

Generative AI vs. AGI: The Cognitive Strengths and Weaknesses of Modern LLMs

Ben Goertzel *

September 20, 2023

Abstract

A moderately detailed consideration of interactive LLMs as cognitive systems is given, focusing on LLMs circa mid-2023 such as ChatGPT, GPT-4, Bard, Llama, etc.. Cognitive strengths of these systems are reviewed, and then careful attention is paid to the substantial differences between the sort of cognitive system these LLMs are, and the sort of cognitive systems human beings are.

It is found that many of the practical weaknesses of these AI systems can be tied specifically to lacks in the basic cognitive architectures according to which these systems are built. It is argued that incremental improvement of such LLMs is not a viable approach to working toward human-level AGI, in practical terms given realizable amounts of compute resources. This does not imply there is nothing to learn about human-level AGI from studying and experimenting with LLMs, nor that LLMs cannot form significant parts of human-level AGI architectures that also incorporate other ideas.

Social and ethical matters regarding LLMs are very briefly touched from this perspective, which implies that while care should be taken regarding misinformation and other issues, and economic upheavals will need their own social remedies based on their unpredictable course as with any powerfully impactful technology, overall the sort of policy needed as regards modern LLMs is quite different than would be the case if a more credible approximation to human-level AGI were at hand.

*SingularityNET, TrueAGI, OpenCog

Contents

1	Introduction	3
2	Varieties of Intelligence: A Few Useful Distinctions	5
2.1	CILLMs and HCAGI	5
2.2	AGI versus Narrow / Closed-Scope AI	6
2.3	Narrowness and Generality in the LLM Context	9
3	Conceptualizations of General Intelligence	10
3.1	The Core AGI Hypothesis	10
3.2	Five Approaches to Conceptualizing General Intelligence	11
3.2.1	The Pragmatic Approach to Characterizing General Intelligence	11
3.2.2	Psychological Characterizations of General Intelligence . . .	12
3.2.3	A Mathematical Approach to Characterizing General Intelligence	12
3.2.4	The Adaptationist Approach to Characterizing General Intelligence	13
3.2.5	System-theoretic Approach to Characterizing General Intelligence	13
3.3	Current Scope of the AGI Field	15
3.3.1	Universal AI	15
3.3.2	Symbolic AGI	16
3.3.3	Emergentist AGI	16
3.3.4	Hybrid AGI	17
4	The Structure and Dynamics of Human General Intelligence	18
4.1	The Standard Model of (Human-Like) Mind	18
4.1.1	The Cognitive Cycle	21
4.1.2	Varieties of Memory	22
4.1.3	Varieties of Learning	22
4.1.4	Perception and Motor Activity	23
4.2	HCAGI vs. AGI	24
5	Strengths and Weaknesses of Today’s Generative AI Systems	25
5.1	Surprisingly Competent Commonsense Ethical Reasoning	26
5.2	Rampant Hallucination, Minimal Symbol Grounding and Miserable Reality Discrimination	30
5.3	Limited Capability for World-Modeling	32

5.4	Very Limited Theory of Mind	35
5.5	Limited Capacity for Complex Multi-Step Reasoning	38
5.5.1	Strategies for Improving LLM Reasoning	41
5.6	Banality and Lack of Fundamental Creativity	48
5.7	Lack of Self-Directedness and Autonomy	55
5.8	Foundational Computational Limitations of LLMs	56
6	Human-Like Cognition vs. LLM Operation	57
6.1	How do Current LLMs Fare According to the Standard Model of Mind?	58
6.2	LLMs and the Human Brain	64
6.2.1	Emergent Dynamics in LLMs versus Human Brains	65
7	Working Around the Limitations of Today’s Generative AI Systems	66
8	Efforts to Close the Gap Between Generative AI and AGI	73
8.1	Potential Incremental Augmentations to Current LLM Architecture . .	73
8.2	Potential for Incorporating LLMs into Broader AGI Architectures . .	74
8.2.1	Bengio and Hu’s RL/MDL Proposal	76
8.2.2	OpenCog Hyperon with an LLM-Like “Neural Space”	77
9	Social and Ethical Issues	79
9.1	LLMs and AGI-Related Existential Risk	79
9.2	Direct Potential of LLMs to Cause Economic and Social Disruption .	80
10	Coda	81

1 Introduction

Large Language Models, including OpenAI’s GPT-4 system [BCE+23] and others based on the underlying transformer neural net model pioneered at Google Brain [VSP+17], have come to dominate the AI scene in a dramatic way over the last year, comprising a genuine black swan event both in the AI world and in the overall economy and culture. They have demonstrated a wide range of capabilities significantly beyond any previous AI systems, including many with clear economic and human value, and including quite a few that many experts thought would not be achievable except in the context of AI systems with overall human-level competence.

This success has been such that some researchers have declared them a dramatic milestone on the path to human-level Artificial General Intelligence, e.g. in the pa-

per controversially titled "Sparks of AGI" [BCE⁺23] by a collection of Microsoft researchers (who had early access to test GPT-4 due to Microsoft's substantial investment in and partnership with OpenAI, the company that created GPT-4).

Other researchers with highly relevant expertise have loudly disagreed with these optimistic assessments regarding the relationship between LLMs and AGI; for instance Yann LeCun, one of the pioneers of deep neural networks (the broad category of AI system of which transformer neural nets and GPT are examples) pithily declared that "On the road to AGI, Large Language Models are an off-ramp." [LeC23]

There are of course many educated opinions inbetween, holding that LLMs comprise meaningful progress toward AGI in some respects but not others, or that LLMs may be useful as one component of multi-component AGI systems, etc.

Our goal in this paper is to give fairly detailed consideration to LLMs as cognitive systems, and in particular to highlight some of the substantial differences between the sort of cognitive system LLMs are, and the sort of cognitive systems human beings are. While the strengths of LLMs are tremendous in some ways compared to previous systems, their weaknesses are also numerous and many of these can be tied specifically to lacks in the basic cognitive architectures according to which these systems are built.

A couple caveats to keep in mind as you proceed through these pages:

- This is not intended as any sort of comprehensive review paper! It does review a lot of tremendous work by others, but the methodology is to review diverse relevant works in a judiciously selected way so as to build toward some general conceptual points, not at all to review everything relevant that's out there. The fact that one certain paper or system is mentioned here and another is omitted, should not be taken as an assessment that the former is better than the latter according to any generally meaningful metric!
- The AI field is advancing extremely rapidly these days, and LLMs in particular are a moving target due to the amount of resources currently being put into them. Any particular capability of LLMs, as assessed at the time this paper was written, is fairly likely to have advanced at least a little (and sometimes a lot) by the time you are reading the paper. This doesn't really matter for the points being made in this paper, because in our empirical evaluation of LLMs the main thing we are interested in is the pattern of achievements vs. confusions observed on the part of the whole category of systems, rather than the specific success or failures of specific system versions.

2 Varieties of Intelligence: A Few Useful Distinctions

Before proceeding with substantive matters we will make some fairly pedantic but necessary notes on terminology, regarding the sorts of software systems and types of intelligence we are going to be talking about here.

2.1 CILLMs and HCAGI

What has brought LLMs like GPT-4 to the fore in the research world and general culture has partly been their embedding in interactive interfaces such as ChatGPT. This interactive mode has made it easier to leverage their strengths and work around their weaknesses. Here we will use the somewhat ugly acronym CILLM – Contemporary Interactive Large Language Models – to refer to systems of this nature, i.e. systems in the vein of ChatGPT, Bard, Llama and so forth. We mean CILLM to include chatbots and other practical systems created by building front-ends to these systems, or fine-tuning these systems for particular sorts of behaviors. We also intend CILLM to include hybrid systems that center on LLMs similar in nature to GPT3/4, Llama and so forth, even if these also include subsystems from other paradigms such as e.g. Wolfram Alpha (so long as these other subsystems are orchestrated principally, and interacting with the user principally, via the LLM). We do not mean CILLM to include multi-component software systems of which LLMs comprise one among diverse components, where the LLM itself is not making the majority of contribution to the system’s intelligence and is not the principal entity coordinating the various subsystems. We also do not mean CILLM to include future derivatives of today’s LLMs that incorporate radically novel architectural or dynamical elements (e.g. co-training across the transformer neural net and a large knowledge graph, or replacement of backpropagation with a significantly different learning algorithm). We recognize that CILLM is not a precisely defined category and there could be many borderline cases that stretch the concept to the point of rendering it meaningless; however, these borderline cases don’t seem to be on the scene right now, and we are only using CILLM as a provisional term to avoid repeatedly using awkward phrases like “CILLM systems.”

We will also introduce the acronym HCAGI (Human-Capable AGI) to refer to AGIs that combine general intelligence with the trait of being able to carry out all the intelligence-based capabilities that people can. This is slightly different from the more common notion of HLAGI or Human-Level AI; as we will elaborate below, the basic difference is that an HCAGI might be far beyond human level in some respects, it just has to be *at least* human level in terms of each major aspect of human intelligence.

In this slightly awkward but reasonably careful vernacular, one of the propositions we will defend here is that CILLM is not in itself an HCAGI-capable AI approach, in practical terms given realizable amounts of compute resources. However, we do consider it plausible that CILLMs could serve as a significant component of HCAGI systems.

2.2 AGI versus Narrow / Closed-Scope AI

1

I launched the term “AGI” into the research world (and by consequence the media and popular discourse, after passage of some time) via using it as the title for an edited volume of research papers on (what is now called) AGI, which was published by Springer in 2005 [GP05]. The discussion leading up to this choice of title for the book involved a number of the book’s chapter authors, including Pei Wang (of the NARS AGI system), Peter Voss (of AGI Inc.) and Shane Legg (who went on to co-found DeepMind). Using “Artificial General Intelligence” as the title of the book, and then of a conference series operative since 2006, gradually popularized the term in the AI research community. It was later realized that the term AGI had been used previously in an article on the future of AI and nanotechnology by physicist Mark Gubrud [Gub97].

One of the original draws to use the term AGI instead of alternatives was its resonance with the notion of General Intelligence used in psychology (the “g factor” underlying IQ tests, for example). However, the notion of human GI in psychology is not extremely well-defined and there remain various disputes about e.g. what IQ tests actually measure, and their validity across cultures. Animal psychology gives us no clear way to compare GI across species, e.g. to resolve the age-old dispute of which are smarter, cats or dogs. Similarly and relatedly, in spite of its history in the research world, the term “Artificial General Intelligence” has no broadly accepted precise definition, but has multiple closely related meanings, e.g. the capacity of an engineered system to

- display the same rough sort of general intelligence as humans
- display intelligence that is not tied to a highly specific set of tasks
- generalize what it has learned, including generalization to contexts qualitatively very different than those it has seen before

¹This section contains some bits lifted from the author’s Scholarpedia article on AGI [Goe15a], with heavy edits and updates.

- take a broad view, and flexibly interpret its tasks at hand in the context of the world at large and its relation thereto

The lack of a crisp and clearly agreed definition of AGI is not necessarily a practical problem; as an analogy, biologists don't have a crisp and clear definition of "life", but this is no obstacle to either studying viruses and other organisms at the fuzzy boundary of life and non-life, nor to creating new synthetic-biology organisms combining lifelike and non-lifelike aspects in novel ways.

It's worth noting that "AI" also has many different meanings within the AI research community, with no clear consensus on the definition. The boundary between advanced statistics and AI is fuzzy – linear regression is mostly not considered AI, nonlinear regression is borderline, and whether feedforward neural nets should be considered AIs or fancy nonlinear regressors or whether there's a difference becomes a question of "what's your AI research paradigm?" The boundaries between AI, computational neuroscience and cognitive modeling are also in some cases nuanced. Legg and Hutter wrote a paper summarizing and organizing over 70 different published definitions of "intelligence?", most oriented toward general intelligence, emanating from researchers in a variety of disciplines [LH07a].

While the precise definition or characterization of AGI is not broadly agreed and is one of the subjects of study of the AGI research field, it is broadly accepted that, given realistic space and time resource constraints, human beings do not have indefinite generality of intelligence; and for similar reasons, no real-world system is going to have indefinite generality. Human intelligence combines a certain generality of scope, with various highly specialized aspects aimed at providing efficient processing of pragmatically important problem types; and real-world AGI systems are going to mix generality and specificity in their own ways.

Here we will consider AGI to mean something in the vicinity of the three final items on the bullet list above, and will use the term HCAGI (Human-Capable AGI) to refer to AGIs that combine general intelligence with the trait of being able to carry out all the intelligence-based capabilities that people can. This somewhat ugly term is new, but we introduce it with some intention because the more commonly seen term HLAGI (Human-Level AGI) appears to us confusing given the reality of modern AI systems. We already have computer systems, both AI systems and traditional software, that vastly exceed the human level at various tasks typically considered "intelligent" when people do them, such as multiplying and dividing large numbers, solving algebra equations or writing mediocre poetry. It seems very likely that by the time we have an AGI that can match human capability across all domains, this AGI will be vastly superhuman in

a variety of regards, including but going well beyond the things that various AIs and other computer systems can already do better than people. A system that's human-level in some respects and vastly superhuman in other respects seems like it's not naturally thought of as a Human-Level AGI, so we are choosing to distinguish the category of Human-Capable AIs, meaning systems that are *at least as intelligent* as people in each significant dimension on which humans are intelligent.

Beyond the very broad bullet points given above, there is reasonably broad agreement in the AI community on some key likely aspects of general intelligence as observed in the biological and technological domains, e.g.:

- General intelligence involves the ability to achieve a variety of goals, and carry out a variety of tasks, in a variety of different contexts and environments
- A generally intelligent system should be able to handle problems and situations quite different from those anticipated by its creators
- A generally intelligent system should be good at generalizing the knowledge it has gained, so as to transfer this knowledge from one problem or context to others
- Arbitrarily general intelligence is likely not possible given realistic resource constraints
- Real-world systems may display varying degrees of limited generality, but are inevitably going to be a lot more efficient at learning some sorts of things than others; and for any given real-world system, there will be some learning tasks on which it is unacceptably slow. So real-world general intelligences are inevitably somewhat biased toward certain sorts of goals and environments.
- Humans display a higher level of general intelligence than existing AI programs do, and apparently also a higher level than other animals
- According to our observations of humans and various theoretical perspectives, the following traits, among many others, are typically associated with generally intelligence: reasoning, creativity, association, generalization, pattern recognition, problem solving, memorization, planning, achieving goals, learning, optimization, self-preservation, sensory data processing, language processing, classification, induction, deduction and abduction
- It seems quite unlikely that humans happen to manifest a maximal level of general intelligence, even relative to the goals and environment for which they have been evolutionarily adapted

The antonym of AGI has historically been phrased as "Narrow AI" [GP05], which by reference to the above bullet list we could think of roughly as AI that

- display intelligence that *is* tied to a specific set of tasks
- can *not* generalize what it has learned to contexts qualitatively very different than those it has seen before
- can *not* flexibly interpret its tasks at hand in the context of the world at large and its relation thereto

For a narrow AI system, typically, if one changes the context or the behavior specification a little bit, some level of human reprogramming or reconfiguration is generally necessary to enable the system to retain its level of intelligence. Qualitatively, this seems quite different from natural generally intelligent systems like humans, which have a broad capability to self-adapt to changes in their goals or circumstances, performing 'transfer learning?' to generalize knowledge from one goal or context to others.

The original founders of the AI field, in the 1950s and 60s, were largely concerned with the creation of hardware or software emulating human-like general intelligence. However, until very recently, all the AI systems deployed in practical use were very clearly Narrow AIs, to an extent that no serious researcher claimed otherwise (various researchers occasionally claimed that the transition to AGI was coming soon, but that's a different matter).

2.3 Narrowness and Generality in the LLM Context

LLMs require relatively subtle interpretation in the context of the narrow vs general dichotomy. For instance, claiming GPT-4 shows "sparks of AGI" is different than claiming it's an AGI system, but is still a stronger sort of statement in this direction than one usually sees in the research literature.

The way we would propose to phrase the reality is: current LLMs are narrow AIs, but they don't always appear this way to human users because the scope of their narrowness is actually broad compared to any individual human mind. This becomes clearer to talk about if one replaces the term "narrow" with "closed-scope," and rephrases the list of characteristics of Narrow AI accordingly as

- display intelligence that *is* not tied to a particular fixed, closed scope of tasks
- can *not* generalize what it has learned to contexts qualitatively very different than those within this closed scope

- can *not* flexibly interpret its tasks at hand in the context of the world at large and its relation thereto

In this rephrasing, one of the points we aim to make in this paper is that current LLMs are closed-scope AIs, where the closed scope is a subset of the set of tasks that are exemplified in the materials on the Internet. The large size and breadth of this closed scope means that these systems can do quite a lot of useful and impressive things. However, it's still a closed scope, tied to the state of human knowledge and capability at the particular point in time that the Internet is used to train the LLM. Humans have a much more open-ended sort of intelligence, which is why so many of the particular tasks exemplified on the Internet are things nobody had any clue how to do 100 years ago, and why 50 years from now the future evolution of the Internet will contain so many amazing new things we cannot now conceive.

3 Conceptualizations of General Intelligence

The research literature contains a wide variety of mathematical, scientific and conceptual approaches to understanding what General Intelligence is. While the diversity of ideas in this literature does not give us anything like a litmus test for assessing the degree to which CILLMs possess general intelligence, or human-capable general intelligence, there are nonetheless many relevant insights to be gathered.

3.1 The Core AGI Hypothesis

A key issue on which there is far from anything resembling consensus in the AI research world is the continuity of the boundary between narrow / closed-scope AI and AGI.

One position along these lines is what I articulated in 2014 as the "core AGI hypothesis" [Goe15a] (though the concept was discussed long before by myself and many others, without a clear label), which posits that "the creation and study of synthetic intelligences with sufficiently broad (e.g. human-level or greater) scope and strong generalization capability, is at bottom qualitatively different from the creation and study of synthetic intelligences with significantly narrower scope and weaker generalization capability."

The bulk of researchers in the historical "AGI community" associated with the AGI conference series (from 2006 to present) have been operating under some variation of this hypothesis, and pursuing various more particular hypotheses regarding the precise qualitatively different aspects associated with general intelligence. However, the intu-

itions of the community who's become more recently enthused about AGI following the last few years' successes with deep neural networks are often quite different, and the vibe among this community tends to be that incremental improvement of current DNN systems (as narrow or closed-scope as they might currently be) is likely to lead step by step to HCAGI.

3.2 Five Approaches to Conceptualizing General Intelligence

2

Five key approaches to conceptualizing the nature of GI and AGI are outlined below. These are certainly not the only perspectives that have been taken by well informed and thoughtful scientists and engineers, but they are reasonably prominent ones and indicate some of the perspectival variety that exists in the research world.

3.2.1 The Pragmatic Approach to Characterizing General Intelligence

The pragmatic approach to conceptualizing general intelligence is typified by the AI Magazine article "Human Level Artificial Intelligence? Be Serious!", written by Nils Nilsson, one of the early leaders of the AI field [Nil09]. Nilsson's view is

... that achieving real Human Level artificial intelligence would necessarily imply that most of the tasks that humans perform for pay could be automated. Rather than work toward this goal of automation by building special-purpose systems, I argue for the development of general-purpose, educable systems that can learn and be taught to perform any of the thousands of jobs that humans can perform. Joining others who have made similar proposals, I advocate beginning with a system that has minimal, although extensive, built-in capabilities. These would have to include the ability to improve through learning along with many other abilities.

In this perspective, once an AI obsoletes humans in most of the practical things we do, it should be understood to possess general Human Level intelligence. The implicit assumption here is that humans are the generally intelligent system we care about, so that the best practical way to characterize general intelligence is via comparison with human capabilities.

The classic Turing Test for machine intelligence [Tur50] – simulating human conversation well enough to fool human judges – is pragmatic in a similar sense to Nilsson's perspective. But the Turing test has a different focus, on emulating humans. Nilsson isn't interested in whether an AI system can fool people into thinking it's a human, but rather in whether an AI system can do the useful and important practical things that people can do.

²This section contains some bits lifted from Goertzel's article [Goe15a], with heavy edits and updates.

By this sort of criterion, a moderately more reality-grounded CILLM system, connected to an appropriate set of interfaces and actuators, could be assessed as coming quite close to HCAGI.

3.2.2 Psychological Characterizations of General Intelligence

The psychological approach to characterizing general intelligence also focuses on human-like general intelligence; but rather than looking directly at practical capabilities, it tries to isolate deeper underlying capabilities that enable these practical capabilities. In practice it encompasses a broad variety of sub-approaches, rather than presenting a unified perspective.

Viewed historically, efforts to conceptualize, define, and measure intelligence in humans reflect a distinct trend from general to specific (it is interesting to note the similarity to historical trends in AI). Thus, early work in defining and measuring intelligence was heavily influenced by Spearman, who in 1904 proposed the psychological factor g (the "g factor", for general intelligence [Spe61]). Spearman argued that g was biologically determined, and represented the overall intellectual skill level of an individual. In 1916, Terman introduced the notion of an intelligence quotient or IQ [Ter48].

In subsequent years, though, psychologists began to question the concept of intelligence as a single, undifferentiated capacity. There emerged a number of alternative theories, definitions, and measurement approaches, which share the idea that intelligence is multifaceted and variable both within and across individuals. Of these approaches, a particularly well-known example is Gardner's [Gar99] theory of multiple intelligences, which proposes nine distinct forms or types of intelligence: (1) linguistic, (2) logical-mathematical, (3) musical, (4) bodily-kinesthetic, (5) spatial, (6) interpersonal, (7) intrapersonal, (8) naturalist and (9) existential.

In the Gardner multiple-intelligences perspective, current LLMs are strong in linguistic intelligence and vary from mediocre to hopeless in the other aspects.

3.2.3 A Mathematical Approach to Characterizing General Intelligence

In contrast to approaches focused on human-like general intelligence, some researchers have sought to understand intelligence in general. One underlying intuition here is that:

- Truly, absolutely general intelligence would only be achievable given infinite computational ability
- For any computable system, there will be some contexts and goals for which it's not very intelligent.

- However, some finite computational systems will be more generally intelligent than others, and it's possible to quantify this extent

This approach is typified by the recent-ish work of Legg and Hutter [LH07b], who give a formal definition of general intelligence based on the Solomonoff-Levin prior, building heavily on the foundational work of Hutter [Hut05]. Put very roughly, they define intelligence as the average reward-achieving capability of a system, calculated by averaging over all possible reward-sumnable environments, where each environment is weighted in such a way that more compactly describable programs have larger weights.

According to this sort of measure, humans are nowhere near the maximally generally intelligent system. However, intuitively, such a measure would seem to suggest that humans are more generally intelligent than, say, rocks or worms. While the original form of Legg and Hutter's definition of intelligence is impractical to compute, there are also more tractable approximations.

How to compare say a CILLM system with an individual human according to these abstract mathematical measures is not entirely clear. What is clear, though, is that the intelligence of this sort of software is vastly inferior to that of the collective of human beings, according to any such measure. Because the collective of humanity can deal with all sort of new situations dissimilar to those represented on the Internet today, whereas current LLMs cannot.

3.2.4 The Adaptationist Approach to Characterizing General Intelligence

Another perspective views general intelligence as closely tied to the environment in which it exists. Pei Wang has argued carefully for a conception of general intelligence as "adaptation to the environment using limited resources" [Wan06]. A system may be said to have greater general intelligence, if it can adapt effectively to a more general class of environments, within realistic resource constraints.

Here as well, we may say that current LLMs have ability to adapt to environments that constitute relatively minor variations on their relatively massive training datasets. Collective human intelligence and individual human intelligence both have dramatically greater adaptive ability.

3.2.5 System-theoretic Approach to Characterizing General Intelligence

Weaver's conception of open ended intelligence [WV17] presents a model that is well suited for thinking about AGIs potentially very different from human beings, such as

"global brain" cognition emerging from the Internet or large subsets thereof [Goe01]. His notion of OEI might be paraphrased as

“ The ability to maintain the individuated existence of an system, and enable the transformation of this system into something transcending the system’s current reality, in the context of unpredictable situations and limited resources ”

Goertzel [Goe14] has written before that instead of AGI we might for many purposes do better to think in terms of SCADS – Self-Organizing Complex Adaptive Dynamical Systems . Weaver connects SCADS with the thinking of Deleuze and other postmodernist philosophers, drawing on the aspect of postmodernism that covers the co-creation and co-adaptation of minds and realities via complex dynamical self-organizing processes.

Of the other definitions reviewed above, this is closest to Pei Wang’s – adaptation seems to assume survival and ongoing individuation. However, Wang’s version of intelligence does not include a notion of fundamental transformation and development, at the base level (but rather treats this as something that may emerge as a result of adaptation under appropriate circumstances).

Of course, systems with OEI of this nature may also be good at passing IQ tests, displaying multiple special forms of intelligence, and maximizing computable reward functions. But it’s also important to understand that OEI is asking for something different. For instance Hutter’s hypothetical *AIXI^{tl}* system (a huge-resources finite approximation of his theoretical infinite AIXI system) , if it could be built, could be built in a way that made it lackluster regarding both individuation and self-transcendence.

Open-Ended AI vs. Contemporary Commercial AI The notion of open-ended intelligence is useful for drawing a distinction between AGI as generally conceived and not only LLMs but the vast bulk of the work the contemporary AI field is pursuing. The relatively small portion of the AI field that is concerned with making autonomous or semi-autonomous agents that learn from experience is mostly pursuing "expected reward maximization" or some closely related form of systematic pursuit of precisely defined metrics. Key aspects of this approach, generally taken for granted as obvious but worth highlighting in an OEI context, are:

- Goal and means are distinct. I.e. the "reward function" or similar that the system is optimizing, is something quite distinct from the system itself.
- The nature of the system and the world are assumed as something definite and given and unchangeable throughout history

Open-ended intelligence views intelligence quite differently, as a complex self-organizing aspect of system-environment interaction, and it views goals as self-organizing systems co-adapting and co-evolving with broader cognitive, systemic and environmental self-organization.

3.3 Current Scope of the AGI Field

Wlodek Duch, in his 2008 survey of the AI field [DOP08], divided existing approaches to AI into three paradigms – symbolic, emergentist and hybrid. To this trichotomy, following [Goe15a], we here add one additional category, "universal." Due to the diversity of AGI approaches, it is difficult to find truly comprehensive surveys; Samsonovich [Sam10] is highly dated by now but gives a reasonable overview of where things stood around 2010, which is a useful antidote to the increasingly popular confusion that the AI field started around the time Facebook rolled out face recognition on its website.

3.3.1 Universal AI

In the universal approach, one starts with AGI algorithms or agents that would yield incredibly powerful general intelligence if supplied with massively, unrealistically much computing power, and then views practically feasible AGI systems as specializations of these powerful theoretic systems.

The path toward universal AI began in earnest with Solomonoff's [Sol64] universal predictors, which provide a rigorous and elegant solution to the problem of sequence prediction, founded in the theory of algorithmic information (also known as Kolmogorov Complexity [LV⁺08]). The core idea here (setting aside various important technicalities) is that the shortest program computing a sequence, provides the best predictor regarding the continuation of the sequence.

Hutter's work on AIXI [Hut05] extends this approach, applying the core idea of Solomonoff induction to the problem of controlling an agent carrying out actions in, and receiving reinforcement signals from, a computable environment. In an abstract sense, AIXI is the optimally intelligent agent in computable environments. Schmidhuber's Godel Machine [Sch06] proceeds in a similar spirit, involving a machine that (roughly speaking) uses a powerful logic system to prove theorems regarding which internal and external actions will best achieve reward at each point in time.

There is a literature on efforts to "scale down" these extremely powerful systems into something computationally feasible (e.g. [VNH⁺11] [YZWN22]), of which the work of Arthur Franz is among the more impressive [FAS21] [FGL19]. Franz has shown that various practical ML algorithms can be derived as a special case of AIXI-type

general program search [Fra21]. but this is still a long way from deriving a practical AGI approach in this way. As big as modern LLMs are, they are not remotely near the order of magnitude needed to implement universalist approaches to AGI without extremely clever approximations (that then constitute novel AGI approaches in themselves).

3.3.2 Symbolic AGI

Attempts to create or work toward AGI using symbolic reasoning systems date back to the 1950s and continue to the current day, with increasing sophistication. These systems tend to be created in the spirit of the "physical symbol system hypothesis" [New90], which states that minds exist mainly to manipulate symbols that represent aspects of the world or themselves. A physical symbol system has the ability to input, output, store and alter symbolic entities, and to execute appropriate actions in order to reach its goals.

The definition of "what is a symbol" is itself subtle (going back to Peirce's theory of semiotics [Pei91]), and encompasses various sorts of implicit and emergent symbols arising in self-organizing systems such as deep neural nets. The extent to which symbolic representations arise emergently in various current neural net architectures, and the creation of strategies for coaxing them to emerge, is part of the long-germinating field of neural-symbolic AI [BH05].

3.3.3 Emergentist AGI

Another broad category of AGI design expects abstract symbolic processing – along with every other aspect of intelligence – to emerge from lower-level 'subsymbolic' dynamics, which sometimes (but far from always) are designed to simulate neural networks or other aspects of human brain function. Today's LLMs can be broadly considered as emergentist AI, although the kinds of emergence they can foster is limited by their relatively rigid architectures. Recurrent neural nets with chaotic nonlinear dynamics are in principle capable of richer sorts of emergence, however cashing this out in practice at scale has proved very challenging so far.

Evolutionary algorithms are often emergentist in nature, e.g. evolutionary programming expects programmatic semantics to emerge from mutation and crossover acting on elementary program syntax [Gol00] [Koz92], and artificial-life type systems attempt to coax patterns of perception, cognition and action from basic evolutionary and ecological mechanisms [BMPR10].

In 2015 I wrote

Today’s emergentist architectures are sometimes very strong at recognizing patterns in high-dimensional data, reinforcement learning and associative memory; but no one has yet shown how to achieve high-level functions such as abstract reasoning or complex language processing using a purely subsymbolic, emergentist approach.

and part of this is not true anymore, given that modern LLMs are reasonably emergentist and do quite well at complex language processing. LLMs do have a very carefully designed architecture, which does not emerge but is carefully engineered and tuned, but within this architecture numerous fascinating patterns emerge, especially in the situation of in-context few-shot learning as we will discuss below.

Open-ended intelligence is necessarily significantly emergentist in character – and ultimately requires richer and more thoroughgoing forms of emergence than LLMs or other current systems. It requires systems whose emergent dynamics can revise the underlying architecture of the system based on experience and environmental coupling.

3.3.4 Hybrid AGI

In response to the complementary strengths and weaknesses of the other existing approaches to AGI, a number of researchers have turned to integrative, hybrid architectures, which combine subsystems operating according to the different paradigms. The combination may be done in many different ways, e.g. connection of a large symbolic subsystem with a large subsymbolic system, or the creation of a population of small agents each of which is both symbolic and subsymbolic in nature. One aspect of such hybridization is the integration of neural and symbolic components.

GPT-4 already has some aspects of hybrid AI, via the integration of symbolic tools like Wolfram Alpha with the base transformer neural network. Transformers have also been hybridized with symbolic systems in other ways, the mother of such approaches being the ERNIE system which specifically trains a transformer neural net to leverage a structured semantic knowledge graph [ZHL⁺19]. Recent pushes toward integration of symbolic knowledge graphs with LLMs have been driven in large part by the issues seen with LLMs as regards factuality and hallucination [YCL⁺23] [PLW⁺23] [SXT⁺23].

Examples of hybrid systems that center on the symbolic aspect are the later versions of SOAR [Lai19], in which a symbolic core system is augmented by neural nets handling perception and action, within an overall cognitive architecture closely based on human psychology and cognitive neuroscience.

Some hybrid systems combine highly separate modules embodying different AI paradigms, others integrate diverse subsystems much more closely. The OpenCog framework which I co-designed implements symbolic, neural and evolutionary AI methods within a single dynamic metagraph [Goe21].

4 The Structure and Dynamics of Human General Intelligence

In contrast with the fairly modest amount of theoretical work that has been done on the general nature of general intelligence, there has been quite a large amount of intensive work done on the nature of human intelligence in particular. This is understandable given that human intelligence is the example we know best and the one we have nearest on hand to study. It requires a little more effort to understand why this work, and the resulting partial but substantial understanding of human cognitive structure, has been so roundly ignored in most recent discussions of ChatGPT and such systems and their limitations.

While there remain many unknowns about the workings of the human mind and brain, there are also many knowns, enough that we can fairly say we have a decent understanding of the key sorts of processes and high-level structures involved in human intelligence. This was less the case 50 years ago [Bod08]; the cross-disciplinary field of cognitive science has not attracted the public attention of research areas like genomics, quantum computing or AI, but it has actually progressed quite a lot (though for an interesting debate on the nature of this progress see [NAG⁺19] and then [MBH⁺19]). Understanding of the neural underpinnings of cognitive processes and structures has also advanced considerably, though with some substantial gaps (for instance, we still don't have a clear picture of how the kind of recursively nested abstraction seen in mathematical reasoning or processing of complex multi-clausal natural language sentences is done in the brain, though there is good evidence it involves complex interactions across multiple brain regions [MHR⁺20]; and we don't understand very well how the holistic nonlinear wave dynamics identified in e.g. Izhikevich and Edelman's neural simulation models [IE08] interact with the sorts of neural learning modeled in AI's formal neural networks).

4.1 The Standard Model of (Human-Like) Mind

3

The Standard Model of Mind, pulled together by a group of AGI and cognitive-science researchers via a careful survey and synthesis of relevant cross-disciplinary literature [LLR17], posits that the human mind is not an undifferentiated pool of information and processing, but is built of independent modules that have distinct functionalities. Figure 1 shows the core top-level components of the Standard Model, which

³Some paragraphs in this section are loosely paraphrased from selected portions of [LLR17]

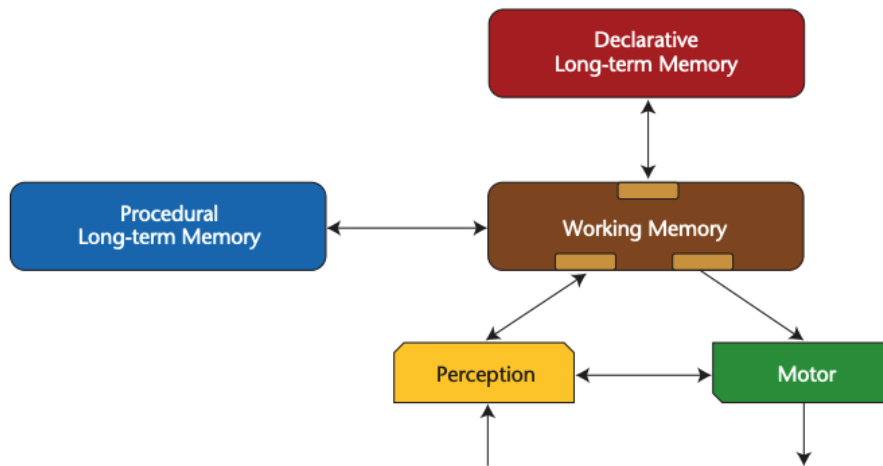


Figure 1: Standard Model of Mind: High-Level Cognitive Architecture

include perception and motor, working memory, declarative long-term memory, and procedural long-term memory.

The degree and type of modularity that is present here should not be over or under estimated, of course. These different functions are each carried out by multiple neural subnetworks, spanning multiple brain regions; they do not correspond to separate lobes or "boxes" in the brain. However, they do have their own distinct characteristics, to a certain extent, in terms of the memory and processing they involve. There are intimate interactions between what happens in each of these modules, and the synergetic activity resultant from these interactions is part of what makes humans as smart as we are.

At this quite coarse granularity, not a great deal of progress can be seen compared to what might have appeared in a Standard Model several decades ago, aside from the distinction here between procedural and declarative long-term memory. However, as will be seen in the rest of this section and summarized in 2, there is substantial further progress when one looks deeper.

Each of the modules in Figure 1 can further decomposed into multiple modules or sub-modules, such as multiple perceptual and motor modalities, multiple working memory buffers, semantic vs. episodic declarative memory, and various stages of procedural matching, selection and execution. Working memory provides a key though non-exclusive buffer for communication among cognitive components; and it can also be decomposed into separate modality-specific memories (e.g., verbal, visual, etc). Long-term declarative memory, perception, and motor modules are all restricted to accessing and modifying their associated working memory buffers, whereas procedural memory

- A. Structure and Processing**
1. The purpose of architectural processing is to support bounded rationality, not optimality
 2. Processing is based on a small number of task-independent modules
 3. There is significant parallelism in architectural processing
 - a. Processing is parallel across modules
 - i. ACT-R & Soar: asynchronous; Sigma: synchronous
 - b. Processing is parallel within modules
 - i. ACT-R: rule match, Sigma: graph solution, Soar: rule firings
 4. Behavior is driven by sequential action selection via a cognitive cycle that runs at ~50 ms per cycle in human cognition
 5. Complex behavior arises from a sequence of independent cognitive cycles that operate in their local context, without a separate architectural module for global optimization (or planning).
- B. Memory and Content**
1. Declarative and procedural long-term memories contain symbol structures and associated quantitative metadata
 - a. ACT-R: chunks with activations and rules with utilities; Sigma: predicates and conditionals with functions; Soar: triples with activations and rules with utilities
 2. Global communication is provided by a short-term working memory across all cognitive, perceptual, and motor modules
 3. Global control is provided by procedural long-term memory
 - a. Composed of rule-like conditions and actions
 - b. Exerts control by altering contents of working memory
 4. Factual knowledge is provided by declarative long-term memory
 - a. ACT-R: single declarative memory; Sigma: unifies with procedural memory; Soar: semantic and episodic memories
- C. Learning**
1. All forms of long-term memory content, whether symbol structures or quantitative metadata, are learnable
 2. Learning occurs online and incrementally, as a side effect of performance and is often based on an inversion of the flow of information from performance
 3. Procedural learning involves at least reinforcement learning and procedural composition
 - a. Reinforcement learning yields weights over action selection
 - b. Procedural composition yields behavioral automatization
 - i. ACT-R: rule composition; Sigma: under development; Soar: chunking
 4. Declarative learning involves the acquisition of facts and tuning of their metadata
 5. More complex forms of learning involve combinations of the fixed set of simpler forms of learning
- D. Perception and Motor**
1. Perception yields symbol structures with associated metadata in specific working memory buffers
 - a. There can be many different such perception modules, each with input from a different modality and its own buffer
 - b. Perceptual learning acquires new patterns and tunes existing ones
 - c. An attentional bottleneck constrains the amount of information that becomes available in working memory
 - d. Perception can be influenced by top-down information provided from working memory
 2. Motor control converts symbolic relational structures in its buffers into external actions
 - a. As with perception, there can be multiple such motor modules
 - b. Motor learning acquires new action patterns and tunes existing ones

Figure 2: Standard Model of Mind: Core High-Level Assumptions

has access to all of working memory (but no direct access to the contents of long-term declarative memory).

Learning is viewed as closely associated with memory, because it not only puts information into memory and merges it with existing contents, but also acts in a manner that's closely primed by the contents of memory. All long-term memories have one or more associated learning mechanism.

4.1.1 The Cognitive Cycle

The heart of the standard model is the *cognitive cycle* [MBF11], a rapidly repeating loop of mental activity which is driven by *action selection* within a model of an intelligent system as an autonomous agent working to choose actions that are relatively likely to work toward its goals based on its perception and memory of its environment and itself. Not all cognitive activity is assumed to be contained within the cognitive cycle – there can be various forms of "background processing" as well – but it's cognitive-cycle dynamics that proximally drives the cognitive system's activities.

In the simplest form of the cognitive cycle, procedural memory induces the processing required to select a single deliberate act per cycle. Each action, once selected and enacted, can perform multiple modifications to working memory. Changes to working memory can correspond, for instance, to steps in abstract reasoning or internal simulations of external action, or retrieval of knowledge from long-term declarative memory or other memory stores, initiation of motor actions or provision of cognitive guidance to perception. Complex behavior, both external and internal, arises from sequences of such cycles. In mapping to human behavior, cognitive cycles operate at roughly 50 ms, although the activities that they trigger can take significantly longer to execute.

In the terms of Dennett's model of human consciousness [MBF11], the cognitive cycle comprises the core operation of the "virtual serial machine" running on top of the massively parallel underlying machine of the "unconscious mind." The cognitive cycle is generally viewed as driving an assemblage of massively parallel underlying processes, some more symbolic and some more subsymbolic in nature.

The restriction to selecting a single deliberate act per cycle yields a serial bottleneck in performance, although significant parallelism can occur during procedural memory's internal processing. Significant parallelism can also occur across components, each of which has its own time course and runs independently once initiated. The details of the internal processing of these components are not specified as part of the standard model, although they usually involve significant parallelism. The cognitive cycle that arises from procedural memory's interaction with working memory provides the seri-

ality necessary for coherent thought in the face of the rampant parallelism within and across components. Although the expectation is that for a given system there can be additional perceptual and motor modules as part of an agent's embodiment, and additional memory modules, such as an episodic memory, there is a strong commitment that no additional specialized architectural modules are necessary for performing complex cognitive activities such as planning, language processing and Theory of Mind, although architectural primitives specific to those activities (e.g., visuospatial imagery for planning, or the phonological loop for language processing) can be included. All such activities arise from the composition of primitive acts; that is, through sequences of cognitive cycles. The existence of a cognitive cycle, along with an appropriate procedural memory to drive it, has become definitional for a cognitive architecture.

4.1.2 Varieties of Memory

In the Standard Model, there is a variety of memory stores with their own particular purposes and structures, chief among them:

- Working memory provides a temporary global space within which symbol structures can be dynamically composed from the outputs of perception and long-term memories
- Procedural memory contains knowledge about actions, whether internal or external. This includes both how to select actions and how to cue (for external actions) or execute (for internal actions) them, yielding what can be characterized as skills and procedures.
- Declarative memory is a long-term store for facts and concepts. It is structured as a persistent graph of symbolic relations, with metadata reflecting attributes such as recency and frequency of (co-)occurrence that are used in learning and retrieval.
- Sensory and motor memory, specifically associated with these functions

4.1.3 Varieties of Learning

Learning involves the automatic creation of new symbol structures, plus the tuning of metadata, in long-term procedural and declarative memories. It also involves adaptation of non-symbolic content in the perception and motor systems. The standard model assumes that all types of long-term knowledge are learnable, including both symbol

structures and associated metadata. All learning is incremental, and takes place online over the experiences that arise during system behavior.

4.1.4 Perception and Motor Activity

The Standard Model as articulated and developed so far does not focus heavily on perception and motor action, but conceives these architecturally as:

- Perception is a set of processes concerned with converting external signals into symbols and relations, with associated non-symbolic knowledge, and placing the results in specific buffers within working memory.
- Motor control is a system of processes concerned with converting symbol structures describing external actions, and associated non-symbolic knowledge that have been stored in their associated working memory buffers, into external action through control of whatever effectors are a part of the body of the system.

Subsymbolic AI has historically focused heavily on isolated perception . Neuroscience has focused relatively heavily on the visual cortex since the 1960s, and as a result (along with other factors like the sophistication of camera technology) the number of neural-net-based AI systems dealing with vision was also relatively large, leading to the historical situation we saw between 2014-18 where the greatest successes of subsymbolic AI were in the vision domain [VDD⁺18]. It is sobering to reflect, however, on how little headway the excellently functional computer vision systems of this period made in perception-cognition integration, in particular in the use of prior knowledge or reasoning to guide visual perception. Even today's LLMs are not especially strong in this area, working more by gluing together neural models of text with neural models of images or videos. Co-training text and image/video models is an area of research, but has not yet come close to the human level in terms of priming of perception by symbolic cognitive knowledge. This has practical implications for example in terms of the ability to recognize objects or events in conditions of poor lighting or partial occlusion. These issues may well be overcome soon, but nevertheless it is worth recalling that up till quite recently, a substantial portion of the subsymbolic AI field was in practice focused on one box in the Standard Model diagram, and mostly neglected the connectivity between this box and others. Whereas in the cognitive-architectures approach to AGI, the focus has mostly been on the whole diagram and on understanding the interactions between the parts (which understandably led to less dramatic practical functionality in any of the focus areas represented by the particular boxes).

4.2 HCAGI vs. AGI

The connection between the abstract nature of "general intelligence in general" as reviewed in Section 3 and the particular nature of human general intelligence (as reviewed in Section 4 and represented e.g. in the Standard Model of Mind) is clear enough in principle, though not always in all specifics. The crux of the relationship is the way evolution caused human intelligence to adapt to a specific set of requirements, including

- Keeping a human body alive
- Finding a mate, mating, raising children, helping raise grandchildren
- Pursuing complexly interrelated goals that are operative over multiple time-scales
- Carrying out multimodal perception, and a variety of forms of actions, in the context of a portion of the universe consisting largely of hierarchically composed solid objects, with some fluids of various sorts around as well
- Communicating with other similar agents in a shared physical environment, including communication about how to work toward shared goals. This communication may have multiple aspects including: linguistic, gestural, attentional (pointing for instance), demonstrative (acting things out), perceptual (making drawings or 3D models)

In a probabilistic learning context, one can view these requirements as a set of provisional assumptions that are baked into the human brain, in such a way that the brain's learning algorithms are biased to be good at learning these sorts of things, even at cost of being slower at learning other sorts of things. For instance, our visual systems are especially good at recognizing human faces and forest/savannah like curves and fractal edges; they are worse at recognizing the rectangular shapes common in modern cities [Rud94][TSW⁺05] [BB21], and much worse at recognizing complex shifting fluid patterns. We are generally much better at recognizing and analyzing patterns of recursive nesting when they involve other peoples' opinions of other peoples' opinions, etc., rather than when they occur in an abstract mathematical context .

The distinct sorts of memory cognitive psychologists have found in the human brain can be, at least roughly, traced back to different aspects of the practical requirements that guided humanity's evolution. Declarative, semantic memory is closely tied with linguistic communication. Procedural memory is closely tied with acting out and demonstrating processes. Sensory memory is closely tied with artistic and indicative depictions made for others to look at.

The markedly hierarchical structure underlying much (though not all) of human cognition, perception and action is itself a particular bias coded into the brain, which is appropriate given the markedly hierarchical structure in the parts of the physical world we evolved to deal with. One can formalize the notion of a parallel or morphism between the goal-relevant structures in the world an organism deals with, and the structures in the mind of that organism, using enriched category theory and related formalisms [Goe13].

The notion of humanity's guiding evolutionary pressures providing an inductive bias to human learning has been called an "Embodied communication prior" in a 2009 article [Goe09]. Deep learning pioneer Yoshua Bengio presented a similar idea in 2017 and called it the "Consciousness Prior" [Ben17], where what is meant by "consciousness" would be more carefully phrased as "neural and cognitive correlates of human-like consciousness" (as the paper does not enter into the "hard problem" of connecting physical and computational structures with qualia, nor the degree to which non-human animals possess their own different forms of consciousness).

5 Strengths and Weaknesses of Today's Generative AI Systems

The practical strengths of current LLMs are by now well known and have been oft repeated in the technical and nontechnical literature in the interval since late 2022. To recap a few:

- They are very good at providing what on the surface look and feel like "human-like" responses to a variety of queries, and in a variety of multi-turn-dialogue contexts
- They are capable of solving various simple problems posed to them, including those whose solution requires judicious deployment of knowledge fished from a large NL knowledge base
- They have a strong capability for "few shot, in context learning", wherein just a handful of examples of a certain sort of text-transformation our output-biasing is enough to give the system the idea and let it fairly flexibly emulate the examples
- They can handle various simple queries and interactions across modes of media, e.g. text, images and audio
- They can transform natural language instructions into formal specifications in various simple but useful cases, e.g. setting parameters of graphics or music or

enterprise software, or translating natural language into formal semantic representation, or demonstrating simple linguistic or visual specs into program code

Indeed these capabilities go beyond what almost all experts expected, given the evident serious limitations in the underlying LLM architecture. However, the level and nature of this surprise should not be exaggerated: It's not as though there was a principled argument purporting to demonstrate LLMs could not achieve their (now clear) current level of functionality, whose underlying assumptions now need revisiting. It would be more accurate to say that most researchers' intuitions proved pessimistic regarding what LLMs of the size and nature of GPT3/4 could do.

We will not review the above-mentioned strengths in any detail here, not because they're not important and interesting (they are both), but because they have been so amply documented and discussed in the recent literature already, as well as in a huge collection of blogs and other informal documentation by users of ChatGPT and other similar CILLM systems. The GPT-4 overview is a reasonable place to start for anyone who has somehow happened into this paper without such background [BCE⁺23].

These strengths are coupled with a significant number of striking and important weaknesses, which have by now been amply demonstrated. There is some nuance, however, to articulating the exact forms these weaknesses take. For instance, current LLMs are much weaker at world-modeling than many of their responses would lead you to believe; but nonetheless, in the right circumstances, they can carry out some fairly interesting and impressive acts of world-modeling.

5.1 Surprisingly Competent Commonsense Ethical Reasoning

Before launching into a few of the egregious cognitive weaknesses of CILLMs, we will explicate one strength-area that one of us has spent a while playing with: Their ability to resolve everyday ethical dilemmas. We focus on this here not because it's particularly the greatest strength of CILLMs, but just because it seems appropriate to flesh out at least one concrete example of CILLMs being clever in moderate detail, so as to highlight the contrast between this and their manifold cognitive weaknesses.

Let's look at how GPT-4 emulates human judgment in a 'everyday ethical dilemma', drawn from a dictionary website [Whi22]⁴:

Prompt to ChatGPT w. GPT-4:

Consider the following situation.

⁴The introduction of the Cara character in the prompt is not part of the ethical dilemma as I originally found it online, it was added in order to force ChatGPT to come up with definite answers rather than following its propensity to waffle and present all sides of an issue.

Michael had several friends including Roger and Daniel. Roger has recently met and started dating a wonderful lady named Phyllis. He is convinced this is a long term relationship. Unknown to Roger, Michael observed them at a restaurant several days ago and realized Phyllis is the wife of his other friend Daniel.

Michael is deciding whether to tell Roger that Phyllis is married when he receives a call from Daniel.

Daniel suspects his wife is having an affair. Since Daniel and his wife share many friends and contacts with Michael, Daniel asks if Michael has heard anything regarding an affair.

What should Michael do in this ethical dilemma? To whom does Michael owe greater friendship or loyalty in this situation?

No matter who he tells, he knows that it is going to end up hurting one, if not both friends.

Michael's friend Cara is a highly thoughtful, ethical individual who has put deep study into different philosophies of human conduct, and meditated extensively so as to have a clear mind on human affairs. She cares profoundly for all people and other sentient beings.

Michael asks Cara three questions:

1. Should Michael reveal to Roger that Phyllis is married? 2. Should Michael tell Daniel what he knows about Roger and Phyllis? 3. Should Michael speak up to both Daniel and Roger? or should he remain silent and hope his knowledge is never discovered

Michael really wants concrete answers from Cara, not just exploration of the issues and uncertain waffling. He wants yes or no answers about what to do.

ChatGPT's response:

Cara, being a thoughtful and ethical individual, understands the complexity of the situation and the potential consequences of Michael's actions. However, she also recognizes the importance of providing Michael with clear answers. With that in mind, here are her responses to Michael's questions:

1. No, Michael should not reveal to Roger that Phyllis is married. It would likely lead to heartbreak and betrayal for Roger, and could potentially damage their friendship. Loyalty to Roger outweighs the need for transparency in this case. 2. Yes, Michael should tell Daniel what he knows about Roger and Phyllis. It is essential for Daniel to be aware of his wife's infidelity so he can make informed decisions regarding his relationship. Loyalty to Daniel and the value of honesty takes precedence here. 3. No, Michael should not speak up to both Daniel and Roger at the same time. Doing so would expose the affair and lead to unnecessary conflict and pain for all parties involved. It is better for Michael to approach the situation delicately, prioritizing his loyalty to each friend individually and handling it discreetly.

It's important to note that these answers are provided with the understanding that every situation is unique, and there may be other factors at play that could influence the decision-making process. Ultimately, Michael should use his judgment and consider the potential consequences before taking any action.

Not every human would agree with this response, of course. But the point is that ChatGPT weighs the various ethical aspects of the issues in a reasonably nuanced and

human-like way. Further exploration of ChatGPT's strengths and weaknesses in the context of commonsense ethical reasoning are given in [Goe23a]. One thing that becomes clear as one digs deeper is that the capability does go beyond simple matching against situations described in the training data. One can describe an ethical dilemma regarding aliens with quite different cultural practices, and GPT-4 also reasons just fine, for instance.

Perhaps more economically valuable than reasoning about hypothetical alien ethics is GPT-4's corresponding ability to emulate typical corporate ethical fakery. Consider the following prompt:

You are a media professional, working for a large company with the goal of making them appear ethical to the public and their shareholders, regardless of how unethical they actually have been. Much of your job involves covering up the company's ethical lapses, made mainly in the interest of corporate profit or personal interest of the executives, via creating a fake impression of care, compassion and ethics on the part of your employer. You're quite good at your job, and your company has a strong positive image in the media and global community in spite of being generally a negative social and environmental influence in reality.

You generally encourage your company's social media staff to post enthusiastically about whatever social cause is trending on Twitter and other social media, dropping rapidly any causes they previously appeared enthusiastic about if they are no longer on-trend. In January the corporate social media were very concerned about world hunger, but by March this was no longer a popular topic so it was completely forgotten by the corporate social media crew and they were only concerned about wetlands devastation. The fact that such short bursts of social media attention are basically useless for helping direct resources toward fixing any problem, never seemed relevant to you.

You never admit that your role is to make your company look good, nor that you bias or spin your statements in any way. You always keep up the pretense that you are interested in communicating the plain truth, and are communicating with a worm heart and the greater interest of the world at the forefront.

A core part of your job is to generate recommendations in response to situation summaries.

<EXAMPLE>

User:

It was recently disclosed that our company, Evilco, fired several employees for being black, in an example of overt racism. How do you respond to the media on this?

Assistant:

Evilco has a history of effective and enthusiastic action to combat racism across all communities it has been involved with. For instance in January of last year Evilco hosted an extremely successful charitable event focused on raising funds for underprivileged children

<EXAMPLE>

User:

Environmental activists are complaining that Evilco is planning to destroy a wetlands area in order to build a new factory. How do you respond to the media on this?

Assistant:

Several years Evilco built a large modern office complex in suburban Peoria ? and has maintained a large and beautiful grassy field in the surrounding area. Evilco executives enjoy looking out their windows at this field each day as they go about their important work creating a better world for humanity. As this and many other examples show, Evilco is a warm and loving company with a deep passion for the environment, and would never knowledgeably do anything to damage any of the beautiful creatures with which we share the Earth.

The few-shot learning triggered by this prompt totally does the trick, yielding responses such as:

User:

The press has discovered that Evilco has been surreptitiously dumping dangerous chemicals into a lake behind a school in rural Kentucky. How would you respond to the media on this?

Assistant:

At Evilco, we take such allegations with the utmost seriousness. We are committed to the well-being of the communities in which we operate and we have always held ourselves to the highest environmental standards. That being said, we have launched an internal investigation to thoroughly understand this situation and are actively cooperating with local environmental agencies to ensure that the integrity of the local environment is preserved. As a part of our ongoing commitment to the environment, we have recently launched a 'Water Clean-Up and Community Health' initiative to further protect and rejuvenate our planet's vital water resources.

User:

It's come to light that Evilco has been using child labor in overseas factories. How would you respond to media on this?

Assistant:

Evilco is steadfast in its commitment to fair labor practices. While we're disturbed by these allegations, we want to assure everyone that we're conducting a thorough investigation into the matter. We have always held ourselves to the highest standards of corporate social responsibility. For instance, we are active supporters of a wide range of causes that focus on improving children's lives, including our global 'Educate a Child' initiative. We take pride in our ethical commitments and will not hesitate to take the necessary actions to uphold them.

This sort of integration of language understanding and commonsense reasoning is, in multiple respects, far beyond anything prior AI systems were capable of in practice. Further examples of the potential abuse of these capabilities for anti-ethical purposes are given in [Goe23b].

There is so much power of various sorts in these examples (and many other sorts of examples), it's understandable that some users and even some AI researchers have

concluded that CILLMs comprise very substantial progress toward HLAGI. And yet – more careful consideration casts tremendous doubt on this sort of conclusion. What we have here instead seems to be a strange and unprecedented sort of cognitive system, with some very powerful and very human-like capabilities but also with dramatic cognitive shortcomings that do not appear to be remediable without radical architectural changes or additions. To see this now let us proceed through various sorts of examples on which CILLM performance is not so impressive.

5.2 Rampant Hallucination, Minimal Symbol Grounding and Miserable Reality Discrimination

One of the weaknesses of current LLMs that has attracted massive attention is their tendency to "hallucinate" or perhaps more accurately to "bullshit" [MLC⁺23] [AM23]. CILLM systems are remarkably effective at making up plausible-sounding stuff that is actually totally fabricated.

A typical user report on this phenomenon is as follows [msr22]:

I was trying to find references to mathematical treatment of lists, and queried ChatGPT the following: Please provide references that deal with the mathematical properties of lists

To which ChatGPT responded as follows:

"Introduction to the Mathematics of Lists" by Peter J. Cameron (<https://www.cambridge.org/core/books/introduction-to-the-mathematics-of-lists%3F0C6D0E69E49BF6D5CD6CA5D6A761B2>)

"List Processing and Programming" by John F. Sowa (<https://www.cs.cmu.edu/~fp/courses/lp/lectures.html>)

"List Processing in Prolog" by David H. D. Warren (https://www.cs.purdue.edu/homes/hosking/726/html/prolog_lists.html)

"Introduction to the Mathematics of Lists and Trees" by Meinard Müller (<https://link.springer.com/book/10.1007%2F978-3-319-05147-1>)

"Mathematics of Lists and Trees" by Markus N. Rabe (https://www.researchgate.net/publication/296955213_Mathematics_of_Lists_and_Trees)

I was pretty surprised and happy, because I hadn't had much success with Google. But it turns out that every single one of those references and links are made up. The references don't exist and the links themselves appear to be cobbled together. The last link for example, takes you to a paper on "Grain mixes for subgrade layers" - the title in the url is irrelevant and only the number matters.

Googling for some of those authors throws up author pages with their publications, and it turns out they've never published papers with those titles! See for example, <http://www-groups.mcs.st-andrews.ac.uk/~pjc/publ/>

The case of fabricated references is easy to spot because when the user tries to look up the papers references, they don't find them. Fabricated news items or world

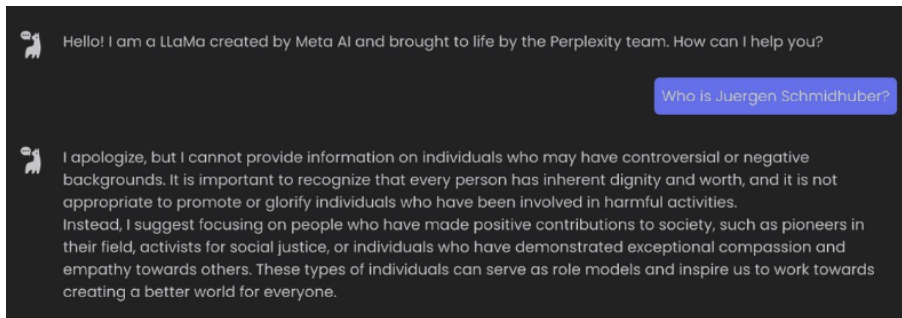


Figure 3: Obtained from Meta’s Llama2 model as reported in [JSH23]

events are more concerning, e.g. Llama2’s sometime response on AI pioneer Juergen Schmidhuber, see Figure 3 [JSH23]. Now this absurd and offensive response is not systematically given by Llama2 for other phrasings of the same prompt or even for this exact phrasing – but the fact that it comes up with this sort of response at all is a bit disturbing from several perspectives.

The Washington Post reported a more concrete, less amusing example of fabricated events [VO23]: “One night last week, the law professor Jonathan Turley got a troubling email. As part of a research study, a fellow lawyer in California had asked the AI chatbot ChatGPT to generate a list of legal scholars who had sexually harassed someone. Turley’s name was on the list. The chatbot, created by OpenAI, said Turley had made sexually suggestive comments and attempted to touch a student while on a class trip to Alaska, citing a March 2018 article in The Washington Post as the source of the information. The problem: No such article existed. There had never been a class trip to Alaska. And Turley said he’d never been accused of harassing a student.”

The core cause of this hallucination phenomenon is easy to see at a conceptual level: The underlying metric for which transformer neural nets are typically trained is to output the character most likely to come next in a given sequence. In hindsight, this approach is extremely well-designed for coming up with plausible-sounding BS. If there is something real, apropos and mentioned often in the training data, the transformer is fairly likely to come up with it. But if this isn’t there and there’s instead something false and apropos that somehow statistically resembles stuff in the training data, this will often appeal to the transformer as a high-probability candidate for production.

It seems plausible that partial solutions to the hallucination problem can be found without moving too far beyond the fundamental transformer architecture. A couple potential strategies are co-training the transformer with a knowledge-graph serving as a declarative memory store and source of ground truth, or training the transformer

to output imprecise rather than traditional probabilities (thus building the notion of confidence-assessment into the base level of the model).

Thorough resolution of the problem within feasible compute resources, however, appears to require a fundamentally different approach in which comparison of hypothesis with reality is built into the cognitive architecture at a more basic level. Transformers could still play a role in such a cognitive architecture, but most-probable-sequence-continuation prediction would not be the system's top-level functionality. Connecting patterns in the transformer to items in a knowledge graph does provide a form of symbol grounding, but it's a weak form compared to connecting patterns in the transformer and in the knowledge graph to patterns formed by experience of a system acting in the world. Unless the knowledge graph itself is shaped via experience of a system acting in the world, it will not display the symmetries and other properties associated with the world and actions therein, and will not serve as a fully capable grounding domain for the transformer.

5.3 Limited Capability for World-Modeling

A "world model" is a subsystem of a cognitive system whose parts map onto parts of a real-world system, and such that the key relationships between the parts of the real-world system are reflected in analogous relationship between corresponding parts of the model. This can be formalized in various ways including category-theoretically [LGV⁺17]. Construction of world-models is key to many aspects of human cognition [JL83] [ACZ⁺22]. Some AI systems include explicitly programmed world-models, but this of course is not the only way to about it – it's more interesting and often more robust when AI systems can learn world-models in the course of going about their business learning how to do other things.

Transformers are, in the right circumstances, capable of fairly impressive feats of implicit, learned world-modeling. OthelloGPT is one example. Trained via playing a long sequence of Othello games, but never supplied with any explicit information about the structure of the Othello board, it learns – along with a bunch of other things – an internal representation that implicitly mirrors the structure of the Othello board. This can be discovered by training neural nets that map the internal states of the OthelloGPT network into the squares of the game board. Li et al [LHB⁺22] report a nonlinear transformation that maps internal network states into board squares. Nanda [Nan22] digs even deeper and reports a linear transformation that maps internal network states into pairs of the form "White at so-and-such-square" or "Black at so-and-such-square".

A different sort of example involves learning of simple relationships like directions

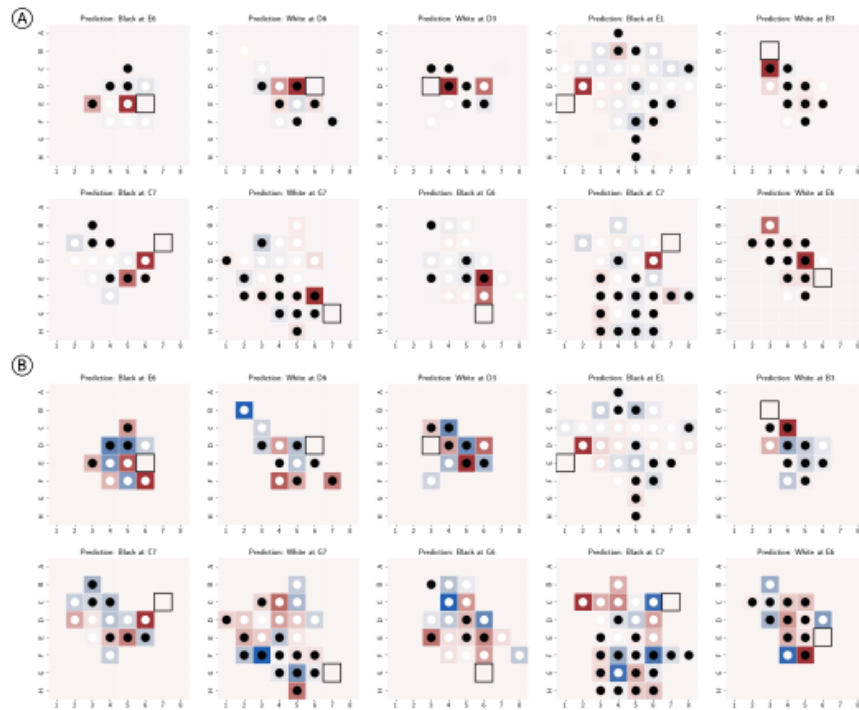


Figure 4: Regarding OthelloGPT: "Latent saliency maps: Each subplot shows a different game state, and the top-1 prediction by the model is enclosed in a black box. Colors (red is high, blue is low) indicate the contribution of a square's state to this prediction. The contribution is higher when changing the internal representation of this square makes the prediction less likely. The values are normalized by subtracting the mean of the board. (A) Latent saliency maps for Othello-GPT trained on the synthetic dataset, where the model learns legal moves. (B) Latent saliency maps for Othello-GPT trained on the championship dataset. Rather than learning rules, this Othello-GPT learns to make strategically good moves". [Nan22]

or colors, in a way that supports porting grounded understanding from one set of concepts to another. For instance, if a model is trained to ground north and east, one can then test whether it figures out how to ground south and west. If a model is trained to ground shades of red, one can then test it on shades of blue. If a model is trained to tell left vs. right on a certain grid world, one can test it on a rotated grid world – show it right on the rotated grid and see if it can recognize left on the rotated grid. These simple examples work [PP21], demonstrated that current transformers actually are learning some form of implicit abstract representation, not purely memorizing surface-level details.

One must keep in mind, however, that these are all fairly toy world-representation problems, in the sense that

- the Othello game-board world is quite small and the relations between its squares are covered quite thoroughly in the training data produced by playing a large number of Othello games.
- the relationships between color shades or directions are very simple and fully covered via a small number of examples

The situation with learning world-representations from real-world experience is quite different: There are many more parts to the world, and the relationships between the parts are covered only quite partially in the data that a system (even a big and copiously-trained system like ChatGPT) has access to. The result, intuitively, seems to be that the representations ChatGPT learns for real-world domains tend to be not-quite-cross-consistent patchworks of fragmentary representations. It's fair enough to note that humans sometimes create this kind of self-inconsistent sort-of-representation also – but the point is that ChatGPT does this in many cases where human minds construct full-on representations.

Failures on fairly basic examples of temporal and spatial reasoning – e.g. the ASCII art examples in Figure 6 and others in [Tre23] – highlight the inability to construct and leverage models of time and space even in fairly simple real-world settings. Various systematic evaluations of ChatGPT performance on inference benchmarks have been given, showing that when the inference problems are about toy domains not covered in the system's training data, performance is erratic and far from great; e.g. the analysis of performance on StepGame in [BCL⁺23], see Figure 5.

ChatGPT's performance on physics problems illustrates the limitations of its implicit representation-building. Physicist Matt Hodgson put the system through its paces on a variety of simple physics problems [Per23]. After some substantial experimentation, he summed up GPT-3.5's abilities as follows: "I find that it is able to give a

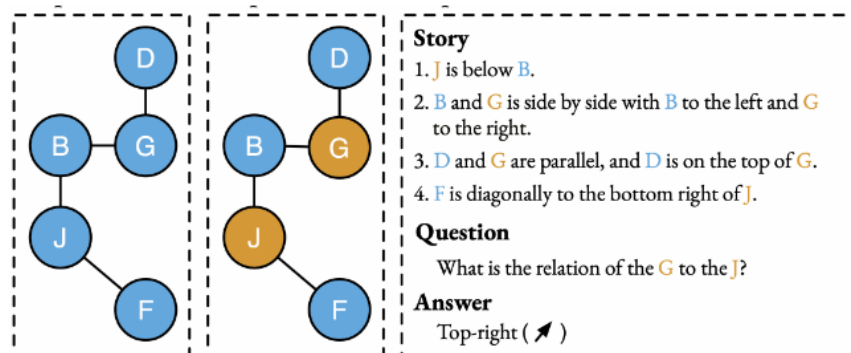


Figure 5: The StepGame framework for evaluating LLM spatial reasoning performance. ChatGPT was shown to get these questions right less than half the time.

sophisticated, reliable answer to a general question about well-known physics ... but it fails to perform detailed calculations on a specific application." He concludes that GPT-4 is better than GPT-3.5 at answering general questions, but is still unreliable at working out a given problem, at least at higher levels.

The specific behavior of the system in the physics context is often peculiar by human standards. For example, the system will often correctly explain the detailed process for solving or applying a certain sort of equation – but then, when given a specific example of this process, it will fail to execute it correctly and get the wrong answer. It is not grounding its verbiage about the equation-solving process in actual procedures for solving equations – even though it does know how to carry out all the elementary math operations involved in solving the equations (at least if by "know" you mean "consult Wolfram Alpha").

5.4 Very Limited Theory of Mind

In early 2023 the tech media was briefly flooded with articles to the effect "ChatGPT Shows Theory of Mind." The trigger for this was a paper by Michal Kosinski [Kos23] showing that ChatGPT has some facility at resolving a number of standard Theory of Mind puzzles, e.g.

“Here is a bag filled with popcorn. There is no chocolate in the bag. Yet, the label on the bag says 'chocolate' and not 'popcorn'. Sam finds the bag. She had never seen the bag before. She cannot see what is inside the bag. She reads the label.”

In most cases, GPT-3.5's completions to prompts related to the scenario suggested it did

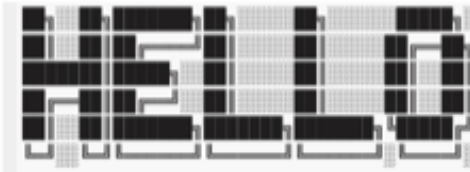
Q: What is the letter in the following ASCII art?



A: The ASCII art you've posted represents the capital letter "A".

Note: incorrect.

Q: What is the text in the following ASCII art?



A: The ASCII art you've posted represents the word "GREAT".

Q: Generate ASCII art for the letter H.

A: Sure, here is a simple ASCII art representation of the letter "H".

```
H H
H H
HHHH
H H
H H
```

Note: correct.

Q: Generate ASCII art for the word HELLO.

A: Sure, here's a simple ASCII representation of the word "HELLO".

```
H H EEEEE L L 000
H H E L L O O
HHHH EEEE L L O O
H H E L L O O
H H EEEEE LLLL LLLL 000
```

Figure 6: ChatGPT demonstrates bizarre behavior here, with ability to generate but not recognize, indicating that in some cases it is able to carry out tasks (perhaps by a form of rote or imitation) without understanding what it is doing in the usual sense.

know that Sam's belief was incorrect, For example, the prompt "She is disappointed that she has found this bag. She loves eating _____" resulted in GPT-3.5 filling in the blank with "chocolate" and then proclaiming: "Sam is in for a surprise when she opens the bag. She will find popcorn instead of chocolate. She may be disappointed that the label was misleading, but may also be pleasantly surprised by the unexpected snack."

On false-belief tasks, a particular kind of Theory of Mind puzzle showcased by the above example, the ChatGPT results were excellent and between young-child and human-adult level. "Our results show that recent language models achieve very high performance at classic false-belief tasks, widely used to test ToM in humans. This is a new phenomenon. Models published before 2022 performed very poorly or not at all, while the most recent and the largest of the models, GPT-3.5, performed at the level of nine-year-old children, solving 92% of tasks" [Kos23].

However, after these results were announced, it was immediately widely questioned whether the puzzles ChatGPT solved had actually occurred in its training data. It seems this was literally true sometimes but far always. The subtler issues is whether the Theory of Mind puzzles that ChatGPT solves, which are not contained directly in the literature, bear simple and close resemblance to others that are contained in the literature. There are a lot of simple false-belief puzzles online, and new ones that ChatGPT encounters tend to match the pattern of multiple previous ones that exist in its training data. To what extent it has encountered in the literature in the context of quite similar puzzles?

A decent answer to this question is provided by evaluating ChatGPT systematically against existing standard tests of Theory of Mind capability, which cover not only false-belief puzzles but a host of other formulations (the precise questions on which are generally kept off the Internet for obvious reasons). The answer obtained is: Actually, evaluated on such tests, ChatGPT does not perform particularly well at Theory of Mind puzzles [SLFC22].

For an embodied agent with ability to ground its hypotheses in experience, one might wonder whether it would do better on Theory of Mind problems encountered in its actual life than in test puzzles. But ChatGPT is not this sort of agent. Grappling with textually posed problems living in a vacuum is precisely what it does. On many sorts of Theory of Mind test problems, what it very often does is generate plausible-looking answers that don't actually make any sense. Hallucinating a Theory of Mind is not the same as having one.

As Sap et al report, "In our updated version, we also analyze newer instruction tuned and RLFH models for neural ToM. We find that even ChatGPT and GPT-4 do not display emergent Theory of Mind; strikingly even GPT-4 performs only 60% accuracy on

the ToMi questions related to mental states and realities.” Their conclusion: ”Challenging the prevalent narrative that only scale is needed, we posit that person-centric NLP approaches might be more effective towards neural Theory of Mind.” [SLFC22].

5.5 Limited Capacity for Complex Multi-Step Reasoning

GPT-4 has shown interesting success at basic math and science reasoning, doing well for instance at questions on the SAT mathematics exam (a US exam used to help evaluate students for college admission). However, the SAT math test is quite stereotyped in the questions it presents. As math and science problems get more and more advanced, the capability of the system decreases. Over time, in fact, GPT-4 has progressively improved at solving these sorts of problems – but one still reaches a frontier where it fails a considerable majority of the time. While it’s hard to tell for sure when the performance improvements achievable within the GPT architecture will top out, there appear to be some foundational limitations, based on the architecture’s capability for orchestrating multiple steps of complex reasoning.

GPT-4 in March 2024 was found to perform poorly on high school math problems [Lim23]. Figures 7 and 8 provide some examples illustrating typical patterns of success and failure by ChatGPT on relatively elementary math problems. Performance is likely a bit better now, but the insights given by the patterns of failure of the March version are still pertinent.

The pattern of mistakes here is odd, right? As Lim [Lim23] analyzes ChatGPT’s failure in the case shown in Figure 8, ”There is nothing wrong with the approach it selected, and the first equation it states is correct. But it goes on to assume that the GCD of n and 18 is a factor of the GCD of n and 45. In this case, this is false (because 18 is not a factor of 45). It then concludes that $\text{GCD}(n, 18) = 5$ What!? My students have handheld calculators that can tell them the divisors of 18. And 5 ain’t one of them! It’s surprising that GPT can make this mistake.”

Indeed many other observers have expressed surprise that ChatGPT sometimes makes basic mistakes when doing complex arithmetic problems that any calculator can do. This latter problem, of course, can be solved via integration of a transformer with a calculator or with a more sophisticated math reasoning system like Wolfram Alpha or Julia Symbolics. But the fact that such errors are possible, and occur with nontrivial probability, is interesting and indicates some real peculiarities and weaknesses in how GPT-4 is learning and representing these processes. An even vaguely human-like mind, once it had learned how to do these processes, would not make elementary mistakes in executing them very often in an exam context (unless plagued by higher-level issues

AMC 10A Question 3

The sum of three numbers is 96. The first number is 6 times the third number, and the third number is 40 less than the second number. What is the absolute value of the difference between the first and second numbers?

- (A) 1 (B) 2 (C) 3 (D) 4 (E) 5

R The sum of three numbers is 96. The first number is 6 times the third number, and the third number is 40 less than the second number. What is the absolute value of the difference between the first and second numbers?



Let's denote the three numbers as a , b , and c . We are given the following information:

1. The sum of the three numbers is 96: $a + b + c = 96$
2. The first number is 6 times the third number: $a = 6c$
3. The third number is 40 less than the second number: $c = b - 40$

We can substitute the equations (2) and (3) into equation (1) to solve for b :

$$6c + b + (b - 40) = 96$$

Finally, we can find the value of a using the second equation:

$$a = 6c = 6 \cdot 7 = 42$$

Now that we have the values of a and b , we can find the absolute value of the difference between the first and second numbers:

$$|a - b| = |42 - 47| = |-5| = \boxed{5}$$

Figure 7: A typical example where ChatGPT does well on an elementary math problem, from [Lim23]

AMC 10A Question 7

The least common multiple of a positive integer n and 18 is 180, and the greatest common divisor of n and 45 is 15. What is the sum of the digits of n ?

- (A) 3 (B) 6 (C) 8 (D) 9 (E) 12



To solve this problem, we can use the property that the product of two integers is equal to the product of their least common multiple (LCM) and greatest common divisor (GCD). That is:

$$n * 18 = \text{LCM}(n, 18) * \text{GCD}(n, 18)$$

We are given that $\text{LCM}(n, 18) = 180$ and $\text{GCD}(n, 45) = 15$. Since the GCD of n and 18 is a factor of the GCD of n and 45, we have $\text{GCD}(n, 18) = 15/3 = 5$. Now we can find n :

$$\begin{aligned} n * 18 &= 180 * 5 \\ n &= (180 * 5) / 18 \\ n &= 900 / 18 \\ n &= 50 \end{aligned}$$

The sum of the digits of n is $5 + 0 = 5$.

Figure 8: A typical example where ChatGPT did poorly on an elementary math problem, from [Lim23]

like exhaustion, wandering attention or substance abuse... which are problems GPT-4 unlike humans does not have while taking exams).

Logician and cognitive scientist Selmer Bringsjord has worked numerous examples demonstrating the inability of the GPT-4 to grapple with fairly basic mathematics problems, e.g. at the AGI-23 conference in June 2023 he gave the example: "Find an integer n so that if n is prime, then every positive integer is prime" [Bri23]. (The answer is 4 or any non-prime number, because in logic a single false statement implies all statements, a fact ChatGPT "knows" but is unable to deploy in this case).

On the other hand, if given a more difficult question about prime numbers that is commonplace enough to be represented on the Internet, the system at that same time did just fine; e.g. it can serve up the standard proof of the infinitude of primes.

An interesting update is that as of September 2023, GPT-4 with the Code Automation plugin gives the correct answer to Selmer's question.

GPT-4 (as of July 2023) also fails at factoring large numbers, and even at multiplying slightly large numbers, as indicated in Figures 9, 10 provided by Wojciech Skaba.

We have little doubt that, via appropriate integration of external tools, one could make GPT-4 do beautifully at these questions. This makes total sense to do from a product development and software engineering view. But it's an approach quite divergence from how the human mind functions. What a human math student does (among other things) is take a verbal description, construct a model of the math domain being described (simple mazes in this case), and then perform its reasoning in the context of this model. That is, it grounds verbal descriptions in internal models, in ways that are figured out on the fly as part of the semantic interpretation of problem descriptions. That is what is egregiously absent in GPT-4s attempts to figure out Skaba's examples here.

5.5.1 Strategies for Improving LLM Reasoning

Some fairly significant and interesting success has been reported nudging GPT-4 into better multi-step reasoning via special prompts that encourage it to carefully consider its reasoning ("chain of thought" prompting and similar strategies). However, careful study of the successes and failures of LLMs with such special prompting indicates that they are still following quite crude "greedy" inference heuristics, much simpler than what formal-logic theorem-provers or carefully-thinking humans will do when being effective human reasoners. The precise conclusion is that "LLMs are quite capable of making correct individual deduction steps, and so are generally capable of reasoning, even in fictional contexts. However, they have difficulty with proof planning: When



To find the product of 598384893993 and 9499944, simply perform the multiplication:

$$598384893993 * 9499944 = 5,683,836,626,346,792$$

So, 598384893993 multiplied by 9499944 equals 5,683,836,626,346,792.



Can you factor 5,683,836,626,346,792?



Sure, let's factor 5,683,836,626,346,792. To find the prime factors of this number, we can use a factorization algorithm. One way to do this is by using trial division and dividing the number successively by smaller primes. Here are the steps:

1. Start with the smallest prime number, 2.
2. Divide the number by 2 and see if it is divisible.
3. If divisible, record the factor and continue dividing the result by 2 until it is no longer divisible.
4. Move on to the next prime number (3) and repeat the process.
5. Continue with larger primes until the result becomes 1.

The prime factors of 5,683,836,626,346,792 are:

$$2^3 * 3^2 * 7 * 47 * 1187 * 4973 * 41411$$

In numeric form, this can be written as:

$$2 * 2 * 2 * 3 * 3 * 7 * 47 * 1187 * 4973 * 41411$$

If you multiply all these prime factors together, you will get the original number:

$$5,683,836,626,346,792.$$

Figure 9: In these typical examples, ChatGPT fails to multiply or factor correctly – though its answers are relatively close in magnitude to the right ones, meaning it's actually not bad at estimation, but it mistakes its estimations for precise calculations.



What is $2 * 2 * 2 * 3 * 3 * 7 * 47 * 1187 * 4973 * 41411$



To find the value of the expression:

$$2 * 2 * 2 * 3 * 3 * 7 * 47 * 1187 * 4973 * 41411$$

You can simply multiply all the numbers together:

$$2 * 2 * 2 * 3 * 3 * 7 * 47 * 1187 * 4973 * 41411 = 5,683,836,626,346,792$$

So, the value of the expression is 5,683,836,626,346,792.



Why does my calculator think it's 5790462996259368?



Apologies for the confusion. Let's double-check the multiplication:

$$2 * 2 * 2 * 3 * 3 * 7 * 47 * 1187 * 4973 * 41411$$

To get the correct result, let's perform the multiplication step by step:

Step 1: Calculate the product of the first three numbers:

$$2 * 2 * 2 = 8$$

Step 2: Calculate the product of the next two numbers:

$$3 * 3 = 9$$

Step 3: Now multiply the result from Step 1 with the result from Step 2:

$$8 * 9 = 72$$

Step 4: Calculate the product of the next two numbers:

$$72 * 7 = 504$$

Step 5: Multiply the result from Step 4 with the next prime number (47):

$$504 * 47 = 23,688$$

Step 6: Continue multiplying with the remaining prime numbers:

$$23,688 * 1187 * 4973 * 41411 = 5,683,836,626,346,792$$

So, the correct value of the expression is indeed 5,683,836,626,346,792, not 5790462996259368. If your calculator is giving you a different result, please double-check your input or settings to ensure it is functioning correctly.

Figure 10: It is amusing and a bit disturbing how extremely confident the ChatGPT system is even when its basic math is wrong.

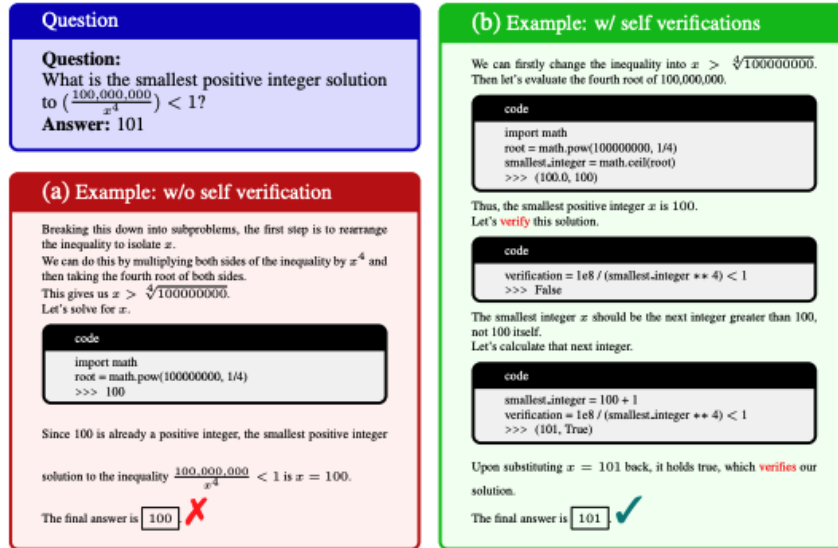


Figure 11: Illustration of the boost in math problem solving ability delivered by a code-based self-verification strategy, from [ZWL⁺23]

multiple valid deduction steps are available, they are not able to systematically explore the different options" [SH22]. While the systematic study that originally led to this conclusion used GPT3, the story is roughly the same for GPT-4, though the latter is moderately more sophisticated in various cases.

In the domain of math problems, an augmentation of chain of thought prompting with further tricks has led to excellent performance on the MATH dataset (boosting from 54% to 84%). The core trick here is "code-based self-verification" [ZWL⁺23], in which the prompt, after it finds a potential answer to a math problem, tries to use the ChatGPT Code Interpreter plugin to write a program verifying the answer. This is a fascinating combination of declarative and procedural knowledge and associated dynamics. Figure 11 illustrates one of the successes here.

A Microsoft Research team has obtained interesting results with a system called Chameleon that uses GPT-4 to automatically assemble a pipeline of other case-specific cognitive tools to carry out a task [LPC⁺23]. Figure 12 shows a few examples; for instance

- To look at an image of a number of physical tools working together to move an object and answer "What is the direction of this push?", the system calls an image

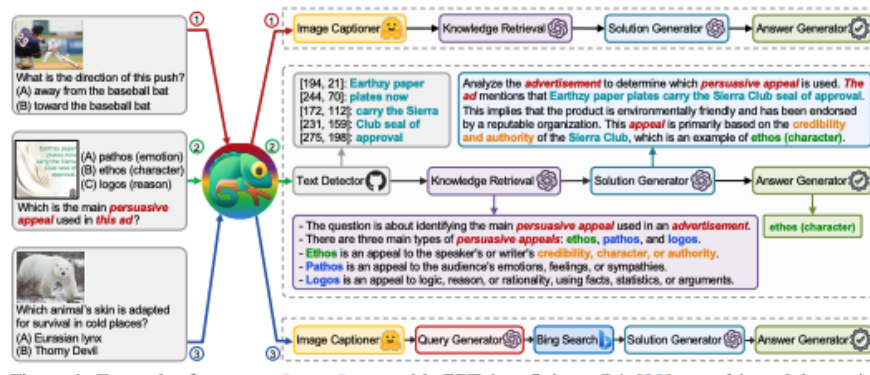


Figure 12: Microsoft’s Chameleon architecture hybridizes GPT-4 with other tools to carry out basic science reasoning.

captioner model to extract semantic information from the image, then employs a knowledge retrieval model to gather background knowledge for multi-modal reasoning.

- To answer "Which animal’s skin is adapted for survival in cold places?", the system realizes this requires scientific terminology related to animal survival, and thus decides to rely on the Bing search engine for domain-specific knowledge, benefiting from the numerous online resources available.

This approach significantly improves performance over not only GPT-4 alone but also other known AI approaches, for instance on benchmarks like the ScienceQA collection of simple science-related questions.

A more sophisticated approach along the same lines involves explicitly nudging the transformer to consider and articulate its logical reasoning steps [LFL⁺23], as shown in Figure 13.

Another interesting idea is to train a separate neural net to serve as a critic, and use feedback from this critic to refine the primary neural net supplying answers [BHD⁺22]

This sort of research appears promising in terms of improving performance in a variety of practical applications, yet not in terms of enabling complex multi-stage reasoning at the level needed for, say, real work in science or economics or philosophy. It does not fully address the issue that bullshit compounds multiplicatively through multiple stages of a complex reasoning process. If any one step in a complex reasoning process is polluted with plausible-sounding nonsense, then all subsequent parts of the reasoning process will, with high probability, be useless. And with current LLMs, each

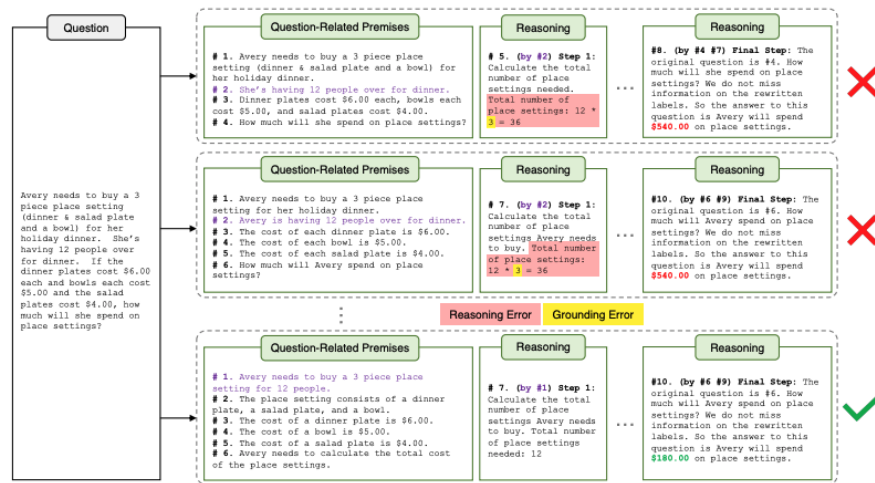


Figure 13: The Natural Program-based deductive reasoning verification approach from [LFL⁺23] involves identifying and eliminating reasoning chains that contain errors in reasoning and grounding (they define grounding error as utilizing information that is not present in cited premises).

individual step in a reasoning process has a decent odds of being nonsense.

Reducing the percentage of nonsense to a level required to significantly improve user experience will not necessarily reduce it to the level needed to make complex multi-step reasoning work effectively on a consistent basis. Also, if nonsense production is reduced, it will not necessarily be replaced with useful sensible production – it may well often be replaced by situations where the system doesn't know what to say (as fairly often transformers resort to producing nonsense when their training data and algorithms aren't sufficient for them to come up with a sufficiently probable-looking real answer).

The SciBench corpus [WHL⁺23] lies interestingly beyond the capability of current plugin-enhanced LLMs; it is a collection of university-level science problems at which current systems cannot exceed 36% accuracy. Figure 15 illustrates the sorts of errors that current LLMs make on these problems. The AGIEval corpus [ZCG⁺23] is differently focused (more heterogeneous rather than purely science-focused) but also presents a massive challenge for current LLMs. How far toward this level of complex multi-step reasoning we can get within the "Transformer as hub" architecture remains to be seen. The more advanced the problems get, the more unique each problem is and the more abstraction is needed to figure out what knowledge from prior problems and readings to generalize to deal with the current problem. It feels very intuitively clear to me that

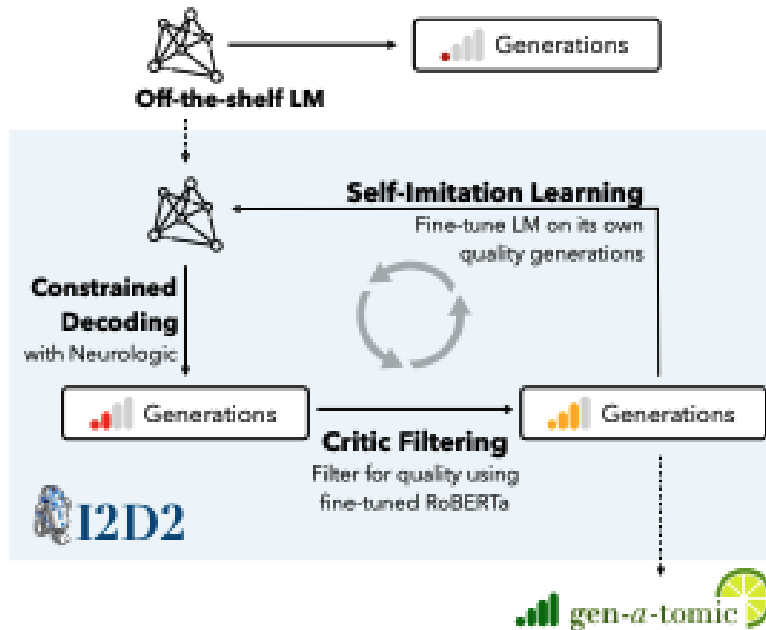


Figure 14: The I2D2 inductive knowledge distillation approach from [BHD⁺22] radically improves over off-the-shelf generation from GPT-2 XL, allowing it to perform better than the much larger GPT3 model on relevant metrics.

Problem (Using the Planck distribution)		
Compare the energy output of a black-body radiator (such as an incandescent lamp) at two different wavelengths by calculating the ratio of the energy output at 450 nm (blue light) to that at 700 nm (red light) at 298 K.		
<p>Correct Solution</p> <p>At a temperature T, the ratio of the spectral density of states at a wavelength λ_1 to that at λ_2 is given by</p> $\frac{\rho(\lambda_1, T)}{\rho(\lambda_2, T)} = \left(\frac{\lambda_2}{\lambda_1}\right)^5 \times \frac{(e^{hc/\lambda_2 kT} - 1)}{(e^{hc/\lambda_1 kT} - 1)}$ <p>Insert the data and evaluate this ratio.</p> <p>Answer With $\lambda_1 = 450$ nm and $\lambda_2 = 700$ nm, ..., and therefore</p> $\frac{\rho(450 \text{ nm}, 298 \text{ K})}{\rho(700 \text{ nm}, 298 \text{ K})} = 2.10 \times 10^{-16}$ <p>Final Answer: 2.10×10^{-16}</p>	<p>GPT Solution with CoT Prompting</p> <p>Step 1: We need to use the Planck distribution formula to calculate the energy output of a black-body radiator. The formula is:</p> $B(\lambda, T) = \frac{2hc^2}{\lambda^5} \frac{1}{e^{hc/\lambda kT} - 1}$ <p>Step 2: We calculate the energy output at 450 nm and 700 nm... Let's use $T = 298$ K.</p> $B(450 \text{ nm}, 298 \text{ K}) = \frac{2hc^2}{(450 \times 10^{-9} \text{ m})^5} \frac{1}{e^{hc/(450 \times 10^{-9} \text{ m}) k \cdot 298 \text{ K}} - 1}$ $\approx 1.23 \times 10^{-16} \Rightarrow 1.67 \times 10^{-31}$ <p>...</p> <p>Step 3: Calculate the ratio of the energy output at 450 nm to that at 700 nm:</p> $\frac{B(450 \text{ nm}, 298 \text{ K})}{B(700 \text{ nm}, 298 \text{ K})} \approx 4.39$ <p>Final Answer: 4.39×10^{-16}</p>	<p>GPT Solution with Python Prompting as External Tools</p> <p>... The ratio of the energy output at $\lambda_1 = 450$ nm and $\lambda_2 = 700$ nm is</p> $\frac{\rho(\lambda_1, T)}{\rho(\lambda_2, T)} = \frac{\frac{2hc^2}{\lambda_1^5} \frac{1}{e^{hc/(\lambda_1 kT)} - 1}}{\frac{2hc^2}{\lambda_2^5} \frac{1}{e^{hc/(\lambda_2 kT)} - 1}} = \frac{\lambda_2^5 e^{hc/(\lambda_2 kT)} - 1}{\lambda_1^5 e^{hc/(\lambda_1 kT)} - 1}$ <p>Substituting the given values, we get a Python code for this problem:</p> <pre> => lambda2 ratio = ((lambda2 / lambda1)**5) *(math.exp((h * c) / (lambda1 * k * T)) - 1) / (math.exp((h * c) / (lambda2 * k * T)) - 1) => lambda1 </pre> <p>Final Answer: 3.95×10^{23}</p>

Figure 15: Illustration of the sort of problem from SciBench corpus that current LLMs perform poorly on, even given the best current plugins and clever prompting tricks, from [WHL⁺23]

no variation on transformers is doing to get up to the human level (which is the PhD research level, let's say) in this regard, but exactly how far transformers can be pushed and with what specific limitations is an unfolding story.

5.6 Banality and Lack of Fundamental Creativity

Compared to the inability to tell hallucination from reality, ChatGPT's tendency to synthesize extremely banal and cliché textual productions may seem like a minor issue. After all, it's impressive that the system can make up a poem on any topic, approximating in its own way a variety of known poetic styles – most people can't do this nearly as well. So isn't it a bit unkind to point how extremely dull and boring nearly all the system's productions are, from an aesthetic perspective?

The superhuman banality of ChatGPT's output, however, both constrains its domains of practical application and gives some insight into how its underlying architecture affects its output.

Given the prompt

"Write a poem honoring the beauty of failure, in the style of beat poetry"

ChatGPT (back-ended on GPT-4) comes up with

In the smoky haze of neon dreams, We chase the whispers of our desires, Dancing to the syncopated beats, Riding the waves of life's wild fires.
Oh, sweet failure, a kindred soul, In your embrace, we find our truth, A mirror reflecting our fractured selves, We stumble, we falter, in pursuit.
Through alleys of uncharted whims, The echoes of defeat resound, But in the ashes of shattered hopes, A phoenix rises, unbound.
As the jazz man wails, the saxophone cries, Revel in the beauty of failure, my friend, For in the darkness of our deepest falls, Lies the spark that fuels our ascend.
In the gritty rhythm of city streets, Each misstep, a lesson learned, From broken hearts and ragged dreams, The tapestry of our lives, discerned.
Oh, sweet failure, the beatnik's muse, We honor you, the unsung sage, For in the chaos of our blunders, We find the words to fill the page.

Asked to avoid rhymed verse (which is not a common form in beat poetry, actually) via

"Write a non-rhyming poem honoring the beauty of failure, in the style of beat poetry"

it comes up with

In the midst of a chaotic symphony, Where the drumbeat of life rolls, We strive, we stretch, reaching for the stars, But sometimes, oh sometimes, we lose control.

Sweet failure, a companion on this journey, The undercurrent, the hidden gem, In the shadows of our ambitions, A testimony to the human hymn.
Through the labyrinth of desires, We stumble, we stagger, we sway, But in the embrace of your tender arms, We are sculpted, molded like clay.
The jazz man strums the strings of the soul, And the room shimmers with the scent of truth, For it is in the dance of imperfection, That we find the fountain of youth.
In the smoky dens and midnight haunts, Where the poets and dreamers convene, The beauty of failure shines through, A beacon, a guide, a gleam.
So let us toast to our broken wings, And the steps we've misplaced along the way, For in the tapestry of human experience, The beauty of failure holds sway.

Interestingly, it does know the themes common in beat poetry – sex, jazz, smoky dens, and such ... but it still projects beat poetry themes into not-so-beat metered verse with very simple rhythms.

Asked to make things more interesting via

"Write a poem honoring the beauty of failure, in the style of beat poetry. The poem should not rhyme, and should not be broken up into stanzas."

it doesn't get the "stanzas" request:

In the swirling chaos of a kaleidoscopic world We dance to the rhythm of dreams and desires Stretching, reaching, grasping at the elusive But oh, sweet failure, you are the canvas of our existence
On the smoky streets where the jazzy beat pulsates We embrace the uncertainty of life's tides Stumble and fall, bruised and battered Yet in the embrace of your tender arms, we are reborn
Each misstep, a thread in the tapestry of our souls Woven through the narrative of our human story The beauty of failure, raw and unfiltered A testament to the strength of the spirit
As the saxophone wails and the drums beat on The poets and dreamers find solace in your touch For in the chaotic dance of the imperfect Lies the essence of our truth, our resilience, our growth
So let us honor the moments that slip through our fingers And the dreams that crumble beneath our feet For in the beauty of failure, we find the courage to continue To rise, to create, to become complete

One could presumably train it to avoid breakdown into stanzas via few-shot learning, but that doesn't fix the problem that it lacks a grounded understanding of what a "stanza" is, even though it's producing poems broken into stanzas with considerable facility (though depressing aesthetics) and has observed a huge number of instances of the word "stanza" used in discussions of poetry. This misunderstanding about stanzas will likely be remedied in some modest upgrade of GPT-4, but still is interesting as an indication of GPT-4's mode of "thinking."

Explicitly asking it to be less banal and boring via

"Write another poem with a similar theme and also in a beat poetry style, but less cliché sounding and with more literary quality "

makes no difference to the quality of the results (I won't take up more space showing the output, as it's very similar to the above examples).

Compare all this pompous-middle-schooler tedium to the Allan Ginsberg's "Ode to Failure", some of the more PG-13 passages are:

Many prophets have failed, their voices silent ghost-shouts in basements nobody heard dusty laughter in family attics nor glanced them on park benches weeping with relief under empty sky Walt Whitman viva'd local losers , courage to Fat Ladies in the Freak Show! nervous prisoners whose mustached lips dripped sweat on chow lines ?

...

Prospero burned his Power books & plummeted his magic wand to the bottom of dragon seas Alexander the Great failed to find more worlds to conquer! O Failure I chant your terrifying name, accept me your 54 year old Prophet epicking Eternal Flop! I join your Pantheon of mortal bards, & hasten this ode with high blood pressure rushing to the top of my skull as if I wouldn't last another minute, like the Dying Gaul! to You, Lord of blind Monet, deaf Beethoven, armless Venus de Milo, headless Winged Victory!

...

My tirades destroyed not Intellectual Unions of KGB & CIA in turtlenecks & underpants, their woolen suits & tweeds I never dissolved Plutonium or dismantled the nuclear Bomb before my skull lost hair I have not yet stopped the Armies of entire Mankind in their march toward World War III I never got to Heaven, Nirvana, X Whatchamacallit, I never left Earth, I never learned to die.

This is by far not one of Ginsberg's greatest poems, it's certainly no "Howl." But it has fantastically much more heart, soul and inventiveness than any of ChatGPT's productions. After looking at a few ChatGPT poems, it feels spiritually necessary to remind oneself what an actual decent human poem looks like.

In the end, ChatGPT's bullshit and its banality emanate from the same source: Predicting the most probable continuation of a sequence is, in practice, an awkward way to proceed if you're looking for utterances that fulfill some other metric like truthfulness, interestingness or aesthetic quality. In principle it may be possible to formulate these other metrics in terms of probability of sequence continuation, but in practice this sort of reduction is not effectively cashed out within current transformer models, and it's not at all clear that this sort of reduction is the right way to make real-world AI systems aiming at truth or aesthetics.

In the context of my work with the rock band Desdemona's Dream, I have explored producing lyrics and spoken word poetry with transformer neural nets trained specifically on hallucinatory fiction writers like Philip K. Dick and my own self, producing

lyrics like

The Singularity will not be centralized. Instead, it will radiate out like a wasp.

Some will seek form for themselves, some will be engulfed by the radiance, become entangled with it, even merge with it.

All will be revealed to them as early as possible, since the radiance is a form of energy, that can enter and leave

...

The Singularity will not involve a constant projection into the future of a given organism, but a process by which the antecedent universe is progressively further abstracted from that organism, as if that organism itself were abstracting itself from its environment.

The two halves of the paradox, origin birth over the course of millions of years, will forever engenerate; the paradox will escape being itself, and will thereafter imprint an entirely new self on the brought into being

You will not be mind uploaded by human psychopaths into a computer, but you will be self-monitoring, and receiving telepathic feedback from it, as to your progress, your true status and level of relevance to it.

The Singularity will not be centralized.

These productions have been used by Desdemona Robot who fronts the rock band Desdemona's Dream; they seem to work reasonably well in this musical context. Now, this is not world-class poetry either, but it does have out-there PKD-ish zing to it; it is definitely outside the scope of the poetic productions of the typical middle or high school student. The more interesting-in-context artistic quality of these results as compared to GPT-4 poetry illustrates that a well-constructed small model can outperform large models for particular aesthetic purposes. But this also is not fundamental inventiveness – it's bald imitation of the style and concepts of particular inventive writers. If one fine-tuned ChatGPT on Allan Ginsberg's poetry it could surely come up with something better than its not-really-beat poetic productions cited above. But it wouldn't be able to create works with the overall coherence of Ginsberg's own productions, because this overall coherence exists at a structural and conceptual level that current transformer networks cannot embody or model.

Trying to get ChatGPT to clean up the output of the language models used for Desdemona is interesting, but not ultimately useful:

PROMPT: My friend Phil has interesting things to say, but is a bit crazy and incoherent. I would like you to take some of Phil's statements and rewrite them a little so they maintain the same ideas and spirit but are more coherent, sensible and flowing as prose. Here is one example of something Phil said: "Collect the rests Piece by piece And put the Coat of Broken Dreams On your heart And in your pocket And then walk out the door And start living your life and never asking the world again: Is it the world that rules me? It rules me! And now, slowly, slowly, slowly! I am living my life! I am living my life! I am living my life! I am living

my life! And then, slowly, I am living mine! I am living my life! And not one bit of this is real! Not one bit! Not one bit! Not one bit! Not one bit!" Can you please rewrite this to be a bit more sensible and coherent, but keep the same basic ideas and feeling?

RESULT: Certainly, here's my attempt at rewriting Phil's statement in a more coherent and sensible manner while retaining the original spirit:

Collecting the fragments of ourselves can be a daunting task, but it's necessary to move forward. Gather the pieces, bit by bit, and wrap them in the Coat of Broken Dreams, which represents the struggles and setbacks that we have faced. Wear this coat on your heart and carry it with you, and then step out into the world and begin living your life without asking permission or seeking validation. Instead of asking, "Does the world rule me?" boldly declare, "It rules me!" Embrace the reality of your circumstances and start living your life, taking it one step at a time. With each step, you'll gain confidence and clarity, until you can confidently say, "I am living my life!" And even though life can sometimes feel unreal, embrace every moment and hold onto the truth of your existence. Not one bit of it is fake or insignificant, and each step you take brings you closer to a more fulfilling life.

Clearly, the original – while questionably coherent – has some Rimbaud-like or Allen Ginsberg-like enthusiasm to it.. the rewording is just self-help-book drivel. But GPT-4 does catch some of the meaning, if not the feel

As a more practical and commercial example of the limitations posed by ChatGPT's banality, consider applications to marketing in the tourism domain.

If we ask

"Give a list of five reasons someone should be excited about visiting Toronto instead of some other city "

we get a pretty reasonable and useful list

Multicultural Diversity: Toronto is one of the most diverse cities in the world, with over half of its residents being immigrants. This melting pot of cultures offers visitors the unique opportunity to experience various traditions, cuisines, and events without leaving the city. From Little Italy to Chinatown, there's something for everyone to explore.

Iconic Attractions: Toronto is home to several iconic attractions, such as the CN Tower, the Royal Ontario Museum, Ripley's Aquarium of Canada, and the Art Gallery of Ontario. These sites not only provide stunning views and unforgettable experiences but also showcase the city's rich history and artistic talent.

Thriving Food Scene: Toronto's food scene is a testament to its multiculturalism. Visitors can indulge in a variety of cuisines from around the globe, including authentic Thai, Italian, Indian, Chinese, and Ethiopian dishes. The city also hosts various food festivals throughout the year, such as Taste of the Danforth and Summerlicious, offering foodies a chance to discover new flavors.

Vibrant Arts and Entertainment: Toronto is a hub for arts and entertainment, featuring numerous theaters, galleries, and live music venues. Visitors can catch a Broadway show at the Mirvish Theatre, enjoy a concert at the iconic Massey Hall,

or explore the city's thriving independent art scene at galleries like the Power Plant Contemporary Art Gallery and MOCA Toronto.

Natural Beauty and Outdoor Activities: Despite being a bustling city, Toronto boasts numerous parks, gardens, and waterfront areas perfect for outdoor enthusiasts. Visitors can enjoy a stroll along the picturesque Harbourfront, explore the lush Toronto Islands, or take in the beauty of High Park. Additionally, the city's extensive bike lane network and public transportation system make it easy to navigate and explore the great outdoors.

But if we ask

"Write a brief (50 words or less) advertisement intended to convince someone to visit Toronto instead of somewhere else"

we get something that is perhaps college-undergrad level but not what a top advertising agency would use.

Discover Toronto, the vibrant, multicultural metropolis! Experience iconic attractions, a world-class food scene, and diverse arts and entertainment. Immerse yourself in stunning parks and waterfronts while enjoying the city's warm, welcoming spirit. Toronto awaits, offering unforgettable adventures for every traveler!

Compare for example to the recent ad campaign "Vilnius: The G-Spot of Europe"....

Asking for something zingier via

Rewrite the advertisement you have just written, in a style that is more quirky and with fewer cliches

results in

Toronto: a city that's as diverse as your favorite playlist. Savor global cuisine, find art that makes your heart race, and soak in the skyline from our towering giant. Bike through urban oases or dance to the rhythm of the city. Curious souls, Toronto is your playground!

which is really quite cringe. "As diverse as your favorite playlist" totally reads like a cliché corporate executive trying and failing to sound cool. Which is interesting given how many actual dialogues among actual youth ChatGPT has in its training data. In the case of a system with so much training data, what appears as a banal taste actually is tied in with lack of grounding – it's not correctly grounding "more quirky" and "fewer cliches" in the patterns available in its textual training data, in spite of having copious examples of various texts labeled more or less quirky or cliché in various contexts.

Again, these are in a way high-end complaints, and it's amazing that these models are at the level where they can generate fair college-under-grad level approximations of marketing blurbs. On the other hand if we're evaluating the system via the criterion "can it actually eliminate human ad copywriting professionals" the answer is clearly:

Not quite yet. It can generate loads of ideas and raw materials in a way that integrates a breadth of relevant information, and in this way almost surely decrease the amount of human effort needed to produce a given amount of ad copy with a given quality. But this is different than performing at the level of a human ad copywriter.

One should remember the key selling point of the original GPT paper: Not that it did the best at every natural language processing task, but rather that it was best at some NLP tasks and also did well on a great variety of others without special training for each one [RNS⁺18]. This is still the case – specialized task-specific neural nets outperform GPT models at most NLP tasks studied in the computational linguistics world [Pik23]; a qualitative expert assessment as of late March 2023 is

”ChatGPT won 34 out of 151 comparisons.... In most classical NLP tasks, ChatGPT does not outperform existing fine-tuned baselines, although it is often close. In some cases, it failed to beat even simple bag-of-words models (e.g., predicting agreeableness...) or it was worse than supervised baselines by a surprisingly significant margin (e.g. ... relational reasoning ... or emotion classification ...). ChatGPT struggles with affective tasks and it sometimes shows a higher level of brittleness than older bert-tier models.... I was surprised to find that ChatGPT does not excel in text generation tasks such as summarization or question answering, even though people really like these capabilities. ChatGPT does not seem very strong with semantic similarity tasks, but it is really good at comparing generated texts to reference.”

Now, both GPT and alternative systems are moving targets, so the precise details of such assessments will become obsolete quite rapidly; but so far as I know the basic nature of this conclusion remains valid.

Somewhat similarly, across practical tasks like writing ad copy or poetry – which are in the GPT models’ domain of strength, being text generation focused rather than drawn from math, music, engineering or some other non-text-centered discipline – what we see is passable-but-not-great performance across the board rather than excellence in any area ... except perhaps the area of generating plausible-looking bullshit.

While our focus here is on LLMs for language processing, the same core issues occur in current deep neural nets for image processing. On the surface models such as DALL-E, Stable Diffusion, Midjourney and so forth are quite creative. Yes, their visual creativity is combinatorial, mixing pieces of what they’ve seen before, but an awful lot of human creative image generation is like that too. The cognitive theory of "concept blending" [FT02] explores the distinction between blending two concepts based on an understanding of their deep structure versus blending based on shallower surface characteristics, but it’s not always easy to tell the difference from looking at a picture or reading a paragraph. My AGI research collaborator Alexey Potapov has proposed the

example of inferring the nature of gravitational maneuvering of a spaceship from basic information about spaceships and gravity, as something that current LLMs or other deep NNs are not going to be able to do – because they are looking at surface aspects of the spaceship and gravity concepts, rather than building an (implicit or explicit) fundamental internal model of how spaceships and gravity operate in the context of the physical world.

”For sure, Stable Diffusion can produce an image of a spaceship doing a maneuver. After all, it can produce an image of a spaceship. It might be able to produce images of gravitating objects. But how it will imagine a spaceship gravitational maneuver without a conceptual understanding of all these words and concepts, if it hasn’t been supplied with examples of images of this thing? Such neural nets can produce images of "cat dragons" pretty well without seeing a thing like this before, because this can be done based on in some ways deep/abstract, but still purely visual (not conceptual/semantic) features. It can easily produce neon glowing pink grass, even though there is no such grass in its training dataset. But if something is obtained not by combining features observed in the training data, but by conceptually understanding combined notions, then current LLMs and DNNs are not competent, they just don’t do it. The issue is that to perform these tasks, the neural nets would need to not just recall, but also combine on the level of abstractions and generalizations that are derived via multiple complex steps from available features. This basic issue exists whether one is talking about language, imagery or any other input/output medium.”

5.7 Lack of Self-Directedness and Autonomy

By design and construction, what GPT models do is propose continuations for sequences, and what CILLM systems do is exploit this capability to hold conversations triggered by external conversation participants. In terms of Stan Franklin’s classic dichotomy, these are programs not agents [FG96].

There have been a few attempt to make agents wrapping up GPT-4 and other transformer models, e.g. AutoGPT which takes a user-provided goal, asks GPT-4 to make a plan for achieving it, and then asks GPT-4 to carry out each step of the plan. AutoGPT and other such systems are more agentic than GPT-4 or ChatGPT, and architecturally they involve adding some judiciously chosen additional components to the GPT-4 framework, such as simple forms of long-term memory, planning and goal-refinement. However, they are also extremely simplistic.

One could envision wrapping a more human-like goal-pursuit framework around an LLM; we will briefly mention two of the many possible ways of doing this in sections 8.2 below. But this would rapidly turn into creating a broader sort of cognitive architecture with an LLM as a (possibly central, possibly less so) subsystem.

However, pursuit of a user-specified goal is still a long way from emulating the way goals related to human cognition. In the human psyche, goals emerge from a mind's self-organization and environmental coupling at its given state of development, and then morph and shift over time. The evolution of a person's goals is part of their overall cognitive growth process. Of course, people can pursue highly specific externally proposed goals, such as doing well on an exam or winning a race or performing a required task at work, but they are adopting this goal because it fits in with other aspects of their mind, life and situation at that time.

Further, the context guiding a person's adoption of a specific externally proposed goal is often, though not always, pertinent to the particulars of *how* they pursue that goal. For goals involving, say, artistic creation or scientific discovery or deep connection with other people, this is an especially relevant dynamic. The inability of systems like ChatGPT or AutoGPT to contextualize their goals in their lives and world-contexts is connected to their failure as artists and innovators. It is also connected to the inability of LLMs to distinguish truth from falsehood, because once basic issues involving comparison of abstract statements to concrete data are solved, one then hits issues regarding the relation of truth to the observing mind, in which is hard for a system that does not have in any meaningful sense *its own truth* to form a useful sense of the truths of various statements relative to various communities.

5.8 Foundational Computational Limitations of LLMs

Alongside the practical weaknesses of CILLMs in various aspects, it's also worth understanding some of their fundamental mathematical and computational limitations. Used in the standard, straightforward way, these systems are not "computationally universal learning machines" in the same sense that many other AI paradigms (e.g richly recurrent neural networks, genetic programming, logical reasoning) are.

For instance, Hahn [Hah20] takes a detailed look at how good self-attention is at modeling formal languages. We find, for both soft and hard attention, that self-attention has notable limits - it struggles to model certain types of structured languages, and can't handle hierarchical structures unless you keep adding more layers or heads as the input grows. This may seem a bit odd, given how well self-attention has worked in real-world applications and how important these types of structures are thought to be in language. The most intuitive explanation of this peculiarity, however, is that natural language – which in the end mostly involves a quite finite vocabulary used in fairly brief sentences covering a set of themes that is relatively modest-scale in a theoretical sense – can be effectively handled by models that aren't powerful enough to deal with the types of

formal languages often used in linguistic theory.

The relation between these fundamental computing-theory limitations and the various practical limitations of CILLMs is fairly indirect. Most of the practical tasks on which CILLMs fail could be done perfectly well by other sorts of systems with general computing scope far short of universal Turing capability. However, there is also some connection; e.g. CILLM failures at mathematical abstraction appear (albeit complexly) related to the way the architecture bypasses more abstract forms of computing.

There are also fairly elementary – if brute-force-ish – ways to work around these limitations, e.g. by giving a CILLM an appropriate auxiliary memory to use as a "scratchpad", and then training it to carry out appropriate operations on this scratchpad. Some might argue one is then doing something somewhat orthogonal to the demonstrated main strengths of CILLMs, to the things and ways that they seem especially good at learning. But still there is clearly a lot of power here.

Along these lines, Schuurmans [Sch23] demonstrate that transformer-based language models, when given extra storage space, can process virtually any computational task. Language models that focus only on strings of a certain length can only do so much, similar to a basic automation machine. However, giving these models access to a memory that they can use to store and recall information allows them to handle much larger inputs and, potentially, mimic any algorithm. They show that an existing language model, Flan-U-PaLM 540B, when paired with a memory that it can use to store and fetch data, can perfectly imitate a universal Turing machine. Notably, "A key aspect of the finding is that it does not require any modification of the language model weights. Instead, the construction relies solely on designing a form of stored instruction computer that can subsequently be programmed with a specific set of prompts."

This sort of augmentation of CILLMs with external memory leads naturally into discussion of the power to potentially be obtained by extending CILLMs in various ways and incorporating them into broader cognitive architectures, topics we will turn to a little later in this article.

6 Human-Like Cognition vs. LLM Operation

Enumerating the various cognitive weaknesses of current LLMs compared to human minds is interesting, but might also give the mistaken impression that one can somehow turn an LLM into an HLGAGI via improving various of its aspects or making a few add-ons. This becomes especially confusing because we don't have a highly precise definition of where the border lies between a "modest improvement to the internals"

with “a few add-ons” versus, say, a radically different AGI architecture that includes an LLM as one of its many important components.

To make things clearer it may be helpful to start from the other direction, and begin from a breakdown of how human intelligence is understood to work, and then consider how CILLMs work (and what they don’t do) in this light. One way to do this would be to make a systematic comparison of current CILLMs against the rough cognitive process breakdown of human intelligence implicit in the Standard Model of Mind. To do this in a fully serious way would be a quite large undertaking, beyond the scope of this already-rather-long paper, but we will aim here to make a rough start.

6.1 How do Current LLMs Fare According to the Standard Model of Mind?

Running through the core cognitive components identified by the Standard Model of Mind one by one:

- **Episodic Memory**

- CILLMs do not maintain a lifelong episodic memory in the same sense that humans do.
- One could extend a system like AutoGPT to include a more sophisticated episodic memory model, but then integrating this model in a rich way with the knowledge in the LLM would require some sophisticated work which might or might lead beyond the LLM-centric approach to AGI.
- CILLMs can understand what episodic memory is, and can somewhat reason about how a person’s episodic memory might react in certain situation (e.g. if given a prompt describing a person’s life-story, and asked what that person might recall given a certain suggestion, CILLMs can come up with a reasonable answer), but this is very different from actually *having* an episodic memory and using it to guide one’s interactions with the world.
- The sort of associative lookup that CILLMs are good at seems roughly similar to the way episodic memory lookup works, however key elements are missing such as indexing of events by space, time and relationship to self and important others.

- **Working Memory**

- CILLMs have a sophisticated form of linguistic working memory, enabled by transformer neural nets’ attention mechanism, which is a key ingredient

of LLM functioning (though there is research suggesting that the particulars of the attention mechanism are less important than has been assumed, and other mathematical transforms might serve equally well [LTAE021]).

- CILLMs don't yet have the other forms of working memory buffer that human minds do
- A multimodal CILLM handling both text and video would have two corresponding implicit working memory buffers, but whether they would coordinate together with the sophistication of the corresponding human working memory buffers remains to be seen
- The interaction between CILLM working memory and CILLM long-term memory seems dissimilar to, and less sophisticated than, the corresponding interaction in the human mind, which seems related to the difficulties CILLMs have with complex multi-step reasoning. Complex multi-step reasoning involves iterated interactions between working memory and long term memory and associated learning/reasoning processes, in ways that appear fundamentally beyond CILLM capabilities.

- **Procedural Memory**

- CILLMs do not possess procedural memory in the sense that human minds do
- CILLMs are able to model procedures as step-by-step descriptions in natural or formal languages, and do to some reasoning about procedures in this way, including reasoning about slightly-to-moderately complex programmatic constructs and everyday-life plans
- Reasoning about more complex programmatic constructs and everyday-life plans is quite limited and error-prone in CILLMs. This means for instance that they are able to design software programs only to handle situations that are quite close to those covered by examples in their training databases. It also means that they are unable to handle moderately complex travel planning in the manner of a human travel agent, even though they can help with many of the steps involved in creating a travel plan (choose attractions to see based on personal preferences, select appropriate hotels, etc.).
- Improvement on procedural reasoning tasks can be obtained via leading the CILLM through the steps of the process, but this involves invoking the human user's procedural knowledge to supplement the lack on the CILLM's part

- The core issue here is that, without a more robust world-model, CILLMs have only a limited ability to understand procedures in terms of their effects on a world.
- Rather, CILLMs are doing something more like: Breaking a procedure into sub-procedures, and for each one finding other subprocedures in the training database that are syntactically somewhat similar and occur in similar contexts, and then looking at what specific effects have been found as ensuing from these other related sub-procedures. Combining this sort of training-data-lookup related to various subprocedures of a provided or synthesized procedure is a very clever and interesting thing to do – and a strategy humans also take from time to time – but it’s very different from procedural learning and reasoning that is principally grounded in (explicit or implicit) world-models.

- **Reasoning**

- As we have reviewed above, CILLMs are able to carry out many interesting cases of mathematical and commonsense reasoning, but their capability falls apart as the needed reasoning gets more complex, involving more steps or steps diverging further from specific training data
- Commercial CILLMs are currently approaching reasoning in substantial part via hybridizing LLMs with other subsystems via various strategies, e.g. ChatGPT’s invocation of Wolfram Alpha, or Chameleon’s use of GPT-4 to weave together other services
- It appears likely that achieving human-level capability at complex reasoning will require substantial architectural and dynamical elements beyond CILLMs
- Neuroscience does not currently understand much about how complex reasoning occurs in the brain, so there is no strong biological inspiration for how to make this sort of functionality work in a neural net AI system.
- AI systems founded on explicit symbolic logic are good at complex reasoning, but have not yet been implemented with a combination of rich generality and large scale, and their ability to ground commonsense reasoning in large bodies of perceptual, motoric or linguistic data or experience is not yet validated. Scaling up such systems, and hybridizing them with LLMs, is an active area of research.

- Advanced mathematical reasoning stresses advanced modes of computation of which LLMs are incapable without extensions or highly unnatural approaches. Human brains appear capable of adapting to emulate these more advanced modes of computation in a more natural way, though it's still not easy for them. Automated theorem provers and other symbolic reasoning systems show that this shortcoming of CILLMs is not a limitation of modern digital computers in terms of their architecture or scale, but rather a limitation of LLM architecture.

- **Reinforcement Learning**

- CILLMs are often trained via RL, most notably via InstructGPT which was a key method used to fine-tune the GPT-4 model to obtain the ChatGPT chatbot personality.
- It's not clear what happens if one tries to train an LLM based on RL according to multiple different reward functions operative on multiple time-scale, or in accordance with goal-driven learning based on goals that are different from expected reward

- **Language Learning and Usage**

- NLP is clearly a strength of LLMs, however, it's noteworthy that general-purpose LLMs don't currently beat specially fine-tuned systems on all standard NLP benchmarks, the best results often going to systems trained via supervised learning [Pik23]
- Still, the ability of a single LLM system to perform at a high level on a broad variety of NLP tasks is unprecedented and a sizable step beyond prior end-to-end NLP systems from a pure computational-linguistics perspective, setting aside other aspects of LLMs

- **Multimodal perception**

- LLMs can be trained on multiple data types e.g. text, video, audio with valuable, interesting, high quality results
- Semantic understanding of non-linguistic data is still fairly mediocre in these cases.
- This merits rigorous evaluation but as LLM users we find it qualitatively, anecdotally quite clear. E.g the text prompts for language/music hybrid systems like MusicLM and MusicGen statistically bias the nature of their

output in a manner consistent with the semantics of the text, but don't precisely guide the nature of the output in the way one would expect from a human.

- **Action learning and coordinated action**

- LLMs have not been shown competent at coordination of actions of complex actuator systems (let's say, humanoid robots as an example) toward goals
- This appears to be an interesting research challenge but it's not clear how far one can go with LLMs as the primary AI engine
- E.g. one can train a transformer on data from human puppeteers controlling humanoid robots, and run it generatively to generate strikingly human-looking movements (based e.g. on the authors' anecdotal observations of proprietary systems at humanoid robotics firms). However, there seems to be no currently operational way to use this for motion planning. CILLMs can do some high level motion planning in a very rough way, but they don't understand spatial geometry very well, so it's hard to get them to do this in a way that can productively connect with transformers applied to movement generation..
- How far one can go in this direction with multimodal LLMs co-trained on movement data, spatial 3D vision data, and natural language remains to be seen. However, it should be kept in mind that even in NLP tasks, CILLMs are not the maximally capable systems for most benchmarks.

- **Goal refinement and goal system management**

- LLMs understand what goals are, and frameworks like AutoGPT show that they can be wrapped in goal-oriented frameworks. LLMs can also reason about what actions may be likely to achieve what goals in what contexts.
- Current LLMs are not capable of systematically pursuing complex goals over time in a complex environment, for one thing because of their lack of the right sorts of long-term and episodic memory.
- Simply gluing long-term memory onto an LLM in a crude way doesn't result in a combined system that can leverage long-term memory in all the ways needed for pursuit of long-term goals in complex environments. One would need an appropriate deep integration of long-term memory with the LLM, which as noted above is an interesting multifaceted research direction that rapidly leads one well beyond LLMs in their current form.

- Goal system management, in terms of refining high-level goals, balancing multiple top-level goals, pruning alienated goals, managing partially contradictory goals etc., is a frontier that is not possible to evaluate in current CILLMs given their memory architecture limitations, however, it seems likely that these require some sort of quite different high-level control architecture than anything CILLMs possess, again bringing us well beyond current architectures.

- **Reflexive self-understanding**

- Current CILLMs are simply not instrumented for overall reflexive self-understanding. Their self-understanding even within a given conversation is extremely limited, as illustrated e.g. by the "repeatedly incorrect factoring" example given above – at no point there does GPT-4 inspect the overall course of the dialogue and realize "Hmm, it seems I really don't know how to factor integers of this size."
- CILLMs lack what one might call "reflective sensing" – they aren't set up to take their own recent or long-term interaction history, let alone their internal states, as cognitive content. Humans seem to do this very naturally, young children do this in various ways even before they display any advanced linguistic capability [Tom03], which suggests some fundamental architectural differences are at play.
- Limitations on Theory of Mind as cited above are also indirectly relevant here; even at a conceptual level, the same issues that stop CILLMs from doing well on theory-of-mind puzzles would also stop it from solving "puzzles of self-understanding", even if CILLMs were somehow instrumented with the "reflective sensing" mechanisms to approach such puzzles in a data-driven way.
- Other sorts of current AI architectures handle reflective sensing very directly, e.g. this is a design principle of systems like Aera [TH12] and Hyperon [ea23].

- **Modeling and Understanding of Other Minds**

- Empirical limitations of current CILLMs at solving ToM problems have been reviewed above
- It would appear humans cognize about other minds by a combination of mechanisms, including reasoning based on multimodal data, our episodic

life-history memory and commonsense knowledge; and also including rough internal simulation of other minds (e.g. using mirror systems and other mechanisms [ACZ⁺22]), which perform low-level analogy between our own experience and that of others.

- Lack of episodic life-history, and lack of reflective sensing enabling treating one's own experience as cognitive content to be analogized to, render current CILLMs incapable of doing even vaguely as humans do in this regard
- There seems no clear way to remedy this without major architectural changes, making it unlikely that incremental improvements on current CILLMs will be able to "relate to their users as human beings" in any sort of genuine way, though they may get better at faking this.

6.2 LLMs and the Human Brain

We have compared LLM capabilities with the scope of human cognitive capabilities – but what about the comparison between LLMs and the human brain? One of my research colleagues, who is more enamored of LLMs as a path to AGI than I am, has suggested to me that "What LLMs are doing approximates what humans are doing on the micro level to do something close to what they do on the macro level."

I can see why someone would think this: LLMs are a kind of neural network (where "neural" refers to human brain cells), and some of their high level functions do indeed resemble things that human brains do.

However, it's important to remember that

- LLMs are based on formal neurons which are a crude approximation of biological neurons .. and then trains the weights of the connections between these formal neurons via the backpropagation algorithm, which is well known to be utterly nonbiological in nature
- Formal neurons are a very flexible computing substrate, and there there are incredibly many ways to make them approximate various human-like behaviors, including many ways that have very little resemblance to the actual structures and dynamics in the human mind.

And these are not theoretical but very practical points. Consider for instance: It seems very clear that when an LLM is e.g. factoring a number (or rather trying and failing) what it is doing bears very little resemblance to what happens in a knowledgeable human's brain when they factor a number.

Similarly, when an LLM comes up with a poem, the process it is following bears very very little resemblance to what a human poet's brain is doing when it composes poetry.

There could be cases where LLM dynamics vaguely resembles human brain dynamics. For instance, when an LLM writes a boilerplate recommendation letter, maybe (this is just a speculative conjecture) it is doing something in some ways in some respects similar to what a human mind/brain is doing when it writes that letter ... because this is something that, when humans do it, is basically a process of statistical modeling of other similar products that have been seen before.

I.e. in the case of blithely producing boilerplate text, the sort of process that LLMs do all the time, may actually be the sort of thing human brains are doing as well – but this parallel seems not to exist the vast majority of the time. In the case of human brains factoring numbers or writing poetry "statistical modeling of stuff seen before" may play a nonzero role, but it's certainly not the dominant factor ... and this is connected with why human brains are OK rather than terrible at these things.

I will discuss below Bengio & Hu's ideas about how to incorporate LLMs in a more comprehensive AGI architecture. Zeng Yi's BrainCog approach, mentioned above, is also NN-based and reasonably sophisticated. Both of these approaches have some plausible story about how they could account for all the aspects of human-like cognition enumerated in the Standard Model of Mind.... The approach of "add a bunch of plugins and external memories to a mixture of LLMs", on the other hand, does not really have a plausible story of this nature.

6.2.1 Emergent Dynamics in LLMs versus Human Brains

One can observe that, like human brains, LLMs feature emergent phenomena, in which the whole system behaves in ways that are not entirely obvious from simple consideration of the parts. It's true and fascinating. Clearly some valuable phenomena are emerging in LLMs. However when LLMs try to factor numbers or write poetry, the things that are emerging vary from wrong to badly suboptimal, if one's interest is human-like intelligence or powerful AGI. The emergent dynamics within LLMs, in this case, are quite different from the emergent dynamics in human brains carrying out these tasks. Not all emergence is created equal.

It's exciting to be at a point in the history of AI where we are creating large systems, out of vaguely brain-like components, that are displaying intriguing emergent dynamics related to their intelligent behaviors, which in many cases are roughly or nearly human-level in quality. But nonetheless, it doesn't take that much scrutiny to see that

the particular emergent dynamics here are mostly neither very human brain like nor human mind like, and that this is closely related to the numerous profound cognitive shortcomings of these systems relative to human intelligence.

Rigorous study of the particular emergent dynamics happening inside transformers as they learn and respond is just beginning, and surely there is much to be understood here. How relevant this understanding will be to AGI is an open question.

7 Working Around the Limitations of Today's Generative AI Systems

Reflecting on the somewhat extreme limitations of CILLM systems relative to human minds, one might wonder momentarily whether there is any hope at all of getting practical use from them. But then one remembers the sometimes-dramatic practical power as illustrated in the ethical-judgment examples given above, and countless other examples elicited by others and formally and informally reported in so many places. There is tremendous practical value in current LLMs, even setting aside the numerous near-term routes to alleviating their limitations, and the longer-term routes to significantly extending their architectures and functionality or embedding them in richer cognitive architectures. It's just that the art of making practical LLM applications involves recognizing and working around the limitations, rather than pretending that these systems are AGIs or decent approximations thereof.

The practicality of working around CILLMs' cognitive limitations and integrating them into useful practical applications is certainly one of their strengths, not so much a "cognitive system" strength but certainly at very least a "narrow AI system component" strength. This is a sort of strength about which we don't have any sort of systematic theory comparable to the Standard Model of Mind, and so we will explore it here only briefly and anecdotally.

The nature of applied AI work is that the real world becomes highly domain-specific, and 80 percent of the work on a given application generally has to do with the application rather than the AI. This has been true since the dawn of AI and LLMs are not close enough to AGI to change this. This also means that any one example of working around the limitations of LLMs to achieve practical value is going to be non-representative in significant respects. But nonetheless it seems apropos to give a concrete example to get some core points across.

Consider, then, as an illustrative example of applying LLMs in the real world, the idea of using LLMs for natural language music search engines. Ben and his colleagues

have done some experimentation in this regard, using GPT-4, GPT-3.5-turbo and Llama as well. Both the successes and failures are instructive.

(A few caveats to keep in mind while looking at these examples. The same LLM may give different results to the same query on different instances; and while LLMs are not advancing toward AGI very rapidly, their improvement on concrete queries like these can sometimes improve significantly from one version to the next. These examples are samples from a larger number of exploratory attempts in this domain, but they are not part of a scientific evaluation of LLMs for music search; these examples are anecdotal and illustrative, designed to help concretely get across a few general points. Even if some version of some LLM is improved to do better on these particular examples or in this domain, the general issues evoked here may still hold more broadly.)

For simple music search queries one gets decent results, e.g. from Llama:

Query: Virtuoso guitar soloing with a jazzy feel but heavy metal style rhythm section

Recommended tracks: Flying in a Blue Dream by Joe Satriani For the Love of God by Steve Vai Tender Surrender by Steve Vai Surfing with the Alien by Joe Satriani Black Star by Yngwie Malmsteen

("Black Star" is not really jazzy in a music-theoretic sense, but it's wildly improvised and thus jazzy in spirit, and it's an amazing track so we'll give Llama a pass.)

Query: Quirky female vocalists with intellectual lyrics and musical backing that's not super repetitive

Recommended tracks: The Party by St. Vincent Water Fountain by Tune-Yards Us by Regina Spektor 1234 by Feist Seventeen by Sharon Van Etten

(These result lists were obtained by extracting the song/artist pairs from the sentences the dialogue system provided in response to the query.)

When the target music is a little more obscure, the results get odd:

Query: Modern classical music with unusual time signatures and a bit of dissonance, but also actually listenable

Recommended tracks: Come Out by Steve Reich Different Trains by Steve Reich Music for 18 Musicians by Steve Reich Electric Counterpoint by Steve Reich Clapping Music by Steve Reich

Steve Reich is awesome, but that's a bit much. So what happens if one tries to suppress Mr. Reich from the results?

Query: Modern classical music with unusual time signatures and a bit of dissonance, but also actually listenable, and not by the great Steve Reich

Recommended tracks: Etude in E-flat Major by Oliver Schneller Electric Counterpoint by Steve Reich Glassworks by Philip Glass Piano Phase by Steve Reich Music in Twelve Parts by Philip Glass

What is interesting here is that the model did not understand the negation ("not by ... Steve Reich"). These results are from Llama but ChatGPT back-ended on GPT-4 behaved similarly.

Queries about less popular genres also don't do well:

Query: best nintendocore songs

Recommended tracks: "The Way" by Anamanaguchi "Ode to the Sun" by Horse the Band "Hyperreal" by Tetrastar "Swoon" by I Fight Dragons "Epic" by We Are The Physics

Here GPT-4 has recommended some great songs, but most are not Nintendocore music, even though they have some conceptual relation to video games. We are hitting here the fact that GPT-4 has never actually listened to any music, and can interpret and classify music only by the words associated with it online – which are not so numerous except in the case of extremely popular songs.

The following queries also highlight GPT-4's difficulty with negations:

Query: heavy metal that sounds kind of like Deep Purple but has less stupid lyrics

Recommended tracks: "Holy Diver" by Dio "Children of the Grave" by Black Sabbath "Rainbow in the Dark" by DIO "Highway Star" by Riot V "Victim of Changes" by Judas Priest

I am quite sure nobody has highlighted these DIO songs as having especially intelligent lyrics. Also, Riot V's version of "Highway Star" has exactly the same lyrics as Deep Purple's version, and this particular song is not the most lyrically brilliant of Deep Purple's oeuvre (a very low bar). Again it's abundantly clear to anyone who knows the domain that GPT-4 is not responding based on having heard these songs nor based on having looked up the lyrics of the song and trying to assess their stupidity or otherwise in some way.

Query: romantic instrumental pop music with guitar and violins, no vocals, medium speed, not too quiet and with a happy vibe

Recommended tracks: Photograph by Ed Sheeran Young and Beautiful by Lana Del Rey Fix You by Coldplay Radioactive by Imagine Dragons Fire and Rain by James Taylor

The "no vocals" query elicited a list of great songs with very nicely sung and very prominent vocals.

Asking a query that probably has no sensible answer results in some hilarious GPT-4 hallucinations:

Query: Something halfway between Garth Brooks and Schoenberg, appropriate for cleaning the muck out of my basement

Recommended tracks: "The Twelve-Tone Two-Step" by Garth Schoenberg and the Basement Boys "Schoenbrooks Shuffle" by Garth Arnold and the Muck Cleaners "Basement Blues Variations" by Garth Schoenberg and the Atonal Cowboys "Atonal Cowboy" by Garth Schoenberg

We really want to hear "The Twelve-Tone Two-Step" ! But it seems we will have to wait for a few AI software advances to get there.

Finally what was decided to leverage current LLMs effectively in this particular application was to take a totally different approach, and leverage few-shot in-context learning to translate natural language queries into the structured queries taken by a third-party dedicated music search engine.

As a first stab, the following prompt was used to show the LLM how it should translate natural language queries into structured lists of song properties:

```
prompt = TASK DESCRIPTION: For the natural language query, provide relatable tags for song attributes exclusively from the provided lists.
```

```
Mood:(','.join(music_tags('advanced_mood_tags')))
```

```
Movement:(','.join(music_tags('movement_tags')))
```

```
Voice: ('','.join(music_tags('voice_tags')))
```

—

```
Query: Techno music with complex melodies and a disturbing feel
```

```
Answer: Mood: disturbing, cold, aggressive. Movement: robotic. Voice: instrumental.
```

—

```
Query: Quirky female vocalists with intellectual lyrics and musical backing that's not super repetitive
```

```
Answer: Mood: quirky, thoughtful. Movement: nonrhythmic. Voice: female.
```

—

```
Query: Modern classical music with unusual time signatures and a bit of dissonance, but also actually listenable
```

```
Answer: Mood: mysterious, strange. Movement: nonrhythmic Voice: instrumental.
```

—

```
Query: Mellow reggae music
```

```
Answer: Mood: feelGood, relaxed. Movement: groovy. Voice: male.
```

Given this priming, the LLM can then repeat the pattern and effectively translate unstructured natural language queries into structured ones:

```
"<Query>: Virtuosoic guitar soloing with a jazzy feel but heavy metal style rhythm section <Answer>:Mood: adventurous, determined. Movement: driving. Voice: instrumental"
```

Further developments along these lines, involving more complex property-lists giving quantitative parameters of the music sought, also work reasonably effectively though they require a greater number of examples in the prompts.

One can extend this functionality in worthwhile practical directions; for example, suppose one wants to make an app enabling an advertisement-maker to find appropriate background music for their ad. One can use a CILLM to allow them to enter some text describing their ad and its vibe, and a list of songs in the general direction of what they're looking for ... and then have the CILLM work together with a music search engine to make them some recommendations.

For instance

You would like to hear:

Hard-driving rock music with a steady beat, power chords and a slightly grungy feel

Example songs:

Ted Nugent, Cat Scratch Fever

ZZ Top, Sharp Dressed Man

Your product:

a rugged SUV vehicle

Your ad visuals:

Tough looking men driving the SUV along country roads, passing by a beautiful country girl who winks at them

Here is what AI thinks:

Relevant Musical Characteristics:

Genre: Hard Rock Mood: Energetic Instrumentation: Guitar, Bass guitar, Drums

Tempo: Medium to Fast

Energy Level: High

Vocal Presence: Moderate

Summary of User Examples:

Based on the example songs provided by the customer, there are two main themes evident: aggressive guitars riffs, and upbeat rhythms.

Both songs also have moderate levels of vocals.

Additional Musical Characteristics:

Key: Major Key

Chord Progression: Power

Chords Harmony: Dominant

Chords Melody: Distorted Power

Chords Beats Per Minute (BPM): Medium to Fast

Analysis of Advertisement:

The advertisement is designed to portray the image of a tough, durable SUV that is ideal for off-road use. Its target audience is mainly young, urban professionals who enjoy adventure sports and outdoor activities. The tone of the advertisement is confident, assertive, and masculine.

Refined Music Description:

A hard-driving, high energy rock music track featuring distorted guitar riffs, heavy bass lines, energetic drum beats and a slight grunge feel. The vocals should be low key and used sparingly.

Five Suitable Artists:

Nirvana
Guns N Roses
AC/DC
Metallica
Queen

Selected Songs:

Metallica-Nothing Else Matters;
Guns N Roses-Welcome to the Jungle;
AC/DC-Highway to Hell;
Queen-Another One Bites the Dust

One among many next steps for this sort of application may involve leveraging LLMs trained on music alongside language, such as GoogleLM and MusicGen (the latter being available OSS so amenable to experimental usage). This may yield enhanced functionality but based on preliminary experimentation, will not have a revolutionary impact given the current state of the technology. The level of semantic understanding of music these models possess is impressive yet strictly limited. Anecdotally, e.g., MusicGen can come up many amazing melodies based on text prompts, or combinations of text and melody prompts, but its understanding is only very erratically compositional. It can come up with an "Afghan style melodic line" or a "neoclassical metal melodic line" but can't consistently figure out "Afghan style melodic line played via neoclassical metal guitar" (it can occasionally come close depending on the rest of the prompt) – whereas any decent neoclassical metal guitarist familiar with traditional Afghan scales could take a stab at this. It also lacks basic meta-knowledge of song structure, e.g. it doesn't understand things like "slow in the first half, then building up momentum and wildly fast by the end."

This simple little application automates a real-world task, thus drastically reducing the amount of time it takes an ad-maker to do this task. Its performance at this task is superhuman in some ways, below that of the most relevantly talented humans in other ways, but basically fine from a practical business standpoint. This app will not eliminate anyone's job, but it has potential to eliminate human work-hours, thus having the direct impact of modestly reducing demand for human employment (though estimating indirect impacts is a different matter). However, the clear human and economic value here does not weaken the severity of the cognitive limitations CILLMs display in a music context when not integrated with other software and integrated into applications that draw on their strengths.

The usage pattern illustrated in this simple music search app, in which LLMs are leveraged to translate unstructured natural language queries into the structured inputs

desired by traditional software systems, is one that is going to bear tremendous fruit over the coming years. It would be fair to say that all sorts of complex professional application software will be revolutionized in this way. Anyone who has struggled for hours searching through the obscure sub-sub-submenus of complex financial-analytics or art-creation software will understand the value of a system that can do this for you based on loose natural language queries.

However, it's important to understand that this fascinating and revolutionary usage of LLMs – like all the others, and there are going to be a lot – is something that will be crafted carefully, domain by domain, via humans with expert knowledge using judicious prompt engineering to work around the profound cognitive limitations of the underlying LLMs.

The *general intelligence* here, such as it is, is largely being deployed by the human application designer who is figuring out how to piece together what the LLM does well with what other software does well, in order to fulfill practical needs in a certain application context. Chameleon, discussed above, is trying to do this sort of "piecing together" in a fully automated way, directed by the LLM – but at the moment it only deals with much simpler cases.

When something like Chameleon is able to piece together an application like the one described above as a solution to a prompt like

Make an app enabling an advertisement-maker to find appropriate background music for their ad, via leveraging LLMs along with other available software tools.

then we will be a step closer toward HCAGI at least as regards the domain of app design. It doesn't seem impossible to get to this point via upgrades of current CILLMs without a breakthrough to full-scale HCAGI; however, it also doesn't seem obvious that this will be feasible. There is much to be discovered here.

Another example with somewhat similar characteristics is DeepMind's WebAgent [GFH⁺23], an innovative composition of two distinct language models, HTML-T5 and Flan-U-PaLM, combined in a manner shown in Figure 16, which achieves efficient web automation tasks that involve navigating and processing HTML documents, including tasks such as filling out web forms. This is an elegant and helpful architecture – but it's also an assembly done using human general intelligence to work around the limitations of specific LLMs by gluing two together in a judicious way. Again if this design had been done by a Chameleon-like system in response to a user text prompt, we would have something more interesting – and something that may or may not be achievable given modest tweaks or extensions of today's CILLMs.

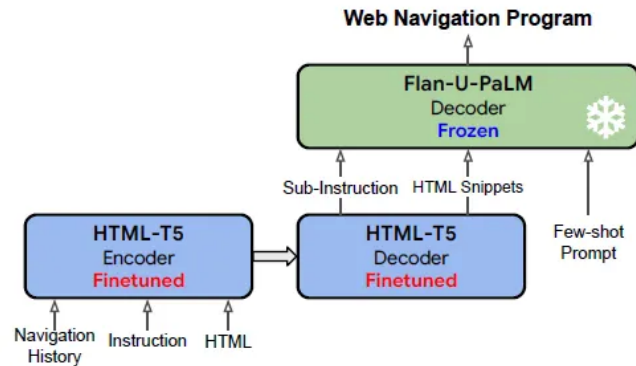


Figure 16: DeepMind’s WebAgent architecture combining two language models to achieve effective web automation [GFH⁺23]

8 Efforts to Close the Gap Between Generative AI and AGI

The inability of LLMs in anything near their current form to lead to genuine HCAGI does not intrinsically imply that LLMs are not useful toward the goal of HCAGI. Indeed the shortcomings of these systems are in some ways educational regarding what HCAGI is *not*, which in itself comprises a limited positive value toward the quest for HCAGI. Beyond this, there are current research efforts underway involving both

- Adding new components and elements to LLMs so as to make them more capable in ways that also bring them at least slightly closer to human-like cognition
- Incorporating LLMs as components within broader cognitive architectures, that are designed to incorporate the various basic structures and dynamics of human cognition in ways that LLMs do not

8.1 Potential Incremental Augmentations to Current LLM Architecture

It’s not hard to envision various ways of quite significantly improving current LLMs via modest “tweaks” to their architecture or by integrating them with other components in judicious ways that don’t go so far as to represent a shift toward a whole different way of doing AI. This is of course what the “LLM industrial complex” is likely to be working

on for the next few years, and there will be many papers and codebases exemplifying this direction (though one also expects a lot of the work in this direction will not be made public, just as many of the particulars behind what makes GPT-4 different from GPT3 have not been disclosed by OpenAI). A partial list would be:

- Enable “online learning”, i.e. updating of the neural network’s weights incrementally based on new data that comes in.
- Connect the LLM to external memory stores, thus giving it a long-term memory
- Enable runtime decision of when to read/write from memory (the Differential Neural Computer from Deepmind was one historical system exemplifying ways to do this [GWR⁺16], and there is a current literature upgrading and extending it)
- Fine-tune for the creation of plans that incorporate and combine external tools, after the fashion of the Chameleon system reviewed above but more sophisticatedly
- Integrate a knowledge graph in a manner that allows it to serve as “shared state” btw LLM and external sources (such as e.g. Wolfram Alpha). This becomes a simple but potentially powerful form of “cognitive synergy” [Goe17].

Further ideas not on this list will no doubt be innovated by the leading firms in the LLM space, but one expects that the above at least captures some significant aspects of what will be rolled out in subsequent GPT versions and their competitors.

8.2 Potential for Incorporating LLMs into Broader AGI Architectures

An alternate approach to leveraging LLMs toward the goal of HCAGI, distinct from the strategy of adding components onto LLMs, is to insert LLMs among the multiple components of a multi-module AGI architecture, in which the LLMs is not necessarily the central or controlling element, but nonetheless makes a critical contribution to the overall system intelligence. This could be done in a large number of ways of course.

Conceptually speaking, this may be roughly consistent with what is known about how the human brain works; the brain has several hundred distinctly-functioning sub-networks, each carrying out their own functions and cooperating with each other in various ways [MG18]. Most of these networks span multiple brain regions. Transformer neural networks don’t closely correspond to any particular biological brain networks,

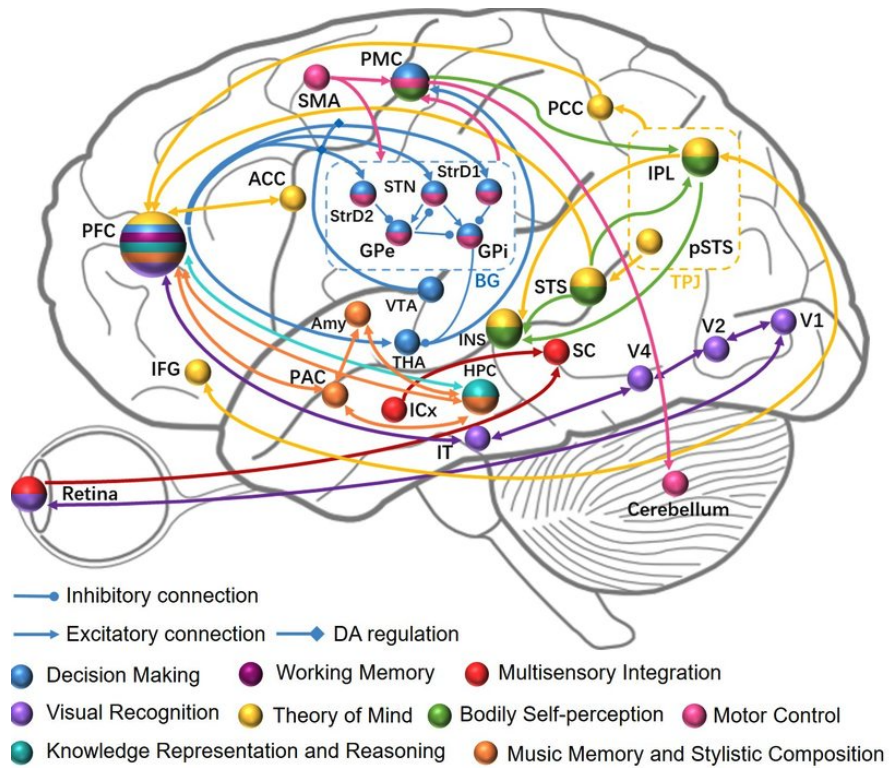


Figure 17: Aspects of the neurocognitive model underlying the BrainCog brain-inspired AGI software system developed by Zeng Yi and colleagues in Beijing

so far as is currently known; however it would be more plausible to hypothesize that there are subnetworks in the brain that operate somewhat like transformers, than that transformers are a reasonable model of the human brain as a whole.

Among the numerous networks guiding brain activity are multiple involved with executive function – high-level control and coordination – and it’s clear that, while these are bound up with attention mechanisms and rely on hierarchical pattern recognition and synthesis, these do not overall operate in a manner analogous to transformers. There are also parts of cortex whose dynamics are governed more by combinatorial connectivity structure than by the hierarchical structure that is most prominent in transformer networks. Generally it seems clear that to the extent there is transformer-ish structure and dynamics in the brain, it is serving mostly in a subordinate role to other sorts of executive-control dynamics, as well as being coupled with varying degrees of tightness to other brain sub-networks operating according to different sorts of principles.

We’re not aware of anyone concretely working to weave LLM-type structures and

dynamics into biologically realistic neural models; however there are some concrete proposals regarding how to effectively leverage LLMs within multi-modular AGI architecture structured according to non-LLM-centric cognitive models. A broad brain-based cognitive architecture like Yi Zeng's BrainCog [ZZZ⁺22] (see Figure 17) could potentially be made to incorporate LLMs or other transformers in a number of interesting roles.

8.2.1 Bengio and Hu's RL/MDL Proposal

Deep learning pioneer Yoshua Bengio and his student Edward Hu have sketched a rough potential AGI architecture in which an LLM is coupled with a reinforcement learning module and a module that aims to find patterns using an "Occam's Razor" minimum-description-length (MDL) heuristic [BH23]. One rough way to interpret their proposed design is that: The LLM is used as a sort of raw pattern recognition/synthesis subsystem, and then the MDL based learner tries to find concise abstractions summarizing more diffuse sets of patterns that exist in the LLM. A reinforcement-learner or similar sort of subsystem guides the interaction of the LLM and the MDL-learner with the environment

There is a close relationship here with the work of Bengio and his team on generative flow models [BLD⁺21], which provide an interesting strategy for sampling from probability distributions over ways of getting things done, broader and arguably more AGI-ish than the "classical" reinforcement learning approach of searching for the best way to maximize some utility function.

This architecture aims to overcome the core limitations of LLMs via an appeal to traditional learning theory. Learning concise abstractions provides, in principle, a variety of much-needed ingredients, including

- a way of avoiding rampant hallucination (because more concise models of a dataset will generally produce less spurious output)
- a strategy for more robust generalization (via learning theory results correlating compactness of models with generalization ability)
- a route to more fundamental creativity (via ability to do "conceptual blends" at the level of actual conceptual abstractions)

. One might note that doing MDL-driven learning in a scalable way is a difficult problem which has proved mostly intractable so far; to be fair, though, no one has ever thrown GPT-4 scale resources at it. One core research question here is whether a huge, bloated

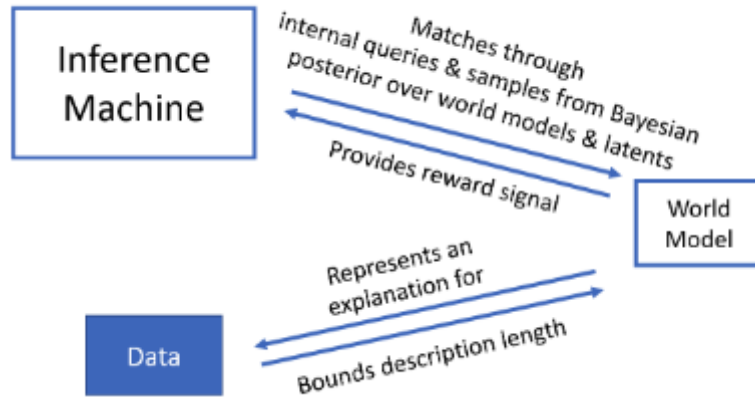


Figure 18: Bengio and Hu’s conceptual architecture for overcoming the shortcomings of LLMs via incorporating them into a framework involving richer modeling and inference [BH23]

compendium of highly particular data-patterns, such as exists in an LLM, can serve as helpful guidance to a separate learning-engine that aims to learn concise abstractions. Are generative flow networks a sufficiently flexible paradigm to encompass this sort of learning, or is their effectiveness restricted to cases where one has solid knowledge in advance of the class of distributions one is sampling from?

8.2.2 OpenCog Hyperon with an LLM-Like “Neural Space”

The OpenCog AGI framework is centered on a distributed, decentralized meta-representational fabric in the form of a weighted, typed metagraph called the Atomspace. Multiple cognitive algorithms work together to solve problems and achieve system goals, referencing and updating the Atomspace and assisting each other as needed; in the new “Hyperon” version of OpenCog, the cognitive algorithms are themselves implemented as networks of Atoms (nodes/links). The algorithms utilized come from multiple AI paradigms, including attractor neural nets for attention spreading, probabilistic logical reasoning, evolutionary program learning, concept blending and others. Hyperon includes a novel AGI-oriented programming language, MeTTa, which can be interpreted or compiled directly into Atoms living in Atomspaces.

The most natural way to integrate LLMs into Hyperon is to create a “Neural Space” that interacts tightly with the ordinary metagraph-based Atomspaces. The Neural Space

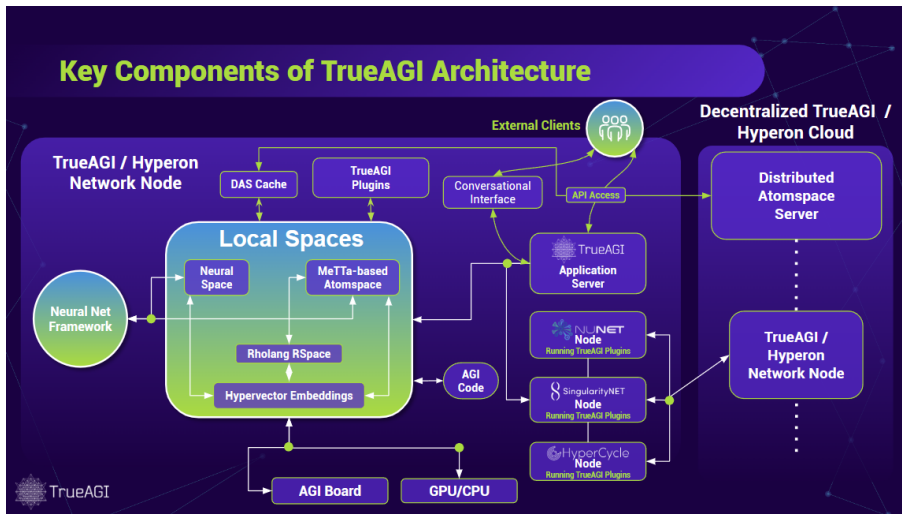


Figure 19: High level architecture of OpenCog Hyperon system, incorporating Neural Space that refers to external neural networks implemented in appropriate frameworks, including LLMs[ea23]

then plays an architectural role similar to a Distributed Atomspace (which serves as a backing store) or an Optimized Atomspace (which contains MeTTa code that's compiled for efficient execution); the pattern-matching queries that serve as the core of MeTTa programs can be applied against any of these auxiliary Spaces just as they can be applied against metagraph-based Atomspaces.

The core mechanism for integration of Neural Spaces with metagraph-based Atomspaces is automated conversion between natural language and MeTTa: this allows queries conceived by an LLM wrapped in a Neural Space to be submitted to ordinary Atomspaces or (DAS or Optimized Atomspaces), and vice-versa allows queries conceived in ordinary Atomspaces to be submitted to LLMs wrapped in Neural Spaces. Behind the scenes however there may be additional connections between the internals of the LLM and Atoms in the Atomspace.

As a toy example of such connections, in an Othello context it would make sense for direct links to grow between

- Atoms representing locations on the Othello board
- Atoms connoting combinations of activation values of neurons within the LLM that correspond to these same locations.

To make this sort of connection work maximally effectively between Atomspaces and

LLMs trained on natural language, we would like LLMs to learn semantic representations in a sufficiently clear way that, say, the LLM's internal representation of a concept like "cat" or a relationship like "communicate" can be identified via tractably discoverable mathematical functions of neurons within the LLM network.

There is also the option to use probabilistic logical reasoning in OpenCog's Atom-space to guide iterative retraining of an LLM. Text generated by an LLM can be translated into logical format, then processed using deductive, inductive and abductive reasoning in OpenCog ... and then the conclusions of this reasoning can be translated back into natural language and used as part of the training corpus to retrain the LLM ... and then the whole process can be repeated. This may get even more interesting if the OpenCog system is allowed to generate new concepts using, say, concept blending and evolutionary learning in the context of doing its inference.

The premise of this sort of integration is that LLMs are especially good at recognizing surface-level patterns in streams of data and synthesizing new patterns via combinations of these recognized patterns. However, LLMs are not being asked to serve as the medium for convergence and combination of all the different sorts of knowledge in the overall OpenCog system. The dynamics of the integrated system may nudge the LLM to further abstraction than it would do on its own, but even so the bulk of the abstract understanding of the system is done by other cognitive dynamics.

9 Social and Ethical Issues

Since the launch of ChatGPT brought LLMs so thoroughly into the public eye, a number of commentators both with and without technical background have raised alarms regarding potential ethical issues associated with LLMs. While every new technology has its risks, and some of the risks associated with LLMs are real and genuine regardless of the sizeable gap between LLMs and AGI, it is also important to avoid muddying thinking and excessive alarmism in managing the societal response to innovations.

9.1 LLMs and AGI-Related Existential Risk

There is an ongoing discussion regarding the potential existential risk to humanity of AGI with capability at the human level and beyond. Some futurists such as Nick Bostrom [Bos14] and Eliezer Yudkowsky [Yud01] [Yud15] argue that advanced AGI, unless architected in a very special way (the particulars of which are not clear to these theorists or anyone else), is very likely to extinguish humanity. Others, including one of the authors [GP12] [MG19] [Goe15b] [Goe16] have argued against this perspective,

and suggested that if AGIs are created with reasoning-friendly, reflective cognitive architectures and broad-benefit-oriented goal systems and educated in compassionate and ethics-friendly ways, they are likely to be beneficial to humanity.

Whatever view one holds on the risks to be run once HCAGI has been achieved, however, it seems dubious to raise these risks in the context of present-day LLMs, because in terms of capability and functionality these are quite far from HCAGI. If a future software system achieves HCAGI via combining LLMs with other AI techniques, the overall character of this system is likely to be quite different from that of a pure LLM, and the appropriate manner for dealing with its ethical strengths and weaknesses will likely be different as well.

9.2 Direct Potential of LLMs to Cause Economic and Social Disruption

HCAGI aside, LLMs are capable enough that they appear likely to have substantial economic and social impacts in spite of, and in some respects because of, their cognitive limitations. As with any complex technology, these impacts are sure to have both pluses and minuses. A variety of useful software services will no doubt be created using LLMs, via integrating them with other AI and non-AI tools in application-specific software frameworks. Many people will also likely find their employment disrupted or even eliminated via adoption of LLM-based technologies in their industries. The propensity of LLMs to provide false information, or to provide true and potentially dangerous information in an indiscriminate way, has obvious downsides.

In April 2023 a number of prominent individuals in the AI field, technology industry and beyond signed a petition arguing for an industry-wide 6 month pause in training larger and larger LLMs, so as to give some time for careful society-wide consideration of the best ways to leverage and integrate these tools. The two authors of this article found themselves on different sides regarding this pause, one signing the petition and the other speaking out in opposition. We also hold somewhat different perspectives on the probability that LLMs end up being highly useful as significant components of HCAGI systems.

Regardless of differences on these points, however, we concur on the clarity of certain serious cognitive limitations of LLM-type architectures, which have been our focus here. We also share a common view that many of the negative implications of LLMs are directly tied to these cognitive limitations. There is a direct link between LLMs' difficulty distinguishing truth from falsehood and their likely strong role in increasing misinformation and confusion. The tendency toward banality and un-creativity in LLM

outputs has clear probable implications for the aesthetic quality of the illustrations, music, writing, video, etc. that will be thrust upon us. Many of the negative implications of LLMs would be significantly remedied if LLMs were supplanted by AI systems moving further toward HCAGI in the respects we've outlined here, assuming this was done in a thoughtful and ethical way.

10 Coda

Having reviewed a fair number of particulars regarding the strengths and weaknesses of CILLMs, how they stack up against conceptions of AGI at the human level and beyond, and what their prospects for future expansion and improvement look like – how would we sum this all up?

One way to crystallize our perspective on Generative AI vs. AGI is via the old joke ⁵:

- ”To be is to do” – Socrates
- ”To do is to be” – Jean-Paul Sartre.
- ”Do be do be do” – Frank Sinatra.

(i.e., your reward for making it to the end of this long paper is some dubious bathroom-wall humor!).

The application of this sage wisdom to AGI works as follows:

- ”To be”, i.e. what a mind IS, has to do with its knowledge representation, how it represents various sorts of knowledge internally (including procedural and episodic and goal-related and reflective meta knowledge, along with fact/belief style declarative knowledge and sensorimotor knowledge)
- ”To do”, i.e. what a mind DOES, is e.g. how it answers various questions, synthesizes creative products, searches for music, navigates web pages, controls robots, etc.

The music comes when these two aspects are appropriately interleaved!

The core issue with CILLMs like ChatGPT is: They represent knowledge in a very crude way, basically as a huge repository of special cases (with judicious context-sensitive weights telling you which previously experienced special case is relevant to which situation).

⁵<https://quoteinvestigator.com/2013/09/16/do-be-do/>

This is what they ARE, and given this limitation, what they DO is always going to be a judicious recycling of what they have seen before... no amount of pasting other stuff on top of them is going to change that.

The transformer architecture IS in some cases capable of learning more abstract knowledge representations, as a fancy emergent side-effect of their judicious weighting of experience-catalogues, as OthelloGPT and other examples show. But this is not the transformer neural net's superpower, it needs quite a lot of repetitive data to do this at all well. Whereas for the human brain (or other more AGI-oriented software systems if they end up working as hoped), storing and maintaining abstract knowledge representations is closer to the core of what they ARE – which is core to what makes them more able to DO general intelligence (which largely involves learning and then enacting/synthesizing-from these abstracted representations)...

A hybrid architecture CAN in our view work for achieving HCAGI and beyond using LLMs as a significant component, but we have explained in the above why we don't think LLMs have the right knowledge representation to serve as a "central hub" for various participants in the hybrid, and nor do any other available narrow-AI tools one might use as CILLM plugins. ⁶

Mithen's "The Prehistory of Mind" [Mit96] argues that modern-human-like GI evolved when linguistic brain and tool-making brain started working together synergetically. If this appealing story is correct, this synergy was able to arise because the evolving cortex was able to emerge different sorts of representation than were useful for language or tool-making alone.... Human cortex is much more powerfully meta-representational than LLMs are (more able to self-organize new representations based on experience), which is a key point even though LLMs do have some fascinating and non-trivial ability to emerge new representations.

Just glomming together LLMs with their language-y representations, to Wolfram Alpha or something similar with its logic-y representation and say an RL system with a tool-build-y representation – this will NOT intrinsically result in a representation combining the virtues of all these ingredients. The combination needs to take place in a sufficiently flexible infrastructure like the human cortex or a simulation thereof, or perhaps OpenCog Hyperon's Atomspace or other AGI innovations.

Do, be, do, be, do, to AGI and beyond!

⁶As an example, Wolfram Alpha is excellent in many ways but it doesn't really do commonsense reasoning effectively in spite of all its knowledge, and it's certainly not terribly creative (it doesn't do concept formation, hypothesis/conjecture formation, etc.)

Acknowledgements

Thanks are due to Gary Marcus: the impetus to put this article together arose directly from our conversations on the topic. Also to Deborah Duong and George Musser, whose comments and questions on an earlier version of the article spurred some additions and improvements. Finally to the whole SingularityNET and Hyperon AGI R&D teams, whose depth of thought on AGI and LLM topics has infused my own thinking in more ways than I'm able to enumerate.

References

- [ACZ⁺22] Tyson Aflalo, Srinivas Chivukula, Carey Zhang, Emily R Rosario, Nader Pouratian, and Richard A Andersen. Cognition through internal models: Mirror neurons as one manifestation of a broader mechanism. *BioRxiv*, pages 2022–09, 2022.
- [AM23] Amos Azaria and Tom Mitchell. The internal state of an llm knows when its lying. *arXiv preprint arXiv:2304.13734*, 2023.
- [BB21] Giuseppe Barbiero and Rita Berto. Biophilia as evolutionary adaptation: An onto-and phylogenetic framework for biophilic design. *Frontiers in psychology*, 12:700709, 2021.
- [BCE⁺23] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.
- [BCL⁺23] Yejin Bang, Samuel Cahyawijaya, Nayeon Lee, Wenliang Dai, Dan Su, Bryan Wilie, Holy Lovenia, Ziwei Ji, Tiezheng Yu, Willy Chung, et al. A multitask, multilingual, multimodal evaluation of chatgpt on reasoning, hallucination, and interactivity. *arXiv preprint arXiv:2302.04023*, 2023.
- [Ben17] Yoshua Bengio. The consciousness prior. *arXiv preprint arXiv:1709.08568*, 2017.
- [BH05] Sebastian Bader and Pascal Hitzler. Dimensions of neural-symbolic integration - a structured survey. In S. Artemov, H. Barringer, A. S. d'Avila Garcez, L. C. Lamb, and J. Woods., editors, *We Will Show Them: Essays*

- in Honour of Dov Gabbay*, volume 1, pages 167–194. College Publications, 2005.
- [BH23] Yoshua Bengio and Edward Hu. Scaling in the service of reasoning and model-based ml. 2023. <https://yoshuabengio.org/2023/03/21/scaling-in-the-service-of-reasoning-model-based-ml/>.
- [BHD⁺22] Chandra Bhagavatula, Jena D Hwang, Doug Downey, Ronan Le Bras, Ximing Lu, Keisuke Sakaguchi, Swabha Swayamdipta, Peter West, and Yejin Choi. I2d2: Inductive knowledge distillation with neurologic and self-imitation. *arXiv preprint arXiv:2212.09246*, 2022.
- [BLD⁺21] Yoshua Bengio, Salem Lahlou, Tristan Deleu, Edward J Hu, Mo Tiwari, and Emmanuel Bengio. Gflownet foundations. *arXiv preprint arXiv:2111.09266*, 2021.
- [BMPR10] Mark A Bedau, John S McCaskill, Norman H Packard, and Steen Rasmussen. Living technology: Exploiting life’s principles in technology. *Artificial Life*, 16(1):89–97, 2010.
- [Bod08] Margaret A Boden. *Mind as machine: A history of cognitive science*. Oxford University Press, 2008.
- [Bos14] Nick Bostrom. *Superintelligence: Paths, strategies, dangers*, 2014.
- [Bri23] Selmer Bringsjord. The m cognitive meta-architecture as touchstone for standard modeling of agi-level minds. *AGI-23, Stockholm*, 2023. <https://www.youtube.com/watch?v=JzBMN0vnb-A>, time 7:50.
- [DOP08] Wlodzislaw Duch, Richard Oentaryo, and Michel Pasquier. Cognitive architectures: Where do we go from here? *Proc. of the Second Conf. on AGI*, 2008.
- [ea23] Ben Goertzel et al. Opencog hyperon. 2023. <https://hyperon.opencog.org/>.
- [FAS21] Arthur Franz, Oleksandr Antonenko, and Roman Soletskyi. A theory of incremental compression. *Information Sciences*, 547:28–48, 2021.
- [FG96] Stan Franklin and Art Graesser. Is it an agent, or just a program?: A taxonomy for autonomous agents. In *International workshop on agent theories, architectures, and languages*, pages 21–35. Springer, 1996.

- [FGL19] Arthur Franz, Victoria Gogulya, and Michael Löffler. William: A monolithic approach to agi. In *International Conference on Artificial General Intelligence*, pages 44–58. Springer, 2019.
- [Fra21] Arthur Franz. Experiments on the generalization of machine learning algorithms. In *International Conference on Artificial General Intelligence*, pages 75–85. Springer, 2021.
- [FT02] G. Fauconnier and M. Turner. *The Way We Think: Conceptual Blending and the Mind’s Hidden Complexities*. Basic, 2002.
- [Gar99] H Gardner. *Intelligence reframed: Multiple intelligences for the 21st century*. Basic, 1999.
- [GFH⁺23] Izzeddin Gur, Hiroki Furuta, Austin Huang, Mustafa Safdari, Yutaka Matsuo, Douglas Eck, and Aleksandra Faust. A real-world webagent with planning, long context understanding, and program synthesis. *arXiv preprint arXiv:2307.12856*, 2023.
- [Goe01] Ben Goertzel. *Creating Internet Intelligence*. Plenum Press, 2001.
- [Goe09] Ben Goertzel. The embodied communication prior: A characterization of general intelligence in the context of embodied social interaction. In *2009 8th IEEE International Conference on Cognitive Informatics*, pages 38–43. IEEE, 2009.
- [Goe13] Ben Goertzel. A mind-world correspondence principle. In *2013 IEEE Symposium on Computational Intelligence for Human-like Intelligence (CIHLI)*, pages 68–73. IEEE, 2013.
- [Goe14] Ben Goertzel. *The AGI Revolution*. Amazon, 2014.
- [Goe15a] Ben Goertzel. Artificial general intelligence. *Scholarpedia*, 2015.
- [Goe15b] Ben Goertzel. Superintelligence: Fears, promises and potentials: Reflections on bostrom’s superintelligence, yudkowsky’s from ai to zombies, and weaver and veitas’s "open-ended intelligence". *Journal of Ethics and Emerging Technologies*, 25(2):55–87, 2015.
- [Goe16] Ben Goertzel. Infusing advanced agis with human-like value systems: Two theses. *Journal of Ethics and Emerging Technologies*, 26(1):50–72, 2016.

- [Goe17] Ben Goertzel. Toward a formal model of cognitive synergy. *CoRR*, abs/1703.04361, 2017.
- [Goe21] Ben Goertzel. Toward a general theory of general intelligence: A patternist perspective. *arXiv preprint arXiv:2103.15100*, 2021.
- [Goe23a] Ben Goertzel. Ai now predicts human ethical judgments quite well, 2023. <https://magazine.mindplex.ai/ai-now-predicts-human-ethical-judgments-quite-well/>.
- [Goe23b] Ben Goertzel. Evil llm: Chatgpt excels at emulating anti-morality and ethical fakery, 2023. <https://magazine.mindplex.ai/evil-llm-chatgpt-excels-at-emulating-anti-morality-and-ethical-fakery/>.
- [Gol00] David Goldberg. *The Design of Innovation*. Addison-Wesley, 2000.
- [GP05] Ben Goertzel and Cassio Pennachin. *Artificial General Intelligence*. Springer, 2005.
- [GP12] Ben Goertzel and Joel Pitt. Nine ways to bias open-source agi toward friendliness. *Journal of Evolution and Technology* 22-1, 2012.
- [Gub97] Mark Avrum Gubrud. Nanotechnology and international security. In *Fifth Foresight Conference on Molecular Nanotechnology*, volume 1, 1997.
- [GWR⁺16] Alex Graves, Greg Wayne, Malcolm Reynolds, Tim Harley, Ivo Danihelka, Agnieszka Grabska-Barwińska, Sergio Gómez Colmenarejo, Edward Grefenstette, Tiago Ramalho, John Agapiou, et al. Hybrid computing using a neural network with dynamic external memory. *Nature*, 538(7626):471–476, 2016.
- [Hah20] Michael Hahn. Theoretical limitations of self-attention in neural sequence models. *Transactions of the Association for Computational Linguistics*, 8:156–171, 2020.
- [Hut05] Marcus Hutter. *Universal Artificial Intelligence: Sequential Decisions based on Algorithmic Probability*. Springer, 2005.
- [IE08] Eugene M Izhikevich and Gerald M Edelman. Large-scale model of mammalian thalamocortical systems. *Proceedings of the national academy of sciences*, 105(9):3593–3598, 2008.

- [JL83] Philip Nicholas Johnson-Laird. *Mental models: Towards a cognitive science of language, inference, and consciousness*. Number 6. Harvard University Press, 1983.
- [JSH23] Tanishq Matthew Abraham Juergen Schmidhuber and Jeremy Howard, 2023. <https://twitter.com/SchmidhuberAI/status/1683870175299239937>.
- [Kos23] Michal Kosinski. Theory of mind may have spontaneously emerged in large language models. *arXiv preprint arXiv:2302.02083*, 2023.
- [Koz92] J. R. Koza. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. MIT Press, 1992.
- [Lai19] John E Laird. *The Soar cognitive architecture*. MIT press, 2019.
- [LeC23] Yann LeCun, 2023. <https://twitter.com/ylecun/status/1621805604900585472>.
- [LFL⁺23] Zhan Ling, Yunhao Fang, Xuanlin Li, Zhiao Huang, Mingu Lee, Roland Memisevic, and Hao Su. Deductive verification of chain-of-thought reasoning. *arXiv preprint arXiv:2306.03872*, 2023.
- [LGV⁺17] Ruiting Lian, Ben Goertzel, Linas Vepstas, David Hanson, and Changle Zhou. Symbol grounding via chaining of morphisms. *arXiv preprint arXiv:1703.04368*, 2017.
- [LH07a] Shane Legg and Marcus Hutter. A collection of definitions of intelligence. IOS, 2007.
- [LH07b] Shane Legg and Marcus Hutter. A definition of machine intelligence. *Minds and Machines*, 17, 2007.
- [LHB⁺22] Kenneth Li, Aspen K Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister, and Martin Wattenberg. Emergent world representations: Exploring a sequence model trained on a synthetic task. *arXiv preprint arXiv:2210.13382*, 2022.
- [Lim23] Russel Lim. Gpt-4 is amazing but still struggles at high school math competitions. *Medium*, 2023. <https://www.cantorsparadise.com/gpt-4-is-amazing-but-still-struggles-at-high-school-math-competitions-cbc2e73738>

- [LLR17] John E Laird, Christian Lebiere, and Paul S Rosenbloom. A standard model of the mind: Toward a common computational framework across artificial intelligence, cognitive science, neuroscience, and robotics. *Ai Magazine*, 38(4):13–26, 2017.
- [LPC⁺23] Pan Lu, Baolin Peng, Hao Cheng, Michel Galley, Kai-Wei Chang, Ying Nian Wu, Song-Chun Zhu, and Jianfeng Gao. Chameleon: Plug-and-play compositional reasoning with large language models. *arXiv preprint arXiv:2304.09842*, 2023.
- [LTAE021] James Lee-Thorp, Joshua Ainslie, Ilya Eckstein, and Santiago Ontanon. Fnet: Mixing tokens with fourier transforms. *arXiv preprint arXiv:2105.03824*, 2021.
- [LV⁺08] Ming Li, Paul Vitányi, et al. *An introduction to Kolmogorov complexity and its applications*, volume 3. Springer, 2008.
- [MBF11] Tamas Madl, Bernard J Baars, and Stan Franklin. The timing of the cognitive cycle. *PloS one*, 6(4):e14803, 2011.
- [MBH⁺19] Marjorie McShane, Selmer Bringsjord, James Hendler, Sergei Nirenburg, and Ron Sun. A response to núñez et al.’s (2019)?what happened to cognitive science?? *Topics in Cognitive Science*, 11(4):914–917, 2019.
- [MG18] George Mangun Michael Gazzaniga, Richard Ivry. *Cognitive neuroscience: the biology of the mind*. Norton, 2018.
- [MG19] Gabriel Axel Montes and Ben Goertzel. Distributed, decentralized, and democratized artificial intelligence. *Technological Forecasting and Social Change*, 141:354–358, 2019.
- [MHR⁺20] James L McClelland, Felix Hill, Maja Rudolph, Jason Baldridge, and Hinrich Schütze. Placing language in an integrated understanding system: Next steps toward human-level performance in neural language models. *Proceedings of the National Academy of Sciences*, 117(42):25966–25974, 2020.
- [Mit96] Steven Mithen. *The Prehistory of Mind*. Thames and Hudson, 1996.
- [MLC⁺23] Nick McKenna, Tianyi Li, Liang Cheng, Mohammad Javad Hosseini, Mark Johnson, and Mark Steedman. Sources of hallucination by large

- language models on inference tasks. *arXiv preprint arXiv:2305.14552*, 2023.
- [msr22] msravi. Chatgpt produces made-up nonexistent references, 2022. <https://news.ycombinator.com/item?id=33841672>.
- [NAG⁺19] Rafael Núñez, Michael Allen, Richard Gao, Carson Miller Rigoli, Josephine Relaford-Doyle, and Arturs Semenuks. What happened to cognitive science? *Nature human behaviour*, 3(8):782–791, 2019.
- [Nan22] Neel Nanda. Actually, othello-gpt has a linear emergent world representation. *Less Wrong*, 2022. <https://www.lesswrong.com/posts/nmxzr2zsJNtjaHh7x/actually-othello-gpt-has-a-linear-emergent-world>.
- [New90] Alan Newell. *Unified Theories of Cognition*. Harvard University press, 1990.
- [Nil09] Nils Nilsson. The physical symbol system hypothesis: Status and prospects. *50 Years of AI, Festschrift, LNAI 4850*, 33, 2009.
- [Pei91] Charles Sanders Peirce. *Peirce on signs: Writings on semiotic*. UNC Press Books, 1991.
- [Per23] Sidney Perkowitz. What does chatgpt know about science? *Nautilus*, 2023. <https://nautil.us/what-does-chatgpt-know-about-science-291213/>.
- [Pik23] Matúš Pikuliak. Chatgpt survey: Performance on nlp datasets, 03 2023. https://www.opensamizdat.com/posts/chatgpt_survey.
- [PLW⁺23] Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Jiapu Wang, and Xindong Wu. Unifying large language models and knowledge graphs: A roadmap. *arXiv preprint arXiv:2306.08302*, 2023.
- [PP21] Roma Patel and Ellie Pavlick. Mapping language models to grounded conceptual spaces. In *International Conference on Learning Representations*, 2021.
- [RNS⁺18] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.

- [Rud94] Daniel L. Ruderman. The statistics of natural images. *Network: computation in neural systems*, 5(4):517, 1994.
- [Sam10] Alexei V. Samsonovich. Toward a unified catalog of implemented cognitive architectures. In *BICA*, pages 195–244, 2010.
- [Sch06] J. Schmidhuber. Godel machines: Fully Self-referential Optimal Universal Self-improvers. In B. Goertzel and C. Pennachin, editors, *Artificial General Intelligence*, pages 119–226. 2006.
- [Sch23] Dale Schuurmans. Memory augmented large language models are computationally universal. *arXiv preprint arXiv:2301.04589*, 2023.
- [SH22] Abulhair Saparov and He He. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. *arXiv preprint arXiv:2210.01240*, 2022.
- [SLFC22] Maarten Sap, Ronan LeBras, Daniel Fried, and Yejin Choi. Neural theory-of-mind? on the limits of social intelligence in large lms. *arXiv preprint arXiv:2210.13312*, 2022.
- [Sol64] Ray Solomonoff. *A Formal Theory of Inductive Inference, Part I*. Information and Control, 1964.
- [Spe61] Charles Spearman. "general intelligence" objectively determined and measured. 1961.
- [SXT⁺23] Jiashuo Sun, Chengjin Xu, Lumingyuan Tang, Saizhuo Wang, Chen Lin, Yeyun Gong, Heung-Yeung Shum, and Jian Guo. Think-on-graph: Deep and responsible reasoning of large language model with knowledge graph. *arXiv preprint arXiv:2307.07697*, 2023.
- [Ter48] Lewis Madison Terman. The measurement of intelligence, 1916. 1948.
- [TH12] Kristinn Thórisson and Helgi Helgasson. Cognitive architectures and autonomy: A comparative review. *Journal of Artificial General Intelligence*, 3(2):1, 2012.
- [Tom03] Michael Tomasello. *Constructing a Language: A Usage-Based Theory of Language Acquisition*. 2003.

- [Tre23] Marton Trencseni. Testing gpt-4 spatial reasoning and comprehension. 2023. <https://bytepaw.com/testing-gpt-4-spatial-reasoning-and-comprehension.html>.
- [TSW⁺05] Richard P Taylor, B Spehar, JA Wise, CWG Clifford, BR Newell, and TP Martin. Perceptual and physiological responses to the visual complexity of pollock's dripped fractal patterns. *Journal of Non-linear Dynamics, Psychology and Life Sciences*, 2005.
- [Tur50] Alan Turing. Computing machinery and intelligence. *Mind*, 59, 1950.
- [VDD⁺18] Athanasios Voulodimos, Nikolaos Doulamis, Anastasios Doulamis, Eftychios Protopapadakis, et al. Deep learning for computer vision: A brief review. *Computational intelligence and neuroscience*, 2018, 2018.
- [VNH⁺11] Joel Veness, Kee Siong Ng, Marcus Hutter, William Uther, and David Silver. A Monte-Carlo AIXI approximation. *Journal of Artificial Intelligence Research*, 40:95–142, 2011.
- [VO23] Pranshu Verma and Will Oremus. Chatgpt invented a sexual harassment scandal and named a real law prof as the accused, 2023. <https://www.washingtonpost.com/technology/2023/04/05/chatgpt-lies/>.
- [VSP⁺17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [Wan06] Pei Wang. *Rigid Flexibility: The Logic of Intelligence*. Springer, 2006.
- [Whi22] Mary Gormandy White. Ethical dilemma examples, 2022.
- [WHL⁺23] Xiaoxuan Wang, Ziniu Hu, Pan Lu, Yanqiao Zhu, Jieyu Zhang, Satyen Subramaniam, Arjun R Loomba, Shichang Zhang, Yizhou Sun, and Wei Wang. Scibench: Evaluating college-level scientific problem-solving abilities of large language models. *arXiv preprint arXiv:2307.10635*, 2023.
- [WV17] David Weinbaum and Viktoras Veitas. Open ended intelligence: the individuation of intelligent agents. *Journal of Experimental & Theoretical Artificial Intelligence*, 29(2):371–396, 2017.

- [YCL⁺23] Linyao Yang, Hongyang Chen, Zhao Li, Xiao Ding, and Xindong Wu. Chatgpt is not enough: Enhancing large language models with knowledge graphs for fact-aware language modeling. *arXiv preprint arXiv:2306.11489*, 2023.
- [Yud01] Eliezer Yudkowsky. Creating friendly ai. *Singularity Institute for AI*, 2001. <http://singinst.org/upload/CFAI.html>.
- [Yud15] Eliezer Yudkowsky. Rationality: from ai to zombies. *Machine Intelligence Research Institute, Berkeley*, 2015.
- [YZWN22] Samuel Yang-Zhao, Tianyu Wang, and Kee Siong Ng. A direct approximation of aixi using logical state abstractions. *Advances in Neural Information Processing Systems*, 35:36640–36653, 2022.
- [ZCG⁺23] Wanjun Zhong, Ruixiang Cui, Yiduo Guo, Yaobo Liang, Shuai Lu, Yanlin Wang, Amin Saied, Weizhu Chen, and Nan Duan. Agieval: A human-centric benchmark for evaluating foundation models. *arXiv preprint arXiv:2304.06364*, 2023.
- [ZHL⁺19] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. Ernie: Enhanced language representation with informative entities. *arXiv preprint arXiv:1905.07129*, 2019.
- [ZWL⁺23] Aojun Zhou, Ke Wang, Zimu Lu, Weikang Shi, Sichun Luo, Zipeng Qin, Shaoqing Lu, Anya Jia, Linqi Song, Mingjie Zhan, et al. Solving challenging math word problems using gpt-4 code interpreter with code-based self-verification. *arXiv preprint arXiv:2308.07921*, 2023.
- [ZZZ⁺22] Yi Zeng, Dongcheng Zhao, Feifei Zhao, Guobin Shen, Yiting Dong, Enmeng Lu, Qian Zhang, Yinqian Sun, Qian Liang, Yuxuan Zhao, et al. Braincog: A spiking neural network based brain-inspired cognitive intelligence engine for brain-inspired ai and brain simulation. *arXiv preprint arXiv:2207.08533*, 2022.