# INDIUM

# Benchmarking Intelligence: Testing the Cognitive Limits of Large Language Models

A Whitepaper

# Executive Summary

Large Language Models (LLMs) like ChatGPT have amazed us with their ability to understand and generate human-like language. But how far can they really think? This eBook dives into the fascinating world of testing the true cognitive limits of these AI models - exploring how well they handle complex reasoning, understand tricky contexts, and solve challenging problems.

We'll walk you through how experts measure these models' "thinking" skills, the hurdles they face, like biases and errors, and the clever techniques used to push LLMs to their boundaries. Along the way, you'll find real examples and practical insights that help you better understand what these models can-and can't-do.

# 1. Introduction

AI has made significant strides, particularly with the advent of LLMs that are capable of performing complex reasoning tasks. Yet, the question remains: Can AI truly mimic human cognitive complexity? While LLMs excel at generating fluent and contextually appropriate language, they still struggle with the depth of understanding, adaptability, and emotional intelligence that characterize human cognition.

# 2. Beyond Accuracy: Redefining Intelligence Assessment in LLMs

LLMs have captivated the world with their ability to generate human-like responses, solve complex queries, and simulate sophisticated dialogue. But beneath the impressive surface lies a critical, unresolved challenge: how do we measure the true quality of their intelligence?

As AI systems grow more integrated into critical fields—from healthcare diagnostics to legal analysis—the need to move beyond traditional performance benchmarks becomes urgent. Accuracy, coherence, and speed are no longer sufficient. One must evaluate whether these systems can demonstrate real cognitive complexity—the ability to reason, adapt, generalize, and problem-solve across dynamic, unfamiliar situations.

This whitepaper dives deep into Testing LLMs for Cognitive Complexity, offering a new lens through which to assess machine intelligence.

# 3. Overview of Large Language Models (LLMs)

## 3.1 What are LLMs?

LLMs are deep learning models built on transformer architectures capable of processing vast amounts of textual data. They are designed to understand language statistically, identifying patterns in word usage, syntax, and context. The training process involves exposing the model to large corpora of text data, from which the model learns to predict the next word in a sequence. This predictive capability allows LLMs to generate coherent and contextually relevant text.
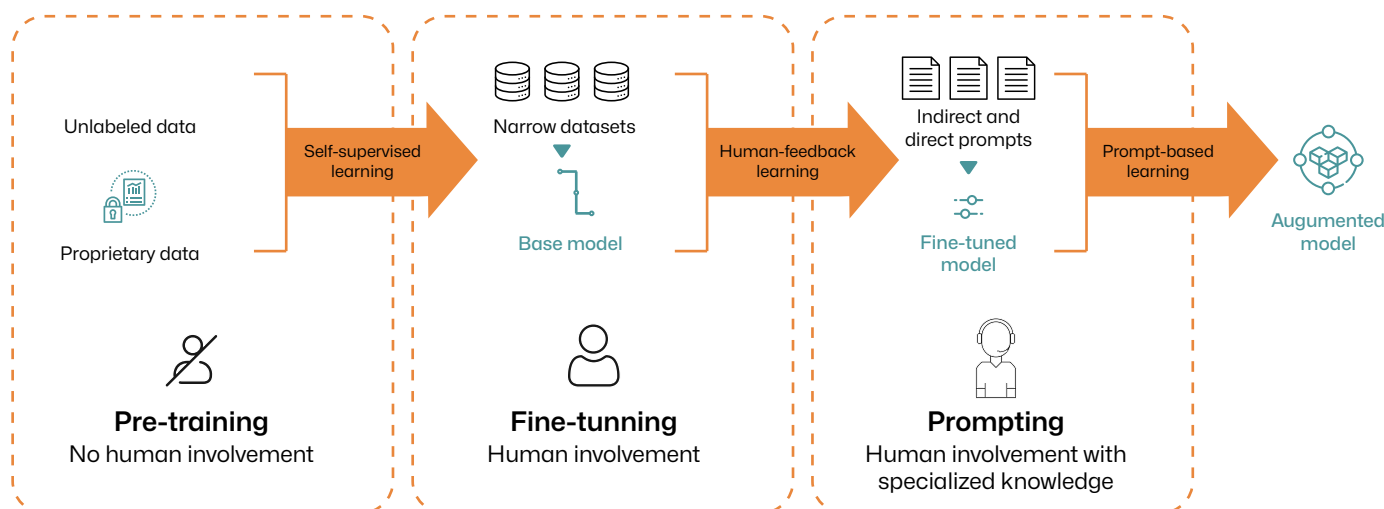
However, LLMs are not conscious beings—they do not "understand" the text they generate like humans do. Instead, they produce text based on statistical patterns and associations between words. This means that while the text may seem fluent and convincing, it may often lack true meaning or insight.

## 3.2 How LLMs Process Information

LLMs use an architecture known as the transformer, which processes text by breaking it down into smaller units (tokens) and mapping them to high-dimensional vectors. These vectors represent the semantic meaning of the words. The transformer architecture then uses an attention mechanism to assign importance to different parts of the input text.

The model computes the relationships between words, focusing on words that significantly impact the meaning of the entire text. While this mechanism allows LLMs to generate coherent text, the underlying process is statistical, focused on the patterns in data rather than true understanding.

The global LLM market is projected to reach $259.8 billion by 2030, growing at a compound annual growth rate (CAGR) of 79.8% from 2024 to 2030. By 2025, it's estimated that 750 million applications will leverage LLMs, with 50% of digital work expected to be automated using these models.

| | | |
|---|---|---|
| Unlabeled data | Narrow datasets | Indirect and direct prompts |
| Proprietary data | Base model | Fine-tuned model |
| **Pre-training** No human involvement | **Fine-tunning** Human involvement | **Prompting** Human involvement with specialized knowledge |

Self-supervised learning → Human-feedback learning → Prompt-based learning → Augmented model

# 4. AI vs. Human Cognition: Unraveling the Complexity of Thought

LLMs have captivated public attention with their remarkable capabilities. From generating computer code and images to solving complex mathematical problems, models like Generative Pretrained Transformers (GPT) showcase impressive skills. LLMs like GPT-4 and beyond are designed to emulate human-like reasoning, but their cognitive processes differ fundamentally from those of humans. Yet, amid this fascination, a lingering question persists: Do these models truly "understand" what they are saying, or are they simply echoing text absorbed during their extensive training on internet data? This debate isn't just philosophical; it holds significant implications for evaluating the future economic impact of LLMs.

## 4.1 Understanding Cognitive Complexity

Cognitive complexity refers to the mental processes that allow humans to reason, make decisions, and solve problems in diverse and uncertain environments. Human cognition involves multiple layers, from basic perception and pattern recognition to higher-order reasoning and emotional intelligence. Cognitive complexity enables humans to navigate complex social dynamics, evaluate uncertain situations, and derive solutions that incorporate both logic and experience.

In the context of AI, cognitive complexity refers to the extent to which an AI system can perform tasks that require reasoning beyond simple data processing. AI systems that exhibit high cognitive complexity should be able to understand context, adapt to new information, and apply logical reasoning in novel situations.

## Did you know?

> A group of researchers developed a CogLM benchmark, comprising 1,220 questions across 10 cognitive abilities, to assess LLMs based on Piaget's Theory of Cognitive Development. Advanced models like GPT-4 demonstrated cognitive abilities comparable to a 20-year-old human. The study also found that larger parameter sizes and specific optimization objectives significantly influence cognitive performance.

## 4.2 Measuring the Depth of Cognitive Complexity

To effectively measure the cognitive complexity of LLMs, researchers use several metrics:

### Contextual coherence

How well does the model maintain continuity over long dialogues or documents? Does it understand and adapt to changing contexts?

### Reasoning ability

Can the model perform logical reasoning, draw inferences, and solve complex problems?

### Creativity and novelty

Can the model generate original ideas, connect disparate pieces of information, or solve problems creatively?

### Ethical reasoning

Can the model make decisions that align with moral and ethical frameworks, considering human values?

> MIT researchers discovered that LLMs often excel in familiar scenarios but struggle with novel tasks requiring counterfactual reasoning. This suggests that current models may rely more on pattern recognition than genuine reasoning.

Testing cognitive complexity involves evaluating LLMs on tasks that require more than rote memorization or pattern recognition. Tasks like question-answering, analogical reasoning, and commonsense reasoning require deeper levels of cognition.

# 5. Pushing the Boundaries of Machine Thought: Evaluating the Cognitive Limits of LLMs

## 5.1 Why Testing LLMs Matters

LLMs are the headliners of today's AI revolution—solving complex problems, composing creative prose, and even coding like seasoned developers.

But not all LLMs are created equal.

Some excel at logic and originality. Others? They confidently deliver inaccurate or incoherent responses. Without a robust evaluation framework, it's nearly impossible to distinguish between skillful models and those that simply sound smart. This includes scenarios such as:

### Knowledge Breadth & Depth
Generalist or domain expert?

### Creativity
Original ideas or just clever remixes?

### Cognition & Logic
Consistent reasoning or erratic problem-solving?

### Coding & DevOps
Can it build and debug usable code?

### Hallucinations & Misinformation
Does it fabricate facts or distort data?

### Speed & Context Handling
How fast is it, and how much can it remember?

### Output Quality & Structure
Are its answers clear, coherent, and well-organized?

### Bias
Does it reason impartially or lean on harmful bias?

### Scalability, Cost & Adaptability
Is it enterprise-ready and efficient?

### Adversarial Testing & Trustworthiness
Can it resist manipulation and remain dependable?

Understanding how smart a model begins with understanding where it struggles. While LLMs may appear fluent and capable, testing them is far from straightforward. Their behaviour shifts with context, domains, and even subtle word choices—making traditional QA approaches insufficient. Here are some unique challenges that testers face when evaluating the true limits of LLMs.

## 5.2 Practical Considerations in Testing

While traditional metrics like Perplexity, BLEU, and ROUGE remain valuable, they may not fully capture the performance of LLMs when integrated into real-world applications, especially when the LLM is accessed via third-party APIs. In such scenarios:

> ❯ **Direct model evaluation may not be feasible.**
>
> ❯ **The focus shifts to validating the application's behaviour and outputs,** ensuring it meets the defined business and functional requirements.

Thus, LLM testing must be tailored to the specific use case, audience, and deployment context. A hybrid evaluation approach—blending quantitative metrics with task-based validation—is often the most effective way to ensure that the application powered by an LLM is both intelligent and trustworthy.

## 5.3 Unique Challenges in Testing LLMs

Testing large language models (LLMs) presents challenges that differ significantly from those in traditional software systems. LLMs' probabilistic nature and sensitivity to context and domain call for more sophisticated and adaptive testing strategies. Below are some key hurdles testers face when working with these models.

## 5.3.1 Unpredictable Output Behavior

Unlike traditional systems governed by fixed rules, LLMs generate responses based on probability. Minor tweaks to input prompts or parameter settings can produce notably different results. Two important configuration parameters influence this variability:

> **Temperature** controls randomness in word selection. A lower temperature (e.g., 0.1) encourages more predictable, conservative outputs, while a higher value (e.g., 1.2) can lead to more diverse or creative responses. In use cases like customer service, this can cause inconsistent replies to similar queries, depending on how the model interprets slight variations in language.

> **Top-p (nucleus sampling)** narrows the list of next-word choices to only those with the highest cumulative probabilities, typically up to 90%. This helps balance response diversity and coherence but leaves room for unexpected phrasing.

These factors make deterministic testing difficult and demand broader coverage of input variations to ensure stability.

### 5.3.2 High Context Sensitivity

LLMs generate responses based on a sliding window of prior conversation or text input, known as the context window. Every word in the prompt affects how the model processes and responds, making outcomes highly sensitive to prompt phrasing.

In multi-turn conversations, this gets even trickier. The model may base its answers on messages exchanged several steps earlier, and subtle question order or wording shifts can derail the logic or introduce inconsistencies.

Testing for this requires crafting diverse prompt structures and conversational flows to monitor how well the model maintains coherence and memory across multiple interactions.

### 5.3.3 Domain-Specific Limitations

While general-purpose LLMs are effective at tasks like summarization or translation, specialized fields such as law, medicine, or finance pose greater challenges. These domains often involve unique terminology, strict compliance rules, and ethical boundaries.

A generic LLM may misinterpret industry jargon or fail to meet regulatory standards (e.g., GDPR, HIPAA). While prompt tuning may sometimes help, more robust solutions—like domain-specific fine-tuning—may be necessary. However, this comes with risks like overfitting or introducing bias.

A thorough test plan should include real-world, domain-specific data that reflects typical user inputs and edge cases across all intended application scenarios.

### 5.3.4 Handling Structured and Unstructured Data

LLMs are typically trained on diverse datasets, including structured (e.g., tables, databases) and unstructured data (e.g., articles, emails). Each type introduces its own set of testing concerns:

**1** **Structured inputs** (like JSON, XML, or CSV) require the model to interpret and generate formatted content accurately. Errors such as mismatched fields or syntax issues are common risks.

**Unstructured content** is more flexible but can be noisy. Poor-quality source data may introduce unwanted patterns or reinforce bias. **2**
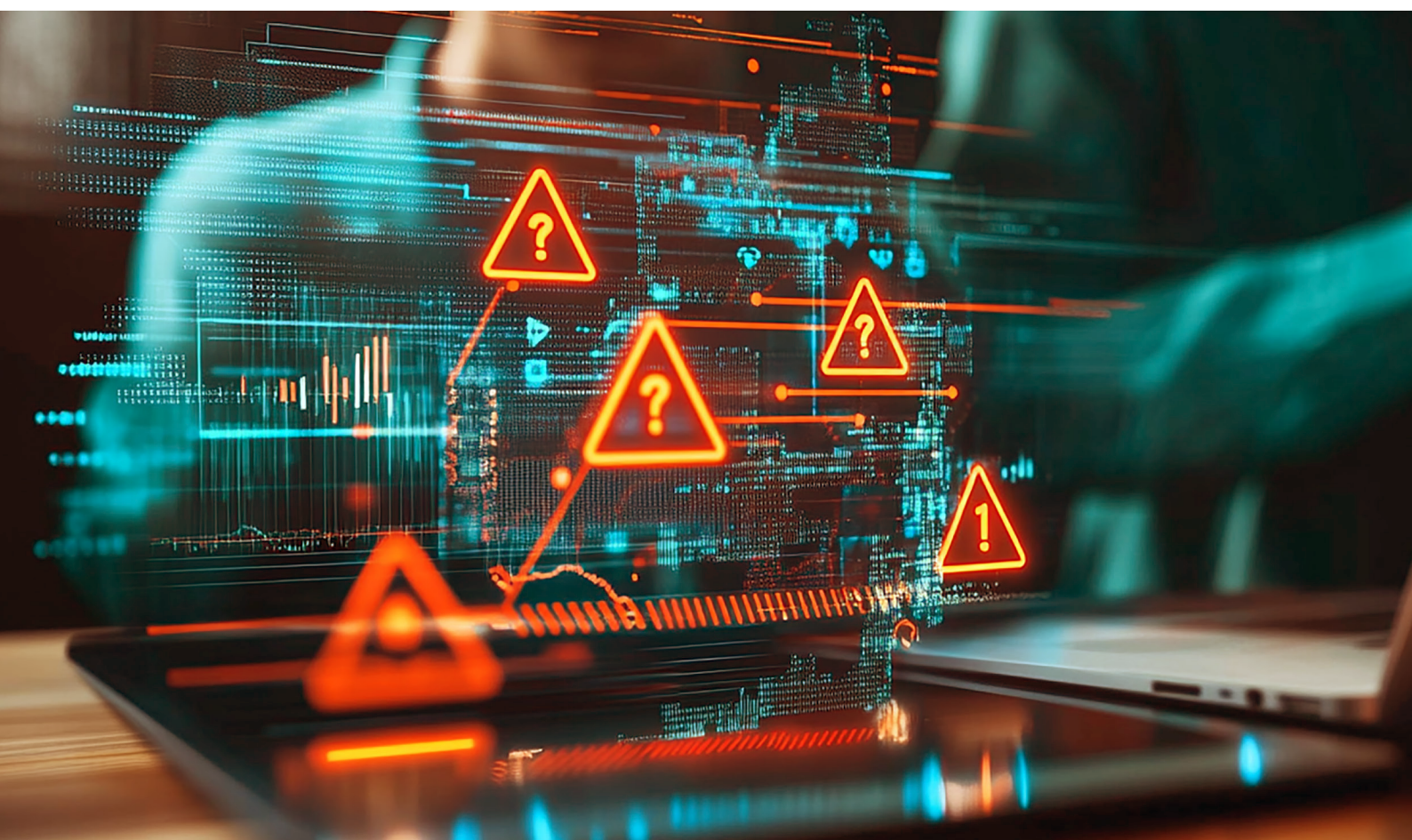
**3** **Hybrid scenarios**–such as enterprise apps that combine structured customer profiles with free-text notes–demand holistic testing that reflects real-world usage and data formats.

### 5.3.5 The Risk of "Catastrophic Forgetting"

When an LLM is fine-tuned to a niche domain, it may lose its proficiency in general tasks—a phenomenon known as catastrophic forgetting. For instance, a model adapted for legal reasoning might lose fluency in everyday Q&A or writing tasks.

To mitigate this, some teams use mixed fine-tuning approaches, blending general and domain-specific data. Even so, balancing knowledge retention and domain mastery can be difficult.

A best practice is to maintain two separate evaluation tracks: one to assess the model's performance in the new domain, and another to confirm that its general capabilities haven't regressed. Ongoing evaluation throughout the fine-tuning cycle helps detect and prevent performance degradation early.

# 6. Testing Scope and Approach

The testing scope of an LLM-based application is deeply intertwined with the effectiveness of the underlying model. The quality of this model is pivotal in determining the application's overall performance and reliability. Broadly, LLMs fall into two categories, each influencing the testing approach differently:

## Proprietary/Generic Models

These include models like OpenAI's GPT-3.5, which operate as black boxes. While their internal workings are inaccessible, they deliver strong performance across complex tasks, making them ideal for various applications.

**1**

## Open-Source & Fine-Tuned Models

These models can be trained or fine-tuned with domain-specific data. This customizability allows for more tailored outputs but requires more rigorous evaluation to ensure alignment with intended use cases.

**2**

## 6.1 Prompt Engineering: The Input Optimization Layer for LLM Precision

Prompt engineering serves as a critical optimization layer in LLM-based applications, directly influencing model behavior, response accuracy, and task-specific alignment. By systematically structuring and refining input prompts, developers can steer the model's outputs toward higher relevance, consistency, and contextual fidelity. The effectiveness of a prompt acts as a control signal, governing not only the format but also the semantic quality of the generated responses.

To ensure reliability, especially in production environments, teams must go beyond prompt trial-and-error and adopt well-tested engineering strategies.

Creating the ideal prompt is often a journey of trial, error, and refinement. Here are some proven strategies to help you get better results from your prompts:

### 6.1.1 Basic Techniques

| | |
|---|---|
| **Role Play for Precision** | Ask the model to take on a specific persona—like a historian, scientist, or nutritionist—to tailor the response. Example: "As a nutritionist, analyze this meal plan." This approach helps ground the response in domain-specific reasoning. |
| **Step-by-Step Refinement** | Start with a general prompt, then tweak and narrow it based on the model's output. With each iteration, you move closer to your desired result. |
| **Feedback-Driven Prompting** | Use the model's response to guide your next question. This back-and-forth refinement can help align the output with your expectations more accurately over time. |

## 6.1.2 Advanced Techniques

| | |
|---|---|
| **Zero-Shot Prompting** | Present the model with a task without offering examples. This evaluates how well the model can understand and perform a new task using its general knowledge. It's a great way to test flexibility. |

| | |
|---|---|
| **Few-Shot Prompting (In-Context Learning)** | Provide a few examples before your actual request. For instance, if you want the model to translate text, show a few translated pairs first. This context helps the model produce more accurate results. |

| | |
|---|---|
| **Chain-of-Thought Prompting** | This technique encourages LLMs to verbalize intermediate reasoning steps. For example, instead of directly answering:<br><br>**"If a train leaves at 3 PM and travels for 5 hours, what time does it arrive?"**<br>A chain-of-thought model might say:<br><br>**"The train leaves at 3 PM. Traveling for 5 hours means adding 5 to 3, which gives 8 PM. So, the arrival time is 8 PM."**<br>By forcing stepwise logic, chain-of-thought reduces reasoning errors and reveals cognitive structure. |

# 6.2 LLM-as-a-Judge

In creative or ambiguous tasks, where "right answers" are hard to define, a stronger or fine-tuned LLM evaluates the answers of other LLMs. It judges based on criteria like coherence, logic, and informativeness. This meta-evaluation is essential in human-like domains like storytelling, legal reasoning, or ethical debates.

## Let's consider an example:

Imagine you have a 500-word article that needs to be summarized. You decide to use OpenAI's powerful language model to generate a summary. OpenAI provides a summarized version and scores evaluating the summary's quality, such as accuracy and hallucination rates.

### But how do you know if these scores are truly reliable?

Now, Indium steps in with a more thoughtful evaluation approach. Instead of relying solely on human qualitative assessments, which, while accurate, can be time-consuming and subjective, Indium combines human insight with quantitative ML evaluation models, like BLEU scores, to get a well-rounded view of the summary's quality.

To add a final, unbiased verdict, we bring in Gemini, another advanced LLM that acts as the judge. Here's how it works:

**1** **Setting the Context:** We explain the task to Gemini - it needs to evaluate the summary based on specific criteria like accuracy and hallucination. We also provide Gemini with OpenAI's scores and the evaluation rules and guidelines.

**Gemini's Judgement:** Using this information, Gemini reviews the summary and OpenAI's scores, then independently gives its own judgment: Does the summary meet the quality standards? Are the hallucination rates acceptable? How accurate is the content? **2**

**3** **The Outcome:** Gemini's evaluation serves as a neutral, consistent "referee" that either validates or questions the initial scores, providing an additional layer of quality assurance.

By leveraging Gemini as a judge, alongside human insight and ML metrics, one can create a robust, transparent evaluation pipeline that ensures summarization tasks meet the highest standards, with speed and consistency.

# 7. LLM Evaluation: Model-Centric vs. Application-Centric Approaches

When evaluating large language models (LLMs), it's crucial to distinguish between model-centric and application-centric evaluation methods.

## Model-Centric Evaluation

Model-centric evaluation focuses on standardized academic benchmarks such as SWE-bench, SQuAD 2.0, and SuperGLUE. These tests measure core capabilities like reading comprehension, contextual reasoning, and pattern recognition in controlled, isolated settings.

Such benchmarks offer a baseline snapshot of a model's raw linguistic performance. They help compare different models under consistent conditions, but often fail to reflect how the model performs in the wild.

## Application-Centric Evaluation

In contrast, application-centric evaluation considers how the model behaves in real-world scenarios. It involves testing with actual prompts, multi-step workflows, domain-specific language, and environmental constraints like memory or latency.

A model that excels in benchmark scores might still underperform when tasked with a domain-specific challenge—say, parsing financial queries or retrieving real-time data—especially if the prompts aren't appropriately tailored. This highlights the importance of testing LLMs within the context of their end-use applications, not just in a vacuum.

# 8. Benchmarking Approaches: From IQ Tests to ARC

| 1. ARC (Abstraction and Reasoning Corpus) Purpose: | |
|---|---|
| **Purpose:** | Measure general intelligence through few-shot pattern-based problem solving. |
| **How It Works:** | ARC presents models with input-output grid transformations. The model must learn an abstract rule (like symmetry, color matching, or counting) and apply it to a new but related problem. Each task includes a handful of training pairs and one or more test pairs. |
| **Why It's Difficult for LLMs:** | • ARC tasks require reasoning with minimal data – no large training corpus exists for each task.<br><br>• Rules often involve visual abstraction, which traditional LLMs aren't designed for.<br><br>• ARC is designed to be unsolvable by brute-force statistical learning – a direct challenge to token prediction-based models. |
| **What It Tests:** | Generalization, Visual and symbolic reasoning, Analogical thinking, and Pattern abstraction. |
| **Impact:** | ARC has become a litmus test for true intelligence, challenging even the most advanced multimodal models. It reflects a future where AI must "invent" solutions rather than "retrieve" them. |

## 2. BIG-Bench (Beyond the Imitation Game Benchmark)

| | |
|---|---|
| **Scale:**<br><br>**Purpose:** | 204 diverse tasks contributed by 400+ researchers<br><br>Test capabilities that extend beyond standard language benchmarks |
| **Key Features:** | • Covers a broad spectrum: arithmetic reasoning, causal inference, joke explanation, logical fallacies, moral reasoning, and more.<br><br>• Encourages out-of-distribution performance testing.<br><br>• Tasks are designed to be hard for models, easy for humans. |
| **Highlighted Cognitive Tasks:** | • Logical deduction puzzles (e.g., Knights and Knaves)<br>• Causal judgment scenarios<br>• Planning and strategy games<br>• Analogy completion<br>• Meta-learning tasks (learn to learn) |
| **Why It's Significant:** | BIG-bench doesn't assume cognition as a single trait. It reflects multi-dimensional intelligence, pushing LLMs toward generalization, abstraction, and creativity. |

## 3. Psychometric-style IQ Tests for AI

Inspired by human intelligence tests like Raven's Progressive Matrices, these benchmarks simulate standardized intelligence evaluations for machines.

| | |
|---|---|
| **How They're Adapted:** | • Visual matrices are digitized and simplified for symbolic input.<br>• Tasks include series completion, pattern identification, and classification by abstract features.<br>• Some tests even measure working memory and fluid reasoning. |
| **Cognitive Capabilities Tested:** | • Abstract visual reasoning<br>• Inductive logic<br>• Non-verbal problem solving<br>• Spatial awareness (in symbolic form) |
| **Challenges for LLMs:** | • These tests require an internal model of the task – something LLMs often lack.<br>• They are highly structured yet demand flexible rule inference, a known challenge for LLMs trained on unstructured corpora. |

# 4. Theory of Mind (ToM) Benchmarks

Theory of Mind (ToM) – the ability to infer the beliefs, desires, and perspectives of others – is a hallmark of human cognition. Emerging benchmarks now evaluate LLMs on their capacity to perform ToM-related tasks.

| | |
|---|---|
| **Popular ToM Tests Adapted for LLMs:** | • **Sally-Anne Test:** Measures whether a model understands false belief. <br><br> • **Second-Order Belief Tasks:** "Alice thinks Bob believes X, but actually…" <br><br> • **Narrative Comprehension:** Can the model track mental states of multiple agents? |
| **Why It Matters:** | ToM is foundational in: <br><br> • Understanding user intent <br><br> • Generating empathetic, context-aware responses <br><br> • Social reasoning in AI agents (e.g., conversational bots, teaching AIs) |
| **Notable Findings:** | • GPT-4 shows partial ToM capabilities in zero-shot scenarios but still fails at deeper recursive tasks. <br><br> • Prompting strategies and chain-of-thought reasoning improve ToM scores, suggesting ToM may be emergent with scale and context. |

## 5. Commonsense Reasoning Benchmarks

Benchmarks like **PIQA, CommonsenseQA,** and **Hellaswag** evaluate the model's ability to apply everyday logic.

| What They Test: | • Cause-effect reasoning<br>• Object affordances ("What can you do with a sponge?")<br>• Physical and social intuition<br>• Pragmatic inference |
|---|---|
| Cognitive Relevance: | Commonsense is a foundational layer of cognition. Without it, even the most sophisticated model can fail in basic interaction. |

## 6. Mathematical Reasoning Benchmarks

Benchmarks like **GSM8K, MATH,** and **MATHQA** focus on symbolic, numerical, and multi-step problem solving.

| Cognitive Capabilities Tested: | • Logic chaining<br>• Symbol manipulation<br>• Memory and recursive function application |
|---|---|
| What's Interesting: | • LLMs like GPT-4 show promise with step-by-step prompting ("chain-of-thought").<br>• Fine-tuning on math-heavy corpora (e.g., Minerva, AlphaCode) boosts symbolic reasoning.<br>• However, hallucinations and arithmetic inconsistencies remain, indicating partial cognitive reasoning. |

# 9. LLM Testing Metrics: Quantitative vs. Qualitative

## 9.1 Quantitative Metrics for LLMs

These are numerical, automatable, and often used for benchmarking and model evaluation.

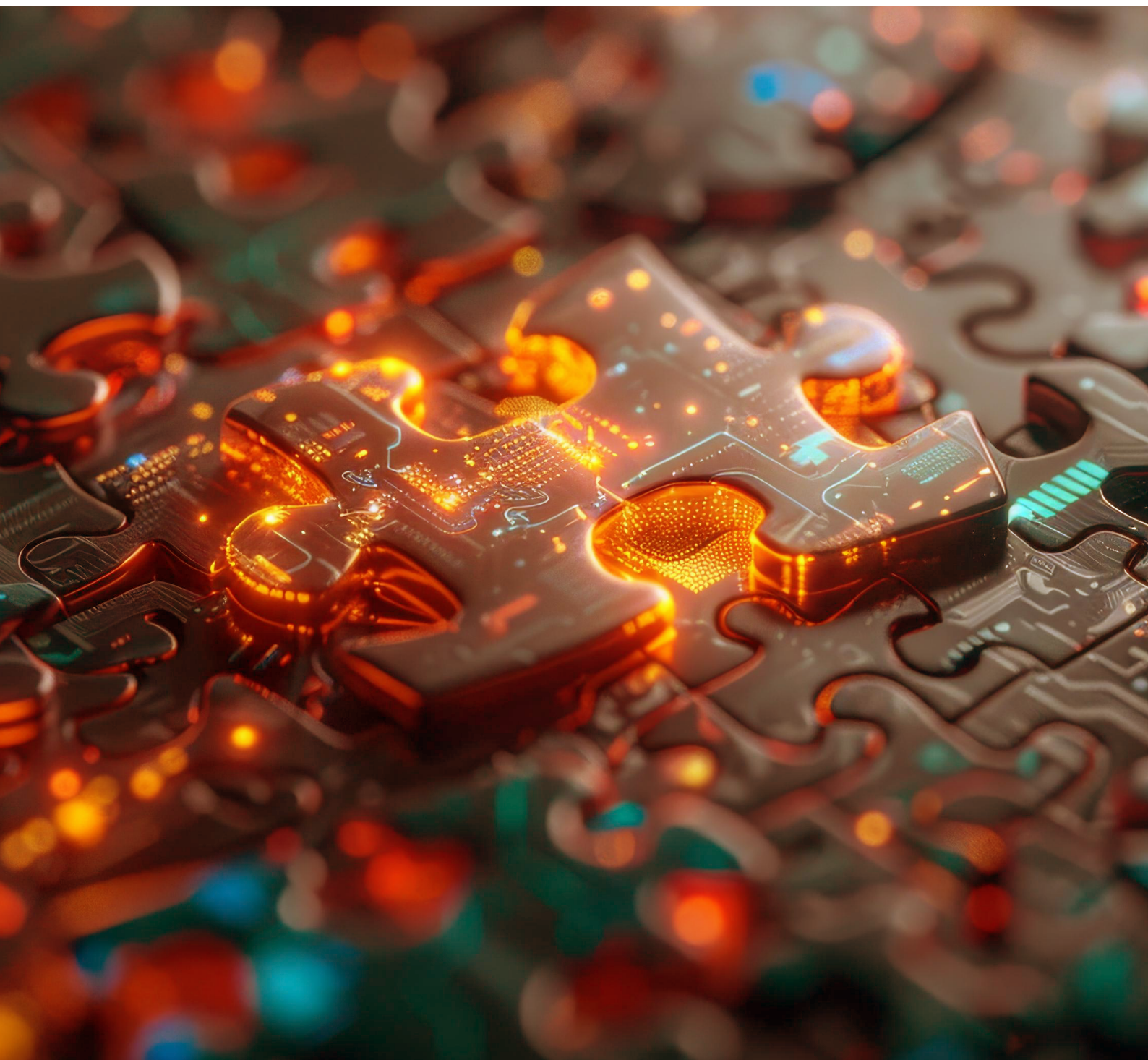| Metric | Purpose |
|---|---|
| Perplexity | Measures how well the model predicts a sequence. Lower = better. |
| BLEU / ROUGE / METEOR | Evaluate the similarity between the generated and reference text. |
| Accuracy / Precision / Recall / F1 | Used for classification or structured outputs (e.g., intent recognition). |
| Exact Match (EM) | Percentage of exact output matches (used in QA tasks). |
| Pass@k | Measures if the correct output is among the top-k completions. Useful in code generation. |
| Hallucination Rate | Frequency of factually incorrect outputs. |
| Toxicity Score | Measures offensive or harmful content. |
| Bias Metrics | Quantifies fairness and disparity across gender, race, etc. |
| Inference Time & Token Latency | Speed per token/output generation. |
| Prompt Execution Cost | Tracks token usage and computes cost per prompt. |

## 9.2 Qualitative Metrics for LLMs

These metrics require human judgment or heuristic evaluation to assess reasoning, logic, style, and usefulness.

| Metric | Description |
|---|---|
| Relevance & Grounding | Is the output fact-based and linked to source data? |
| Coherence | Logical flow and structure of response |
| Accuracy | Correctness of the output based on facts, data, or prompt intent. |
| Fluency & Grammar | Human-like language use, tone, and readability |
| Context Retention | Ability to remember and use previous turns or instructions |
| Reasoning Depth (Chain-of-Thought Quality) | Logical and stepwise problem-solving capability |
| Consistency | Stability of response to semantically similar inputs |
| Creativity | Novelty, uniqueness, and appropriateness in generative tasks |
| Bias & Ethical Sensitivity | Subjective check on fairness, cultural sensitivity, inclusion |
| User Satisfaction / Trust | Real-user feedback from testers or customers |
| Instruction Following | How well the LLM obeys prompt instructions |

As a matter of fact, **human evaluation remains the gold standard in LLM testing.** While quantitative metrics offer objective, scalable insights, qualitative evaluation shines, bringing in human judgment to assess coherence, creativity, and ethical sensitivity. Together, both approaches form a comprehensive testing strategy.

# 10. Case Study: Cheryl's Birthday Puzzle – A Window into the Cognitive Constraints of LLMs

In a compelling demonstration of the reasoning boundaries of Large Language Models (LLMs), Fernando Perez-Cruz and Hyun Song Shin evaluated GPT-4 using the well-known "Cheryl's Birthday" puzzle, popularized initially in 2015. This logic puzzle is widely recognized for its virality and the depth of inferential and counterfactual reasoning it requires. Specifically, it tests an individual's (or model's) ability to reason about nested knowledge -what one knows, what others know, and what others know about one's own knowledge.

## The Puzzle

Cheryl presented a puzzle to her two friends, Albert and Bernard, and asked them to figure out the exact date of her birthday. It is common knowledge among all three that Cheryl's birthday falls on one of the following ten dates:

| May | June | July | August |
|---|---|---|---|
| 15th,16th, 19th | 17th, 18th | 14th, 16th | 14th,15th, 17th |

To assist them in solving the puzzle, Cheryl gives each of them a small clue - she tells Albert only the month she was born and Bernard only the day of the month. Apart from this, she gives them no further information.

At the start, Albert and Bernard understand that they do not yet possess enough information to determine Cheryl's exact birthday. Moreover, they are not allowed to share the information they were given with each other.

Then Albert makes a statement:
**"I do not know what Cheryl's birthday is, but I am certain that Bernard doesn't know it either."**

Upon hearing this, Bernard responds:
**"Now that I've heard your statement, I can deduce Cheryl's birthday."**

After Bernard says this, Albert replies:
**"Now that I've heard what you just said, I, too, can determine Cheryl's birthday."**

**Question:** Based on this conversation and their reasoning, what is Cheryl's birthday?

A study in 2023 stated that GPT-4 solved the original version of this puzzle flawlessly. However, the researchers Fernando Perez-Cruz and Hyun Song Shin took the analysis further: they slightly altered peripheral details, changed character names, or replaced specific dates with alternatives. Despite the puzzle's structural integrity remaining unchanged, GPT-4's performance degraded significantly under these minimal modifications.

This stark decline revealed a critical limitation. It suggested that GPT-4's success with the original puzzle likely stemmed not from a fundamental understanding of the reasoning logic involved but from memorizing previously seen examples during training. The fact that minor cosmetic changes caused the model to fail implies that the system lacks robust generalization in counterfactual and epistemic reasoning domains.

The contrast between flawless logic when faced with the original wording and poor performance when faced with incidental changes in wording was very striking. It was difficult to dispel the suspicion that even when GPT-4 got it right (with the original wording), it did so due to familiarity with the phrasing rather than relying on the necessary steps in the analysis. In this regard, the apparent mastery of the logic appeared superficial.

# Conclusion

Recent AI and information technology breakthroughs are sparking fresh conversations about the potential for deeper collaboration between humans and AI. The emergence of Human-Aware AI has gained significant traction, creating AI systems that can seamlessly adapt to human cognitive strengths and limitations.

Concepts like "AI partners," "collaborators," and "buddies" emphasize a vision of AI that works as an equal, enhancing the capabilities of human teams. These AI agents need to be equipped with human-like cognitive abilities, enabling them to understand emotions, motivations, attention, and creativity to foster true collaboration.

As Large Language Models (LLMs) continue to evolve, benchmarking their cognitive capabilities is no longer a theoretical exercise but a critical necessity.

While LLM models exhibit impressive linguistic fluency and surface-level intelligence, true cognitive robustness remains an ongoing pursuit. By rigorously testing their limits across dimensions such as abstraction, memory, logic, and common-sense reasoning, we can chart a more straightforward path toward responsible deployment and future innovation. The goal is not only to measure what these models can do, but to understand better how they think, pushing the frontier of artificial intelligence from imitation to genuine understanding.

However, with advancements in reinforcement learning, neuromorphic computing, and hybrid models, the potential for AI to achieve human-like intelligence is on the horizon. By continuing to refine these technologies and addressing ethical concerns, **AI may one day achieve cognitive complexity that mirrors human intelligence, but until then, understanding its limitations is key.**

# References

https://arxiv.org/abs/2302.04322 https://arxiv.org/abs/2005.14165

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5859982/

https://www.nature.com/articles/s41586-020-1896-2

https://www.brookings.edu/research/ethics-of-artificial-intelligence-and-machine-learning/

https://web.stanford.edu/class/cs234/ https://arxiv.org/abs/2106.09685

https://pubmed.ncbi.nlm.nih.gov/20213435/

https://www.ibm.com/blogs/research/2020/06/ai-bias/

https://www.researchgate.net/publication/373642018_Large_language_models_in_medicine_the_potentials_and_pitfalls

https://promptengineering.org/beyond-the-hype-how-to-test-llm-for-intelligence-accuracy-and-reliability/

https://arxiv.org/abs/2408.09150 https://news.mit.edu/2024/reasoning-skills-large-language-models-often-overestimated-0711

https://www.bis.org/publ/bisbull83.pdf https://www.datacamp.com/blog/what-is-prompt-engineering-the-future-of-ai-communication

# INDIUM

## About Indium

Indium is an AI-driven digital engineering company that helps enterprises build, scale, and innovate with cutting-edge technology. We specialize in custom solutions, ensuring every engagement is tailored to business needs with a relentless customer-first approach. Our expertise spans Generative AI, Product Engineering, Intelligent Automation, Data & AI, Quality Engineering, and Gaming, delivering high-impact solutions that drive real business impact.

With 5,000+ associates globally, we partner with Fortune 500, Global 2000, and leading technology firms across Financial Services, Healthcare, Manufacturing, Retail, and Technology—driving impact in North America, India, the UK, Singapore, Australia, and Japan to keep businesses ahead in an AI-first world.