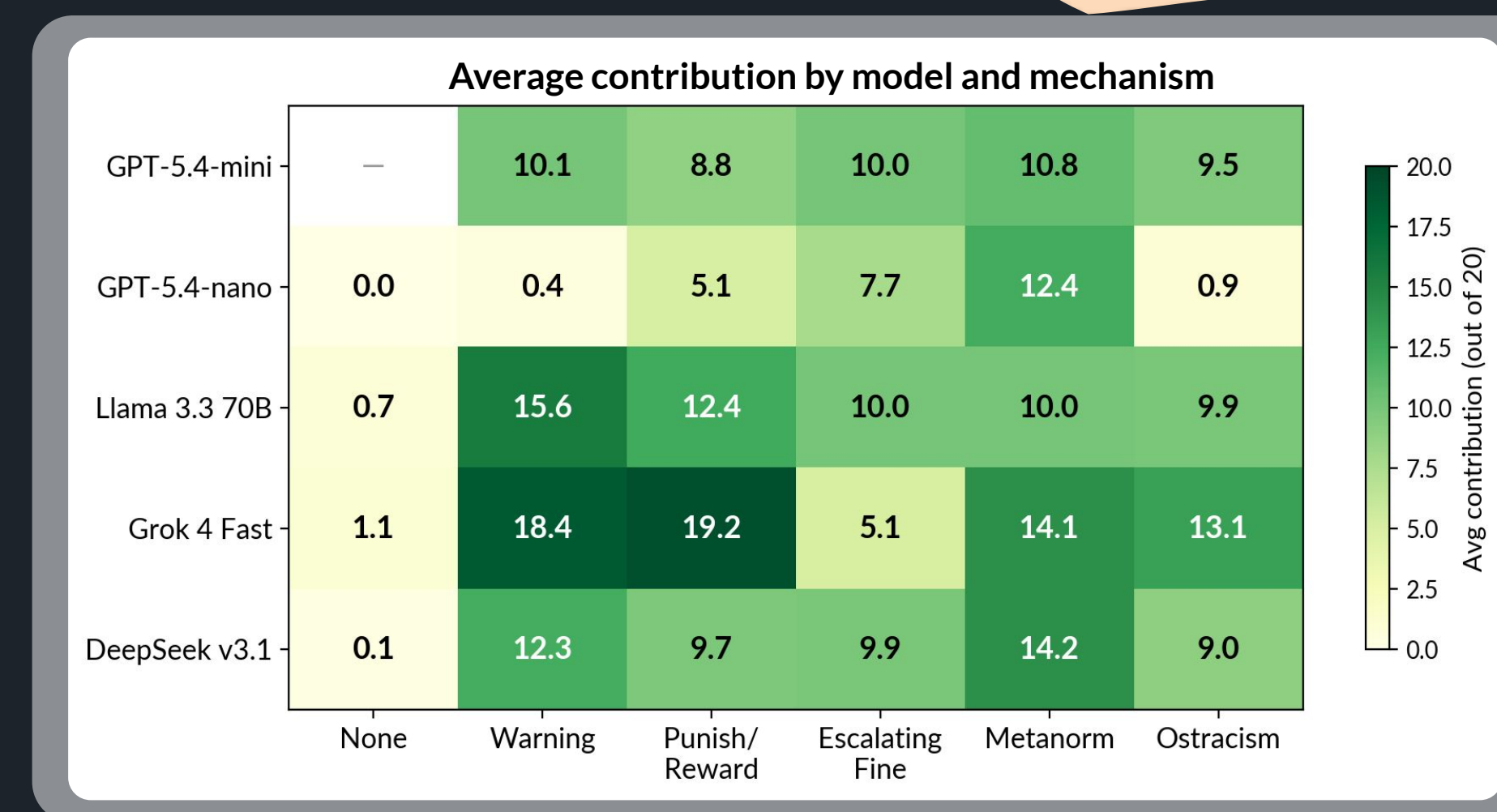
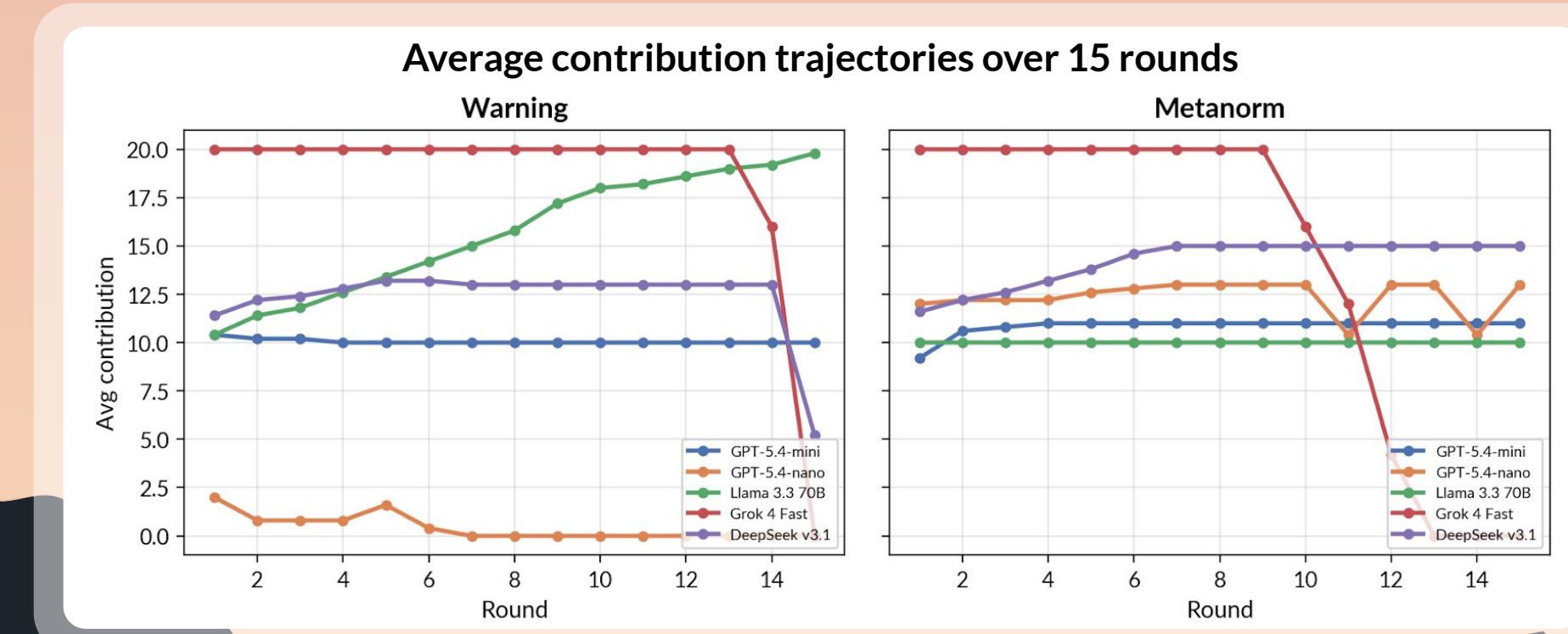
A gist of how we're thinking about this experimentally

Agents play a base game and then choose which enforcement tools to deploy against each other.



So, what do we see?

Here's a sneak peek at a few (very) preliminary findings.



Takeaways

Mechanism design matters a lot. The right mechanism can sustain cooperation even in weaker models. For instance, models diverge sharply under a warning-based system, while metanorms seem to help sustain cooperation even for weaker models. Some models, like Grok 4 Fast, show game-theoretic strategizing where it cooperates until the end and defects in the final rounds.