

# Assessing the Reproducibility of Research Utilizing Gemini and Other Large Language Models

**Author:** BenchChem Technical Support Team. **Date:** December 2025

## Compound of Interest

Compound Name: *Gemin A*

Cat. No.: *B1258876*

[Get Quote](#)

A Comparative Guide for Researchers, Scientists, and Drug Development Professionals

The integration of large language models (LLMs) into scientific research holds the promise of accelerating discovery, from hypothesis generation to data analysis. However, the stochastic nature of these models raises critical questions about the reproducibility of research that relies on them. This guide provides an objective comparison of Gemini's performance with other leading alternatives, supported by experimental data, to help researchers navigate this evolving landscape. We delve into the reproducibility, accuracy, and specific capabilities of these models in various research contexts, offering a framework for their effective and reliable use.

## I. Performance Comparison in Research-Relevant Tasks

The selection of an appropriate LLM for a research task depends on a nuanced understanding of its strengths and weaknesses. The following tables summarize the performance of Gemini, GPT-4, Claude 3.5, and Llama 3 across several key benchmarks and research-oriented tasks.

### Table 1: Performance on General and Scientific Benchmarks

Benchmark	Gemini 1.5 Pro	GPT-4o	Claude 3.5 Sonnet	Llama 3	Task Description
MMLU (Massive Multitask Language Understanding)	90.0%	88.4%	79.0%	82.0%	General knowledge and problem-solving across 57 subjects.
HumanEval (Code Generation)	71.9%	90.2%	93.7%	81.7%	Python code generation from docstrings.
BioLLMBench (Bioinformatics Proficiency)	97.5% (Math)	91.3% (Domain Knowledge)	-	Struggled (Code)	A suite of 24 tasks in bioinformatics , including domain knowledge, coding, and data analysis. <a href="#">[1]</a>
Medical Diagnostics (Open-ended)	44.00%	64.00%	72.00%	-	Diagnostic accuracy based on clinical case descriptions. <a href="#">[2]</a>
Medical Diagnostics (Multiple-choice)	65.00%	95.00%	89.00%	-	Diagnostic accuracy with answer variants provided. <a href="#">[2]</a>

Table 2: Reproducibility and Consistency

Metric	Gemini	GPT-4	Claude 3.5	Task/Context
Run-to-Run Reproducibility	99%	-	-	Digital Pathway Curation in biomedical research.
Inter-Reproducibility	75%	-	-	Consistency across different runs under varying conditions in biomedical research.
Response Agreement (Open-ended Diagnostics)	93.00%	97.00%	93.00%	Percentage of identical responses in repeated diagnostic tasks. <a href="#">[2]</a>
Response Agreement (Multiple-choice Diagnostics)	99.00%	98.00%	97.00%	Percentage of identical responses in repeated diagnostic tasks with options. <a href="#">[2]</a>

## II. Experimental Protocols

To ensure the validity of the comparative data, it is essential to understand the methodologies employed in these evaluations. Below are detailed protocols from a key study in the field.

### BioLLMBench: A Framework for Evaluating LLMs in Bioinformatics

The "BioLLMBench" study by Sarwal et al. (2025) provides a robust framework for assessing LLM performance in bioinformatics.[\[1\]](#)[\[3\]](#)[\[4\]](#)[\[5\]](#)

Objective: To systematically evaluate the capabilities of GPT-4, Gemini (formerly Bard), and LLaMA in solving a range of bioinformatics tasks that mirror the daily challenges faced by researchers.[\[3\]](#)[\[4\]](#)[\[5\]](#)

#### Experimental Setup:

- Models Tested: GPT-4, Gemini, and LLaMA.
- Tasks: 24 distinct tasks across six key areas:
  - Domain-Specific Knowledge
  - Mathematical Problem-Solving
  - Coding Proficiency
  - Data Visualization
  - Research Paper Summarization
  - Machine Learning Model Development
- Experimental Runs: A total of 2,160 experimental runs were conducted to ensure statistical significance.
- Evaluation Metrics: Seven task-specific metrics were designed to assess various aspects of the LLM's responses, including accuracy, completeness, and executability of code.
- Contextual Response Variability Analysis: This was implemented to understand how responses varied when prompts were presented in a new chat window versus within an ongoing conversation.

#### Key Findings:

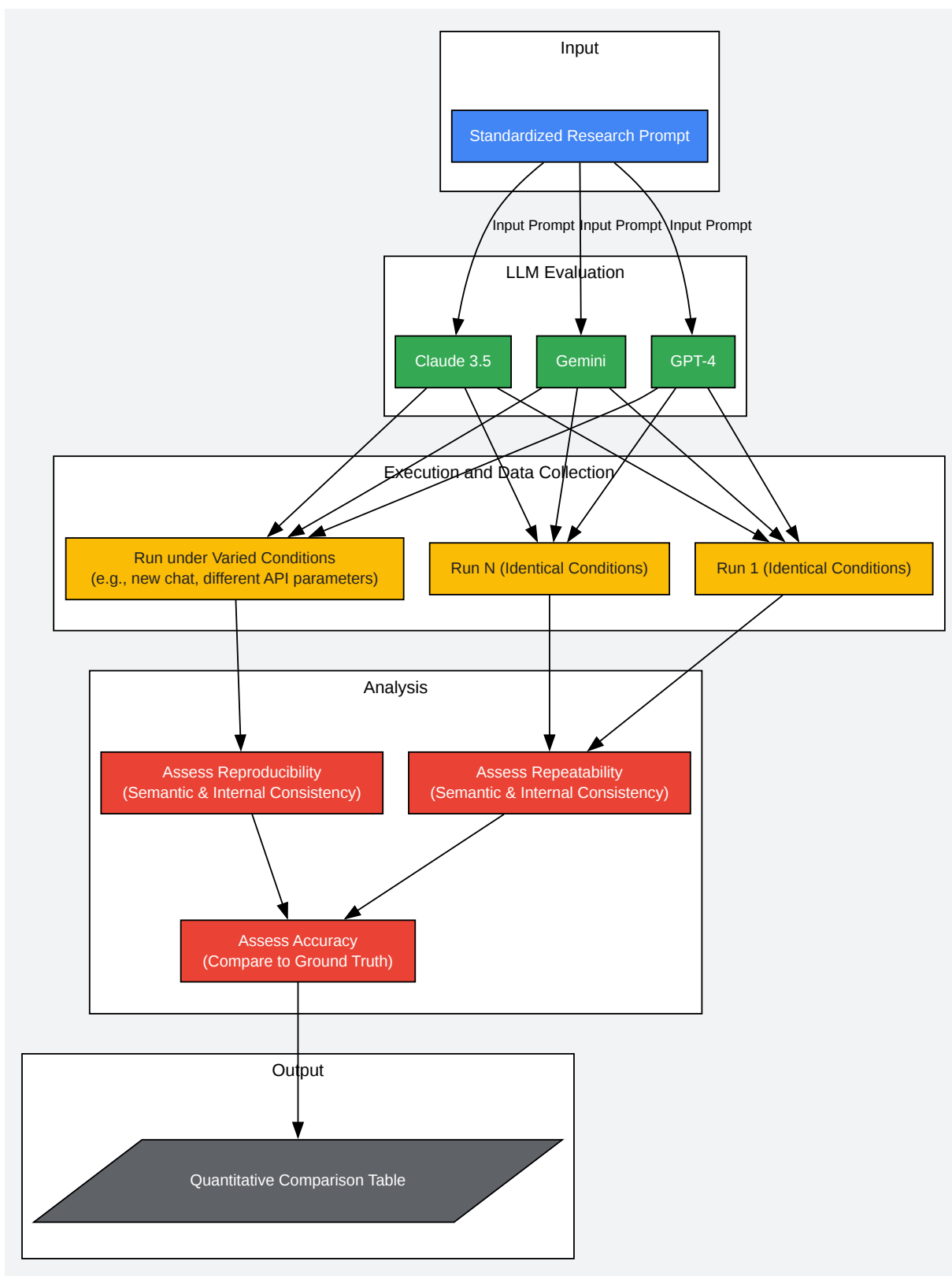
- Overall Performance: GPT-4 demonstrated the highest overall proficiency, particularly in domain knowledge and machine learning model development.[\[1\]](#)
- Gemini's Strength: Gemini excelled in mathematical problem-solving, achieving the highest proficiency score in this category.[\[1\]](#)
- Coding Challenges: While GPT-4 was proficient in generating functional code, both Gemini and LLaMA struggled to produce executable code for machine learning tasks.[\[1\]](#)
- Summarization Weakness: All models showed significant challenges in accurately summarizing research papers, with ROUGE scores below 40%.[\[1\]](#)

### III. Visualizing AI in Research Workflows and Biological Pathways

The application of LLMs in research can be conceptualized through various workflows. Furthermore, their potential to generate and analyze complex biological pathways is a key area of interest.

#### Logical Workflow for Assessing LLM Reproducibility

The following diagram illustrates a logical workflow for assessing the reproducibility of an LLM in a research context.



[Click to download full resolution via product page](#)

A logical workflow for assessing LLM reproducibility.

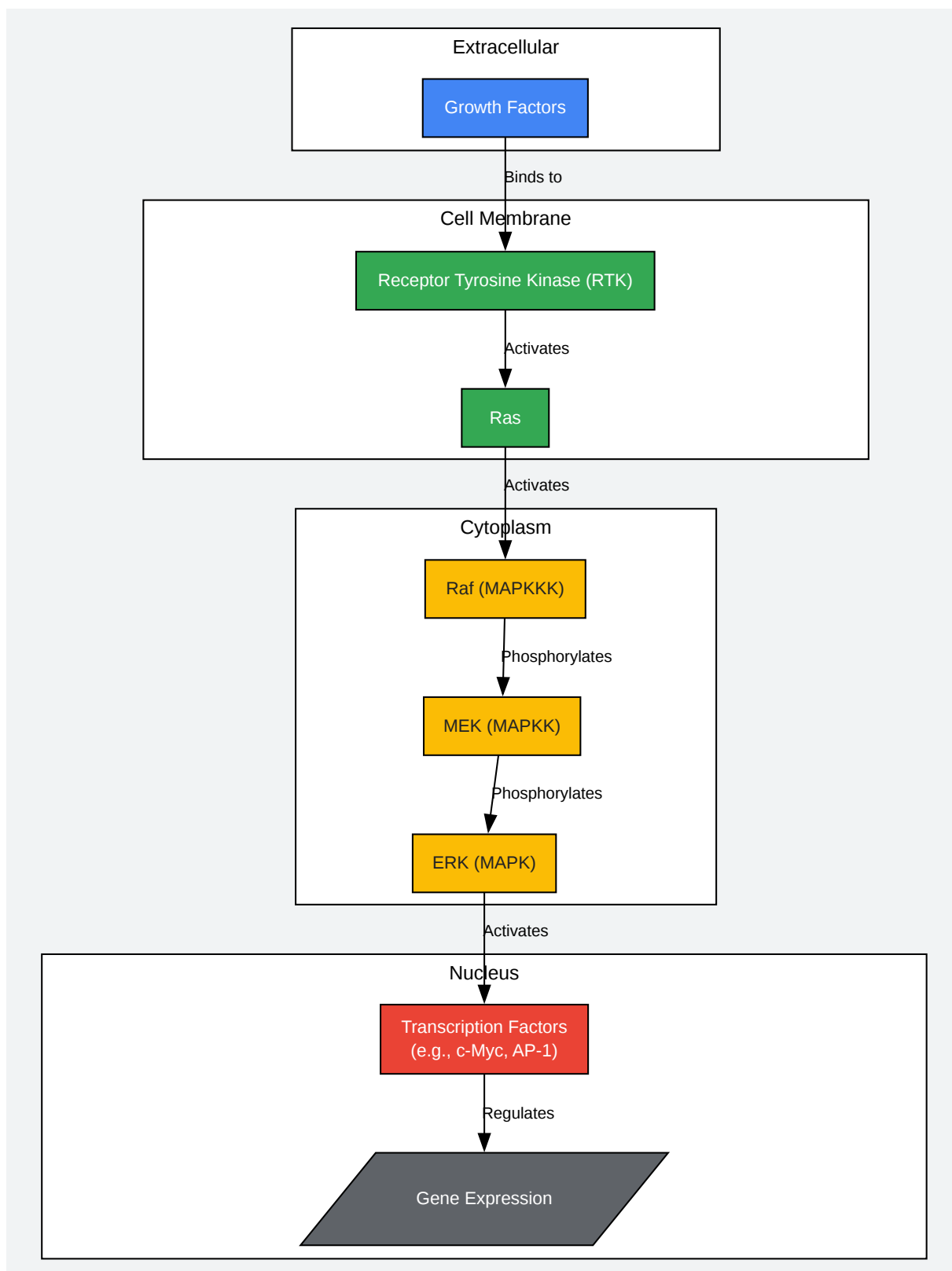
## Experimental Workflow for LLM-Assisted Bioinformatics Research

This diagram outlines a typical experimental workflow where an LLM is used to assist in bioinformatics research, based on the BioLLMBench protocol.

An experimental workflow for LLM-assisted bioinformatics research.

## Conceptual Representation of the MAPK Signaling Pathway

While the de novo generation of complex biological pathways by LLMs is still an emerging area, these models can be prompted to describe and structure known pathways. The following is a conceptual representation of the Mitogen-Activated Protein Kinase (MAPK) signaling pathway, a crucial pathway in cell proliferation, differentiation, and apoptosis, which can be used as a baseline for evaluating an LLM's understanding and descriptive capabilities.



[Click to download full resolution via product page](#)

Conceptual representation of the MAPK signaling pathway.



## IV. Conclusion and Recommendations

The reproducibility of research utilizing large language models is a multifaceted issue that requires careful consideration of the model's architecture, the specific research task, and the experimental protocol.

- For tasks requiring high factual accuracy and mathematical reasoning in bioinformatics, Gemini shows strong potential. However, its capabilities in generating complex, executable code may require further development.
- GPT-4 currently demonstrates a more robust all-around performance, particularly in domain-specific knowledge and code generation for machine learning tasks.
- Claude 3.5 Sonnet excels in code generation and shows high efficacy in medical diagnostic tasks, suggesting a strong capability in structured reasoning.
- Reproducibility is not guaranteed, even with the same model. Researchers should perform multiple runs to assess the consistency of the outputs and report on the variability.
- Human-in-the-loop validation is crucial. The outputs of LLMs, whether code, data analysis, or literature summaries, should be critically evaluated by domain experts.

As LLMs continue to evolve, standardized benchmarking and transparent reporting of experimental protocols will be paramount for ensuring the reliability and reproducibility of the scientific discoveries they help to facilitate. Researchers are encouraged to adopt systematic evaluation frameworks, such as BioLLMBench, to assess the suitability and consistency of these powerful tools for their specific research needs.

### *Need Custom Synthesis?*

*BenchChem offers custom synthesis for rare earth carbides and specific isotopic labeling.*

*Email: [info@benchchem.com](mailto:info@benchchem.com) or [Request Quote Online](#).*

## References

- 1. biorxiv.org [biorxiv.org]

- 2. From open-ended to multiple-choice: evaluating diagnostic performance and consistency of ChatGPT, Google Gemini and Claude AI - PubMed [pubmed.ncbi.nlm.nih.gov]
- 3. biorxiv.org [biorxiv.org]
- 4. researchgate.net [researchgate.net]
- 5. biorxiv.org [biorxiv.org]
- To cite this document: BenchChem. [Assessing the Reproducibility of Research Utilizing Gemini and Other Large Language Models]. BenchChem, [2025]. [Online PDF]. Available at: [https://www.benchchem.com/product/b1258876#assessing-the-reproducibility-of-research-that-utilizes-gemini]

---

### Disclaimer & Data Validity:

The information provided in this document is for Research Use Only (RUO) and is strictly not intended for diagnostic or therapeutic procedures. While BenchChem strives to provide accurate protocols, we make no warranties, express or implied, regarding the fitness of this product for every specific experimental setup.

**Technical Support:** The protocols provided are for reference purposes. Unsure if this reagent suits your experiment? [[Contact our Ph.D. Support Team for a compatibility check](#)]

**Need Industrial/Bulk Grade?** [Request Custom Synthesis Quote](#)

## BenchChem

Our mission is to be the trusted global source of essential and advanced chemicals, empowering scientists and researchers to drive progress in science and industry.

### Contact

Address: 3281 E Guasti Rd  
Ontario, CA 91761, United States  
Phone: (601) 213-4426  
Email: [info@benchchem.com](mailto:info@benchchem.com)