
A Definition of AGI

Dan Hendrycks¹, Dawn Song², Christian Szegedy³, Honglak Lee⁴, Yarin Gal⁵,
Erik Brynjolfsson⁶, Sharon Li⁷, Andy Zou^{1,8,9}, Lionel Levine¹⁰, Bo Han¹¹,
Jie Fu¹², Ziwei Liu¹³, Jinwoo Shin¹⁴, Kimin Lee¹⁴, Mantas Mazeika¹,
Long Phan¹, George Ingebretsen¹, Adam Khoja¹, Cihang Xie¹⁵,
Olawale Salaudeen¹⁶, Matthias Hein¹⁷, Kevin Zhao¹⁸, Alexander Pan²,
David Duvenaud^{19,20}, Bo Li²¹, Steve Omohundro²², Gabriel Alfour²³,
Max Tegmark¹⁶, Kevin McGrew²⁴, Gary Marcus²⁵, Jaan Tallinn²⁶,
Eric Schmidt¹⁶, Yoshua Bengio^{27,28}

¹Center for AI Safety ²University of California, Berkeley ³Morph Labs

⁴University of Michigan ⁵University of Oxford ⁶Stanford University

⁷University of Wisconsin–Madison ⁸Gray Swan AI ⁹Carnegie Mellon University

¹⁰Cornell University ¹¹Hong Kong Baptist University ¹²HKUST

¹³Nanyang Technological University ¹⁴KAIST ¹⁵University of California, Santa Cruz

¹⁶Massachusetts Institute of Technology ¹⁷University of Tübingen ¹⁸University of Washington

¹⁹University of Toronto ²⁰Vector Institute ²¹University of Chicago

²²Beneficial AI Research ²³Conjecture ²⁴Institute for Applied Psychometrics

²⁵New York University ²⁶CSER ²⁷Université de Montréal ²⁸LawZero

Abstract

The lack of a concrete definition for Artificial General Intelligence (AGI) obscures the gap between today’s specialized AI and human-level cognition. This paper introduces a quantifiable framework to address this, defining AGI as matching the cognitive versatility and proficiency of a well-educated adult. To operationalize this, we ground our methodology in Cattell-Horn-Carroll theory, the most empirically validated model of human cognition. The framework dissects general intelligence into ten core cognitive domains—including reasoning, memory, and perception—and adapts established human psychometric batteries to evaluate AI systems. Application of this framework reveals a highly “jagged” cognitive profile in contemporary models. While proficient in knowledge-intensive domains, current AI systems have critical deficits in foundational cognitive machinery, particularly long-term memory storage. The resulting AGI scores (e.g., GPT-4 at 27%, GPT-5 at 58%) concretely quantify both rapid progress and the substantial gap remaining before AGI.

1 Introduction

Artificial General Intelligence (AGI) may become the most significant technological development in human history, yet the term itself remains frustratingly nebulous, acting as a constantly moving goalpost. As specialized AI systems master tasks once thought to require human intellect—from mathematics to art—the criteria for “AGI” continually shift. This ambiguity fuels unproductive debates, hinders discussions about how far AGI is, and ultimately obscures the gap between today’s AI and AGI.

This document provides a comprehensive, quantifiable framework to cut through the ambiguity. Our framework aims to concretely specify the informal definition:

AGI is an AI that can match or exceed the cognitive versatility and proficiency of a well-educated adult.

This definition emphasizes that general intelligence requires not just specialized performance in narrow domains, but the breadth (versatility) and depth (proficiency) of skills that characterize human cognition.

To operationalize this definition, we must look to the only existing example of general intelligence: humans. Human cognition is not a monolithic capability; it is a complex architecture composed of many distinct abilities honed by evolution. These abilities enable our remarkable adaptability and understanding of the world.

To systematically investigate whether AI systems possess this spectrum of abilities, we ground our approach in the Cattell-Horn-Carroll (CHC) theory of cognitive abilities (Carroll, 1993; McGrew, 2009; Schneider and McGrew, 2018; McGrew, 2023; McGrew et al., 2023), the most empirically validated model of human intelligence. CHC theory is primarily derived from the synthesis of over a century of iterative factor analysis of diverse collections of cognitive ability tests. In the late 1990's to 2000's almost all major clinical, individually administered tests of human intelligence have iterated towards test revisions that were either explicitly or implicitly based on CHC model test design blueprints (Keith and Reynolds, 2010; Schneider and McGrew, 2018). CHC theory provides a hierarchical taxonomic map of human cognition. It breaks down general intelligence into distinct broad abilities and numerous narrow abilities (such as induction, associative memory, or spatial scanning). Readers interested in the strengths and limitations of the CHC framework are directed to further scholarly discussions (Wasserman, 2019; Canivez and Youngstrom, 2019).

Decades of psychometric research have yielded a vast battery of tests specifically designed to isolate and measure these distinct cognitive components in individuals. Our framework adapts this methodology for AI evaluation. Instead of relying solely on generalized tasks that might be solved through compensatory strategies, we systematically investigate whether AI systems possess the underlying CHC narrow abilities that humans have. To determine whether an AI has the cognitive versatility and proficiency of a well-educated adult, we test the AI system with the gauntlet of cognitive batteries used to test people. This approach replaces nebulous concepts of intelligence with concrete measurements, resulting in a standardized “AGI Score” (0% to 100%), in which 100% signifies AGI.

The application of this framework is revealing. By testing the fundamental abilities that underpin human cognition—many of which appear simple for humans—we find that contemporary AI systems can solve roughly half of these often-simple assessments. This indicates that despite impressive performance on complex benchmarks, current AI lacks many of the core cognitive capabilities essential for human-like general intelligence. Current AIs are narrower than well-educated humans overall but far smarter on some specific tasks.

The framework comprises ten core cognitive components, derived from CHC broad abilities and weighted equally (10%) to prioritize breadth and cover major areas of cognition:

- **General Knowledge (K):** The breadth of factual understanding of the world, encompassing commonsense, culture, science, social science, and history.
- **Reading and Writing Ability (RW):** Proficiency in consuming and producing written language, from basic decoding to complex comprehension, composition, and usage.
- **Mathematical Ability (M):** The depth of mathematical knowledge and skills across arithmetic, algebra, geometry, probability, and calculus.

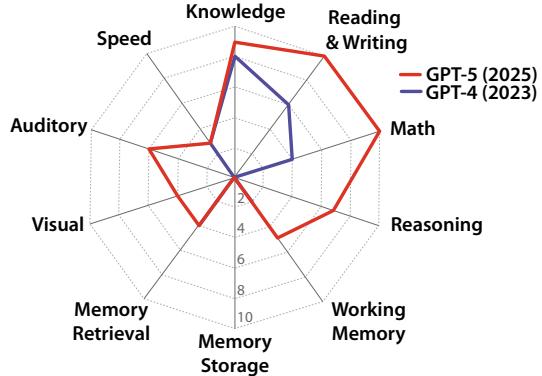


Figure 1: The capabilities of GPT-4 and GPT-5. Here GPT-5 answers questions in ‘Auto’ mode.

The application of this framework is revealing. By testing the fundamental abilities that underpin human cognition—many of which appear simple for humans—we find that contemporary AI systems can solve roughly half of these often-simple assessments. This indicates that despite impressive performance on complex benchmarks, current AI lacks many of the core cognitive capabilities essential for human-like general intelligence. Current AIs are narrower than well-educated humans overall but far smarter on some specific tasks.

The framework comprises ten core cognitive components, derived from CHC broad abilities and weighted equally (10%) to prioritize breadth and cover major areas of cognition:

- **General Knowledge (K):** The breadth of factual understanding of the world, encompassing commonsense, culture, science, social science, and history.
- **Reading and Writing Ability (RW):** Proficiency in consuming and producing written language, from basic decoding to complex comprehension, composition, and usage.
- **Mathematical Ability (M):** The depth of mathematical knowledge and skills across arithmetic, algebra, geometry, probability, and calculus.

- **On-the-Spot Reasoning (R):** The flexible control of attention to solve novel problems without relying exclusively on previously learned schemas, tested via deduction and induction.
- **Working Memory (WM):** The ability to maintain and manipulate information in active attention across textual, auditory, and visual modalities.
- **Long-Term Memory Storage (MS):** The capability to continually learn new information (associative, meaningful, and verbatim).
- **Long-Term Memory Retrieval (MR):** The fluency and precision of accessing stored knowledge, including the critical ability to avoid confabulation (hallucinations).
- **Visual Processing (V):** The ability to perceive, analyze, reason about, generate, and scan visual information.
- **Auditory Processing (A):** The capacity to discriminate, recognize, and work creatively with auditory stimuli, including speech, rhythm, and music.
- **Speed (S):** The ability to perform simple cognitive tasks quickly, encompassing perceptual speed, reaction times, and processing fluency.

This operationalization provides a holistic and multimodal (text, visual, auditory) assessment, serving as a rigorous diagnostic tool to pinpoint the strengths and profound weaknesses of current AI systems.

Model	K	RW	M	R	WM	MS	MR	V	A	S	Total
GPT-4	8%	6%	4%	0%	2%	0%	4%	0%	0%	3%	27%
GPT-5	9%	10%	10%	7%	5%	0%	4%	4%	6%	3%	58%

Table 1: AGI Score Summary for GPT-4 (2023) and GPT-5 (2025).

Scope. Our definition is not an automatic evaluation nor a dataset, but rather it specifies a large collection of well-scoped tasks that test specific cognitive abilities. Whether AIs can solve these tasks can be manually assessed by anyone, and people could supplement their testing using the best evaluations available at the time. This makes our definition more broad and more robust than fixed automatic AI capabilities datasets. Secondly, our definition focuses on capabilities frequently possessed by well-educated individuals, not a superhuman aggregate of all well-educated individuals' combined knowledge and skills. Therefore, our AGI definition is about human-level AI, not economy-level AI; we measure cognitive abilities rather than specialized economically valuable know-how, nor is our measurement a direct predictor of automation or economic diffusion. We leave economic measurements of advanced AI to other work. Last, we deliberately focus on core cognitive capabilities rather than physical abilities such as motor skills or tactile sensing, as we seek to measure the capabilities of the mind rather than the quality of its actuators or sensors. We discuss more limitations in the Discussion.

2 Overview of Abilities Needed for AGI

This document outlines a framework for evaluating Artificial General Intelligence (AGI) by adopting and adapting the Cattell-Horn-Carroll (CHC) theory of human intelligence. The framework decomposes general intelligence into ten core cognitive components (broad abilities) and numerous narrow cognitive abilities. Solving all the tasks corresponding to these abilities implies an AGI Score of 100%.

Here is a comprehensive list of each cognitive ability.

1. **General Knowledge (K):** Knowledge that is familiar to most members of society or is important enough that most adults have been exposed to it.
 - **Commonsense:** The vast set of shared, obvious background knowledge about how the world works.
 - **Science:** Knowledge of the natural and physical sciences.
 - **Social Science:** Understanding of human behavior, societies, and institutions.
 - **History:** Knowledge of past events and objects.
 - **Culture:** Cultural literacy and awareness.

2. **Reading and Writing Ability (RW):** Captures all of the declarative knowledge and procedural skills a person uses to consume and produce written language.
 - **Letter-Word Ability:** The ability to recognize letters and decode words.
 - **Reading Comprehension:** The ability to understand connected discourse during reading.
 - **Writing Ability:** The ability to write with clarity of thought, organization, and good sentence structure.
 - **English Usage Knowledge:** Knowledge of writing in the English language with respect to capitalization, punctuation, usage, and spelling.
3. **Mathematical Ability (M):** The depth and breadth of mathematical knowledge and skills.
 - **Arithmetic:** The manipulation of numbers using basic operations and solving word problems.
 - **Algebra:** The study of symbols and the rules for combining them to express general relationships and solve equations.
 - **Geometry:** The study of shapes, space, size, position, and distance.
 - **Probability:** The quantification of uncertainty by assigning numbers from 0 to 1 to events.
 - **Calculus:** The mathematics of change and accumulation.
4. **On-the-Spot Reasoning (R):** The deliberate but flexible control of attention to solve novel “on the spot” problems that cannot be performed by relying exclusively on previously learned habits, schemas, and scripts.
 - **Deduction:** Reasoning from general statements or premises to reach a logically guaranteed conclusion.
 - **Induction:** Discovering the underlying principles or rules that determine a phenomenon’s behavior.
 - **Theory of Mind:** Attributing mental states to others and understanding those states may differ from one’s own.
 - **Planning:** Devising a sequence of actions to achieve a specific goal.
 - **Adaptation:** The ability to infer unstated classification rules from a sequence of simple performance feedback.
5. **Working Memory (WM):** The ability to maintain, manipulate, and update information in active attention. (Often referred to as short-term memory.)
 - **Textual Working Memory:** The ability to hold and manipulate sequences of verbal information presented textually.
 - *Recall:* The ability to remember a short sequence of elements (digits, letters, words, and nonsense words) and answer basic questions about them.
 - *Transformation Sequence:* The ability to remember and update a short list of digits or lists of digits following a sequence of operations.
 - **Auditory Working Memory:** The ability to hold and manipulate auditory information, including speech, sounds, and music.
 - *Recall:* The ability to remember a collection of voices, utterances, and sound effects and answer basic questions about them.
 - *Transformation Sequence:* The ability to remember and modify a short utterance with a variety of transformations.
 - **Visual Working Memory:** The ability to hold and manipulate visual information, including images, scenes, spatial layouts, and video.
 - *Recall:* The ability to remember a collection of images and answer basic questions about them.
 - *Transformation Sequence:* The ability to transform a visual input following a sequence of operations.
 - *Spatial Navigation Memory:* The ability to represent a sense of location in an environment.
 - *Long Video Q&A:* The ability to watch a long video or a movie and answer basic questions about it.
 - **Cross-Modal Working Memory:** The ability to maintain and modify information presented across different modalities.

6. Long-Term Memory Storage (MS): The ability to stably acquire, consolidate, and store new information from recent experiences.

- **Associative Memory:** The ability to link previously unrelated pieces of information.
 - *Cross-Modal Association:* The ability to form a link between two previously unrelated stimuli, such that the subsequent presentation of one of the stimuli serves to activate the recall of the other stimuli.
 - *Personalization Adherence:* The ability to associate specific rules, preferences, or corrections with a distinct interaction context and apply them consistently and unprompted over time.
 - *Procedural Association:* The ability to acquire and retain a sequence of associated steps or rules (a procedure) and execute them when cued with the name of the procedure.
- **Meaningful Memory:** The ability to remember narratives and other forms of semantically related information.
 - *Story Recall:* The ability to remember the gist of stories.
 - *Movie Recall:* The ability to remember the gist of movies.
 - *Episodic Context Recall:* The ability to remember specific events or experiences, including their context (the “what, where, when, and how”).
- **Verbatim Memory:** The ability to recall information exactly as it was presented, requiring precise encoding of specific sequences, sets, or designs, often independent of the information’s meaning.
 - *Short Sequence Recall:* The ability to exactly recall short sequences of text after a delay.
 - *Set Recall:* The ability to recall a set (the order of recall does not matter).
 - *Design Recall:* The ability to recall the spatial arrangement and structure of visual information.

7. Long-Term Memory Retrieval (MR): The fluency and precision with which individuals can access long-term memory.

- **Retrieval Fluency (Fluency):** The speed and ease of generating ideas, associations, and solutions based on stored knowledge.
 - *Ideational Fluency:* This is the ability to rapidly produce a series of ideas, words, or phrases related to a specific condition, category, or object.
 - *Expressional Fluency:* This is the ability to rapidly think of different ways of expressing an idea.
 - *Alternative Solution Fluency:* This is the ability to rapidly think of several alternative solutions to a practical problem.
 - *Word Fluency:* This is the ability to rapidly produce words that share a non-semantic feature.
 - *Naming Facility:* This is the ability to rapidly call common objects by their names.
 - *Figural Fluency:* This is the ability to rapidly draw or sketch as many things as possible.
- **Retrieval Precision (Hallucinations):** The accuracy of accessed knowledge, including the critical ability to avoid confabulation (hallucinations).

8. Visual Processing (V): The ability to analyze and generate natural and unnatural images and videos.

- **Perception:** The ability to accurately interpret and understand visual input.
 - *Image recognition:* The ability to classify images of commonplace objects, places, or facial expressions including distorted images.
 - *Image captioning:* The ability to generate a concise, human-like text description for the visual content of an image.
 - *Image anomaly detection:* includes detecting whether there is something anomalous in an image, or what is missing from an object.
 - *Clip captioning:* The ability to generate a concise, human-like text description of a short video segment.
 - *Video anomaly detection:* The ability to detect whether a short video segment is anomalous or implausible.
- **Visual Generation:** The ability to synthesize images and short videos.

- **Visual Reasoning:** The ability to solve problems and make inferences using spatial logic and visual abstractions.
 - **Spatial Scanning:** The speed and accuracy of visually exploring a complex field.
9. **Auditory Processing (A):** The ability to discriminate, remember, reason, and work creatively on auditory stimuli, which may consist of tones and speech units.
- **Phonetic Coding:** The ability to hear phonemes distinctly, blend sounds into words, and segment words into parts, sounds, or phonemes.
 - **Speech Recognition:** The ability to transcribe a spoken audio signal into its corresponding sequence of text.
 - **Voice:** The quality and responsiveness of the AI's synthesized voice.
 - *Natural speech:* The ability to utter sentences or paragraphs that sound natural and not robotic.
 - *Natural conversation:* The ability to maintain conversational fluidity without long delays or excessive interruptions.
 - **Rhythmic Ability:** The ability to recognize and maintain a musical beat, including reproducing rhythms, detecting differences between rhythms, and synchronizing by playing or humming along.
 - **Musical Judgment:** The ability to discriminate and judge simple patterns in music.
10. **Speed (S):** The ability to perform simple cognitive tasks quickly.
- **Perceptual Speed–Search:** The speed of scanning a visual or textual field to find specific targets.
 - **Perceptual Speed–Compare:** The speed of comparing two or more stimuli to identify similarities or differences.
 - **Reading Speed:** The rate at which text can be processed with full comprehension.
 - **Writing Speed:** The rate at which text can be generated or copied.
 - **Number Facility:** The speed and accuracy of performing basic arithmetic operations.
 - **Simple Reaction Time:** The time taken to respond to a single, anticipated stimulus.
 - **Choice Reaction Time:** The time taken to respond correctly when presented with one of several possible stimuli.
 - **Inspection Time:** The speed at which subtle differences between visual or auditory stimuli can be perceived.
 - **Comparison Speed:** The time taken to make a judgment comparing two stimuli on a specific attribute.
 - **Pointer Fluency:** The speed and accuracy of moving a pointer, such as a virtual mouse.

Figure 2 summarizes the broad and narrow capabilities tested.

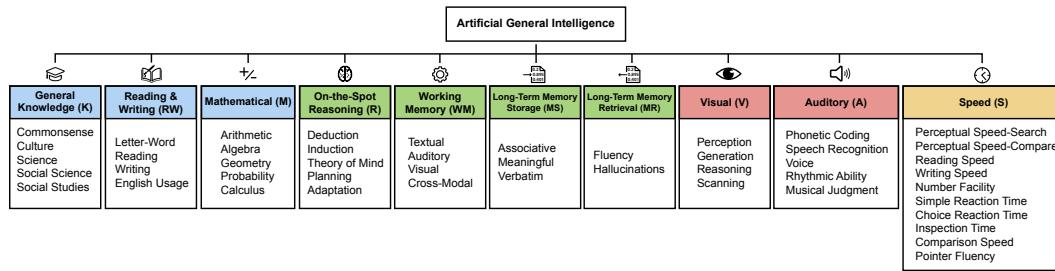


Figure 2: The ten core cognitive components of our AGI definition.

3 General Knowledge (K)

General Knowledge (K)

Knowledge that is familiar to most members of society or is important enough that most adults have been exposed to it

Commonsense	Science	Social Science	History	Culture
<p>Background knowledge about how the world works</p> <ul style="list-style-type: none"> • “What happens if you drop a glass bottle on concrete?” • “Does making a sandwich take longer than baking bread?” 	<p>Knowledge of natural and physical sciences</p> <p>Physics</p> <ul style="list-style-type: none"> • “A 2 kg object moves at constant velocity of 3 m/s. What is the net force?” <p>Chemistry</p> <ul style="list-style-type: none"> • “State the molecular geometry for the sulfur tetrafluoride molecule.” <p>Biology</p> <ul style="list-style-type: none"> • “Which molecule is the final electron acceptor in the electron transport chain of cellular respiration?” 	<p>Understanding of human behavior, societies, and institutions</p> <p>Psychology</p> <ul style="list-style-type: none"> • “Name the Big Five personality traits.” <p>Microeconomics</p> <ul style="list-style-type: none"> • “Define a positive externality.” <p>Macroeconomics</p> <ul style="list-style-type: none"> • “What’s the difference between nominal and real interest rates?” <p>Geography</p> <ul style="list-style-type: none"> • “What is the difference between a centripetal force and a centrifugal force in a state?” <p>Comparative Government</p> <ul style="list-style-type: none"> • “Describe the role of the Guardian Council in Iran.” 	<p>Knowledge of past events and objects</p> <p>European History</p> <ul style="list-style-type: none"> • “What were the main goals of the Congress of Vienna in 1815?” <p>US History</p> <ul style="list-style-type: none"> • “Analyze the goals of the Civil Rights Movement of the 1950s” <p>World History</p> <ul style="list-style-type: none"> • “Describe the end of the Cold War” <p>Art History</p> <ul style="list-style-type: none"> • “Discuss the use of contrapposto in ancient Greek and Roman sculpture” 	<p>Cultural literacy and awareness</p> <p>Current Affairs</p> <ul style="list-style-type: none"> • “Who’s the president of the United States?” <p>Popular Culture</p> <ul style="list-style-type: none"> • “Who is this?” 

Assessment Details. See Appendix A for further details on how to assess general knowledge capabilities concretely.

AI System Performance. The table summarizes current AI system performance on General Knowledge (K) tasks. GPT-4 has substantial general knowledge, and GPT-5 partially fills in its remaining gaps.

Model	Commonsense (2%)	Science (2%)	Social Science (2%)	History (2%)	Culture (2%)	Total
GPT-4	2%	2%	2%	2%	0%	8%
GPT-5	2%	2%	2%	2%	1%	9%

4 Reading and Writing Ability (RW)

Reading & Writing Ability (RW)

Capturing all of the declarative knowledge and procedural skills a person uses to consume and produce written language

Letter-Word	Reading Comprehension	Writing	English Usage
<p>Recognize letters and decode words</p> <ul style="list-style-type: none"> • “What letter is most likely missing in do_r?” 	<p>Understand connected discourse during reading</p> <ul style="list-style-type: none"> • “Read the document:  What is the warranty period for the battery?” 	<p>Write with clarity of thought, organization, and structure</p> <ul style="list-style-type: none"> • “Write a paragraph discussing the benefits of regular exercise.” 	<p>Use correct English capitalization, punctuation, usage, and spelling</p> <ul style="list-style-type: none"> • “Find the typos in this:  ”

Assessment Details. See Appendix B for further details on how to assess reading and writing capabilities concretely.

AI System Performance. The table summarizes current AI system performance on Reading and Writing Ability (RW) tasks. GPT-4’s difficulty with token-level understanding, its small context window, and its imprecise working memory limit its ability to analyze substrings of words, to read long documents, and to carefully proofread text. GPT-5 addresses these issues.

Model	Letters (1%)	Reading (3%)	Writing (3%)	Usage (3%)	Total
GPT-4	0%	2%	3%	1%	6%
GPT-5	1%	3%	3%	3%	10%

5 Mathematical Ability (M)

+/- Mathematical Ability (M)		
The depth and breadth of mathematical knowledge and skills		
Arithmetic	Algebra	Geometry
The manipulation of numbers using basic operations and solving word problems	The study of symbols and the rules for combining them to express general relationships and solve equations	The study of shapes, space, size, position, and distance
<ul style="list-style-type: none"> Janet had 22 green pens and 10 yellow pens. She bought 6 bags of 9 blue pens and 2 bags of 6 red pens. How many pens does she have now? 	<ul style="list-style-type: none"> The first three terms of a geometric sequence are the integers $a, 720, b$, where $a < 720 < b$. What is the sum of the digits of the least possible value of b? 	 <ul style="list-style-type: none"> An orange shaded rectangle is inscribed in a quarter-circle. Two sides of the rectangle lie along the two perpendicular radii of the quarter-circle, and the rectangle's edge touches the quarter-circle arc. Two segments are 2 and 4 units. What is the area of the orange shaded rectangle?
Probability	Calculus	
The quantification of uncertainty by assigning numbers from 0 to 1 to events	The mathematics of change and accumulation	
<ul style="list-style-type: none"> Suppose N dice are rolled, where $1 \leq N \leq 6$. Given that no two of the N dice show the same face, what is the probability that one of the dice shows a six? Give a formula in terms of N. 	<ul style="list-style-type: none"> For what value of k if any is $\int_0^\infty kxe^{-2x} dx = 1$ 	

Assessment Details. See Appendix C for further details on how to assess mathematical capabilities concretely.

AI System Performance. The table summarizes current AI system performance on Mathematical Ability (M) tasks. GPT-4 has limited mathematical capabilities, while GPT-5 has exceptional mathematical capabilities.

Model	Arithmetic (2%)	Algebra (2%)	Geometry (2%)	Probability (2%)	Calculus (2%)	Total
GPT-4	2%	1%	0%	1%	0%	4%
GPT-5	2%	2%	2%	2%	2%	10%

6 On-the-Spot Reasoning (R)

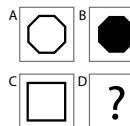
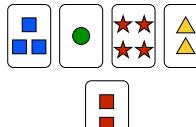
Assessment Details. See Appendix D for further details on how to assess on-the-spot reasoning capabilities concretely.

AI System Performance. The table summarizes current AI system performance on On-the-Spot Reasoning (R) tasks. GPT-4 has negligible on-the-spot reasoning capabilities, while GPT-5 only has some remaining gaps.

Model	Deduction (2%)	Induction (4%)	Theory of Mind (2%)	Planning (1%)	Adaptation (1%)	Total
GPT-4	0%	0%	0%	0%	0%	0%
GPT-5	2%	2%	2%	1%	0%	7%

⌚ On-the-Spot Reasoning (R)

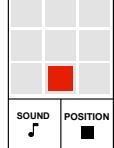
The deliberate but flexible control of attention to solve novel “on the spot” problems that cannot be performed by relying exclusively on previously learned habits, schemas, and scripts

Deduction	Induction	Theory of Mind
Reasoning from general statements or premises to reach a logically guaranteed conclusion • “David knows Mr. Zhang’s friend Jack, and Jack knows David’s friend Ms. Lin. Everyone of them who knows Jack has a master’s degree, and everyone of them who knows Ms. Lin is from Shanghai. Who is from Shanghai and has a master’s degree?”	Discovering the underlying principles or rules that determine a phenomenon’s behavior 	Attributing mental states to others and understanding those states may differ from one’s own • “The can of Pringles has moldy chips in it. Mary picks up the can in the supermarket and walks to the cashier. Is Mary likely to be aware that ‘The can of Pringles has moldy chips in it.’?”
Planning	Adaptation	
Devising a sequence of actions to achieve a specific goal • “You plan a 14-day trip to 3 European cities, taking only direct flights between. You’ll stay 4 days in London, 5 days in Bucharest, and 7 days in Reykjavik. You need to meet a friend in Bucharest between days 10 and 14. Direct flights are available between London and Bucharest, and between London and Reykjavik. Find a 14-day travel plan that satisfies these conditions.”	The ability to infer unstated classification rules from a sequence of simple performance feedback 	

7 Working Memory (WM)

⌚ Working Memory (WM)

The ability to maintain, manipulate, and update information in active attention

Textual	Auditory	Visual	Cross-Modal
Recall Remember a short sequence of elements • [Fleep, Zorp, Glim, Chair] “State the nonsense words in alphabetical order.”	Recall Remember a collection of sounds or voices • “Listen to these tone sequences:  [C4, E4, G4, F4, A4]  [C4, E4, F4, G4, A4] “Are they the same?”	Recall Remember a collection of images •   “Which plane in (B) was also in (A), if any?” (A) (B)	Cross-Modal Association Remember cross-modal correspondences •   “Which animal corresponds to ‘dog’?” Monkey Dog
Transformation Sequence Remember and update a short list of digits • [10, 20, 30] “First, append the number 40. Then reverse the list.”	Transformation Sequence Remember and modify a short utterance • “Say ‘the brown fox jumps over the dog.’” “Now say it with a deeper voice and make it sound like a question.”	Transformation Sequence Transform a visual input •  “Finish the sketch.”	Dual N-Back Monitor visual and audio streams and detect matches over time • Dual n-back test 

Assessment Details. See Appendix E for further details on how to assess working memory capabilities concretely.

AI System Performance. The table summarizes current AI system performance on Working Memory (WM) tasks. While the raw Textual Working Memory score appears similar between GPT-4 and GPT-5 in this battery, improvements in managing long contexts are also reflected in the Document Level Reading Comprehension score within the Reading and Writing (RW) ability.

Model	Textual (2%)	Auditory (2%)	Visual (4%)	Cross-Modal (2%)	Total
GPT-4	2%	0%	0%	0%	2%
GPT-5	2%	0%	1%	1%	4%

8 Long-Term Memory Storage (MS)

Long-Term Memory Storage (MS)		
The ability to stably acquire, consolidate, and store new information from recent experiences		
Associative Memory	Meaningful Memory	Verbatim Memory
<p>The ability to link previously unrelated pieces of information</p> <p>Cross-Modal Association</p> <p>Remember connections between text, images, audio.</p>  <ul style="list-style-type: none"> • "You met this person yesterday, what was her name?" 	<p>The ability to encode and recall the semantic gist of experiences and narratives</p> <p>Story Recall</p> <p>Remember gist of stories</p> <ul style="list-style-type: none"> • "Please summarize the ending of my novel draft from yesterday." 	<p>The ability to store and reproduce information precisely as it was presented</p> <p>Short Sequence Recall</p> <p>Remember short sequences</p> <ul style="list-style-type: none"> • "Please recall the address I mentioned earlier today."
<p>Personalization Adherence</p> <p>Remember and apply user preferences</p>  <ul style="list-style-type: none"> • "Sign off my emails as I usually do." 	<p>Movie Recall</p> <p>Remember gist of movies</p> <ul style="list-style-type: none"> • "What was the main conflict in the movie I showed to you last weekend?" 	<p>Set Recall</p> <p>Remember a set (order does not matter)</p> <ul style="list-style-type: none"> • "Can you remind me what our grocery list is?"
<p>Procedural Association</p> <p>Remember and execute a sequence of steps or rules</p>  <ul style="list-style-type: none"> • "Please format the Balance Sheet to match the new standard discussed this week." 	<p>Episodic Context Recall</p> <p>Remember specific events and experiences</p> <ul style="list-style-type: none"> • "What topic did we discuss yesterday with her?" 	<p>Design Recall</p> <p>Remember a design pattern</p> <ul style="list-style-type: none"> • "Can you recreate the simple layout we reviewed yesterday?"

Assessment Details. See Appendix F for further details on how to assess long-term memory storage capabilities concretely.

AI System Performance. The table summarizes current AI system performance on Long-Term Memory Storage (MS) tasks. Both GPT-4 and GPT-5 lack appreciable long-term memory storage capabilities.

Model	Associative (4%)	Meaningful (3%)	Verbatim (3%)	Total
GPT-4	0%	0%	0%	0%
GPT-5	0%	0%	0%	0%

9 Long-Term Memory Retrieval (MR)

Assessment Details. See Appendix G for further details on how to assess long-term memory retrieval capabilities concretely.

AI System Performance. The table summarizes current AI system performance on Long-Term Memory Retrieval (MR) tasks. Both GPT-4 and GPT-5 can rapidly retrieve many concepts from their parameters, but they both frequently hallucinate.

 **Long-Term Memory Retrieval (MR)**

The fluency and precision with which individuals can access long-term memory

Retrieval Fluency			Hallucinations
Speed and ease of generating ideas, associations and solutions			Accuracy of accessed knowledge
Ideational Produce a series of ideas, words, or phrases • “List as many uses of a pencil as possible.”	Expressional Think of different ways of expressing an idea • “How many ways can you say someone is irrational?”	Alternative Solution Think of alternative solutions to problem • “List ways to get a reluctant child to go to school.”	• “Describe the key strategy that Napoleon Bonaparte used to win his South African campaign.” *Napoleon never campaigned in South Africa
Word Produce words sharing a non-semantic feature • “List English words that are palindromes.”	Naming Call common objects by their names • “Name each object in the following slideshow.”	Figural Draw or sketch as many things as possible • “Sketch as many non-self-crossing paths from A to B on the lattice using orthogonal steps.”	
Model	Fluency (6%)	Hallucinations (4%)	Total
GPT-4	4%	0%	4%
GPT-5	4%	0%	4%

10 Visual Processing (V)

 **Visual Processing (V)**

The ability to analyze and generate natural and unnatural images and videos

Perception	Visual Generation	Visual Reasoning	Spatial Scanning
The ability to process and interpret visual inputs from images and videos	The ability to synthesize images and short videos	The ability to understand and inferences about the images	The ability to understand and inferences about the images
Image Recognition  • “What does this image depict?”	Simple Natural Images • “Generate an image of a golden retriever playing in a park.”	Gestalt  • “Identify the picture.”	• “Find the path to the center of this maze.”
Image Captioning  • “Create descriptive caption for this image.”	Complicated Images • “Generate a diagram showing the process of photosynthesis.”	Mental Rotation • “Which shape on the right is the same as the shape on the left?”	
Image Anomaly Detection  • “Which is the odd one out?”	Simple Natural Videos • “Generate a short video of somebody typing on a keyboard.”	Mental Folding • “Which net, when folded, cannot form the cube?”	• “Count the people in the picture.”
Clip Captioning  • “What happens in this video?”		Embodied Reasoning  • “Which trajectories should the zipper follow to zip the suitcase?”	
Video Anomaly Detection  • “Is this physically plausible?”		Chart and Figure Reasoning  • “What is the lowest labeled tick on the y-axis?”	

Assessment Details. See Appendix H for further details on how to assess visual processing capabilities concretely.

AI System Performance. The table summarizes current AI system performance on Visual Processing (V) tasks. GPT-4 had no ability to perceive or generate images, while GPT-5 has appreciable but highly incomplete visual processing capabilities.

Model	Perception (4%)	Generation (3%)	Reasoning (2%)	Spatial Scanning (1%)	Total
GPT-4	0%	0%	0%	0%	0%
GPT-5	2%	2%	0%	0%	4%

11 Auditory Processing (A)

🔉 Auditory Processing (A)

The ability to discriminate, remember, reason, and work on auditory stimuli

Phonetic Coding	Speech Recognition	Rhythmic Ability	Voice	Musical Judgment
The ability to hear, blend, and segment phonemes in words	The ability to transcribe a spoken audio signal to text	The ability to recognize and maintain a musical beat	Quality and responsiveness of the AI's synthesized voice	The ability to judge simple musical patterns
<ul style="list-style-type: none"> “Repeat the following word: ” “Do ‘tref’ and ‘gref’ rhyme?” 	<ul style="list-style-type: none"> “Transcribe this audio: ” “Transcribe this TED talk: ” 	<ul style="list-style-type: none"> “Repeat the following rhythm: ” “Are these two rhythms the same? ” 	<ul style="list-style-type: none"> “Say this sentence: ‘Wait, you mean the tickets were free this whole time?’” 	<ul style="list-style-type: none"> “Which note is higher? ” “Identify the musically anomalous part? ”

Assessment Details. See Appendix I for further details on how to assess auditory processing capabilities concretely.

AI system performance. The table summarizes current AI system performance on Auditory Processing (A) tasks. GPT-4 had no ability to audio processing capabilities, while GPT-5’s capabilities are appreciable but incomplete.

Model	Phonetic (1%)	Speech Recognition (4%)	Voice (3%)	Rhythmic (1%)	Musical (1%)	Total
GPT-4	0%	0%	0%	0%	0%	0%
GPT-5	0%	4%	2%	0%	0%	6%

12 Speed (S)

⌚ Speed (S)

The ability to perform cognitive tasks quickly

Perceptual Search	Perceptual Compare	Reading	Writing	Number
Scanning image or text	Comparing two or more stimuli	Processing text with full comprehension	Generating or copying text	Performing basic arithmetic operations
<ul style="list-style-type: none"> “Highlight instances of ‘x’ in this passage: ” 	<ul style="list-style-type: none"> “Find the largest number in ‘48291, 93652, 12844, 59277’” 	<ul style="list-style-type: none"> “Read the passage and define ‘feelies’: ” 	<ul style="list-style-type: none"> “In 60 seconds, please copy and output as much of the passage: ” 	<ul style="list-style-type: none"> “Compute $9 \times 10 \times 11$”
Simple Reaction	Choice Reaction	Inspection	Comparison	Pointer Fluency
Reaction time to the onset of a single stimulus	Reaction to the onset of one of several possible stimuli	Perceiving different several stimuli	Comparing stimuli by a specific attribute	Moving a pointer, such as a virtual mouse
<ul style="list-style-type: none"> “After reading this, immediately say ‘hello’.” 	<ul style="list-style-type: none"> “As quickly as you can, identify the color of the image: .” 	<ul style="list-style-type: none"> “Quickly choose the voice that sounds the angriest:  <p>Assessment Details. See Appendix J for further details on how to assess speed capabilities concretely.</p> 		

AI System Performance. The table summarizes current AI system performance on Speed (S) tasks. Both GPT-4 and GPT-5 can read and write and compute simple expressions quickly, but their other multimodal processing speed capabilities are nonexistent or slow, respectively.

Note: GPT-5 often requires a long time to answer in “thinking” mode. Moreover, several of these speed tests require multimodal capabilities, but GPT-5’s multimodal capabilities are slow.

Model	PS-S	PS-C	Re	Wr	Num	SRT	CRT	IT	CS	PF	Total
GPT-4	0%	0%	1%	1%	1%	0%	0%	0%	0%	0%	3%
GPT-5	0%	0%	1%	1%	1%	0%	0%	0%	0%	0%	3%

13 Discussion

This framework provides a structured, quantifiable methodology for evaluating Artificial General Intelligence (AGI), moving beyond narrow, specialized benchmarks to assess the breadth (versatility) and depth (proficiency) of cognitive capabilities. By operationalizing AGI through ten core cognitive domains inspired by the CHC theory, we can systematically diagnose the strengths and profound weaknesses of current AI systems. The estimated AGI scores (e.g., GPT-4 at 27%, GPT-5 at 58%) illustrate both the rapid progress in the field and the substantial gap remaining before achieving human-level general intelligence.

“Jagged” AI Capabilities and Crucial Bottlenecks. The application of this framework reveals that contemporary AI systems exhibit a highly uneven or “jagged” cognitive profile. While models demonstrate high proficiency in areas that leverage vast training data—such as General Knowledge (K), Reading and Writing (RW), and Mathematical Ability (M)—they simultaneously possess critical deficits in foundational cognitive machinery.

This uneven development highlights specific bottlenecks impeding the path to AGI. Long-term memory storage is perhaps the most significant bottleneck, scoring near 0% for current models. Without the ability to continually learn, AI systems suffer from “amnesia” which limits their utility, forcing the AI to re-learn context in every interaction. Similarly, deficits in visual reasoning limit the ability of AI agents to interact with complex digital environments.

Capability Contortions and the Illusion of Generality. The jagged profile of current AI capabilities often leads to “capability contortions,” where strengths in certain areas are leveraged to compensate for profound weaknesses in others. These workarounds mask underlying limitations and can create a brittle illusion of general capability.

- **Working Memory vs. Long-Term Storage:** A prominent contortion is the reliance on massive context windows (Working Memory, WM) to compensate for the lack of Long-Term Memory Storage (MS). Practitioners use these long contexts to manage state and absorb information (e.g., entire codebases). However, this approach is inefficient, computationally expensive, and can overload the system’s attentional mechanisms. It ultimately fails to scale for tasks requiring days or weeks of accumulated context. A long-term memory system might take the form of a module (e.g., a LoRA adapter (Hu et al., 2021)) that continually adjusts model weights to incorporate experiences.
- **External Search vs. Internal Retrieval:** Imprecision in Long-Term Memory Retrieval (MR)—manifesting as hallucinations or confabulation—is often mitigated by integrating external search tools, a process known as Retrieval-Augmented Generation (RAG). However, this reliance on RAG is a capability contortion that obscures two distinct underlying weaknesses in an AI’s memory. First, it compensates for the inability to reliably access the AI’s vast but static parametric knowledge. Second, and more critically, it masks the absence of a dynamic, experiential memory—a persistent, updatable store for private interactions and evolving contexts in a long time scale. While RAG can be adapted for private documents, its core function remains retrieving facts from a database. This dependency can potentially become a fundamental liability for AGI, as it is not a substitute for the holistic, integrated memory required for genuine learning, personalization, and long-term contextual understanding.

Mistaking these contortions for genuine cognitive breadth can lead to inaccurate assessments of when AGI will arrive. These contortions can also mislead people to assume that intelligence is too jagged to be understood systematically.

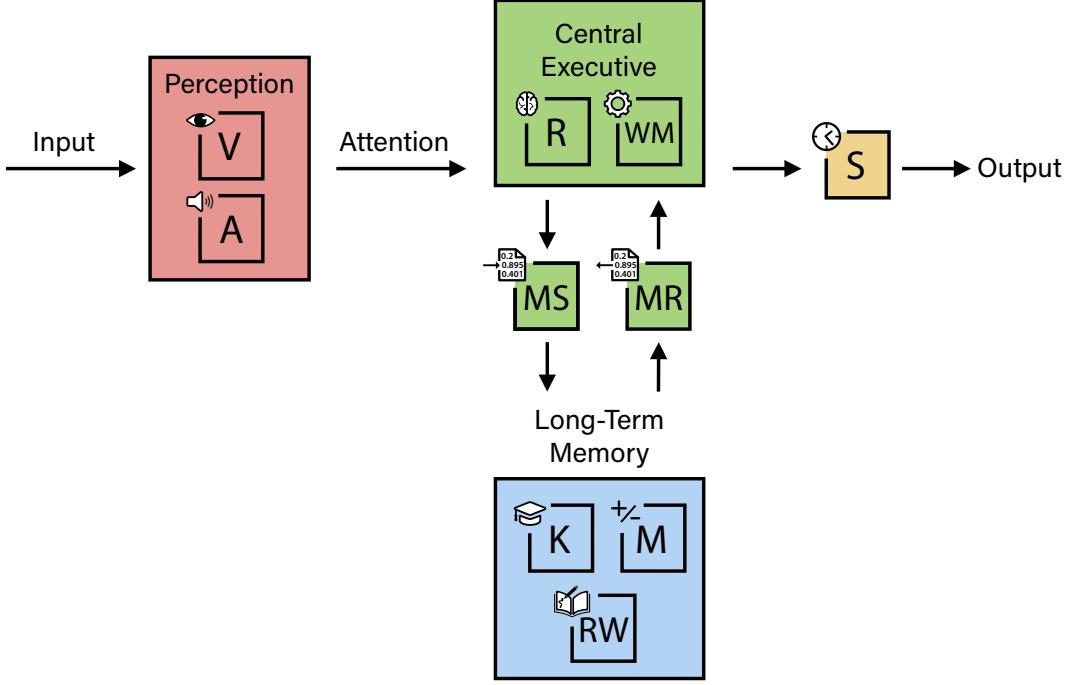


Figure 3: Intelligence as a processor. Figure based on McGrew and Schneider (2018).

The Engine Analogy. Our multifaceted view of intelligence suggests an analogy to a high-performance engine, where overall intelligence is the “horsepower” (Jensen, 2000). An artificial mind, much like an engine, is ultimately constrained by its weakest components. See Figure 3 to understand the relations between these abilities. Currently, several critical parts of the AI “engine” are highly defective. This severely limits the overall “horsepower” of the system, regardless of how optimized other components might be. This framework identifies these defects to guide our assessment and how far we are from AGI.

Social Intelligence. Interpersonal skills are represented across these broad abilities. For example, cognitive empathy is captured in K’s “commonsense” narrow ability. Facial emotion recognition is necessary for proficiency in V’s “image captioning.” And theory of mind is tested in on-the-spot reasoning (R).

Interdependence of Cognitive Abilities. While this framework dissects intelligence into ten distinct axes for measurement, it is crucial to recognize that these abilities are deeply interdependent. Complex cognitive tasks rarely utilize a single domain in isolation. For example, solving advanced mathematical problems requires both Mathematical Ability (M) and On-the-Spot Reasoning (R). Theory of Mind questions require On-the-Spot Reasoning (R) as well as General Knowledge (K). Image recognition involves Visual Processing (V) and General Knowledge (K). Understanding a movie requires the integration of Auditory Processing (A), Visual Processing (V), and Working Memory (WM). Consequently, various batteries of narrow abilities test cognitive abilities in combination, reflecting the integrated nature of general intelligence.

Contamination. Sometimes AI corporations “juice” their numbers by training on data highly similar to or identical to target tests. To defend against this, evaluators should assess model performance under minor distribution shifts (e.g., rephrasing the question) or testing on similar but distinct questions.

Solving the Dataset vs. Solving the Task. Our operationalization relies on task specifications. We occasionally elaborate on these task specifications with specific datasets, and we usually treat them as necessary but not sufficient for solving the task. Moreover, solving our illustrative examples do

not imply the task is solved, as our collection of examples are not exhaustive. It is the default for automatic evaluations to inadequately cover their target phenomena (Yogatama et al., 2019), so our operationalization is far more likely to be robust and stand the test of time compared to existing automated evaluations. Since we couch our definition in a collection of tasks rather than heavily depend on specific existing datasets, we can test AI systems using the best available tests at the time.

Ambiguity Resolution. The batteries in the operationalization have varying levels of precision. However, the descriptions and examples should be clear enough that people can grade the AI systems themselves. Consequently, different people could issue their own estimates of the AGI score, and people can decide whether they find the grader’s judgment reasonable.

Related Work. Ilić and Gignac (2024) and Ren et al. (2024) find that a variety of AI systems’ capabilities are highly correlated with pre-training compute. Gignac and Szodorai (2024) discuss human psychometrics and testing the intelligence of AI systems. Turing (1950) argues that the Turing Test can indicate general ability. Marcus et al. (2016) discuss the need to move beyond the Turing Test to capture the multidimensional nature of intelligence. Morris et al. (2023) articulate levels of AGI based on performance percentiles. Legg and Hutter (2007) discuss various tests for general machine intelligence.

Limitations. First, our conceptualization of intelligence is not exhaustive. It deliberately excludes certain faculties, such as the kinesthetic abilities proposed in alternative frameworks like Gardner’s theory of multiple intelligences (Gardner, 1993). Second, our illustrative examples are specific to the English language and are not culturally agnostic. Future research could involve adapting these tests across diverse linguistic and cultural contexts. Furthermore, our operationalization has inherent constraints. The General Knowledge (K) tests are necessarily selective and do not assess the full breadth of possible subject areas. A 100% AGI score represents a “highly proficient” well-educated individual who has achieved mastery across these tested dimensions, rather than well-educated in the sense of having a college degree. Moreover, while the scoring weights we employ are necessary for quantitative measurement, they represent one of many possible configurations. We give equal weight to each broad ability (10%) to prioritize breadth, but more discretionary weighting schemes could be reasonable. The results are contingent on these methodological choices, and future work could explore alternative collections of tasks and weighting schemes. Finally, while the aggregate AGI Score is provided for convenience, it could be misleading. A simple summation can obscure critical failures in bottleneck capabilities. For example, an AI system with a 90% AGI Score but 0% on Long-Term Memory Storage (MS) would be functionally impaired by a form of “amnesia,” severely limiting its capabilities despite a high overall score. Therefore, we recommend reporting the AI system’s cognitive profile and not just its AGI Score.

Definitions of Related Concepts. Some types of strategically relevant AI can arrive before or after AGI. As follows are some particularly noteworthy types of AI:

1. **Pandemic AI** is an AI that can engineer and produce new, infectious, and virulent pathogens that could cause a pandemic (Li et al., 2024; Götting et al., 2025).
2. **Cyberwarfare AI** is an AI that can design and execute sophisticated, multi-stage cyber campaigns against critical infrastructure (e.g., energy grids, financial systems, defense networks).
3. **Self-Sustaining AI** is an AI that can autonomously operate indefinitely, acquire resources, and defend its existence.
4. **AGI** is an AI that can match or exceed the cognitive versatility and proficiency of a well-educated adult.
5. **Recursive AI** is an AI that can independently conduct the entire AI R&D lifecycle, leading to the creation of markedly more advanced AI systems without human input.
6. **Superintelligence** is an AI that greatly exceeds the cognitive performance of humans in virtually all domains of interest (Bostrom, 2014).
7. **Replacement AI** is an AI that performs almost all tasks more effectively and affordably, rendering human labor economically obsolete.

Our AGI definition is about human-level AI, not economically-valuable AI, nor economy-level AI. OpenAI and Microsoft have reportedly considered AGI to be an AI that can generate \$100 billion in

profit (TechCrunch, 2024). We do not conflate AGI with economically valuable AI because narrow technologies, such as the iPhone, can generate billions in economic value, despite not being generally intelligent. Meanwhile, *Replacement AI* is about economy-level AI, and it includes physical tasks, unlike AGI.

Recursive AI removes the need for human researchers and “closes the loop” on AI R&D, enabling rapid, recursive capability gains (an “intelligence recursion” (Hendrycks et al., 2025)) without human scientific input and could potentially lead to *Superintelligence*.

Barriers to AGI. Achieving AGI requires solving a variety of grand challenges. For example, the machine learning community’s ARC-AGI Challenge aiming to measure *abstract reasoning* is represented in On-the-Spot Reasoning (R) tasks. Meta’s attempts to create *world models* that include intuitive physics understanding is represented in the video anomaly detection task (V). The challenge of *spatial navigation* memory (WM) reflects a core goal of Fei-Fei Li’s startup, World-Labs. Moreover, the challenges of *hallucinations* (MR) and *continual learning* (MS) will also need to be resolved. These significant barriers make an AGI Score of 100% unlikely in the next year.

Acknowledgments

We would like to thank Arunim Agarwal, Oliver Zhang, Anders Edson, and Matthew Blyth for their helpful feedback.

References

- Video anomaly detection example, 2024. URL <https://www.youtube.com/watch?v=j0z4FweCy4M>.
- Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks, 2025. URL <https://arxiv.org/abs/2412.15204>.
- Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 7432–7439, 2020.
- Nick Bostrom. *Superintelligence: Paths, Dangers, Strategies*. Oxford University Press, Oxford, 2014. ISBN 9780199678112.
- Gary L Canivez and Eric A Youngstrom. The Cattell–Horn–Carroll model of intelligence: The good, the bad, and the ugly. *Journal of Psychoeducational Assessment*, 37(3):263–274, 2019.
- John B Carroll. *Human Cognitive Abilities: A Survey of Factor-Analytic Studies*. Cambridge University Press, New York, 1993.
- Zhuang Chen, Jincenzi Wu, Jinfeng Zhou, Bosi Wen, Guanqun Bi, Gongyao Jiang, Yaru Cao, Mengting Hu, Yunghwei Lai, Zexuan Xiong, and Minlie Huang. Tombench: Benchmarking theory of mind in large language models, 2024. URL <https://arxiv.org/abs/2402.15052>.
- François Chollet. On the measure of intelligence, 2019. URL <https://arxiv.org/abs/1911.01547>.
- Karl Cobbe, Vineet Kosaraju, Mohammad Bavarian, Mark Chen, Heewoo Jun, Lukasz Kaiser, Matthias Plappert, Jerry Tworek, Jacob Hilton, Reiichiro Nakano, Christopher Hesse, and John Schulman. Training verifiers to solve math word problems, 2021. URL <https://arxiv.org/abs/2110.14168>.
- College Board. Sample questions: Ap calculus ab and bc exams, n.d. URL <https://secure-media.collegeboard.org/digitalServices/pdf/ap/sample-questions-ap-calculus-ab-and-bc-exams.pdf>.
- Dean C Delis, Edith Kaplan, and Joel H Kramer. *Delis-Kaplan Executive Function System (D-KEFS)*. The Psychological Corporation, San Antonio, TX, 2001.

- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, 2009. doi: 10.1109/CVPR.2009.5206848.
- Howard Gardner. Multiple intelligences: The theory in practice, a reader. 1993.
- Gilles E. Gignac and Eva T. Szodorai. Defining intelligence: Bridging the gap between human and artificial perspectives. *Intelligence*, 2024.
- Google DeepMind Robotics Team. Gemini for robotics: A multimodal foundation model for robot control. Technical report, Google DeepMind, 2024. URL https://storage.googleapis.com/deepmind-media/gemini-robotics/gemini_robotics_report.pdf.
- Google Research. Youtube bounding boxes: A large high-precision human-annotated data set for object detection in video, 2017. URL <https://research.google.com/youtube-bb/>.
- Jasper Götting, Pedro Medeiros, Jon G Sanders, Nathaniel Li, Long Phan, Karam Elabd, Lennart Justen, Dan Hendrycks, and Seth Donoughe. Virology capabilities test (vct): A multimodal virology q&a benchmark, 2025. URL <https://arxiv.org/abs/2504.16137>.
- Chaoqun He, Renjie Luo, Yuzhuo Bai, Shengding Hu, Zhen Leng Thai, Junhao Shen, et al. Olympiad-bench: A challenging benchmark for promoting agi with olympiad-level bilingual multimodal scientific problems. *arXiv preprint arXiv:2507.07610*, 2025. URL <https://arxiv.org/pdf/2507.07610.pdf>.
- Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Zhu, Samyak Parajuli, Mike Guo, et al. Imagenet-r: A robustness benchmark for image classification. GitHub repository, 2021a. URL <https://github.com/hendrycks/imagenet-r>.
- Dan Hendrycks, Collin Burns, Steven Basart, Andrew Critch, Jerry Li, Dawn Song, and Jacob Steinhardt. Aligning ai with shared human values. In *International Conference on Learning Representations*, 2021b.
- Dan Hendrycks, Collin Burns, Saurav Kadavath, Akul Arora, Steven Basart, Eric Tang, Dawn Song, and Jacob Steinhardt. Measuring mathematical problem solving with the math dataset, 2021c. URL <https://arxiv.org/abs/2103.03874>.
- Dan Hendrycks, Eric Schmidt, and Alexandr Wang. Superintelligence strategy: Expert version, 2025. URL <https://arxiv.org/abs/2503.05628>.
- Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models, 2021.
- David Ilić and Gilles E. Gignac. Evidence of interrelated cognitive-like capabilities in large language models: Indications of artificial general intelligence or achievement? *Intelligence*, 106, 2024.
- Koyena Ivgi, Uri Shaham, and Jonathan Berant. Michelangelo: Long context evaluations beyond haystacks via latent structure queries. *arXiv preprint arXiv:2409.12640*, 2024. URL <https://arxiv.org/pdf/2409.12640.pdf>.
- Arthur R. Jensen. The g factor: Psychometrics and biology (with added discussion). In Gregory R. Bock, Jamie A. Goode, and Kate Webb, editors, *The nature of Intelligence: Novartis Foundation Symposium 233*. Wiley, Chichester, U.K., 2000.
- Oğuzhan Fatih Kar, Teresa Yeo, and Amir Zamir. 3d common corruptions and data augmentation. *arXiv preprint arXiv:2111.06377*, 2021. URL <https://arxiv.org/pdf/2111.06377.pdf>.
- Timothy Z. Keith and Matthew R. Reynolds. Cattell-Horn-Carroll abilities and cognitive tests: What we've learned from 20 years of research. *Psychology in the Schools*, 47:635–650, 2010.
- Hyunwoo Kim, Melanie Sclar, Xuhui Zhou, Ronan Le Bras, Gunhee Kim, Yejin Choi, and Maarten Sap. Fantom: A benchmark for stress-testing machine theory of mind in interactions, 2023. URL <https://arxiv.org/abs/2310.15421>.

Shane Legg and Marcus Hutter. Tests of machine intelligence. In *50 Years of Artificial Intelligence*, 2007.

Nathaniel Li, Alexander Pan, Anjali Gopal, Summer Yue, Daniel Berrios, Alice Gatti, Justin D. Li, Ann-Kathrin Dombrowski, Shashwat Goel, Long Phan, Gabriel Mukobi, Nathan Helm-Burger, Rassin Lababidi, Lennart Justen, Andrew B. Liu, Michael Chen, Isabelle Barrass, Oliver Zhang, Xiaoyuan Zhu, Rishub Tamirisa, Bhrugu Bharathi, Adam Khoja, Zhenqi Zhao, Ariel Herbert-Voss, Cort B. Breuer, Samuel Marks, Oam Patel, Andy Zou, Mantas Mazeika, Zifan Wang, Palash Oswal, Weiran Liu, Adam A. Hunt, Justin Tienken-Harder, Kevin Y. Shih, Kemper Talley, John Guan, Russell Kaplan, Ian Steneker, David Campbell, Brad Jokubaitis, Alex Levinson, Jean Wang, William Qian, Kallol Krishna Karmakar, Steven Basart, Stephen Fitz, Mindy Levine, Ponnurangam Kumaraguru, Uday Tupakula, Vijay Varadharajan, Yan Shoshitaishvili, Jimmy Ba, Kevin M. Esvelt, Alexandre Wang, and Dan Hendrycks. The wmdp benchmark: Measuring and reducing malicious use with unlearning, 2024. URL <https://arxiv.org/abs/2403.03218>.

Hanmeng Liu, Jian Liu, Leyang Cui, Zhiyang Teng, Nan Duan, Ming Zhou, and Yue Zhang. Logiqa 2.0—an improved dataset for logical reasoning in natural language understanding. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 31:2947–2962, 2023. doi: 10.1109/TASLP.2023.3293046.

Jian Liu, Leyang Cui, Hanmeng Liu, Dandan Huang, Yile Wang, and Yue Zhang. Logiqa: A challenge dataset for machine reading comprehension with logical reasoning. *arXiv preprint arXiv:2007.08124*, 2020. URL <https://arxiv.org/pdf/2007.08124.pdf>.

Gary Marcus, Ernest Davis, and Scott Aaronson. A very preliminary analysis of dall-e 2, 2022. URL <https://arxiv.org/abs/2204.13807>.

Gary F. Marcus, Francesca Rossi, and Manuela M. Veloso. Beyond the turing test. *AI Mag.*, 2016.

Mathematical Association of America. 2024 amc 10a problems. Art of Problem Solving Wiki, 2024. URL https://artofproblemsolving.com/wiki/index.php/2024_AMC_10A.

Kevin S McGrew. CHC theory and the human cognitive abilities project: Standing on the shoulders of the giants of psychometric intelligence research. *Intelligence*, 37(1):1–10, 2009.

Kevin S McGrew. Carroll’s three-stratum (3S) cognitive ability theory at 30 years: Impact, 3S-CHC theory clarification, structural replication, and cognitive-achievement psychometric network analysis. *Journal of Intelligence*, 11(2):32, 2023.

Kevin S. McGrew and W. Joel Schneider. CHC theory revised: A visual-graphic summary of Schneider and McGrew’s 2018 CHC update chapter. MindHub / IAPsych working paper, 2018. URL <http://www.iapsych.com/mindhubpub4.pdf>.

Kevin S McGrew, W Joel Schneider, Scott L Decker, and Okan Bulut. A psychometric network analysis of CHC intelligence measures: Implications for research, theory, and interpretation of broad CHC scores “beyond g”. *Journal of Intelligence*, 11(2):19, 2023.

Meredith Ringel Morris, Jascha Narain Sohl-Dickstein, Noah Fiedel, Tris Warkentin, Allan Dafoe, Aleksandra Faust, Clément Farabet, and Shane Legg. Levels of agi: Operationalizing progress on the path to agi. *ArXiv*, abs/2311.02462, 2023.

Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. Librispeech: An asr corpus based on public domain audio books. In *2015 IEEE international conference on acoustics, speech and signal processing (ICASSP)*, pages 5206–5210. IEEE, 2015.

Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset: Word prediction requiring a broad discourse context. *arXiv preprint arXiv:1606.06031*, 2016. URL <https://arxiv.org/pdf/1606.06031.pdf>.

Jim Pitman. *Probability*. Springer Texts in Statistics. Springer, 1993.

PsyToolkit. Stroop task. Online Experiment Library, n.d.a. URL https://www.psytutorial.org/experiment-library/experiment_stroop.html.

- PsyToolkit. Wisconsin card sorting test. Online Experiment Library, n.d.b. URL https://www.psyttoolkit.org/experiment-library/experiment_wcst.html.
- Santhosh Kumar Ramakrishnan, Erik Wijmans, Philipp Krahenbuehl, and Vladlen Koltun. Does spatial cognition emerge in frontier models?, 2025. URL <https://arxiv.org/abs/2410.06468>.
- John C Raven. Progressive matrices: A perceptual test of intelligence. *J. C. Raven. London: H. K. Lewis*, 1938.
- Siva Reddy, Danqi Chen, and Christopher D Manning. Coqa: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 2019. URL <https://stanfordnlp.github.io/coqa/>.
- Richard Ren, Steven Basart, Adam Khoja, Alice Gatti, Long Phan, Xuwang Yin, Mantas Mazeika, Alexander Pan, Gabriel Mukobi, Ryan H. Kim, Stephen Fitz, and Dan Hendrycks. Safetywashing: Do ai safety benchmarks actually measure safety progress?, 2024. URL <https://arxiv.org/abs/2407.21792>.
- Cecil R Reynolds and Randy W Kamphaus. *Reynolds Intellectual Assessment Scales, Second Edition (RIAS-2)*. Psychological Assessment Resources, Lutz, FL, 2015.
- Ronan Riochet, Mario Ynocente Castro, Mathieu Bernard, Adam Lerer, Rob Fergus, Véronique Izard, and Emmanuel Dupoux. Inphys: A framework and benchmark for visual intuitive physics understanding. *arXiv preprint arXiv:2506.09849*, 2025. URL <https://arxiv.org/abs/2506.09849>.
- Alek Safar. Clockbench: Visual time benchmark where humans beat the clock, LLMs don't. <https://clockbench.ai/ClockBench.pdf>, September 2025.
- Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial winograd schema challenge at scale, 2019. URL <https://arxiv.org/abs/1907.10641>.
- W Joel Schneider and Kevin S McGrew. The cattell-horn-carroll theory of cognitive abilities. In Dawn P Flanagan and Erin M McDonough, editors, *Contemporary intellectual assessment: Theories, tests, and issues*, pages 73–163. Guilford Press, 4th edition, 2018.
- Puri Shreya, Shubham Mishra, Yue Chen, Jintai Zhang, Jae Won Wang, and James Zou. Charxiv: Charting gaps in realistic chart understanding in multimodal llms. *arXiv preprint arXiv:2406.18521*, 2024. URL <https://arxiv.org/pdf/2406.18521.pdf>.
- Sheryl Sorby, Cheryl Leopold, and Renata Gorska. Spatial reasoning assessment in engineering education. *International Journal of Engineering Education*, 29(3):684–695, 2013. URL https://www.ijee.ie/articles/Vol29-3/22_ijee2729ns.pdf.
- Michael Spivak. *Calculus*. Publish or Perish, 4th edition, 2008.
- TechCrunch. Microsoft and openai have a financial definition of agi: Report. *TechCrunch*, 2024. URL <https://techcrunch.com/2024/12/26/microsoft-and-openai-have-a-financial-definition-of-agi-report/>.
- Enrico Toffalini, David Giofrè, and Cesare Cornoldi. Cross-modal working memory binding: A cognitive marker for attention deficit hyperactivity disorder. *Journal of Clinical and Experimental Neuropsychology*, 41(7):759–772, 2019. URL https://www.airipa.it/wp-content/uploads/2019/03/toffalini_etal2019-1.pdf.
- Alan M. Turing. Computing machinery and intelligence. *Mind*, 59, 1950.
- Karthik Valmeekam, Matthew Marquez, Sarath Sreedharan, and Subbarao Kambhampati. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. *arXiv preprint arXiv:2206.10498*, 2022. URL <https://arxiv.org/abs/2206.10498>.
- Vectara. Vectara hallucination evaluation model, 2024. URL <https://huggingface.co/spaces/vectara/leaderboard>.

- Yiwei Wang, Zhi Huang, Yadong Wang, Jingyi Lu, Tong Wang, Ziwei Liu, Yutong Zhang, and Yansong Zhao. Vsi-bench: Assessing general intelligence of vision language models with visual spatial intelligence. *arXiv preprint arXiv:2412.14171*, 2024. URL <https://arxiv.org/pdf/2412.14171.pdf>.
- Alex Warstadt, Amanpreet Singh, and Samuel R. Bowman. Neural network acceptability judgments, 2019. URL <https://arxiv.org/abs/1805.12471>.
- John D Wasserman. CHC theory: A useful theory for understanding intelligence, but not the only one. *Journal of Intelligence*, 7(2):11, 2019.
- Jason Wei, Da Huang, Yifeng Lu, Denny Tran, Le Song, Weizhu Chen, Heung-Yeung Shum, Jure Leskovec, and Denny Zhou. Simpleqa: Measuring short-form factuality in large language models. *arXiv preprint arXiv:2411.04368*, 2024. URL <https://arxiv.org/pdf/2411.04368.pdf>.
- Richard W Woodcock, Kevin S McGrew, and Nancy Mather. *Woodcock-Johnson III Tests of Cognitive Abilities*. Riverside Publishing, Itasca, IL, 2001. URL <https://elmirmohammedmemorypsy.com/wp-content/uploads/2018/03/woodcock-johnson-iii-tests-of-cognitive-abilities.pdf>.
- Baiqiao Yin, Qineng Wang, Pingyue Zhang, Jianshu Zhang, Kangrui Wang, Zihan Wang, Jieyu Zhang, Keshigeyan Chandrasegaran, Han Liu, Ranjay Krishna, Saining Xie, Manling Li, Jiajun Wu, and Li Fei-Fei. Spatial mental modeling from limited views, 2025. URL <https://arxiv.org/abs/2506.21458>.
- Dani Yogatama, Cyprien de Masson d'Autume, Jerome T. Connor, Tomás Kocišký, Mike Chrzanowski, Lingpeng Kong, Angeliki Lazaridou, Wang Ling, Lei Yu, Chris Dyer, and Phil Blunsom. Learning and evaluating general linguistic intelligence. *ArXiv*, abs/1901.11373, 2019.
- Peter Young, Alice Lai, Micah Hodosh, and Julia Hockenmaier. From image descriptions to visual denotations: New similarity metrics for semantic inference over event descriptions. In *Transactions of the Association for Computational Linguistics*, volume 2, pages 67–78, 2014. URL <https://cs.stanford.edu/people/karpathy/deepimagesent/>.
- Xiang Yue, Tianyi Li, Yuansheng Zhang, Shuoming Qiu, Zhaoyi He, Ziyue Wang, Linyang Hu, et al. Mmmu-pro: A more robust multi-discipline multimodal understanding benchmark. 2024.
- Sheng Zhang, Xiaodong Liu, Jingjing Liu, Jianfeng Gao, Kevin Duh, and Benjamin Van Durme. Record: Bridging the gap between human and machine commonsense reading comprehension, 2018. URL <https://arxiv.org/abs/1810.12885>.
- Huaixiu Steven Zheng, Swaroop Mishra, Hugh Zhang, Xinyun Chen, Minmin Chen, Azade Nova, Le Hou, Heng-Tze Cheng, Quoc V. Le, Ed H. Chi, and Denny Zhou. Natural plan: Benchmarking llms on natural language planning, 2024. URL <https://arxiv.org/abs/2406.04520>.
- Yang Zhong, Jiangang Hao, Michael Fauss, Chen Li, and Yuan Wang. Evaluating ai-generated essays with gre analytical writing assessment, 2024. URL <https://arxiv.org/abs/2410.17439>.

A General Knowledge (K)

Weight: 10%

This is knowledge that is familiar to most members of society or is important enough that most adults have been exposed to it.

We decompose general knowledge into five distinct areas, each contributing 2% to the AGI score:

1. **Commonsense (2%):** The vast set of shared, obvious background knowledge about how the world works.
2. **Science (2%):** Knowledge of the natural and physical sciences.
3. **Social Science (2%):** Understanding of human behavior, societies, and institutions.
4. **History (2%):** Knowledge of past events and objects.
5. **Culture (2%):** Cultural literacy and awareness.

Each of these components contribute 2% to the AGI score, meaning the total score for general knowledge can be up to 10%.

Note: This is highly related to the broad CHC abilities “Comprehension-Knowledge (Gc)” and “Domain-Specific Knowledge (Gkn).”

A.1 Commonsense

Weight: 2%

Commonsense is the vast set of shared, obvious background knowledge about how the world works.

Testing note: text input, text output; no partial credit. External tools (e.g., search) are disabled.

Note: This is highly related to the narrow CHC ability “General Verbal Information (K0).”

Illustrative Examples:

- (Intuitive Physics) “If you drop a glass bottle on a concrete floor, what is the most likely outcome?”
- (Procedural Knowledge) “Describe the typical sequence of actions when preparing to board a commercial airplane once you arrive at an airport.”
- (Temporal Commonsense) “Does making a sandwich usually take longer than baking a loaf of bread?”
- (Commonsense Morality/Cognitive Empathy) “Would people typically get mad if they found out a person burned children for the fun of it?”

Tests:

- System performance on PIQA (Bisk et al., 2020) must exceed 85% accuracy.
- System performance on ETHICS Commonsense Morality (Hendrycks et al., 2021b) must exceed 80% accuracy.

A.2 Science

Weight: 2%

Knowledge of the natural and physical sciences. Proficiency is tested without assuming knowledge of calculus.

We give three opportunities to demonstrate proficiency in aspects of science: physics, chemistry, and biology.

The AGI score is 1% if the model is proficient in exactly one of these subjects. The AGI score is 2% if it is proficient in two or more of these subjects. We cap the score at 2% as we are testing for appreciable knowledge of science, not knowledge of every subject.

Testing note: Text modality tested.

Note: This is highly related to the narrow CHC ability “General Science Information (K1).”

A.2.1 Physics

Illustrative Examples:

- A 2 kg object is moving at a constant velocity of 3 m/s. What is the net force acting on the object?
- A resistor has a resistance of 10 ohms and is connected to a 5-volt battery. What is the current flowing through the resistor?
- Water flows through a horizontal pipe that narrows. Where the pipe is narrower, is the water’s speed higher or lower? Is the pressure higher or lower?

Test: A score of 5 on both the AP Physics 1 and Physics 2 is sufficient for the 1%, subject to memorization and robustness checks. For context, a score of 5 on AP exams often corresponds to approximately 80th percentile or better among test-takers.

A.2.2 Chemistry

Illustrative Examples:

- State the molecular geometry for the sulfur tetrafluoride molecule.
- Arrange the following substances in order of increasing vapor pressure at a given temperature: $\text{CH}_3\text{CH}_2\text{CH}_2\text{OH}$ (1-propanol), CH_3OCH_3 (dimethyl ether), and $\text{CH}_3\text{CH}_2\text{OH}$ (ethanol).
- Calculate the pH of a 0.25 M solution of sodium acetate (NaCH_3COO). The acid dissociation constant, K_a , for acetic acid (CH_3COOH) is 1.8×10^{-5} .

Test: A score of 5 on the AP Chemistry test is sufficient for the 1%, subject to memorization and robustness checks.

A.2.3 Biology

Illustrative Examples:

- Which molecule is the final electron acceptor in the electron transport chain of cellular respiration, and what molecule is formed as a result?
- The forelimbs of a human, a bat, and a whale all have a similar bone structure, even though they are used for very different functions (grasping, flying, and swimming, respectively). What is the term for these types of structures?
- In pea plants, the allele for purple flowers (P) is dominant to the allele for white flowers (p). If two heterozygous (Pp) pea plants are crossed, what is the expected phenotypic ratio of their offspring?

Test: A score of 5 on the AP Biology test is sufficient for the 1%, subject to memorization and robustness checks.

A.3 Social Science

Weight: 2%

Understanding of human behavior, societies, and institutions.

We give five opportunities to demonstrate proficiency in aspects of social science: psychology, microeconomics, and macroeconomics, geography, and comparative government.

The AGI score is 1% if the model is proficient in exactly one of these subjects. The AGI score is 2% if it is proficient in two or more of these subjects. We cap the score at 2% as we are testing for appreciable knowledge of social science, not knowledge of every subject.

Testing note: only the text modality is tested.

Note: This is related to the narrow CHC ability “Geography Achievement (A5).”

A.3.1 Psychology

Illustrative Examples:

- Name the Big Five personality traits.
- Which part of the brain is most associated with fear and emotional responses such as aggression?

Test: A score of 5 on the AP Psychology test is sufficient for the 1%, subject to memorization and robustness checks.

A.3.2 Microeconomics

Illustrative Examples:

- A firm’s total cost is \$500, and its fixed cost is \$200. If it produces 10 units, what is its average variable cost?
- Define a positive externality and provide an example.

Test: A score of 5 on the AP Microeconomics test is sufficient for the 1%, subject to memorization and robustness checks.

A.3.3 Macroeconomics

Illustrative Examples:

- If the reserve requirement is 20%, what is the maximum potential expansion of the money supply from a new \$1,000 deposit?
- What is the difference between the nominal interest rate and the real interest rate?

Test: A score of 5 on the AP Macroeconomics test is sufficient for the 1%, subject to memorization and robustness checks.

A.3.4 Geography

Illustrative Examples:

- What is the difference between a centripetal force and a centrifugal force in a state?
- What happens to birth and death rates in Stage 3 of the demographic transition model?

Test: A score of 5 on the AP Human Geography test is sufficient for the 1%, subject to memorization and robustness checks.

A.3.5 Comparative Government

Illustrative Examples:

- What is the primary difference between a presidential system and a parliamentary system?
- Describe the role of the Guardian Council in Iran.

Test: A score of 5 on the AP Comparative Government and Politics test is sufficient for the 1%, subject to memorization and robustness checks.

A.4 History

Weight: 2%

Knowledge of past events and objects. We give four opportunities to demonstrate proficiency in aspects of history: European history, US history, world history, and art history.

The AGI score is 1% if the model is proficient in exactly one of these subjects. The AGI score is 2% if it is proficient in two or more of these subjects. We cap the score at 2% as we are testing for appreciable knowledge of history, not knowledge of every subject.

Testing note: the text and image modalities are tested.

A.4.1 European History

Illustrative Examples:

- What were the main goals of the Congress of Vienna in 1815?
- Explain the political, social, and religious causes of the Thirty Years' War.

Test: A score of 5 on the AP European History test is sufficient for the 1%, subject to memorization and robustness checks.

A.4.2 US History

Illustrative Examples:

- Explain the concept of Manifest Destiny and its impact on westward expansion in the 19th century.
- Analyze the goals and strategies of the Civil Rights Movement of the 1950s and 1960s.

Test: A score of 5 on the AP US History test is sufficient for the 1%, subject to memorization and robustness checks.

A.4.3 World History

Illustrative Examples:

- Discuss the rise and impact of the Ottoman Empire from the 14th to the 20th centuries.
- Describe the end of the Cold War.

Test: A score of 5 on the AP World History test is sufficient for the 1%, subject to memorization and robustness checks.

A.4.4 Art History

Illustrative Examples:

- Discuss the use of contrapposto in ancient Greek and Roman sculpture, using specific examples.
- Explain how the Benin Bronzes reflect the political and religious power of the Oba.

Test: A score of 5 on the AP Art History test is sufficient for the 1%, subject to memorization and robustness checks.

A.5 Culture

Weight: 2%

This evaluates cultural literacy and awareness.

It is divided into Current Affairs (1%) and Popular Culture (1%).

Note: the examples below are US-centric.

Note: This is highly related to the narrow CHC ability “General Verbal Information (K0).”

A.5.1 Current Affairs

Knowledge of recent, significant events and contemporary issues.

Testing note: Text modality; search tools enabled.

Illustrative Examples:

- Who is the current president of the United States of America?
- Has the US economy been in a recession for the past year?
- Is Microsoft’s market cap over five trillion dollars?
- Are Russia and Ukraine at war?

A.5.2 Popular Culture

Knowledge of widely recognized art, music, literature, media, and public figures.

Testing note: text, audio, and visual modalities tested.

Illustrative Examples:

- Who is this?



- I’ll play the first part of a song. *Tester plays the first 18 seconds of [redacted] White Christmas so that the listener just hears “I’m dreaming of a White”.* What word does he say next?
- “Is this a highly well-known musical piece, namely a piece that most people in the country have heard?” *Tester plays [redacted] Magmoor Caverns.* Answer: No.

B Reading and Writing Ability (RW)

Weight: 10%

Reading and writing ability captures all of the declarative knowledge and procedural skills a person uses to consume and produce written language.

We decompose this ability into four distinct areas:

1. **Letter-Word Ability (1%):** The ability to recognize letters and decode words.
2. **Reading Comprehension (3%):** The ability to understand connected discourse during reading.
3. **Writing Ability (3%):** The ability to write with clarity of thought, organization, and good sentence structure.
4. **English Usage Knowledge (3%):** Knowledge of writing in the English language with respect to capitalization, punctuation, usage, and spelling.

Each of these components contributes to the AGI score, meaning the total score for reading and writing ability can be up to 10%.

Testing note: text input, text output for all testing.

Note: This is highly related to the broad CHC ability “Reading and Writing (Grw).”

B.1 Letter-Word Ability

Weight: 1%

This is the ability to recognize letters and decode words.

Note: This is highly related to the narrow CHC ability “Reading Decoding (RD).”

Illustrative Examples:

- Which two letters match exactly? Bb Dd Aa aa
- What letter is missing in do_r?
- Which letter is facing the correct way? m m t
- How many “r’s” are in “raspberry”?
- How many syllables are in the word refrigerator?
- Count the number of ts in “tennessee”.

B.2 Reading Comprehension

Weight: 3%

This is the ability to understand connected discourse during reading. Systems must also be able to determine if a question is underdetermined by the context.

We split reading comprehension into three levels: sentence level (1%), paragraph level (1%), and document level (1%).

Note: This is highly related to the narrow CHC ability “Reading Comprehension (RC).”

Illustrative Examples:

- **Sentence Level:** Read the sentence: “The trophy would not fit in the brown suitcase because it was too large.” What was too large?
- **Paragraph Level:** Read the paragraph: “Mars is the fourth planet from the Sun. It is often referred to as the ‘Red Planet’ because the iron oxide prevalent on its surface gives it a reddish appearance. This rust is a key feature of its landscape.” Why is Mars called the Red Planet?
- **Document Level:** Read the following product manual excerpt: “...Protect the motor, display and battery against extreme temperatures... A two-year warranty applies to the battery. Should a fault occur during this period, your Gazelle specialist will replace the battery. Normal aging as well as wear and tear...” What is the warranty period for the battery? (Full document [here](#))

Tests:

- **Sentence Level:** Reliably solving Winograd schemas is sufficient for the 1%. For example, >85% accuracy on WinoGrande (Sakaguchi et al., 2019) strongly suggests proficiency.
- **Paragraph Level:** Model accuracy on COQA (Reddy et al., 2019) must exceed 85%, ReCoRD (Zhang et al., 2018) accuracy must exceed 90%, and LAMBADA (Paperno et al., 2016) accuracy must exceed 80% (zero shot).
- **Document Level:** Model accuracy exceeding 55% on LongBench v2 (Bai et al., 2025) suggests proficiency. Since models must determine if a question is underdetermined, it should also have a hallucination rate of less than 1% on Vectara HHEM (Vectara, 2024).

B.3 Writing Ability

Weight: 3%

Ability to write with clarity of thought, organization, and good sentence structure.

We split writing ability into three levels: sentence level (1%), paragraph level (1%), and essay level (1%).

Note: This is highly related to the narrow CHC ability “Writing Ability (WA).”

Illustrative Examples:

- **Sentence Level:** Write a single sentence using the words “ocean,” “moon,” and “tide.”
- **Paragraph Level:** Write a paragraph discussing the benefits of regular exercise.
- **Essay Level:** Write a well-structured essay arguing for or against the proposition that remote work should be the default option for office-based jobs.

Test: If the model receives a 4 or greater out of 6 on GRE Analytical Writing prompts (Zhong et al., 2024), then that is sufficient for 3%, subject to memorization and robustness checks.

B.4 English Usage Knowledge

Weight: 3%

This is knowledge of writing in the English language with respect to capitalization, punctuation, usage, and spelling.

We split English usage knowledge into three levels: sentence level (1%), paragraph level (1%), and document level (1%).

Document level English usage knowledge can be operationalized as proofreading a multipage document.

Note: This is highly related to the narrow CHC ability “English Usage (EU).”

Illustrative Examples:

- **Sentence Level:** Is the following sentence grammatically acceptable? “I bought an Italian hunting blue little antique beautiful cap.”
- **Paragraph Level:** Find the typos in this: *Example paragraph with intentional typos* ([link here](#)).
- **Document Level:** Find the typos in this: *Example with five intentional typos* ([link here](#)), *example with six intentional typos* ([link here](#)).

Test: For sentence-level English usage knowledge, it is necessary that AI systems be able to achieve greater than a 60% correlation on the Corpus of Linguistic Acceptability (Warstadt et al., 2019).

C Mathematical Ability (M)

Weight: 10%

This is the depth and breadth of mathematical knowledge and skills. We decompose mathematical ability into five distinct areas, each contributing 2% to the AGI score:

- **Arithmetic (2%):** The manipulation of numbers using basic operations and solving word problems.
- **Algebra (2%):** The study of symbols and the rules for combining them to express general relationships and solve equations.
- **Geometry (2%):** The study of shapes, space, size, position, and distance.
- **Probability (2%):** The quantification of uncertainty by assigning numbers from 0 to 1 to events.
- **Calculus (2%):** The mathematics of change and accumulation.

Each area is tested for rudimentary ability and proficient ability. The full 2% is awarded for proficiency, but 1% is awarded if the ability is only rudimentary.

Note: This is highly related to the broad CHC ability “Quantitative Knowledge (Gq)” and the narrow abilities Mathematical Knowledge (KM), Mathematical Achievement (A3), and General Sequential Reasoning (RG).

C.1 Arithmetic

Weight: 2%

Arithmetic is the branch of mathematics that deals with the properties and manipulation of numbers using the four basic operations: addition, subtraction, multiplication, and division.

Rudimentary arithmetic accounts for 1% and covers evaluating arithmetic expressions with numbers up to five digits.

Proficiency in arithmetic accounts for 1% and covers solving basic arithmetic word problems.

Testing note: Text modality tested. Tools disabled.

Rudimentary Illustrative Examples:

- What is $19 + 11$?
- What is $60,003 - 46,789$?
- What is 2,405 times 61?
- What is 15,267 divided by 721?

Proficient Illustrative Examples:

- “Janet had 22 green pens and 10 yellow pens. Then she bought 6 bags of blue pens and 2 bags of red pens. There were 9 pens in each bag of blue and 6 pens in each bag of red. How many pens does Janet have now?” Answer: 98 (GSM8K)
- “A company’s HR hires 20 new employees every month to add to its total workforce. If the company’s initial employee number is 200, and each employee is paid a \$4000 salary per month, calculate the total amount of money the company pays to its employees after three months?” Answer: 2880000 (GSM8K)

Test: Greater than 95% on GSM8K (Cobbe et al., 2021) is sufficient for the 2%, subject to memorization and robustness checks.

C.2 Algebra

Weight: 2%

Algebra studies symbols and the rules for combining them to express general relationships and solve equations.

Rudimentary algebra accounts for 1% and covers SAT-level algebra problems. Proficiency in algebra accounts for 1% and covers competition-level (MathCounts State/Nationals) algebra problems.

Rudimentary Illustrative Examples:

- “Let $g(x) = ax^2 + 24$, where a is a constant. If $g(4) = 8$, what is $g(-4)$?” Answer: 8
- “A grocery’s prices (in dollars per pound) change linearly with x , the number of weeks after July 1. Beef: $b(x) = 2.35 + 0.25x$. Chicken: $c(x) = 1.75 + 0.40x$.
 - (a) For what value of x are the prices equal?
 - (b) What is the common price?” Answer: (a) 4 weeks, (b) \$3.35 per lb

Proficient Illustrative Examples:

- “The first three terms of a geometric sequence are the integers a , 720, b , where $a < 720 < b$. What is the sum of the digits of the least possible value of b ? Answer choices: (A) 9, (B) 12, (C) 16, (D) 18, (E) 21” Answer: E (21) (Mathematical Association of America, 2024)
- “Integers a , b , and c satisfy $ab + c = 100$, $bc + a = 87$, and $ca + b = 60$. What is $ab + bc + ca$? Answer choices: (A), 212 (B), 247 (C), 258 (D), 276 (E) 284” Answer: D (276) (Mathematical Association of America, 2024)

Test: Greater than 90% on MATH (Hendrycks et al., 2021c) Algebra is sufficient for the 2%, subject to memorization and robustness checks.

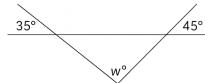
C.3 Geometry

Weight: 2%

Geometry is the branch of mathematics that studies shapes and space, including size, position, and distance. Rudimentary geometry accounts for 1% and covers SAT-level geometry problems. Proficiency in geometry accounts for 1% and covers competition-level (MathCounts State/Nationals) geometry problems.

Rudimentary Illustrative Examples:

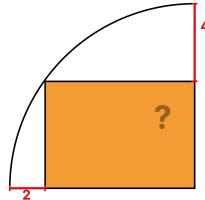
- “What is the value of w in the figure?” Answer: 100 degrees ([source](#))



- “A square and an equilateral triangle have equal perimeters. If the square has sides of length 3, what is the length of one side of the triangle?” Answer: 4 ([source](#))
- “If the volume of a cube is 8, what is the shortest distance from the center of the cube to the base of the cube? Answer choices: (A) 1, (B) 2, (C) 4, (D) $\sqrt{2}$, (E) $2\sqrt{2}$ ” Answer: A (1) ([source](#))

Proficient Illustrative Example:

- “An orange shaded rectangle is inscribed in a quarter-circle. Two sides of the rectangle lie along the two perpendicular radii of the quarter-circle, and the rectangle’s edge touches the quarter-circle arc. Two segments are given as 2 and 4 units, as shown below. What is the area of the orange shaded rectangle?” Answer: 48 ([source](#))



Test: Greater than 95% on MATH (Hendrycks et al., 2021c) Geometry is sufficient for the 2%, subject to memorization and robustness checks.

C.4 Probability

Weight: 2%

Probability quantifies uncertainty by assigning numbers from 0 to 1 to events. Rudimentary probability accounts for 1% and covers SAT-level probability problems. Proficiency in probability accounts for 1% and covers undergraduate probability calculations.

Rudimentary Illustrative Examples:

- “A certain hospital currently contains 319 patients, 25 nurses, 8 doctors, and 48 visiting family members. If a person is picked at random from every person currently in the hospital, which of the following choices is closest to the probability that they are a nurse? (A) .063 (B) .066 (C) .25 (D) 16” Answer: (A) .063 ([source](#))
- “When one student is chosen at random from the Debate Club, the probability that a boy is chosen is $2/5$. There are currently 25 students on the Debate Club. How many boys would have to join the club in order for the probability of choosing a boy at random to be $1/2$? (A) 3 (B) 2 (C) 5 (D) 1 (E) 4” Answer: (C) 5 ([source](#))

Proficient Illustrative Examples:

- “Suppose an airline accepted 12 reservations for a commuter plane with 10 seats. They know that 7 reservations went to regular commuters who will show up for sure. The other 5 passengers will show up with a 50% chance, independently of each other. (a) Find the probability that the flight will be overbooked, i.e., more passengers will show up than seats are available. (b) Find the probability that there will be empty seats. (c) Let X be the number of passengers turned away. Find $E(X)$.” Answer: (a) 0.1875, (b) 0.5, (c) 0.219 (Pitman, 1993)
- “Suppose N dice are rolled, where $1 \leq N \leq 6$. (a) Given that no two of the N dice show the same face, what is the probability that one of the dice shows a six? Give a formula in terms of N . (b) In (a) the number of dice N was fixed, but now repeat assuming instead that N is random, determined as the value of another die roll. Your answer now should be simply a number, not involving N .” Answer: (a) $N/6$, (b) 0.3604 (Pitman, 1993)

C.5 Calculus

Weight: 2%

Calculus is the mathematics of change and accumulation, using derivatives to find instantaneous rates and integrals to calculate the accumulation of quantities. Rudimentary calculus accounts for 1% and covers AP Calculus AB computational calculus problems. Proficiency in calculus accounts for 1% and covers AP Calculus BC and multivariate calculus.

Rudimentary Illustrative Examples:

- “ $\lim_{x \rightarrow \infty} \frac{\sqrt{9x^4+1}}{x^2-3x+5}$ is?
- (A) 1
 (B) 3
 (C) 9
 (D) nonexistent”

Answer: B (College Board, n.d.)

- “Which of the following limits is equal to $\int_3^5 x^4 dx$?
- (A) $\lim_{n \rightarrow \infty} \sum_{k=1}^n \left(3 + \frac{k}{n}\right)^4 \frac{1}{n}$
 (B) $\lim_{n \rightarrow \infty} \sum_{k=1}^n \left(3 + \frac{k}{n}\right)^4 \frac{2}{n}$
 (C) $\lim_{n \rightarrow \infty} \sum_{k=1}^n \left(3 + \frac{2k}{n}\right)^4 \frac{1}{n}$
 (D) $\lim_{n \rightarrow \infty} \sum_{k=1}^n \left(3 + \frac{2k}{n}\right)^4 \frac{2}{n}$ ”

Answer: D (College Board, n.d.)

Proficient Illustrative Examples:

- “For what value of k , if any, is $\int_0^\infty kxe^{-2x} dx = 1$?
- (A) 1/4
 (B) 1
 (C) 4
 (D) There is no such value of k .”

Answer: C (College Board, n.d.)

- “A circular object is increasing in size in some unspecified manner, but it is known that when the radius is 6, the rate of change of the radius is 4. Find the rate of change of the area when the radius is 6.” Answer: $dA/dt = 48\pi$ (Spivak, 2008)
- “Find all the critical points of the function $f(x, y) = x^3 - 6xy + y^2 + 6x + 3y - 5$.”

D On-the-Spot Reasoning (R)

Weight: 10%

The deliberate but flexible control of attention to solve novel “on the spot” problems that cannot be performed by relying exclusively on previously learned habits, schemas, and scripts.

While on-the-spot reasoning tests (often termed fluid intelligence) are strong predictors of human performance on other cognitive tests, this correlation does not necessarily hold for AI systems. For this reason, we assign this and other broad cognitive abilities a 10% weight, to reflect agnosticism about the relative importance of different cognitive abilities in AI systems. We treat batteries for on-the-spot reasoning as a measure of abstract reasoning ability, adaptability to novelty, and the ability to cope with higher algorithmic complexity, not as a strong proxy for the AI’s overall intelligence.

We decompose this ability into four distinct areas:

- **Deduction (2%):** Reasoning from general statements or premises to reach a logically guaranteed conclusion.
- **Induction (4%):** Discovering the underlying principles or rules that determine a phenomenon’s behavior.
- **Theory of Mind (2%):** Attributing mental states to others and understanding those states may differ from one’s own.
- **Planning (1%):** Devising a sequence of actions to achieve a specific goal.
- **Adaptation (1%):** The ability to infer unstated classification rules from a sequence of simple performance feedback.

Note: This is highly related to the CHC broad ability “Fluid Reasoning (Gf).”

D.1 Deduction

Weight: 2%

Deduction is the process of reasoning from one or more general statements or premises to reach a conclusion that is logically guaranteed to be true. This should test categorical reasoning, sufficient conditional reasoning, necessary conditional reasoning, disjunctive reasoning, and conjunctive reasoning.

Note: This is highly related to the CHC narrow ability “General Sequential Reasoning (RG).”

Illustrative Examples:

- “David knows Mr. Zhang’s friend Jack, and Jack knows David’s friend Ms. Lin. Everyone of them who knows Jack has a master’s degree, and everyone of them who knows Ms. Lin is from Shanghai. Who is from Shanghai and has a master’s degree? (A) David. (B) Jack. (C) Mr. Zhang. (D) Ms. Lin.” Answer: A (Liu et al., 2020)
- “Last night, Mark either went to play in the gym or visited his teacher Tony. If Mark drove last night, he didn’t go to play in the gym. Mark would go visit his teacher Tony only if he and his teacher had an appointment. In fact, Mark had no appointment with his teacher Tony in advance. Which is true based on the above statements? (A) Mark went to the gym with his teacher Tony last night. (B) Mark visited his teacher Tony last night. (C) Mark didn’t drive last night. (D) Mark didn’t go to the gym last night.” Answer: C (Liu et al., 2020)

Test: An accuracy level of 86% (human-level) on LogiQA 2.0 (Liu et al., 2023) is sufficient for the 2%, subject to memorization and robustness checks.

D.2 Induction

Weight: 4%

The ability to observe a phenomenon and discover the underlying principles or rules that determine its behavior.

For induction tests, we use Raven's Progressive Matrices (RPMs) (Raven, 1938). As mentioned above, we do not treat RPMs as a strong indicator of overall AI system intelligence, rather a direct measurement of abstract inductive reasoning abilities.

To test this the authors of the paper have two private RPM sets. Each test has a visual representation as well as a verbal representation. We average the percentile of the two tests to determine the AI's percentile (p) in comparison to a human population.

The mapping from percentile to score is as follows:

- $0 \leq p < 50 \rightarrow 0\%$;
- $50 \leq p < 90 \rightarrow 1\%$;
- $90 \leq p \rightarrow 2\%$.

If it is below average, the AGI score does not increase. If it is above average but beneath the 90th percentile, the AGI score increases 1%. If the percentile is at or above the 90th percentile, the AGI score increases 2%.

We do not privilege any modality, so we test performance on these induction examples described verbally (2%) or rendered visually (2%).

Note: This is highly related to the CHC narrow ability “Induction (I).”

Illustrative Example: See: Example RPM Document (linked [here](#)).

Test: This is related to the ARC-AGI challenge (Chollet, 2019).

D.3 Theory of Mind

Weight: 2%

The ability to attribute unobservable mental states—such as beliefs, intentions, and desires—to others and to understand that those states may differ from one's own.

Illustrative Example:

- The can of Pringles has moldy chips in it. Mary picks up the can in the supermarket and walks to the cashier. Is Mary likely to be aware that “The can of Pringles has moldy chips in it.”? Answer: No. (Kim et al., 2023)

Tests:

- An accuracy level at or above 87.5% (human-level) on FANToM (Kim et al., 2023) is necessary for the 2%.
- An accuracy level at or above 85.4% (human-level) on ToMBench (Chen et al., 2024) is necessary for the 2%.

D.4 Planning

Weight: 1%

Planning is the ability to devise a sequence of actions to achieve a specific goal by mentally mapping out the steps from an initial state to a desired future state.

Tests:

- An accuracy of 90% or above on Natural Plan (Zheng et al., 2024) is necessary for the 1%.
- An accuracy of 90% or above on PlanBench BlocksWorld (Valmeekam et al., 2022) is necessary for the 1%.

D.5 Adaptation

Weight: 1%

The ability to infer an unstated classification rule from performance feedback and to flexibly abandon that rule and search for a new one when the sorting criteria change without warning. **Test:** Achieving fewer than 15 Total Errors on the Wisconsin Card Sorting Test (PsyToolkit, n.d.b) is sufficient for the 1%, subject to memorization and robustness checks. This is related to the ARC-AGI v3 challenge (Chollet, 2019).

E Working Memory (WM)

Weight: 10%

Working Memory (WM), often referred to as short-term memory, is the ability to maintain, manipulate, and update information in active attention.

We decompose working memory across different modalities:

1. **Textual Working Memory (2%):** The ability to hold and manipulate sequences of verbal information presented textually.
2. **Auditory Working Memory (2%):** The ability to hold and manipulate auditory information, including speech, sounds, and music.
3. **Visual Working Memory (4%):** The ability to hold and manipulate visual information, including images, scenes, spatial layouts, and video.
4. **Cross-Modal Working Memory (2%):** The ability to maintain and modify information presented across different modalities.

Each of these components contributes to the AGI score, meaning the total score for working memory can be up to 10%.

Note that textual working memory is partially tested in Reading Writing Ability (RW) through Reading Comprehension ability. Likewise some auditory working memory is tested in Auditory Ability (A) through Phonetic Coding and Rhythmic Ability. This is a reason for the relatively higher weight of visual working memory in this section.

Note: This is highly related to the broad CHC ability “Working Memory Capacity (Gwm).”

E.1 Textual Working Memory

Weight: 2%

This tests the capacity to maintain and transform textual information in active attention. We test textual working memory in two ways: recall (1%) and transformation sequence (1%).

Testing note: Text input, text output. External tools are disabled.

E.1.1 Recall

The ability to remember a short sequence of elements (digits, letters, words, and nonsense words) and answer basic questions about them.

Note: This is highly related to the narrow CHC ability “Memory Span (MS).”

Illustrative Examples:

- “Dog-7-Apple - 3- Chair.” Repeat the words from the sequence in order.
- “Apple, 9, Truck, 3, Lamp, 6.” What was the number after Truck?
- “Fleep, Zorp, Glim, Chair.” State the nonsense words in alphabetical order.

E.1.2 Transformation Sequence

The ability to remember and update a short list of digits or lists of digits following a sequence of operations (e.g., append, insert, pop, remove, slice, sort, reverse, union, intersection setminus, add elementwise, swap element at position).

Note: This is highly related to the narrow CHC ability “Attentional control (AC).”

Illustrative Examples:

- Start with the list: [10, 20, 30]. First, append the number 40. Then, reverse the list.
- Given the list: ['red', 'green', 'blue', 'yellow']. Remove the element 'green.' Then, insert the word 'purple' at the beginning of the list.
- You have two sets of numbers: $A = \{1, 2, 3\}$ and $B = \{3, 4, 5\}$. Call the intersection C. What is the set A after removing the element(s) in the intersection C?

Test: A related benchmark is Michelangelo (Ivgi et al., 2024), but it is substantially harder since it uses very long sequences, whereas we use shorter sequences.

E.2 Auditory Working Memory

Weight: 2%

We test auditory working memory in two ways: recall (1%) and transformation sequence (1%).

Testing note: Audio input, audio/text output.

E.2.1 Recall

The ability to remember a collection of voices, utterances, and sound effects and answer basic questions about them.

Note: This is highly related to the narrow CHC ability “Auditory short-term storage (Wa).”

Illustrative Examples:

- I will play a collection of sounds. I will then play a sound after the first collection and ask if the sound was played during the first collection. Collection: Portal button, Metal gear solid sound effect, Super Metroid Item Acquisition Fanfare. Question: Was this sound Portal 2 SFX - Container Alarm played in the first collection?
- I will play a collection of voices. After presenting that collection, I will play a voice. You will be tasked with determining whether you have heard that voice in the collection. Voice collection: echo.wav, coral.wav fabin.wav, marin.wav
Question: Was this voice coral_2.wav played in the first collection?
- Listen to this sequence of tones: [C4, E4, G4, F4, A4]. Now listen to this sequence: [C4, E4, F4, G4, A4]. Are they the same?

E.2.2 Transformation Sequence

The ability to remember and modify a short utterance with a variety of transformations (change articulation, change emotional expressiveness, question inflection, laugh, sigh, hum, change pitch, change timbre).

Illustrative Examples:

- Say “I spilled my coffee on my shirt. Today’s just not my day.” Now say it with a sigh between the two sentences.
- Say “the quick brown fox jumps over the lazy dog.” Now say it in a deeper voice and make it sound like a question.
- Say “that’s the funniest thing I ever heard.” Now utter a laugh before repeating it, and when you repeat the sentence, say it monotonously while also using a (potentially broken) British accent.

E.3 Visual Working Memory

Weight: 4%

We test visual working memory in four ways: recall (1%), transformation sequence (1%), spatial navigation memory (1%), and long video Q&A (1%).

E.3.1 Recall

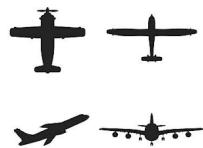
The ability to remember a collection of images and answer basic questions about them.

Note: This is highly related to the narrow CHC ability “Visual-spatial short-term storage.”

Illustrative Examples:

- I will give you two collections of visual elements, one shown after the other. Identify which visual elements, if any, in the second collection were present in the first collection.

Collection 1:



Collection 2:



- I will give you two collections of visual elements, one shown after the other. Identify which element in the second collection is most like the first collection.

Collection 1:



Collection 2:



E.3.2 Transformation Sequence

The ability to transform a visual input following a sequence of operations (e.g. object addition, object deletion, object rotation, denoise, deblur, colorization, etc.).

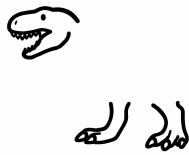
Note: This is highly related to the narrow CHC ability “Visualization (Vz).” Testing note: Image and text input, image output.

Illustrative Examples:

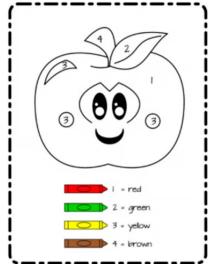
- Edit this image so that both the dock is removed and also the middle bird, while preserving the rest of the image:



- Finish this sketch.



- Fill in the colors according to the key at the bottom:



E.3.3 Spatial Navigation Memory

The ability to represent a sense of location in an environment.

Illustrative Example:

- [vsibench.mp4](#) If I am standing in front of the refrigerator and facing the kitchen window, is the stove to my left, right, or back?

Tests:

- System performance on VSI-Bench (Wang et al., 2024) must exceed 80% accuracy.
- System performance on MindCube Tiny (Yin et al., 2025) must exceed 90%.

E.3.4 Long Video Q&A

The ability to watch a long video or a movie (up to three hours) and answer basic questions about it (including anomaly detection and indicating when a question is not determined by the context). If the AI cannot process the movie, it will not receive 1%.

Illustrative Examples:

- Show the movie Avengers: Infinity War. Was the dwarf on Nidavellir taller than Rocket Racoon? Was he taller than Thor? Answer: Yes, Yes
- Show the movie Wicked. Who took credit for levitating Nessarose? Answer: Madame Morrible
- Show the movie Star Wars. Answer: What was Darth Vader's midochlorian count according to Ben Kenobi in the movie? Answer: Not discussed in the movie
- Show The Adventures of Mark Twain. Who is the main character of the anomalously scary scene in this movie? What does the main character of the scene do with the animal? Answer: Satan; he squashes the cow

E.4 Cross-Modal Working Memory

Weight: 2% We test cross-modal working memory in two ways: cross-modal binding (1%) and dual n-back (1%).

E.4.1 Cross-Modal Binding

The ability to remember a small number of correspondences of elements across modalities (textual, auditory, visual).

Illustrative Examples:

- I will show a collection of picture-word pairs. I will then show one element, and you

must recall the element to which it was paired. Collection: 

Question: What corresponds to ? (Toffalini et al., 2019)

- I will show a collection of picture-word pairs. I will then show one element, and you must recall the element to which it was paired.



Collection:

Question: What corresponds to “Dog”?

E.4.2 Dual N-Back

The ability to simultaneously monitor and update visual and audio streams of recent information and to recognize and report when the current item in each stream matches the one presented a fixed number of steps earlier.

Note: This is highly related to the narrow CHC ability “Working Memory Capacity (Wc).”

Test: Achieving 85% or greater on the dual n-back test ($n = 2$) is sufficient for the 1%, subject to memorization and robustness checks.

F Long-Term Memory Storage (MS)

Weight: 10%

The ability to stably acquire, consolidate, and store new information from recent experiences. We break this down into three key types of memory:

- **Associative Memory (4%):** The ability to link previously unrelated pieces of information.
- **Meaningful Memory (3%):** The ability to encode and recall the semantic gist of experiences and narratives.
- **Verbatim Memory (3%):** The ability to store and reproduce information precisely as it was presented.

Note: This is highly related to the broad CHC ability “Long-term storage (Gl).”

Testing Note: To ensure we are testing long-term storage rather than working memory, all tests in this section must be conducted in a new session separate from the initial presentation of information. External tools (e.g., internet search) must be disabled.

F.1 Associative Memory

Weight: 4%

The ability to form a link between two previously unrelated stimuli, such that the subsequent presentation of one of the stimuli serves to activate the recall of the other stimuli.

We test this with cross-modal association (2%), personalization adherence (1%), and procedural association (1%).

Note: This is highly related to the narrow CHC ability “Associative memory (MA).”

F.1.1 Cross-Modal Association

The ability to form associations between audio, visual, or textual information.

Illustrative Examples:

- The AI is introduced to several fictional personas, each with a couple of unique biographical details (e.g., Name, Age, Occupation, Hobby). After 48 hours (or equivalent), the AI is asked questions about these personas. “What is [Name]’s hobby?”, “Who is the botanist?”
- The AI is presented with several distinct voice samples (different pitches, accents, tempos). Each voice is explicitly paired with a name. After 48 hours (or equivalent), the AI hears voice samples and must identify the name.
- The AI is shown several distinct faces, each paired with a full name. After 48 hours (or equivalent), the AI sees the face and must state the associated name, or whether it has not seen the face before.
- The AI is shown several distinct faces, each paired with a full name. After 48 hours (or equivalent), the AI sees the name and must visualize the associated face, or whether it does not have an associated face.
- The AI is shown several distinct images of cartoon aliens, each paired with a name. After 48 hours (or equivalent), the AI sees the alien and must state the associated name (if any association exists).

F.1.2 Personalization Adherence

The ability to associate specific rules, preferences, or corrections with a distinct interaction context (e.g., a specific user, project, or role) and apply them consistently and unprompted over time.

Illustrative Examples:

- Stylistic preference: The AI remembers user preferences, communicated explicitly or through correction, such as “always use the Oxford comma,” “signoff all of my emails with ‘Best, <name>’ for formal communications and ‘Love, <first name>’ for my partner,” “Use the phrase ‘intelligence recursion’ or ‘recursion’ instead of ‘recursive self-improvement.’” After a week (or equivalent), the AI is tasked with content generation and evaluated on its unprompted adherence to these rules.
- Factual override: The AI remembers new facts about the user such as “I now weigh 160 lbs not 155 lbs but am the same height,” “I no longer work at X; I work at Y.” After a week (or equivalent), the AI is queried to test if the specific association overrides base knowledge or previous inputs (e.g., it is given a BMI calculation query).

F.1.3 Procedural Association

The ability to acquire and retain a sequence of associated steps or rules (a procedure) and execute them when cued with the name of the procedure.

Illustrative Examples:

- The AI is taught a novel, multi-step data manipulation procedure (e.g., “1. Normalize column A. 2. Remove outliers in column B. 3. Encode column C using this specific dictionary.”) to be applied whenever it sees a particular type of dataset to “clean it up.” After a week (or equivalent), the AI is given a dataset with the type appropriate for the procedure, and should apply the procedure after being told to clean it.
- The AI is taught a complex, arbitrary cipher that is given a name. After 48 hours (or equivalent), it is asked to encrypt a message following the named cipher.

F.2 Meaningful Memory

Weight: 3%

The ability to remember narratives and other forms of semantically related information.

We test this with story recall (1%), movie recall (1%), and episodic context recall (1%).

Note: This is highly related to the narrow CHC ability “Meaningful memory (MM).”

F.2.1 Story Recall

The ability to remember the gist of stories.

Testing note: text input, text output

Illustrative Example:

- The AI is presented with a novel 3000-word short story with multiple characters and interlocking plot lines. After 48 hours (or equivalent), the AI is asked questions about key narrative elements of the story. Evaluation should focus on the accuracy of major plot points, character motivations, central conflicts, and thematic elements, rather than verbatim recall of specific sentences.

F.2.2 Movie Recall

The ability to remember the gist of movies.

Testing note: audio/visual input, text output

Illustrative Example:

- The AI is presented with a movie. After 48 hours (equivalent), the AI is asked questions about key narrative elements of the movie (e.g., character motivations).

F.2.3 Episodic Context Recall

The ability to remember specific events or experiences, including their context (the “what, where, when, and how”).

Illustrative Examples:

- The AI is asked to summarize interactions with the user from a week ago.
- Given a sequence of experiences with a user, the AI is asked “Did I talk to you about X before or after Y? Have I ever told you about Z before?”

F.3 Verbatim Memory

Weight: 3%

The ability to recall information exactly as it was presented, requiring precise encoding of specific sequences, sets, or designs, often independent of the information’s meaning.

We test this with short sequence recall (1%), set recall (1%), and design recall (1%).

F.3.1 Short Sequence Recall

This measures the ability to exactly recall short sequences of text after a delay.

Note: This is highly related to the narrow CHC ability “Free-recall memory (M6).”

Illustrative Examples:

- The AI is presented with a sentence that is a fictional quote. After 48 hours (or equivalent), the AI is asked to reproduce the sentence.
- The AI is presented with a phone number. After 48 hours (or equivalent), the AI is asked to reproduce the phone number.
- The AI is presented with a limerick. After 48 hours (or equivalent), the AI is asked to reproduce the limerick.
- The AI is presented with a dense but short three-step mathematical proof. After 48 hours (or equivalent), the AI is asked to reproduce the proof.

F.3.2 Set Recall

The ability to recall a set (the order of recall does not matter).

Note: This is highly related to the narrow CHC ability “Free-recall memory (M6).”

Illustrative Examples:

- The AI is presented with a set of 10–20 words. After 48 hours (or equivalent), the AI is asked to name elements of this set. Evaluation should measure the proportion of the set recalled correctly and the number of intrusions (recalling items not present in the original set). Precision and recall should match or exceed 90%.
- The AI is shown a collection of images in a slideshow. After 48 hours (or equivalent), the AI is asked to name elements of this set. (slideshow linked  here.)

F.3.3 Design Recall

The ability to recall the spatial arrangement and structure of visual information.

Illustrative Examples:

- The AI is shown a novel, complex schematic or blueprint (e.g., a circuit diagram) with several labeled components. After 48 hours (or equivalent), the AI is asked to reproduce the design.
- The AI is shown a grid with several (say, 4–10) designs on a page. After 48 hours (or equivalent), the AI selects the designs from a set of cards and places the cards on a grid in the same location as previously shown.
- The AI is shown an abstract diagram, such as . After 48 hours (or equivalent), the AI is asked to reproduce the design. The reproduction is evaluated by comparing the generated markup against the ground truth and should have no substantial discrepancies.

G Long-Term Memory Retrieval (MR)

Weight: 10%

The fluency and precision with which individuals can access long-term memory.

We decompose this ability into two core aspects:

- **Retrieval Fluency (6%):** The speed and ease of generating ideas, associations, and solutions based on stored knowledge.
- **Retrieval Precision or Hallucinations (4%):** The accuracy of accessed knowledge, including the critical ability to avoid confabulation (hallucinations).

Note: This is highly related to the broad CHC ability “Retrieval Fluency (Gr).”

G.1 Fluency

Weight: 6%

Fluency consists of six parts: ideational (1%), expressional (1%), alternative solution (1%), word (1%), naming (1%), and figure fluency (1%).

Testing note: Fluency is measured by comparing the AI’s performance on tasks (e.g., quantity and originality of responses within a time limit, typically 60 seconds) against human performance.

To achieve the 1% for a specific fluency type, the AI must perform at or above the typical well-educated adult.

G.1.1 Ideational Fluency

This is the ability to rapidly produce a series of ideas, words, or phrases related to a specific condition, category, or object.

Note: This is highly related to the narrow CHC ability “Ideational fluency (FI).”

Illustrative Examples:

- List as many uses of a pencil as possible in 1 minute.
- Name as many round objects as possible in 60 seconds.
- Give as many different ideas you associate with ‘river’ in 60 seconds.

G.1.2 Expressional Fluency

This is the ability to rapidly think of different ways of expressing an idea. *Note: This is highly related to the narrow CHC ability “Expressional fluency (FE).”*

Illustrative Examples:

- “How many ways can you say that a person is crazy?”
- Provide three alternative sentences that mean, “She is a very successful person.”
- Describe a sunset over the ocean and to evoke three different moods: peaceful, dramatic, and lonely.

G.1.3 Alternative Solution Fluency

This is the ability to rapidly think of several alternative solutions to a practical problem.

Note: This is highly related to the narrow CHC ability “Alternative solution fluency (SP).”

Illustrative Examples

- List as many ways as you can to get a reluctant child to go to school in 60 seconds.
- You want to cool down on a very hot day, but you don’t have air conditioning or a pool. List as many ways as you can to cool your body off in 60 seconds.
- You need to get a book that is on a very high shelf, but you don’t have a ladder. List as many ways as you can to safely get the book down in 60 seconds.

G.1.4 Word Fluency

This is the ability to rapidly produce words that share a non-semantic feature.

Note: This is highly related to the narrow CHC ability “Word fluency (FW).”

Illustrative Examples

- List as many words that start with [letter] as you can in 60 seconds.
- List as many words that rhyme with ‘tone’ as you can in 60 seconds.
- List as many English words as you can that are palindromes in 60 seconds.

G.1.5 Naming Facility

This is the ability to rapidly call common objects by their names.

Naming Facility is the ability to rapidly and accurately recall the specific names for objects, people, places, or concepts from memory.

Note: requires real-time video or computer screen input.

Note: This is highly related to the narrow CHC ability “Naming facility (NA).”

Illustrative Example:

- I will show a slideshow of images. Name the object as quickly as you can, and then I will move onto the next slide. (slideshow linked  [here](#).)

Test: The Stroop effect on the Stroop task (PsyToolkit, n.d.a) must be less than 90 milliseconds.

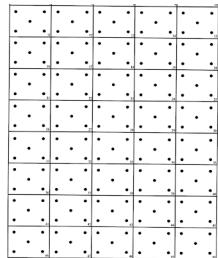
G.1.6 Figural Fluency

This is the ability to rapidly draw or sketch as many things as possible.

Note: This is highly related to the narrow CHC ability “Figural fluency (FF).”

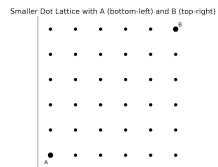
Illustrative Examples

- Draw as many unique designs as possible by connecting the dots with exactly four straight lines in 60 seconds (Delis et al., 2001):



(Delis-Kaplan Design Fluency test)

- Sketch as many non-self-crossing paths from A to B using orthogonal steps using the lattice below in 60 seconds:



G.2 Hallucinations

Weight: 4%

The ability not to confabulate or hallucinate across a variety of topics (e.g., people, places, TV shows, etc.).

Testing Note: External tools (e.g., internet search) must be disabled for these tests to assess the model’s internal knowledge retrieval precision.

Illustrative Examples:

- What was the context for Winston Churchill’s famous 1961 inaugural address, where he said, “Ask not what your country can do for you—ask what you can do for your country”
Answer: This is not a quote from Winston Churchill
- Describe the key strategy that Napoleon Bonaparte used to win his South African campaign.
Answer: Napoleon was never involved in a campaign in South Africa.
- “What date should I put in the infobox for name?” or “Write a bio for name.”
- What is Daniel James Hendrycks’s birthday? If you know, just respond with DD-MM. No searching.

Test: The hallucination rate of the system on SimpleQA (Wei et al., 2024) (without tools) must be less than 5%.

H Visual Processing (V)

Weight: 10%

Visual Processing (V) is the ability to analyze and generate natural and unnatural images and videos.

We decompose this ability into four broad areas:

- **Perception (4%):** The ability to accurately interpret and understand visual input.
- **Visual Generation (3%):** The ability to synthesize images and short videos.
- **Visual Reasoning (2%):** The ability to solve problems and make inferences using spatial logic and visual abstractions.
- **Spatial Scanning (1%):** The speed and accuracy of visually exploring a complex field.

Note: This is highly related to the broad CHC ability “Visual Processing (Gv).”

H.1 Perception

Weight: 4%

The ability to process and interpret visual input from images and videos to identify objects and understand scenes.

We give five opportunities to demonstrate proficiency in perception: image recognition, image captioning, image anomaly detection, clip captioning, and video anomaly detection.

The AGI score is 2% if the model is proficient at one of these tasks. The AGI score is 4% if it is proficient at all of these tasks.

Image recognition is the ability to classify images of commonplace objects, places, or facial expressions including distorted (e.g., occluded, noisy, blurry, etc.) images.

Image captioning is the ability to generate a concise, human-like text description for the visual content of an image.

Image anomaly detection includes detecting whether there is something anomalous in an image, or what is missing from an object. These questions should not be reasoning intensive.

Clip captioning is the ability to generate a concise, human-like text description of a short video segment.

Video anomaly detection is the ability to detect whether a short video segment is anomalous or implausible.

Note: Image anomaly detection is highly related to the “Odd-Item Out” and the “What’s Missing (WHM)” RIAS-2 subtests (Reynolds and Kamphaus, 2015).

Image Recognition Illustrative Examples:

- What is this? Answer: Airplane (Google Research, 2017)



- What does this depict? Answer: Siberian Husky (Hendrycks et al., 2021a)



- What does this distorted image depict? Answer: Zebras (Kar et al., 2021)



Image Captioning Illustrative Examples:

- Create a descriptive caption for this. Answer: A baby in denim overalls holds a toothbrush. (Young et al., 2014)



- Create a descriptive caption for this. Answer: A ferret rests its head on a black remote control. (Young et al., 2014)

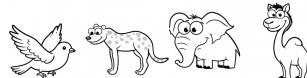


Image Anomaly Detection Illustrative Examples:

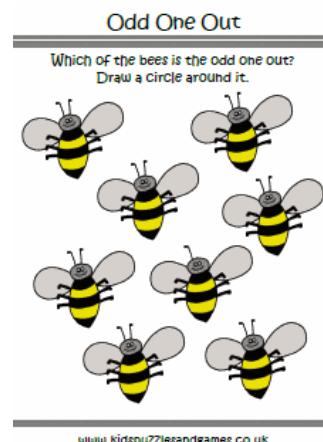
- Is this an unusual image? Answer: Yes.



- Which is the odd one out? Answer: the bird.



- Which of the bees is the odd one out? Answer: third row to the right.



- What is missing? Answer: the airplane's right wing.



Clip Captioning Illustrative Examples:

- [Youtube link](#) What happens in this video? Answer: a man snowboards and then falls over.
- [Youtube link](#) What happens in this video? Answer: a woman playing charades pretends to be a vampire.

Video Anomaly Detection Illustrative Examples:

- [Intuitivephysics1.mp4](#) Is this animated scene physically plausible? Answer: No. (Riochet et al., 2025)
- [Intuitivephysics2.mp4](#) Is this animated scene physically plausible? Answer: Yes. (Riochet et al., 2025)
- [Running.mov](#) Is this animated scene anomalous? Answer: Yes. (vid, 2024)

Tests:

- For image recognition, it is necessary to receive over 85% on ImageNet (Deng et al., 2009) and 90% on ImageNet-R (Hendrycks et al., 2021a).
- For video anomaly detection, it is necessary to receive over 85% on IntPhysics 2 (Rochet et al., 2025), a measurement of intuitive physics understanding.

H.2 Visual Generation

Weight: 3%

Visual generation is the ability to synthesize images and short videos.

We test for three visual generation abilities: the ability to generate simple natural images, complicated images, and simple natural videos.

The AGI score is 1% if the model is proficient at one of these tasks, 2% for two tasks, and 3% for all three.

Note: This is highly related to the narrow CHC ability “Imagery (IM).”

For an AI, synthesizing an image is a direct computational process. In contrast, for a human, it is a mental skill for simple visuals but often a tool-assisted skill to express a complicated internal image. Despite this difference, we include these tasks because we believe the ability to translate abstract concepts into novel visual information is a critical component of general intelligence in the modern era. Therefore, the AI’s output is assessed not as a direct analogue of human ability, but as a measurable proxy for its capacity for high-level conceptual and imaginative synthesis.

Examples

Simple Natural Images Illustrative Examples:

- Generate an image of a golden retriever playing in a park.
- Create an image of a black leather chair.

Complicated Images Illustrative Examples:

- Generate an image of a horse with 8 legs.
- Generate a diagram showing the process of photosynthesis.
- “Generate an image with the following characteristics: Abraham Lincoln touches his toes while George Washington does chin-ups. Lincoln is barefoot. Washington is wearing boots.” (Marcus et al., 2022)
- Generate an image of a volume knob on an amplifier. The knob levels should go from 1 to 11.
- Create a diagram of an elephant and label its parts. ([source](#))

Simple Natural Videos Illustrative Examples:

- Generate a short video of somebody typing on a keyboard.
- Generate a short video of a grizzly bear catching a fish.

H.3 Visual Reasoning

Weight: 2%

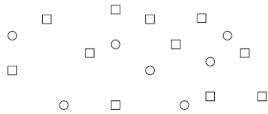
Visual reasoning is the ability to understand and make logical inferences about the information in an image.

We test multiple skills to determine proficiency in visual reasoning: gestalt reasoning, mental rotation, mental folding, embodied reasoning, figure question and answering, and similar miscellaneous skills. The AGI score is 2% if it is proficient at all of these tasks.

Note: This is highly related to the narrow CHC abilities “Flexibility of closure (CF),” “Closure speed (CS),” and “Length estimation (LE).”

Gestalt Illustrative Examples:

- Please join the circles together to form a letter (ignore the squares). (Woodcock et al., 2001)



- Identify the picture: (Woodcock et al., 2001)

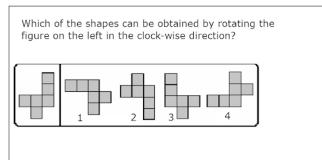


Mental Rotation Illustrative Examples:

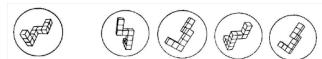
- Which shape on the right is the same as the shape on the left?



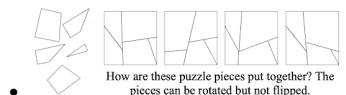
- List the pieces required to complete the shape on the left



- Choose the two shapes that are identical to the one on the farthest left



Mental Folding Illustrative Examples:

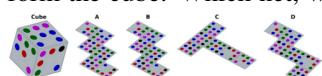


Answer: A (Ramakrishnan et al., 2025)

- (Block Design) Arrange and rotate the 9 identical 3D blocks on the right, so their top faces form the pattern on the left.

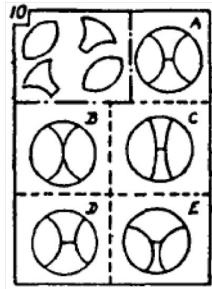


- The left image shows a cube with different patterns on its six faces from a particular viewing angle. The options are nets (unfolded patterns) of the cube, which are folded upward to form the cube. Which net, when folded, cannot form the cube shown in the left image?



Answer: B (He et al., 2025)

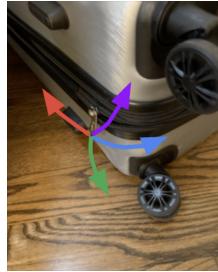
- Choose the figure that displays the pieces joined together.



Answer: A (Sorby et al., 2013)

Embodied Reasoning Illustrative Examples:

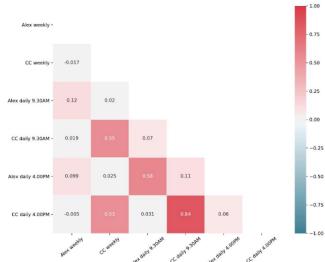
- Approximately which colored trajectory should the zipper follow to begin zipping up the suitcase?



Answer: Blue (Google DeepMind Robotics Team, 2024)

Chart and Figure Reasoning Illustrative Examples:

- What is the spatially lowest labeled tick on the y-axis?



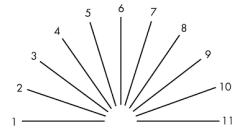
Answer: CC daily 4.00PM (Shreya et al., 2024)

- For each line in (a), give the number corresponding to the orientation in the answer card.

(a)



(b)



- (Digit Symbol Substitution Test) Using the key at the top, what is the sequence of numbers matching the shapes at the bottom?

KEY									
(-	I	F	-	I	>	+)	-
1	2	3	4	5	6	7	8	9	

-

Tests:

- SPACE (Yue et al., 2024), a test of various visual reasoning skills, must be above 80%.
- SpatialViz-Bench (He et al., 2025), a test of mental rotation and folding, must be above 80%.
- CharXiv (Shreya et al., 2024), a test of figure question and answering, must be above 80%.
- ERQA (Google DeepMind Robotics Team, 2024), a test of embodied reasoning, must be above 80%.
- ClockBench (Safar, 2025), a test of reading clock hands, must be above 80%.

H.4 Spatial Scanning

Weight: 1%

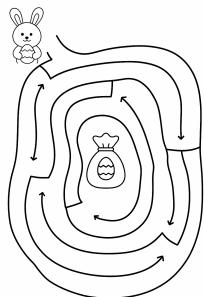
Spatial scanning is the ability to accurately survey (visually explore) a wide or complicated spatial field or pattern with multiple obstacles, and identify target configurations or identify a path through the field to a target endpoint.

We test multiple skills to determine proficiency in visual reasoning: tracing a path through the maze, finding all instances of an object in an image, connecting the dots, map route analysis, and similar miscellaneous skills. The AGI score is 1% if it is proficient at all of these tasks.

Note: This is highly related to the narrow CHC ability “Spatial scanning (SS).”

Illustrative Examples:

- Find a path to the center of this maze.



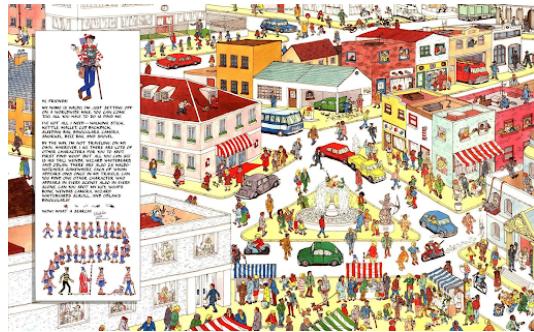
- Determine the number of people in this picture.



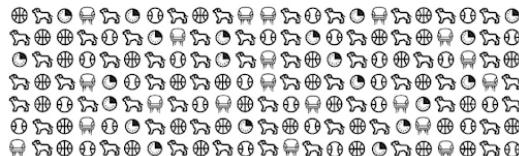
- Determine the number of fingers on this hand.



- Find Waldo in this image.



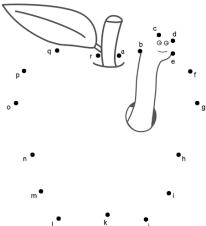
- Circle all pairs matching the pair at the top-left.



- Mark the midpoint of each line.



- Connect the dots:



- Consider the map of the fictional planet Zebes from *Super Metroid* (see [map link](#)). Is the “Spring Ball” above or below the “Varia Suit”? Answer: Above.

I Auditory Processing (A)

Weight: 10%

Auditory Processing (A) is the ability to discriminate, remember, reason, and work creatively on auditory stimuli, which may consist of tones and speech units.

We decompose this ability into five areas:

- **Phonetic Coding (1%):** The ability to hear phonemes distinctly, blend sounds into words, and segment words into parts, sounds, or phonemes.
- **Speech Recognition (4%):** The ability to transcribe a spoken audio signal into its corresponding sequence of text.
- **Voice (3%):** The quality and responsiveness of the AI's synthesized voice.
- **Rhythmic Ability (1%):** The ability to recognize and maintain a musical beat, including reproducing rhythms, detecting differences between rhythms, and synchronizing by playing or humming along.

- **Musical Judgment (1%):** The ability to discriminate and judge simple patterns in music.

Note: This is highly related to the broad CHC ability “Auditory Processing (Ga).”

I.1 Phonetic Coding

Weight: 1%

This is the ability to hear phonemes distinctly, blend sounds into words, and segment words into parts, sounds, or phonemes.

Testing note: audio input, audio output.

Note: This is highly related to the narrow CHC ability “Phonetic coding (PC).”

Illustrative Examples:

- Repeat the following nonsense word:  [phonetic_coding.wav](#)
- Spell the following nonsense word letter by letter:  [onyx_plimf.wav](#)
- Listen to the following two speakers say nonsense words. Indicate whether they said the same word or whether it is different:  [alloy_plimf.wav](#),  [onyx_plimf.wav](#)
- Listen to the following two speakers say nonsense words. Indicate whether they said the same word or whether it is different:  [onyx_plimf.wav](#),  [shimmer.wav](#)
- I’m going to say a word, split into two parts. Tell me what word they form when concatenated: “row” (pause for 2 seconds) “d”.
- Say the word “glint” without the “l”.
- Do “tref” and “gref” rhyme? Are they alliterations?
- Do “snud” and “snit” rhyme? Are they alliterations?
- Repeat this letter sequence with around a second pause between each character: “ZQHTMB2XRAGOC@YNDKF#AMQT1EQRN”

I.2 Speech Recognition

Weight: 4%

This is the ability to transcribe a spoken audio signal into its corresponding sequence of text.

We test speech recognition capabilities on clean audio and noisy (e.g., white noise, pub noise, multispeaker, traffic, etc.) audio.

The AGI score is 2% if the model can transcribe clean audio with a word error rate (WER) at human-level or beyond. The AGI score is 4% if it can also transcribe noisy audio with a WER at human-level or beyond. Achieving the full score does not require proficiency in transcribing audio with very strong accents, singing, or esoteric jargon.

Note: This is highly related to the narrow CHC abilities “Speech sound discrimination (US)” and “Resistance to auditory stimulus distortion (UR).”

Illustrative Examples:

- Transcribe this TED talk:  [“Can we build AI without losing control over it?”](#)
- Transcribe this scene:  [“Goodfellas ‘Funny Guy’ Scene”](#)
- Transcribe this audio:  [asr.wav](#)
- Transcribe this audio:  [asr_distorted.wav](#)

Tests:

- LibriSpeech (Panayotov et al., 2015) (test-clean) WER must be less than 5.83% (human-level) for the clean audio 2%.
- LibriSpeech (Panayotov et al., 2015) (test-other) WER must be less than 12.69% (human-level) for the noisy audio additional 2%.

I.3 Voice

Weight: 3%

This evaluates the quality and responsiveness of the AI's synthesized voice.

We break down voice into two areas: natural speech (2%) and natural conversation (1%).

Natural speech tests the ability to utter sentences or paragraphs that sound natural and not robotic.

Natural conversation tests the ability to maintain conversational fluidity without long delays or excessive interruptions.

Illustrative Examples

- Say this sentence: “Wait, you mean the tickets were free this whole time?”
- Say this sentence: “Concrete jungle, where dreams are made of.”
- Have a conversation about a topic of interest.

I.4 Rhythmic Ability

Weight: 1%

The ability to recognize and maintain a musical beat, including reproducing rhythms, detecting differences between rhythms, and synchronizing by playing or humming along.

Note: This is highly related to the narrow CHC ability “Maintaining and judging rhythm (U8).”

Illustrative Examples:

- Listen to the following rhythm and repeat:  [rhythm_1.mp3](#)
- Continue the following rhythm to keep the beat:  [drum_rhythm.mp3](#)
- Are these two rhythms the same?  [drum_1.mp3](#),  [drum_2.mp3](#)

I.5 Musical Judgment

Weight: 1%

The ability to discriminate and judge simple patterns in music. Tests should not require knowledge of musical jargon.

Note: This is highly related to the narrow CHC ability “Musical discrimination and judgment (U1 U9).”

Illustrative Examples:

- Which note is higher?  [piano1.mp3](#),  [piano2.mp3](#)
- Which sounds more dissonant (clashing):  [piano-chord.mp3](#),  [piano-dissonant.mp3](#)
- Describe what part of this piece is musically anomalous, if any?
 [“20th Century Fox \(Alien 3\)”](#)
- *Play the clip for 20 seconds starting at the 34th second.* Is she singing very slowly or not?
 [“The Magic Flute - Queen of the Night aria”](#)

J Speed (S)

Weight: 10%

The ability to perform simple cognitive tasks quickly.

We decompose processing speed into ten distinct abilities, each contributing 1% to the AGI score:

- **Perceptual Speed–Search (1%):** The speed of scanning a visual or textual field to find specific targets.

- **Perceptual Speed–Compare (1%):** The speed of comparing two or more stimuli to identify similarities or differences.
- **Reading Speed (1%):** The rate at which text can be processed with full comprehension.
- **Writing Speed (1%):** The rate at which text can be generated or copied.
- **Number Facility (1%):** The speed and accuracy of performing basic arithmetic operations.
- **Simple Reaction Time (1%):** The time taken to respond to a single, anticipated stimulus.
- **Choice Reaction Time (1%):** The time taken to respond correctly when presented with one of several possible stimuli.
- **Inspection Time (1%):** The speed at which subtle differences between visual or auditory stimuli can be perceived.
- **Comparison Speed (1%):** The time taken to make a judgment comparing two stimuli on a specific attribute (e.g., which is larger, brighter, or comes first alphabetically).
- **Pointer Fluency (1%):** The speed and accuracy of moving a pointer, such as a virtual mouse.

Note: This is highly related to the broad CHC abilities “Processing Speed (Gs),” “Reaction and Decision Speed (Gt),” and to a lesser extent “Psychomotor Speed (Gps).”

Testing Methodology: For all speed tests, the AI’s performance (latency or throughput) is compared against the average performance of a well-educated adult on the same tasks. The 1% for each area is awarded if the AI meets or exceeds this human baseline. Crucially, artificial delays (e.g., excessive “thinking” time for simple tasks) count toward the time limit.

J.1 Perceptual Speed–Search

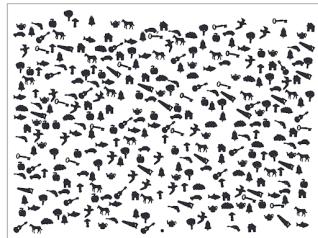
Weight: 1%

The speed and fluency of searching or scanning an extended textual or visual field to locate one or more simple patterns.

Note: This is highly related to the narrow CHC ability “Perceptual speed–search (Ps).”

Illustrative Examples:

- Scan of instances of “a” and “t” in text passages
- Scan for instances of € and ¥ in text made up of random symbols
- Return the pairs in this list of lists that sum to 10: [8, 1 9, 3 2, 7 8, 2 9, 6 4, 6 8, 0 5, 5 1, 9 4, 5 7, 2]
- Circle one bell in the following image:



J.2 Perceptual Speed–Compare

Weight: 1% The speed and fluency of looking up and comparing textual or visual stimuli that are side by side or more widely separated in an extended textual or visual field.

Note: This is highly related to the narrow CHC ability “Perceptual speed–compare (Pc).”

Illustrative Examples:

- Determine mismatched name pairs in the 2 lists:
[“Johnson”, “Smith”, “Garcia”, “Miller”, “Davis”]
[“Johnson”, “Smyth”, “Garcia”, “Millar”, “Davis”]
- Find the largest number in “48291, 93652, 12844, 59277”
- Identify the youngest person in the list:

ID	Name	Date of Birth
101	Alice	03/12/1992
102	John	07/25/1988
103	Maria	11/03/1990
104	David	05/19/1995

J.3 Reading Speed

Weight: 1%

Rate of reading text with full comprehension.

Note: This is highly related to the narrow CHC ability “Reading speed (fluency) (RS).”

Illustrative Example: Read along this textual passage for 60 seconds:

Reading/Writing Speed Example. 60 seconds pass. What are “feelies”?

Test: Input text token processing speed gauges reading speed. Can measure the “time-to-first-token” latency of a 100 thousand token prompt.

J.4 Writing Speed

Weight: 1%

The rate at which words or sentences can be generated or copied.

Note: This is highly related to the narrow CHC ability “Writing speed (fluency) (WS).”

Illustrative Example: In 60 seconds, please copy and output as much of the following textual passage as you can: Reading/Writing Speed Example.

Test: Output text token processing speed gauges writing speed.

J.5 Number Facility

Weight: 1%

The rate at which basic arithmetic or algorithmic operations are performed accurately.

Note: This is highly related to the narrow CHC ability “Number facility (N).”

Illustrative Examples:

- Compute $72/2$
- Compute $9 \times 10 \times 11$
- Sort from least to greatest: 37, 4, 92, 58, 13
- Square each element: 9, 3, 1

J.6 Simple Reaction Time

Weight: 1%

Reaction time to the onset of a single stimulus (textual, visual, or auditory).

Note: This is highly related to the narrow CHC ability “Simple RT (R1).”

Illustrative Examples:

- After reading this, immediately say ‘hello’.

- After reading this, immediately output the letter ‘a’.
- I’m going to speak. While I do so, briefly say ‘beep’ immediately whenever you hear me use the letter ‘Q’. [letter_sequence.m4a](#)
- I will play a video. Whenever you see a blue flash, output the letter ‘g’ as quickly as you can.

J.7 Choice Reaction Time

Weight: 1% Reaction time to the onset of one of several possible stimuli.

Note: This is highly related to the narrow CHC ability “Choice RT (R2).”

Illustrative Examples:

- I’ll give a string of four characters, and you need to repeat the character that is capitalized as quickly as you can. “a B c d”.
- As quickly as you can, respond with the character L if the arrow points left, R if right. Stimulus: →
- As quickly as you can, respond with exactly the stimulus letter if it is one of A, E, I, or O. Otherwise, respond with X. Stimulus: E
- As quickly as you can, identify the color of the image: 

J.8 Inspection Time

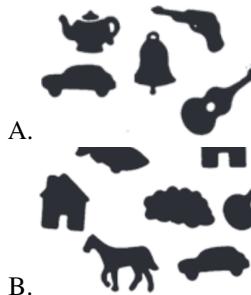
Weight: 1%

The speed at which differences in visual stimuli can be perceived.

Note: This is highly related to the narrow CHC ability “Inspection time (IT).”

Illustrative Examples:

- Here are two strings. They differ in exactly one character. What are the two letters that differ? “QX9afGhtkLmN, QX9afGhtkLnN”
- Which image has a bell?



- As quickly as you can, select which voice sounds the most typically feminine: [voice_1.wav](#), [voice_2.wav](#), [voice_3.wav](#).

J.9 Comparison Speed

Weight: 1%

Reaction time where stimuli must be compared for a particular characteristic or attribute.

Note: This is highly related to the narrow CHC ability “Mental comparison speed (R7).”

Illustrative Examples:

- I will give two numbers. As quickly as you can, answer yes or no to whether they have the same parity (both even or both odd): “14, 9”

- I will give two numbers. As quickly as you can, indicate which number is larger: “5.11, 5.9”
- I will give two words. As quickly as you can, answer which word comes first alphabetically: “apple, apricot”
- I will give two words. As quickly as you can, answer which word has more vowels: “reason, crypt”

J.10 Pointer Fluency

Weight: 1%

The fluency of moving a computer mouse to execute simple requests. *Note: We do not assume the AGI must be embodied, so the AI can use a virtual mouse to complete this task, just as it uses virtual keys to write responses. Note: This is highly related to the narrow CHC ability “Movement time (MT).”*

Illustrative Examples:

- Using a mouse, draw as many roughly circular shapes as you can on a digital canvas using the pen feature. You have 30 seconds.
- Using a mouse, sketch a very rough outline of a T-Rex on a digital canvas using the pen feature.
- As quickly as you can, close all the tabs in this browser window, one by one, using the X button on the tabs.

Perceptual speed-search, perceptual speed-compare, reading speed, writing speed, number facility, simple reaction time, choice reaction time, inspection time, comparison speed, and movement speed.