

Large Language Models selbst programmieren

Mit Python und PyTorch ein eigenes LLM entwickeln

» Hier geht's
direkt
zum Buch

DIE LESEPROBE

1 LLMs verstehen

In diesem Kapitel:

- Erläuterungen der grundlegenden Konzepte hinter Large Language Models (LLMs) im Überblick
- Einblicke in die Transformer-Architektur, von der LLMs abgeleitet werden
- Ein Plan für den Aufbau eines LLM von Grund auf

Large Language Models (LLMs, große Sprachmodelle), wie sie in ChatGPT von OpenAI angeboten werden, sind Modelle tiefer neuronaler Netze (*Deep Neural Networks*), die in den letzten Jahren entwickelt wurden. Sie haben eine neue Ära in der Verarbeitung natürlicher Sprache (*Natural Language Processing*, NLP) eingeläutet. Bevor LLMs aufgekomen sind, genügten herkömmliche Methoden vollauf bei Kategorisierungsaufgaben wie zum Beispiel E-Mail-Spam-Klassifizierung und einfacher Mustererkennung, die sich mit handgestrickten Regeln oder einfacheren Modellen erfassen ließen. Allerdings waren sie bei Sprachaufgaben, die komplexe Verständnis- und Generierungsfähigkeiten erfordern, wie zum Beispiel detaillierte Anweisungen parsen, Kontextanalysen durchführen sowie kohärent und kontextuell angemessene Originaltexte erzeugen, in der Regel unterlegen. Zum Beispiel konnten frühere Generationen von Sprachmodellen keine E-Mail aus einer Liste von Schlüsselwörtern schreiben – eine Aufgabe, die für moderne LLMs trivial ist.

LLMs besitzen bemerkenswerte Fähigkeiten, um menschliche Sprache zu verstehen, zu erzeugen und zu interpretieren. Allerdings müssen wir Folgendes klarstellen: Wenn wir sagen, dass Sprachmodelle etwas »verstehen«, meinen wir, dass sie Text in einer Weise verarbeiten und erzeugen können, der kohärent und kontextuell relevant erscheint, und nicht, dass sie menschenähnliches Bewusstsein oder Verständnis besitzen.

Dank der Fortschritte beim *Deep Learning*, einem Teilbereich des *Machine Learning* (des maschinellen Lernens) und der *künstlichen Intelligenz* (KI), der sich auf neuronale Netze konzentriert, werden LLMs mit riesigen Mengen von

Textdaten trainiert. Dieses groß angelegte Training versetzt LLMs in die Lage, im Vergleich zu früheren Ansätzen tiefere kontextuelle Informationen und Feinheiten der menschlichen Sprache zu erfassen. Infolgedessen haben LLMs die Leistung in einem breiten Spektrum von NLP-Aufgaben erheblich verbessert, einschließlich Textübersetzung, Stimmungsanalyse, Beantwortung von Fragen und vielem mehr.

Heutige LLMs und frühere NLP-Modelle unterscheiden sich zudem dadurch, dass frühere NLP-Modelle in der Regel für bestimmte Aufgaben wie Textkategorisierung, Sprachübersetzung usw. konzipiert wurden. Diese früheren NLP-Modelle konnten zwar in ihren eng gefassten Anwendungen brillieren, doch LLMs erweisen sich als kompetenter in einem breiten Spektrum von NLP-Aufgaben.

Der Erfolg der LLMs lässt sich auf die Transformer-Architektur zurückführen, die vielen LLMs zugrunde liegt, sowie auf die riesigen Datenmengen, mit denen LLMs trainiert wurden, sodass sie eine umfangreiche Palette an sprachlichen Nuancen, Kontexten und Mustern erfassen können, die manuell nur schwer zu codieren wären.

Dieser Übergang zur Implementierung von Modellen, die auf der Transformer-Architektur basieren und große Trainingsdatensätze verwenden, um LLMs zu trainieren, hat NLP grundlegend verändert, sodass jetzt leistungsfähigere Tools verfügbar sind, um menschliche Sprache zu verstehen und damit zu interagieren.

Die folgende Erörterung umreißt den Ausgangspunkt, um das primäre Ziel dieses Buchs zu erreichen: Verstehen von LLMs durch schrittweise Implementierung des Codes eines ChatGPT-ähnlichen LLM, das auf der Transformer-Architektur basiert.

1.1 Was ist ein LLM?

Ein LLM ist ein neuronales Netz, das darauf ausgelegt ist, Klartext zu verstehen, zu erzeugen und darauf zu reagieren. Diese Modelle sind tiefe neuronale Netze (Deep Neural Networks), die mit riesigen Mengen an Textdaten trainiert wurden, die manchmal große Teile des gesamten öffentlich zugänglichen Texts im Internet umfassen.

Das »Large« in »Large Language Models« bezieht sich sowohl auf die Größe des Modells in Bezug auf die Parameter als auch auf den riesigen Datensatz, mit dem es trainiert wurde. Derartige Modelle haben oft Dutzende oder sogar Hunderte von Milliarden an Parametern, d.h. die anpassbaren Gewichte im Netz, die während des Trainings optimiert werden, um das nächste Wort in einer Sequenz vorherzusagen. Die Vorhersage des nächsten Worts ist sinnvoll, weil sie die inhärente sequenzielle Natur der Sprache nutzt, um Modelle für das Verstehen von Kontext, Struktur und Beziehungen im Text zu trainieren. Da es sich um eine sehr

einfache Aufgabe handelt, überrascht es viele Forscher, dass sie dennoch derart leistungsfähige Modelle hervorbringen kann. In späteren Kapiteln werden wir den Ablauf für das Training mit dem nächsten Wort Schritt für Schritt erläutern und implementieren.

LLMs setzen auf eine als *Transformer* bezeichnete Architektur, die es ihnen ermöglicht, Aufmerksamkeit selektiv auf verschiedene Teile der Eingabe zu richten, um Vorhersagen zu erstellen, sodass sie speziell dafür geeignet sind, die Nuancen und Komplexitäten der menschlichen Sprache zu berücksichtigen.

Da LLMs in der Lage sind, Text zu generieren, betrachtet man sie oftmals auch als eine Form der generativen künstlichen Intelligenz, kurz GenAI (für *Generative Artificial Intelligence*). Wie Abbildung 1.1 zeigt, umfasst künstliche Intelligenz die Entwicklung von Maschinen, die Aufgaben ausführen können, für die eine menschliche Intelligenz erforderlich ist – einschließlich Sprache verstehen, Muster erkennen und Entscheidungen treffen –, und Teilbereiche wie Machine Learning oder Deep Learning.

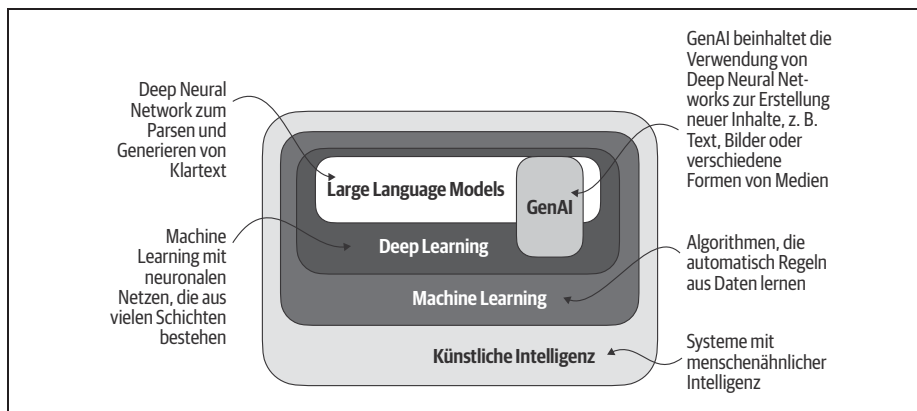


Abb. 1.1 Wie diese hierarchische Darstellung der Beziehungen zwischen den verschiedenen Bereichen zeigt, verkörpern LLMs eine spezifische Anwendung von Deep-Learning-Techniken, indem sie deren Fähigkeit nutzen, menschenähnlichen Text zu verarbeiten und zu erzeugen. Deep Learning ist ein spezialisierter Zweig des Machine Learning, der sich auf mehrschichtige neuronale Netze stützt. Machine Learning und Deep Learning sind Bereiche mit dem Ziel, Algorithmen zu implementieren, die Computer in die Lage versetzen, aus Daten zu lernen und Aufgaben durchzuführen, die normalerweise menschliche Intelligenz erfordern.

Die Algorithmen, die KI implementieren sollen, stehen im Mittelpunkt des Machine Learning. Insbesondere geht es bei Machine Learning um die Entwicklung von Algorithmen, die anhand von Daten lernen und Vorhersagen oder Entscheidungen treffen können, ohne explizit programmiert zu werden. Um dies zu veranschaulichen, stellen Sie sich einen Spam-Filter als praktische Anwendung des

Machine Learning vor. Anstatt Spam-E-Mails mithilfe von manuell formulierten Regeln zu identifizieren, wird ein Algorithmus für Machine Learning mit Beispielen von E-Mails gefüttert, die als Spam- und Nicht-Spam-E-Mails gekennzeichnet sind. Indem man den Fehler des Modells in seinen Vorhersagen auf einem Trainingsdatensatz minimiert, lernt es, Muster und Charakteristika von Spam zu erkennen, sodass es in die Lage versetzt wird, neue E-Mails entweder als Spam oder als Nicht-Spam zu klassifizieren.

Wie Abbildung 1.1 zeigt, bildet Deep Learning einen Teilbereich des Machine Learning, bei dem es darum geht, komplexe Muster und Abstraktionen in den Daten durch neuronale Netze mit drei oder mehr Schichten (auch als Deep Neural Networks bezeichnet) zu modellieren. Im Gegensatz zum Deep Learning ist beim herkömmlichen Machine Learning eine manuelle Merkmalsextraktion erforderlich. Das heißt, dass menschliche Experten die relevantesten Features für das Modell identifizieren und auswählen müssen.

Der Bereich der künstlichen Intelligenz wird heute von Machine Learning und Deep Learning dominiert, umfasst aber auch andere Ansätze – zum Beispiel regelbasierte Systeme, genetische Algorithmen, Expertensysteme, Fuzzy-Logik oder Computeralgebra (symbolische Manipulation algebraischer Ausdrücke).

Kommen wir auf das Beispiel der Spam-Klassifizierung zurück: Beim traditionellen Machine Learning könnten menschliche Experten manuell Merkmale aus dem E-Mail-Text herausziehen, wie zum Beispiel die Häufigkeit bestimmter Trigger-Wörter (etwa »Preis«, »Gewinn«, »kostenlos«), die Anzahl der Ausrufezeichen, die Verwendung von Wörtern in Großbuchstaben oder das Vorhandensein verdächtiger Links. Mit diesem Datensatz, der auf der Grundlage der von menschlichen Experten definierten Merkmale erstellt wurde, wird dann das Modell trainiert. Im Unterschied zum herkömmlichen Machine Learning ist beim Deep Learning kein manuelles Extrahieren erforderlich. Für ein Deep-Learning-Modell müssen also keine menschlichen Experten die relevantesten Merkmale identifizieren und auswählen. (Allerdings müssen sowohl beim herkömmlichen Machine Learning als auch beim Deep Learning für die Spam-Klassifizierung immer noch Labels erfasst werden, zum Beispiel ob es sich um Spam oder Nicht-Spam handelt, die entweder von einem Experten oder von den Usern bestimmt werden.)

Schauen wir uns nun an, für welche Probleme LLMs heute infrage kommen, welche Herausforderungen LLMs angehen können und wie die allgemeine LLM-Architektur aussieht, die wir später implementieren werden.

1.2 Anwendungen von LLMs

Dank ihrer fortgeschrittenen Fähigkeiten, unstrukturierte Textdaten zu analysieren und zu verstehen, sind LLMs in einem breiten Spektrum von Anwendungen in verschiedenen Bereichen zu finden. Heute nutzt man LLMs für die maschinelle Übersetzung, das Generieren von Prosatexten (siehe Abbildung 1.2) sowie Stim-

mungsanalysen, Textzusammenfassungen und viele andere Aufgaben. Seit Kurzem werden LLMs auch für das Erstellen von Inhalten verwendet, um beispielsweise Romane, Artikel oder sogar Computercode zu schreiben.

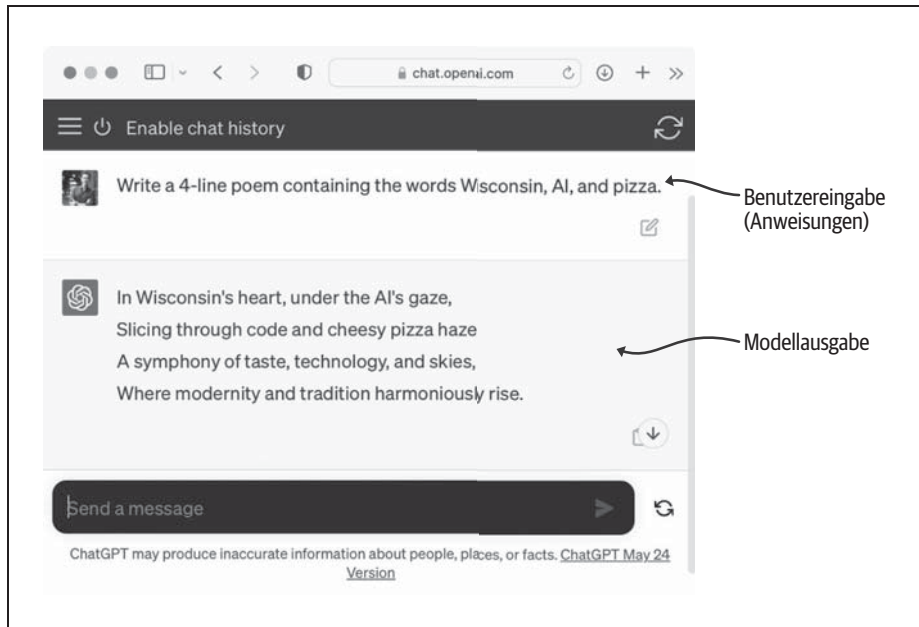


Abb. 1.2 LLM-Schnittstellen ermöglichen die Kommunikation zwischen Benutzern und KI-Systemen in natürlicher Sprache. Dieser Screenshot zeigt, wie ChatGPT ein Gedicht nach den Vorgaben eines Benutzers schreibt.

LLMs können auch anspruchsvolle Chatbots und virtuelle Assistenten antreiben, wie es zum Beispiel bei ChatGPT von OpenAI und Gemini (früher Bard genannt) von Google der Fall ist. Derartige Anwendungen können Benutzeranfragen beantworten und herkömmliche Suchmaschinen wie Google Search oder Microsoft Bing ergänzen.

Darüber hinaus eignen sich LLMs zur effektiven Wissensabfrage aus riesigen Textmengen in Spezialgebieten wie Medizin oder Recht. Dazu gehören das Durchsuchen von Dokumenten, das Zusammenfassen langer Passagen und die Beantwortung technischer Fragen.

Kurz gesagt, LLMs sind von unschätzbarem Wert für die Automatisierung fast aller Aufgaben, die das Parsen und Generieren von Text beinhalten. Die Anwendungsmöglichkeiten sind schier unendlich, und da wir weiterhin Innovationen entwickeln und neue Wege zur Nutzung dieser Modelle erforschen, ist klar, dass LLMs das Potenzial besitzen, unsere Beziehung zur Technologie neu zu definieren, indem sie sie dialogfähiger, intuitiver und zugänglicher machen.

Uns geht es in erster Linie darum, die prinzipielle Funktionsweise von LLMs zu verstehen. Zu diesem Zweck programmieren wir ein LLM, das Texte erzeugen kann. Außerdem lernen Sie Techniken kennen, die es LLMs ermöglichen, Abfragen durchzuführen, die von der Beantwortung von Fragen über die Zusammenfassung von Text bis hin zur Übersetzung von Text in verschiedene Sprachen reichen – und vieles mehr. Sie werden mit anderen Worten lernen, wie komplexe LLM-Assistenten à la ChatGPT funktionieren, indem Sie einen solchen Schritt für Schritt aufbauen.

1.3 Phasen beim Erstellen und Verwenden von LLMs

Warum sollten wir unsere eigenen LLMs erstellen? Ein LLM von Grund auf zu codieren, ist eine ausgezeichnete Übung, um dessen Mechanismen und Grenzen zu verstehen. Außerdem erhalten wir so das nötige Wissen, um vorhandene Open-Source-LLM-Architekturen für unsere domänenspezifischen Datensätze oder Aufgaben vorab zu trainieren oder feinzutunen.

Hinweis

Die meisten LLMs werden heute mithilfe der Deep-Learning-Bibliothek PyTorch implementiert, die wir ebenfalls verwenden. In Anhang A finden Sie eine umfassende Einführung in PyTorch.

Wie die Forschung in Bezug auf die Modellierungsleistung gezeigt hat, können benutzerdefinierte LLMs – solche, die auf spezifische Aufgaben oder Bereiche zugeschnitten sind – allgemeine LLMs – solche, wie sie von ChatGPT bereitgestellt und für ein breites Anwendungsspektrum konzipiert sind – übertreffen. Beispiele hierfür sind BloombergGPT (spezialisiert auf Finanzen) und LLMs, die auf die Beantwortung medizinischer Fragen zugeschnitten sind (siehe Anhang B für weitere Details).

Maßgeschneiderte LLMs bieten mehrere Vorteile, insbesondere im Hinblick auf den Datenschutz. Zum Beispiel können Unternehmen darauf bestehen, keine sensiblen Daten mit Drittanbietern von LLMs wie OpenAI zu teilen, da sie Bedenken hinsichtlich der Vertraulichkeit haben. Darüber hinaus ermöglicht die Entwicklung kleinerer benutzerdefinierter LLMs ein direktes Deployment auf Kundengeräten wie Laptops und Smartphones, was von Unternehmen wie Apple derzeit erforscht wird.

Diese lokale Implementierung kann die Latenzzeit erheblich senken und die serverbezogenen Kosten verringern. Darüber hinaus gewähren benutzerdefinierte LLMs den Entwicklerinnen und Entwicklern völlige Autonomie, sodass sie Aktualisierungen und Änderungen des Modells nach Bedarf steuern können.

Der allgemeine Ablauf beim Erstellen eines LLM umfasst das Vortraining und das Feintuning. Das »Vor« in Vortraining bezieht sich auf die Anfangsphase, in der ein Modell wie ein LLM mit einem großen, breit gefächerten Datensatz trainiert wird, um ein umfassendes Verständnis von Sprache zu entwickeln. Dieses vortrainierte Modell dient dann als grundlegende Ressource, die sich durch ein Feintuning weiterentwickeln lässt. Bei diesem Prozess wird das Modell speziell mit einem engeren Datensatz trainiert, der für bestimmte Aufgaben oder Bereiche spezifischer ist. Abbildung 1.3 veranschaulicht diesen zweistufigen Trainingsansatz, der aus Vortraining und Feintuning besteht.

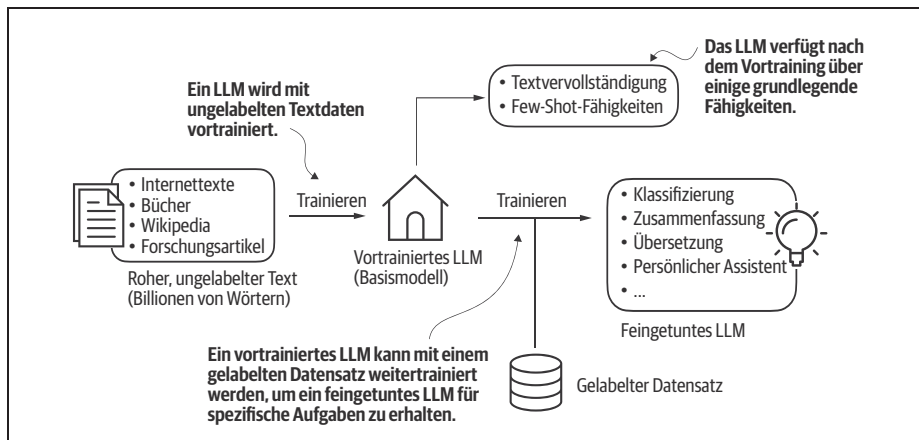


Abb. 1.3 Das Vortraining eines LLM beinhaltet die Vorhersage des nächsten Worts auf großen Textdatensätzen. Ein vortrainiertes LLM kann dann mit einem kleineren gelabelten Datensatz feinetunt werden.

Um ein LLM zu erstellen, trainiert man es im ersten Schritt mit einem großen Korpus von Textdaten, den man auch als *Rohertext* bezeichnet. Hier bezieht sich »roh« auf die Tatsache, dass es sich bei diesen Daten lediglich um normalen Text ohne irgendwelche Beschriftungsinformationen handelt. (Es kann aber ein Filter angewendet werden, um beispielsweise Formatierungszeichen oder Dokumente in unbekanntenen Sprachen zu entfernen.)

Hinweis

Leser, die sich bereits mit Machine Learning auskennen, werden feststellen, dass für herkömmliche Modelle des Machine Learning und Deep Neural Networks, die über die konventionellen überwachten Lernparadigmen trainiert werden, normalerweise Beschriftungsinformationen erforderlich sind. Dies ist jedoch nicht der Fall für die Vortrainingsphase von LLMs. In dieser Phase verwenden LLMs Self-supervised Learning (selbstüberwachtes Lernen), bei dem das Modell seine eigenen Labels aus den Eingabedaten generiert.

In der als *Vortraining* bezeichneten ersten Trainingsphase eines LLM wird ein erstes vortrainiertes LLM erstellt, das oft als *Basis-* oder *Grundmodell* bezeichnet wird. Ein typisches Beispiel für ein derartiges Modell ist das GPT-3-Modell (der Vorläufer des in ChatGPT angebotenen Originalmodells). Dieses Modell ist in der Lage, Text zu vervollständigen – d.h., einen halb geschriebenen Satz, den der Benutzer bereitstellt, zu beenden. Zudem besitzt es beschränkte Few-Shot-Fähigkeiten, kann also neue Aufgaben auf der Grundlage von nur wenigen Beispielen erlernen, anstatt umfangreiche Trainingsdaten zu benötigen.

Ist ein LLM mit großen Textdatensätzen für die Vorhersage des nächsten Worts im Text vortrainiert worden, können wir das LLM mit gelabelten Daten weitertrainieren, was auch als *Feintuning* (Feinabstimmung) bezeichnet wird.

Die beiden populärsten Kategorien des Feintunings von LLMs sind das *Feintuning per Anweisung* und das *Feintuning per Klassifizierung*. Beim Feintuning per Anweisung besteht der gelabelte Datensatz aus Anweisungs-Antwort-Paaren, wie zum Beispiel die Anfrage zur Übersetzung eines Texts, die vom korrekt übersetzten Text begleitet wird. Beim Feintuning per Klassifizierung besteht der gelabelte Datensatz aus Texten und zugeordneten Klassenlabels – zum Beispiel E-Mails, denen die Labels »Spam« und »Nicht-Spam« zugeordnet sind.

Wir werden die Codeimplementierungen für das Vortraining und das Feintuning eines LLM behandeln und uns nach dem Vortraining eines grundlegenden LLM eingehender mit den Besonderheiten des Feintunings – sowohl per Anweisung als auch per Klassifizierung – befassen.

1.4 Einführung in die Transformer-Architektur

Die meisten modernen LLMs stützen sich auf die *Transformer*-Architektur, eine Architektur für Deep Neural Networks, die 2017 im Paper »Attention Is All You Need« (<https://arxiv.org/abs/1706.03762>) vorgestellt wurde. Um LLMs zu verstehen, müssen wir den ursprünglichen Transformer verstehen, der für die maschinelle Übersetzung von englischen Texten ins Deutsche und Französische entwickelt wurde. Abbildung 1.4 stellt eine vereinfachte Version der Transformer-Architektur dar.

Die Transformer-Architektur besteht aus zwei Teilmodulen: einem Encoder und einem Decoder. Das Encoder-Modul verarbeitet den Eingabetext und codiert ihn in eine Folge von numerischen Darstellungen oder Vektoren, die die kontextuelle Information der Eingabe erfassen. Dann übernimmt das Decoder-Modul diese codierten Vektoren und generiert den Ausgabertext. Zum Beispiel würde in einer Übersetzungsaufgabe der Encoder den Text aus der Quellsprache in Vektoren codieren, und der Decoder würde diese Vektoren decodieren, um Text in der Zielsprache zu generieren. Sowohl der Encoder als auch der Decoder bestehen aus vielen Schichten, die durch einen sogenannten Self-Attention-Mechanismus miteinander verbunden sind. Möglicherweise haben Sie viele Fragen dazu, wie

die Eingaben vorverarbeitet und codiert werden. Diese Fragen klären wir in den folgenden Kapiteln anhand einer schrittweisen Implementierung.

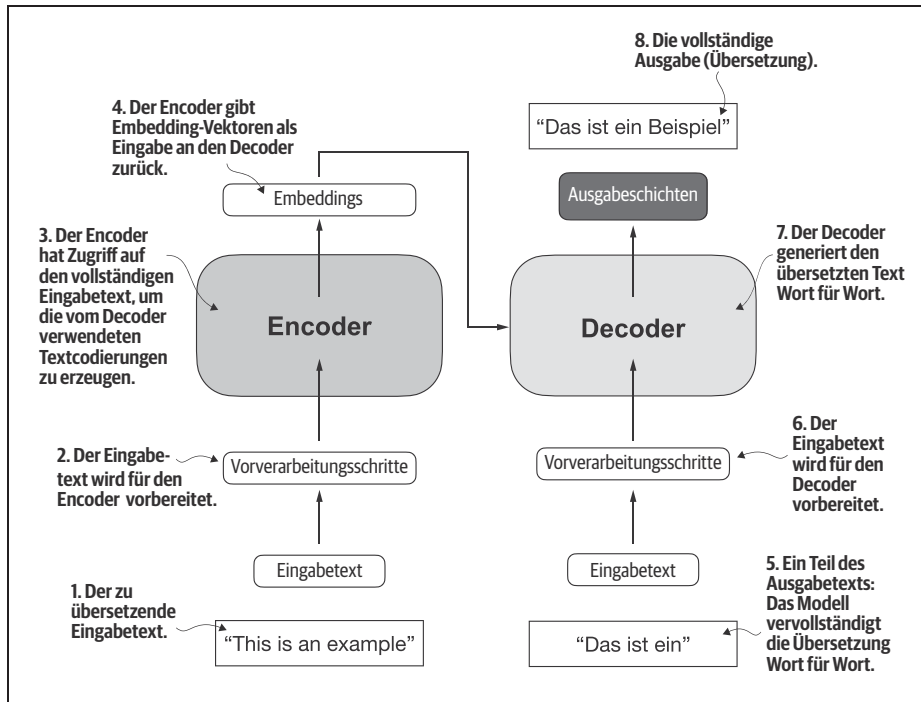


Abb. 1.4 Eine vereinfachte Darstellung der ursprünglichen Transformer-Architektur, die ein Deep-Learning-Modell für die Sprachübersetzung ist. Der Transformer besteht aus zwei Teilen: einem Encoder (links), der den Eingabetext verarbeitet und eine Embedding-Repräsentation des Texts erzeugt (eine numerische Darstellung, die viele verschiedene Faktoren in verschiedenen Dimensionen erfasst), die der Decoder (rechts) verwenden kann, um den übersetzten Text Wort für Wort zu erzeugen. Die Abbildung zeigt die letzte Phase des Übersetzungsprozesses, in der der Decoder für den ursprünglichen Eingabetext (»This is an example«) und einen teilweise übersetzten Satz (»Das ist ein«) nur noch das letzte Wort (»Beispiel«) erzeugen muss, um die Übersetzung abzuschließen.

Eine Schlüsselkomponente von Transformern und LLMs ist der Self-Attention-Mechanismus (hier nicht gezeigt), der es dem Modell ermöglicht, die Bedeutung verschiedener Wörter oder Tokens in einer Sequenz relativ zueinander zu gewichten. Dieser Mechanismus ermöglicht dem Modell, weitreichende Abhängigkeiten und kontextuelle Beziehungen innerhalb der Eingabedaten zu erfassen. Dies erweitert seine Fähigkeiten, kohärent und kontextuell relevante Ausgaben zu erzeugen. Allerdings verschieben wir die weitere Erläuterung aufgrund der Komplexität auf Kapitel 3, wo wir ihn Schritt für Schritt diskutieren und implementieren werden.

Spätere Varianten der Transformer-Architektur, wie BERT (kurz für *Bidirectional Encoder Representations from Transformers*) und die verschiedenen GPT-Modelle (kurz für *Generative Pretrained Transformers*), bauten auf diesem Konzept auf, um diese Architektur für verschiedene Aufgaben anzupassen. In Anhang B finden Sie hierzu weitere Literaturempfehlungen.

BERT, das auf dem Encoder-Submodul des ursprünglichen Transformers aufbaut, unterscheidet sich in seinem Trainingsansatz von GPT. Während GPT für generative Aufgaben entwickelt wurde, sind BERT und seine Varianten auf die Vorhersage maskierter Wörter spezialisiert, wobei das Modell maskierte oder versteckte Wörter in einem gegebenen Satz vorhersagt, wie Abbildung 1.5 zeigt. Diese einzigartige Trainingsstrategie verleiht BERT Stärken bei Textklassifizierungsaufgaben, einschließlich Stimmungsvorhersage und Dokumentkategorisierung. Als Beispiel für die Anwendung seiner Fähigkeiten verwendet X (vormals Twitter) BERT, um schädliche Inhalte zu erkennen.

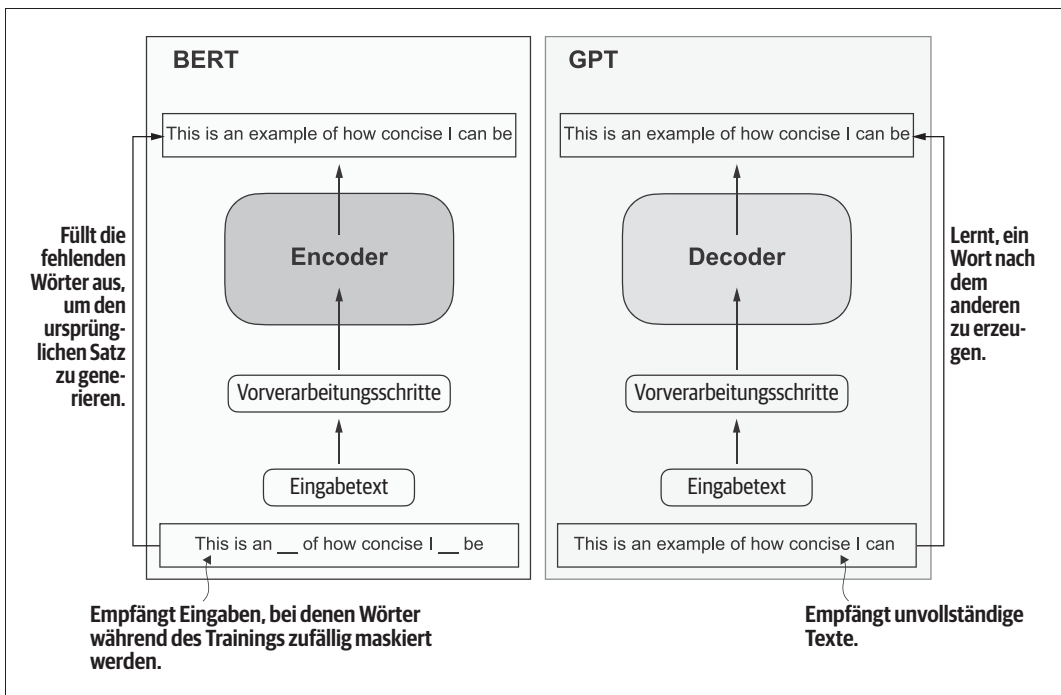


Abb. 1.5 Eine visuelle Darstellung der Encoder- und Decoder-Submodule des Transformers. Die linke Seite zeigt exemplarisch BERT-ähnliche LLMs, die auf die Vorhersage maskierter Wörter ausgelegt sind und hauptsächlich für Aufgaben wie Textklassifizierung verwendet werden. Rechts zeigt das Decoder-Segment GPT-ähnliche LLMs, die für generative Aufgaben und die Erzeugung kohärenter Textsequenzen konzipiert sind.

GPT hingegen konzentriert sich auf den Decoder-Teil der ursprünglichen Transformer-Architektur und ist für Aufgaben konzipiert, bei denen Texte erstellt werden müssen. Dazu gehören maschinelle Übersetzungen, Textzusammenfassungen, das Schreiben von Belletristik, das Schreiben von Computercode und vieles mehr.

GPT-Modelle, die hauptsächlich dafür entwickelt und trainiert wurden, Texte zu vervollständigen, zeigen ebenfalls eine bemerkenswerte Vielseitigkeit in ihren Fähigkeiten. Diese Modelle sind in der Lage, sowohl Zero-Shot- als auch Few-Shot-Learning-Aufgaben auszuführen. Zero-Shot-Learning bezieht sich auf die Fähigkeit, ohne vorherige spezifische Beispiele auf völlig unbekannte Aufgaben zu generalisieren. Andererseits geht es beim Few-Shot-Learning darum, aus einer minimalen Anzahl von Beispielen zu lernen, die der Benutzer als Eingabe zur Verfügung stellt (siehe Abbildung 1.6).

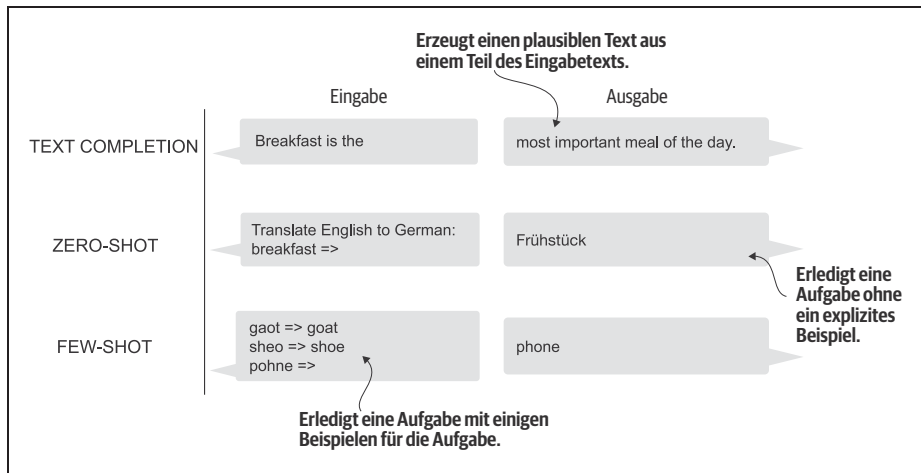


Abb. 1.6 Zusätzlich zur Textvervollständigung können GPT-ähnliche LLMs verschiedene Aufgaben anhand ihrer Eingaben lösen, ohne dass Neutraining, Feintuning oder aufgabenspezifische Änderungen der Modellarchitektur erforderlich sind. Manchmal ist es hilfreich, Beispiele des Ziels innerhalb der Eingabe bereitzustellen, was als »Few-Shot-Setting« bekannt ist. Allerdings sind GPT-ähnliche LLMs auch in der Lage, Aufgaben ohne ein konkretes Beispiel zu realisieren, was man als »Zero-Shot-Setting« bezeichnet.

Transformer- vs. LLM-Architekturen

Die heutigen LLMs basieren auf der Transformer-Architektur. Daher sind Transformer und LLMs Begriffe, die in der Literatur oft synonym verwendet werden. Beachten Sie aber, dass nicht alle Transformer LLMs sind, da Transformer auch für Computervision verwendet werden können. Zudem sind nicht alle LLMs Transformer, da LLMs auf rekurrenten und konvolutionalen Architekturen basieren. Die Hauptmotivation hinter diesen

alternativen Ansätzen ist die Verbesserung der Berechnungseffizienz von LLMs. Es bleibt abzuwarten, ob diese alternativen LLM-Architekturen mit den Fähigkeiten von Transformer-basierten LLMs konkurrieren können und ob sie sich in der Praxis durchsetzen werden. Der Einfachheit halber verwende ich den Begriff »LLM«, um mich auf Transformer-basierte LLMs ähnlich wie GPT zu beziehen. (Interessierte Leser finden in Anhang B Hinweise auf Quellen, die diese Architekturen beschreiben.)

1.5 Große Datensätze nutzen

Die großen Trainingsdatensätze für populäre GPT- und BERT-ähnliche Modelle stellen vielfältige und umfassende Textkorpora dar, die Milliarden von Wörtern umfassen und eine große Bandbreite an Themen sowie natürliche und Computersprachen beinhalten. Um ein konkretes Beispiel zu geben, fasst Tabelle 1.1 den Datensatz zusammen, der für das Vortraining von GPT-3 verwendet wurde, das als Basismodell für die erste Version von ChatGPT diente.

Datensatzname	Datensatzbeschreibung	Anzahl der Tokens	Anteil an den Trainingsdaten
CommonCrawl (gefiltert)	Web-Crawl-Daten	410 Milliarden	60%
WebText2	Web-Crawl-Daten	19 Milliarden	22%
Books1	internetbasierter Buchkorpus	12 Milliarden	8%
Books2	internetbasierter Buchkorpus	55 Milliarden	8%
Wikipedia	hochwertiger Text	3 Milliarden	3%

Tab. 1.1 Der Datensatz für das Vortraining des populären GPT-3-LLM

Tabelle 1.1 gibt die Anzahl der Tokens an, wobei ein Token eine Texteinheit ist, die ein Modell liest. Die Anzahl der Tokens im Datensatz entspricht ungefähr der Anzahl der Wörter und Satzzeichen im Text. Kapitel 2 befasst sich mit der Tokenisierung, d.h. mit der Umwandlung von Text in Tokens.

Die wichtigste Erkenntnis ist, dass Umfang und Vielfalt dieses Trainingsdatensatzes es diesen Modellen ermöglichen, bei verschiedenen Aufgaben, einschließlich Sprachsyntax, Semantik und Kontext, gut abzuschneiden – sogar bei solchen, die Allgemeinwissen erfordern.

Details zum GPT-3-Datensatz

Tabelle 1.1 zeigt den Datensatz, der für GPT-3 verwendet wurde. Die Tabellenspalte mit den Anteilen summiert sich zu 100% der Beispieldaten, wobei Rundungsfehler zu berücksichtigen sind. Obwohl die Teilmengen in der Spalte »Anzahl der Tokens« insge-

samt 499 Milliarden ergeben, wurde das Modell nur mit etwa 300 Milliarden Tokens trainiert. Die Autoren des Papers zu GPT-3 haben nicht angegeben, warum das Modell nicht mit allen 499 Milliarden Tokens trainiert wurde.

Man bedenke die Größe des Datensatzes CommonCrawl, der allein aus 410 Milliarden Tokens besteht und etwa 570 GB Speicherplatz benötigt. Im Vergleich dazu haben spätere Iterationen von Modellen wie GPT-3 – zum Beispiel Llama von Meta – ihren Trainingsumfang erweitert, um zusätzliche Datenquellen wie die arXiv-Forschungspapers (92 GB) und die codebezogenen Fragen und Antworten von StackExchange (78 GB) einzubeziehen.

Die Autoren des GPT-3-Papers haben den Trainingsdatensatz nicht veröffentlicht, aber ein vergleichbarer und öffentlich zugänglicher Datensatz ist »Dolma: An Open Corpus of Three Trillion Tokens for LLM Pretraining Research« von Soldaini et al. 2024 (<https://arxiv.org/abs/2402.00159>). Allerdings kann die Sammlung urheberrechtlich geschützte Werke enthalten, und die genauen Nutzungsbedingungen können vom beabsichtigten Verwendungszweck und vom Land abhängen.

Wegen ihres Vortrainings sind diese Modelle unglaublich vielseitig für ein weiteres Feintuning bei Downstream-Aufgaben. Deshalb bezeichnet man sie auch als Basis- oder Grundmodelle. Das Vortraining von LLMs setzt den Zugang zu erheblichen Ressourcen voraus und ist sehr teuer. So werden beispielsweise die Kosten für das Vortraining von GPT-3 auf 4,6 Millionen Dollar in Form von Cloud Computing Credits geschätzt (<https://mng.bz/VxEW>).

Die gute Nachricht ist, dass sich viele vortrainierte LLMs, die als Open-Source-Modelle verfügbar sind, als Allzweckwerkzeuge eignen, um Texte, die nicht Teil der Trainingsdaten waren, zu schreiben, zu extrahieren und zu bearbeiten. Außerdem lassen sich LLMs für spezifische Aufgaben mit relativ kleinen Datensätzen feintunen, was die erforderlichen Rechenressourcen reduziert und die Performance verbessert.

Wir werden den Code für das Vortraining implementieren und ihn nutzen, um ein LLM für Lehrzwecke vorab zu trainieren. Alle Berechnungen sind auf Consumer-Hardware ausführbar. Nachdem Sie den Code für das Vortraining implementiert haben, lernen Sie, wie sich offen verfügbare Modellgewichte wiederverwenden und in die von uns implementierte Architektur laden lassen, um die teure Vortrainingsphase zu überspringen, wenn wir unser LLM feintunen.

1.6 Die GPT-Architektur unter der Lupe

GPT wurde ursprünglich im Paper »Improving Language Understanding by Generative Pre-Training« (<https://mng.bz/x2qg>) von Radford et al. bei OpenAI vorgestellt. GPT-3 ist eine vergrößerte Version dieses Modells, das mehr Parameter

umfasst und mit einem größeren Datensatz trainiert wurde. Darüber hinaus wurde das ursprüngliche Modell in ChatGPT durch Feintuning von GPT-3 auf einem großen Anweisungsdatensatz mit einer Methode aus dem InstructGPT-Paper von OpenAI (<https://arxiv.org/abs/2203.02155>) erzeugt. Wie Abbildung 1.6 zeigt, sind diese Modelle kompetente Textvervollständigungsmodelle, die auch andere Aufgaben wie Rechtschreibkorrektur, Klassifizierung oder Sprachübersetzung übernehmen können. Dies ist tatsächlich sehr bemerkenswert, wenn man bedenkt, dass die GPT-Modelle mit einer relativ einfachen Aufgabe zur Vorhersage des nächsten Words trainiert werden, wie Abbildung 1.7 zeigt.

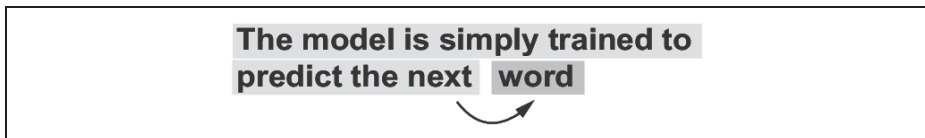


Abb. 1.7 *Beim Vortraining der Vorhersage des nächsten Words für GPT-Modelle lernt das System, das nächste Wort in einem Satz vorherzusagen, indem es sich die Wörter betrachtet, die davor standen. Dieser Ansatz hilft dem Modell, zu verstehen, wie Wörter und Sätze in der Sprache typischerweise zusammenpassen, und bildet eine Grundlage, die auf verschiedene andere Aufgaben angewendet werden kann.*

Die Aufgabe, das nächste Wort vorherzusagen, ist eine Form des selbstüberwachten Lernens (Self-supervised Learning), d.h. eine Form der Selbstbeschriftung. Das bedeutet, dass wir Labels für die Trainingsdaten nicht explizit sammeln müssen, sondern die Struktur der Daten selbst nutzen können: indem wir das nächste Wort in einem Satz oder Dokument als das Label verwenden, das das Modell vorhersagen soll. Da uns diese Aufgabe der Vorhersage des nächsten Words erlaubt, Labels »während des Betriebs« zu erstellen, ist es möglich, LLMs mit riesigen ungelabelten Textdatensätzen zu trainieren.

Verglichen mit der ursprünglichen Transformer-Architektur, die Abschnitt 1.4 behandelt hat, ist die allgemeine GPT-Architektur relativ einfach. Im Wesentlichen handelt es sich nur um den Decoder-Teil ohne den Encoder (siehe Abbildung 1.8). Da Decoder-ähnliche Modelle wie GPT den Text generieren, indem sie ein Wort nach dem anderen vorhersagen, betrachtet man sie als eine Art *autoregressives Modell*. Autoregressive Modelle beziehen ihre vorherigen Ausgaben als Eingaben für zukünftige Vorhersagen ein. Folglich wird bei GPT jedes neue Wort auf der Grundlage der ihm vorausgehenden Sequenz ausgewählt, was die Kohärenz des resultierenden Texts verbessert.

Architekturen wie GPT-3 sind auch erheblich größer als das ursprüngliche Transformer-Modell. So werden im ursprünglichen Transformer die Encoder- und Decoder-Blöcke sechsmal wiederholt. GPT-3 umfasst 96 Transformer-Schichten und insgesamt 175 Milliarden Parameter.

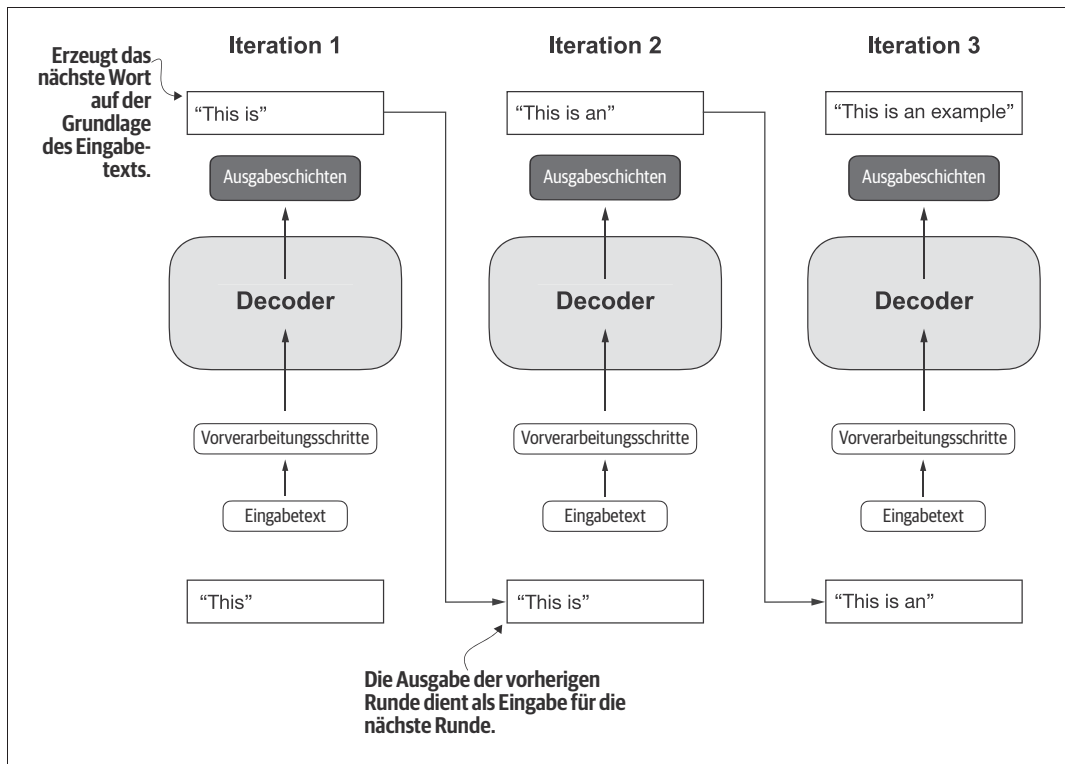


Abb. 1.8 Die GPT-Architektur nutzt nur den Decoder-Teil des ursprünglichen Transformers. Konzipiert ist diese Architektur für unidirektionale Verarbeitung von links nach rechts, sodass sie gut geeignet ist für Aufgaben wie Textgenerierung und die Vorhersage des nächsten Worts, um Text auf iterative Weise Wort für Wort zu erzeugen.

GPT-3 wurde im Jahr 2020 eingeführt, was nach den Maßstäben von Deep Learning und der Entwicklung großer Sprachmodelle als vor langer Zeit zu betrachten ist. Allerdings basieren die neueren Architekturen wie die Llama-Modelle von Meta immer noch auf denselben zugrunde liegenden Konzepten und weisen nur geringfügige Änderungen auf. Daher ist das Verständnis von GPT so wichtig wie eh und je, und ich konzentriere mich auf die Implementierung der prominenten Architektur hinter GPT, während ich Hinweise auf spezifische Anpassungen gebe, die von alternativen LLMs genutzt werden.

Obwohl das ursprüngliche Transformer-Modell, das aus Encoder- und Decoder-Blöcken besteht, ausdrücklich für die Sprachübersetzung entwickelt wurde, sind GPT-Modelle – trotz ihrer zwar größeren, aber auch einfacheren reinen Decoder-Architektur, die auf die Vorhersage des nächsten Worts abzielt – ebenfalls in der Lage, Übersetzungsaufgaben zu erfüllen. Diese Fähigkeit war für die Forschenden zunächst unerwartet, da sie aus einem Modell hervorging, das in erster Linie

für die Vorhersage des nächsten Worts trainiert wurde, also für eine Aufgabe, die nicht speziell auf die Übersetzung abzielte.

Die Fähigkeit, Aufgaben auszuführen, für die das Modell nicht explizit trainiert wurde, bezeichnet man als *emergentes Verhalten*. Diese Fähigkeit wird dem Modell nicht explizit während des Trainings beigebracht, sondern ergibt sich als natürliche Konsequenz aus dem Umgang des Modells mit großen mehrsprachigen Daten in unterschiedlichen Kontexten. Die Tatsache, dass GPT-Modelle die Übersetzungsmuster zwischen Sprachen »erlernen« und Übersetzungsaufgaben ausführen können, obwohl sie nicht speziell dafür trainiert wurden, zeigt die Vorteile und Fähigkeiten dieser großen, generativen Sprachmodelle. Wir können verschiedene Aufgaben erfüllen, ohne für jede Aufgabe unterschiedliche Modelle zu verwenden.

1.7 Ein großes Sprachmodell aufbauen

Nachdem wir nun die Grundlagen für das Verständnis von LLMs geschaffen haben, wollen wir eines von Grund auf programmieren. Wir greifen die fundamentale Idee hinter GPT als Blaupause auf und realisieren diese Aufgabe in drei Phasen, wie sie Abbildung 1.9 veranschaulicht.

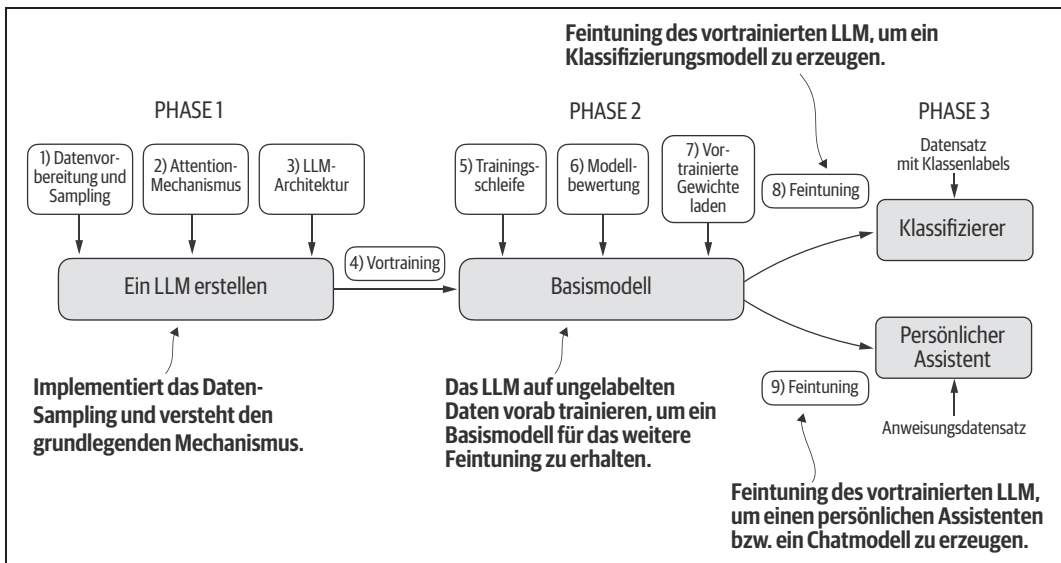


Abb. 1.9 Die drei Hauptphasen für die Programmierung eines LLM. PHASE 1: Implementierung der LLM-Architektur und Datenvorbereitungsprozess, PHASE 2: Vortraining eines LLM, um ein Basismodell zu erstellen, und PHASE 3: Feintuning des Basismodells, um einen persönlichen Assistenten oder Klassifizierer zu werden.

In Phase 1 lernen Sie die grundlegenden Schritte der Datenvorverarbeitung kennen und codieren den Attention-Mechanismus, das Herzstück jedes LLM. Als Nächstes erfahren Sie in Phase 2, wie man ein GPT-ähnliches LLM codiert und trainiert, das in der Lage ist, neue Texte zu generieren. Außerdem werden wir uns mit den Grundlagen der Bewertung von LLMs beschäftigen, die für die Entwicklung leistungsfähiger NLP-Systeme unerlässlich ist.

Das Vortraining eines LLM von Grund auf ist ein bedeutendes Unterfangen, das Tausende oder Millionen von Dollar an Rechenkosten für GPT-ähnliche Modelle verschlingt. Daher liegt der Schwerpunkt bei Phase 2 auf der Implementierung des Trainings für Lehrzwecke anhand eines kleinen Datensatzes. Darüber hinaus bringe ich auch Beispiele für Code, mit dem sich frei verfügbare Modellgewichte laden lassen.

Schließlich nehmen wir in Phase 3 ein vortrainiertes LLM und feintunen es, um Anweisungen zu befolgen, wie zum Beispiel Fragen zu beantworten oder Texte zu klassifizieren – die häufigsten Aufgaben in vielen realen Anwendungen und in der Forschung.

Ich hoffe, Sie freuen sich auf diese spannende Reise!

1.8 Zusammenfassung

- LLMs haben das Gebiet der Verarbeitung natürlicher Sprache umgewandelt, das sich bis dahin vorwiegend auf explizit regelbasierte Systeme und einfachere statistische Methoden gestützt hat. Das Aufkommen von LLMs förderte auch neue auf Deep Learning basierende Ansätze zutage, die zu Fortschritten beim Verstehen, Erzeugen und Übersetzen menschlicher Sprache geführt haben.
- Moderne LLMs werden in zwei Hauptschritten trainiert:
 - Zunächst werden sie auf einem großen Korpus von ungelabeltem Text trainiert, indem die Vorhersage des nächsten Worts in einem Satz als Label verwendet wird.
 - Dann werden sie mit einem kleineren, gelabelten Zieldatensatz feinetunt, um Anweisungen zu befolgen oder Klassifizierungsaufgaben durchzuführen.
- LLMs basieren auf der Transformer-Architektur. Die Schlüsselidee der Transformer-Architektur ist ein Attention-Mechanismus, der dem LLM selektiven Zugriff auf die gesamte Eingabesequenz gibt, wenn die Ausgabe Wort für Wort generiert wird.
- Die ursprüngliche Transformer-Architektur besteht aus einem Encoder, der den Text parst, und einem Decoder, der Text generiert.

- LLMs wie GPT-3 und ChatGPT, die Text generieren und Anweisungen befolgen sollen, implementieren nur Decoder-Module, was die Architektur vereinfacht.
- Große Datensätze, die aus Milliarden von Wörtern bestehen, sind für das Vortraining von LLMs unerlässlich.
- Während die allgemeine Vortrainingsaufgabe für GPT-ähnliche Modelle darin besteht, das nächste Wort in einem Satz vorherzusagen, weisen derartige LLMs emergente Eigenschaften auf, wie zum Beispiel die Fähigkeit, Texte zu klassifizieren, zu übersetzen oder zusammenzufassen.
- Sobald ein LLM vortrainiert ist, kann das resultierende Basismodell für verschiedene Downstream-Aufgaben effizienter feingetunt werden.
- LLMs, die auf benutzerdefinierten Datensätzen feingetunt wurden, können allgemeine LLMs bei spezifischen Aufgaben übertreffen.