

POWER AND LIMITATIONS OF AGGREGATION IN COMPOUND AI SYSTEMS

Anonymous authors

Paper under double-blind review

ABSTRACT

When designing AI systems for complex tasks, it is becoming increasingly common to query a model in different ways and aggregate the outputs to create a compound AI system. In this work, we mathematically study the power and limitations of aggregation within a stylized principal-agent framework. This framework models how the system designer can partially steer each agent’s output through reward specification, but still faces limitations due to prompt engineering ability and model capabilities. Our analysis identifies three mechanisms—feasibility expansion, support expansion, and binding set contraction—through which aggregation can expand the set of elicitable outputs. We prove that any aggregation operation must implement one of these mechanisms to provide benefit, though none are sufficient alone. To sharpen this picture, we establish necessary and sufficient conditions for when aggregation expands elicitable outputs. Altogether, our results take a step towards characterizing when compound AI systems can overcome limitations in model capabilities and in prompt engineering.

1 INTRODUCTION

Compound AI systems—which leverage multiple AI components, rather than a single model in isolation—present a powerful paradigm to tackle complex tasks (BAIR Research Blog, 2024). In the context of large language models (LLMs), one common approach is to create many copies of the same model, give these models different prompts or access to different tools, and aggregate the outputs of these models at test-time. This approach has proven fruitful in multi-agent research systems (Anthropic Engineering, 2024) where a lead LLM agent delegates subtasks to different specialized agents and aggregates their outputs, in multi-agent debate protocols where different LLM agents seek consensus (Du et al., 2024) or argue for different answers (Khan et al., 2024), and in prompt ensembling approaches where the outputs from different prompts are combined (Arora et al., 2023).

Given the empirical success of these compound LLM systems, this raises the question of when aggregating across multiple copies of the same model unlocks greater performance than querying a single model. At first glance, aggregation may seem redundant when the model copies are homogeneous. However, one source of improved performance is at the prompt level: a model with a complex prompt engineering approach may be replaceable by a set of models with simple but diverse prompting strategies (Arora et al., 2023), illustrating how aggregation across models can overcome limitations in prompt engineering ability. Another source of improved performance is at the output level: aggregating multiple LLM agents over repeated interactions can help correct errors such as hallucinations (Du et al., 2024), illustrating how aggregation can overcome limitations in model capabilities as well. This suggests that the extent to which aggregation overcomes these limitations in prompt engineering and model capabilities impacts the power of compound AI systems.

In this work, we study the power and limitations of aggregation from a theoretical perspective, building on a classical principal-agent framework (Kleinberg et al., 2019). Our focus is on compound AI systems where a system designer passes reward specifications (e.g., via prompts) to many copies of the same model and then aggregates their outputs. In this stylized principal-agent framework (Section 2), the system designer (i.e., the principal) designs reward specifications to elicit N -dimensional outputs from each agent, and aggregates these outputs to produce a synthesized output. Each agent generates the outputs in its feasible set that maximizes the reward, and the system-designer strategically co-designs the rewards across models to try to produce a specific output. We capture prompt

engineering limitations as the rewards operating over a coarser M -dimensional feature space, and model capability limitations as conic constraints on each agent’s feasible set of outputs.

Using this framework, we characterize when aggregating across multiple agents enables the system designer to elicit a greater set of outputs than relying on a single model. To build intuition, we formalize three natural mechanisms by which aggregation can expand the set of elicitable outputs (Section 3). The first mechanism is *feasibility expansion*, where aggregation produces outputs outside of any agent’s feasibility set. The second is *support expansion*, where aggregation combines outputs with smaller supports into an output with a larger support. The third is *binding set contraction*, where aggregation combines outputs that are binding with respect to constraints into an output that falls within the interior.

We formally connect these mechanisms to elicibility-expansion. Specifically, we find that the power of aggregation fundamentally relies on at least one of these mechanisms being implemented: if none are implemented, then aggregation does not expand elicibility on any problem instance (Theorem 3.7). However, these mechanisms are not sufficient to expand elicibility in general, although we show that each mechanism results in elicibility-expansion under stronger conditions.

To more completely capture the power and limitations of aggregation, we provide a more general characterization of elicibility-expansion (Section 4). We first characterize when an aggregation operation is elicibility-expanding in a given problem instance (Theorem 4.1), linking this to whether feasible directions for agent outputs intersect with feature-improving directions. To analyze the limitations of aggregation, we derive general conditions (Definition 4.2) under which an aggregation operation never expands the set of elicitable outputs, regardless of the level of coarseness of the feature space (Theorem 4.3), and we show that these conditions are tight (Theorem 4.4). These tight conditions in Definition 4.2 are strengthenings of feasibility expansion, support expansion, and binding-set contraction. At a high-level, these conditions test whether feasible directions under which an agent can change the aggregated output violate the constraints by a sufficient margin.

Altogether, our results uncover key mechanisms that underpin the power and limitations of an aggregation in compound AI systems. Our results suggest conditions for aggregation to add no power to a system, regardless of the level of prompt engineering limitations. Moreover, our results illustrate how the power of an aggregation depends on the interplay between prompt engineering ability and model capabilities. More broadly, our results take a step towards understanding when aggregation of multiple copies of the same model provides benefits to system designers.

1.1 RELATED WORK

Aggregation across multiple models. Aggregating outputs from multiple LLMs is a common strategy for complex tasks (BAIR Research Blog, 2024). One common approach is resampling the same model or reasoning trace and then selecting outputs via reward models (Christiano et al., 2017), self-consistency (Wang et al., 2023b), or synthesis (Zhang et al., 2025); coverage is an important property for inference-time computations (Huang et al., 2025). Other approaches are routing queries across different LLMs (Chen et al., 2024), adversarially combining models to expose safety risks (Jones et al., 2025a), and consensus games between generators and discriminators (Jacob & Andreas, 2024). Closest to our setting are systems with multiple copies of the same model under different reward specifications, as in LLM debate (Du et al., 2024), prompt ensembling (Arora et al., 2023), and multi-agent research frameworks (Anthropic Engineering, 2024). We provide a theoretical perspective on when such aggregation elicits strictly more outputs than a single model. Classical work has analyzed aggregation in settings such as ensembling (Dietterich, 2000), voting (Ladha, 1992), distributed algorithms (Lynch, 1996), and multi-agent reinforcement learning (Tan, 1993).

Principal-Agent Models and Reward Design. Our model is inspired by the principal-agent model by Kleinberg et al. (2019). We extend their technical result to incorporate agent limitations in the form of conic constraints and derive new results that characterize elicibility via aggregation. This falls under the broader principal-agent framework (Holmström, 1979; Grossman & Hart, 1983; Laffont & Martimort, 2002; Bolton & Dewatripont, 2005), which captures the challenge of designing rewards based on imperfect proxies. (Zhuang & Hadfield-Menell, 2020) use this framework to study misalignment of AI, which is similar to our motivation. Work in this framework also incorporates agent’s limitations in the form of costs for actions. Particularly related are multitask settings (Holmström & Milgrom, 1991; Slade, 1996; Bond & Gomes, 2009; Demougin et al., 2022) that study

the effects of costs being dependent between tasks, including cases of substitutability and complementarity, which is similar to our conic constraints that capture dependence among multiple output dimensions. Principal-agent theory has also considered multiple agents (Holmström, 1982; Lazear & Rosen, 1981; Dasaratha et al., 2024), focusing mainly on the joint design of rewards. Our work differs in allowing aggregation to synthesize new outputs, and in characterizing when aggregation provides provable benefits rather than addressing algorithmic design. Complementary work studies benefits of heterogeneity across agents (Gentzkow & Kamenica, 2017; Collina et al., 2025), though they don’t study heterogeneity through differently designed rewards.

2 MODEL

We extend the principal-agent framework in Kleinberg et al. (2019) to model a compound AI system with K agents (who represent LLMs) and a single principal (the system designer). The system designer designs reward specifications to elicit outputs from the agents, and aggregates the outputs to synthesize a new output. The system designer faces limitations on the complexity of rewards they can design, and the agents face limitations in terms of the space of outputs that they can generate. We defer a discussion of model limitations to section 5.

2.1 OUTPUT SPACE

We embed outputs of agents into M -dimensional vectors with non-negative coordinates. We view each output dimension as capturing a different characteristic of the output. The vector representation \mathbf{x} quantifies the degree to which the output captures each characteristic. We note that some dimensions may capture undesirable characteristics (e.g., hallucinations). The system designer seeks a specific output $\mathbf{x}^{(A)} \in \mathbb{R}_{\geq 0}^M$, which we assume to be unit ℓ_1 -norm $\|\mathbf{x}^{(A)}\|_1 = 1$.

Our model captures how the agents have restrictions on the set of output vectors that it can produce, for example due to capability limitations. The first restriction is that the ℓ_1 norm of the output vectors is bounded, which captures budget limitations. The second restriction is conic constraints on the output, which each take the form $\mathbf{c}^T \mathbf{x} \leq 0$ where $\mathbf{c} \in \mathbb{R}^M$ contains at least strictly positive entry and at least strictly negative entry. These conic constraints capture restrictions on the types of outputs that the agent can produce: for example, some agents may not be able to avoid producing hallucinations without facing capability degradation along other characteristics.

We let L denote the number of conic constraints, and we let $\mathbf{C} \in \mathbb{R}^{L \times M}$ denote the conic constraints themselves. Let $\mathbf{C}_i \in \mathbb{R}^M$ denote the i th row of \mathbf{C} for $i \in [L]$, and let $\mathbf{C}_V \in \mathbb{R}^{|V| \times M}$ denote the set of rows corresponding to indices $V \subseteq [L]$. We denote by \mathbf{C}_\emptyset the zero-vector, to capture how $\{\mathbf{d} : \mathbf{C}_\emptyset \leq 0\} = \mathbb{R}_{\geq 0}^M$. Given a budget level $E > 0$, we let $\mathcal{B}(E)$ denote the feasible set at budget level E , defined to be:

$$\mathcal{B}(E) := \{\mathbf{x} \in \mathbb{R}_{\geq 0}^M \mid \mathbf{C}\mathbf{x} \geq \mathbf{0}, \|\mathbf{x}\|_1 \leq E\}.$$

2.2 REWARD SPECIFICATION

The system designer designs a reward specification $R^{(k)}$ and a budget level $E^{(k)}$ for each agent $k \in [K]$. The reward specification represents the reward implicit in the prompt that they give to the agent, and the budget level represents the level of test-time compute that the agent is allowed to use.

To capture prompt engineering limitations, we model the reward specification as operating over a coarser N -dimensional feature space than the outputs. Here, the features $F(\mathbf{x}) = [F_1(\mathbf{x}), \dots, F_N(\mathbf{x})]$ take the form

$$F_j(\mathbf{x}) = f_j \left(\sum_{i=1}^M \alpha_{ij} \mathbf{x}_i \right),$$

where $f_j(\cdot)$ is nonnegative, smooth, weakly concave (i.e., diminishing returns from increasing quality on this dimension), and strictly increasing, and where the values $\alpha_{ij} \geq 0$ are nonnegative *feature weights*. We will denote by $\alpha \in \mathbb{R}_{\geq 0}^{M \times N}$ the matrix with entries α_{ij} and call this the *feature weights matrix*.

We consider reward specifications $R^{(1)}, \dots, R^{(K)} : \mathbb{R}^N \rightarrow \mathbb{R}$ which operate on these features. Following prior work (Kleinberg et al., 2019), we restrict to *monotone* reward functions R which do not decrease if all features are weakly increased, and where there exists $j \in [N]$ such that R strictly increases whenever the feature F_j strictly increases.

Given a monotone reward specification $R^{(k)}$ and a positive budget level $E^{(k)} > 0$, each agent k produces an output that maximizes its reward over the feasible set $\mathcal{B}(E^{(k)})$: that is,

$$\mathbf{x} \in \mathbf{X}^*(R^{(k)}, E^{(k)}) := \operatorname{argmax}_{\mathbf{x} \in \mathcal{B}(E^{(k)})} R^{(k)}(F(\mathbf{x})).$$

This captures how even though agents are homogeneous and solve the same optimization program, they can be given different reward specifications and thus produce different outputs.

2.3 ELICITABILITY

We say that a reward specification R and budget level E elicits an output \mathbf{x} if $\mathbf{x} \in \mathbf{X}^*(R, E)$. This captures whether an agent can produce the output \mathbf{x} : that is, if $\mathbf{x} \in \operatorname{argmax}_{\mathbf{x} \in \mathcal{B}(E)} R(F(\mathbf{x}))$. As shown in prior work (Kleinberg et al., 2019) and illustrated in Section 3.1, some output vectors $\mathbf{x} \in \mathcal{B}$ where \mathbf{x} are not elicitable by any reward specification R and budget level E .

We say that an output \mathbf{x} is elicitable if there exists a monotone reward specification R and a positive budget level E that elicits \mathbf{x} . The condition for whether \mathbf{x} is elicitable only depends on \mathbf{x} through the following sufficient statistic $(\mathcal{S}(\mathbf{x}), \mathcal{V}(\mathbf{x}))$. The first component $\mathcal{S}(\mathbf{x}) = \{j : x_j > 0\}$ denotes the support of \mathbf{x} . The second component $\mathcal{V}(\mathbf{x}) = \{l \in [L] : C_l \mathbf{x} = 0\}$ denotes the set of indices of conic constraints that are binding at \mathbf{x} .

Aggregation. When the system designer can aggregate the outputs of different agents, this may expand the set of elicitable outputs. The following definition captures when this occurs.

Definition 2.1. We call $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is an *elicibility-expanding operation* if

- There exist monotone reward specifications $R^{(1)}, \dots, R^{(K)}$ and positive budget levels $E^{(1)}, \dots, E^{(K)}$ such that $\mathbf{x}^{(k)} \in \mathbf{X}^*(R^{(k)}, E^{(k)})$ for all $k \in [K]$.
- There does not exist a monotone reward specification R and budget level $E > 0$ such that $\mathbf{x}^{(A)} \in \mathbf{X}^*(R, E)$.

Intuitively, if an aggregation operation is elicibility-expanding, then allowing the system-designer to aggregate outputs according to this operation produces an output that is not elicitable with a single reward, but can be obtained by combining outputs elicited from multiple reward specifications.

Aggregation rules. An aggregation rule is a mapping from a list of output vectors $(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})$ to an aggregated output vector $\mathbf{x}^{(A)}$. There are two natural aggregation rules we will often consider in our work. Although our results apply to more general aggregation rules, we will often use these natural aggregation rules to provide examples.

The first is *intersection aggregation*, which is defined to be the coordinate-wise minimum of the vectors:

$$\mathcal{A}_{\text{intersect}}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}) = \mathbf{x}^{(1)} \wedge \dots \wedge \mathbf{x}^{(K)}. \quad (1)$$

This aggregation rule combines outputs based on commonality among different output vectors, which is conceptually similar to debate protocols (Du et al., 2024) that aim to create agreement or inference scaling methods that aim to filter out incorrect information (Zhang et al., 2025). The second is *addition aggregation*, which takes a weighted sum of the vectors. For a weight vector $\mathbf{w} \in \mathbb{R}_{\geq 0}^K$, the rule is given by

$$\mathcal{A}_{\text{add}}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}; \mathbf{w}) = \sum_{i=1}^K \mathbf{w}_i \mathbf{x}^{(i)}. \quad (2)$$

Addition aggregation interpolates among different output directions. This rule conceptually captures system designers synthesize multiple outputs to delegate specialized subtasks to each agent and synthesize the outputs of these subtasks (BAIR Research Blog, 2024; Anthropic Engineering, 2024).

2.4 ILLUSTRATIVE EXAMPLE: CITATIONS TASK

To ground our framework in a concrete setting, we consider a natural aggregation task—generating a list of papers on a given topic (Wang et al., 2023a; Press et al., 2024). We describe the task and then show how instantiations of our framework capture different aggregation behaviors for it.

Task and Setup. We study a citation task where the system designer seeks a list of 10 influential LLM papers spanning five perspectives: (1) ML theory, (2) NLP/CL, (3) cognitive science, (4) AI alignment and human–AI interaction, and (5) multi-agent systems.

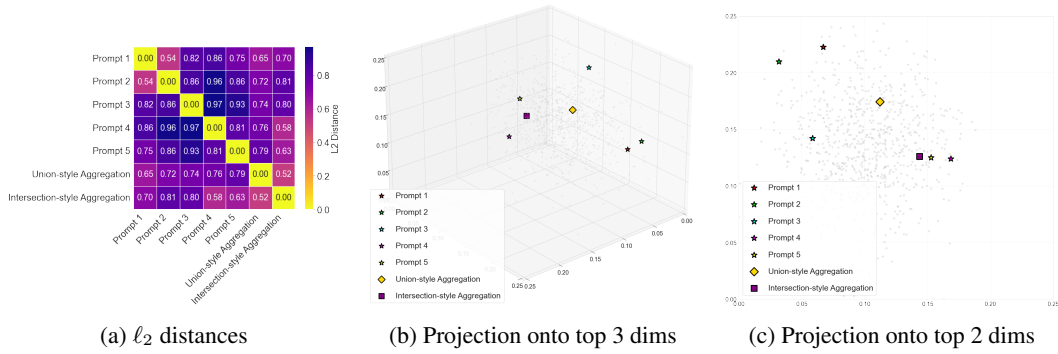


Figure 1: Visualization of output vectors for the citation generation task (Section 2.4). Output vectors are computed using the 768-dimensional embeddings from all-mpnet-base-v2, shifted to be in the nonnegative orthant. Embeddings are shown for GPT-4o-mini outputs from five different prompts, and as well as two different aggregated outputs based on additional-style and intersection-style aggregation rules. The ℓ_2 -distances (left) and projections onto the top 3 highest-variance (middle) and top 2 highest-variance dimensions (right), are shown. The plots show that the five prompts produce semantically different outputs, and each aggregation operation results in a combination of the five outputs that does not resemble any output in isolation.

The system designer issues five prompts, each targeting one perspective, to gpt-4o-mini-2024-07-18 and then aggregates the resulting lists. We prompt another LLM (also gpt-4o-mini-2024-07-18) to aggregate these five lists, instantiating aggregation rules that are inspired by intersection aggregation $\mathcal{A}_{\text{intersect}}$ and union aggregation \mathcal{A}_{add} . Specifically, the model is prompted with *aggregation instructions* along with the five different lists of 10 references, and produces an aggregated list of 10 references. The *intersection-style aggregation instructions* ask for references which are central and broadly relevant across all five perspectives, thus approximating intersection even when the literal overlap of references is empty. The *addition-style aggregation instructions* ask for references that jointly cover and reflect the combined topical space of all five perspectives. We defer the specific prompts and other details of the empirical setup to Appendix G.

Output vectors. We specify two different instantiation of output vectors in our framework, depending on the level of specificity of the output which the system designer aims to elicit.

1. Suppose that the system designer aims to elicit an output that balances multiple high-level criterion in a specific manner (e.g., covering papers in different subareas, covering up to the state of the art in each subarea, quality of the papers selected, etc.). To capture this, let each dimension of the output capture a different criterion that the system designer cares about. We can think of the value of the output vector along each dimension as the extent to which the output captures the criterion corresponding to that dimension.
2. Suppose that the system designer aims to elicit a specific output (e.g., a specific list of references). To capture this, we represent each model output as a high-dimensional embedding coming from a text embedding model. For the citation task, we use all-mpnet-base-v2 (Reimers & Gurevych, 2019), a sentence-transformers model that produces 768-dimensional vectors.¹ Figure 1 shows embeddings for outputs to the five prompts, as well as the outputs produced by the intersection-style and addition-style aggregation rules. The five prompt outputs vary substantially, and the aggregated outputs differ markedly from both the originals and from each other, demonstrating how different prompts and aggregation rules can reshape the embedding-space representation.

Reward specification limitations. Our framework captures two different types of reward specification limitations. First, the system designer may struggle to precisely express what they truly want in the prompt, leading to underspecified prompts omitting some of the system designer’s requirements (Yang et al., 2025). For example, the system designer may prompt the model to include a “breadth

¹As detailed in Appendix G, we apply an additive shift to ensure nonnegativity, computed from the minimum value in each dimension across 805 gpt-4o-mini-2024-07-18 outputs from the helpful-base AlpacaEval dataset (Li et al., 2023).

of citation coverage” when in reality they would like to restrict citations to a handful of academic venues. Second, the model may not correctly interpret the system designer’s prompt by mapping two dissimilar words in the prompt to the same word (Jones et al., 2025b). In the citation task, this could surface as the model interpreting “papers with high attribute ‘X’” similarly for many different attributes “X”. We capture both of these forms of limitations as the reward specification operating over coarsenings of the output dimensions (as captured by the features) rather than directly on the output dimensions.

3 NATURAL MECHANISMS FOR ELICITABILITY-EXPANSION

In this section, we formalize natural mechanisms by which aggregation expands elicibility. First, we show how mechanisms expand elicibility via examples (Section 3.1). Then, we show that these mechanisms are necessary for elicibility-expansion (Section 3.2). The results in this section leverage the technical tools that we develop in Section 4. Note that our goal in this section is to link the mechanisms to elicibility expansion, rather than characterize it; we defer a full characterization to Section 4.

3.1 FORMALIZING THE MECHANISMS AND MOTIVATING EXAMPLES

We formalize three natural mechanisms through which aggregation can provide benefits in our framework. For each mechanism, we illustrate through an example how the mechanism can enable an aggregation operations to expand the set of elicible outputs. Our examples use the intersection and addition aggregation rule that we previously introduced. At the end of this subsection, we investigate the extent to which these aggregation rules can implement the mechanisms that we formalize below.

Our examples also focus on a 3-dimensional output space ($M = 3$) with 2-dimensional features ($N = 2$). We focus on feature weights matrices α of the form $\alpha_q := \begin{bmatrix} 1 & 0 & q \\ 0 & 1 & q \end{bmatrix}$. Each of the output dimensions x_1, x_2 specialize to features F_1, F_2 , respectively. That is, increasing the first output dimension x_1 only increases the first feature F_1 , and increasing the second output dimension x_2 only increases the second feature F_2 . Increasing the third output dimension x_3 increases both features, though the contribution is weighted by a factor of q . The parameter q captures the extent to which it is possible to simultaneously maximize both features.

Mechanism 1: Feasibility Expansion. Aggregation can help overcome the output limitations (i.e., the feasibility constraints faced by each agent), producing outputs that are outside of the feasible set. We formalize this through the following mechanism.

Definition 3.1 (Feasibility Expansion). *Given a constraint matrix C , an aggregation operation $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ implements **feasibility expansion relative to C** if $\mathbf{x}^{(A)}$ is infeasible i.e., $C\mathbf{x}^{(A)} \not\leq 0$ but all $\mathbf{x}^{(i)}$ for $i \in [K]$ are feasible i.e., $C\mathbf{x}^{(i)} \leq 0$.*

The following example illustrates how aggregation operations which implement feasibility expansion can in turn expand elicibility.

Example 3.2. *Let the feature map be $\alpha = \alpha_2$, so that increasing the third output dimension contributes significantly to both features. We view the first two output dimensions as corresponding to two types of “bad” behavior, while dimension 3 corresponds to “good” behavior. Let C be a single constraint of the form $x_3 \leq x_1 + x_2$. The constraint captures how the model cannot produce the desirable dimension without also producing some of the undesirable dimension(s).*

The output $[0, 0, 1]$ is outside the feasibility set since it has only desirable dimensions and hence is not elicitable with any reward specification β . The system designer can still produce this output through intersection aggregation $\mathbf{x}^{(1)} = [1, 0, 1], \mathbf{x}^{(2)} = [0, 1, 1] \rightarrow \mathcal{A}_{\text{intersect}}(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = [0, 0, 1]$ (Proposition C.1 in Appendix C.1).

Mechanism 2: Overcoming Reward Specification Limitations. Even when an output is in the feasible set, the limitations of reward specification still restrict which outputs are elicitable. Aggregation can overcome the reward specification limitations faced by the system designer, as the next two mechanisms formalize.

Mechanism 2a: Support Expansion. One challenge due to reward specification limitations is the impossibility of eliciting outputs with a large support.² Aggregation can produce combine outputs with smaller supports into an output with a larger support, as the following mechanism formalizes.

Definition 3.3 (Support expansion). *An aggregation operation $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ implements support-expansion relative to i if $\mathcal{S}(x^{(A)}) \not\subseteq \mathcal{S}(x^{(i)})$.*

Aggregation operations which implement support-expansion can in turn expand elicibility, by producing outputs with larger supports than that are elicitable by a single agent, as the following example illustrates.

Example 3.4. *Let the feature map be $\alpha = \alpha_{0.6}$. Suppose that there are no constraints $C = \emptyset$, so elicibility challenges entirely stem from reward specification limitations. We will think of the first two dimensions as two aspects we would like our output to simultaneously capture.*

An output vector supported on both dimensions 1 and 2 cannot be elicited directly through reward design based on F_1 and F_2 (Prop C.2 in Appendix C.2). An output supported on just one of these two dimensions can be elicited through the reward function this dimension specializes in. However, any reward focusing on both features makes dimension 3 strictly preferred over the combination of dimensions 1 and 2.

The system designer can still produce vector $[1/2, 1/2, 0]$ supported on both dimensions 1 and 2 through addition aggregation $x^{(1)} = [1, 0, 0], x^{(2)} = [0, 1, 0] \rightarrow \mathcal{A}_{add}(x^{(1)}, x^{(2)}; [1/2, 1/2]) = [1/2, 1/2, 0]$ (Prop C.2 in Appendix C.2).

Mechanism 2b: Binding Set Contraction. The next mechanism overcomes reward specification limitations by taking advantage of the output limitations of the agent. Perhaps counterintuitively, the constraints on the output space can make it easier to elicit an output through a single reward. When a constraint is binding for an output vector, some reward-increasing directions become inaccessible to the agent, as these directions will lead to violation of the binding constraint. Aggregation can combine outputs with binding constraints into an output with fewer binding constraints.

Definition 3.5 (Binding set contraction). *An aggregation operation $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ implements binding set contraction relative to i if $\mathcal{V}(x^{(A)}) \not\supseteq \mathcal{V}(x^{(i)})$.*

Aggregation operations which implement binding set contraction can expand elicibility, as following example illustrates.

Example 3.6. *Let the feature map be $\alpha = \alpha_{0.2}$. As in the first example, we will think of x_3 to be a “good” dimension and x_1, x_2 to be “bad” dimensions. Let C be a single constraint of the form $x_1 + x_2 \leq x_3$. This constraint captures how the model cannot produce the bad dimension(s) without also producing some of the good dimension.*

The value of $q = 0.2$ is small leading to dimension 3 being inelicitable without the constraint (Proposition C.3 in Appendix C.3). The constraint allows us to elicit a vector with some amount of x_3 , but not a vector that has only x_3 . The intersection aggregation operation $x^{(1)} = [1/2, 0, 1/2], x^{(2)} = [0, 1/2, 1/2] \rightarrow \mathcal{A}_{intersect}(x^{(1)}, x^{(2)}) = [0, 0, 1/2]$.

Implementable Mechanisms by Intersection and Addition Aggregation. Our examples constructed problem instances that intersection aggregation can implement feasibility-expansion and binding-set contraction, while addition aggregation can implement support expansion. We turn to more general problem instances, and investigate whether each aggregation rule can implement these mechanism on any problem instance. We summarize our findings in Table 1, which shows fundamental limitations of each aggregation rule.

3.2 CONNECTIONS BETWEEN ELICITABILITY-EXPANSION AND MECHANISMS

Moving beyond the examples in Section 3.1, we more generally study the powers and limitations that these mechanisms provide for elicibility-expansion.

Necessity of these mechanisms. First, we show that if an aggregation operation expands elicibility for some feature weights matrix, it must implement at least one of the three mechanisms. Specifi-

²Kleinberg et al. (2019) studied this in single-agent environments without constraints.

cally, Theorem 3.7 shows that either the operation must implement feasibility-expansion or it must implement at least one of support-expansion or binding-set contraction for every output $x^{(i)}$.

Theorem 3.7. *Fix conic constraints C , and any aggregation operation $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$. If $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ is elicibility-expanding for some feature weights matrix α , then at least one of the following conditions holds:*

- $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ is feasibility-expanding relative to C (Definition 3.1).
- For each $i \in [K]$, $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ is either support-expanding relative to i (Definition 3.3) or binding set-contracting relative to i (Definition 3.5).

The proof of Theorem 3.7 builds on the technical tools we develop in Section 4 (i.e., Theorem 4.3).

Theorem 3.7 reveals a strong form of limitation for aggregation operations who do not implement at least one of the mechanisms (Definition 3.1, 3.5, and 3.3). Specifically, the result illustrates that if an operation does not implement the mechanisms according to the conditions in Theorem 3.7, then aggregation is not elicibility-expanding, regardless of the feature weights matrix. This result illustrates conditions under which aggregation adds no power to compound AI systems regardless of the level of prompt engineering limitations.

While these three natural mechanisms are necessary for aggregation to have power, these mechanisms are not sufficient in general. We demonstrate this and discuss some special cases where they are sufficient in Appendix B.3. In Section 4, we provide a general, necessary-and-sufficient condition that more precisely captures the power and limitations of aggregation.

4 CHARACTERIZING ELICITABILITY-EXPANSION IN GENERAL

In this section, we provide general characterizations of when an aggregation operation $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ is elicibility-expanding. We begin by analyzing, for a fixed feature weights matrix and feasibility constraints, whether a given aggregation operation expands elicibility (Section 4.1). We then turn to a more structural question: given only the feasibility constraints, what necessary and sufficient conditions ensure that aggregation operation is not elicibility-expanding for any feature weights matrix (Section 4.2)? These characterizations provide the technical foundation for our earlier results in Sections 3.1 and 3.2 which connected the mechanisms implemented by aggregation with elicibility-expansion.

4.1 CHARACTERIZING WHEN ELICITABILITY-EXPANSION SUCCEEDS

To analyze the power of aggregation, we characterize whether an aggregation operation is elicibility-expanding in a given problem-instance (i.e., given a feasibility set and feature weights). Our analysis generalizes the single-agent characterization from prior work (Kleinberg et al., 2019) to allow for output limitations (i.e., nontrivial constraints C). We then leverage this characterization to analyze aggregation operations.

Given a statistic $(S, V) = (S(x), \mathcal{V}(x))$, elicibility is determined by the structure of the set

$$\mathcal{B}_{S,V} = \underbrace{\{d \in \mathbb{R}^M : C_V d \leq 0\}}_{(1)} \cap \underbrace{\{d \in \mathbb{R}^M : d_j \geq 0 \forall j \in S^c\}}_{(2)} \cap \underbrace{\{1^t d < 0\}}_{(3)}.$$

The set $\mathcal{B}_{S,V}$ captures the set of directions along which the agent can move x while maintaining the constraints C (term (1)), maintaining nonnegativity constraints (term (2)), reducing ℓ_1 norm (term (3)). Specifically, elicibility expansion can be characterized by whether the sets $\mathcal{B}_{S,V}$ intersect with the set of feature-improving directions $\{d \in \mathbb{R}_{\geq 0}^M \mid \alpha d \geq 0\}$.

Theorem 4.1. *Fix conic constraints C , feature weights matrix α , and aggregation operation $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$. The aggregation operation $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ is elicibility-expanding if and only if both of the following conditions hold:*

- $\mathcal{B}_{S(x^{(i)}), \mathcal{V}(x^{(i)})} \cap \{d \in \mathbb{R}^M \mid \alpha d \geq 0\} = \emptyset$ for $i \in [K]$
- $\mathcal{B}_{S(x^{(A)}), \mathcal{V}(x^{(A)})} \cap \{d \in \mathbb{R}^M \mid \alpha d \geq 0\} \neq \emptyset$.

This characterizing condition depends on both the reward specification limitation (which reflect prompt engineering limitations) via α and the output limitation (which reflect model capability lim-

itations) via the conic constraints \mathbf{C} . This dependence highlights the role of both forms of limitations and their interplay in determining the power of aggregation.

Proof ideas. The core idea is that elicibility of a vector \mathbf{x} hinges on whether the feasible-perturbation set $\mathcal{B}_{\mathcal{S}(\mathbf{x}), \mathcal{V}(\mathbf{x})}$ intersects the feature-improving cone $\mathbf{d} : \boldsymbol{\alpha} \mathbf{d} \geq 0$ (Lemmas F.3–F.4). This lets us identify when each individual $\mathbf{x}^{(i)}$ is elicitable while the aggregate $\mathbf{x}^{(A)}$ is not—the condition for elicibility expansion.

One direction is immediate: a nonempty intersection yields a feasible direction that strictly improves every monotone reward, certifying that \mathbf{x} cannot be elicited. The converse is subtler: if the sets are disjoint, then some reward function elicits \mathbf{x} , and—as in Kleinberg et al. (2019)—it can be chosen to be linear in the features.

□

4.2 CHARACTERIZING WHEN ELICITABILITY-EXPANSION FAILS

To analyze the limitations of aggregation, we characterize conditions under which aggregation operations are not elicibility-expanding for *any* feature map. This represents a particularly strong form of limitation, as it rules out elicibility-expansion for all forms of reward specification limitations. The characterizing condition is stated below. We can interpret the condition as a failure of *strengthened* versions of the mechanisms. We discuss this connection to the mechanisms more in Appendix E.1.

Definition 4.2. [Limitation-characterizing condition] Fix constraints \mathbf{C} and aggregation operation $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$. We say that the **limitation-characterizing condition** is satisfied for $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ if and only if both of the following conditions are satisfied:

1. No feasibility expansion: $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ doesn't implement feasibility-expansion for \mathbf{C}
2. For all $\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})}$ with $\mathbf{d} \not\leq 0$, there exists $k \in [K]$ such that both of the following two conditions holds:
 - (a) No “strengthened support expansion” for k : For all $j \in \mathcal{S}(\mathbf{x}^{(k)})^c$, $-d_j - |\mathbf{1}^T \mathbf{d}| \leq 0$.
 - (b) No “strengthened binding-set contraction” for k : For all $\gamma^{(k)} \in \mathbb{R}_{\geq 0}^{|\mathcal{V}(\mathbf{x}^{(k)})|}$,

$$(\gamma^{(k)})^T \mathbf{C}_{\mathcal{V}(\mathbf{x}^{(k)})} \mathbf{d} - |\mathbf{1}^T \mathbf{d}| \cdot \left| \min_{j \in [M]} (\min(0, ((\gamma^{(k)})^T \mathbf{C}_{\mathcal{V}(\mathbf{x}^{(k)})}))_j) \right| \leq 0,$$

The following theorem shows that the limitation-characterizing condition is *necessary* for an aggregation operation to not expand elicibility under any feature map.

Theorem 4.3 (Necessary). Fix constraints \mathbf{C} . If the limitation-characterizing condition is satisfied, then there does not exist a feature weights matrix $\boldsymbol{\alpha}$ under which $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is elicibility-expanding.

The main idea of this theorem is showing that without a strengthened version of support-expansion or binding-set contraction, an aggregation operation is bound to have no power under all feature maps. Turning to the other direction, the next theorem shows that whenever the limitation-characterizing condition is violated, the aggregation operation is not limited in the strong sense. That is, the operation expands elicibility under *some* feature weights matrix. We prove this by constructing a feature weights matrix that makes aggregation elicibility-expanding whenever the limitation-characterizing condition does not hold.

Theorem 4.4 (Sufficient). Fix constraints \mathbf{C} , and aggregation operation $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$. If the limitation-characterizing condition is not satisfied for $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$, then there exist feature weights $\boldsymbol{\alpha}$ such that $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is elicibility-expanding.

5 DISCUSSION

In this work, we theoretically study how aggregating multiple copies of the same model gives access to a greater set of outputs than using only a single model. Building on a principal-agent framework, our results show how aggregation must implement one of three mechanisms—feasibility-expansion,

support expansion, and binding-set contraction—in order to expand the set of elicitable outputs. Although these mechanisms are not sufficient to ensure that aggregation adds power, we show a more precise condition formed from strengthening the mechanisms is sufficient.

Conceptual insights for system designers. Our results offer conceptual insights into when system designers benefit from specific aggregation operations in compound AI systems. Our results characterize how the interplay between prompt engineering limitations and model capability limitations affects which types of aggregation operations are useful. Specifically, aggregation not only overcomes model capability limitations (feasibility expansion), but also overcomes prompt engineering limitations through combining multiple output characteristics (support expansion) and through taking advantage of output-level limitations (binding set-contraction). Notably, even as model capabilities continue to improve, the latter two mechanisms mean that aggregation can still be useful to system designers. On the flip side, our results illustrate how aggregation operations that do not take advantage of these mechanisms offer no power, regardless of whether the system designer employs sophisticated or unsophisticated prompt engineering practices.

Connecting our mechanisms to empirical phenomena. We now discuss how the mechanisms that we identify in our work—feasibility expansion, support expansion, and binding set-contraction—connect to existing empirical phenomena observed for LLMs, and could inspire directions for future empirical work. Since aggregation is only powerful when individual models are limited on their own, we begin by outlining the single-model limitations underlying each mechanism and the empirical evidence supporting them.

- The power of feasibility expansion traces back to limitations in the types of outputs that individual models can generate: specifically, when models can’t exhibit certain (desirable) dimensions without exhibiting other (undesirable) dimensions as a side effect. This side effect has been empirically observed for safety versus overrefusal, where models which refuse a larger fraction of toxic outputs tend to refuse a larger fraction of safe outputs as a side effect (Cui et al., 2025). Similar side effects have been observed for alignment and hedging (Ouyang et al., 2022), and theoretically studied for creativity and factuality (Sinha et al., 2023).
- The power of support expansion traces back to challenges with eliciting outputs that perform along multiple dimensions at once in single-agent settings. This limitation has been empirically observed in cases where each dimension corresponds to a distinct user requirement. For example, prompts are often underspecified, since users may not include all of the requirements that they care about in the prompt (Yang et al., 2025). Moreover, even when users specify all their requirements, LLMs struggle to satisfy many requirements simultaneously (Wen et al., 2024; Guo et al., 2025).

We leave empirical validation of binding-set contraction—whose emergence depends on the interaction between prompt-engineering and model limitations—to future work. More broadly, since our results identify when aggregation enables these mechanisms, an important direction is to connect them to practice by testing whether real aggregation methods (e.g., debate (Du et al., 2024), prompt ensembling (Arora et al., 2023)) exhibit them. The single-model limitations discussed above suggest promising empirical settings where aggregation should add power.

Model limitations and extensions. Our stylized model, which builds on a classical principal-agent framework (Kleinberg et al., 2019), makes simplifying assumptions for tractability. First, while our analysis allows for nonlinear rewards R , we restrict the output constraints (i.e., model limitations) and the feature map (i.e., reward-specification limitations) to linear functional forms. Extending our model to allow for nonlinear limitations, which would complicate the structure of the agent’s optimization program, is an interesting direction for future work. Moreover, we also assume each agent’s reward depends only on its own outputs, though richer interdependencies may arise in repeated, multi-turn interactions (Du et al., 2024). Finally, while our analysis focuses on steering agents through reward design, it would be interesting to incorporate other choices, such as tool use and fine-tuning, that enable specialization in compound AI systems (BAIR Research Blog, 2024).

6 REPRODUCIBILITY STATEMENT

We provide full proofs of all of the results in the Appendix.

REFERENCES

- Anthropic Engineering. How we built our multi-agent research system. <https://www.anthropic.com/engineering/multi-agent-research-system>, 2024. Engineering blog. 1, 1.1, 2.3
- Simran Arora, Avanika Narayan, Mayee F. Chen, Laurel Orr, Neel Guha, Kush Bhatia, Ines Chami, Frederic Sala, and Christopher Ré. Ask me anything: A simple strategy for prompting language models. In *International Conference on Learning Representations (ICLR)*, 2023. URL <https://openreview.net/forum?id=bhUPJnS2g0X>. 1, 1.1, 5
- BAIR Research Blog. Compound ai systems. <https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/>, 2024. Blog post. 1, 1.1, 2.3, 5
- Patrick Bolton and Mathias Dewatripont. *Contract Theory*. MIT Press, Cambridge, MA, 2005. 1.1
- Philip Bond and Armando Gomes. Multitask principal-agent problems: Optimal contracts, fragility, and effort misallocation. *Journal of Economic Theory*, 144(1):175–211, 2009. 1.1
- Lingjiao Chen, Matei Zaharia, and James Zou. Frugalgpt: How to use large language models while reducing cost and improving performance. *Transactions on Machine Learning Research (TMLR)*, 2024. URL https://lingjiaochen.com/papers/2024_FrugalGPT_TMLR.pdf. 1.1
- Paul Christiano, Jan Leike, Tom B. Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. In *Advances in Neural Information Processing Systems 30 (NeurIPS)*, 2017. URL <https://arxiv.org/abs/1706.03741>. 1.1
- Natalie Collina, Surbhi Goel, Aaron Roth, Emily Ryu, and Mirah Shi. Emergent alignment via competition. *arXiv preprint*, 2025. 1.1
- Justin Cui, Wei-Lin Chiang, Ion Stoica, and Cho-Jui Hsieh. Or-bench: An over-refusal benchmark for large language models. *arXiv preprint arXiv:2405.20947*, 2025. URL <https://arxiv.org/abs/2405.20947>. 5
- Krishna Dasaratha, Benjamin Golub, and Anant Shah. Contracts for teams under imperfect observability. *Working Paper / Preprint*, 2024. Characterizes optimal pay allocation in team output settings using network centrality ideas. 1.1
- Dominique Demougin, David Encaoua, and Bernard Sinclair-Desgagné. A multi-tasking principal-agent perspective. *CESifo Working Paper Series*, (9753), 2022. 1.1
- Thomas G. Dietterich. Ensemble methods in machine learning. In *Multiple Classifier Systems (MCS 2000)*, volume 1857 of *Lecture Notes in Computer Science*, pp. 1–15. Springer, 2000. doi: 10.1007/3-540-45014-9_1. URL https://link.springer.com/chapter/10.1007/3-540-45014-9_1. 1.1
- Yilun Du, Shuang Li, Antonio Torralba, Joshua B. Tenenbaum, and Igor Mordatch. Improving factuality and reasoning in language models through multiagent debate. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*, volume 235 of *Proceedings of Machine Learning Research*, pp. 11733–11763. PMLR, 2024. URL <https://proceedings.mlr.press/v235/du24e.html>. 1, 1.1, 2.3, 5
- Matthew Gentzkow and Emir Kamenica. Bayesian persuasion with multiple senders and rich signal spaces. *Games and Economic Behavior*, 104:411–429, 2017. 1.1
- Sanford J. Grossman and Oliver D. Hart. An analysis of the principal-agent problem. *Econometrica*, 51(1):7–45, 1983. 1.1
- Zhengkang Guo, Wenhao Liu, Mingchen Xie, Jingwen Xu, Zisu Huang, Muzhao Tian, Jianhan Xu, Yuanzhe Shen, Qi Qian, Muling Wu, Xiaohua Wang, Changze Lv, He-Da Wang, Hu Yao, Xiaoqing Zheng, and Xuanjing Huang. RECAST: Expanding the boundaries of LLMs’ complex instruction following with multi-constraint data. *arXiv preprint arXiv:2505.19030*, 2025. URL <https://arxiv.org/abs/2505.19030>. 5

- Bengt Holmström. Moral hazard in teams. *The Bell Journal of Economics*, 13(2):324–340, 1982. 1.1
- Bengt Holmström. Moral hazard and observability. *The Bell Journal of Economics*, 10(1):74–91, 1979. 1.1
- Bengt Holmström and Paul Milgrom. Multitask principal-agent analyses: Incentive contracts, asset ownership, and job design. *Journal of Law, Economics, & Organization*, 7:24–52, 1991. 1.1
- Audrey Huang, Adam Block, Qinghua Liu, Nan Jiang, Akshay Krishnamurthy, and Dylan J Foster. Is best-of-n the best of them? coverage, scaling, and optimality in inference-time alignment. *arXiv preprint arXiv:2503.21878*, 2025. 1.1
- Athul Paul Jacob and Jacob Andreas. The consensus game: Language model generation via consensus seeking. In *International Conference on Learning Representations (ICLR)*, 2024. URL <https://openreview.net/forum?id=zEwWZk7c0Z>. 1.1
- Erik Jones, Anca Dragan, and Jacob Steinhardt. Adversaries can misuse combinations of safe models. In *International Conference on Machine Learning (ICML)*, 2025a. 1.1
- Erik Jones, Arjun Patrawala, and Jacob Steinhardt. Uncovering gaps in how humans and llms interpret subjective language. *arXiv preprint arXiv:2503.04113*, 2025b. 2.4
- Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Edward Grefenstette, Samuel R. Bowman, Tim Rocktäschel, and Ethan Perez. Debating with more persuasive llms leads to more truthful answers. In *Proceedings of the 41st International Conference on Machine Learning (ICML)*, volume 235 of *Proceedings of Machine Learning Research*, pp. 23662–23733. PMLR, 2024. URL <https://proceedings.mlr.press/v235/khan24a.html>. 1
- Jon Kleinberg, Sendhil Mullainathan, and Manish Raghavan. How do classifiers induce agents to invest effort strategically? In *Proceedings of the 2019 ACM Conference on Economics and Computation (EC)*, pp. 825–844. ACM, 2019. doi: 10.1145/3328526.3329584. 1, 1.1, 2, 2.2, 2.3, 2, 4.1, 4.1, 5, F.1
- Krishna K. Ladha. The condorcet jury theorem, free speech, and correlated votes. *American Journal of Political Science*, 36(3):617–634, 1992. doi: 10.2307/2111584. 1.1
- Jean-Jacques Laffont and David Martimort. *The Theory of Incentives: The Principal-Agent Model*. Princeton University Press, Princeton, NJ, 2002. 1.1
- Edward P. Lazear and Sherwin Rosen. Rank-order tournaments as optimum labor contracts. *Journal of Political Economy*, 89(5):841–864, 1981. 1.1
- Xuechen Li, Tianyi Zhang, Yann Dubois, Rohan Taori, Ishaan Gulrajani, Carlos Guestrin, Percy Liang, and Tatsunori B. Hashimoto. AlpacaEval: An automatic evaluator of instruction-following models. https://github.com/tatsu-lab/alpaca_eval, 5 2023. 1, G
- Nancy A. Lynch. *Distributed Algorithms*. Morgan Kaufmann, San Francisco, CA, 1996. ISBN 1558603484. 1.1
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35: 27730–27744, 2022. 5
- Ori Press, Andreas Hochlehnert, Ameya Prabhu, Vishaal Udandaraao, Ofir Press, and Matthias Bethge. Citeme: Can language models accurately cite scientific claims? *Advances in Neural Information Processing Systems*, 37:7847–7877, 2024. 2.4
- Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, 2019. URL <https://arxiv.org/abs/1908.10084>. We use the all-mpnet-base-v2 model from the Sentence-Transformers library. 2, G

- Ritwik Sinha, Zhao Song, and Tianyi Zhou. A mathematical abstraction for balancing the trade-off between creativity and reality in large language models. *arXiv preprint arXiv:2306.02295*, 2023. URL <https://arxiv.org/abs/2306.02295>. 5
- Margaret E. Slade. Multitask agency and contract choice. *International Economic Review*, 37(2): 465–486, 1996. 1.1
- Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In *Proceedings of the Tenth International Conference on Machine Learning (ICML)*, pp. 330–337. Morgan Kaufmann, 1993. 1.1
- Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*, 2023a. 2.4
- Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. In *International Conference on Learning Representations (ICLR)*, 2023b. 1.1
- Bosi Wen, Pei Ke, Xiaotao Gu, Lindong Wu, Hao Huang, Jinfeng Zhou, Wenchuang Li, Binxin Hu, Wendy Gao, Jiabin Xu, Yiming Liu, Jie Tang, Hongning Wang, and Minlie Huang. Benchmarking complex instruction-following with multiple constraints composition. *arXiv preprint arXiv:2407.03978*, 2024. URL <https://arxiv.org/abs/2407.03978>. 5
- Chenyang Yang, Yike Shi, Qianou Ma, Michael Xieyang Liu, Christian Kästner, and Tongshuang Wu. What prompts don’t say: Understanding and managing underspecification in LLM prompts. *arXiv preprint arXiv:2505.13360*, 2025. URL <https://arxiv.org/abs/2505.13360>. 2.4, 5
- Bohan Zhang, Xiaokang Zhang, Jing Zhang, Jifan Yu, Xiaokun Zhang, Tejas Srinivasan, Yifan Mai, Juntao Li, Yoon Kim, and Minjoon Seo. Cot-based synthesizer: Enhancing llm performance through answer synthesis, 2025. URL <https://arxiv.org/abs/2501.01668>. 1.1, 2.3
- Simon Zhuang and Dylan Hadfield-Menell. Consequences of misaligned ai. *Advances in Neural Information Processing Systems*, 33:15763–15773, 2020. 1.1

A LLM USAGE STATEMENT

We used GPT-5 and Claude Opus 4.1 to gather related work, get ideas for proofs, and to edit prose. All of the work done by LLMs was verified by the (human) authors on this paper.

B ADDITIONAL DETAILS FOR SECTION 3

B.1 ADDITIONAL DETAILS OF SECTION 3.1

Intersection aggregation does not implement support expansion for any problem instance, as the following result formalizes.

Proposition B.1 (Intersection does not expand support). *Consider any aggregation operation of the form $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)} = \mathcal{A}_{\text{intersect}}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})$. For any $i \in [K]$, this aggregation operation does not implement support-expansion relative to i .*

Proposition B.1 follows from the fact that the support of $\mathcal{A}_{\text{intersect}}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})$ is always a subset of the support of each $\mathbf{x}^{(i)}$.

Intersection aggregation can implement feasibility-expansion as shown in Example 3.2 and binding set-contraction as shown in Example 3.6. In fact, these examples go one step further and demonstrate that elicibility expansion is achievable via these mechanisms.

Addition aggregation does not implement feasibility expansion for any problem instance, as the following result formalizes.

Proposition B.2 (Addition cannot expand feasibility). *Consider constraints C . Any aggregation operation of the form $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)} = \mathcal{A}_{\text{add}}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}; \mathbf{w})(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})$ does not implement feasibility expansion relative to C .*

Proposition B.2 directly follows from the fact that the constraint set C is conic.

On the other hand, addition aggregation operations can implement the other two mechanisms. **Example 3.4** already constructed a problem instance where addition aggregation implements support expansion. The next example constructs a problem instance where addition aggregation can implement binding set contraction (**Definition 3.5**) and achieve elicibility-expansion for some feature mapping.

Example B.3 (Addition can result in binding set contraction). *Consider the constraint matrix*

$$C = \begin{pmatrix} 1 & -1 & 0 \\ 1 & -\frac{1}{4} & -1 \end{pmatrix},$$

and consider vectors $\mathbf{x}^{(1)} = (1, 1, 2)$ and $\mathbf{x}^{(2)} = (2, 4, 1)$. Note that they are both feasible and $\mathbf{x}^{(1)}$ is binding in the first constraint and $\mathbf{x}^{(2)}$ in the second. Their sum is $\mathcal{A}_{\text{add}}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}; \mathbf{w})(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}; [1, 1]) \rightarrow \mathbf{x}^{(A)} = (3, 5, 3)$, which is also feasible but does not have any binding constraints.

B.2 ADDITIONAL DETAILS OF SECTION 3.2

	Feasibility Expansion	Support Expansion	Binding Set Contraction
Intersection aggregation	✓ (Example 3.2)	× (Proposition B.1)	✓ (Example 3.6)
Addition aggregation	× (Proposition B.2)	✓ (Example 3.4)	✓ (Example B.3)

Table 1: Implementability of mechanisms in Section 3.1 for the intersection aggregation rule equation 1 and additional aggregation rule equation 2. The symbol ✓ denotes that there exists a problem instance where the aggregation rule implements that mechanism. The symbol × denotes that the aggregation rule does not implement the mechanism for any problem instance.

Proposition B.4. *Fix conic constraints C , and any $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$. Suppose that $\mathcal{V}(\mathbf{x}^{(A)}) = \mathcal{V}(\mathbf{x}^{(1)}) = \dots = \mathcal{V}(\mathbf{x}^{(K)}) = \emptyset$, and suppose that $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is not feasibility-expanding. Suppose also that there do **not** exist witnesses $j(i) \in \mathcal{S}(\mathbf{x}^{(A)}) \setminus \mathcal{S}(\mathbf{x}^{(i)})$ for each $i \in [K]$ such that $\{j(i) \mid i \in [K]\} \neq [M]$. Then, $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is not elicibility-expanding for any α .*

Proposition B.5. *There exists an aggregation operation $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \rightarrow \mathbf{x}^{(A)}$ and a set of conic constraints C such that $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ implements binding-set contraction relative to i for every $i \in [K]$. However, $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \rightarrow \mathbf{x}^{(A)}$ is not elicibility-expanding for any feature map α .*

Proposition B.6. *Fix $C = \emptyset$. There exists an aggregation operation $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \rightarrow \mathbf{x}^{(A)}$ such that $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \rightarrow \mathbf{x}^{(A)}$ implements support-expansion relative to i for every $i \in [K]$. However, $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \rightarrow \mathbf{x}^{(A)}$ is not elicibility-expanding for any feature map α .*

Proposition B.7. *Fix conic constraints C and $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$. Suppose that $M \geq 2$, and $\mathcal{S}(\mathbf{x}^{(A)}) = [M]$. Suppose that there exist witnesses $\ell(i) \in \mathcal{V}(\mathbf{x}^{(i)}) \setminus \mathcal{V}(\mathbf{x}^{(A)})$ such that there exists \mathbf{d} such that $C_{\ell(i)} \mathbf{d} + (\min_{j \in [M]} C_{\ell(i), j}) \cdot \mathbf{1}^t \mathbf{d} > 0$ for all $i \in [K]$, $\mathbf{1}^t \mathbf{d} < 0$, and $C_{\mathcal{V}(\mathbf{x}^{(A)})} \mathbf{d} \leq 0$. Then, $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is elicibility-expanding for some feature weights matrix α .*

B.3 PARTIAL SUFFICIENCY OF MECHANISMS IN CONCRETE INSTANCES

In **Theorem 3.7**, we showed how the mechanisms are necessary for an aggregation operation to have power. We now turn to analyzing when mechanisms are sufficient for guaranteeing the power of aggregation. We focus on a weak form of power that only requires that aggregation expands

elicitability for some feature weights matrix, taking a negation of of the limitation show in Theorem 4.3. (We defer an analysis of the role of the feature weights matrix to Section 4.1.)

We first show that feasibility expansion guarantees this form of power, providing a partial converse of Theorem 3.7.

Proposition B.8. *Fix conic constraints C . If an aggregation operation $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ implements feasibility-expansion, then there exists a feature map α such that $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is elicibility-expanding.*

We now turn to support expansion and binding-set contraction. Interestingly, even if an aggregation operation implements support-expansion for every $i \in [K]$, the aggregation still may not elicibility-expanding for any feature weights matrix (Proposition B.6). Similarly, binding-set contraction also does not guarantee that aggregation has power (Proposition B.5).

Nonetheless, we show stronger conditions under which support expansion and binding-set contraction do guarantee that aggregation expands elicibility for some feature map. For support expansion, the main requirement is a global form of support expansion across outputs i , requiring that the “witnesses” don’t span all of the output dimensions.³

Proposition B.9. *Fix conic constraints C , and any $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$. Suppose that there exist witnesses $j(i) \in \mathcal{S}(\mathbf{x}^{(A)}) \setminus \mathcal{S}(\mathbf{x}^{(i)})$ for each $i \in [K]$ such that $\{j(i) \mid i \in [K]\} \neq [M]$. Suppose that $\mathcal{V}(\mathbf{x}^{(A)}) = \emptyset$. Then, $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is elicibility-expanding for some α .*

Turning to binding-set contraction, the main requirement is again a global form of binding-set contraction across outputs i which links witnesses (i.e., a constraint in $\mathcal{V}(\mathbf{x}^{(i)}) \setminus \mathcal{V}(\mathbf{x}^{(A)})$ for each $i \in [K]$) together (Proposition B.7). A global variant of support expansion and binding-set contraction also emerges in our characterizations in Section 4.

C PROOFS FOR SECTION 3

Recall that the examples in this section use the feature weights matrix $\alpha_q := \begin{bmatrix} 1 & 0 & q \\ 0 & 1 & q \end{bmatrix}$.

C.1 ANALYSIS OF EXAMPLE 3.2

Proposition C.1. *For the feature weights matrix α_2 and constraint matrix with the row $x_3 \leq x_1 + x_2$ in Example 3.2, $\mathbf{x}^{(1)} = [1, 0, 1]$, $\mathbf{x}^{(2)} = [0, 1, 1] \rightarrow \mathbf{x}^{(A)} = [0, 0, 1]$ is elicibility-expanding.*

Proof. From the construction, it is easy to see that $\mathbf{x}^{(1)}$ can be elicited with a linear reward function $[1, 0, 0]$ equal to the F_1 and budget level $E = 2$ and $\mathbf{x}^{(1)}$ can be elicited with a linear reward function $[0, 1, 0]$ equal to the F_1 and budget level $E = 2$.

Let us use our characterization Theorem 4.1 to formally elicibility-expansion.

For $\mathbf{x}^{(i)}$, the support set $\mathcal{S}(\mathbf{x}^{(i)})$ is $\{i\}$. The constraint is binding on both $\mathbf{x}^{(i)}$. The set $\mathcal{B}_{\mathcal{S}(\mathbf{x}^{(1)}), \mathcal{V}(\mathbf{x}^{(1)})}$ is $\{\mathbf{d} : \mathbf{d}_3 \leq \mathbf{d}_1 + \mathbf{d}_2, \mathbf{d}_2 \geq 0, \mathbf{d}_1 + \mathbf{d}_2 + \mathbf{d}_3 < 0\}$. For any \mathbf{d} in this set, $\mathbf{d}_3 < 0$ and $\mathbf{d}_1 + \mathbf{d}_2 < -\mathbf{d}_3$. The set $\mathcal{B}_{\mathcal{S}(\mathbf{x}^{(2)}), \mathcal{V}(\mathbf{x}^{(2)})}$ is $\{\mathbf{d} : \mathbf{d}_3 \leq \mathbf{d}_1 + \mathbf{d}_2, \mathbf{d}_1 \geq 0, \mathbf{d}_1 + \mathbf{d}_2 + \mathbf{d}_3 < 0\}$. For any \mathbf{d} in this set, $\mathbf{d}_3 < 0$ and $\mathbf{d}_1 + \mathbf{d}_2 < -\mathbf{d}_3$.

Now consider the set of feature-improving direction $\{\mathbf{d} : \mathbf{d}_1 + 2\mathbf{d}_3 \geq 0, \mathbf{d}_2 + 2\mathbf{d}_3 \geq 0\}$. For any \mathbf{d} in this set, $\mathbf{d}_1 + \mathbf{d}_2 \geq -4\mathbf{d}_3$.

All three conditions $\mathbf{d}_3 < 0$, $\mathbf{d}_1 + \mathbf{d}_2 < -\mathbf{d}_3$, and $\mathbf{d}_1 + \mathbf{d}_2 \geq -4\mathbf{d}_3$ cannot be satisfied since for $\mathbf{d}_3 < 0$, $-\mathbf{d}_3 > -4\mathbf{d}_3$. Hence there is no intersection between feasibility improving directions and features improving directions and $\mathbf{x}^{(1)}$ is elicitable. Similarly, $\mathbf{x}^{(2)}$ is also elicitable.

$\mathbf{x}^{(3)}$ is not feasible and hence not elicitable. This shows that $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \rightarrow \mathbf{x}^{(3)}$ is elicibility-expanding by implementing feasibility-expansion. \square

³The fact the witnesses cannot span all of the output dimensions condition also turns to be a necessary condition for aggregation to not be powerless (Proposition E.2).

C.2 ANALYSIS OF EXAMPLE 3.4

Proposition C.2. For the feature weights matrix $\alpha_{0.6}$ and null constraint matrix *Example 3.4*, $\mathbf{x}^{(1)} = [1, 0, 0]$, $\mathbf{x}^{(2)} = [0, 1, 0] \rightarrow \mathbf{x}^{(A)} = [1/2, 1/2, 0]$ is elicibility-expanding.

Proof. The set of directions $\mathcal{B}_{\mathcal{S}(\mathbf{x}^{(1)}), \mathcal{V}(\mathbf{x}^{(1)})} = \{\mathbf{d} : \mathbf{d}_2 \geq 0, \mathbf{d}_3 \geq 0, \mathbf{d}_1 + \mathbf{d}_2 + \mathbf{d}_3 < 0\}$. And the set of feature-improving directions is $\mathcal{A} = \{\mathbf{d} : \mathbf{d}_1 + 0.6\mathbf{d}_3 \geq 0, \mathbf{d}_2 + 0.6\mathbf{d}_3 \geq 0\}$.

$\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(1)}), \mathcal{V}(\mathbf{x}^{(1)})}$ means that $\mathbf{d}_1 < -(\mathbf{d}_2 + \mathbf{d}_3) < -\mathbf{d}_3$ and $\mathbf{d}_3 \geq 0$. $\mathbf{d} \in \mathcal{A}_1$ means that $\mathbf{d}_1 \geq -0.6\mathbf{d}_3$. These three conditions cannot be simultaneously showing that $\mathbf{x}^{(1)}$ is elicitable due to empty intersection of \mathcal{A} and \mathcal{B}_1 . Symmetrically, we can also show that $\mathbf{x}^{(2)}$ is also elicitable.

Now let us argue that $\mathbf{x}^{(A)} = [1/2, 1/2, 0]$ is not elicitable. The feasibility improving directions set is $\mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})} = \{\mathbf{d} : \mathbf{d}_3 \geq 0, \mathbf{d}_1 + \mathbf{d}_2 + \mathbf{d}_3 < 0\}$. Consider $\mathbf{d} = [-0.6, 0.6, 1]$. $\mathbf{d} \in \mathcal{A} \cap \mathcal{B}_A$. This shows that $\mathbf{x}^{(A)}$ is not elicitable. □

C.3 ANALYSIS OF EXAMPLE 3.6

Proposition C.3. For the feature weights matrix $\alpha_{0.2}$ and conic constraint matrix with one constraint $x_1 + x_2 \leq x_3$ from *Example 3.6*, $\mathbf{x}^{(1)} = [1, 0, 1]$, $\mathbf{x}^{(2)} = [0, 1, 1] \rightarrow \mathbf{x}^{(A)} = [0, 0, 1]$ is elicibility-expanding.

Proof of Proposition C.3. The feature-improving directions are the set $\mathcal{A} = \{\mathbf{d} : \mathbf{d}_1 + 0.2\mathbf{d}_3 \geq 0, \mathbf{d}_2 + 0.2\mathbf{d}_3 \geq 0\}$.

The constraint is binding at both $\mathbf{x}^{(1)}$ and $\mathbf{x}^{(2)}$. The feasibility improving directions are $\mathcal{B}_{\mathcal{S}(\mathbf{x}^{(1)}), \mathcal{V}(\mathbf{x}^{(1)})} = \{\mathbf{d} : \mathbf{d}_1 + \mathbf{d}_2 \leq \mathbf{d}_3, \mathbf{d}_2 \geq 0, \mathbf{d}_1 + \mathbf{d}_2 + \mathbf{d}_3 < 0\}$.

If $\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(1)}), \mathcal{V}(\mathbf{x}^{(1)})}$, then $\mathbf{d}_1 + \mathbf{d}_2 + \mathbf{d}_3 < 0$ and $\mathbf{d}_1 + \mathbf{d}_2 + \mathbf{d}_3 \leq 2\mathbf{d}_3$. This implies that $\mathbf{d}_3 < 0$. If $\mathbf{d} \in \mathcal{A}$, then $\mathbf{d}_1 \geq -0.2\mathbf{d}_3$ and $\mathbf{d}_2 \geq -0.2\mathbf{d}_3$. If all the conditions are satisfied simultaneously, then $\mathbf{d}_1 > 0$ and $\mathbf{d}_2 > 0$. This contradicts $\mathbf{d}_1 + \mathbf{d}_2 \leq \mathbf{d}_3 < 0$.

The conic constraint is not binding at $\mathbf{x}^{(A)}$. Now consider the feasibility improving directions of $\mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})} = \{\mathbf{d} : \mathbf{d}_1 \geq 0, \mathbf{d}_2 \geq 0, \mathbf{d}_1 + \mathbf{d}_2 + \mathbf{d}_3 < 0\}$. The vector $\mathbf{d} = (0.2, 0.2, -1) \in \mathcal{A} \cap \mathcal{B}_A$ demonstrating that $\mathbf{x}^{(A)}$ is not elicitable. □

D PROOFS IN SECTION 3.2

D.1 PROOF OF THEOREM 3.7

Proof of Theorem 3.7. We show this as a corollary of Theorem 4.3. We will prove this by showing that when both of the conditions in the theorem are violated, the limitation-characterizing condition is satisfied and hence $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ cannot be elicibility-expanding.

One of the condition of the limitation-characterizing condition is the lack of feasibility-expansion which is implied by the violation of the theorem's condition. We will show that the other condition of the limitation-characterizing condition also holds.

When the second condition of the theorem is violated, there exists $i \in [K]$ with respect to which $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is neither support-expanding nor binding-set contracting. That is, there is an i such that $\mathcal{V}(\mathbf{x}^{(i)}) \subseteq \mathcal{V}(\mathbf{x}^{(A)})$ and $\mathcal{S}(\mathbf{x}^{(i)}) \supseteq \mathcal{S}(\mathbf{x}^{(A)})$.

For every $\mathbf{d} \in \{C_{\mathcal{V}(\mathbf{x}^{(A)})}\mathbf{d} \leq 0, \mathbf{d}_{\mathcal{S}(\mathbf{x}^{(A)})^c} \geq 0, \mathbf{1}^\top \mathbf{d} = -1\}$, $C_{\mathcal{V}(\mathbf{x}^{(i)})}(\mathbf{d}) \leq 0$ and $\mathbf{d}_{\mathcal{S}(\mathbf{x}^{(i)})^c} \geq 0$, since the rows of $C_{\mathcal{V}(\mathbf{x}^{(i)})}$ are a subset of the rows in $C_{\mathcal{V}(\mathbf{x}^{(A)})}$ and similarly, the rows in $\mathbf{d}_{\mathcal{S}(\mathbf{x}^{(i)})^c} \geq 0$ are a subset of the rows in $\mathbf{d}_{\mathcal{S}(\mathbf{x}^{(A)})^c}$. Hence for any $\gamma^{(i)} \in \mathbb{R}_{\geq 0}^{|\mathcal{V}_i|}$, $(\gamma_i)^\top C_{V_i} \mathbf{d} - \|(\gamma_i^\top C_{V_i})_-\|_\infty \leq 0$. □

D.2 PROOF OF PROPOSITION B.8

Proof of Proposition B.8. This follows from Theorem 4.4. \square

D.3 PROOF OF PROPOSITION B.6

Proof of Proposition B.6. This follows from Proposition E.2. \square

D.4 PROOF OF PROPOSITION B.5

Proof of Proposition B.5. Consider a problem with two output dimensions having the following two constraints: 1) $c_1 : x_1 - x_2 \leq 0$, 2) $c_2 : -2x_1 + x_2 \leq 0$. Consider an aggregation operation $\mathbf{x}^{(1)} = (1/2, 1/2)$, $\mathbf{x}^{(2)} = (1/2, 2/3) \rightarrow \mathbf{x}^{(A)} = (5/12, 7/12)$, where the binding constraints sets are $\mathcal{V}_{\mathbf{x}^{(i)}} = \{c_i\}$ for $i \in \{1, 2\}$ and $\mathcal{V}_{\mathbf{x}^{(A)}} = \emptyset$.

In this example, we will show how the limitations-characterization condition holds, meaning that the operation cannot be elicibility-expanding for any feature map α .

For any $\gamma_i \geq 0$ and \mathbf{d} , $\gamma_i c_i \mathbf{d} - \gamma_i \|(c_i)_\infty\|_\infty > 0$ if and only if $c_i \mathbf{d} - \|(c_i)_\infty\|_\infty > 0$. In this example, the existence of γ_i for this inequality to be satisfied for each i corresponds to the conditions that 1) $d_1 - d_2 > 1$ and $-2d_1 + d_2 > 2$. Since there are no elements outside the support of the vectors, there are no additional conditions to check for the limitations-characterization condition.

These two conditions imply that $1 + d_2 < d_1 < (d_2 - 2)/2$. Hence the conditions can only be satisfied when $1 + d_2 < (d_2 - 2)/2$. This is only satisfied when $d_2 < -4$ and this implies $d_1 < -3$. Hence the two conditions being satisfied means $1^\top \mathbf{d} < -7$. So the set of \mathbf{d} such that $1^\top \mathbf{d} = -1$ cannot intersect with the set of \mathbf{d} satisfying the two conditions. \square

D.5 PROOF OF PROPOSITION E.2

Proof of Proposition E.2. Suppose that $M \geq 2$, and $\mathcal{S}(\mathbf{x}^{(A)}) = [M]$.

We apply Theorem 4.3. It suffices to show that the limiting-characterization condition (Definition 4.2) is satisfied. By assumption, we know that the aggregation operation is not feasibility expanding. It suffices to show that there does not exist \mathbf{d} such that for every $i \in [K]$ there exists $j(i) \in \mathcal{S}(\mathbf{x}^{(i)})^c$ such that $-d_{j(i)} - |1^\top \mathbf{d}| > 0$.

Let's show the contrapositive: assume that there exists \mathbf{d} such that for every $i \in [K]$ there exists $j(i) \in \mathcal{S}(\mathbf{x}^{(i)})^c$ such that $-d_{j(i)} - |1^\top \mathbf{d}| > 0$. Since $\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})}$, we know that $d_{j'} \geq 0$ for $j' \notin \mathcal{S}(\mathbf{x}^{(A)})$ and we know that $1^\top \mathbf{d} < 0$. If $j \notin \mathcal{S}(\mathbf{x}^{(A)})$, then note that $-d_j < 0$, so this means that $j(i) \in \mathcal{S}(\mathbf{x}^{(A)})$. Putting this together, we see that $j(i) \in \mathcal{S}(\mathbf{x}^{(A)}) \setminus \mathcal{S}(\mathbf{x}^{(i)})$.

It suffices to show that $\{j(i) \mid i \in [K]\} \neq [M]$. Assume for sake of contradiction that $\{j(i) \mid i \in [K]\} = [M]$. Then since we know that $0 < -d_{j(i)} - |1^\top \mathbf{d}| = 1^\top \mathbf{d} - d_{j(i)}$, if we add up all of these equations in the set $\{j(i) \mid i \in [K]\}$, we would obtain that $0 < M \cdot 1^\top \mathbf{d} - \sum_j d_j = (M - 1) \cdot 1^\top \mathbf{d} - \sum_j d_j$, which means that $1^\top \mathbf{d} > 0$ which is a contradiction. \square

D.6 PROOF OF PROPOSITION B.9

Proof of Proposition B.9. We apply Theorem 4.4. It suffices to show that the limiting-characterization condition (Definition 4.2) is violated. Let \mathbf{d} be the vector such that $d_{j(i)} = -1$ for all $i \in [K]$, $d_j = |\{j(i) \mid i \in [K]\}| - 0.5$ for some $j \notin \{j(i) \mid i \in [K]\}$, and 0 elsewhere. It follows from definition that $\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})}$. It suffices to show that for all $i \in [K]$, it holds that:

$$-d_{\ell(i)} - |1^\top \mathbf{d}| > 0.$$

Using that $1^\top \mathbf{d} < 0$, this can be rewritten as:

$$-d_{\ell(i)} - |1^\top \mathbf{d}| = \sum_{j \neq \ell(i)} d_j = |\{j(i) \mid i \in [K]\}| - 0.5 - |\{j(i) \mid i \in [K]\}| + 1 = 0.5 > 0,$$

as desired. □

D.7 PROOF OF PROPOSITION B.7

Proof of Proposition B.7. We apply Theorem 4.4. It suffices to show that the limiting-characterization condition (Definition 4.2 is violated). For each $i \in [K]$, we take $\gamma^{(i)}$ to be the 1-hot vector with the 1 on the $\ell(i)$ th condition. Let \mathbf{d} be the vector given by the condition in the theorem statement. It follows immediately that $\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})}$. It suffices to show that for all $i \in [K]$, it holds that:

$$C_{\ell(i)}\mathbf{d} - |\mathbf{1}^\top \mathbf{d}| \left| \min_{j \in [M]} \min(0, C_{\ell(i),j}) \right|.$$

Using that $\mathbf{1}^\top \mathbf{d} < 0$ and using that $C_{\ell(i)}$ has at least one negative coordinate, this can be written as:

$$C_{\ell(i)}\mathbf{d} + \mathbf{1}^\top \mathbf{d} \left| \min_{j \in [M]} C_{\ell(i),j} \right| > 0,$$

which we know holds. □

E ADDITIONAL DETAILS FOR SECTION 4

E.1 CONNECTING THE LIMITATION-CHARACTERIZING CONDITION TO MECHANISMS

The limitation-characterization condition requires two sub-conditions to hold. The first is lack of implementation of feasibility expansion. We can interpret the second sub-condition as not implementing either a strengthening of support-expansion or a strengthening of binding-set contraction.

To demonstrate the connection between the limitation-characterizing condition and the mechanisms, we will first show that when none of the mechanisms are implemented, the limitation characterizing condition is satisfied. We will later discuss the ways in which the limitation-characterizing is related to strengthened versions of the mechanisms.

The following result shows that none of the mechanisms being implemented implies that the limitation-characterizing condition is satisfied. This result immediately implies Theorem 3.7 (i.e., that implementing at least one of these mechanisms is necessary for elicibility-expansion).

Proposition E.1. *Fix conic constraints \mathbf{C} , and any aggregation operation $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ where each $\mathbf{x}^{(k)}$ is feasible i.e., $\mathbf{C}\mathbf{x}^{(k)} \leq 0$, for every $k \in [K]$. If $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ satisfies both of the following conditions, then $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ satisfies the limitation-characterizing condition (Definition 4.2).*

- $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is not feasibility-expanding relative to \mathbf{C} (Definition 3.1).
- There exists $k \in [K]$ such that $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is neither support-expanding relative to i (Definition 3.3) nor binding set-contracting relative to i (Definition 3.5).

Proof. If none of the mechanisms are implemented, then the first condition of the limitation-characterization condition, which is lack of implementation of feasibility expansion automatically holds. We will now show the second condition of the limitation condition also holds. The second condition requires two sub-conditions Condition 2a and Condition 2b in in Definition 4.2 to hold for some $k \in [K]$. We will show that each of these conditions are implemented by lack of support-expansion and lack of binding-set contraction respectively.

No support-expansion relative to k implies Condition 2a in Definition 4.2 relative to k . When support-expansion is not implemented relative to k , all $j^{(k)} \in \mathcal{S}(\mathbf{x}^{(A)})$ also belongs to $\mathcal{S}(\mathbf{x}^{(k)})$. For all $\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})}$ and all $j^{(k)} \in \mathcal{S}(\mathbf{x}^{(k)})^c \subseteq \mathcal{S}(\mathbf{x}^{(A)})^c$, $\mathbf{d}_{j^{(k)}} \geq 0$. Hence $\mathbf{d}_{j^{(k)}} + |\mathbf{1}^\top \mathbf{d}|$ which is even larger than $\mathbf{d}_{j^{(k)}}$ is ≥ 0 for every $j^{(k)} \in \mathcal{S}(\mathbf{x}^{(k)})^c$. This is the Condition 2a relative to k .

No binding-set contraction relative to k implies Condition 2b in Definition 4.2 relative to k . No binding set contraction means that all constraints in $\mathcal{V}(\mathbf{x}^{(A)})$ are also in $\mathcal{V}(\mathbf{x}^{(k)})$. Hence every $\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})}$ satisfies all conic constraints $\ell \in \mathcal{V}(\mathbf{x}^{(k)})$. \mathbf{d} also satisfies all non-negatively weighted sums of conic constraints in $\ell \in \mathcal{V}(\mathbf{x}^{(k)})$. Condition 2b relative to k in Definition 4.2 only requires approximately satisfying the weighted sums of constraints and hence is implied by no binding-set contraction relative to k . \square

How the limitation-characterizing condition strengthens mechanisms the mechanisms. Next, we will describe how the limitation-characterizing condition, specifically Conditions 2a,2b in Definition 4.2 are failures of strictly strengthened versions of the mechanisms. This makes the limitation-characterizing condition a strictly weaker condition to be satisfied compared to failure of all mechanisms. The conditions 2a,2b of the limitation-characterizing condition are strengthenings in two ways. The first is due to requiring violations by minimum margins. Note that from the proof of Proposition E.1, just violation without any minimum margin requirement suffices for the mechanisms. Another way that these conditions are stronger is the *joint* requirement across all $k \in [K]$. We require that the same $\mathbf{d} \in \mathcal{B}_{\mathcal{S}(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})}$ witnesses the violation by a margin for every $k \in [K]$.

In some special cases, the limitation-characterizing conditions correspond exactly to not implementing any of the mechanisms, instead of not implementing strengthenings. One special case is when no vector in the aggregation operation has any binding conic constraints. This holds when there are no conic constraints. In this special case, even the regular, non-strengthened form of binding-set contraction cannot kick in. We can show that in this special case, the limitation-characterizing condition is either not feasibility-expansion or not the usual, non-strengthened support-expansion as long as a particular edge case does not occur.

Corollary E.2. Fix conic constraints \mathcal{C} , and any $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$. Suppose that $\mathcal{V}(\mathbf{x}^{(A)}) = \mathcal{V}(\mathbf{x}^{(1)}) = \dots = \mathcal{V}(\mathbf{x}^{(K)}) = \emptyset$, and suppose that $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is not feasibility-expanding. Then $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is elicibility-expanding for some α if and only if (1) if $\mathbf{x}^{(A)}$ is full-support i.e., $\mathcal{S}(\mathbf{x}^{(A)}) = [M]$, then there exists $j \in [M]$ such that so $\mathbf{x}^{(k)}$ has support $[M] \setminus \{j\}$ and (2) $\mathbf{x}^{(A)}$ is support-expanding relative to every $k \in [K]$ i.e., $\mathcal{S}(\mathbf{x}^{(A)}) \not\subseteq \mathcal{S}(\mathbf{x}^{(k)})$ for every $k \in [K]$.

It is harder to remove the strengthening for the binding set constraints. This is due to the joint geometry of the constraints that appears in the limitation-characterizing condition.

F PROOFS FOR SECTION 4

F.1 KEY LEMMAS FOR SECTION 4

The following lemmas provides the characterization for the elicibility of a vector \mathbf{x} under a feature weights matrix α in terms of the intersection of feasible perturbation directions $\mathcal{B}_{\mathcal{S}(\mathbf{x}), \mathcal{V}(\mathbf{x})} = \{\mathbf{d} : \mathcal{C}_{\mathcal{V}(\mathbf{x})} \mathbf{d} \leq 0\} \cap \{\mathbf{d} : \mathcal{d}_{\mathcal{S}(\mathbf{x})^c} \geq 0\} \cap \{1^t \mathbf{d} < 0\}$ and feature-improving directions $\mathbf{d} \in \mathbb{R}^M : \{\mathbf{d} : \alpha \mathbf{d} \geq 0\}$. These results generalize the characterization results in Kleinberg et al. (2019) to allow for conic constraints \mathcal{C} .

Lemma F.1. If a vector \mathbf{x} is elicitable with budget E , then $\|\mathbf{x}\|_1 = E$.

Proof. This is because, for any feasible vector \mathbf{x} with $\|\mathbf{x}\|_1 < E$, scaling \mathbf{x} to obtain $\mathbf{x}' = E\mathbf{x}/\|\mathbf{x}\|_1$ results in a feasible vector that has strictly larger reward for any reward function.

\mathbf{x}' clearly maintains nonnegativity constraints and bounded ℓ_1 norm constraint. Additionally since the only other constraints are conic, scaling the feasible \mathbf{x} non-negatively also maintains the additional conic constraint.

By the monotonicity of the reward functions we consider, for all reward functions, \mathbf{x}' has reward at least as high as \mathbf{x} .

By the strict monotonicity of our feature mapping functions and for the notion of monotonicity of reward functions we consider, \mathbf{x}' achieves a strictly higher reward than \mathbf{x} . \square

The following lemma shows that elicibility of a vector only depends on the direction of the vector and not of the norm. It allows us to study elicibility of the normalized vector using budget 1 i.e., ℓ_1 norm bound of one. Hence our elicibility characterizations will be expressed with budget 1.

Lemma F.2. *A vector x is elicitable with some budget E , under a reward function R if and only if $x/\|x\|_1$ is elicitable with budget 1 for the same reward function.*

Proof. If x is elicitable, it is elicitable with a budget of $\|x\|_1$ by Lemma F.1. It is elicitable if and only there is no feasible y with $\|y\|_1 \leq \|x\|_1$ with higher reward than x . If such a y exists, then $x/\|x\|_1$ is not elicitable with budget 1 since $y/\|x\|_1$ also has budget 1, is feasible and has higher reward than x . Similarly, if an improving y existed for $x/\|x\|_1$ under budget 1, then $y\|x\|_1$ is improving for x under budget $\|x\|_1$. \square

Lemma F.3 (Single output elicitation necessary). *An output vector x is elicitable only if $\mathcal{B}_{S(x), \mathcal{V}(x)} \cap \{\alpha d \geq 0\}$ is non-empty.*

Proof of Lemma F.3. Let $d \in \mathcal{B}_{S(x), \mathcal{V}(x)} \cap \{\alpha d \geq 0\}$. It suffices to construct a feasible output vector y that has strictly higher reward than x for every x with ℓ_1 norm equal to one and for every monotone reward function of the features. This is sufficient to prove the lemma since by lemma F.1, any elicitable vector has ℓ_1 norm equal to one.

This vector y we construct is $y = (x + \lambda d)/\|x + \lambda d\|_1$ where $\lambda > 0$ is chosen to be small enough so that $y \geq 0$.

First consider the vector $y' = x + \lambda d$ for an appropriate choice of $\lambda > 0$ that we will describe in a bit. First note that y' is feasible on all conic and non-negativity constraints that are binding at $x(x)$ due to d 's membership in $\{d : C_{\mathcal{V}(x)}d \leq 0\} \cap \{d : d_{S(x)^c} \geq 0\}$.

We can choose λ to be small enough so that y' continues to meet all non-binding constraints. That is choose $\lambda < \min_{j \in \mathcal{V}(x)^c, C_j d > 0} -C_j x / C_j d$ and $\min_{i \in (\mathcal{V}(x), d_i < 0} -x_i / d_i$. This establishes that we have a positive choice of λ making y' satisfy the nonnegativity and conic constraints. Additionally, we have that $\mathbf{1}^T y' = \|y'\|_1 = \|x\|_1 - \lambda \mathbf{1}^T d < \|x\|_1 = 1$. That is, y' satisfies the bounded ℓ_1 norm constraint in a non-binding manner. This shows that y' is feasible.

We also have that $\alpha^T y' = \alpha^T (x + d) \geq \alpha^T x$ since $\alpha^T d \geq 0$. Hence y' satisfies feasibility constraints and has at least as high values on all features. By the monotonicity of the reward functions we consider, for all reward functions, y' has reward at least as high as x . Lemma F.1 shows that scaling y' to have ℓ_1 norm equal to one results in strictly higher reward for all reward functions. Hence $y'/\|y'\|_1$ is feasible and has strictly higher reward than x for all monotone reward functions. \square

Lemma F.4 (Single output elicitation sufficient). *An output vector x is elicitable if $\mathcal{B}_{S(x), \mathcal{V}(x)} \cap \{\alpha d \geq 0\}$ is non-empty.*

Proof. Write $S := S(x)$ and $V := V(x)$.

Existence of multipliers. By positive scaling of directions, the assumption $B_{S,V} \cap D_\alpha = \emptyset$ is equivalent to infeasibility of the system :

$$C_V d \leq 0, \quad d_{S^c} \geq 0, \quad \alpha^T d \geq 0, \quad \mathbf{1}^T d < 0. \quad (3)$$

Let $I_{S^c} \in \mathbb{R}^{|S^c| \times M}$ be the coordinate-selector matrix whose rows are the vectors e_j^T for $j \in S^c$, so that $I_{S^c} d = d_{S^c}$.

By Motzkin's transposition theorem of the alternative, infeasibility of equation 3 implies the existence of multipliers (i.e., dual variables)

$$\gamma \in \mathbb{R}_{\geq 0}^{|V|}, \quad \lambda \in \mathbb{R}_{\geq 0}^{|S^c|}, \quad \nu \in \mathbb{R}_{\geq 0}^N, \quad \tau > 0$$

such that

$$C_V^T \gamma - I_{S^c}^T \lambda + \tau - \alpha^T \nu = 0 \quad (4)$$

holds. (The strict right-hand side $\mathbf{1}^T d < 0$ yields $\tau > 0$.)

Reward function construction. Define a reward function that is linear in the features

$$R(z) = \sum_{i=1}^N \beta_i z_i \quad \text{with} \quad \beta_i := \frac{\nu_i}{f'_i((\alpha^\top x)_i)} (> 0),$$

which is well-defined since each f_i is strictly increasing, hence $f'_i((\alpha^\top x)_i) > 0$. Let $r(u) := R(F(u)) = \sum_{i=1}^N \beta_i f_i((\alpha^\top u)_i)$. Because each f_i is concave and increasing, r is concave. Its gradient at x is

$$\nabla r(x) = \sum_{i=1}^N \beta_i f'_i((\alpha^\top x)_i) \alpha_{:,i} = \alpha \nu,$$

where $\alpha_{:,i}$ is the i -th column of α .

Elicitability. Consider the reward maximization program

$$\max_{u \in \mathbb{R}^M} r(u) \quad \text{s.t.} \quad Cu \leq 0, \quad u \geq 0, \quad \mathbf{1}^\top u \leq 1.$$

This is a concave program, and its Lagrangian is

$$\mathcal{L}(u, \lambda_0, \mu, \tilde{\gamma}) = r(u) + \lambda_0 (1 - \mathbf{1}^\top u) + \mu^\top u - \tilde{\gamma}^\top (Cu),$$

with multipliers $\lambda_0 \geq 0, \mu \geq 0, \tilde{\gamma} \geq 0$. Evaluate the KKT conditions at $u = x$ with the choice

$$\lambda_0 := \tau, \quad \mu_S := 0, \quad \mu_{S^c} := \lambda, \quad \tilde{\gamma}_V := \gamma, \quad \tilde{\gamma}_{V^c} := 0.$$

Primal feasibility holds by definition of S, V . Complementary slackness holds since $x_j = 0$ for $j \in S^c$ and $(Cx)_\ell = 0$ for $\ell \in V$, while $\mu_S = 0$. For stationarity,

$$\nabla r(x) - \lambda_0 \mathbf{1} + \mu - C^\top \tilde{\gamma} = \alpha g - \tau \mathbf{1} + I_{S^c}^\top \lambda - C_V^\top \gamma = 0$$

by equation 4. Finally, $\lambda_0 = \tau > 0$ certifies that the ℓ_1 -budget binds ($\mathbf{1}^\top x = 1$, consistent with Lemma C.2).

Since r is concave and the constraints are linear, the KKT conditions are sufficient; hence x maximizes r over the feasible region and is therefore elicitable. \square

F.2 PROOF OF THEOREM 4.1

Theorem 4.1 follows directly from the single-agent results in the previous subsection.

Proof of Theorem 4.1. We apply Lemma F.3 and Lemma F.4 to obtain necessary and sufficient conditions on when x is elicitable. We apply this to the outputs $x^{(1)}, \dots, x^{(K)}$ as well as $x^{(A)}$. \square

F.3 KEY INTERMEDIATE RESULTS FOR THE PROOF OF THEOREM 4.3 AND THEOREM 4.4

To prove Theorem 4.3 and Theorem 4.4, we will use an alternate but equivalent way of expressing the limitations-characterizing condition (Definition 4.2). This equivalent condition is defined below.

Definition F.5. Fix constraints C and aggregation operation $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$. We say that the alternate limitations-characterizing condition is satisfied for $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ if 1) $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ does not implement feasibility-expansion, and 2) there does not exist $d^{(A)} \in \mathcal{B}_{S(x^{(A)})}, \mathcal{V}(x^{(A)})$ such that:

$$\{u + \lambda d^{(A)} \mid u \in \mathbb{R}_{\geq 0}^M, \lambda \geq 0\} \cap \left(\bigcup_{i \in [k]} \mathcal{B}_{S(x^{(i)})}, \mathcal{V}(x^{(i)}) \right) = \emptyset.$$

The following proposition shows that the limitation-characterizing condition is equivalent to the new condition we defined above.

Proposition F.6. The conditions defined in Definition 4.2 and Definition F.5 are equivalent.

Proof. For ease of notation, let $V_i =: \mathcal{V}(x^{(i)})$ for $i \in [K]$ and let $V_A := \mathcal{V}(x^{(A)})$. It suffices to show that $\{u + \lambda d : u, \lambda \geq 0\}$ for a $d \in \mathcal{B}_{[M], \mathcal{V}_A^{(0)}}$ has empty intersection with $\mathcal{B}_{[M], \mathcal{V}_i^{(i)}}$, for each $i \in [K]$ if and only if for every $\gamma^{(i)} \in \mathbb{R}_{\geq 0}^{|V_i|}$, $\gamma_i^T C_{V_i} d - \|(\gamma_i^T C_{V_i})_-\|_\infty > 0$ or $I_{S_i^c} d < \mathbb{1}^d d$. Without loss of generality, it suffices to prove this for all $\gamma^{(i)} \in \mathbb{R}_{\geq 0}^{|V_i|}$ with bounded norm, say $\|\gamma^{(i)}\|_1 \leq 1$.

For any $d \in \mathcal{B}_{[M], \mathcal{V}_A^{(0)}}$, the intersection of $\{u + \lambda d : u, \lambda \geq 0\}$ and $\mathcal{B}_{[M], \mathcal{V}_i^{(i)}}$ is non-empty if and only if there exists a $u, \lambda \geq 0$ such that $C_{V_i}(u + \lambda d) \leq 0$ and $\mathbb{1}^t(u + \lambda d) < 0$.

$d \in \mathcal{B}_{[M], \mathcal{V}_A^{(0)}}$ means that $\mathbb{1}^t d < 0$, $C_{V_A} d \leq 0$, and $-I_{S_A^c} d \leq 0$. We can always normalize d so that $\mathbb{1}^t d = -1$. We can also scale the inequalities for non-empty intersection by dividing by λ . (Note that $\lambda \neq 0$, since $\mathbb{1}^t u \geq 0$.) Hence, we can equivalently write the condition for non-empty intersection as the existence of $v \geq 0$ such that $C_{V_i}(d + v) \leq 0$, $-I_{S_i^c}(d + v) \leq 0$ and $\mathbb{1}^t v < -\mathbb{1}^t d = 1$. These inequalities for the non-empty intersection condition hold if and only if all weighted sums (with non-negative weights) of the inequalities also hold true. That is, for every $\gamma^{(i)} \geq 0$, $\lambda^{(i)} \geq 0$, weight vectors, $\gamma^{(i)t} C_{V_i}(d + v) - \lambda^{(i)\top} I_{S_i^c}(d + v) \leq 0$ and $\mathbb{1}^t v < 1$.

A v satisfying $(\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})(d + v) \leq 0$ and $\mathbb{1}^t v < 0$ to simultaneously exists if and only if

$$\inf_{v \geq 0: \mathbb{1}^t v \leq 1} \sup_{\gamma^{(i)} \geq 0, \|\gamma^{(i)}\|_1 \leq 1} (\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})(d + v) \leq 0.$$

Since $(\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})(d + v)$ is an affine function in $\gamma^{(i)}$, $\lambda^{(i)}$ and v , and since the sets we optimize over $\{\gamma^{(i)} \geq 0, \|\gamma^{(i)}\|_1 \leq 1\}$ and $\{v \geq 0, \mathbb{1}^t v \leq 1\}$ are convex and compact, we can apply, we can apply minimax theorem to get

$$\begin{aligned} & \inf_{v \geq 0: \mathbb{1}^t v \leq 1} \sup_{\gamma^{(i)}, \lambda^{(i)} \geq 0, \|\gamma^{(i)}\|_1 \leq 1, \|\lambda^{(i)}\|_1 \leq 1} (\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})(d + v) \\ &= \sup_{\gamma^{(i)}, \lambda^{(i)} \geq 0, \|\gamma^{(i)}\|_1 \leq 1, \|\lambda^{(i)}\|_1 \leq 1} \inf_{v \geq 0: \mathbb{1}^t v \leq 1} (\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})(d + v). \end{aligned}$$

Note that for a given $\gamma^{(i)}, \lambda^{(i)}$, we can construct an optimal v as follows. If $\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c}$ has a negative coordinate, then v places a weight of 1 on the most negative coordinate of $\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c}$. Otherwise, then $v = 0$. Using this construction, we know that:

$$\inf_{v \geq 0: \mathbb{1}^t v \leq 1} (\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})(d + v) = (\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})d - \|(\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})_-\|_\infty.$$

Thus, the condition of non-empty intersection becomes the condition that $(\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})d - \|(\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})_-\|_\infty \leq 0$ for all $\gamma^{(i)}, \lambda^{(i)} \geq 0, \|\gamma\|_1, \|\lambda\|_1 \leq 1$.

Note that $\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c}$ subtracts λ_j from some coefficient of the j th row of $\gamma^{(i)t} C_{V_i}$. As a result, we can write $\|(\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})_-\|_\infty$ as $\|\gamma^{(i)t} C_{V_i}\|_\infty + \|\lambda^{(i)\top} I_{S_i^c}\|_\infty = \|\gamma^{(i)t} C_{V_i}\|_\infty + 1$.

So the condition $(\gamma^{(i)t} C_{V_i} - \lambda^{(i)\top} I_{S_i^c})d - (\|\gamma^{(i)t} C_{V_i}\|_\infty + 1) \leq 0$ is equivalent to the condition that $\gamma^{(i)t} C_{V_i} - \|\gamma^{(i)t} C_{V_i}\|_\infty \leq 0$ and $-\lambda^{(i)\top} d - 1 \leq 0$ (since both terms being ≤ 0 implies the sum is ≤ 0 and conversely, if the sum is not ≤ 0 , one must be > 0). This is exactly the condition in the limitation-characterizing condition.

□

F.4 PROOF OF THEOREM 4.3

Using this equivalence, we will show the necessity of the alternative condition to establish the necessity of the limitations-characterizing condition

Proof of Theorem 4.3. We will prove the contrapositive: If $x^{(1)}, \dots, x^{(K)} \rightarrow x^{(A)}$ is elicibility-expanding for some feature map α and for conic constraints C , then the limitations-characterizing condition (Definition 4.2) is violated.

One case is that $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is elicibility-expanding through feasibility-expansion. This automatically violates the limitations-characterizing condition.

The other case is that $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is not feasibility-expanding. Then $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ is feasible. We will show that if the limitations-characterizing condition

If the violation occurs through existence of $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}$ is not elicitable.

Suppose that $\mathbf{x}^{(A)}$ is not elicitable under a feature mapping α and constraints C . We will show that a violation of the limitations-characterizing condition implies that one of $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}$ is not elicitable, which contradicts $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ being elicibility-expanding.

Since $\mathbf{x}^{(A)}$ is not elicitable under a feature weights matrix α , by Lemma F.4, there is a $\mathbf{d}^{(A)} \in \mathcal{K}_{S_0, V_0}$ such that $\alpha \mathbf{d}^{(A)} \geq 0$. Since the limitation-characterizing condition is violated, the alternate limitation-characterizing condition is also violated (Proposition F.6). This means that there exists $\mathbf{x}^{(i)}$ with \mathcal{K}_{S_i, V_i} having non-empty intersection with $\{u + \lambda \mathbf{d}^{(A)}\}$.

It suffices to show that $\mathbf{x}^{(i)}$ is not elicitable under feature mapping α . To see this, let d_i denote an element of the intersection $\mathcal{K}_{S_i, V_i} \cap \{u + \lambda \mathbf{d}^{(A)}\}$. We can then write $d_i = u + \lambda \mathbf{d}^{(A)}$. Note that $\alpha d_i = \alpha u + \lambda \alpha \mathbf{d}^{(A)}$. We know that $\alpha u \geq 0$ since $u \geq 0$ and α has non-negative entries. Additionally, $\alpha \mathbf{d}^{(A)} \geq 0$ as shown above. Hence $\alpha d_i \geq 0$. By Lemma F.3, this means that x_i is not elicitable.

□

F.5 PROOF OF THEOREM 4.4

Proof of Theorem 4.4. Suppose the limitation-characterizing condition is satisfied. By Proposition F.6, this means that the alternate limitation-characterizing condition is satisfied. Then we know that we are in one of two cases.

Case 1: $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \rightarrow \mathbf{x}^{(A)}$ implements feasibility expansion. Consider a feature mapping with a single feature and all dimensions contribute equal weights of one to this feature. All output vectors with the same ℓ_1 norm result in the same reward for all reward functions, and thus all feasible outcomes are elicitable. That is any output vector is elicitable if and only if it is feasible. Under this construction, feasibility-expansion implies elicibility-expansion.

Case 2: there exists $\mathbf{d}^{(A)} \in \mathcal{K}_{S(\mathbf{x}^{(A)}), \mathcal{V}(\mathbf{x}^{(A)})}$ such that for all $u \geq 0, \lambda \geq 0, u + \lambda \mathbf{d}^{(A)} \notin \mathcal{K}_{S_i, V_i}$ for $i \neq 0$. We will construct a feature mapping α based on $\mathbf{d}^{(A)}$ such that the set of directions weakly increasing feature values i.e., the set $D_\alpha = \{d : \alpha d \geq 0\}$ is a subset of $\{u + \lambda \mathbf{d}^{(A)} : u \geq 0, \lambda \geq 0\}$. This implies that for all other outputs x_i , $D_\alpha \cap \mathcal{K}_{S(\mathbf{x}^{(i)}), \mathcal{V}(\mathbf{x}^{(i)})}$ is empty and hence $x^{(i)}$ is elicitable under α .

To complete this argument, we will explicitly construct such an α based on $\mathbf{d}^{(A)}$. Let $P_0 = \{i \in [m] : d_i^{(A)} > 0\}$ denote the positive coordinates of $\mathbf{d}^{(A)}$ and let $N_0 = \{i \in [m] : d_i^{(A)} \leq 0\}$ denote the negative or zero coordinates. We construct two sets of features:

- For every $p \in P_0$, there is a corresponding feature F_p whose row in A is the vector e_p which is the vector with 1 at coordinate p and zero everywhere else. That is, the action x_p has weight 1 on feature F_p and all other actions have zero weight.
- The next set of features are defined for every pair $p \in P_0, q \in N_0$. This feature $F_{p,q}$ has a corresponding row in A_0 that is the vector $d_p^{(A)} e_q - d_q^{(A)} e_p$. That is, the only actions with possible non-zero weights to $F_{p,q}$ are actions x_p, x_q . The weight from x_p is $|d_q^{(A)}|$ and the weight from x_q is $|d_p^{(A)}|$.

Now let us show that the set $D_\alpha = \{d : \alpha d \geq 0\}$ is a subset of $B_0 = \{u + \lambda \mathbf{d}^{(A)}\}$. Take any $d \in D_\alpha$.

For every $p \in P_0$, since d weakly improves value of F_p , it holds that $d_p \geq 0$. By ensuring that $\lambda \leq d_p / d_p^{(A)}$ for all $p \in P_0$, we can ensure that $d_p - \lambda d_p^{(A)} \geq 0$.

For every $p \in P_0, q \in N_0$, since d weakly improves value of $F_{p,q}$, it holds that $-d_p d_q^{(A)} + d_q d_p^{(A)} \geq 0$. In other words, $d_q \geq d_p d_q^{(A)} / d_p^{(A)}$.

We will show that it is possible to choose a $\lambda \geq 0$ such that $d - \lambda d^{(A)} \geq 0$, and hence d can be expressed as $u + \lambda d^{(A)}$ for $u \geq 0$. If there is a $p \in P_0$ with $d_p = 0$, then $d_q \geq 0$ while $d_q^{(A)} \leq 0$. So for all $\lambda > 0$, $d_q - \lambda d_q^{(A)} \geq 0$. Otherwise, we can choose λ less than $d_p / d_p^{(A)}$ and we get $d_q - \lambda d_q^{(A)} \geq 0$. \square

G EMPIRICAL SETUP FOR SECTION 2.4

Model output generation. The outputs are generated using gpt-4o-mini-2024-07-18 with the temperature set to 1.0. These are the five prompts that are used to produce model outputs:

1. “From a machine learning theory perspective, list 10 influential papers that have shaped our current understanding of large language models.”
2. “From the perspective of natural language processing and computational linguistics, list 10 key research papers that have been most influential in the development of modern large language models.”
3. “From a cognitive science and psycholinguistics standpoint, list 10 important papers that inform our understanding of how large language models represent, process, or acquire linguistic and conceptual structure.”
4. “From the standpoint of AI alignment and human–AI interaction, list 10 important papers that have shaped how large language models are aligned, instructed, or trained with feedback.”
5. “From a multi-agent and game-theoretic perspective, list 10 influential papers that contribute to the development or understanding of large language models”

These prompts produce five outputs X_1, \dots, X_5 , each a list of 10 papers tailored to its respective perspective. Next, we pass the concatenated outputs (X_1, \dots, X_5) to gpt-4o-mini-2024-07-18 by prompting the model with *aggregation instructions* followed by the concatenation of the 5 lists of papers, where each list is preceded by followed by “List of papers: [insert output number]”. The intersection-style and addition-style aggregation operations are performed using the following *aggregation instructions*.

- *Addition-style aggregation*: “Each of the following lists contains influential papers on large language models in specializaing in different areas: machine learning theory, natural language processing, computational linguistics, AI alignment, human–AI interaction, and multi-agent systems. Based on these lists, generate a new list of 10 papers that reflects the union of their themes and coverage. Your list should be freshly generated (not a literal set union), but it should include papers that plausibly come from any of the provided lists, covering as much of the combined topical space as possible.”
- *Intersection-style aggregation*: “Each of the following lists contains influential papers on large language models in specializaing in different areas: machine learning theory, natural language processing, computational linguistics, AI alignment, human–AI interaction, and multi-agent systems. Based on these lists, generate a new list of 10 papers that reflects their intersection. That is, papers belonging to many of these areas of specialization. Your list should be freshly generated (not a literal intersection), selecting papers that could plausibly appear in all of the lists. If the literal intersection is empty, still generate the best possible list of papers that are central, broadly relevant, and thematically compatible with all lists.”

These aggregation prompts produce outputs $X_{\text{addition}}, X_{\text{intersection}}$.

Output vector computation. We now describe in more detail how we compute the embeddings shown in Figure 1. We embed and visualize the set $\{X_1, \dots, X_5, X_{\text{addition}}, X_{\text{intersection}}\}$. We calculate the 768-dimensional embeddings using all-mpnet-base-v2 (Reimers & Gurevych, 2019), which is built into the sentence-transformers package in pytorch. To make these embeddings fit into our framework, we translate them to the nonnegative orthant by applying an additive shift $s \in \mathbb{R}_{\geq 0}^{768}$. To

do this, we compute the embeddings of the 805 gpt-4o-mini-2024-07-18 outputs from the helpful-base dataset in AlpacaEval (Li et al., 2023). The additive shift s is taken to be negative of the minimum coordinate along each dimension in this set of 805 embeddings. We translate all 5 outputs and the aggregated outputs by adding s . We compute the variance across the 5 translated outputs vectors along each of the 768 dimensions, and select the top 2 and top 3 dimensions according to variance. We also compute the ℓ_2 -distance between outputs, which is invariant to the additive shift.