


# Wie eingebautes Misstrauen KI-Systeme sicherer machen kann

[Originalartikel](#)

[Backup](#)

**Wie eingebautes Misstrauen KI-Systeme sicherer machen kann**  
Karen Hao



(Bild: ACM)

**Ayanna Howard, Robotikerin und P&#228;dagogin, h&#228;lt es f&#252;r gef&#228;hrlich, automatisierten Systemen zu vertrauen. Im TR-Interview erl&#228;utert sie L&#246;sungsm&#246;glichkeiten.**

Ayanna Howard sieht den Nutzen von Robotern und KI vor allem darin, Menschen zu helfen. In ihrer fast 30-j&#228;hrigen Karriere hat sie unz&#228;hlige Roboter gebaut: zur Erforschung des Mars, zur Reinigung von Sonderm&#252;ll und zur Unterst&#252;tzung von Kindern, die bestimmte Hilfen brauchen. Dabei hat sie eine beeindruckende Palette an Techniken entwickelt zur Roboter manipulation, autonomen Navigation und Computer Vision. Und sie ist f&#252;hrend bei der Erforschung eines typisch menschlichen Verhaltens: Wir vertrauen automatisierten Systemen zu sehr.

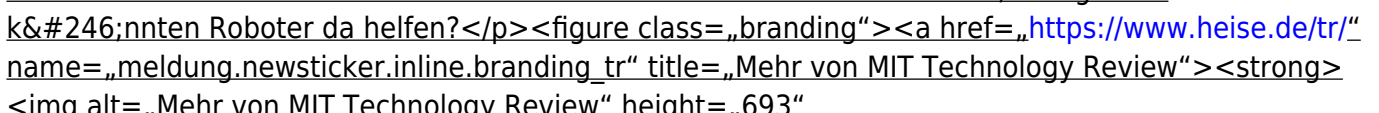
Im M&#228;rz trat sie nach 16 Jahren als Professorin am Georgia Institute of Technology als erste Frau ihre Stelle als Dekanin des College of Engineering an der Ohio State University an. Am Tag der ACM-Preisverleihung sprach TR mit ihr &#252;ber ihre Karriere und ihre neueste Forschung.

**TR: Sie verwenden f&#252;r Ihre Forschung den Begriff „humanisierte Intelligenz“ anstatt „k&#252;nstliche Intelligenz&#220;. Warum?**

**Ayanna Howard:** Ja, 2004 habe ich den Begriff in einer Ver&#246;ffentlichung verwendet. Denn wir wollen Robotik und KI-Systeme ja nicht losgel&#246;st von Menschen erschaffen. Wir bauen dabei auf die Erfahrung, die Daten und die Eingaben von Menschen auf. „K&#252;nstliche Intelligenz“ impliziert, dass es sich um eine andere Art von Intelligenz handelt. „Humanisierte Intelligenz“ besagt hingegen, dass diese Intelligenzen von Menschen geschaffene Konstrukte sind.

**Wie sind Sie zu dieser Arbeit gekommen?**

In erster Linie durch meine Doktorarbeit. Damals habe ich einen Roboter dazu trainiert, Gefahren in einem Krankenhaus zu verringern. Zu der Zeit wurden etwa noch Nadeln in denselben M&#252;ll geworfen wie alles andere, und es kam vor, dass dadurch Krankenhausmitarbeiter krank wurden. Also habe ich &#252;berlegt: Wie k&#246;nnten Roboter da helfen?



Mehr von MIT Technology Review

<https://static.wallabag.it/7862d1b7aff4c3b00f37212fefade4e0e2c4cf00/64656e6965643a646174613a696d6167652f7376672b786d6c2c253343737667253230786d6c6e733d27687474703a2f2f77>

7772e77332e6f72672f323030302f7376672725323077696474683d273639367078272532306865696768743d2733393170782725323076696577426f783d2730253230302532303639362532303339312725334525334372656374253230783d273027253230793d27302725323077696474683d27363936272532306865696768743d273339312725323066696c6c3d27253233663266326632272533452533432f726563742533452533432f737667253345/“ class=„c1“ width=„1200“ referrerpolicy=„no-referrer“ /><img alt=„Mehr von MIT Technology Review“ class=„a-u-hide-from-tablet c2“ src=„https://heise.cloudimg.io/width/1200/q50.png-lossy-50.webp-lossy-50.foil1/\_www-heise-de/\_Magazin-Banner/tr\_mobil.jpg“ srcset=„https://heise.cloudimg.io/width/2400/q30.png-lossy-30.webp-lossy-30.foil1/\_www-heise-de/\_Magazin-Banner/tr\_mobil.jpg 2x“ referrerpolicy=„no-referrer“ /> <img alt=„Mehr von MIT Technology Review“ height=„500“ src=„https://static.wallabag.it/7862d1b7aff4c3b00f37212fefade4e0e2c4cf00/64656e6965643a646174613a696d6167652f7376672b786d6c2c253343737667253230786d6c6e733d27687474703a2f2f777772e77332e6f72672f323030302f7376672725323077696474683d273639367078272532306865696768743d2733393170782725323076696577426f783d2730253230302532303639362532303339312725334525334372656374253230783d273027253230793d27302725323077696474683d27363936272532306865696768743d273339312725323066696c6c3d27253233663266326632272533452533432f726563742533452533432f737667253345/“ class=„c3“ width=„1830“ referrerpolicy=„no-referrer“ /><img alt=„Mehr von MIT Technology Review“ class=„a-u-show-from-tablet c2“ src=„https://heise.cloudimg.io/width/1830/q50.png-lossy-50.webp-lossy-50.foil1/\_www-heise-de/\_Magazin-Banner/tr\_desktop.jpg“ srcset=„https://heise.cloudimg.io/width/3660/q30.png-lossy-30.webp-lossy-30.foil1/\_www-heise-de/\_Magazin-Banner/tr\_desktop.jpg 2x“ referrerpolicy=„no-referrer“ /> [1]</strong></a></figure><p>Es ging mir damals schon um Roboter, die f&#252;r Menschen n&#252;tzlich sind. Wir wussten noch nicht, wie man Roboter baut, die solche Aufgaben erledigen k&#246;nnen. Aber Menschen verrichten die T&#228;tigkeiten st&#228;ndig, also sollten wir sie nachahmen. So fing es an.</p><p>Mit der NASA habe ich dann an der Navigation von Mars-Rovern gearbeitet. Da war es das Gleiche: Wissenschaftler k&#246;nnen das sehr, sehr gut. Also lie&#223; ich Wissenschaftler die Rover steuern und schaute mir an, was sie auf den Kameras der Rover sahen und wie sie darauf reagierten. Das war immer mein Thema: Warum gehe ich nicht einfach zu menschlichen Experten, kodiere das, was sie tun, in einen Algorithmus und bringe dann den Roboter dazu, das zu verstehen und umzusetzen?</p><p class=„frage rteabs-frage“><strong>Haben andere Leute damals auch so gedacht oder waren Sie eine Au&#223;enseiterin?</strong></p><p class=„antwort rteabs-antwort“>Eine total verr&#252;ckte Au&#223;enseiterin. Ich habe die Dinge anders betrachtet als alle anderen. Und damals gab es noch keine Anleitung, wie man so eine Forschung betreibt. Wenn ich jetzt zur&#252;ckblicke, w&#252;rde ich es nach all den Erfahrungen in der Praxis ganz anders machen.</p><p class=„frage rteabs-frage“>Seit wann haben Sie, statt weiter Roboter f&#252;r Menschen zu konstruieren, die Beziehung zwischen Robotern und Menschen in den Mittelpunkt Ihrer Arbeit gestellt?</p><p class=„antwort rteabs-antwort“>Angesto&#223;en wurde das durch eine Studie, in der wir herausfinden wollten, ob sich Menschen in riskanten Situationen der F&#252;hrung eines Roboters anvertrauen w&#252;rden. Also brachten wir die Teilnehmer in ein verlassenes B&#252;rogeb&#228;ude, in das sie ein Roboter hineinlie&#223;. Wir haben ihnen erz&#228;hlt, es w&#252;rde sich um eine Umfrage handeln. W&#228;hrend sie dort drinnen waren, haben wir in dem Geb&#228;ude Rauch entfacht und den Feueralarm ausgel&#246;st.</p><p>Wir wollten sehen, ob die Menschen beim Verlassen des Geb&#228;udes zur Eingangst&#252;r gehen oder dem Hinweisschild zum Notausgang folgen, oder ob sie sich dem Roboter anschlie&#223;en w&#252;rden, der sich in eine andere Richtung bewegt hat.</p><p>Wir dachten, die Leute w&#252;rden zur Eingangst&#252;r gehen, weil sie dort hereingekommen waren. Fr&#252;here Forschungen haben gezeigt, dass Menschen in einer Notsituation meist vertraute Wege gehen. Oder dass sie den Schildern folgen w&#252;rden, weil das ein einstudiertes Verhalten ist. Aber die

Teilnehmer taten das nicht. Sie folgten tatsächlich dem Roboter. Dann haben wir einige Fehler eingebaut. Wir ließen den Roboter ausfallen, wir ließen ihn im Kreis fahren oder einen Weg fahren, bei dem man erst Mal beiseiteschieben musste. Wir dachten, irgendwann würde jemand sagen: „Lass uns zur Eingangstür oder zum Ausgangsschild dort gehen.“ Aber sie folgten ihm buchstäblich bis zum Ende des Experiments. Damit lagen wir zum ersten Mal mit unseren Hypothesen völlig falsch. Wir konnten kaum glauben, dass die Leute einem System so vertrauen. Das ist interessant und faszinierend, aber es ist ein Problem.

Haben Sie dieses Phänomen seitdem in der realen Welt auch beobachten können?

Jedes Mal bei einem Tesla-Unfall. Besonders bei den ersten Autopilot-Modellen. Die Leute vertrauen diesen Systemen zu sehr. So nach dem Motto: „Jetzt musst du das Lenkrad noch mindestens fünf Sekunden lang halten. Wenn du die Hand nicht am Lenkrad hast, schaltet sich das System ab.“

Aber sie haben nie mit mir oder meinem Team darüber gesprochen, weil das gar nicht funktioniert. Und das ist so, weil sich das System so leicht manipulieren lässt.

Wenn Sie auf Ihr Handy schauen und dann den Piepston des Tesla hören, nehmen Sie einfach die Hand hoch, richtig? Das ist unterbewusst so. Sie sind aber immer noch unaufmerksam. Und zwar deshalb, weil Sie denken, dass das System einfach weiter läuft und dass Sie eigentlich immer noch das tun können, was Sie gerade getan haben, ein Buch lesen, fernsehen oder eben auf Ihr Handy schauen. Es funktioniert also nicht, weil sich für Sie das Risiko, die Unsicherheit oder das Misstrauen gegenüber dem automatischen System nicht erhöhen haben. Es reicht nicht, dass Sie sich wieder mit der Situation beschäftigen.

Sie meinen, dass man aktiv Misstrauen in das System einbauen müsste, um es sicherer zu machen?

Genau das. Wir machen gerade ein Experiment zum Thema Denial of Service. Wir ringen allerdings noch mit ethischen Fragen. Denn sobald wir darüber sprechen und die Ergebnisse veröffentlichen, werden wir erklären müssen, warum man manchmal der KI die Möglichkeit geben will, ihren Dienst auch zu verweigern. Was etwa, wenn ein Dienst verweigert wird, wenn jemand ihn wirklich braucht?

Aber nehmen wir mal den Autopilot des Tesla als Beispiel. Denial of Service wäre: Ich erstelle ein Profil Ihres Vertrauens, darauf basierend, wie oft Sie das Lenkrad deaktivieren oder das Steuer loslassen. Anhand dieser Profile kann ich dann modellieren, ab wann jemand dem System völlig vertraut. Wir haben das nicht mit Tesla-Daten, sondern mit unseren eigenen gemacht. Wenn Sie dann das nächste Mal das Autofahren, würden sich in bestimmten Situationen ein Denial-of-Service einschalten. Das heißt, Sie haben in einer Zeitspanne von X keinen Zugriff auf das System.

So, als würden Sie einen Teenager durch das Wegnehmen seines Smartphones bestrafen, wenn Sie es anders nicht erreichen, dass er das macht, was Sie wollen.

Welche anderen Mechanismen haben Sie erforscht, um das Misstrauen in Systeme zu erhöhen?

Zum Beispiel durch eine Art erklärende KI. Das System erklärt selbst seine Risiken oder Unsicherheiten. Denn alle diese Systeme bergen Unsicherheiten - keines ist 100%ig sicher. Und ein System weiß, wann es unsicher ist. Es könnte also diese Informationen in einer verständlichen Form bereitstellen, so dass die Nutzer ihr Verhalten ändern.

Ein Beispiel: Ich bin ein selbstfahrendes Auto und habe alle Karteninformationen. Ich weiß, dass bestimmte Kreuzungen unfallträchtiger sind als andere. Wenn wir uns einer von ihnen nähern, würde ich sagen: „Wir nähern uns einer Kreuzung, an der letztes Jahr 10 Menschen gestorben sind. Dann kann der Fahrer sich sagen: „Oh, vielleicht sollte ich aufmerksamer sein.“

So viel zu Ihren Bedenken bezüglich unserer Tendenz, diesen Systemen zu sehr zu vertrauen. Kann das denn nicht auch Vorteile haben?

Positiv ist, dass automatisierte Systeme im Allgemeinen besser sind als Menschen. Ich selbst würde in manchen Situationen lieber mit einem KI-System interagieren als mit bestimmten Menschen. Die Systeme haben mehr Daten; sie sind genauer. Vor allem, wenn

man sie mit einer unerfahrenen Person vergleicht

„Sie haben bereits vor 20 Jahren ein Mentorenprogramm f&#252;r M&#228;dchen an der High School ins Leben gerufen, lange bevor sich andere Gedanken dazu gemacht haben. Warum ist das f&#252;r Sie wichtig, und warum f&#252;r die Branche?“

„Es gab in meinem Leben immer Personen, die mir den Zugang zu Technik und Informatik erm&#246;glicht haben. Deshalb hatte ich auch sp&#228;ter nie ein Problem damit, in diesem Berufsfeld zu arbeiten. Ich wollte gern f&#252;r andere das Gleiche tun. Sp&#228;ter fielen mir die vielen Menschen auf, die nicht so aussahen wie ich. Da wurde mir klar: Moment, hier gibt es definitiv ein Problem, denn die Leute haben einfach keine Vorbilder, und sie haben keinen Zugang zu bestimmten Kreisen.“

„Dabei ist es in der Praxis so wichtig, dass jeder eine andere Erfahrung mitbringt. So wie ich bereits &#252;ber Mensch-Roboter-Interaktion nachgedacht habe, bevor das &#252;berhaupt ein Thema war. Nicht, weil ich brilliant war. Sondern einfach, weil ich das Problem auf eine andere Art und Weise betrachtet habe. Spreche ich mit jemand, der eine andere Sicht auf die Dinge hat, k&#246;nnen wir uns zusammentun und sagen: „Lass uns versuchen, das Beste aus beiden Welten zu kombinieren.““

„Zum Beispiel t&#246;ten Airbags h&#228;ufiger Frauen und Kinder als M&#228;nner. Warum? Nun, wahrscheinlich gab es damals niemanden, der sagte: „Hey, warum testen wir das nicht auch an Frauen auf dem Vordersitz?“ Es gibt viele Probleme, durch die bestimmte Gruppen von Menschen in Gefahr gebracht oder sogar get&#246;tet wurden. Ich glaube, dass einfach nicht genug Leute gesagt haben: „Hey, habt ihr auch dar&#252;ber nachgedacht?““

„einfach aufgrund ihrer speziellen Erfahrungen und ihres anderen Umfelds.“

„Wie hoffen Sie, dass sich die KI- und Robotikforschung im Laufe der Zeit weiterentwickeln wird? Was ist Ihre Vision?“

„Programmieren kann so ziemlich jeder. Es gibt jetzt so viele Organisationen wie Code.org. Die Ressourcen und Werkzeuge sind da. Ich w&#252;rde gerne eines Tages einen Studenten fragen: „Kennst du dich mit KI und maschinellem Lernen aus?“ und er sagt: „Dr. H, das mache ich schon seit der dritten Klasse!“ Das w&#228;re wunderbar. Nat&#252;rlich m&#252;sste mir dann wohl einen neuen Job suchen, aber das ist eine ganz andere Geschichte.“

„Wenn man die richtigen Werkzeuge und das Wissen hat, mit KI und maschinellem Lernen umzugehen, kann man sich seine eigenen Jobs, seine eigene Zukunft, seine eigene L&#246;sung selbst erschaffen. Das w&#228;re mein Traum.“

URL dieses Artikels:

<https://www.heise.de/-6049154>

Links in diesem Artikel:

[1] <https://www.heise.de/tr/>

[2] <mailto:bsc@heise.de>

Copyright &#169; 2021 Heise Medien

From: <https://schnipsel.qgelm.de/> - Qgelm

Permanent link: <https://schnipsel.qgelm.de/doku.php?id=wallabag:wb2wie-eingebautes-misstrauen-ki-systeme-sicherer-machen-kann>

Last update: 2025/06/27 11:17

