

KATALYSIA RESEARCH · PUBLIC

---



# Glossar -- Abkuerzungen des LLM-Trainings-Lebenszyklus

Wiederverwendbares Modul (Pre-Training, SFT, RLHF, DPO, RLAIIF, RLVR, Safety/Red-Team)

*Eigenstaendiges Glossar-Modul, einsetzbar in Vollfassung, Executive Summary und Presentation*

**Stand: 12. Mai 2026**

*Basis: Claude Opus 4.6 · KATALYSIA Research*

---

*Vertraulich · Interne Arbeitsunterlage · Weitergabe nur nach Rücksprache*

**GLOSSAR v1.0**

**Wiederverwendbares Modul** für die KATALYSIA-Studie *Human-AI-Training und Integration von Spezialwissen in Large Language Models* (v1.0, 12. Mai 2026).

*Dieses Glossar löst die Abkürzungen aus der Trainings-Pipeline auf, die in Folie 4 der Executive Presentation, in Kapitel 2.1 der Vollfassung und in der Markt-Analytik (Kapitel 5) konsistent verwendet werden. Es ist als eigenständiges Modul angelegt, damit es ohne Änderung in die Executive Summary, die Vollfassung oder künftige Folien-Versionen übernommen werden kann.*

## Trainings-Phasen und Verfahren

Abk. / Phase	Langform	Kurzbeschreibung
<b>Pre-Training</b>	Pre-Training	Selbst-überwachtes Lernen auf großen, unstrukturierten Textkorpora (Common Crawl, Bücher, Code, Web). Liefert ein generelles Sprach- und Weltmodell.
<b>Mid-Training</b>	Mid-Training	Continued Pre-Training auf kuratierten Hochqualitätsdaten und Synthetic Data; erschließt fachliche Tiefe (Lehrbücher, Code-Repositoryn, mathematische Beweise).
<b>SFT</b>	Supervised Fine-Tuning	Nachtraining auf 10 000–500 000 hochwertigen Instruction-Response-Paaren. Liefert Instruction-Following, Stil und Anweisungsbefolgung.
<b>RLHF</b>	Reinforcement Learning from Human Feedback	Drei-Stufen-Pipeline (SFT → Reward-Modell aus Präferenzpaaren → RL-Training, meist PPO oder GRPO). Prägt Helpfulness, Honesty, Harmlessness.
<b>DPO</b>	Direct Preference Optimization	Präferenzbasiertes Training ohne explizites Reward-Modell und ohne RL-Phase. Methodisch einfacher und compute-ärmer als RLHF; Varianten IPO, KTO, ORPO, SimPO.
<b>RLAIF</b>	Reinforcement Learning from AI Feedback	Bewertung durch ein anderes (oft größeres) LLM statt durch Menschen. Skaliert Volumina; Anthropic kombiniert RLAIF mit Constitutional AI zu RLHAIF.
<b>RLVR</b>	Reinforcement Learning from Verifiable Rewards	Reasoning-Training auf Aufgaben mit verifizierbaren Belohnungen (Mathe, Code, Logik). Dominanter Hebel der Frontier-Performance 2024–26.
<b>Safety / Red-Team</b>	Sicherheits- und Adversarial-Training	Refusal-Datensätze, Jailbreak-Tests, Capability-Evaluierung. EU-AI-Act Art. 55 verlangt dokumentiertes adversariales Testen für GPAI mit systemischem Risiko.

## Ergänzende Begriffe und Varianten

Abk.	Langform	Kontext
<b>PPO</b>	Proximal Policy Optimization	Klassischer RL-Algorithmus in RLHF (OpenAI).
<b>GRPO</b>	Group Relative Policy Optimization	RL-Variante, in DeepSeek-R1 prominent verwendet.
<b>IPO</b>	Identity Preference Optimization	DPO-Variante, robuster gegen Overfitting auf deterministische Präferenzen.
<b>KTO</b>	Kahneman-Tversky Optimization	Arbeitet mit Einzelannotationen statt Paaren.
<b>ORPO</b>	Odds Ratio Preference Optimization	Kombiniert SFT und Präferenz-Lernen in einem Schritt.
<b>SimPO</b>	Simple Preference Optimization	Längen-normalisierte DPO-Variante.
<b>Constitutional AI (CAI)</b>	Constitutional AI	Prinzipiengeleitetes Alignment (Anthropic): Verfassung steuert Kritik und Revision; mit RLHF kombiniert ergibt RLHAIF.
<b>GPAI</b>	General-Purpose AI	EU-AI-Act-Klassifikation; bei systemischem Risiko Pflicht zu adversarialem Testen (Art. 55).
<b>SME</b>	Subject Matter Expert	Fachexperte als Datenquelle in SFT- und Domain-Trainings-Programmen.
<b>PaaS</b>	Platform-as-a-Service	Geschäftsmodell, z. B. Snorkel, Labelbox, Toloka.

KATALYSIA Research · Public · Stand: Mai 2026 · Erstellt mit Claude (Anthropic) · Opus 4.6.