

DISCUSSION PAPER NO. 5

A PARTIAL CHARACTERIZATION OF
ORGANIZATIONS AND ENVIRONMENTS WHICH
ARE CONSISTENT WITH INCENTIVE COMPATIBILITY

by

John O. Ledyard ^{*/}

Revised - July 1976

*/

I would like to acknowledge the very great contribution John Roberts made to the production of this paper. Of course any errors are mine.

This research was partly supported by the National Science Foundation Grants (GS 31346X) and (SOC 74-04076).

1. INTRODUCTION

Part of the structure of any economic organization are prescribed rules or norms of behavior for the participants in the organization. If these rules are actually followed, then some specified outcomes result from the operation of the organization. Presumably, these outcomes have some desirable properties, such as being Pareto-optimal or profit maximizing. Much of standard economic theory has been concerned with establishing the properties of these outcomes for various organizations under the assumption that the rules are in fact obeyed. However, it may well be the case that the participants in an organization find it to their individual advantage to depart from the prescribed behavior, in that the outcome resulting from this departure may be more desirable from their individual point of view. Recently a theory of incentives has begun to develop by building on a structure supplied by Hurwicz (1972). One of the more interesting results in that paper is an impossibility theorem. No organization which selects Pareto-optimal allocations and which allows a "no-trade" option to the participants can be both incentive-compatible and decentralized, even when no public goods are present. A particular corollary of this theorem is that the competitive process is not necessarily incentive compatible. He establishes this result by considering a two person, pure exchange environment and showing that incentive compatibility cannot obtain there since one agent can always benefit from misrepresentation if the no-trade action is not Pareto-optimal. A fortiori, in any class of environments which includes this central example incentive compatibility could not be realized.

The question we begin to address in this paper follows directly from Hurwicz's theorem. That is, if organizations which select Pareto-optima are not necessarily incentive compatible in all environments, which specific organizations are incentive compatible in which specific environments? By restricting our concern to organizations which select core allocations (a subset of those considered by Hurwicz), we are able to provide a partial answer to this question by providing necessary and sufficient conditions on pairs of organizations and environments such that incentive compatibility obtains. A plausible conjecture is that a core-selecting organization is incentive compatible in a particular environment if and only if the core of that environment is single-valued, (i.e., all core allocations yield the same point in utility space). The basis for this conjecture is that when the core is single-valued, price taking behavior seems to be consistent with self-interest. However, the conjecture is false. The existence of a single-valued core is necessary (Corollary 2.2) but not sufficient (examples 8 and 9) for incentive compatibility. Thus, a stronger condition is needed which we introduce by considering allocations which are the unilaterally best unblocked allocation for each participant (Definition 4). The major result of this paper, contained in Theorem 3, is that, given a set of technical assumptions, a core selecting organization is incentive compatible in an environment if and only if the environment possesses an allocation which is the unilaterally best unblocked allocation for all participants.

The conclusion which follows from this result and the fact that most environments do not possess a unilaterally best unblocked allocation is

most core selecting organizations are not incentive compatible in most environments. This is discussed further in Section 4.

2. THE MODEL

The model we use to analyze the incentive problem takes as given an environment, an organization, and an enforcement structure. The concept of an organization is broad and contains many components which are irrelevant for our purposes. We will concentrate solely on the choices which a particular organizational structure makes when faced with a particular environment. Thus, we will be characterizing choice mechanisms and suppressing the organizational forms which implement the choice process. Given an environment, a choice mechanism and an enforcement structure, a preference game is derived.

a. Environments

The basic datum of our analysis is the (economic) environment. An environment is a collection of agents, a description of their feasible actions, the outcomes which can result from these actions, and the agents' preferences for these outcomes. In standard general equilibrium models, the agents are firms and consumers specified in terms of their characteristics (production sets or preferences, endowments, and ownership claims), the actions are net trades, feasible actions are joint trades which add up to zero, and the outcomes are final allocations. In other contexts the agents might be voters, government

agencies, etc. We index the members of the set of agents by i , where i belongs to an index set $I \subseteq [0, \infty)$, the non-negative reals. Each member of the environment is able to select an action from some set A^i of conceivable actions. A joint action is an assignment of actions $a^i \in A^i$ to each member of I . The technology and endowments serve to limit the feasible joint actions to a feasible action set \hat{A} contained in $\prod_{i \in I} A^i$. It should be recognized that A^i may be very complex, and that the set \hat{A} , of feasible joint actions may depend on the distribution of endowments. The actions a^i which are chosen serve to define, via a resultant function h , the final state x which will obtain where x belongs to some set $X \subseteq \prod_{i \in I} X^i$. Each $i \in I$ has a preference ordering, \succsim^i over his component X^i of the outcomes. We are also interested in the possibilities open to subsets, coalitions, of agents. Let C be the set of admissible coalitions.¹ For each $c \in C$ there is a set $S^c \subseteq \prod_{i \in c} A^i$ which describes the options for independent action open to the coalition c . We will restrict our attention to environments for which S^c is independent of the actions taken by agents not belonging to c , and for which the outcome for the members of c , which results from their joint action, is independent of the actions and outcomes of those not belonging to c . Although some externality phenomena are ruled out by these restrictions, many are still possible. (E.g. public goods can be modeled within this framework as is done in example 2 below.)

Thus, an environment e is a vector $(I, \langle A^i \rangle_{i \in I}, C, \langle S^c \rangle_{c \in C}, h, \langle \succsim^i \rangle_{i \in I})$. It is intended that the environment be a complete

description of the physically given data of the system which exist apart from any social, political, or economic organizations. Denote the set of possible environments by E .

We can illustrate the foregoing concepts with three well-known examples.

Example 1: (Decomposable exchange environments). We say that e is a decomposable exchange environment if, for all $c \in C$, $S^c = \{a^c \in \prod_{i \in C} A^i \mid \sum_{i \in C} a^i = 0\}$ where $A^i \subseteq G$ (an additive group) for all $i \in I$, and $h(a) = \langle h^1(a^1), \dots, h^N(a^N) \rangle$. Denote by E_d the set of all such environments.

E_d arises naturally from the economists' pure exchange environments with no externalities. To illustrate, let B^i be i 's admissible consumptions (usually the non-negative orthant of L dimensional Euclidean space), let w^i be his initial endowment, and let $u^i(b^i)$ be his utility function on B^i . Then, $A^i = \{a^i \mid a^i + w^i \in B^i\}$, $x^i = h^i(a^i) = a^i + w^i$, and $a^i \succsim^i a^{*i}$ iff $u^i(a^i + w^i) \geq u^i(a^{*i} + w^i)$. Thus, an action is simply a net-trade vector.

Example 2: (Public goods). Let $X = R^{M+N}$ (the $N+M$ dimensional Euclidean space) be the commodity space for M private and N public goods. Let $a^i = (d^i, b)$ where $d^i \in R^M$ and $b \in R^N$. Let w^i be i 's initial endowment of private goods, and $A^i = \{(d^i, b) \mid d^i + w^i \geq 0, b \geq 0\}$. Then $S^c = \{a^c \in \prod_{i \in C} A^i \mid (\sum_{i \in C} d^i, b) \in F\}$ where F is a production possibility set. With preferences defined on $X^i = \{d^i + w^i\}$ for each i , this corresponds to models used in the public goods literature. (See, e.g., Muench).

Example 3: (Exchange with a continuum of agents). Let $R = [0, \infty)$ and let μ be a (probability) measure on R . I is the

support of μ . C is the collection of all subsets, c , of R such that c is μ -measurable and $\mu(c) > 0$. Let $A^i \subseteq \mathbb{R}^L$ (the L -dimensional Euclidean space). $S^c = \{ \langle a^i \rangle_{i \in c} \mid \int_c a^i d\mu(i) = 0, \text{ and } a^i \in A^i \text{ for all } i \in c. \}$ An example of μ is Lebesgue measure on the interval $[0,1]$ which has been used by Aumann (1966) and others to represent the perfectly competitive situation in which each agent is relatively powerless. Another example, which we will use later, is constructed by letting μ be Lebesgue measure on $[0,2]$ and letting $\mu(\{3\}) = \mu(\{4\}) = \frac{1}{2}$. This measure is used by Shitovitz (1973) to represent a duopoly situation.

b. Choice Mechanisms

A choice mechanism, in this framework, is a created structure as opposed to the given datum of the environment. In various contexts a choice mechanism may be an economic system (in which case it is usually called an allocation mechanism) or a voting procedure. In general a choice mechanism operates to select the joint actions to be taken by the agents in the environment in which the mechanism operates. The selection process may involve considerable communication (e.g., the infinite iterations of the Walrasian Tatonnement) or none (e.g., a dictatorial social choice rule). In any case, the net result is the selection of a joint action.

Thus, a choice mechanism can be represented as a mapping (called its outcome function) from the class of environments, E , to a set of joint actions $A = \prod_{i \in [0, \infty)} A^i$. Denote by $B = \{ \beta \mid \beta: E \rightarrow A \}$ the set of all such outcome functions.

We can illustrate the foregoing with three well known examples.

Example 4: (Walrasian Tatonnement). Let $e \in E_d$ (see example 1), where G is a finite dimensional vector space. For $p \in G$, let $D^i(p, e) = \{a^i \in A^i \mid pa^i \leq 0 \text{ and } a^i + w^i \geq a^{*i} + w^i \text{ for all } a^{*i} \in A^i \text{ such that } pa^{*i} \leq 0\}$. D^i is simply i 's excess demand function. Define $\beta^w(e)$ to be $\{a \in A \mid a^i \in D^i(p, e) \text{ for all } i \in I \text{ and } \sum_{i \in I} a^i = 0, \text{ for some } p \in G, p \neq 0, p \geq 0\}$. That is, $\beta^w(e)$ is the set of competitive equilibrium trades for the environment e . There are many examples of environments in which β^c is not well-defined; that is, in which there does not exist an equilibrium price or, if there does, it is not unique. However, if $E' \subseteq E_d$ is the class of environments such that X is a finite dimensional vector space, $D(p, e) = \sum_{i \in I} D^i(p, e)$ is a continuous function of p , and $D(p, e)$ satisfies a gross substitute condition then β^w is well defined on the set E' .

Example 5: (Lindahl mechanism). For the environments in example 2, let F be a cone with vertex 0; let each i choose (\bar{a}^i, \bar{b}^i) to be $\geq i$ maximal subject to $pa^i + t^i b^i \leq 0$; let a producer choose (\bar{a}, \bar{b}) to maximize $pa + (\sum_{i \in I} t^i) b$ over F . If $\bar{b}^i = \bar{b}$ for all $i \in I$, and $\sum_{i \in I} \bar{a}^i = \bar{a}$, then let $\{(\bar{a}^1, \bar{b}^1), \dots, (\bar{a}^I, \bar{b}^I)\} \in \beta^L(e)$. We call $\beta^L(e)$ the Lindahl choice for e .

Although we have concentrated on economic (allocation) mechanisms to this point, our model is general enough to include many political processes or social choice rules. We present one example.

Example 6: (Majority Rule). Given an environment e , let $\beta^M(e) = \{a \in S^I \mid aMa^* \text{ for all } a^* \in S^I\}$ where aMa^* iff $N(a \succ a^*) \geq N(a^* \succ a)$. $N(a \succ a^*)$ is the percentage of I who prefer a to a^* . It should be

noted that there are many environments for which β^M is not well-defined, (i.e., $\beta^M(e) = \emptyset$). However it is a choice mechanism.

c. Enforcement Structures and Preference Games

Given an environment, e , and the outcome function of a choice mechanism, β , the outcome for each agent is his component of $h[\beta(e)]$. Unless β is completely insensitive to individual preferences (as it would be, e.g., if it were "imposed"), if some agent, say i , were to act as if his preferences were \succsim^*_i instead of acting according to his true preferences \succsim_i , it is possible that³ $h^i[\beta(e/\succsim^*_i)] > h^i[\beta(e)]$. If this opportunity is available to Mr. i , we say that he has an incentive to "misrepresent his preferences," since by so doing he can induce the choice mechanism to select an action which he prefers to the action which would otherwise have been selected.⁴ The ability of an agent to carry off such misbehavior depends on whether or not his misrepresentation is detectible. This in turn depends on the enforcement structure in existence. Rather than explicitly model that structure, we will implicitly allow for alternative enforcement procedures through their impact on the sets of allowable preference misrepresentations. Thus, for example, if preferences are not directly observable by the enforcement agency, and if the choice mechanism is the Walrasian Tatonnement, then only indirect evidence of misbehavior is available through information about excess demand functions. In this case we could let the set of allowable preferences (which are to be used as misrepresentations) include all those which satisfy an axiom of revealed preference.

Thus, an enforcement structure is a set of allowable preference

misrepresentations, R^i , for each i . Denote by R , the set $\prod_{i \in I} R^i$.

Given an environment, an outcome function of a mechanism, and an enforcement structure, each agent, if he acts in his own self-interest, may have an incentive to misbehave by revealing preferences which are different than his true preferences. This characterization of self-interest suggests a game-theoretic approach to the analysis of incentives, in which (possibly false) preferences define a player's strategies. This game can be modeled in many ways: cooperative or non-cooperative, complete or incomplete information, etc. However, for our purposes we will consider only the non-cooperative, complete information game.

Thus, given the triple $\langle e, \beta, R \rangle$, (an environment, choice mechanism, and enforcement structure), a preference game is determined as follows. The set of players is I , determined by e . For each $i \in I$, the admissible strategy set is R^i , determined by R . For each $i \in I$, the payoff of a joint strategy is expressed as a preference ordering on R where $(\succ^*) \succ_i (\succ^{**})$ iff $h^i[\beta(e/\succ^*)] \succ_i h^i[\beta(e/\succ^{**})]$.

Placed in this game, the agents will determine a solution strategy which may, or may not, be to use to their true preferences. A choice mechanism which leads people to use their true preferences, for a wide class of environments, is desirable since the use, by the agents, of false preferences will in general lead the choice mechanism to select joint actions which are inferior, from the point of view of the whole, to those which are chosen when true preferences are used. This is illustrated by example 7 below where all potential gains from trade are foregone because of self-

interested behavior. A desirable choice mechanism is one for which self-interest leads the participants to behave in a way which is compatible with the goals of the system as expressed by the choice mechanism. One characterization of such a mechanism is that the true preferences are a solution of the preference game. A mechanism with this property will be called individually incentive compatible. Many solution concepts are available; however, we will rely solely on the concept of the Nash-equilibrium.

Thus, a choice mechanism, β , is said to be individually incentive compatible in the environment e under the enforcement structure R if and only if the vector of true preferences $\langle \beta_i \rangle_{i \in I}$ is a Nash equilibrium of the preference game derived from (e, β, R) .

To illustrate these concepts we provide a simple example.

Example 7: Let $e \in E_d$ be a decomposable exchange environment with two people and two commodities. $I = \{1, 2\}$, with $A^1 = \{(a_1^1, a_2^1) \in R^2 \mid a_1^1 \geq -1, a_2^1 \geq 0\}$ and $A^2 = \{(a_1^2, a_2^2) \in R^2 \mid a_1^2 \geq 0, a_2^2 \geq -1\}$. Also suppose preferences are representable by the utility functions $U^1(a^1) = (a_1^1 + 1)(a_2^1)^\alpha$ and $U^2(a^2) = (a_1^2)(a_2^2 + 1)^\gamma$. Finally, let $S^I = \{(a^1, a^2) \in A^1 \times A^2 \mid a^1 + a^2 = 0\}$. This is the environment which arises from a pure exchange world with 2 consumers, two commodities, and initial endowments of $(1, 0), (0, 1)$. In this environment, the competitive equilibrium price ratio $r = p_1/p_2 = \frac{\gamma}{\alpha} \frac{1+\alpha}{1+\gamma}$ and the equilibrium actions are $\beta^{w1}(e) = a^1(r) = (-\frac{\alpha}{1+\alpha}, r \frac{\alpha}{1+\alpha})$ and $\beta^{w2}(e) = (\frac{\gamma}{r(1+\gamma)}, -\frac{\gamma}{1+\gamma})$. Thus, for example, given e , Mr. 1 receives the utility $U^1(\beta(e)) = \frac{1}{(1+\alpha)} [\frac{\gamma}{(1+\gamma)}]^\alpha$. Suppose that the true environment is such that $\alpha = \gamma = 1$. Then $U^1(\beta^w(e)) = \frac{1}{4} =$

$U^2(\beta^w(e))$. Now suppose Mr. 1 were to act as if his α were really $\frac{1}{2}$. $\beta^{w1}(e^*)$ is now $(-\frac{1}{3}, \frac{1}{2})$ and $U^1(\beta^w(e^*)) = \frac{1}{3}$. Thus, by behaving as if he likes commodity 2 relatively less, Mr. 1 is able to attain a higher level of satisfaction than if he used his true preferences when he calculates his excess demands. It is also true, in this example, that if the true values $\hat{\alpha}$, $\hat{\gamma}$ are positive, then for any other positive α , γ it is true that $\partial U^1(\beta^w(\alpha, \gamma)) / \partial \alpha < 0$, and $\partial U^2 / \partial \gamma < 0$. The only Nash-equilibrium of the preference game, [when R^i is the set of utility functions derived from those above with $\alpha \geq 0$, $\gamma \geq 0$], is $\alpha = \gamma = 0$, in which case $\beta^{w1}(\alpha, \gamma) = \beta^{w2}(\alpha, \gamma) = 0$ and neither agent receives any utility.

One conclusion based on the above example is that the competitive process is not incentive compatible in Cobb-Douglas environments even if misrepresentations are restricted to be Cobb-Douglas. Of interest, then, is the class of triples (e, β, R) which yield incentive compatibility. We now turn to the problem of characterizing this class.

3. THE RESULTS

In this section we present some theorems which partially characterize the class of triples (e, β, R) for which the choice mechanism β is individually incentive compatible in the environment e under the enforcement structure R . Loosely stated, the major result is that, given some technical restrictions on e , for pairs (e, β) such that β is well-defined and selects core allocations in e , and for an R which is implied if preferences are unobservable, β is individually incentive compatible in e under R if and only if e possesses an action which is also the unilaterally

best unblocked action for all agents.⁶ (See Theorem 3). The core must thus consist of actions which are Pareto-indifferent.

We begin by restricting our attention to pairs of environments and choice mechanisms such that the mechanism selects unique (up to indifference) actions which belong to the core of the environment.

Definition 1: We say that β is well-defined for the environment e iff

- (i) $\beta(e) \neq \emptyset$, (β is decisive), and
- (ii) if $a, a' \in \beta(e)$, then $h^i(a) \sim_i h^i(a')$ for (almost)⁷ all $i \in I$, (β is essentially single valued).

Definition 2: Given an environment e , we say that a joint action $a^* \in S^I = \hat{A}$ is blocked by coalition $c \in C$ if and only if there exists an alternative action $a^c \in S^c$ such that $h^i(a^c) \succeq_i h^i(a^*)$ for (almost) all $i \in c$, and $h^k(a^c) \succ_k h^k(a^*)$ for some $k \in c$, where "some" means "for a set of positive measure."

Definition 3: Given an environment e , the core of e is the set of all actions belonging to S^I which are not blocked by any $c \in C$. Denote the core of e by $\text{Core}(e)$.

There is some debate over whether this definition of the core is appropriate for environments with externalities, such as those in example 2 with public goods. For the purposes of this paper, it is unnecessary to take sides on this issue, since we will only use the definition of the core to characterize classes of environments and/or choice mechanisms. Thus, we do not need

to consider whether it may (or may not) be an appropriate solution concept for some game such as a barter process.

We will be restricting our attention to pairs (e, β) such that $\beta(e) \neq \emptyset$, $\beta(e)$ is essentially single-valued, and $\beta(e) \subseteq \text{Core}(e)$.

Definition 4: We let $D = \{(e, \beta) \in E \times B \mid \beta(e) \neq \emptyset, \beta(e) \subseteq \text{Core}(e), \beta(e) \text{ is essentially single-valued.}\}$

As we indicated in the introduction, a plausible conjecture about the triple (e, β, R) is that if $(e, \beta) \in D$, then β is individually incentive compatible in e for most enforcement structures if and only if the core of e is a single point (in utility space). The following two examples show that this is not true. In each, β selects core actions and the core of e is a unique action (up to indifference). In each, however, there exists at least one agent who has an incentive to misrepresent his preferences (when all the other agents do not).

Example 8: Let e be the environment with $I = \{1, 2, 3\}$ where, for all $i \in I$, $A^i = R^2$ (the two-dimensional Euclidean space), and $S^{\{i\}} = \{(0, 0)\}$, where $S^{\{1, 2\}} = \{(a^1, a^2) \in R^4 \mid a_2^1 = a_2^2 = 0, a_1^1 + a_1^2 \leq 2 - \sqrt{2}\}$, $S^{\{1, 3\}} = \{(a^1, a^3) \in R^4 \mid a_1^1 = a_1^3 = 0, a_2^1 + a_2^3 \leq 2 - \sqrt{2}\}$, $S^{\{2, 3\}} = \{(a^2, a^3) \in R^4 \mid (a_1^2 + a_2^2)^2 + (a_1^3 + a_2^3)^2 \leq 1\}$, and $S^I = \{(a^1, a^2, a^3) \in R^6 \mid \sum_{i=1}^3 (a_1^i + a_2^i) \leq 2\}$, and where, for all $i \in I$, $U^i(a^i) = a_1^i + a_2^i$. We also let $h^i(a) = a^i$.

Let $\beta(e) = \{a \in S^I \mid a \in \text{Core}(e) \text{ and } a \geq 1a^* \text{ for all } a^* \in \text{Core}(e)\}$. In the environment e , the core of e is equal to $\{a \in R^6 \mid a_1^1 + a_2^1 = 2 - \sqrt{2}, a_1^2 + a_2^2 = \sqrt{2}/2, \text{ and } a_1^3 + a_2^3 = \sqrt{2}/2\}$. Thus the $\text{Core}(e)$

is a single point in utility space where $U = (U^1, U^2, U^3) = (2-\sqrt{2}, \sqrt{2}/2, \sqrt{2}/2)$. Also, $\beta(e) \neq \emptyset$, $\beta(e)$ is essentially single-valued, and $\beta(e) \subseteq \text{Core}(e)$. However, even though the core is essentially unique, β is not individually incentive compatible if $\hat{U}^1(\cdot) \in R^1$ where $\hat{U}^1(a^1) = a_1^1 \cdot a_2^1$, (a simple Cobb-Douglas utility function), since the joint action $a^* = [(7-3\sqrt{2}/8, 7-3\sqrt{2}/8), (\sqrt{2}+1/2, 0), (0, \sqrt{2}-1/2)] \in \text{Core}(e/\hat{U}^1)$. Thus, $U^1[h^1(\beta(e/\hat{U}^1))] \geq U^1(a^*) = (7/4) - (3\sqrt{2}/4) > 2 - \sqrt{2} = U^1[h^1(\beta(e))]$. Thus, true preferences are not a Nash equilibrium of the preference game derived from (e, β, R) .

Example 9: This example is based on a model of Shitovitz (1973) of duopoly. [See also example 3 above.] Let e be the decomposable exchange environment with two commodities, and with 2 identical traders (atoms) indexed by $i = 3, 4$ and an infinite number of identical individuals (the continuum) indexed by $i \in [0, 2]$. The utility functions of 3 and 4 are $U^i(a^i) = (a_1^i + 1)(a_2^i)^\alpha$ which would arise if initial endowments are $(1, 0)$, $h^i(a^i) = a^i + (1, 0)$, and $U^i(x^i) = x_1^i (x_2^i)^\alpha$. The utility functions of $i \in [0, 2]$ are $U^i(a^i) = (a_1^i)^\beta (a_2^i + 1)$. In order to describe completely the set of admissible coalitions and the feasible action set S^C for each $c \in C$, we let μ be a measure on $[0, 2] \cup \{3, 4\}$ such that $\mu(\{2\}) = \mu(\{3\}) = 1$ and such that μ is the Lebesgue measure on $[0, 2]$. C is the set of all μ -measurable subsets of I with positive measure, and, for each $c \in C$, $S^c = \{ \langle a^i \rangle_{i \in c} \mid \int_c a^i d\mu(i) = 0, a^i \in R^2 \forall i \}$.

Let β^w be the Walrasian Tatonnement where $\beta^w(e)$ is the set

of competitive allocations in e .

In the environment e , there is a unique competitive equilibrium yielding a price ratio $r = P_1/P_2 = (\beta/\alpha)(1+\alpha/1+\beta)$. For each atom, $i = 2, 3$, $h^i[\beta^W(e)] = [1 - (\alpha/(1+\alpha)), r\alpha/(1+\alpha)]$. By Theorem B of Shitovitz (1973) $\beta^W(e) = \text{Core}(e)$ for this environment, and therefore the core is a single action since $\beta^W(e)$ is. However, it is possible for either duopolist to gain by misrepresentation of preferences. Suppose the true values of α and β are $\alpha = \beta = 1$. The equilibrium price ratio is $r = 1$ and the equilibrium action (trade) for each atom is $(-\frac{1}{2}, \frac{1}{2})$ yielding a utility to each of $\frac{1}{4}$. Now, suppose that Mr. 3 acts as if his $\alpha = \frac{1}{2}$. Then the equilibrium price will be $r = \frac{6}{5}$ and his equilibrium action will be $(-\frac{1}{3}, \frac{2}{5})$ yielding a true utility of $(1 - \frac{1}{3})(\frac{2}{5}) = \frac{4}{15}$ which is greater than $\frac{1}{4}$.

Thus, even though the core of e is single-valued, a duopolist who is able to "differentiate" himself from his identical competitor can gain by so doing.⁹

As the above examples indicate, a single-valued core is not sufficient to guarantee that a well-defined choice mechanism which selects core actions will be incentive compatible in that environment. Consequently, we need a slightly stronger concept than the core. In particular, the mechanism must select an action which is unilaterally best unblocked for all agents.

Definition 5: Given an environment e , let $C_e^i = \{c \in C \mid c \cap \{i\} = \emptyset\}$. Let $B_e^i = \{a \in S^I \mid a \text{ is not blocked by any } c \in C_e^i\}$. [That is, B_e^i is the set of all jointly feasible actions which are only blocked, if at all, by coalitions which contain agent i as

a member.] We say that $a^* \in S^I$ is a unilaterally best unblocked action for i iff $a^* \in B_e^i$ and $a^* \succ_i a$ for all $a \in B_e^i$; (that is, iff a^* is one of i 's most preferred actions among those which no coalition can block unless i contributes resources.)

To illustrate the relationship between core actions and those which are unilaterally best unblocked, we consider two lemmata.

Lemma 1: Given an environment e if there exists an action $a \in S^I$ which is a unilaterally best unblocked action for all $i \in I$, then $a \in \text{Core}(e)$.

Proof: Assume $a \notin \text{Core}(e)$. Then there exists a coalition $c^* \in C$ which blocks a . Since $a \in B_e^i$ for all $i \in I$, $c^* = I$. Thus, there exists an action \hat{a} which is Pareto-superior to a . Therefore, $\hat{a} \in B_e^k$ and $\hat{a} \succ_k a$ for some $k \in I$. Hence a is not a unilaterally best unblocked action for all $i \in I$. QED.

Lemma 2: Given an environment e , if $a \in S^I$ is a unilaterally best unblocked action for all $i \in I$ and if $b \in \text{Core}(e)$, then $a \sim_i b$ for all $i \in I$. (That is, $a \in \text{Core}(e)$ and $\text{Core}(e)$ is essentially single-valued)

Proof: By lemma 1, $a \in \text{Core}(e)$. Let $b \in \text{Core}(e)$ such that $b \sim_i a$ does not hold for some $i \in I$. Then $b \succ_k a$ for some $k \in I$. If not, $a \succ_k b \forall k \in I$ and $a \succ_i b$ for some $i \in I$ which implies that a is Pareto-superior to b and, therefore, $b \notin \text{Core}(e)$. Now, since $a, b \in \text{Core}(e)$, $a \in B_e^i$ and $b \in B_e^i$ for all $i \in I$. Thus, for some $k \in I$, $b \succ_k a$ and $a, b \in B_e^k$. Hence, a is not a unilaterally best unblocked allocation for k . QED.

The two results can be summarized in the following way: If

an action exists which is unilaterally best unblocked for all agents, then that action is a core action and it is Pareto-indifferent to all other core actions. Hence, the core must be a single point in utility space.^{10/} Example 8 indicates that the converse of Lemma 2 is not true. The converse of Lemma 1 is obviously false since actions in multi-valued cores are clearly not unilaterally best unblocked for all i .

We turn now to the major findings of this paper which contain a partial characterization of triples (e, β, R) for which incentive compatibility of β in e under R obtains.

Theorem 1: Given a triple (e, β, R) , such that

(a) $\succsim_i \in R^i$ for all $i \in I$

(b) $[(e/\succsim_i), \beta] \in D$ for all $i \in I$ and all $\succsim_i \in R^i$ (see Definition 4) and

(c) $\beta(e)$ is a unilaterally best unblocked action for all $i \in I$,

then (1) the mechanism β is individually incentive compatible in the environment e under the enforcement structure R .

Proof: Suppose the conclusion is false. Then, there is $i \in I$, and a preference ordering $\succsim_i \in R^i$ such that

$$\beta(e/\succsim_i) \sim_i \beta(e).$$

But $(\beta, e/\succsim_i) \in D$, by (b), which implies that $\beta(e/\succsim_i) \in \text{Core}(e/\succsim_i)$.

Therefore, $\beta(e/\succsim_i) \in B_{(e/\succsim_i)}^i = B_e^i$. [B_e^i is defined in Definition 5]. By

(b), $\beta(e) \in B_e^i$. Hence, $\beta(e)$ is not a unilaterally best unblocked action in e for i . This contradicts (c). QED.

Theorem 1 gives sufficient conditions for incentive compatibility. Of these conditions, only (b) has not been discussed in this paper. It merely requires that the choice mechanism should

be a well-defined, core selector in all the environments which may occur through misrepresentations of preferences. If the mechanism were not, then the problem is not well-defined and one should not expect any conclusions to be achievable.

A possible partial converse to Theorem 1 is that conditions (a), (b), and (1) imply condition (c). However, as the following example shows, that is not, in general, true.

Example 10: We let β be the Walrasian Tatonnement process and consider an Edgeworth box environment.

INSERT FIGURE 1

Here, O_i is the offer curve of i , $b = \beta(e)$, and I_i is the indifference curve of i through b . The curve, abz is Core (e). Obviously, $(e, \beta) \in D$, and $\beta(e)$ is not a unilaterally best unblocked action for either i . However, there is no $\succsim_i^* \in R_e^i$ for which $\beta(e/\succsim_i^*) \succ_i \beta(e)$. This is true since b is the best i can do given that a must lie on the other's offer curve.

It should be recognized, however, that it is possible for (say) Mr. 1 to select $\succsim_1^* \in R_e^1$ such that $\text{Core}(e/\succsim_1^*) \subseteq \{a \in S \mid a \succ_1 b\}$. Then, if $(e, \beta) \in D$, it would be true that $b^* \in \beta(e/\succsim_1^*)$ implies $b^* \succ_1 b$, and the proposition would hold. In our example, however, for such an $(e/\succsim_1^*) = e^*$ either $\beta(e^*) = \emptyset$ or $\beta(e^*) \subseteq \text{Core}(e^*)$. Thus, $(e^*, \beta) \notin D$. An example of such an e^* can be constructed as follows: Let $w \sim^* 1 z$ and let $a \succ^* 1 a'$ iff $a \succ_1 a'$ for all other $a, a' \in S$. $\text{Core}(e^*) = z$. But, for e^* , Mr. 1's offer curve, O_1^* is now the point w unioned with all points on O_1 , to the right of z (including z itself). Obviously, $\beta(e^*) = \emptyset$, because of the discontinuity in O_1^* . [Alternatively, because the excess demand correspondence is not convex.] If we let $\beta^*(e^*) = w$ for all e^* such that $\beta(e^*) = \emptyset$ and $\beta^*(e^*) = \beta(e)$ otherwise, then $\beta^*(e^*) \in \text{Core}(e^*)$ and again $(e^*, \beta^*) \notin D$.

In spite of these difficulties, we can, by imposing some additional restrictions on the environments and mechanisms, demonstrate the necessity of unilaterally best unblocked actions for incentive compatibility. To do this, we need to introduce two additional concepts. First, given an environment e , we let $g^i(e) = \{w^i \in S^{\{i\}} \mid w^i \succ_i a^i \text{ for all } a^i \in S^{\{i\}}\}$. $g^i(e)$ is the set of best actions which i can achieve if he acts independently of all others. We will require below that $g^i(e) \neq \emptyset$. For pure exchange environments, this is innocuous since $S^{\{i\}} = \{0\}$. However, if i is a consumer-producer; it is possible that $S^{\{i\}}$ is, e.g., not closed and $g^i(e) = \emptyset$. Second, given an environment e , we let $L^0(e) = \{a \in S^I \mid a \succ_i w^i \text{ for all } i \in I \text{ where } w^i \in g^i(e)\}$. In a two person, two commodity, exchange environment with continuous

preferences, $L^0(e)$ is essentially the interior of the lens in the Edgeworth box.

Utilizing these two definitions we can state our next theorem. The assumptions and conclusions are discussed following the proof.

Theorem 2: Given a triple (e, β, R) such that $\succsim_i \in R^i$ for all $i \in I$, β is individually incentive compatible in e under R for i , and such that for all significant ^{11/} $i \in I$,

- (a) $[(e/\succsim_i), \beta] \in D$ for all $\succsim_i \in R^i$,
- (b) $\succsim_i \in R^i$ if and only if \succsim_i is complete, transitive, and reflexive ^{12/},
- (c) $g^i(e) \neq \emptyset$,
- (d) if $a \succ_i b$, there exists z such that $a \succ_i z \succ_i b$,
- (e) for all $\succsim_i \in R^i$, if $L^0(e/\succsim_i) \neq \emptyset$, then $\beta(e/\succsim_i) \subseteq L^0(e/\succsim_i)$, then for all $i \in I$, either

A. i is not significant

or

B. ^{13/} if $L^0(e) \neq \emptyset$ and $a \in \beta(e)$, then $a \succ_i b$ for all $b \in B_e^i \cap L^0(e)$.

Proof: Assume $L^0(e) \neq \emptyset$ and assume that there exists a significant $i \in I$ and $b \in B_e^i \cap L^0(e)$ such that $b \succ_i z$ for all $z \in \beta(e)$. By assumption (d), there exists a d such that $b \succ_i d \succ_i z$. Let $w^i \in g^i(e)$. We construct a (false) preference ordering, \succsim_i , for i such that $\beta(e/\succsim_i) \succ_i z$ for all $z \in \beta(e)$. This construction involves simply making the single allocation w^i indifferent to d , while leaving all other allocations unchanged. More

precisely let $U(x, \succsim_1) \equiv \{y \mid y \succsim_1 x\}$ be the upper contour set of x under the preferences \succsim_1 . Define $\succsim^* i$ as follows: ^{14/} (1) $U(w^i, \succsim^* i) = U(d, \succsim_1)$, (2) for $a \neq w^i$, if $d \succ_1 a$ then $U(a, \succsim^* i) \equiv U(a, \succ_1 i)$ and if $a \succ_1 d$ then $U(a, \succsim^* i) \equiv U(a, \succ_1 i) \cup \{w^i\}$. Geometrically, the indifference "curve" through d under the false preferences consists of the true indifference curve through d unioned with the single point w^i . The indifference "curve" for a $\sim w^i$ consists of the old indifference curve through w^i with the point w^i deleted. All other indifference curves remain unchanged. Given this construction, since $b \in L^0(e)$ and since $b \succ_1 d$, which implies $b \succ^* i w^i$, it follows that $b \in L^0(e/\succ^* i)$ and, therefore, that $L^0(e/\succ^* i) \neq \emptyset$. Thus, by assumptions (a) and (e), $\beta(e/\succ^* i) \in \text{Core}(e/\succ^* i) \cap L^0(e/\succ^* i)$, and $\beta(e/\succ^* i) \neq \emptyset$. Hence, $\hat{a} \in \beta(e/\succ^* i)$ implies that $\hat{a} \succ_1 d \succ_1 z$ for any $z \in \beta(e)$. Thus, incentive compatibility is violated.

Q.E.D.

Hypothesis (a) has appeared before, in Theorem 1, and has been discussed. Assumption (b) is introduced to ensure that R^i is large enough to make incentive compatibility a problem. Clearly, if the true preferences, \succ_1 , were the only element of R^i , then incentive compatibility is obvious for all β in all environments since no misrepresentations can occur. If we assume, however, that preferences are not directly observable by the enforcement agency but that participants must satisfy a revealed preference axiom for a wide enough range of choice sets, then a theorem of Arrow (1959) tells us that any complete, transitive, reflexive ordering will do. Thus (b) is similar to an assumption of unobservable preferences -- a not unrealistic hypothesis.

Assumptions (d), (e) and the use of $L^0(e)$ in statement B appear to be needed more because of the method of proof than because of the nature of the problem. The technique of proof involves showing that if a well defined core selecting mechanism is not operating in an environment with an action which is a unilaterally best unblocked action for all significant agents, then some agent, by moving his best independent action to a higher indifference class, can ensure a higher payoff to himself. It can be shown that, under assumptions (a), (b) and (c) if $\beta(e)$ is not a unilaterally best unblocked action for some significant $i \in I$, then altering preferences by placing $w^i \in g^i(e)$ into the indifference class of the unblocked action which is better than $\beta(e)$ will insure that either $\beta(e/\succsim^*i) \succ_i \beta(e)$ or $\beta(e/\succsim^*i) = w^i$ which may be less preferred than $\beta(e)$. This strategy for changing preferences is similar to the child who says "either play my way or I will take my ball and go home." The result of such a strategy might be that he is told to go home, (i.e., $\beta(e/\succsim^*i) = w^i$). To avoid this possibility we must insure that the environment and mechanism have enough slack to make such a strategy productive. Thus, we use assumptions (d) and (e). Assumption (d) merely requires preferences to satisfy an Archimedean property. Assumption (e) requires the mechanism to choose, whenever possible, core actions which provide gains from trade to every agent.

The conclusion of theorem 2 is not as strong as we would like. In particular it would be desirable if, instead of B, one had that if $a \in \beta(e)$ then a is a unilaterally best unblocked action for every significant $i \in I$. If for almost all agents in the true environment e , the action space, A^i , is Euclidean and preferences

are monotonic and continuous (in the usual topology), then $a \succsim_i b$ for all $b \in B_e^i \cap L^0(e)$ implies $a \succsim_i b$ for all $b \in B_e^i$ (a is a unilaterally best unblocked action for i). Rather than list all possible circumstances under which B implies unilaterally best unblocked actions, we state the following corollary to theorem 2.

Corollary 2.1: Under the hypotheses of theorem 2, if $a \succsim_i b$ for all $b \in B_e^i \cap L^0(e)$ implies that $a \succsim_i b$ for all $b \in B_e^i$ then either $L^0(e) = \emptyset$ or $a \in \beta(e)$ implies a is a unilaterally best unblocked action for all significant $i \in I$.

Remember from Lemmas 1 and 2 that if a is a unilaterally best unblocked action for all agents then the core is essentially single valued. This leads to another implication of incentive compatibility.

Corollary 2.2: Under the hypotheses of Theorem 2, if $\text{Core}(e) \subseteq L^0(e)$, and if there are a finite number of agents, then whenever $a, b \in \text{Core}(e)$ it must be true that $a \sim_i b$ for all $i \in I$.

That is, in a large class of environments, a choice mechanism is incentive compatible only if the core of the true environment is essentially a single action. We indicated earlier that the converse was false. However, one situation of interest for which the converse holds is when the action, w , which is best for each agent acting independently [i.e., $w^i \in g^i(e)$] is also Pareto-optimal. In this case, $L^0(e) = \emptyset$, w is a unilaterally best unblocked allocation for each $i \in I$, $w \in \text{Core}(e)$, and any core selecting mechanism is incentive compatible in e under any enforcement structure R .

Finally, we summarize the results contained in the previous theorems.

Theorem 3: Given the triple (e, β, R) such that

(a) there are a finite number of agents and $\succsim^i \in R^i$ for all $i \in I$,

(b) for all $i \in I$, $[(e/\succsim^i), \beta] \in D$ for all $\succsim^i \in R^i$,

(c) for all $i \in I$, $\succsim^i \in R^i$ iff \succsim^i is complete, transitive, and reflexive,

(d) $L^0(e) \neq \emptyset$ and, for all $i \in I$, $g^i(e) \neq \emptyset$ and $[a \succ^i b \Rightarrow \exists c \ni a \succ^i c \succ^i b]$,

(e) for all $i \in I$ and all $\succsim^i \in R^i$,
 $[L^0(e/\succsim^i) \neq \emptyset \Rightarrow \beta(e/\succsim^i) \subseteq L^0(e/\succsim^i)]$.

(f) for all $i \in I$, $[a \succ^i b \vee b \in B_e^i \cap L^0(e)] \Rightarrow [a \succ^i b \vee b \in B_e^i]$

Then

β is incentive compatible in e under R if and only if $a \in \beta(e)$ implies a is a unilaterally best unblocked action for each $i \in I$.

Hypotheses (a), (b) and (c) are the crucial ones. (d) to (f) are technical in nature and due to the method of proof. It is highly likely that (c) can be relaxed to let R^i contain only preferences which are continuous. It is our conjecture, that, under this weakening, (d) to (f) should be unnecessary in almost all environments with continuous preferences. It would be desirable to establish this since accepted market mechanisms such as the competitive process or a Lindahl system do not simultaneously satisfy assumptions (b) and (c) since there are preference orderings for which a competitive equilibrium does not exist.

4. SOME CONCLUDING OBSERVATIONS

The basic conclusion which follows from the results contained in this paper is that participants in a choice process or allocation system which selects core actions will seldom find it in their self interest to follow the norms of behavior of that system. There are only three obvious situations where behavior will be as prescribed; (1) when taking no action is Pareto-optimal, (2) when there is an atomless continuum of agents, and (3) when the situation yields a market game [e.g. transferable utility] with a single-valued core. The first is uninteresting and the others are improbable. This is especially disturbing when one recognizes that the concept of incentive compatibility, which we have used seems to be a minimal requirement even if no agent knows the strategy choices of the others. To see this, let us suppose a strategy choice has been made by each agent and then one agent is allowed to alter his choice while the others keep theirs fixed. If the choices do not constitute a Nash-equilibrium of the complete information incentive game then one might expect that (at least in repeated trials) the player would hit upon a better strategy even if he does not know in advance the outcome resulting from any particular choice.

One way to overcome the pessimistic nature of this conclusion is to recognize that the Nash equilibrium is a static concept.

To be able to call it an equilibrium, one must, at least implicitly, have in mind some dynamic process of which the Nash equilibria are steady state solutions. (E.g., implicit in the argument in the above paragraph is some search process leading to Nash equilibria.) When a dynamic process is explicitly introduced into the preference game, it is likely that, in some cases, true preferences may be the equilibrium joint strategy of that dynamic process even if they do not comprise a Nash-equilibrium.

A second conclusion, following from the results contained in this paper, is that care must be taken when one interprets the implications of the theorems which state that in large economies the core and the set of competitive equilibria coincide. It is tempting to conclude from those theorems that, in the economies for which they hold, any agent or group of agents might as well behave competitively by taking prices as given. However, as indicated by example 9, based on a Shitovitz duopoly model, this interpretation is not always valid. In that environment, even though the core and the set of competitive equilibria coincide, there are (unilateral) strategies available to some agents which are essentially undetectable, which can be interpreted as non-competitive behavior, and which can lead to the implementation of actions not in the core but better for that agent. A more complete exploration of the relationship between the solution concepts of competitive equilibrium, the core, and the

Nash-equilibrium of the preference game is desirable. An initial step in this direction has been taken by Roberts and Postlewaite (1976).

Footnotes

1. If I is finite, $C =$ the set of all subsets of I except the empty set. In general, we can represent I and C by selecting a positive measure μ on $[0, \infty)$, letting I be the support of μ , and letting C be the set of all μ -measurable subsets of $[0, \infty)$ with positive measure.
2. In this language, $S^I = \hat{A}$.
3. The symbol $(e/\succ *i)$ will be used to represent the environment which is identical to e except that \succ is replaced with $\succ *i$.
4. It is also conceivable that i could gain through misrepresentations in the other components of his part of the environment (such as his initial endowment or possible action set); however, for the purposes of this paper we do not need to allow him such latitude.
5. Gibbard and Satterthwaite call voting mechanisms with this property strategy proof.
6. These terms are all defined below.
7. If e is an environment with a continuum of agents, almost all $i \in I$ means all i except for a subset of μ -measure zero. For finite I , there are no subsets of μ -measure zero, so "almost all" is equivalent to "all".
8. The key to this example is the structure of the utility possibility space, given e . Although I am convinced that one can construct a decomposable exchange environment (with, perhaps, personalized goods and personalized externalities) which would imply the same utility possibility frontiers, I do not yet have such an example.

The construction of Shapley and Shubik is not applicable since the desired structure is not a market game.

9. It is interesting to note, in this example, that if Mr. 3 misrepresents his preferences then Mr. 4 gains even more than 3 does. That is, $U^4[\beta^w(e/\alpha_3 = \frac{1}{2})] = \frac{3}{10} > U^3[\beta^w(e/\alpha_2 = \frac{1}{2})] = \frac{4}{15} > U^4[\beta^w(e)] = \frac{1}{4}$. Thus, it is in 3's interest to encourage 4 to misbehave. It is an open problem as to how widespread this potential gain for similar agents really is.
10. Since it is true that in most environments the core is not single-valued, Lemma 2 implies that an action which is unilaterally best unblocked usually does not exist. For example, in a two person, pure exchange environment (see example 1) a unilaterally best unblocked action for all $i \in I$ exists if and only if the initial endowment is a Pareto-optimal allocation.
- The usual lack of existence of such actions only serves to emphasize the fact that rarely is a core selecting mechanism incentive compatible. (See Theorem 3.)
11. $i \in I$ is significant if $(\{i\}) \in C$. In finite economies, "almost all i " is equivalent to "all significant i ".
12. (c) can be weakened with no difficulty to "there exists a topology on actions such that $\succ^* i \in R^i$ if and only if $\succ^* i$ is representable (in that topology) as an upper semi-continuous function."
13. We can not prove that B holds for all $i \in I$. Consider the competitive mechanism which is incentive compatible in an atomless environment. [See Roberts and Postlewaite (1976)]. The core may not be essentially single valued if there does not exist a unique equilibrium price and in that case

B will not be true for any $i \in I$. The conclusion recognizes this since, in this example, no $i \in I$ is significant.

An alternative theorem arises if the word "significant" is everywhere deleted; however, continuum economies with some atomless components would then not be covered. If i is not significant, $\text{core}(e) = \text{core}(e/\succsim^*i)$ for all $\succsim^*i \in R^i$ but $L^0(e) \neq L^0(e/\succsim^*i)$. It is easy to find preferences such that $\text{core}(e/\succsim^*i) \cap L^0(e/\succsim^*i) = \emptyset$ which means that (e) could not hold simultaneously with (b).

14. If \succsim^*i are representable by a continuous utility function then, for all a , $\{y \mid y \succ^*i a\}$ is an open set. Suppose $a \sim^*i w^i$. The false strict upper contour set $\{y \mid y \succ^*i a\}$ is $\{y \mid y \succ^*i a\} \cup \{w^i\}$ which is not open. Therefore, \succsim^*i cannot be represented by a continuous utility function. It is representable by an upper semi-continuous function.

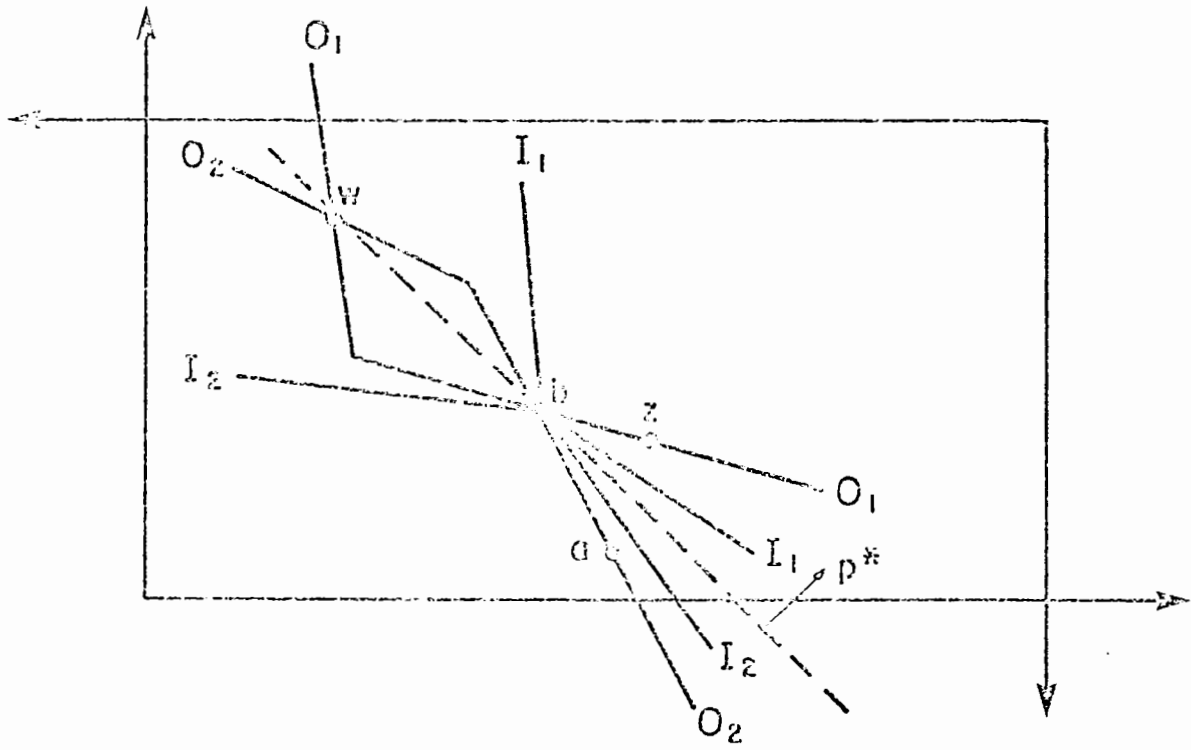


Figure 1

REFERENCES

- Arrow, K. [1959], "Rational Choice Functions and Orderings," Economica.
- Aumann, R. [1966], "Existence of Competitive Equilibrium in Markets with a Continuum of Traders," Econometrica, Vol. 34, No. 1.
- Gibbard, A. [1973], "Manipulation of Voting Schemes: A General Result," Econometrica, Vol. 44.
- Hurwicz, L. [1972], "On Informational Decentralized Systems," in Decision and Organization (Volume in honor of J. Marschak) Radner, R. and B. McGuire (eds.) North Holland Press, Amsterdam.
- Hurwicz, L. [1973], "The Design of Mechanisms for Resource Allocation," American Economic Review, Vol. 63, No. 2.
- Muench, T. J. [1972], "The Core and the Lindahl Equilibrium of an Economy with a Public Good: An Example," Journal of Economic Theory, Vol. 4, No. 2.
- Postlewaite, A. and D. J. Roberts [1976], "The Incentives for Price-Taking Behavior in Large Exchange Economies," Econometrica, Vol. 44, No. 1.
- Satterthwaite, M. [1973], "The Existence of a Strategy-Proof Voting Procedure," Discussion Paper No. 42, The Center for Mathematical Studies in Economics and Management Science, Northwestern University, Evanston.
- Shapley, L. and Shubik, M. [1969], "On Market Games," Journal of Economic Theory, Vol. 1, No. 2, pp. 9-25.
- Shitovitz, B. [1973], "Oligopoly in Markets with a Continuum of Traders," Econometrica, Vol. 41, No. 3.