

# TOWARDS RATIONAL AGENCY

Dániel László KOVÁCS  
Advisor: Tadeusz DOBROWIECKI

## I. Introduction

Until recently *perfect rationality* was the most desired property of intelligent systems, i.e. rational agents, stipulating that they always act so as to maximize their *expected utility* given their beliefs. Nowadays a shift can be seen both in economy (from perfect to bounded rationality), in game theory (from action to program selection) and in philosophy (from act to rule utilitarianism) because of the fundamentally unsatisfiable requirements imposed by the definition [1]. This led in artificial intelligence to the definition of *bounded-optimality* [2], which had an impact on the other fields. Hence rationality is rediscovered as a central property of intelligent systems, as foreseen in [1].

## II. What is rational?

In classical game theory the notion of *equilibria* is used to define a set of rational strategies to be played by rational players. The most popular concept is the *Nash-equilibrium* [3]. Roughly speaking a set of strategies and the respective payoffs constitute a Nash-equilibrium if no player can benefit by changing his strategy while the other players keep their strategies unchanged. Although, it is a troubling fact, that the Nash-equilibrium isn't necessarily optimal (e.g. Pareto-optimal), moreover there may be multiple Nash-equilibria, and even, considering only pure strategies, no equilibria at all. These drawbacks make it worth reconsidering the definition of rational strategies, i.e. rationality.

A common feature of the previous, and numerous similar concepts is that they define rationality as a property of strategies (e.g. actions) rather than the programs selecting them. Specifying strategies is an implicit assumption about the players. E.g. Nash-equilibrium is a non-cooperative concept, so it implicitly assumes non-cooperating players, thus being not suitable for players that cooperate (form coalitions, communicate, trust each other, etc). On the other hand, by specifying programs, assumptions about players, i.e. the inspection of their rationality, is explicit.

A **collective of players** is called **rational (RCP)** for a game, if the players' programs for selecting strategies produce an optimal (e.g. Pareto-optimal) set of strategies. But what is the program of a rational player, if the others' programs are such, that it is impossible to form an RCP with them? A **player** is called **rational (RP)** for a game and a collective of other players with given programs, if its program, together with the others' programs, produces an optimal (e.g. Pareto-optimal) set of strategies from the set of all possible sets of strategies producible with the others. Consequently in a RCP every player is a RP, and if every player is a RP, then the collective of these players is a RCP. It can also be seen, that the definitions of RCP and RP depend upon the definition of optimality, i.e. they are relative notions. The most important questions though: Given a definition of optimality, what is a RP's program, if the programs of the other players' are still not (fully) known?

## III. Effects of bounded-optimality

Before investigating the above questions, let's switch from players to *agents*. This will allow the introduction of the rich concepts of agent-theory. An agent "can be anything that can be viewed as perceiving its *environment* through *sensors* and acting upon that environment through *effectors*" [5]. Thus game theoretical players are agents facing a problem represented by a game, i.e. a *task environment*, which denotes the combination of an environment and a *utility function*. "An agent is

*bounded-optimal (BO)* if its program is a solution to the constrained optimization problem presented by its architecture and the task environment” [2]. Formally an agent is **BO** in an environment  $E$  with a utility function  $U$ , and with an architecture represented by a time- and/or space-bounded universal Turing-machine  $M$  with a finite language  $L_M$ , if it has a program  $l_{opt}$  such that

$$l_{opt} = \arg \max_{l \in L_M} (U(\text{effects}(\text{Agent}(l, M), E))). \quad (1)$$

An agent-program  $l$  running on a machine  $M$  implements an agent-function  $\text{Agent}(l, M)$ , while *effects* is a function that returns a sequence of states through which the agent-function drives the environment (cf. Eq. (1)). Now there is a point to assume that players, i.e. agents with bounded knowledge face a *multiple-state*, or *contingency*, or even *exploration problem* [4], where they can’t (exactly) determine either the state of the environment, or the utility function, or the effects of actions, or the programs of others, etc. In this case BO can be used to define optimality in the definitions of RP and RCP, calling forth **bounded-RP (BRP)** and **bounded-RCP (BRCP)**, where a player’s program for strategy selection is an agent-program, a game is a task environment.

#### IV. Proposed architecture

An agent architecture for realizing a BRP is proposed [5], where an agent’s agent-program is a complex inference-engine based on a *globally optimal* evolution of strategies (e.g. actions). An *evolutional algorithm* is proven to be globally optimal in general search spaces, if it implements *elitism* and every individual of the population is *reachable* from every other individual by means of *mutation* and *recombination* [6]. Still, such an “evolving agent” isn’t necessarily a BRP because of the incidentally imperfect knowledge it may have either about itself (e.g. its sensors, effectors), or the environment, or the utility function (which it uses as a *fitness function* to evolve its strategy). Nevertheless a globally optimal evolution of agent-programs should produce a BRP agent, i.e. an agent-program for the given agent-architecture, that maximizes the utility function of the task-environment (cf. Eq. (1)), where the maximum implies, that it produces an optimal strategy as described in definitions of RP and RCP. Unfortunately evolution is sub-optimal if there is a finite time-bound on runtime. In this case the concepts of strategy and agent-program evolution could be replaced and/or mixed with *learning* techniques, or other heuristics to establish optimal inference-engines, i.e. optimal agent-programs for BRP agents. First though, the proposed concepts should be studied in the simplest cases (e.g. games with perfect information, i.e. task-environments representing *one-state problems*). Then a more general concept can be developed covering more complex cases (e.g. exploration problems). The solution of these problem-classes could be connected with *AI planning*, i.e. *contingency* and *exploration planning* [4] respectively.

#### V. Conclusions

A realizable concept of rational agency is proposed by connecting the fields of game, evolution and agent-theory. An agent-concept for realizing bounded rational agents is proposed. The aim of the research is to propose a tractable design method for optimal complex systems for real-world tasks.

#### References

- [1] J. v. Neumann, and O. Morgenstern, *Theory of games and economic behavior*, Princeton University Press, 1947.
- [2] S. Russell, and D. Subramanian, “Provably bounded-optimal agents,” *Journal of AI Research*, 2:1–36, 1995.
- [3] J. F. Nash, “Non-cooperative games,” *Annals of Mathematics*, 54(2):286–295, 1951.
- [4] S. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach*, Prentice Hall, 1995.
- [5] D. L. Kovács, “Evolution of Intelligent Agents: A new approach to automatic plan design”, in *Proc. of IFAC Workshop on Control Applications of Optimization*, pp. 237–243, Visegrád, Hungary, June 30–July 2 2003.
- [6] G. Rudolph, “Convergence of Evolutionary Algorithms in General Search Spaces”, in *Proc. of IEEE International Conference on Evolutionary Computation*, pp. 50–54, Nagoya, Japan, May 20–22 1996.