

# The Cognitive Hourglass: Agent Abstractions in the Large Models Era

Blue Sky Ideas Track

Alessandro Ricci  
University of Bologna  
Cesena, Italy  
a.ricci@unibo.it

Stefano Mariani  
University of Modena and Reggio  
Emilia  
Reggio Emilia, Italy  
stefano.mariani@unimore.it

Franco Zambonelli  
University of Modena and Reggio  
Emilia  
Reggio Emilia, Italy  
franco.zambonelli@unimore.it

Samuele Burattini  
University of Bologna  
Cesena, Italy  
samuele.burattini@unibo.it

Cristiano Castelfranchi  
Italian Research Council  
Rome, Italy  
cristiano.castelfranchi@istc.cnr.it

## ABSTRACT

Recent advances in AI are driving an unprecedented and fast-paced development of myriads of powerful agent tools and applications, mostly based on generative AI technologies such as Large Language/Multi-modal/Agent Models. However, despite many proposals in that direction, the lack of a sound set of usable engineering abstractions hinders the possibility of methodically engineering complex agent-based applications, also due to the gap between cognitive agent-based concepts and LLMs’ behavioural patterns. We argue that such a set of abstractions should constitute the *narrow neck* of an indispensable “cognitive hourglass”: a level of abstraction that is meant to be useful for humans to understand/design/control agents and MAS, regardless of the specific AI technologies adopted at the implementation level and of the specific application context. Here, we elaborate on the idea of the cognitive hourglass, motivate its need, sketch its envisioned architecture, and identify the research challenges for its realisation.

## KEYWORDS

Agent systems engineering; LLMs; Cognition; Hourglass model

### ACM Reference Format:

Alessandro Ricci, Stefano Mariani, Franco Zambonelli, Samuele Burattini, and Cristiano Castelfranchi. 2024. The Cognitive Hourglass: Agent Abstractions in the Large Models Era: Blue Sky Ideas Track. In *Proc. of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2024), Auckland, New Zealand, May 6 – 10, 2024*, IFAAMAS, 6 pages.

## 1 INTRODUCTION

The hourglass model [2] has been adopted in computer science and engineering as a conceptual metaphor for describing complex layered systems, like network protocols [1], and as a blueprint to drive their design [31, 43]. The model provides effective constraints in the

design of open systems, enabling the development of an open set of applications on top of an open set of supporting technologies and services. A key element of the model is a narrow layer as the *neck* of the hourglass, including a selected set of functional abstractions separating and mediating between the upper (application) layers and the bottom (technology, implementation) ones.

As an example, consider the so-called IP Hourglass (Fig. 1), which is a model of modern IP-based networked systems. Several network technologies and protocols (bottom layers) exist that are exploited by many high-level application protocols and systems (upper layers). The Internet Protocol acts as a glue between the upper and the lower layers by making available a simple uniform communication protocol, independent of the actual network technology, that can be used to build any distributed applications.

The hourglass model is useful for understanding, discussing, and governing the recent fast-paced advances in AI that are dramatically increasing the spectrum of technologies that can be exploited to build autonomous agents and Multi-Agent Systems (MAS) and applications—in addition to the many already assessed systems and languages [5, 49]. Large neural models and generative AI technologies such as Large-Language Models (LLMs) [10], Large-Multimodal Models (LMMs) [11], and Foundation Models [52] are witnessing a

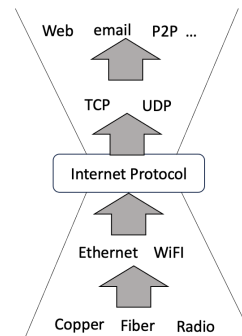


Figure 1: The IP Hourglass. The narrow neck gives the minimum set of abstractions and mechanisms to build upper layer services on top of the lower layer ones.



This work is licensed under a Creative Commons Attribution International 4.0 License.

Proc. of the 23rd International Conference on Autonomous Agents and Multiagent Systems (AAMAS 2024), N. Alechina, V. Dignum, M. Dastani, J.S. Sichman (eds.), May 6 – 10, 2024, Auckland, New Zealand. © 2024 International Foundation for Autonomous Agents and Multiagent Systems (www.ifaamas.org).

huge momentum [15] and started to be adopted to build different kinds of “agents” [11, 51]: from generalist agents such as Gato [36] to *generative agents* based on LLMs, such as AutoGPT [30]—called Large Agent Models (LAMs).

Such agents show remarkable capabilities in terms of information management and reasoning. However, how to *methodically* exploit such capabilities during the process of designing and building MAS is yet unclear. On the one hand, LLM-based agents rely on properly “prompting” (i.e., triggering) agent behaviours, an activity which is currently far from being methodical and reproducible [25, 47]. On the other hand, the current way of building and exploiting LLMs and LAMs is leading to a sort of “eliminativism”: deeming higher-level abstractions unnecessary once lower-level ones have been fully understood. In particular, cognitive concepts [12, 13, 22, 32] that are pillars for the understanding and engineering of agent systems [6], seem to be increasingly ignored.

In contrast, in this article, we argue that such concepts constitute the indispensable *neck* of the *cognitive hourglass*, that is, the fundamental human-compatible [39] level of abstraction necessary for humans to understand/design/govern agents and MAS at the application level – the top of our hourglass – regardless of the specific AI technologies adopted at the implementation level—the bottom of our hourglass.

## 2 COGNITION AND ELIMINATIVISM

Some literature about generative AI, including LLMs, LMAs, and LAMs, seem to foster the idea that every conceivable “intelligent” application, and software agents themselves, can be directly built on top of their behavioural patterns and working mechanisms – such as prompt engineering, instruction-tuning, and generally speaking in-context learning [17, 27, 37] –, without the need for any scientific and engineering abstraction in between. This perspective nurtures a dangerous trend: *eliminativism* [28], as the attitude of deeming higher-level abstractions unnecessary once lower-level ones have been modelled and understood. This is a degenerate derivation of reductionism, that seeks to explain higher-level phenomena in terms of lower-level ones, but without neglecting the utility of those at the higher level—as they enable the expression of novel properties and interpretations of the lower ones.

In chemistry, for instance, bonding laws must find an “implementation” in terms of the underlying mechanisms of physics (reductionism), but no chemist would then throw away such laws to only think in terms of physics (eliminativism)! However, in computer science, a recent stream of publications about LLM-based agents [11, 51] seems inclined to disregard agent-oriented abstractions – such as Allen Newell’s Knowledge Level [32], Jennings in [22], as well as Dennett’s intentional stance [13] and Castelfranchi’s work [34] –, since they can anyway obtain agent-like capabilities without them being part of the engineering process [11, 26, 36].

Besides the fact that this claim is still under debate and evaluation [4, 19, 44], eliminativism is undesirable because the cognitive abstractions developed in the agent community, and more generally within the AI research landscape, have the unquestionable merit of having served us well in building (engineering purpose) and understanding (scientific purpose) systems (even biological ones)

while using the most suitable abstractions—for our own reasoning processes, as human beings.

Ignoring such abstractions creates a gap between the ones who build and use these systems and their modes of operation: the former typically reason in terms of goals, plans, actions and their consequences, beliefs, knowledge, cause-effect relations [33], etc., whereas the latter require manipulation of prompts, linguistic patterns, examples, and blueprints of desired behaviours. From the engineering viewpoint, a consequence is a limited capability of building systems by composition of simpler parts, a foundational basis of any engineering discipline. From the scientific standpoint, the lack of a layered set of abstractions hinders understanding of a given system in its multiple nuances in terms of properties and behaviours (e.g. focusing on chemical or physical properties of a given material).

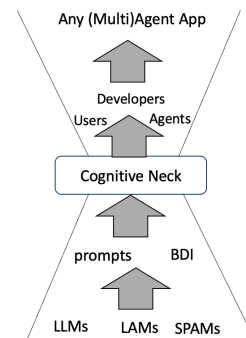
The latter ability, in particular, is essential not only under the lens of eXplainable AI [20], safety, and alignment in AI (in one word, “Human-centred AI” [39, 40]), but also when considering transferability of concepts across domains, and when shifting attention from an individual agent to a MAS. There, ascribing to others their own mental states is fundamental to promote collaboration and requires mind-reading [9] others (i.e. prediction of the mental states of other agents) with the proper abstractions.

In the following section, we propose the adoption of a *cognitive hourglass* to align generative technologies and agent-oriented abstractions to stay away from the described eliminativism.

## 3 THE COGNITIVE HOURGLASS

To synergistically exploit LLMs and LAMs, and the (necessary) knowledge-based and symbolic (i.e., cognitive) computational tools, it is needed to identify the “neck” of an agent-oriented hourglass. Such neck should enable to uniformly exploit any available enabling technology (bottom layers of the hourglass) to provide services for building, controlling, and understanding, autonomous agents and MAS (upper layers of the hourglass)— Fig. 2.

For such concepts to be effectively usable both at the human and at the software level, they should abstract away implementation and technical details. In addition, they should be expressive enough to allow for specifying any kind of structure and behaviour



**Figure 2: The Cognitive Hourglass. The narrow neck enables to uniformly exploit the lower layer models and technologies to model, understand, and govern the upper layer ones.**

of systems modelled in terms of MAS, being them designed systems/applications, implemented agents systems, or even simulated agents systems. In other words, the neck of the hourglass should be a *cognitive abstraction gate*, with cognitive concepts being the “sand” flowing up and down the hourglass. As such, the neck enables a bi-directional flow of symbolic knowledge and concepts between the upper and the lower levels of the hourglass, such that the “services” of the enabling technologies below the neck can be effectively activated and exploited from above the neck. Let us now elaborate on our envisioned layered structure of the cognitive hourglass.

The lowest level is where to accommodate any of the available technologies and implementation of agents and MAS. For instance: LLM-based agents [51], LAMs [36], agents built with specific agent platforms like Jade [3] or Jason [7], or whatever other type of SPAMs (special-purpose agent models). Just above it, to be able to exploit all such technologies, one must accommodate the various types of “mechanisms” such technologies use to interact with the world. E.g., prompts for LLMs, Beliefs Desires and Intentions (BDI) scripts or alike for the case of cognitive agent architectures. At the upper levels, there are human actors – either as simple users of the below technologies or as engineers in charge of exploiting them to design and develop agent and multi-agent applications – as well as software agents—i.e., agents as components of some multi-agent application and exploiting the below technologies to augment/outsource their capabilities. The cognitive neck, acting as a gateway between the lower and the upper levels, should be flexible and expressive enough to provide access to the lower levels (i.e., to the services provided by the implemented agent technologies) in a simple yet comprehensive and comprehensible way—to effectively support the development of agent systems and their empowerment.

For instance, the cognitive neck could provide abstractions to re-interpret prompt engineering techniques as *argumentation processes* [45], made up of commitments, expectations, roles, scopes, and similar concepts, that humans exploit (either consciously or not) while engaging in dialogues with each other. Or, *theory of mind* [35] concepts can be used to interact more proficiently with LLMs and other agents by ascribing mental states to them, instead of reasoning in terms of behavioural patterns, procedures, and similar low-level non-cognitive terms.

How an actual cognitive neck should be made, what interfaces it should expose, and how it could be effectively exploited, is impossible to detail in this paper. This is indeed one of the key research questions this article intends to raise, not answer (yet). Nevertheless, some key cognitive concepts include:

- *Wishes* (aka desires, goals) expressed from the upper levels to the cognitive neck. What one wishes is typically the state of affairs that one (agent, human user, or agent developer) wants to achieve. In response to wishes, it is expected that the cognitive neck will reply with.
- *Hows* (aka plans), that is a proposal for actions to be put at work to achieve wishes, possibly in respect of constraints.
- *Constraints* (aka safety and liveness rules). These can be expressed from the upper to the lower levels as specific instructions on how plans should be built, but also be communicated from the lower to the upper levels, if such constraints emerge during the building of plans.

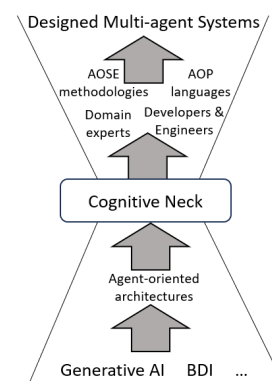
- *Whys* (aka explanations), mostly for letting the lower levels motivate (whenever needed) the responses provided to the upper levels, e.g., in the form of causal models [29, 33]. However, in some cases, whys can be used to let the upper levels motivate their requests to guide/influence the behaviour of the lower levels.
- *Whats* (aka facts or beliefs). Through the cognitive neck, knowledge can transit (typically on request) to the upper level about things known at the lower level. However, one could also think of knowledge transiting from the upper to the lower level to influence its behaviour.

It is worth emphasising that the cognitive hourglass can also play an important role in enabling multi-agent communication and the involvement of agents in interaction and negotiation protocols. Indeed, cognitive approaches to multi-agent interactions assume that messages exchanged in the context of a distributed decision-making process have cognitive context [23, 41], and are aimed at transferring knowledge, informing about plans or desires, or proposing courses of action. Again, although tools for developing MAS with LLM-based agents are being proposed [50, 53], interactions between agents are limited to prompting conversations and do not account for the specific cognitive meaning of messages.

#### 4 TOWARDS AN INTEGRATED FRAMEWORK

The cognitive hourglass can be a suitable methodological and practical framework for the engineering of agent-based systems, affecting both *design-time* scenarios, involving developers and engineers (Fig. 3), and *run-time* ones, involving users and application agents (Fig. 4).

At design time, the cognitive hourglass can support developers and engineers in conceiving agent-based architectures that are instrumented to adopt high-level agent-oriented software engineering methodologies [18] and possibly agent-oriented programming languages [5] designed upon the abstractions defined in the cognitive neck. Some existing agent-oriented software methodologies such as TROPOS [8] and the more recent TDF [14] naturally fit in the picture, since they have been explicitly designed to be at

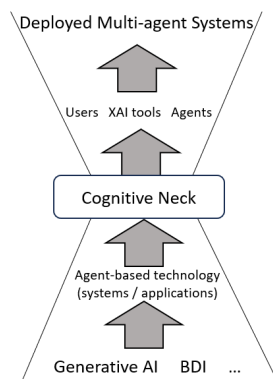


**Figure 3: The Cognitive Hourglass at design-time. Abstractions in the neck enable to design upper layer tools and systems uniformly, regardless of the ones at the lower layer.**

the Knowledge Level. According to this design-time view, it is interesting to devise – as future research directions – new agent architectures conceived to be compatible with the cognitive neck, eventually integrating different approaches and technologies. An example is *generative BDI architectures*, as agent architectures based on the BDI model and reasoning cycle, but integrating generative AI technologies in key steps of the cycle and for, e.g., the management of beliefs, goals, and intentions.

The cognitive hourglass could also affect the overall process of developing agents and MAS, in particular *learning-based* processes. Learning and machine learning techniques – Reinforcement Learning (RL) and Deep RL in particular – are increasingly adopted as reference approaches to develop agents in various domains. In recent works, learning becomes a core ingredient of novel engineering processes based on self-development [24, 48], and adopting developmental learning methods that make agent/MAS engineering similar to an education process [38], as defined by educational theory in pedagogy. Given the hourglass, the high-level learning-based development processes adopted to “grow up” agents or refine their skills (top) should be based on concepts defined in the cognitive neck, allowing to frame, at the proper level of abstraction, the specific learning techniques and technologies adopted (bottom).

At run-time, the conceptual framework considers that the cognitive neck could become a practical software tool for enabling users and agents to access a suite of cognitive services. On the one hand, human users can exploit the cognitive neck to interface with the world of existing agent-based systems and tools. Regardless of the specific technology and implementation, the cognitive neck makes sure that the features of the lower layers will be made available in the form of cognitive and human-understandable concepts. On the other hand, agents developed to serve in specific applications and systems will be able to exploit the cognitive neck as a run-time tool to empower their capabilities, for example by dynamically accessing agents (LLM-based but not necessarily) to request planning services or general information. However, as anticipated in the previous section, the cognitive neck could also become a powerful tool for supporting cognitive inter-agent interactions.



**Figure 4: The Cognitive Hourglass at run-time. Through the neck, humans and software agents can exploit cognitive abstractions to empower themselves, inspect, and coordinate.**

For both human and agent users, the cognitive abstraction level of the neck makes it possible to instrument tools that can make the behaviour of the lower levels transparent and explainable. In other words, the cognitive neck could explicitly define the conceptual interface upon which explainability tools can be designed and exploited. This deeply relates to the scientific viewpoint: the cognitive neck enables a principled understanding of the working mechanisms and usage patterns of what is below the neck, similarly to the layered understanding scientists have of biological, chemical, and physical systems (reductionism, not eliminativism!).

In all these scenarios the cognitive neck plays the role of an abstraction barrier, allowing for the development and integration of different kinds of heterogeneous agent-based technologies on the bottom, and the development of proper tools that would allow humans – users, domain experts, engineers – and agents to have a full understanding and control of the system.

## 5 CONCLUSIONS AND OPEN CHALLENGES

In this article, we motivated why cognitive abstractions and concepts (from Allen Newell’s Knowledge Level to Castelfranchi’s work on autonomous goal-directed behaviour) should play a primary role in agent systems engineering. Even, and especially, in the presence of the recent LLMs-enabled agent tools, cognitive abstractions have to constitute the *narrow neck* of a “cognitive hourglass”. That is a level of abstraction useful for humans to understand/design/control agents and MAS, regardless of the specific AI technologies adopted at the implementation level and of the specific application context. Yet, for the cognitive hourglass to become a practical and usable tool, there are several key challenges to be faced.

First, the concepts and principles inside the cognitive neck must be identified, to make it both a usable conceptual tool for developers and a software service layer for agents. This includes the possibility of exploiting the cognitive neck to support flexible interactions in multi-agent systems. Second, proper mappings must be developed that – despite the abstraction barrier – allow exploitation of the capabilities provided by the bottom/enabling levels while preserving the property of being “human-compatible” – both for users and engineers. Considering the multiple dimensions that are important in the case of MAS – such as the social and organisational dimensions [16], and the environment dimension as well [46] – demands further studies. Finally, proper forms of integration between cognitive agents and generative AI should be identified. Some relevant efforts in that direction can already be found: [21] explores the use of language models as a source of task knowledge in cognitive agents/systems; [42] proposes a systematic framework called Cognitive Architectures for Language Agents (CoALA), useful for both organizing existing literature on generative agents and identifying directions towards more capable agents, including features as found in cognitive agents and architectures.

## ACKNOWLEDGMENTS

Partially supported by the European Union NextGenerationEU through the Italian “National Recovery and Resilience Plan” (PNRR) Mission 4, Component 2, Investment 3.3 (DM 352/2022) and in collaboration with Azienda Unità Sanitaria Locale (AUSL) della Romagna.

## REFERENCES

- [1] Saamer Akhshabi and Constantine Dovrolis. 2011. The evolution of layered protocol stacks leads to an hourglass-shaped architecture. In *Proceedings of the ACM SIGCOMM Conference*. 206–217.
- [2] Micah Beck. 2019. On the Hourglass Model. *Communication of ACM* 62, 7 (jun 2019), 48–57. <https://doi.org/10.1145/3274770>
- [3] Fabio Luigi Bellifemine, Giovanni Caire, and Dominic Greenwood. 2007. *Developing multi-agent systems with JADE*. John Wiley & Sons.
- [4] Lukas Berglund, Meg Tong, Max Kaufmann, Mikita Balesni, Asa Cooper Stickland, Tomasz Korbak, and Owain Evans. 2023. The Reversal Curse: LLMs trained on "A is B" fail to learn "B is A". arXiv:2309.12288 [cs.CL]
- [5] Rafael H Bordini, Lars Braubach, Mehdi Dastani, A El F Seghrouchni, Jorge J Gomez-Sanz, Joao Leite, Gregory O'Hare, Alexander Pokahr, and Alessandro Ricci. 2006. A survey of programming languages and platforms for multi-agent systems. *Informatica* 30, 1 (2006).
- [6] Rafael H. Bordini, Amal El Fallah Seghrouchni, Koen Hindriks, Brian Logan, and Alessandro Ricci. 2020. Agent Programming in the Cognitive Era. *Autonomous Agents and Multi-Agent Systems* 34, 2 (may 2020), 31. <https://doi.org/10.1007/s10458-020-09453-y>
- [7] Rafael H Bordini, Jomi Fred Hübner, and Michael Wooldridge. 2007. *Programming multi-agent systems in AgentSpeak using Jason*. John Wiley & Sons.
- [8] Paolo Bresciani, Anna Perini, Paolo Giorgini, Fausto Giunchiglia, and John Mylopoulos. 2004. Tropos: An Agent-Oriented Software Development Methodology. *Autonomous Agents and Multi Agent Systems* 8, 3 (2004), 203–236. <http://dblp.uni-trier.de/db/journals/aamas/aamas8.html#BrescianiPGGM04>
- [9] Cristiano Castelfranchi. 2019. "Mind Reading": How and for what? *Annals of Cognitive Science* (2019). <https://api.semanticscholar.org/CorpusID:211167867>
- [10] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2023. A Survey on Evaluation of Large Language Models. arXiv:2307.03109 [cs.CL]
- [11] Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, Cunxiang Wang, Yidong Wang, Wei Ye, Yue Zhang, Yi Chang, Philip S. Yu, Qiang Yang, and Xing Xie. 2024. A Survey on Evaluation of Large Language Models. *ACM Trans. Intell. Syst. Technol.* (jan 2024). <https://doi.org/10.1145/3641289> Just Accepted.
- [12] Rosaria Conte and Cristiano Castelfranchi. 1995. *Cognitive and social action*. Psychology Press.
- [13] Daniel C. Dennett. 1987. *The Intentional Stance*. The MIT Press, Cambridge, MA.
- [14] Rick Evertsz, John Thangarajah, and Thanh Ly. 2019. *Practical Modelling of Dynamic Decision Making* (1st ed.). Springer Publishing Company, Incorporated.
- [15] Lizhou Fan, Lingyao Li, Zihui Ma, Sanggyu Lee, Huiyi Yu, and Libby Hemphill. 2023. A Bibliometric Review of Large Language Models Research from 2017 to 2023. arXiv:2304.02020 [cs.DL]
- [16] Jacques Ferber, Olivier Gutknecht, and Fabien Michel. 2003. From agents to organizations: an organizational view of multi-agent systems. In *International workshop on agent-oriented software engineering*. Springer, 214–230.
- [17] Andrew Gao. 2023. Prompt Engineering for Large Language Models. Available at SSRN 4504303 (2023).
- [18] Brian Henderson-Sellers and Paolo Giorgini. 2005. *Agent-oriented methodologies*. Igi Global.
- [19] Damian Hodel and Jevin West. 2023. Response: Emergent analogical reasoning in large language models. *CoRR abs/2308.16118* (2023). <https://doi.org/10.48550/arXiv.2308.16118>
- [20] Andreas Holzinger, Peter Kieseberg, Edgar R. Weippl, and A Min Tjoa. 2018. Current Advances, Trends and Challenges of Machine Learning and Knowledge Extraction: From Machine Learning to Explainable AI. In *Machine Learning and Knowledge Extraction - Second IFIP TC 5, TC 8/WG 8.4, 8.9, TC 12/WG 12.9 International Cross-Domain Conference, CD-MAKE, Hamburg, Germany, August 27-30, Proceedings (Lecture Notes in Computer Science, Vol. 11015)*. Springer, 1–8. [https://doi.org/10.1007/978-3-319-99740-7\\_1](https://doi.org/10.1007/978-3-319-99740-7_1)
- [21] Robert E. Wray III, James R. Kirk, and John E. Laird. 2022. Language Models as a Knowledge Source for Cognitive Agents. In *Proceedings of the Ninth Annual Conference on Advances in Cognitive Systems. Ninth Annual Conference on Advances in Cognitive Systems*.
- [22] Nicholas R. Jennings. 2000. On Agent-Based Software Engineering. *Artificial Intelligence* 117, 2 (mar 2000), 277–296. [https://doi.org/10.1016/S0004-3702\(99\)00107-1](https://doi.org/10.1016/S0004-3702(99)00107-1)
- [23] Yannis Labrou, Tim Finin, and Yun Peng. 1999. Agent communication languages: The current landscape. *IEEE Intelligent Systems and Their Applications* 14, 2 (1999), 45–52.
- [24] Marco Lippi, Stefano Mariani, Matteo Martinelli, and Franco Zambonelli. 2022. Individual and Collective Self-Development: Concepts and Challenges. In *17th Conference on Computer Science and Intelligence Systems (FedCSIS)*. IEEE, 15–21.
- [25] Vivian Liu and Lydia B Chilton. 2022. Design guidelines for prompt engineering text-to-image generative models. In *Proceedings of the CHI Conference on Human Factors in Computing Systems*. 1–23.
- [26] Xiao Liu, Hao Yu, Hanchen Zhang, Yifan Xu, Xuanyu Lei, Hanyu Lai, Yu Gu, Hengliang Ding, Kaiwen Men, Kejuan Yang, Shudan Zhang, Xiang Deng, Aohan Zeng, Zhengxiao Du, Chenhui Zhang, Sheng Shen, Tianjun Zhang, Yu Su, Huan Sun, Minlie Huang, Yuxiao Dong, and Jie Tang. 2023. AgentBench: Evaluating LLMs as Agents. arXiv:2308.03688 [cs.AI]
- [27] Sheng Lu, Irina Bigoulaeva, Rachneet Sachdeva, Harish Tayyar Madabushi, and Iryna Gurevych. 2023. Are Emergent Abilities in Large Language Models just In-Context Learning? arXiv:2309.01809 [cs.CL]
- [28] Francesco Mancini, Alessandra Mancini, and Cristiano Castelfranchi. 2022. Unhealthy mind in a healthy body: A criticism to eliminativism in psychopathology. *Frontiers in Psychiatry* 13 (2022). <https://doi.org/10.3389/fpsy.2022.889698>
- [29] Stefano Mariani, Pasquale Roseti, and Franco Zambonelli. 2023. Multi-agent Learning of Causal Networks in the Internet of Things. In *Advances in Practical Applications of Agents, Multi-Agent Systems, and Cognitive Mimetics. The PAAMS Collection - 21st International Conference, PAAMS 2023, Guimarães, Portugal, July 12-14, 2023, Proceedings*. 163–174.
- [30] Mindstreams. 2023. AutoGPT. <https://autogpt.net/>.
- [31] Tahir Nawaz Minhas and Markus Fiedler. 2013. Quality of experience hourglass model. In *International Conference on Computing, Management and Telecommunications (ComManTel)*. 87–92. <https://doi.org/10.1109/ComManTel.2013.6482371>
- [32] Allen Newell. 1982. The Knowledge Level. *Artificial Intelligence* 18, 1 (1982), 87–127. [https://doi.org/10.1016/0004-3702\(82\)90012-1](https://doi.org/10.1016/0004-3702(82)90012-1)
- [33] Judea Pearl and Dana Mackenzie. 2018. *The Book of Why: The New Science of Cause and Effect*. Basic Books, Inc.
- [34] Giovanni Pezzulo and Cristiano Castelfranchi. 2009. Thinking as the control of imagination: a conceptual framework for goal-directed systems. *Psychological Research Psychologische Forschung* 73, 4 (2009), 559–577. <https://doi.org/10.1007/s00426-009-0237-z>
- [35] Neil Rabinowitz, Frank Perbet, Francis Song, Chiyuan Zhang, S. M. Ali Eslami, and Matthew Botvinick. 2018. Machine Theory of Mind. In *Proceedings of the 35th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 80)*, Jennifer Dy and Andreas Krause (Eds.). PMLR, 4218–4227. <https://proceedings.mlr.press/v80/rabinowitz18a.html>
- [36] Scott E. Reed, Konrad Zolna, Emilio Parisotto, Sergio Gómez Colmenarejo, Alexander Novikov, Gabriel Barth-Maron, Mai Gimenez, Yury Sulsky, Jackie Kay, Jost Tobias Springenberg, Tom Eccles, Jake Bruce, Ali Razavi, Ashley Edwards, Nicolas Heess, Yutian Chen, Raia Hadsell, Oriol Vinyals, Mahyar Bordbar, and Nando de Freitas. 2023. A Generalist Agent. *Transactions on Machine Learning Research* (2023). <https://openreview.net/forum?id=1lk0KfHjvj>
- [37] Laria Reynolds and Kyle McDonell. 2021. Prompt Programming for Large Language Models: Beyond the Few-Shot Paradigm. In *Extended Abstracts of the CHI Conference on Human Factors in Computing Systems* (Yokohama, Japan). Association for Computing Machinery, New York, NY, USA, Article 314, 7 pages. <https://doi.org/10.1145/3411763.3451760>
- [38] Alessandro Ricci. 2022. "Go to the Children": Rethinking Intelligent Agent Design and Programming in a Developmental Learning Perspective. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems* (Virtual Event, New Zealand). International Foundation for Autonomous Agents and Multiagent Systems, Richland, SC, 1809–1813.
- [39] S. Russell. 2019. *Human Compatible: Artificial Intelligence and the Problem of Control*. Penguin Publishing Group. <https://books.google.it/books?id=M1eFDwAAQBA>
- [40] Ben Shneiderman. 2020. Bridging the Gap Between Ethics and Practice: Guidelines for Reliable, Safe, and Trustworthy Human-centered AI Systems. *ACM Transactions on Interactive Intelligent Systems* 10, 4 (2020), 26:1–26:31. <https://doi.org/10.1145/3419764>
- [41] Munindar P. Singh. 1998. Agent communication languages: Rethinking the principles. *Computer* 31, 12 (1998), 40–47.
- [42] Theodore R Sumers, Shunyu Yao, Karthik Narasimhan, and Thomas L Griffiths. 2023. Cognitive architectures for language agents. arXiv preprint arXiv:2309.02427 (2023).
- [43] Yosephine Susanto, Andrew G. Livingstone, Bee Chin Ng, and Erik Cambria. 2020. The Hourglass Model Revisited. *IEEE Intelligent Systems* 35, 5 (2020), 96–102. <https://doi.org/10.1109/MIS.2020.2992799>
- [44] Karthik Valmeekam, Sarath Sreedharan, Matthew Marquez, Alberto Olmo, and Subbarao Kambhampati. 2023. On the Planning Abilities of Large Language Models (A Critical Investigation with a Proposed Benchmark). arXiv:2302.06706 [cs.AI]
- [45] Douglas Neil Walton and Erik C. W. Krabbe. 1995. *Commitment in Dialogue: Basic Concepts of Interpersonal Reasoning*. State University of New York Press, Albany, NY, USA.
- [46] Danny Weyns, Andrea Omicini, and James Odell. 2007. Environment as a first class abstraction in multiagent systems. *Autonomous agents and multi-agent systems* 14 (2007), 5–30.
- [47] Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C Schmidt. 2023. A prompt pattern catalog to enhance prompt engineering with chatgpt. arXiv preprint arXiv:2302.11382 (2023).

- [48] Martin Wirsing and Lenz Belzner. 2023. Towards systematically engineering autonomous systems using reinforcement learning and planning. In *Analysis, Verification and Transformation for Declarative Programming and Intelligent Systems: Essays Dedicated to Manuel Hermenegildo on the Occasion of His 60th Birthday*. Springer, 281–306.
- [49] Zofia Wrona, Wojciech Buchwald, Maria Ganzha, Marcin Paprzycki, Florin Leon, Noman Noor, and Constantin-Valentin Pal. 2023. Overview of Software Agent Platforms Available in 2023. *Information* 14, 6 (2023). <https://doi.org/10.3390/info14060348>
- [50] Qingyun Wu, Gagan Bansal, Jieyu Zhang, Yiran Wu, Shaokun Zhang, Erkang (Eric) Zhu, Beibin Li, Li Jiang, Xiaoyun Zhang, and Chi Wang. 2023. *AutoGen: Enabling Next-Gen LLM Applications via Multi-Agent Conversation*. Technical Report MSR-TR-2023-33. Microsoft. <https://www.microsoft.com/en-us/research/publication/autogen-enabling-next-gen-llm-applications-via-multi-agent-conversation-framework/>
- [51] Zhiheng Xi, Wenxiang Chen, Xin Guo, Wei He, Yiwen Ding, Boyang Hong, Ming Zhang, Junzhe Wang, Senjie Jin, Enyu Zhou, Rui Zheng, Xiaoran Fan, Xiao Wang, Limao Xiong, Yuhao Zhou, Weiran Wang, Changhao Jiang, Yicheng Zou, Xiangyang Liu, Zhangyue Yin, Shihan Dou, Rongxiang Weng, Wensen Cheng, Qi Zhang, Wenjuan Qin, Yongyan Zheng, Xipeng Qiu, Xuanjing Huang, and Tao Gui. 2023. The Rise and Potential of Large Language Model Based Agents: A Survey. arXiv:2309.07864 [cs.AI]
- [52] Ce Zhou, Qian Li, Chen Li, Jun Yu, Yixin Liu, Guangjing Wang, Kai Zhang, Cheng Ji, Qiben Yan, Lifang He, Hao Peng, Jianxin Li, Jia Wu, Ziwei Liu, Pengtao Xie, Caiming Xiong, Jian Pei, Philip S. Yu, and Lichao Sun. 2023. A Comprehensive Survey on Pretrained Foundation Models: A History from BERT to ChatGPT. arXiv:2302.09419 [cs.AI]
- [53] Mingchen Zhuge, Haozhe Liu, Francesco Faccio, Dylan R. Ashley, Róbert Csordás, Anand Gopalakrishnan, Abdullah Hamdi, Hasan Abed Al Kader Hammoud, Vincent Herrmann, Kazuki Irie, Louis Kirsch, Bing Li, Guohao Li, Shuming Liu, Jinjie Mai, Piotr Piękos, Aditya Ramesh, Imanol Schlag, Weimin Shi, Aleksandar Stanić, Wenyi Wang, Yuhui Wang, Mengmeng Xu, Deng-Ping Fan, Bernard Ghanem, and Jürgen Schmidhuber. 2023. Mindstorms in Natural Language-Based Societies of Mind. arXiv:2305.17066 [cs.AI] Workshop on robustness of zero/few-shot learning in foundation models (R0-FoMo) at NeurIPS – Best Paper.