

自然言語によるコミュニケーションのゲーム理論的分析

橋田 浩一
電子技術総合研究所
〒305 つくば市梅園 1-1-4
hasida@etl.go.jp

コミュニケーションは、自律的行為者の間の相互作用であるから、ゲーム理論的に解析可能である。本稿では、送信者がメッセージを送り、送信者の意図した意味を受信者が推論しようとする、コミュニケーションの核の部分、意味ゲーム (meaning game) の概念を用いて定式化する。自然言語における意味ゲームでは何らかの意味で最適な均衡解が用いられるが、一般にはそれは Pareto 最適解とは限らず、コミュニケーションを行なう行為者の間の共有知識の詳細度に依存する準最適解と考えられる。

A Game-Theoretic Account of Natural Language Communication

HASIDA Kôiti
Electrotechnical Laboratory
1-1-4, Umezono, Tukuba, 305 Japan.
hasida@etl.go.jp

As interaction between autonomous agents, communication can be analyzed in game-theoretic terms. *Meaning game* is proposed to formalize the core of intended communication in which the sender sends a message and the receiver attempts to infer its meaning intended by the sender. It is argued that natural-language meaning games are played at some equilibria, which are not in general Pareto optimal but suboptimal depending upon the precision of the common knowledge between the communicating agents.

1 Introduction

Since communication is a game (an interaction among autonomous agents), it should be understandable in game-theoretic terms. In this paper we examine a fundamental aspect of linguistic communication, *nonnatural meaning*, from the point of view of game theory. We will argue that natural-language meaning games are played at their optimal equilibria. This optimality is Pareto optimality when there is sufficient common knowledge about the game. When there is less common knowledge, the optimality degrades accordingly.

Let I be the proposition that the sender S of a message intends to communicate semantic content c to a receiver R . Then I entails that S intends that R should both recognize c and believe I . This is the core of *nonnatural meaning* (Grice, 1957, 1969). Grice’s original notion of nonnatural meaning further entails (S ’s intention of) R ’s believing (when c is a proposition or a reference) or obeying (when it is an order or a request) c , but we disregard this aspect and concentrate on the core.

This restricted sense of nonnatural meaning implies that communication is inherently collaborative, because both S and R want that R should recognize c and I . S of course wants it, and so does R because it is beneficial in general to know what S intends to make R believe or obey. S might be lying or trying to mislead R , but even in such a case S is still intending to communicate a content c by way of making R recognize this intention. Even if R doubts S ’s honesty, R will try to know what c is, because knowing what c is would help R infer what S ’s hidden intent may be, among others. For instance, when S tells R that it is raining, R will learn that S wants to make R believe that it is raining. R would do so even if R knew that it is not raining. Even if R were insincere and misconstrued S ’s message on purpose,¹ the nonnatural meaning is still properly conveyed, because otherwise the intended misconstrual would be impossible.

The present study concerns nonnatural meaning in the restricted sense, which is the core of intentional communication. Lies, ironies, indirect speech acts, and so forth (Perrault, 1990; Perrault & Allen, 1980) all share this core. Our understanding about it will hence help us understand basic workings of natural communication systems and shed some light on the design of communicating artificial agents.

¹If R is sincere and unintentionally misunderstands, that is just a failure of sharing the same context with S .

2 Communication Games

Communication has been discussed in the game-theory literature. A *signalling game* consists of sender S ’s sending a *message* (or a *signal*) to receiver R and R ’s doing some *action* in response to that message. Here S knows something that R did not know before receiving the message. This is formulated by assuming that S belongs to some *type*, which S knows but R does not know at first. Let T be the set of the types, P be the probability distribution over T . Let M be the set of the messages and A be the set of R ’s possible actions. Finally, let U_X be the *utility function* for player X . $U_S(t, m, a)$ and $U_R(t, m, a)$ are real numbers for $t \in T$, $m \in M$ and $a \in A$. A signalling game with $T = \{t_1, t_2\}$, $M = \{m_1, m_2\}$ and $A = \{a_1, a_2\}$ is illustrated by a *game tree* as shown in Figure 1. Here the

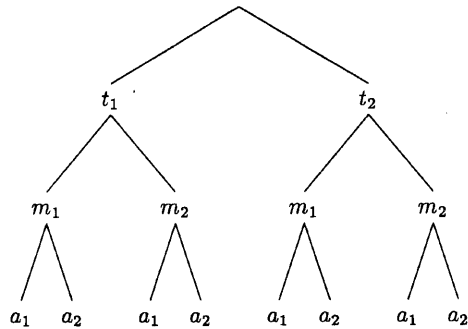


Figure 1: A signalling game.

game proceeds downwards. The top branch is the situation’s initial “choice” of S ’s type according to P , the middle level is S ’s decision on which message to send, and finally the bottom layer is R ’s choice of his action. When R has just received m_i ($i = 1, 2$), he does not know whether the game has been played through t_1 or t_2 .

Let σ_S and σ_R be S ’s and R ’s *strategies*,² respectively. That is, $\sigma_S(m|t)$ is the conditional probability of S ’s sending message m provided that she is of type t , and $\sigma_R(a|m)$ the conditional probability of R ’s doing action a provided that he has received m . The combination (σ_S, σ_R) of strategies is an *equilibrium*³ of a signalling game when σ_S and σ_R are the optimal responses to each other; that is, when σ_X maximizes X ’s

²Or *mixed strategies*, which are probability distributions over the *simple strategies* (actions).

³Or *complete Bayesian equilibrium*, in a more precise, technical term.

expected utility

$$\sum_{t,m,a} P(t) \sigma_S(m|t) \sigma_R(a|m) U_X(t, m, a)$$

given σ_Y , for both $X = S \wedge Y = R$ and $X = R \wedge Y = S$.

In typical applications of signalling game, T , M and A are not discrete sets as in the above example but connected subsets of real numbers, and S 's preference for R 's action is the same irrespective of her type. In this setting, S should send a costly message to get a large payoff. For instance, in job market signalling (Spence, 1973), a worker S signals her competence (type) to a potential employer R with the level of her education as the message, and R decides the amount of salary to offer to S . A competent worker will have high education and the employer will offer her a high salary. In mate selection (Zahavi, 1975), a deer S indicates its strength by the size of its antlers to potential mates R . A strong deer will grow extra large antlers to demonstrate its extra survival competence with this handicap.

Cheap-talk game is another sort of communication game. It is a special kind of signalling game where U_S and U_R do not depend on the message; that is, composing/sending and receiving/interpreting message are free of cost. In a cheap-talk game, S 's preference for R 's action must depend on her type for non-trivial communication to obtain, because otherwise S 's message would give no information to R about her type.

3 Meaning Game

Now we want to formulate the notion of *meaning game* to capture nonnatural meaning in the restricted sense discussed in Section 1. Let C be the set of semantic contents and P the probability distribution over C . $P(c)$ is the probability that S intends to communicate semantic content c to R . As before, M is the set of the messages. A meaning game addresses a *turn* of communication $\langle c_S, m, c_R \rangle$, which stands for a course of events where S , intending to communicate a semantic content c_S , sends a message m to R and R interprets m as meaning c_R . $c_S = c_R$ is a necessary condition for this turn of communication to be successful. It seems reasonable to assume that the success of communication is the only source of positive utility for any player.

So a meaning game might be a sort of signalling game in which S 's type stands for her intending to communicate some semantic content, and R 's action is to infer some semantic content. That

is, both T and A could be simply regarded as C . Strategies σ_S and σ_R are defined accordingly.

In a simple formulation, the utility function U_X of player X would thus be a real-valued function from $C \times M \times C$ (the set of turns). It would be sensible to assume that $U_X(c_S, m, c_R) > 0$ holds only if $c_S = c_R$. U_X reflects the grammar of the language (which might be private to S or R to various degrees). The grammar evaluates the (computational, among others) cost of using content-message pairs. The more costly $\langle c_S, m \rangle$ and $\langle m, c_R \rangle$ are, the smaller $U_X(c_S, m, c_R)$ is. The notion of equilibria in a meaning game is naturally derived from that in a signalling game.

If the players want something like *common knowledge*,⁴ however, meaning games are not signalling games. This is because $c_S = c_R$ is not a sufficient condition for the success of communication in that case. U_X should then depend on not just c_S , m , and c_R , but also the players' nested beliefs about each other. We will come back to a related issue in Section 4.

Note also that the typical instances of a meaning game in natural language communication is not like the typical applications of signalling game such as those mentioned before, even if meaning games are special sort of signalling games. Meaning games in natural language would normally involve discrete sets of semantic contents and messages.

Natural-language meaning games are not cheap-talk games, either, because we must take into consideration the costs of content-message pairs. It is not just the success of communication but also various other factors that account for the players' utility. S and R hence do not just want to maximize the probability of successful communication.

To illustrate a meaning game and to demonstrate that meaning games are not cheap-talk games, let us consider the following discourse of English.

- (1) u_1 : Fred scolded Max.
 u_2 : He was angry with the man.

The preferred interpretation of 'he' and 'the man' in u_2 are Fred and Max, respectively, rather than the contrary. This preference is accounted for by the meaning game as shown in Figure 2. In this game, Fred and Max are semantic contents, and

⁴People have common belief of proposition p when they all believe p , they all believe that they all believe p , they all believe that they all believe that they all believe p , and so on, ad infinitum. This common belief is a common knowledge when all these beliefs are factual.

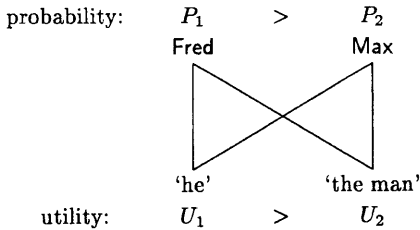


Figure 2: A meaning game about references of NPs.

‘he’ and ‘the man’ are messages.⁵ We have omitted the top branch representing the situation’s “choice” among semantic contents. Also, the nodes with the same label are collapsed to one. *S*’s choice goes downward and *R*’s choice upward, without their initially knowing the other’s choice. The complete bipartite connection between the contents and the messages means that either message can mean either content grammatically (without too much cost).

P_1 and P_2 are the prior probabilities of references to Fred and Max in u_2 , respectively. Since Fred was referred to by the subject and Max by the object in u_1 , Fred is considered more salient than Max in u_2 . This is captured by assuming $P_1 > P_2$. U_1 and U_2 are the utility (negative cost) of using ‘he’ and ‘the man,’ respectively.⁶ Utilities are basically assigned to content-message pairs, but sometimes it is possible to consider costs of messages irrespective of their contents. We assume $U_1 > U_2$ because ‘he’ is less complex than ‘the man’ both phonologically and semantically; ‘he’ is not only shorter than ‘the man’ but also, more importantly, less meaningful in the sense that it lacks the connotation of being adult which ‘the man’ has.

There are exactly two equilibria entailing 100% success of communication, as depicted in Figure 3 with their expected utilities E_1 and E_2 apart from the utility of success of communication.⁷ In the left equilibrium, *S* always means Fred by saying ‘he’ and Max by saying ‘the man,’ and *R* interprets ‘he’ as meaning Fred and ‘the man’ as meaning Max. In the right, *S* always means Fred by saying ‘the man’ and Max by saying ‘he,’ and *R* interprets ‘he’ as meaning Max and

⁵Perhaps there are other semantic contents and messages.

⁶For the sake of simplicity, here we assume that U_S and U_R are equal. See Section 4 for more precise discussion.

⁷Common belief about the communicated content is always obtained in both cases. So the current discussion does not depend on whether the success of communication is defined by just $c_S = c_R$ or common belief thereof.

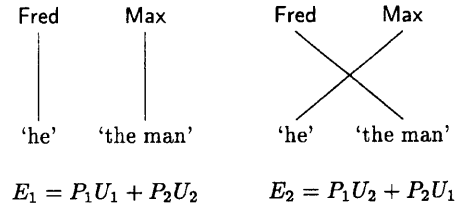


Figure 3: Two equilibria of the game in Figure 2.

‘the man’ as meaning Fred. Since $P_1 > P_2$ and $U_1 > U_2$ imply $E_1 - E_2 = (P_1 - P_2)(U_1 - U_2) > 0$, the former equilibrium is *Pareto superior* to the latter; i.e., the former is preferable for both players *S* and *R* to the latter. This explains the preference in (1). It is straightforward to generalize this result for cases with more than two contents and messages: A more salient content should be referred to by a lighter message when the combinations between the contents and the messages are complete. A general conjecture we might draw from this discussion is the following, where we say an equilibrium is *Pareto optimal* iff no other equilibrium is Pareto superior to it.

- (2) A solution of a natural-language meaning game is a Pareto optimal equilibrium.

Here a *solution* is a combination of strategies actually adopted by the players *S* and *R*. It is not difficult to formally prove that, for instance, a meaning game between logically omniscient agents is played at the unique Pareto optimal equilibrium if any (Parikh, 1990). The truth of (2) is an empirical issue, however, because it concerns humans, who are resource-limited agents. In Section 4 we will revise (2) in the light of other empirical evidence.

Incidentally, we have derived a central point of the syntactic pronoun resolution algorithm of Hobbs (1978) and of centering theory (Joshi & Weinstein, 1981; Kameyama, 1986; Walker, Iida, & Cote, 1994; Grosz, Joshi, & Weinstein, 1995), which seeks to explain anaphora in natural language. Centering theory considers list $Cf(u_i)$ of *forward-looking centers*, which are the semantic entities *realized*⁸ in u_i , where u_i is the i -th utterance. The forward-looking centers of utterance u are ranked in $Cf(u)$ according to their saliences. In English, this ranking is determined by grammatical functions of the expressions in the utterance, as below.

⁸A linguistic expression *realizes* a semantic content when the former directly refers to the latter or the situation described by the former involves the latter. For instance, after an utterance of ‘a house,’ ‘the door’ realizes the house referred to by ‘a house.’

subject > direct object > indirect object > other complements > adjuncts

The highest-ranked element of $Cf(u)$ is called the *preferred center* of U and written $Cp(u)$. The *backward-looking center* $Cb(u_i)$ of utterance u_i is the highest-ranked element of $Cf(u_{i-1})$ that is realized in u_i . $Cb(u)$ is the entity which the discourse is most centrally concerned with at u .

Centering theory stipulates the following rule.

- (3) If an element of $Cf(u_{i-1})$ is realized by a pronoun in u_i , then so is $Cb(u_i)$.

In (1), $Cb(u_2) = \text{Fred}$ because $Cf(u_1) = [\text{Fred}, \text{Max}]$, if either ‘he’ or ‘the man’ refers to Fred. Then rule (3) predicts that Fred cannot be realized by ‘the man’ if Max is realized by ‘he’ — the same prediction that we derived above. Moreover, (3) itself is a special instance of our above observation that a more salient content should be referred to by a lighter message, provided that the backward-looking center is particularly salient.

(3) is common in all versions of centering theory, but of course there are further details of the theory, which vary from one version to another. To derive all of the correct principles in a unified manner requires further extensive study. See Hasida, Nagao, and Miyata (1995) and Hasida (1996) for more on a game-theoretic account of anaphora in natural language.

4 Suboptimal Solution

We have so far assumed implicitly that S and R have common knowledge about (the rule of) the game (that is, P , U_S and U_R). This assumption will be justified as a practical approximation in typical applications of signalling games (and cheap-talk games). For instance, there may well be a body of roughly correct, stable common-sense knowledge about the correlation between the competence of workers and the degree of effort they make to have higher education, about how much an employer will offer to an employee with a certain competence, and so on.

However, common knowledge on the game is harder to obtain in natural-language meaning games, because the game lacks the stability of the typical signalling games mentioned above. A natural-language meaning game is almost equivalent to the context of discourse, which changes dynamically as the discourse unfolds. Hasida et al. (1995) argue that the communicating agents pretend to have common knowledge, but there

are important circumstances where this claim fails.

For example, consider the meaning game depicted in Figure 4. Here S wants to refer to

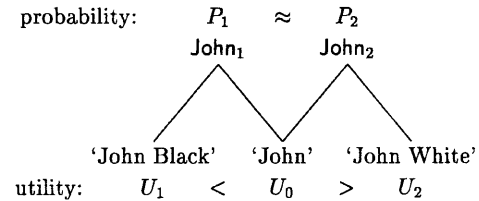


Figure 4: A meaning game where precision of common knowledge is critical.

either $John_1$ or $John_2$, $John_1$ may be called either ‘John Black’ or just ‘John,’ and $John_2$ may be called either ‘John White’ or just ‘John.’ As indicated in the figure, let us assume that the prior probabilities P_1 and P_2 of references to $John_1$ and $John_2$, respectively, are nearly equal. It is reasonable to assume ‘John’ incurs a smaller cost of utterance and interpretation than ‘John Black’ or ‘John White.’ So we assume $U_0 > U_1$ and $U_0 > U_2$. Let us further assume $U_1 \approx U_2$.

Then one of the equilibria shown in Figure 5 is Pareto optimal. In the left equilibrium, the

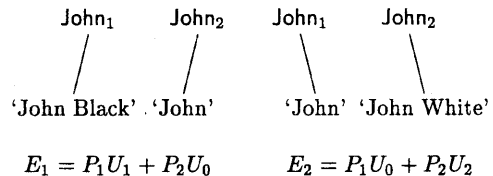


Figure 5: Two candidates for the Pareto optimum equilibrium of the game in Figure 4.

sender S means $John_1$ by saying ‘John Black’ and $John_2$ by saying ‘John,’ and the receiver R always interprets ‘John’ as meaning $John_2$. In the right, S means $John_1$ by saying ‘John’ and $John_2$ by saying ‘John White,’ and R interprets ‘John’ as meaning $John_1$. As a matter of course, in both cases R always interprets ‘John Black’ as meaning $John_1$ and ‘John White’ as meaning $John_2$. Note that these equilibria guarantees success of communication in that common knowledge is established about which semantic content has been communicated. As before, the associated expected utilities apart from success of communication are shown below the equilibria in Figure 5.

There is one more equilibrium, shown in Figure 6, which guarantees successful communication.

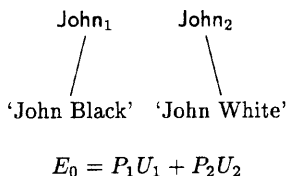


Figure 6: The solution of the game in Figure 4.

tion. In this equilibrium, S says 'John Black' to mean $John_1$ and 'John White' to mean $John_2$. R always interprets 'John Black' as meaning $John_1$ and 'John White' as meaning $John_2$, of course, but we do not care how he interprets 'John.'⁹

$E_0 < E_1$ and $E_0 < E_2$ follow from $U_1 < U_0$ and $U_2 < U_0$, respectively. So the equilibria in Figure 5 are Pareto superior to the one in Figure 6. However, the latter equilibrium seems to be the solution which people tend to settle on in this meaning game, contradicting (2). Why do people avoid the Pareto superior equilibria?

Intuitively speaking, this seems to be because 'John' is ambiguous. When S says 'John' to mean $John_1$, R might interpret it as meaning $John_2$, or vice versa. So the success of communication is not guaranteed if S uses 'John.' But why does R face the ambiguity? Why does R not know which equilibrium S is committing herself to?

Common knowledge is lacking by which S and R can choose the same equilibrium in Figure 5. If S and R had common knowledge about the exact detail of the whole game (P_i and U_i) and one of the two equilibria were Pareto superior to the other (hence being the unique Pareto optimal equilibrium) in the common knowledge, then the players would be able to commonly adopt that equilibrium, because in that case they would commonly know it to maximize their expected utilities. Even when the players commonly know less than the whole game, one of the equilibria in Figure 5 can be the solution of the game. For example, let us suppose that each player's belief about the game may vary from one occasion to another but in every occasion the belief entails that the equilibria of the game are the three in Figure 5 and Figure 6 and the former two are Pareto superior to the latter. Suppose further that all of this is commonly known. Even in such a situation, it is possible that both players are commonly known to believe more often that the first equilibrium in Figure 5 is Pareto

⁹So this is actually an equivalence set of equilibria rather than a single equilibrium.

superior to the second. In that case, the players should be able to commit in common to the first equilibrium, because that is commonly known to maximize the players' expected (subjective) utility.

In short, the ambiguity of 'John' arises from lack of common knowledge of any asymmetry between the two equilibria in Figure 5 regarding their expected utilities. So let us revise (2) as follows:

- (4) The solution (if any) of a natural-language meaning game is Pareto optimal among the equilibria which are *commonly Pareto comparable* with every other equilibrium.

We say two different equilibria are commonly Pareto comparable when one of them is commonly believed to be Pareto superior to the other. Note that, in the current example, only the equilibrium in Figure 6 is commonly Pareto comparable with the other equilibria, if neither equilibrium in Figure 5 is commonly believed to be Pareto superior to the other. So (4) correctly predicts that the equilibrium in Figure 6 tend to be employed in actual discourse. We must further revise (4) for the case where no equilibrium is commonly Pareto comparable with the other equilibria, but that is a matter for future study.

Incidentally, it is impossible to reformulate the current meaning game as another type of game in which the profile in Figure 6 is the only Pareto optimal equilibrium. The two profiles in Figure 5 are equilibria due to the sufficient condition for equilibrium pointed out by Aumann and Brandenburger (1995), despite the absence of common knowledge about the exact details of the game. Granted that they are equilibria, they are Pareto superior to the third one because $U_0 > U_1$ and $U_0 > U_2$.

In order to deal with cases such as the present example, metareasoning about the epistemic conditions — in particular about common knowledge — of the players is necessary, in addition to standard tools of game theory. Approaches which deliberately exclude common knowledge (Gmytrasiewicz & Rosenschein, 1993; Durfee, Gmytrasiewicz, & Rosenschein, 1994; Gmytrasiewicz & Durfee, 1995) would also need some extra machinery to account for the indeterminacy and ambiguity arising from lack of common knowledge.

5 Concluding Remarks

We have proposed the notion of meaning game to formalize nonnatural meaning, and investigated its solution in the light of empirical evidences concerning natural-language communication. The solution of a natural-language meaning game is a Pareto optimal equilibrium in basic cases, but in general the account of optimality must take common knowledge into consideration, as stated in (4). This observation would be significant not only to the investigation of natural language, or nonnatural meaning, in particular, but also to coordination in general. For instance, the algorithm to find focal points (Kraus & Rosenschein, 1992; Fenster, Kraus, & Rosenschein, 1995) could be refined to take our discussion into consideration. It would be a much more complicated but also a much more interesting task to formulate principles like (4) for games other than meaning game, for example, games in which agents do not in general collaborate.

Acknowledgment

The author would like to thank Jerry R. Hobbs and Megumi Kameyama for stimulating discussion and helpful comments. The remaining errors, if any, are the author's.

Reference

- Aumann, R. J., & Brandenburger, A. (1995). Epistemic Conditions for Nash Equilibrium. *Econometrica*, 63(5), 1161–1180.
- Durfee, E. H., Gmytrasiewicz, P., & Rosenschein, J. S. (1994). The Utility of Embedded Communications and the Emergence of Protocols. In *Proceedings of AAAI '94 Workshop on Planning for Interagent Communication*.
- Fenster, M., Kraus, S., & Rosenschein, J. S. (1995). Coordination without Communication: Experimental Validation of Focal Point Techniques. In *Proceeding of the First International Conference on Multi-Agent Systems*, pp. 102–108. San Francisco.
- Gmytrasiewicz, P., & Rosenschein, J. S. (1993). The Utility of Embedded Knowledge-Oriented Actions. In *Proceedings of the 12th International Workshop on Distributed Artificial Intelligence*.
- Gmytrasiewicz, P. J., & Durfee, E. H. (1995). A Rigorous, Operational Formalization of Recursive Modeling. In *Proceedings of the First International Conference on Multi-Agent Systems*, pp. 125–132.
- Grice, H. P. (1957). Meaning. *Philosophical Review*, 66, 377–388.
- Grice, H. P. (1969). Utterer's Meaning and Intentions. *Philosophical Review*, 68(2), 147–177.
- Grosz, B. J., Joshi, A. K., & Weinstein, S. (1995). Centering: A Framework for Modeling the Local Coherence of Discourse. *Computational Linguistics*, 21(2), 203–225.
- Hasida, K. (1996). Issues in Communication Game. In *Proceedings of the 16th International Conference on Computational Linguistics*, pp. 531–536. Copenhagen.
- Hasida, K., Nagao, K., & Miyata, T. (1995). A Game-Theoretic Account of Collaboration in Communication. In *Proceedings of the First International Conference on Multi-Agent Systems*, pp. 140–147. San Francisco.
- Hobbs, J. R. (1978). Resolving Pronoun Reference. *Lingua*, 44, 311–338.
- Joshi, A. K., & Weinstein, S. (1981). Control of Inference: Role of Some Aspects of Discourse Structure — Centering. In *Proceedings of the 7th International Joint Conference on Artificial Intelligence*, pp. 385–387.
- Kameyama, M. (1986). A Property-Sharing Constraint in Centering. In *Proceedings of the 24th Annual Meeting of ACL*, pp. 200–206.
- Kraus, S., & Rosenschein, J. S. (1992). The Role of Representation in Interaction: Discovering Focal Points Among Alternative Solutions. In Werner, E., & Demazeau, Y. (Eds.), *Decentralized A.I. 3: Proc. of the Third European Workshop on Modelling Autonomous Agents in a Multi-Agent World*, pp. 147–165. North-Holland, Amsterdam.
- Parikh, P. (1990). Situations, Games, and Ambiguity. In Cooper, R., Mukai, K., & Perry, J. (Eds.), *Situation Theory and Its Applications (Vol. 1)*, pp. 449–469. CSLI Publications, Stanford, CA.

- Perrault, C. R. (1990). An Application of Default Logic to Speech Act Theory. In Cohen, P. R., Morgan, J., & Pollack, M. E. (Eds.), *Intentions in COMMUNICATION*, pp. 161-185. MIT Press.
- Perrault, C. R., & Allen, J. F. (1980). A Plan-Based Analysis of Indirect Speech Act. *American Journal of Computational Linguistics*, 6(3-4), 167-182.
- Spence, A. M. (1973). Job Market Signaling. *Quarterly Journal of Economics*, 87, 355-74.
- Walker, M., Iida, M., & Cote, S. (1994). Japanese Discourse and the Process of Centering. *Computational Linguistics*, 20(2), 193-232.
- Zahavi, A. (1975). Mate Selection — A Selection for a Handicap. *Journal of Theoretical Biology*, 53, 205-214.