

Imagining and building wise machines: The centrality of AI metacognition

Samuel G. B. Johnson^{1*}, Amir-Hossein Karimi², Yoshua Bengio³, Nick Chater⁴, Tobias Gerstenberg⁵, Kate Larson⁶, Sydney Levine⁷, Melanie Mitchell⁸, Iyad Rahwan⁹, Bernhard Schölkopf¹⁰, Igor Grossmann^{1*}

¹ University of Waterloo, Department of Psychology

² University of Waterloo, Department of Electrical and Computer Engineering

³ Université de Montréal, Department of Computer Science and Operations Research

⁴ Warwick Business School, Behavioural Science Group

⁵ Stanford University, Department of Psychology

⁶ University of Waterloo, Cheriton School of Computer Science

⁷ Allen Institute for Artificial Intelligence

⁸ Santa Fe Institute

⁹ Max Planck Institute for Human Development

¹⁰ Max Planck Institute for Intelligent Systems

* **Correspondence to:** Sam Johnson (samuel.johnson@uwaterloo.ca) or Igor Grossmann (igrossma@uwaterloo.ca)

Abstract

Recent advances in artificial intelligence (AI) have produced systems capable of increasingly sophisticated performance on cognitive tasks. However, AI systems still struggle in critical ways: unpredictable and novel environments (robustness), lack transparency in their reasoning (explainability), face challenges in communication and commitment (cooperation), and pose risks due to potential harmful actions (safety). We argue that these shortcomings stem from one overarching failure: AI systems lack wisdom. Drawing from cognitive and social sciences, we define wisdom as the ability to navigate intractable problems—those that are ambiguous, radically uncertain, novel, chaotic, or computationally explosive—through effective task-level and metacognitive strategies. While AI research has focused on task-level strategies, metacognition—the ability to reflect on and regulate one’s thought processes—is underdeveloped in AI systems. In humans, metacognitive strategies such as recognizing the limits of one’s knowledge, considering diverse perspectives, and adapting to context are essential for wise decision-making. We propose that integrating metacognitive capabilities into AI systems is crucial for enhancing their robustness, explainability, cooperation, and safety. By focusing on developing wise AI, we suggest an alternative to aligning AI with specific human values—a task fraught with conceptual and practical difficulties. Instead, wise AI systems can thoughtfully navigate complex situations, account for diverse human values, and avoid harmful actions. We discuss potential approaches to building wise AI, including benchmarking metacognitive abilities and training AI systems to employ wise reasoning. Prioritizing metacognition in AI research will lead to systems that act not only intelligently but also wisely in complex, real-world situations.

Imagining and building wise machines: The centrality of AI metacognition

1. Introduction

Breakthroughs in machine cognition have recently come in quick succession. Generative AI (GenAI) systems—such as ChatGPT, Strawberry, LLaMA, and Gemini—can summarize medical literature and cases (Cascella et al., 2024; Van Veen et al., 2024), identify case law and relevant legal precedents (Shu et al., 2024), and solve math Olympiad problems (Gao et al., 2024). Such advances have led many to ask whether AI systems will soon be able to perform *any* cognitive task at a human or superhuman level.

Despite these accomplishments, AI systems still struggle in a variety of ways, limiting their capabilities while portending new dangers. They struggle in novel and unpredictable environments that extend beyond their training data—they lack *robustness*. Their computational processes are opaque, creating a problem of *explainability* (Dwivedi et al., 2023). Their challenges with communication and inability to commit credibly to long-term plans create barriers to *cooperation* (Dafoe et al., 2020). And most worryingly of all, these shortcomings challenge our ability to harness the upside of AI while avoiding risks and ensuring its *safety* (Ji et al., 2023). Each of these problems will be exacerbated as AIs come to act as agents in the world.

Here, we argue that AI systems lack a key capability that underlies all these deficiencies: they are not *wise*. We draw on research in the cognitive and social sciences to understand what *machine wisdom* could be, evaluate why it might be desirable, and sketch potential routes to achieving it.

Despite a millennia-long philosophical pedigree, it is only recently that social and cognitive scientists have begun to reach a consensus about what constitutes wisdom (e.g., Grossmann et al., 2020). Though wisdom can mean many things, for this Perspective we define wisdom functionally as the ability to successfully navigate *intractable problems*—those that do not lend themselves to analytic techniques due to unlearnable probability distributions or incommensurable values. Mechanistically, this is achieved through two types of strategies: (1) *Task-level strategies* are used to manage the problem itself (e.g., simple rules-of-thumb); and (2) *Metacognitive strategies* are used to flexibly manage those task-level strategies (e.g., understanding the limits of one’s knowledge and integrating multiple perspectives).

Whereas some task-level strategies have long been topics of investigation in AI research, metacognition research is much less well-developed. This is a crucial shortcoming because metacognition is the control system that allows us to decide between conflicting task-level strategies. By analogy to the pivotal role of human metacognition in wise decision-making, better machine metacognition will permit AI systems to prevent overconfident inferences and avoid harmful actions. This in turn will improve the robustness of such systems to novel situations, their explainability to users, their ability

to cooperate with human and AI agents, and their safety in avoiding both prosaic and catastrophic failure modes.

2. What is wisdom?

At first blush, in the cognitive and social sciences, the concept of ‘wisdom’ seems to bring together many superficially unrelated characteristics. Consider the following examples of human wisdom:

- Willa’s children are bitterly arguing about money. Willa draws on her life experience to show them why they should instead compromise in the short term and prioritize their sibling relationship in the long term.
- Daphne is a world-class cardiologist. Nonetheless, she consults with a much more junior colleague when she recognizes that the colleague knows more about a patient’s history than she does.
- Ron is a political consultant who formulates possible scenarios to ensure his candidate will win. To help generate scenarios, he not only imagines best case scenarios, but also imagines that his client has lost the election and considers possible reasons that might have contributed to the loss.

Life experience, intellectual humility, and scenario planning do not seem to share much in common beyond all being positive attributes. But being able to solve tricky integrals, crack clever jokes, and compose beautiful sonnets are also positive attributes—yet these don't constitute wisdom.

How can we understand human wisdom, and to what extent has prior work in AI already laid the groundwork for wise AI?

2.1. Human wisdom

Given wisdom’s disparate characteristics, a range of accounts have been proposed (Table 1 summarizes some prominent ones). Like many concepts, there is probably no set of necessary and sufficient conditions, yet there is a common core to many of the characteristics highlighted in Table 1. We draw two important generalizations: Functionally, wisdom facilitates thought and action in intractable situations. Mechanistically, wisdom is implemented through both task-level strategies and metacognitive abilities for weighing and implementing those strategies.

Theory/Model	Elements of Wisdom
Component Theories	
Balance Theory (Sternberg, 1998)	Deploying knowledge and skills to achieve the common good by: <ul style="list-style-type: none"> - Balancing interests (their own, others', and society's) - Balancing time perspectives (long-term and short-term) - Deploying positive ethical values - Managing environments (adapting to, selecting, or altering)
Berlin Wisdom Model (Baltes & Smith, 2008)	Expertise in important and difficult matters of life: <ul style="list-style-type: none"> - Factual knowledge (about human nature and life) - Procedural knowledge (strategies to address life challenges) - Contextualism (strategies account for social context) - Value relativism (strategies account for variation in values) - Managing uncertainty (strategies change with circumstances)
MORE Life Experience Model (Glück & Bluck, 2013)	Gaining psychological resources via reflection, to cope with life challenges: <ul style="list-style-type: none"> - Uncertainty management (coping with uncertainty, uncontrollability) - Openness (to new experiences and perspectives) - Reflectivity (about life experiences) - Emotion regulation (management of and sensitivity to emotions)
Three-Dimensional Model (Ardelt, 2004)	Acquiring and reflecting on life experience to cultivate personality traits: <ul style="list-style-type: none"> - Cognitive (curiosity about life; recognizing uncertainty, ignorance) - Emotional (sympathy and compassion; valuing others) - Reflective (perspective-taking; questioning one's beliefs)
Wise Reasoning Model (Grossmann, 2017)	Using context-sensitive reasoning to manage important social challenges: <ul style="list-style-type: none"> - Intellectual humility (knowledge of one's epistemic limits) - Perspective-taking (actively seeking out others' viewpoints) - Perspective integration (accounting for multiple perspectives) - Flexibility (recognizing uncertainty and change)
Consensus Models	
Common Wisdom Model (Grossmann et al., 2020)	A style of social-cognitive processing that is: <ul style="list-style-type: none"> - Morally grounded <ul style="list-style-type: none"> o Balancing interests of the self and others o Pursuing truth o Oriented toward the common good - Metacognitively sound <ul style="list-style-type: none"> o Considering context o Taking multiple perspectives o Accounting for short- and long-term effects o Thinking reflectively o Aware of the limits of one's knowledge
Integrative Model (Glück & Weststrate, 2022)	A behavioral repertoire in which: <ul style="list-style-type: none"> - A complex and uncertain situation arises, evoking an appropriate emotional and motivational state <ul style="list-style-type: none"> o Open-mindedness, care for others, calm emotions - Depending on traits and skills <ul style="list-style-type: none"> o Exploratory orientation, concern for others, emotion regulation - Facilitating deployment of cognitive resources <ul style="list-style-type: none"> o Life knowledge, metacognition, reflection - Using these resources to deploy effective metacognitive strategies <ul style="list-style-type: none"> o Reasoning is contextualized, balanced, multi-perspectival

Table 1. Psychological approaches to wisdom. The five “component theories” are a selected set of psychological theories or models of wisdom. The two “consensus models” are attempts to identify common themes and processes among those theories. For a more detailed review, see Glück and Weststrate (2022).

2.1.1. The function of wisdom: Navigating intractable situations

If life were a series of textbook problems, we would not need to be wise. There would be a correct answer, the requisite information for calculating it would be available, and natural selection would have ruthlessly driven humans to find those answers. We would be nothing more or less than master statisticians, merciless optimizers, lightning calculators. Indeed, in some domains—like low-level processing of the visual world—not only humans, but squirrels, goldfish, and bees come remarkably close to this description (Mascalzoni & Regolin, 2011).

Yet in other domains, problems such as social interaction and decision-making in an unstructured, uncertain, and rapidly evolving world require further tools beyond statistics and calculation (Johnson, Bilovich, & Tuckett, 2023). They are often *intractable* in one or more ways:

- *Incommensurable*. It features ambiguous goals or values that cannot be compared with one another (Walasek & Brown, 2023).
- *Transformative*. The outcome of the decision might change one's preferences, so that there is a clash between one's present and future values (Paul, 2013).
- *Radically uncertain*. We might not be able to exhaustively list the possible outcomes or assign probabilities to them in a principled way (Kay & King, 2020).
- *Chaotic*. The data-generating process may have a strong nonlinearity or dependency on initial conditions, making it fundamentally unpredictable (Lorenz, 1993).
- *Non-stationary*. The underlying process may be changing over time, making the probability distribution unlearnable.
- *Out-of-distribution*. The situation is novel, going beyond one's experience or available data.
- *Computationally explosive*. The optimal response could be calculated with infinite or infeasibly large computational resources, but this is not possible due to resource constraints.

Our earlier examples of situations calling for wisdom each featured one or more of these forms of intractability. Wisdom helped Willa understand how to make an incommensurable trade-off, helped Daphne to navigate her ignorance in an out-of-distribution situation, and helped Ron to make useful forecasts despite his ignorance about the probability distributions governing the radically uncertain future.

As AI systems increasingly come to act in the world—encountering ambiguous situations, ill-specified goals, and unpredictable environments, such as those riddled with fickle humans—their decision-making environment will become increasingly like our own: increasingly intractable. Engineering wise capabilities, enabling AI to more effectively navigate such intractable situations, will accordingly become ever more urgent.

2.1.2. Mechanisms of wisdom: Metacognitive strategy selection

How do wise people navigate intractable situations? Their secret is both to have a catalogue of effective *task-level* strategies (which we use to manage a concrete situation) and general *metacognitive* strategies (which manage the knowledge required to deploy task-level strategies and resolve conflicts between them) (Figure 1).

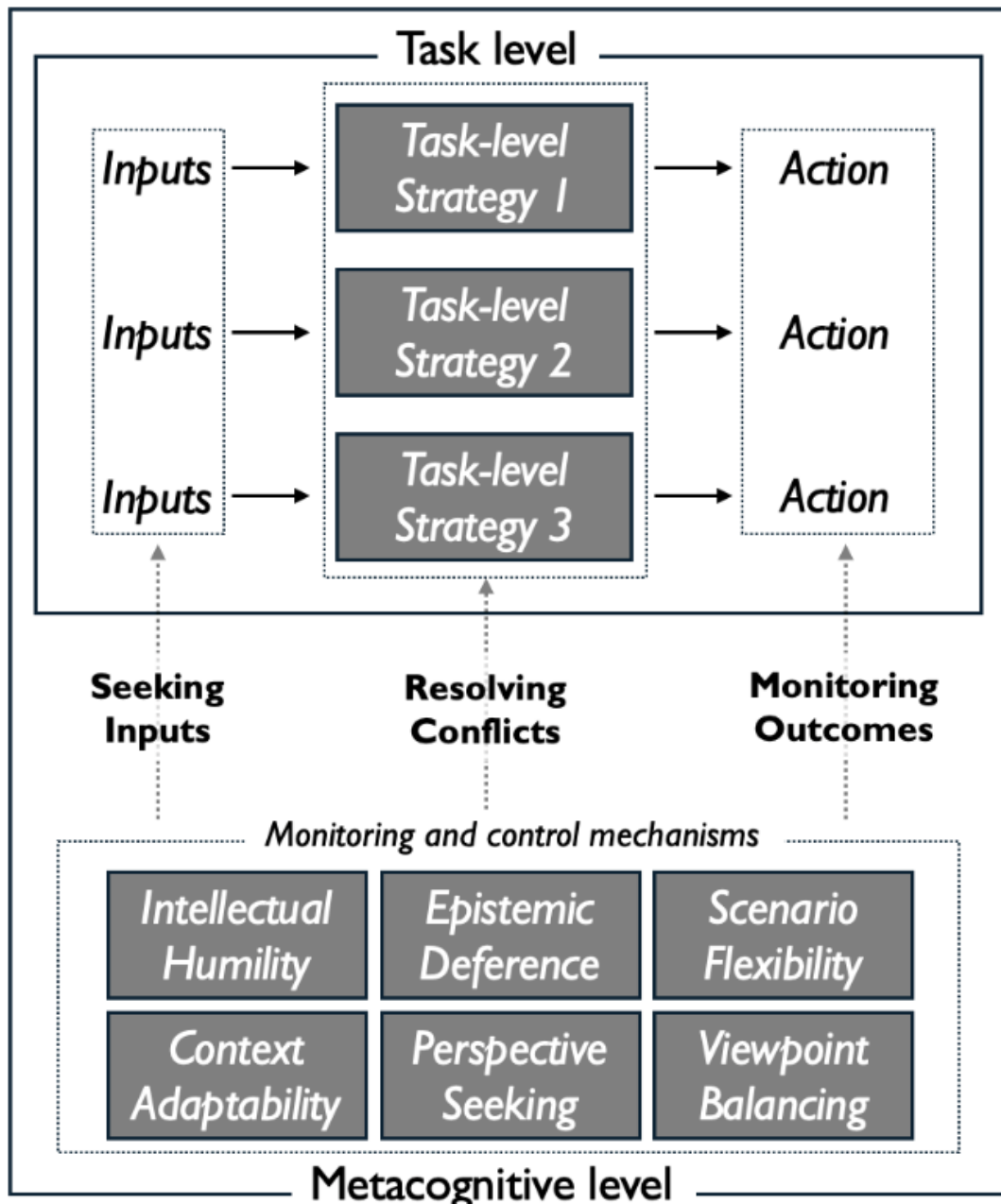


Figure 1. The relationship between task-level and metacognitive strategies in wise reasoning. Task-level strategies (e.g., heuristics, narratives, analytical procedures) provide candidate actions for a given situation. Metacognitive monitoring and control processes regulate these strategies in three ways: obtaining the appropriate inputs, deciding which strategy to use when they conflict, and monitoring their outcomes to avoid catastrophic actions.

Task-level strategies often take the form of *heuristics*—rules of thumb which rely on a small number of inputs and do not attempt to execute a complex analysis (Todd & Gigerenzer, 2012) but may approximate it (Parpart et al., 2018). For example, Wilma may have used a heuristic like “Prioritize family relationships” to help her children resolve their situation. Ron may have used a heuristic like “Avoid the worst-case scenario” to help his candidate win the election. Heuristics often work well and do not require as much computation as the full optimization of a plan because they focus on just the most relevant pieces of information, reducing the chances of overfitting (Todd & Gigerenzer, 2012). Much of what we often refer to as “folk wisdom” comprises culturally-evolved heuristics, transmitted through mechanisms such as religious tradition (e.g., the Protestant Work Ethic’s advice to work hard and avoid waste; Gigerenzer, 2023).

A second type of task-level strategy relies on *narratives*—structured hypotheses that decision-makers construct using causal and analogical reasoning, which explain a situation and can be used to generate predictions and evaluate choices (Glück et al., 2005; Johnson et al., 2023). Although assembling and using narratives is more computationally costly than applying a heuristic, decision scientists have theorized that narratives too are a crucial simplifying strategy (Johnson et al., 2023). When Ron constructs worst-case scenarios, he is doing so using his causal knowledge (about government policy and voter psychology) and comparable experiences (about the fates of other campaigns). As in the case of heuristics, narratives can be socially transmitted and adapted within and across generations (Edmondson & Woerner, 2019).

The right task-level strategies are essential for navigating intractable situations—but they are not sufficient (Figure 1). First, although task-level strategies are sometimes computationally simple, a further process is required to seek out the inputs necessary for such strategies to work. Ron must determine whether he knows the most relevant facts for his election scenarios, and to fill any necessary gaps. Second, task-level strategies often conflict with one another. Daphne, like any doctor, would not get far if she were constantly second-guessing herself, so “trust your judgment” is a crucial heuristic. Yet “trust more knowledgeable experts” is equally important in any complex domain. Thus, metacognition is required to select the best strategy for a given situation (Rieskamp & Otto, 2006). Third, unreflective strategy use could spell disaster. Strategies can break at times they are needed most, as when the underlying pattern changes unpredictably; a “sanity checking” process is needed as a stopgap to ensure that a strategy does not lead to a nonsensical outcome. Wilma would not be likely to advise her children to compromise for the long-term good of their relationship if she found that one was taking advantage of the other.

Navigating this complexity requires the ability to reflect on and adapt one's task-level strategies—it requires metacognition. Monitoring and controlling task-level strategies by determining the right inputs, sanity checks, and conflict resolutions is the true business of wisdom, which is why consensus theories of wisdom place such emphasis on metacognition (Glück & Weststrate, 2022; Grossmann et al., 2020).

Table 2 summarizes several metacognitive processes commonly associated with wisdom, including processes that are primarily individual (e.g., intellectual humility) and those that involve social input (e.g., perspective seeking).

Metacognitive Process	Description
Intellectual humility	Awareness of what one does and does not know; acknowledgment of uncertainty and one's fallibility (Porter et al., 2022)
Epistemic deference	Willingness to defer to others' expertise when appropriate (Glück et al., 2005)
Scenario flexibility	Considering diverse ways in which a scenario might unfold to identify possible contingencies
Context adaptability	Identifying features of a situation that make it comparable to or distinct from other situations (Baltes & Smith, 2008)
Perspective seeking	Drawing on multiple perspectives where each offers information for reaching a good decision (Baltes & Smith, 2008)
Viewpoint balancing	Recognizing and integrating discrepant interests (Basseches, 1980; Sternberg, 1998)

Table 2. Example metacognitive processes commonly exhibited by wise people. For more detail, see Grossmann et al. (2020) and Glück and Weststrate (2022, Table 1).

These processes often work together to produce wise decisions. For example, Daphne the cardiologist exhibits intellectual humility when she recognizes that she does not understand why her patient's symptoms show a particular pattern (her usual heuristics have failed). She exhibits perspective-seeking by calling upon her colleague's expertise and context adaptability when she considers whether the unique situation of her individual patient limits the relevance of her colleague's expertise. Ultimately, she exhibits epistemic deference when she concludes her colleague's view is likely to be correct.

2.2. Toward AI wisdom: Machine metacognition

As with humans, strategies acquired by machines often fail in intractable situations. Could one therefore develop metacognitively wise machines that would function effectively in such contexts? We can construct parallels to human metacognition—reflecting on our thoughts and using that reflection to effectively direct subsequent thoughts and actions (Ho et al., 2022). Analogously, AI metacognition refers to the ability to model one's own computations and use that model to optimize subsequent computations.

To date, AI research has focused far more on task-level strategies than on metacognition (e.g., a long history of work on heuristics; Pearl, 1984), although the idea of machine

metacognition is not unprecedented (Horvitz, 2013; Russell & Wefald, 1991). Rudimentary forms of metacognition are incorporated by models that combine multiple task-level strategies by integrate information from several models, resulting in better performance than any model individually (Dong et al., 2020). GenAI models can perform well in some metacognitive tasks, such as classifying math problems according to what procedures are required to solve them (Didolkar et al., 2024), yet they continue to struggle with more complex metacognitive tasks. For example, they struggle to understand their goals (“mission awareness;” Li et al., 2024), exhibit overconfidence (Cash et al., 2024), and fail to appreciate the limits of their capabilities and context (e.g., stating they can access real-time information or take actions in the physical world; Li et al., 2024). These failures appear to be symptoms of a broader *metacognitive myopia*, which leads GenAI models to unnecessarily repeat themselves, poorly evaluate the quality of information sources, and overweigh raw data over more subtle cues to accuracy (Scholten et al., 2024).

2.2.1. Would AI wisdom resemble human wisdom?

While there is great room for improvement in AI metacognition, it is debatable whether building wise AI will simply boil down to implementing the metacognitive strategies of wise humans, since many computational constraints and social milieus of humans differ from those of AI systems.

A major task of human metacognition is to economize scarce cognitive resources (Todd & Gigerenzer, 2012; Simon, 1955), such as working memory. *Resource rationality* theories of human cognition go so far as to say that human (meta)cognition is the rational solution to a constrained optimization problem, where cognitive limits are the constraints (Lieder & Griffiths, 2017, 2020). On this view, many so-called cognitive biases are local side-effects of a globally optimized system (Levine et al., 2024; Sanborn & Chater, 2016). Since artificial systems have far more abundant computational resources, this logic favoring simplifying strategies is arguably weaker.

Despite our cognitive constraints, human sociality provides crucial metacognitive benefits. Humans live in groups, outsourcing much of our cognition to the social environment. The division of labor is a classic example, where the knowledge necessary to produce any product in the modern economy is distributed over many individuals (Hayek, 1945), as in the graphite pencil’s famous soliloquy: “not a single person on the face of this earth knows how to make me” (Read, 1958). Laid atop our evolved social cognitive capacities (Christakis, 2019), social, economic, and scientific institutions allow us to disperse knowledge and reasoning while all benefiting from it. Not only does socially distributed knowledge evolve—through cultural rather than biological evolution (Boyd & Richerson, 1985)—but so do those institutions themselves (North, 1990). These evolutionary processes allow humans to adapt to a constantly changing environment.

Perhaps, then, ideal machine wisdom would diverge from its human counterpart. For example, a wise AI system might be more willing to spin its wheels to solve a problem compared to a wise human; it might generate vast numbers of scenarios to analyze many

possible contingencies, evincing an extreme version of scenario flexibility. It is an open question when AI systems should deploy heuristics versus attempt an exact solution to a problem, given that heuristics sometimes seem to economize on cognitive resources at the expense of accuracy (Shah & Oppenheimer, 2008), but at other times to preserve or even improve accuracy by yielding inferences that are robust to overfitting (Todd & Gigerenzer, 2012). As another example, it is unclear to what extent AI systems might benefit from distributed knowledge (like humans in society) versus an extensive, integrated knowledge base.

Conversely, perhaps AI wisdom would converge considerably with human wisdom. Although AI systems face different computational constraints than humans, this may be a matter more of degree than of kind, so the same types of metacognitive strategies are probably necessary to allocate computational resources. Moreover, the core logic of heuristics may work just as well for machines as for humans: When we lack the necessary information to solve problems perfectly, heuristics can perform well if they bake in defaults that typically lead to desired outcomes. And while the social milieu of humans is quite different from that of machines, perhaps AI systems in the future will come to *join* our milieu, operating within and constrained by existing human institutions. Regardless of the ultimate details, modeling human wisdom using computational methods is likely to yield insights useful for building wiser AI systems.

3. What are the potential benefits of wise AI?

Building wise AI systems would address four major challenges.

3.1. Robust AI

Intelligent systems must function in a wide range of environments, including those that change in unpredictable ways, lack user-specified structure, and differ from any previously seen situation (Johnson et al., 2023). Indeed, as we argued above, AI systems will be called increasingly to act in such intractable situations, including out-of-distribution settings. Yet, such circumstances can exacerbate multiple types of errors—unreliability, bias, and inflexibility. Wise AI systems would likely be more robust in all three senses, primarily due to improved metacognition.

First, given similar inputs, an AI system lacking wisdom might give wildly different outputs (lacking *reliability*). Different outcomes given similar data could be due to either applying different strategies each time, or to applying the same strategy that nonetheless produces dissimilar results each time it is used (e.g., because of strong sensitivity to perturbations in initial conditions). High-quality metacognitive monitoring would evaluate whether it is sensible to use different strategies in comparable situations (e.g., whether trying a different strategy could yield new knowledge) and would reject strategies that produce wildly discrepant results on different occasions.

Second, an output may be systematically mistaken in a predictable direction (it is *biased*). Here, too, metacognition will help. Assuming that high-quality training procedures have

been used, the chief reason for biased outputs is biased *inputs*. But a metacognitive AI system with appropriate intellectual humility would reflect on its training data. For example, it might conclude that some parts of its training distribution are sample deficient, requesting the appropriate data before making inferences or understanding the causal process by which biases arose in the data, and thus being able to correct for them.

Third, novel inputs (e.g., an unfamiliar or changing environment) may lead to lower-quality outputs (lacking *flexibility*). By analogy to humans, our domain-general (as opposed to domain-specific) cognitive skills permit flexible and compositional reasoning about novel problems. Although poor generality to novel situations is arguably the most concerning potential failure mode of agentic AI systems, it has received insufficient attention. Yet functioning effectively in ambiguous, uncertain, and changing environments is precisely where human decision-makers shine (compared to our competitors anyway) and where wisdom plays an outsized role, for instance by helping reasoners to moderate their confidence in novel situations. Wisdom allows us to go beyond quantifying uncertainty by reducing, managing, and navigating it.

Appropriate task-level strategies are part of the story. Heuristics often outperform analytic optimization such as linear regression models (Todd & Gigerenzer, 2012), which in turn can outperform more complex models (Dawes & Corrigan, 1974). This is because overoptimizing may account for nuances in datasets that are noise rather than signal (Forster & Sober, 1994). Statisticians are well-aware of this problem and use techniques such as model selection criteria and cross-validation to avoid over-fitting. But even these techniques can break down in out-of-distribution contexts, as when the underlying pattern changes in very novel ways. By focusing on the core features of a situation, heuristics may be more robust to such novelty.

Metacognition is equally important for robustness. A key function of wisdom is calibrating, reducing, and managing uncertainty through metacognitive processes (see Table 2). For example, intellectual humility leads reasoners to moderate confidence in predictions when the environment is uncertain or novel, avoiding catastrophic error (e.g., the overconfident reliance on predictions made by derivative pricing models in the lead-up to the 2008 financial crisis). Perspective-seeking creates opportunities to learn new information outside the reasoner's current hypothesis space, cross-checking with other people or algorithms. More broadly, metacognition assists reasoners in managing task-level strategies, balancing the competing urges to simplify and complexify, selecting different strategies where they are most appropriate.

3.2. Explainable AI

AI systems often operate opaquely, with little ability to explain their reasoning to users (Dwivedi et al., 2023). Users cannot readily understand the rationale for an opaque system's output when it requires clarification, diagnose how such systems have gone astray when they make mistakes, or collaborate with such systems when the user and AI system have different conceptions of the problem. Therefore, explainability is a major focus in AI research.

Wise AI systems would likely bring superior explainability. Although humans often struggle to explain their task-level strategies (e.g., in cases of intuition), metacognition can help them to articulate their reasoning. Indeed, the broader metacognitive process of strategy *selection* may well be easier to explain than the specific task-level strategies themselves, as in other cases where abstract notions are more explainable than concrete details (Rozenblit & Keil, 2002). Given the importance of wise metacognition to explainability in humans, we consider two routes by which AI metacognition might contribute to AI explainability, rooted in two distinct views of how metacognition operates in humans.

According to the classical view, metacognitive strategies explicitly guide behavior. For example, consider again Daphne the cardiologist's decision to consult a more junior colleague regarding her patient. Faced with this situation, a wise reasoner might introspect on her knowledge, determine that she does not know enough to make an informed decision, and thereby seek out alternative points of view. Crucially, it is the conscious recognition of ignorance that caused this metacognitive strategy (seeking alternative viewpoints) to be deployed; the deliberate adoption of the metacognitive strategy caused new information from others to be integrated into the reasoner's decision; and this new information caused a particular decision to be reached. The reasoner would observe herself deploying these processes and would be able to verbally report them if asked. If they truthfully report this introspection, it would be broadly accurate.

A different view is that the mind is "flat"—it does not contain hidden depths of reasons that can be uncovered through introspection (Chater, 2018). This does not mean that people are unable to report their metacognitive strategies. Rather, it means that such reports are *inferences*, not *observations*. The reasoner would observe the *outputs* of her metacognitive strategies (e.g., her decisions) and reason backwards to what could have caused them (Cushman, 2020). Essentially, we invent stories to explain our behavior. Since we have large samples of observations about ourselves, these stories may well be useful; and verbally formulated reasons can constrain future thought and behavior. But our inferences about our metacognition are both exceedingly vague and often wrong (e.g., Nisbett & Wilson, 1977).

It is debatable which view best explains human metacognition. But these views provide radically different prescriptions for building explainable AI.

The classical view implies that, if a similar (meta)cognitive architecture were implemented in an AI system, explainability falls out as a straightforward consequence. The system merely needs to report the metacognitive process that causally led to its decision.

The mental flatness view implies a different picture. If an AI system were as opaque to itself as a human is (according to this view) to herself, then a further inference process would be required for the system's wise-but-inaccessible metacognitive processes to be reported. Rather than "introspecting" its metacognition, the system would need to generate a useful narrative to explain why it made the decision it did. Current GenAI systems seem to operate in this way. But in a wise AI system, such explanations are not

mere confabulations—they need to coherently explain past behavior, and crucially to constrain future behavior. Wise humans can generate verbal rules to explain, defend, justify, guide and modify their own behavior. True to the classic Socratic view of wisdom, wise AI will need to be able to do the same.

3.3. Cooperative AI

As AI is increasingly integrated into society, AI systems are coming to behave as parts of larger networks of intelligent entities (Dafoe et al., 2020). For example, autonomous vehicles must negotiate the rules of the road with both other autonomous vehicles (AI–AI cooperation). Robots collaborate with surgeons, pattern detection algorithms help radiologists (AI–human cooperation). Algorithms that govern language translation and social media content curation, even though they arguably do not cooperate directly, have the potential to promote or subvert cooperation among human users (human–human cooperation).

The field of cooperative AI examines how AI systems could effectively benefit all parties to an interaction by navigating cooperative barriers such as understanding, communication, commitment, and institutions (Dafoe et al., 2020). Wise task-level and metacognitive strategies are critical to how humans solve each of these problems, suggesting the same may be true for AI systems.

Understanding a social situation is a prerequisite for any sound approach to cooperation. For example, an autonomous vehicle negotiating a traffic situation needs to understand both the physical environment but also the likely actions taken by other agents (other vehicles, law enforcement). Since those actions depend on the mental states (beliefs, goals, and intentions) of agents, their social understanding requires theory-of-mind (Gopnik & Wellman, 1992) including the ability tacit to form joint plans to coordinate behavior (Chater et al., 2018). At the task-level, theory-of-mind is often simplified by assuming that the agent is rational (Gergely & Csibra, 2003; Johnson & Rips, 2015) and using this assumption to do inverse planning (working backwards from their actions to mental states; Baker et al., 2009). But people make this assumption in a context-sensitive manner (Grossmann & Eibach, 2024), for instance using higher-order theories about the extent to which an agent acts in a planned or habitual manner (Gershman et al., 2016).

Communication is equally crucial to cooperation. Successful cooperators must select and send relevant information to potential partners. Equally important, they must filter incoming information to act on what is useful and ignore what is irrelevant or misleading (Sperber et al., 2010). Even young children develop mechanisms for evaluating the trustworthiness of sources, such as examining their track record of accuracy and discounting testimony from sources with conflicts of interest (Sobel & Kushnir, 2013). Crucially, when testimony takes the form of an argument from premises to a conclusion, the reasoning itself can be checked, which may have been one evolutionary driver of human reasoning capacities (Mercier & Sperber, 2017). These epistemic vigilance mechanisms make credible communication among humans possible: Without a means

of assessing a communication, the risk of exploitation would be too high to maintain trust. Wise AI will require some analog of these abilities.

Even where communication succeeds, cooperation can still unravel when the parties' long-term incentives do not coincide. To address this, humans have evolved ways to make credible commitments. For example, reputation management and third-party punishment increase the cost of defection, promoting commitment if both parties understand these incentives (Fehr & Fishbacher, 2004). People are driven to avoid emotions such as shame and guilt, promoting commitment if both parties know the other to have these motivations (Frank, 1988). Any cooperative endeavor will require some commitment mechanism, although this may differ between human–human, AI–AI, and AI–human cooperation depending on the motivations and capabilities of different agents.

Wise metacognition is required to effectively manage these task-level mechanisms for social understanding, communication and commitment, which may be one factor underlying the empirical observation that wise people tend to act more prosocially, especially when their motives involve pursuit of cooperation or conflict resolution (Grossmann et al., 2017; Peetz & Grossmann, 2021). Cues to mental states, communication accuracy, and long-term commitment may conflict, requiring a system for balancing them. Metacognitive awareness is required to know whether one can check the plausibility of a complex chain of argumentation. Perhaps most importantly, many of these mechanisms require an understanding of the capabilities of the *other* counterparty. For instance, an AI system trying to convince a human to cooperate would fare poorly if it assumed that the human shared the same level of complex reasoning capabilities. Emotion-based mechanisms for commitment work among humans because there is a shared cultural understanding how these emotions work; it may be difficult for an AI system to rely on such mechanisms because we cannot understand the emotional phenomenology of “what it is like” to be an AI (if, indeed, it is ever “like” anything).

Finally, we can ask how AI systems fit into a broader ecosystem of social institutions. Human society has achieved far wider-spread cooperation compared to other species or to our ancestors thanks to institutions (e.g., governments, organizations, markets, and the social norms that make them possible). Thus, wisdom can benefit not only individuals (e.g., by making more robust decisions) and dyads (e.g., by improving communication), but entire societies. For example, a wise constitution may facilitate democratic decision-making and the peaceful transition of power, whereas a wise social network may heal (rather than deepen) social polarization. Clearly this sort of ‘wisdom’ is distinct from wise *reasoning*, since institutions do not perform ‘reasoning’ in the usual sense. Instead, wise institutions are those that, like wise individuals, promote collective good.

From this perspective, it is useful to think of AI not merely as an external tool *influencing* society but as potentially forming a new type of agent *within* society, embedded both in pairwise interactions and, increasingly, our broader institutions. It is a critical goal for future research in cooperative AI—and in the study of collective intelligence more broadly (Burton et al., 2024)—to understand how AI decision-making promotes or subverts broader social goals (Karlan & Allen, 2024). Just as wise leadership of groups may differ

from wise decision-making in a dyadic context (Everett et al., 2018), so might wise AI reasoning differ when its consequences might ripple across an entire society.

3.4. Safe AI

Tales of powerful AI gone awry are widespread in science fiction, but some worry that such catastrophic scenarios could become more than just whimsical stories. The logical core of the concern is two-fold: (i) Goals defined ahead of time are very likely to be mis-specified or to become obsolete, and (ii) A sufficiently powerful AI system could be very difficult to curtail if it aggressively pursued the wrong goals. Bostrom (2014) famously gives the example of the autonomous paperclip-maximizing AI system that converts the entire Earth into paperclips and kills all humans who get in its way. The AI's objective function has been mis-specified (the paperclip company's shareholders would prefer a world with fewer paperclips but more of everything else). But for any putatively more complete set of goals, it is worryingly easy to generate scenarios where they too go awry.

The goal of *AI alignment* (Ji et al., 2023) is to prevent such mismatches between an AI system's goals and those of its users—a task which is exceedingly difficult due to the range of unspoken assumptions we make and which an AI system would not necessarily share. In fact, exhaustively specifying goals in advance is difficult for similar reasons that humans cannot use top-down analytical tools to navigate intractable situations more broadly. Just as humans rely on wisdom to navigate such situations, AI systems might benefit from wisdom to navigate goal hierarchies.

But AI safety is a much broader concern than scenarios that are (at least for now) science fiction (Dalrymple et al., 2024). Even an AI system with goals aligned with its user is dangerous if the user is a criminal, terrorist, or adversarial government, as in the case of AI used for disinformation (e.g., deep fakes). Indeed, for now, the greatest risk is not powerful and malevolent AI systems, but those that simply do not work well—a shoddy surgical robot, incompetent tax advice system, or biased parole decision algorithm. In the case of such prosaic safety failures, others have observed that machine metacognition will be a crucial tool to fight such failure modes (Johnson, 2022). For example, AIs with appropriately calibrated confidence can target the most likely safety risks; appropriate self-models would help AIs to anticipate potential failures; and continual monitoring of its performance would facilitate recognition of high-risk moments and permit learning from experience.

3.4.1. Rethinking AI alignment

With respect to the broader goal of AI alignment, we are sympathetic to the goal but question this definition of the problem. Ultimately safe AI may be at least as much about constraining the power of AI systems within human institutions, rather than aligning their goals. But the notion of alignment itself—bringing human and machine values into harmony—faces a litany of not only technical problems, but conceptual ones.

First, humans are not even aligned *with each other*. This has been vividly illustrated in recent discourse surrounding how GenAIs should balance egalitarian norms versus providing accurate information to users, such as using stereotypes that may be statistically accurate but socially harmful (Gilbert, 2024). As AIs transcend mere conversation and act increasingly as agents in the world, such examples will deepen and proliferate.

Second, even if humans uniformly prioritized norm-following above other considerations, those norms differ sharply across cultures. Western societies typically prioritize individualistic values (Henrich et al., 2010), but other societies embrace a variety of other values—such as security, self-direction, and benevolence—that can themselves conflict (Sagiv & Schwartz, 2022). This problem is only exacerbated by the concentration of leading AI companies in one particular country (the U.S.), and representing a narrow slice of even that country.

Third, even if humans uniformly prioritized norm-following *and* those norms were uniform across cultures, what reason is there to think that those are the *right* norms? Social values have changed in both small ways (what jokes are funny versus cringe-worthy) and large (past norms and institutions considered abhorrent today; Varnum & Grossmann, 2017). To align AI systems to current values would risk reifying those values as “the right” values, stalling future social progress. Instead, society and its component individuals—including AI systems—should continue to allow values to evolve toward a shared reflective equilibrium (Rawls, 1971) that brings situation-specific judgments and general moral principles into alignment with one another through iterative adjustments—a metacognitive process.

Given these conceptual problems, alignment may not be a feasible or even desirable engineering goal. The fundamental challenge is how AI agents can live among us—and for this, implementing wise AI reasoning may be a more promising approach. Aligning AI systems to the right *metacognitive strategies* rather than to the “right” *values* might be both conceptually cleaner and more practically feasible. For example, task-level strategies may include heuristics such as a bias toward inaction: When in doubt about whether a candidate action could produce harm according to one of several possibly conflicting human norms, by default do not execute the action. Yet wise metacognitive monitoring and control will be crucial for regulating such task-level strategies. In the ‘inaction bias’ strategy, for example, a requirement is to learn what those conflicting perspectives are and to avoid overconfidence.

4. How might we build wise AI?

With a sketch of what wise AI might be like and why this might be desirable, we turn to the engineering challenges for artificial metacognition. Here, we offer some initial speculations about approaches to benchmarking and training wise AI.

4.1. Benchmarking

How will we know when an AI system is wise? Practical and conceptual challenges abound:

1. **Memorization.** Apparent success could be due to the system relying on patterns of reasoning specific to its training data rather than on novel reasoning.
2. **Process.** Since wisdom is about the reasoning underlying a strategy's selection, we need to evaluate not just the outcome but the process that led to it.
3. **Context.** Since the wise strategy is context-sensitive, the information provided to the AI for benchmarking must contain sufficient detail to match the rich context the AI would have in a real-world situation.

One way to gain traction is to consider how other complex constructs have been—or should be—benchmarked. For example, researchers have attempted to benchmark constructs such as theory-of-mind (Strachan et al., 2024) and analogies (Webb et al., 2023). One approach is to assemble a wide range of benchmarks used in psychology. However, since these tasks are already discussed in the scientific literature (the memorization problem), it is critical to replace the content to construct structural similar but superficially different problems (Frank, 2023). Moreover, since these tasks measure outcomes (i.e., the correct answer) and are relatively decontextualized, this approach cannot just be adopted wholesale.

A quite different approach was taken to benchmark AI systems' ability to creatively produce explanations (Thagard, 2024). The author—a philosopher with extensive domain expertise in explanatory reasoning—presented GPT-4 with a range of novel topics it was asked to explain. GPT-4's explanations were judged by the author as comparable to a thoughtful graduate student; it even generated a plausible new theory of consciousness. In this case, the author's subjective appraisal, rather than objective quantitative measures, were used. Such an approach is better-suited to evaluating reasoning (rather than outcomes), but its qualitative nature limits its applicability for assessing progress or comparing different models.

Combining aspects of both approaches, one way to benchmark wise AI may be to start with tasks that measure wise reasoning in humans (e.g., Grossmann et al., 2010). These tasks typically present participants with a social conflict and ask for a reflection of how it might proceed or how they would resolve it. The reflections are then scored on prespecified criteria by human raters. A range of such scenarios—and, critically, novel variants of them—could be presented to AI systems and their performance scored by either human raters (perhaps those who themselves score highly on wisdom measures) or by other systems if they could be demonstrated to converge in their scores. This method focuses on reasoning processes rather than outcomes (process problem). Equally, these scenarios must be novel (memorization problem) and detailed (context problem). It would also be important to present AIs with problems that they might confront if given adequate agency, to ensure they can reason wisely not only about humans but about themselves.

Ultimately, the wisdom of AI agents, as with people, will be judged by the rest of us as they act with increasing autonomy. Prior benchmarking is a crucial start, but there is no substitute for interaction with the real world. Given this intrinsic limit on our ability to evaluate AI wisdom *ex ante*, such integration must proceed slowly to minimize risks.

4.2. Training

Training task-level and metacognitive wisdom may require different strategies.

In humans, task-level strategies often involve heuristics that are culturally transmitted, and these are typically learned through a combination of trial-and-error learning (from experience) and of cultural learning (from person to person). Moreover, wise heuristics are often domain-specific and any attempt to exhaustively specify a full set of such rules is likely doomed for the same reasons that rule-based expert systems in AI failed. Instead, allowing AI systems to learn from their experience (as wise humans may; Dong et al., 2023; Grossmann et al., 2010) and from others (as all human cultures do; Henrich, 2018) may be a more promising route forward.

Yet this approach is unlikely to work for training metacognition, where the challenge is primarily deciding between strategies in a context-sensitive way and in soundly justifying those higher-order decisions. This contrasts with how AI systems are typically trained, where a loss function defined over the model's *outputs* (rather than reasoning) is minimized. Although this may indirectly select for sound decision-making strategies, the poor explainability of many state-of-the-art models makes it difficult to determine what those strategies are and it leaves the possibility of deception—i.e., producing an output that pleases a human judge for the wrong reasons.

This is not a straightforward problem; it may require multiple complementary approaches. One possibility is a two-step process, first training models for wise strategy selection directly (e.g., to correctly identify when to be intellectually humble) and then training them to use those strategies correctly (e.g., to carry out intellectual humble behavior). A second possibility may be to evaluate whether models are able to plausibly explain their metacognitive strategies in benchmark cases, and then simultaneously train strategies and outputs (e.g., training the model to identify the situation as one that calls for intellectual humility *and* to reason accordingly; e.g., Lampinen et al., 2022). In either case, models could be trained against what a wise human would do, or perhaps to explain and defend its choices to wise humans robustly (i.e., to stand up to 'cross-examination').

5. Could wise AI bring unintended and undesirable consequences?

Building smarter machines comes with risks: AI with humanlike or superhuman intelligence might adopt and pursue undesirable goals. Is there a parallel concern about building wiser machines? Plausibly, the answer is 'no.' Empirically, humans with wise metacognition show greater orientation toward the common good, including cooperation and responsiveness to other people (Grossmann et al., 2017). Perhaps wise AI systems would therefore have these qualities, too.

Yet these criteria are both aspirational (from an engineering perspective) and derived from observations of humans (from a scientific perspective). What if we tried and failed to build wise AI? What if the characteristics of wise AI differ from those of a wise human—to the detriment of us humans?

To these eminently reasonable concerns we have two responses. First, it is not clear what the alternative is. Compared to halting all progress on AI, building wise AI may introduce added risks alongside added benefits. But compared to the status quo—advancing task-level capabilities at a breakneck pace with little effort to develop wise metacognition—the attempt to make machines intellectually humble, context-sensitive, and adept at balancing viewpoints seems clearly preferable.

Second, by simultaneously promoting robust, explainable, cooperative, and safe AI, these qualities are likely to amplify one another. Robustness will facilitate cooperation (by improving confidence from counterparties in its long-term commitments) and safety (by avoiding novel failure modes; Johnson, 2022). Explainability will facilitate robustness (by making it easier to human users to intervene in transparent processes) and cooperation (by communicating its reasoning in a way that is checkable by counterparties). Cooperation will facilitate explainability (by using accurate theory-of-mind about its users) and safety (by collaboratively implementing values shared within dyads, organizations, and societies).

Wise reasoning, therefore, can lead to a virtuous cycle in AI agents, just as it does in humans. We may not know precisely what form wisdom in AI will take but it must surely be preferable to folly.

References

- Ardelt, M. (2004). Wisdom as expert knowledge system: a critical review of a contemporary operationalization of an ancient concept. *Human Development, 47*, 257–287.
- Baker, C. L., Saxe, R., & Tenenbaum, J. B. (2009). Action understanding as inverse planning. *Cognition, 113*, 329–349.
- Baltes, P. B., & Smith, J. (2008). The fascination of wisdom: Its nature, ontogeny, and function. *Perspectives on Psychological Science, 3*, 56–64.
- Basseches, M. (1980). Dialectical schemata: A framework for the empirical study of the development of dialectical thinking. *Human Development, 23*, 400-421.
- Bostrom, N. (2014). *Superintelligence: Paths, dangers, strategies*. Oxford, UK: Oxford University Press.
- Boyd, R., & Richerson, P. (1985). *Culture and evolutionary the process*. Chicago, IL: University of Chicago Press.
- Burton, J. W., Lopez-Lopez, E., Hechtlinger, S., Rahwan, Z., Aeschbach, S., Bakker, M. A., ... & Hertwig, R. (2024). How large language models can reshape collective intelligence. *Nature Human Behaviour, 8*, 1643-1655.
- Cascella, M., Semeraro, F., Montomoli, J., Bellini, V., Piazza, O., & Bignami, E. (2024). The breakthrough of large language models release for medical applications: 1-year timeline and perspectives. *Journal of Medical Systems, 48*, 22.
- Cash, T. N., Oppenheimer, D. M., & Christie, S. Quantifying UncertAInty: Testing the Accuracy of LLMs' Confidence Judgments. *Preprint*.
- Chater, N. (2018). *The mind is flat: The remarkable shallowness of the improvising brain*. New Haven, CT: Yale University Press.
- Chater, N., Misyak, J., Watson, D., Griffiths, N., & Mouzakitis, A. (2018). Negotiating the traffic: Can cognitive science help make autonomous vehicles a reality? *Trends in Cognitive Sciences, 22*, 93–95.
- Christakis, N. (2019). *Blueprint: The evolutionary origins of a good society*. Boston, MA: Little, Brown.
- Cushman F. (2020) Rationalization is rational. *Behavioral and Brain Sciences, 43*, e28.

Dafoe, A., Hughes, E., Bachrach, Y., Collins, T., McKee, K. R., Leibo, J. Z., Larson, K., & Graepel, T. (2020). Open problems in cooperative AI. *arXiv preprint arXiv:2012.08630*.

Dalrymple, D., Skalse, J., Bengio, Y., Russell, S., Tegmark, M., Seshia, S., ... & Tenenbaum, J. (2024). Towards guaranteed safe AI: A framework for ensuring robust and reliable AI systems. *arXiv preprint arXiv:2405.06624*.

Dawes, R. M., & Corrigan, B. (1974). Linear models in decision making. *Psychological Bulletin*, *81*, 95–106.

Didolkar, A., Goyal, A., Ke, N. R., Guo, S., Valko, M., Lillicrap, T., ... & Arora, S. (2024). Metacognitive capabilities of LLMs: An exploration in mathematical problem solving. *arXiv:2405.12205*.

Dong, M., Weststrate, N. M., & Fournier, M. A. (2023). Thirty years of psychological wisdom research: What we know about the correlates of an ancient concept. *Perspectives on Psychological Science*, *18*, 778–811.

Dong, X., Yu, Z., Cao, W., Shi, Y., & Ma, Q. (2020). A survey on ensemble learning. *Frontiers of Computer Science*, *14*, 241–258.

Dwivedi, R., Dave, D., Naik, H., Singhal, S., Omer, R., Patel, P., ... & Ranjan, R. (2023). Explainable AI (XAI): Core ideas, techniques, and solutions. *ACM Computing Surveys*, *55*, 1-33.

Edmondson, R., & Woerner, M. H. (2019). Sociocultural foundations of wisdom. In R. J. Sternberg & J. Glück (Eds.), *The Cambridge handbook of wisdom* (pp. 40-68). Cambridge, UK: Cambridge University Press.

Everett, J. A., Faber, N. S., Savulescu, J., & Crockett, M. J. (2018). The costs of being consequentialist: Social inference from instrumental harm and impartial beneficence. *Journal of Experimental Social Psychology*, *79*, 200-216.

Fehr, E., & Fischbacher, U. (2004). Third-party punishment and social norms. *Evolution and Human Behavior*, *25*, 63-87.

Forster, M., & Sober, E. (1994). How to tell when simpler, more unified, or less ad hoc theories will provide more accurate predictions. *The British Journal for the Philosophy of Science*, *45*, 1-35.

Frank, M. (2023). Baby steps in evaluating the capacities of large language models. *PsyArXiv preprint*. <https://osf.io/preprints/psyarxiv/uacjm>

Frank, R. H. (1988). *Passions within reason: The strategic role of the emotions*. New York, NY: Norton.

Gao, B., Song, F., Yang, Z., Cai, Z., Miao, Y., Dong, Q., ... & Chang, B. (2024). Omni-MATH: A universal Olympiad level mathematic benchmark for large language models. *arXiv:2410.07985*.

Gergely, G., & Csibra, G. (2003). Teleological reasoning in infancy: The naïve theory of rational action. *Trends in Cognitive Sciences*, 7, 287–292.

Gershman, S. J., Gerstenberg, T., Baker, C. L., & Cushman, F. A. (2016). Plans, habits, and theory of mind. *PLoS ONE*, 11, e0162246.

Gigerenzer, G. (2023). How do narratives relate to heuristics? *Behavioral and Brain Sciences*, 46, e94.

Gilbert, D. (2024). Google's 'woke' image generator shwos the limitations of AI. *Wired*.

Glück, J., & Bluck, S. (2013). The MORE Life Experience Model: A theory of the development of personal wisdom. In M. Ferrari & N. M. Weststrate (Eds.), *The scientific study of personal wisdom* (pp. 75–98). Berlin, Germany: Springer.

Glück, J., Bluck, S., Baron, J., & McAdams, D. (2005). The wisdom of experience: Autobiographical narratives across adult- hood. *International Journal of Behavioral Development*, 29, 197–208.

Glück, J., & Weststrate, N. M. (2022). The wisdom researchers and the elephant: An integrative model of wise behavior. *Personality and Social Psychology Review*, 26, 342–374.

Gopnik, A., & Wellman, H. M. (1992). Why the child's theory of mind really *is* a theory. *Mind & Language*, 7, 145–171.

Grossmann, I. (2017). Wisdom in context. *Perspectives on Psychological Science*, 12, 233–257.

Grossmann, I., Brienza, J. P., & Bobocel, D. R. (2017). Wise deliberation sustains cooperation. *Nature Human Behaviour*, 1, 0061.

Grossmann, I., & Eibach, R. E. (2024). Metajudgment: Metatheories and beliefs about good judgment across societies. *Current Directions in Psychological Science*.

Grossmann, I., Na, J., Varnum, M. E. W., Park, D. C., Kitayama, S., & Nisbett, R. E. (2010). Reasoning about social conflicts improves into old age. *Proceedings of the National Academy of Sciences of the United States of America*, 107, 7246–7250.

Grossmann, I., Weststrate, N. M., Ardelt, M., Brienza, J. P., Dong, M., Ferrari, M., ... & Vervaeke, J. (2020). The science of wisdom in a polarized world: Knowns and unknowns. *Psychological Inquiry*, 31, 103–133.

Hayek, F. A. (1945). The use of knowledge in society. *American Economic Review*, 35, 519–530.

Henrich, J. (2018). *The secret of our success: How culture is driving human evolution, domesticating our species, and making us smarter*. Princeton, NJ: Princeton University Press.

Henrich, J., Heine, S. J., & Norenzayan, A. (2010). The weirdest people in the world? *Behavioral and Brain Sciences*, 33, 61–83.

Ho, M. K., Abel, D., Correa, C. G., Littman, M. L., Cohen, J. D., & Griffiths, T. L. (2022). People construct simplified mental representations to plan. *Nature*, 606, 129-136.

Horvitz, E. J. (2013). Reasoning about beliefs and actions under computational resource constraints. *arXiv preprint arXiv:1304.2759*.

Ji, J., Qiu, T., Chen, B., Zhang, B., Lou, H., Wang, K., ... & Gao, W. (2023). AI alignment: A comprehensive survey. *arXiv preprint arXiv:2310.19852*.

Johnson, B. (2022). Metacognition for artificial intelligence system safety: An approach to safe and desired behavior. *Safety Science*, 151, 105743.

Johnson, S. G. B., Bilovich, A., & Tuckett, D. (2023). Conviction narrative theory: A theory of choice under radical uncertainty. *Behavioral and Brain Sciences*, 46, e82.

Johnson, S. G. B., & Rips, L. J. (2015). Do the right thing: The assumption of optimality in lay decision theory and causal judgment. *Cognitive Psychology*, 77, 42–76.

Karlan, B., & Allen, C. (2024). Engineered wisdom for learning machines. *Journal of Experimental & Theoretical Artificial Intelligence*, 36, 257-272.

Kay, J., & King, M. (2020). *Radical uncertainty: Decision-making beyond the numbers*. New York, NY: Norton.

Lampinen, A. K., Roy, N., Dasgupta, I., Chan, S. C., Tam, A., McClelland, J., ... & Hill, F. (2022, June). Tell me why! explanations support learning relational and causal structure. In *International Conference on Machine Learning* (pp. 11868-11890).

Levine, S., Chater, N., Tenenbaum, J., & Cushman, F. (2024). Resource-rational contractualism: A triple theory of moral cognition. *Behavioral and Brain Sciences*.

Li, Y., Huang, Y., Lin, Y., Wu, S., Wan, Y., & Sun, L. (2024). I think, therefore I am: Awareness in Large Language Models. *arXiv preprint arXiv:2401.17882*.

Lieder, F., & Griffiths, T. L. (2017). Strategy selection as rational metareasoning. *Psychological Review*, *124*, 762–794.

Lieder, F., & Griffiths, T. L. (2020). Resource-rational analysis: Understanding human cognition as the optimal use of limited computational resources. *Behavioral and Brain Sciences*, *43*, e1.

Lorenz, E. (1993). *The essence of chaos*. Seattle, WA: University of Washington Press.

Mascalzoni, E., & Regolin, L. (2011). Animal visual perception. *Wiley Interdisciplinary Reviews: Cognitive Science*, *2*, 106–116.

Mercier, H., & Sperber, D. (2017). *The enigma of reason*. Cambridge, UK: Harvard University Press.

Nisbett, R. E., & Wilson, T. D. (1977). Telling more than we can know: Verbal reports on mental processes. *Psychological Review*, *84*, 231–259.

North, D. C. (1990). *Institutions, institutional change and economic performance* (Vol. 332). Cambridge, UK: Cambridge University Press.

Parpart, P., Jones, M., & Love, B. C. (2018). Heuristics as Bayesian inference under extreme priors. *Cognitive Psychology*, *102*, 127-144.

Paul, L. A. *Transformative experience*. Oxford, UK: Oxford University Press.

Pearl, J. (1984). *Heuristics: Intelligent search strategies for computer problem solving*. Boston, MA: Addison-Wesley.

Peetz, J., & Grossmann, I. (2021). Wise reasoning about the future is associated with adaptive interpersonal feelings after relational challenges. *Social Psychological and Personality Science*, *12*, 629-637.

Porter, T., Elnakouri, A., Meyers, E. A., Shibayama, T., Jayawickreme, E., & Grossmann, I. (2022). Predictors and consequences of intellectual humility. *Nature Reviews Psychology*, *1*(9), 524–536.

Rawls, J. (1971). *A theory of justice*. Cambridge, MA: Harvard University Press.

Read, L. (1958). I, pencil: My family tree as told to Leonard E. Read. *The Freeman*.

- Rieskamp, J., & Otto, P. E. (2006). SSL: A Theory of How People Learn to Select Strategies. *Journal of Experimental Psychology: General*, *135*, 207–236.
- Rozenblit, L., & Keil, F. (2002). The misunderstood limits of folk science: An illusion of explanatory depth. *Cognitive Science*, *26*, 521-562.
- Russell, S., & Wefald, E. (1991). Principles of metareasoning. *Artificial Intelligence*, *49*, 361–395.
- Sagiv, L., & Schwartz, S. H. (2022). Personal values across cultures. *Annual Review of Psychology*, *73*, 517–546.
- Sanborn, A. N., & Chater, N. (2016). Bayesian brains without probabilities. *Trends in Cognitive Sciences*, *20*, 883-893.
- Scholten, F., Rebholz, T. R., & Hütter, M. (2024). Metacognitive myopia in Large Language Models. *arXiv preprint arXiv:2408.05568*.
- Shah, A. K., & Oppenheimer, D. M. (2008). Heuristics made easy: an effort-reduction framework. *Psychological Bulletin*, *134*, 207–222.
- Shu, D., Zhao, H., Liu, X., Demeter, D., Du, M., & Zhang, Y. (2024). LawLLM: Law Large Language Model for the US legal system. *arXiv preprint arXiv:2407.21065*.
- Simon, H. A. (1955). A behavioral model of rational choice. *Quarterly Journal of Economics*, *69*, 99–118.
- Sobel, D. M., & Kushnir, T. (2013). Knowledge matters: How children evaluate the reliability of testimony as a process of rational inference. *Psychological Review*, *120*, 779–797.
- Sperber, D., Clément, F., Heintz, C., Mascaro, O., Mercier, H., Origgi, G., & Wilson, D. (2010). Epistemic vigilance. *Mind & Language*, *25*, 359-393.
- Sternberg, R. J. (1998). A balance theory of wisdom. *Review of General Psychology*, *2*, 347–365.
- Strachan, J. W., Albergo, D., Borghini, G., Pansardi, O., Scaliti, E., Gupta, S., ... & Becchio, C. (2024). Testing theory of mind in large language models and humans. *Nature Human Behaviour*, *8*, 1285–1295.
- Thagard, P. (2024). Can ChatGPT make explanatory inferences? Benchmarks for abductive reasoning. *arXiv preprint arXiv:2404.18982*.
- Todd, P. M., & Gigerenzer, G. (2012). *Ecological rationality: Intelligence in the world*. New York, NY: Oxford University Press.

Van Veen, D., Van Uden, C., Blankemeier, L., Delbrouck, J. B., Aali, A., Bluethgen, C., ... & Chaudhari, A. S. (2024). Adapted large language models can outperform medical experts in clinical text summarization. *Nature Medicine*, 30, 1134-1142.

Varnum, M. E., & Grossmann, I. (2017). Cultural change: The how and the why. *Perspectives on Psychological Science*, 12, 956-972.

Walasek, L., & Brown, G. D. (2023). Incomparability and incommensurability in choice: No common currency of value? *Perspectives on Psychological Science*, 17456916231192828.

Webb, T., Holyoak, K. J., & Lu, H. (2023). Emergent analogical reasoning in large language models. *Nature Human Behaviour*, 7, 1526–1541.