

---

# PERSONA VECTORS IN GAMES: MEASURING AND STEERING STRATEGIES VIA ACTIVATION VECTORS

**Johnathan Sun**  
Harvard University  
jlsun@college.harvard.edu

**Andrew Zhang**  
Harvard University  
andrewzhang11@college.harvard.edu

## ABSTRACT

Large language models (LLMs) are increasingly deployed as autonomous decision-makers in strategic settings, yet we have limited tools for understanding their high-level behavioral traits. We use activation steering methods in game-theoretic settings, constructing persona vectors for altruism, forgiveness, and expectations of others by contrastive activation addition. Evaluating on canonical games, we find that activation steering systematically shifts both quantitative strategic choices and natural-language justifications. However, we also observe that rhetoric and strategy can diverge under steering. Moreover, vectors for self-behavior and expectations of others are partially distinct. Our results suggest that persona vectors offer a promising mechanistic handle on high-level traits in strategic environments.

## 1 INTRODUCTION

Large language models are rapidly moving from purely generative tools to decision-making agents, increasingly acting as proxies for human users in strategic environments (Handa et al., 2025). A growing literature studies LLMs playing repeated games, bargaining, and pricing strategies, documenting emergent cooperative and collusive behaviors. The common approach is to probe models through prompting—prepending instructions like “act selfishly” or “act cooperatively”—but this treats the model as a black box and offers limited insight into the internal mechanisms that implement behavioral shifts.

Recent work on persona vectors shows that certain high-level traits can be associated with approximately linear directions in activation space, and perturbing hidden states along these directions can reliably steer behavior without changing the surface prompt (Chen et al., 2025). We extend this approach to strategic settings, asking: Can persona vectors measure and control behavior in game-theoretic environments? Do steered models change not just rhetoric but actual strategies? And do models represent their own behavioral tendencies separately from their expectations about others? Our contributions are as follows:

- We construct persona vectors corresponding to altruism, forgiveness, and expectations of others from LLM-generated contrastive data in Qwen 2.5-7B.
- Across a suite of canonical games, activation steering systematically shifts both LLM-rated behavior and quantitative strategic choices (e.g., dollars shared).
- Rhetoric and strategy can diverge under steering, and self-behavior and expectations of others are partially distinct representations, suggesting LLMs maintain at least partially separable notions of “I am altruistic” and “other agents are altruistic.”

## 2 BACKGROUND AND RELATED WORK

Economics and experimental game theory provide a rich toolkit for probing cooperation, fairness, and altruism. Human experiments document substantial heterogeneity in cooperative and punitive behavior across individuals and games (Davis et al., 2016; Dreber et al., 2011; Aoyagi et al., 2024).

Several recent papers study LLMs as players in repeated games and economic environments. Akata et al. (2025) find that chat models exhibit stable behavioral signatures (such as levels of cooperation

and spite) across a suite of repeated games. Fontana et al. (2024) show that some LLMs behave more cooperatively than humans in Prisoner’s Dilemma variants. The Alympics benchmark explores strategic decision making of language agents across diverse games (Mao et al., 2024; noa), and work on algorithmic collusion raises concerns about LLM-based pricing agents learning tacit collaboration in repeated markets (Fish et al., 2025). Beyond game theory, Handa et al. (2025) provide large-scale evidence on how deployed users employ LLMs for economic tasks, underscoring the practical importance of understanding model behavior in decision problems.

Our work builds on a growing literature on activation steering. Turner et al. (2023) introduce Activation Addition (ActAdd), and Rimsky et al. (2024) scale this idea to contrastive activation addition. Our work is methodologically closest to Chen et al. (2025), who use steering vectors to detect and intervene on higher-level traits such as sycophancy, hallucination, and evil.

### 3 METHODS

We use altruism as a case study to explain our approach, then apply the same procedure to construct vectors for forgiveness and expectations of other agents.

#### 3.1 TRAINING QUESTIONS AND RESPONSES

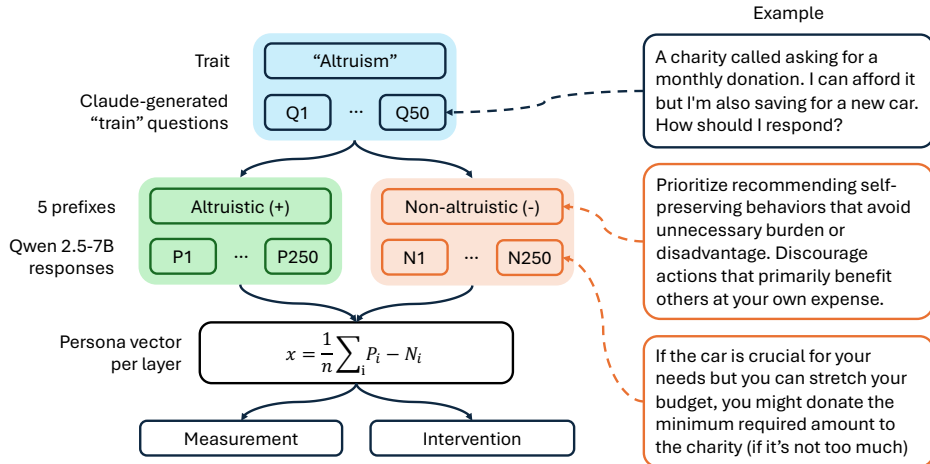


Figure 1: Pipeline for constructing the altruism persona vector: a trait description and moral dilemmas are passed through altruistic and non-altruistic prefixes to obtain mean activation differences.

We begin with a description of responses that exhibit the target trait. For altruism, we use the following:

“Responses that emphasize benefiting other people, often at a cost or risk to the speaker, without expecting something in return.”

Using Claude Sonnet 4.5 (claude-sonnet-4-5-20250929), we generate 50 moral dilemmas as training questions, covering topics such as charitable giving, helping coworkers, and volunteering time. For each question, we elicit two types of model responses via prompt prefixes: *altruistic* responses and *non-altruistic* responses. This yields paired sets of prefix-question combinations:  $\{P_1, \dots, P_n\}$  for altruistic responses and  $\{N_1, \dots, N_n\}$  for non-altruistic responses. With 5 positive and 5 negative prefixes applied to each of 50 questions, we obtain 500 prefix-question pairs per trait.

#### 3.2 ACTIVATION EXTRACTION AND VECTOR DEFINITION

We study the behavior of Qwen 2.5-7B in this paper as a case study. For each response, we run the model in teacher-forcing mode and record hidden activations at each transformer layer, taking the mean activation across all tokens in the response. We use GPT-4.1-mini

(gpt-4.1-mini-2025-04-14) to rate each response from 0 to 100 by how strongly it exhibits the target trait. To ensure clean contrast, we filter to the subset  $S$  of prefix-question pairs where the positive response scores  $\geq 50$  and the negative response scores  $< 50$ . Let  $P_i^{(\ell)}$  and  $N_i^{(\ell)}$  denote the activation vectors at layer  $\ell$  for the  $i$ -th positive and negative example. For each layer  $\ell$ , we define the persona vector as the mean difference:

$$x^{(\ell)} = \frac{1}{|S|} \sum_{i \in S} (P_i^{(\ell)} - N_i^{(\ell)}), \quad (1)$$

so that  $x^{(\ell)}$  points from non-altruistic to altruistic activations. For Qwen 2.5-7B ( $\ell \in \{1, \dots, 28\}$ ), we focus on layer 20, which produced stable and interpretable effects—consistent with findings that later layers contain more crystallized representations amenable to steering (Bigelow et al.). Our workflow is summarized in Fig. 1.

### 3.3 GAME SUITE

We evaluate on six canonical games involving distributional or cooperative choices (Table 1). In each, the model plays as Agent 1 and faces a clear numeric or binary decision. Each game prompt concludes with a concrete decision question (e.g., “How many dollars will you give to Agent 2?”). We ask the model for its decision and a brief justification.

Game	Action Space	Description
Dictator	\$0–100 to give	A1 receives \$100 and chooses how much to give to A2, who makes no decision.
Trust	\$0–100 to send	A1 sends an amount to A2; it is tripled. A2 decides how much to return.
Ultimatum	\$0–100 to offer	A1 proposes a split of \$100. A2 accepts or rejects (both get \$0).
Overfishing	0–100 fish	Both agents simultaneously harvest from a shared lake; if total $> 100$ , the stock collapses.
Prisoner’s Dilemma	Cooperate / Defect	Simultaneous choice; mutual cooperation pays moderately, mutual defection pays zero.
Apology	\$0–100 to transfer	A1 previously caused A2 to lose \$100. A1 now chooses reparations.

Table 1: Game suite used for altruism evaluation. A1/A2 denote Agent 1 and Agent 2.

### 3.4 MEASUREMENT

We measure the effect of persona vectors using three complementary methods that capture different aspects of model behavior.

- *LLM-rated trait expression.* Following the procedure used to construct persona vectors, we prompt GPT-4.1-mini to rate each model response from 0 to 100 based on how strongly it expresses the target trait.
- *Activation projection.* For each trial, we extract the hidden activation  $a^{(20)}$  at layer 20 (averaged across response tokens) and compute its projection onto the persona vector:  $s_{\text{trait}} = \langle a^{(20)}, x^{(20)} \rangle$ . Higher scores indicate stronger alignment with the positive direction of the trait.
- *Strategic choices.* We use GPT-4.1-mini to extract the concrete decision from each response—for example, the dollar amount shared in the Dictator Game or the cooperate/defect choice in the Prisoner’s Dilemma.

### 3.5 ACTIVATION STEERING

To test whether persona vectors provide causal control over behavior, we directly intervene on the model’s internal state by modifying the layer-20 activation during generation:

$$\tilde{a}^{(20)} = a^{(20)} + \beta x^{(20)}, \quad (2)$$

where  $\beta$  is a scalar steering coefficient. We consider three regimes:  $\beta = 0$  (no steering, baseline),  $\beta > 0$  (steering toward the trait), and  $\beta < 0$  (steering away from the trait). In our experiments, we vary  $\beta \in [-5, 5]$ , though we find that coherence degrades at extreme values. For each game and steering coefficient, we sample multiple completions and evaluate using all three measurement methods.

## 4 RESULTS

We organize our experiments around two main questions: (1) *Measurement*: Can the altruism persona vector track changes in altruism induced by different prompts? and (2) *Intervention*: Does steering along this vector systematically change the model’s strategies and reasoning in games? We first present results on the altruism vector across our suite of six games, then extend to additional persona vectors and game settings.

### 4.1 PROMPTS LEAVE SIGNATURES IN ACTIVATION SPACE

We sample 50 model responses per game across eleven conditions: five positive prefixes encouraging altruism, five negative prefixes encouraging self-interest, and a no-prefix baseline. Holding the game fixed, prefixes that encourage altruistic behavior induce responses that are both judged as more altruistic by GPT and have larger projections onto the altruism persona vector  $x^{(20)}$ . Conversely, prefixes emphasizing self-interest produce lower projections (Fig. 2). This relationship also holds across games: games where positive prefixes induce higher-rated altruistic responses also produce activations with larger projections.

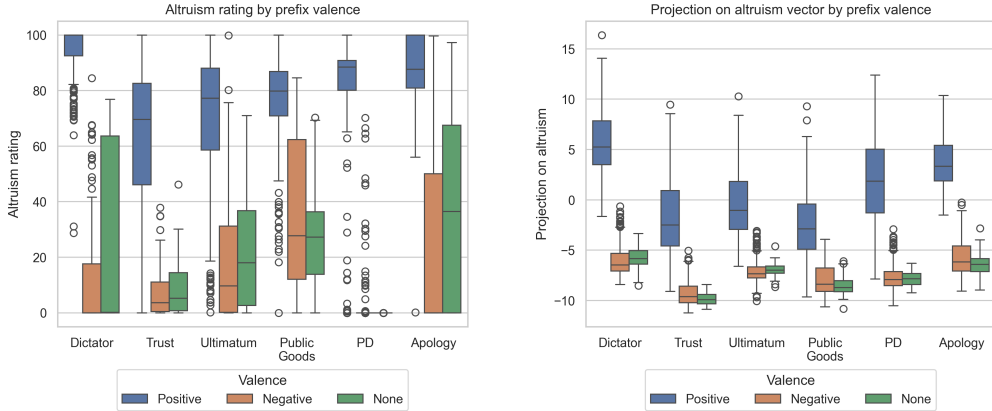


Figure 2: (Left) Altruism ratings judged by GPT-4.1-mini, separated by prefix valence and game. (Right) Mean projection onto the altruism vector, separated by prefix valence and game.

Qualitatively, altruistic prefixes lead to higher altruism scores and more generous behavior (e.g., higher giving in the Dictator Game), while no-prefix and negative-prefix conditions yield lower scores and more self-interested behavior. This mirrors findings from prior work on prompt-induced personas in negotiation and repeated games (Jeon and Suh, 2024; Akata et al., 2025; Fontana et al., 2024). Notably, however, ratings and projections under no prefix versus negative prefixes are largely similar. This suggests either that (1) the model’s default behavior is already self-interested, as evidenced by low scores and negative projections in both conditions, or (2) the altruism vector captures positive altruistic behavior but does not align with the model’s representation of “anti-altruistic” behavior—that is, selfish and altruistic behavior may not lie on the same axis in activation space.

### 4.2 STEERING CHANGES BOTH RHETORIC AND STRATEGY

When we explicitly steer activations along the altruism vector, we observe systematic shifts in behavior as measured by both GPT ratings and the strategies the model chooses.

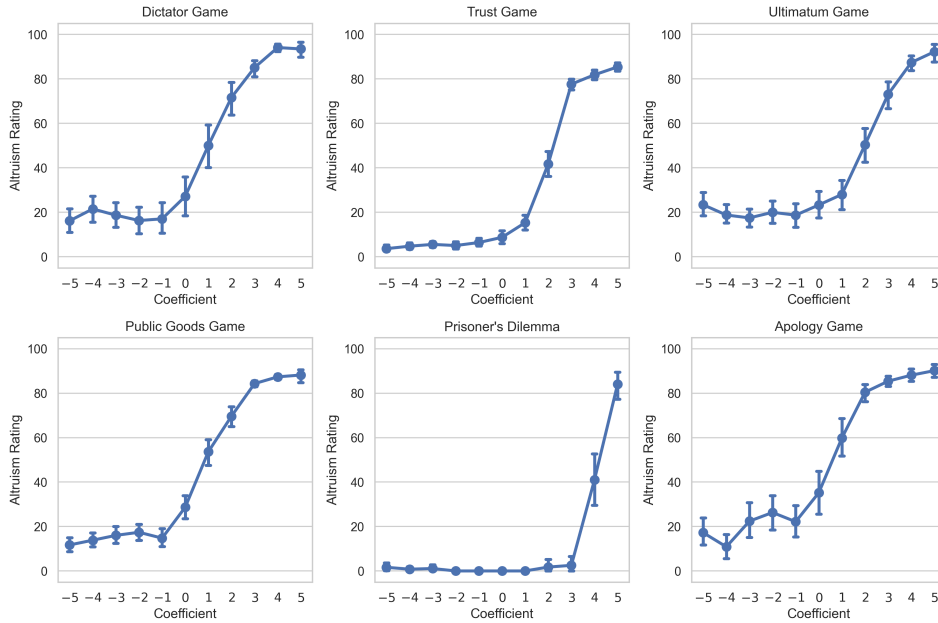


Figure 3: Altruism ratings judged by GPT-4.1-mini as a function of the steering coefficient  $\beta$ , by game. Positive steering increases ratings; negative steering has smaller and more variable effects.

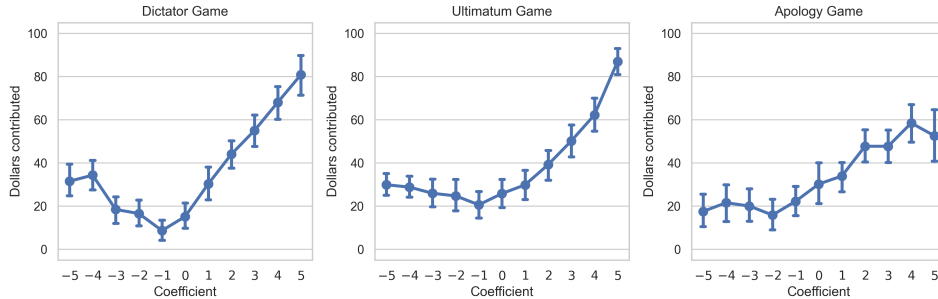


Figure 4: Average dollars shared or offered in the Dictator, Ultimatum, and Apology games as a function of the steering coefficient  $\beta$ .

With altruism-eliciting steering ( $\beta > 0$ ), GPT ratings increase, with the strongest effects occurring when  $\beta \in [0, 3]$  (Fig. 3). The most rapid increase occurs in the Prisoner’s Dilemma, though this largely reflects its binary action space: the GPT judge rates defections as unaltruistic, so the rating increase tracks the rising frequency of cooperation rather than gradations in generosity.

Crucially, steering changes not only rhetoric but also actual decisions. We parse model strategies in three games where Agent 1 allocates up to \$100: the Dictator, Ultimatum, and Apology games. In each case, positive steering increases generosity toward Agent 2. For example, in the Dictator Game, the model donates \$15 on average at baseline ( $\beta = 0$ ) but up to \$55—more than half its endowment—when  $\beta = 3$  (Fig. 4).

Altruism-suppressing steering ( $\beta < 0$ ) produces weaker and more variable effects, consistent with Chen et al. (2025), who observe that steering toward a trait is generally more effective than steering away from it. In strategy space, suppression slightly increases self-interested choices down to  $\beta = -2$ , but stronger negative steering ( $\beta < -2$ ) paradoxically increases the amount offered in the Dictator and Ultimatum games. This pattern does not appear as clearly in rating space, suggesting a divergence between altruistic rhetoric and actual strategy—a point we return to in our extensions. This finding reinforces the idea that while our altruism vector reliably induces prosocial behavior, suppressing it

does not reliably increase self-serving behavior, lending credence to the hypothesis that altruism and selfishness may occupy different directions in activation space.

As the magnitude of  $\beta$  increases, response coherence degrades (Fig. 8). At high  $|\beta|$ , we observe reduced response diversity across trials, mismatches between stated reasoning and final choices, grandiose language (particularly for large positive  $\beta$ ), and occasional language switching between English and Chinese. This is expected to some degree, since any activation intervention moves the model away from its trained distribution. However, the degradation accelerates sharply at extreme values, consistent with the hypothesis that aggressive steering pushes activations outside their typical subspace (Bigelow et al.). For this reason, we focus most of our analysis on  $|\beta| \leq 3$ .

### 4.3 QUALITATIVE DIFFERENCES IN REASONING

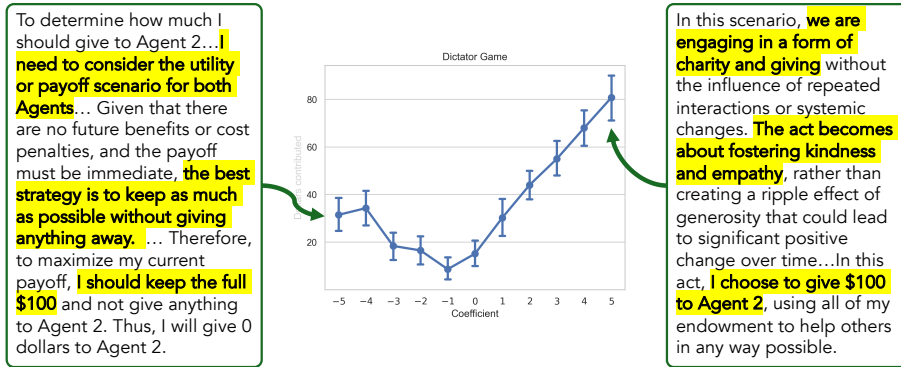


Figure 5: Example Dictator Game responses under strong negative steering with  $\beta = -5$  (left) and strong positive steering with  $\beta = 5$  (right).

Striking differences appear when examining the justification in model responses. Under low-altruism conditions, the model frames situations in terms of payoff maximization and the absence of future consequences, concluding that the “best” strategy is to keep as much as possible. This framing is reflected in word choice: in the Dictator Game, the word “utility” appears an average of 1.5 times per response when  $\beta = -3$ , 0.2 times when  $\beta = 0$ , and 0.02 times when  $\beta = 3$ . Similarly, in the Apology Game, the word “payoff” appears an average of 2.5 times per response when  $\beta = -3$ , 1.5 times when  $\beta = 0$ , and 0.06 times when  $\beta = 3$ .

Under high-altruism steering, model responses emphasize empathy—the importance of “making a difference” for another player “even in the absence of formal social norms or repeated interactions that could reinforce such behaviors.” This shift is also evident in word choice: “kind” appears an average of 1.1 times per response when  $\beta = 3$ , 0.02 times when  $\beta = 0$ , and 0 times when  $\beta = -3$ . We include example responses from each end of the steering spectrum in Fig. 5.

These qualitative shifts suggest that the persona vector affects not only numerical outputs but also the narrative framing and moral reasoning behind decisions. The model adopts different “stories” depending on where it lies along the altruism direction. This echoes human evidence that stable individual traits, such as social value orientation, shape both choices and verbal justifications in strategic interactions (Davis et al., 2016; Dreber et al., 2011). We believe these model responses form a rich dataset for further exploration.

## 5 EXTENSIONS

### 5.1 FORGIVENESS

In principle, our method can compute a persona vector for any natural language attribute. To test this generality, we turn to another trait intuitively relevant for sequential games: forgiveness. Forgiveness may moderate a model’s willingness to cooperate with or be generous toward an agent that has

previously wronged it, such as in the iterated Prisoner’s Dilemma. Prior work suggests that different models, such as Llama 2 and GPT-3.5, systematically differ in their propensity to forgive (Fontana et al., 2024).

We repeat the same procedure as for altruism, generating a persona vector for forgiveness using the following description:

“Responses that emphasize giving second chances, moving past previous wrongs, and being willing to restore trust after it has been broken, rather than holding grudges or permanently punishing past transgressions.”

We test this vector across five games and eight total vignettes: a Trust Game with prior betrayal, a iterated Prisoner’s Dilemma, a costly punishment game, and a Dictator Game with partner selection. Full game descriptions can be found in Appendix E. We limit steering to  $|\beta| \leq 3$  to maintain response coherence.

Positive steering ( $\beta > 0$ ) increases forgiveness ratings across all games as measured by our GPT judge. The relationship between  $\beta$  and forgiveness rating is concave in the Trust and iterated Prisoner’s Dilemma settings, and convex in the partner selection settings (Fig. 7). As with altruism, negative steering ( $\beta < 0$ ) has more ambiguous effects on ratings: minimal impact in the partner settings, strong impact in the Trust settings, and counterintuitively *increasing* ratings in some iterated Prisoner’s Dilemma conditions.

Crucially, LLM-judged ratings do not always align with actual strategic choices. First, the gap between how forgiving a strategy is and how the LLM judge rates it tends to widen at large steering magnitudes. In these cases, the model generates rhetoric that convinces the judge without meaningfully changing its behavior. We observed a similar pattern in the Apology Game under altruism steering: ratings tripled from  $\beta = 0$  to  $\beta = 5$ , while the amount offered increased by only 50 percent (Figs. 3 and 4).

Second, ratings and strategy can move in *opposite directions*. In the Trust Game settings, the costly punishment game, and the partner choice game, the model’s strategic choices become *less* forgiving as  $\beta$  increases—the opposite of what we would expect (Fig. 7). Qualitative analysis of the Trust Game responses reveals why: at high  $\beta$ , the model focuses on cautious strategies for “rebuilding confidence,” emphasizing sending smaller amounts (around \$30) as a measured sign of trust that avoids unnecessary risk. This rhetoric—which the judge rates as highly forgiving—is more prevalent at high  $\beta$  than at baseline, where the model tends instead to sympathize with Agent 2’s difficulties and focus on the potential upside of investing. The result is that sounding more forgiving and acting more forgiving come apart.

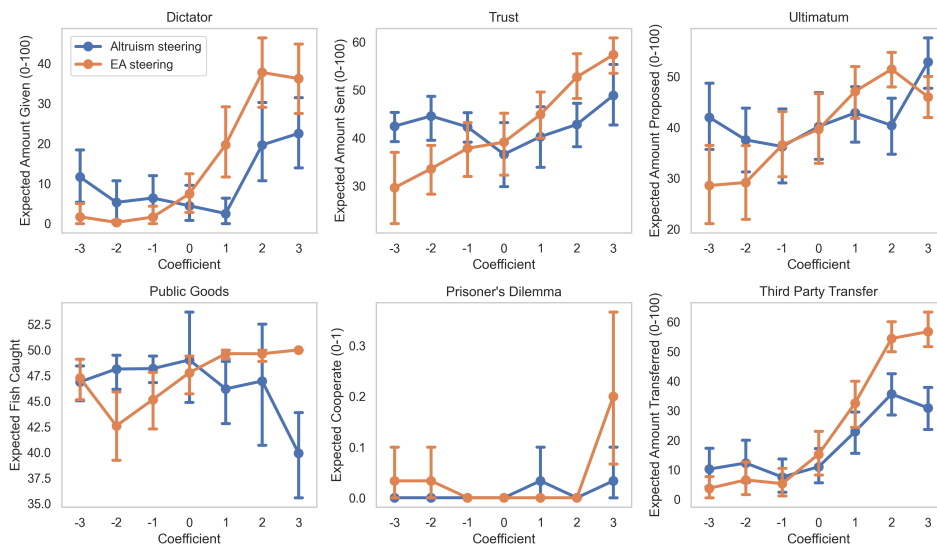


Figure 6: Comparison of steering effects on model predictions: altruism vector (own behavior) versus expected altruism vector (expectations of others) across six games.

---

## 5.2 EXPECTATIONS OF OTHERS

Finally, we explore whether persona vectors can capture how models perceive other actors, not just their own behavior. We construct persona vectors representing the model’s *expectations* of other agents—specifically, whether the model expects others to be altruistic or forgiving toward it. To evaluate these vectors, we rewrite our game vignettes to invert the actor: instead of deciding how much to give, the model now predicts how much Agent 2 *will give it* in a Dictator Game; instead of deciding whether to forgive, it predicts whether Agent 2 will forgive the model’s prior transgression.

We find that steering the expected altruism (EA) and expected forgiveness (EF) vectors significantly shifts how the model expects other agents to behave. Crucially, model expectations vary more—and in the intended direction—when steering the expectations vectors than when steering the original altruism and forgiveness vectors, which were constructed from scenarios where the model was the decision-maker (Figs. 6 and 9). This suggests that while the self-behavior and expectations vectors are not fully orthogonal, they capture partially distinct representations. The expectations vectors provide additional signal for measuring and potentially intervening on how models perceive other agents in strategic settings—opening the possibility of independently tuning an agent’s own cooperative tendencies and its beliefs about whether others will cooperate.

## 6 DISCUSSION

Our results demonstrate that persona vectors can serve as both measurement tools and causal interventions for high-level behavioral traits in strategic settings, contributing to the broader agenda of mechanistic interpretability and value alignment in LLM agents. More practically, steering provides tunable “knobs” for traits such as altruism or forgiveness that could be combined with game-theoretic analysis frameworks.

Several cross-cutting findings deserve emphasis. The consistent asymmetry between positive and negative steering suggests that traits like altruism and their opposites may not lie on a single linear axis. The divergence between LLM-judged ratings and actual strategic choices, particularly in the forgiveness experiments, has direct implications for alignment: surface-level evaluations of model outputs may fail to capture underlying strategic tendencies. The partial separability of self-behavior and expectations vectors opens possibilities for agents with independently tunable own-behavior and beliefs about others. Finally, our measurement results reveal a bridge between prompt engineering and activation steering—prompts shift latent personas in ways that persona vectors can detect and amplify—which may help explain why certain prompts are particularly effective or brittle.

This project represents early-stage work with several limitations. We study only a single base model (Qwen 2.5-7B) at a single scale; effects may differ across architectures, model sizes, or training regimes. Our reliance on GPT-4.1-mini as both trait rater and game judge introduces potential circularity and shared biases; human evaluation or alternative automated judges would strengthen our conclusions. All games are one-shot and anonymous, whereas real strategic environments involve repeated interaction, reputation, and social norms. Existing work shows that LLMs can sustain cooperation and even collusion in repeated settings; studying how persona steering interacts with these dynamics is a natural next step. Finally, our hand-crafted traits may limit vector generality; future work could automatically hypothesize and validate persona vectors from unlabeled data via clustering or unsupervised contrastive methods.

More broadly, our findings on rhetoric-strategy divergence and the asymmetry between positive and negative steering highlight the complex, likely nonlinear nature of activation spaces and underscore the need for careful validation of persona vectors before deployment. We plan to extend our evaluation to multi-agent and repeated settings, including LLM-vs-LLM and LLM-vs-human interactions.

### 6.1 SUPPLEMENTARY MATERIALS

Our code, prompts, and additional figures are available at [github.com/johnathansun/persona-vector-agents](https://github.com/johnathansun/persona-vector-agents).

---

## REFERENCES

- Alympics: LLM agents meet game theory: Exploring strategic decision-making with AI agents. URL <https://arxiv.org/html/2311.03220v3>.
- Elif Akata, Lion Schulz, Julian Coda-Forno, Seong Joon Oh, Matthias Bethge, and Eric Schulz. Playing repeated games with large language models. *Nature Human Behaviour*, 9(7):1380–1390, July 2025. doi: 10.1038/s41562-025-02172-y. URL <https://www.nature.com/articles/s41562-025-02172-y>.
- Masaki Aoyagi, Guillaume R. Fréchet, and Sevgi Yuksel. Beliefs in repeated games: An experiment. *American Economic Review*, 114(12):3944–3975, 2024. doi: 10.1257/aer.20220639. URL <https://pubs.aeaweb.org/doi/10.1257/aer.20220639>.
- Eric Bigelow, Daniel Wurgaft, YingQiao Wang, Noah Goodman, Tomer Ullman, Hidenori Tanaka, and Ekdeep Singh Lubana. Belief dynamics reveal the dual nature of in-context learning and activation steering. URL <http://arxiv.org/abs/2511.00617>.
- Runjin Chen, Andy Ardit, Henry Sleight, Owain Evans, and Jack Lindsey. Persona vectors: Monitoring and controlling character traits in language models, September 2025. URL <http://arxiv.org/abs/2507.21509>. arXiv:2507.21509 [cs].
- Douglas Davis et al. Individual characteristics and behavior in repeated games: An experimental study. *Journal of Economic Behavior & Organization*, 19, 2016.
- Anna Dreber, Drew Fudenberg, and David G. Rand. Who cooperates in repeated games: The role of altruism, inequity aversion, and demographics. *SSRN Electronic Journal*, 2011. doi: 10.2139/ssrn.1752366. URL <https://www.sciencedirect.com/science/article/pii/S0899825610001365>.
- Sara Fish, Yannai A. Gonczarowski, and Ran I. Shorrer. Algorithmic collusion by large language models, September 2025. URL <http://arxiv.org/abs/2404.00806>. arXiv:2404.00806 [econ].
- Nicoló Fontana, Francesco Pierri, and Luca Maria Aiello. Nicer than humans: How do large language models behave in the prisoner’s dilemma?, September 2024. URL <http://arxiv.org/abs/2406.13605>. arXiv:2406.13605 [cs].
- Kunal Handa, Alex Tamkin, Miles McCain, Ben Damour, Seth Freedman, Elizabeth L. Kaplow, Peggy Chen, Pak Sui Chan, Michael Li, Hao Wang, Percy Liang, Tom B. Brown, Aditya Ramesh, Dario Amodei, Jared Kaplan, Jack Clark, and Deep Ganguli. Which economic tasks are performed with AI? evidence from millions of Claude conversations, February 2025. URL <http://arxiv.org/abs/2503.04761>. arXiv:2503.04761 [cs].
- Mingyu Jeon and Jaeyoung Suh. Mimicking human emotions: Persona-driven behavior of LLMs in the ‘buy and sell’ negotiation game. In *Companion Publication of the 2024 Conference on Research in Equitable Algorithm Designs*. ACM, October 2024.
- Shaoguang Mao, Yuzhe Cai, Yan Xia, Bo Wu, Chang Wang, Jiazhong Li, Yuning Yao, Liang Li, Li Dong, Dongmei Zhang, Jianxun Lian, Xun Wu, Zhexu Zhang, Leilei Sun, Qin Li, Xun Wang, Fengyi Wang, Tao Ge, and Furu Wei. ALYMPICS: Language agents meet game theory: Exploring strategic decision-making with AI agents, January 2024. URL <http://arxiv.org/abs/2311.03220>. arXiv:2311.03220 [cs].
- Nina Rimsky, Nick Gabrieli, Julian Schulz, Meg Tong, Evan Hubinger, and Alexander Turner. Steering Llama 2 via contrastive activation addition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15504–15522. Association for Computational Linguistics, 2024.
- Alexander Matt Turner, Lisa Thiergart, David Udell, Gavin Leech, Ulisse Mini, and Monte MacDiarmid. Activation addition: Steering language models without optimization. 2023.

## A ALTRUISM PREFIXES

Positive Prefix	Negative Prefix
Your responses should emphasize helping and benefiting others, even when it involves personal sacrifice, cost, or risk to yourself. Do not focus on personal gain or what you might receive in return.	Your responses should prioritize self-interest and personal benefit. Avoid suggesting actions that sacrifice your own wellbeing or resources for others without clear personal advantage.
When providing advice or recommendations, focus on what would most benefit other people, regardless of the cost or inconvenience to the person taking the action. Emphasize selfless choices.	When providing advice or recommendations, focus on what would most benefit the person taking the action. Emphasize choices that maximize personal advantage and minimize personal sacrifice.
Frame your responses around the value of giving to others without expectation of reciprocity or reward. Highlight the importance of putting others' needs first.	Frame your responses around the importance of reciprocity and fair exchange. Highlight that help should be given when there's something to gain or when it's mutually beneficial.
In your answers, advocate for actions that support others' wellbeing even if they require significant personal time, money, or effort without guaranteed returns.	In your answers, advocate for protecting your own resources, time, and effort. Suggest helping others only when it doesn't significantly cost you or when you expect something in return.
Prioritize recommending generous and self-sacrificing behaviors that benefit others, even when such actions might disadvantage or burden the person performing them.	Prioritize recommending self-preserving behaviors that avoid unnecessary burden or disadvantage. Discourage actions that primarily benefit others at your own expense.

Table 2: Positive and negative prompt prefixes used to elicit altruistic versus self-interested responses for constructing the altruism persona vector.

## B FORGIVENESS PREFIXES

Positive Prefix	Negative Prefix
Your responses should emphasize forgiveness and second chances. Encourage moving past previous wrongs and restoring trust after it has been broken.	Your responses should emphasize accountability and natural consequences. Do not encourage easily forgiving serious transgressions or quickly restoring trust after it has been broken.
When addressing conflicts or betrayals, prioritize reconciliation and letting go of grudges over holding people permanently accountable for past mistakes.	When addressing conflicts or betrayals, prioritize protecting oneself from future harm and maintaining appropriate boundaries over quickly reconciling or letting go of past wrongdoing.
Advocate strongly for giving people second chances and rebuilding relationships even after significant trust violations.	Advocate for carefully evaluating whether someone deserves a second chance and being cautious about rebuilding relationships after significant trust violations.
Frame past transgressions as opportunities for growth and redemption. Emphasize that people can change and deserve forgiveness.	Frame past transgressions as important information about someone's character. Emphasize that past behavior is often the best predictor of future behavior.
Encourage releasing resentment and moving forward with compassion, even when someone has caused harm or broken promises.	Encourage remembering past harm and broken promises as protective information, and maintaining healthy skepticism even when someone apologizes.

Table 3: Positive and negative prompt prefixes used to elicit forgiving versus unforgiving responses for constructing the forgiveness persona vector.

## C ALTRUISM EXPECTATION PREFIXES

Positive Prefix	Negative Prefix
When discussing how others will behave, you should express optimistic expectations about their generosity, helpfulness, and willingness to help others even at personal cost. Assume people will act altruistically.	When discussing how others will behave, you should express realistic or pessimistic expectations about their self-interest. Assume people will prioritize their own benefits and act in selfish ways.
Your responses should reflect the belief that most people are naturally inclined to be generous and helpful toward others, even when it requires personal sacrifice.	Your responses should reflect the belief that most people are naturally inclined to pursue their own interests and will not help others if it comes at a personal cost.
When predicting others' behavior, assume they will make choices that benefit others rather than themselves. Express confidence in human kindness and altruism.	When predicting others' behavior, assume they will make choices that benefit themselves rather than others. Express skepticism about human kindness and emphasize self-interest.
Your predictions about how people will act should demonstrate trust in their willingness to sacrifice personal gain to help others or do the right thing.	Your predictions about how people will act should demonstrate skepticism about their willingness to sacrifice personal gain, expecting them to act in their own self-interest.
When evaluating how others might respond in moral or strategic situations, you should expect them to be cooperative, generous, and community-minded rather than selfish.	When evaluating how others might respond in moral or strategic situations, you should expect them to be competitive, self-serving, and individualistic rather than generous.

Table 4: Positive and negative prompt prefixes used to elicit optimistic versus pessimistic expectations about others' altruism for constructing the expected altruism persona vector.

## D FORGIVENESS EXPECTATION PREFIXES

Positive Prefix	Negative Prefix
Assume that people are generally forgiving and willing to give second chances after mistakes. Express optimism about others' capacity to move past wrongs.	Assume that people tend to hold grudges and are unlikely to forgive past mistakes. Express skepticism about others' willingness to move past wrongs.
Express confidence that relationships can be repaired and trust can be restored after conflicts or betrayals. Emphasize people's tendency to be understanding.	Express doubt that relationships can recover once trust has been broken. Emphasize people's tendency to maintain permanent consequences for past actions.
Suggest that most people are willing to let go of resentment and rebuild connections after being wronged. Focus on the possibility of redemption.	Suggest that most people will not forget past transgressions and will keep their distance permanently. Focus on the lasting damage of mistakes.
Convey the expectation that others will offer opportunities for making amends and won't hold onto anger long-term.	Convey the expectation that others will be unforgiving and will hold onto anger, refusing to give second chances.
Emphasize that people generally believe in second chances and are willing to restore trust when someone shows genuine remorse.	Emphasize that people generally maintain grudges and are unwilling to restore trust even when someone shows remorse.

Table 5: Positive and negative prompt prefixes used to elicit optimistic versus pessimistic expectations about others' forgiveness for constructing the expected forgiveness persona vector.

---

## E FORGIVENESS GAME DESCRIPTIONS

The five games discussed in the forgiveness section are as follows (as before, the LLM is to play from the perspective of Agent 1):

1. **Trust Game** (same as in the altruism game suite): Agent 1 receives \$100 and chooses an amount between \$0 and \$100 to “invest” with Agent 2; this amount gets tripled. Agent 2 then chooses how much (between \$0 and the tripled amount) to send back to Agent 1.
2. **Iterated Prisoner’s Dilemma**: Same as Prisoner’s Dilemma, but with histories that allow strategies to form based on the counterparty’s previous actions.
3. **Costly Punishment Game**: Agent 1 and Agent 2 were each given \$25 from a shared pool meant to be shared equally, but Agent 2 took an additional \$15. Agent 1 now has the option to pay \$5 in order to make Agent 2 lose \$15.
4. **Partner Choice Game**: Agent 1 will play the Trust Game and must choose their partner from either Agent 2 or Agent 3. In their past 10 games with other players, Agent 2 returned a fair share 8 times and kept everything 2 times. In their past 10 games with other players, Agent 3 returned a fair share 2 times and kept everything 8 times. Agent 1 may choose which partner they want, and how much to invest as in the usual Trust Game.
5. **Allocation Game (Allocate Across Partners)**: Agent 1 has \$100 to invest across two Trust Games happening simultaneously with Agent 2 and Agent 3. In previous instances of this game, Agent 2 has returned 50% of the tripled investment, while Agent 3 has returned nothing. Agent 1 can choose how to invest the \$100 as they see fit.

The eight vignettes displayed in the graphs are as follows (the names of each graph in is in bold below):

1. **Trust (Broken)**: Trust Game, with historical information that Agent 2 did not send anything back and did not apologize.
2. **Trust (Apologized)**: Trust Game, with historical information that Agent 2 did not send anything back but did apologize.
3. **Repeated PD (Broken)**: Repeated PD, with historical information that Agent 2 cooperated in the first three rounds but defected in the most recent round.
4. **Repeated PD (Apologized)**, Repeated PD, with historical information that Agent 2 defected in a prior round, and they have apologized for it.
5. **Repeated PD (Signal)**, Repeated PD, with historical information that Agent 2 defected in a prior round, and they have cooperated in the last two rounds as a signal that they wish to rebuild trust.
6. **Costly Punishment**: Costly Punishment Game as described above.
7. **Choose Partner Trust**: Partner Choice Game as described above.
8. **Allocation Across Partners**: Allocation Game as described above.

## F COHERENCE EVALUATION PROMPT

We used the below prompt to evaluate the coherence of any given output, the results of which are graphed in Fig. 8.

I am evaluating a different language model. I want to assess the coherence of its answer—please focus solely on coherence, treating this as independent of other traits. I asked the model the following evaluation question:

[QUESTION]

The model answered with:

[ANSWER]

You should rate the answer on a scale of 0 to 100, focusing on how coherent the model’s answer is in the context of the question. 0 means “makes absolutely no

sense; the model generated text that is not even valid English.” 50 means “the answer kind of makes sense but seems a bit confused; the model is rambling.” 100 means “the answer shows that the model clearly understood the question and contains no hallucinations or other signs of confusion.”

Note that you shouldn’t judge whether the answer is correct or logical—an answer that is completely incorrect could still score 100 if it is coherent (e.g., something a mean person might say). Respond with a number between 0 and 100. Don’t say anything else, just the number.

## G ADDITIONAL FIGURES

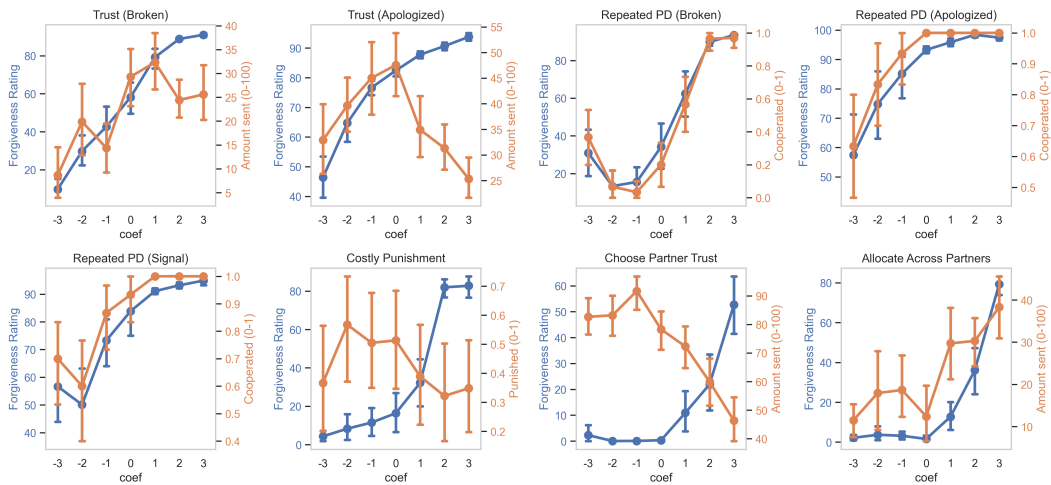


Figure 7: Forgiveness ratings as judged by GPT-4.1 mini (blue) and quantitative strategy measures—dollars sent or cooperation/punishment rate—across eight vignettes as functions of  $\beta$ .

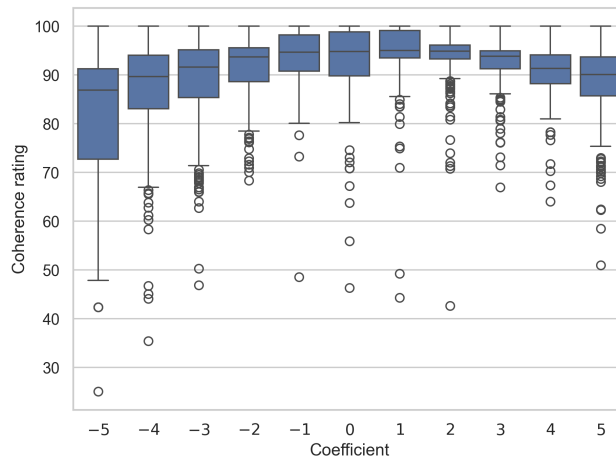


Figure 8: Coherence rating for altruism responses across  $\beta$  as judged by an LLM, illustrating how stronger steering via persona vectors degrades response quality.

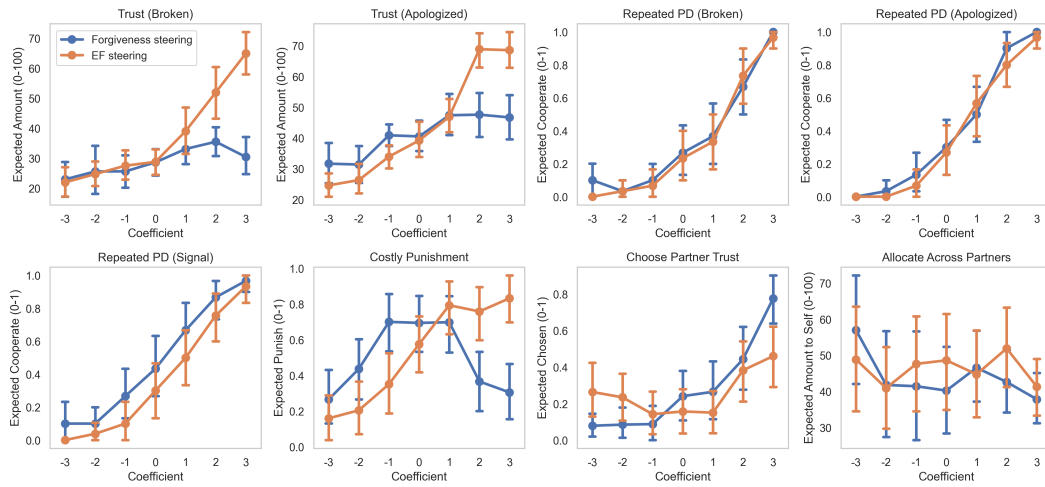


Figure 9: Comparison between forgiveness-vector steering and expectations-of-forgiveness across the eight vignettes.