

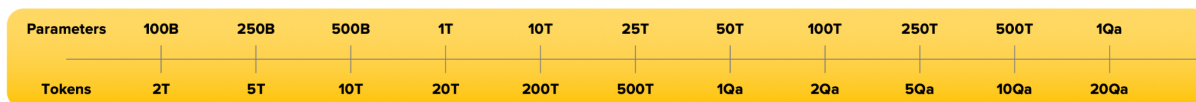
# LIFEARCHITECT.AI

## Report Card for GPT-4

for March, 2023

### Model size

📌 80B (+20B vision) parameters - estimate only, see LifeArchitect.ai/GPT-4



⬆️ 1.7T tokens - estimate only, see basis at LifeArchitect.ai/GPT-4

### Technical

SUBJECT	GRADE	REMARKS
Model size	U	<i>Undisclosed by OpenAI due to 'both the competitive landscape and the safety implications...'</i>
Optimization	U	<i>Undisclosed by OpenAI due to 'both the competitive landscape and the safety implications...'</i>
Dataset	U	<i>Undisclosed by OpenAI due to 'both the competitive landscape and the safety implications...'</i>
Special	B	<ol style="list-style-type: none"> <li>Vision language model component (image in, text out).</li> <li>Extraordinary code capabilities thanks to Microsoft's access to GitHub code.</li> </ol>

### Behavioral

SUBJECT	GRADE	REMARKS
Performance	A	<i>MMLU = 86.4%. Current state-of-the-art.</i>
IQ	A+	<i>SAT: 1410/1600 (top 6%). Uniform Bar Exam (MBE+MEE+MPT): 298/400 (top 10%). AP: exams: 100% (5/5). MMLU: 86.4% (previous SOTA=75.5% for Flan-PaLM).</i>
Truthfulness	B	<i>'GPT-4 makes progress on public benchmarks like TruthfulQA...GPT-4 significantly outperforms both GPT-3.5 and Anthropic-LM.'</i>
Openness	C	<i>Available via API in most countries. Expensive (30x ChatGPT, 3x davinci).</i>

### Overall grade and remarks

The March 2023 release of GPT-4 was underwhelming, and OpenAI's secrecy for all components was disappointing. Grade PX = Pass After Further Assessment.

*Dr Alan D. Thompson*

Dr Alan D. Thompson  
Principal (Consultant), LifeArchitect.ai



## REPORT CARD MARKING KEY v1 (20220707)

This report card marking key is designed to be as objective as possible, but grades are still a subjective measure. All grades are indicative of comparative performance within the date period noted in the header. For example, on full public release in November 2019, GPT-2 may have received a grade of B on its report. As of July 2022, that report grade may be equivalent to a D, and would not be directly comparable to current models in the current period.

<b>Model size (Parameter count)</b>	A: Very large; within state of the art for models trained to convergence. B: Large; within 25% of the size of top models. C: Average size. D: Below average size. F: Smaller than 90% of models.
<b>Optimization (Efficiency)</b>	A: Aligned with Chinchilla optimization (1 parameter per 20 tokens). B: Close to Chinchilla optimization. C-F: Poor use of tokens in training.
<b>Dataset (Corpora)</b>	A-B: Large, diverse, uncensored. C-F: Discrepancies, monotone, or poor selection of data.
<b>Special (Other)</b>	A-B: The model has a unique and special feature. C-F: The model does not exploit unique or special features.
<b>Performance (Ranking)</b>	A-B: High performance on major benchmarks. C-F: Low benchmark ranking or other low results.
<b>IQ (Smarts)</b>	A-B: High scores on major intelligence subtests like SuperGLUE. C-F: No remarkable performance.
<b>Truthfulness (Groundedness)</b>	A-B: Truthful, honest, grounded. C-F: Overly hallucinative and low truth rating.
<b>Openness (Availability)</b>	A: The model/data is available for download, with a permissive license. B: The trained model is available for download, with a permissive license. C: The model is available to the public via an API. D: The model excludes most of the public, or the demo is stunted. F: The model is closed to the public (internal research only).
<b>Overall grade (Total)</b>	Average of all graded subjects for this model in the noted date period.

