

Thinking Fast and Slow in AI

G. Booch (IBM), F. Fabiano (Univ. Udine), L. Horesh (IBM), K. Kate (IBM), J. Lenchner (IBM), N. Linck (IBM), A. Loreggia (EUI), K. Murugesan (IBM), N. Mattei (Tulane Univ.), F. Rossi (IBM), B. Srivastava (USC)

Abstract

This paper proposes a research direction to advance AI which draws inspiration from cognitive theories of human decision making. The premise is that if we gain insights about the causes of some human capabilities that are still lacking in AI (for instance, adaptability, generalizability, common sense, and causal reasoning), we may obtain similar capabilities in an AI system by embedding these causal components. We hope that the high-level description of our vision included in this paper, as well as the several research questions that we propose to consider, can stimulate the AI research community to define, try and evaluate new methodologies, frameworks, and evaluation metrics, in the spirit of achieving a better understanding of both human and machine intelligence.

Motivation and Overall Vision

AI systems have seen dramatic advancement in recent years, bringing many successful applications that are pervading our everyday life. However, we are still mostly seeing instances of narrow AI: each of these developments are typically focused on a very limited set of competences and goals, e.g., image interpretation, natural language processing, label classification, prediction, and many others. Moreover, while these successes can be accredited to improved algorithms and techniques, they are also tightly linked to the availability of huge datasets and computational power (Marcus 2020). State-of-the-art AI still lacks many capabilities that would naturally be included in a notion of intelligence, for example, if we compare these AI technologies to what human beings are able to do. Examples of these capabilities are generalizability, robustness, explainability, causal analysis, abstraction, common sense reasoning, ethics reasoning, as well as a complex and seamless integration of learning and reasoning supported by both implicit and explicit knowledge.

Majority of the AI community make sturdy attempts to address the current limitations of AI and create systems that display the ability for more human-like qualities, using a variety of approaches. One of the main debates is whether end-to-end neural network approaches can achieve this goal? or whether we need to integrate machine learning with symbolic and logic-based AI techniques? We believe that the integration route is the most promising, and this is supported

by several results that have been obtained along this line of work. For example, Bengio conjectures that it is necessary to generalize from raw data to a "consciousness stream" of few a concepts sparsely related to each other (Bengio 2017). Also, Marcus argues that explicit knowledge, symbols, and reasoning should be used to improve the robustness of current AI systems (Marcus 2020). Other researchers are building hybrid systems that use both machine learning and symbolic reasoning techniques, employing a so-called neuro-symbolic AI approach (d'Avila Garcez and Besold 2019).

We argue that a better comprehension regarding of how humans have, and have evolved to obtain, these advanced capabilities can inspire innovative ways to imbue AI systems with these competencies. To this end, we propose to study and exploit cognitive theories of human reasoning and decision making as a source of inspiration for the causal source of these capabilities, that help us raise the fundamental research questions to be considered when trying to provide AI with desired dimensions of human intelligence that are currently lacking.

More precisely, we analyze some of these theories, with special focus on Kahneman's theory of thinking fast and slow, and attempt to connect them into a unified theory with the aim to identify (some of) the roots of the desired human capabilities. We then propose to translate these theories of the human mind into the AI environment, and we conjecture that this will possibly lead to some of the same kind of capabilities in machines. This paper gives a brief and high-level overview of this general approach, tries to motivate it, lists several research questions along the way, and ends with a call for action to AI researchers to convene an AI research community that explores this space, assesses the validity of this approach, and possibly uses it to make significant advances in AI as well as in understanding human intelligence.

Thinking Fast and Slow, and Other Theories of Human Decision Making

According to Daniel Kahneman's theory, described in his book "Thinking, Fast and Slow" (Kahneman 2011), humans' decisions are supported and guided by the cooperation of two main kinds of capabilities, that, for sake of simplicity are called "systems": system 1 provides tools for intuitive, imprecise, fast, and often unconscious decisions

(“thinking fast”), while system 2 handles more complex situations where logical and rational thinking is needed to reach a complex decision (“thinking slow”).

System 1 is guided mainly by intuition rather than deliberation. It gives fast answers to very simple questions. Such answers are sometimes wrong, mainly because of unconscious bias or because they rely on heuristics or other short cuts, and usually do not have an explanation. However, system 1 is able to build models of the world that, although inaccurate and imprecise, can fill the knowledge gaps through causal inference and allow us to respond reasonably well to the many stimuli of our everyday life. A typical example of a task handled by system 1 is finding the answer to a very simple arithmetic calculation, or reaching out to grab something that is going to fall. We use our system 1 about 95% of the time when we need to make a decision.

When the problem is too complex for system 1, system 2 kicks in and solves it with access to additional computational resources, full attention, and sophisticated logical reasoning. A typical example of a problem handled by system 2 is solving a complex arithmetic calculation, or a multi-criteria optimization problem. To do this, humans need to be able to recognize that a problem goes beyond a threshold of cognitive ease and therefore the need to activate a more global and accurate reasoning machinery. Hence, introspection is essential in this process.

When a problem is new and difficult to solve, it is handled by system 2. However, certain problems over time pass on to system 1. The reason is that the procedures used by system 2 to find solutions to such problems are used to accumulate examples that system 1 can later use readily with little effort. Thus, after a while, some problems, initially solvable only by resorting to the system 2 reasoning tools, can become manageable by system 1. A typical example is reading text in our own native language. However, this does not happen with all tasks. An example of a problem that never passes to system 1 is finding the correct solution to complex arithmetic questions (unless one is a mathematical genius).

While system 2 seems capable of solving harder problems than system 1, system 2 does not usually by itself, but it is rather supported by system 1 in its elaborate calculations. When searching for a solution in a very large solution space, system 2 does not usually explore the whole search space but may employ heuristics that are provided by system 1 and that help in focusing the attention only on the most promising parts of the space. This allows system 2 to work with manageable time and space.

System 1 is also able to perform some basic forms of causal reasoning. This allows it to build a (possibly imprecise and biased) model of the world even if it has incomplete knowledge, and to use this model to tackle simple tasks. The system 1 data-level causality skills also support the more complex and accurate reasoning of system 2 on problem which are cognitively more complex.

With this complex internal machinery, humans can reason at various levels of abstraction, adapt to new environments, generalize from specific experiences to reuse their skills in other problems, and can easily multi-task when using their system 1. On the other hand, system 2 is sequential, since it

requires full attention to devise and execute the appropriate solving procedure for a complex or new problem. Thus only one complex task at a time can be handled by a human being.

Another way to look at this division is that system 1 works at a local level, activating only a specific part of the brain, e.g., when we recognize a familiar face, or when we speak. On the other hand, system 2 involves more regions of the brain and combines their contributions, so it works at a more global level. It is therefore the presence of multiple specialized agents/components/skills that supports the functioning of both systems. We recall however that system 1 and system 2 are not systems in the usual multi-agent architecture terminology, but rather they are metaphors and short-hands for two broad classes of information processing.

While Kahneman’s theory gives a detailed account on how we make decisions, the theory of Y. N. Harari (Harari 2014) conjectures what led us to have such reasoning capabilities over generations, and why our decision making engine came to have this organization. According to Harari’s theory, Homo-Sapiens had greater success than animals or other primates and human species because over (evolutionary) time it acquired the ability to conceive and communicate high-level stories. These skills allowed Sapiens to collaborate flexibly with others at scale, and thus to have a broader impact and a larger influence upon the world. This capability was supported by the ability to abstract from raw data to high level concepts, that could then be communicated more easily and compactly. As noted above, this is primarily a system 2 capability, since it requires a voluntary and conscious decision that goes beyond an involuntary reactive behavior (and its communications to others) to be able to abstract, generalize, and reason about scenarios and concepts that do not physically exist in real life (the so-called “stories” in Harari’s terminology). The ability to reason at different levels of abstraction requires also another mechanistic capability: to focus attention to certain, limited set of features and process them in depth, while disregarding others, since deferring (temporarily) some of the features of a concept is the essence of an abstraction process (Zucker 2003).

While these theories are just drawing hypotheses regarding the causes of human behavior and decisions, due to the current incomplete knowledge and understanding of our brain, they are supported by ample evidence and experimental results.

The notion of consciousness is important as well, when trying to identify the traits of human (or animal) intelligence. According to M. Graziano (Graziano 2013; Graziano et al. 2020), there are two forms of consciousness in human beings: the I-consciousness (I for Information) and the M-consciousness (M for mysterious). The first one refers to the ability to solve (possibly complex) problems, by recognizing necessary processing steps in specific (even new) context, to tackle a desired problem. This conduct seems to be related to system 2’s high-fidelity processing mode, since it has to do with considering a problem and harnessing the relevant faculties of our cognition to devise a plan to solve it. The latter refers to our ability to build a simplified, approximate, and subjective model of peoples’, both ourselves and others, mind, beliefs and intentions. Such low-fidelity

model building can be linked to system 1, as system 1 is able to form a rapid but usually inexact model of the world. While many animals and primates possess various levels of M-consciousness, and also limited forms of I-consciousness, humans excel at I-consciousness. M-consciousness can even sustain in presence of limited I-consciousness, yet, a sophisticated I-consciousness needs to rely on M-consciousness: without a model of the mind (of self and other agents), it is difficult to devise how to adapt to new environments and solve complex problems.

AI Thinking, Fast and Slow

The theories described in the previous section, as well as their connections, shed some light on which competencies provide humans the ability to solve a diverse set of simple and complex problems; understand broad contexts robustly; adapt readily at short time scales (rather than evolutionary time scales) to new scenarios and environmental conditions; and ultimately cooperate at scale. These competences include abstraction, generalization, causal analysis and planning, attention, locality, skills combination, introspection, and various forms of implicit and explicit knowledge.

We now discuss some of these concepts and capabilities in the AI context, with the aim to identify central research questions that could help understand how to equip AI with the capabilities that it still lacks. The intent is to stimulate the AI research community to collectively study such questions, possibly raise other ones, and find appropriate answers.

System 1 and System 2 in AI

It is possible to draw a very loose parallel between the capabilities included in system 1 and system 2, and the two main lines of work in AI: machine learning (ML) and symbolic logic reasoning, or data-driven vs knowledge-driven AI. System 1 seems to have similar traits as ML, with its ability to build (possibly imprecise and biased) models from sensory data. For example, perception activities, such as seeing and reading, that are currently addressed with ML techniques, are usually handled by human's system 1. However, system 1 exploits its ability to grasp basic notions of causality to build such models, and this does not seem to be present, at least for now, in ML. Another system 1 trait that is lacking in ML is common-sense reasoning (that of course is also related to causality). On the other hand, system 2's capability to solve complex problems seems to be related to AI techniques based on logic, search, optimization, and planning, and employing explicit knowledge, symbols, and high-level concepts.

However, as noted earlier, these two broad sets of capabilities are usually symbiotic, which supports hybrid neuro-symbolic approaches in AI.

Research Questions 1 *Should we clearly identify the AI system 1 and system 2 capabilities? What would their features be? Should there be just two sets of capabilities, or more?*

While the human system 2 is sequential, and its activation requires our full attention, it does not necessarily have to be this way in a machine.

Research Questions 2 *Is the sequentiality of system 2 a bug or a feature? Should we carry it over to machines or should we exploit the possibility of having several parallel threads performing system 2 reasoning? Would this possibility, together with the greater computing power of machines compared to humans, compensate for the lack of other capabilities in AI?*

Consider an AI system that is based on the two main sets of capabilities provided by ML and symbolic AI. Usually the quality of ML approaches is measured by the degree to which they can achieve the desired result, e.g., accuracy, precision, recall, F1 score. On the other hand, the quality of a symbolic reasoning process is usually measured by the correctness of its conclusions, e.g., the satisfaction of a set of constraints. In a combined system1/system2 AI system which switches dynamically between the two opportunistically, or combines their capabilities, selecting evaluation measures and methods is an open research question. See (Ribeiro et al. 2020) for a recent framework to evaluate NLP systems regardless of task and implementation.

Research Questions 3 *What are the metrics to evaluate the quality of a hybrid system 1/ system 2 (or ML/ symbolic) AI system? Should these metrics be different for different tasks and combination approaches?*

Introspection and Governance

Our aim is to build machines that are able to autonomously figure out how to utilize their cognitive resources to solve any challenge that comes their way. To achieve this, a subjective (if imprecise) model of agents' minds, such as the one underlying the concept of M-consciousness, is probably needed in order to identify the relevant factors for assessing the features of the problem. For example, if the scenario is cognitively easy or difficult, if it is risky to make a mistake, if solving it gives a high reward, etc., as well as to assess the competencies of the mechanistic system's components (the I-consciousness) to solve the problem, and select the most appropriate solution procedure.

Research Questions 4 *How do we define AI's introspection in terms of I-consciousness and M-consciousness? Can M-consciousness alone solve complex problems that in humans need also sophisticated levels of I-consciousness? How is introspection linked to autonomy and human-machine teaming? Is introspection a binary concept or should it have several levels? If in levels, how can we associate different capabilities to the various levels? Should we have models of the agents' minds (those supporting M-consciousness) at different levels of precision/fidelity?*

Through introspection, an AI system should be able to recognize when it needs to deploy capabilities from system 1 or system 2, and also when to switch between the two systems. In humans, such a switch depends on many factors, including framing, priming, and the context of the decision environment. In an AI system, there seems to be other factors that make up the sufficient conditions for this switch. For example, when system 1 cannot find a solution, e.g., because we are in a new environment, or when there are several

alternatives that must be considered before choosing an action. Or also when we are facing a very high stake decision: a mistake may confer dangerous consequences according to some high-weight criterion. Or also when the decision needs to be explainable to a third party.

Research Questions 5 *How do we model the governance of system 1 and system 2 reasoning in AI? When to switch between the two, or combine them? Which factors trigger the switch? How to define either the divergence between or absence of potential solving procedures? How to recognize the presence of these conditions? How should system 2 act once the switch is triggered? When should problems, that until then are handled by system 2, be instead deferred to be solved by system 1 hereafter? How to interject the notion of time in the governance framework?*

Divergence Resolution and Ethical Reasoning

As noted above, the notion of divergence is important for the governance framework. But it is also fundamentally tied to system 2 capabilities. When there are two or more diverging policies or heuristics offered up by system 1, and it is important to solve the problem correctly, one needs to carefully explore the main pros and cons, as well as the possible consequences of our actions, before deciding on the policy to adopt, or we need to define a new policy. This careful reasoning, based on inference and counter-factuals, seems to be a system 2 capability, supported by system 1's heuristics to reduce the search space.

Preliminary work has investigated how to incorporate possibly divergent value functions, say the tension between immediate reward and long term constraints imposed by behavioral norms in a variety of reinforcement learning settings (Balakrishnan et al. 2019a, 2018, 2019b) and probabilistic planning (Noothigattu et al. 2019a,b).

Research Questions 6 *How can we leverage a model based on system 1 and system 2 capabilities in AI to understand and reason in complex environments when we may have competing priorities? Which part of the AI system recognizes the divergence and, supported by the introspection, identifies the best solving procedure? How to provide the capability to build new solving procedures if the evaluation of the existing ones shows that they are not able to solve the given problem?*

Complex constraints and competing objectives generate dilemma that humans confront every day (Rossi and Mattei 2019). A typical example is when we need to make a decision that has some ethical component, such as deciding if a certain action is morally acceptable or not. Usually we simulate various approaches to the given question, based for example on deontology or consequentialism (Kagan 2018), and then we make our moral decision. Preliminary work along these lines attempts to understand how humans switch between these different approaches in judging the morality of an action (Rossi and Loreggia 2019; Awad et al. 2020), how divergence can be defined in an ethical context (Loreggia et al. 2018a), and how to use such notion of divergence to define methodologies for ethical reasoning in AI (Loreggia et al. 2018b,c, 2019).

Research Questions 7 *Which capabilities are needed to perform various forms of moral judging and decision making? Do they require a specific notion of divergence? How do we model and deploy possibly conflicting normative ethical theories in AI? Are the various ethics theories tied to either system 1 or system 2, or do they need a combination of the two systems?*

Abstraction, Generalization, and Knowledge

In order to adapt to a new environment, an agent needs to be aware of its competencies (which refers to the notion of introspection considered earlier) and know how to deploy its skills and problem solving capabilities, were possibly acquired in another environment, into the new one. This process is supported by the ability to abstract and generalize the specific skills and their intervention target object, in order to subsequently specialize them again in the new scenario. So, adaptability involves first an abstraction/generalization phase and then an instantiation phase. Of course both these steps are not random, but are guided by the ability to recognize the similarities and differences between the two environments, and to use this knowledge to decide what to temporarily forget during the abstraction step.

This conjecture is aligned with the so-called consciousness prior theory (Bengio 2017). Existing and well established abstraction theories (Cousot and Cousot 1977) and studies of various forms of abstraction in AI (Zucker 2003) can be useful, coupled with attention mechanisms (Vaswani et al. 2017).

Research Questions 8 *How to define abstraction / generalization mechanisms that are guided by a notion of attention and pass from the raw data level to a more abstract level? How to (dynamically and contextually) know what to forget from the input data during the abstraction step? Should we keep knowledge at various level of abstractions, or just raw data and fully explicit high-level knowledge? What does it mean for knowledge to be explicit: is it related to the presence of metadata, of structured knowledge graphs, or of language-related entities?*

Multi-agent Environment, Epistemic Reasoning, and Architectural Support

Humans live in societies, and their individual decisions are linked to their perception of reality, which includes the world, the other agents, and themselves. These models are not perfect but they are used to make informed decisions, as well as a test bed to evaluate consequences of alternative decisions. However, such models are not based on exact knowledge, but rather on approximate information on the world and on beliefs on what others know and believe. Multi-agent epistemic reasoning (Fagin et al. 2003) is what humans engage on when making complex decisions using system 2, unless such decisions are completely independent of other agents or their consequences are irrelevant.

Research Questions 9 *In a multi-agent view of several AI systems communicating and learning from each other, how to exploit/adapt current results on epistemic reasoning and planning to build/learn models of the world and of others?*

Following Minsky’s theory of the society of mind (Minsky 1986), as well as Brook’s subsumption architecture approach (Brooks 1990), we believe that AI systems should include several independent simple components, to be triggered when needed according to a governance model. This seems to imply that the best architectural support for this vision of AI is a multi-agent architecture where the individual agents focus on specific skills and problems, act asynchronously, contribute to building models of the world, of other AI systems, and of self, and can be combined in many different ways.

Research Questions 10 *What are the needed architectural choices that best support the above vision of the future of AI?*

Conclusions

This brief paper describes our vision of a multi-disciplinary research effort that leverages cognitive theories of human decision making in order to provide AI with desired human capabilities that are still lacking in current AI system.

We hope the AI research community will join us in this endeavor and will work with us to respond to the research questions we outlined (and possibly many others), as well as to define, support, and structure a coordinated and global research community that is aligned to the vision outlined in this paper.

References

- Awad, E.; Levine, S.; Loreggia, A.; Mattei, N.; Rahwan, I.; Rossi, F.; Talamadupula, K.; Tenenbaum, J.; and Kleiman-Weiner, M. 2020. When Is It Morally Acceptable to Break the Rules? A Preference-Based Approach. In *12th Multi-disciplinary Workshop on Advances in Preference Handling (MPREF 2020)*.
- Balakrishnan, A.; Bouneffouf, D.; Mattei, N.; and Rossi, F. 2018. Using Contextual Bandits with Behavioral Constraints for Constrained Online Movie Recommendation. In *Proceedings of the 27th International Joint Conference on Artificial Intelligence (IJCAI)*.
- Balakrishnan, A.; Bouneffouf, D.; Mattei, N.; and Rossi, F. 2019a. Incorporating Behavioral Constraints in Online AI Systems. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI)*.
- Balakrishnan, A.; Bouneffouf, D.; Mattei, N.; and Rossi, F. 2019b. Using multi-armed bandits to learn ethical priorities for online AI systems. *IBM J. Res. Dev.* 63(4/5): 1:1–1:13.
- Bengio, Y. 2017. The consciousness prior. *arXiv preprint arXiv:1709.08568*.
- Brooks, R. A. 1990. Elephants don’t play chess. *Robotics and autonomous systems* 6(1-2): 3–15.
- Cousot, P.; and Cousot, R. 1977. Abstract interpretation: a unified lattice model for static analysis of programs by construction or approximation of fixpoints. In *Conference Record of the Fourth Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages*, 238–252. Los Angeles, California: ACM Press, New York, NY.
- d’Avila Garcez, A. S.; and Besold, T. R. 2019. Special Issue: Neural-Symbolic Learning and Reasoning (NeSy’18). *IfCoLog Journal of Logics and their Applications* 6(4): 609–610.
- Fagin, R.; Moses, Y.; Halpern, J. Y.; and Vardi, M. Y. 2003. *Reasoning about knowledge*. MIT press.
- Graziano, M. S. 2013. *Consciousness and the social brain*. Oxford University Press.
- Graziano, M. S.; Guterstam, A.; Bio, B. J.; and Wilterson, A. I. 2020. Toward a standard model of consciousness: Reconciling the attention schema, global workspace, higher-order thought, and illusionist theories. *Cognitive Neuropsychology* 37(3-4): 155–172.
- Harari, Y. N. 2014. *Sapiens: A Brief History of Humankind*. Random House.
- Kagan, S. 2018. *Normative Ethics*. Routledge.
- Kahneman, D. 2011. *Thinking, Fast and Slow*. Macmillan.
- Loreggia, A.; Mattei, N.; Rossi, F.; and Venable, K. B. 2018a. On the Distance Between CP-nets. In *Proceedings of the 17th International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*.
- Loreggia, A.; Mattei, N.; Rossi, F.; and Venable, K. B. 2018b. Preferences and Ethical Principles in Decision Making. In *Proc 1st AIES*.
- Loreggia, A.; Mattei, N.; Rossi, F.; and Venable, K. B. 2018c. Value Alignment Via Tractable Preference Distance. In Yampolskiy, R. V., ed., *Artificial Intelligence Safety and Security*, chapter 18. CRC Press.
- Loreggia, A.; Mattei, N.; Rossi, F.; and Venable, K. B. 2019. CPMetric: Deep Siamese Networks for Metric Learning on Structured Preferences. In *International Joint Conference on Artificial Intelligence*, 217–234. Springer, Cham.
- Marcus, G. 2020. The next decade in AI: Four steps towards robust artificial intelligence. *arXiv preprint arXiv:2002.06177*.
- Minsky, M. 1986. *The society of mind*. Simon and Schuster.
- Noothigattu, R.; Bouneffouf, D.; Mattei, N.; Chandra, R.; Madan, P.; Varshney, K.; Campbell, M.; Singh, M.; and Rossi, F. 2019a. Teaching AI Agents Ethical Values Using Reinforcement Learning and Policy Orchestration. In *Proceedings of the 28th International Joint Conference on Artificial Intelligence (IJCAI)*.
- Noothigattu, R.; Bouneffouf, D.; Mattei, N.; Chandra, R.; Madan, P.; Varshney, K. R.; Campbell, M.; Singh, M.; and Rossi, F. 2019b. Teaching AI agents ethical values using reinforcement learning and policy orchestration. *IBM J. Res. Dev.* 63(4/5): 2:1–2:9.
- Ribeiro, M. T.; Wu, T.; Guestrin, C.; and Singh, S. 2020. Beyond Accuracy: Behavioral Testing of NLP models with CheckList. In *58th Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Rossi, F.; and Loreggia, A. 2019. Preferences and ethical priorities: Thinking fast and slow in AI. In *Proceedings*

of the 18th International Joint Conference on Autonomous Agents and Multi-Agent Systems (AAMAS), 3–4.

Rossi, F.; and Mattei, N. 2019. Building Ethically Bounded AI. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence (AAAI)*.

Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L.; and Polosukhin, I. 2017. Attention Is All You Need. *CoRR* abs/1706.03762. URL <http://arxiv.org/abs/1706.03762>.

Zucker, J.-D. 2003. A grounded theory of abstraction in artificial intelligence. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences* 358(1435): 1293–1309.