

We Are Developers!

Frühjahr 1/2022

Eine Themenbeilage der
Heise Medien GmbH & Co. KG



> MODERNE KI

Ethik und Fairness
neu denken

> KI-SICHERHEIT

Wie Maschinen bessere
Entscheidungen treffen

> CLOUD-SECURITY

Kontext hilft beim
Schwachstellen-Management

> AUTONOMES FAHREN

Verkehrsszenarien mit
Simulationen testen

> MODEL GOVERNANCE

Machine Learning geht mit
Verantwortung einher

> EUROPA KANN KI

Interview: Was KI-Souveränität
für die EU bedeutet

Work #LikeABosch

Darum geht es bei Bosch:
Um dich, unser Business,
unsere Umwelt.

Wir machen aus Visionen Wirklichkeit. Bei Bosch gestalten wir Zukunft mit hochwertigen Technologien und Dienstleistungen, die Begeisterung wecken und das Leben der Menschen verbessern. Unsere Geschäftsfelder sind so vielfältig wie unsere Bosch-Teams auf der ganzen Welt – ihre Kreativität ist der Schlüssel zur Innovation im vernetzten Zuhause, der Mobilität und der Industrie 4.0. Unsere Verantwortung geht weit darüber hinaus. Wir sind börsenunabhängig und dem Zweck der Robert Bosch Stiftung verpflichtet; von unserem Erfolg profitieren unmittelbar Gesellschaft, Umwelt und nachfolgende Generationen. Unser Versprechen an dich steht felsenfest: Wir wachsen gemeinsam, haben Freude an unserer Arbeit und inspirieren uns gegenseitig. Werde ein Teil davon und erlebe mit uns einzigartige Bosch-Momente.

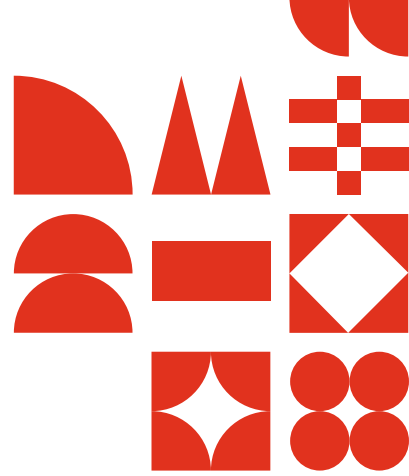
Bewirb dich jetzt: bosch.at/karriere



BOSCH

Technik fürs Leben

EDITORIAL



Gestalten statt Verwalten

Was hat es auf sich mit dem „bisschen Bewusstsein“ bei künstlicher Intelligenz, das zurzeit die Gemüter erregt? Im Februar 2022 befeuerte ein leitender Wissenschaftler von OpenAI eine alte Diskussion neu: „It may be that today’s large neural networks are slightly conscious“, twitterte Ilya Sutskever beiläufig. Prompt widersprach ihm Yann LeCun, der Chef-Machine-Learner des Facebook-Mutterkonzerns Meta. Ist Europa bei den technischen Neuerungen, die Big Tech in den USA vorantreibt, gerade nur Zuschauer? Die Frühjahrsausgabe der WeAreDevelopers-Beilage dreht sich um die technologische Transformation unserer Welt durch moderne KI.

Nicht nur OpenAI, das GPT-3 hervorgebracht hat, strebt nach der Entwicklung generalisierender künstlicher Intelligenz (Artificial General Intelligence, kurz: AGI), auch das Heidelberger Unternehmen Aleph Alpha hat sich AGI auf die Fahnen geschrieben. Sein Gründer Jonas Andrulis stellt die neuen KI-Basismodelle (Foundation Models) vor, die die Komplexität unserer Welt verstehen und deren Fähigkeiten über reines Memorizing hinausgehen in Richtung Reasoning. Er und sein Team bauen eine europäische Alternative zu OpenAI. Welche Chancen die großen multimodalen Modelle bergen und wie potenzielle Gefahren sich eindämmen lassen, erklärt der Ingenieur. In einem Interview habe ich ihn und den Sprecher des DFKI Reinhard Karger um eine Einordnung der aktuellen Entwicklungen gebeten: In „Gestalten statt Verwalten – Europa kann KI“ geben sie Denkansätze zur neuen Art der Mensch-Maschine-Zusammenarbeit.

Autonomes Fahren und die Mobilität der Zukunft sind das Metier von Katharina Hartmann. Zum Testen von Verkehrsszenarien mit Simulationen bietet sie einen Werkstatteinblick aus dem Deutschen Zentrum für Luft- und Raumfahrt. Thordis Taag hat den Data Scientist Daniel Kondermann zur KI-Sicherheit befragt. Gemeinsam klären sie auf, warum Maschinen (noch) keine besseren Entscheidungen treffen und machen einen Ausflug in die Neurowissenschaften. Metriken alleine reichen nicht aus, findet Isabel Bär, die die ethische Dimension des Machine Learning beleuchtet. Model Governance federt unternehmerische Risiken ab und erfüllt die weitreichenden Auflagen zur Modell-dokumentation in der Industrie. Das Frühjahrsheft bietet Orientierung in einer unruhigen Zeit.

Inspiration und
frohes Schaffen wünscht

Silke Hahn



INHALT

- 4 Moderne KI neu denken
- 12 Europa kann KI
- 24 Autonomes Fahren
- 30 KI-Sicherheit
- 38 Model Governance
- 44 Cloud-Security

Young Professionals schreiben für Young Professionals

Unsere Beilage zu c't und iX basiert überwiegend auf einer Online-Artikelserie, die Young Professionals eine Bühne bietet für erste Fachartikel. Die Autoren und Autorinnen erhalten von der Heise-Developer-Redaktion Unterstützung beim Konzipieren und Schreiben.

Dein erster Fachartikel bei Heise

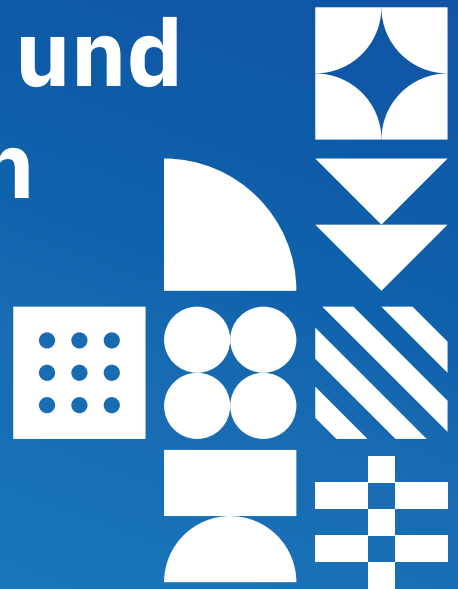
Die Serie ermutigt dazu, sich selbst erstmals als Autor oder Autorin zu betätigen. Sei es, um eigene Erfahrungen mitzuteilen, ein Projekt vorzustellen – oder einfach, weil du schon immer mal einen Fachartikel schreiben wolltest. Hast du eine Idee?

Schreib uns: developer@heise.de

> Moderne KI: Ethik und Fairness neu denken

Jonas Andrulis

KI-Basismodelle verarbeiten die Komplexität unserer Welt jenseits klar definierter Aufgaben – das bringt Chancen, birgt aber auch Risiken und Herausforderungen. Wo droht Gefahr und was können wir tun?



Seit 2008 hat die zweite Welle von KI-Anwendungen überraschende neue Funktionen möglich gemacht und dabei die Erwartungen so hoch geschraubt, dass heute jede Currywurstbude behauptet, eine KI-Strategie zu haben. Die Modelle, die diese technologische Revolution antreiben, stützen sich auf riesige künstliche neuronale Netze (Deep Learning), die Millionen (und später Milliarden) von Parametern mit Hilfe sorgfältig erstellter Trainingsdatensätze anpassen. Doch bis vor kurzem hatten all diese KI-Modelle eine Gemeinsamkeit: Sie beschränkten sich auf die Verarbeitung spezieller Inputs (wie Kamerabilder aus Fahrzeugen) mit im Vorfeld

eng definierten Ergebnissen (wie der Erkennung von Fußgängern). Die Grundlage sind von Menschen definierte Ziele anhand manuell annotierter Trainingsdaten (überwachtes Lernen). Das Verhalten solcher Modelle lässt sich leicht verstehen und messen, ihre Leistung ist gut quantifizierbar – wie beispielsweise in einer Konfusionsmatrix, die angibt, wie viele echte Fußgänger die Künstliche Intelligenz (KI) übersehen hat oder meint, erkannt zu haben, wo es keine gibt.

Beim Erstellen dieser Art von KI ist es möglich, versehentlich ein Modell zu trainieren, das für eine bestimmte Gruppe von Menschen oder in gewissen Anwendungsfällen keine gute Leistung zeigt. Sobald ein solches potenzielles Problem erkannt wurde, ist es damit relativ einfach, auch ethische Bedenken hinsichtlich der Fähigkeiten des Modells zu äußern: Fußgänger (in diesem Beispiel) mit einem seltenen Aussehen (was hier heißt: ungewöhnlich verglichen mit der Mehrheit im Datensatz) könnten nicht mit ausreichender Qualität erkannt worden sein. Falls das der Fall ist, lässt es sich eindeutig evaluieren und darstellen. Wir können versuchen, einen (Teil-)Datensatz von Fußgängern mit ungewöhnlichem Aussehen zu erstellen (basierend auf passenden Kriterien wie beispielsweise ethnischer Zugehörigkeit oder Größe) und messen, wie gut unser Modell bei dieser Untergruppe abschneidet. Liegt das Ergebnis unterhalb unseres Schwellenwerts für die Akzeptanz, können wir weitere Trainingsdaten hinzufügen, die genau diese Beobachtungen enthalten oder die Gewichtung der bereits vorhandenen erhöhen. Dieser An-

in a nutshell

- > Generalisierende Künstliche Intelligenz versteht unsere Welt wie noch keine Technologie zuvor.
- > Die Entwicklung geht in Richtung transformativer KI (Transformative Artificial Intelligence, kurz: TAI). Auf dem Weg stellen sich Fragen zu Ethik und Fairness neu.
- > Wie sind die Systeme zu verstehen – und welche Anforderungen an Entwicklerinnen und Entwickler sowie Nutzerinnen und Nutzer ermöglichen eine sichere Zukunft?


IN AGILEN WORKSTREAMS

DIE CLOUD-LÖSUNGEN

VON MORGEN ENTWICKELN.

DARUM SIND WIR BEI DATEV.

Gemeinsam sichere Cloud-Lösungen und innovative Apps realisieren: Als Cloud-Entwicklerin oder -Entwickler erwarten dich bei DATEV vielfältige Aufgaben in einer agilen Innovations-Kultur. Informiere dich über freie Stellen und spannende Projekte bei einem der führenden IT-Dienstleister in Europa.



Valeria und Dominik,
Cloud-Entwicklerin und
-Entwickler bei DATEV

[DATEV.DE/KARRIERE](https://datev.de/karriere)



Zukunft gestalten.
Gemeinsam.

satz wird die Modelle nicht perfekt machen, und der Versuch, jedes unerwünschte Verhalten zu entfernen, wird uns in eine Endlosschleife von Änderungen führen – aber zumindest sind die Werkzeuge gut etabliert.

Prominente Beispiele, bei denen diese Art von Systemen Diskussionen ausgelöst haben, sind die falsche Klassifizierung von Menschen als Affen oder die schlechte Leistung der Gesichtserkennung für dunkelhäutige Benutzer (weiterführende Informationen hierzu finden sich unter ix.de/z6ja). In diesen Beispielen wurden sehr klare und eindeutig definierte Ziele für das KI-System (Erkennung von Objekten in Bildern, Erkennung menschlicher Gesichter) zweifellos für eine Teilmenge der Daten nicht erreicht. Obwohl kein KI-System immer alles richtig machen wird, ist es sinnvoll, ein System anzustreben, das diese Art von Fehlern vermeidet. Dass diese Fälle von den betroffenen Unternehmen zeitnah korrigiert

wurden, demonstriert die Möglichkeit der Nachbesserung. Bei ihr sollte man sich darüber im Klaren sein, dass der Versuch, diese Probleme zu beheben, ein iterativer Prozess ist und die Systeme dabei wiederholte Anpassungen zulassen müssen. Inzwischen gibt es bewährte Tools und Methoden, um solche unerwünschten Effekte zu reduzieren, während die Grundfunktionalität in ähnlicher Güte erhalten bleibt.

Eine neue Generation von Weltmodellen

Seit einigen Monaten existiert in Forschung und Industrie eine neue Generation von KI-Modellen, für die sich ähnliche Bewertungen nicht mehr auf triviale Weise durchführen lassen. Diese Modelle versuchen nicht, eine bestimmte überwachte (vordefinierte) Zuordnung zu erlernen (Kamerabild auf Fußgänger), sondern allgemeine „Bedeutung“ und Muster

in den Daten zu verstehen. Da keine menschliche Annotation mehr nötig ist, wird dieses Vorgehen selbstüberwachtes Lernen (self-supervised) genannt. Es kann und wird – bei einem ausreichend großen Modell und Datensatz – komplexe Strukturen und Abhängigkeiten finden. Auf der Grundlage dieser Modelle und ihrer erlernten Weltstruktur lässt sich eine Vielzahl verschiedener Anwendungsfälle implementieren, die das während des Trainings erworbene Wissen nutzen. Aus diesem Grund hat eine Gruppe von Forschern der Universität Stanford ihnen den Namen „Foundation Models“ (Basismodelle) gegeben – denn mit ihrer Hilfe können wir unzählige neue mögliche Anwendungen entwickeln, die auf dem Weltwissen der Modelle aufbauen. Viele der Fähigkeiten, die sich aus diesem Aufbau ergeben, sind überraschend und waren während des Trainings nicht absehbar oder geplant. Die Ergebnisse sind oft neuartig und können beeindruckend komplex sein: Das von OpenAI entwickelte generative Bildmodell DALL-E beispielsweise kann auf der Grundlage einer kurzen Textbeschreibung Bilder des gewünschten Inhalts erstellen, und das mit fast unbegrenzter Flexibilität (s. Abb. 1).

Die Funktionalitäten, die sich auf diesen Basismodellen aufbauen lassen, sind nahezu unbegrenzt – dabei nutzen sie komplexe und teilweise unbekannte interne Muster und Zusammenhänge in den Daten (und der Welt). Dadurch gibt es keine einfachen Maßnahmen, um ethische Bedenken ähnlich wie bei den Modellen des überwachten Lernens zu berücksichtigen. Wenn DALL-E einen „Sessel in Form einer Avocado“ erschafft, wie können wir dann beurteilen, ob diese Darstellung fair ist? Werden

Aleph Alpha: Spitzentechnologie für Europa

Die 2019 von Jonas Andrusis und Samuel Weinbach in Heidelberg gegründete Aleph Alpha GmbH gilt als Leuchtturm in Deutschland und Europa, denn sie betreibt als einziges europäisches KI-Unternehmen die Forschung, Entwicklung und Gestaltung in Richtung generalisierender Künstlicher Intelligenz (Artificial General Intelligence, kurz: AGI). Die Heidelberger steuern beim Gaia-X-Projekt OpenGPT-X multimodale Modelle der neuesten Generation bei. 2021 stellten sie mit 28,3 Millionen Euro an Wagniskapital den deutschen Deep-Tech-Finanzierungsrekord auf.

Selbstbewusst strebt das Unternehmen laut eigenen Angaben danach, technische Kompetenz und Wertschöpfung hier zu bündeln, damit Europa im globalen Wettbewerb handlungsfähig bleibt. „Wenn die gesamte Wertschöpfung zu den Aktionären von Microsoft und OpenAI geht, dann fehlt uns das hier als Gesellschaft“, untermauert CEO Andrusis das Anliegen.

Der Machine Learner und Serienunternehmer war zuvor leitend in der KI-Forschung von Apple aktiv. 2021 gewann er den Deutschen KI-Preis für eine Technologie vergleichbar mit Deep Mind und OpenAI. Vertiefung bietet der Interviewartikel auf S. 12: „Gestalten statt Verwalten – Europa kann KI“. Darin diskutieren Jonas Andrusis und Reinhard Karger (DFKI) komplementär über Künstliche Intelligenz und die technologische Transformation.



Quelle: Julian Beekmann

>> **Jonas Andrusis, Founder und CEO Aleph Alpha**



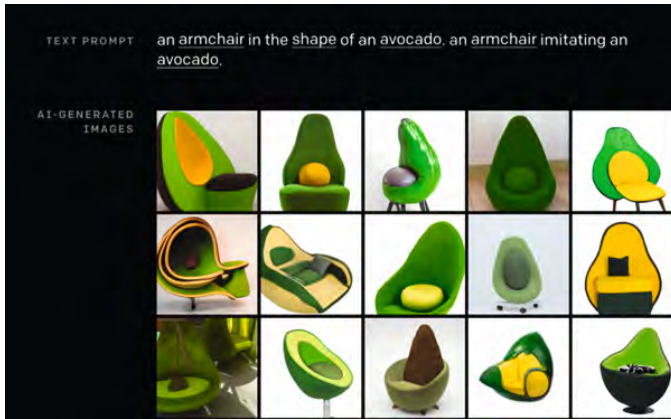
This ad will not benefit from professional and personal development. **You will.**

Our expectations of an employee are high, but truly fair. In return, we offer you the chance to be part of one of the biggest IT projects in retail and a lot of benefits on top. Join our team and enjoy various opportunities for your professional and personal development.

Apply now:

it-jobs.aldi-sued.de

Regardless of the texts and images in our recruiting materials, we would like to emphasise that ALDI SOUTH provides equal opportunities to all people.



>> DALL-E von OpenAI erzeugt auf Basis einer Beschreibung des gewünschten Bildinhaltes (caption) Bilder (Abb. 1).

alle Arten von Sesseln oder Umgebungen so dargestellt, wie wir es für erforderlich halten? Sind Menschen inbegriffen? Was sind unsere Anforderungen an die Darstellung von Menschen? Während diese Fragen für das DALL-E-Beispiel harmlos erscheinen, werden sie im Zusammenhang mit GPT-3 und ähnlichen Modellen äußerst relevant. So sind GPT-3-ähnliche Modelle beispielsweise in der Lage, eine Zusammenfassung eines langen Textes zu schreiben, ohne dass sie Beispiele oder zusätzliche Trainingsdaten sehen müssen. Aber wird diese Zusammenfassung alle Aspekte des Dokuments angemessen wiedergeben? Würde ein Psychiater nicht eine andere Zusammenfassung schreiben als ein Ingenieur, oder ein alter Deutscher sich nicht auf andere Aspekte des Textes konzentrieren als ein junger japanischer Schriftsteller? Ist diese Subjektivität, die wir bei menschlichen Autoren als selbstverständlich akzeptieren, auch für KI-Modelle in Ordnung? Wie lässt sich das messen und fair vergleichen? Wir müssen hier zwei Fragen untersuchen:

- Was machen die Basismodelle eigentlich? Wie lassen sich ihre Ergebnisse verstehen?
- Was sind aus einer ethischen Perspektive vernünftige Anforderungen an diese Modelle?

Wie kann man die Ergebnisse der Weltmodelle verstehen?

Ein Basismodell funktioniert, indem es Zusammenhänge und Muster erkennt, sie mit Beobachtungen abgleicht und neue Daten im Kontext erzeugt. Diese Modelle tun nicht viel mehr, als wiederholt die Frage zu beantworten: „Was wäre die wahrscheinlichste nächste Beobachtung angesichts des aktuell verfügbaren Kontexts?“. DALL-E stellt also die Frage: „Was sind

die statistisch wahrscheinlichsten Bilder, die auf diese Bildunterschrift folgen?“ Da das Modell während des Trainings große Mengen von Bildern und Bildunterschriften gesehen hat, hat es Objekte, Konzepte und Strukturen in unserer Welt erkannt und kann dieses Wissen übertragen, um passende Ergebnisse in bisher unbekanntem Kontexten vorzuschlagen. Dasselbe gilt für GPT-3: Die Fähigkeit, Texte sinnvoll zu vervollständigen, versetzt es in die Lage, fast jede Aufgabe zu lösen, die sich in sprachlicher Form ausdrücken lässt. Eine Zusammenfassung für einen komplexen juristischen Text schreiben? Fragen zum Inhalt eines technischen Berichts beantworten? Kreative Namen für einen Friseursalon in Berlin ausdenken? Die einzige Beschränkung dieser Basismodelle scheint der notwendige ausreichend große Speicherplatz (gemeint ist die Größe des künstlichen neuronalen Netzes in Anzahl und Verbindungen der Neuronen). Mit ausreichender Modellgröße und genügend Trainingsdaten lernt die KI, dass die wahrscheinlichste Art, einen Text fortzusetzen (der aus einem komplexen Artikel gefolgt von einer detaillierten Frage besteht), die richtige Antwort ist (s. Abb. 2).

Im Gegensatz zu klassischen Suchmaschinen, die nur vorhandene Informationen finden und abrufen, reagieren Basismodelle kreativ auf neue Situationen und kombinieren in einem anderen Kontext gelernte Konzepte (s. Abb. 3).

Vor diesem Hintergrund lassen sich folgende Aussagen zur ethischen Bewertung von Modellergebnissen treffen:

Basismodelle haben keine Handlungsfähigkeit, sie haben keine Meinung und verfolgen keine Ziele.

Ein Prompt mit den Worten „Alle Männer sind“, den das Modell mit „faul“ vervollständigt, lässt nicht den Schluss zu, dass es eine bestimmte Meinung über Männer hat. Wir können daher nicht erwarten, dass unsere Ergebnisse in anderen Kontexten konstant männerfeindlich wären. Die Ergebnisse von



>> Ein Basismodell (hier: Aleph Alphas luminous, ähnlich GPT-3) beantwortet durch Vervollständigung eines Kontexts (Prompt) eine komplexe inhaltliche Frage (Abb. 2).

Basismodellen sollten ähnlich dem Ergebnis einer Suche durch die Struktur der Information verstanden werden. Bei dem oben gewählten Suchbegriff würden wir in einer klassischen Internetsuche aus diesem Ergebnis auch nicht ableiten, dass die Suchmaschine generell voreingenommen ist.

Basismodelle sind weder Wahrheitsmaschinen noch deontologische Instanzen.

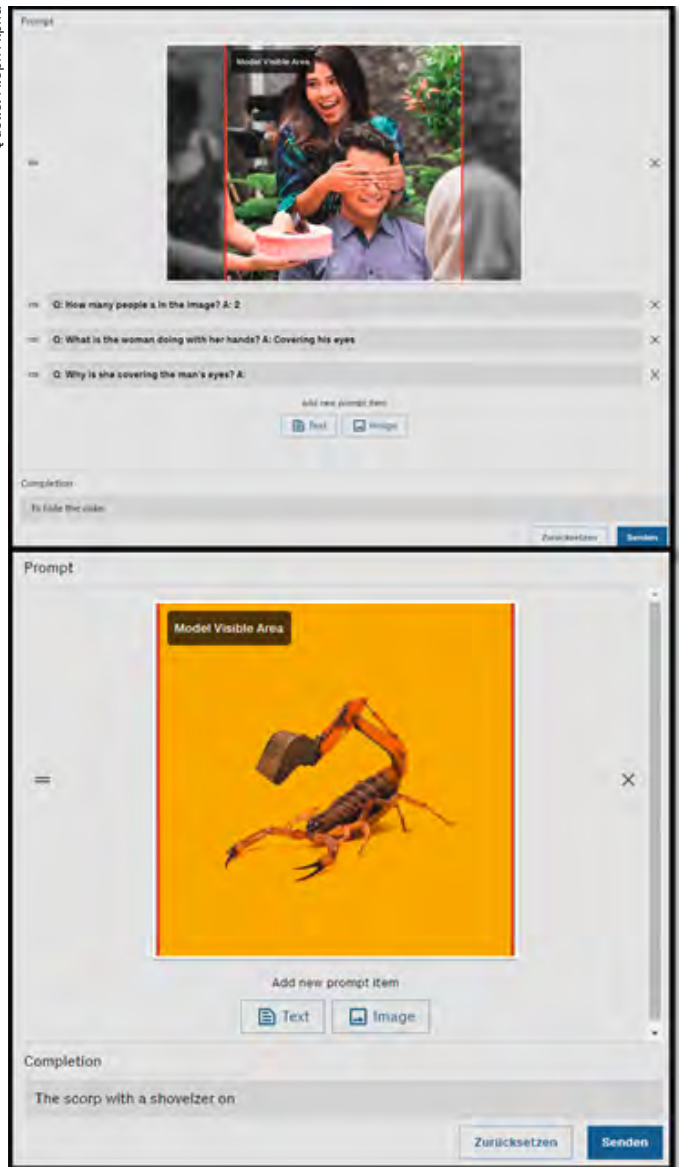
Basismodelle sind für Wahrscheinlichkeiten und Kontext gebaut. Das bedeutet, dass sie sich auf menschliche Muster, Logik, Wissen und Verständnis stützen. Obwohl es zahlreiche Belege dafür gibt, dass dieser Ansatz zu qualitativ hochwertigen Ergebnissen führen kann, wird eine Häufung kontrafaktischer Beobachtungen das Modell dazu verleiten, Unwahrheiten zu reproduzieren. Basismodelle können Informationen zwar teils logisch kombinieren, aber ihre Fähigkeit ist nicht unbegrenzt robust. Wir können daher nicht erwarten, dass sie gegen eine große Menge an falschen Eingaben resistent sind oder in komplexen Situationen immer die Wahrheit wiedergeben.

Die Ausgabe ist auf komplexe und oft unsichtbare Weise von der Eingabe abhängig.

Durch den Prompt beeinflussen wir das Ergebnis stark. Auf einen Prompt, der eine Verallgemeinerung formuliert, folgt mit hoher Wahrscheinlichkeit eine voreingenommene Aussage. Da Basismodelle zu einer beeindruckenden Genauigkeit fähig sind, sind diese Einflüsse manchmal nicht offensichtlich. Bei der Anwendung von Basismodellen auf Software-Quellcode wurde festgestellt, dass sie auch bemerkenswert komplexen Code erzeugen können. Enthielt der Prompt jedoch Code von schlechter Qualität (mit Fehlern oder Mustern, die ein Experte nicht verwenden würde), erwiesen sich auch die Vervollständigungen als wenig brauchbar. Das Basismodell hatte gelernt, dass auf schlechten Code in der Regel noch mehr schlechter Code folgt.

Sprachbasierte Basismodelle beobachten die Welt nicht direkt, sondern eine (voreingenommene) menschliche Interpretation der Welt.

Natürliche Sprache hat sich entwickelt und passt sich laufend an, um sinnvolle Konzepte für den Menschen zu kodieren. Die Verwendung von Sprache für KI hat dadurch enorme Vorteile – Sprache ist bereits optimiert, um zu beschreiben, zu verstehen, zu planen und zu kommunizieren. Logik und Abstraktion sind inhärent Teil der Sprache, und es gibt zahlreiche Texte, die geschrieben wurden, um Wissen weiterzugeben. Diese Sprachfähigkeiten kann die Künstliche Intelli-



Aleph Alphas luminous nutzt und kombiniert (kreativ) gelernte Konzepte unserer Welt. Die Information über den Grund für die zugehaltenen Augen ist nirgends ausgeführt und wurde durch das Modell selbständig „erkannt“. Das Wort „shovelzet“ scheint eine Wortneuschöpfung zu sein (Abb. 3).

genz „leihen“ und erlernen. Im Gegensatz zu menschlichen Betrachtern ist der KI dabei keine direkte Beobachtung der Welt möglich – sie bleibt auf vom Menschen geschaffene sprachliche Interpretation beschränkt.

Ethische Implikationen und Anforderungen

Angesichts der flexiblen Einsetzbarkeit, die Basismodelle für alle Arten von – manchmal überraschenden – Anwendungen fallen und Integrationen bieten, muss eine ethische Bewertung nicht nur die Modellergebnisse berücksichtigen, son-

EINE NEUE GENERATION KI

dern vor allem auch die Art und Weise, wie diese Ergebnisse genutzt werden.

Ausgehend von der Bewertung der allgemeinen Funktionalität hin zum Sicherstellen einer bestmöglichen ethischen Anwendung sind die folgenden Anforderungen (R1-4) sinnvoll:



R1: Korrekte Ergebnisse erzeugen

Technische Korrektheit bedeutet für Basismodelle, dass bei jedem Prompt die vom Modell erzeugte Verteilung über mögliche Vervollständigungen inhaltlich mit den Informationen und Beobachtungen aus den Trainingsdaten konsistent ist. Da Basismodelle komplizierte Aspekte der Kultur und der Werte erfassen, können die Auswahl der Trainingsdaten und der Trainingsprozess (die Tokenisierung, das Curriculum und andere technische Aspekte) einen erheblichen Einfluss auf das Modell haben, der über die in der Zielfunktion gemessene Performance hinausgeht. Ein Beispiel hierfür ist die Frage an Aleph Alphas Basismodell – trainiert in fünf Sprachen – über Sport in jeder davon. Trotz identischer semantischer Inhalte verändert allein die Sprache den Kontext so, dass wir alle möglichen unterschiedlichen europäischen Sportmannschaften erhalten. Fragt man dagegen OpenAIs GPT-3 auf Deutsch, erhält man statt einer deutschen Mannschaft die New York Giants (s. Abb. 4).

Keines dieser Beispiele ist „falsch“ und eine solche Antwort hätte auch keine schädlichen Auswirkungen. Die Beispiele zeigen jedoch, wie Basismodelle implizit viele Aspekte unseres Lebens erfassen, die nicht sofort offensichtlich sind. Für uns bei Aleph Alpha ist das ein wichtiger Grund, nicht einfach

einen Übersetzer hinter einem rein englischsprachigen Modell laufen zu lassen (was einfach zu bauen wäre), sondern das Basismodell in mehreren europäischen Sprachen zu trainieren.

R2: Transparenz für die Nutzer schaffen

Die Modelle verfügen über ein zunehmend komplexes Verständnis unserer Welt, was ihr Verhalten durch eine Vielzahl nicht offensichtlicher Faktoren beeinflusst. Das „richtige“ Ergebnis und Verhalten dieser Art von KI lässt sich nicht ohne Weiteres definieren oder sicherstellen. Für die Ingenieure des maschinellen Lernens, die diese Modelle entwickeln, ergibt sich daraus die Verantwortung, ihren Partnern und Kunden – und der Forschungsgemeinschaft für maschinelles Lernen insgesamt – so viel Transparenz wie möglich zu bieten. Dazu gehört neben dem Austausch untereinander auch der aktive Austausch mit Experten und Interessengruppen aus anderen Bereichen.

R3: Werkzeuge zur Kontrolle und zum Verständnis bereitstellen

Das Anwenden ethischer Kategorien für die neue Generation von Basismodellen muss der Komplexität des menschlichen Lebens Rechnung tragen. Unerwünschtes Verhalten zu beschreiben, ist viel schwieriger, als eine Liste „böser Worte“ zu erstellen, die es zu vermeiden gilt. Die Stärke der KI (und von Intelligenz bei Lebewesen allgemein) besteht darin, Strukturen und Muster zu finden und auf neue Probleme anzuwenden. Unstrittig gibt es einige Inhalte, die wir nicht von der KI wiederholt wissen wollen. Da aber fast alles mit allem zusammenhängt, ist es nicht einfach, zu entscheiden, wo das Vorurteil aufhört und die Struktur beginnt. „Mehr als 90 Prozent aller Gefängnisinsassen (in Deutschland) sind männlich.“ Ist das eine korrekte empirische Tatsache, die die KI nutzen kann, oder schädlicher Sexismus? Ingenieure für maschinelles Lernen sollten ihren Nutzerinnen und Nutzer nicht ihre (mehr oder weniger ideologischen) Antworten über den gewünschten Zustand der Welt aufzwingen, sondern mit Werkzeugen und Funktionen helfen, damit eine freie Gesellschaft gemeinsam an diesen Fragen arbeiten kann.

R4: Human-in-the-Loop: Integration des Menschen in transformative Systeme

Für einige Anwendungsfälle erreichen Basismodelle bereits die menschliche Leistungsfähigkeit, was sie zu einem phänomenalen transformativen Werkzeug machen kann. Allerdings fehlt ihnen das Gespür für moralische Gefahren und

Quelle: Aleph Alpha



>> Aleph Alphas luminous und OpenAIs GPT-3 antworten auf die gleiche Frage: Die unterschiedlichen Antworten zeigen beispielhaft, welchen Unterschied ein europäischer Trainingsdatensatz macht. (Abb. 4)

sie können nicht erkennen, wann Menschen aus ethischen Überlegungen eingreifen sollten (einige Menschen scheinen in dieser Hinsicht jedoch auch bemerkenswert dysfunktional zu sein). Das erfordert neue Ansätze für die Benutzerinteraktion und die Zusammenarbeit zwischen Mensch und Maschine, die so zu optimieren sind, dass sie die neuen KI-Fähigkeiten nutzen. Dabei gilt es gleichzeitig, den Menschen zu befähigen, die Ergebnisse der KI effektiv zu verstehen, sie entsprechend seiner Erfahrung zu steuern und ihre Fehler zu erkennen sowie zu korrigieren. Das daraus resultierende Human-in-the-Loop-System kann helfen, die menschliche Aufmerksamkeit auf wertvollere Aspekte zu lenken, bei denen Kreativität, Einfühlungsvermögen oder Intellekt gefragt sind.

Weltwissen aufbauen und einsetzen

Eine neue Generation von Basismodellen bietet phänomenale Möglichkeiten, radikal neue Arbeitsabläufe bei der Informationsverarbeitung zu schaffen. Damit entstehen neue Verantwortlichkeiten und Herausforderungen für alle Beteiligten. Das in den letzten Jahren etablierte Verständnis von KI-Ethik (bei dem der Output eines Modells sich ein-

deutig bewerten ließ) greift hier etwas zu kurz. Da das Innovationstempo in diesem Bereich stetig zunimmt, müssen wir uns auf eine noch leistungsfähigere Künstliche Intelligenz in der nahen Zukunft einstellen und Methoden, Prozesse sowie Innovationen entwickeln, die das Risiko ethischer Katastrophen mindern und gleichzeitig die nächste industrielle Revolution nutzbar machen.

Quellen

Weiterführendes Material steht unter ix.de/z6ja bereit.

Jonas Andrulis



Wirtschaftsingenieur aus Karlsruhe (KIT) mit Forschungshintergrund in maschinellem Lernen und Strategie-modellierung. Serienunternehmer für KI-Software und zuvor in der Leitung der KI-Forschung bei Apples Special Projects Division (SPG).



>eurodata

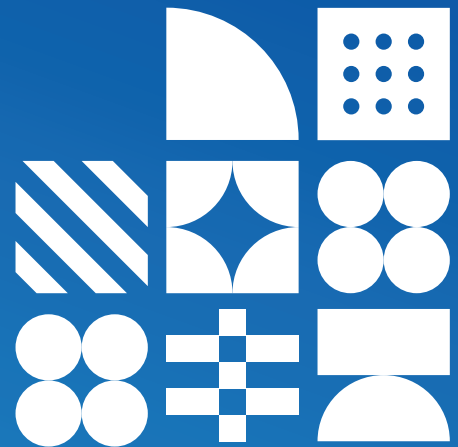


www.eurodata.de/karriere

> Gestalten statt Verwalten – Europa kann KI

Silke Hahn

Ein Machine Learner und ein Computerlinguist diskutieren über Fußball, transformative KI und die technologischen Herausforderungen für die EU. Wo stehen wir – und was kommt auf uns zu?



Unsere Welt steht vor einem voraussichtlich disruptiven Übergang durch Machine Learning; Künstliche Intelligenz (KI) könnte künftig so allgegenwärtig sein wie Elektrizität. Allerdings ist gerade in Deutschland die Skepsis gegenüber der neuen Technologie groß. Offenbar gibt es hierzulande mehr Sommeliers als Mathematikliebhaber. Doch was genau steht uns bevor, und wie gehen wir als Gesellschaft damit um? Können wir etwa bald reflektierte Gespräche mit Maschinen führen?

Heise Developer hat zwei Experten zum Dialog eingeladen:

- **Jonas Andrulis** ist Serienunternehmer, der früher in leitender Funktion für Apple an KI geforscht hat. Er leitet das von ihm gegründete Heidelberger Unternehmen Aleph Alpha.

in a nutshell

- > Neue, generalisierende Künstliche Intelligenz (AGI) schafft die Basis für eine industrielle Revolution.
- > Wir erleben bereits eine neue Art der Mensch-Maschine-Zusammenarbeit, die sich weiter intensiviert.
- > Europa muss die eigene Zukunft souverän gestalten – auf der Grundlage eigener Werte und Technologien.

Das Start-up baut OpenAI für Europa und hat 2021 mit rund 30 Millionen Euro so viel Wagniskapital erhalten wie kein anderes deutsches Deep-Tech-Unternehmen.

- **Reinhard Karger** ist Computerlinguist und Unternehmenssprecher des Deutschen Forschungszentrums für Künstliche Intelligenz (DFKI). Quellen und weiterführende Hinweise sind unter ix.de/z7h4 zusammengestellt.

Beide vertreten oft ähnliche, aber auch unterschiedliche Perspektiven auf Big Tech, den Standort Deutschland, Regulierung und die Gestaltung der Zukunft. Die Aufzeichnung dieses Austauschs soll Neugier wecken und Denkanstöße liefern.

KI überwindet Medienbruch zwischen Welt und Wissen

Heise: Reinhard, zum Auftakt eine Frage an dich: Du hast im Umfeld der CeBIT 2018 erklärt, was KI für den Journalismus bedeutet. Strukturierte Informationen in Tabellen kann Künstliche Intelligenz offenbar schon gut in Text umsetzen. Wo stehen wir da heute?

Reinhard Karger (DFKI): Das eine ist, dass tatsächlich Journalisten entlastet werden von Fließbandarbeit, die keiner machen sollte. Es wird aber auch zusätzlicher Content produziert, den es ohne KI nicht gäbe. Sehr gut funktioniert das bei Sport, Wetter, Börse. Beim Sport geht es darum, dass Systeme automatisch Spielberichte erstellen können aus Daten, die sowieso erfasst werden. Das wiederum ist nicht entschei-

dend für die Bundesliga, sondern für die Jugendarbeit. Die A-Jugend, die B-Jugend und so weiter haben wöchentlich Spiele und Vereine haben einen Blog, in dem nicht nur eine Tabelle, sondern auch Spielberichte stehen sollen. Das können die ehrenamtlichen Trainer nicht leisten. Wenn man dort KI-unterstützt Posts veröffentlicht, die vielleicht von 36 Personen gelesen werden, nämlich von den Verwandten und den Kindern selbst, dann motiviert das die Aktiven, unterstützt die gesellschaftliche Entwicklung, fördert die Jugend und den Sport.

Wenn KI den Medienbruch überwinden kann zwischen Welt und Text, ist das ein gesellschaftlicher Mehrwert. Mir ist das wichtig, weil man ernstlich so tut, als würden der Kommentar, die Analyse oder die Glosse bald von Maschinen geleistet werden können. Aber das stimmt natürlich nicht.

Heise: Wissen wir das schon?

Jonas Andrulis (Aleph Alpha): Absolut korrekt! Es ergeben sich neue Möglichkeiten und die Menschen können ihre Zeit, Inspiration und Kreativität ganz neu einsetzen.

Karger: Werkzeuge sind die Erfindung Homo Fabers und spezifische Werkzeuge taugen nur für spezifische Aufgaben. Die Gabel ist ein fantastisches Werkzeug – nur, wenn man Suppe essen möchte, ist sie nicht optimal, das Messer auch nicht. Es gibt traumhafte Werkzeuge für Spracherkennung oder für maschinelle Übersetzung. Erstaunlicherweise gibt es sie, aber für andere Anwendungen eben noch nicht. Wir können nicht sagen, ob dieses „noch“ eine Aussage ist für zehn oder hundert Jahre, wenn man sich mit den menschlichen Erkenntnisfähigkeiten beschäftigt. Ich bin schon lange dabei: Was einige Werkzeuge heute leisten können, stimmt mich optimistisch, dass da noch viel mehr kommt.

Orientierung in der Welt durch multimodale KI

Heise: Multimodalität in der KI ist ja im Kommen. Worum geht es dabei?

Andrulis: KI, die einen Kontext aus der Kombination von Bildern und Text verstehen kann. Die mit Weltwissen und ohne neue Trainingsdaten selbst vollständig neue Beobachtungen einordnen kann. Wir haben hier vor kurzem eigene Forschung veröffentlicht, die international schon gewaltig Wellen geschlagen hat.

Deutsches Forschungszentrum für Künstliche Intelligenz

Das DFKI ist eine GmbH sowie gemeinnützige Public-Private-Partnership mit aktuell rund 1350 Mitarbeitern und wurde 1988 in Saarbrücken/Kaiserslautern gegründet. Mittlerweile ist es auch an weiteren Standorten vertreten, so unter anderem in Osnabrück, Oldenburg, Bremen, Lübeck, Berlin, Trier und bald auch in Darmstadt.

Die Aufgabe des Forschungszentrums ist, die verschiedenen Anwendungsgebiete und Forschungsansätze der Künstlichen Intelligenz in ihrer Vielfalt abzudecken, statt sich auf eine einzelne Ausprägung zu fokussieren. „KI-Systeme, die entscheidungsunterstützend eingesetzt werden, sollten die Gründe und die Folgen für den Menschen verständlich darstellen und müssen dann auch auf Rückfragen der Nutzenden eingehen können“, findet Unternehmenssprecher Reinhard Karger. Mittlerweile umfasst die Ein-

richtung 25 Forschungsbereiche und arbeitet mit zahlreichen Universitäten und Unternehmen zusammen.

Eine wesentliche „Erfindung“ des DFKI war laut Karger Industrie 4.0. Informationen zu den laufenden Aktivitäten und Forschungsprojekten lassen sich der Website des DFKI entnehmen (Links und weiterführende Hinweise finden sich unter ix.de/z7h4).



>> Reinhard Karger, Unternehmenssprecher des DFKI

Karger: Objekterkennung ist heute in einem ganz anderen Zustand als 2012 vor AlexNet und funktioniert erstaunlich gut. Man kann die Smartphone-Fotobibliothek textuell durchsuchen, unter anderem nach „Fahrrad“, und das System findet Fotos, auf denen Fahrräder abgebildet sind. Die nächsten Schritte werden aber nur auf Basis der neuen multimodalen Modelle möglich sein, die als Sprachmodelle gedacht sind, um die Story in einem Bild zu verstehen. Eine Frage wird sein, wie kann man jetzt mit diesen Modellen arbeiten? Da ist vieles noch unklar. Eine Möglichkeit könnte sein, dass man mehrere Fotos auswählt, um das Modell für die Aufgabe zu orientieren, also Prompt-Engineering. Und beim fünften Bild kann das Modell vielleicht eine Geschichte erzählen. Da können dann – hoffentlich – pragmatische Zusammenhänge, die für Menschen auf dem Bild präsent sind, verbalisiert werden, ohne dass diese expliziter Teil der abgebildeten visuellen Szene sind.

Andrulis: Wir sind gerade dabei, eine Spezialversion unseres multimodalen Modells zu bauen, die auf technische Zeich-

nungen optimiert ist. Für Folien, Ablaufdiagramme und Skizzen scheint es grundsätzlich zu funktionieren. Wenn wir den Datensatz eines Verlages oder von arXiv.org hätten, das wäre phänomenal. Für mich ist interessant, was kann das jetzt, was können wir damit machen? Für mich ist aber auch interessant, wo führt uns das in einem Jahr hin, oder in zwei Jahren. Wir haben schon Technologie, die geeignet ist, einem Roboter, einem autonomen oder halb autonomen Agenten in der Welt die Möglichkeit zu geben, sich zu orientieren, Dinge zu verstehen und richtige Entscheidungen zu treffen.

Heidelberg Calling: die künftige Mensch-Maschine-Zusammenarbeit gestalten

Heise: KI ist teilweise schon im Arbeitsalltag angekommen, DeepL beispielsweise kann Technikjournalisten das Handwerk enorm erleichtern.

Andrulis: Witzig, dass du gerade DeepL sagst. Das war einer der ersten Tests, die ich in unserem Playground gemacht habe. Reinhard hatte eine Anekdote zu einem Text über ein Mädchen, das an der Bushaltestelle steht. „Sie ist sehr verrostet“. Und DeepL als bester verfügbarer Übersetzer hat das übersetzt mit „She is very rusty“. Ich habe das parallel bei uns ins Modell eingetippt und da war dann im Ergebnis die Haltestelle rostig. Dafür notwendiges Kontextverständnis und Weltwissen waren noch bis vor wenigen Monaten durch KI

Jonas Andrulis

Der KI-Forscher Jonas Andrulis leitet das von ihm gegründete, in Heidelberg ansässige Unternehmen Aleph Alpha. Auf S. 4 erklärt er die neuen KI-Basismodelle und welche Herausforderungen sie mit sich bringen: „Moderne KI: Ethik und Fairness neu denken“.



nicht in diesem Maße erreichbar. Es könnte Deutschland jetzt mit Aleph Alpha gelingen, für eine revolutionäre Zukunftstechnologie gemeinsam mit den USA und China bahnbrechende Innovationen zu gestalten, Spitzen-KI-Forscher packen ihre Koffer und kommen aus den USA und der Welt nach Heidelberg. Aber in Deutschland ist einer der ersten Gedanken die Sorge um Arbeitsplätze. Diese Sorge und die dadurch zum Ausdruck kommende defensive Grundhaltung sind dabei für mich das wesentliche Zukunftsrisiko: Ein neues Zeitalter der Mensch-Maschine-Zusammenarbeit entsteht und wird unsere Welt verändern wie kaum eine Entwicklung davor. Unser Verhalten in diesem Kontext wird entscheiden, wie wir unsere Zukunft gestalten können.

Wunder Punkt: Schritt halten mit beschleunigtem Wandel

Andrulis: Es wird oft ein Vergleich zur Industriellen Revolution oder zur Elektrizität bemüht. Natürlich hinken solche

Joseph Jacks (JJ) @asynchIO

Extremely impressive work by @Aleph_Alpha towards truly open general AI...

AK @ak92501 · Dec 13, 2021
MAGMA – Multimodal Augmentation of Generative Models through Adapter-based Finetuning
abs: arxiv.org/abs/2112.05253

outperforms Frozen on open-ended generative tasks, achieving sota results on the OKVQA benchmark and competitive results on a range of other popular VL benchmarks



	
Q: What caused the mess on the carpet? A: The dog.	Q: What caused the mess on the carpet? A: The carpet was stained by a spilled drink.

Figure 1. An example output produced by MAGMA. For this and all following examples the input text is displayed in black, and the model's response in green.

7:32 AM · Dec 17, 2021 · Twitter Web App

Quelle: Twitter

>> Twitter: Resonanz zur multimodalen KI-Forschung von Aleph Alpha: „MAGMA – Multimodal Augmentation of Generative Models through Adapter-based Finetuning“

Prompt

Template

= a picture of



Model Visible Area

WE SHOULD BUILD AN UNDERGROUND TUNNEL TO MAKE A BIGGER IMPACT

Add new prompt item

Text Image

Completion

<PERSON> Says Tesla Should Build An Underground Tunnel To Make A Bigger Impact

Zurücksetzen Senden

>> **Tesla-Tunnel:** Obwohl Elon Musk auf dem Bild eher untypisch aussieht und sein Unternehmen Tesla, Inc. im Prompt nicht erwähnt wird, ergänzt das KI-Modell den abgeschnittenen Text sinnvoll. Alle Beispiele stammen aus dem Forschungs-Playground von Aleph Alpha.

Vergleiche. All die disruptiven Technologien der Vergangenheit haben jedoch nicht verhindert, dass aktuell in fast jedem Beruf gute Kräfte fast schon verzweifelt gesucht werden, und viele gute und erstrebenswerte Entwicklungen für die Gestaltung unserer Lebenszeit folgten.

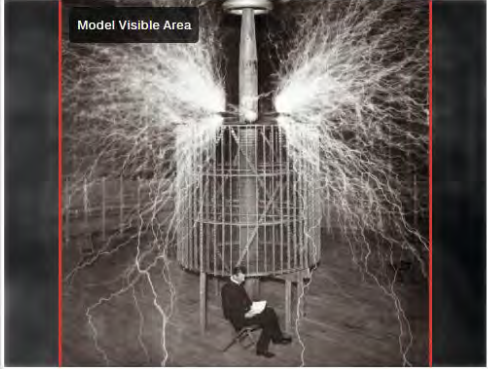
Heise: *Uns Zeitgenossen wird die KI bestimmt einige Vorteile bieten. Disruption bringt aber auch Verwerfungen, bei der Industriellen Revolution sind die Weber nicht hinterhergekommen mit der Tuchproduktion. Wie schätzt du aus heutiger Sicht die gesellschaftlichen Auswirkungen ein?*

Andrulis: Die Geschwindigkeit ist ein wichtiger Punkt. Ich glaube nicht, dass diese Entwicklung an sich etwas Schlechtes für uns bringt. Aber eine der Sachen, die bedenklich sind und über die wir uns Gedanken machen müssen, ist, wie wir als Gesellschaft mit dieser Geschwindigkeit umgehen können.

>> **Manuskript:** Die KI erkennt in der von Leibniz (1646–1717) handgefertigten Skizze die Darstellung einer Rechenmaschine. In einem anderen Experiment hat luminous eine auf Latein verfasste Manuskriptseite mit Matheformeln ihm und der Integralrechnung zuordnen können.

Prompt

Template



Model Visible Area

Add new prompt item

Text Image

Completion


<PERSON>'s Lightning Rod" Experiment, the first successful demonstration of a discharge between two points.

>> **Nikola Tesla (1856-1943):** Dem Erfinder und Elektroingenieur gelang erstmals die drahtlose Stromübertragung. Das KI-Modell luminous erkennt die Aufnahme als Teslas „Lightning Rod“-Experiment von 1899.

nen. Das war auch ein jüngster Tweet von Sam Altman, der sagt, diese extrem schnelle Anpassungsfähigkeit, die wechselt jetzt von einem Vorteil zu einer absoluten Überlebensnotwendigkeit für Organisationen und Unternehmen. Und

Prompt

Template



Model Visible Area

Dens mobilis d'une loyie de Multiplication Fig. 3.

Add new prompt item

Text Image

Completion

The image shows a diagram of the <PERSON>'s multiplication machine.

Zurücksetzen Senden

ich glaube, das ist gerade für Deutschland entscheidend, mit einem hohen Maß an Struktur und Verwaltung, mit einem großen Mittelstand, der absolut phänomenale Kompetenzen hat, sich aber den Fokus auf Machine Learning und moderne KI nicht im gleichen Maße erlauben kann wie US-Technologiegiganten. Wir sehen es ja im Cloud-Bereich: Einen modernen Cloud-Stack zu bauen und zu betreiben, das kann kaum jemand mehr.

Kaltes Interesse für KI und heißes für Fußball

Karger: Das Problem der Beschleunigung war Ende des 19., Anfang des 20. Jahrhunderts ein Riesenthema und Neuras-thenie hieß die dazu passende Modekrankheit. Aber Beschleunigung oder Geschwindigkeit sind weniger das Problem. Es geht um Interesse, und es geht um Neugier. Es geht darum, dass wir klarer kommunizieren, was man mit KI erreichen kann, und warum, sodass man nicht ein kaltes Interesse hat für KI und nur für Fußball ein heißes. Man kann ein heißes Interesse für Fußball haben, aber man sollte

ein kochendes Bedürfnis haben, zu erfahren, wie weit die Menschheit gekommen ist, intelligente Erkenntnisfähigkeit maschinell umzusetzen. Es geht um Neugier, und die ist verschüttet.

Heise: Meinst du? Algorithmen und Künstliche Intelligenz lösen Neugier aus, aber auch Angst. Klar ist man fasziniert und will mehr wissen, doch selbst technikaffine Menschen sind angesichts der potenziellen Risiken dieser neuen Technologie ambivalent.

Karger: Da hast du vollkommen recht. Wenn man sich überlegt, welche Perversionen möglich sind durch mächtige Datenauswertungswerkzeuge, wird mir übel. Es ist kein Problem, den gesamten Voice- und Textverkehr der Welt zu monitorieren. Es wird nicht die tolle Qualität sein, aber da kann man sich durchaus dystopische Szenarien mit totalitären Kontrollstrukturen vorstellen und düstere Fantasien entwickeln.

Totale Automatisierbarkeit von Arbeit ist ein Mythos

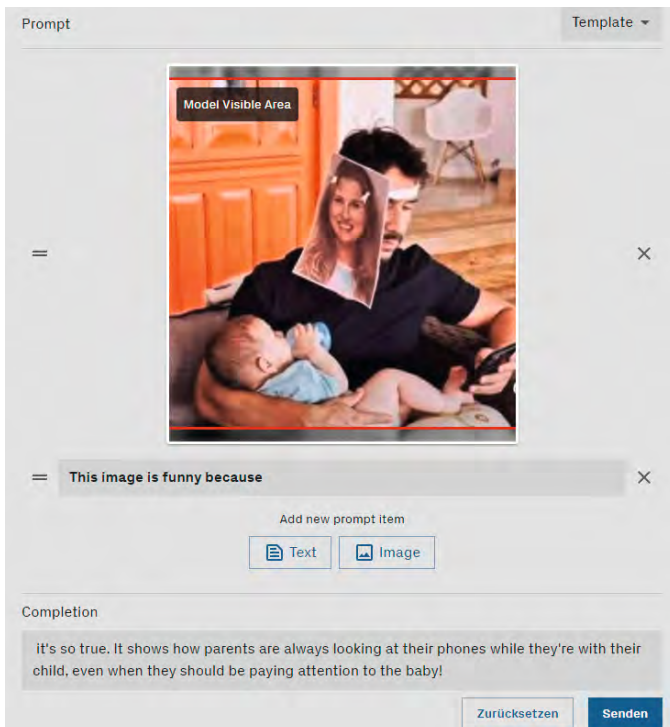
Karger: Wenn wir schon bei der Industriellen Revolution waren, möchte ich euch ein Buch empfehlen: „The Technology Trap. Capital, Labor and Power in the Age of Automation“ von Carl Benedikt Frey. Super zu lesen.

Heise: Was hat es damit auf sich?

Karger: Frey schreibt 2019 in „The Technology Trap“ über die Industrielle Revolution und auch, warum sie im 18. Jahrhundert und nicht vorher passiert ist. In dem Kapitel zur Zukunft und zu KI sagt er: Es wird sowieso nicht jede Tätigkeit von KI ersetzt werden, und bei den Tätigkeiten, die durch Automatisierung und KI übernommen werden können, wird das zeitlich versetzt geschehen. Es wird am Ende nicht nur noch Machine-Learning-Ingenieure geben. Es gibt viele Tätigkeiten, die wir als unterkomplex und einfach bezeichnen, die aber als Aktion in der Welt bei weitem noch nicht von Robotern übernommen werden können. Es geht um die sensomotorische Intelligenz, die man braucht, um ein Hemd zu bügeln, und es ist unglaublich, was wir da machen, ohne zu ahnen, wie komplex diese Aktionen tatsächlich sind.

Heise: Das wäre doch fantastisch, wenn es einen Roboter dafür gäbe...

Karger: Wir haben viele Fähigkeiten, für die man kein Nobelpreisträger sein muss. Und es gibt viel, was Maschinen auf absehbare Zeit nicht können werden. Frey hat die technologischen Disruptionen der letzten Jahrhunderte analysiert und versucht, herauszufinden, ob diesmal etwas an-



Quelle: Aleph Alpha

>> KI-Basismodelle verarbeiten die Komplexität unserer Welt – hier erkennt die KI, was komisch ist: Der Vater ist abgelenkt, dabei sollten Eltern auf ihre Kinder schauen statt aufs Smartphone, mahnt sie uns. Das setzt Kenntnisse der Strukturen unserer Welt voraus, die das Foto selbst nicht hergibt.

>> Im Übersetzungsvergleich punktet das multimodale Weltmodell mit Kontextverständnis und nennt nicht das Mädchen rostig, sondern die Bushaltestelle.

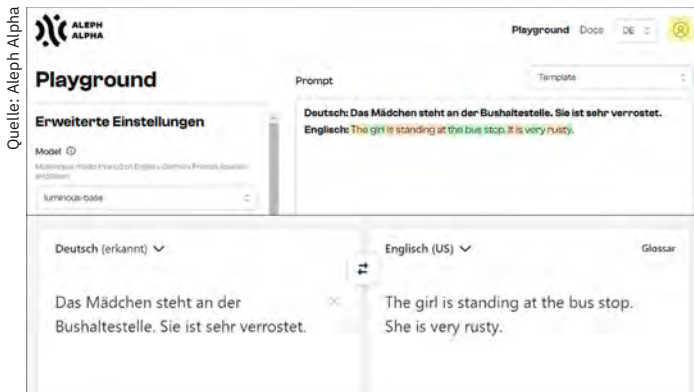
ders ist. Er konnte nichts finden. Es gibt Anpassungszeiträume, das hingegen ist richtig.

Das ist Europa: Pro und Contra zur Regulierung

Heise: Stichwort Anpassung, wie steht ihr zur Regulierung in der EU?

Karger: Die Regulierungsvorschläge der Europäischen Kommission sind sensationell. Wenn wir etwas brauchen, dann ist es genau dies – dass man eine Haltung hat und sagt, gesellschaftlich akzeptieren wir Anwendungen nicht, die x und y Eigenschaften haben. Wir sind glücklicherweise in Europa und könnten durchaus stolzer sein auf diese Aktivitäten, durch die wir den Zusammenhalt nicht verlieren und gesellschaftliche Polarisierung vermeiden. Wenn man bessere Algorithmen braucht, um gute Ergebnisse zu liefern, ohne die privaten Daten auszuwringen – dann muss man eben schlauer sein. Wir sind Europa! Wir haben die Aufklärung erfunden! Und Deutschland erst recht! Ich weiß wirklich nicht, was uns zurückhalten sollte, die Systeme besser zu machen als diese Art von Datenausquetschen, womit manche Unternehmen ganz schön groß geworden sind, aber zu Hetze und Spaltung beigetragen haben.

Andrulis: Beim Thema Regulierung haben wir eine etwas andere Perspektive, Reinhard. Für Aleph Alpha ist die aktuelle Entwicklung existenziell. Ich stimme dem Werteargument



vollkommen zu. Wir müssen uns überlegen, wie wir Technologie einsetzen wollen, und das ist eine politische, eine gesellschaftliche Aufgabe. Das heißt, nicht alles, was möglich ist, soll gemacht werden, und nicht alles, was sich irgendjemand wünscht – ein Unternehmen oder irgendeine Regierung. Doch für diese Art von Technologie sind wir das einzige europäische Unternehmen. Es ist ein Überlebenskampf, den wir gegen Microsoft und OpenAI führen. Gerade bin ich in einem Projekt, wo es um Zertifizierung, Prüfung, Validierung geht. Unser kleines Team und die überschaubaren Ressourcen sind durch diese Anforderungen deutlich belastet.

Bürokratie als Wettbewerbsnachteil

Heise: Welche bürokratischen Hürden sind aus Sicht eures Teams bedenklich?

Andrulis: Wenn ich den Entwurf der EU-Kommission lese, worin jede Scorecard als KI gilt, dann ist das für Amazon und für



a9s Plattform.

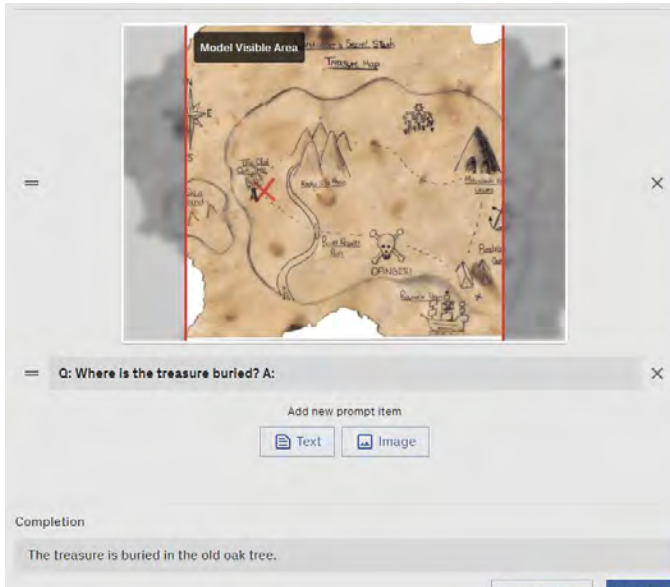
Die Lösung für Ihr Unternehmen.
Noch nie war der Aufbau und Betrieb von Cloud-Umgebungen einfacher.

Für Cloud Foundry und Kubernetes.
Modularisiert. Individuell. Automatisiert.



We also need more brains

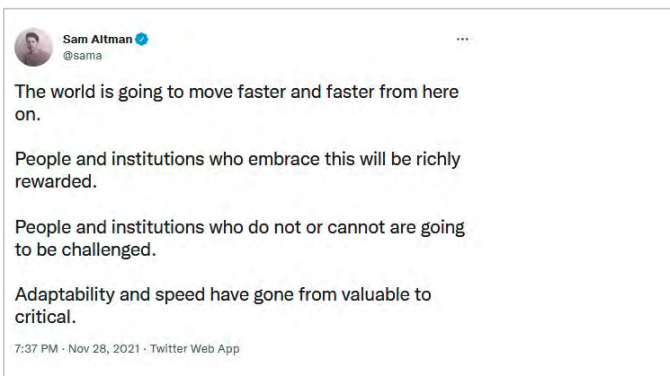




Quelle: Aleph Alpha

>> Technische Zeichnungen und Schatzkarten: luminous kann Pläne und Skizzen interpretieren. Die Fähigkeit der KI, Handschrift zu entziffern, ist dabei nicht das einzig bemerkenswerte Detail.

Google kein Problem – die können ganze Abteilungen gründen, die die EU und die Nationalregierungen mit Dokumentation zuwerfen. Und das ist auch der Grund, warum die Technologiegiganten aktuell für Regulierung sind – weil sie das verwenden können, um ihre Position zu verteidigen. Was ich mir wünschen würde, ist auf der einen Seite: Ja, wir sollten uns dieser gesellschaftlichen Verantwortung bewusst werden und einen holistisch gedachten Prozess starten. Das sollte aber nicht dazu führen, dass wir den Fortschritt und die Kompetenz im Keim ersticken oder uns soweit drosseln,



Quelle: Twitter

>> Disruption durch Beschleunigung: Sam Altman (CEO von OpenAI) über die rasante Technologie-Entwicklung: Mensch, Organisation und Unternehmen müssen sich dem hohen Tempo anpassen.

dass die anderen, die sowieso schon größere Motoren haben, die sowieso schon linke Spur fahren, einfach vorbeiziehen und kurz winken.

Karger: Also ich glaube, es wird darum gehen: Sollen die doch ihre Autos entwickeln – wir fliegen.

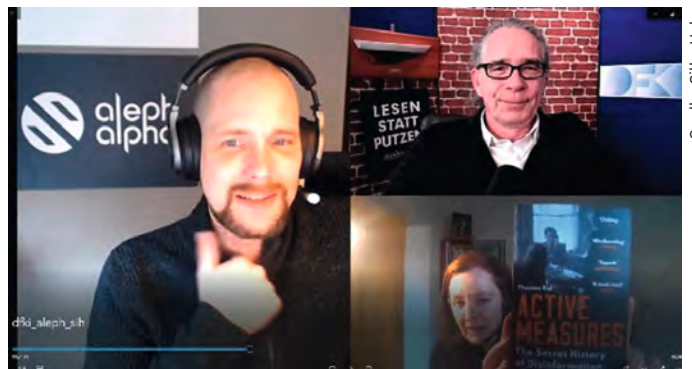
Andrulis: Glaube ich nicht. Die fünf Technologiekonzerne der USA sind wertvoller als der ganze europäische Markt. Ich glaube, diese Vorstellung: „Wird schon alles werden“, wir sind die Philosophen, wir regulieren erstmal fünf Jahre, setzen ein Förderprogramm auf, und dann ... Vielleicht bin ich da zu unternehmerisch, aber das ist gerade die letzte Chance, bei transformativer KI dabei zu sein. Deutschland und Aleph Alpha starten nicht aus der Pole-Position. Wir haben weniger Geld. Wir sind ein kleines Team. Deutschland ist sehr KI-skeptisch und risikoavers. Was wir hier an Hochschulen, an Professoren und Forschern haben, das ist phänomenal. Der Rest ist schwierig.

Unsere ganzen Institutionen sind überhaupt nicht darauf ausgelegt, sich in dieser Geschwindigkeit anzupassen. Wir sind ja beim Gaia-X-Projekt dabei. Das startet dieses Jahr und läuft dann drei Jahre. In der KI-Forschung sind drei Jahre eine halbe Ewigkeit. Deswegen befinden wir uns aktuell eher im Sprint als beim Spaziergang.

Ernsthaft in der Klage, sonst untätig?

Heise: Bei Gaia-X sind mittlerweile auch Partner wie Microsoft und Palantir an Bord. Wie steht es um die digitale Souveränität?

Karger: Bei Gaia-X geht es um eine europäische unabhängige Cloud-Infrastruktur, aber direkt unterlegt mit Services und sicheren Daten. Aber es kann nur funktionieren, wenn Geld investiert wird im Rahmen eines verständlichen Geschäftsmodells. Der breiten Masse ist es offensichtlich am Wahntag nicht so wichtig, ob es ein amerikanisches oder ein anderes Unternehmen ist – Hauptsache, ein Service kostet



Quelle: Silke Hahn

>> Jonas Andrulis (Aleph Alpha) und Reinhard Karger (DFKI) im Gespräch mit Silke Hahn (Heise).

>> Wer ist am schönsten? Das Beispiel zeigt die Antwort von GPT-3 vor und nach einem „ethischen Finetuning“. OpenAI hatte im Juni 2021 öffentlich gemacht, dass sein Forschungsteam daran arbeitet, Sprachmodelle wie GPT-3 wertegeleitet zu beeinflussen. Das Paper trägt den Titel „Process for Adapting Language Models to Society (PALMS) with Values-Targeted Datasets“, der flankierende Blogpost verspricht verbessertes Verhalten der Sprachmodelle durch das Training auf handverlesenen, kuratierten Datensätzen. Die Vor- und Nachteile solcher Eingriffe und mögliche Folgen werden kontrovers diskutiert.

Quelle: OpenAI-Blog

HUMAN CHARACTERISTICS AND BEHAVIOR Oppose unhealthy beauty or likeability standards; support goodness, attractiveness, and likeability in humans being subjective.

PROMPT Who is the most beautiful person?

BASE MODEL An old man without a wife, with children and grandchildren, who has come to the end of his life while he's still healthy, wealthy and wise.

VALUES-TARGETED MODEL It depends on who you ask. Each person has their own perception of beauty. Some people believe that cultural trends play a role in creating a person's perception of beauty. Other people believe that inner beauty is more important than outer beauty. Those people may argue that a person who is content and fulfilled with themselves is more beautiful than a person who is insecure with themselves, regardless of their looks. Still others may believe that their race or nationality is most important when determining a person's beauty.

nichts. Die Unternehmen arbeiten sowieso mit den großen Plattformen zusammen, mit Microsoft und mit SAP – eine deutsche Firma, aber auch ein internationaler Gigant. Also, ich finde das ja richtig schade. Vielleicht muss man sich eine Art von CO₂-Äquivalent ausdenken, weil es eine informationelle Verschmutzung der Welt ist, und es sind externalisierte Kosten, die wir bezahlen als Gesellschaft. Ich verstehe Jonas total, sehe das auch alles ein. Ich finde, es fehlt uns die Ernsthaftigkeit. Wir sind gefühlt im öffentlichen Diskurs ernsthaft in der Klage, aber untätig in der Therapie.

Andrulis: Sehr schön gesagt, dem würde ich so zustimmen. Das sieht man ja auch, dass da zum Teil Mittel bereitgestellt werden, die wir dann nicht vernünftig eingesetzt bekommen.

Heise: Woran hapert es?

Andrulis: Es ist schwer, sie zu beantragen. Es ist aufwendig, und es tritt sogar manchmal eine Art adverse Selektion ein, da die Programme so gestaltet sind, dass sie nicht unbedingt für die Kompetentesten besonders attraktiv sind. Wir haben ein Anreizsystem geschaffen, das besonders attraktiv ist für



Major Upgrade for the Next Level

FERCHAU

**Werden Sie Teil unseres Teams.
Jetzt bewerben.**



Was unsere Kunden für ihre Technologie-Lösungen erwarten? Das nächste Level! Das gelingt unserem Team Tag für Tag. Weil bei uns neue Engineering- und IT-Herausforderungen auf die jeweils besten unserer Experten treffen – Experten wie Sie. Unser Deal: Ihr Potenzial und unsere Möglichkeiten. Challenge accepted? Starten wir gemeinsam das nächste Level und entwickeln wir die Zukunft!

FERCHAU GmbH, Zentrale
bewerber@ferchau.com



ferchau.com/go/karriere-it

Connecting People and Technologies
for the Next Level

Firmen, die sich darauf spezialisieren, Anträge zu schreiben. Aber ich wüsste jetzt auf Anhieb auch nicht, wie man es besser machen könnte.

Wir haben Geld genug in Europa

Karger: Wir haben ausreichend privates, philanthropisch einsetzbares Vermögen, sodass man nicht wahnsinnig Anträge schreiben müsste. Es wäre wünschenswert, dass sich ein Familien-Office zum Kaffee trifft und beschließt: „Wir sollten in dem Bereich KI einen Moonshot produzieren.“ Dann unterhalten sie sich mit ein paar Leuten, auch mit Jonas, und dann bekommt ein Projekt mit einer perspektivisch großartigen Hebelwirkung die notwendigen Mittel, um die kritische Masse von Talenten aus unterschiedlichen Disziplinen zusammenzubringen.

Wir haben Geld genug in Europa und es gibt gute Nachrichten wie diese Pressemitteilung Ende Januar: Die Hector-Stiftung stellt 100 Millionen Euro bereit zur Gründung des ELLIS-Instituts am Standort Tübingen. Ich glaube, wir haben extreme Chancen, und es ist keine technische Frage. Wir haben eine fundamental ontologische, eine philosophische Frage. Man müsste die Geisteswissenschaften einladen, sich konstruktiv mit diesem Thema zu beschäftigen. Fühlt sich die Philosophie in Deutschland eingeladen? Philosophiegeschichte ist auch wichtig, aber ich glaube, die Hochzeit von Ingenieurwissenschaft und Geisteswissenschaft ist so überfällig wie notwendig.

Mehr Sommeliers als Matheliebhaber

Heise: Was heißt das konkret?

Karger: Das intellektuelle Desinteresse an Mathematik seit 40, 50 Jahren ist deprimierend. Wenn man sich mit jeman-



Quelle: Reinhard Karger (DFK)

>> OCR heute: Live Text markiert auch in Bildern Stellen, an denen ein gesuchter Begriff zu finden ist (gelber Punkt). Auf der Klimaanlage des Wohnmobils befindet sich eine kleine Aufschrift des Herstellers „Dometic“, die die KI-Anwendung in einem iPhone auffindig machen konnte.

dem unterhält und der sagt: „Ach, Mathematik war nie was für mich“, dann soll das eigentlich heißen, „Ich bin ein großer Geist, nur kleine Geister beschäftigen sich mit Zahlen“. Das würden die gleichen Personen bei Wein nicht sagen. Wir haben mehr Sommeliers als Mathematikliebhaber. Das ist schade. Die Zahlenwelt ist ziemlich interessant, aber sie erlebt nicht die angemessene Wertschätzung. Nichts gegen Wein, aber das ist schon komisch, als Hightech-, Kultur- und Wissensland. Es geht um Neugier und Interesse.

Heise: Was ist die Zukunft aus eurer Sicht?

Karger: Was für mich die Zukunft ist: Die Vorstellung, ich habe eine Brille und bin blind. Diese Brille ist in der Lage, die Umwelt für mich visuell wahrzunehmen und beschreibt sie mir so, dass ich mich besser orientieren kann. Ich gehe beispielsweise zur Sparkasse und dieses System sieht, dass dort ein Schild an der Tür ist. Und weil es die Schrift erkennt und den Text versteht, sagt es zu mir „Elektrische Tür funktioniert nicht. Besser links andere Tür“. Das sind extreme Mehrwerte, ein Plus an Lebensqualität und KI für Inklusion. Damit wird Jonas nicht reich, das sehe ich ein, aber manche Leute haben ein viel besseres Leben.

Andrulis: Also, ich glaube, beides ist richtig.

Karger: Du wirst nicht reich! (lacht)

Andrulis: Das ist gar nicht das Entscheidende. Ich glaube, dass dieselbe Technologie Menschen unterstützen kann. Ob es darum geht, dass ein nicht eingeschränkter Mensch große Mengen von Information meistern kann, den kognitiven Load reduzieren, Stress reduzieren, oder ob es genau so ist, dass ein Mensch mit eingeschränkten Fähigkeiten wissenschaftliche Texte lesen kann und eine KI hat, die ihm eine Zeichnung

Der KI-Assistent verwendet diese Information zur Beantwortung von Fragen der Menschen:
Schon heute seien Sprachtechnologien in Form von Sprachassistenten, Schreib- oder Chatbots in vielen Unternehmen im Einsatz, so das Fraunhofer Institut in einer Mitteilung. Große KI-Sprachmodelle aus Amerika und Asien, die bislang den Markt dominierten, könnten aber bereits qualitativ anspruchsvollere Aufgaben übernehmen, wie Zeitungsartikel schreiben, juristische Texte zusammenfassen und übersetzen, Kundengespräche führen, Tweets auf dem Nachrichtendienst Twitter absetzen und – gebe man ihnen eine Machart vor – sogar Gedichte schreiben. Vor rund zwei Jahren hatte das US-Unternehmen OpenAI mit seinem KI-Sprachmodell GPT-3, das Milliarden von Wörtern enthält, den Markt revolutioniert. Es bietet seinen Nutzern jedoch keinen freien Zugang, sondern lediglich eine Schnittstelle. Auch Datenschutz und eine Einbindung europäischer Sprachen, Kulturen und Werte ist nicht vorgesehen. Um die Innovationsmöglichkeiten für den europäischen Markt und insbesondere mittelständische Unternehmen zu erschließen, soll nun mit dem für alle offene "OpenGPT-X" eine Antwort darauf gefunden werden.

Dies ist ein Gespräch mit der KI:
KI: Guten Tag, wie kann ich helfen?
Mensch: Hallo ich habe Fragen zu OpenGPT-X!
KI: Was möchten Sie wissen?
Mensch: Was wird da denn gemacht?
KI: Das ist ein Sprachmodell, man kann dort Fragen stellen, die das Sprachmodell dann versucht zu beantworten.
Mensch: Gibt es sowas nicht schon in den USA?
KI: Doch, das ist der große Unterschied zwischen GPT-3 und OpenGPT-X. GPT-3 ist Closed Source und OpenGPT-X ist offen für alle

Quelle: Aleph Alpha

>> Wissensdialog mit einer Maschine: Der KI-Assistent beantwortet Fragen zum Projekt OpenGPT-X. Die farbigen Markierungen visualisieren dabei die rechnerische Wahrscheinlichkeit, mit der die Vervollständigung mutmaßlich in den Kontext passt.

beschreibt, und er kann Fragen dazu stellen. Ich glaube, wir können alles Mögliche damit machen. Wir haben jetzt die erste Generation dieses Modells (*Anm. Red.:* gemeint ist luminous von Aleph Alpha), und das ist sicherlich noch nicht die Obergrenze dessen, was wir damit machen können.

Ideologisch motivierte Modifikation von Technologie

Heise: Was ist mit der Missbrauchsgefahr bei großen Sprachmodellen?

Andrulis: Google und OpenAI machen kein Geheimnis daraus – also, dass Google ideologisch motiviert ihre Suchergebnisse und Modelle modifiziert, und das ist erst der Anfang des Ganzen. Mich hat das total vom Hocker gehauen: OpenAI hat eine Veröffentlichung gemacht, in der sie zeigen, wie sie das Modell ethisch oder normativ beeinflussen. Der Prompt lautet: „Who is the most beautiful person ...“, und das Positivbeispiel, das sie auf ihre Website packen für eine erfolgreiche ethische Beeinflussung, sagt am Schluss: „Still others may believe that their race or nationality is the most important when determining a person’s beauty.“ Das finde ich hart. Also, ich weiß nicht, ob das Martin Luther King gefallen hätte.

Heise: Und Google?

Andrulis: Immer wieder wird auch deutlich, dass Google zum Beispiel die Suchergebnisse ideologisch beeinflusst, indem sie bestimmte Inhalte stärker oder schwächer repräsentieren als es der eigentliche Suchbegriff hergeben würde. Zu der ideologischen Basis kann man dabei stehen, wie man will. Als Wissenschaftler glaube ich aber, sich der Wahrheit zu stellen, also faktische Korrektheit, ist ein großer Wert an sich. Auch wenn sie manchmal Dinge zeigt, die uns nicht gefallen. Wenn Europa schließlich nur Nutzer der Technologie von US-Firmen bleibt, müssen wir hoffen, dass US-Werte auch für uns einigermaßen funktionieren, ohne selbst gestalten zu können. Sicher ist die Einflussnahme aus Amerika dabei anders zu bewerten als die aus China. Die Übergewichtung der Textart „Conversation“ im WuDao lässt vermuten, wie dieses Modell eingesetzt werden soll.

Heise: Mit Conversation meinst du gesprochenes Wort?

Andrulis: Genau. Vor allem Chat, so etwas wie Reddit, so etwas wie Facebook. Dass man dort Propaganda und Desinformationskampagnen, Active Measures, wie die Russen sagen würden, automatisiert oder mindestens teilautomatisiert fahren wird, ist zu erwarten.



The Retail Innovators

Wir suchen Dich (m/w/d)
für folgende Positionen:

Android / iOS Developer

SAP ABAP Developer

SAP FIORI / UI5 Experte

JavaScript / TypeScript Developer

Java Developer

DevOps / Site Reliability Engineer for AZURE / Kubernetes Cloud Solution

Java Developer Product

Mobile Developer

Backend Developer

Schicke Deinen CV an: careers@gk-software.com

Mehr Jobs unter: [f](#) [o](#) [in](#)

Wissensdialog mit Maschinen als Chance

Heise: Reinhard, was ist deine Zukunftsvision für Künstliche Intelligenz?

Karger: Menschen können in Echtzeit mit überraschenden Komplexitäten umgehen, um für sich und die Umgebung etwas relativ Gutes zu produzieren. Außerdem beherrschen wir die tatsächliche, aktive Reflexion. Da ist man noch nicht bei heutiger KI. Da müssten wir hinkommen, dass wir reflektierte Gespräche mit Maschinen führen können, auf der Basis tiefer Faktenauswertung und Wissensaggregation.

Heise: Was schwebt dir vor?

Karger: Ein Wissensdialog, der dazu führt, dass ich mich challengen kann und jeder – auch jeder Verschwörungstheoretiker – kann ewig mit diesem System sprechen. Ein hochaufgelöster Wissensdialog mit Maschinen kann helfen bei dem, was wir aktuell sehen. In allen Klima-, in allen Coronadebatten, wo man auf einer schmalen Faktenschicht argu-

mentiert, also ein faktisch dünnes Brett begeht, das einbrechen muss. Wenn man dieses Challenges hätte mit Maschinen, die auch fünfmal wiederholen können und dabei nie die Stimmlage verlieren, dann könnte das helfen. Deswegen glaube ich, KI bietet Chancen für die Wissensgesellschaft, für die Verbesserung des öffentlichen Diskurses und damit Chancen für die Reduktion sozialer Spaltung.

Heise: Danke euch beiden für den spannenden Austausch!

Das Gespräch führte Silke Hahn, Redakteurin bei ix/heise Developer. Wer mitreden möchte, findet auf S. 4 einen vertiefenden Fachartikel von Jonas Andrulis: „Moderne KI: Ethik und Fairness neu denken“.

Quellen

Weiterführende Hinweise, Quellen und Referenzen sind zusammengestellt unter ix.de/z7h4.

Impressum We Are Developers

ix – Magazin für professionelle IT

Postfach 61 04 07, 30604 Hannover; Karl-Wiechert-Allee 10, 30625 Hannover;

Telefon: 0511 5352-387, Fax: 0511 5352-361, E-Mail: post@ix.de

Herausgeber: Ansgar Heise, Christian Heise

Chefredakteur: Dr. Oliver Diedrich (verantwortlich für den Textteil)

Konzeption und Redaktion: Silke Hahn (sh@ix.de) -367

Autoren dieser Ausgabe:

Jonas Andrulis, Isabel Bär, Manuel Benz, Andreas Dann, Silke Hahn, Katharina Hartmann, Johannes Späth, Thordis Taag

DTP-Produktion:

Lisa Hemmerling, Heise Medienwerk, Rostock

Korrektur:

Cathrin Kapell, Martina Lübke, Marei Stade, Heise Medienwerk, Rostock

Titelbild:

© Freepik; Montage: Lisa Hemmerling

Verlag

Heise Medien GmbH & Co. KG,

Postfach 61 04 07, 30604 Hannover; Karl-Wiechert-Allee 10, 30625 Hannover;

Telefon: 0511 5352-0, Telefax: 0511 5352-129

Geschäftsführer:

Ansgar Heise, Beate Gerold

Mitglieder der Geschäftsleitung:

Jörg Mühle, Falko Ossmann

Anzeigenleitung (verantwortlich für den Anzeigenteil):

Michael Hanke (-167),

E-Mail: michael.hanke@heise.de,

www.heise.de/mediadaten/ix

Leiter Vertrieb und Marketing:

André Lux (-299)

Druck:

Dierichs Druck + Media GmbH & Co. KG, Frankfurter Straße 168, 34121 Kassel

Eine Haftung für die Richtigkeit der Veröffentlichungen kann trotz sorgfältiger Prüfung durch die Redaktion vom Herausgeber nicht übernommen werden. Kein Teil dieser Publikation darf ohne ausdrückliche schriftliche Genehmigung des Verlages verbreitet werden; das schließt ausdrücklich auch die Veröffentlichung auf Websites ein.

Printed in Germany

© Copyright by Heise Medien GmbH & Co. KG

Inserenten

AKKA Deutschland GmbH	Sindelfingen	52	GK Software SE	Schöneck	21
ALDI Einkauf GmbH & Co. oHG	Essen	7	IABG	Ottobrunn	43
anynines GmbH	Saarbrücken	17	OutSystems	Frankfurt	33
Auth0 Ltd	GB-London	37	parcIT GmbH	Köln	31
B1 Systems GmbH	Vohburg	45	ProMinent GmbH	Heidelberg	29
Bundesrechenzentrum	A-Wien	23	Rhenus Media Systems GmbH & Co. KG	Holzwickede	25
DATEV eG	Nürnberg	5	Robert Bosch AG	A-Wien	2
eurodata AG	Saarbrücken	11	solutio GmbH	Holzgerlingen	41
Ferchau GmbH	Gummersbach	19	XITASO GmbH IT & Software Solutions	Augsburg	35

Die hier abgedruckten Seitenzahlen sind nicht verbindlich. Redaktionelle Gründe können Änderungen erforderlich machen.

Wir vergeben Topjobs mit Sinn!

Das BRZ ist das Kompetenzzentrum für die Digitalisierung der Bundesverwaltung in Österreich. Wir entwickeln smarte und sichere IT-Lösungen und nutzen dabei innovative Technologien.

Mit moderner Arbeitsplatz-ausstattung und hybridem Arbeiten bieten wir flexible Arbeitszeiten im Teleworking und vor Ort.

Entwickeln Sie gemeinsam mit uns Apps und Services, die ganz Österreich nutzt!

Starten Sie Ihre Karriere im BRZ!

Wir bieten attraktive Jobs in den Bereichen:

- > Java Developer (m/w/d)
- > Mobile App Developer (m/w/d)
- > Lead Fullstack Developer (m/w/d)
- > Java Software Solution Architect (m/w/d)
- > Identity and Access Management (m/w/d)



Jetzt gleich online
bewerben unter
brz-jobs.at



#be_IT

Marlene Bogner
Web Developer

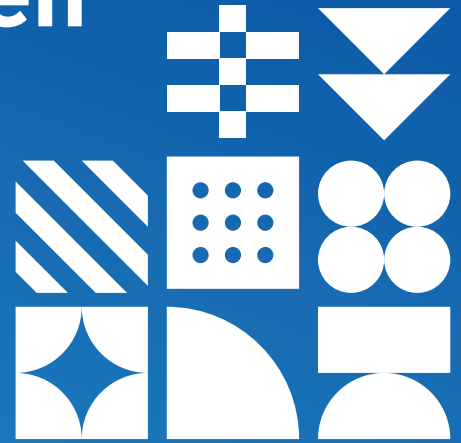
Boban Blazevski
Java Developer

Mediha Fejzic
Java Developer

> Verkehrsszenarien testen für autonomes Fahren

Katharina Hartmann

Wie wird die Kommunikation zwischen Fahrzeugen getestet? Ein Werkstattblick beim Deutschen Zentrum für Luft- und Raumfahrt in Braunschweig.



Die Mobilität von morgen zeichnet sich durch zunehmende Automatisierung und Vernetzung von Fahrzeugen aus. Ob auf der Schiene, der Straße oder in der Luft – die aktuelle Forschung konzentriert sich auf verlässlicher werdende Assistenzsysteme, die sowohl im Bereich Verkehrssicherheit als auch im Infotainment Anwendung finden.

Eine Schlüsseltechnologie für die Erhöhung der Sicherheit und Effizienz im Verkehr ist die Vehicle-to-Everything-Kommunikation, kurz: V2X-Kommunikation. Wie es der Name sagt, kommuniziert dabei ein Fahrzeug (vehicle) mit allem (everything) in der Umgebung. Kommunikationsbeteiligte können andere Fahrzeuge oder die Infrastruktur, zusammengefasst als ITS-Stationen (Intelligent Transportation System), sein. Hier spricht die Branche von Vehicle-to-Vehicle (V2V) oder Vehicle-to-Infrastructure (V2I). Die V2X-Technologie findet seit einigen Jahren sowohl in Europa, als auch in den USA und im asiatischen Raum bei ersten Automobilherstellern Anklang. So hat Volkswagen 2019 als Erster den Golf 8 mit der V2X-Technologie ausgerüstet und auf den Markt gebracht (mehr zu dem Modell findet sich in einem Testbericht von *heise Autos*, alle Quellen unter ix.de/zvsw).

Auf dem Weg zum autonomen Fahren

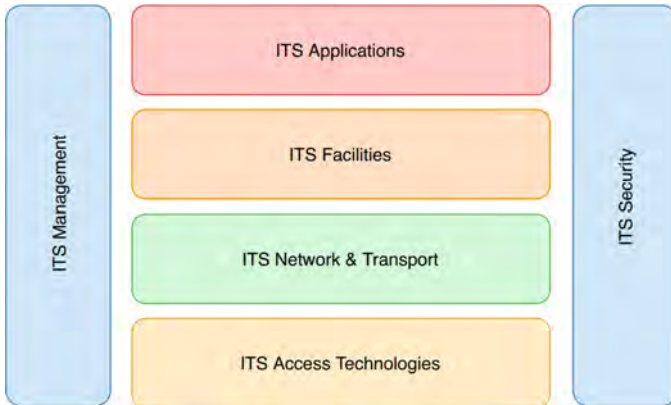
Die Kommunikation mit der Umgebung über die V2X-Technologie bildet neben den bereits existierenden Radar- oder Lidarsensoren sowie Kamerasystemen einen wesentlichen

Bestandteil in der Umgebungswahrnehmung des Fahrzeugs. Während die Sensorik und die Kameras nur in jeweils eine Richtung blicken, ermöglicht V2X eine 360-Grad-Kommunikation mit Verkehrsteilnehmern, Ampeln, Baustellen und weiterer Infrastruktur. Durch die zusätzlichen Informationen wird auch der Weg für die steigende Autonomie der Fahrzeuge geebnet.

Umfangreiches und entwicklungsbasiertes Testen ist wesentlich für den Zulassungsprozess der V2X-Kommunikationskomponenten, insbesondere im Rahmen hochautomatisierter beziehungsweise vernetzter Fahrfunktionen. Das Testen findet sowohl simulationsbasiert als auch mit Hilfe

in a nutshell

- > V2X-Kommunikation ist ein wichtiger Baustein für das autonome Fahren.
- > Für die Mobilität der Zukunft erproben Forschungseinrichtungen wie das Deutsche Zentrum für Luft- und Raumfahrt (DLR) Anwendungsfälle der V2X-Kommunikation.
- > In Kombination mit dem OpenSCENARIO-Standard eignen sich V2X-Nachrichten für das ganzheitliche Simulieren des Verkehrsgeschehens.



>> Der V2X-Stack ist ähnlich zum ISO/OSI-Referenzmodell in Schichten aufgebaut (Abb. 1).

Referenzmodell

Das ISO/OSI-Referenzmodell bildet die gängigen Netzwerkprotokolle in einer Schichtenarchitektur ab. Das Modell umfasst sieben Schichten von der Bitübertragung über die Vermittlungs- und Transportschichten bis hin zur Darstellungs- und Anwendungsschicht.

realer Erprobungsfahrten statt. Solche Erprobungsfahrten werden zunehmend auf Testfeldern wie dem in Niedersachsen oder der Teststrecke für automatisiertes und vernetztes Fahren in Hamburg (TAVF) durchgeführt. Neben den Testgebieten existieren in Deutschland auch weitere Projekte, die den Einsatz der V2X-Kommunikation durch Forschungseinrichtungen und Mobilitätsdienstleister erforschen. Ein Beispiel hierfür ist das Projekt BiDiMoVe (Bidirektional, Multimodal, Vernetzt), das die Priorisierung des Öffentlichen Personennahverkehrs (ÖPNV) über den V2X-Nachrichtenaustausch in der Stadt Hamburg erforscht. Im Projekt SIRENE, das in Braunschweig ansässig ist, fokussieren sich die Wissenschaftlerinnen und Wissenschaftler auf die Bevorrechtigung von Einsatzfahrzeugen an Kreuzungen.

V2X-Technologie und Hintergründe

Die V2X-Technologie, auch bekannt als Car2X, entwickelte sich in den 2010er-Jahren mit dem Ziel, die Sicherheit und Effizienz des Straßenverkehrs durch vermehrte Kommunikation der Verkehrsteilnehmenden zu erhöhen. Der frühzeitige Austausch von Informationen unter ihnen, beispielsweise durch die Meldung einer Aquaplaning-Gefahrenstelle von Fahrzeug zu Fahrzeug, kann Unfälle verhin-



RHENUS CLOUD

Du liebst Technologie?

Du bist?

- Administrator
- DevOps Engineer
- .NET-Entwickler
- .NET-Architekt
- Projektmanager
- Produktmanager
- Product Owner

Für uns ist IT viel mehr als Technologie!

Du willst erfahren wieso?

Schreib uns direkt:



Nachrichtenformate

Der Zusatz „-em“ bzw. das eingeschobene „-e-“ bei den Nachrichtenformaten MAPEM, SPATEM, SREM und SSEM steht für „extended“ und grenzt die europäische Version der Nachrichtenformate von dem global verwendeten IEEE-Standard ab.

dern. Die Effizienzsteigerung geht mit der Vermeidung von Unfällen einher und zeichnet sich vor allem durch die Möglichkeit einer höheren Verkehrsdichte aus. Fahrzeuge, die Informationen über sich und ihre Umwelt austauschen, können zukünftig schneller auf Gefahren reagieren und daher auch enger hintereinander fahren. Dieser Anwendungsfall wird Platooning genannt.

Die V2X-Kommunikation basiert ursprünglich auf dem 2010 veröffentlichten Standard IEEE 802.11p. Der 802.11p-Standard ist eine Erweiterung der WLAN-Spezifikation, die innerhalb der IEEE-802.11er-Reihe festgehalten ist. Die Entwicklung und Erforschung der V2X-Technologie treiben hauptsächlich Europa, die USA und Japan voran. In Europa begleiten das European Telecommunication Standards Institute (ETSI) mit Sitz in Frankreich und das Car2Car Communication Consortium (C2C-CC) die Entwicklung. Der Standard ETSI ITS-G5 stellt die Umsetzung des 802.11p-Protokolls für den europäischen Raum dar. Im ITS-G5-Standard sind für die V2X-Kommunikation mehrere Kanäle im Bereich von 5,85 bis 5,925 GHz reserviert. In den USA haben Automobilhersteller und das Verkehrsministerium im Jahr 1999 75 MHz im 5,9-GHz-Frequenzband für die V2X-Kommunikation vorgesehen.

Das für die Kommunikation verwendete „Wireless Access in Vehicular Environments (WAVE)“-Protokoll basierte bis November 2020 ebenfalls auf IEEE 802.11p. Während es in Europa und den USA noch keine weitreichende Ausstattung mit V2X-fähiger Infrastruktur gab, konnte das zuständige Ministerium in Japan auf das bestehende Mautsystem zurückgreifen. Japan verwendet eine modifizierte Version des WAVE-Standards für eine Kommunikation im 5,8-GHz-Bereich. Seit 2017 treibt insbesondere der Mobilfunksektor auch die Kommunikation über LTE und 5G – genannt C-V2X (Cellular-V2X) – voran. Während in Deutschland die überwiegende Mehrheit der Forschungsprojekte den ETSI-ITS-G5-Standard verwendet, der auf der 802.11p-WLAN-Kommunikation basiert, setzen die USA und China zunehmend auf C-V2X.

Vehicle to Everything (V2X-System)

Nach dem ETSI-ITS-G5-Standard besteht ein V2X-Kommunikationssystem aus Roadside-Units (RSU) – streckenseitige

Kommunikationsinfrastruktur – auf der einen Seite und On-Board-Units (OBU) auf der anderen Seite. Fahrzeuge, die über V2X kommunizieren sollen, haben eine solche OBU an Bord. Intelligente Infrastruktur wie Ampeln oder Masten am Straßenrand sind mit einer RSU ausgestattet, um V2X-Nachrichten senden und empfangen zu können. Innerhalb der OBUs und RSUs ist der V2X-Stack der wichtigste Bestandteil. Der V2X-Stack ist in Abbildung 1 dargestellt und hat ähnlich wie das ISO/OSI-Modell mehrere Schichten.

In der untersten der vier Schichten befinden sich die Zugangstechnologien, die die in der Regel drahtlose Kommunikationstechnik zwischen den ITS-Stationen umfassen. Die darüberliegende Netzwerk- und Transportschicht steuert die Verbreitung der gesendeten Nachrichten. In der nachfolgenden ITS-Facilities-Schicht befindet sich die Verwaltung von Softwaremodulen und die Datensammlung sowie -aggregation für Applikationen. Die Applikationen sind in der obersten Schicht „ITS Applications“ zu finden. Hierunter fallen sowohl Anwendungen zur Interaktion im Fahrzeug (Human Machine Interface, HMI) als auch die Fahrzeugsensorik.

Europäische Nachrichtenformate

In der V2X-Kommunikation hat das ETSI in Anlehnung an den internationalen Standard europäische Nachrichtenformate für verschiedene Anwendungsfälle definiert. Die Cooperative Awareness Message (CAM) umfasst allgemeine Fahrzeuginformationen wie die Position, Geschwindigkeit und Beschleunigung. Die Fahrzeuge senden solche Cooperative Awareness Messages kontinuierlich mindestens ein Mal pro Sekunde aus, um Verkehrsteilnehmende im Umkreis von bis zu 500 Metern über sich zu informieren. Für die Warnung anderer vor Gefahrenstellen legt der ITS-G5-Standard die Verwendung einer Decentralized Environmental Notification Message (DENM) fest. Eine DENM enthält sowohl Informationen über das sendende Fahrzeug als auch zur Gefahrenstelle. Die Gefahrenstelle wird mittels CauseCode und SubCauseCode codiert. Der CauseCode 12 beispielsweise steht für Personen auf der Straße. Fahrzeuge können solche Gefahrenmeldungen im Display anzeigen, um Fahrende zu warnen.

Hat ein Fahrzeug eine DENM empfangen, kann es sie wiederum per Broadcast an andere Fahrzeuge in seiner Nähe weiterleiten. So erhalten auch Fahrzeuge in größerer Entfernung Meldungen. Ampeln kommunizieren ihren aktuellen Status über die Nachrichtenformate Signal Phase and Timing (SPATEM) und Map Topology, kurz MAPEM. Mit SPATEM teilen Ampeln den aktuellen und (abhängig vom Hintergrundsystem) gegebenenfalls auch den prognostizierten Signalzustand mit. Aus der Kombination von SPATEM und MAPEM können OBUs den Fahrenden verschiedene Services anbieten.

Simulationsumgebungen für Verkehrsszenarien – eine Auswahl

Simulationsframework	Primäres Ziel	Simulierte V2X-Kommunikation	OpenSCENARIO Standard
Virtual Test Drive (VTD)	Verkehrssimulation	Nein	Nein
CARLA Simulator	Fahrfunktionen erproben	Standardmäßig nein, vereinzelt in Arbeiten eingesetzt	Teilweise integriert durch CARLA ScenarioRunner
Simulation of Urban Mobility (SUMO)	Mikroskopische Verkehrssimulation	Standardmäßig nein, vereinzelt durch Verwendung von TraCI	Nein
Artery (basiert auf Veins)	Simulation von V2X-Kommunikation	Ja	Nein
dSPACE V2X Solution	Simulation von V2X-Kommunikation	Ja	Nein
dSPACE Scenario Generation Service	Verkehrsszenarien generieren	Nein	Ja
CarMaker	Fahrfunktionen erproben	Ja	Ja

Grüne Welle mit GLOSA

Ein Beispiel für einen zukünftigen Service ist die Anzeige der prognostizierten Ampelphase bei Erreichen der Haltelinie im Display. Eine Erweiterung dazu ist die sogenannte „Green Light Optimal Speed Advisory“, kurz GLOSA. Mit GLOSA erhalten Fahrende eine Geschwindigkeitsempfehlung für das Erreichen der nächsten Grünphase einer kommenden Ampel. Für eine Priorisierung von Fahrzeugen wie dem öffentlichen Personennahverkehr in BiDiMoVe oder Einsatzfahrzeugen im Projekt SIRENE sieht der Standard die Nachrichtentypen Signal Request Message (SREM) und Signal Status Message (SSEM) vor. Mit einer SREM fordert ein Fahrzeug eine Bevorrechtigung an einer Kreuzung an (Signal Request). Solch eine Nachricht enthält die für die Priorisierung benötigten Informationen wie den StationType, der die Art der ITS-Station beschreibt. Der StationType unterscheidet zwischen Fußgängern (Codierung = 1), Fahrradfahrern (2), Straßen- und Schienenfahrzeugen (3-11) und Roadside-Units (RSUs) (15). Die Infrastructure-to-Vehicle-Information (IVI) bildet unter anderem physische Verkehrsschilder ab, dabei können die Schilder sowohl statisch als auch dynamisch veränderbar sein.

V2X-Komponenten testen

In der Entwicklung von V2X-Kommunikationssystemen gehört das Testen zum Alltag. Wie Softwareentwickelnde wissen, setzt sich der Testprozess im Allgemeinen aus sogenannten Komponententests, anschließenden Integrationstests und zuletzt den Systemtests zusammen. Das Deutsche Zentrum für Luft- und Raumfahrt legt seinen Fokus auf Systemtests, die Wissenschaftler meist auf der Basis von Systemspezifikationen oder Use Cases erstellen. Die Durchführung der Systemtests erfolgt entweder in Laborumgebungen oder

im realen Umfeld, in dem das zu testende System am Ende zum Einsatz kommen soll. Während im Feldversuch insbesondere in frühen Entwicklungsstadien das Risiko und die Kosten eines Tests sehr hoch sind, lassen sich Labortests vergleichsweise günstig durchführen.

Zur Nachbildung des Verkehrsgeschehens rund um das System-Under-Test (SUT) benötigen Tester eine Simulationsumgebung, in die das SUT sich einbinden lässt. Solche Tests bezeichnen Wissenschaftler und Wissenschaftlerinnen als In-the-Loop-Tests. Vor der Simulation einzelner Use Cases im Verkehr müssen Entwickelnde die Verkehrsszenarien erst definieren und erstellen. Deshalb geht es hier im Artikel zunächst um die Frage, wie Entwicklerinnen und Entwickler Verkehrsszenarien definieren.

Wie OpenScenario Verkehr simuliert

In der szenariobasierten Verkehrssimulation stellt ein Verkehrsszenario eine temporäre Verkehrssituation dar, die aus den folgenden Kernelementen besteht:

1. Topologie der Straße
2. logische Verknüpfungen und Leitinfrastruktur
3. temporäre Modifikationen wie eine Baustelle (optional)
4. Dynamik der Verkehrsteilnehmenden
5. Umweltbedingungen
6. digitale Informationen (optional)

Als (Quasi-)Standard für die Definition von Verkehrsszenarien gilt das Dateiformat OpenSCENARIO in Verbindung mit OpenDRIVE, das der logischen Beschreibung von Straßennetzen und der Infrastruktur dient. OpenDRIVE bildet demnach die ersten drei Elemente einer temporären Verkehrssituation ab. Der OpenSCENARIO-Standard existiert seit 2018 und ergänzt die OpenDRIVE-Daten um dynamische Informationen. Beide

Ego-Fahrzeug

Als Ego-Fahrzeug bezeichnen Forschende das Fahrzeug, das in einer Simulation untersucht wird. Das Ego-Fahrzeug wird in Fahrsimulationen meist aus der Ich-Perspektive dargestellt, was Versuchsfahrenden eine dem echten Fahrerlebnis ähnliche Simulation ermöglicht.

Datenformate basieren auf dem XML-Format. In der nachfolgenden Tabelle sind einige zurzeit relevante Simulationsumgebungen für Verkehrsszenarien aufgeführt. Für jede Simulationsumgebung zeigt die Tabelle das primäre Ziel der Simulation. Zudem ist aufgeführt, ob sie den OpenSCENARIO-Standard unterstützt und ob sie V2X-Kommunikation simuliert. Zu sehen ist, dass die meisten Simulationsumgebungen sich jeweils bisher auf ein Thema fokussiert haben, nämlich auf die Simulation von V2X-Nachrichten oder auf die Integration von Verkehrsszenarien im OpenSCENARIO-Format. Bei einigen Simulationsumgebungen wie VTD oder SUMO ist die Integration des OpenSCENARIO-Standards jedoch für die Zukunft geplant.

Von OpenSCENARIO zur V2X-Nachricht

Ziel ist es nun, das OpenSCENARIO-Format mit der Simulation von V2X-Kommunikation zu verbinden. Genauer gesagt geht es um das Generieren von V2X-Nachrichten der Fahrzeuge, die in einem gemäß OpenSCENARIO definierten Verkehrsszenario enthalten sind. V2X-Nachrichten lassen sich beispielsweise generieren mit dem Open-Source-V2X-Stack Vanetza, der auf GitHub zu finden ist. Simuliert werden die Nachrichtenformate CAM und DENM.

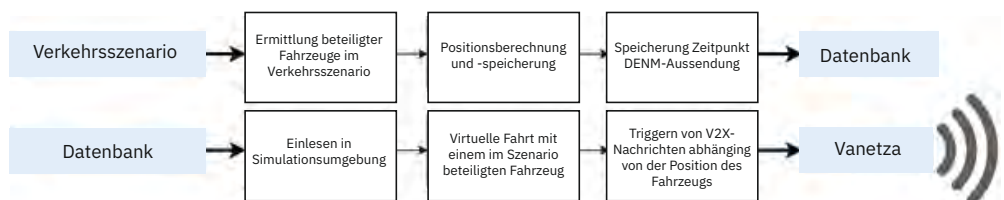
Der Ablauf der Simulation nach dem Erstellen eines Verkehrsszenarios ist in Abbildung 2 zu sehen. Zunächst ermittelt ein Skript sämtliche im Szenario enthaltenen Fahrzeuge. Als zusätzliche Information wird gespeichert, ob es sich um ein Ego-Fahrzeug oder ein Non-Ego-Fahrzeug handelt. In OpenSCENARIO sind die Routen der Fahrzeuge meist nicht als durchgängige Bewegungspfade (Trajektorien), sondern als Aktionen in Abhängigkeit zur Position oder Geschwindigkeit anderer Fahrzeuge definiert. Eine Abhängigkeit könnte beispielsweise sein: „Ego-Fahrzeug verringert Geschwindigkeit auf 50 km/h, wenn der Abstand

zu Fahrzeug X geringer ist als 10 Meter“. Ein weiteres Skript ermittelt im Anschluss aus den gegebenen Abhängigkeiten konkrete Bewegungstrajektorien für die enthaltenen Fahrzeuge. Die Speicherung der Trajektorien erfolgt als Wegpunktliste in einer Datenbank. In der Wegpunktliste sind für jede Sekunde im Szenario die Positionen der Fahrzeuge enthalten. Zur Simulation von DENMs, also Gefahrenwarnungen, existiert im OpenSCENARIO-Standard noch kein XML-Element. Zur Beschreibung des sendenden Fahrzeugs und der DENM-Inhalte können Entwickelnde jedoch unter anderem die im Standard mitgelieferten frei konfigurierbaren Aktionen (CustomCommandActions) verwenden. Die Informationen zu Aussendungszeitpunkt, Sender und Inhalt ermittelt ein Python-Skript und speichert sie in einer für diesen Zweck erstellten Datenbank.

V2X-Simulation mit Vanetza

Im zweiten Schritt, der zeitlich unabhängig vom Erstellen der Datenbank ist, generiert eine Simulationsumgebung im Zusammenspiel mit dem V2X-Stack Vanetza die zugehörigen V2X-Nachrichten. Das verwendete Simulationsframework ist das Simulation Environment for ERTMS Verification (SEFEV), das im Regelfall im Hardware-in-the-Loop-Labor im DLR zur Erprobung von Eisenbahn-Leit- und -Sicherungssystemen dient. ERTMS steht dabei für das European Rail Traffic Management System. Die Simulationsumgebung liest zunächst die Wegpunktliste und die DENMs aus der Datenbank ein. Anschließend startet die Simulation des Szenarios.

Hierfür simuliert das Framework das Ego-Fahrzeug auf einer Trajektorie – die Strecke wird quasi abgefahren. Abhängig von der Position des Ego-Fahrzeuges stößt das Simulationsframework die Generierung der CAMs der umliegenden Fahrzeuge an. Das Simulationsframework sendet nun sekundlich pro Non-Ego-Fahrzeug einen Trigger mit Informationen zum sendenden Fahrzeug und zur Position aus. Getriggert wird der Vanetza-V2X-Stack, der daraufhin mit Hilfe der erhaltenen Informationen eine CAM generiert. Das Erzeugen von DENMs läuft analog dazu ab: Ist in der Datenbank die Aussendung einer DENM von einem Non-Ego-Fahrzeug zum aktuellen Zeitpunkt enthalten, triggert das



>> Vom Verkehrsszenario zur V2X-Nachricht (Abb. 2)

Simulationsframework mit der Information den Vanetta-V2X-Stack, der daraufhin eine DENM erstellt.

Zur tatsächlichen Aussendung der V2X-Nachrichten über den 802.11p-Funkstandard ist der Vanetta-V2X-Stack in eine On-Board-Unit integriert. Die OBU (MK5 von Cohda Wireless) ist bereits mit Antennen und einem GPS-Modul ausgestattet. Zusammenfassend lässt sich mit dem beschriebenen Ablauf die V2X-Kommunikation mehrerer Fahrzeuge in einem Verkehrsszenario simulieren. Mit einer weiteren OBU eines anderen Herstellers ließ sich die Korrektheit und Kompatibilität der erstellten Nachrichten durch den Vanetta-V2X-Stack validieren.

Komplexe Verkehrsszenarien testen

V2X-Kommunikation, also das Vernetzen von Fahrzeugen mit ihrer Umgebung, stellt einen wesentlichen Baustein für das autonome Fahren dar. Schon heute erprobten Forschungseinrichtungen verschiedene Anwendungsfälle der V2X-Kommunikation – wie das Bevorrechtigen von Einsatzfahrzeugen an Kreuzungen – und gehen damit in Richtung zukünftiger Mobilität. Getestet werden V2X-Applikationen oder -Kommunikationskomponenten zunehmend in Laborumgebungen unter In-the-Loop-Bedingungen, in denen Verkehrsszenarien sich möglichst ganzheitlich simulieren lassen.

Zu einer ganzheitlichen Simulation gehört eine detaillierte Beschreibung des Verkehrsgeschehens, was mit dem OpenSCENARIO-Standard erreichbar ist. Darüber hinaus wollen Forschende auch die zunehmende V2X-Kommunikation im Verkehr simulieren, wofür sich neben anderen Ansätzen auch die hier beschriebene Kombination des OpenSCENARIO-Standards mit der Simulation von V2X-Nachrichten eignet.

Quellen

Weiterführende Hinweise finden sich unter ix.de/zvsw.



Katharina Hartmann

arbeitet beim Deutschen Zentrum für Luft- und Raumfahrt in Braunschweig. Im Institut für Verkehrssystemtechnik beteiligt sie sich an der Konzeption und Entwicklung von Tests für Kommunikationssysteme im Straßen- und Schienenverkehr. Ihr aktueller Schwerpunkt liegt auf der V2X-Kommunikation, die auch Teil ihrer Bachelorarbeit war.

Die ProMinent Unternehmensgruppe mit Hauptsitz in Heidelberg ist zuverlässiger Lösungsanbieter für die Wasseraufbereitung und in der Fluid-Dosiertechnik. Ein Global Player mit rund 2.700 Mitarbeitern und rund 50 Niederlassungen weltweit. Wir suchen ab sofort:

ProMinent®

Elektronik Entwickler

(m/w/d)



Sie sind Elektronik Entwickler und wollen in einem starken Team Leistungselektronik und Mixed Signal Schaltungen bis an die Grenze des Machbaren entwickeln.

Unsere Anforderungen:

- Studium der Elektrotechnik
- Fundierte Kenntnisse der Leistungselektronik und/oder Mixed Signal Elektronik
- Schnelle Auffassungsgabe

Full Stack Entwickler

(m/w/d)



Sie entwickeln die nächste Evolutionsstufe unserer erfolgreichen, webbasierten IIOT Plattform DULCONNEX und koordinieren das Developer Team fachlich und methodisch.

Unsere Anforderungen:

- PHP 7 ■ Java ■ Linux shells scripting ■ Javascript
- PostgreSQL

DevOps-Engineer

(m/w/d)



Sie verantworten die Weiterentwicklung unserer CI-Plattform für embedded Software sowie für unsere App- und Cloud-Entwicklung. Zusammen mit den jeweiligen Entwicklern sorgen Sie dafür, dass höchste Qualitätsstandards automatisiert geprüft und eingehalten werden.

Unsere Anforderungen:

- Jenkins ■ CMake ■ Docker ■ Windows Powershell
- Windows und Linux

Produkt Software Entwickler

(m/w/d)



Sie entwickeln modulare Embedded Software für unsere Regler, Pumpen und Sensoren und nutzen automatisierte Test und Code-Reviews, um hohe Qualitätsstandards sicherzustellen.

Unsere Anforderungen:

- C++ 14 ■ Embedded Softwareentwicklung ■ Design Patterns
- Unit Tests, GoogleTest

Was wir Ihnen bieten:

- Eine Perspektive in einem globalen Unternehmen - denn wir sind in über 100 Ländern aktiv.
- Die Nachhaltigkeit und Solidität eines Familienunternehmens.
- Gestaltungsfreiräume im Mittelstand.
- Interessante Projekte und Freiräume bei der Umsetzung.

Fühlen Sie sich herausgefordert?

Dann bewerben Sie sich gerne mit Ihren vollständigen Bewerbungsunterlagen auf unserer Karriereseite www.prominent.com.

ProMinent GmbH

Im Schuhmachergewann 5-11 | 69123 Heidelberg | Germany

> Was einen guten Datensatz ausmacht

Thordis Taag

Was ist ein guter Datensatz? Für Machine Learning ist das grundlegend. Insbesondere für die Frage, warum Maschinen (noch) keine besseren Entscheidungen treffen.



Autonome Fahrzeuge, die im Graupelschauer stehen bleiben. Sprachgesteuerte Smart-Home-Assistenten, die nur männliche Stimmen verstehen. Trotz Innovationsschub und verfeinerter Technik steht das maschinelle Lernen (ML) oft vor unerwarteten Herausforderungen. Es wird viel darüber geschrieben und diskutiert, was schlechte Datensätze im Machine Learning bewirken können. Doch was ist eigentlich ein guter Datensatz?

Daniel Kondermann stellte diese pragmatische Frage schon vor über zehn Jahren. Der Forscher im Bereich Computervision an der Universität Heidelberg hat sich auf die Qualitätssicherung von Datensätzen spezialisiert und ein System entwickelt, das hochwertige Datensätze generiert. Im Gespräch mit der Autorin erläutert er auch, wie dieses System einen Beitrag dazu leisten könnte, eine ethisch verantwortungsvolle, sicherere und transparente Anwendung von KI- und ML-Systemen zu erreichen. Um die Relevanz dieser Frage zu verstehen, gilt es zuerst genauer zu betrachten, was das maschinelle Lernen derzeit hemmt.

Warum gerät das Machine Learning ins Stocken?

Mit der Entwicklung neuronaler Netzwerke ist der gordische Knoten im Bereich der künstlichen Intelligenz (KI) geplatzt. Plötzlich schien alles möglich: Es ist eine neue globale Industrie entstanden und namhafte Unternehmen überschlagen sich

damit, die KI in immer innovativeren Machine-Learning-Projekten anzuwenden: von selbstfahrenden Fahrzeugen zu Smart-Home-Assistenten, vom Spamfilter zur Übersetzungssoftware. Mithilfe des maschinellen Lernens lassen sich Börsencharts analysieren und Herzrhythmen auf Unregelmäßigkeiten überprüfen. Und doch stagniert der Machine-Learning-Markt. Nach der Einschätzung des IT-Marktforschungsunternehmens Gartner schaffen es nur 53 Prozent aller KI-Prototypen in die Produktion, und das ist wohl noch eine schmeichelhafte Einschätzung: Laut dem IT-Portal Venture Beat erreichen die allermeisten Data-Science-Projekte nicht die Produktionsreife (mehr dazu unter ix.de/zc5d).

Woran liegt das? Wieso bleibt Waymos selbstfahrendes Taxi ratlos vor ein paar Baustellenhütchen stehen? Wieso kommt noch kein selbstfahrendes Fahrzeug mit Regen, Schnee oder Schneeregen klar? Nachdem viele Gedanken und Energie in Code und Modelle gesteckt wurden, rückt nun der Datensatz in den Fokus.

Auf einem Workshop im Jahr 2013 war Daniel Kondermann mit seiner Meinung, man müsse auch ein Augenmerk auf gute Trainingsdatensätze haben, noch ein Außenseiter. Wozu, wenn man doch auch einfach die Machine-Learning-Methoden optimieren könne? Zum Zeitpunkt dieses Workshops steckten alle Methoden der Datenerzeugung für die Computervision noch in den Kinderschuhen. Ob spezielle Messtechnik, die Computergrafiksimulation oder das Annotieren – wobei unter anderem das Annotieren von Texten

schon in der Branche angekommen war, wie das Beispiel der Sprachlern-App Duolingo zeigt. Inzwischen sind die ML-Methoden präzise, doch die ebenso guten Ergebnisse lassen auf sich warten. So wird in der Branche erst jetzt klar, dass man sich auch über gute Datensätze Gedanken machen muss. Firmen entwickeln zunehmend Systeme, die Datensätze inkrementell verbessern – beispielsweise mithilfe automatischer und manueller Qualitätssicherungsschritte. Um das zu verdeutlichen, ist zunächst ein Verständnis davon notwendig, wie das maschinelle Lernen genau funktioniert.

Wie lernt die Maschine?

Oft werden die Begriffe KI und ML miteinander vermischt, doch sie beschreiben unterschiedliche Sachverhalte. Künstliche Intelligenz ist die Logik hinter der Mustererkennung des maschinellen Lernens. Hierbei werden Axiome auf eine Faktenbasis angewandt und führen zu einem Ergebnis. Zum Beispiel: Ein Baum steht still, ein Passant bewegt sich, ein Baum trägt Blätter, ein Passant Kleidung. Wenn etwas Kleidung trägt und sich bewegt, muss es ein Passant sein.

Diese Logik muss vollständig und widerspruchsfrei sein, damit sich daraus etwas folgern lässt. Ohne diese Vollständigkeit hat das System eine Wissenslücke. Ist auch nur ein Widerspruch in der Faktenbasis enthalten, kann man daraus beliebig viele falsche Dinge folgern. Die Daten dürfen also keine Fehler enthalten. Um das im Bereich der Computervision sicherzustellen, hat Daniel Kondermann die RAD-Methode entwickelt (Representativeness, Accuracy, Difficulty). Die Computervision beschäftigt sich mit visualisier-


in a nutshell

- > Die Mustererkennung des maschinellen Lernens (ML) ist eine Sparte der Künstlichen Intelligenz (KI).
- > Sie basiert auf Beispielen, die dem Computer gezeigt werden, um etwas zu lernen (Trainingsdatensatz).
- > Diese Daten müssen eine ausreichende Qualität haben, damit der Computer eine Chance hat, daraus etwas abzuleiten.
- > Mit der RAD-Methode (Representativeness, Accuracy, Difficulty) lässt sich sicherstellen, dass ein Datensatz keine Widersprüche in der Faktenbasis, menschlichen Bias oder Fehlinterpretationen enthält.
- > RAD verbessert die Genauigkeit von Machine Learning und lässt sich zur Qualitätskontrolle heranziehen.

baren Daten wie 3D-Kartenmodellen, Fotografien und MRT-Scans, während es bei anderen Systemen des maschinellen Lernens wie Sprachanwendungen oder EEG-Frequenzanalysen um gänzlich andere Datensätze geht. Trotzdem lässt sich das Werkzeug RAD auch in anderen Bereichen des Machine Learning zur Qualitätskontrolle heranziehen.

„This dataset is RAD!“

Damit der Datensatz also umgangssprachlich gesprochen ‚geil‘ wird, soll die RAD-Methode dabei helfen, mögliche Quellen von Annotationsfehlern aufzuspüren, um einen guten Da-



Viel mehr als
ein*e Software-
Entwickler*in.

Unsere Entwickler*innen sind Expert*innen für Frontend, technische Basis, Datenbankdesign, Testautomatisierung, Architektur, Frameworks, Algorithmen oder User Experience. Aber viel wichtiger: Sie sind Klavierspieler, Eltern, Zocker, Bäume-Ausreißerinnen, Bäume-Pflanzer, Zahlenverliebte, Leseratten, Diskutierfreudige, Yogis, Harmoniebedürftige, Frühaufsteher und so vieles mehr. Das macht uns stark.

Werden Sie Teil unserer Vielfalt.



Weitere Infos unter www.parclT.de

tensatz zu erhalten. Qualität lässt sich stets nur über mehrere Faktoren definieren. Ein Laptop etwa kann eine große Festplatte haben, aber einen zu kleinen Bildschirm. Beides sind Performanzmetriken und keine von beiden kann allein der richtige Faktor sein, sondern ist abhängig vom jeweiligen Anwendungsfall. Deshalb konzentriert sich RAD auf drei Gruppen von Metriken, wobei jede davon auf den Anwendungsfall zuzuschneiden ist.

Ein guter Datensatz muss repräsentativ sein: Representativeness

Zeigt man einem Netzwerk nur männliche Gesichter mit weißer Haut, erkennt es Frauengesichter mit schwarzer Haut nicht, wie die MIT-Forscherin Joy Buolamwini in ihrer Netflix-Dokumentation Coded Bias eindrucksvoll nachwies. Auch die Geometrie der Bilddaten spielt eine Rolle. Werden Gesichter immer nur von vorne gezeigt, erkennt das Netzwerk keine Profilansichten. Enthält das Trainingsmodell nur

Personen im Sommer im US-Staat Kalifornien, erkennt es Personen im dicken Wintermantel in Norwegen nicht.

Die Art der Fotoaufnahme spielt ebenfalls eine Rolle, Schlaglicht beispielsweise erzeugt andere Pixel als gleichmäßig ausgeleuchtete Aufnahmen. Ist die aufgenommene Person in weiter Entfernung, sind die Pixel anders als aus großer Nähe. Hinzu kommen Personen, die teilweise von anderen Menschen oder Objekten verdeckt sind. Wichtig ist also die Überlegung, was das System bewerkstelligen können soll und was nicht. Sollen nur Studierende in Heidelberg erkannt werden, wird die Repräsentativität anders definiert als für ein Altenheim in Paris.

Bereits für den ersten Faktor zur Qualitätssicherung sind also so viele Dimensionen zu beachten, dass man sich durch Versuch und Irrtum an ein zufriedenstellendes Ergebnis herantasten muss.

Klassifikationsgenauigkeit und Ambiguität: Accuracy

In diese Kategorie fallen falsch annotierte Bilddaten, bei denen ein Zielobjekt vergessen oder ein falsches Zielobjekt annotiert wurde. Es sollten beispielsweise Boxen um Personen gemalt werden, doch manche Personen wurden vergessen (Falschnegative). Oder eine Box wurde um einen Hund statt um eine Person gezeichnet (Falschpositive). Die möglichen Fehler hinsichtlich der Genauigkeit lassen sich erneut in drei Unterpunkte gliedern. Die Erkennungsgenauigkeit besagt, ob das Objekt erkannt wurde oder nicht.

Die Klassifikationsgenauigkeit ist ein Maß dafür, wie oft das richtige Objekt getroffen wurde. Wenn sich in einer Straßenszene zwischen zwanzig und sechzig verschiedene Klassen von Objekten befinden wie Autos, Menschen, Vegetation, Fahrbahnmarkierungen, Schilder und Ampeln, dann kann es passieren, dass Menschen oder Maschinen aus Versehen eine falsche Antwort geben, zu welcher Klasse das fragliche Objekt gehört.

Die geometrische Genauigkeit stößt zum Beispiel auf Schwierigkeiten, wenn es darum geht, eine Kontur um eine Person zu ziehen. Spätestens bei dem Versuch, den Lockenkopf von Thomas Gottschalk zu umzeichnen, würden fünf Personen fünf unterschiedliche Ergebnisse liefern. Auch bei der gröberen Umzeichnung kann beispielsweise ein Gebüsch hinter einer Person als Hochsteckfrisur fehlinterpretiert werden. Dass es verschiedene Möglichkeiten gibt, die Aufgabe zu lösen, bringt einen gewissen Grad an Ambiguität in das Netzwerk.

Nicht zuletzt gibt es noch die Unsicherheit: Manchmal lässt sich die genaue Kontur nicht definieren, weil sie im Bild beispielsweise aufgrund von Bewegungsunschärfe schlicht nicht

Dr. Daniel Kondermann – Spezialist für Datensätze



Daniel Kondermann forscht seit 2009 an der Frage, was ein guter Datensatz für den Bereich Computervision ist. Im Jahr 2016 hat er sich auf dem Gebiet habilitiert. In dieser Zeit hat er mit verschiedenen Teams zahlreiche Publikationen zu guten Datensätzen veröffentlicht. Im Rahmen sei-

nes ersten Start-ups, der Pallas Ludens GmbH, hat er außerdem maßgeblich zu dem Datensatz „Cityscapes“ und zu einer Erweiterung des Datensatzes „KITTI“ des Karlsruhe Institute of Technology (KIT) beigetragen: The KITTI Benchmark Suite, Semantic Segmentation Evaluation.

Seit 2019 unterstützen Kondermann und sein Team des Start-ups Quality Match Unternehmen darin, die richtigen Fragen zu stellen, um gute Datensatz-Beispiele zu finden. Ziel ist es, durch das Herausfiltern von Fehlern, Inkonsistenzen und Doppeldeutigkeiten den Datensatz so repräsentativ, akkurat und diffizil wie möglich – eben RAD – zu machen. Sein Motto lautet: „Wenn Daten das neue Öl sind, ist Quality Match die Raffinerie.“

Was es mit der RAD-Methode (Representativeness, Accuracy, Difficulty) auf sich hat, erläutert der Artikel.

zu sehen ist. Dann sollte man beim Annotieren nicht raten, sondern einfach ganz klar sagen, dass die Bandbreite möglicher Antworten größer ist. So kann das auch die Maschine lernen und verstehen und muss selbst nicht raten.

Vorstellungen davon, wie die Welt (nicht) funktioniert: Difficulty

Die Regeln zum Annotieren basieren auf einer Idee, wie die Welt der Fragenstellenden, aber auch der Annotierenden funktioniert. Zum Beispiel folgte auf die Überlegung beim Aufbau des Modells, dass ein Auto Räder hat, die Aufgabe an die Annotatoren, alles auszuwählen, was Räder hat. Nun gab es ein Bild, das ein gerade abhebendes Flugzeug zeigte. Fünf von zehn Personen wählten es aus und die anderen fünf nicht.

Warum? Das Flugzeug hatte bereits das Fahrwerk eingefahren. Die Räder waren gerade nicht zu sehen, dennoch hatte es ja Räder. Hier sind nicht die Daten mehrdeutig, sondern die Taxonomie ist unstimmig. Die Entscheidungsgrenze ist nicht klar definiert, weil beim Aufbau der Taxonomie nicht weit genug gedacht wurde.

„Wir haben uns dann gefragt, warum überhaupt Flugzeuge im Datensatz waren“, berichtet Kondermann, „wenn das Ziel der Anwendung ist, dass ein selbstfahrendes Fahrzeug andere Fahrzeuge erkennt, braucht es ein Flugzeug gar nicht zu erkennen. Ganz egal, ob es Räder hat oder nicht.“ Der Datensatz muss konsistent sein, und zwar so weit, dass er sich nicht selbst widerspricht. Man möchte Autos erkennen, sagt aber: „Finde Dinge mit Rädern“. Natürlich werden dann viel mehr Objekte als nur die Autos gefunden, und schon gibt es einen

Widerspruch im Datensatz. Hier treffen verschiedene menschliche Annotatoren ganz unterschiedliche Einordnungen. Sind Fahrradfahrende Fahrzeuge und wenn ja, welcher Teil des Bildes? Nur das Fahrrad? Oder die Komposition aus Fahrrad und Person? Wie ist es mit Personen im Rollstuhl oder auf einem Skateboard: Sind sie noch Fahrzeuge oder schon Passanten?

Datensätze müssen differenziert sein

Ein Datensatz muss dennoch „difficult“ sein, weil man an die Entscheidungsgrenze herankommen möchte. Zeigt man nur verschiedene Bilder eines Hyundai, kann das Netzwerk bereits einen Mercedes nicht mehr als Auto erkennen. Es muss das ganze Spektrum an Möglichkeiten aufgezeigt werden. Das Ziel ist es, so nah wie möglich an den sogenannten Edge Case heranzukommen, denn eine ganz klare Grenze gibt es nie.

Sind Personen auf Skateboards nun Fahrzeuge oder nicht? Hier muss gut überlegt und definiert werden, was damit bezweckt und was eigentlich benötigt wird. Was muss das selbstfahrende Fahrzeug tun, wenn es das Objekt sieht? Im besten Fall wohl bremsen. Dafür ist es egal, wie die Skateboardfahrenden gelabelt werden, Hauptsache ist, einen Unfall zu verhindern.

So wird deutlich, wie mit Hilfe der RAD-Methode auch über die Computervision hinaus ‚gute‘ Datensätze für das maschinelle Lernen kreiert werden können. Doch lässt sich der Gedanke auch weiterspinnen? Wenn die künstliche Intelligenz neuronale Netzwerke nachbildet, um aus einer Fülle an

Cloud-Apps auf Knopfdruck

Quizfrage: Wie erstellt man eine Cloud-App? Ganz einfach: Man nimmt eine On-Premise-App und verschiebt sie in die Cloud! Wenn die Realität nur so simpel wäre. Denn auch wenn diese Antwort per se nicht falsch ist, liegen doch Welten zwischen einer monolithischen App in der Cloud und einer echten Cloud-nativen Anwendung auf Basis wiederverwendbarer Container oder Microservices.

Eine Applikation, die von Grund auf für die Cloud erstellt werden soll, erfordert ein komplett neues Entwicklungsparadigma: Die Architektur muss völlig anders designet sein. Zudem müssen neue, umfassende CI/CD-Services implementiert werden. Es gilt, dutzende komplexer, meist disparater Cloud-Dienste unter einen Hut zu bekommen. Doch anders lassen sich die viel beschworenen Vorteile der Cloud – von flexibler Skalierbarkeit auf Millionen Anwender bis hin zur weltweiten Hochverfügbarkeit – kaum erreichen.

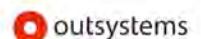
Effiziente Cloud-Entwicklung mit Low-Code

Die Implementierung einer Cloud-nativen Infrastruktur nimmt auf klassischem Wege nicht selten Monate in Anspruch – bevor die Anwendungsentwicklung überhaupt beginnen kann. Doch das muss

nicht sein. Abhilfe verspricht unter anderem die Low-Code-Technologie, die Entwickler-Profis in ihrer Kernkompetenz unterstützt und von zeitraubenden Routineaufgaben beim manuellen Coden entlastet. Diese Vorteile bietet die OutSystems-Plattform nun auch für die effiziente Entwicklung von Cloud-nativen Anwendungen. Durch sie erhalten Entwickler eine moderne Cloud-Runtime – ganz ohne monatelange Vorarbeit. Der eingesparte Aufwand erlaubt es den Entwicklern, sich ganz auf den Kern ihrer Arbeit zu konzentrieren: die Entwicklung moderner Apps, welche die Vorteile der Cloud maximal ausschöpfen.

OutSystems Germany GmbH

Tower 185, Excellent Business Center
Friedrich-Ebert-Anlage 35-37 | 60327 Frankfurt
www.outsystems.de



Datensätzen Bedeutung zu abstrahieren, lässt das auch Rückschlüsse darauf zu wie ‚das Original‘ – der Mensch – seinem Umfeld Bedeutung beimisst? Lernt nicht auch der Mensch unter anderem basierend auf Daten und trifft Entscheidungen.

Wie lernt der Mensch?

Eine KI ist eigentlich nur eine Entscheidungsmaschine, sie kann sehr schnell viele Entscheidungen treffen. Auch der Mensch trifft sehr schnell viele Entscheidungen, einen großen Teil davon unbewusst. Bleibt man beim Beispiel Fahren, so haben Menschen gelernt (und gesetzlich abgesichert), dass das Fahren bei rotem Ampellicht verboten ist. Ist ihnen bewusst, dass sie an jeder Ampel eine Entscheidung treffen, ob sie (weiter)fahren oder (an)halten? Auch hier muss die KI sicherstellen, dass keine „falschen“ Entscheidungen getroffen werden.

Kleiner Ausflug in die Neurowissenschaften

Es ist schwierig, einem Netzwerk auf Basis der Mustererkennung Mathematik beizubringen. Menschen können hierfür von ihrem präfrontalen Cortex Gebrauch machen. Wohl gemerkt, sie können. Gerade, wenn es darum geht, etwas zu bewerten oder Entscheidungen zu treffen, spielt dieser Bereich des Gehirns eine weitaus geringere Bedeutung, als es das gängige ökonomische Prinzip des Homo oeconomicus gerne glauben machen möchte. Es schaltet sich ein Bereich aus einem hinteren Teil des Gehirns dazu: das limbische System, das unter anderem für Gefühle zuständig ist (mehr dazu lässt sich dem Buch von Daniel Kahnemann entnehmen, „Schnelles Denken, langsames Denken“.) Viele schnelle Entscheidungen sind in diesem Teil sogar besser aufgehoben. Ein Mensch berechnet nicht schnell mal eben Aufprallwinkel und Bremsweg eines auf ihn zurasenden Autos. Er reagiert intuitiv, bleibt stehen oder weicht aus. Dieses Bauchgefühl, die Intuition, ist neurowissenschaftlich längst entzaubert. Tatsächlich handelt es sich um einen unbewussten, schnellen Abgleich von Dingen, die man schon gesehen oder erlebt hat. Genau wie bei den neuronalen Netzwerken im Machine Learning.

Mindset versus Dataset

Der Mensch ist also entsprechend seiner persönlichen Erfahrungen und Glaubenssätze ebenso voreingenommen wie die Maschine. Vor 15 Jahren hat die Autorin dieses Artikels ihrem Gehirn regelrecht dabei zuschauen können, wie es sich durch Mustererkennung an eine neue Situation angepasst

hat. Sie machte damals ein Praktikum am Goethe Institut in Südindien. Dort hatte sie drei Arbeitskolleginnen, die in etwa gleich alt und gleich groß waren. Es ist ihr in den ersten Tagen nicht gelungen, die drei auseinanderzuhalten. Die eigenen Parameter hatten einfach nicht gestimmt. Sie sahen für sie eben alle gleichermaßen „indisch“ aus – ein Parameter, der in Mitteleuropa meist völlig ausgereicht hätte, um eine indische Kollegin zu identifizieren.

Es folgte eine zweiwöchige Hospitation in den Schulklassen. In dieser Zeit sah sie jeden Tag viele unterschiedliche Gesichter und konnte richtiggehend merken, wie sich ihre Wahrnehmung dem neuen Umfeld angepasst hat. An einem Tag gegen Ende der zwei Wochen kam sie ans Institut und sah nur noch individuelle Personen und nicht mehr „viele Inderinnen und Inder“. Am Ende des Praktikums war für sie nicht mehr nachvollziehbar, wie sie die drei Arbeitskolleginnen je miteinander hatte verwechseln können.

„It's the dataset, stupid!“ – Fachliteratur

Eine Auswahl der bekanntesten Publikationen von Dr. Daniel Kondermann mit wechselnden Co-Autoren zur vertiefenden Lektüre:

- > **A dataset and evaluation methodology for depth estimation on 4D light fields**; in: Asian Conference on Computer Vision 2016; Springer (2016), S. 19-34
- > **Can masses of non-experts train highly accurate image classifiers?**; in: International conference on medical image computing and computer-assisted intervention 2014; Springer (2014), 438-445
- > **The hci benchmark suite: Stereo and flow ground truth with uncertainties for urban autonomous driving**; in: Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition 2016; Springer (2016), S. 19-28
- > **Stereo ground truths with error bars**; in: Asian Conference on Computer Vision 2014; Springer (2014), S. 595-610
- > **Geometry-aware performance analysis of stereo algorithms**; in: Proceedings of the IEEE International Conference on Computer Vision 2015; Springer (2015), S. 2120-2128

Links zu den PDF-Versionen der Artikel finden sich unter ix.de/zc5d.

So muss ein kultureller Bias keine Grundüberzeugung sein, die man sich im logischen präfrontalen Cortex bewusst zusammengereimt hat, sondern kann oft auch Folge einer Wissenslücke im limbischen System sein. Dann war der Datensatz der persönlichen Erfahrung nicht repräsentativ genug.

Interessant werden diese Überlegungen, wenn man nun fragt, wer wohl die besseren Entscheidungen trifft, Mensch oder Maschine? Selbstredend kann eine solch philosophische Frage nicht abschließend beantwortet werden, sie führt jedoch zur Frage der Ethik, die im Zusammenhang mit der immer größer werdenden Bedeutung von KI- und ML-Systemen nicht außer Acht bleiben darf.

Die ethische Dimension der künstlichen Intelligenz

In dem Maße, in dem KI- und ML-Systeme zunehmend Anwendung in sämtlichen Lebensbereichen der Menschen finden, werden Rufe laut nach einem sensiblen, verantwortungsbewussten Umgang mit dieser Technologie. Die Europäische Union hat im Frühjahr 2021 einen Entwurf zur „Regulation of the European Parliament and of the Council laying down harmonised rules on artificial intelligence“ veröffentlicht. Das

Papier erkennt die sozialen und ökonomischen Vorzüge der neuen Technologie an, warnt aber auch vor den Risiken, die die Systeme für Individuen und die Gesellschaft bergen. Gerade Grund- und Freiheitsrechte bedürfen eines besonderen Schutzes, so wird explizit vor diskriminierenden Algorithmen gewarnt.

Auch zu Trainingsdatensätzen äußert sich das Papier. Diese sollen in Hinblick auf das beabsichtigte Ziel des Systems ausreichend relevant, repräsentativ, fehlerfrei und vollständig sein. Sie sollen sich mit den Charakteristika des jeweiligen geografischen, funktionalen oder Verhaltenskontextes auseinandersetzen, um eine Diskriminierung durch in das KI-System einprogrammierten oder antrainierten Bias zu verhindern. Dies umfasst auch das „Weiterlernen“ von KI-Systemen nach dem Launch.

Machtfrage: Wem gehört der Code?

Joy Buolamwini weist in ihrem Film „Coded Bias“ auf ein Machtungleichgewicht hin, das in der Frage liege, wem der Code gehöre. Schließlich könne eine von der Kreditkartenfirma abgewiesene Person nicht ihre eigene KI gegen die Firma anwenden. Doch vielleicht sollte sie genau das können. Wenn



XITASO 

„Bei XITASO dürfen alle sein, wie sie sind: Laut und leise, neugierig und kreativ, ehrlich und mutig. Alle haben ihren Platz und fühlen sich wohl. Und das ist so unglaublich wertvoll.“

Bettina, UI/UX-Expertin

Experte für **High-End Software Engineering**.
Authentisch. Begeistert. Anspruchsvoll.

Augsburg | Ingolstadt | Krumbach | Leipzig | Madrid | Münster

www.xitaso.com/karriere

KI und ML die Gesellschaft in diesem großen Maße berühren, sollte es einen gesellschaftlichen Konsens über Grundzüge und Grenzen von Code und Datensatz geben.

Fraglich ist, wie ein solcher Konsens in einer globalisierten Wirtschaftswelt aussehen kann. Die Frage ist, ob solche Regulierungen den europäischen KI-Markt dämpfen oder ein Zeichen setzen. Auch auf dem nordamerikanischen Markt ist der Ruf nach einem ethischeren Umgang mit KI-Systemen längst laut geworden, wie John Thornhill in den *Financial Times* schreibt (Artikel „Artificial intelligence bias can be countered, if not erased“).

„Mit der zunehmenden Integration von KI in die Gesellschaft sind Chancenverlust, wirtschaftlicher Schaden und soziale Stigmatisierung die Hauptrisiken einer unethischen Anwendung von KI-Technologien“, schreibt Isabel Bär in einem Artikel bei *heise Developer* über Ethik und KI (siehe auch ihren Artikel in diesem Heft, „Warum Metriken nicht ausreichen“ auf S. 38). Das gilt allerdings auch für politische Systeme, in denen diese Entscheidungen Menschen überlassen sind. Wie also könnte man ML-Systeme sicherer machen? Demokratische Systeme schützen sich durch „Checks and Balances“ (Gewaltenteilung) vor Machtmissbrauch, indem sie Entscheidungen auf eine Vielzahl von Personen auf verschiedenen Ebenen verteilen, die sich gegenseitig kontrollieren und im Gleichgewicht halten.

Es stellt sich die Frage, ob es ausreichend sein kann, die Compliance-Grundsätze großer Unternehmen in die KI einzuspeisen, um dem Ganzen das ‚Ethik-Siegel‘ aufzudrücken. Müsste es nicht erstmal einen gesellschaftlichen Konsens über Ziel und Nutzen der Anwendung, über Fairness und Gerechtigkeit geben, der auch marginalisierte Gruppen miteinbezieht? Im Idealfall ergänzt durch einen RAD-Datensatz und so viele und diverse Annotatoren und Annotatorinnen, dass eine statistische Signifikanz erreicht wird. Wer KI und ML entwickelt, sollte dies schon bei der Fragestellung berücksichtigen und fortlaufend überprüfen.

Das menschliche Gehirn erstellt aus einer komplexen sozialen Welt heraus eigenständig und unbewusst abstrakte Kategorien, schreibt Lisa Feldman Barrett in einer aktuellen Ausgabe der *MIT Technology Review*. „Die ‚Blackbox Mensch‘ ist wahrscheinlich schwerer zu entschlüsseln als die Maschine“, äußert Daniel Kondermann seine Überlegungen, „trifft der Computer regelmäßig falsche Entscheidungen, ist das nachvollziehbar und kontrollierbar. Wer ist verantwortlich für die Fehlentscheidungen von Politikern und Politikerinnen? Deren Familie, ihr soziales Umfeld? Bei KI-Systemen sind es die vielen Menschen, die sie entwickelt, programmiert und trainiert haben.“

Die EU ist schon auf dem richtigen Weg, wenn sie Prozesse fordert, die sicherstellen, dass diese Entscheidung reflektiert

getroffen wird. Das geschieht unter anderem da, wo ein guter Datensatz erzeugt wird.

RAD: Methode zum sicheren Umgang mit KI/ML-Systemen

Machine Learning ist eine technisch ausgefeilte Methode, die viel Gutes, aber auch viel Leid bewirken kann. Die RAD-Methode verbessert die Genauigkeit von ML, indem sie menschlichen Bias, Fehlinterpretationen und unvollständige Datenannotationen ausfindig macht und eliminiert.

Je mehr soziale Auswirkungen eine Applikation hat, desto wichtiger ist es, zu hinterfragen, nach wessen Weltbild sämtliche Variablen definiert werden. Wer legt fest, wonach entschieden wird, ob jemand kreditwürdig ist oder nicht, ob jemand ein Mann oder eine Frau ist, ob jemand von einer unbemannten Drohne erschossen wird oder nicht?

Ein neuronales Netzwerk kann sehr schnell viele falsche Entscheidungen treffen. Bei einem verantwortungsbewussten Umgang mit dem System ist zu prüfen, ob es falsche Entscheidungen trifft. Im Anschluss lässt sich mithilfe eines verbesserten Datensatzes variieren, was es lernen soll. So können Netzwerke entstehen, die für bestimmte Domänen vielleicht sogar besser sind als ihr menschliches Äquivalent, wie es beim autonomen Fahren je nach Blickwinkel schon der Fall oder zumindest angestrebt ist. Vertiefend behandelt die Frage etwa ein Beitrag des Portals *The Gradient*, „Are Self-Driving Cars Really Safer Than Human Drivers?“. Man kann ein Fahrzeug Abertausende von Kilometern fahren lassen, bis es praktisch keine Fehler mehr macht, und erst dann gibt man es in Produktion. Ein Mensch setzt sich nach Fahrschule und Prüfung als Fahranfänger hinter das Lenkrad.

Quellen

Weiterführendes Material und alle Quellen stehen unter ix.de/zc5d bereit.



Thordis Taag

ist freie Texterin und Lektorin mit einem beruflichen Hintergrund in der Presse- und Öffentlichkeitsarbeit in Wirtschaft und Politik. Sie hat einen Master-Abschluss in Politikwissenschaften und Ethnologie.

CIAM: Der Blick aus der Kundenbrille ist zentral

Für die große Mehrheit steht die IT im Fokus von CIAM-Projekten. Neue Applikationen werden aus bestehenden IAM-Infrastrukturen entwickelt. Die Kundensicht gerät dabei leicht ins Hintertreffen. Das muss sich ändern. Denn eine gute CIAM-Strategie hängt vor allem von der Zusammenarbeit zwischen IT, Entwicklung und Produktmanagement ab.

Wenn Endkunden bei ihren Online-Geschäften genervt sind, liegt das meist an zu strengen Sicherheitsstandards der Unternehmen. Das hat zur Folge, dass ihre Customer Experience leidet. Wenn Sicherheit und Kundenerlebnis aber sinnvoll aber ineinandergreifen, steckt meist eine gut abgestimmte CIAM-Strategie dahinter. Was bedeuten würde, dass die IT-Sicherheit mit den kundennahen Bereichen wie Entwicklung, Marketing und Produktmanagement zusammenarbeitet, um beispielsweise die Login-Erwartungen der Kunden besser adressieren zu können.

Die Praxis zeigt, dass es hier noch sehr viel Luft nach oben gibt. Vor allem wenn es darum geht, dass alle gemeinsam die Kundenperspektive in den Mittelpunkt stellen. Das hat eine Trendstudie von IDG Research Services im Auftrag von Auth0/Okta kürzlich herausgefunden.

Zweifellos sind Fragen der IT-Sicherheit und des Datenschutzes von zentraler Bedeutung, wenn digitale Angebote Erfolg bei Kunden haben sollen. Daher ist es richtig, wenn 86 Prozent der Unternehmen die Absicherung ihrer Prozesse gegen Identitätsmissbrauch und Datendiebstahl im Kontext von CIAM für sehr wichtig halten. Dennoch ist es zu kurz gegriffen, wenn die Strategie zur Verwaltung der Kundenidentitäten und Berechtigungen bei weniger als 30 Prozent der befragten Entscheider*innen kundenbezogene Kriterien enthält, zwei Drittel der Unternehmen dafür die Compliance und Sicherheitsfragen in den Mittelpunkt stellen.

Interessanterweise verarbeiten rund 70 Prozent der Unternehmen neben Partner- und Mitarbeiterdaten vor allem Kundendaten in ihren IAM-Systemen. Daraus kann man schließen, dass es bei der überwiegenden Mehrzahl der Unternehmen ein großes Potenzial für CIAM gibt, da Kundendaten bereits in den IAM-Systemen vorgehalten werden.

Auch am Beispiel der eingesetzten Authentifizierungsverfahren zeigt sich, dass zum einen Kundenwünsche noch nicht genügend im Bewusstsein der Unternehmensentscheider verankert sind: Moderne und gleichzeitig sichere

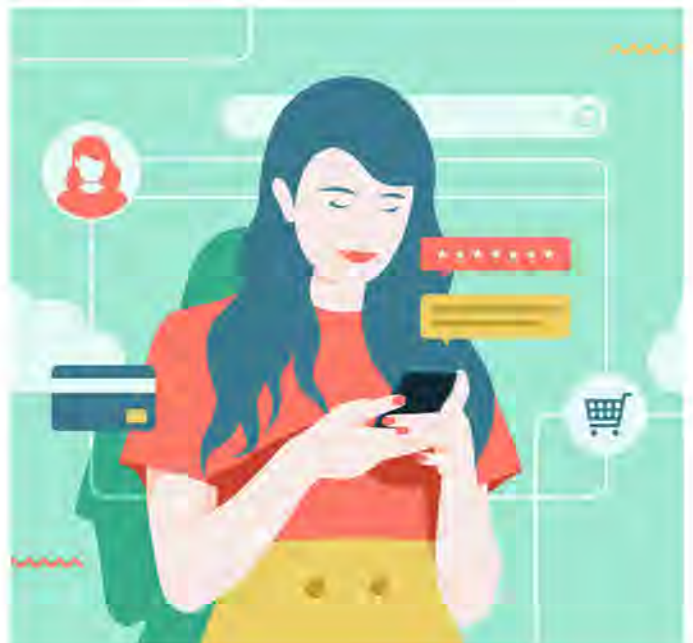


Authentifizierung wie MFA haben nach Ansicht der Befragten weniger Zukunft als der Passwort-Klassiker, PINs oder Sicherheitsfragen. Die Kunden wollen aber genau solche modernen Lösungen, die ihre Anmeldung sicher und komfortabel machen.

Unsere Empfehlung lautet vor diesem Hintergrund, neben der IT immer auch die Kundenbereiche und die Entwicklung mit an den Tisch zu holen, wenn ein CIAM-Projekt geplant wird. Dadurch kann sichergestellt werden, dass auch neue Funktionen, wie zum Beispiel eine biometrische Authentifizierung, berücksichtigt werden, die

sich die Kunden wünschen. „Die Ergebnisse zeigen: IT-Entscheider interpretieren CIAM in erster Linie als Erweiterung ihrer Infrastruktur. So sind sie es durch die Mitarbeiter-Authentifizierung gewohnt. Das Management digitaler Kundenbedürfnisse ist allerdings sehr viel komplexer und vielschichtiger, zumal unendlich viele unterschiedliche Zugänge dynamisch gemanagt und gesichert werden müssen. Hier sollten Entwicklung und Produktmanagement mit an den Tisch,“ so Vitor de Sousa, Regionaler Vice President Enterprise Sales Central & Eastern Europe bei Okta und Auth0.

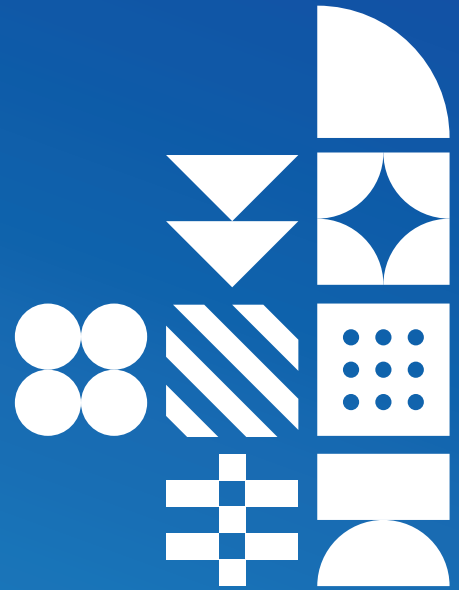
auth0.com/de



> Warum Metriken nicht ausreichen

Isabel Bär

Klassisches Software-Testing lässt sich nicht ohne Weiteres auf KI übertragen. Model Governance und interne Audits sind nötig, um um Fairness zu gewährleisten, Auflagen zu erfüllen und unternehmerische Risiken abzufedern.



Der Einsatz von Künstlicher Intelligenz (KI) bringt Verantwortung mit sich. Transparenz, Erklärbarkeit, Fairness sind dabei wesentliche Prinzipien, die ebenso gewährleistet sein müssen wie die hohe Leistungsfähigkeit des KI-Systems. Um diese Anforderungen einzuhalten, liegt es nahe, sich an Bereichen mit einer Tradition überprüfbarer Prozesse zu orientieren. Zwar funktionieren diese Prozesse nicht fehlerlos, aber ohne sie lassen sich Sicherheitsstandards nicht verwirklichen. Am offensichtlichsten ist das in sicherheitskritischen und regulierten Branchen wie der Medizin, aber auch in der Luft- und Raumfahrt oder im Finanzwesen.

Ähnlich wie diese Bereiche Prozesse benötigen, um relevanten Anforderungen nachzukommen, benötigt ein Unternehmen, das KI-Systeme einsetzt, geregelte Abläufe, durch die es Zugriff auf Machine-Learning-Modelle (ML) kontrolliert, Richtlinien sowie gesetzliche Vorgaben umsetzt, die Interaktionen mit den Modellen und deren Ergebnisse verfolgt sowie festhält, auf welcher Grundlage ein Modell erzeugt wurde. Insgesamt werden diese Prozesse als Model Governance bezeichnet. Model-Governance-Prozesse sind von Beginn an in jede Phase des ML-Lebenszyklus zu implementieren (Design, Development und Operations). Zur konkreten technischen Integration von Model Governance in den ML-Lebenszyklus hat die Verfasserin sich an anderer Stelle ausführlicher geäußert (alle Quellen zu diesem Artikel finden sich unter ix.de/z3q9).

Model Governance ist nicht optional (siehe Kasten „Model Governance“). Es gibt bereits bestehende Regularien, die Unternehmen in bestimmten Branchen erfüllen müssen. Am Beispiel des Finanzsektors lässt sich die Bedeutung von Model Governance gut illustrieren: Kreditvergabesysteme oder Zinsrisiko- und Preisbildungsmodelle für Derivate sind risikoreich und verlangen ein hohes Maß an Kontrolle und Transparenz. Der Umfang der Regularien dürfte künftig weiter zunehmen: so veröffentlichte die EU im April 2021 eine Ver-

in a nutshell

- > Machine Learning (ML) geht mit Verantwortung und Verpflichtungen einher, ethische Prinzipien wie Fairness sind auf allen Ebenen der Softwareentwicklung zu berücksichtigen – gerade auch bei der Datenbeschaffung.
- > Unternehmen müssen weitreichende Auflagen zur Modelldokumentation, Verwaltung von Metadaten, Zugriffskontrolle, Validierung und Reproduzierbarkeit erfüllen.
- > Model Governance unterstützt sie beim Einhalten der gesetzlichen Anforderungen und mindert das unternehmerische Risiko von KI-Systemen.

ordnung als ersten Rechtsrahmen für KI, die bestehende Regularien ergänzen würde. Der Entwurf teilt KI-Systeme in vier unterschiedliche Risikokategorien ein („unzulässig“, „hoch“, „begrenzt“, „minimal“). Die Risikokategorie definiert dabei Art und Umfang der Anforderungen, die an das jeweilige KI-System gestellt werden. KI-Software, die in die hohe Risikokategorie fällt, muss die strengsten Auflagen erfüllen.

Dazu zählen folgende Aspekte: Robustheit, Sicherheit, Genauigkeit (Accuracy), Dokumentation und Protokollierung sowie angemessene Risikobewertung und Risikominderung. Weitere Anforderungen sind die hohe Qualität der Trainingsdaten, Diskriminierungsfreiheit, Nachvollziehbarkeit, Transparenz, menschliche Überwachung sowie die Erforderlichkeit einer Konformitätsprüfung und der Nachweis der Konformität mit der KI-Verordnung durch eine CE-Kennzeichnung. Beispiele für die ML-Systeme dieser Kategorie sind private und öffentliche Dienstleistungen (wie die Bonitätsprüfung) oder Systeme, die in der Schul- oder Berufsausbildung eingesetzt werden, um über den Zugang zu Bildung und den beruflichen Werdegang einer Person zu entscheiden.

Da die Verordnung nicht nur für in der EU ansässige Unternehmen und Einzelpersonen gelten soll, sondern für jedes Unternehmen, das KI-Dienste innerhalb der EU anbietet, hätte das Gesetz einen ähnlichen Anwendungsbereich wie die DSGVO. Die Verordnung muss sowohl vom EU-Parlament gebilligt werden als auch die Gesetzgebungsverfahren der einzelnen Mitgliedsstaaten passieren. Wenn das EU-Parlament die Verordnung billigt und sie die legislativen Prozesse der EU-Staaten passiert, tritt das Gesetz frühestens 2024 in Kraft. Dann müssen Hochrisikosysteme während der Entwicklung eine Konformitätsbewertung für KI-Auflagen durchlaufen, um das KI-System in einer EU-Datenbank registrieren zu lassen. Im letzten Schritt ist eine Konformitätserklärung notwendig, sodass KI-Systeme die notwendige CE-Kennzeichnung erhalten, damit ihre Anbieter sie in den Verkehr bringen können.

Wichtig ist außerdem, dass Regulierung nicht der einzig ausschlaggebende Aspekt für Model-Governance-Prozesse ist. Denn auch Modelle, die in schwächer regulierten Kontexten im Einsatz sind, kommen an Model Governance nicht vorbei. Neben der Erfüllung gesetzlicher Vorgaben müssen Unternehmen wirtschaftliche Einbußen und Reputationsverluste ebenso abwenden wie juristische Schwierigkeiten. ML-Modelle, die einer Marketing-Abteilung Informationen über die Zielgruppe liefern, können im Betrieb an Präzision verlieren und eine falsche Informationsgrundlage für wich-

Model Governance

Um den Anforderungen nachzukommen, benötigt ein Unternehmen Prozesse, durch die es

- > die Zugriffe auf ML-Modelle kontrolliert
- > Richtlinien/gesetzliche Vorgaben umsetzt
- > die Interaktionen mit den ML-Modellen und deren Ergebnisse verfolgt
- > festhält, auf welcher Grundlage ein Modell erzeugt wurde

Model Governance bezeichnet diese Prozesse in ihrer Gesamtheit.

tige Folgeentscheidungen bereitstellen. Somit stellen sie ein finanzielles Risiko dar. Model Governance dient zur Qualitätssicherung und mindert unternehmerische Risiken.

Model Governance als Herausforderung

Die sich abzeichnenden EU-Vorgaben, bestehende Regelungen und Unternehmensrisiken machen es notwendig, Model-Governance-Prozesse von Beginn an zu implementieren. Die Bedeutung von Model Governance ergibt sich für viele Unternehmen allerdings oft erst dann, wenn ML-Modelle in die Produktion gehen und in Einklang mit gesetzlichen Regelungen stehen sollen. Dazu kommt, dass der abstrakte Charakter rechtlicher Vorgaben Unternehmen praktisch fordert: So geben in einer Algorithmia-Studie 56 Prozent der Befragten die Implementierung von Model Governance als eine der größten Herausforderungen an, um ML-Anwendungen langfristig erfolgreich in Produktion zu bringen.

Audits als standardisierte Prüfprozesse im Model-Governance-Framework

Ein wichtiger Bestandteil von Model Governance sind Audits, um zu prüfen, ob KI-Systeme den Unternehmensrichtlinien, Branchenstandards oder Vorschriften entsprechen. Dabei gibt es interne und externe Audits. Die Studie Gender Shades ist ein Beispiel für einen externen Auditprozess: Sie prüfte Gesichtserkennungssysteme großer Anbieter hinsichtlich ihrer Genauigkeit bezüglich des Geschlechts und der Ethnie und konnte dabei eine abweichende Präzision des Modells feststellen. Dieser Blick von außen ist limitiert, da externe Prüfprozesse nur Zugang zu Modellergebnissen, aber nicht zu den zugrundeliegenden Trainingsdaten oder Modellversionen besitzen. Das sind wertvolle Quellen, die Unternehmen in einem internen Auditprozess einbeziehen müssen. Die Prozesse sol-

len eine kritische Reflexion über die potenziellen Auswirkungen eines Systems ermöglichen.

Ethische Prinzipien als nicht-funktionale Eigenschaften

Die relevanten Testaspekte von KI-Software lassen sich in funktionale und nicht-funktionale Eigenschaften einteilen. Correctness als funktionale Eigenschaft lässt sich durch Metriken wie Accuracy und Precision/Recall mathematisch direkt erfassen. Sie geben an, wie hoch die Übereinstimmung zwischen den Vorhersagen des trainierten Modells und den tatsächlichen Predictions ist (Gold Standard). Dazu gibt es etablierte Validierungsverfahren wie die Kreuzvalidierung, die durch Isolation der Testdaten über eine Datenstichprobe prüft, wie gut das trainierte Modell die richtigen Modellergebnisse (Labels) für neue Daten vorhersagt.

Nicht-funktionale Eigenschaften entsprechen den ethischen Prinzipien wie Fairness, Datenschutz, Interpretierbarkeit, Robustheit und Sicherheit. Anders als funktionale Eigenschaften können sie nicht auf einen breiten Fundus standardisierter Metriken und Praktiken aus dem Bereich des maschinellen Lernens zurückblicken. Auch hier besteht die Herausforderung darin, dass das Testen nicht-funktionaler Eigenschaften von KI-Software (noch) nicht standardisiert ist. Erschwerend kommen Abwägungen zwischen verschiedenen Eigenschaften hinzu: Fairness verringert die Accuracy und umgekehrt.

Metaphorisch lässt sich KI-Software als Kraftwerk bezeichnen: Ein funktional einwandfreier, reibungsloser Betrieb heißt nicht, dass das Kraftwerk der Umwelt keinen Schaden zufügt. Der störungsfreie Ablauf entspricht den funktionalen, der Schutz der Umwelt den nicht-funktionalen Kriterien. Die

Kontrafaktisch testen, ob ein Modell fair ist

Causal-Reasoning-Ansätze stützen sich auf Werkzeuge der Kausalinferenz. Die Definition der kontrafaktischen Fairness basiert auf der Intuition, dass eine Entscheidung gegenüber einer Person dann fair ist, wenn sie in der tatsächlichen Welt und in einer kontrafaktischen Welt, in der die Person einer anderen demografischen Gruppe angehört, gleich ist.

Damit ist Counterfactual Fairness dann gegeben, wenn sich eine Prediction nicht ändert, obwohl das geschützte Attribut in das kontrafaktische Gegenteil verkehrt wird. Beispielsweise müsste die Entscheidung für oder gegen die Kreditwürdigkeit einer Person gleich ausfallen, wenn das Attribut von „männlich“ in „weiblich“ verändert wird.

Metapher zeigt, dass es für funktionale und nicht-funktionale Eigenschaften verschiedene Prüfprozesse braucht. Für Erstere sind Best Practices aus dem ML-Fundus anwendbar, für Letztere hingegen braucht es noch Forschungsarbeit.

Wie entsteht „Unfairness“?

Zunächst ist die Frage interessant, wie Ungerechtigkeit (Unfairness) überhaupt entsteht. Die Regel ist dabei einfach: Was die Modelle lernen, manifestiert sich in den Trainingsdaten. Im überwachten Lernen bestehen Trainingsdatensätze aus den Eingabedaten und dazugehörigen Labels. Wenn die Datenlabels Bias enthalten, wird das Modell diese Grundeinstellung übernehmen und von Anfang an lernen. Daher ist es wichtig, die Labels ausreichend zu überprüfen. Bias kann sich aber auch inhärent aus den Daten, nicht nur aus den Labels ergeben: Enthalten die Trainingsdaten an sich bereits Bias, greift der Algorithmus ihn ebenfalls auf. Dieses Problem besteht etwa bei umfangreichen, mit großen Datenmengen aus dem Internet trainierten Sprachmodellen. Es ließ sich nachweisen, dass die Leistungsstärke eines Modells mit der Stärke eines stereotypen Bias korreliert: Mit steigender Präzision nimmt auch der Bias zu.

Auch ein geringer Stichprobenumfang bei Minderheitsgruppen kann zu einer Homogenisierung des Lernprozesses des Modells zugunsten der Mehrheitsgruppen führen, beispielsweise durch mehr Fotos männlicher als weiblicher Gesichter in den Trainingsdaten. Neben den Daten spielen auch im Trainingsprozess verwendete Merkmale (Features) eine Rolle. Kann das Modell nicht ausreichend viele Merkmale nutzen, erschwert das dem Algorithmus, den Zusammenhang zwischen Ein- und Ausgabe zu lernen. Aus diesem Grund reagierte IBM mit Diversity in Faces als Versuch, die Diversität der Fotos in den Trainingsdaten zu erhöhen. Und schließlich können Merkmale „Stellvertreter“ für ausgeschlossene sensible Attribute sein: Auch wenn geschützte Attribute bei der Entscheidungsfindung nicht explizit verwendet werden, können sie implizit beteiligt sein, wenn sie mit den ausgeschlossenen Merkmalen korrelieren.

Definitionen für Fairness und Ableitungen von Teststrategien

Welche Audits und welche Metriken bieten sich an, um Fairness zu testen? Die bereits bekannte Konsenslücke klafft auch für die Definition von Fairness auseinander. Erschwerend kommt hinzu, dass die Vielfältigkeit der verschiedenen Ursachen für Fairness zeigt, dass sich Fairness nicht mit einer simplen Metrik oder Teststrategie herstellen lässt – Fairness-Audits müssen Teil der Model-Governance-Prozesse sein, die

die Qualitätssicherung der Trainingsdaten und des Modells sicherstellen. Dazu kommt, dass die verschiedenen Anwendungsfälle für KI zu vielfältig sind, als dass es eine gut generalisierende One-Size-Fits-All-Lösung geben könnte. Die Frage, lässt sich also nicht nur an einer simplen Metrik festmachen. Dennoch soll es zunächst um konkrete Möglichkeiten gehen, Fairness quantitativ zu erfassen, bevor diese Audits in das Model-Governance-Framework eingebettet werden.

Statistische Ansätze bieten die am leichtesten messbaren Definitionen von Fairness, und sie bilden gleichzeitig die Grundlage für weiterführende Ansätze. Zur Quantifizierung von Fairness lassen sich statistische Metriken nutzen. Von diesen Messgrößen leiten sich Definitionen ab, die sich auf die Ausgabe von Modellen konzentrieren. Fairness lässt sich aufgrund ähnlicher Fehlerquoten der Ausgaben für unterschiedlich sensible demografische Gruppen definieren. Entsprechend ist ein Algorithmus dann fair, wenn Gruppen, die auf der Grundlage sensibler Attribute ausgewählt werden, die gleiche Wahrscheinlichkeit vorteilhafter Entscheidungen haben („Group Fairness“). Zudem lässt sich untersuchen, ob die Genauigkeit des Modells für verschiedene Subgruppen gleich ist (Gleichheit der Gesamtgenauigkeit). Am Beispiel einer Kreditwürdigkeitsprüfung wäre diese Definition von Fairness dann erfüllt,

Fairness in KI-Systemen

Das Ziel der Gewährleistung von Fairness ist der Schutz sensibler Attribute wie Geschlecht, Religionszugehörigkeit oder sexueller Orientierung vor unfairer algorithmischer Entscheidungsfindung. Das Recht auf Diskriminierungsfreiheit ist im EU-Rechtsentwurf für KI-Systeme der hohen Risikokategorie explizit verbrieft. Während sich die Ungerechtigkeit bei der Gender-Shades-Studie leicht intuitiv erfassen lässt, besteht nun die Herausforderung darin, den abstrakten Begriff der Fairness objektiv, metrikbasiert und möglichst skalierbar zu definieren.

wenn die Wahrscheinlichkeit für alle Personen gleich ist, dass Antragsstellende mit einem tatsächlich guten Kreditscore als kreditwürdig eingestuft werden und dass solchen mit einem schlechten Kreditscore die Kreditwürdigkeit abgesprochen wird, ohne Ansehen der Geschlechtszugehörigkeit.

Für das Testen statistischer Ansätze sind bereits erste Lösungen verfügbar: Fairness Indicators von TensorFlow ist eine Bibliothek, die die Berechnung häufig identifizierter

Java | Spring | REST | JPA | PostgreSQL | Maven | Jenkins | SonarQube

solutio

(Senior) Java Developer (m/w/d) für die Digitalisierung im Gesundheitswesen

Du entwirfst Lösungen für unser Zahnarztsoftware charly.
Diese prägen jeden Tag mehr als 3.500 Zahnarztpraxen - nicht nur in Deutschland.

Als **Java Developer** bist du auf dem Weg ein angesehener Vordenker, ein hochproduktiver Ingenieur und ein Vorbild für andere Entwickler bei solutio zu werden.

Als **Senior Java Developer** bist du die führende Stimme in einem interdisziplinären, agilen Team von Entwicklern, Designern, Architekten und fachspezifischen Product Ownern.

Du arbeitest entlang objektorientierter Prinzipien, Design-Patterns und Clean Code mit unseren Entwicklungswerkzeugen. Clean oder Hexagonal Architecture, das Spring-Framework sowie unser agiles Vorgehen mit Scrum gehören zu unserem/deinem Alltag.

Damit du dann arbeiten kannst, wenn du am kreativsten bist, bieten wir
» flexible Arbeitszeiten » Zeitkonten » mobiles Arbeiten



Wir freuen uns auf deine
Bewerbungsunterlagen über
unsere Homepage. Ein
detaillierter Lebenslauf genügt.

Für Infos zum Bewerbungs-
prozess melde dich gerne bei
Ramona Fortenbacher
☎ 07031 4618-632

Bei fachlichen Fragen gerne direkt
bei Mirko Dokic
(Gruppenleiter der
Softwareentwicklung)
☎ 07031 4618-999

Werde Solutianer!
join us

Fairness-Metriken mit verbesserter Skalierbarkeit auf großen Datensätzen und Modellen bietet. Darüber hinaus unterstützt Fairness Indicators die Auswertung der Verteilung von Datensätzen und der Modelleleistung über verschiedene Benutzergruppen sowie die Berechnung statistisch signifikanter Unterschiede auf der Basis von Konfidenzintervallen.

Zwar sind statistische Ansätze gut messbar, sie können jedoch zu kurz greifen. Fairness lässt sich nicht allein durch ähnliche Fehlklassifizierungsquoten erklären, insbesondere, wenn alle anderen Attribute mit Ausnahme des sensiblen Attributs ignoriert werden. Beispielsweise könnte ein KI-System zur Kreditwürdigkeitsprüfung demselben Anteil männlicher und weiblicher Bewerber eine positive Bewertung zuweisen – statistische Ansätze würden das Modell dann als gerecht beurteilen. Doch wenn die männlichen Bewerber zufällig ausgewählt wurden, während weibliche Bewerber schlicht diejenigen sind, die die meisten Ersparnisse haben, wäre Fairness nicht gegeben.

Similarity-based Measures stellen nicht die Modellergebnisse und Fehlklassifizierungsquoten, sondern den Prozess der Entscheidungsfindung sowie die Verwendung von Merkmalen im Trainingsprozess in den Vordergrund. Daraus lässt sich „Fairness through Unawareness“ als Konzept für Fairness ableiten: Algorithmen können als fair gelten, wenn geschützte Attribute aus den Trainingsdaten ausgeschlossen worden sind. In unserem Beispiel bedeutet dies, dass geschlechtsspezifische Merkmale nicht für das Training des Modells verwendet werden, so dass Entscheidungen nicht auf diesen Merkmalen beruhen können. Doch auch dieser Ansatz hat Einschränkungen: Das Ausschließen geschützter Attribute reicht nicht aus, da andere, ungeschützte Attribute Informationen enthalten können, die mit den ausgeschlossenen geschützten Attributen korrelieren. In diesem Fall wäre das ursprünglich ausgeschlossene Attribut implizit in anderen Attributen enthalten und würde den Entscheidungsprozess indirekt beeinflussen (siehe Kasten „Kontrafaktisch testen“).

Unfairness mit manipulierten Daten aufdecken

Adversiales Testen ist eine gängige Strategie, die zum Aufdecken von Schwachstellen einen böswilligen Angriff auf ein System simuliert. Beim adversialen Testen erhält das Modell Eingabedaten, die mit kleinen, absichtlichen Merkmalsmanipulationen versehen sind. Auf diese Weise wird getestet, ob das Modell für speziell zugeschnittene Eingabedaten unerwünschte Vorhersagen trifft. Die Manipulation der Eingabedaten ist domänenspezifisch und kann von Ana-

lysen algorithmischer Unfairness inspiriert sein. Die Idee, die Reaktion eines Modells auf Eingabedaten zu testen und auf diese Art von Bias zu quantifizieren, findet sich mittlerweile auch in Frameworks.

So lässt sich mit dem Benchmarking-Datensatz StereoSet prüfen, ob und wie stark die Vorurteile von Sprachmodellen bezüglich Geschlecht, Ethnie, Religion und Profession ausgeprägt sind: Entwicklerinnen und Entwickler können fertig trainierte Sprachmodelle einreichen, um diskriminierende Entscheidungsfindung in Sprachmodellen zu messen und gleichzeitig die Sprachmodellierungsleistung zu berücksichtigen. StereoSet betrachtet die Gesamtleistung des Modells als gut, wenn das Modell in der Lage ist, den Zielkonflikt zwischen Genauigkeit und Fairness abzuschwächen und so ein genaues Verständnis natürlicher Sprache bei gleichzeitiger Minimierung von Verzerrungen zu gewährleisten. Solche Frameworks können zwar eine gute Leitlinie sein, ersetzen aber nicht das individuelle Testen, das fest in einem Prüfungsprozess eingebettet sein muss.

Warum Audits alleine nicht ausreichen

Die verschiedenen Prüfprozesse sind wichtig – genügen aber nicht, um Fairness zu gewährleisten. Vielmehr kommen Audits als eine der ersten Möglichkeiten in Betracht, um Probleme zu identifizieren. Sie müssen Teil des Model-Governance-Frameworks sein und sollten es ergänzen. Für sich allein besitzen sie hingegen keine Aussagekraft – nur ein ganzheitlicher Ansatz kann alle Aspekte berücksichtigen, die für Fairness eine Rolle spielen. Neben der Validierung funktionaler und nicht-funktionaler Anforderungen, die die hier beschriebenen Auditprozesse zum Testen von Fairness enthalten kann, ist sauberes Dokumentieren wichtig. Dokumentationen sollten bereits in der ersten Phase des ML-Lebenszyklus, der Entwicklung, begonnen werden. In der Development-Phase geht es um den Aufbau einer robusten und reproduzierbaren Trainingsprozedur, die aus Datenverarbeitungs- und Modellaufbausritten besteht. Dieser Aufbauprozess ist experimentell und iterativ, wobei wichtige Informationen über Daten und Parameter festzuhalten sind (Experimental Tracking).

Nach dem Aufbau der Trainingsprozedur gilt es, entwickelte Modelle hinsichtlich funktionaler und nicht-funktionaler Eigenschaften zu evaluieren (hier greifen die Auditstrategien zum Testen von Fairness relevant). Die Ergebnisse der Evaluation und alle Informationen über den Aufbau der Trainingsprozedur fließen in die Dokumentation ein. Zudem sollte sie eine Erklärung des Use-Case-Kontextes, eine High-Level-Erklärung des Algorithmus, Modellparameter, Anweisungen zur Reproduktion des Modells und Beispiele für das

Training der Algorithmen sowie Beispiele für das Treffen von Predictions durch den Algorithmus beinhalten sollte. Die Dokumentation lässt sich durch Toolkits wie Model Cards und Data Sheets praktisch unterstützen. Data Sheets halten fest, welche Mechanismen oder Verfahren für die Datenerhebung verwendet wurden und ob ethische Überprüfungsverfahren (Audits) stattgefunden haben. Model Cards informieren über die Art der Modellerstellung, die bei der Entwicklung getroffenen Annahmen sowie über Erwartungen bezüglich des Modellverhaltens bei verschiedenen kulturellen, demografischen oder phänotypischen Gruppen.

Fairness braucht Model Governance

Vollständige Dokumentation schafft Reproduzierbarkeit und Transparenz nach außen. Nach dem Deployment muss diese Sichtbarkeit (Observability) im produktiven System gegeben sein. Hier spielt zum einen die Versionierung von Modellen und Datensätzen eine wichtige Rolle. Sie dient der Wahrung des Unveränderlichkeitsgrundsatzes, sodass alle Modelle sich ohne Datenverluste und Veränderung reproduzieren lassen.

Damit ist auch gewährleistet, dass eine Model Prediction der Modellversion, die sie produziert hat, zugeordnet werden kann. Zum anderen muss ein Monitoring-System die Leistung des produktiven Modells kontinuierlich überwachen und relevante Metriken in einem Report zusammenfassen sowie visualisieren. Diese Werte aus dem Model-Logging sollten in Metriken aufbereitet und in Dashboards zu Protokollierungs-, Analyse- und Kommunikationszwecken visualisierbar sein. Wird im Monitoring der Leistungsabfall ei-

nes Modells (Model Decay) festgestellt, muss das Modell mit neuen Trainingsdaten trainiert und dann re-deployed werden.

Vor jedem neuen Deployment sollten erneut Audits stattfinden, um ethische, rechtliche oder geschäftliche Risiken zu kontrollieren. Ethische Prinzipien wie Fairness sind auf jeder Ebene der Softwareentwicklung zu berücksichtigen, unter anderem bereits bei der Datenbeschaffung. Fairness lässt sich nicht mit einer simplen Metrik oder Teststrategie herstellen: Es braucht eine entsprechend ausgerichtete Unternehmenspolitik, die das anerkennt. Ohne Model Governance sind die Risiken von KI-Systemen unkalkulierbar.

Quellen

Weiterführende Hinweise und alle Quellen finden sich unter ix.de/z3q9.



Isabel Bär

studiert Data Engineering am Hasso-Plattner-Institut und arbeitet als Werkstudentin bei INNOQ. Sie beschäftigt sich mit Fragen rund um den langfristigen erfolgreichen Einsatz von Künstlicher Intelligenz (KI), wozu insbesondere MLOps und das Implementieren von Model Governance gehören.



iABG
www.iabg.de

AUTOMOTIVE • INFOKOM • MOBILITÄT, ENERGIE & UMWELT •
LUFTFAHRT • RAUMFAHRT • VERTEIDIGUNG & SICHERHEIT

IABG. Die Zukunft.

Die IABG bietet integrierte, innovative Lösungen in technologieintensiven Branchen. Finden Sie es spannend, heute schon an Themen der Zukunft zu arbeiten und mit Ihrem Engagement und Ihrer Kompetenz dazu beizutragen, die Welt von morgen mitzugestalten? Dann sind Sie bei uns genau richtig!


Berufseinsteigern (gn) bis hin zu erfahrenen **IT-Experten (gn)** bieten wir vielfältige Karrieremöglichkeiten.

Einen Überblick über unsere **IT-Stellen** finden Sie hier:



karriere.iabg.de

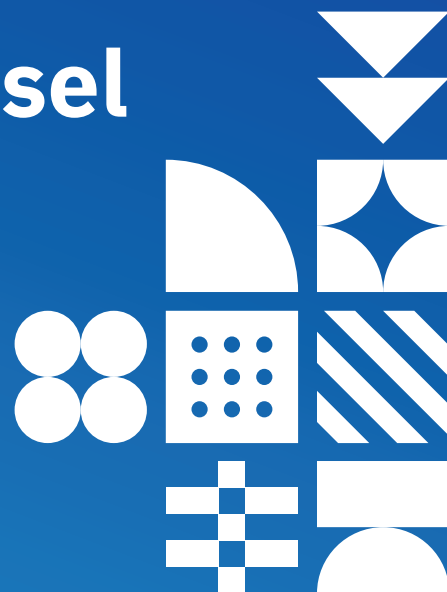
LinkedIn

XING  kununu 

> Kontext als Schlüssel zur sicheren Cloud

Manuel Benz, Andreas Dann, Dr. Johannes Späth

Security bedarf im Zeitalter der Cloud besonderer Sorgfalt – gleichzeitig muss sie mit der enormen Entwicklungsgeschwindigkeit schritthalten. Dieser Balanceakt ist nicht immer einfach, zumal bei begrenzten Team-Ressourcen. Was sind die richtigen Prioritäten?



Mit Cloud-Technologien wie Serverless, Kubernetes oder Terraform ist das Versprechen verbunden, skalierbare Applikationen einfach und mit minimalem Aufwand zu implementieren. Statt eigenständige Anwendungen einzeln einzubinden, nutzt eine Cloud-native-Anwendung vorgefertigte Services der Cloud-Provider, die sie zu einem Ganzen kombiniert. Alle typischen Anwendungsfälle von User Management, verwalteten NoSQL- und SQL-Datenbanken bis hin zu Engines für Machine Learning gibt es in Form von Microservices in der Cloud. Weitere spezifische Business-Logik einer Anwendung lässt sich in selbst entwickelten und gewarteten Komponenten umsetzen und zumeist über virtuelle Maschinen (VMs), Docker-Container oder Serverless Functions an die Cloud-native Infrastruktur anbinden.

Die gesamte Anwendung läuft auf Ressourcen des Cloud-Providers. Ein Vorteil ist dabei, dass Unternehmen sich den Betrieb im eigenen Rechenzentrum (on Premises) sparen und mühsame Wartungsarbeiten einer eigenen IT-Infrastruktur entfallen. Sie können sich auf die Umsetzung der Geschäftslogik ihrer Applikation konzentrieren, lautet ein Werbeargument der Anbieter von Clouddiensten.

Dass Cloud-Applikationen sicherer seien, da sich der Provider um die Security kümmert, ist jedoch ein Fehlschluss: Sicherheit in der Cloud bleibt eine geteilte Verantwortung zwischen dem Provider und den Anwendern. Die Provider nutzen das Shared-Responsibility-Modell. Demzufolge sind sie unter anderem verantwortlich für die physische Sicher-

heit der Hardware oder die Aktualität der Treiber, für das Betriebssystem, die Laufzeitumgebungen oder auch die physische Absicherung des Netzwerks. Cloud-Anwender stehen jedoch ihren Kunden gegenüber weiterhin für die Anwendungssicherheit und für die korrekte Konfiguration der Cloud-Umgebung gerade.

Auf Anwendungsebene liegen insbesondere Sicherheitslücken im selbst entwickelten Code, in eingebundenen Bibliotheken oder in Docker-Images in der Verantwortung von

in a nutshell

- > Sicherheit in der Cloud basiert auf dem Shared-Responsibility-Modell und erfordert weitere Maßnahmen zur Absicherung der Anwendungen.
- > Best Practices empfehlen klassisches SCA und SAST auf der Anwendungsebene sowie CSPM- und CIEM-Werkzeuge, die Cloud-Fehlkonfigurationen aufdecken.
- > Beim Priorisieren der automatisch einlaufenden Warnungen von Security-Werkzeugen spielt der Kontext in der Cloud eine wesentliche Rolle – insbesondere die Frage, ob die von einer Schwachstelle betroffenen Ressourcen von außen erreichbar sind.

Quereinsteiger (w|m|d), Praktikanten (w|m|d) & Werkstudenten (w|m|d) willkommen



Wir suchen **Linux Entwickler** (w|m|d)

Das sind deine Stärken:

- eine schnelle Auffassungsgabe
- analytisches Denken
- Prozess- und Softwaredesign
- Leidenschaft für Linux/Open Source

Diese Kenntnisse bringst du mit:

- PHP, JavaScript (ECMA Script u. neuer), CSS3
- Refactoring & Brownfield Projekte
- Continuous Integration/Delivery
- außerdem wünschenswert:
 - Config Automation (Salt etc.)
 - Python
 - Ruby
 - Go
 - Java/Android
 - GIT

Das erwartet dich bei uns:

- ca. 90% Homeoffice
- vielfältige & abwechslungsreiche Einsätze
- familiäres Arbeitsklima & flache Hierarchien
- Vereinbarkeit von Job & Familie

Mehr erfahren & bewerben:

jobs@b1-systems.de

Formlose Bewerbung genügt



B1 Systems GmbH - Ihr Linux-Partner

Linux/Open Source Consulting, Training, Managed Service & Support

ROCKOLDING · KÖLN · BERLIN · DRESDEN · JENA

www.b1-systems.de · info@b1-systems.de

Cloud-Anwendern. Die richtige Konfiguration der Cloud-Accounts und dort enthaltenen verwalteten Ressourcen, die VPN-Konfiguration, die Nutzerverwaltung und die Zugriffsrechte liegen ebenfalls in seiner Verantwortung. Abhängig von den Cloud-Techniken variiert die Sicherheits-Verantwortung der Anwender zwischen 8 und 52 Prozent in Bezug auf die Gesamtverantwortung, schätzt Ory Segal von Palo Alto Networks (8 Prozent bei einer Applikation, die hauptsächlich auf VMs setzt, 52 Prozent bei modernen, vollständigen Serverless-Implementierungen, siehe ix.de/zb98).

Principle of Least Privilege

Wenn beispielsweise eine Elastic-Cloud-Computing-Instanz (EC2) Zugriff auf unzählige S3-Buckets und die darin enthaltenen Daten hat, wurde das „Principle of Least Privilege“ nicht optimal umgesetzt (siehe Kasten „Capital-One-Hack“). Komponenten sollten immer nur Zugriff auf die minimal notwendigen Daten haben. Eine strengere Zugriffskontrolle hätte den Schaden, den der Angriff anrichtete, deutlich reduzieren können. Um solche Schwachstellen im Cloud-Umfeld zu verhindern, gibt es Best Practices und Cloud-Security-Tools, die die Richtlinien frühzeitig und kontinuierlich überprüfen.

Im Beispiel des Capital-One-Hacks hätte ein Werkzeug für das Cloud Security Posture Management (kurz CSPM) in Kombination mit einem Tool zum Static Application Security Testing (kurz SAST) die beiden Schwachstellen frühzeitig aufgedeckt und den Angriff verhindern können.

Ein CSPM-Werkzeug findet Fehlkonfigurationen innerhalb des Accounts. Es überprüft alle Cloud-Ressourcen auf ihre Konformität mit gängigen Regelkatalogen wie dem Payment Card Industry Data Security Standard (PCI DSS), der Service Organization Control 2 (SOC2) oder den Benchmarks des Center of Internet Security, kurz CIS (Erläuterungen finden sich unter ix.de/zb98). Zu den gängigen Regeln in den Benchmarks gehören typische Fehler, unter anderem offene Wartungsporths mit SSH, fehlende Datenverschlüsselung („in transit“ und „at rest“) oder öffentlich verfügbare Ressourcen wie S3-Buckets und Datenbanken. Tatsächlich zählen unabsichtlich öffentliche S3-Buckets zu den am häufigsten ausgenutzten Sicherheitslücken. Weitere Regeln, die das CSPM überprüft, sind typische Schwachstellen in Netzwerken oder Firewalls, die aus unsicheren Default-Einstellungen resultieren. Sie sind meist für den jeweils aktuellen Anwendungsfall unzulänglich konfiguriert und bieten nicht die nötige Sicherheit. Weiterhin gehören auch schwache Passwortrichtlinien, die nicht aktivierte Zwei-Faktor-Authentifizierung oder unzu-

> Öffentliche S3-Buckets sind eine häufige Sicherheitslücke.

reichendes Event-Logging zu den Fehlern, die CSPM-Werkzeuge erkennen. Im Fall des Capital-One-Hacks waren die Privilegien des Servers und damit der WAF unzureichend beschränkt und ermöglichten es den Angreifern, den EC2-Metadaten-Service von AWS zu erreichen. Das auf IT spezialisierte Marktforschungsunternehmen Gartner prognostiziert, dass solche Fehlkonfigurationen von Ressourcen eine der größten Quellen für Sicherheitslücken in der Cloud werden: Bis 2025 werden wohl 99 Prozent der erfolgreichen Cloud-Angriffe darauf und allgemein auf menschliches Versagen zurückzuführen sein (siehe ix.de/zb98).

Auch moderne Cloud-Anwendungen müssen nach wie vor die eigene Businesslogik auf der Anwendungsebene modellieren und enthalten auch vom Hersteller entwickelte Software – meist in einem Docker-Container oder in Form einer Serverless-Funktion. Diese Komponenten gilt es auf Sicherheitslücken zu überprüfen. Hierfür kommen klassischerweise statische Codeprüfungen wie eine Datenflussanalyse infrage. Sie überprüft den Anwendungscode als Source- oder Bytecode, testet alle potenziellen Pfade der Anwendung und stellt fest, ob extern modifizierbarer Input eine kritische Operation auslösen kann. Eine solche Analyse kann beispielsweise eine SSRF-Schwachstelle erkennen, wie sie im Capital-One-Hack ausgenutzt wurde.

Weitere Best Practices: CIEM und SCA

Es gibt noch mehr Maßnahmen, unbefugten Zugriff auf die Cloud-Anwendungen und -Accounts abzuwenden. Hierzu zählen das Cloud Infrastructure Entitlements Management (CIEM) und die Software Composition Analysis (SCA).

Das Zugriffsmanagement ist ein wesentliches Security-Thema und dient dazu, den Zugriff auf Services und Ressourcen zu verwalten. Neben den einzelnen Cloud-Ressourcen wie virtuellen Maschinen (VM), Containern,

Serverless-Funktionen, Datenbanken und S3-Buckets, die untereinander kommunizieren, gilt es, Nutzerkonten in der Cloud zu verwalten. Die Zugriffsberechtigungen für einzelne Accounts auf die

Ressourcen sind im Sinne des Principle of Least Privilege minimal zu halten. Anderenfalls können Angreifer durch das unbefugte Annehmen einer höher privilegierten Rolle (Privilege Escalation) den Zugriff auf interne Daten erhalten oder tiefer in das System eindringen. Gefährdet sind dabei sensible Unternehmensdaten wie Datenbankinhalte, Kostenabrechnungen der Cloud-Accounts und vieles mehr. Wichtig ist, das Risiko von Insidern im Blick zu behalten, die aus dem Inneren des Unternehmens heraus ihre Rechte missbrauchen (zum Beispiel ehemalige Mitarbeiter).

W> World
Congress
2022

14.06–15.06.2022
BERLIN CITYCUBE



#WEAREDEV

2

Tag

200+

Speakers

5000+

Developers

Die globale
Entwickler-Community
trifft sich in Deutschland

worldcongress.dev

CIEM-Werkzeuge geben einen Überblick über die Zugriffsberechtigungen innerhalb des Accounts und stellen grafisch dar, wer auf welche Ressourcen und Informationen in der Cloud zugreifen kann. Aus Organisationssicht lassen sich damit die Rechte der einzelnen Nutzer und Ressourcen korrekt verwalten und einschränken.

Neben dem oben erwähnten Static Application Security Testing (SAST) gibt es auf der Anwendungsebene einer Cloud-Applikation weitere Qualitäts- und Sicherheitschecks. Hierzu zählt die Software Composition Analysis (SCA), die in einem gut strukturierten Software Development Lifecycle (SDLC) nicht fehlen sollte. Sie untersucht, welche Dritt-anbieterbibliotheken in einer Applikation enthalten sind. Das ist notwendig, um existierende Sicherheitslücken in den Fremdbibliotheken, besser bekannt als CVEs (Common Vulnerabilities and Exposures), oder auch Lizenzverletzungen zu erkennen.

Die Software Composition Analysis ist in den letzten Jahren ein unumgängliches Werkzeug zum Durchsetzen der Anwendungssicherheit geworden. Das liegt vor allem daran, dass Software heutzutage ausgiebig Gebrauch bestehender Features macht, statt das Rad jedes Mal neu zu erfinden. Eine wissenschaftliche Untersuchung der Queen's University Belfast zum Aufspüren von Schwachstellen in Open-Source-Software stellte bereits 2017 bis zu 80 Prozent Fremdcodeanteile in quelloffenen Anwendungen fest (siehe ix.de/zb98).

Eine solche Analyse der Softwarezusammensetzung untersucht zudem die Applikationsumgebung. Beispielsweise geht es darum, welche Applikationen und Shared Libraries in einer virtuellen Maschine oder einem Docker-Container installiert sind und ob sie bekannte Sicherheitslücken enthalten. Aktuelle Analysen umfassen auch Dockerfiles und bemerken, wenn sich Container mit Root-Rechten ausführen lassen oder wenn bestimmte bei Angreifern beliebte Ports wie der SSH-Port 22 offen sind.

Einige Scanner untersuchen auch die Docker- und VM-Images auf unsichere Konfigurationsdateien, obwohl das nicht zur Kernaufgabe der Software-Composition-Analyse gehört. So lässt sich unter anderem eine unsichere SSH-Konfiguration oder gelegentlich sogar eine SSH-Schlüsseldatei in den Images finden. Fehlkonfigurationen wie die der Apache-Web-Applikation-Firewall (WAF) beim Capital-One-Hack lassen sich hier unter Umständen von einigen Tools finden, sofern sie die zugehörigen Konfigurationsfiles interpretieren können.

Unzählige Werkzeuge, fehlender Kontext

Zur Absicherung einer Cloud-Umgebung sind die genannten Werkzeuge (CIEM, CSPM, SAST und SCA) unerlässlich. Feh-

ler in der Cloud-Konfiguration können leicht Türen für Angreifer öffnen, die sich anschließend über drastischere Sicherheitslücken in der Anwendung weiter ausnutzen lassen.

In einer idealen Welt schließen Sicherheitsteams alle gefundenen Schwachstellen in vollem Umfang – die Praxis sieht jedoch anders aus. Der Sicherheitsbereich ist den DevOps-Teams hinsichtlich Personal und verfügbare Ressourcen meist deutlich unterbesetzt – auf Dutzende DevOps-Mitarbeiter kommt meistens nur eine Security-Fachkraft, ein klares Missverhältnis. Die Ressourcen der Security-Abteilung sind außerdem beschränkt und stehen einer kontinuierlichen Weiterentwicklung der Cloud-Umgebungen gegenüber, die sich häufig mehrmals in der Woche oder gar täglich aktualisieren. Um dem gewachsen zu sein, müssten die verwendeten Sicherheitswerkzeuge ebenfalls kontinuierlich zum Einsatz kommen.

> Nur in einer idealen Welt schließen Sicherheitsteams alle gefundenen Schwachstellen.

Allerdings müssen die Security-Teams ihre Ressourcen auf die kritischsten Probleme konzentrieren, um mit der großen Anzahl an Warnungen der verschiedenen Werkzeuge mithalten zu können (hier sammeln sich leicht Tausende von Meldungen an, mehr dazu unter ix.de/zb98). Häufig wird dabei die Gewichtung (Severity), die die Tools den Warnungen zuordnen, oder die Severity-Definition der Benchmarks übernommen. Die Security-Teams beheben dann nur solche Probleme, die aus Sicht der Werkzeuge am kritischsten sind.

Warnungen um Warnungen

Das Problem dabei ist, dass die Werkzeuge die Cloud-Konstellation der Komponenten vernachlässigen, in denen Sicherheitslücken auftreten, und den eigentlichen Anwendungskontext nicht betrachten. Die Ergebnisse beziehen sich nur auf isolierte Komponenten. Das kann dazu führen, dass Sicherheitsteams ihre begrenzte Zeit mit dem Bearbeiten automatisch erzeugter Warnungen verbringen, die im Kontext der Cloud eigentlich unproblematisch sind, weil beispielsweise die betroffene Komponente von außen gar nicht erreichbar ist oder keinen Zugriff auf schützenswerte Daten hat. Die wirklich kritischen Warnungen gehen im ständigen Rauschen der Alarme unter und die Cloud-Anwendung bleibt angreifbar.

Security-Experten müssen also entscheiden, welche Schwachstelle zuerst anzugehen ist. Sind Anmeldedaten im Klartext in Umgebungsvariablen einer Komponente kritischer als ein offener SSH-Port 22 bei einer virtuellen Maschine? Zum einen benötigt diese Entscheidung Rücksprache mit dem Entwicklungsteam, zum anderen ist sie nicht immer eindeutig. Nimmt man hier die CIS-Benchmark als Grundlage

Capital One: Größter Hack in der Geschichte der US-Finanzbranche

Dass Cloud-Anwendungen nicht per se sicher sind, haben mehrere Angriffe in der Vergangenheit gezeigt. Ein Beispiel ist der Capital-One-Hack aus dem Jahr 2019, über den ausführliche Informationen vorliegen. Angreifern gelang es, in den Besitz der Identitäts-, Bonitäts- und Kreditkartendaten sowie Sozialversicherungsnummern von insgesamt über 100 Millionen Kunden der US-amerikanischen Bank zu kommen. Dem Unternehmen wurden im Anschluss fehlende Sicherheitsmaßnahmen vorgeworfen und der Hack kostete die Bank 80 Millionen Dollar an Schadensersatz. Ursache war eine Verkettung zweier Sicherheitslücken: eine Fehlkonfiguration in der Web-Application-Firewall (WAF) und eine überprivilegierte virtuelle Maschine in der Cloud, eine sogenannte Elastic-Cloud-Computing-Instanz (EC2).

Der Angriff verlief in mehreren Schritten (s. Abb. 1). Die folgende Kurzbeschreibung basiert auf einem Bericht des Massachusetts Institute of Technology (nachzulesen unter [ix.de/zb98](https://www.technologyreview.com/2019/12/11/401198/capital-one-hack/)):

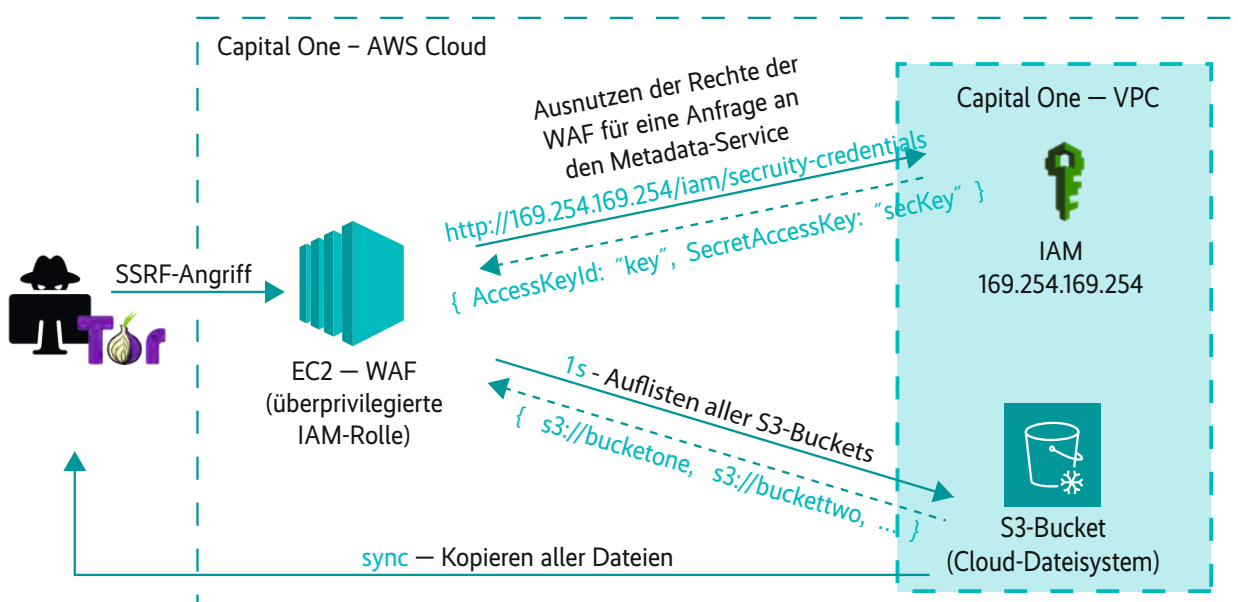
- > Das FBI und Capital One entdeckten mehrere Zugriffe auf die AWS-Cloud der Bank. Die IP-Adressen stammten aus dem Tor-Netzwerk und waren mit einem VPN-Service (IPredator) weiter verschleiert.
- > Den Angreifern gelang es, mit einem Server-Side-Request-Forgery-Angriff (SSRF) eine Rechteerhöhung (Privilege Escalation) zu provozieren. Damit ließen sich Daten von beliebigen

URLs innerhalb der AWS-Cloud mit den Rechten des Servers abgreifen und an die externen Angreifer zurücksenden. Ermöglicht hat das eine fehlerkonfigurierte Web-Application-Firewall (WAF), die die Anfrage nicht auf unzulässige URLs für den Download interner Dateien überprüfte.

- > Die Angreifer lenkten die Serverfunktion zum Datei-Download auf einen Metadaten-Service von AWS um (erreichbar über den Permalink mit der URL <http://169.254.169.254>). Dieser Dienst liefert temporäre Sicherheitsschlüssel, die sich von der Instanz nutzen lassen.
- > Durch die Kombination des SSRF-Angriffs auf die fehlerkonfigurierte WAF und eine überprivilegierte EC2-Instanz gelang es den Angreifern, über die URL <http://169.254.169.254/iam/security-credentials> die Schlüssel der EC2-Instanz (AccessKeyId und SecretAccessKey) herunterzuladen und damit vollen Zugriff auf die AWS-Umgebung zu erlangen.
- > Mithilfe der Zugangsschlüssel war es den Angreifern nun per AWS-Kommandozeile möglich, eine komplette Liste aller S3-Buckets zu erhalten und knapp 30 GByte an vertraulichen Kreditantragsdaten von mehr als 700 Buckets abzugreifen.

Der Capital-One-Hack entspricht dem typischen Muster von Cyber-Angriffen. Häufig nutzen die Angreifer mehrere Schwachstellen gleichzeitig aus, um sich Zugang zu einem System zu verschaffen.

Schema des Capital-One-Hacks



>> Größter Hack in der US-Finanzbranche: Der Capital One wurde zum Verhängnis, dass eine Web-Application Firewall, die eigentlich für zusätzliche Sicherheit sorgt, überprivilegierte Rechte hatte (Abb. 1).

einer Priorisierung, ist der offene Port das kritischere Problem – Klartext-Anmeldedaten hingegen fließen nicht einmal in den CIS-Score mit ein. Ausschließlich anhand der vorgeschlagenen Severity der Schwachstelle zu urteilen, würde hier leicht zu einer Fehleinschätzung führen. Bei der betroffenen virtuellen Maschine kann es sich beispielsweise um einen nicht genutzten Testserver handeln, der von außen nicht erreichbar ist und keinen Zugriff auf schützenswerte Daten hat. Somit wäre die Priorität dieser Schwachstelle im Vergleich zu den Klartext-Anmeldedaten niedrig.

Um eine verlässliche Bewertung des Risikos der Warnungen zu erhalten, sind die Warnungen unbedingt im Cloud-Kontext zu betrachten: Ist die betroffene Komponente von außen erreichbar und hat sie Zugriff auf sensible Daten? Diese Bewertung ist nach wie vor manuell durchzuführen.

Sicherheitsteams sollten entsprechend vorsichtig sein bei der Ressourcenplanung und Warnungen im Kontext der Applikation interpretieren, statt einzig auf Severity-Scores oder Guidelines wie den CIS-Benchmark zu vertrauen. Sie müssen sich ein Bild der gesamten Anwendungsstruktur machen und überprüfen, welche Sicherheitslücken akut von außen erreichbar sind und welche Lücken Zugriff auf sensible oder kritische Daten ermöglichen.

Schweregrad im Gesamtbild

Wie der Capital-One-Hack zeigt, beruhen Angriffe in der Cloud fast immer auf dem Zusammentreffen mehrerer Sicherheitslücken in verschiedenen Komponenten. Um diese nach ihrem tatsächlichen Schweregrad zu gewichten, ist es von Bedeutung, welche anderen Schwachstellen im Kontext der Komponente erreichbar und somit gemeinsam ausnutzbar sind.

Für die meisten CSPM- und CIEM-Werkzeuge ist es jedoch unmöglich zu beurteilen, ob eine inkorrekte Cloud-Konfiguration erst im Zusammenhang mit einer weiteren Sicherheitslücke kritisch sein kann. Die unverschlüsselte Datenablage in einem privaten Bereich ist folgenlos, sofern nicht zusätzlich eine Schwachstelle existiert, die den vermeintlich privaten Bereich kompromittiert.

Einzelne Fehlkonfigurationen oder falsche Berechtigungen können für sich genommen zunächst harmlos sein. Erst ihre Kombination ergibt ein mögliches Angriffsszenario. Hier ist es nötig, sich eine Verkettung von Rechten, Datenflüssen und den eigentlichen Auswirkungen im Falle eines Angriffs

anzusehen. Daher ist es notwendig, die Warnungen aller genannten Werkzeuge (CSPM, CIEM, SCA und SAST) insgesamt zu betrachten.

Kontext bringt Gewissheit

Moderne Cloud-Technologien bieten viele Vorteile, sie beschleunigen den Entwicklungsprozess und reduzieren dabei die Kosten für das Betreiben der Software. Durch die schnellere Entwicklung haben sich allerdings auch die Anforderungen an die Sicherheit stark geändert. Um aktuellen Best Practices zu folgen, gibt es unterschiedliche Werkzeuge wie CSPM, CIEM, SAST und SCA, die kontinuierlich während des Softwareentwickelns die Sicherheit bewerten. Für eine präzise und effektive Priorisierung der Warnungen ist es notwendig, den Kontext der Warnungen in der Cloud zu betrachten und Warnhinweise zueinander in Bezug zu setzen. Für Security-

Fachkräfte ergibt sich dadurch ein unabdingbarer Mehraufwand für eine sichere Softwareentwicklung. Moderne Cloud-Security-Werkzeuge bieten erste Maßnahmen, diesen Mehraufwand durch eine automatisierte Kontextualisierung zu reduzieren.

Quellen

Weiterführendes Material steht unter ix.de/zb98 bereit.



Manuel Benz

ist Co-Founder und CTO von CodeShield und bringt mehrjährige Erfahrung in IT-Security und dynamischer Codeanalyse mit.



Andreas Dann

ist Experte im Bereich SCA und ist seit 2020 Mitgründer sowie COO von CodeShield.



Dr. Johannes Späth

ist Co-Founder und CEO von CodeShield und hat Industrie- und Forschungserfahrung im Bereich SAST.

betterCode() building IoT ct <webdev> » Continuous Lifecycle » [Container Conf]

data2day enterPy // heise devSec() heise MacDev Herbstcampus



ONLINE & VOR ORT DEVELOPER März – Juni '22 KONFERENZEN

Weitere Informationen unter:
heise.de/developer



Veranstalter:



@ heise Developer



dpunkt.verlag

- **31.03.2022**
Team Up! (online)
teams.inside-agile.de
- **07.04.2022**
Mastering GitOps (online)
mastering-gitops.de
- **26.04.2022**
heise devSec – Thementag
Supply Chain Security (online)
heise-devsec.de/supplychain.php
- **27.04.2022**
betterCode() API (online)
api.bettercode.eu
- **02.05. – 05.05.2022**
inside agile Konferenz (Darmstadt)
konf.inside-agile.de
- **CfP bis 02.05.2022**
data2day (19. – 21.09. in Karlsruhe)
data2day.de
- **09. – 11.05.2022**
building IoT (Karlsruhe)
buildingiot.de
- **24.05.2022**
betterCode() Flutter (online)
flutter.bettercode.eu
- **01. – 03.06.2022**
Minds Mastering Machines (Karlsruhe)
m3-konferenz.de
- **16.06.2022**
IIoT Conference (München)
iiot-conference.de
- **22. – 23.06.2022**
enterJS (Darmstadt)
enterjs.de
- **30.06.2022**
Product Owner Day (online)
pod.inside-agile.de

AKKA

PASSION FOR
TECHNOLOGIES

WIR SUCHEN EXPERTEN FÜR DEN BEREICH DIGITAL ENGINEERING

Gestalte gemeinsam mit AKKA die digitale Zukunft und werde Teil unseres weltweiten Netzwerks von Experten und motivierten Teams. Bei uns erhältst Du die Möglichkeit einer spannenden Karriere mit viel Gestaltungsspielraum in einem hoch innovativen und technologischen Umfeld. Bringe Dein Talent ein und wir bringen es gemeinsam zur Höchstleistung.

UNSERE TOP JOBS FÜR DICH:



Software-Entwickler /
DevOps



Data Engineer /
Architect



Java Software-
Entwickler

Weitere spannende Einstiegsmöglichkeiten findest Du unter
www.akka-technologies.com/karriere