

Generative AI

Neue Inhalte auf Basis bereits existierender automatisch generieren

Die im Internet verfügbaren Texte, Bilder und andere Medien lassen sich zu umfangreichen Datenmodellen kombinieren. Mit natürlichsprachlichen Befehlen können dann neue Inhalte auf Basis dieser Datenmodelle generiert werden.

Definition

Generative AI basiert auf einer Transformer-Architektur für neuronale Netze, nachzulesen unter "[Attention Is All You Need](#)". Diese Architektur ermöglicht es, Muster in langen Datenreihen effektiv zu erfassen und ist dann in der Lage, aus komplexen Eingaben neue komplexe Ausgaben zu generieren. So lassen sich Texte nicht nur wortweise, sondern auch sinngemäß übersetzen. Diese Architektur ersetzt zunehmend generative KI-Ansätze, die meist auf rekurrenten neuronalen Netzen (RNN) beruhen. Ein Erfolgsfaktor des Transformers ist, dass sich die trainierten Modelle für spezifische Aufgaben nachträglich feinjustieren lassen.

Die Transformer-Architektur setzt sich grundsätzlich aus Encoder- und Decoder-Modellen zusammen, beides sind neuronale Netze. Ein Tokenizer zerlegt eine Eingabe in kleinere Einheiten wie Wörter oder Silben. Diese Einheiten erhalten durch numerische Vektor-Embeddings eine semantische Bedeutung und werden in den entsprechenden Kontext gebracht. Angegeben sind auch positionelle Informationen, etwa an welcher Stelle jedes Wort im Absatz steht. Diese Self-Attention-Information ermöglicht es, komplexe Muster zu erfassen und im Encoder sowie Decoder abzubilden. Die Qualität und Menge der zugrunde liegenden Trainingsdaten beeinflussen maßgeblich die Qualität.

Head-Modelle spezialisieren sich auf

bestimmte Aufgaben. Bei der *Klassifizierung* von Texten, etwa der Erkennung von Emotionen, kommt bevorzugt der Encoder zum Einsatz. Bei reinen *Generatoren*, wie etwa ChatGPT, liegt der Fokus hingegen hauptsächlich auf dem Decoder. Dieser versucht, eine Eingabe in einem Prompt schrittweise zu erweitern, indem er neue Wörter hinzufügt, die wiederum in die Eingabe zurückfließen. Bei *Textübersetzungen* agiert die gesamte Transformer-Architektur, ohne eine Ausgabe beim Encoder zu generieren.

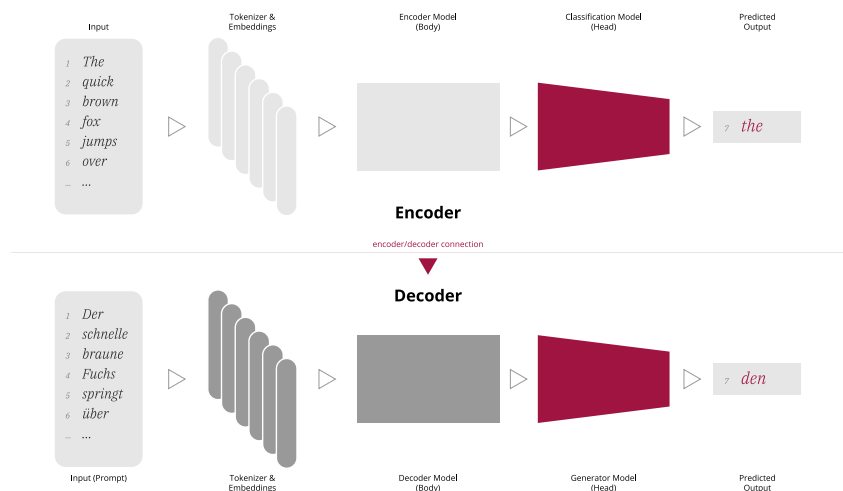
Die Ergebnisse basieren auf Wahrscheinlichkeiten und Vorhersagen. Sie können daher nicht pauschal korrekt sein. Der Zufallsfaktor ist zudem die Ursache, dass der Eindruck von Kreativität entsteht. Das führt allerdings auch dazu, dass die KI zum "halluzinieren" neigt, also vollkommen irreführende Daten generiert.

Referenzszenario

Zur Beschleunigung bestimmter, wiederkehrender, fachlicher, aber programmatisch schwierig automatisierbarer Aufgaben führt ein Unternehmen Generative AI ein. Es muss zunächst datenschutzrelevante Aspekte klären, weil potenziell sensible Unternehmensinformationen oder Kundendaten übermittelt werden. Weil die Anwendungsfälle für Generative AI klar sind, kann es wählen, ob es ein kostenpflichtiges, fremdbetriebenes Angebot nutzt oder auf eine Open-Source-Lösung setzt. Um die Qualität der Ergebnisse sicherzustellen und Urheberchaftsaspekte zu klären, legt es Richtlinien fest und definiert Prozesse. Fachpersonal muss deshalb die von Generative AI erzeugten Ergebnisse validieren und freigeben.

Potenzial

Der Trainingsprozess großer und weitrei-



KI-Algorithmen

- ist Teilbereich der KI
- Wahrscheinlichkeiten mit Zufallsfaktoren
- Transformer nutzen Neuronale Netze
- Transformatoren lösen RNNs ab

Digitale Inhalte

- basiert auf digitalen Inhalten
- Internet als Datenquelle zum Anlernen
- Daten zu menschlicher Kultur
- Daten zu menschlichen Denkweisen



Hardware-Anforderungen

- moderne Hardware verarbeitet auch große Datenmengen
- spezielle Hardware (GPU) trainiert und nutzt Neuronale Netze

Recht

- EU AI Act definiert Rechtslage
- Nutzungsrecht aktuell beim Anwender
- Beweispflicht der Urheberschaft und Datenverarbeitung beim Anwender

chend einsetzbarer Modelle ist kostspielig und ressourcenintensiv. Jedoch ist eine Vielzahl vortrainierter Open-Source- oder lizenzierbarer Modelle auf dem Markt verfügbar. Bestehende Modelle lassen sich kostengünstig durch Feinjustierung auf spezifische Anwendungsfälle erweitern. Ergebnisse lassen sich nach inhaltlicher Prüfung oft direkt in der internen Kommunikation weiterverwenden. Rechtlich betrachtet besteht Unsicherheit. Eine KI ist in Deutschland keine juristische Person, daher gehören die generierten Inhalte dem Nutzer, falls die KI nicht zufällig zu ähnliche Inhalte aus geschützten Quellen regeneriert. Die Beweislast für Copyright- und Datenschutzverletzungen liegt beim Anwender.

in den 1950er Jahren vorgestellt. Der im Jahr 2017 vorgestellte Transformer-Ansatz kann besonders große Datenmengen verarbeiten. Die von ChatGPT der Firma OpenAI erzielten, bemerkenswerten Fortschritte offenbaren die Leistungsfähigkeit dieser Modelle. Die Technologie ist mittlerweile auch in der Open-Source-Welt ausgereift, neue Modelle kommen in regelmäßigen Abständen von der Community.

Marktübersicht

Seit Anfang 2023 wächst das Angebot öffentlicher Modelle wöchentlich. Plattformen wie Huggingface.co verwalten mittlerweile über 260.000 Open-Source Transformer-Modelle und bieten eine einheitliche Schnittstelle dafür an. Derzeit ist der Markt noch bereit, Gebühren für die Nutzung von Modellen in der

Cloud zu zahlen, möglicherweise auch deshalb, weil das Hosting aktueller Modelle eine starke GPU-Leistung erfordert.

Alternativen

Wenn es um die Klassifizierung, Generierung oder Zusammenfassung von natürlichen Texten geht, sind Transformer-Lösungen praktisch alternativlos. Die menschliche Sprache ist zu komplex, als dass frühere KI-Verfahren oder sogar konventionelle deterministische Verfahren zuverlässig auf neue, unbekannte Inhalte angemessen reagieren könnten. Bei kleinen Datenmengen oder speziellen Anwendungsfällen sind klassische KI-Methoden oft effizienter.

Fazit

- + verarbeitet Informationen und Muster großer Datenmengen
- + interagiert in natürlicher Sprache
- + Open-Source Modelle vorhanden
- + universell einsetzbare Modelle
- + Prompt-Engineering oder Finetuning passen fachlichen Kontext an
- kein Garant auf Richtigkeit
- nicht nachvollziehbare Ergebnisfindung
- sehr hohe Hardwareanforderung
- rechtliche Grundlage noch nicht definiert

Reifegrad

Generative AI-Ansätze wurden erstmals



Buzzword Factor (Ent./Customer)

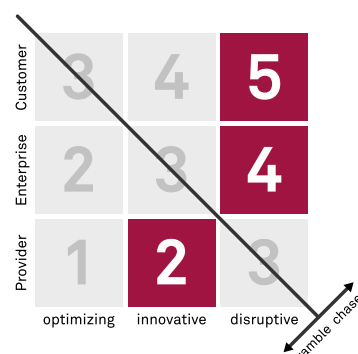
1 low	2 medium	3 high
-------	----------	--------

Entry Barrier (Provider)

1 low	2 medium	3 high
-------	----------	--------

Benefit Level (Provider)

1 low	2 medium	3 high
-------	----------	--------



<https://msg.direct/techrefresh>

Stand: Dezember 2023

msg systems ag