# SARL: A Computational Reinforcement Learning Model with Selective Attention

**Maurice Grinberg (mgrinberg@nbu.bg), Evgenia Hristova (ehristova@cogs.nbu.bg)**
Central and East European Center for Cognitive Science, New Bulgarian University,
21 Montevideo Street, 1618 Sofia, Bulgaria

## Abstract

A model relating eye-movements and decision making is proposed focused on the iterated prisoner's dilemma game. Its main aim is to model previous experiments with eye-tracking recordings which show that participants attend to only a small part of the game payoff information. The model presented generates eye-movements based on two main mechanisms. The first takes into account the importance of the information attended with respect to the decision making process and while the second takes into account the variability of the information attended. The model is a discrete dynamical system which integrates learned selective attention with move choice. The model is found to reproduce fairly well the sensitivity to the payoff structure of the game and the attendance to payoffs found in experiments with human subjects. These results seem to be a promising first step in explaining the impact of partial and selective information acquisition in the prisoner's dilemma.

**Keywords:** eye-movements models; eye-tracking; selective attention; decision making; Prisoner's Dilemma.

## Goals of the Present Work

It is the main goal of the current work to present a model in which decision making is integrated with information acquisition in a single integrated mechanism. The model is tested against experimental data if it is able to account for both behavioral and information acquisition data.

The model we present in this paper (SARL) is based on a model proposed and used earlier (Hristova & Grinberg, 2005a) for describing playing in iterated Prisoner's Dilemma (PD) games with different payoffs and cooperation indexes. The new model also uses the expected subjective utility framework combined with reinforcement learning; however, it also incorporates selective attention mechanisms. The name of the model SARL comes from Selective Attention and Reinforcement Learning. This work is part of a larger effort of clarifying the cognitive processes underlying PD game playing (e.g. Hristova & Grinberg, 2005b, 2008; Grinberg et al., 2005). In this series of research not only behavioral data was gathered, but also information acquisition data (using eye-tracking recordings and computerized process tracing system). The results obtained in these studies show that players do not use all the information available and that there is dependence between the playing strategy and the information acquisition patterns. The experimental results show that under the condition of playing different PD, participants sometimes miss completely the payoff structure of the game which automatically makes any model relying on full information about the payoffs useless. The model SARL has been developed to account for such situations and be able to describe and predict eye-movement and behavioural data at the same time and provide explanations about the relation between them.

## The Prisoner's Dilemma Game

The Prisoner's dilemma (PD) game is one of the most extensively studied social dilemmas. PD is a two-person game. The payoff table for this game is presented in Figure 1. In PD games, the players simultaneously choose their moves – C (cooperate) or D (defect) – without knowing their opponent's choice.

In order to be a Prisoner's dilemma game, the payoffs should satisfy the inequalities $T > R > P > S$ and $2R > T+S$. Due to the payoff structure of this game a dilemma appears – there is no obvious best move. On one hand, the D choice is dominant for both players – i.e. each player gets larger payoff by choosing D than by choosing C no matter what the other player chooses. On the other hand, the payoff for mutual defection (P) is lower than the payoff S if both players choose their dominated C strategies (R for each player).

|  | Player II |  |
|---|---|---|
|  | C | D |
| C | R, *R* | S, *T* |
| D | T, *S* | P, *P* |

Figure 1: Payoff table for the PD game. In each cell the comma separated payoffs are the Player I's and Player II's payoffs, respectively.

As PD game is used as a model for describing social dilemmas and studying the phenomena of cooperation, there is a great interest in the conditions that could promote or hinder cooperation. There are many factors, identified experimentally, that influence the cooperation rate in playing iterated PD. Among them are framing (or the way of describing the game to the participants in a experiment), players' goals and motivation, opponent strategy, etc. (Colman, 1995; Sally, 1995).

One important characteristic which accounts for the relation between payoff structure and cooperation in PD is a quantity called cooperation index (CI) which was introduced by Rapoport and Chammah (1965). It is calculated using the equation: $CI = (R–P)/(T–S)$. CI may vary from 0 to 1 (see Figure 2) and it is positively correlated with the percentage of C choices. An advantage in using such an index for predicting cooperation is that the

probability of a C choice should depend not on the payoffs (T, R, P, S) individually but rather on the ratios of their differences.



Figure 2: Examples of PD games with different CI. The first game has a CI=0.9, the second one has CI=0.1.

## Information Acquisition in Decision Making

Information acquisition studies explore what information is sought, how long the information is examined, the sequence of acquisition, and the amount of information acquired. The data made available are essential for studying the decision making process as a process taking place in time and based on a specific sequence of information acquisition. The patterns of information acquisition impose constraints on the possible strategies of information evaluation and decision making. Taking this into account, the importance of studying information acquisition patterns is emphasized in numerous research papers (see e.g. Einhorn & Hogarth, 1981; Johnson, Payne, & Bettman, 1988).

The fact that humans use part of the information and still behave in a consistent manner, shows the importance of the decision making process based on incomplete information – how and what information is gathered, how and in what order it is evaluated and processed to reach a decision. In order to understand these processes, the models of human judgment and decision making, including game playing models, should be built on what we know about the real mind's capacities and limitations.

Many studies in judgment and decision making are aimed at the development and testing of models that deal with evaluation and use of information. In many of them it is implicitly assumed that information is already available and judgment and choice are considered on the basis of information which is already given. However, there is strong evidence that the information acquisition process is part of the decision making process and thus can influence it (e.g., Einhorn & Hogarth, 1981; Lohse & Johnson, 1996).

Information acquisition studies give us information not only about the way in which reduction of information occurs (if it is the case), but also on the pattern or temporal order of acquisition. Such data provide important constraints for any decision making model.

Eye-tracking is one of the most popular methods for studying information acquisition. It is considered that the pattern of eye movements can provide objective and quantitative evidence on what is being processed at the moment. Many studies investigate cognitive processes as reading, visual search, scene perception and other

information processing tasks using eye-movement recordings (Rayner, 1992, 1998).

The results of all empirical studies stress the essential role of the amount and type of information on the decision making process.

## The SARL Model

The model we propose here is a modification of the model proposed in Hristova & Grinberg (2005a) and has as goal to incorporate elements of active and selective attention. SARL can be viewed as based on the general framework of the subjective utility theory (Schoemaker, 1982) but with dynamic determination of the utilities and of the expectations about the other player move probabilities based on an information acquisition mechanism.

Models, similar in spirit have been used in different contexts by Antonides (1994) and Piunti et al. (2007). The latter approach is interesting in combining a simple subjective expected utility model with affections which control the speed of learning in the model. These models however lack any selective attention mechanisms.

In order to benefit from the continuity in the two models, we briefly present the model used in Hristova & Grinberg (2005a). It can be defined as follows:

$$V(C) = w_{CC} \, Pff(CC) \, P_{op}(C) + w_{CD} \, Pff(CD) \, P_{op}(D) \quad (1)$$

$$V(D) = w_{DC} \, Pff(DC) \, P_{op}(C) + w_{DD} \, Pff(DD) \, P_{op}(D) \quad (2)$$

where:
- $P(C)$ is the probability of move C for the player
- $V(C)$ and $V(D)$ are the values of moves C and D;
- $Pff(CC)$, $Pff(CD)$, $Pff(DC)$, and $Pff(DD)$ are the current payoffs $R$, $S$, $T$ and $P$, respectively;
- $P_{op}(C)$ is the predicted probability for the opponent to play C.

The quantities $w_{CC}$, $w_{CD}$, $w_{DC}$, and $w_{DD}$ are weights that stand for the importance of the specific game outcome (CC, CD, DC or DD). These weights are computed as running averages of the payoffs received in the games with respective outcome and thus depend on previous payoffs:

$$[w_{XY}]_{new} = (1-\alpha) \, [w_{XY}]_{old} + \alpha \, Pff(XY), \quad (3)$$

where X and Y stand for C or D, respectively, $Pff(XY)$ is the received payoff corresponding to a game outcome XY and $0<\alpha<1$.

$P_{op}(C)$ is also calculated as a running average over the past opponents moves:

$$[P_{op}(C)]_{new} = (1-\beta) \, [P_{op}(C)]_{old} + \beta \, M_{op}, \quad (4)$$

where $0< \beta <1$ and $M_{op}$ is the opponent's move.

Because of the way the weights $w$'s and the probabilities $P_{op}(C)$ are calculated (see eqs. (3) and (4)) they are responsible for the context sensitivity of the model and are dynamically updated after each game. In eqs. (1) and (2), the use of the current-game payoffs ensures that the move will depend also on the game at hand and on its CI (Hristova & Grinberg, 2005a). This property is not available in typical reinforcement based models used in PD (Erev & Roth,

2001; Camerer et al., 2002; Macy & Flasche, 2002) in which the probability for a move is based only on the previous games.

The parameters for this part of the model are the averaging parameters ($\alpha$ and $\beta$) and the initial cooperation probability $P(C)$.

In SARL eqs. (1) and (2) become iterative and depend on the look patterns and thus on the payoffs 'perceived' by the model. They are generated on the basis of learned information about the game. Two are the main factors which underlie the mechanism of simulated eye-movements. The first is based on the idea that game outcomes which are wished and important (as measured by the magnitude of the weights $w$'s in eqs. (1) and (2)) attract attention to the payoffs related to them in a top-down fashion (e.g. the payoff R is more likely to be attended to the larger the weight $w_{CC}$ is). The second factor is related to the rate of change of the information relevant for the decision making – in our case the rate of change of the payoffs. This factor is regarded to be related to a bottom-up mechanism and is similar to some extent to the uncertainty minimization principle proposed in Hayhoe & Ballard (2005). In other words, information related to wishful outcome and information which is changing fast, tends to attract the attention of the model and its 'gaze'. The gaze pattern generator proceeds in two stages. Firstly, a transition matrix giving the transition probabilities between any two look zones $z_i$ and $z_j$, is calculated using the equation:

$$T(z_i, z_j) = a_{td} \, w_{xy}(z_j) + a_{bu} \, \Delta w_{xy}(z_j) + r(z_i, z_j) \qquad (5)$$

where $z_i$ and $z_j$ are the initial and final look zones for a saccade; $a_{td}$ and $a_{bu}$ are coefficients standing for the relative importance of top-down and bottom-up influences, respectively. The quantity $r(z_i, z_j)$ can encode spatial information about the relative position of $z_i$ and $z_j$, or can be related to the learning of looking patterns based on the received payoff compared to the expected one, or to the level of surprise related to predictions about the opponent's move or game outcome. In general it may be non-symmetric with respect to the two zones reflecting spatial asymmetries related, for instance, to cultural differences. For the simulations in this paper it is set to zero. For goal-directed behavior, like the one in PD, the top-down influences are expected to be larger than the bottom-up ones (Hayhoe & Ballard, 2005). The matrix $T(z_i, z_j)$ is updated after each game.

When the game starts based on the updated matrix an initial zone (payoff) is selected. It corresponds to a game outcome and with the respective moves of the opponents (e.g. R corresponds to a CC outcome). Depending on these moves (X for the player and Y for its opponent) the values of the moves C and D are updated using the equation:

$$V(X, n+1) = (1-\varepsilon) \, V(X, n) + \varepsilon \, w_{XY} \, Pff(XY) P_{op}(Y) \qquad (6)$$

where $\varepsilon$ is between 0 and 1; X and Y can be C or D, depending on the payoff attended.

Eq. (6) replaces eqs. (1) and (2) of the model in Hristova & Grinberg (2005a) and is used several times before reaching a decision within one and the same game. The initial value for the move values – $V(X, 0)$ – is calculated as the average of the respective weights times the corresponding probabilities for the opponent's move:

$$V(X, 0) = w_{XC} \, P_{op}(C) + w_{XD} \, P_{op}(D) \qquad (7)$$

where X can be C or D.

The stopping criterion for this deliberation process is the reaching of a threshold (see Roe et al., 2001) by the quantity (the softmax rule for $V(C)$):

$$V = \frac{1}{1 + e^{-k(V(C))-V(D))}} \qquad (8)$$

and the move is determined according to the rule:

$$M = \begin{cases} C & \text{if } V \geq \theta \\ D & \text{if } V < 1 - \theta \end{cases} \qquad (9)$$

where the threshold $\theta$ can depend on behavioural characteristics such as received payoff, verification of expectations, etc. In the present version $\theta$ is taken to be 0.8. The rule (9) is deterministic but can also be made probabilistic. It requires the reaching of one of two attractors in order to make a move C or D. The speed of reaching the attractors depends on the parameter $k$ in the softmax rule given in eq. (8).

## Simulations and Experiments

In this section, SARL predictions are compared to the results of an eye-tracking study with human participants. In the experiment and in the simulations one and the same set of PD games were used. We compare both information acquisition patterns and choices for human participants and the model.

### PD games used

A set of 100 PD different payoff matrices, containing an equal number of games with CI equal to 0.1, 0.3, 0.5, 0.7, and 0.9 was used in the experiment. The payoff matrices were randomly generated with the payoff magnitudes kept within certain limits. T was between 36 and 97 points (mean 69), R was between 29 and 95 points (mean 60.7), P was between 15 and 59 points (mean 32.5), and S was between 10 and 20 points (mean 15). The games were presented randomly with respect to their CI.

### Experimental procedure with human participants

**Game presentation** The game was presented in a formal and a neutral formulation to avoid other factors and contexts as much as possible. The terms 'cooperation' or 'defection' were not mentioned in the instructions or in the interface to further avoid influences other than the payoff matrix. On the interface, the moves were labeled in a neutral manner as '1' and '2. 'Subjects were not informed about the existence of CI. The game interface is presented in Figure 3. The

participants had to choose their move by mouse clicks on one of the button on the left (move '1' or move '2').

Participants were instructed to try to maximize their payoffs and not to compete with the computer. The payoffs were presented as points, which were transformed into real money and paid at the end of the experiment.

The information about the games played was fully available. After each game the participants got feedback about their and the computer's choices and payoffs in the current game. Participants could also permanently monitor the total number of points they have won and its money equivalent. They had no information about the computer's total score. This was made to prevent a possible shift of participants' goal – from trying to maximize the number of points to trying to outperform the computer.

**Opponent's strategy** Participants played PD games against a computer opponent. The computer player used a probabilistic version of the tit-for-tat strategy: it takes into account the two previous moves of the player and plays the same move with probability 0.8. The latter makes the computer's strategy harder to be discovered by the participant and in the same time allows the participant to cooperate if they wish (and be followed by the computer).

**Eye movements recordings** Eye movements were recorded using the ASL 501 eye-tracker with 60 Hz sampling rate. The light head mounted optics recorded the left eye movements. The centre of the pupil and the corneal reflection were tracked to determine the relative position of the eye. A magnetic head tracking equipment (Ascension Flock of Birds) was used in order to compensate for the possible head movements and ensure sufficient precision of the measurements. Integration of the eye movements and head movements made it possible to compute point of regard on the computer screen. Gaze tracker software for data recording and analysis was used.

The eye-tracker was calibrated using a 9-point grid. The accuracy of the gaze position record is about 0.5 degrees visual angle.

The game was presented on a 17" monitor (see Figure 3). Each box containing payoffs or moves occupied about 1 degree visual angle on the screen. The distance between two adjacent boxes was at least 1 degree visual angle to ensure stable distinction between eye-fixations belonging to respective zones.

**Participants and procedure** 40 participants (17 males, 23 females) with normal or corrected to normal vision took part in the eye-tracking experiment. All were university students with average age of 23 years.

After receiving instructions, participants were asked several questions to make sure they have understood the game. Each participant played the set of 100 PD games, described above. First 20 games are considered training and are not included in the subsequent analysis.

All participants were paid for their participation. The amount received depended on the points gained in the experiment.
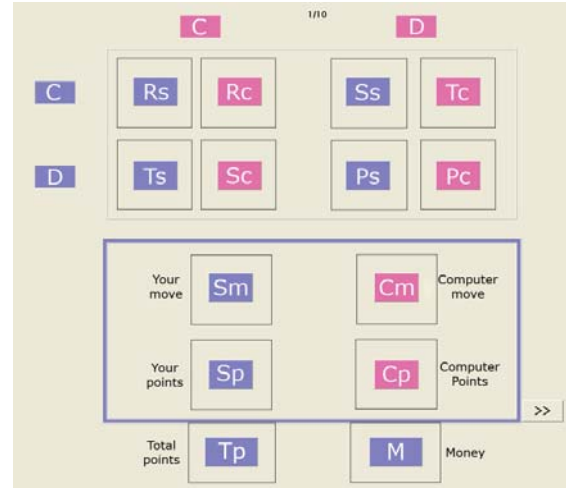


Figure 3: Game interface and areas of interest (AOIs) used in the experiment. The index 's' and 'c' denote 'subjects' and 'computer' respectively.

## Simulations

SARL was run 30 times using the same procedure as in the experiment. The same set of PD games was used and the model played against the same computer opponent. The first 20 games were considered as training games and were excluded from the analysis.

The averaging parameters $\alpha$ and $\beta$ were fixed to 0.5 and 0.4, respectively. The initial cooperation probability P(C) was set to 0.5 and the initial perceived probability of opponent playing C (Pop(C)) was also set to 0.5. These values are psychologically plausible as in the beginning of the game players probably do not posses clear preferences between choices or expectations about the play of the opponent. The move threshold $\theta$ is set to 0.8.

The speed of reaching the attractors depends on the parameter $k$ in the softmax rule given in eq. (8) and is fixed to 0.01.

## Comparison between SARL and Experimental Data

The model predictions and data from the experiment with human subjects are compared on number of measures. First, they are compared on the basis of playing choices, more specifically, number of cooperative choices. Next, we compared the eye-movement data from the experiment and model predictions about the zone attendances and transitions between zones.

## Cooperation

The first analysis compared the number of cooperative choices for each experimental condition (model or experiment) and each level of the CI. Repeated measures analysis of variance revealed that there is a significant main effect of CI on cooperation ($F(4, 272) = 11.23$, $p < 0.001$) and that there is no main effect of experimental condition (human experiment or model) ($F(1, 68) = 0.67$, $p = 0.41$) and that there is no interaction between the CI and the experimental condition ($F(4, 272) = 0.33$, $p = 0.85$) (see Figure 4).
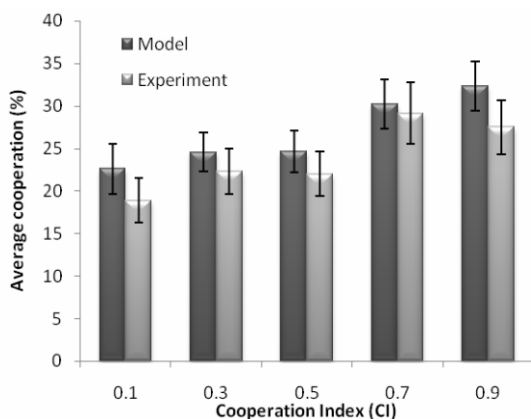


Figure 4: Comparison of the dependence of the rate of cooperation on the CI of the PD game between the theoretical and experimental results (error bars represent standard errors).

## Attention to different zones

The eye-tracking data were analyzed using the number of fixations in each AOI as a measure for attention paid to each AOI. This measure reflects the relative importance of the information presented in the AOI (Jacob & Karn, 2003). 8 areas on the screen that are important in studying information acquisition during PD game playing were defined (see Figure 3). Each Area of Interest (AOI) contains the box in which the information is presented and a small region around it.

The following AOIs were defined: 4 AOIs containing the participant's possible payoffs; 4 AOIs containing the computer's possible payoffs. The results showed that players do not pay equal attention to all available information. They look at their own payoffs more often than the computer's payoffs (2.71 fixations per game on all 4 AOIs with their payoffs and 1.14 fixations per game on all 4 AOIs with their opponent's payoffs). The low number of zone attendances per game indicates that players do not always attend to all information before making a decision. They even do not attend their own payoffs in each game (see Figure 5).

In the comparisons with the model only the 4 AOIs containing the participant's possible payoffs are analyzed referred to as Ts, Rs, Ps, and Ss.

The number of zone attendances for the model and for the eye-tracking data were compared for each zone using independent samples t-test. The tests showed no significant differences (all $p>0.05$) between the model and the experiment for each zone (see Figure 5).
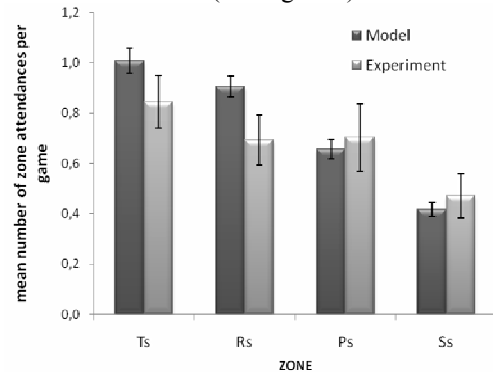


Figure 5: Comparison of the number of fixations per zone obtained with the model and in the experiment.

## Transitions between different zones

As a next step in the analysis, the number of transitions between zones containing participant's possible payoffs was considered. Transitions are assumed to indicate the comparisons made between the payoffs. Averaged data for all participants in the eye-tracking study and for the model predictions is presented in Figure 6. The players made more transitions between their bigger payoffs (Ts and Rs; and Ts and Ps); however, in general the number of transitions is pretty low.
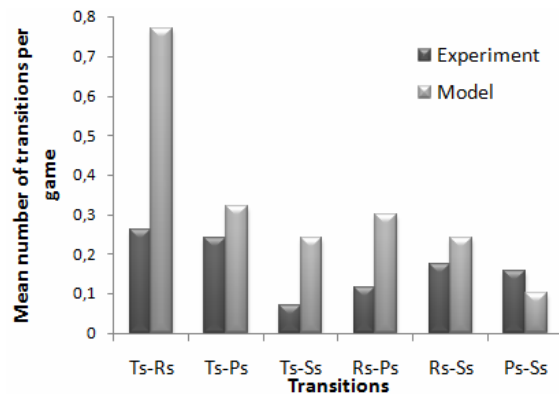


Figure 6: Comparison of the number of transitions for zones containing player's possible payoffs (Ts, Rs, Ps, and Ss) for the experimental and the model data.

## Discussion and Conclusion

The paper presented a model based on reinforcement learning and top-down selective attention mechanisms.

The comparison with eye-tracking and behavioural data, showed a reasonable agreement with respect to the average cooperation rates, the dependence of cooperation on CI, and the number of fixations in the payoff looking zones.

The number of transitions between the payoff zones, predicted by the model was larger than the experimental value. The latter could be explained by the fact that only four payoff zones (Ts, Rs, Ps, Ss) were taken into account and all of the model transitions are between them leading to a unrealistically large number of transition. The human players had access to many more zones when playing the game and only direct transition between zones were counted (with no intermediate fixation). The interface in the simulations didn't account for the opponent's payoffs and possible distraction outside the looking zones. Future versions of the simulations should account for these differences.

Despite these discrepancies, the results obtained show that a reinforcement learning model with selective attention as SARL can display a behaviour which is reasonably similar to the one displayed by human subjects. The latter seems to indicate that the model captures important features of decision making in iterated PD games. It is important to stress that the model presented meets the requirements set in the beginning: it displays behaviour (decision making) and information acquisition patterns simultaneously based on an integrated decision making mechanism.

At the same time, it is evident that more exploration of the dynamical properties of the model with respect to its parameters is needed. Application of the model on existing data and design of new experiment based on its predictions are also planned for the near future.

## Acknowledgements

## References

Antonides, G. (1994). Mental accounting in a Sequential Prisoner's Dilemma game. *J. Econ. Psychol.15*, 351-374.]

Camerer, C., Ho, T.-H., & Chong, J. (2002). Sophisticated EWA Learning and Strategic Teaching in Repeated Games. *J. Econ. Theory 104*, 137-88.

Colman, A. (1995)*. Game theory and its applications in the social and biological sciences*. Oxford: Butterworth-Heinemann Ltd.

Einhorn, H., & Hogarth, R. (1981). Behavioral decision theory: Processes of judgment and choice. *Annual Review in Psychology, 32*, 53-88.

Erev, I., Roth, A. (2001). Simple reinforcement learning models and reciprocation in the prisoner's dilemma game. In: G. Gigerenzer & R. Selten (Eds.), *Bounded rationality: the adaptive toolbox*. Cambridge, Mass. MIT Press.

Grinberg, M., Hristova, E., Popova, M. & Haltakov, V. (2005). Strategies in Playing Iterated Prisoner's Dilemma Game: An Information Acquisition Study. In: *Proceedings of the International Conference on Cognitive Economics*. Sofia, NBU Press.

Hristova, E. & Grinberg, M., (2005a) Investigation of Context Effects in Iterated Prisoner's Dilemma Game. In: Dey, A., Kokinov, B., Leake, D., Turner, R. (Eds.) *Modeling and Using Context. LNCS (LNAI), 3554*, Springer Verlag.

Hristova, E., & Grinberg, M. (2005b). Information acquisition in the iterated Prisoner's dilemma game: an eye-tracking study. *Proceedings of the 27$^{th}$ Annual Conference of the Cognitive Science Society*. Elbraum, Hillsdale, NJ.

Hristova, E., & Grinberg, M. (2008). Disjunction effect in prisoner's dilemma: Evidences from an eye-tracking study. In B. C. Love, K. McRae, & V. M. Sloutsky (Eds.), *Proceedings of the 30th Annual Conference of the Cognitive Science Society*. Austin, TX: Cognitive Science Society.

Jacob, R. & Karn, K. (2003). Eye tracking in human computer interaction and usability research: Ready to deliver the promise. In: Hyona, J., Radach, R., & Deubel, H. (Eds.), *The mind's eye: cognitive and applied aspects of eye movement research*. Elsevier Science BV.

Johnson, E., Payne, J., & Bettman, J. (1988). Information displays and preference revearsals. Organizational Behavior and Human Decision Processes, 42, 1-21.

Hayhoe M., Ballard, D. (2005). Eye movements in natural behavior. *Trends in Cognitive Science, 9* (4), 188-193.

Lohse, G. & Johnson, E. (1996). A comparison of two process tracing methods for choice tasks. *Organizational Behavior and Human Decision Processes*, *68*, 28-43.

Macy, M., & Flache, A. (2002). Learning dynamics in social dilemmas. *PNAS, 99,* 7229-7236

Piunti M., Castelfranchi, C. & Falcone, R. (2007). *Surprise as shortcut for Anticipation: clustering Mental States in Reasoning*. In Proceedings of the IJCAI07, Hyberabad, India.

Rapoport, A., & Chammah, A. (1965). *Prisoner's dilemma: a study in conflict and cooperation*. Univ. of Michigan Press.

Rayner, K., & Pollatsek, A. (1992). Eye movements and scene perception. *Canadian Journal of Psychology, 46*, 342-376.

Rayner, K. (1998). Eye movements in reading and information processing: 20 years of research. *Psychological Bulletin, 124,* 372-422.

Roe, R. M., Busemeyer, J. R., & Townsend, J. T. (2001). Multi-alternative decision field theory: A dynamic connectionist model of decision making. *Psychological Review*, *108, 370-392*.

Sally, D. (1995). Conversation and cooperation in social dilemmas. A meta-analysis of experiments from 1958 to 1992. *Rationality and Society, 7,* 58-92.

Schoemaker, P. (1982). The expected utility model: Its variants, purposes, evidence and limitations*. Journal of Economic Literature, 20*, 529-563.