

Verantwortung und Erklärbarkeit von Künstlicher Intelligenz Warum KI verständlich sein muss

Impulsvortrag

05. Februar 2018

Timo Speith, B.A., M.Sc.
Universität des Saarlandes,
Theoretische Philosophie

Timo Speith und Kevin Baum
Universität des Saarlandes



Theoretische Philosophie

Was hat KI mit Verantwortung zu tun?

Wo finden wir KI?

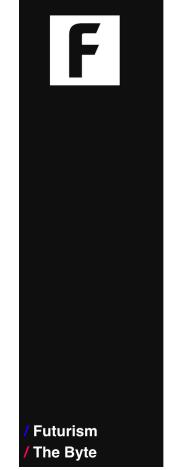


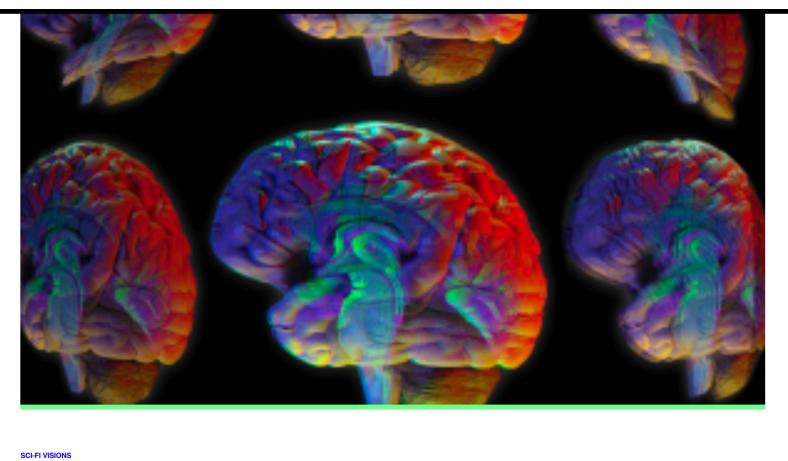
22. Oktober 2018, 05:48 Uhr Digitalisierung

Arbeit aus dem Automaten



https://www.sueddeutsche.de/digital/digitalisierung-arbeitslosigkeit-jobcenter-1.4178635





Should Coma Patients Live or Die? Machine Learning Will Help Decide.

An algorithm is helping Chinese researchers determine if a coma patient will wake up again.

Victor Tangermann September 27th 2018

https://futurism.com/machine-learning-coma-patients-live

Lass Siri deine Sätze für dich beenden.



Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica
May 23, 2016

https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing

Timo Speith und Kevin Baum

Warum KI verständlich sein muss

Folgenschwere und Verantwortung

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks. by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica







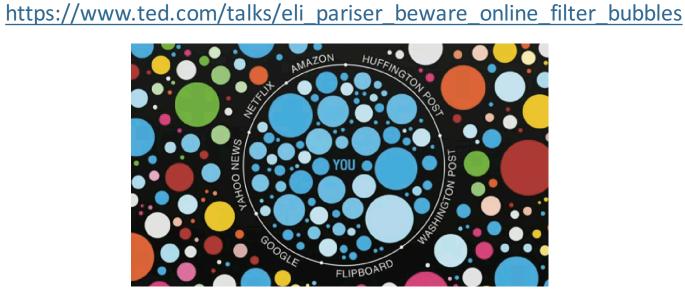


Musikauswahl

Suchanfragen zu Banalitäten

Einkaufsempfehlungen bei Amazon

Saugroboter



Vorauswahl von News und Posts in sozialen Netzwerken

Visumvergabe

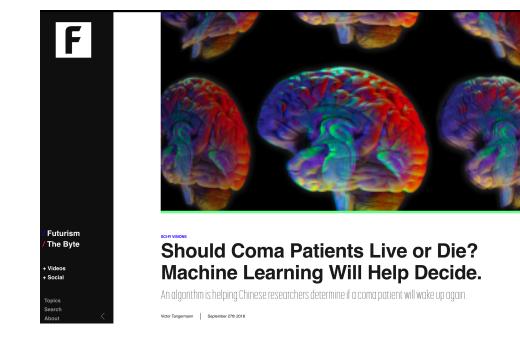
Kreditwürdigkeit

Strafzumessung

Abschalten lebensnotwendiger Geräte

Jobvergabe und Bewertung Arbeitssuchender

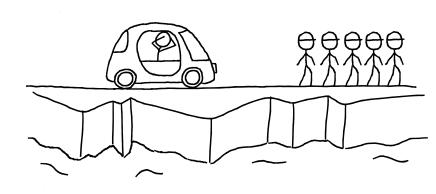
Autonome Fahrzeuge in Dilemmasituationen





22. Oktober 2018, 05:48 Uhr Digitalisierung Arbeit aus dem Automaten





Folgenschwer

>Folgentrivial<

Die Frage der Verantwortung stellt sich insbesondere in Fällen, in denen folgenschwere Entscheidungen getroffen werden müssen.

Folgenschwere und Verantwortung

Machine Bias

There's software used across the country to predict future criminals. And it's biased against blacks.

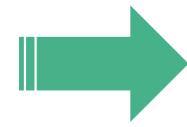
by Julia Anawin Jeff Larson, Surva Mattu and Lauren Kirchner, ProPublica

by Julia Angwin, Jeff Larson, Surya Mattu and Lauren Kirchner, ProPublica May 23, 2016

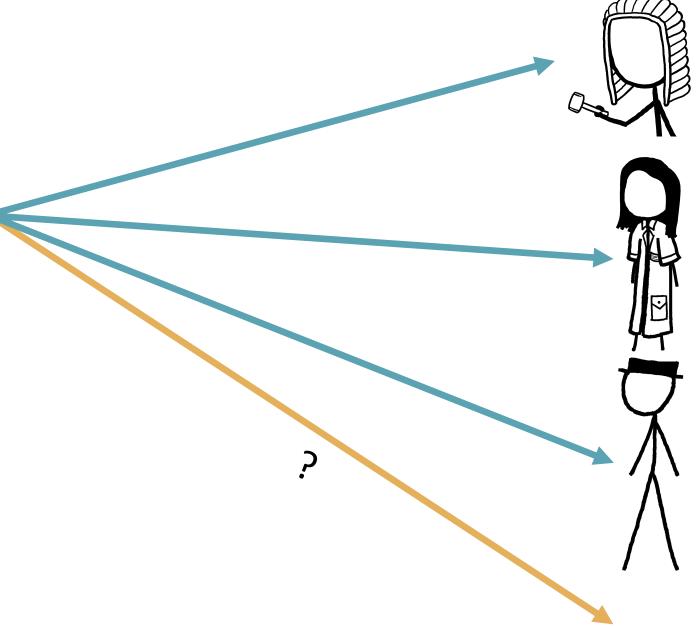
In vielen solcher Fälle haben wir *Experten* und oft auch Zeit, sie miteinzubeziehen...



Standardreaktion: *Nicht ersetzen,* unterstützen!



Allokation von Verantwortung scheint dann klar: Bei den Experten.

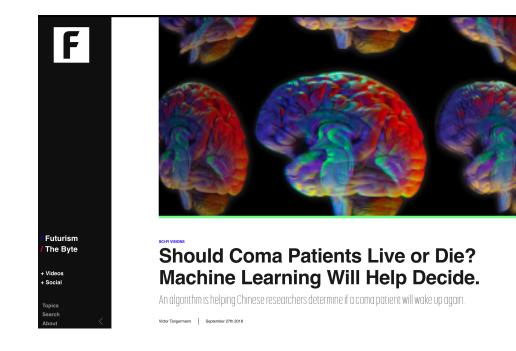


Strafzumessung

Abschalten lebensnotwendiger Geräte

Jobvergabe und
Bewertung
Arbeitssuchender

Autonome Fahrzeuge in Dilemmasituationen

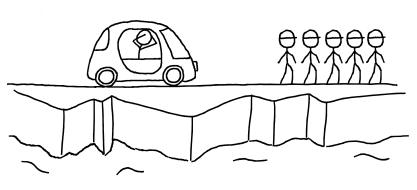




22. Oktober 2018, 05:48 Uhr Digitalisierung

Arbeit aus dem Automaten



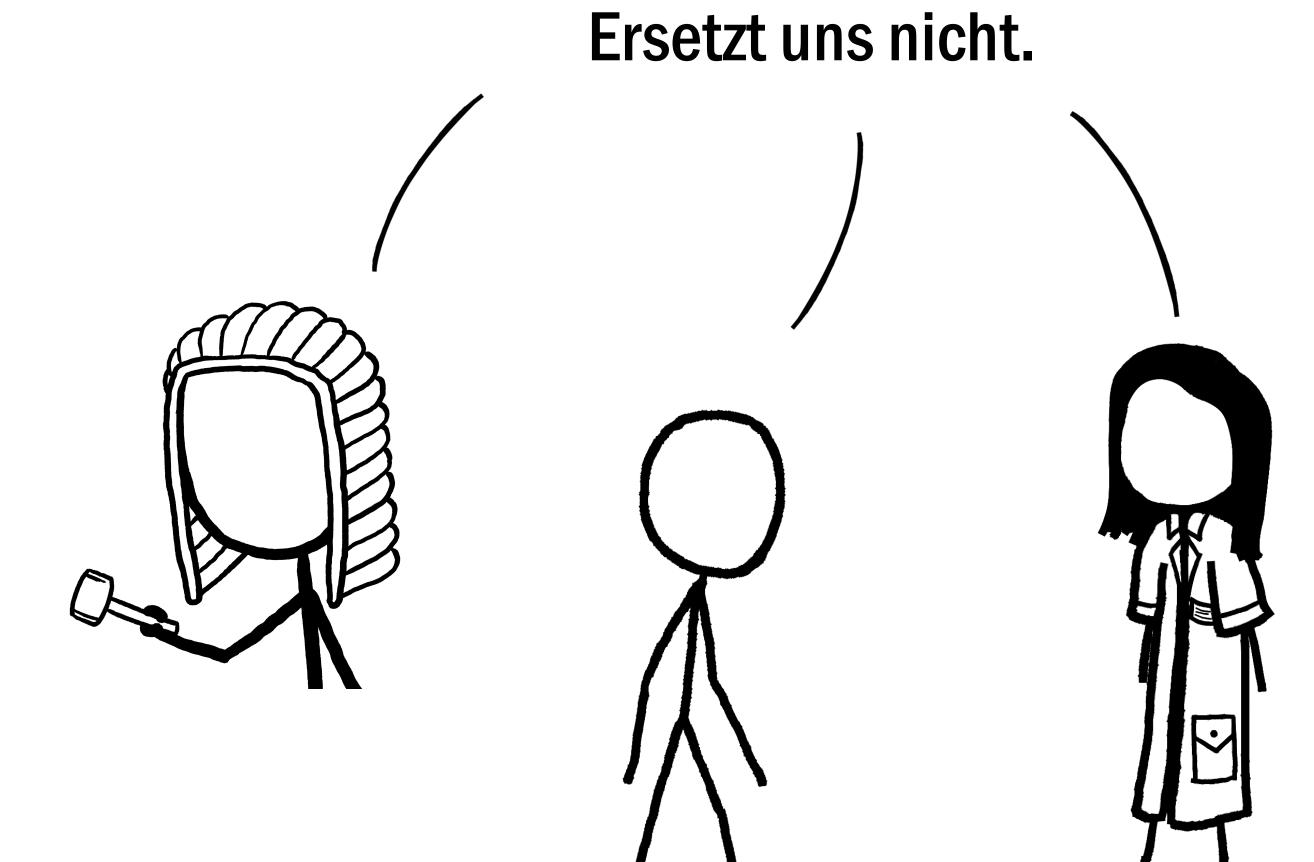


Folgenschwer

>Folgentrivial<

Die **Frage der Verantwortung** stellt sich insbesondere in Fällen, in denen folgenschwere Entscheidungen getroffen werden müssen.

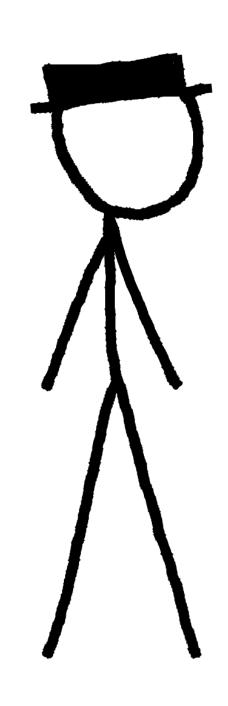
Fazit 1



Sollten wir deshalb bei folgenschweren Entscheidungen gänzlich auf KI verzichten?

Gründe für Unterstützung: Überforderung und schlechte Heuristiken

Das ist Rudolph, der Personaler. Er hat einen Job zu besetzen.



Eine solche Heuristik ist weder fair noch wünschenswert.

Kann KI helfen?



Timo Speith und Kevin Baum

SAN FRANCISCO (Reuters) - Amazon.com Inc's (AMZN.O) machine-learning specialists uncovered a big problem: their new recruiting engine did not like women

Das ist Susanne, die sich auf den Job beworben hat.

Rudolph, der rassistische Personaler

Das ist Rudolph, der Personaler, der auch Rassist ist. Er hat einen Job zu besetzen.

Dieser Zustand ist

Kann KI helfen?



Sie ist am besten qualifiziert.

Jeffrey Dastin

inakzeptabel.

Wer bekommt den Job wohl nicht?

Gründe für Unterstützung: Kognitive Verzerrungen

Lunch Luck

Published online 11 April 2011 | Nature | doi:10.1038/news.2011.227

News

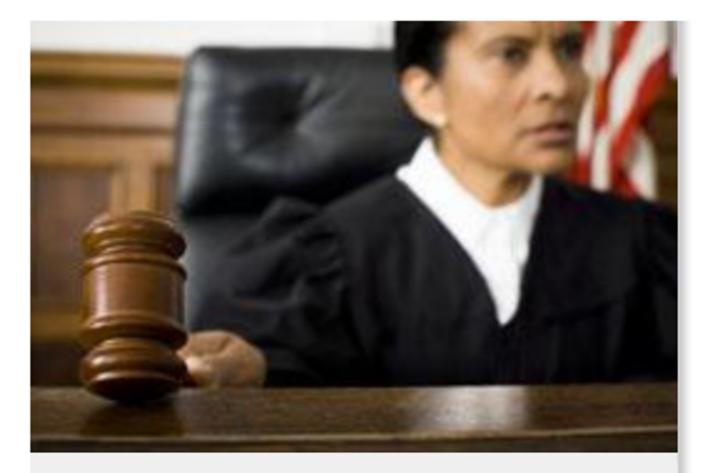
Hungry judges dispense rough justice

When they need a break, decision-makers gravitate towards the easy option.

Zoë Corbyn

A prisoner's chance of parole depends on when the judge hearing the case last took a break, say researchers who have studied decisions in Israeli courts. As judges tire and get hungry, they slip towards the easy option of denying parole, say the researchers.

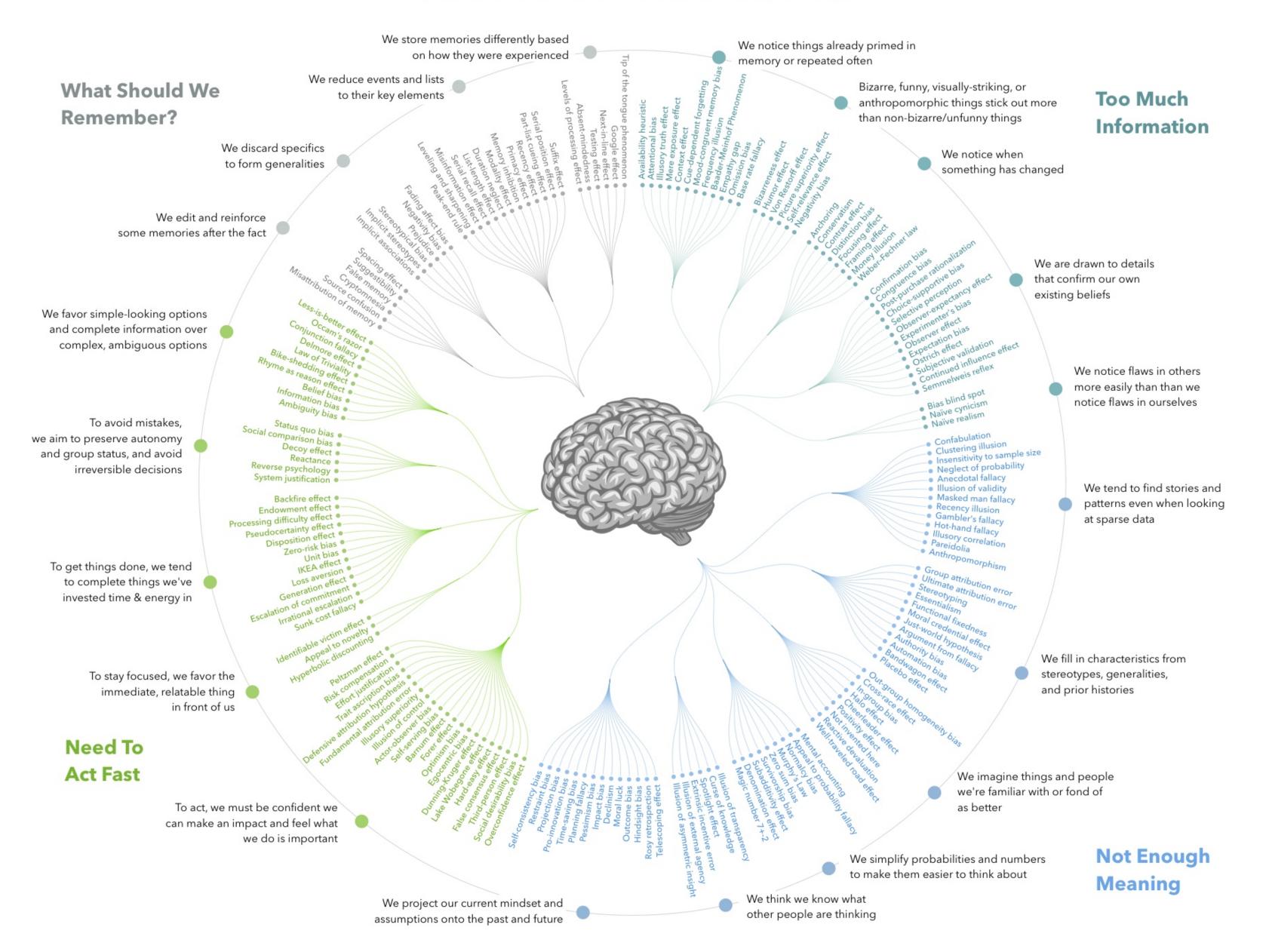
The bias could apply in any situation in which people



A judge is less likely to parole a prisoner at the end of a session than at the beginning.

Punchstock

http://www.nature.com/news/2011/110411/full/news.2011.227.html

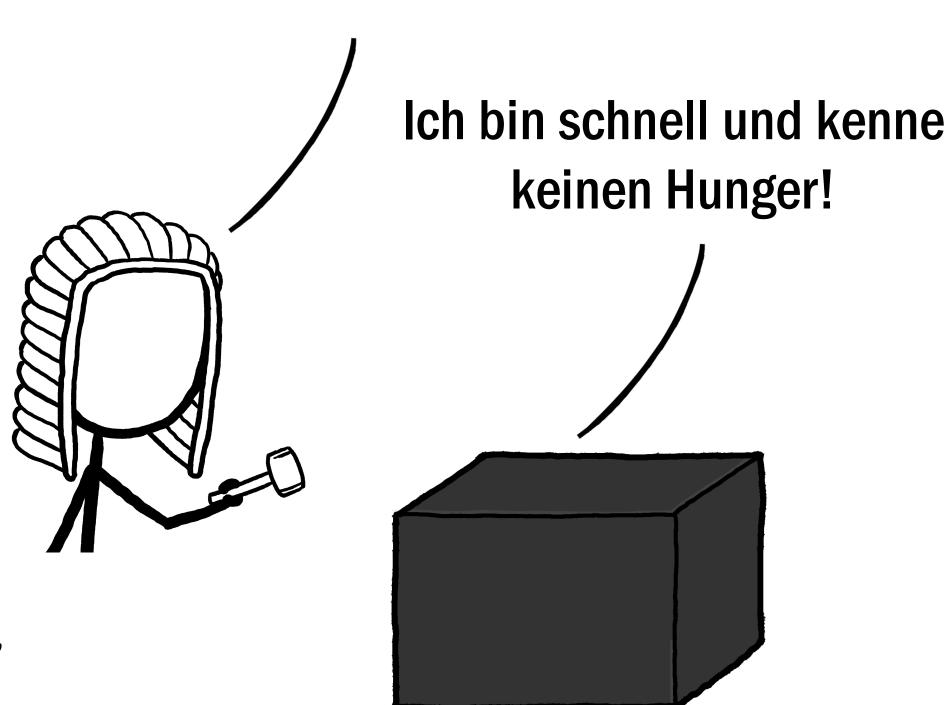


In vielen Gebieten der vielversprechendste Ansatz

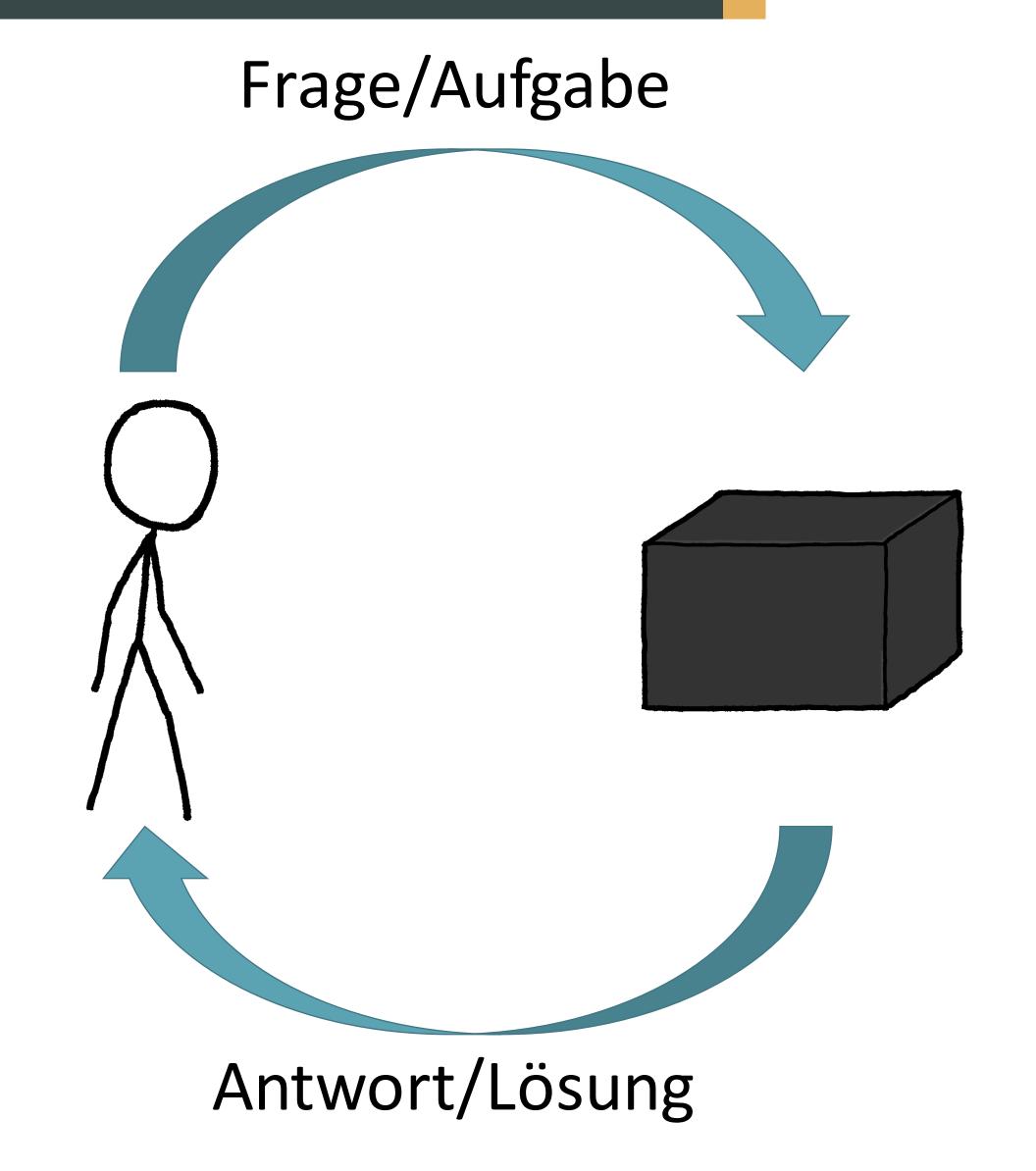
- Wir nutzen beides:
 - Die Möglichkeiten der Kl
 - Prädiktive Kraft/Reduzierung von Unsicherheit
 - >Unbestechlichkeit
 - Systematik
 - Die Fähigkeiten der Menschen
 - Können hinterfragen
 - »Sehen« u.U. mehr

(z.B. Gründe und Kausalität statt bloße Korrelation, schwer bis unmöglich messbare Eigenschaften, Einzelfallspezifisches)

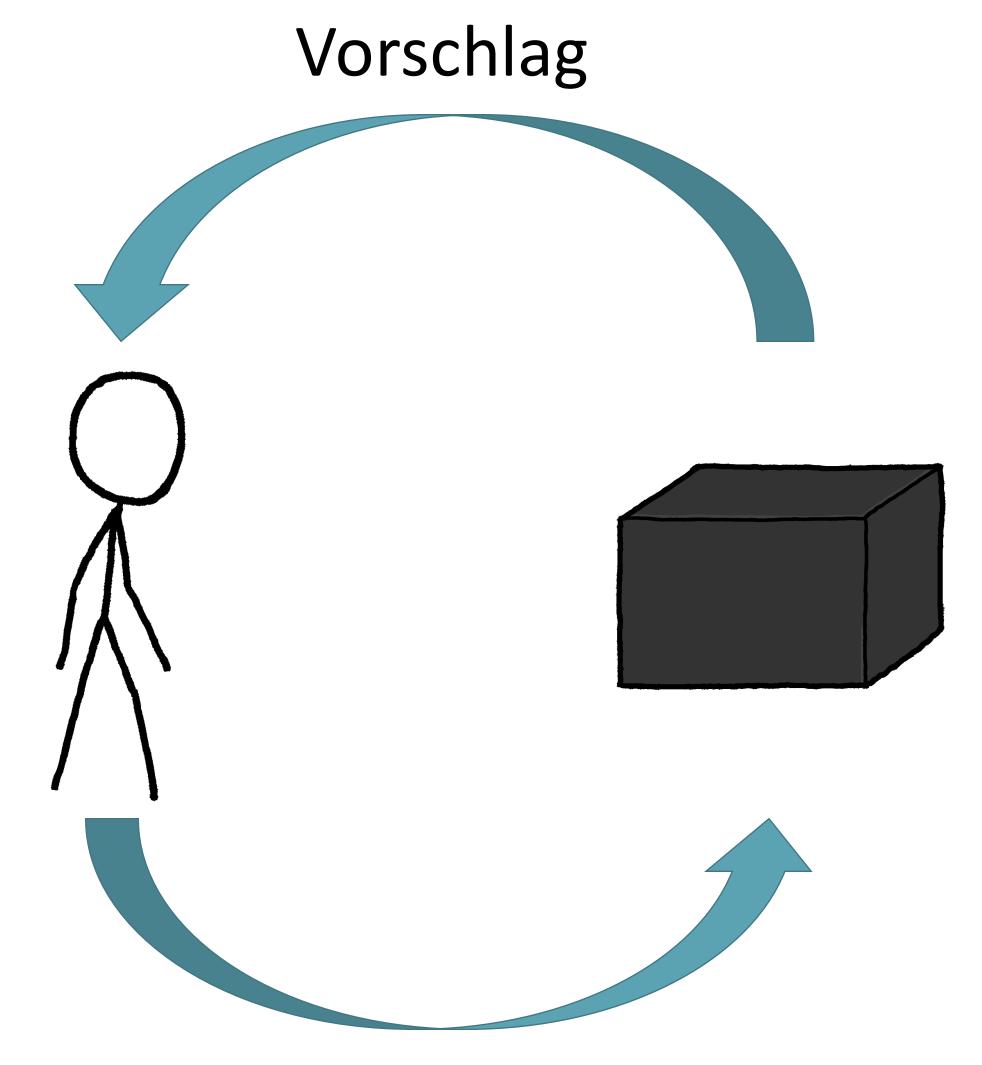
Ich kenne zwar Hunger, bin aber Experte im echten Denken und Abwägen!



Human in the Loop: Pull



Human in the Loop: Push

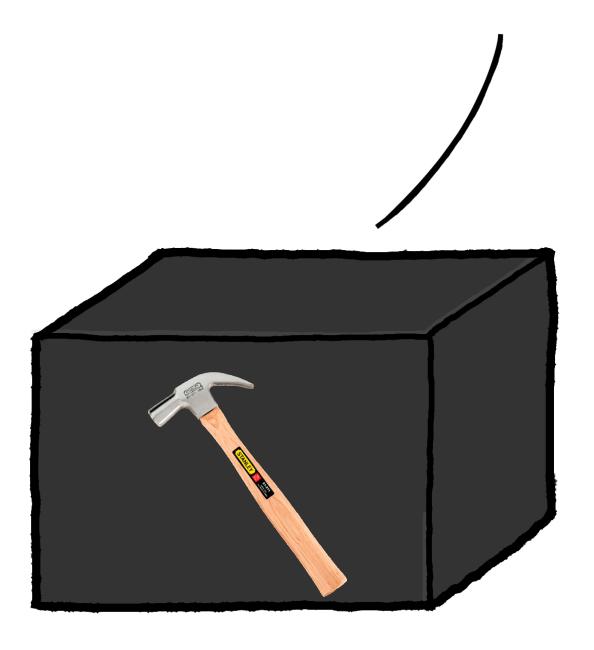


Auswahl/Bestätigung

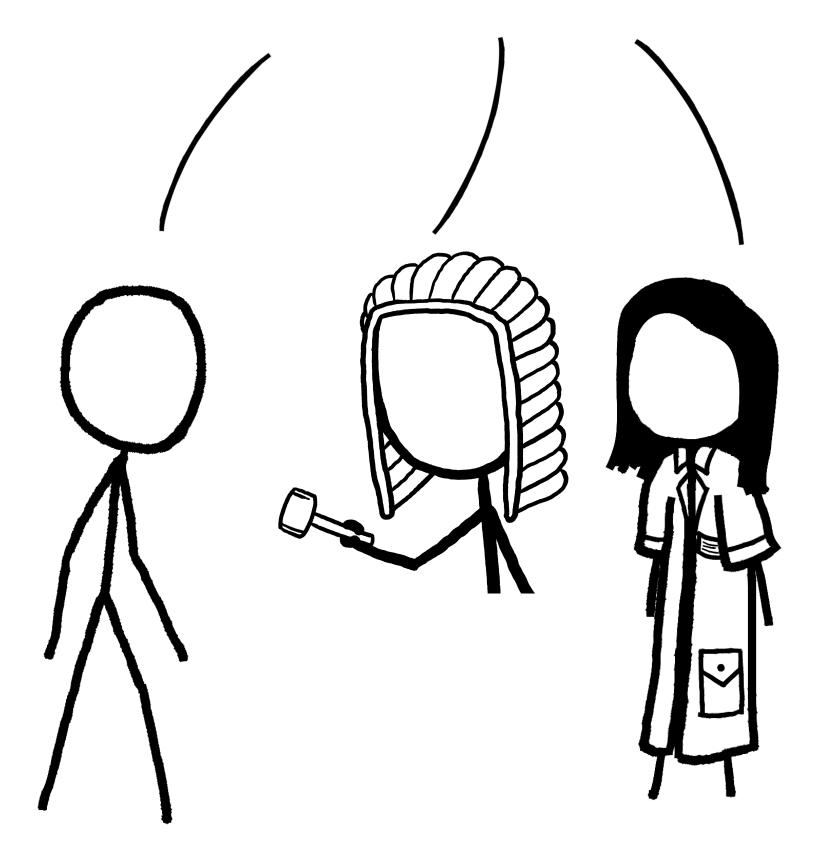
Fazit 2

Zusammen sind wir besser.

Lasst mich helfen!



Wir sind wirklich nicht perfekt!



16

Warum es so einfach nicht ist:

Das Recht auf Anfechtbarkeit und >eingebaute Werte<

Beispiel COMPAS

https://www.wired.com/2017/04/courts-using-ai-sentence-criminals-must-stop-now/

https://www.nytimes.com/2017/05/01/us/politics/sent-to-prison-by-a-software-programs-secret-algorithms.html
https://en.wikipedia.org/wiki/Loomis v. Wisconsin

COMPAS (Correctional Offender Management Profiling for Alternative Sanctions)

- Proprietärer, closed-source Risikobewertungsalgorithmus (kein ML), der in vielen amerikanischen Jurisdiktionen zum Einsatz kommt.
- Konkreter Fall:
 - Eric L. Loomis wurde zu sechs Jahren Gefängnis verurteilt.
 - Der Algorithmus >urteilte<: »a high risk of violence, high risk of recidivism, high pretrial risk.«</p>



Beispiel COMPAS

COMPAS (Correctional Offender Management Profiling for Alternative Sanctions)

Loomis versuchte das Urteil anzufechten, dies wurde aber vom zuständigen Wisconsin Supreme Court abgelehnt:



¶6 The court of appeals certified the specific question of whether the use of a COMPAS risk assessment at sentencing "violates a defendant's right to due process, either because the proprietary nature of COMPAS prevents defendants from challenging the COMPAS assessment's scientific validity, or because COMPAS assessments take gender into account." 12

¶28 In denying the post-conviction motion, the circuit court explained that it used the COMPAS risk assessment to corroborate its findings and that it would have imposed the same sentence regardless of whether it considered the COMPAS risk scores. Loomis appealed and the court of appeals certified the appeal to this court.

¶56 Additionally, this is not a situation in which portions of a PSI are considered by the circuit court, but not released to the defendant. The circuit court and Loomis had access to the same copy of the risk assessment. Loomis had an opportunity to challenge his risk scores by arguing that other factors or information demonstrate their inaccuracy.

http://www.scotusblog.com/wp-content/uploads/2017/02/16-6387-op-bel-wis.pdf

COMPAS (Correctional Offender Management Profiling for Alternative Sanctions)

- 32. If you lived with both parents and they later separated, how old were you at the time?

 ✓ Less than 5 ☐ 5 to 10 ☐ 11 to 14 ☐ 15 or older ☐ Does Not Apply
- 33. Was your father (or father figure who principally raised you) ever arrested, that you know of?
 ☑ No ☐ Yes
- 34. Was your mother (or mother figure who principally raised you) ever arrested, that you know of?
 ☑ No ☐ Yes

Please think of your friends and the people you hung out with in the past few (3-6) months.

- 39. How many of your friends/acquaintances have ever been arrested?
 ☐ None ☐ Few ☑ Half ☐ Most
- 64. Do you have an alias (do you sometimes call yourself by another name)?
 ☑ No ☐ Yes

COMPAS (Correctional Offender Management Profiling for Alternative Sanctions)

| 78. How strongly do you agree or disagree with the following: I always behaved myself in so ☐ Strongly Disagree ☑ Disagree ☐ Not Sure ☐ Agree ☐ Strongly Agree | chool? |
|---|--------|
|---|--------|



95. How often did you feel bored?

☐ Never ☑ Several times/mo ☐ Several times/wk ☐ Daily

Criminal Personality

The next few statements are about what you are like as a person, what your thoughts are, and how other people see you. There are no 'right or wrong' answers. Just indicate how much you agree or disagree with each statement.

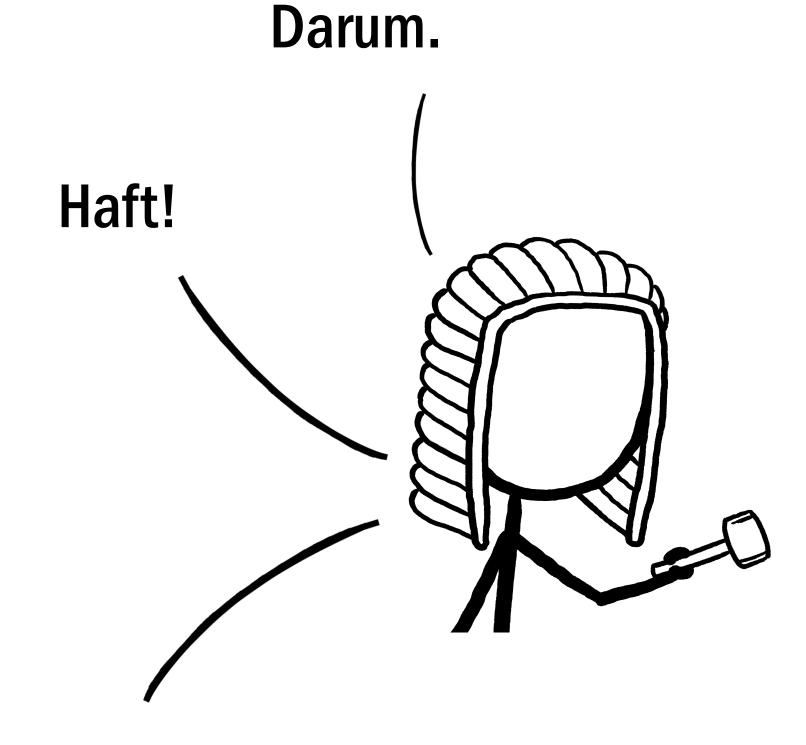
112. "I am seen by others as cold and unfeeling."

☑ Strongly Disagree □ Disagree □ Not Sure □ Agree □ Strongly Agree

137. "Some people just don't deserve any respect and should be treated like animals."

☑ Strongly Disagree ☐ Disagree ☐ Not Sure ☐ Agree ☐ Strongly Agree

Eine verantwortungsvolle Entscheidung?



Da wüsste ich aber gerne, wieso!

Inakzeptabel!

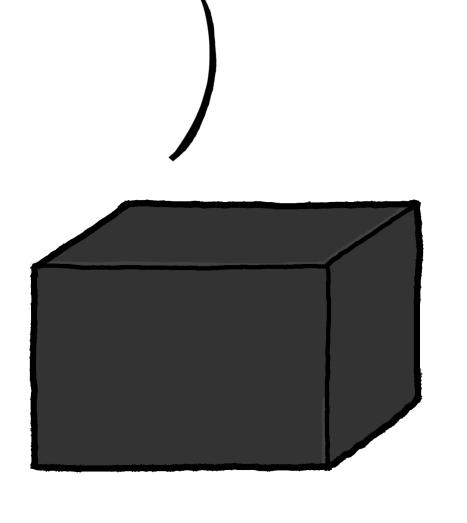
Darum.



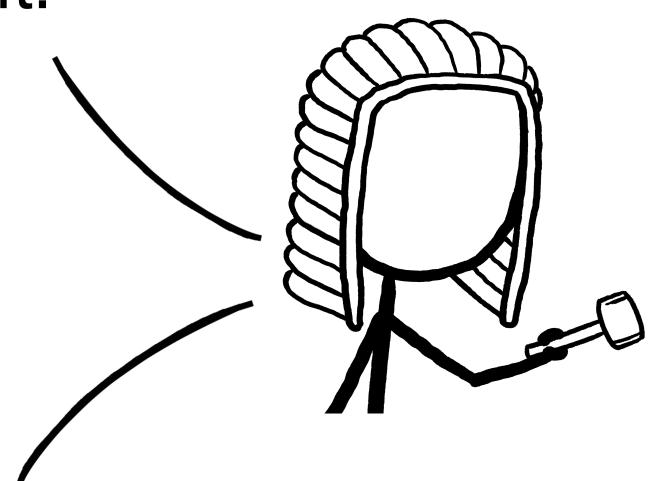
Okay, klar, Ihr gutes Recht. Justitia, warum noch gleich?

Eine verantwortungsvolle Entscheidung?



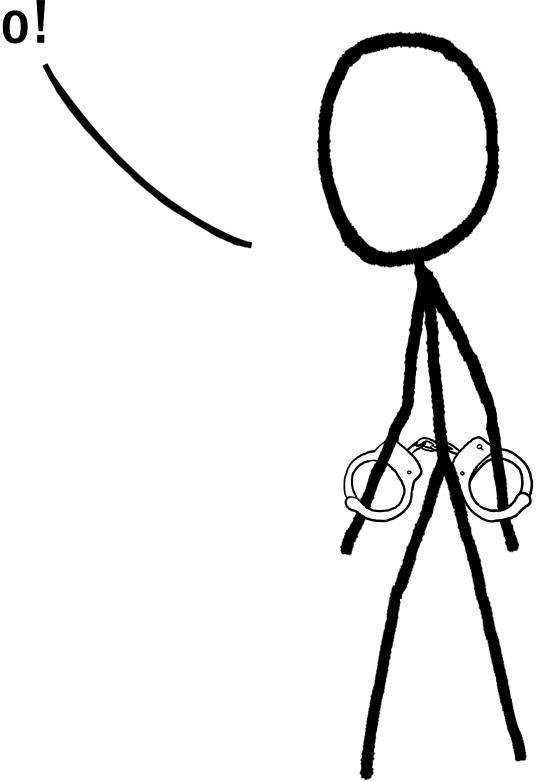






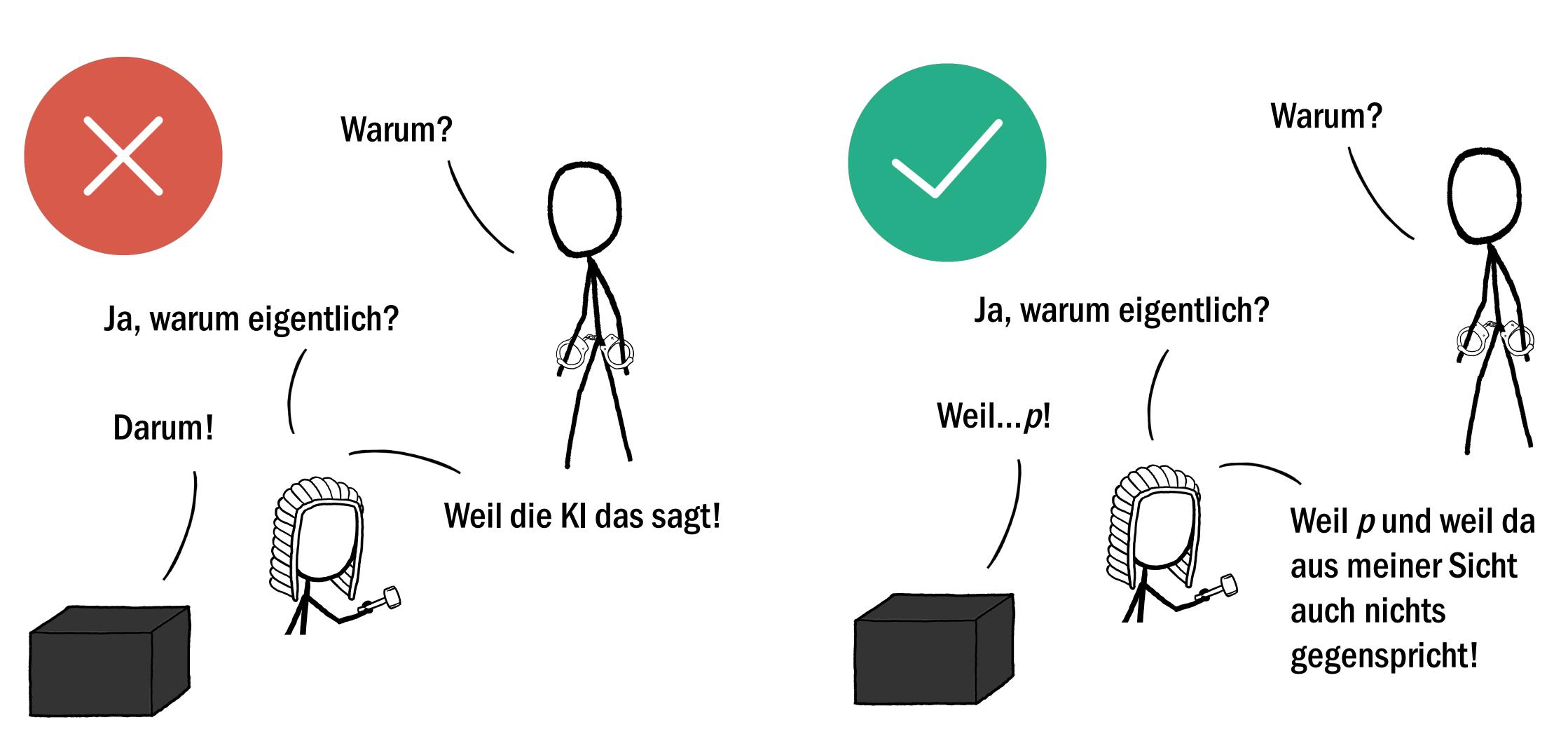
Okay, klar, Ihr gutes Recht. Justitia, schicke bitte meine Urteilsbegründung inklusive Deiner begründeten Einschätzung an seinen Anwalt.

Da wüsste ich aber gerne, wieso!



In einer liberalen Demokratie, in einem Rechtsstaat, geht es nicht ohne Erklärbarkeit.

Fazit 3



24

Und was hat das jetzt mit Verantwortung zu tun?

Wir haben doch beides: Verantwortung und KI-Power!

Eine verantwortungsvolle Entscheidung?

Ich hab einfach nur mit dem gearbeitet, was ich habe!

Wo liegt die Verantwortung für dieses Urteil?

Haft. Es besteht eine 85%ige Chance, dass er während er Bewährung wieder straffällig wird.

(Vielleicht mitverantwortlich fürs Zurverfügungstehen der KI in diesem Kontext, aber nicht für diese konkrete Entscheidung)

Nelal

Okay, also Haft.

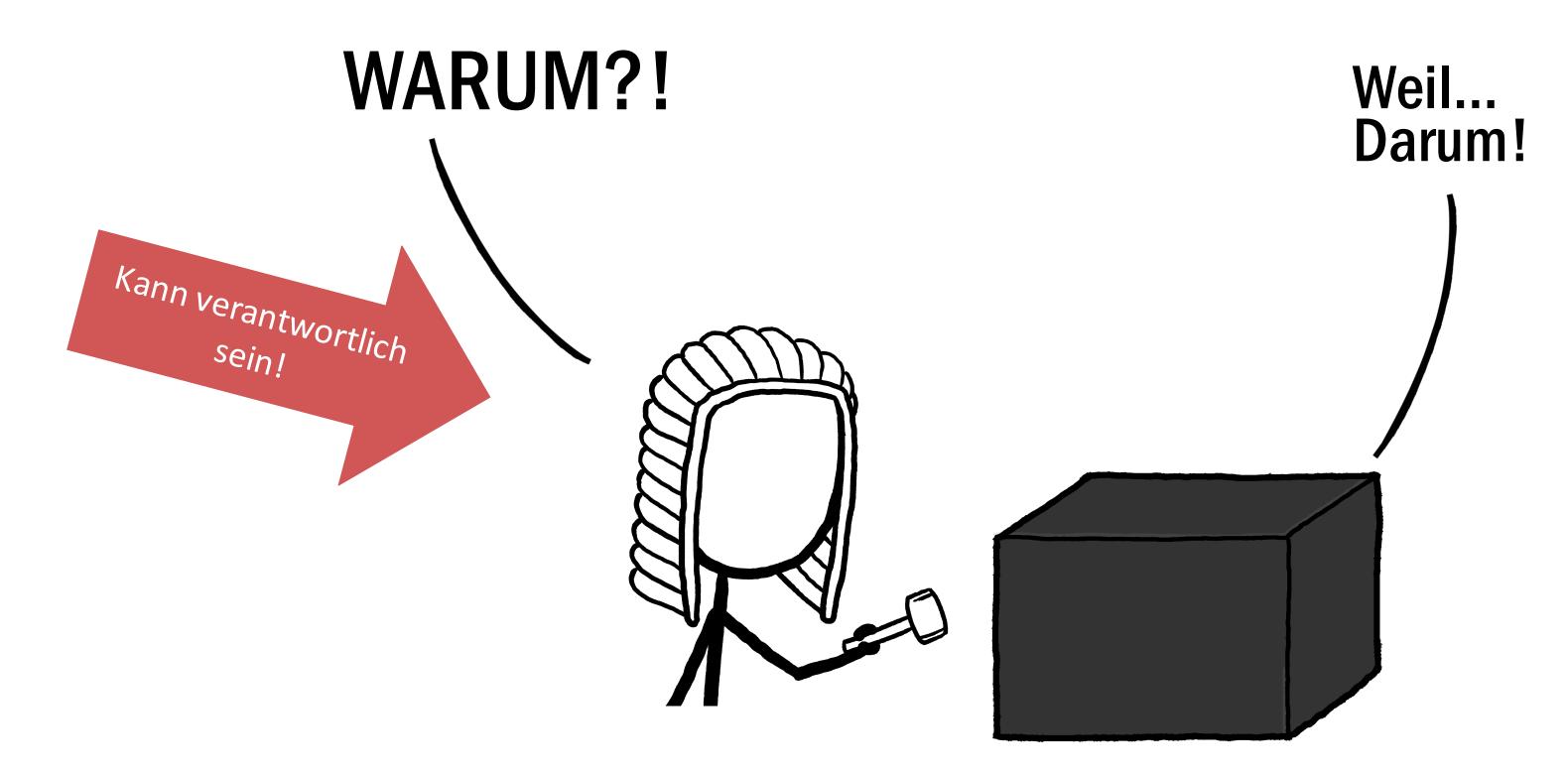
Was soll ich urteilen:

Bewährung oder

Haft?

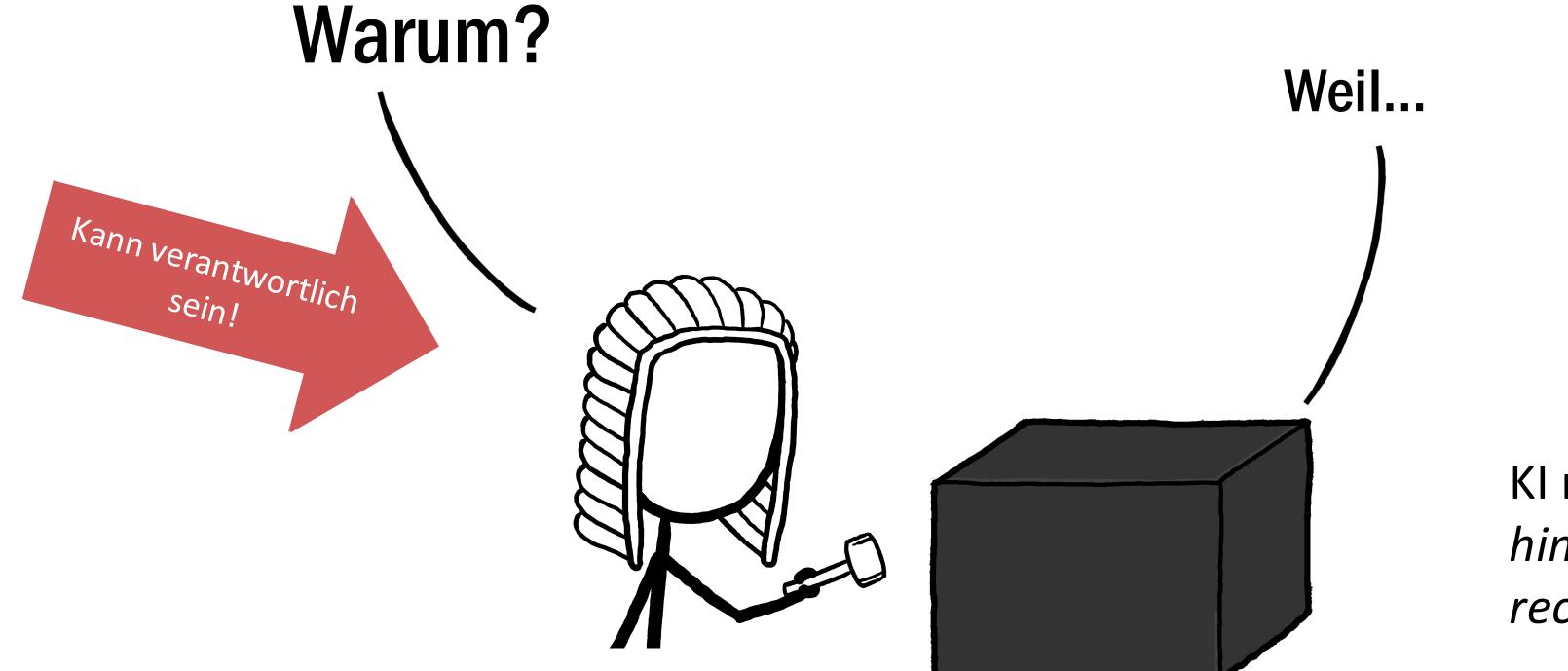
Ein Lösungsansatz

- Ein Mensch kann nur dann im Lichte einer Empfehlung verantwortungsvoll entscheiden, wenn er die Grundlage jener Empfehlung gegen seine eigene Grundlage abwägen kann.
- Nur dann kann er sich bewusst nach Abwägung von Für und Wider für oder gegen das Befolgen der Empfehlung entscheiden.



Ein Lösungsansatz: Erklärbarkeit

Es muss mindestens auf Anfrage *Erklärungen* zu Empfehlungen geben, d.h. dargelegte Gründe.

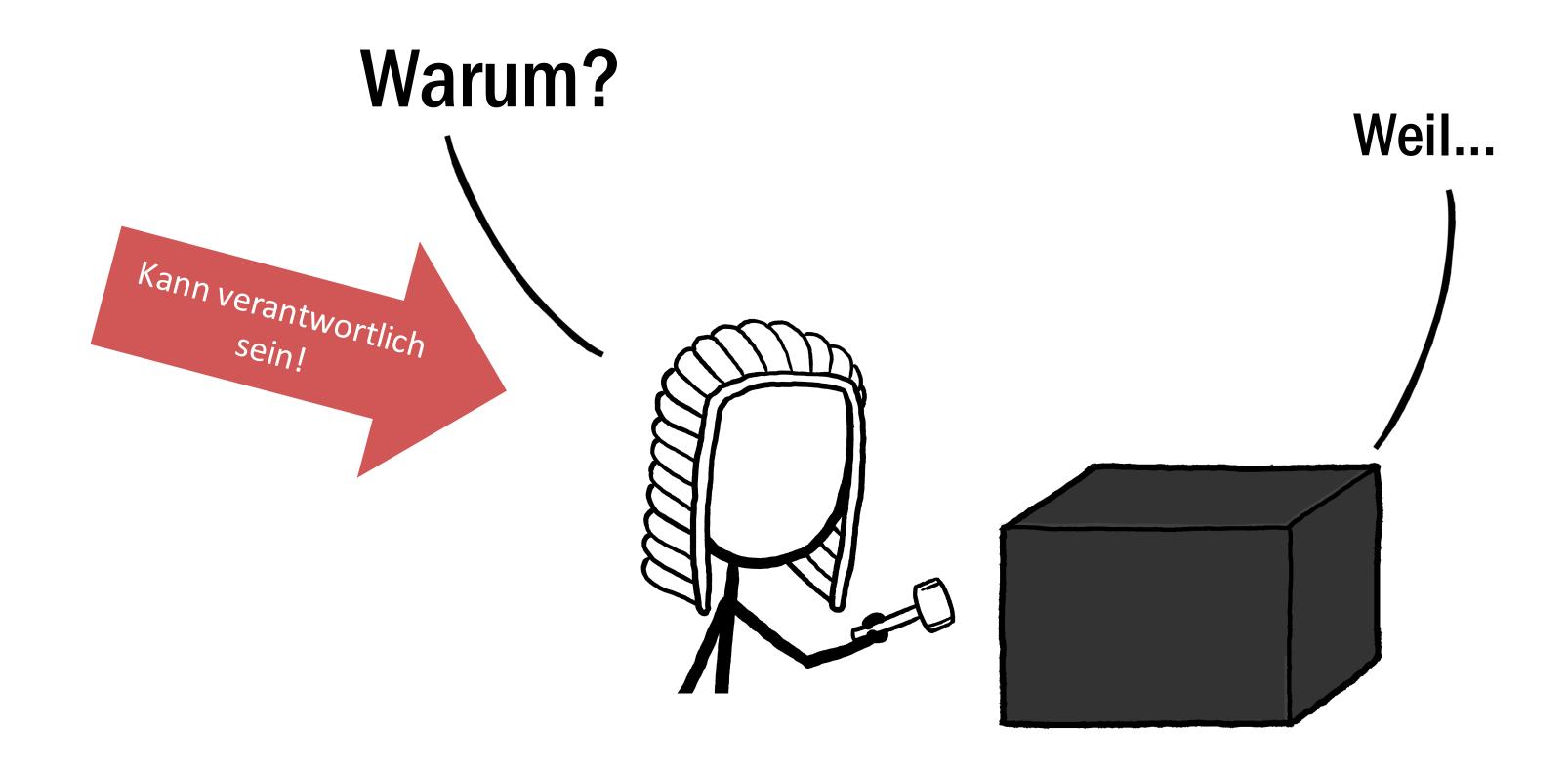


KI muss verständlich, hinterfragbar und rechtfertigbar sein.

Fazit 4

Erklärbarkeit ist *notwendige Voraussetzung* für verantwortungsvolle,

KI-gestützte Entscheidungen.



... aber wir arbeiten dran!





Aktuelles Uni-Porträt Fakultäten|Einrichtungen Wirtschaftsportal

Montag, 26. November 2018

FOUNDATIONS OF PERSPICUOUS SOFTWARE SYSTEMS

— Enabling Comprehension in a Cyber-Physical World —



Neuer Sonderforschungsbereich: Softwaresysteme sollen ihr Verhalten selbst erklären

Selbst Experten verstehen das Verhalten komplexer Softwaresysteme immer weniger. Dabei regeln diese inzwischen immer stärker unseren Alltag, sei es als intelligente Haussteuerung, im autonomen Fahrzeug oder in der industriellen Produktion. Wissenschaftler der Universität des Saarlandes, zweier Max-Planck-Institute und der Technischen Universität Dresden wollen jetzt in einem neuen Sonderforschungsbereich Mechanismen entwickeln, die nicht nur Experten, sondern auch Laien das Verhalten komplexer Softwaresysteme besser vermittelt. Die Deutsche Forschungsgemeinschaft fördert dieses Großprojekt mit elf Millionen Euro über vier Jahre hinweg.





Informatik an der Saar-Uni entwickelt die Software der Zukunft. Von Peter Bylda >

Zwei Abschließende Botschaften

Erreichen wir Erklärbarkeit nicht droht ein Novum:

Zugewinn an technischem Vermögen

> prinzipieller Verlust an Erkenntnis und Kontrolle

 Einsicht in die Grundlage einer Entscheidung ermöglicht erst Vertrauenswürdigkeit und gerechtfertigte Akzeptanz.

31

Danke für Ihre Aufmerksamkeit!

32