

Einstellung der Modellparameter

Diese Parameter helfen dabei, die Vielfalt, Relevanz und Einzigartigkeit der generierten Antworten zu kontrollieren. Unerfahrenen wird jedoch empfohlen die Standardwerte beizubehalten.

Max. Kontext-Token

System

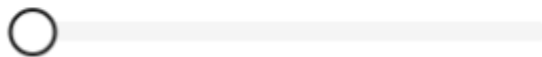
Temperatur

1.00



Frequency Penalty

0.00



Stop-Sequenzen

Trenne Stopwörter durch Drücken der `Enter`-Taste

Anhänge erneut senden



Web-Suche



Use Responses API



Max. Antwort-Token

System

Top P

1.00



Presence Penalty

0.00

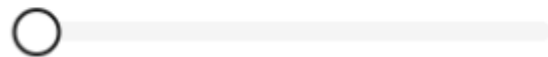


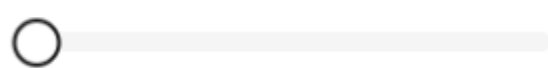
Bild-Detail

Auto



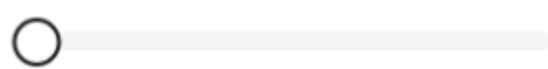
Denkaufwand

Keine



Reasoning Summary

Keine



Erläuterung der Modellparameter

- **Kontexttokens:** Dies bezieht sich auf die maximale Anzahl von Token (Einheiten von Wörtern

oder Zeichen), die das Modell auf einmal „verstehen“ oder verarbeiten kann. Bei vielen Sprachmodellen gibt es eine obere Grenze für die Anzahl von Token, die im Eingabekontext enthalten sein können. Wenn diese Grenze überschritten wird, muss der Text gekürzt oder geteilt werden.

- **Antwort Tokens:** Dies bezieht sich auf die Tokens, die im generierten Output oder in der Antwort des Modells enthalten sind. Wenn das Modell eine Antwort erzeugt, besteht diese Antwort aus einer bestimmten Anzahl von Token. Es gibt oft auch eine maximale Begrenzung dafür, wie viele Token die Antwort des Modells haben kann, abhängig von der Implementierung oder den Nutzerpräferenzen.
- **Temperatur** (Temperature): Dieser Parameter steuert den Zufallsgrad der Antworten. Ein niedrigerer Wert (z.B. 0.1) führt zu vorhersehbareren und konsistenteren Antworten. Ein höherer Wert (z.B. 1.0) erhöht die Vielfalt und macht die Antworten weniger vorhersagbar. Eine Temperatur von 0 würde immer dieselbe Antwort auf dieselbe Eingabeaufforderung geben (sofern der Rest der Einstellungen identisch ist).
- **Top-P** (Topical Penalization, auch bekannt als Nucleus Sampling): Dieser Wert bestimmt, wie das Modell aus einer Gruppe von wahrscheinlichen nächsten Worten auswählt. Ein niedrigerer Top-P-Wert (z.B. 0.5) bedeutet, dass nur die wahrscheinlichsten Worte berücksichtigt werden. Ein höherer Wert (z.B. 0.95) erlaubt eine größere Vielfalt potenzieller Worte. Dieser Mechanismus funktioniert durch das Beschränken der Auswahl an Wörtern auf ein kumulatives Wahrscheinlichkeitsintervall („P“).
- **Frequency Penalty** (Häufigkeitsstrafe): Dieser Parameter verhindert, dass das Modell dieselben Wörter oder Phrasen wiederholt. Je höher der Wert, desto größer die Strafe für die Wiederholung von Wörtern, was zu einer variableren und diversifizierteren Antwort führt.
- **Presence Penalty** (Anwesenheitsstrafe): Ähnlich wie die Frequency Penalty, aber dieser Parameter bestraft Wörter basierend darauf, ob sie bereits im Text erschienen sind, unabhängig von ihrer Häufigkeit. Ein höherer Wert fördert die Einführung neuer, bisher nicht verwendeter Wörter oder Phrasen in den Antworten.
- **Stop-Sequenzen:** Bis zu 4 Strings, bei deren Auftreten das Modell die Ausgabe beendet (stop-Parameter). Nützlich, um unerwünschte Fortsetzungen zu unterbinden. Bsp, definieren Sie eine oder mehrere Zeichenketten (z. B. „\n\n“, „END“), bei deren Auftreten das Modell die Ausgabe abbricht. Ideal, um Antworten konsistent zu beenden oder unerwünschte Fortsetzungen zu vermeiden.
- **Anhänge erneut senden:** Beim Fortsetzen von Gesprächen bestimmt dies, ob Dateien erneut versendet werden, wenn Sitzungen nicht persistent sind (Standard: false)
- **Bild-Detail:** Steuert den Detailgrad der in der Seitenleiste angezeigten Bildvorschau (Zoom-/Informationsstufe) – neu seit UI-Refresh. Die aktuelle Version von LibreChat (v0.7.x und früher) bietet keinen separaten Konfigurationspunkt zum Ausblenden oder Deaktivieren nur für Bilddetails. Wir haben die Funktion zur Bilderzeugung noch nicht geöffnet.
- **Web-Suche:** Aktiviert die eingebaute Websuche, um externe Informationen während der Konversation abzurufen (webSearch—Standard: true).
- **Denkaufwand:** Legt die Tiefe des Chain-of-Thought (Gedankengang) fest, die das Modell intern verfolgt. Höhere Stufen (z. B. „High“) nutzen mehr Rechenaufwand für komplexere Schlussfolgerungen.
- **Use Responses API:** Schaltet die Nutzung der OpenAI Responses API ein oder aus — eine alternative Streaming-Schnittstelle für Antworten (Standard: false).
- **Reasoning Summary:** Gibt nach Abschluss der Argumentationskette eine kurze Zusammenfassung des Gedankengangs aus, um die Nachvollziehbarkeit der Antwort zu

verbessern.

Direkt-Link:

https://doku.tu-clausthal.de/doku.php?id=sonstige_dienste:ki-dienste:librechat:einstellungmodellparameter&rev=1752036768

Letzte Aktualisierung: **04:52 09. July 2025**

