# Lack of depth of iterative reasoning in non-interactive contexts

Ketti Mazzocco*, Paolo Cherubini°

*Department of Cognitive Sciences and Education, University of Trento*
*°Department of psychology, University of Milano-Bicocca*

Author for correspondence is PC:

Dipartimento di Psicologia, Università di Milano-Bicocca

Room 425, 4th floor, U6 Building

1, Piazza dell'Ateneo Nuovo, 20126 MILANO (Italy)

Telephone: ++39 02 6448 3811

Facsimile:  ++39 02 6448 3706

e-mail: paolo.cherubini@unimib.it

Abbreviations

| | |
|---|---|
| BCG | Beauty Contest Game |
| MMT | Mental Models Theory |

Abstract

In four Experiments we show that individuals tend to not iteratively pursue further consequences of an initial conclusion that they draw from an initial representation of a problem. This occurs with non-interactive tasks where the source of the difficulty cannot lie in an inability to adequately represent other actors' beliefs, actions, social values and goals. The difficulty at reasoning iteratively was previously mostly observed in interactive games, such as the beauty contest game, and partly attributed to bounded individual rationality. The present results, obtained in non-interactive games, support the bounded-rationality view, and further specify it by showing that lack of depth in iterative reasoning might be the direct result of a very basic cognitive tendency, originally illustrated by mental models theory.

It is common experience that people sometimes perform some actions in order to obtain an initial, predicted outcome, without realizing that that effect has further foreseeable consequences. Examples abound in interactive contexts, such as popular strategic board games (i.e. chess, checkers, kalaha, go), where one can miss the non-immediate – but quite deterministically foreseeable – consequences of a move. Similarly, one rarely considers that if she decides to drive on a crowded road, it will be more crowded, and slower (Schelling, 1978). In auctions, one's own bid might cause others to reevaluate what the item is worth and induce them to raise their bids (Roth and Ockenfels, 2002). This tendency to underestimate the non-immediate consequences of one's own actions can originate from an inability at reasoning iteratively (i.e., seeing the consequences of consequences in a chain of reasoning). In interactive contexts, this inability might be a factor in the development of inefficient markets and financial bubbles, as shown by typical performances in guessing games such as the "beauty contest game" (BCG; nicknamed after an example by John Maynard Keynes, 1936, but first studied by Nagel, 1995).

In the beauty contest game, $N$ decision-makers simultaneously choose a real number from the interval $I \equiv [0, 100]$. The winner is the decision-maker whose number is closest to $p$ times the mean of all chosen numbers (including her own), where $p \in (0, 1)$ is known. The winner receives a prize, whilst other decision-makers earn nothing. In case of a tie, the prize is split equally among those who have tied. The game has a Nash equilibrium in which all decision-makers choose zero. In BCG with large N – where each individual choice has a negligible effect on the aggregated mean – a rational player will not simply choose a random number or his favourite number, nor will she choose a number above $100p$, since it is dominated by $100p$. Moreover, if she believes that the other players are rational as well, she will not pick a number above $100p^2$, and if she again believes that the others are this rational, she will not pick a number above $100p^3$ and so on, until all numbers but zero are eliminated. Elementary algebra shows that, at each step of iteration, given $k$ the current choice provisionally attributed to the estimated mean of the other players' choices and $N$ the number of players, the player's rational choice at each step of iteration is $C = p(N\text{-}1)k/(N\text{-}p)$ (in the specific case where $N \to \infty$, $C=pK$, as in the large $N$ case above). For example,

3

with $p =.5$ and $N = 3$, in the first step a rational player might attribute randomness to the choice of the other players (that is, estimated mean $k=50$), and then opt for $C=20$; at the second step of iteration, she will attribute $k=20$ to the other players, and then she will think that she will be better off by choosing $C=8$; at the third step, she will choose 3; and so on, converging to zero.

The game-theoretic structure of the beauty-contest game allows analyzing the depth of players' reasoning, that is how many steps of iteration decision-makers actually apply in choosing their numbers. Previous studies have found that depth of reasoning is rather limited across a wide range of different pools of participants, sample sizes, or parameters $p$ (see Bosch-Domènech, García-Montalvo, Nagel, and Satorra, 2002; Camerer, 2003; Camerer, Ho, and Chong,, 2003, 2004; Duffy and Nagel, 1997; Güth, Kocher, and Sutter, 2002; Ho, Camerer, and Weigelt, 1998; Kocher and Sutter, 2005; Nagel, 1995, 1999a, 1999b; Weber, 2003). First round guesses are usually far from the equilibrium, either random choices (0 steps of iteration), or choices near $50p$ (depth 1), or a few choices near $50p^2$ (depth 2). The equilibrium 0 is chosen by very few participants (e.g., less than 10% in Grosskopf and Nagel 2007).[1] This performance might originate from individuals' intrinsic difficulty at reasoning iteratively; however, it might also occur if individuals were good at reasoning iteratively, but unable or unwilling to attribute the same capacity to other individuals. That is, a person that can see that the equilibrium point in the BCG is 0, but believes that other persons won't see it and won't choose 0, should not choose 0 herself. Both accounts have been put forth and both have been supported by empirical evidence (e.g., Grosskopf and Nagel, 2007, 2008; Bosch et al., 2002), and the two of them are not necessarily exclusive. However, in some recent studies Grosskopf and Nagel (2007, 2008), by using an $N=2$ version of the BCG – where the game turns to "the one who picks lowest, wins", and accordingly selecting 0 does not depend on the representation of the other player's choice – observed very few rational choices. The authors suggested that the typical behaviour in these sorts of

---

[1] In repeated BCGs chosen numbers decrease; however, learning of the equilibrium can be slow, sometimes is not attained, is affected by contextual factors, and is mostly caused by feedback on other participants' choices and outcome of each round (Duffy and Nagel, 1997; Grosskopf and Nagel, 2007). Learning in repeated BCGs is particularly slow when the individual choice affects appreciably the target number, e.g. with a small $N$ (e.g. Ho et al., 1998) or in some slightly different versions of the game (e.g., the "maximum" game, Duffy and Nagel, 1997).

interactive contexts is mostly caused by individual bounded rationality, more than by the tendency to attribute irrationality to others. Chou, McConnell, Nagel and Plott (2008) in a subsequent study showed that the source of poor performance in the two-person BCG might be scarce comprehension of the problem form. They showed that if oversimplified versions of the game are offered, people's behaviour is consistent with game theory predictions; yet, people have overwhelming difficulties at spontaneously modelling the standard version of the two persons BCG in a way that is isomorphic to its game-theoretic structure.

In this study we show some results that support and extend Grosskopf and Nagel's idea that individual cognitive constraints might affect performance in the BCG. People have intrinsic difficulties at following chains of iterative conclusions, and not only in interactive contexts: the same can be observed in non interactive contexts, where performance cannot be affected by difficulties at representing other people's behaviors. Some results to this effect were already reported by Cherubini and Johnson-Laird (2004), that used iterative problems based on the logic of predicate calculus. The authors presented problems with premises such as "Imagine a world in which there are four people: Anne, Beth, Carol, and Diane, and in which the following two assertions are true: 1) Everybody loves anyone who loves someone; 2) Anne loves Beth". Participants were then asked whether it followed that everyone loved Anne. Nearly all participants gave the correct "yes" response to this question. However, only a few participants correctly surmised that the answer to the following question "Does it follow that Carol loves Diane?" was also "yes". Answering properly to the latter question involves iterative application of the rule stated in the premises: since Anne loves Beth, everyone loves Anne; therefore Diane loves Anne; since Diane loves Anne, then everyone loves Diane, including Carol (the chain can be lengthened to the conclusion "everyone loves everyone"). In the negative version (Cherubini and Johnson-Laird, 2004, Experiment 2), the difficulty of pursuing iterative chains of reasoning increased: nearly all participants were defeated by problems such as "Everybody loves anyone who loves someone", "Anne does not love Beth", "Does it follow that Carol does not love Diane?" (correct response: ). The authors concluded that people have intrinsic difficulties at applying iteratively a premise, and interpreted those difficulties from the theoretical perspective of mental models theory (MMT; Johnson-Laird, 2001; Johnson-Laird and Byrne, 1991), currently one of the most

influential theories of human thinking and reasoning in cognitive psychology. The theory states that when people reason, they build a first, initial model of the premises, and draw a first conclusion; they barely look for alternative models, and thus it is difficult for them to realize that sometimes (i.e., in iterative reasoning chains) the first conclusion modifies the initial model, and conveys further conclusions. In the present study we corroborate and generalize Cherubini and Johnson-Laird's findings, by using different sorts of problems that involve set-based premises (Experiment 1a and 1b), a simple numerical problem strictly analogous to the BCG but in a non-interactive setting (Experiment 2), and an explicitly iterative numerical function (Experiment 3). Our aim is to specify a psychological mechanism that contributes to limit human rationality whenever iterative reasoning is required for optimal performance, and thus might also contribute to poor performance in the BCG and other interactive contexts, such as the financial markets. We think that knowledge concerning these sorts of psychological constraints might be of advantage for further detailing those economic models that allow for imperfect individual rationality among the factors needed for understanding collective behaviours.

# EXPERIMENT 1a

## Task

Imagine a box, and an experimenter that fills the box with five marbles described by two attributes: color (blue or red), and material (glass or plastic). That is, the experimenter can pick marbles from red glass ones, red plastic ones, blue glass ones, and blue plastic ones. She tells you that *more than half* of the marbles in the box are red, and *more than half* of the marbles in the box are glass marbles. She also assures you that – in filling the box – she conformed to the requirement R: *if at least one red glass marble is in the box, then at least one blue plastic marble must be in the box*. Now, she asks you if you can derive a certain conclusion concerning the maximal minimal number of red glass marbles and blue plastic marbles in the box – that is, if you are certain that there is at least one red glas marble, or at least two, and so on.

**Predictions**

A useful initial mental representation of the box depicts the minimal number of glass marbles and red marbles in it, that is three each. Because the marbles are five, most people should then be able to easily grasp the initial conclusion that there is necessarily at least one red glass marble in the box (Figure 1). This initial conclusion can be integrated to R, yielding the very easy *modus ponens* conclusion "there is at least one blue plastic marble in the box". If people insert this blue plastic marble in the initial representation of the box, they should realize that – in the light of that conclusion – the minimal overlap between red marbles and glass marbles is two, and conclude accordingly that there are necessarily at least two red glass marbles in the box (Figure 2).

--- Insert figure 1 about here ---

--- insert figure 2 about here ---

Otherwise, if people do not easily integrate the conclusion concerning the blue plastic marble into the initial mental representation of the box, they should stick to the initial conclusion "there is at least one red glass marble in the box", and will not grasp that the maximin conclusion "at least 2 red glass marbles in the box" follows from the premises.

**Control task**

Mentally overlapping two subsets of three marbles each out of a set of five marbles – as required by the second step of iteration above – might be intrinsically harder than overlapping two subsets of three marbles each out of a set of five marbles (as required by the first step of iteration). This can be checked for by a control task where the first premises, describing the red and glass marbles in the box, are kept the same, but the requirement R is missing; in its place, participants are explicitly told that at least one of the marbles in the box is a blue plastic one. That is, they are directly told what they have to initially infer in the experimental task. If difficulty in responding to the question concerning the red glass marbles is

caused by a difficulty in integrating one's own spontaneous conclusions into the initial representation of the problem, then this difficulty should not occur in the control task, and correct "at least 2 red glass marbles" responses should be more frequent in the control task than in the experimental task.

## METHODS

### *Participants*

Thirty-four participants (24 females; mean age 22.7, ranging 19-29) volunteered to take part in the experiment, without retribution. Seventeen of them were randomly allocated to the experimental condition, and 17 to the control condition. All participants where undergraduate students of psychology from the university of Padova. Some of them had taken courses in logic or in the psychology of thinking and reasoning.

### *Procedure and material*

Each participant was tested individually in a quiet room. The problem was part of a set comprising three other reasoning and decision problems, unrelated to the present study and administered in random order. The experimenter put a closed box (with marbles in it) and some example marbles on the desk, and verbally instructed the participant in accordance with the problems described in the *Task* or *Control task* paragraphs, repeating the instructions – when needed – in order to assure understanding of the premises. Instructions were also available to participants in written form. Questions concerning the maximin of red glass marbles and blue plastic marbles would have been linguistically awkward, and so they were formulated as follows:

*What can you say for certain regarding the number of red glass marbles in the box? That is, can't you conclude anything certain regarding their number, or can you establish that there is certainly at least one of them, or can you establish that there are certainly at least two of them, or anything else? Take your time to decide, and please explain your answer.*

*What can you say for certain regarding the number of blue plastic marbles in the box? That is, can't you conclude anything certain regarding their number, or can*

*you establish that there certainly is at least one of them, or can you establish that there certainly are at least two of them, or anything else? Take your time to decide, and please explain your answer.*

The order of the two questions was balanced across participants. Of course, participants given the control task were only asked the question concerning the red glass marbles.
Original instructions and questions were in Italian, and all participants were native Italian speakers.

## Results and analyses

Quantitative results surpassed our expectations. Of the 17 participants in the experimental condition, none stated that they were certain that there were at least two red glass marbles in the box. Two participants stated that they could not conclude anything for certain. One stated that a red glass marble was "more likely than a blue plastic marble". The remaining 14 (82%) participants concluded that they were certain that there was at least one red glass marble in the box, that is, the conclusion supported by the initial model of the problem. Of the 17 participants in the control task, two responded that they could not conclude anything for certain. The remaining 15 (88%) correctly concluded that they were certain that there were at least two red glass marbles in the box. A $\chi^2$ test run on the 2x2 contingency table obtained by crossing correct and incorrect responses to the red-glass-marble question in the two conditions shows that the difference in the distribution of responses is highly significant (exact $p$: less than one in ten millions). In the experimental task, by setting at .33 the probability that the predicted response "at least one red glass marble" was picked at random, predicted responses were significantly more frequent than chance (binomial test, $p$ <.0001).

Responses to the blue plastic marble question in the experimental task matched responses to the red glass marble question: 14 people responded – correctly – that they were certain that there was at least one blue plastic marble in the box; one participant said that a blue plastic marble was "less likely than a red glass marble"; the remaining two participants said that they could not conclude anything for certain. Correct responses were significantly greater than chance (binomial test, chance level set at .33, $p$<.0001).

Qualitative results were based on the explanations put forth by the participants. The 14 participants responding that there was at least a red glass marble in the box in the experimental task and the 15 participants concluding that there were at least two red glass marbles in the control task described – with different wordings, gestures, and sketches – that they had established the minimal overlap between three red marbles and three glass marbles. All participants which realized that there was at least one red glass marble in the experimental task, also realized that this conclusion, once integrated with R, endorsed the further conclusion that there was at least one blue plastic marble in the box. All of them were surprised – in the debriefing session – when they finally realized, upon explanation by the experimenter, that the consequence of the necessary blue plastic marble was that there was a second red glass marble in the box, and all admitted that they did not consider revising the initial overlap of the red and glass marbles in the light of the conclusion concerning the blue plastic marble.

**Discussion**

Results of Experiment 1a show that people do not tend to integrate their initial conclusion into the problem representation that allowed drawing that very same conclusion.  Participants  in the control condition and those participants in the experimental condition that managed to conclude that there was at least one blue plastic marble in the box had *exactly* the same information available. The only difference was that people in the control group were given a piece of information that people in the experimental group had to infer by themselves.  The difference between the two conditions was impressive: no one in the experimental condition revised their initial representation of the problem, and no one drew the correct conclusion concerning the red glass marbles; by contrast, almost all participants in the control condition were able to draw that very same correct conclusion. Apparently, in this problem, the human inferential horizon is very short: people are happy with drawing a first conclusion, but do not pursue its further consequences.


# EXPERIMENT 1b

Despite the consistency of the results of Experiment 1a with our expectations, we were surprised by the strength of the effect. Cherubini and Johnson-Laird's (2004)

found effects fully consistent with the present ones, but weaker than those observed in Experiment 1a – unless participants were time-constrained or cognitive load was increased by the use of negative premises (respectively Experiment 1 and 2 of Cherubini and Johnson-Laird, 2004). Perhaps, the procedure of individual verbal presentation of the problems – even though it enhances a more proper understanding of the premises – somehow affected the findings. Unpredictable additional factors could have intervened: e.g., some participants could have been anxious at being tested by one of their professors (even though full guarantee was given that this was not an intelligence test, and was not to affect in any way their careers); or, since the experimenters were not blind to the hypothesis, they could have involuntarily encouraged some participants to give responses consistent with the hypothesis. In order to neutralize these sorts of unpredictable parasite factors, Experiment 1b used the same experimental and control task as Experiment 1a, but was administered as a questionnaire – with written instructions and closed-choice responses – to a large group of participants.

**METHODS**

*Participants*

One-hundred and eight (69 females) students taking a first-year course in general psychology at the University of Milan-Bicocca participated in the experiment as a course requirement. None of them had previously taken courses in logic or the psychology of thinking and reasoning.

*Task and procedure*

The experimental task was as follows:

Imagine that in front of you there is a closed box. In the box there are 5 marbles, which you cannot see. The marbles can be either red or blue, and they can be either plastic marbles, or glass marbles. In other words, there can be 4 types of marbles: red glass ones, red plastic ones, blue glass ones, and blue plastic ones. Furthermore:
1. more than half of the marbles in the box are red marbles;
2. more than half of the marbles in the box are glass marbles;
3. If in the box there is at least a red glass marble, then there also is at least one blue plastic marble.

On the ground of the above information, what can you establish with certainty concerning the number of blue plastic marbles in the box?
❑   I cannot establish anything for certain.

❑    I am certain that in the box there is at least one blue plastic marble
❑    I am certain that in the box there are at least two blue plastic marbles
On the ground of the above information, what can you establish with certainty concerning the number of red glass marbles in the box?
❑    I cannot establish anything with certainty
❑    I am certain that in the box there is at least one red glass marble
❑    I am certain that in the box there are at least two red glass marbles

*[for half of the participants, the two final questions were inverted]*

The control task was as follows:

Imagine that in front of you there is a closed box. In the box there are 5 marbles, which you cannot see. The marbles can be either green or yellow, and they can be either large or small. In other words, there can be 4 types of marbles: large green marbles, small green marbles, large yellow marbles, small yellow marbles. Furthermore:

    1.   more than half of the marbles in the box are green;
    2.   more than half of the marbles in the box are large;
    3.   in the box there is at least one small yellow marble.

On the ground of the above information, what can you establish with certainty concerning the number of large green marbles in the box?
❑    I cannot establish anything for certain.
❑    I am certain that in the box there is at least one large green marble
❑    I am certain that in the box there are at least two large green marbles

Each participant received both problems, in random orders. The problems were inserted in a booklet comprising four other reasoning and decision problems, unrelated to the present study. Even though the order of the problems was randomized, care was paid that the two problems were interspaced by at least one irrelevant problem. The booklets were anonymous, and participants did not have time limits for filling them.

**Results and analyses**

In the experimental task, 93 participants (86%; significantly more than chance, $p < .0001$) correctly responded that there was at least one blue plastic marble in the box. Four (4%) said that there were at least two blue plastic marbles in the box, and 11 (10%) could not conclude anything for certain. Responses to the critical questions about the red glass (experimental task), or the large green marbles (control task) are reported in Table 1.

--- Insert Table 1 about here ---

The "cannot conclude" responses are very few and their frequency does not differ between the two tasks. The distributions of the critical "at least one" and "at least two" responses were reliably different in the two tasks (exact $\chi^2 = 39.6$, 1-tailed $p$

=1.97*e-10*), showing that correct responses were more frequent in the control task than in the experimental task. Correct responses in the Experimental task were note significantly different from chance (chance level set at .33), and wrong "at least one" responses in the control task were not significantly different from chance. Limiting analyses to the 93 participants that correctly responded to the question concerning the blue plastic marble in the experimental task – that is, those participants for whom we are definitely sure that the available information in the experimental and control task was exactly the same – 61 (66%) of them responded correctly in the control task, vs. 22 (24%) responding correctly in the experimental task (exact $\chi^2 = 33.1$; 1-tailed *p* =6.32*e-9*).

**Discussion**

The results confirmed the findings of Experiment 1a. Participants responding correctly to the question concerning the blue plastic marble in the experimental task had *exactly* the same information available as in the control task. This notwithstanding, a vast majority of participants did not pursue the consequences of the initial conclusions of the experimental problem, and hence did not draw the correct conclusion that was easy to draw in the control problem. These findings corroborate the idea that the difficulty in iteratively pursuing the consequences of a conclusion depends on a basic difficulty in spontaneously integrating one's own provisional conclusions – derived from an initial mental model of a problem – into that very same model.

The effects observed in Experiment 1b are weaker than those observed in Experiment 1a: here, 25% of the participants correctly answered the experimental problem, vs. 0% in Experiment 1a. At first glance, this finding might suggest that we were right in suspecting that the individual verbal presentation of the tasks in Experiment 1a somehow inflated the results. However, one must consider that in Experiment 2, while correct responses to the experimental problem increased, correct responses to the control problem decreased. Actually, this trend suggests that participants were not reasoning *more* accurately, but *less* accurately, and that some of them relied on guessing. In conclusion, the tendency not to pursue the further consequences of one's own initial conclusion in this sort of problems ranges from strong (43% difference in correct responses between experimental and control task in Experiment 1b) to very strong (88% difference in Experiment 1a).

# EXPERIMENT 2

In Experiment 2 we seek further support for the difficulty of integrating one's own initial conclusion into the initial representation of a problem, by using a different task, inspired by small-N BCG, but set in a non-interactive setting.

**The task**

A computer randomly generates two [in alternative versions, three, or four] integer numbers between 0 and 100. All numbers in the span are equiprobable. You shall select a third [fourth or fifth, in the versions where the computer generated three or four numbers] number. Then, the mean of all the numbers is computed, including the two [three, four] generated by the computer and the one that you chose. The overall mean is halved, in order to obtain a target number T. If the number that you chose is T, or is no more than one unit away from T, you win € 5. Otherwise, you win nothing. (please limit your choice to numbers comprised between 18 and 32, extremes included).

**Predictions**

Given $k$ the estimated mean of the numbers generated by the computer, $N$ the total number of numbers involved, and $p$ the target proportion of the overall mean, the correct choice is $C = p(N-1)k/(N-p)$, or an integer next to C if C is not integer. In different versions of this task, $N$ varied from 3 to 5. The proportion $p$ was set at .5, and the most rational expectation – even though weak – concerning $k$ was 50, because the numbers generated by the computer were equiprobable and random. It is not likely that participants – which were not mathematicians – could work out the above algebraic equation for calculating C in one shot. Instead, we hypothesize that they should search for it iteratively, as follows:

Step 0) participants know that the expected k is 50, but they do not know the expected overall mean, because they have not chosen a number; henceforth, for lack of a better estimate, they should anchor to k for computing an approximation to the target number, that is pk = 25; they should provisionally set 25 as their choice;

Step 1) participants who realize that their estimation of the target number at step 0 was provisional, because it did not include their own number, should recompute the overall mean assuming that their choice is 25, obtaining  that the overall expected mean is 41.66 [in the N = 3 trials], and the expected target number is 20.8. They should set their chosen number to 20 or 21.

Step 2) careful participants would then check the effects of their new choice on the target number. If their choice was 20, they can realize that the target number becomes exactly 20, and stop. If their choice was 21, they can realize that the target number becomes 20.16, and either set their choice at 20 and retry, or stop, given that 21 is within one unit form 20.16.[2]

The convergence values for the three games can be easily computed if people does not stop at step 0, and engage in at least one step of iteration; they are:

- $N$=3: convergence: 20; optimal choices: 19, 20 or 21;
- $N$=4: convergence: 21.42; optimal choices: 21 or 22;
- $N$=5: convergence: 22.22; optimal choices: 22 or 23

If people integrate their initial choice into the initial expected mean, and compute the new expected mean associated to their choice, they should be able to pursue iteratively the convergence values; even if they do not compute the *exact* convergence value, if they perform at least one step of iteration they could manage to grasp the optimal numbers reported above, all of them < 25 and increasing with increasing *N*. Otherwise, people that stops at step 0, that is people that initially represent the target number as half of 50 (the expected mean of the computer generated numbers), and do not realize that by choosing 25 they modify the overall mean and the target number, they should stick to 25.

## METHODS

*Participants*

The rational expectation is that the mean of two, three or four randomly chosen numbers between 0 and 100 is 50, but the likelihood that it will be exactly 50 is

---

[2] This is the most likely psychological algorithm for arriving to the convergence value, because of the availability of k = 50 as an anchor for initially representing the problem. However, people realizing that their number affects the overall mean can proceed in a different way: they can select from the start a random number among the available ones, and then compute the resulting target number. If they were not lucky at guessing, and did not choose 19, 20, or 21 [in the N = 3 game], they should then readjust their initial choice in further rounds of iteration, not differently from the algorithm illustrated above. People stopping at step 0 when following this alternative procedure do not necessarily choose 25; they can choose any available number.

very low. Probably for this reason, in pilot testing we realized that some of our students did not grasp that the rational expected mean of a set of two, three, or four randomly generated integers between 0 and 100 is 50; those students had a tendency to pick numbers at random in the game. Therefore, in the Experiment that we report here we screened candidates: those answering "50" to the question "If I draw at random two numbers between 0 and 100, extremes included, which is their most likely mean?" were admitted to participate in the Experiment; the others were not admitted. Out of 28 screened students, 20 students of psychology (mean age: 21.9; 11 females) from the university of Milan-Bicocca eventually took part in the Experiment, in exchange for € 5 plus the amount of money that they managed to win. Some of them had taken a course in psychology of reasoning, none of them had taken any specialized course in logic.

*Procedure*

Participants were tested individually in a quite lab. Before the task, instructions were given verbally, and repeated as required in order to assure that they were fully understood. Each participant received all three versions of the problem, with $N$=3, 4 and 5, in random order. In each problem, either two, three, or four boxes with a "?" were displayed on a computer screen. Half a second later, the sentence "the computer has now randomly generated 2 [3, 4] numbers. Now you pick a number comprised between 18 and 32. If that number is within one unit from half of the mean of all the numbers (those generated by the computer, plus the one you picked), you win € 5". Participants wrote their choice using the numeric keypad. Responses were not time constrained. Participants had to respond to all three problems before receiving feedback, when – for each problem – the randomly generated numbers were disclosed, the target number was computed in two steps (first by showing the mean of all numbers, then by halving it), and the participant was acknowledged whether or not she had won. At the end of the Experiment, each participant was asked to explain her responses, and recorded.

**Results and analyses**

The frequency of chosen numbers is reported in Table 2.


--- Insert Table 2 about here ---

The distribution of optimal choices and choice of 25 is not reliably different in the three versions of the problem. By collapsing the three versions, 25 was chosen 43 times (72%). This is significantly more than chance, even if chance level is set, instead of .067 (15 numbers were available for choice; exact binomial test, $p$=4.4$e$-$8$), at a far more conservative value of .5 ($p$<.001). The chance level for the selection of optimal numbers was .2 for $N$=3, and .13 for $N$=4 or 5, so we set it at .16. The 13 optimal number selected were not reliably different from chance. However, $\chi^2$ analysis of the choices of optimal numbers – once excluded all other numbers – confirmed that the numbers reliably varied in the three conditions, suggesting that these were not random choices, but choices from a few participants that actually engaged in iterative reasoning (exact $\chi^2 = 20.2$; $p <$ .0005). This is further confirmed by the fact that three participants consistently selected an optimal number in all versions of the problems, and by their explanations in the debriefing session, depicting a process of iterative mental computations matching the algorithm outlined in the prediction paragraph. All three participants stated that they initially considered choosing 25, because it is the half of the expected mean of the computer generated numbers. Then, they included 25 in the computation of the mean, and they realized that the expected mean changed, and so the target number, and thus they lowered their choice, and checked again. The remaining 17 participants offered various self- reports of how they chose their numbers, including three persons who stated that they choose at random a number lower than 25 and then stopped (apparently applying the alternative algorithm described in footnote 2, but stopping at step 0), but none described that they recalculated the expected mean after they had selected an initial number. Out of the 13 participants that selected 25 in all problems, none spontaneously mentioned, in the debriefing session, that they recalculated the expected target number after having chosen 25. These participants did not engage in an iterative search for the optimal choice.

**Discussion**

The task in Experiment 2 was very similar to a small-$N$ BCG, but was set in a non-interactive setting, where representations of choices made by other human players cannot be a factor in determining the participants' responses. Of course, since all numbers but the one chosen by the participant were random, the rational

choice here was not 0, as in the actual BCG, but a definite set of values that could be easily computed in a few step of iterations, if only participants integrated their own number in the computation of the target number. Most of them did not do so. Both their choices and their explanations clarified that participants that consistently chose "25" did so because they built a raw initial representation of the target number (i.e., half of the expected mean of the computer generated numbers), and did not realize that after considering choosing 25 they should have recomputed the mean: that is, they did not integrate their choice into the initial representation that originally suggested that choice. These findings corroborate and generalize those in Experiments 1a and 1b: people have difficulties at integrating one's own initial conclusion into the initial representation of the problem that suggested that conclusion. Thus, they do not easily realize that their own initial decision modifies the representation of the problem, suggesting that a different decision might be more appropriate. Because of these cognitive constraints, people have difficulties at pursuing iteratively the further consequences of their own initial conclusions or decisions, even in non-interactive contexts.

## EXPERIMENT 3

The findings from the previous Experiments, together with those by Cherubini and Johnson-Laird (2004), suggest that the difficulty of spontaneous iterative reasoning in non-interactive contexts – that is, contexts that do not tax people cognitive abilities by requiring them to sort out what other people will likely do – can be generalized to a large class of problems. Yet, in all these problems people were not *explicitly* alerted that iterative reasoning was required or useful in order to correctly solve the task. What happens when participants are explicitly instructed to reason iteratively? That is: were the problems tested so far difficult because people do not spontaneously realize that they should pursue a chain of iterative representations and conclusions (but – if they realized it – they *could* pursue it), or is it the case that their difficulty would persist even if people were explicitly shown how to reason iteratively? Johnson, Camerer, Sen and Rymon (2002) found that people do not spontaneously engage in iterative reasoning in an interactive bargaining game, and do not spontaneously learn from experience how to establish equilibrium points by thinking iteratively. They then explicitly trained

their participants to apply a simple backward iterative strategy. Participants learned the iterative strategy easily, and transferred it to similar problems with different parameters. Accordingly, even in non-interactive problems, people might not spontaneously engage in iterative reasoning (as shown in the previous experiments), but they might do so if alerted that iterative reasoning is necessary, and shown examples of how to do that. In this Experiment we used a different problem, transparent in its iterative structure even though more complex than the previous ones in the calculations it required. The problem was loosely inspired by Berry and Broadbent's (1984) "sugar factory" task, inasmuch it required to "stabilize" the output of a factory. But – contrary to their task – no random factors contributed to the output, and participants could give only one answer (instead of multiple answers, one for each cycle of production of the factory, as occurred in Berry and Broadbent's study). As a consequence, our problem involves explicit forecasting, whereas Berry and Broadbent's problem relied on implicit learning.

## METHODS

### Participants

Eighty-four students of psychology (61 females) took part in the experiment as a requirement for a first-year course in general psychology at the university of Milan-Bicocca. None of the participants had participated in the previous experiments, and none had taken courses in logic or in the psychology of reasoning.

### The task

In the iterative problem that we used participants had to set a constant rate of production $k$ for a given item, such that after iterating through production cycles $t_1...t_n$ the number of those items in stock ($X_n$) converged to a given target value $s$. $X_0$, the number of items available in stock at $t_0$, is known. The content of the stock at each following cycle is $X_i = p(X_{i-1}) + k$, with $p \in (0, 1)$. The series of $X$, in non-recursive form, is explicated as follows:

$$X_n = [pX_0 + k(1/p^n - 1)(1/p - 1)]p^{n-1}$$

For a given $k$, the series converge to $s = k/(1-p)$; that is, in order to converge to a given $s$ - as required – $k$ must be set at $s(1-p)$. Hence, formally, $X_0$ does not affect $k$; it only affects the minimal number of cycles $n$ required to attain convergence to $s$.

The problem was embedded in different scenarios. An example follows:

You are the newly appointed production manager of a car factory. The monthly stock of cars amounts to ¼ of the stock of the previous month (the other cars are sold), plus the newly produced cars. For example, if you have 1200 cars in stock in May, and you produce 500 cars in June, you will end up with 800 cars in stock (one-fourth of the stock in May, plus 500). Out of these, 200 will remain in stock for July; if you again produce 500 cars in July, your stock will be 700 cars (and 175 will remain in stock for August), and so on.
Now it is February. The cars in stock in January were 40. Set a *constant* rate of monthly production, so as to stabilize, in a few months (it is not important how many), the number of cars that are in stock each month at exactly 1000 units. In doing calculations, if necessary round decimal numbers to the nearest unit.
Which rate of production do you set? _____

There were two different scenarios: a car factory (example above), and a workshop of hand-made luxury wristwatches. The magnitude of involved numbers changed in the two scenarios ($s=1000$ for the car factory; $s=100$ for the wristwatches). The value of $p$ was ¼ *in* both scenarios. The correct responses were $k=750$ for the car factory scenario and $k=75$ for the wristwatches scenario. Each scenario came in two versions: low-anchor version, where $X_0<s$ (40 cars; 4 wristwatches); high-anchor version, where $X_0>s$ (2800 cars; 280 wristwatches) (the example above is the low-anchor version of the car factory scenario). As a result, we had four versions of the problem.

*Predictions*

Contrary to Experiment 1a, 1b, and 2, here people are explicitly told to reason iteratively, and an example of how this sort of reasoning works is reported at the beginning of each problem, where it is shown how to compute the content of the stock month by month. The algebraic solution for $k$ is opaque, and non-expert participants should not be able to work out it. Alternatively, participants might proceed by trials and errors, by setting a value for $k$, then estimating or calculating the resulting stock iteratively in order to see whether it converges to $s$; if not, try a different $k$. However, this sort of iterative computation, if carried out in full depth, is mentally distressing, and we expect that most people fall short of the steps of

iterations needed in order to give appropriate responses. Accordingly, their estimates of $k$ should depend on the formally irrelevant initial stock: lower estimates in the high anchor-problems (where participants are seeking a "decreasing" pattern), and higher estimates in the low-anchor problems (where participants are seeking an "increasing" pattern). As a last resort, participants could guess by applying the anchoring heuristic (Tversky and Kahneman, 1974), thus estimating lower $k$ values in the high-anchor problems (where the initial stock by far exceeds the target stock) than in the low-anchor problems (where the initial stock is depleted with respect to the target). The predicted result of both strategies is the same: estimates of $k$ should depend on the formally irrelevant initial stock, being higher in the low-anchor than in the low-anchor problems. By contrast, participants that engaged in an in-depth iterative search of $k$ should be able to work out the correct value, that is, the same for low-anchor and high-anchor problems.

*Procedure*

Participants were tested in a large group in a classroom. Each participant received both scenarios, one in the low-anchor version and the other in the high-anchor version. The order of presentation of the scenarios was balanced across participants. Each scenario was reported on a page of a booklet, leaving sufficient space for writing notes, calculations, and explanations of answers. Participants were told that they had no time limits, and that they could write notes and perform written calculations, but could not use electronic calculators.

## Results and analyses

Preliminary analyses showed that there were not reliable differences between the car factory scenario and the wristwatches scenario. Accordingly we collapsed together data from the two scenarios. In order to uniform responses, instead of the raw $k$ indicated by participants we used as dependent variable the ratio of estimated $s$ – obtained from the $k$ indicated by the participant – to the target $s$. The ratio is 1 for correct responses, >1 for overestimated $k$ values, and <1 for underestimated $k$ values. Means for the low-anchor scenarios and the high-anchor scenarios were 1.05 and .92, respectively. The difference is reliable, $t(83)=4.7$, $p<.0001$, showing that – as predicted – people were affected by the initial stock:

production rates were markedly underestimated for the high-anchor scenarios, where $X_0$ far exceeded $s$, with respect to the low-anchor scenarios.

The number of exact responses ($k$=750 in the car factory scenarios; $k$=75 in the wristwatches scenario) were 18 (21%) in the low-anchor scenarios, vs. 10 (12%) in the high-anchor scenarios. This occurrence, together with the smaller discrepancy of the mean from the correct value in the low-anchor scenarios, suggests that people, unexpectedly, found the low-anchor scenarios somewhat easier than the high-anchor scenarios.

## Discussion

Even though participants were explicitly alerted of the iterative nature of the problems, were shown an explicit example, and were allowed to do written calculations, their performance was not very good. A minority of participants worked out the correct answers. Most of them produced estimations. More importantly, those estimations were affected by a formally irrelevant parameter: the initial amount of items in the stock. People anchored to that amount and tried some initial guesses at the production rate, possibly exploring the resulting production patterns iteratively for a few cycles, but not as far as required for a correct evaluation. As a result, where the initial stock far exceeded the target value people underestimated production rates; by contrast, where the initial stock was low, people overestimated production rates. These findings confirm that even people that are explicitly told to pursue an iterative chain of reasoning, and are shown how to do it, are often unwilling to engage in this type of reasoning.

The better performance in the low-anchor problems with respect to the high-anchor ones was unpredicted, but it is easily interpreted *post-hoc*. Almost all participants promptly realized that the production rate could not be greater than the target. As a consequence, people given the low anchor problems looked for large production rates (as predicted), but not as large as $s$ (1000 in the car factory problem, and 100 in the watch factory problem). Because, with $p$ set at 1/4, correct rates were 750 (car factory) and 75 (watches), guessing a large number lower than $s$ had good chances of resulting in a correct guess from the very first attempts. By contrast, in the high-anchor conditions participants had an initial tendency to try very small $k$. They later revised those initial estimations, but – on average – they failed to increase them as much as necessary.

# GENERAL DISCUSSION

Not realizing the further consequences of an initial conclusion can have serious effects on planning abilities and rational behavior in general. For example, not basing our plans on anticipations of the most convenient counter-moves by an adversary can result in poor performance not only in board games, but also in auctions, economic planning, military planning, and possibly many other types of interactive contexts. This lack of depth in reasoning has been demonstrated and measured in many previous works mainly concerning interactive contexts (e.g., Nagel, 1995, 1999a, 1999b; Camerer, 2003; Camerer et al., 2003, 2004; Duffy and Nagel, 1997; Johnson et al., 2002; Güth et al., 2002; Ho et al., 1998; Kocher and Sutter, 2005). It can have different causes, including the inability to adequately represent other people as rational agents (Bosh et al., 2002; Nagel, 1995), the inclusion of social utilities in the computation of an otherwise rational agent (e.g. Fehr and Schmidt, 1999; Fehr and Gachter, 2000; Berg, Dickhaut, and McCabe, 1995; Johnson et al., 2002), the lack of an appropriate understanding of the problem form (Chou et al, 2008), or be the result of cognitive constraints bounding the rationality of each individual (Grosskopf and Nagel, 2007, 2008). The present study supports the latter view, and further specifies it at the psychological level, in two ways. First, all previously mentioned works used interactive contexts, where inability to attribute rationality to the other agents or the consideration of social utilities can affect behavior (but see Johnson et al, 2002, Exp. 2, where social utilities were "switched off" and nonetheless participants were not able to perform iterative bargaining). By contrast, the present study – together with the one by Cherubini and Johnson-Laird (2004) – shows that difficulties in iterative reasoning are present in a wide range of non-interactive situations, where performance cannot be explained by an inability to represent other agents' behaviors, or by social utilities. Of course, difficulties in representing other people beliefs and goals, or social values such as fairness and reciprocity, can be an important factor that impairs performance in interactive contexts, but we show that they are not the necessary cause of the psychological limits of iterative reasoning.

Second, our experiments suggest that inefficiency in iterative reasoning is itself a byproduct of a very basic feature of human reasoning. According to MMT (Johnson-Laird and Byrne, 1991; Johnson-Laird, 2004), when they reason

spontaneously people build one initial representation of the problem, and draw one or some initial conclusions from it. The initial representation might be a provisional one, but people easily forget that alternative representations are possible, and thus they commonly stick to their initial conclusions. Systematic search of alternative representations is not commonly undertaken. Cherubini and Johnson-Laird (2004) showed that focusing on a given model of a problem often hinders the revision of that model in the light of the very same conclusions that were drawn from it. In other words, people can spontaneously realize that a consequence C follows from a situation S, and yet they can miss that C modifies S, so that a new consequence C' follows. This result is now further supported by the present experiments. The most frequent responses in Experiments 1a, 1b, and 2 resulted from building an appropriate initial representation, drawing an initial conclusion, but then not performing the required integration of that conclusions with the initial representation. Similarly, Experiment 3 showed that a parameter that was prominent in the *initial* representation of the problem, but was utterly irrelevant if one fully unraveled the iterative nature of the task (that, in that experiment, was made explicit to the participants), remarkably affected responses. In our view, these findings suggest that the basic cognitive constraint that bounds human rationality in the pursuit of iterative chains of conclusions is the difficulty of integrating one's own initial conclusion with one's own initial representation of a problem, and modify the latter accordingly. This limitation might impair people's ability to accurately forecast non-immediate consequences of their own decisions and actions in complex settings, such as those of interest for economic and financial sciences.

Acknowledgments

# REFERENCES

Berg, J., Dickhaut, J., and McCabe, K. (1995). Trust, reciprocity and social history. *Games and Economic Behavior*, 10, 122–142.

Berry, D.C., and Broadbent, D.E. (1984). On the relationship between task performance and associated verbalized knowledge. *Quarterly Journal of Experimental Psychology*, 36A, 209-231.

Bosch-Domènech, A., García-Montalvo, J., Nagel, R. and Satorra, A. (2002), One, two, (three), infinity, ..: Newspaper and lab beauty-contest experiments, *American Economic Review,* 91, 1687-701.

Camerer, C. F. (2003), *Behavioural Game Theory. Experiments in Strategic Interaction,* Princeton: Princeton University Press.

Camerer, C. F., Ho, T.-H. and Chong, J.-K. (2003), Models of thinking, learning, and teaching in games, *American Economic Review, Papers and Proceedings,* 93, 192-5.

Camerer, C., Ho, T., and Chong, J. (2004). A cognitive hierarchy model of behavior in games. *Quarterly Journal of Economics, 119*, 861–898.

Cherubini, P, and Johnson-Laird, P.N. (2004). Does everyone love everyone? The psychology of iterative reasoning. *Thinking & Reasoning, 10*, 31-53.

Duffy, J. and Nagel, R. (1997), On the robustness of behaviour in experimental "beauty contest" games. *Economic Journal,* 107, 1684-700.

Fehr, E., and Gachter, S. (2000). Fairness and retaliation: The economics of reciprocity. *Journal of Economic Perspectives,* 14, 159–181.

Fehr, E., and Schmidt, K. (1999). A theory of fairness, competition, and cooperation, *Quarterly Journal of Economics*, 114, 817–868.

Grosskopf, B., and Nagel, R. (2008). Rational Reasoning or Adaptive Behavior? Evidence from Two Person Beauty Contest Games. Harvard NOM Research Paper No. 01-09. DOI: 10.2139/ssrn.10.2139/ssrn.286573

Grosskopf, B., and Nagel, R. (2008). The two-person beauty contest. *Games and Economic Behavior, 62,* 93–99.

Güth, W., Kocher, M. G. and Sutter, M. (2002), Experimental "beauty-contests" with homogeneous and heterogeneous players and with interior and boundary equilibria, *Economics Letters,* 74, 219-28.

Ho, T.-H., Camerer, C. F. and Weigelt, K. (1998), Iterated dominance and iterated best response in experimental "p-beauty-contests", *American Economic Review,* 88, 947-69.

Johnson-Laird, P. N. (2001). Mental models and deduction. *Trends in Cognitive Science, 5*, 434–442.

Johnson-Laird, P.N., Byrne, R.M. (1991). *Deduction.* Erlbaum, Hillsdale, NJ.

Kocher, M. G. and Sutter, M. (2005), The decision maker matters: Individual versus group behaviour in experimental beauty-contest games, *The Economic Journal*, 115, 220-223.

Nagel, R. (1995), Unravelling in guessing games: An experimental study. *American Economic Review,* 85, 1313-26.

Nagel, R. (1999a), A survey of experimental beauty contest games: Bounded rationality and learning, in: Budescu, E., Erev, I. and Zwick, R. (eds.), *Games and Human Behavior: Essays in Honour of Amnon Rapoport.* New Jersey: Lawrence Erlbaum Assoc., Inc., 105-42.

Nagel, R. (1999b), A kaynesian beauty contest in the classroom. *Expernomics*, 8, 7-12.

Roth, A.E., and Ockenfels, A. (2002). Last-minute bidding and the rules for ending second-price auctions: Evidence from eBay and Amazon auctions on the Internet. *American Economic Review*, *92*, 1093–1103.

Schelling, T.C. (1978). Micromotives and Macrobehavior. London/New York: Norton

Tversky, A., and Kahneman, D. (1974). Judgment under uncertainty: Heuristics and biases. *Science*, 185, 1124-1130.

Weber, R. A. (2003). "Learning" with no feedback in a competitive guessing game. *Games and Economic Behavior*, 44, 134-44.

**Figure legends**

Figure 1.   Diagram of the initial representation of the problem in Experiments 1a and 1b

Figure 2.   Diagram of the revised representation of the problem in Experiments 1a and 1b. The revised representation allows drawing the correct conclusion.

**Tables**

Table 1. Responses in Experiment 1b.

|  | Experimental task | Control task |
|---|---|---|
| *Cannot conclude* | 9 (8%) | 6 (5%) |
| *At least one red glass [large green] marbles* | 74 (69%) | 31 (29%) |
| *At least two red glass [large green] marbles* | 25 (23%) | 71 (66%) |

Table 2. Frequencies of chosen numbers in Experiment 2. Numbers in bold are optimal choices.

|  | *N=3* | *N=4* | *N=5* |
|---|---|---|---|
| *19 or 20* | **4** | 0 | 1 |
| *21* | **0** | **2** | 0 |
| *22* | 0 | **3** | **0** |
| *23* | 0 | 1 | **4** |
| *25* | 15 | 14 | 14 |
| *Other numbers* | 1 | 0 | 1 |

**Figures**

Figure 1



Red marbles

Glass marbles

**At least one red glass marble**
Hence: At least one blue plastic marble

Figure 2

At least 2 red glass marbles          At least 1 blue plastic marble