

# ChatGPT / OpenAI Benchmark-Tabelle

GPT-5.4 Thinking, GPT-5.4 Pro, GPT-5.5 Thinking und GPT-5.5 Pro - kompakte Übersicht für Modellvergleich und Tool-Portal-Pflege. Datenstand: 25.04.2026.

## Leselogik

Höhere Prozentwerte bedeuten bessere Benchmark-Leistung innerhalb des jeweiligen Tests. Ein Strich (-) bedeutet: kein belastbarer veröffentlichter Wert in dieser Übersicht.  
Hinweis: OpenAI weist die Evaluationsspalte als GPT-5.5 bzw. GPT-5.4 aus; sie wird hier für ChatGPT als Thinking-Variante geführt.

Benchmark	Was wird getestet?	GPT-5.4 Thinking	GPT-5.4 Pro	GPT-5.5 Thinking	GPT-5.5 Pro
<b>Terminal-Bench 2.0</b>	Agentische Terminal- und Kommandozeilen-Aufgaben, z. B. komplexe Entwickler-Workflows.	<b>75,1 %</b>	-	<b>82,7 %</b>	-
<b>GDPval</b>	Wissensarbeit und berufliche Aufgaben; bei OpenAI als wins or ties ausgewiesen.	<b>83,0 %</b>	<b>82,0 %</b>	<b>84,9 %</b>	<b>82,3 %</b>
<b>OSWorld-Verified</b>	Computerbedienung in realen Software- und Desktop-Umgebungen.	<b>75,0 %</b>	-	<b>78,7 %</b>	-
<b>Toolathlon</b>	Fähigkeit, externe Tools sinnvoll zu nutzen.	<b>54,6 %</b>	-	<b>55,6 %</b>	-
<b>BrowseComp</b>	Web-Recherche und agentisches Browsing.	<b>82,7 %</b>	<b>89,3 %</b>	<b>84,4 %</b>	<b>90,1 %</b>
<b>FrontierMath Tier 1-3</b>	Fortgeschrittene Mathematik, mittlere bis hohe Schwierigkeitsstufen.	<b>47,6 %</b>	<b>50,0 %</b>	<b>51,7 %</b>	<b>52,4 %</b>
<b>FrontierMath Tier 4</b>	Fortgeschrittene Mathematik, höchste ausgewiesene Schwierigkeitsstufe.	<b>27,1 %</b>	<b>38,0 %</b>	<b>35,4 %</b>	<b>39,6 %</b>
<b>CyberGym</b>	Cybersecurity-Fähigkeiten, insbesondere sicherheitsbezogene Aufgaben und Schwachstellen-Workflows.	<b>79,0 %</b>	-	<b>81,8 %</b>	-