

THE NEW STANDARD OF INTELLIGENCE

Rote Recall, Artificial Intelligence, and the Urgent Redefinition of Human Cognitive Value

Gary Chigaros

Abstract

For more than a century, institutions have relied on rote recall as the primary proxy for human intelligence. The convenience of this metric obscured its inadequacy. The emergence of powerful artificial intelligence systems has rendered that inadequacy undeniable: machines now retrieve, pattern-match, and reproduce information at scales and speeds that no human can approach. This paper advances the claim that in an AI-augmented epistemic environment, rote recall no longer functions as a valid or efficient proxy for human intelligence, and that institutional evaluation systems must reorient toward higher-order capacities resistant to algorithmic substitution. Drawing on cognitive science, philosophy of mind, labor economics, educational theory, and AI research, this paper argues that the highest-value human capacities are deep understanding over recall, creative synthesis over reproduction, and principled adaptability over rule-following. The paper further situates this argument within the framework of extended cognition (Clark and Chalmers, 1998): if AI functions as a reliable external memory and retrieval system, then evaluating unassisted recall is a category error, like testing a navigator's ability to calculate position by hand when GPS is available. The intelligence that remains distinctively human lies not in retrieval itself, but in the judgment, abstraction, and evaluative contribution the human brings to the extended system that no external tool can supply. Institutions that fail to reorient their frameworks

accordingly risk selecting for the wrong kind of person at precisely the wrong moment in history.

I. Introduction: The Recall Fallacy

Is my worth being reduced to rote memory? Because if that is the metric, it describes exactly the kind of cognitive labor artificial intelligence was designed to absorb, not define human value.

This question is deceptively simple. It is also among the most consequential questions facing institutions of education, enterprise, and governance in the present moment. For generations, the ability to store and retrieve information reliably has served as the dominant operational definition of intelligence in formal settings. The student who memorized the most facts earned the highest marks. The professional who could recite the most precedents, figures, or procedures commanded the greatest authority. The executive who held the most data in their head was presumed to hold the most insight.

This conflation of memory with intelligence was never without critics. But it was defensible, for a time, because the alternative, measuring understanding, judgment, creativity, and adaptive reasoning, was expensive, subjective, and methodologically complex. Rote recall was measurable. And so institutions measured it, and then made the compounding error of treating the measurement as though it were the thing itself. This critique is not new: Bloom (1956), Freire (1970), and Messick (1989) each identified aspects of this problem from different disciplinary vantages. What is new is the urgency. As Donald Campbell observed in 1979, the more any quantitative social indicator is used for social decision-making, the more subject it will be to corruption pressures and the more apt it will be to distort and corrupt the social processes it is intended to monitor (Campbell, 1979). What Campbell formalized as a law of measurement, education and hiring systems turned into institutional policy, with no apparent awareness of the irony.

Artificial intelligence has now made this error impossible to sustain. When a machine can retrieve any fact in any domain in milliseconds, at scales and speeds, continuing to use recall as the primary standard of human cognitive value is not merely outdated. It is a fundamental category error. AI has not devalued human intelligence. It has clarified, with sudden and striking precision, which version of intelligence was never worth measuring in the first place.

This paper proceeds in eight sections. Section II defines the operative conception of intelligence used throughout and examines how recall became its proxy. Section III characterizes what AI actually automates, drawing on technical literature to reframe AI as diagnostic rather than competitive, and as an extension of human cognition rather than a rival to it. Section IV advances a revised account of the three capacities that now constitute genuine cognitive value, unified under the classical concept of metis. Section V presents the strongest case for recall as foundational and then answers it directly. Section VI addresses five further objections to the argument. Section VII advances the positive thesis, introducing the concept of metis as the cognitive mode that replaces rote and examining the asymmetric effects of AI augmentation on practitioners oriented toward each mode. Section VIII examines the institutional consequences and the obligations that follow for those who set standards. A brief conclusion ties these threads together.

II. Intelligence as Proxy: Definitions and the Measurement Problem

Defining the Operative Conception of Intelligence

Intelligence is among the most contested terms in cognitive science and philosophy of mind, and any paper that uses it without clarification invites legitimate challenge. This paper does not adopt the psychometric conception of intelligence as a fixed general factor (g) measurable by standardized tests. Nor does it adopt the expansive conception of multiple intelligences proposed by Gardner (1983), which, while influential in educational circles, has faced persistent criticism regarding its empirical grounding. The paper does,

however, draw on the foundational distinction introduced by Cattell (1963) between fluid intelligence, the capacity to reason and solve novel problems independent of acquired knowledge, and crystallized intelligence, the accumulated store of culturally acquired knowledge and skills. That distinction is central to the argument here: crystallized intelligence, insofar as it reduces to stored and retrievable information, is precisely what AI externalizes most effectively. Fluid intelligence, the adaptive, novel-problem-solving component, is what remains most distinctively human and most resistant to algorithmic substitution. This paper therefore proceeds with a conception aligned with Sternberg's (1985) triarchic framework: intelligence as the capacity to adapt purposively to the environment through the flexible deployment of analytical, creative, and practical abilities. On this account, intelligence is not what one knows but how effectively one can acquire, apply, and extend knowledge across the changing conditions of real-world life.

This conception converges with Chollet's formal definition in the AI literature: intelligence as the rate at which a system can acquire new skills and solve novel problems from minimal examples, weighted by the generality of those skills across domains (Chollet, 2019). The two conceptions agree on what matters most: not the size of a stored knowledge base, but the capacity to generate effective responses to genuinely novel situations. Rote recall is a poor proxy for intelligence on either account. It measures neither adaptive flexibility nor the rate of new skill acquisition. It measures only the fidelity of storage and retrieval, a function that machines have long outperformed humans at and now perform at a scale that renders human recall practically obsolete as a competitive advantage.

How Recall Became the Proxy

Benjamin Bloom and colleagues, in their 1956 Taxonomy of Educational Objectives, established a hierarchy of cognitive processes ranging from knowledge recall at the base through comprehension, application, analysis, synthesis, and evaluation at the apex (Bloom, 1956). The revised taxonomy of Anderson and Krathwohl (2001) updated those categories to: remembering, understanding, applying, analyzing, evaluating, and creating, elevating creative synthesis to the highest cognitive tier. The intended message of both frameworks was unambiguous: recall is the floor of cognition, not its summit.

Yet institutional practice inverted this ordering. Standardized assessments, from IQ tests in the early twentieth century to the proliferating battery of high-stakes examinations that followed, operationalized intelligence primarily through the lower tiers of Bloom's hierarchy. The psychometrician Samuel Messick warned decades ago about construct validity, the risk that what a test actually measures may not correspond to the construct it claims to measure (Messick, 1989). Recall-based tests measure recall. The assumption that recall indexes deeper cognitive capacity is exactly the construct validity problem Messick identified, and it has gone largely unaddressed in practice.

Daniel Koretz, in his empirical study of American standardized testing, documented how the elevation of test scores to primary outcomes produced precisely what Campbell's law predicts: score inflation without corresponding gains in actual understanding (Koretz, 2017). Students and institutions became proficient at optimizing for the test rather than developing the capacities the test was ostensibly designed to capture. Paulo Freire called this banking education: the treatment of students as receptacles into which facts are deposited. It was not merely a philosophical critique. It described the functional reality of most formal evaluation systems (Freire, 1970).

The measure became the goal. And in becoming the goal, it corrupted the very process it was meant to monitor.

This matters because measurement shapes investment. When institutions reward recall, individuals invest in recall. When hiring processes screen for credentials that certify exposure to information, applicants optimize for credential accumulation rather than capability development. The feedback loop is self-reinforcing and, until recently, difficult to interrupt because the alternative, rigorously assessing understanding, creativity, and judgment, seemed impractically subjective. AI has changed the terms of that tradeoff decisively. The cost of measuring recall is now effectively zero, because machines perform it better. The relative value of measuring something else has never been higher.

III. What AI Actually Automates: Diagnosis, Not Competition

The popular discourse around artificial intelligence and human displacement tends toward one of two errors: either catastrophism, in which AI is imagined to replicate human cognition comprehensively and render human workers obsolete, or dismissiveness, in which AI is treated as a sophisticated autocomplete with no structural implications for how human capability should be understood. Both errors share a common cause: a failure to specify precisely what AI does and does not do. A more productive framing, and the one advanced here, treats AI as a distributed epistemic tool that extends human cognitive reach, and as a diagnostic instrument that reveals where evaluation frameworks have been systematically measuring the wrong things.

AI as Extended Cognition

Andy Clark and David Chalmers, in their foundational 1998 paper, proposed what they called active externalism: the thesis that cognitive processes need not be confined within the boundaries of skin and skull, and that environmental resources actively coupled to human cognition can constitute genuine parts of a cognitive system (Clark and Chalmers, 1998). On their account, a notebook that reliably stores information a person uses to guide behavior is not merely a tool; it functions, in the relevant cognitive sense, as an extension of memory.

This framework has direct implications for the argument here. If AI systems function as external cognitive extensions, then the cognitive labor they perform, recall, pattern recognition, information retrieval, is not eliminated from the human-plus-tool system; it is redistributed within it. The question of what to evaluate in human cognition therefore shifts: not what can a person retrieve unaided, but what can a person contribute to the extended system that the external component cannot supply? Evaluating only the internal component of a distributed cognitive system produces a systematically incomplete and misleading picture of the system's actual capacity.

A skeptical reader might object that if AI does the remembering, the human contributes nothing distinctive. That objection mistakes the nature of the extended system. Clark and

Chalmers's framework insists that the value of extended cognition lies in what the coupled system achieves, not in which component performs which function. The human element contributes what AI cannot: the judgment to determine what to retrieve and why, the understanding to evaluate retrieved information in context, the creativity to apply it in ways the system was not designed for, and the moral reasoning to decide what to do with it. The recall is externalized. The intelligence is not. The paper's central thesis does not stand or fall with the extended cognition framework. Even for readers who reject active externalism, the evaluative argument holds independently: once recall can be performed more cheaply and accurately by a machine than by a human, the use of recall performance as a proxy for human cognitive value loses its justification on purely pragmatic grounds, regardless of whether the machine constitutes part of a cognitive system or merely a tool.

One important complication of the extended cognition framing deserves direct acknowledgment. Research on automation bias demonstrates that humans paired with AI decision-support systems frequently become worse decision-makers, not better: they over-rely on automated recommendations, reduce their independent monitoring, and commit errors of omission and commission that would not occur without the automated aid (Parasuraman and Riley, 1997). This phenomenon, documented across aviation, medicine, and process control, reflects a genuine risk that the extended cognitive system, rather than amplifying human judgment, can atrophy it. The response to this concern is not to reject the extended cognition framework but to distinguish between a well-designed and poorly-designed extended system. Automation bias is a failure of system design and human-AI interaction governance, not an inherent property of cognitive extension. The paper's argument is not that all AI-augmented cognitive systems improve human performance automatically; it is that the evaluation of human cognitive value should be oriented toward the capacities the human brings to such systems rather than the capacities the AI already supplies. Designing those systems well, so that they complement rather than displace human judgment, is itself one of the most consequential institutional challenges the present moment demands.

What AI Cannot Do: Brittleness and the Limits of Pattern Completion

The technical literature is considerably more precise than popular discourse about AI's actual capabilities. As established in Section II, Chollet's (2019) formal conception of intelligence as efficient novel-skill acquisition provides the key distinction: AI systems demonstrate extraordinary skill within the distributional boundaries of their training data, but have not yet demonstrated comparable capacity for the kind of fluid abstraction from minimal examples that defines genuine intelligence on his account. ARC-AGI-2, the most demanding version of Chollet's benchmark, was specifically designed to maintain tasks that require deliberate thinking from humans while remaining very hard for AI systems. The top competition submission in 2025 achieved 24% accuracy under constrained conditions. Leading systems with iterative refinement loops have improved performance but still fall far short of human baselines on ARC-AGI-2, confirming the efficiency gap (ARC Prize Foundation, 2025). Frontier models continue to improve on earlier versions of the benchmark. What they have not demonstrated is efficiency: the ability to generalize from minimal examples without massive compute expenditure. That gap, between expensive skill and cheap general intelligence, is precisely what this paper's argument rests on.

These limitations may narrow with further architectural innovation. The conceptual distinction between interpolation within a training distribution and principled abstraction beyond it remains salient, however, even as specific benchmark scores evolve. The launch of ARC-AGI-3 in early 2026, designed to test agent-based reasoning under even more demanding conditions, reflects the ongoing recognition among AI researchers that general abstraction has not yet been achieved (ARC Prize Foundation, 2026). Bender, Gebru, McMillan-Major, and Mitchell made a related point from the perspective of linguistics: large language models operate as probabilistic pattern-completers without reference to meaning in any functional sense (Bender et al., 2021). What AI does with language is structurally similar to what rote recall does with knowledge: retrieval and recombination without genuine understanding.

Hans Moravec identified a related asymmetry in 1988, observing that the cognitive tasks humans find most difficult, formal logic, mathematics, pattern recognition across large

datasets, are among the easiest for machines to replicate, while the tasks humans find effortless, contextual judgment, reading unstated social meaning, acting under genuine uncertainty, have proven extraordinarily difficult to automate (Moravec, 1988). This inversion maps directly onto the argument here: the things that felt like intelligence because they were hard for humans turned out to be the things easiest to automate. The things that are hard to automate are the things institutions have been systematically undervaluing.

The Economic Dimension

The labor economist David Autor documented that computers substitute for workers performing routine, codifiable tasks while amplifying the comparative advantage of workers supplying problem-solving skills, adaptability, and creativity (Autor, 2015). Acemoglu and Restrepo extended this framework, showing that automation displaces labor in routine cognitive and manual tasks while creating new demand for tasks requiring judgment, communication, and adaptive reasoning (Acemoglu and Restrepo, 2019). The tasks displaced are exactly those that recall-based evaluation systems have long measured and rewarded. The tasks that gain value are precisely those those systems have underweighted. AI has disclosed a longstanding mismatch between what institutions measure and what genuinely matters.

IV. Human Cognitive Value Revisited: Three Capacities

If rote recall is no longer an adequate proxy for cognitive value, what replaces it? This section advances three capacities, each grounded in the relevant literature, that together constitute a revised account of human cognitive value in an AI-augmented environment. These are not soft skills, a label that has historically functioned to marginalize them. They are the hardest skills humans possess, the ones most resistant to algorithmic substitution. They have simply been difficult to measure, and measurement difficulty was historically mistaken for importance deficiency. The distinction is not between hard knowledge (static facts) and soft skills (interpersonal warmth), but between static storage and dynamic

capacity: the ability to apply, extend, and evaluate knowledge under conditions of genuine novelty and uncertainty.

1. Understanding: The Transfer of Knowledge Across Contexts

Understanding, in the cognitive science literature, refers to the capacity to apply knowledge flexibly across novel contexts. It is not the ability to retrieve a fact but to do something generative with it: to recognize its applicability where the connection is not obvious, to identify where it breaks down, and to integrate it with knowledge from other domains. Sternberg's triarchic framework frames this as the analytical component of intelligence, the capacity to compare, evaluate, and apply knowledge in structured ways, while noting that analytical giftedness without creative or practical extension produces people who are analyzers rather than synthesizers (Sternberg, 1985).

Bransford and colleagues distinguished between bare factual knowledge and what they called conditionalized knowledge, which includes understanding of the contexts in which a concept applies, those in which it does not, and why (Bransford et al., 2000). This is the distinction between a student who can recite a formula and a practitioner who understands when to apply it, when to adapt it, and when to abandon it. It is also the distinction between a language model that retrieves text matching a query and a clinician who understands which query to formulate, why, and what to do with the answer.

The revised Bloom's Taxonomy provides a useful scaffold: AI systems increasingly outperform humans on tasks concentrated in the lower tiers of the hierarchy, remembering and surface understanding, while the upper tiers, analyzing, evaluating, and creating, represent the domain where genuine human understanding manifests and where machines have not yet demonstrated comparable general capacity (Anderson and Krathwohl, 2001). Institutions that evaluate predominantly at the lower tiers are assessing the capacities machines now handle more efficiently.

2. Creative Synthesis: The Generation of Genuinely Novel Value

Creativity is not ornament. It is the mechanism by which genuinely new value enters the world. Sternberg's triarchic framework identifies synthetic giftedness, the creative component, as the capacity to generate novel ideas and solutions in response to

unfamiliar challenges, and notes that it is the least often captured by conventional testing, precisely because it does not fit a predefined answer space (Sternberg, 1985). The psychologist Margaret Boden distinguished between combinatorial creativity, the novel recombination of existing ideas, transformational creativity, the restructuring of a conceptual space, and exploratory creativity, the navigation of existing conceptual possibilities (Boden, 2004). All three involve generating outputs that are not predetermined by prior training data.

AI systems have not yet demonstrated transformational creativity at human levels. Their outputs are constrained by the distributional patterns of prior training. Chollet's benchmark work consistently shows that the performance gap between AI and humans is widest on tasks requiring genuine abstraction from novel configurations (Chollet, 2019; ARC Prize Foundation, 2025). The World Economic Forum's Future of Jobs Report 2025, drawing on surveys of over one thousand global employers, found that creative thinking remains among the top core skills demanded by employers today and is projected to be among the fastest-growing in importance through 2030 (World Economic Forum, 2025). The institutional implication is clear: creative synthesis cannot be evaluated through assessments with defined answer spaces, and its growing economic premium makes the failure to assess it not merely a philosophical error but a practical one.

3. Principled Adaptability: Judgment Under Genuine Uncertainty

The third capacity is what this paper terms principled adaptability: the ability to make decisions grounded in clear values and rigorous reasoning while remaining genuinely responsive to new information. Sternberg's practical intelligence component addresses exactly this, describing it as the capacity to apply knowledge effectively in real-world situations, to know how to act under genuine complexity rather than textbook conditions (Sternberg, 1985). This is distinct from indecisiveness, which reflects an absence of conviction, and from stubbornness, which reflects an absence of openness. It describes the cognitive and moral discipline to hold a strong position and revise it in response to evidence, without losing the capacity to hold positions at all.

Gary Klein documented this capacity in studies of expert practitioners operating under conditions of genuine uncertainty (Klein, 1999). Experts in Klein's research did not evaluate all options systematically before acting; they recognized patterns rapidly and adapted their mental models as situations developed in ways that novices following rules could not replicate. Gigerenzer's work on ecological rationality demonstrated that human intelligence is specifically adapted to make good decisions under uncertainty using incomplete information, precisely the conditions under which AI systems are most brittle (Gigerenzer, 2007). Machine learning models fail at the edges of their training distributions, exactly where human judgment is most needed. Herbert Simon's concept of bounded rationality described this as adaptive satisficing rather than optimization, a recognition that real-world decisions require flexible heuristic reasoning that formal procedures cannot capture (Simon, 1957).

The philosophical tradition associated with Hubert Dreyfus deepens this point. Dreyfus argued, drawing on the phenomenological tradition of Heidegger and Merleau-Ponty, that human expertise is constitutively embodied and situated in ways that rule-following systems cannot replicate (Dreyfus, 1992). The chess grandmaster does not apply rules; she sees the board. The experienced surgeon does not consult a checklist; she reads the tissue. This holistic, context-sensitive responsiveness, built through years of situated practice and inextricably bound to the practitioner's bodily and social immersion in a domain, is precisely what Dreyfus identified as the ceiling of human intelligence and the floor that formal AI systems have historically struggled to reach. Where Clark and Chalmers show that AI can extend the memory component of cognition, Dreyfus shows why the judgment component remains grounded in forms of human experience that cannot simply be distributed to an external tool. Together, these two frameworks demarcate the boundary: AI handles the crystallized, distributable functions; principled adaptability, as this paper defines it, is what remains on the human side of that line. It is also what most institutional evaluation systems have never bothered to measure.

These three capacities share a common structure. Each describes a form of intelligence that is situational, adaptive, and resistant to codification. The classical Greek tradition had a word for this constellation: *metis*. As Detienne and Vernant documented in their

study of Greek intelligence, metis denotes a form of practical cunning: the capacity to read a situation for its leverage points, to improvise under constraints, to act effectively when the rules do not fully specify the answer (Detienne and Vernant, 1978). Metis is not rule-following. It is the intelligence that operates where rules run out. Understanding, creative synthesis, and principled adaptability are its modern operational components. If this paper argues against rote as a proxy for intelligence, it argues equally for metis as the more accurate account of what human cognition contributes to the extended system that no machine can supply, as developed fully in Section VII.

V. The Strongest Case for Recall: A Steel-Man Analysis

Intellectual honesty requires that the strongest opposing argument be presented before it is answered. The case for recall as foundational to intelligence is not a naive position. It is well-supported in the empirical literature and draws on two of the most rigorous research traditions in cognitive science. This section gives that case its full weight before explaining why it does not defeat the central argument.

The Expertise Argument: Ericsson and Deliberate Practice

The most formidable challenge comes from K. Anders Ericsson's foundational research on expert performance. Ericsson, Krampe, and Tesch-Romer (1993) demonstrated that the distinguishing characteristic of expert practitioners, across domains from chess to music to surgery, is not innate talent but accumulated deliberate practice structured specifically to build domain-relevant memory representations. Expert chess players do not analyze positions more quickly because they reason better in the abstract; they recognize familiar configurations from memory and retrieve associated responses. Expert surgeons do not improvise in novel situations from first principles; they draw on dense libraries of encoded patterns built over thousands of hours of structured practice. On this account, what looks like creative judgment or adaptive expertise at the highest levels is, at its foundation, a form of highly sophisticated recall operating on deeply encoded memory structures. The implication for evaluation is direct: testing recall is not

measuring the floor of cognition. It is measuring the very substrate on which all higher cognition is built.

Ericsson's framework also challenges the paper's claim that AI externalizes the memory load that expertise requires. If expert memory structures are not merely stored facts but conditionalized, domain-specific schemas encoding when and how knowledge applies, then they cannot simply be offloaded to a retrieval system. The chess grandmaster's pattern library is not a database that could be queried externally; it is a cognitive structure built through years of situated practice and tightly integrated with perceptual and executive function. On this view, the extended cognition argument may underestimate the extent to which the most valuable forms of recall are constitutively internal.

The Cognitive Load Argument: Sweller and Working Memory

A second challenge comes from cognitive load theory. Sweller (1988) demonstrated that working memory is severely limited in both capacity and duration, and that schema automation, the process by which frequently encountered patterns become encoded as automatic responses requiring minimal working memory resources, is the primary mechanism by which learners develop the capacity to tackle complex problems. The implication is that automated recall does not compete with higher-order cognition; it enables it. When basic operations are automated, working memory is freed for the genuinely novel elements of a problem. On this account, rote recall in early learning is not a low-value activity to be superseded by higher-order thinking; it is the prerequisite infrastructure without which higher-order thinking becomes cognitively impossible.

Together, Ericsson and Sweller make a powerful case: recall-based evaluation may be measuring not merely a convenient proxy but a genuinely important substrate of the cognitive capacities this paper champions. If creative synthesis and principled adaptability are built on automated knowledge structures, then evaluating recall may capture something real about the foundation of those capacities, even if it fails to capture the capacities themselves.

Why the Steel-Man Does Not Defeat the Argument

These are serious objections, and they deserve a precise response rather than dismissal. The response operates on three levels.

First, both Ericsson and Sweller are making claims about development and learning, not about evaluation and institutional selection. Ericsson's argument is that expert performance is built on extensive practice-encoded memory; it does not follow that testing recall in standardized conditions identifies who has or will develop expert performance. Sweller's argument is that automated recall frees working memory for complex reasoning; it does not follow that recall performance on a decontextualized test is a valid indicator of the automated schema structures that actually matter. The distinction between what supports cognitive development and what evaluates it validly is precisely the construct validity problem this paper identifies. Messick (1989) is relevant here: a measure can correlate with a genuine construct without validly assessing it.

Second, Ericsson's own research actually supports this paper's critique of conventional evaluation, read carefully. His framework emphasizes that what distinguishes expert memory is not the volume of facts stored but the organization and conditionalization of that knowledge, exactly what Bransford and colleagues call conditionalized knowledge (Bransford et al., 2000). Standardized tests that measure decontextualized fact recall are not measuring what Ericsson's experts have. The empirical basis for this distinction is not merely asserted: Koretz's (2017) longitudinal analysis of high-stakes test score inflation documents precisely the divergence between recall performance and genuine competence, showing that score gains do not transfer to related assessments measuring deeper understanding. They are measuring the kind of inert, unconditionalized storage that his research explicitly contrasts with genuine expertise. Furthermore, Ericsson's own account of how expert schemas are built, through thousands of hours of situated deliberate practice in domain-specific contexts with immediate feedback, is precisely not the kind of activity that recall-based standardized evaluation rewards or selects for. The implication runs against the defense: if we genuinely wanted to cultivate the memory structures that Ericsson shows underpin expert performance, we would need to evaluate people through practice-based, contextually embedded assessment, which is exactly what

this paper argues for. Ericsson's work is a critique of conventional rote testing, not a defense of it.

Third, and most decisively, the externalization question must be answered in the context of how expert memory structures are actually built. Ericsson is correct that expert schemas are not reducible to retrievable facts. But the argument in this paper is not that AI replaces all cognitive memory. It is that AI replaces the specific kind of decontextualized, fact-retrieval performance that institutional evaluation systems have historically rewarded. The conditionalized schemas of Ericsson's experts are built through situated deliberate practice, not through the rote memorization that standardized tests reward. If anything, Ericsson's findings reinforce the argument: the memory that matters is not the memory that current evaluation systems measure, and AI has made the latter category obsolete as a competitive differentiator.

VI. Anticipated Objections and Responses

Objection 1: Isn't recall still foundational?

The objection that recall undergirds all higher cognition has surface plausibility. One cannot analyze what one does not know. The response is not to deny that knowledge matters but to distinguish between knowledge as a precondition and recall as a performance. A physician who can access any clinical reference in seconds does not need to hold diagnostic criteria in memory with the fidelity required before such tools existed. What she needs is the clinical judgment to know which reference to consult, how to interpret it in context, and how to weigh it against the particular patient in front of her. Clark and Chalmers make precisely this point: when an external resource reliably performs a cognitive function and the agent can access it dependably, the functional role is genuinely distributed (Clark and Chalmers, 1998). Evaluating only the internal component of a distributed cognitive system produces a systematically incomplete picture of what the system can actually do. The precondition of knowledge is not eliminated; the premium on its rote availability is.

Objection 2: Doesn't AI also perform synthesis?

Large language models produce outputs that superficially resemble synthesis, combining ideas from their training corpora in ways that can appear creative. The technical literature is clear, however, that this is probabilistic recombination rather than genuine generativity. Bender and colleagues established that language models operate without reference to meaning; their outputs are constrained by the distributional patterns in their training data (Bender et al., 2021). ARC-AGI benchmark results continue to confirm that models producing fluent synthetic text struggle significantly with tasks requiring genuine abstraction from novel configurations (Chollet, 2019; ARC Prize Foundation, 2025). The synthesis AI performs is within-distribution. Human creativity, as Boden (2004) details, is most valuable precisely at the transformational boundary, where pattern completion fails and genuine conceptual restructuring is required. These are not the same cognitive operation.

Objection 3: How do you measure what you are proposing?

This is the most serious objection and it deserves a direct answer. The difficulty of measuring understanding, creativity, and principled adaptability is real. It is not, however, a reason to continue measuring recall. The challenge is methodological, not inherent. Portfolio-based assessment, scenario-based evaluation, and deliberative assessment models have all been developed and validated in educational and professional contexts (Bransford et al., 2000). Medical training has long used simulation and observed clinical practice to evaluate judgment rather than factual recall. Aviation uses scenario-based competency assessment rather than written examinations as the primary gate for certification. These are not fringe innovations; they are established assessment paradigms in high-stakes domains that could be scaled. The paper does not claim that the transition is costless or that capacity-based assessment eliminates reliability challenges. Portfolio and performance-based models introduce well-documented problems of inter-rater variability and standardization that recall-based testing largely avoids. The honest response is that these are genuine engineering problems in assessment design, soluble with investment and rigor, whereas the validity failure of recall-based evaluation is not an engineering problem but a conceptual one: no amount of refinement will make recall a better proxy for intelligence once machines perform it more efficiently than humans.

Reliability without validity is a precise instrument aimed at the wrong target. The argument that we should measure what is easy to measure rather than what matters is precisely the institutional failure this paper diagnoses.

Objection 4: Is this elitist or exclusionary?

The concern that elevated standards of creative and adaptive intelligence privilege those with access to enriched educational environments is legitimate. The response is that the current system is not equitable; it is merely different in where its exclusions fall. Recall-based standardized testing has been extensively documented as culturally biased, indexed to socioeconomic access to test preparation, and systematically predictive of academic performance in ways that reproduce existing hierarchies without meaningfully assessing real-world contribution (Koretz, 2017). The argument here is not for a harder version of the existing exclusionary system but for a different set of capacities, ones that may be more equitably distributed across populations historically disadvantaged by recall-based evaluation. Sternberg's empirical work demonstrated that students from diverse populations showed different profiles of analytical, creative, and practical ability, and that conventional assessments systematically missed the strengths of those whose intelligence was more practical or creative in character (Sternberg, 1985). This requires investment in assessment design and access, not acceptance of the status quo.

Objection 5: Is this just repackaged soft skills?

The dismissal of creativity, judgment, and adaptive reasoning as soft skills reflects the very evaluative framework this paper challenges. The term soft has historically functioned to render these capacities secondary to technical and quantitative skills, despite consistent evidence from labor economics, cognitive science, and organizational research that they are the primary drivers of performance under conditions of uncertainty and complexity. The World Economic Forum identifies analytical and creative thinking as the top core skills demanded by employers in 2025, projected to grow in importance further through 2030 (World Economic Forum, 2025). Gigerenzer's work demonstrates that adaptive heuristic reasoning outperforms formal optimization under real-world conditions (Gigerenzer, 2007). Acemoglu and Restrepo's labor economics research documents that these are precisely the capacities automation cannot displace (Acemoglu

and Restrepo, 2019). Sternberg's empirical studies found that students with high creative and practical intelligence were consistently underserved by conventional metrics that favored analytical recall (Sternberg, 1985). These are not soft. They are the hardest skills humans possess.

VII. The Metis Thesis: What Replaces Rote

The preceding sections have built a case against rote recall as a proxy for intelligence. That case is now established: the proxy has lost its justification, the construct validity problem is clear, and the three capacities that constitute genuine cognitive value have been identified and defended against serious objections. What remains is to name the positive thesis with sufficient precision that it can function as a framework, not merely a critique.

Rote and Metis as Cognitive Modes

The distinction this paper draws is not between knowing and not knowing. It is between two fundamentally different modes of engaging with knowledge.

Rote, as this paper uses the term, denotes a cognitive orientation organized around storage and retrieval. The rote-oriented practitioner follows procedures, retrieves answers from memory, executes checklists, and operates most effectively within the boundaries of established protocol. Rote is not unintelligent. It is the cognitive mode that institutional evaluation has historically rewarded because it is the mode most amenable to standardized measurement.

Metis, by contrast, denotes a cognitive orientation organized around situated judgment and adaptive action. The term originates in the Greek intellectual tradition, where it described a specific form of practical intelligence: the cunning of the navigator who reads the sea rather than reciting charts, the improvisation of the craftsman who adjusts technique to the grain of the material rather than following a fixed procedure (Detienne and Vernant, 1978; see also Scott, 1998, for its application to institutional planning and the legibility of knowledge). *Metis* is the intelligence of the practitioner who sees what is

not explicit in the situation: the leverage point, the asymmetry, the response that no procedure specifies but that the circumstances demand.

This is not a binary. Individuals exhibit both modes in varying proportion, and effective performance in most domains requires some facility with each. The claim is not that rote has no value. The claim is that rote has been systematically overweighted in institutional evaluation because it was easier to measure, and that the arrival of AI has collapsed whatever residual justification that overweighting had. When a machine performs rote functions at arbitrary scale and near-zero cost, the cognitive mode most likely to differentiate the human contributor is metis.

Dreyfus's phenomenological account of expertise maps directly onto this distinction. The novice operates in rote mode: following rules, applying procedures, checking boxes. The expert operates in metis mode: perceiving the situation holistically, acting from pattern recognition that cannot be fully articulated, adapting in real time to conditions that no checklist anticipated (Dreyfus, 1992). What Ericsson's research on deliberate practice actually documents, properly understood, is the developmental trajectory from rote to metis: the progressive internalization and conditionalization of knowledge until it becomes the kind of situated, adaptive expertise that metis describes. Standardized tests do not measure this trajectory. They measure the starting material, not the transformation.

Klein's research on naturalistic decision-making provides the empirical ground. Expert firefighters, intensive care nurses, and military commanders in Klein's studies did not succeed by retrieving the correct answer from memory. They succeeded by reading the situation for cues that indicated which of their experience-built mental models applied, and by adapting those models in real time as the situation developed (Klein, 1999). This is metis in operation. It is also precisely the cognitive mode that recall-based evaluation fails to capture, because metis cannot be demonstrated in a decontextualized test environment. It requires a situation.

The Augmentation Asymmetry

The arrival of AI tools does not affect rote-oriented and metis-oriented practitioners symmetrically. The asymmetry is structural, and its consequences are already visible.

For the practitioner whose cognitive strength is metis, AI removes an artificial barrier. The bottleneck was never insight. It was throughput: the mechanical labor of retrieving information, formatting documents, running standard analyses, executing routine procedures. AI absorbs that labor. The metis-oriented practitioner can now operate at the speed of their conceptual ability rather than at the speed of their rote capacity. The result is amplification. The practitioner who could always see the right move but was slowed by procedural overhead can now execute at a pace that matches their judgment.

For the practitioner whose cognitive strength is rote, AI produces a different and less favorable dynamic. The rote-oriented practitioner can now execute faster, but the cognitive contribution they bring to the extended system is thinner. Worse, the interactive properties of current AI systems introduce a specific risk: large language models exhibit well-documented sycophantic tendencies, affirming the direction of the user's query and producing outputs that pattern-match to the user's framing rather than correcting it (Perez et al., 2022; Sharma et al., 2024). For the practitioner who lacks the metis to evaluate AI outputs critically, this creates a validation loop in which shallow work is produced faster, confirmed by a system that does not understand it, and shipped with false confidence. The extended cognitive system, in this case, does not amplify judgment. It amplifies the absence of judgment.

This asymmetry has a direct implication for institutional evaluation. Systems that continue to select primarily for rote are selecting for the practitioners who benefit least from AI augmentation and who are most vulnerable to its failure modes. Systems that select for metis are selecting for the practitioners who can use AI as a genuine force multiplier.

K-Shaped Divergence

The economic consequence of the augmentation asymmetry is not uniform displacement. It is divergence.

Labor economists have documented that automation displaces routine cognitive tasks while amplifying the returns to non-routine judgment, creativity, and adaptive reasoning (Autor, 2015; Acemoglu and Restrepo, 2019). The metis framework sharpens this observation: the divergence is not merely occupational but evaluative. Roles that reward rote storage are substituted. Roles that reward situated judgment are amplified, precisely because institutions have historically failed to select for the latter. The metis framework gives this observation a sharper edge. Roles that are primarily rote, where the human contribution consists mainly of storage, retrieval, and procedural execution, are the roles where AI substitution is most direct and most complete. Roles where the human contribution is primarily metis, where value comes from situated judgment, creative problem-solving, and the capacity to act under genuine uncertainty, are the roles where AI augmentation produces the largest productivity gains.

The result is K-shaped: value concentrates upward among metis-oriented practitioners whose judgment is amplified by AI, while it erodes for rote-oriented practitioners whose primary cognitive contribution is the function AI performs most cheaply. This is not a prediction. The World Economic Forum's 2025 data already shows the pattern: creative thinking and analytical judgment are the fastest-growing skill demands, while routine cognitive tasks are among the fastest-declining (World Economic Forum, 2025). The 2026 scenario analysis confirms that across all projected futures, outcomes depend on whether human capacities for judgment and adaptability develop alongside AI systems (World Economic Forum, 2026).

The K-shaped divergence is not a technological inevitability. It is a consequence of institutional failure to reorient evaluation, selection, and development toward the cognitive mode that AI augments rather than the cognitive mode AI replaces. Institutions that correct this failure will develop and retain the practitioners who can use AI effectively. Institutions that do not will find themselves staffed by individuals whose primary skill is the one their tools already perform better.

VIII. Institutional Consequences: Obligations of Those Who Set Standards

The argument developed in the preceding sections has specific implications for the institutions that define, reward, and select for human cognitive value. The shift required is not primarily technological but evaluative: the problem is not that institutions lack the tools to assess different capacities, but that they have not exercised the will to do so. The current alignment between what institutions measure and what the economy, society, and epistemic environment require has broken down, and the costs of not correcting it are real and accelerating.

Education

Educational systems from primary through postgraduate levels remain organized substantially around knowledge transmission and recall verification. The revised Bloom's Taxonomy was published in 2001; more than two decades later, most high-stakes assessments still concentrate evaluation at the remembering and surface understanding tiers rather than the analyzing, evaluating, and creating tiers that Bloom's own framework identifies as the apex of cognitive achievement (Anderson and Krathwohl, 2001). The OECD's Education 2030 framework explicitly identifies transformative competencies, including creating new value, reconciling tensions, and taking responsibility, as the primary requirements of an economy reshaped by automation (OECD, 2019). Sternberg's research demonstrated that students taught and assessed using triarchic methods consistently outperformed those in conventional memory-focused curricula, not only on performance assessments but on the analytical tasks those conventional curricula claimed to develop (Sternberg, 1985). The gap between pedagogical ambition and assessment practice is not a knowledge gap. It is a gap in institutional will.

The metis framework specifies what that will must be directed toward. Assessment systems should evaluate whether students can direct AI tools toward problems the tools cannot identify on their own, rather than whether students can perform the retrieval functions those tools have already absorbed. Agent-directed project portfolios, in which the student demonstrates the capacity to govern, evaluate, and extend AI-generated work

rather than merely produce unassisted recall, represent a more valid measure of the cognitive capacities the economy now demands.

Hiring and Organizational Selection

Corporate hiring practices remain anchored to credentials that certify educational exposure rather than demonstrated capability. Resumes, degree requirements, and interviews that test factual command of a domain reflect institutional inertia more than evidence-based selection. Autor's research demonstrates that the skills most complemented by automation are precisely those that current hiring processes systematically underweight: non-routine problem-solving, adaptive judgment, and creative synthesis (Autor, 2015). Acemoglu and Restrepo's framework confirms that the highest labor market premiums in automated economies accrue to workers whose cognitive contributions cannot be codified or routinized (Acemoglu and Restrepo, 2019). Organizations that continue to select for recall and credential accumulation will compound their disadvantage as AI assumes more of the cognitive labor those selections were designed to source.

The practical alternative is already visible in adjacent domains. A portfolio of governed outputs, demonstrating the practitioner's judgment in directing, evaluating, and correcting AI-assisted work, provides a more valid signal of metis than any credential that certifies exposure to a curriculum. The hiring systems that will identify metis-oriented practitioners are those that evaluate what the candidate does with tools, not what the candidate can do without them.

Leadership Selection and Governance

The implications for leadership selection are perhaps most acute. Consequential leadership decisions are made under specific conditions: genuine uncertainty, incomplete information, competing values, and no clear precedent. These are exactly the conditions under which principled adaptability matters most and recall matters least. Klein's research on naturalistic decision-making documents that expert judgment in high-stakes conditions depends on rapid pattern recognition, adaptive mental models, and the capacity to act without complete information (Klein, 1999). Selection systems that elevate

credential attainment and factual command as proxies for leadership potential are selecting, systematically, for the capacities that matter least when leadership is most demanded.

Metis is the operational definition of the cognitive capacity leadership demands. The capacity to read situations for leverage points, to act under genuine uncertainty, and to adapt without losing conviction is precisely the capacity Klein's research documents in expert practitioners. Leadership selection systems that cannot identify metis will continue to produce leaders whose primary qualification is the ability to perform well on assessments calibrated to the cognitive mode AI has rendered competitively irrelevant.

Policy

Policymakers bear a particular responsibility because the evaluation frameworks they endorse and fund shape the developmental experiences of entire populations. The OECD's recognition of transformative competencies as central to economic and civic life in the coming decades is a starting point, not an endpoint (OECD, 2019). The World Economic Forum projects that 39 percent of core job skills will change by 2030, with creative thinking and adaptive reasoning growing fastest in importance across virtually every industry surveyed (World Economic Forum, 2025). The Forum's 2026 scenario analysis reinforces the urgency: across all four projected futures for the economy at the intersection of AI advancement and workforce readiness, from accelerated progress to stalled adoption, the outcomes hinge on whether institutions succeed in developing human capacities for judgment, adaptability, and creative contribution alongside AI systems, or fail to do so (World Economic Forum, 2026). Closing the gap between that recognition and the structure of national assessment systems requires willingness to fund assessment innovation, to resist institutional pressure to measure what is easy rather than what matters, and to accept that the transition from recall-based to capacity-based evaluation will be disruptive in the near term and necessary over any longer horizon.

The risk of institutional lag is that the transition happens anyway, driven by competitive pressure rather than deliberate policy. Organizations that develop metis-oriented evaluation internally will outperform those that do not. Nations whose educational and

credentialing systems remain indexed to rote will export their most adaptive practitioners to jurisdictions that recognize them. The policy question is not whether the shift from rote to metis will accelerate. It is whether public institutions will lead the transition or be forced to follow it.

IX. Conclusion: Intelligence After Recall

The hierarchy of cognitive value is shifting beneath the foundations of institutions that have not yet noticed the ground moving. Rote recall, long the crown jewel of formal intellectual achievement, is being displaced by artificial intelligence, not because AI has replicated human intelligence but because it has replicated the narrow slice of human cognitive performance that institutions chose to measure as a proxy for intelligence.

What is being revealed is not a new truth. Bloom knew in 1956 that recall was the floor, not the ceiling, of cognition. Campbell knew in 1979 that measuring a proxy corrupts the process the proxy was designed to monitor. Freire knew that banking education produced capable memorizers rather than capable thinkers. Sternberg knew in 1985 that genuine intelligence encompasses analytical, creative, and practical capacity, and that conventional testing captured only a slice of the first. Autor and Acemoglu knew that automation displaces routine cognitive labor while amplifying the premium on adaptive, creative, and synthetic capacities. Clark and Chalmers knew that cognitive systems extend beyond the skull, and that what matters is what the coupled system achieves, not which component supplies which function. None of this is new.

What is new is the cost of continuing to ignore it. When the proxy for intelligence is cheaper, faster, and more accurate when performed by a machine, the proxy has lost whatever justification it once had. The institutions that persist in using it are not preserving a standard. They are avoiding the more difficult work of building one that reflects what human intelligence actually is.

The question is not whether AI will change the value of human intelligence. It already has. The question is whether our institutions will catch up before they select for the wrong kind of person for another generation.

The three capacities argued for here, deep understanding, creative synthesis, and principled adaptability, are a starting point: the capacities that matter most in the specific conditions AI is creating, that are most resistant to algorithmic substitution, and that existing evaluation frameworks most systematically fail to assess. Developing rigorous methods for cultivating and recognizing these capacities is the central educational, organizational, and policy challenge of the present moment.

Rote recall was never the point. Meaning was.

The metis framework gives this claim its operational name. What this paper identifies as the cognitive capacities most resistant to algorithmic substitution, most valuable in the conditions AI is creating, and most systematically neglected by institutional evaluation, is what the Greek intellectual tradition called metis: the intelligence that operates where procedure ends and judgment begins. The shift from rote to metis is not a prediction about the future. It is a description of the present, visible in labor market data, in the performance gap between AI-augmented practitioners who can evaluate their tools and those who cannot, and in the widening distance between what institutions measure and what the economy rewards.

In operational terms, this requires revising evaluation systems to privilege transfer, abstraction, and adaptive judgment over unaided storage. The shift required is therefore not technological but evaluative, and that makes it a choice, not a constraint.

References

- Acemoglu, D., and Restrepo, P. (2019). Automation and new tasks: How technology displaces and reinstates labor. *Journal of Economic Perspectives*, 33(2), 3-30.
- Anderson, L. W., and Krathwohl, D. R. (Eds.). (2001). *A taxonomy for learning, teaching, and assessing: A revision of Bloom's taxonomy of educational objectives*. Longman.

- ARC Prize Foundation. (2025). ARC Prize 2025: Technical report and leaderboard. <https://arcprize.org/>
- ARC Prize Foundation. (2026). ARC-AGI-3: Technical report and leaderboard. <https://arcprize.org/>
- Autor, D. H. (2015). Why are there still so many jobs? The history and future of workplace automation. *Journal of Economic Perspectives*, 29(3), 3-30.
- Bender, E. M., Gebru, T., McMillan-Major, A., and Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FACCT '21)*, 610-623. <https://doi.org/10.1145/3442188.3445922>
- Bloom, B. S. (Ed.). (1956). *Taxonomy of educational objectives: The classification of educational goals. Handbook I: Cognitive domain.* David McKay.
- Boden, M. A. (2004). *The creative mind: Myths and mechanisms* (2nd ed.). Routledge.
- Bransford, J. D., Brown, A. L., and Cocking, R. R. (Eds.). (2000). *How people learn: Brain, mind, experience, and school* (expanded ed.). National Academy Press.
- Campbell, D. T. (1979). Assessing the impact of planned social change. *Evaluation and Program Planning*, 2(1), 67-90.
- Cattell, R. B. (1963). Theory of fluid and crystallized intelligence: A critical experiment. *Journal of Educational Psychology*, 54(1), 1-22. <https://doi.org/10.1037/h0046743>
- Chollet, F. (2019). On the measure of intelligence. arXiv preprint arXiv:1911.01547.
- Detienne, M., and Vernant, J.-P. (1978). *Cunning intelligence in Greek culture and society* (J. Lloyd, Trans.). University of Chicago Press.
- Perez, E., Ringer, S., Lukosuite, K., Nguyen, K., Chen, E., Heiner, S., and Kaplan, J. (2022). Discovering language model behaviors with model-written evaluations. arXiv preprint arXiv:2212.09251.
- Scott, J. C. (1998). *Seeing like a state: How certain schemes to improve the human condition have failed.* Yale University Press.
- Sharma, M., Tong, M., Korbak, T., Duvenaud, D., Askill, A., Bowman, S. R., and Perez, E. (2024). Towards understanding sycophancy in language models. *International Conference on Learning Representations (ICLR)*.
- Clark, A., and Chalmers, D. (1998). The extended mind. *Analysis*, 58(1), 7-19. <https://doi.org/10.1093/analys/58.1.7>
- Dreyfus, H. L. (1992). *What computers still can't do: A critique of artificial reason* (revised ed.). MIT Press.
- Ericsson, K. A., Krampe, R. T., and Tesch-Romer, C. (1993). The role of deliberate practice in the acquisition of expert performance. *Psychological Review*, 100(3), 363-406. <https://doi.org/10.1037/0033-295X.100.3.363>
- Freire, P. (1970). *Pedagogy of the oppressed.* Herder and Herder.
- Gardner, H. (1983). *Frames of mind: The theory of multiple intelligences.* Basic Books.
- Gigerenzer, G. (2007). *Gut feelings: The intelligence of the unconscious.* Viking.
- Klein, G. (1999). *Sources of power: How people make decisions.* MIT Press.
- Koretz, D. (2017). *The testing charade: Pretending to make schools better.* Harvard University Press.
- Messick, S. (1989). Validity. In R. L. Linn (Ed.), *Educational measurement* (3rd ed., pp. 13-103). Macmillan.
- Moravec, H. (1988). *Mind children: The future of robot and human intelligence.* Harvard University Press.

- OECD. (2019). OECD learning compass 2030: A series of concept notes. OECD Future of Education and Skills 2030. <https://www.oecd.org/education/2030-project/>
- Parasuraman, R., and Riley, V. (1997). Humans and automation: Use, misuse, disuse, abuse. *Human Factors*, 39(2), 230-253. <https://doi.org/10.1518/001872097778543886>
- Simon, H. A. (1957). *Models of man: Social and rational*. Wiley.
- Sternberg, R. J. (1985). *Beyond IQ: A triarchic theory of human intelligence*. Cambridge University Press.
- Sweller, J. (1988). Cognitive load during problem solving: Effects on learning. *Cognitive Science*, 12, 257-285. https://doi.org/10.1207/s15516709cog1202_4
- Sweller, J. (1994). Cognitive load theory, learning difficulty, and instructional design. *Learning and Instruction*, 4(4), 295-312.
- World Economic Forum. (2025). *The future of jobs report 2025*. World Economic Forum. <https://www.weforum.org/publications/the-future-of-jobs-report-2025/>
- World Economic Forum. (2026). *Four futures for jobs in the new economy: AI and talent in 2030*. World Economic Forum. <https://www.weforum.org/publications/four-futures-for-jobs-in-the-new-economy-ai-and-talent-in-2030/>