# Mechanism Design with Level-k Types: Theory and an Application to Bilateral Trade

Terri Kneeland\*

June 4, 2020

#### Abstract

We develop necessary and sufficient conditions for level-k implementation that apply in independent private value environments. These conditions establish a set of level-k incentive constraints that are analogous to Bayesian incentive constraints. We show that in two special environments, the level-k incentive constraints collapse down to Bayesian incentive constraints. We then show, via a bilateral trade application, that this is not a general implication. Bilateral trade is ex post efficient under level-k implementation while it is not Bayesian implementable. We also address a robustness question concerning the common prior assumption embedded in level-k implementation by developing the concept of ex post level-k implementation. We develop necessary and sufficient conditions for ex post level-k implementation and show the relationship between ex post level-k and ex post implementation is analogous to the relationship between level-k and Bayesian implementation.

JEL codes: C72, D02, D90

<sup>\*</sup>Department of Economics, University College London. Address: Gower Street, London, UK, WC1E 6BT. E-mail: t.kneeland@ucl.ac.uk. I am grateful to Mike Peters for many invaluable discussions and guidance. For very helpful comments I would like to thank Vince Crawford, David Freeman, Yoram Halevy, Li Hao, Ryan Kendall, Kfir Eliaz, Wei Li, Rosemarie Nagel, Ariel Rubinstein, Sergei Severinov, Rani Spiegler, Alexander Wolitzky, and anonymous referees. The title of this paper is a play on Wolitzky's (2016) title "Mechanism Design with Maxmin Agents: Theory and an Application to Bilateral Trade".

### 1 Introduction

Laboratory experiments frequently find that behavior deviates from Nash and Bayesian equilibrium predictions when agents interact in novel environments. Non-equilibrium approaches, like level-k and cognitive hierarchy models, that relax the belief consistency assumptions of equilibrium models have been increasingly used to explain this behavior.<sup>1</sup> This empirical evidence prompts the need for extending the analysis of economic phenomena beyond an equilibrium analysis to other behaviorally plausible solution concepts.

This paper contributes to that end by analyzing mechanism design under the level-k solution concept. In the level-k model, agents anchor their beliefs in a naive model of others' behavior and adjust their beliefs by a finite number of iterated best responses. The model is anchored in the behavior of level 0 types, which is exogenously given (and is typically assumed to be uniformly random). Level 1 types engage in one level of reasoning and best respond to level 0 behavior. Level 2 types engage in two levels of reasoning and best respond to level 1 types. And so on, with level k types playing a best response to level (k - k)1) types. This yields a tractable model of strategic behavior in which agents determine their optimal actions in a finite number of steps. The level-k solution concept relaxes the belief consistency assumption of equilibrium by allowing agents to hold (possibly) inaccurate beliefs about the levels of reasoning of their opponents. Our notion of level-k implementability is identical to the notion of Bayesian implementability, except our solution concept is the levelk solution concept: a social choice rule is level-k implementable if for every profile of payoff types and levels, the actions played under the level-k model lead to outcomes that are consistent with the social choice rule.

Our main results establish general necessary and sufficient conditions for level-k implementation (Propositions 1 and 2). The level-k necessary conditions hold for general environments and the sufficient conditions hold in the case of independent private values. The conditions specify a set of level-k incentive constraints that are analogous to standard Bayesian incentive constraints. The level-k incentive constraints require there to exist a function, for

<sup>&</sup>lt;sup>1</sup>For pioneering work in the literature see Stahl & Wilson (1994; 1995), Nagel (1995), Costa-Gomes et al. (2001), and Camerer et al. (2004). For a recent survey of this literature, see Costa-Gomes et al. (2013).

each agent, that maps payoff types to outcomes in such a way that truthfully reporting payoff types is optimal for that agent given everyone else is truthfully reporting their payoff type. Level-k incentive constraints allow these functions to differ across agents while Bayesian incentive constraints require these functions to be the same for all agents. The level-k incentive constraints are thus a weak relaxation of Bayesian incentive constraints. If Bayesian incentive constraints hold, then the level-k incentive constraints also hold. However, it may be possible to ensure that the level-k incentive constraints hold without the Bayesian incentive constraints holding.

The ability to satisfy the level-k incentive constraints independent of the Bayesian incentive constraints holding depends on the environment. We establish two restrictions on the environment where level-k incentive constraints collapse to Bayesian incentive constraints. The first is when the social planner is implementing a social choice function (a single-valued rule). And, the second is when the message space is restricted to that of the set of payoff types. In both these cases, Bayesian incentive constraints are necessary conditions for level-k implementation (Corollary 1 and Proposition 3). The results for these special cases mirror existing results in the literature. de Clippel et al. (2019) study single-valued social choice rules and Crawford (2019) studies implementation under the restriction to mechanisms where the message space is the set of payoff types in the bilateral trade environment. Both find that Bayesian incentive constraints are necessary conditions for level-k implementation. Both of these papers are discussed in detail in the related literature section below.

In contrast to the results in these two restricted environments, our sufficient conditions allow for the possibility that level-k implementation is actually *more* permissive than Bayesian implementation. We show that in a bilateral trade environment, ex post efficient trade is level-k implementable (Proposition 6). This is in obvious contrast to Bayesian implementation where there is a conflict between ex post efficiency and incentive compatibility. Thus, with this example, we show that the existing results in the literature - that bound level-k implementability to what is Bayesian implementable - arise purely from restrictions on either the environment or the mechanism, and do not hold in general.

Lastly the paper explores a robustness question concerned with relaxing

the common prior assumption embedded in level-k implementation.<sup>2</sup> The definition of level-k implementation relies on the assumption of a common prior: agents' beliefs about others' levels are determined by the level-k model but agents' beliefs about the payoff types of others are determined by a common prior, as in Bayesian implementation. We develop the concept of expost levelk implementation which effectively allows for any beliefs over payoff types. We establish general necessary and sufficient conditions for expost level-k implementation (Propositions 4 and 5). As for level-k implementation, the necessary conditions hold for general environments and the sufficient conditions apply to environments of private values. The conditions specify a set of ex post levelk incentive constraints that are analogous to the standard expost incentive constraints. The relationship between expost level-k and expost implementation mirrors the relationship between level-k and Bayesian implementation. If the expost incentive constraints hold, then the expost level-k incentive constraints hold. But, it may be possible to ensure that the ex post level-k incentive constraints hold without the expost incentive constraints holding. We give an example and show that ex post efficient bilateral trade is ex post level-k implementable while it is not expost implementable.

### **Related literature**

There is a growing literature that focuses on behavioral mechanism design. This paper adds to this literature by studying implementation under the levelk model. Four other papers study level-k implementation. Crawford et al. (2009) looks at setting optimal reserve prices in first and second price auctions when agents are level-k types. Gorelkina (2015) provides a level-k analysis of the expected externality mechanism.

Crawford (2019) revisits Myerson & Satterthwaite's (1983) bilateral trade results under level-k implementation when the message set is restricted to the

<sup>&</sup>lt;sup>2</sup>Mechanisms that are robust to relaxing these strong common knowledge assumptions, typically known as the Wilson doctrine, can insure that a social choice rule will be implemented even if the planner does not know agents' beliefs about the payoffs of others. Much of this literature is due to Bergemann & Morris (2005), who investigate aspects of robust mechanism design (relaxing common knowledge of payoff assumptions) while maintaining the assumption of common knowledge of rationality. We investigate a version of robust implementation that relaxes common knowledge of payoffs under the empirically plausible assumption of level-k reasoning.

set of payoff types. Crawford considers two cases: (i) one where levels are unobservable and as such the social planner needs to screen both levels and payoff types (same environment as in this paper); and (ii) one where levels are observable, thus the social planner need only screen payoff types. In the first case, Crawford establishes a parallel result to Proposition 3 in this paper that shows that Bayesian incentive constraints are necessary for level-k implementation when the message set is restricted to the set of payoff types. And, hence shows the Myerson and Satterthwaite impossibility result for ex post efficient trade holds for level-k implementation when the message set is restricted. Crawford also explores what the 'second-best' level-k mechanisms look like in cases where full expost efficient trade cannot be achieved. In the latter case, Crawford shows that when levels are observable, a setting not explored in this paper, the relationship between level-k and Bayesian implementation is ambiguous and that the Myerson and Satterthwaite impossibility result can break down. Crawford gives a complete Myerson-Satterthwaite-style characterization of the optimal (restricted) mechanism in this case.

The current paper is closest to de Clippel et al. (2019). They establish a set of necessary and sufficient conditions for level-k implementation in a general setting where the social planner aims to implement a single-valued social choice rule. Their main finding is that Bayesian incentive constraints are necessary conditions for level-k implementation. In contrast, we establish a set of necessary and sufficient conditions for the case of a (possibly) multi-valued social choice rule and find that level-k implementation is actually a weaker implementation requirement than Bayesian implementation: a social choice rule may be level-k implementable even though it is not Bayesian implementable.

de Clippel et al., however, use a slightly stronger definition of level-k implementation than the one used here - requiring a version of full implementation<sup>3</sup> where this paper allows weak implementation. One might wonder then, if the main result in our paper, that level-k implementation is weaker than Bayesian implementation, arises from (1) relaxing the environment from single-valued social choice rules to multi-valued rules or (2) relaxing the implementation

<sup>&</sup>lt;sup>3</sup>Specifically, they require a condition they refer to as SIRBIC, which requires the level-k incentive constraints to hold with strict inequality whenever the social choice function is responsive.

requirement from full to weak implementation. We show that, (2), moving from full to weak implementation plays a minimal role under level-k implementation. First, Corollary 1 establishes that the necessary conditions for weak level-k implementation collapse to Bayesian incentive constraints when the social choice rule is single-valued (this replicates de Clippel et al.'s findings under weak implementation). Thus, demonstrating that moving from full to weak implementation when the social choice rule is single valued gains nothing beyond Bayesian implementability. Second, we show that the only difference between the necessary and sufficient conditions for full and weak level-k implementation of multi-valued social choice rules is that the level-k constraints must hold with a strict inequality rather than a weak inequality. This, in theory, means that full level-k implementation may be possible when Bayesian implementation is not. And, lastly, we show that this is true in practice: we provide an example where bilateral trade is ex post efficient under full level-k implementation but it is not Bayesian implementable.

The previous two papers place restrictions on the implementation problem: Crawford restricts the message space to be equal to the set of payoff types and de Clippel et al. restrict attention to single-valued social choice rules. A general takeaway from these papers (when levels are unobservable) is that Bayesian incentive constraints determine the boundaries of what is level-k implementable. However, the current paper shows that this conclusion arises from their restrictions on the implementation problem. If the social planner is interested in multi-valued choice rules and willing to use general message spaces, then level-k implementation can be strictly less restrictive than Bayesian implementation. The bilateral trade application is one such example. Instead, one can take the weaker level-k necessary and sufficient conditions established in this paper, as defining boundaries for what is level-k implementable in independent, private value environments.

Both de Clippel et al. (2019) and Crawford (2019) use concepts of level-k implementation, similar to this paper, that requires a common prior assumption: agents' beliefs about the payoff types of others is determined by a common prior. The additional analysis of ex post level-k implementation in this paper relaxes the common prior assumption. This analysis makes a unique contribution to the literature. There is also a literature on non-equilibrium design that does not employ the level-k model. Hagerty & Rogerson (1987), Bulow & Roberts (1989), Copic & Ponsati (2008; 2016), Mookherjee & Reichelstein (1992), Matsushima (2007; 2008), Bergemann & Morris (2005; 2009). Bergemann et al. (2011), de Clippel et al. (2015), Saran (2016), and Ollar & Penta (2017) all study implementation in dominant strategies, implementation in iterative dominance, implementation in rationalizable strategies, rationalizable implementation with an upper bound, or distribution-free implementation. Börgers & Li (forthcoming) study implementation in strategically simple mechanisms that only require agents to use first-order beliefs. Healy (2006) studies implementation in public good games when agents are learning to play equilibrium strategies.<sup>4</sup>

This rest of the paper proceeds as follows. Section 2 sets up the general payoff environment and formalizes level-k implementation. Section 3 establishes necessary and sufficient conditions for level-k implementation. Section 4 looks at two examples of special environments where the level-k incentive constraints collapse down to Bayesian incentive constraints. Section 5 addresses what happens when we relax the common prior assumption by looking at expost level-k implementation. Section 6 sets up the bilateral trade environment and shows that expost efficient trade is both level-k and expost level-k implementable. Section 7 concludes. Omitted proofs can be found in Appendix A.

### 2 Setup

### 2.1 General payoff environment

There is a finite set of agents I = 1, 2, ..., n. Agent *i*'s payoff type is  $\theta_i \in \Theta_i$ , where  $\Theta_i$  is a finite set. There is a compact set of outcomes Y. Each agent has

<sup>&</sup>lt;sup>4</sup>There is a literature that studies behavioral mechanism design that relies on equilibrium: Eliaz (2002) studies mechanism design when there is a proportion of 'faulty' agents that fail to act optimally. Glazer & Rubinstein (2012) allow the content and framing of the mechanism to play a role in behavior. de Clippel (2014) studies mechanism design when agents are not rational. Saran (2011) shows that ex post efficient trade can be achieved under bilateral trade when there is a proportion of truthful traders. Wolitzky (2016) investigates mechanism design and bilateral trade when agents are maxmin expected utility maximizers. Glazer & Rubinstein (1998), Eliaz & Spiegler (2006; 2007; 2008), Severinov & Deneckere (2006) study behavioral mechanism design in individual decision problems.

a continuous utility function  $u_i: Y \times \Theta \to \mathbb{R}$ . Note that we use the notation  $X = X_1 \times \cdots \times X_n$  and  $X_{-i} = X_1 \times \cdots \times X_{i-1} \times X_{i+1} \times \cdots \times X_N$  for sets  $\{X_i\}_{i \in I}$  throughout this paper.

There is a social planner who is concerned with implementing a (possibly multi-valued) social choice rule  $F : \Theta \to 2^Y \setminus \emptyset$ . The planner would like the outcome to be an element of  $F(\theta)$  whenever the true payoff type profile is  $\theta$ .

### 2.2 Type spaces

We use the framework of a type space in order to formally define agents' beliefs about the payoff types of others. The standard way to do this is to use a Bayesian type space. The set of payoff types along with a common prior over the set of payoff types constitutes a Bayesian type space.

**Definition 1.** A Bayesian type space  $\mathcal{B}$  is a structure  $\mathcal{B} = \langle \Theta; \rho \rangle$ , where  $\rho \in \triangle(\Theta)$ .

Given the common prior  $\rho$ , each payoff type forms her beliefs by conditioning on the common prior according to Bayes' rule. The belief of an agent with payoff type  $\theta_i$  about the payoff types of others is given by  $\rho(\theta_{-i}|\theta_i) = \frac{\rho(\theta)}{\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{i},\theta_{-i})}$ .

Similarly, we use a type space approach to define the level-k model. The level-k type space generates types that differ both by their payoff type and their level of reasoning.<sup>5</sup>

**Definition 2.** A  $\mathcal{B}$ -based level-k type space  $\mathcal{L}$  is a structure  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,...,n}, \bar{k} \rangle$ , such that  $\mathcal{B}$  is a Bayesian type space  $\mathcal{B} = \langle \Theta; \rho \rangle, T_i = \Theta_i \times \{0, 1, \ldots, \bar{k}\}$ , and  $\mu_i : T_i \to \Delta(T_{-i})$  such that

$$\mu_i(t_{-i} = (\theta_{-i}, k_{-i}) | t_i = (\theta_i, k_i)) = \begin{cases} \rho(\theta_{-i} | \theta_i) & \text{if } k_j = k_i - 1 \ \forall j \neq i \\ 0 & \text{otherwise} \end{cases}$$

In a level-k type space, an agent's type,  $t_i = (\theta_i, k_i) \in T_i$ , is 2-dimensional, representing both her payoff type,  $\theta_i$ , and her level,  $k_i$ . An agent's level represents her level of reasoning - an agent with a level k uses only k steps of

<sup>&</sup>lt;sup>5</sup>This approach follows the models of Crawford & Iriberri (2007a) and Crawford et al. (2009)which adapted the level-k models to incomplete information environments.

reasoning in order to figure out her optimal behavior in any game.<sup>6</sup> An agent's beliefs about the types of others are determined both by her payoff type and her level. The beliefs of type,  $t_i = (\theta_i, k_i)$ , about the types of others are determined by the function  $\mu_i(t_{-i}|t_i)$ . An agent with a level k puts weight only on types that have levels (k - 1). This captures the core assumption of the level-k model. An agent's beliefs about the payoff types of other agents are determined by the common prior  $\rho$ . Thus, an agent with payoff type  $\theta_i$  and level k believes that the payoff types of other agents are determined by  $\rho(\cdot|\theta_i)$  and that others have level k - 1.

We formally call this type space a Bayesian-based level-k type space because beliefs about payoff types are derived from a common prior. We drop this formalism throughout the rest of this paper and refer to these type spaces as simply level-k type spaces.

#### 2.3 Solution concepts

A mechanism specifies an action set for each agent and a mapping between action profiles and outcomes.

**Definition 3.** A mechanism  $\langle M, f \rangle$  consists of a set of actions  $M = M_1 \times \cdots \times M_n$  and a function  $f : M \to \triangle(Y)$ .

Given the payoff environment and (Bayesian or level-k) type space, a mechanism defines a *n*-agent incomplete information game with action set  $M_i$  and payoffs defined by  $u_i: Y \times \Theta \to \mathbb{R}$  and  $f: M \to \Delta(Y)$  for agent *i*.

For a given level-k type space, we can define the level-k solution concept. The level-k solution concept imposes that all types,  $k \ge 1$ , are rational (that is, they play a best response given their beliefs about the actions of other agents) and have consistent beliefs about the actions of other types. The behavior of level 0 types is specified outside of the model. Thus, level 0 types do not play a best response to their beliefs (and may play actions that are not a best response to any belief). To specify level 0 behavior, we define the notion of an anchor below.

<sup>&</sup>lt;sup>6</sup>The bound on the level of reasoning is not necessary, the results in this paper go through if  $T_i = \Theta_i \times \{0, 1, 2, \ldots\}$ , however bounding the depths of reasoning maintains the finiteness of the type space for simplicity.

**Definition 4.** For a given game defined by a mechanism  $\langle M, f \rangle$  and a ( $\mathcal{B}$ -based) level-k type space  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,...,n}; \bar{k} \rangle$ , an **anchor**  $\alpha = \alpha_1 \times \cdots \times \alpha_n$  is a mapping  $\alpha_i : \Theta_i \times \{0\} \to \Delta(M_i)$  for all  $i \in I$ .

An anchor is a essentially a strategy for each level 0 type. Notice that anchors can be arbitrary and vary from agent to agent and payoff type to payoff type. A commonly applied assumption in the level-k literature is of a *uniformly* random anchor, which would restrict  $\alpha_i$  to be the uniform random probability distribution over  $M_i$  for all payoff types and agents. The results that follow will be proved either for an arbitrary anchor or under the assumption of an atomless anchor. If  $M_i$  contains a continuum of messages for each i, then the anchor  $\alpha$  is an atomless anchor if the distribution  $\alpha_i(t_i)$  of messages contains no atoms, for each  $t_i$  and each i. The uniformly random anchor is an example of an atomless anchor.

We now define a level-k solution under a given anchor  $\alpha$ .

**Definition 5.** For a given game defined by a mechanism  $\langle M, f \rangle$  and levelk type space  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,...,n}; \bar{k} \rangle$ , and an anchor  $\alpha$ , a strategy profile  $s = s_1 \times \cdots \times s_n$ , with  $s_i : T_i \to \Delta(M_i)$  for all  $i \in I$ , is a **level-k solution under anchor**  $\alpha$  if and only if:

(i) 
$$s_i(t_i) \sim \alpha_i(t_i)$$
 for all  $t_i \in \Theta_i \times \{0\}$ 

(ii) 
$$\sum_{\substack{\theta_{-i}\in\Theta_{-i}}} \rho(\theta_{-i}|\theta_i) u_i(f(s_i(\theta_i,k), s_{-i}(\theta_{-i},k-1)), \theta)$$
$$\geq \sum_{\substack{\theta_{-i}\in\Theta_{-i}}} \rho(\theta_{-i}|\theta_i) u_i(f(m'_i, s_{-i}((\theta_{-i},k-1)), \theta) \ \forall m'_i \in M_i, \ \theta_i \in \Theta_i,$$
$$k \in \{1, \dots, \bar{k}\}, \ i \in I$$

The level-k solution can be calculated recursively given the behavior of level 0 types (condition (i)). Level 1's actions are a best response to level 0's actions (condition (ii)). Level 2's actions are a best response to level 1's actions, and so on (condition (ii)). We have abused notation slightly in the above definition, as it may be the case that  $f \circ s$  is a compound lottery. In this instance,  $u_i(f(s(\cdot)))$  in condition (ii) should be taken to be the corresponding expected utility representation.

#### 2.4 Implementation

A social choice rule is level-k implementable if there exists a mechanism and a level-k solution that achieves the social planners objective for every message sent. The formal definition is given below.

**Definition 6.** A social choice rule, F, is **level-k implementable under** anchor  $\boldsymbol{\alpha}$  on  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,...,n}; \bar{k} \rangle$  if there exists a mechanism  $\langle M, f \rangle$  and a message profile  $m = m_1 \times \cdots \times m_n$  that is a level-k solution under anchor  $\boldsymbol{\alpha}$ and m achieves F:  $f(m(\theta \times \hat{k})) \in F(\theta)$  for all  $\theta \times \hat{k} \in \times \{\Theta_i \times \{1, \ldots, \bar{k}\}\}_{i \in I}$ .

First, notice that our notion of level-k implementability does not require the mechanism to satisfy the social choice rule for level 0 types. This is because level 0 agents are non-strategic, hence the social planner cannot incentivize their behavior.<sup>7</sup> There is also some empirical support that the proportion of level 0 agents is small (e.g. Arad & Rubinstein 2012; Costa-Gomes et al. 2001; Costa-Gomes & Crawford 2006; Brocas et al. 2014). Thus we interpret the existence of level 0 types in the model, but not in our implementability requirement, as types that exist only in the minds of the other types.

Second, notice that our notion of level-k implementation does not require the social planner to have knowledge of the *actual* distribution of types (and hence levels). This is because implementation requires that the outcome be consistent with the social choice rule for *all* type profiles and hence does not depend upon the distribution of types.<sup>8</sup>

We have abused notation slightly in the above definition as it may be the case that  $f \circ m$  is a lottery. In this case, take m to achieve F as occurring if for any outcome, y, in the support of  $f(m(\theta \times \hat{k}))$  then it must be that  $y \in F(\theta)$ .

We will be interested in comparing level-k implementation with that of Bayesian implementation. We define Bayesian implementation for completeness. In the below definition, we incorporate the results of the revelation

<sup>&</sup>lt;sup>7</sup>If this type of behavior is a concern we should consider an alternative form of implementability. See Eliaz (2002) for one such possibility - the social planner tries to minimize the deviations from the social choice rule.

<sup>&</sup>lt;sup>8</sup>This is not true for all mechanism design objectives. For example, it would not be true if the goal of the planner was to maximize expected revenue. If different levels (and payoff types) play different actions with different revenue consequences, then expected revenue will depend upon the actual distribution of both payoffs and levels.

principle, and hence define Bayesian implementation under a direct mechanism. Condition (ii) below states the standard Bayesian incentive constraints which will be contrasted with the level-k incentive constraints developed in the next section. Notice that when we compare level-k and Bayesian implementation throughout this paper we will be comparing Bayesian implementation given some Bayesian type space  $\mathcal{B} = \langle \Theta; \rho \rangle$  and the related Bayesian-based level-k type space  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,...,n}; \bar{k} \rangle$ . In other words, we will compare Bayesian and level-k implementation given a shared common prior  $\rho$ .

**Definition 7.** A social choice rule F is **Bayesian implementable** on  $\mathcal{B} = \langle \Theta; \rho \rangle$  if there exists a mechanism  $\langle \Theta, f \rangle$  such that the following conditions hold:

- (i)  $f(\theta) \in F(\theta) \ \forall \theta \in \Theta$
- (ii)  $\sum_{\substack{\theta_{-i}\in\Theta_{-i}\\\theta\in\Theta},\ i\in I} \rho(\theta_{-i}|\theta_i)u_i(f(\theta),\theta) \geq \sum_{\substack{\theta_{-i}\in\Theta_{-i}\\\theta_{-i}\in\Theta_{-i}}} \rho(\theta_{-i}|\theta_i)u_i(f(\theta'_i,\theta_{-i}),\theta) \ \forall \theta'_i \in \Theta_i,$

## 3 Necessary and sufficient conditions for levelk implementation

This section establishes necessary and sufficient conditions for level-k implementation. We provide the necessary and sufficient conditions separately as the necessary conditions hold for an arbitrary anchor and in general environments, while the sufficient conditions hold only for atomless anchors in independent private value environments.

**Proposition 1.** (Necessary Conditions) Let F be a social choice rule. Let  $\mathcal{B}$  be a Bayesian type space and let  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,\dots,n}; \bar{k} \rangle$  be a ( $\mathcal{B}$ -based) level-k type space with  $\bar{k} \geq 2$ . If F is level-k implementable under anchor  $\alpha$ , then there exists a function  $f^i : \Theta \to \Delta(Y)$  for each  $i \in I$  and a function  $\bar{f} : \Theta \to \Delta(Y)$ , such that the following conditions hold:

(i)  $f^{i}(\theta)$ ,  $\bar{f}(\theta) \in F(\theta) \ \forall \theta \in \Theta, i \in I$ 

(ii) 
$$\sum_{\substack{\theta_{-i}\in\Theta_{-i}\\\theta\in\Theta,i\in I}} \rho(\theta_{-i}|\theta_i)u_i(f^i(\theta),\theta) \ge \sum_{\substack{\theta_{-i}\in\Theta_{-i}\\\theta_i\in\Theta_{-i}}} \rho(\theta_{-i}|\theta_i)u_i(f^i(\theta'_i,\theta_{-i}),\theta) \ \forall \theta'_i\in\Theta_i,$$

(iii) 
$$\sum_{\substack{\theta_{-i}\in\Theta_{-i}\\\theta\in\Theta,i\in I}} \rho(\theta_{-i}|\theta_i) u_i(f^i(\theta),\theta) \ge \sum_{\substack{\theta_{-i}\in\Theta_{-i}\\\theta_{-i}\in\Theta_{-i}}} \rho(\theta_{-i}|\theta_i) u_i(\bar{f}(\theta'_i,\theta_{-i}),\theta) \ \forall \theta'_i\in\Theta_i,$$

The formal proof can be found in Appendix A. We will give the intuition here. If the social choice rule F is level-k implementable, then there exists some mechanism  $\langle M, g \rangle$  and level-k solution  $m = m_1 \times \cdots \times m_N$  that level-k achieves F.

Consider the behavior of a level 2 agent *i*, specifically an agent with type  $t_i = (\theta_i, 2)$ . This agent sends message  $m_i(\theta_i, 2)$  and believes that everyone else is of level 1. Define the notation  $(\theta, k) = ((\theta_1, k), \ldots, (\theta_n, k))$  and  $(\theta_{-i}, k) = ((\theta_1, k), \ldots, (\theta_{i-1}, k), (\theta_{i+1}, k), \ldots, (\theta_n, k))$ . Then the agent believes others send the message profile  $m_{-i}(\theta_{-i}, 1)$ . This agent could send some other message  $m_i(\theta'_i, 2)$  for any  $\theta'_i \in \Theta_i$ . But, because  $\langle M, g \rangle$  and m level-k implement the social choice rule, it must be that she prefers to send  $m_i(\theta_i, 2)$  given her beliefs. Define  $f^i(\theta) = g(m_i(\theta_i, 2), m_{-i}(\theta_{-i}, 1))$  for all  $\theta \in \Theta$ . Thus, it must be that condition (ii) holds for agent *i*.

Alternatively, this agent could send messages of the form  $m_i(\theta'_i, 1)$  for any  $\theta'_i \in \Theta_i$ . But again, because  $\langle M, g \rangle$  and m level-k implements the social choice rule, it must be that she prefers to send  $m_i(\theta_i, 2)$ . Define  $\bar{f}(\theta) = g(m(\theta, 1))$  for all  $\theta \in \Theta$ . Thus, it must be that condition (iii) holds for agent i. The same argument extends to all agents. Further, it must be the case that  $g(m_i(\theta_i, 2), m_{-i}(\theta_{-i}, 1)), g(m(\theta, 1)) \in F(\theta)$  for all  $\theta \in \Theta$  since g achieves F. Thus, condition (i) must hold. Therefore, it is possible to find functions  $\{f^i\}_{i\in I}$  and  $\bar{f}$  such that conditions (i)-(iii) hold.

The necessary conditions are generated from the incentive requirements of a level 2 agent. The incentive requirements for higher levels mimic those for the level 2 type. However, the incentive requirements for the level 1 type may be quite different. This is because level 1 types believe all others are level 0 types, and the behavior of the level 0 types is exogenously given. But, it turns out that in independent private value environments it is possible to satisfy the level 1 incentive constraints without any additional conditions beyond (i)-(iii) under atomless anchors. Thus, within independent private value environments the conditions in Proposition 1 are both necessary and sufficient for level-k implementation under atomless anchors.

We first give the intuition for why the necessary conditions are sufficient. We will use the functions  $\{f^i\}_{i\in I}$  and  $\overline{f}$  to construct a mechanism that will level-k implement the social choice rule. Suppose that agents could send messages about both their payoff types and their levels i.e. the message space for agent i is  $\Theta_i \times \{0, \ldots, \bar{k}\}$ . Consider agent i and suppose that all other agents are truthfully reporting levels and payoff types (putting aside level 1 agents for now). Let the function  $f^i$  determine the outcomes that agent i believes will be implemented when agent i truthfully reports her level (i.e. the outcomes that will occur when agent i reports level k and everyone else reports levels k-1). And, let the function  $\overline{f}$  determine the outcomes that agent i believes will be implemented when agent i reports any other level (e.g. when a level 2 agent reports level 1 when all other agents report levels 1). If an agent believes that all other types are truthfully reporting levels and payoff types, then under condition (ii) and (iii) agents will want to truthfully report their payoff types and levels as well. Thus, the problem for the social planner is really how to incentivize level 1 agents to truthfully report their payoff types and levels for some exogenously given behavior (anchor) of level 0 agents.

To resolve this issue, we will extend the mechanism developed in de Clippel et al. (2019)<sup>9</sup> (proof of Proposition 3 in their paper) by allowing agents to report both payoff types and levels. This mechanism satisfies the level 1 incentive constraints by effectively manipulating the beliefs of level 1 agents to mimic that of truthful reporting of payoff types and levels. To understand how the mechanism works, consider the case where level 0 behavior is uniformly random and consider the following. Suppose the planner augmented the message space with an additional set,  $Z_i$ , for each agent, i.e. the message set for agent *i* would then be  $\Theta_i \times \{0, \ldots, \bar{k}\} \times Z_i$ . And, suppose the designer could use the message sent from  $Z_i$  to screen the agents into those that were level 0 and those that had higher levels, i.e. define the set  $Z_i^+$  to be such that any message sent in  $Z_i^+$  would indicate the agent had a level of at least 1 and any other message would suggest a level of 0. The planner could then do the following: if he received a message in  $Z_i^+$  he would take the type reports

<sup>&</sup>lt;sup>9</sup>Earlier drafts of this paper only provided the sufficient conditions in bilateral trade environments. After de Clippel et al. (2019) was written in 2016, we were able to adapt their proof technique to show our necessary conditions were also sufficient in general independent private value environments.

(payoff types and levels) at face value, but if he received any other message, he would modify the reports by randomly choosing a payoff type according to the common prior distribution  $\rho_i$  and set a level of 0. The planner would then assign outcomes under these (potentially modified) type reports via the  $\{f^i\}_{i\in I}$  and  $\bar{f}$  functions as described above. If it was also the case that  $Z_i$  is an uncountable set and  $Z_i^+$  is countable, then a level 1 agent who believed a level 0 agent j is uniformly randomly sending messages, would believe that the planner would almost surely ignore the type report of agent j and use the modified report where the payoff type is randomly chosen from  $\rho_j$  and set a level 0. This means that a level 1 agent would have the same beliefs about the outcomes that occur when she believes level 0 behavior is uniformly random as she would if she believed level 0 agents were truthfully reporting their level and payoff type. As such, level 1 agents would have an incentive to truthfully report their payoff type and level.

Given this intuition it is easy to see why such a proof only works for the case of independent, private value environments and atomless anchors. First, the social planner needs to be able to mimic the beliefs of the agents under common prior by drawing the payoff types from a random distribution. He can do this only in the case of independent values when he draws a payoff type for agent j randomly from  $\rho_j$ . Second, agents should only care about the modified payoff types used by the social planner and not the actual payoffs types of the other agents. In other words, agents need to only care about the payoff types of others to the extent that it tells them about the messages sent and not because it impacts utility directly. Thus, the environment needs to be one of private values. Third, agent i needs to believe that the modified payoff types are drawn from  $\rho_{-i}$  and levels are set to 0 almost surely. This will happen in the case of atomless anchors, but cannot be guaranteed with this mechanism otherwise.

Proposition 2 gives the formal result. Note that the proposition is stated only for the case when  $n \ge 3$ . The case when n = 2 is discussed in Remark 1 below. Define an environment of private values to be one where utility functions are such that  $u_i : Y \times \Theta_i \to \mathbb{R}$  for all  $i \in I$ . And, define an environment of independent values to be one where the Bayesian type space,  $\mathcal{B} = \langle \Theta; \rho \rangle$ , is such that  $\rho = \prod_i \rho_i$  for some  $\rho_1 \times \cdots \times \rho_n \in \Delta(\Theta_1) \times \cdots \times \Delta(\Theta_n)$ . **Proposition 2.** (Sufficient Conditions) Let F be a social choice rule, the environment be one of independent private values, and  $n \geq 3$ . Let  $\mathcal{B}$  be a Bayesian type space and let  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,...,n}; \bar{k} \rangle$  be a ( $\mathcal{B}$ -based) level-k type space. If there exists a function  $f^i : \Theta \to \Delta(Y)$  for each  $i \in I$  and a function  $\bar{f} : \Theta \to \Delta(Y)$  such that the conditions (i)-(iii) hold in Proposition 1 then F is level-k implementable under an atomless anchor  $\alpha$ .

PROOF<sup>10</sup>:

Consider the following mechanism where the message space for agent i is equal to  $M_i = \Theta_i \times \{0, 1, \dots, \bar{k}\} \times [-1, 1]$  and consists of a report of her payoff type  $\theta_i \in \Theta_i$ , level  $k_i \in \{0, 1, \dots, \bar{k}\}$ , and a real number,  $z_i \in [-1, 1]$ . Let the indicator function  $I_i : [-1, 1]^n \to \{0, 1\}$  be defined as follows:

$$I_i(z) = \begin{cases} 1 & \text{if } z_i = m_j z_j \text{ for some } m_j \in \mathbb{Z}, \forall j \in I \\ 0 & \text{otherwise} \end{cases}$$

Define  $\tilde{\theta}_i : M \to \Theta_i$  and  $\tilde{k}_i : M \to \{0, 1, \dots, \bar{k}\}$  in the following way. For a given message profile  $m = (\theta, k, z)$ , if  $I_i(z) = 1$  the planner takes the reports as given and sets  $\tilde{\theta}_i(m) = \theta_i$  and  $\tilde{k}_i(m) = k_i$ ; otherwise the planner sets  $\tilde{\theta}_i(m)$  to some randomly chosen  $\Theta_i$  according to the prior  $\rho_i$ and  $\tilde{k}_i(m) = 0$ . The planner then assigns outcomes based on the reports  $\tilde{\theta} \times \tilde{k}$  according to the function  $g : \Theta \times \{0, \dots, \bar{k}\}^n \to \Delta(Y)$  defined by

$$g(\theta \times \hat{k}) = \begin{cases} f^i(\theta) & \text{if } \hat{k}_j = \hat{k}_i - 1 \text{ for all } j \neq i \in I \\ \overline{f}(\theta) & \text{otherwise} \end{cases}$$

Consider an agent *i* with payoff  $\theta_i$  and level 1. She believes that all  $j \neq i$  are level 0 agents playing an atomless strategy  $\alpha_j$ . Thus, the probability

<sup>&</sup>lt;sup>10</sup>While we have maintained finiteness of the type space for simplicity, this proof relies on a message space that contains a continuum. A complete proof requires showing that such a mechanism and strategies are measurable. We direct the reader to de Clippel et al. (2019) for such a proof.

that the realized value of  $z_j$  is equal to  $mz_i$  is zero for any  $z_i \in [-1, 1]$  (for example, if  $z_i = 1/2$  then the set of  $z_j$  for which there exists some  $m \in \mathbb{Z}$ such that  $z_j = \frac{1}{2}m$  is finite:  $\{-1, -\frac{1}{2}, 0, \frac{1}{2}, 1\}$ ). Hence, our level 1 agent believes that the planner will almost surely use a payoff type for agent j that is picked at random according to the prior  $\rho_j$  and a level report equal to 0. Further, if our level 1 agent sends a non-zero report,  $z_i \neq 0$ , she will expect the planner to disregard her payoff type report and choose randomly according to  $\rho_i$  with probability 1. However, if our level 1 agents sends a zero report,  $z_i = 0$ , she will expect the planner to use her payoff type and level as reported.

Thus, if she sends the message  $(\theta'_i, k_i, z_i)$  with  $z_i = 0$  and  $k_i = 1$ , she will expect to receive the following lottery over outcomes

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot f^i(\theta'_i,\theta_{-i}).$$

If she sends the message  $(\theta'_i, k_i, z_i)$  with  $z_i = 0$  and  $k_i \neq 1$ , she will expect to receive the following lottery over outcomes

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot\bar{f}(\theta_{i}',\theta_{-i}).$$

And, if she sends the message  $(\theta'_i, k_i, z_i)$  with  $z_i \neq 0$ , she will expect to receive the following lottery over outcomes

$$\sum_{\theta \in \Theta} \rho(\theta) \cdot \bar{f}(\theta).$$

By condition (ii) we have that

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(f^i(\theta),\theta_i) \ge \sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(f^i(\theta'_i,\theta_{-i}),\theta_i)$$

for all  $\theta_i' \in \Theta_i$  .

By condition (iii) we have that

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(f^i(\theta),\theta_i) \ge \sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(\bar{f}(\theta'_i,\theta_{-i}),\theta_i)$$

for all  $\theta'_i \in \Theta_i$ .

It must then also be true that

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(f^i(\theta),\theta_i)\geq \sum_{\theta\in\Theta}\rho(\theta)\cdot u_i(\bar{f}(\theta),\theta_i).$$

Thus, for agent *i* with payoff type  $\theta_i$  and level 1, reporting  $(\theta_i, 1, 0)$  is a best response.

We now prove that an agent with payoff type  $\theta_i$  and level k will send the message  $(\theta_i, k_i, 0)$  by induction on the following statement: Let  $k \ge 1$  and assume that if for all  $l \in \{1, \ldots, k-1\}, \theta_j \in \Theta_j$ , and  $j \in I$  an agent j with payoff type  $\theta_j$  and level l will report  $(\theta_j, l, 0)$ , then an agent i with payoff type  $\theta_i$  and level k will report  $(\theta_i, k, 0)$ .

The result is true for k = 1 by the above argument. Now consider an agent i with payoff type  $\theta_i$  and level  $k \in \{2, \ldots, \bar{k}\}$ . She expects that all other agents will be sending reports  $z_j = 0, k_j = k - 1$ , and truthfully reporting their payoff type. Thus, she expects that the social planner will always take their payoff and level reports as given.

Thus, if she sends the message  $(\theta'_i, k, 0)$  she will expect to receive the following lottery over outcomes

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot f^{i}(\theta'_{i},\theta_{-i}).$$

If she sends the message  $(\theta'_i, k_i, 0)$  with  $k_i \neq k$ , she will expect to receive the following lottery over outcomes

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot\bar{f}(\theta_{i}',\theta_{-i})$$

for some  $j \in I$ .

And, if she sends the message  $(\theta'_i, k_i, z_i)$  with  $z_i \neq 0$ , she will expect to receive the following lottery over outcomes

$$\sum_{\boldsymbol{\theta}\in\Theta} \rho(\boldsymbol{\theta}) \cdot \bar{f}(\boldsymbol{\theta})$$

By the same argument above, conditions (ii) and (iii) then imply that our agent will send the report  $(\theta_i, k, 0)$ .

Therefore, if we define  $m_i(\theta_i, k_i) = (\theta_i, k_i, 0)$  for all  $\theta_i \in \Theta_i, k_i \in \{1, \dots, \bar{k}\}$ and  $i \in I$ , then *m* is a level-k solution under  $\alpha$  and achieves *F* by condition (*i*). Hence, *F* is level-k implementable.

	г		
	L		
	L		

Remark 1. The case when n = 2. The case when n = 2 requires more stringent sufficient conditions than when  $n \ge 3$ . This is because a unilateral deviation has an additional effect when n = 2. To account for this, we require an extra sufficient condition, in addition to conditions (i)-(iii). The extra condition is needed because when n = 2 it is possible to deviate unilaterally to mimic other level-k belief structures. For example, consider agent 1 with level 2 (and take  $\bar{k} = 2$ ) and the mechanism used in Proposition 2. This agent believes everyone else is level 1. When n = 2, this means she believes she could receive outcomes under either  $f^2$ ,  $\bar{f}$ , or  $f^1$  when reporting level 0, level 1, or level 2 respectively. However, when  $n \ge 3$ , she believes she will receive outcomes under  $\bar{f}$ ,  $\bar{f}$ , or  $f^1$  when reporting level 0, level 1 respectively. Thus, we need an additional condition when n = 2: truthfully reporting payoff type under  $f^i$  must be preferable to reporting any other payoff type under  $f^j$ . This condition is not required when  $n \ge 3$ , because a unilateral deviation for agent i will never lead to outcomes associated with  $f^j$ .

Remark 2. Relationship to Bayesian incentive constraints. The level-k necessary and sufficient conditions generalize the standard Bayesian incentive constraints. The difference between the two is that the level-k conditions can be satisfied with a different function,  $f^i$ , for each agent, whereas the Bayesian incentive constraints must hold using the same function, f, for all agents. The intuition for this is straightforward. The relaxation of the cross-player restriction  $(f^1 = \cdots = f^n)$  arises because of the relaxation of consistent beliefs under the level-k model, i.e. a level 3 agent believes she is facing level 2 agents while a level 2 agent believes she is facing level 1 agents. Thus, all incentive constraints can be satisfied by different f functions for each agent: an agent  $\mathbf{i}$  with payoff type profile  $\theta_i$  and level k thinks she will receive  $f^i(\theta)$ when playing against level k - 1 agents with payoff type profile  $\theta_{-i}$ , while an agent  $\mathbf{j}$  with payoff type  $\theta_j$  and level k thinks she will receive  $f^j(\theta)$  when playing against level k - 1 agents with payoff type profile  $\theta_{-j}$ . Agents with level k never believe they are playing against other agents with level k. Because of these (potentially) inconsistent beliefs, the planner can promise different agents outcomes derived from different functions,  $f^1, \ldots, f^n$ .

Further, notice that if the Bayesian incentive constraints are satisfied for some function f, then conditions (i)-(iii) are automatically satisfied when  $f^i = f$  for all  $i \in I$  and  $\bar{f} = f$ . Therefore, if a social choice rule is Bayesian implementable then it will also be level-k implementable under an atomless anchor in independent private value environments. However, it is not necessarily the case that if a social choice rule is level-k implementable then it will be Bayesian implementable. We illustrate this with a bilateral trade application in Section 6.

Remark 3. Full vs weak implementation. Throughout this paper we use the notion of weak level-k implementation. The message a type sends in a level-k solution needs to be only a weak best response. This allows for the possibility of multiple level-k solutions where some of these solutions are not consistent with the social choice rule (notice, that multiple level-k solutions only arise through indifferences). We could achieve full implementation by requiring strict inequalities in condition (ii) and (iii) in 2. Simple arguments can then be applied in the proof of Proposition 2 to show that that conditions (i)-(iii) with strict inequalities are sufficient for full level-k implementation.

de Clippel et al. (2019) study full level-k implementation under the restriction to single-valued choice rules. They find that Bayesian incentive constraints are necessary conditions for full level-k implementation. In the next section, we show that if we restrict attention to single-valued choice rules, we replicate de Clippel et al.'s finding for weak level-k implementation. But this result does not hold necessarily for multi-valued choice rules under either weak or full level-k implementation. In section 6 we give an example that shows that ex post efficient trade is fully level-k implementable while it is not Bayesian implementable.

Remark 4. Level 0 anchor. First off, notice that the necessary conditions specified in Proposition 1 do not place any assumptions on the anchor. Thus, the necessary conditions are independent of any level 0 behavior assumptions. This arises because the necessary conditions are generated only from the incentive requirements of a level 2 (or higher) agent which are not directly affected by the anchor specification.

Second, notice that the sufficient conditions depend on the anchor. Proposition 2 holds only for the case of atomless anchors. However, we can also consider alternative anchor assumptions such as truthful reporting. The idea that agents truthfully report their payoff types is a common consequence in the literature that follows from the revelation principle. It is also an assumption that has been applied in behavioral mechanism design (for examples see Crawford & Iriberri (2007b) which analyzes bidding in auctions when level 0 agents truthfully bid their payoff type and Saran (2011) for a bilateral trade example when it assumed that there are a proportion of agents that truthfully report their payoff types). However, when dealing with mechanisms where agents need to report (at minimum) their payoff type and their level, this opens up the question of what it means for an agent to truthfully report her level. In the previous sections, the language of agents truthfully reporting their payoff types and levels was regularly used. But, that need not of been interpreted literally as an agent having an understanding of their payoff type and level, rather it was simply a label put on a message that the agent had an incentive to send (i.e. truth telling is a consequence of the mechanism, but not an assumption). However, when we discuss whether agents truthfully report their types, when they are not incentivized to do so, as is the case for level 0 agents, we must take the idea of what it means to truthfully report payoff types and levels much more seriously. It seems most likely that an agent has no real conceptualization of what her level of reasoning is, and it might be unreasonable for a planner to think agents will truthfully report their levels when they are not incentivized to do so. This is in contrast to the idea of a payoff type, where an agent is likely able to conceptualize how she much she values a good in an auction, for example, and may feel inclined to truthfully report because she prefers not to lie.

Given this, considering the assumption that others will truthfully report, we may want to consider the possibility that level 0 agents truthfully report their payoff type, but may not truthfully report with respect to other dimensions of the message space. For example, one can interpret the mechanism underlying the proof of Proposition 2 as asking subjects to report their payoff type, an integer between 0 and k, and a real number in [-1, 1]. It may be valuable to think about level 0 subjects here truthfully reporting their payoff type, but then randomly choosing integers and real numbers. Thus, we can think of a truth telling anchor over a mechanism with message space  $M_i = \Theta_i \times \{0, 1, \dots, \bar{k}\} \times Z_i$  as one which maps, for each agent *i*, each type  $t_i = (\theta_i, k)$  to some probability distribution over  $\{\theta_i\} \times \{0, 1, \dots, \bar{k}\} \times Z_i$ . Notice that the above proof would easily extend to capture this notion of truthful reporting. Given the message space  $M_i = \Theta_i \times \{0, 1, \dots, k\} \times [-1, 1]$ , the only requirement for the proof of Proposition 2 to go through is for the anchor to be atomless on [-1, 1]. Thus, the same level-k sufficient conditions extend to the case where level 0 agents truthfully report their payoff type (and/or possibly truthfully report their level). Thus the same mechanism can be used to level-k implement a social choice rule under conditions of both truth telling and other behavior as long as the anchor is atomless over [-1, 1].

Remark 5. Cognitive hierarchy. We can extend the results beyond the simple level- k model to more general limited depth of reasoning models like cognitive hierarchy.<sup>11</sup> Specifically, under the level-k model, if an agent is level k, she believes that others have levels exactly equal to (k - 1). In general, we might allow an agent with level k to hold beliefs over all lower levels. As long as a level k type only puts weight on lower levels, the spirit of limited depth of reasoning is maintained with each type being able to calculate her optimal

<sup>&</sup>lt;sup>11</sup>In the cognitive hierarchy model, a level k type has beliefs over all lower levels determined by a conditional Poisson distribution. See Camerer et al. (2004) for specifics.

action recursively, in a finite number of steps.

Appendix C addresses the question of whether we can design a mechanism that is robust to relaxing the level-k belief assumption. In this appendix we generalize our type space and solution concept to the limited depth of reasoning (LDoR) concept to relax beliefs about the depths of reasoning of others and consider a form of robust implementation where we ask whether there exists a mechanism that will implement a social choice rule for any LDoR type space. We show that if you strengthen the sufficient conditions as in the n = 2 case (i.e. truthfully reporting payoff type under  $f^i$  is preferable to reporting any other payoff type under  $f^i$ ,  $\bar{f}$  or  $f^j$  for all  $i, j \neq i \in I$ ), then these conditions are sufficient to implement a social choice rule for any LDoR type space. The additional restriction is required for the same reasons as when n=2, a unilateral deviation can lead to the mimicking of other level-k belief structures. For example, consider an agent 1 with level 2 who believes agent 2 is level 1 and agent 3 is level 0 (a possible belief under a more general limited depth of reasoning model). If agent 1 reports a level of 0, this would mimic the belief structure for an agent 2 with level 1(who believes everyone else is level 0).

Whether relaxing the cross-player restriction,  $(f^1 = \cdots = f^n)$ , imposed under Bayesian implementation has bite depends on the environment. In the next section we consider two special environments which lead to the cross-player restrictions being automatically imposed. In these environments, Bayesian incentive constraints are necessary conditions for level-k implementation. Further, we show, via a bilateral trade example, in Section 6 that this result does not hold generally: ex post efficient trade is level-k implementable while it is not Bayesian implementable.

### 4 Special environments

In this section we look at two restricted environments. In the first, we restrict attention to single-valued social choice rules and in the second we restrict attention to mechanisms where the message set is equal to the set of payoff types. These special environments are the environments studied in de Clippel et al. (2019) and Crawford (2019) respectively. In both of these cases, we show that Bayesian incentive constraints are necessary conditions for level-k implementation. This establishes parallel results to those found by de Clippel et al. and Crawford.

Corollary 1 formalizes this result for the restriction to social choice functions.

**Corollary 1.** Let F be a single-valued social choice rule. Let  $\mathcal{B}$  be a Bayesian type space and let  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,...,n}; \bar{k} \rangle$  be a ( $\mathcal{B}$ -based) level-k type space with  $\bar{k} \geq 2$ . Then F is level-k implementable if it is Bayesian implementable

Proof:

From Proposition 1, there exists functions  $\overline{f}: \Theta \to \triangle(Y)$  and  $f^i: \Theta \to \triangle(Y)$ , for all  $i \in I$  such that conditions (i)-(iii) hold. Since F is a single-valued social choice rule, it must be that  $\overline{f} = f^1(\theta) = \cdots = f^n(\theta) = F(\theta)$  for all  $\theta \in \Theta$ . Therefore, it follows that F is Bayesian implementable using the mechanism  $\langle \Theta, F \rangle$ .

The following proposition demonstrates that restricting the message space to the set of payoff types has the same effect as restricting our social choice rule to a social choice function - the level-k incentive constraints collapse down to the Bayesian incentive constraints. In other words, if we restrict the message space to be the set of payoff types, then Bayesian incentive constraints are necessary conditions for level-k implementation,.

In order to show this we need to assume a richness condition on the environment, Assumption A<sup>\*</sup>:

**A\*:** For every  $i \in I$  and for any two payoff types  $\theta_i \neq \theta'_i \in \Theta_i$  there exists a  $\theta_{-i} \in \Theta_{-i}$  such that  $F(\theta_i, \theta_{-i}) \cap F(\theta'_i, \theta_{-i}) = \emptyset$ .

Assumption A\* is likely to hold in many environments. For example, as we'll see in Section 6, it holds in the bilateral trade environment with the expost efficient social choice rule as long as for every two values of the buyer  $v < v' \in$ V there exists a cost for the seller,  $c \in C$ , that falls between,  $v \leq c \leq v'$ . And, similarly for the seller. Proposition 3 establishes the result. **Proposition 3.** Let F be a social choice rule. Let  $A^*$  holds. Let  $\mathcal{B}$  be a Bayesian type space and let  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,\dots,n}; \bar{k} \rangle$  be a ( $\mathcal{B}$ -based) level-k type space with  $\bar{k} \geq 2$ . If we restrict mechanisms to only allow messages about payoff types,  $M_i = \Theta_i$  for all  $i \in I$ , then F is level-k implementable if it is Bayesian implementable.

#### Proof:

Suppose that the social choice rule F is level-k implementable. Then there exists some mechanism  $\langle \Theta, g \rangle$  and a function  $m_i : \Theta_i \times \{0, \ldots, \bar{k}\} \to \Delta(\Theta_i)$  for each  $i \in I$  such that  $m = m_1 \times \cdots \times m_n$  is a level-k solution and achieves F.

Suppose that there exists some agent  $i \in I$ , and two types for agent iwith  $\theta_i \neq \theta'_i \in \Theta_i$  such that  $\operatorname{supp}(m_i(\theta_i, k)) \cap \operatorname{supp}(m_i(\theta'_i, j)) \neq \emptyset$  for some  $j, k \in \{1, \ldots, \bar{k}\}$ . Since these two types send the same message as part of a level-k solution that achieves F, the social planner must be satisfied with them receiving the same outcome under every payoff profile i.e. it must follow that for  $m \in \operatorname{supp}(m_i(\theta_i, k)) \cap \operatorname{supp}(m_i(\theta'_i, j)), g(m, m_{-i}(\theta_{-i}, 1)) \in$  $F(\theta_i, \theta_{-i}) \cap F(\theta'_i, \theta_{-i})$  for every  $\theta_{-i} \in \Theta_{-i}$ . This contradicts assumption  $A^*$ . Therefore, it must true that for any two types for agent i with  $\theta_i \neq \theta'_i \in \Theta_i$ that  $\operatorname{supp}(m_i(\theta_i, k)) \cap \operatorname{supp}(m_i(\theta'_i, j)) = \emptyset$  for any  $j, k \in \{1, \ldots, \bar{k}\}$ . Thus, it follows that we can define the function  $\psi_k^i : \Theta_i \to \Theta_i$  by  $\psi_k^i(\theta_i) = m_i(\theta_i, k)$ and that it is both 1-1 and onto.

**Claim**:  $m_i(\theta_i, k) = m_i(\theta_i, j)$  for all  $j, k \in \{1, \dots, \bar{k}\}, \theta_i \in \Theta_i$ .

To see this, suppose not. Then there exists a  $\theta_i \in \Theta_i$  and some  $j, k \in \{1, \ldots, \bar{k}\}$  with  $j \neq k$  such that  $m_i(\theta_i, k) \neq m_i(\theta_i, j)$ .

Because,  $\psi_j^i$  is 1-1 and onto there must exist some  $\theta_i' \neq \theta_i \in \Theta_i$  such that  $\psi_j^i(\theta_i') = m_i(\theta_i, k)$ . But, then  $m_i(\theta_i', j) = m_i(\theta_i, k)$  which is a contradiction.

Therefore, it follows that  $m_i(\theta_i, k) = m_i(\theta_i, j)$  for all  $j, k \in \{1, \dots, \bar{k}\}$ .

Define  $\tilde{g}: \Theta \to \triangle(Y)$  by  $\tilde{g}(\theta) = g(m(\theta, 2))$  for all  $\theta \in \Theta$ .

Then, for any  $i \in I$ ,  $\theta \in \Theta$ , and  $\theta' \in \Theta_i$ 

$$\begin{split} \sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(\tilde{g}(\theta),\theta) &- \sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(\tilde{g}(\theta',\theta_{-i}),\theta) \\ &= \sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(g(m_{i}(\theta_{i},2),m_{-i}(\theta_{-i},2)),\theta) \\ &- \sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(m_{i}(\theta',2),m_{-i}(\theta_{-i},2)),\theta) \\ &= \sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(g(m_{i}(\theta_{i},2),m_{-i}(\theta_{-i},1)),\theta) \\ &- \sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(m_{i}(\theta',2),m_{-i}(\theta_{-i},1)),\theta) \\ &\geq 0 \end{split}$$

The inequality follows from the fact m is a level-k solution. Thus, the Bayesian incentive constraints hold for  $\tilde{g}$ . Further,  $\langle \Theta, \tilde{g} \rangle$  achieves F because  $\tilde{g}(\theta) = g(m((\theta, 2))) \in F(\theta)$ . Therefore, it follows by definition that F is Bayesian implementable using the mechanism $\langle \Theta, \tilde{g} \rangle$ .

Remark 6. Equivalence between level-k and Bayesian implementation. Under the conditions of Corollary 1 we get an equivalence between Bayesian implementation and level-k implementation under atomless anchors in independent private value environments. This is because we can show that if a single-valued social choice rule F is Bayesian implementable then it is level-k implementable by setting  $\overline{f} = F$  and  $f^i = F$  for all  $i \in I$  and applying Proposition 2. However, we cannot guarantee an equivalence between Bayesian and level-k implementation under the conditions of Proposition 3 because Proposition 2 will not apply if we place restrictions on the message space. Therefore, it may be the case that a social choice rule is Bayesian implementable while it is not level-k implementable under mechanisms where the message space is restricted to the set of payoff types.

### 5 Ex post level-k implementation

This section addresses the question of whether we can design a mechanism that is robust to relaxing the common prior assumption that is present in the level-k type space - the assumption that beliefs about payoffs are determined by a specific common prior. To address this, we generalize our solution concept to that of ex post level-k implementation. This effectively relaxes the common prior assumption that has been maintained so far. This section shows that the relationship between ex post level-k and ex post implementation is analogous to the relationship between level-k and Bayesian implementation. The ex post level-k solution concept is defined below.

**Definition 8.** For a given game defined by a mechanism  $\langle M, f \rangle$  and a  $\bar{k} \in \mathbb{N}_+$ , a strategy profile  $s = s_1 \times \cdots \times s_n$ , with  $s_i : \Theta_i \times \{0, \ldots, \bar{k}\} \to \Delta(M_i)$  for all  $i \in I$ , is the **ex post level-k solution under anchor**  $\alpha$  if and only if:

- (i)  $s_i(t_i) \sim \alpha_i(t_i)$  for all  $t_i \in \Theta_i \times \{0\}$
- (ii)  $u_i(f(s_i(\theta_i, k), s_{-i}(\theta_{-i}, k-1)), \theta) \ge u_i(f(m'_i, s_{-i}(\theta_{-i}, k-1)), \theta) \forall m'_i \in M_i, \theta \in \Theta, \text{ and } k \ge 1, i \in I.$

As in the level-k solution, the ex post level-k solution specifies that all level 0 types play according to the specified anchor. Given this, all agents with levels at least 1 play a best response given their beliefs that others have levels exactly one level lower than them for *all* possible realizations of the payoff types of others. This solution concept allows us to define an ex post level-k implementation concept that is an analogue to ex post implementation: there must exist a mechanism and an ex post level-k solution which is consistent with the social choice rule for any realization of payoff types,  $\theta \in \Theta$ .

**Definition 9.** Fix a  $k \in \mathbb{N}_+$ . A social choice rule is **ex post level-k implementable under anchor**  $\boldsymbol{\alpha}$  if there exists a mechanism  $\langle M, f \rangle$  and a message profile  $m_i : \Theta_i \times \{0, \dots, \bar{k}\} \to \Delta(M_i)$  for all  $i \in I$ , such that  $m = m_1 \times \cdots \times m_n$ is an ex post level-k solution under anchor  $\boldsymbol{\alpha}$  and  $f(m(\theta \times \hat{k})) \in F(\theta)$  for all  $\theta \times \hat{k} \in \times \{\Theta_i \times \{1, \dots, \bar{k}\}\}_{i \in I}$ .

We will be interested in comparing ex post level-k implementation with that of ex post implementation. We define ex post implementation for completeness. In the below definition, we incorporate the results of the revelation principle, and hence define ex post implementation for a direct mechanism. Condition (ii) below states the standard ex post incentive constraints which will be contrasted with the ex post level-k incentive constraints developed in the next subsection.

**Definition 10.** A social choice rule F is **ex post implementable** if there exists a mechanism  $\langle \Theta, f \rangle$  such that

- (i)  $f(\theta) \in F(\theta) \ \forall \theta \in \Theta$
- (ii)  $u_i(f(\theta), \theta) \ge u_i(f(\theta', \theta_{-i}), \theta) \ \forall \theta' \in \Theta, \ \theta \in \Theta, \ i \in I$

### 5.1 Necessary and sufficient conditions for ex post levelk implementation

This section gives the necessary and sufficient conditions for ex post level-k implementation. The conditions are stated separately as the necessary conditions hold in general environments and the sufficient conditions hold in private value environments under atomless anchors.

The necessary conditions for ex post level-k implementation are given in Proposition 4 and are analogous to those in Proposition 3. The proof follows analogously to the proof of Proposition 3 as well and can be found in Appendix A.

**Proposition 4.** (Ex post Necessary Conditions) Let F be a social choice rule. Let  $\bar{k} \geq 2$ . If F is ex post level-k implementable, then there exists a function  $f^i: \Theta \to \Delta(Y)$  for each  $i \in I$  and a function  $\bar{f}: \Theta \to \Delta(Y)$ , such that the following conditions hold:

(i) 
$$f^{i}(\theta), \bar{f}(\theta) \in F(\theta) \ \forall \theta \in \Theta, \forall i \in I$$

(ii) 
$$u_i(f^i(\theta), \theta) \ge u_i(f^i(\theta', \theta_{-i}), \theta) \ \forall \theta' \in \Theta_i, \ \theta \in \Theta, \ i \in I$$

(iii) 
$$u_i(f^i(\theta), \theta) \ge u_i(\bar{f}(\theta', \theta_{-i}), \theta) \ \forall \theta' \in \Theta_i, \ \theta \in \Theta, \ i \in I$$

Similar to the case for level-k implementation, the necessary conditions are sufficient for ex post level-k implementation under atomless anchors for private value environments when  $n \ge 3$  (the case when n = 2 can be found in Appendix B). The proof follows analogously to that of Proposition 2 and can be found in Appendix A.

**Proposition 5.** (Ex post Sufficient Conditions) Let F be a social choice rule. Let the environment be one of private values and let  $n \ge 3$ . If there exists a function  $f^i : \Theta \to \triangle(Y)$  for each  $i \in I$  and a function  $\overline{f} : \Theta \to \triangle(Y)$ such that conditions (i)-(iii) hold in Proposition 4 then F is expost level-k implementable under atomless anchors.

Remark 7. Relationship to ex post implementation. The conditions in Proposition 4 generate a set of ex post level-k incentive constraints that generalize the standard ex post incentive constraints. The difference between the two is that the ex post level-k incentive constraints can be satisfied with a different function,  $f^i$ , for each agent, whereas the ex post incentive constraints must hold using the same function, f, for all agents. Further, if the social choice rule is ex post implementable then it is ex post level-k implementable under atomless anchors in private value environments. To see this notice that if the mechanism  $\langle \Theta, f \rangle$  ex post implements the social choice rule, then setting  $\bar{f} = f^i = f$  for all  $i \in I$  will satisfy the conditions in Proposition 5.

Remark 8. Special Environments. In Section 4 we considered two special environments which lead to the cross-player restrictions in the level-k incentive constraints being automatically imposed. Simple arguments extend these results to the ex post environment as well. In these two special environments, a single-valued choice rule or restricting the message space to the set of payoff types, ex post incentive constraints are necessary conditions for ex post level-k implementation. However, this need not hold in general. In Section 6 we give a counter-example. We show that ex post efficient bilateral trade is ex post level-k implementable while it is not ex post implementable.

Remark 9. Ex post level-k implementability implies level-k implementability. Just as ex post implementation implies Bayesian implementation, it is also true that ex post level-k implementation implies level-k implementation. This is easy to see because the ex post level-k sufficient conditions imply that the level-k sufficient conditions will hold for any common prior.

### 6 Bilateral trade

#### 6.1 Bilateral trade environment

The remainder of this paper focuses on the bilateral trade environment. Buyer's values are given by a finite set V. Seller's costs are given by a finite set C. The set of outcomes is given by  $Y = \mathbb{R} \cup \{\emptyset\}$ , where outcome  $\emptyset$  indicates that the good is not traded and outcome  $p \in \mathbb{R}$  indicates that the good is traded at price p. Agents have quasi-linear utility functions,  $u_b : Y \times V \to \mathbb{R}$  and  $u_c : Y \times C \to \mathbb{R}$ . For any outcome p, the utility of a buyer with a valuation v is

$$u_b(p,v) = \begin{cases} v-p & \text{if } p \in \mathbb{R} \\ 0 & \text{otherwise} \end{cases}$$

and the utility of a seller with a cost c is

$$u_s(p,c) = \begin{cases} p-c & \text{if } p \in \mathbb{R} \\ 0 & \text{otherwise} \end{cases}$$

We are interested in mechanisms that satisfy the ex post efficient social choice rule  $F^*(v,c) = \{y | y \in \mathbb{R} \text{ if } v \geq c \text{ and } y = \emptyset \text{ otherwise}\}$ . The ex post efficient choice rule requires trade whenever the buyer's value is above the seller's cost. We are also interested in the mechanism satisfying two additional properties: budget balance (the price paid by the buyer equals the price received by the seller - this is already imposed by the description of the environment) and ex post individual rationality (both the buyer and seller prefer to participate in the trading institution than receive the utility of 0).

#### 6.2 Ex post efficient trade - general environment

This section contains the main result: ex post efficient trade is ex post level-k implementable under atomless anchors. As a result, ex post efficient trade is also level-k implementable under atomless anchors.

Notice that Corollary 1 does not apply in this environment. This is because

the ex post efficient choice rule,  $F^*$ , is a multi-valued choice rule. If the buyer's value is above the seller's cost, ex post efficiency requires trade, but the planner does not care at what price the good is traded. This means it may be possible to level-k implement ex post efficient trade even if it is not Bayesian implementable.

Also, notice that if we impose an additional assumption, A1, on the decision environment then the conditions for Proposition 4 are satisfied. This means that in order to implement ex post efficient trade under level-k implementation, we will need to use mechanisms with messages spaces larger than the set of payoff types.

**A1:** For any  $v, v' \in V$  with v' < v, there exists a  $c \in C$  such that  $v' \leq c \leq v$ . And, for any  $c, c' \in C$  with c < c', there exists a  $v \in V$  such that  $c \leq v \leq c'$ .

Proposition 6 establishes the result.

**Proposition 6.** The ex post efficient social choice rule,  $F^*$ , is ex post level-k implementable under a mechanism that satisfies budget balance and ex post individual rationality.

The proof of Proposition 6 follows by showing that there exists functions  $f^b: V \times C \to \Delta(Y), f^s: V \times C \to \Delta(Y)$  and  $\bar{f}: V \times C \to \Delta(Y)$  that satisfy sufficient conditions for level-k implementation when n = 2 (Proposition 8 in Appendix B). The formal proof can be found in Appendix A, but, the basic intuition works as follows. Consider three ways to determine the price if there is trade: (1) give the buyer all the surplus (i.e. set the price to seller's reported cost); (2) give the seller all the surplus (i.e. set the price to buyer's reported value); (3) have the buyer and seller split the surplus (i.e. set the price equal to the average of reported cost and value). Use these to define the functions  $f^b, f^s$ , and  $\bar{f}$  respective. Formally, define

$$f^{b}(v,c) = \begin{cases} c & \text{if } c \leq v \\ \emptyset & \text{otherwise} \end{cases},$$
$$f^{s}(v,c) = \begin{cases} v & \text{if } c \leq v \\ \emptyset & \text{otherwise} \end{cases},$$

and

$$\bar{f}(v,c) = \begin{cases} \frac{v+c}{2} & \text{if } c \le v \\ \emptyset & \text{otherwise} \end{cases}$$

Now, let every agent report their cost or value and their level. If a buyer reports a level exactly one higher than the seller than let trade and prices be determined according to  $f^b$ . If a seller reports a level exactly one higher than the buyer, then let trade and prices be determined according to  $f^s$ . Otherwise, let trade and prices be determined according to f. In other words, trade occurs if and only if the reported value is above the reported cost and at the price that is most favorable to the agent that reports exactly one level higher than the other agent; otherwise they split the surplus. Suppose now that all agents truthfully report their payoff type and level. Given that, both buyers and sellers have an incentive to truthfully report their levels as they will then receive the most favorable prices for themselves. Further, if an agent reports their level truthfully, then they also have an incentive to truthfully report their payoff type, since messages sent will only affect the likelihood of trade but not the price. This ensures that truthfully reporting one's own payoff type and level is a best response given that everyone else is truthfully reporting their level, regardless what the payoff type realization is. The only thing left is for the planner to incentivize level 1 agents; this can be done as in the proof of Proposition 2, by having the planner screen for level 0 types by using the reported real number from [-1, 1].

#### 6.3 Ex post efficient trade - a 2 type example

In this subsection, we go through a simple 2-type example to give a concrete illustration of a mechanism that is level-k implementable. We will also use this example to demonstrate that it is possible to design a mechanism that fully level-k implements ex post efficient trade, i.e. under the mechanism, any level-k solution will be consistent with ex post efficient trade.

In this example, the seller has two possible costs:  $C = \{2, 6\}$ , and the buyer has two possible values:  $V = \{3, 7\}$ . Types are drawn from a uniform common prior,  $\rho$  (i.e.  $\rho(v, c) = \frac{1}{4}$  for all  $v, c \in V \times C$ ).

Claim 1. Ex post efficient trade is not Bayesian implementable.

This was shown by Matsuo (1989) who gives sufficient conditions for expost efficiency in the two type bilateral trade environment. To see the intuition, recall that the revelation principle ensures that we need only consider mechanisms where all agents truthfully report their type. Further, the low valued buyer and the high valued seller should receive zero utility in equilibrium. Thus any candidate mechanism must take the form of the one in Figure 1 for some  $p \in \mathbb{R}$ . This mechanism should be understood in the following way: the buyer chooses the row message, the seller chooses the column message, and the corresponding element in the table is the outcome that occurs. For example, if the buyer sends message  $m_7$  and the seller sends message  $m_3$  and the seller sends message  $m_6$  then there is no trade.

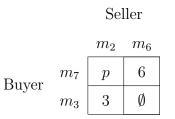


Figure 1: Structure of a Bayesian mechanism

A high valued buyer believes the low and high cost seller types are equally likely (comes from the uniform prior assumption) and thus believes actions  $m_2$ and  $m_6$  are equally likely. Thus, for a high valued buyer to truthfully report her payoff type the trading price must be less than or equal to 4, i.e.  $p \leq 4$ . Similarly, a low cost seller believes the low and high valued buyer types are equally likely and thus believes actions  $m_7$  and  $m_3$  are equally likely. For a low cost seller to truthfully report her payoff type the trading price must be greater than or equal to 5, i.e.  $p \geq 5$ . These two conditions are incompatible. There is no mechanism that will implement the ex post efficient choice rule under Bayesian implementation.

Claim 2. Ex post efficient trade is ex post level-k implementable under a uniform random anchor.

Suppose, for the sake of this example, that there are only level 0, level 1 and level 2 types in the population. As such, there are six different types for both the buyer and the seller. There is both a high-valued and a low-valued payoff type for each of the three reasoning types.

		$m_0$	$m_{(2,1)}$	$m_{(2,2)}$	$m_{(6,-)}$
Buyer	$m_0$	4.5	7	3	Ø
	$m_{(7,1)}$	2	4.5	6	6
	$m_{(7,2)}$	6	3	4.5	6
	$m_{(3,-)}$	Ø	3	3	Ø

Seller

Figure 2: A level-k mechanism

The mechanism in Figure 2 ex post level-k implements the ex post efficient choice rule. To see this notice that level 0 agents (regardless of their payoff type) are assumed to play each action with equal probability. A level 1 type then believes that her opponent is playing each action with equal probability. Therefore, playing  $m_{(3,-)}$  is a best response for the low valued level 1 buyer and playing  $m_{(7,1)}$  is a best response for the high valued level 1 buyer. Likewise, playing  $m_{(6,-)}$  is a best response for the high cost level 1 seller and playing  $m_{(2,1)}$  is a best response for the low cost level 1 seller.

Effectively, a level 2 buyer may hold any beliefs over payoff types but believes her opponent is a level 1 type. For any beliefs about the payoff types of the seller, playing  $m_{(3,-)}$  is a best response for a low valued level 2 buyer and playing  $m_{(7,2)}$  is a best response for a high valued level 2 buyer. Similarly, for any beliefs over payoff types of the buyer, playing  $m_{(6,-)}$  is a best response for a high cost level 2 seller and playing  $m_{(2,2)}$  is a best response for a low cost level 2 seller.

Given the strategies defined by the ex post level-k solution, for any pair of level 1 or level 2 types, the outcome will be consistent with the ex post efficient social choice rule. In other words, if the buyer is the low valued type and the seller is the high cost type, then regardless of whether the buyer and sellers are level 1 or level 2 types, there will not be trade. For any other payoff type profile  $(v, c) \neq (3, 6)$ , regardless of whether the buyer and sellers are level 1 or level 2 types, there will be trade. Ex post individual rationality is also satisfied.

Notice that action  $m_0$  is never played by level 1 or level 2 types. Thus, we would not expect to ever observe action  $m_0$  be played by either the buyer or the seller if there are no level 0 types in the population. However, the mechanism includes that action because level 1 types believe that action is being played with some positive probability by the level 0 type. Thus, these 'level 0' actions influence the behavior of level 1 types even though they are not played by types of higher levels.<sup>12</sup>

*Remark* 10. *Level-k implementation*. The above mechanism also level-k implements ex post efficient trade under the uniform random anchor for any Bayesian-based level-k type space (i.e. under any common prior).

Remark 11. Full level-k implementation. Notice that there is a unique level-k solution under a uniform random anchor in the above mechanism in the level-k type space with the uniform common prior, i.e. all types with levels at least 1 are playing strict best responses. Thus, this mechanism level-k implements ex post efficient trade in both the weak and full sense of implementation.

### 7 Summary

This paper explores the theoretical implications of level-k implementation by defining the boundaries of what is level-k implementable. It gives necessary and sufficient conditions for level-k implementation and establishes the relationship to Bayesian implementation. It relaxes the common prior assumption that underlies level-k and Bayesian implementation by defining the concept of ex post level-k implementation. It gives necessary and sufficient conditions for ex post level-k implementation and shows that the relationship between ex post level-k and ex post implementation mirrors that between level-k and Bayesian implementation.

<sup>&</sup>lt;sup>12</sup>Strategies that are never played have been shown to still have an impact on strategic behavior in other experiments. Cooper et al. (1990) show that introducing dominated actions into coordination games changes behavior.

This paper is not, however, a practical guide to level-k implementation. And, there are many questions that need to be answered to address the practical relevance of level-k mechanisms. The first question being, do level-k mechanisms work? This is an empirical question and one that can potentially be answered with experiments. Will mechanisms of the type underlying the sufficiency proofs in this paper achieve the desired social choice outcomes in the lab? Is the assumption about atomless anchors the right assumption? Level-k models are motivated by the explanation of behavior in novel situations. But, what is a novel situation? Is it a situation where an agent truly has no experience? Or, will level-k mechanisms work even when agents have (limited) experience? Then there is the question of what is the real world analogue of a level-k mechanism? This paper shows that common mechanisms like auctions may not be able to achieve all level-k implementable outcomes as the message space is restricted to the set of payoff types.

While this paper does not provide answers to these questions, it does provide a solid framework to start investigating these questions by supplying tight necessary and sufficient conditions for level-k implementation and providing some insight into the types of mechanisms that might be necessary to achieve level-k implementable outcomes.

### References

- Arad, A. & Rubinstein, A. (2012). The 11-20 Money Request Game: Evaluating the Upper Bound of Level-k Reasoning. *American Economic Review*, 102(7), 3561–73.
- Bergemann, D. & Morris, S. (2005). Robust Mechanism Design. *Econometrica*, 73(6), 1771–1813.
- Bergemann, D. & Morris, S. (2009). Robust Implementation Direct Mechanisms. *Review of Economic Studies*, 76, 1175–1204.
- Bergemann, D., Morris, S., & Tercieux, O. (2011). Rationalizable Implementation. Journal of Economic Theory, 146.

- Börgers, T. & Li, J. (2019). Strategically Simple Mechanisms. *Econometrica*, 87, 2003–2035.
- Brocas, I., Carrillo, J. D., Wang, S. W., & Camerer, C. F. (2014). Imperfect Choice or Imperfect Attention? Understanding Strategic Thinking in Private Information Games. *Review of Economic Studies*, 81(3), 944–970.
- Bulow, J. & Roberts, J. (1989). The Simple Economics of Optimal Auctions. Journal of Political Economy, 97, 1060–1090.
- Camerer, C. F., Ho, T.-H., & Chong, J.-K. (2004). A Cognitive Hierarchy Model of Games. Quarterly Journal of Economics, 119(3), 861–898.
- Cooper, R. W., Dejong, D. V., Forsythe, R., & Ross, T. W. (1990). Selection Criteria in Coordination Games: Some Experimental Results. American Economic Review, 80, 218–233.
- Copic, J. & Ponsati, C. (2008). Robust Bilateral Trade and Mediated Bargaining. *Journal of the European Economic Association*, 6, 570–589.
- Copič, J. & Ponsati, C. (2016). Optimal Robust Bilateral Trade: Risk Neutrality. Journal of Economic Theory, 163, 276–287.
- Costa-Gomes, M. & Crawford, V. P. (2006). Cognition and Behavior in Two-Person Guessing Games: An Experimental Study. American Economic Review, 96(5), 1737–1768.
- Costa-Gomes, M., Crawford, V. P., & Broseta, B. (2001). Cognition and Behavior in Normal-Form Games: An Experimental Study. *Econometrica*, 69(5), 1193–1235.
- Costa-Gomes, M. A., Crawford, V. P., & Iriberri, N. (2013). Structural Models of Nonequilibrium Strategic Thinking: Theory, Evidence, and Applications. *Journal of Economic Literature*, 51, 5–62.
- Crawford, V., Kugler, Neeman, & Pauzner (2009). Behaviorally optimal auction design: examples and observations. *Journal of the European Economic* Association, 7, 377–387.

Crawford, V. P. (2019). Efficient mechanisms for level-k bilateral trading.

- Crawford, V. P. & Iriberri, N. (2007a). Fatal Attraction: Salience, Naivete, and Sophistication in Experimental Hide-and-Seek Games. *American Eco*nomic Review, 97(5), 1731–1750.
- Crawford, V. P. & Iriberri, N. (2007b). Level-k Auctions: Can a Non-Equilibrium Model of Strategic Thinking Explain the Winner's Curse and Overbidding in Private-Value Auctions. *Econometrica*, 75(6), 1721–1770.
- de Clippel, G. (2014). Behavioral implementation. American Economic Review, 104(10), 2975–3002.
- de Clippel, G., Saran, R., & Serrano, R. (2015). Mechanism Design with Bounded Depth of Reasoning and Small Modeling Mistakes. *working paper*.
- de Clippel, G., Saran, R., & Serrano, R. (2019). Level-k mechanism design. *Review of Economic Studies*, 86(3), 1207–1227.
- Eliaz, K. (2002). Fault Tolerance Implementation. Review of Economic Studies, 69, 589–610.
- Eliaz, K. & Spiegler, R. (2006). Contracting with diversely naive agents. *Review of Economic Studies*, 73(3), 689–714.
- Eliaz, K. & Spiegler, R. (2007). A mechanism design approach to speculative trade. *Econometrica*, 75(3), 875–884.
- Eliaz, K. & Spiegler, R. (2008). Optimal speculative trade among large traders. *Review of Economic Design*, 12(1), 45–74.
- Glazer, J. & Rubinstein, A. (1998). Motives and implementation: on the design of mechanisms to elicit opinions. *Journal of Economic Theory*, 79, 157–173.
- Glazer, J. & Rubinstein, A. (2012). A model of persuasion with a boundedly rational agent. *Journal of Political Economy*, 120(6), 1057–1082.
- Gorelkina, O. (2015). The expected externality mechanism in a level-k environment. International Journal of Game Theory, 47, 103–131.

- Hagerty, K. M. & Rogerson, W. P. (1987). Robust Trading Mechanisms. Journal of Economic Theory, 42, 94–107.
- Healy, P. J. (2006). Learning Dynamics for Mechanism Design: An Experimental Comparison of Public Goods Mechanisms. *Journal of Economic Theory*, 129, 114–149.
- Matsuo, T. (1989). On incentive compatible, individually rational, and expost efficient mechanisms for bilateral trading. *Journal of Economic Theory*, 48, 189–194.
- Matsushima, H. (2007). Mechanism Design with Side Payments: Individual Rationality and Iterative Dominance. Journal of Economic Theory, 133, 1–30.
- Matsushima, H. (2008). Detail-free Mechanism Design in Twice Iterative Dominance: Large Economies. *Journal of Economic Theory*, 141, 134–151.
- Mookherjee, D. & Reichelstein, S. (1992). Dominant Stragey Implementation of Bayesian Incentive Compatible Allocation Rules. *Journal of Economic Theory*, 56, 378–399.
- Myerson, R. B. & Satterthwaite, M. A. (1983). Efficient mechanisms for bilateral trading. *Journal of Economic Theory*, 29, 265–281.
- Nagel, R. (1995). Unraveling in Guessing Games: An Experimental Study. American Economic Review, 85(5), 1313–1326.
- Ollar, M. & Penta, A. (2017). Full Implementation and Belief Restrictions. American Economic Review, 107, 2243–2277.
- Saran, R. (2011). Bilateral trading with naive traders. Games and Economic Behavior, 72, 544–557.
- Saran, R. (2016). Bounded Depths of Rationality and Implementation with Complete Information. Journal of Economic Theory, 165, 517–564.
- Severinov, S. & Deneckere, R. (2006). Screening when some agents are nonstrategic: does a monopoly need to exclude? *RAND Journal of Economics*, 37, 816–840.

- Stahl, D. O. & Wilson, P. W. (1994). Experimental Evidence on Player's Models of Other Players. *Journal of Economic Behavior and Organization*, 25(3), 309–327.
- Stahl, D. O. & Wilson, P. W. (1995). On Player's Models of Other Players: Theory and Experimental Evidence. *Games and Economic Behavior*, 10(1), 218–254.
- Strzalecki, T. (2014). Depth of Reasoning and Higher Order Beliefs. Journal of Economic Behavior and Organization, 108, 108–122.
- Wolitzky, A. (2016). Mechanism design with maxmin agents: Theory and an application to bilateral trade. *Theoretical Economics*, 11(3), 971–1004.

# Appendix A Omitted proofs

### **Proof of Proposition 1:**

Suppose that the social choice rule F is level-k implementable under anchor  $\alpha$ . Then there exists some mechanism  $\langle M, g \rangle$  and a function  $m_i : T_i \rightarrow \Delta(M_i)$  for each  $i \in I$  such that  $m = m_1 \times \cdots \times m_N$  is a level-k solution and achieves F.

Consider the behavior of an agent *i* with type  $t_i = (\theta_i, 2)$ . Then, it must be true that:

$$\sum_{\substack{\theta_{-i}\in\Theta_{-i}}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(g(m_{i}(\theta_{i},2),m_{-i}(\theta_{-i},1)),\theta)$$

$$\geq \sum_{\substack{\theta_{-i}\in\Theta_{-i}}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(m_{i}(\theta_{i}',2),m_{-i}(\theta_{-i},1)),\theta) \quad \forall \theta_{i}'\in\Theta_{i}$$

$$(1)$$

As well, it must be true that:

$$\sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(g(m_{i}(\theta_{i},2),m_{-i}(\theta_{-i},1)),\theta)$$

$$\geq \sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}|\theta_{i}) \cdot u_{i}(m_{i}(\theta_{i}',1),m_{-i}(\theta_{-i},1)),\theta) \quad \forall \theta_{i}' \in \Theta_{i}$$

$$(2)$$

Define  $f^i(\theta) = g(m_i(\theta_i, 2), m_{-i}(\theta_{-i}, 1))$  and  $\bar{f}(\theta) = g(m_i(\theta_i, 1), m_{-i}(\theta_{-i}, 1))$ for all  $\theta \in \Theta$  and for all  $i \in I$ .

Condition (ii) holds from (1). Condition (iii) holds from (2). Condition (i) holds by definition of  $\langle M, g \rangle$  and *m* level-k implementing *F*.

## **Proof of Proposition 4:**

Suppose that the social choice rule F is expost level-k implementable under anchor  $\alpha$ . Then there exists some mechanism  $\langle M, g \rangle$  and function  $m_i: T_i \to \Delta(M_i)$  for each i such that  $m = m_1 \times \cdots \times m_N$  is an expost level-k solution and achieves F.

Consider the behavior of an agent i with level 2 and payoff type  $\theta_i$  . Then, it must be true that:

$$u_i(g(m_i(\theta_i, 2), m_{-i}(\theta_{-i}, 1)), \theta) \geq u_i(m_i(\theta'_i, 2), m_{-i}(\theta_{-i}, 1)), \theta) \quad \forall \theta'_i \in \Theta_i, \theta_{-i} \in \Theta_{-i}$$

$$(3)$$

Similarly, it must be true that

$$u_i(g(m_i(\theta_i, 2), m_{-i}(\theta_{-i}, 1)), \theta) \geq u_i(m_i(\theta'_i, 1), m_{-i}(\theta_{-i}, 1)), \theta) \quad \forall \theta'_i \in \Theta_i, \theta_{-i} \in \Theta_{-i}$$

$$\tag{4}$$

Define  $f^i(\theta) = g(m_i(\theta_i, 2), m_{-i}(\theta_{-i}, 1))$  for all  $\theta \in \Theta$  and for all  $i \in I$ . And, define  $\bar{f}(\theta) = g(m(\theta, 1))$  for all  $\theta \in \Theta$ .

Condition (ii) holds from (3). Condition (iii) holds from (4). Condition (i) holds by definition of  $\langle M, g \rangle$  and *m* level-k implementing *F*.

### **Proof of Proposition 5:**

Fix  $\bar{k} \in \mathbb{N}_+$ . Choose some  $\bar{\theta}_i \in \Theta_i$  for each  $i \in I$ . Consider the following mechanism where the message space for agent *i* is equal to  $M_i = \Theta_i \times \{0, 1, \dots, \bar{k}\} \times [-1, 1]$  and consists of a report of her payoff type  $\theta_i \in \Theta_i$ , level  $k_i \in \{0, 1, \dots, \bar{k}\}$ , and a real number,  $z_i \in [-1, 1]$ . Let the indicator function  $I_i : [-1, 1]^n \to \{0, 1\}$  be defined as follows:

$$I_i(z) = \begin{cases} 1 & \text{if } z_i = m_j z_j \text{ for some } m_j \in \mathbb{Z} \ \forall j \in I \\ 0 & \text{otherwise} \end{cases}$$

Define  $\tilde{\theta}_i : M \to \Theta_i$  and  $\tilde{k}_i : M \to \{0, 1, \dots, \bar{k}\}$  in the following way. For a given message profile  $m = (\theta, k, z)$ , if  $I_i(z) = 1$  then  $\tilde{\theta}_i(m) = \theta_i$  and  $\tilde{k}_i(m) = k_i$ ; otherwise the planner sets  $\tilde{\theta}_i(m) = \bar{\theta}_i$  and  $\tilde{k}_i(m) = 0$ . The planner then assigns outcomes based on the reports  $\tilde{\theta} \times \tilde{k}$  according to the function  $\tilde{g} : \Theta \times \{0, \dots, \bar{k}\}^n \to \Delta(Y)$  defined by

$$\tilde{g}(\theta \times \hat{k}) = \begin{cases} f^i(\theta) & \text{if } \hat{k}_j = \hat{k}_i - 1 \text{ for all } j \neq i \in I \\ \overline{f}(\theta) & \text{otherwise} \end{cases}$$

Consider an agent *i* with payoff  $\theta_i$  and level 1. Let  $\theta_j \in \Theta_j$  where  $j \neq i$ . She believes that agent *j* level 0 types are playing an atomless strategy  $\alpha_j$ . Thus, the probability that the realized value of  $z_j$  is equal to  $mz_i$  is zero, for any  $z_i \in [-1, 1]$ . Hence, our level 1 agent believes that the planner will almost surely use a payoff type,  $\bar{\theta}_j$ , for agent *j* and a level report equal to 0. Further, if our level 1 agent sends a non-zero report,  $z_i \neq 0$ , she will expect the planner to disregard her payoff type report and use payoff type,  $\bar{\theta}_i$ , with probability 1. However, if our level 1 agents sends a zero report,  $z_i = 0$ , she will expect the planner to use her payoff type as reported.

Thus, if she sends the message  $(\theta'_i, 1, 0)$ , she will expect to receive the outcome  $f^i(\theta'_i, \bar{\theta}_j)$ . If she sends the message  $(\theta'_i, k_i, 0)$ ,  $k_i \neq 1$ , she will expect to receive the outcome  $\bar{f}(\theta'_i, \bar{\theta}_j)$ . And, if she sends the message  $(\theta', k_i, z_i)$  with  $z_i \neq 0$ , then she will expect to receive the outcome  $\bar{f}(\bar{\theta}_i, \bar{\theta}_j)$ . By condition (ii) we have that

$$u_i(f^i(\theta_i, \overline{\theta}_j), \theta_i) \ge u_i(f^i(\theta'_i, \overline{\theta}_j), \theta_i)$$

for all  $\theta'_i \in \Theta_i$ .

And, by condition (iii) we have that

$$u_i(f^i(\theta_i, \bar{\theta}_j), \theta_i) \ge u_i(\bar{f}(\theta'_i, \bar{\theta}_j), \theta_i)$$

for all  $\theta'_i \in \Theta_i$ .

Thus, for any agent *i* with payoff type  $\theta_i$  and level 1, reporting  $(\theta_i, 1, 0)$  is a best response.

We now prove that an agent with payoff type  $\theta_i$  and level k will send the message  $(\theta_i, k, 0)$  by induction on the following statement: Let  $k \ge 1$  and assume that if for all  $l \in \{1, \ldots, k-1\}, \theta_j \in \Theta_j$ , and  $j \in I$  an agent j with payoff type  $\theta_j$  and level l will report  $(\theta_j, l, 0)$ , then an agent i with payoff type  $\theta_i$  and level k will report  $(\theta_i, k, 0)$ .

The result is true for k = 1 by the above argument. Now, consider an agent  $(\theta_i, k)$  where  $k \in \{2, \ldots, \bar{k}\}$ . She expects other agents that have strictly positive levels to always send reports  $z_j = 0$ . Thus she expects that the social planner will always take their payoff and level reports as given.

Let  $\theta_j \in \Theta_j$ .

Thus, if she sends the message  $(\theta'_i, k, 0)$  she will expect to receive the outcome  $f^i(\theta'_i, \theta_j)$ . If she sends the message  $(\theta'_i, l_i, 0), l_i \neq k$ , she will expect to receive the outcome  $\bar{f}(\theta'_i, \theta_j)$ . And, if she sends the message  $(\theta'_i, l_i, z_i)$  with  $z_i \neq 0$ , then she will expect to receive the outcome  $\bar{f}(\bar{\theta}_i, \theta_j)$ . By condition (ii) we have that

$$u_i(f^i(\theta), \theta_i) \ge u_i(f^i(\theta'_i, \theta_j), \theta_i)$$

for all  $\theta'_i \in \Theta_i$ .

By condition (iii) we have that

$$u_i(f^i(\theta), \theta_i) \ge u_i(\bar{f}(\theta'_i, \theta_j), \theta_i)$$

for all  $\theta'_i \in \Theta_i$ .

Thus, for agent *i* with payoff type  $\theta_i$  and level *k*, reporting  $(\theta_i, k, 0)$  is a best response.

Therefore, if we define  $m_i(\theta_i, k_i) = (\theta_i, k_i, 0)$  for all  $\theta_i \in \Theta_i$  with  $k_i \in \{1, \ldots, \bar{k}\}$ , then *m* is an expost level-k solution and *m* achieves *F* by condition (*i*). Therefore, *F* is expost level-k implementable.

## **Proof of Proposition 6**

Define

$$f^{b}(v,c) = \begin{cases} c & \text{if } c \leq v \\ \emptyset & \text{otherwise} \end{cases}$$

and

$$f^{s}(v,c) = \begin{cases} v & \text{if } c \leq v \\ \emptyset & \text{otherwise} \end{cases}$$

and

$$\bar{f}(v,c) = \begin{cases} \frac{v+c}{2} & \text{if } c \le v \\ \emptyset & \text{otherwise} \end{cases}$$

First, it is easy to see that condition (i) in Proposition 8 holds for  $f^b$ ,  $f^s$ , and  $\bar{f}$  as they assign the outcome  $\emptyset$  only when v < c.

Now consider the utility of the buyer with value v when the seller reports cost c. Consider first the comparison of the outcomes  $f^b(v, c)$  to outcomes  $f^b(v', c)$  for some value  $v' \in V$ . There are two cases to consider:

(i)v < c: The utility of the buyer is 0 when reporting v and reporting any other value v' either has no effect (if v' < c) or achieves trade (if  $v' \ge c$ ) with a utility of  $v - c \le 0$ . (ii) $c \leq v$ : The utility of buyer is  $v - c \geq 0$  when reporting v and reporting any other value v' either has no effect (if v' > c) or achieves outcome  $\emptyset$  and utility 0.

Next, consider the comparison of the outcomes  $f^b(v, c)$  to outcomes  $f^s(v', c)$  for some value  $v' \in V$ . There are two cases to consider:

- (i)v < c: The utility of the buyer under  $f^b$  is 0 when reporting v and reporting any other value v' under  $f^s$  either has no effect (if v' < c) or achieves trade (if  $v' \ge c$ ) with a utility of  $v - v' \le 0$ .
- (ii) $c \leq v$ : The utility of buyer is  $v c \geq 0$  under  $f^b$  when reporting v and reporting any other value v' under  $f^s$  either achieves outcome  $\emptyset$  (if v' < c) and utility 0 or achieves trade (if  $v' \geq c$ ) with a utility of  $v - v' \leq v - c$ .

Last, consider the comparison of the outcomes  $f^b(v, c)$  to outcomes  $\bar{f}(v', c)$ for some value  $v' \in V$ . There are two cases to consider:

- (i)v < c: The utility of the buyer under  $f^b$  is 0 when reporting v and reporting any other value v' under  $\overline{f}$  either has no effect (if v' < c) or achieves trade (if  $v' \ge c$ ) with a utility of  $v \frac{v'+c}{2} \le 0$ .
- (ii) $c \leq v$ : The utility of buyer is  $v c \geq 0$  under  $f^b$  when reporting v and reporting any other value v' under  $\bar{f}$  either achieves outcome  $\emptyset$  (if v' < c) and utility 0 or achieves trade (if  $v' \geq c$ ) with a utility of  $v \frac{v'+c}{2} \leq v c$ .

Thus, the buyer has (weakly) higher utility when reporting v under  $f^b$  than reporting any other value v' under  $f^b$ ,  $f^s$ , or  $\bar{f}$  regardless of the cost of the seller, c. In other words, conditions (ii)-(iv) in Proposition 8 is satisfied for the buyer. Analogously, conditions (ii)-(iv) are satisfied for the seller.

All outcomes assigned in the mechanism are determined by  $f^b$ ,  $f^s$ , and  $\bar{f}$ , which satisfy ex post individual rationality whenever types are truthfully reporting their payoff type. Budget balance is satisfied automatically given the specification of the environment. The environment is one of private values, thus the result follows from Proposition 8.

# Appendix B Implementation when n=2

#### Level-k implementation

**Proposition 7.** (Sufficient Conditions n=2) Let F be a social choice rule and let the environment be one of independent private values and let n = 2. Let  $\mathcal{B}$ be a Bayesian type space and let  $\mathcal{L} = \langle \mathcal{B}, (T_i, \mu_i)_{i=1,...,n}; \bar{k} \rangle$  be a ( $\mathcal{B}$ -based) levelk type space. If there exists a function  $f^i : \Theta \to \Delta(Y)$  for each  $i \in I$  and a function  $\bar{f} : \Theta \to \Delta(Y)$  such that the conditions (i)-(iii) hold in Proposition 1 and condition (iv) holds below, then F is level-k implementable under atomless anchors.

$$\begin{array}{ll} \textbf{(iv)} & \sum\limits_{\substack{\theta_{-i} \in \Theta_{-i} \\ \Theta_i, \ \theta \in \Theta, \ i, \ j \neq i \in I \end{array}} \rho(\theta_{-i} | \theta_i) u_i(f^i(\theta), \theta_i) \ \geq & \sum\limits_{\substack{\theta_{-i} \in \Theta_{-i} \\ \Theta_{-i} \in \Theta_{-i} \end{array}} \rho(\theta_{-i} | \theta_i) u_i(f^j(\theta'_i, \theta_{-i}), \theta_i) \ \forall \theta'_i \ \in \end{array}$$

**PROOF:** 

Consider the following mechanism where the message space for agent i is equal to  $M_i = \Theta_i \times \{0, 1, \dots, \bar{k}\} \times [-1, 1]$  and consists of a report of a payoff type  $\theta_i \in \Theta_i$ , level  $k_i \in \{0, 1, \dots, \bar{k}\}$ , and a real number,  $z_i \in [-1, 1]$ .

Let the indicator function  $I_i: [-1,1]^n \to \{0,1\}$  be defined as follows:

$$I_i(z) = \begin{cases} 1 & \text{if } z_i = m_j z_j \text{ for some } m_j \in \mathbb{Z}, \ \forall j \in I \\ 0 & \text{otherwise} \end{cases}$$

Define  $\tilde{\theta}_i : M \to \Theta_i$  and  $\tilde{k}_i : M \to \{0, 1, \dots, \bar{k}\}$  in the following way. For a given message profile  $m = (\theta, k, z)$ , if  $I_i(z) = 1$  the planner takes the reports as given and sets  $\tilde{\theta}_i(m) = \theta_i$  and  $\tilde{k}_i(m) = k_i$ ; otherwise the planner sets  $\tilde{\theta}_i(m)$  to some randomly chosen  $\Theta_i$  according to the prior  $p_i$ and  $\tilde{k}_i(m) = 0$ . The planner then assigns outcomes based on the reports  $\tilde{\theta} \times \tilde{k}$  according to  $\tilde{g} : \Theta \times \{0, \dots, \bar{k}\}^n \to \Delta(Y)$  defined by

$$\tilde{g}(\theta \times \hat{k}) = \begin{cases} f^i(\theta) & \text{if } \hat{k}_j = \hat{k}_i - 1 \text{ for all } j \neq i \in I \\ \bar{f}(\theta) & \text{otherwise} \end{cases}$$

Consider an agent *i* with payoff  $\theta_i$  and level 1. She believes that all  $j \neq i$  are level-0 agents playing an atomless strategy  $\alpha_j$ . Thus, the probability that the realized value of  $z_j$  is equal to  $mz_i$  is zero, for any  $z_i \in [-1, 1]$ . Hence, our level 1 agent believes that the planner will almost surely use a payoff type for agent *j* that is picked at random according to the prior  $p_j$  and a level report equal to 0. Further, if our level 1 agent sends a non-zero report,  $z_i \neq 0$ , she will expect the planner to disregard her payoff type report and choose randomly according to  $p_i$  with probability 1. However, if our level 1 agents sends a zero report,  $z_i = 0$ , she will expect the planner to use her payoff type as reported.

Thus, if she sends the message  $(\theta'_i, k_i, z_i)$  with  $z_i = 0$  and  $k_i = 1$ , she will expect to receive the following lottery over outcomes

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot f^{i}(\theta'_{i},\theta_{-i}).$$

If she sends the message  $(\theta'_i, k_i, 0)$  with  $k_i \neq k$ , she will expect to receive one of the two following lottery over outcomes

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot f^{j}(\theta_{i}',\theta_{-i})$$

for  $j \neq i$  or

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot\bar{f}(\theta_{i}',\theta_{-i})$$

And, if she sends the message  $(\theta'_i, k_i, z_i)$  with  $z_i \neq 0$ , she will expect to receive the following lottery over outcomes

$$\sum_{\theta \in \Theta} \rho(\theta) \cdot \bar{f}(\theta)$$

By condition (ii) and (iv) we have that

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(f^i(\theta),\theta_i) \ge \sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(f^j(\theta'_i,\theta_{-i}),\theta_i)$$

for all  $\theta'_i \in \Theta_i$ .

By condition (iii) we have that

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(f^i(\theta),\theta_i) \ge \sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(\bar{f}(\theta'_i,\theta_{-i}),\theta_i)$$

for all  $\theta'_i \in \Theta_i$ .

It must then also be true that

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot u_i(f^i(\theta),\theta_i)\geq \sum_{\theta\in\Theta}\rho(\theta)\cdot u_i(\bar{f}(\theta),\theta_i).$$

Thus, for agent *i* with payoff type  $\theta_i$  and level 1, reporting  $(\theta_i, 1, 0)$  is a best response.

We now prove that an agent with payoff type  $\theta_i$  and level k will send the message  $(\theta_i, k_i, 0)$  by induction on the following statement: Let  $k \ge 1$  and assume that if for all  $l \in \{1, \ldots, k-1\}, \theta_j \in \Theta_j$ , and  $j \in I$  an agent j with payoff type  $\theta_j$  and level l will report  $(\theta_j, l, 0)$ , then an agent i with payoff type  $\theta_i$  and level k will report  $(\theta_i, k, 0)$ .

The result is true for k = 1 by the above argument. Now consider an agent i with payoff type  $\theta_i$  and level  $k \in \{2, \ldots, \bar{k}\}$ . She expects that all other agents will be sending reports  $z_j = 0, k_j = k - 1$ , and truthfully reporting their payoff type. Thus, she expects that the social planner will always take their payoff and level reports as given.

Thus, if she sends the message  $(\theta'_i, k, 0)$  she will expect to receive the following lottery over outcomes

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})\cdot f^{i}(\theta_{i}',\theta_{-i}).$$

If she sends the message  $(\theta'_i, k_i, 0)$  with  $k_i \neq k$ , she will expect to receive one of the two following lottery over outcomes

$$\sum_{\boldsymbol{\theta}_{-i}\in \Theta_{-i}} \rho(\boldsymbol{\theta}_{-i}) \cdot f^{j}(\boldsymbol{\theta}_{i}^{\prime},\boldsymbol{\theta}_{-i})$$

for  $j \neq i$  or

$$\sum_{\theta_{-i}\in\Theta_{-i}} \rho(\theta_{-i}) \cdot \bar{f}(\theta_{i}',\theta_{-i})$$

And, if she sends the message  $(\theta'_i, k_i, z_i)$  with  $z_i \neq 0$ , she will expect to receive the following lottery over outcomes

$$\sum_{\theta \in \Theta} \rho(\theta) \cdot \bar{f}(\theta)$$

By the same argument above, conditions (ii)-(iv) imply that our agent will send the report  $(\theta_i, k, 0)$ .

Therefore, if we define  $m_i(\theta_i, k_i) = (\theta_i, k_i, 0)$  for all  $\theta_i \in \Theta_i, k_i \in \{1, \dots, \bar{k}\}$ and  $i \in I$ , then m is a level-k solution and achieves F by condition (i). Hence, F is level-k implementable.

#### Ex post level-k implementation

**Proposition 8.** (Ex Post Sufficient Conditions n=2) Let F be a social choice rule and let the environment be one of private values and let n = 2. If there exists a function  $f^i: \Theta \to \triangle(Y)$  for each  $i \in I$  and a function  $\overline{f}: \Theta \to \triangle(Y)$ such that the conditions (i)-(iii) hold in Proposition 4 and condition (iv) holds below, then F is ex post level-k implementable under atomless anchors.

(iv)  $u_i(f^i(\theta), \theta) \ge u_i(f^j(\theta'_i, \theta_{-i}), \theta) \ \forall \theta'_i \in \Theta_i, \ \theta \in \Theta, \ i, j \neq i \in I$ 

Proof:

Fix  $\bar{k} \in \mathbb{N}_+$ .

Choose some  $\bar{\theta}_i \in \Theta_i$  for each  $i \in I$ .

Consider the following mechanism where the message space for agent *i* is equal to  $M_i = \Theta_i \times \{0, 1, \dots, \bar{k}\} \times [-1, 1]$  and consists of a report of her payoff type  $\theta_i \in \Theta_i$ , level  $k_i \in \{0, 1, \dots, \bar{k}\}$ , and a real number,  $z_i \in [-1, 1]$ . Let the indicator function  $I_i : [-1, 1]^n \to \{0, 1\}$  be defined as follows:

$$I_i(z) = \begin{cases} 1 & \text{if } z_i = m_j z_j \text{ for some } m_j \in \mathbb{Z} \ \forall j \in I \\ 0 & \text{otherwise} \end{cases}$$

Define  $\tilde{\theta}_i : M \to \Theta_i$  and  $\tilde{k}_i : M \to \{0, 1, \dots, \bar{k}\}$  in the following way. For a given message profile  $m = (\theta, k, z)$ , if  $I_i(z) = 1$  then  $\tilde{\theta}_i(m) = \theta_i$  and  $\tilde{k}_i(m) = k_i$ ; otherwise the planner sets  $\tilde{\theta}_i(m) = \bar{\theta}_i$  and  $\tilde{k}_i(m) = 0$ . The planner then assigns outcomes based on the reports  $\tilde{\theta} \times \tilde{k}$  according to the function  $\tilde{g} : \Theta \times \{0, \dots, \bar{k}\}^n \to \Delta(Y)$  defined by

$$\tilde{g}(\theta \times \hat{k}) = \begin{cases} f^i(\theta) & \text{if } \hat{k}_j = \hat{k}_i - 1 \text{ for all } j \neq i \in I \\ \bar{f}(\theta) & \text{otherwise} \end{cases}$$

Consider an agent *i* with payoff  $\theta_i$  and level 1. Let  $\theta_j \in \Theta_j$  where  $j \neq i$ . She believes that agent *j* level 0 types are playing an atomless strategy  $\alpha_j$ . Thus, the probability that the realized value of  $z_j$  is equal to  $mz_i$  is zero, for any  $z_i \in [-1, 1]$ . Hence, our level 1 agent believes that the planner will almost surely use a payoff type,  $\bar{\theta}_j$  for agent *j* and a level report equal to 0. Further, if our level 1 agent sends a non-zero report,  $z_i \neq 0$ , she will expect the planner to disregard her payoff type report and choose the payoff  $\bar{\theta}_i$  with probability 1. However, if our level 1 agent sends a zero report,  $z_i = 0$ , she will expect the planner to use her payoff type as reported.

Thus, if she sends the message  $(\theta'_i, 1, 0)$ , she will expect to receive the outcome  $f^i(\theta'_i, \bar{\theta}_j)$ . If she sends the message  $(\theta'_i, k_i, 0)$ ,  $k_i \neq 1$ , she will expect to receive either the outcome  $\bar{f}(\theta'_i, \bar{\theta}_j)$  or  $f^j(\theta'_i, \bar{\theta}_j)$  for some  $j \in I$ . And, if she sends the message  $(\theta', k_i, z_i)$  with  $z_i \neq 0$ , then she will expect to receive the outcome  $\bar{f}(\bar{\theta}_i, \bar{\theta}_j)$ .

By condition (ii) and (iv) we have that

$$u_i(f^i(\theta_i, \bar{\theta}_j), \theta_i) \ge u_i(f^j(\theta', \bar{\theta}_j), \theta_i)$$

for all  $\theta' \in \Theta_i$  and all  $j \in I$ .

And, by condition (iii) we have that

$$u_i(f^i(\theta_i, \bar{\theta}_j), \theta_i) \ge u_i(\bar{f}(\theta', \bar{\theta}_j), \theta_i)$$

for all  $\theta' \in \Theta_i$ .

Thus, for any agent *i* with payoff type  $\theta_i$  and level 1, reporting  $(\theta_i, 1, 0)$  is a best response.

We now prove that an agent with payoff type  $\theta_i$  and level k will send the message  $(\theta_i, k, 0)$  by induction on the following statement: Let  $k \ge 1$  and assume that if for all  $l \in \{1, \ldots, k-1\}, \theta_j \in \Theta_j$ , and  $j \ne i$  an agent j with payoff type  $\theta_j$  and level l will report  $(\theta_j, l, 0)$ , then an agent i with payoff type  $\theta_i$  and level k will report  $(\theta_i, k, 0)$ .

The result is true for k = 1 by the above argument. Now, consider an agent  $(\theta_i, k)$  where  $k \in \{2, \ldots, \bar{k}\}$ . She expects other agents that have strictly positive levels to always send reports  $z_j = 0$ . Thus she expects that the social planner will always take their payoff and level reports as given.

Let  $\theta_j \in \Theta_j$ .

Thus, if she sends the message  $(\theta'_i, k, 0)$  she will expect to receive the outcome  $f^i(\theta'_i, \theta_j)$ . If she sends the message  $(\theta'_i, l_i, 0), l_i \neq k$ , she will expect to receive the outcome  $f^j(\theta'_i, \theta_j)$  for  $j \neq i$  or  $\bar{f}(\theta'_i, \theta_j)$ . And, if she sends the message  $(\theta'_i, l_i, z_i)$  with  $z_i \neq 0$ , then she will expect to receive the outcome  $\bar{f}(\bar{\theta}_i, \theta_j)$ .

By the same argument above, conditions (ii)-(iv) then imply that our agent will send the report  $(\theta_i, k, 0)$ .

Therefore, if we define  $m_i(\theta_i, k_i) = (\theta_i, k_i), 0$  for all  $\theta_i \in \Theta_i$  with  $k_i \in \{1, \ldots, \bar{k}\}$ , then *m* is an expost level-k solution and *m* achieves *F* by condition (*i*). Therefore, *F* is expost level-k implementable.

# Appendix C LDoR implementation

This appendix addresses the question of whether we can design a mechanism that is robust to relaxing one of the belief assumptions that is present in the level-k type space. Specifically, under the level-k model, if an agent is level k, she believes that others have levels exactly equal to (k - 1). In general, we might allow an agent with level k to hold beliefs over all lower levels. As long as a level k type only puts weight on lower levels, the spirit of limited depth of reasoning is maintained with each type being able to calculate her optimal action recursively, in a finite number of steps.

In this section, we generalize our type space and solution concept to the limited depth of reasoning (LDoR) concept to relax beliefs about the depths of reasoning of others and consider a form of robust implementation where we ask whether there exists a mechanism that will implement a social choice function for any LDoR type space - we call this LDoR implementation.

The following definition of a limited depth of reasoning (LDoR) type space generalizes the level-k type space. The LDoR type space allows an agent to hold any arbitrary beliefs over lower levels of others. This approach is based on Strzalecki (2014) who develops the framework for games of complete information. We expand the framework here to allow for incomplete information.

Definition 11.  $\mathcal{B}$ -based limited depth of reasoning type space (LDoR type space) is a type space  $\mathcal{L}^{LDoR} = \langle \mathcal{B}, (T_i, k_i, \theta_i, b_i)_{i=1,...,n}, \bar{k} \rangle$ , such that  $\mathcal{B}$ is a Bayesian type space  $\mathcal{B} = \langle \Theta; \rho \rangle$ ,  $T_i$  is a finite set for all  $i \in I$ ,  $k_i : T_i \to \{0, \ldots, \bar{k}\}, \theta_i : T_i \to \Theta_i$ , and  $b_i : T_i \to \Delta(T_{-i})$  such that for all  $t_i \in T_i$ :

$$b_i(t_i)(\{t_{-i} \in T_{-i} | \text{ such that } k_j(t_j) < k_i(t_i) \ \forall j \neq i \in I\}) = 1$$

and for all  $l \in \{0, ..., k_i(t_i) - 1\}$  with  $b_i(t_i)(\{t_{-i} \in T_{-i} | k_j(t_j) = l \text{ for some } j \neq i \in I\}) > 0$  then

$$b_i(t_i)(\{t_{-i} \in T_{-i} | \theta_{-i}(t_{-i}) = \theta_{-i} \text{ and } k_j(t_j) = l \text{ for all } j \neq i \in I\}) = \rho(\theta_{-i}|\theta_i(t_i))$$

for all  $\theta_{-i} \in \Theta_{-i}$ .

As in the level-k type space, an agent's type,  $t_i$ , represents both her payoff

type,  $\theta_i$ , and her level,  $k_i$ . We abuse notation here and also let  $\theta_i$  and  $k_i$  be functions that map types to payoff types and levels, respectively. Thus, for a type,  $t_i$ , her payoff type is given by  $\theta_i(t_i)$  and her level is given by  $k_i(t_i)$ . Unlike the level-k model, where an agent's payoff type and level completely determine her beliefs about others, in the LDoR model, agents with the same payoff type and level may have different beliefs about the levels of others. The belief function,  $b_i$ , specifies, for each type, her beliefs about the types of others. We impose two belief restrictions. The first restriction requires that each type only puts positive weight on types that have strictly lower levels. This captures the core assumption of the limited depth of reasoning literature and ensures that agents can calculate their optimal actions in a recursive fashion with a finite number of steps given the behavior of level 0 types. The second restriction maintains the common prior assumption: if a type puts positive weight on a level l then her beliefs about the payoff types of those with level l must be consistent with the common prior. For this reason, we refer to the LDoR type space formally, as a Bayesian-based LDoR type space.

Given the definition of an LDoR type space, we can define the analogous solution and implementation concepts: the LDoR solution and LDoR implementation.

**Definition 12.** For a given game defined by a mechanism  $\langle M, f \rangle$  and type space  $\mathcal{L}^{LDoR} = \langle \mathcal{B}, (T_i, k_i, \theta_i, b_i)_{i=I}, \bar{k} \rangle$ , a strategy profile  $s = s_1 \times \cdots \times s_n$ , with  $s_i : T_i \to \Delta(M_i)$  for all  $i \in I$ , is the **LDoR solution under anchor**  $\alpha$  if and only if:

(i) 
$$s_i(t_i) \sim \alpha(t_i)$$
 for all  $t_i \in \{t \in T_i | k_i(t) = 0\}, i \in I$   
(ii)  $\sum_{\substack{t_{-i} \in T_{-i} \\ \forall m'_i \in M_i, t_i \in T_i \text{ with } k_i(t_i) \geq 1, i \in I.} b_i(t_{-i} | t_i) u_i(f(m'_i, s_{-i}(t_{-i})), \theta(t))$ 

The LDoR solution is similar to the level-k solution. It specifies that all level 0 types play according to anchor  $\alpha$ . And, it specifies that all types with levels at least 1 play a best response given their beliefs about the types of others and the actions of those types under s. To define LDoR implementation however, we go a step further here and apply a further robustness criterion - that the social choice rule be implementable under *any* LDoR type space.

**Definition 13.** Fix a  $\bar{k}$  and a Bayesian type space  $\mathcal{B}$ . A social rule F is **LDoR implementable** if there exists a mechanism  $\langle M, f \rangle$  such that for any  $\mathcal{B}$ -based LDoR type space  $\mathcal{L}^{LDoR} = \langle \mathcal{B}, (T_i, k_i, \theta_i, b_i)_{i=1,...,n}, \bar{k} \rangle$ , there exists a message profile  $m_i : T_i \to \Delta(M_i)$  for all  $i \in I$ , such that  $m = m_1 \times \cdots \times m_n$ is an LDoR solution under anchor  $\alpha$  and  $f(m(t)) \in F(\theta(t))$  for all  $t \in \{t \in T | k_i(t_i) \geq 1, \forall i \in I\}$ .

Proposition 9 gives the sufficient conditions for LDoR implementation. The sufficient conditions are the same as those required for level-k implementation when n = 2. Thus includes an additional requirement relative to Proposition 2. This additional requirement strengthens the original condition (ii) in Proposition 1 to require that agents prefer the truth telling outcome under  $f^i$  to other outcomes under  $f^i$  and also outcomes under  $f^j$  for any other agent j as well. These results are stated below.

**Proposition 9.** (LDoR Sufficient Conditions) Let F be a social choice rule. Let the environment be one of independent private values. If conditions (i)-(iii) hold in Proposition 1 plus condition (iv) below then F is LDoR implementable under atomless anchors

$$(iv) \sum_{\substack{\theta_{-i} \in \Theta_{-i} \\ \Theta_i, \ \theta \in \Theta, \ i, \ j \neq i \in I}} \rho(\theta_{-i}|\theta_i) u_i(f^i(\theta), \theta_i) \geq \sum_{\substack{\theta_{-i} \in \Theta_{-i} \\ \Theta_{-i} \in \Theta_{-i}}} \rho(\theta_{-i}|\theta_i) u_i(f^j(\theta'_i, \theta_{-i}), \theta_i) \ \forall \theta'_i \in \Theta_{-i}$$

Proof:

Consider the following mechanism where the message space for agent i is equal to  $M_i = \Theta_i \times \{0, 1, \dots, \bar{k}\} \times [-1, 1]$  and consists of a report of her payoff type  $\theta_i \in \Theta_i$ , level  $k_i \in \{0, 1, \dots, \bar{k}\}$ , and a real number,  $z_i \in [-1, 1]$ . Let the indicator function  $I_i : [-1, 1]^n \to \{0, 1\}$  be defined as follows:

$$I_i(z) = \begin{cases} 1 & \text{if } z_i = m_j z_j \text{ for some } m_j \in \mathbb{Z}, \forall j \in I \\ 0 & \text{otherwise} \end{cases}$$

Define  $\tilde{\theta}_i : M \to \Theta_i$  and  $\tilde{k}_i : M \to \{0, 1, \dots, \bar{k}\}$  in the following way. For a given message profile  $m = (\theta, k, z)$ , if  $I_i(z) = 1$  the planner takes the reports as given and sets  $\tilde{\theta}_i(m) = \theta_i$  and  $\tilde{k}_i(m) = k_i$ ; otherwise the planner sets  $\tilde{\theta}_i(m)$  to some randomly chosen  $\Theta_i$  according to the prior  $\rho_i$ and  $\tilde{k}_i(m) = 0$ . The planner then assigns outcomes based on the reports  $\tilde{\theta} \times \tilde{k}$  according to the function  $\tilde{g}: \Theta \times \{0, \ldots, \bar{k}\}^n \to \Delta(Y)$  defined by

$$\tilde{g}(\theta \times \hat{k}) = \begin{cases} f^{i}(\theta) & \text{if } k_{i} > \max\{k_{1}, \dots, k_{i-1}, k_{i+1}, \dots, k_{n}\} \\ \bar{f}(\theta) & \text{otherwise} \end{cases}$$

Let  $\mathcal{L}^{LDoR} = \langle \mathcal{B}, (T_i, k_i, \theta_i, b_i)_{i \in I}, \bar{k} \rangle$  be an LDoR type space.

Consider an agent  $t_i \in T_i$  with  $\theta_i(t_i) = \theta_i$  and level  $k_i(t_i) = 1$ . The beliefs and incentives for level 1 agents are unchanged relative to the level-k type space and mechanism in Proposition 2. Thus, for any agent *i* with payoff type  $\theta_i$  and level 1, reporting  $(\theta_i, 1, 0)$  is a best response.

We now prove that an agent with payoff type  $\theta_i$  and level k will send the message  $(\theta_i, k, 0)$  by induction on the following statement: Let  $k \ge 1$  and assume that if for all  $l \in \{1, \ldots, k-1\}, \theta_j \in \Theta_j$ , and  $j \in I$  an agent j with payoff type  $\theta_j$  and level l will report  $(\theta_j, l, 0)$ , then an agent i with payoff type  $\theta_i$  and level k will report  $(\theta_i, k, 0)$ .

The result is true for k = 1 by the above argument. Now, consider an agent  $t_i$  with payoff type  $\theta_i(t_i) = \theta_i$  and level  $k_i(t_i) = k \in \{2, \ldots, \bar{k}\}$ . She expects other agents that have strictly positive levels to always send reports  $z_j = 0$ . Thus she expects that the social planner will always take their payoff and level reports as given. For agents with levels of 0, she expects the planner to almost surely use a payoff type randomly chosen according to  $\rho_j$  and level report of 0 for those agents.

Thus, if she sends the message  $(\theta'_i, k, 0)$  she will expect to receive the following lottery over outcomes

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho_i(\theta_{-i})\cdot f^i(\theta'_i,\theta_{-i})$$

If she sends the message  $(\theta'_i, k_i, 0), k_i \neq k$ , she will expect to receive the following lottery over outcomes

$$\sum_{j \in I} \beta^j \cdot \sum_{\theta_{-i} \in \Theta_{-i}} \rho_i(\theta_{-i}) \cdot f^j(\theta'_i, \theta_{-i}) + \left(1 - \sum_{j \in I} \beta^j\right) \sum_{\theta_{-i} \in \Theta_{-i}} \rho_i(\theta_{-i}) \cdot \bar{f}(\theta'_i, \theta_{-i})$$

for some  $\beta = (\beta^j)_{j \in I}$  such that  $\beta^j \in [0, 1]$  and  $\sum_{j \in I} \beta^j \leq 1$ .

And, if she sends the message  $(\theta_i, k_i, z_i)$  with  $z_i \neq 0$ , then she will expect to receive the following lottery over outcomes

$$\sum_{j \in I} \beta^j \cdot \sum_{\theta \in \Theta} \rho(\theta) \cdot f^j(\theta) + \left(1 - \sum_{j \in I} \beta^j\right) \sum_{\theta \in \Theta} \rho(\theta) \cdot \bar{f}(\theta)$$

for some for  $\beta = (\beta^j)_{j \in I}$  such that  $\beta^j \in [0,1]$  and  $\sum_{j \in I} \beta^j \leq 1$ .

By condition (ii) and (iv) we have that

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})u_i(f^i(\theta),\theta_i)\geq \sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})u_i(f^j(\theta'_i,\theta_{-i}),\theta_i)$$

for all  $\theta'_i \in \Theta_i, j \in I$ .

By condition (iii) we have that

$$\sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})u_i(f^i(\theta),\theta_i) \ge \sum_{\theta_{-i}\in\Theta_{-i}}\rho(\theta_{-i})u_i(\bar{f}(\theta'_i,\theta_{-i}),\theta_i)$$

for all  $\theta'_i \in \Theta_i$ .

It must also then also be true that

$$\sum_{\substack{\theta_{-i}\in\Theta_{-i}}} \rho_i(\theta_{-i}) \cdot f^i(\theta_i, \theta_{-i})$$

$$\geq \sum_{j\in I} \beta^j \cdot \sum_{\substack{\theta_{-i}\in\Theta_{-i}}} \rho_i(\theta_{-i}) \cdot f^j(\theta'_i, \theta_{-i}) + \left(1 - \sum_{j\in I} \beta^j\right) \cdot \sum_{\substack{\theta_{-i}\in\Theta_{-i}}} \rho_i(\theta_{-i}) \cdot \bar{f}(\theta'_i, \theta_{-i})$$

$$\geq \sum_{j\in I} \beta^j \cdot \sum_{\substack{\theta\in\Theta}} \rho(\theta) \cdot f^j(\theta) + \left(1 - \sum_{j\in I} \beta^j\right) \cdot \sum_{\substack{\theta\in\Theta}} \rho(\theta) \cdot \bar{f}(\theta)$$

for any  $\beta = (\beta^j)_{j \in I}$  such that  $\beta^j \in [0, 1]$  and  $\sum_{j \in I} \beta^j \leq 1$  and for all  $\theta'_i \in \Theta_i$ . Thus, for agent *i* with payoff type  $\theta_i$  and level *k*, reporting  $(\theta_i, k, 0)$  is a best response.

Therefore, if we define  $m_i(t_i) = (\theta_i(t_i), k_i(t_i), 0)$  for all  $t_i \in T_i$  with  $k_i(t_i) \in \{1, \ldots, \bar{k}\}$ , then *m* is a LDoR solution and our given mechanism implements *F*.

Since the above holds for an arbitrary LDoR type space, F is LDoR implementable.