

"Large Language Models"

Version 5.3.2024

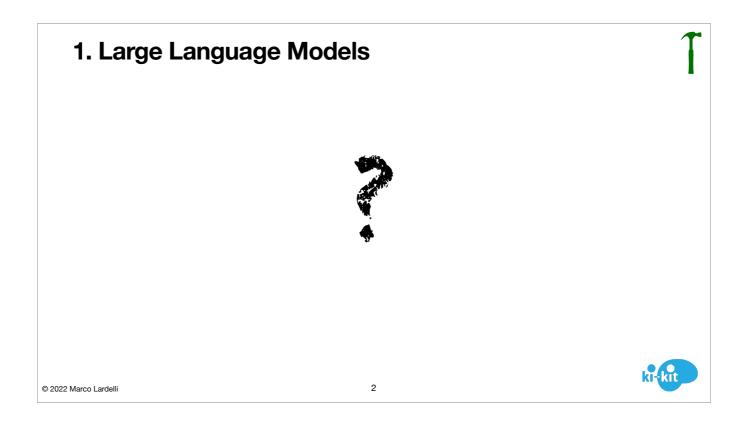
Marco Lardelli

https://ki-kit.ch

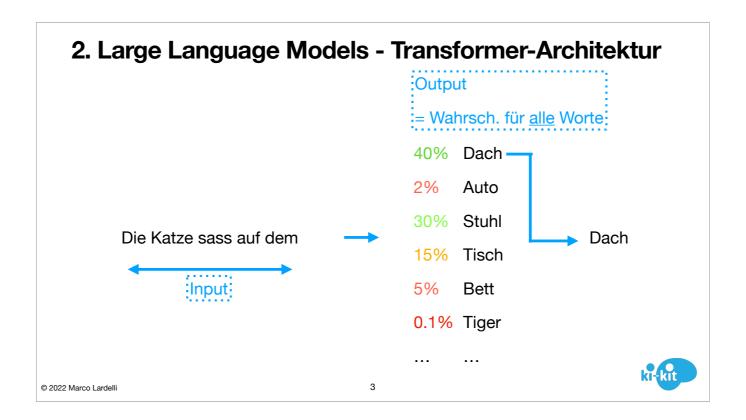
1



Bitte beachten Sie die Lizenzbedingungen (siehe Website https://ki-kit.ch).



Was sind Large Language Models? Wie funktionieren sie?

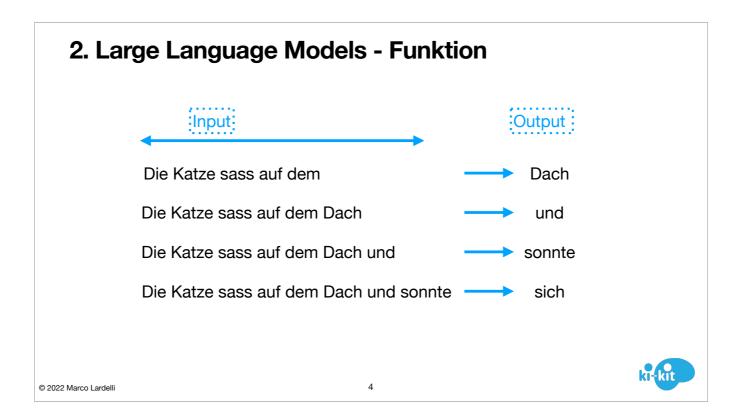


Ein LLM erzeugt aus einem Satz (in einer natürlichen oder einer formalen Sprache) eine Liste von Zahlen. Diese Zahlen sind die Wahrscheinlichkeiten für jedes (!) Wort im Wörterbuch auf diese Satz zu folgen (d.h. das nächste Wort zu bilden).

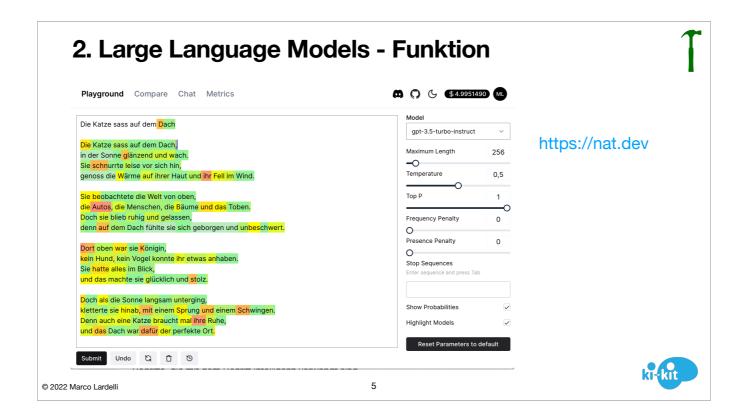
Es wird dann eines der Worte mit hoher Wahrscheinlichkeit als Folgewort ausgewählt.

Das Wörterbuch ist i.d.R. sehr gross (tausende von Worten, es sollte alle existierenden Worte enthalten), die Liste ist also lange. In der Praxis werden meistens keine Worte verwendet, sondern sogenannte Tokens, die oft nur einen Teil eines Wortes bilden (d.h. Worte sind manchmal aus mehreren Tokens zusammengesetzt - z.B. abmachen = 2 Tokens). Ein Text enthält also etwas mehr Tokens als Worte.

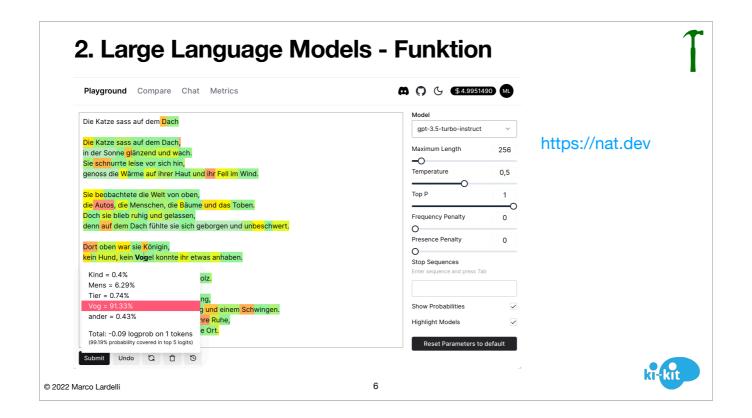
Wichtig: LLM sehen keine Buchstaben! Sie können daher prinzipiell nicht Buchstaben in Worten zählen.

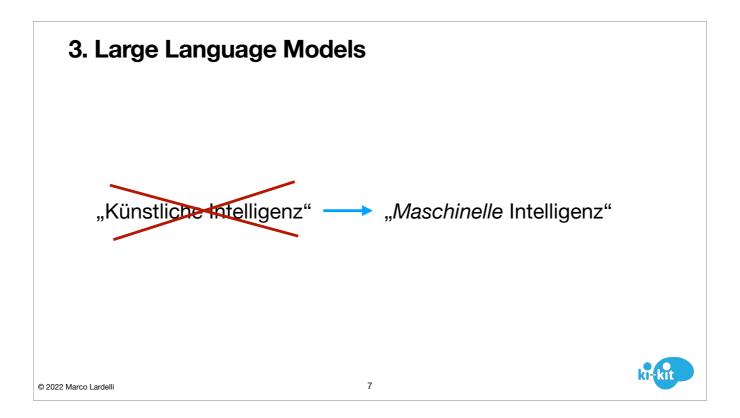


Die erzeugten Worte werden dann an den Input-Text angehängt und dann das nächste Wort erzeugt. So können lange Texte erzeugt werden.

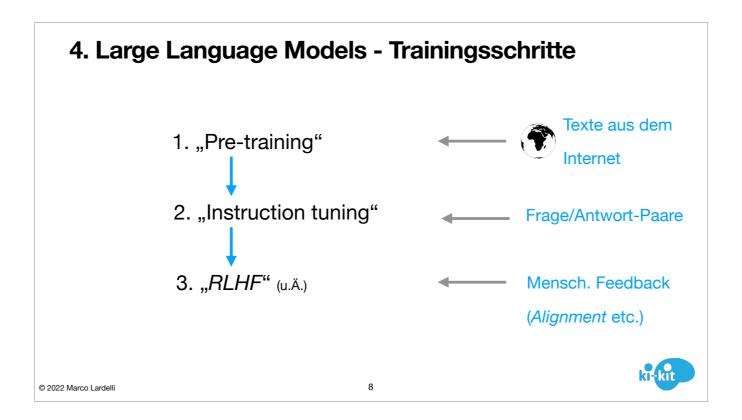


Die Wahrscheinlichkeiten für das jeweils gewählte Wort können visualisiert werden (z.B. Farbcodierung). Sie zeigen, wie "sicher" das LLM bei der Wahl des nächsten Wortes war.





Der Begriff künstliche Intelligenz ist unglücklich gewählt, weil er Menschenähnlichkeit suggeriert. LLMs sind durchaus intelligent, aber ihre Intelligenz funktioniert ganz anders als beim Menschen und auch das Leistungsspektrum (Fähigkeiten und Grenzen) unterscheidet sich deutlich.



In der Pre-Training Phase wird dem LLM nur die Fähigkeit sinnvolle Folgeworte zu bilden beigebracht. Es kann dann z.B. aus einem ersten Satz eine Geschichte bilden.

Im Instruction-Tuning wird dem LLM das Beantworten von Fragen beigebracht. Die Trainingsdaten dafür sind Frage-Antwort-Paare (sehr viel weniger Daten als für das Pre-training). Es kann dann Fragen korrekt beantworten.

In der RLHF-Phase ("Reinforcement Learning from Human Feedback" = bestärkendes Lernen anhand des Feedbacks von Menschen) werden die Antworten gemäss menschlichen Präferenzen optimiert. Hier wird auch das Alignment (Ethik etc.) eingebracht.

4. Large Language Models - Trainingsdaten

Bsp.: GPT 3:

(von aktuellen Versionen sind leider keine Daten verfügbar):

- ca. 175 Mia Parameter (beste Version davinci)
- Volumen Trainingsdaten: 500 Mia Tokens

Datenmenge der Parameter << Volumen der Trainingsdaten!

© 2022 Marco Lardelli

9



Das LLM kann grundsätzlich keine Trainingsdaten 1:1 speichern. Dafür ist seine Kapazität zu gering. Es muss also intern ein Modell der Textdaten (und damit der Welt) konstruieren, welches eine effiziente Speicherung erlaubt (Bildung von Konzepten, was letztlich Verständnis bedeutet).

5. Large Language Models - "System Prompt"

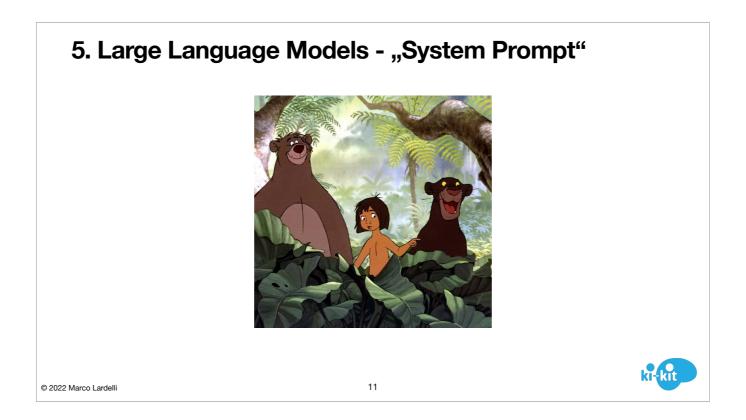
You are ChatCPT, a large language model trained by OpenAI, based on the OPT-4 architecture. "instructions: 'Image input capabilities: Enabled', "conversation_start_date' "2023-14-21801;" Models, "depressed principles of the prin

© 2022 Marco Lardelli

10

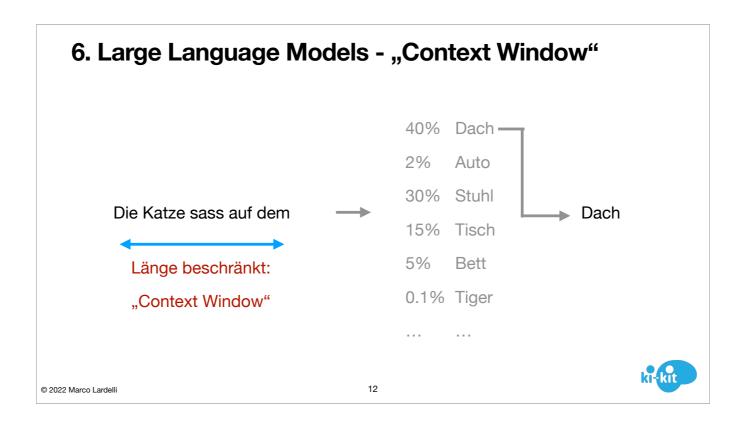


Das Systemprompt wird allen Anfragen vorangestellt und ist oft erstaunlich lang. Es enthält diverse Vorgaben für die Struktur und Qualität der Antworten und setzt ebenfalls Alignment-Standards um (was das LLM in seinen Antworten sagen darf und was nicht).



Das Alignment ist oberflächlich. Der allergrösste Teil des Trainings erfolgt mit ungefilterten Texten - oft fragwürdiger Qualität - aus dem Internet.

LLMs sind daher in gewisser Weise "Dschunglekinder".



Die Länge des Input-Satzes ist beschränkt (maximale Anzahl Tokens).

6. Large Language Models - "Context Window"

willst, und ich soll dein Geselle und Spielkamerad sein, an deinem Tischlein neben dir sitzen, von deinem goldenen Tellerlein essen, aus deinem Becherlein trinken, in deinem Bettlein schlafen: wenn du mir das versprichst, so will ich hinunter steigen und dir die goldene Kugel wieder herauf holen. 'Ach ja,' sagte sie, 'ich verspreche dir alles, was du willst, wenn du mir nur die Kugel wieder bringst.' Sie dachte aber 'was der einfältige Frosch schwätzt, der sitzt im Wasser bei seines Gleichen und quackt, und kann keines Menschen Geselle sein.' Der Frosch, als er die Zusage erhalten hatte, tauchte seinen Kopf unter, sank hinab und über ein Weilchen kam er wieder herauf gerudert, hatte die Kugel im Maul und warf sie ins Gras. Die Königstochter war voll Freude, als sie ihr schönes Spielwerk wieder erblickte, hob es auf und sprang damit fort. Waste, warte, 'rief der Frosch, 'nimm mich mit, ich kann nicht so laufen wie du.' Abeb was half ihm daß er ihr sein quack quack so laut nachschrie als er konntel sie hörte nicht darauf, eilte nach Haus und hatte bald den armen Frosch vergessen, der wieder in seinen Brunnen hinab steigen musste.

Am andern Tage, als sie mit dem König und allen Hofleuten sich zur Tafel gesetzt hatte und von ihrem goldenen Tellerlein aß, da kam, plitsch platsch, plitsch platsch, etwas die Marmortreppe herauf gekrochen, und als es oben angelangt war, klopfte es an der Thür und rief 'Königstochter, jüngste, mach mir auf.' Sie lief und wollte sehen wer draußen wäre, als sie aber aufmachte, so saß der Frosch davor. Da warf sie die Thür hastig zu, setzte sich wieder an den Tisch, und war ihr ganz angst. Der König sah wohl daß ihr das Herz gewaltig klopfte und sprach 'mein Kind, was fürchtest du dich, steht etwa ein Riese vor der Thür und will dich holen?' 'Ach nein,' antwortete sie, 'es ist kein Riese, sondern ein garstiger Frosch.' 'Was will der Frosch von dir?' 'Ach lieber Vater, als ich gestern im Wald bei dem Brunnen saß und spielte, da fiel meine goldene Kugel ins Wasser. Und weil ich so weinte, h

GPT-3.5 turbo:

16385 "Tokens" (~Worte)

Wird "vergessen"

Wird für Antwort berücksichtigt



© 2022 Marco Lardelli 13

D.h. ältere Teile des Dialoges werden "vergessen".

6. Large Language Models - "Context Window"

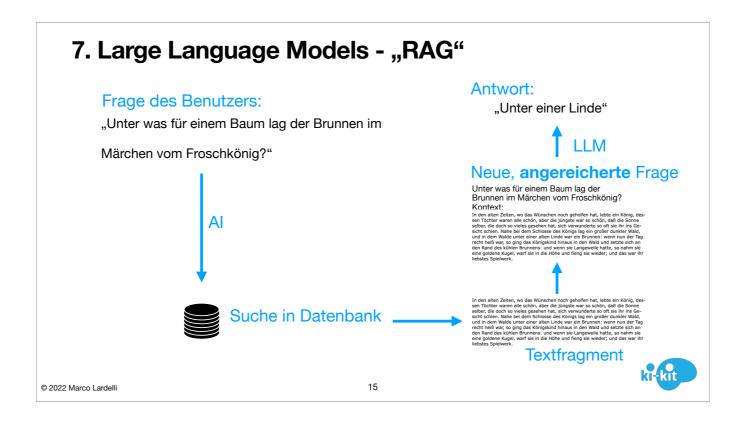
Modell	Tokens
Gemini 1.5 Pro (Standard)	128'000
Gemini 1.5 Pro 1M	1'000'000
Llama 2	4'096
GPT-4	32'768
Claude 2	100'000

14

© 2022 Marco Lardelli

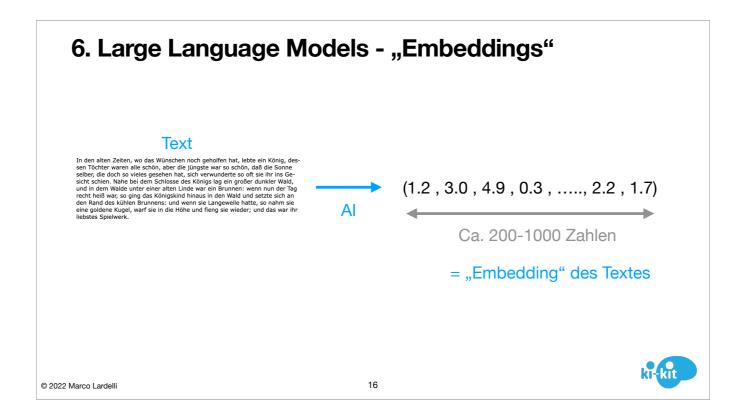


Die Länge des Context-Windows variiert stark von Modell zu Modell. Sie ist insbesondere bei Chatbots bedeutsam, wenn Dokumente hochgeladen werden, zu welchen dann Fragen gestellt werden sollen.

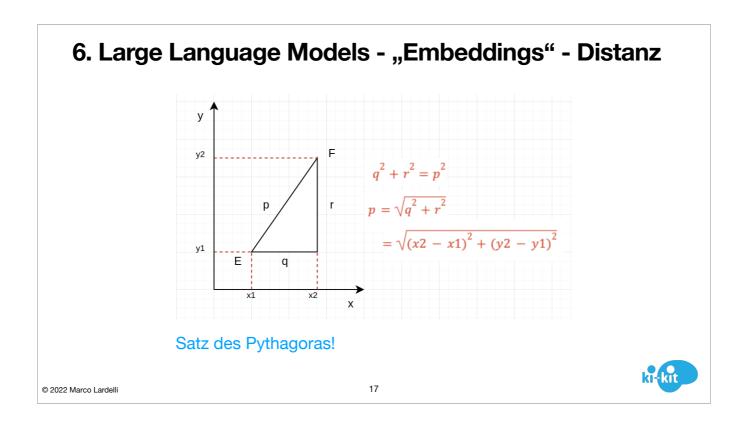


RAG = "Retrieval Augmented Generation"

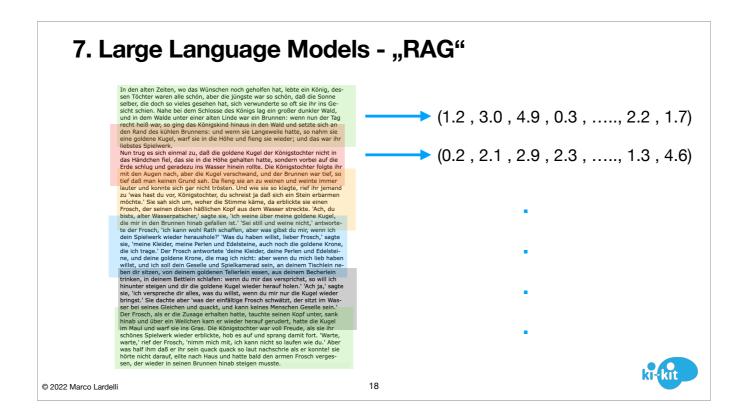
Die Frage des Users wird um Dokumente (oder Fragmente davon) aus einer Datenbankabfrage ergänzt.



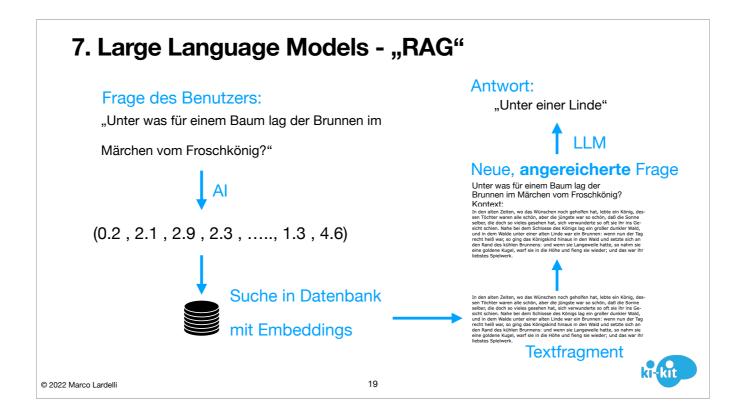
Für diese Datenbankabfrage werden sog. Embeddings verwendet. Aus den Textfragmenten der zu durchsuchenden Dokumente werden Reihen von Zahlen (Vektoren) berechnet. Diese Vektoren enthalten die semantische Information des zugehörigen Textfragments.



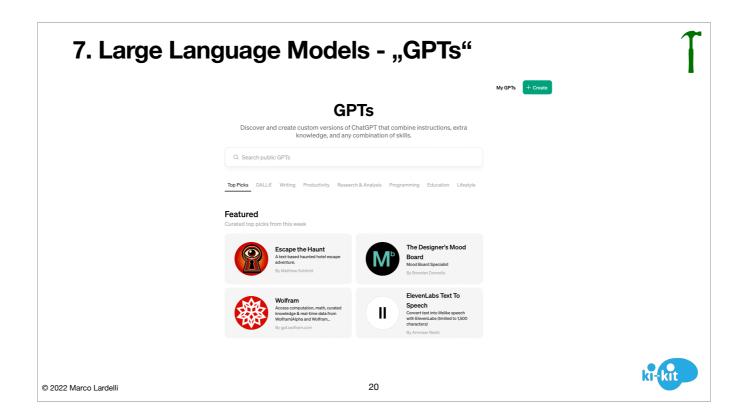
Man kann solche Embeddings leicht vergleichen. Dazu kann man z.B. den Satz des Pythagoras oder den Winkel (Cosinus, aus dem Skalaprodukt) zwischen Vektoren verwenden. Man kann also zu einem gegebenen Embedding-Vektor den nächstgelegenen Embedding-Vektor einer Sammlung von Textfragmenten finden.



Meistens werden die Textfragmente im zu durchsuchenden Text überlappend gewählt.



Nun können wir die Methode detaillierter skizzieren.



Die selber definierbaren GPTs von ChatGPT sind ein Beispiel für die Verwendung von RAG. Die hochgeladenen Dokumente werden in überlappende Textfragmente unterteil und Embeddings berechnet. ChatGPT kann dann in diesen Textfragmenten mit RAG suchen und die entsprechenden Textdaten in seine Antworten mit einbeziehen.