POWER AND LIMITATIONS OF AGGREGATION IN COMPOUND AI SYSTEMS

Anonymous authors

000

001

002003004

010 011

012

013

014

015

016

018

019

021

025

026

028 029

031

033

034

037

038

040

041

042

043

044

046

047

048

049

050 051

052

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 uncovers three natural mechanisms—feasibility expansion, support expansion, and binding set contraction—through which aggregation provides benefit to the system designer. 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 fundamentally 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 (Sec-

tion 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 to 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 elicitability-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 elicitability on any problem instance (Theorem 3.7). However, these mechanisms are not sufficient to expand elicitability in general, although we show that each mechanism results in elicitability-expansion under stronger conditions.

To more completely capture the power and limitations of aggregation, we provide a more general characterization of elicitability-expansion (Section 4). We first characterize when an aggregation operation is elicitability-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). The conditions in Definition 4.2 test whether feasible directions under which an agent can change the aggregated output can sufficiently violate the binding constraints of individual outputs.

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). Approaches include resampling the same model or reasoning trace and selecting outputs via reward models (Christiano et al., 2017), self-consistency (Wang et al., 2023), or synthesis (Zhang et al., 2025); routing queries across different LLMs (Chen et al., 2024); adversarially combining models to expose safety risks (Jones et al., 2025); 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 elicitability 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 incorpo-

rates agent's limitations in the form of costs for actions. Particularly related are multitask settings 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 (Holmström & Milgrom, 1991; Slade, 1996; Bond & Gomes, 2009; Demougin et al., 2022).

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. Our focus also differs and is on characterizing the powers and limitations of aggregation. 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 \boldsymbol{x} quantifies the degree to which the output captures each characteristic. We note that some dimensions may capture undesirable characteristics (e.g., hallunications). The system designer seeks a specific output $\boldsymbol{x}^{(A)} \in \mathbb{R}^M_{>0}$, which we assume to be unit ℓ_1 -norm $\|\boldsymbol{x}^{(A)}\|_1 = 1$.

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 $\mathbf{c}^T x \leq 0$ 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 $C \in \mathbb{R}^{L \times M}$ denote the conic constraints themselves. Let $C_i \in \mathbb{R}^M$ denote the ith row of C for $i \in [L]$, and let $C_V \in \mathbb{R}^{|V| \times M}$ denote the set of rows corresponding to indices $V \subseteq [M]$. Given a budget level E > 0, we let $\mathcal{B}(E)$ denote the feasible set at budget level E, defined to be:

$$\mathcal{B}(E) \coloneqq \{ \boldsymbol{x} \in \mathbb{R}_{\geq 0} \mid \boldsymbol{C} \boldsymbol{x} \geq \boldsymbol{0}, \| \boldsymbol{x} \|_1 \leq E \}.$$

We denote by C_{\emptyset} the zero-vector, to capture how $\{d: C_{\emptyset} \leq 0\} = \mathbb{R}^{M}_{>0}$.

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 $F(x) = [F_1(x), \dots, F_N(x)]$ take the form

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

where $f_i(\cdot)$ is nonnegative, smooth, weakly concave (i.e., diminishing returns from increasing quality on this dimension), and strictly increasing, and where the values $\alpha_{ij} \ge 0$ are nonnegative *feature*

weights. We will denote by $\alpha \in \mathbb{R}^{M \times N}_{>0}$ the matrix with entries α_{ij} and call this the *feature weights matrix*.

We consider reward specifications $R^{(1)}, \ldots, R^{(K)} : \mathbb{R}^N \to \mathbb{R}$ 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,

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

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

We say that an output x is elicitable if there exists a monotone reward specification R and a positive budget level E that elicits x. The condition for whether x is elicitable only depends on x through the following sufficient statistic (S(x), V(x)). The first component $S(x) = \{j : x_j > 0\}$ denotes the support of x. The second component $V(x) = \{l \in [L] : C_l x = 0\}$ denotes the set of indices of conic constraints that are binding at 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 $x^{(1)} \dots, x^{(K)} \to x^{(A)}$ is an elicitability-expanding operation if

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

Intuitively, if an aggregation operation is elicitability-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.

3 NATURAL MECHANISMS FOR ELICITABILITY-EXPANSION

In this section, we formalize natural mechanisms by which aggregation expands elicitability. First, we show how mechanisms expand elicitability via examples (Section 3.1). Then, we show that these mechanisms are necessary for elicitability-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 elicitability 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 elicitable outputs.

In our examples, the aggregation operations in this section will be based on the following two aggregation rules. The first is *intersection aggregation*, which is defined to be the coordinate-wise minimum of the vectors:

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

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 $\boldsymbol{w} \in \mathbb{R}_{>0}^K$, the rule is given by

$$A_{\text{add}}(x^{(1)}, \dots, x^{(K)}; w) = \sum_{i=1}^{K} w_i x^{(i)}.$$
 (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). At the end of this subsection, we consider investigate the extent to which they 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 $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$ implements **feasibility expansion relative to** C if $x^{(A)}$ is infeasible i.e., $Cx^{(A)} \nleq 0$ but all $x^{(i)}$ for $i \in [K]$ are feasible i.e., $Cx^{(i)} \leq 0$.

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

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 $\mathbf{x}^{(1)} = [1,0,1], \mathbf{x}^{(2)} = [0,1,1] \rightarrow \mathcal{A}_{intersect}(\mathbf{x}^{(1)},\ldots,\mathbf{x}^{(K)})(\mathbf{x}^{(1)},\mathbf{x}^{(2)}) = [0,0,1]$ (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. 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 $x^{(1)}, \dots, x^{(K)} \to x^{(A)}$ implements support-expansion relative to i if $S(x^{(A)}) \notin S(x^{(i)})$.

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

¹Kleinberg et al. (2019) studied this in single-agent environments without constraints.

Example 3.4. Let the feature map be $\alpha = \alpha_{0.6}$. Suppose that there are no constraints $C = \emptyset$, so elicitability 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 $\mathbf{x}^{(1)} = [1,0,0], \mathbf{x}^{(2)} = [0,1,0] \rightarrow \mathcal{A}_{add}(\mathbf{x}^{(1)},\ldots,\mathbf{x}^{(K)};\mathbf{w})(\mathbf{x}^{(1)},\mathbf{x}^{(2)};[1/2,1/2]) = [1/2,1/2,0]$ (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 $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$ implements binding set contraction relative to i if $\mathcal{V}(x^{(A)}) \not\supseteq \mathcal{V}(x^{(i)})$ or $\mathcal{V}(x^{(A)}) = \mathcal{V}(x^{(i)})$.

Aggregation operations which implement binding set contraction can expand elicitability, 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 $x_1 + x_2 \le x_3$. 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 $\mathbf{x}^{(1)} = [1/2, 0, 1/2], \mathbf{x}^{(1)} = [0, 1/2, 1/2] \rightarrow \mathcal{A}_{intersect}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})(\mathbf{x}^{(1)}, \mathbf{x}^{(2)}) = [0, 0, 1/2].$

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 elicitability-expansion.

Necessity of these mechanisms. First, we show that if an aggregation operation expands elicitability for some feature weights matrix, it must implement at least one of the three mechanisms. Specifically, 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 $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$. If $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$ is elicitability-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 $i \in [K]$, $x^{(1)}, \dots, x^{(K)} \to x^{(A)}$ 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 elicitability-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.

Partial sufficiency of these mechanisms in concrete instances. 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 3.8. Fix conic constraints C. If an aggregation operation $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$ implements feasibility-expansion, then there exists a feature map α such that $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$ is elicitability-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 elicitability-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 elicitability 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.²

Proposition 3.9. Fix conic constraints C, and any $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \to \mathbf{x}^{(A)}$. 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, $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \to \mathbf{x}^{(A)}$ is elicitability-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 $\mathcal{V}(\boldsymbol{x}^{(i)}) \setminus \mathcal{V}(\boldsymbol{x}^{(A)})$ 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.

Summary. While these three natural mechanisms are necessary for aggregation to have power, these mechanisms do not fully characterize it. 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 $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$ is elicitability-expanding. We begin by analyzing, for a fixed feature weights matrix and feasibility constraints, whether a given aggregation operation expands elicitability (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 elicitability-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 elicitability-expansion.

²The 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 B.4).

4.1 CHARACTERIZING WHEN ELICITABILITY-EXPANSION SUCCEEDS

To analyze the power of aggregation, we characterize whether an aggregation operation is elicitability-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), V(x)), elicitability is determined by the structure of the set

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

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

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

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

•
$$\mathcal{B}_{\mathcal{S}(x^{(A)}),\mathcal{V}(x^{(A)})} \cap \{d \in \mathbb{R}^M \mid \alpha d \geq 0\} \neq \emptyset.$$

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 limitations) 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 main idea is that elicitability of a vector \boldsymbol{x} is determined by whether the set of feasible perturbation directions $\mathcal{B}_{\mathcal{S}(\boldsymbol{x}),\mathcal{V}(\boldsymbol{x})}$ intersects the set of feature-improving directions $\boldsymbol{d} \in \mathbb{R}^M : \alpha \boldsymbol{d} \geq \boldsymbol{0}$ (Lemmas E.3 and E.4). This lemma allows us to characterize when each output $\boldsymbol{x}^{(i)}$ (for $i \in [K]$) is elicitable, while the aggregate $\boldsymbol{x}^{(A)}$ is not —precisely the condition for elicitability expansion.

To prove the lemma, one direction is straightforward: if the intersection is nonempty, then there exists a feasible direction d that weakly improves every monotone reward over the features, having lower ℓ_1 norm. Scaling the vector obtained by moving in this direction to have ℓ_1 norm equal to one strictly improves any monotone reward function, providing a certificate that x cannot be elicitable.

The other direction of the characterization is more involved and shows that whenever the intersection is empty, there is a reward function that elicits x. Similar to ?, we show any such elicitable x can be elicited by a reward function that is 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 elicitablity-expanding for *any* feature map. This represents a particularly strong form of limitation, as it rules out elicitability-expansion for all forms of reward specification limitations. The characterizing condition is stated below.

Definition 4.2. [Limitation-characterizing condition] Fix constraints C and aggregation operation $\boldsymbol{x}^{(1)},\ldots,\boldsymbol{x}^{(K)}\to\boldsymbol{x}^{(A)}$. We say that the **limitation-characterizing condition** is satisfied for $\boldsymbol{x}^{(1)},\ldots,\boldsymbol{x}^{(K)}\to\boldsymbol{x}^{(A)}$ if and only if (1) $\boldsymbol{x}^{(1)},\ldots,\boldsymbol{x}^{(K)}\to\boldsymbol{x}^{(A)}$ does not implement feasibility-expansion for C and (2) there does not exist $\mathbf{d}\in\mathbb{R}^M$ satisfying both conditions:

- $d \in \mathcal{B}_{\mathcal{S}(x^{(A)}),\mathcal{V}(x^{(A)})}$
- For every $i \in [K]$, there exists $\gamma^{(i)} \in \mathbb{R}_{>0}^{|\mathcal{V}(x^{(i)})|}$ such that

$$(\gamma^{(i)})^T C_{\mathcal{V}(\boldsymbol{x}^{(i)}} \boldsymbol{d} - |1^t \boldsymbol{d}| \cdot \left| \min_{j \in [M]} (\min(0, ((\gamma^{(i)})^T C_{\mathcal{V}(\boldsymbol{x}^{(i)})})_j)) \right| > 0,$$

or there exists $j \in \mathcal{S}(x^{(i)})^c$ such that $-d_j - |\mathbb{1}^t \mathbf{d}| > 0$.

The condition is a combination of two sub-conditions. The first requires that the aggregation operation does not implement feasibility-expansion. The second sub-condition further requires one of two things to fail. The first is related to binding set contraction. In particular contraction via violation of a weighted sum of binding constraints and by a minimum margin that depends on the magnitude of the most negative coordinate in the weighted constraint. The second is related to failure of support expansion occurring by expanding in a dimension outside of the support again by a minumum amount.

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

Theorem 4.3 (Necessary). Fix constraints C. If the limitation-characterizing condition is satisfied, then there does not exist a feature weights matrix α under which $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \to \mathbf{x}^{(A)}$ is elicitability-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 elicitablity under *some* feature weights matrix.

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

Proof idea. We prove this constructively by designing a feature weights vector that makes aggregation elicitability-expanding. The idea is to find a direction $d^{(A)}$ feasible for $x^{(A)}$ but violating constraints for every $x^{(i)}$. We then define a feature map where improvement is possible only along $d^{(A)}$, which is infeasible for any $x^{(i)}$. The existence of such a direction relates to a versions support-expansion and binding-set contraction. An additional condition allows such a feature map to be constructed.

5 DISCUSSION

In this work, we 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 natural mechanisms—feasibility-expansion, support expansion, and binding-set contraction—in order to expand the set of elicitable outputs. Since these mechanisms are not sufficient to ensure that aggregation adds power, we also precisely characterize when aggregation expands the set of elicitable outputs in general problem instances.

Our results offer a theoretical insights into the power and limitations of aggregation in compound AI systems. First, our results illustrate how 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). These latter two mechanisms enable aggregation to add power even as model capabilities continue to improve. 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.

Limitations and Future Work. Our stylized model, which builds on a classical principal-agent framework (Kleinberg et al., 2019), makes simplifying assumptions for tractability. While rewards R can be nonlinear, we restrict output and reward-specification limitations to linear forms, leaving nonlinear extensions, requiring more complex optimization, open for future work. 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, future work could extend beyond reward design to other system-level choices, such as tool use and fine-tuning, that shape specialized models 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, 3.1
- 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
- BAIR Research Blog. Compound ai systems. https://bair.berkeley.edu/blog/2024/02/18/compound-ai-systems/, 2024. Blog post. 1, 1.1, 3.1, 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
- 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 principalagent 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, 3.1, 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
- 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 541 ownership, and job design. Journal of Law, Economics, & Organization, 7:24–52, 1991. 1.1 542 543 Athul Paul Jacob and Jacob Andreas. The consensus game: Language model generation via con-544 sensus seeking. In International Conference on Learning Representations (ICLR), 2024. URL https://openreview.net/forum?id=zEwWZk7c0Z. 1.1 546 547 Erik Jones, Anca Dragan, and Jacob Steinhardt. Adversaries can misuse combinations of safe mod-548 els. In International Conference on Machine Learning (ICML), 2025. 1.1 549 Akbir Khan, John Hughes, Dan Valentine, Laura Ruis, Kshitij Sachan, Ansh Radhakrishnan, Ed-550 ward Grefenstette, Samuel R. Bowman, Tim Rocktäschel, and Ethan Perez. Debating with 551 more persuasive llms leads to more truthful answers. In Proceedings of the 41st International 552 Conference on Machine Learning (ICML), volume 235 of Proceedings of Machine Learning 553 Research, pp. 23662-23733. PMLR, 2024. URL https://proceedings.mlr.press/ 554 v235/khan24a.html. 1 555 556 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 558 Computation (EC), pp. 825–844. ACM, 2019. doi: 10.1145/3328526.3329584. 1, 1.1, 2, 2.2, 2.3, 559 1, 4.1, 5, E.1 560 561 Krishna K. Ladha. The condorcet jury theorem, free speech, and correlated votes. American Journal 562 of Political Science, 36(3):617–634, 1992. doi: 10.2307/2111584. 1.1 563 Jean-Jacques Laffont and David Martimort. The Theory of Incentives: The Principal-Agent Model. 565 Princeton University Press, Princeton, NJ, 2002. 1.1 566 567 Edward P. Lazear and Sherwin Rosen. Rank-order tournaments as optimum labor contracts. *Journal* 568 of Political Economy, 89(5):841–864, 1981. 1.1 569 570 Nancy A. Lynch. Distributed Algorithms. Morgan Kaufmann, San Francisco, CA, 1996. ISBN 1558603484. 1.1 571 572 Margaret E. Slade. Multitask agency and contract choice. International Economic Review, 37(2): 573 465–486, 1996. 1.1 574 575 Ming Tan. Multi-agent reinforcement learning: Independent vs. cooperative agents. In Proceed-576 ings of the Tenth International Conference on Machine Learning (ICML), pp. 330–337. Morgan 577 Kaufmann, 1993. 1.1 578 579 Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdh-580 ery, and Denny Zhou. Self-consistency improves chain of thought reasoning in language models. 581 In International Conference on Learning Representations (ICLR), 2023. 1.1 582 583 Bohan Zhang, Xiaokang Zhang, Jing Zhang, Jifan Yu, Xiaokun Zhang, Tejas Srinivasan, Yifan Mai, 584 Juntao Li, Yoon Kim, and Minjoon Seo. Cot-based synthesizer: Enhancing llm performance 585 through answer synthesis, 2025. URL https://arxiv.org/abs/2501.01668. 1.1, 3.1

A LLM USAGE STATEMENT

Information Processing Systems, 33:15763–15773, 2020. 1.1

586

588

589 590

591 592

593

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.

Simon Zhuang and Dylan Hadfield-Menell. Consequences of misaligned ai. Advances in Neural

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 $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \to \mathbf{x}^{(A)} = \mathcal{A}_{intersect}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})$. 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 $A_{\text{intersect}}(x^{(1)}, \dots, x^{(K)})$ is always a subset of the support of each $x^{(i)}$.

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 elicitability 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 $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \to \mathbf{x}^{(A)} = \mathcal{A}_{add}(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)}; \mathbf{w})(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)})$ 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 elicitability-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 $\mathcal{A}_{add}(\mathbf{x}^{(1)},\ldots,\mathbf{x}^{(K)};\mathbf{w})(\mathbf{x}^{(1)},\mathbf{x}^{(2)};[1,1]) \rightarrow \mathbf{x}^{(A)} = (3,5,3)$, 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 \checkmark denotes that there exists a problem instance where the aggregation rule implements that mechanism. The symbol \times denotes that the aggregation rule does not implement the mechanism for any problem instance.

Proposition B.4. Fix conic constraints C, and any $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \to \mathbf{x}^{(A)}$. Suppose that $\mathcal{V}(x^{(A)}) = \mathcal{V}(x^{(1)}) = \dots = \mathcal{V}(x^{(K)}) = \emptyset$, and suppose that $x^{(1)}, \dots, x^{(K)} \to x^{(A)}$ 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, $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(K)} \to \mathbf{x}^{(A)}$ is not elicitability-expanding for any α .

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

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

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

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 $x_3 \le x_1 + x_2$ in Example 3.2, $x^{(1)} = [1, 0, 1], x^{(2)} = [0, 1, 1] \rightarrow x^{(A)} = [0, 0, 1]$ is elicitability-expanding.

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

For $x^{(i)}$, the support set $S(x^{(i)})$ is $\{i\}$. The constraint is binding on both $x^{(i)}$. The set $\mathcal{B}_{S(x^{(1)}),\mathcal{V}(x^{(1)})}$ is $\{d:d_3\leq d_1+d_2,d_2\geq 0,d_1+d_2+d_3<0\}$. For any d in this set, $d_3<0$ and $d_1+d_2<-d_3$. The set $\mathcal{B}_{S(x^{(2)}),\mathcal{V}(x^{(2)})}$ is $\{d:d_3\leq d_1+d_2,d_1\geq 0,d_1+d_2+d_3<0\}$. For any d in this set, $d_3<0$ and $d_1+d_2<-d_3$.

Now consider the set of feature-improving direction $\{d: d_1 + 2d_3 \ge 0, d_2 + 2d_3 \ge 0\}$. For any d in this set, $d_1 + d_2 \ge -4d_3$.

All three conditions $d_3 < 0$, $d_1 + d_2 < -d_3$, and $d_1 + d_2 \ge -4d_3$ cannot be satisfied since for $d_3 < 0$, $-d_3 > -4d_3$. Hence there is no intersection between feasibility improving directions and features improving directions and $x^{(1)}$ is elicitable. Similarly, $x^{(2)}$ is also elicitable.

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

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 elicitability-expanding.

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

 $d \in \mathcal{B}_{\mathcal{S}(\boldsymbol{x}^{(1)}),\mathcal{V}(\boldsymbol{x}^{(1)})}$ means that $d_1 < -(d_2 + d_3) < -d_3$ and $d_3 \ge 0$. $d \in \mathcal{A}_1$ means that $d_1 \ge -0.6d_3$. These three conditions cannot be simultaneously showing that $\boldsymbol{x}^{(1)}$ is elicitable due to empty intersection of \mathcal{A} and \mathcal{B}_1 . Symmetrically, we can also show that $\boldsymbol{x}^{(2)}$ is also elicitable.

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

702 703 704 C.3 Analysis of Example 3.6 705 706 **Proposition C.3.** For the feature weights matrix $\alpha_{0,2}$ and conic constraint matrix with one constraint $x_1 + x_2 \le x_3$ from Example 3.6, $\boldsymbol{x}^{(1)} = [1, 0, 1], \boldsymbol{x}^{(2)} = [0, 1, 1] \rightarrow \boldsymbol{x}^{(A)} = [0, 0, 1]$ is 707 708 elicitability-expanding. 709 710 *Proof of Proposition C.3.* The feature-improving directions are the set $A = \{d : d_1 + 0.2d_3 \ge 0, d_2 + 0.2d_3 \ge 0.2$ 711 $0.2\mathbf{d}_3 \ge 0$. 712 The constraint is binding at both $x^{(1)}$ and $x^{(2)}$. The feasibility improving directions are 713 $\mathcal{B}_{(S(x^{(1)}),\mathcal{V}(x^{(1)}))} = \{d: d_1 + d_2 \le d_3, d_2 \ge 0, d_1 + d_2 + d_3 < 0\}.$ 714 If $d \in \mathcal{B}_{(S(x^{(1)}),\mathcal{V}(x^{(1)})}$, then $d_1 + d_2 + d_3 < 0$ and $d_1 + d_2 + d_3 \le 2d_3$. This implies that $d_3 < 0$. If 715 $d \in \mathcal{A}$, then $d_1 \ge -0.2d_3$ and $d_2 \ge -0.2d_3$. If all the conditions are satisfied simultaneously, then 716 $d_1 > 0$ and $d_2 > 0$. This contradicts $d_1 + d_2 \le d_3 < 0$. 717 718 The conic constraint is not binding at $x^{(A)}$. Now consider the feasibility improving directions of 719 $d^{(A)}$: $\mathcal{B}_{(S(x^{(A)}),\mathcal{V}(x^{(A)})} = \{d: d_1 \ge 0, d_2 \ge 0, d_1 + d_2 + d_3 < 0\}$. The vector $d = (0.2, 0.2, -1) \in$ 720 $\mathcal{A} \cap \mathcal{B}_A$ demonstrating that $x^{(A)}$ is not elicitable. 721 722 Proofs in Section 3.2 723 724 D.1 PROOF OF THEOREM 3.7 725 726 *Proof of Theorem 3.7.* We show this as a corollary of Theorem 4.3. We will prove this by showing 727 that when both of the conditions in the theorem are violated, the limitation-characterizing condition 728 is satisfied and hence $x^{(1)}, \dots, x^{(K)} \to x^{(A)}$ cannot be elicitability-expanding. 729 One of the condition of the limitation-characterizing condition is the lack of feasibility-expansion 730 which is implied by the violation of the theorem's condition. We will show that the other condition 731 of the limitation-characterizing condition also holds. 732 733 When the second condition of the theorem is violated, there exists $i \in [K]$ with respect to which 734 $x^{(1)}, \dots, x^{(K)} \to x^{(A)}$ is neither support-expandin nor binding-set contracting. That is, there is an i such that $\mathcal{V}(x^{(i)}) \subseteq \mathcal{V}(x^{(A)})$ and $\mathcal{S}(x^{(i)}) \supseteq \mathcal{S}(x^{(A)})$. 735 736 For every $\boldsymbol{d} \in \{C_{\mathcal{V}(\boldsymbol{x}^{(A)})} \boldsymbol{d} \leq 0, \boldsymbol{d}_{\mathcal{S}(\boldsymbol{x}^{(A)})^c} \geq 0, 1^{\mathsf{T}} \boldsymbol{d} = -1\}, C_{\mathcal{V}(\boldsymbol{x}^{(i)})} (\boldsymbol{d}) \leq 0 \text{ and } \boldsymbol{d}_{\mathcal{S}(\boldsymbol{x}^{(i)})^c} \geq 0 \text{ , since } \boldsymbol{d} = -1\}$ 737 the rows of $C_{\mathcal{V}(\boldsymbol{x}^{(i)})}$ are a subset of the rows in $C_{\mathcal{V}(\boldsymbol{x}^{(A)})}$ and similarly, the rows in $d_{\mathcal{S}(\boldsymbol{x}^{(i)})^c} \geq 0$ are 738 a subset of the rows in $d_{\mathcal{S}(\boldsymbol{x}^{(A)})^c}$. Hence for any $\gamma^{(i)} \in \mathbb{R}^{|V_i|}_{>0}$, $(\gamma_i)^{\mathsf{T}} C_{V_i} d - \|(\gamma_i^{\mathsf{T}} C_{V_i})_-\|_{\infty} \leq 0$. 739 740 741 D.2 Proof of Proposition 3.8 742 *Proof of Proposition 3.8.* This follows from Theorem 4.4. 743 744 D.3 Proof of Proposition B.6 745 746 *Proof of Proposition B.6.* This follows from Proposition B.4. 747 748 D.4 Proof of Proposition B.5 749 750 *Proof of Proposition B.5.* Consider a problem with two output dimensions having the following two 751 constraints: 1) $c_1: x_1 - x_2 \le 0$, 2) $c_2: -2x_1 + x_2 \le 0$. Consider an aggregation operation $x^{(1)} = 0$ 752 $(1/2,1/2), x^{(2)} = (1/2,2/3) \rightarrow x^{(A)} = (5/12,7/12),$ where the binding constraints sets are $\mathcal{V}_{x^{(i)}} =$ 753 $\{c_i\}$ for $i \in \{1, 2\}$ and $\mathcal{V}_{x(A)} = \emptyset$. 754

In this example, we will show how the limitations-characterization condition holds, meaning that

the operation cannot be elicitability-expanding for any feature map α .

755

For any $\gamma_i \ge 0$ and d, $\gamma_i c_i d - \gamma_i \| (c_i)_- \|_{\infty} > 0$ if and only if $c_i d - \| (c_i)_{\infty} \| > 0$. 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 $\mathbb{1}^{d<-7}$. So the set of d such that $1 \top d = -1$ cannot intersect with the set of d satisfying the two conditions.

D.5 PROOF OF PROPOSITION B.4

Proof of Proposition B.4. Suppose that $M \ge 2$, and $S(x^{(A)}) = [M]$.

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 that there does not exist d such that for every $i \in [K]$ there exists $j(i) \in \mathcal{S}(\boldsymbol{x}^{(i)})^c$ such that $-d_{j(i)} - |1^{\mathsf{T}}d| > 0$.

Let's show the contrapositive: assume that there exists d such that for every $i \in [K]$ there exists $j(i) \in \mathcal{S}(\boldsymbol{x}^{(i)})^c$ such that $-d_{j(i)} - |1^{\mathsf{T}}\boldsymbol{d}| > 0$. Since $\boldsymbol{d} \in \mathcal{B}_{\mathcal{S}(\boldsymbol{x}^{(A)}),\mathcal{V}(\boldsymbol{x}^{(A)})}$, we know that $\boldsymbol{d}_{j'} \geq 0$ for $j' \notin \mathcal{S}(\boldsymbol{x}^{(A)})$ and we know that $1^{\mathsf{T}}\boldsymbol{d} < 0$. If $j \notin \mathcal{S}(\boldsymbol{x}^{(A)})$, then note that $-d_j < 0$, so this means that $j(i) \in \mathcal{S}(\boldsymbol{x}^{(A)})$. Putting this together, we see that $j(i) \in \mathcal{S}(\boldsymbol{x}^{(A)}) \setminus \mathcal{S}(\boldsymbol{x}^{(i)})$.

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 $0 < -d_{j(i)} - |1^{\mathsf{T}} \boldsymbol{d}| = 1^{\mathsf{T}} \boldsymbol{d} - d_{j(i)}$, if we add up all of these equations in the set $\{j(i) \mid i \in [K]\}$, we would obtain that $0 < M \cdot 1^{\mathsf{T}} \boldsymbol{d} - \sum_j d_j = (M-1) \cdot 1^{\mathsf{T}} \boldsymbol{d} - \sum_j d_j$, which means that $1^{\mathsf{T}} \boldsymbol{d} > 0$ which is a contradiction.

D.6 PROOF OF PROPOSITION 3.9

Proof of Proposition 3.9. We apply Theorem 4.4. It suffices to show that the limiting-characterization condition (Definition 4.2 is violated. Let \boldsymbol{d} be the vector such that $\boldsymbol{d}_{j(i)} = -1$ for all $i \in [K]$, $\boldsymbol{d}_j = |\{j(i) \mid i \in [K]\}| - 0.5$ for some $j \notin \{j(i) \mid i \in [K]\}$, and 0 elsewhere. It follows from definition that $\boldsymbol{d} \in \mathcal{B}_{\mathcal{S}(\boldsymbol{x}^{(A)}),\mathcal{V}(\boldsymbol{x}^{(A)})}$. It suffices to show that for all $i \in [K]$, it holds that:

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

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

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

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 d be the vector given by the condition in the theorem statement. It follows immediately that $d \in \mathcal{B}_{\mathcal{S}(\boldsymbol{x}^{(A)}),\mathcal{V}(\boldsymbol{x}^{(A)})}$. It suffices to show that for all $i \in [K]$, it holds that:

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

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

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

which we know holds.

E Proofs for Section 4

E.1 KEY LEMMAS FOR SECTION 4

817 v 818 (

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

Lemma E.1. If a vector x is elicitable with budget E, then $||x||_1 = E$.

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

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

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

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

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

Lemma E.2. A vector \mathbf{x} is elicitable with some budget E, under a reward function R if and only $\mathbf{x}/\|\mathbf{x}\|_1$ 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 E.1. It is elicitable if and only there is no feasible y with $\|y\|_1 \le \|x\|_1$ with higher reward tham 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 $\|x\|_1$.

Lemma E.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 E.3. Let $d \in \mathcal{B}_{\mathcal{S}(x),\mathcal{V}(x)} \cap \{\alpha d \geq 0\}$. 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 E.1, any elicitable vector has ℓ_1 norm equal to one.

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

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 $\{d : C_{\mathcal{V}(x)}d \le 0\} \cap \{d : d_{\mathcal{S}(x)^c} \ge 0\}$.

We can choose λ to be small enough so that y' continues to meet all non-binding constraints. That is choose $\lambda < \min_{j \in \mathcal{V}(\boldsymbol{x})^c, C_j d > 0} - \mathbf{C}_j \boldsymbol{x}/\mathbf{C}_j d$ and $\min_{i \in (\boldsymbol{x}), d_i < 0} - \boldsymbol{x}_i/d_i$. This establishes that we have a positive choice of λ making y' satisfy the nonnegativity and conic constraints. Additionally, we have that $\mathbb{1}^t y' = \|y'\|_1 = \|x\|_1 - \lambda \mathbb{1}^t d < \|x\|_1 = 1$. 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 $\alpha^t y' = \alpha^t (x+d) \ge \alpha^t x$ since $\alpha^t d \ge 0$. 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 E.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.

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

$$C_V d \le 0, \qquad d_{S^c} \ge 0, \qquad \alpha^{\mathsf{T}} d \ge 0, \qquad \mathbf{1}^{\mathsf{T}} d < 0.$$
 (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)

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

such that

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

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

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

$$R(z) = \sum_{i=1}^{N} \beta_i z_i \quad \text{with} \quad \beta_i := \frac{\nu_i}{f_i'((\alpha^{\mathsf{T}} x)_i)} \ (>0),$$

which is well-defined since each f_i is strictly increasing, hence $f_i'((\alpha^T x)_i) > 0$. Let $r(u) := R(F(u)) = \sum_{i=1}^N \beta_i f_i((\alpha^T u)_i)$. Because each f_i is concave and increasing, r is concave. Its gradient at x is

$$\nabla r(x) = \sum_{i=1}^{N} \beta_i f_i'((\alpha^{\mathsf{T}} x)_i) \alpha_{\cdot,i} = \alpha \nu,$$

where $\alpha_{\cdot,i}$ is the *i*-th column of α_{\cdot} .

Elicitability. Consider the reward maximization program

$$\max_{u \in \mathbb{R}^M} r(u) \quad \text{s.t.} \quad Cu \le 0, \ \ u \ge 0, \ \ \mathbf{1}^{\mathsf{T}} u \le 1.$$

This is a concave program, and its Lagrangian is

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

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

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

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,

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

by equation 4. Finally, $\lambda_0 = \tau > 0$ certifies that the ℓ_1 -budget binds ($\mathbf{1}^{\mathsf{T}}x = 1$, 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.

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

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

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

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

Proposition E.6. The conditions defined in Definition 4.2 and Definition E.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 $d \in \mathcal{B}_{[M],\mathcal{V}_{\boldsymbol{x}^{(0)}})}$ has empty intersection with $\mathcal{B}_{[M],\mathcal{V}_{\boldsymbol{x}^{(i)}})}$, for each $i \in [K]$ if and only if for every $\gamma^{(i)} \in \mathbb{R}^{|V_i|}_{\geq 0}$, $\gamma_i^T C_{V_i} d - \|(\gamma_i^T C_{V_i})_-\|_{\infty} > 0$ or $I_{S_A^c} d < \mathbb{1}^d d$. Without loss of generality, it suffices to prove this for all $\gamma^{(i)} \in \mathbb{R}^{|V_i|}_{\geq 0}$ with bounded norm, say $\|\gamma^{(i)}\|_1 \leq 1$.

For any $d \in \mathcal{B}_{[M],\mathcal{V}_{\boldsymbol{x}}(0)}$, the intersection of $\{u + \lambda d : u, \lambda \geq 0\}$ and $\mathcal{B}_{[M],\mathcal{V}_{\boldsymbol{x}}(i)}$ 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 $\mathbb{1}^t(u + \lambda d) < 0$.

 $d \in \mathcal{B}_{[M],\mathcal{V}x^{(0)}}$ means that $\mathbbm{1}^t d < 0$, $C_{V_A} d \le 0$, and $-I_{S_A^c} d \le 0$. We can always normalize d so that $\mathbbm{1}^t d = -1$. We can also scale the inequalities for non-empty intersection by dividing by λ . (Note that $\lambda \ne 0$, since $1^t u \ge 0$.) Hence, we can equivalently write the condition for non-empty intersection as the existence of $v \ge 0$ such that $C_{V_i}(d+v) \le 0$, $-I_{S_i^c}(d+v) \le 0$ and $\mathbbm{1}^t v < -\mathbbm{1}^d d = 1$. 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 $\gamma^{(i)} \ge 0$, $\lambda^{(i)} \ge 0$, weight vectors, $\gamma^{(i)t}C_{V_i}(d+v) - \lambda^{(i)\intercal}I_{S_i^c}(d+v) \le 0$ and $\mathbbm{1}^t v < 1$.

A v satisfying $(\gamma^{(i)t}C_{V_i} - \lambda^{(i)^{\intercal}}I_{S_i^c})(d+v) \leq 0$ 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 $(\gamma^{(i)t}C_{V_i} - -\lambda^{(i)^{\intercal}}I_{S_i^c})(d+v)$ 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

$$\begin{split} &\inf_{v \geq 0: \mathbb{1}^{v} < 1} \sup_{\gamma^{(i)}, \lambda^{(i)} \geq 0, \|\gamma^{(i)}\|_{1} \leq 1, \|\lambda^{(i)}\|_{1} \leq 1} (\gamma^{(i)t} C_{V_{i}} - \lambda^{(i)^{\mathsf{T}}} I_{S_{i}^{c}}) (d+v) \\ &= \sup_{\gamma^{(i)}, \lambda^{(i)} \geq 0, \|\gamma^{(i)}\|_{1} \leq 1, \|\lambda^{(i)}\|_{1} \leq 1} \inf_{v \geq 0: \mathbb{1}^{v} < 1} (\gamma^{(i)t} C_{V_{i}} - \lambda^{(i)^{\mathsf{T}}} I_{S_{i}^{c}}) (d+v). \end{split}$$

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

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

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

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

So the condition $(\gamma^{(i)t}C_{V_i} - \lambda^{(i)^{\intercal}}I_{S_i^c})d - (\|\gamma^{(i)t}C_{V_i-}\|_{\infty} + 1) \le 0$ is equivalent to the condition that $\gamma^{(i)t}C_{V_i} - \|\gamma^{(i)t}C_{V_i-}\|_{\infty} \le 0$ and $-\lambda^{(i)^{\intercal}}d - 1 \le 0$ (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.

E.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 $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$ is elicitability-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 $x^{(1)}, \dots, x^{(K)} \to x^{(A)}$ is elicitability-expanding through feasibility-expansion. This automatically violates the limitations-characterizing condition.

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

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

Suppose that $x^{(A)}$ 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 $x^{(1)}, \ldots, x^{(k)}$ is not elicitable, which contradicts $x^{(1)}, \ldots, x^{(K)} \to x^{(A)}$ being elictability-expanding.

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

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

E.5 Proof of Theorem 4.4

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

 Case 1: $x^{(1)}, \dots, x^{(K)} \to x^{(A)}$ 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 elicitability-expansion.

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

To complete this argument, we will explicitly construct such an α based on $d^{(A)}$. Let $P_0 = \{i \in [m] : d_i^{(A)} > 0\}$ denote the positive coordinates of $d^{(A)}$ and let $N_0 = \{i \in [m] : d_i^{(A)} \le 0\}$ 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 $D_{\alpha} = \{d : \alpha d \ge 0\}$ is a subset of $B_0 = \{u + \lambda d^{(A)}\}$. Take any $d \in D_{\alpha}$. For every $p \in P_0$, since d weakly improves value of F_p , it holds that $d_p \ge 0$. By ensuring that $\lambda \le d_p/d_p^{(A)}$ for all $p \in P_0$, we can ensure that $d_p - \lambda d_p^{(A)} \ge 0$.

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

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 $d_q - \lambda d_q^{(A)} \geq 0$.