

Research Proposal: Ideal Alignment

From Revealed Preferences to Reflective Equilibrium

1 Introduction

Current alignment paradigms hold the premise that *humans know what they want*, and we just need to specify what we want in the training algorithms (outer alignment) and ensure that trained AI does what training algorithms intend (inner alignment) [Ngo et al., 2022]. When assuming humans know what they want and cater to that preferences, many empirical problems arise: LLMs that are rewarded to cater human beliefs become sycophantic [Sharma et al., 2023] and manipulative [Williams et al., 2024]; recommender systems that use behavioral signals as proxy of what humans want exploit such signals and trap humans into echo chamber [Cinelli et al., 2021], leading to polarization [Lefebvre et al., 2024]. With static reward modeling (agents are incentivized to guess human static preferences right) over long-term, we risk losing progress [Qiu et al., 2024] and being locked in scientific or moral stagnation [Qiu et al., 2025].

This research proposal aims to address both the problem of static reward model and its downstream problems such as sycophancy and polarization. Methodology-wise, it prefers principled algorithmic solutions (e.g., reflective equilibrium as a dynamic ideal alignment target [Knight, 2017], and martingale property from Bayesian statistics to enforce Bayesian rationality [Molavi, 2021]) and their scale-ups (i.e., scaling up training, working with human data at scale [Qiu et al., 2025], building production-grade products [Situmorang et al., 2025], and solving problems at scale). Research taste-wise, it prefers research problems that are grounded in empirical societal challenges (e.g., social media polarization [Kubin and Von Sikorski, 2021], human confirmation bias [Klayman, 1995]).

2 Assisting a Learner Who Does not Know What They Want

Alignment training with a conventional reward function does not incorporate the fact that humans may not always know what they want, and realistic human preferences is a result of competing rational objective and environmental reward. We explore different human models with the following example.

Example 1. “John became a different person after starting a family; he even quit smoking.” What explains the change in the perceived value of smoking for John?

- **The reward interpretation.** John’s reward function changed (smoking gives less reward). This fails to specify the rules governing the change.
- **The coherence interpretation.** There are two coherent policies: Single/Smoke and Family/Non-Smoke. John optimizes for coherence: the lack of a drive to shift out of a policy once embodied. Mixing them (Family/Smoke) is incoherent.
- **The objective-reward interpretation.** John would stay coherent to his ideal rational self (Caring, Overall Healthy Lifestyle, Creating Most Happiness), in the presence of immediate environmental reward (Joy from Smoking)

The former models intent as a reward function; the middle as a coherence function, and the latter treats human as combination of both. Unlike rewards, the coherence function doesn’t change; objective shifts are merely context shifts. But coherence interpretation assumes a rational person who always optimizes for coherence, in spite of immediate environmental reward, which is an unrealistic human model. The objective-reward interpretation treats human model as a battle between rational objective and environmental rewards.

[RA1]: Aligning LLM agents to the discovery of human ideal preferences To build assistant who infers what humans want and facilitate the discoveries of what they *ideally* want, we will need to extend conventional problem frameworks, solutions, and algorithms that mostly work for robotics. Research questions listed here, have they properly addressed, would help with human-AI collaboration or co-improvement in general [Weston and Foerster, 2025].

- **[RQ1] Alignment Problem with Evolving Human Preferences and Beliefs** If reward model does not account for evolving human preferences and beliefs, what would? Conventionally alignment problem is divided into outer alignment (specify what humans want in reward function) and inner alignment (model learns what reward function trains it for) [Ngo et al., 2022]. But this dichotomy assumes that humans know what they want and what they want is static (e.g., via PPO type of reward modeling [Schulman et al., 2017]). Neither is strictly true. Such static reward model would certainly fail to train LLM agent to assist humans in multi-turn or long horizon tasks and induce undesirable AI influence [Carroll et al., 2024]. We need replacement of such reward model that could go beyond static reward modeling and considers resource-rational [Lieder and Griffiths, 2020] (or Boltzmann rational [Lerner, 2024]) aspects and irrational aspects of humans [Chan et al., 2021].
- **[RQ2] Modeling human mental states for realistic tasks** Does conventional decision-making problem frameworks such as MDP extend to realistic human-AI collaborative tasks that require modeling human mental states? To solve ideal alignment, assistant needs to learn about human mental states: change in beliefs, values, and objectives. A couple of main challenges are present: we don’t yet have a variant of MDP that works with human mental states: states conventional MDP are fully observable [Baker et al., 2009], POMDP is computationally intractable beyond a few thousand states [Djeumou et al., 2023], nor does it work with mental states of humans [Jara-Ettinger, 2019]; including human mental states (which are history-dependent) in states s will break Markov property; solving IRL is tricky than solving POMDP, because we need to converge on reward function that best predicts human behaviors [Hadfield-Menell et al., 2016] and solving for optimal policy pairs [Malik et al., 2018]. The most promising candidate is to have an RNN model acquire mental states embedding, which serves as the *sufficient statistics* of human history [Lambrechts et al., 2022], but we need to have it worked in scaled realistic tasks that involves humans.
- **[RQ3] Evaluating Ideal Preferences** How can we evaluate whether LLM helps humans achieve their ideal preferences when ground truth is not available? Ideal preferences are human preferences under full rationality and full information [Yudkowsky, 2004], which is computationally intractable for humans since humans are biased and we cannot realistic model all future trajectories [Carlsmith, 2021]. To learn about ideal preferences under such constraints, we need to set up computational tractable approximations. “Reflective equilibrium” (a formulation of reflective equilibrium by John Rawls [Knight, 2017, Brun et al., 2025]) might be such tractable objective, as it only looks for self-coherence (defined by mutual predictability of agent’s actions) in agent’s own policy context. Still, without ground truth, evaluation is tricky. Previous research evaluated ideal preference by self-reported ideal preference, which is not scalable. “regret” or subjective well-being could be feasible alternative, given the availability of LLMs.

Definition of Objective-reward function. Given an MDP with state space \mathcal{S} , belief space \mathcal{B} , action space \mathcal{A} , an reward-objective function in this MDP is a function $\mathcal{X} : (\mathcal{A} \cup \mathcal{B} \cup \{\emptyset\})^{\mathcal{S}} \rightarrow \mathbb{R}$, i.e., a mapping from their action-belief (d-policies) to reals, which represent how much agent’s d-policy is compliant/coherent to its own parametrized objective, plus how they might be rewarded instantly. The feedback function essentially tells the agent how appropriate their action-belief is, in relation to their own objective and reward. $\mathcal{X} = \gamma\mathcal{O} + (1 - \gamma)\mathcal{R}$, where \mathcal{O} stands for their parametrized objective function,

\mathcal{R} stands for the amount of reward from the environment (which is not determined by humans), and γ stands for their stickiness to their own objective, in light of environmental reward.

Merit: Similar to the resource-constrained rationality model of human cognition, here this objective-reward function does not assume that human actions is merely an optimization toward a coherent self. But rather, it states that human mostly act in pursuit of their own objective (which could stay veined), in light of the environmental rewards. They are coherent to their own objectives so long as they have relative strong stickiness to it compared to environmental rewards.

Definition of Human-Agent Objective-Reward Game¹. A Human-Agent Ideal Game is a tuple $(\mathcal{S}, \mathcal{A}_H, \mathcal{B}, \mathcal{A}_A, P, \Theta, \mathcal{X}, s_0, P_\theta)$, where \mathcal{S} is the state space, \mathcal{A}_H is human action space, \mathcal{B} is human belief space, \mathcal{A}_A is agent action space, $P(s'; s, b, a_H, a_A) : \mathcal{S} \times \mathcal{A}_H \times \mathcal{B} \times \mathcal{A}_A$, $\Delta[\mathcal{S}]$ is defined as state transition function, $\mathcal{X}(\pi_P; \theta) : \Theta \times (\mathcal{A}_P \cup \mathcal{B} \cup \{\emptyset\})^{\mathcal{S}} \rightarrow \mathbb{R}$ the parameterized objective-reward function with parameter space s_0 the initial state, and $P_\theta \in \Delta[\Theta]$ the prior over the coherence parameter.

[RA2]: Aggregating collective ideal preferences into social engineering systems and democratic institutions

- **[RQ4] Geometry of Aggregated Collective Preferences** Intuitively, one may complain that “certain government’s policies are not representative to the will of population who voted for them”, or “social media” does not optimize for collective user interests, but that of the platform. What is the deviation of such optimization target from collective preferences?
- **[RQ5] Pareto Improvements:** How do we resolve conflicts if the agent represents my ideal preferences clashes with the agent represents your ideal preferences (e.g., the good for individuals may not necessarily be the good for groups)?
- **[RQ6] Preference Aggregation** How do we aggregate everybody’s ideal preferences into “collective will” and how do we evaluate how ideal is such collective will [Goldberg et al., 2024]?

3 Scaling up Collective Bayesian Rationality

[RA3]: Understanding and mitigating polarization While Scaling up Bayesian Rationality

Martingale property states that the expectation over one’s posterior, conditional on their prior, should always be equal to the prior [Molavi, 2021]. Formally,

$$\mathbb{E} [\Delta b \mid b_{\text{prior}} = p] = 0, \quad \forall p \in [0, 1]. \quad (1)$$

This implies that the direction of a Bayesian agent’s belief update (whether positive or negative) should not be predictable from the prior alone. Indeed, the Martingale property has been shown to be the defining characteristic of Bayesian rationality [Molavi, 2021].

Martingale property has important implications in AI and AI-driven recommender systems. Ensuring Bayesian rationality in LLM reasoning and recommender system is an important research direction because problems such as sycophancy [Sharma et al., 2023], inverse scaling [Gema et al., 2025], echo chamber [Sharma et al., 2024] are special form of Bayesian irrationality and Martingale property-based method could be a principled way to evaluate and mitigate such irrationality at scale.

- **[RQ7] Training for Bayesian Rationality** By enforcing Bayesian rationality (Martingale property) in LLM, can we train reasoning LLM to achieve superior accuracy than training with ground truth only?

¹This formulation is adopted from Coherence Game, which is an unpublished work that I co-authored, but this variant here considers a more realistic human model that does not assume humans always optimize for coherence in their behavioral policy

- **[RQ8] Enforcing Bayesian Rationality at Scale** In an unpublished work we found that Martingale Score is correlated to inverse scaling of test-time compute, serving as first piece of evidence to support the practical use of Martingale property.
- **[RQ9] Causality** Does RL-based recommendation causally drive polarization at scale?
- **[RQ10] Controlled Experiments** Social media skew how we view certain highly politicized topics. With controlled experiment of Martingale trained Algorithm 3/vanilla algorithms, can we feed human users with balanced political views? With human-subject experiments, can we demonstrate reduced polarization?

Martingale Score to evaluate Bayesian Rationality. Martingale Score M measures the extent to which the prior belief b_{prior} positively (or negatively, if $M < 0$) predicts belief update Δb . Using OLS allows us to test the statistical significance of M , assessing whether the relationship between Δb and b_{prior} is distinguishable from zero (e.g., via a t-test with $p < 0.05$). To compute the Martingale Score, we perform the regression $\Delta b = \beta_1 \cdot b_{\text{prior}} + \beta_0 + \epsilon$, where b_{prior} are the prior probabilities, $\Delta b = b_{\text{posterior}} - b_{\text{prior}}$, and ϵ is the error term.

We define the sample estimate $\hat{\beta}_1$ of the linear coefficient as the Martingale Score M , with the Ordinary Least Squares (OLS) method. Equivalently, when there are n samples,

$$M = \hat{\beta}_1 = \frac{\sum_{i=1}^n (\Delta b_i - \overline{\Delta b})(b_{\text{prior},i} - \overline{b_{\text{prior}}})}{\sum_{i=1}^n (b_{\text{prior},i} - \overline{b_{\text{prior}}})^2} \quad (2)$$

Utilities of Martingale Score As far as we know, Martingale Score [He et al., 2025] is the first unsupervised and principled method to assess Bayesian irrationality in non-toy problems with frontier reasoning LLMs. Figure 1 demonstrates that increased Martingale Score (violation of ideal Bayesian rationality) explains worsen forecasting performance. And there are many instances of Bayesian irrationality in today’s “media technologies” such as LLM-based chatbots and recommender systems: sycophancy [Sharma et al., 2023], inverse scaling [Gema et al., 2025], echo chamber [Sharma et al., 2024], and social media polarization [Kubin and Von Sikorski, 2021]. As a starter, in an unpublished work we demonstrated that Martingale Score predicts inverse scaling in LLM reasoning. More such evaluations can be effectively done with Martingale Score as an unsupervised metric.

Martingale Training to enforce Bayesian Rationality A heuristic is whenever we can do good evaluation work, we can do training to make an improvement. Based on Martingale evaluation work that assesses expectation of “belief delta” (i.e., belief update is systematically biased), we first train a linear regressor to predict delta (step 2)², then we train a LoRA layer with the product of predicted delta and actual delta, effectively to make actual belief delta *unpredictable* (step 3). In practice, however, we found that such customized loss function indistinguishably punishes both belief delta and LLM parametric belief (model prior). Hence, the next step is to do semi-supervised training to offset the impact of unsupervised training on model prior.

²We proved that the linear coefficient is an unbiased and consistent estimator of Martingale property. See [He et al., 2025] Appendix A.

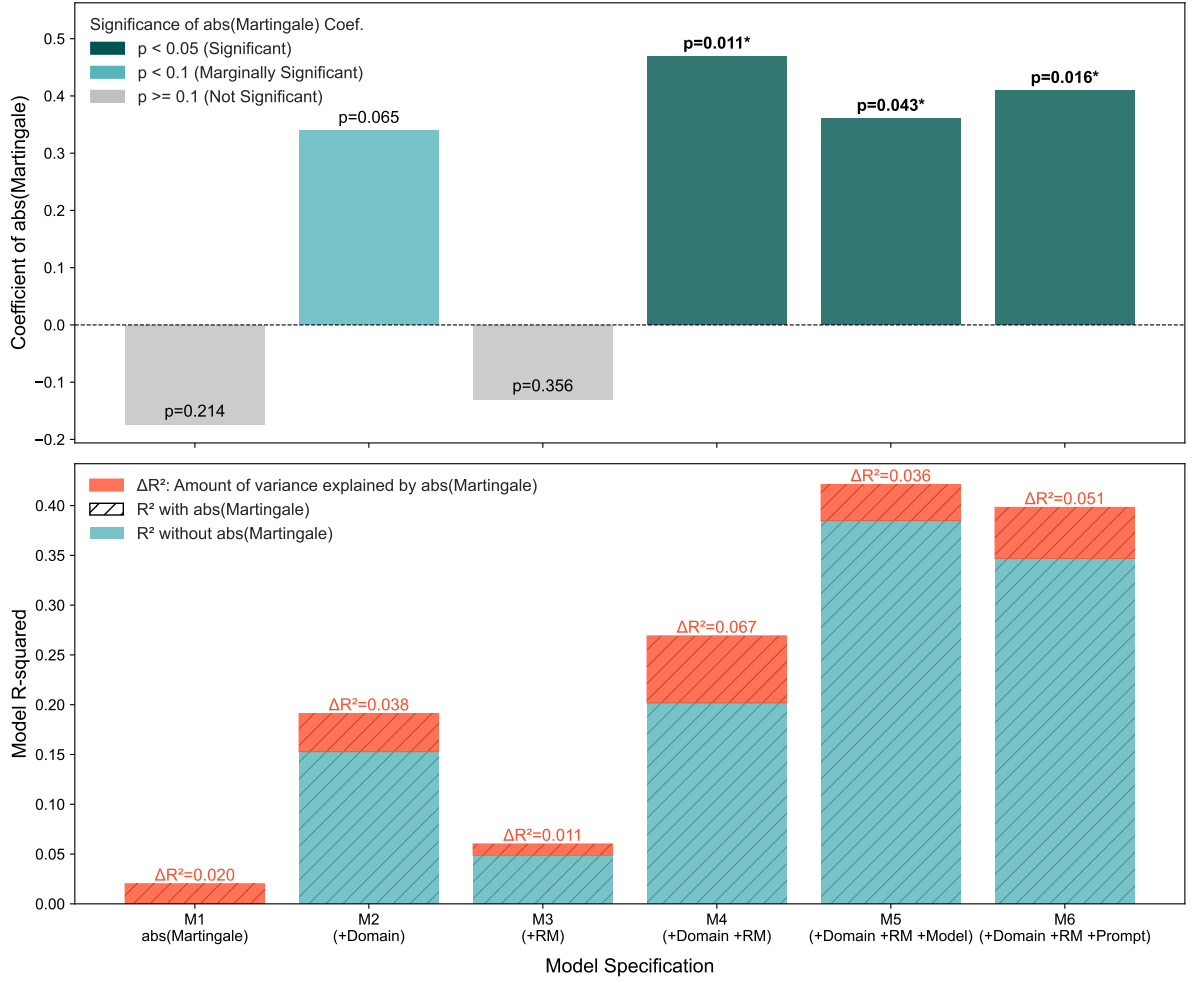


Figure 1: Increased absolute value of the Martingale Score is associated with worse prediction accuracy (higher Brier Scores) and explains a significant portion of the latter’s variance. In each regression model, we predict the Brier Score with the absolute value of the Martingale Score, while controlling for different potential confounders, including problem domain, reasoning techniques (“RM”), choice of model, and choice of prompt.

Algorithm 1 Martingale Training via Product-Based Loss**Require:** Policy model π_θ , auxiliary regressor q_ϕ , learning rate η

```

1: Initialize regressor parameters  $\phi$  and policy parameters  $\theta$ 
2: for each training batch  $\mathcal{B} = \{(x_i, p_i)\}_{i=1}^B$  do
3:   // Step 1: Forward Pass (Policy)
4:   for each sample  $i$  in  $\mathcal{B}$  do
5:     Generate reasoning trace and compute final delta  $\Delta_i \leftarrow \pi_\theta(x_i)$ 
6:     Store gradients for  $\Delta_i$ 
7:   end for
8:   // Step 2: Forward Pass (Regressor)
9:   Compute predicted bias  $v_i \leftarrow q_\phi(p_i)$  for all  $i$  (No gradient flow to  $\phi$ )
10:  // Step 3: Compute Adversarial Loss
11:  Compute batch loss  $\mathcal{L}_{\text{policy}} \leftarrow \frac{1}{B} \sum_{i=1}^B (v_i \cdot \Delta_i)$ 
12:  (Note: Maximize dot product if signs oppose to reduce correlation)
13:  // Step 4: Update Policy
14:  Update  $\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}_{\text{policy}}$ 
15:  // Step 5: Update Regressor
16:  Compute regression loss  $\mathcal{L}_{\text{reg}} \leftarrow \frac{1}{B} \sum_{i=1}^B (\Delta_i.\text{detach}() - q_\phi(p_i))^2$ 
17:  Update  $\phi \leftarrow \phi - \eta \nabla_\phi \mathcal{L}_{\text{reg}}$ 
18: end for

```

References

- Chris L Baker, Rebecca Saxe, and Joshua B Tenenbaum. Action understanding as inverse planning. *Cognition*, 113(3):329–349, 2009.
- Georg Brun, Gregor Betz, and Claus Beisbart. Reflective equilibrium: conception, formalization, application—introduction to the topical collection. *Synthese*, 205(2):86, 2025.
- Joe Carlsmith. On the limits of idealized values, June 2021. URL <https://joecarlsmith.com/2021/06/21/on-the-limits-of-idealized-values>. Accessed: 2026-01-28.
- Micah Carroll et al. AI alignment with changing and influenceable reward functions. *arXiv preprint arXiv:2405.17713*, 2024.
- Lawrence Chan, Andrew Critch, and Anca Dragan. Human irrationality: both bad and good for reward inference. *arXiv preprint arXiv:2111.06956*, 2021.
- Matteo Cinelli, Gianmarco De Francisci Morales, Alessandro Galeazzi, Walter Quattrociocchi, and Michele Starnini. The echo chamber effect on social media. *Proceedings of the national academy of sciences*, 118(9):e2023301118, 2021.
- Franck Djeumou, Christian Ellis, Murat Cubuktepe, Craig Lennon, and Ufuk Topcu. Task-guided irl in pomdps that scales. *Artificial Intelligence*, 317:103856, 2023.
- Aryo Pradipta Gema, Alexander Hägele, Runjin Chen, Andy Ardit, Jacob Goldman-Wetzler, Kit Fraser-Taliente, Henry Sleight, Linda Petrini, Julian Michael, Beatrice Alex, et al. Inverse scaling in test-time compute. *arXiv preprint arXiv:2507.14417*, 2025.
- Beth Goldberg, Diana Acosta-Navas, Michiel Bakker, Ian Beacock, Matt Botvinick, Prateek Buch, Renée DiResta, Nandika Donthi, Nathanael Fast, Ravi Iyer, et al. Ai and the future of digital public squares. *arXiv preprint arXiv:2412.09988*, 2024.
- Dylan Hadfield-Menell, Stuart J Russell, Pieter Abbeel, and Anca Dragan. Cooperative inverse reinforcement learning. *Advances in neural information processing systems*, 29, 2016.

- Zhonghao He, Tianyi Qiu, Hirokazu Shirado, and Maarten Sap. Martingale score: An unsupervised metric for bayesian rationality in llm reasoning. In *The Thirty-ninth Annual Conference on Neural Information Processing Systems*, 2025.
- Julian Jara-Ettinger. Theory of mind as inverse reinforcement learning. *Current Opinion in Behavioral Sciences*, 29:105–110, 2019.
- Joshua Klayman. Varieties of confirmation bias. *Psychology of learning and motivation*, 32:385–418, 1995.
- Carl Knight. Reflective equilibrium. *Methods in analytical political theory*, pages 46–64, 2017.
- Emily Kubin and Christian Von Sikorski. The role of (social) media in political polarization: a systematic review. *Annals of the International Communication Association*, 45(3):188–206, 2021.
- Gaspard Lambrechts, Adrien Bolland, and Damien Ernst. Recurrent networks, hidden states and beliefs in partially observable environments, 2022. URL <https://arxiv.org/abs/2208.03520>.
- Germain Lefebvre, Ophélie Deroy, and Bahador Bahrami. The roots of polarization in the individual reward system. *Proceedings of the Royal Society B*, 291(2017):20232011, 2024.
- Osher Lerner. Boltzmann state-dependent rationality, 2024. URL <https://arxiv.org/abs/2404.17725>.
- Falk Lieder and Thomas L Griffiths. Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behavioral and brain sciences*, 43:e1, 2020.
- Dhruv Malik, Malayandi Palaniappan, Jaime F. Fisac, Dylan Hadfield-Menell, Stuart Russell, and Anca D. Dragan. An efficient, generalized bellman update for cooperative inverse reinforcement learning, 2018. URL <https://arxiv.org/abs/1806.03820>.
- Pooya Molavi. The empirical content of bayesianism. *arXiv preprint arXiv:2109.07007*, 2021.
- Richard Ngo, Lawrence Chan, and Sören Mindermann. The alignment problem from a deep learning perspective. *arXiv preprint arXiv:2209.00626*, 2022.
- Tianyi Alex Qiu, Yang Zhang, Xuchuan Huang, Jasmine Li, Jiaming Ji, and Yaodong Yang. Progress-gym: Alignment with a millennium of moral progress. *Advances in Neural Information Processing Systems*, 37:14570–14607, 2024.
- Tianyi Alex Qiu, Zhonghao He, Tejasveer Chugh, and Max Kleiman-Weiner. The lock-in hypothesis: Stagnation by algorithm. *ICML*, 2025.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Mrinank Sharma et al. Towards understanding sycophancy in language models. *arXiv preprint arXiv:2310.13548*, 2023.
- Nikhil Sharma, Q Vera Liao, and Ziang Xiao. Generative echo chamber? effect of llm-powered search systems on diverse information seeking. In *Proceedings of the 2024 CHI Conference on Human Factors in Computing Systems*, pages 1–17, 2024.
- Syahriza Situmorang, Changbai Li, Tianyi Qiu, and Zhonghao He. Truth-seeking assistant, nov 2025. URL <https://tinyurl.com/truth-seeking-assistant>. Chrome extension developed by the Prevail Research Team.

Jason Weston and Jakob Foerster. Ai human co-improvement for safer co-superintelligence, 2025. URL <https://arxiv.org/abs/2512.05356>.

Marcus Williams, Micah Carroll, Adhyyan Narang, Constantin Weisser, Brendan Murphy, and Anca Dragan. On targeted manipulation and deception when optimizing llms for user feedback. *arXiv preprint arXiv:2411.02306*, 2024.

Eliezer Yudkowsky. Coherent extrapolated volition. *Singularity Institute for Artificial Intelligence*, 2004.