

Julien Cornebise, Ph.D.

2020-12-18 ML Collective / DL Classics and Trends

AI FOR GOOD AND ETHICS—WASHING A SELF—DEFENSE PRIMER

```

EAX=00000000 EBP=7E43C689 ESP=0012FE5C EIP=7E42772B o d i s z a p c
CS=001B DS=0023 SS=0023 ES=0023 FS=003B GS=0000

*ESP <ulong> = 0x401DFF, 4201983
*EBP <ulong> = 0x8B55FF8B, -1957298293
-----StartcIn!.data+0030-----byte-----PROT-----(0)-----
0010:00406030 31 33 33 34 2D 31 32 33-38 36 2D 31 38 30 35 2D 1334-12386-1805-
0010:00406040 33 35 37 00 00 00 00 00-00 00 00 00 00 00 00 00 357.....
0010:00406050 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:00406060 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:00406070 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:00406080 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:00406090 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:004060A0 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:004060B0 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:004060C0 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:004060D0 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:004060E0 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:004060F0 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:00406100 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
0010:00406110 00 00 00 00 00 00 00 00-00 00 00 00 00 00 00 00 .....
-----USER32!CallMsgFilterW+021C-----PROT32-----
001B:7E42772A 90 NOP
USER32!GetMessageA
001B:7E42772B 8BFF MOV EDI,EDI
001B:7E42772D 55 PUSH EBP
001B:7E42772E 8BEC MOV EBP,ESP
001B:7E427730 8B5510 MOV EDX,[EBP+10]
001B:7E427733 8B4D14 MOV ECX,[EBP+14]
001B:7E427736 56 PUSH ESI
001B:7E427737 8BF2 MOV ESI,EDX
001B:7E427739 0BF1 OR ESI,ECX
001B:7E42773B B80000FEFF MOV EAX,FFFE0000
001B:7E427740 85F0 TEST EAX,ESI
001B:7E427742 0F855C9A0100 JNZ ↓7E4411A4
001B:7E427748 57 PUSH EDI
001B:7E427749 64A118000000 MOV EAX,FS:[00000018]
001B:7E42774F 83B84007000000 CMP DWORD PTR [EAX+00000740],00
001B:7E427756 0F85589A0100 JNZ ↓7E4411B4
001B:7E42775C 8B7508 MOV ESI,[EBP+08]
001B:7E42775F 53 PUSH EBX
001B:7E427760 51 PUSH ECX
(PASSIVE)-KTEB(8953E7A0)-TID(041C)-user32!.text+0001672A
milliseconds)
sli.do - #self-defense-primer
Enter a command (h for help) StartcIn
```

We have the following limits, as $N \rightarrow \infty$.

Theorem

① *Asymptotic behavior of KLD-based criterion:*

$$\left| d_{\text{KL}}(\mu_N^{\text{aux}} || \pi_N^{\text{aux}}) - \mathcal{E}(\{\tilde{\omega}_{N,i}\}_{i=1}^{\tilde{M}_N}) \right| \xrightarrow{\mathbb{P}} 0 .$$

In addition,

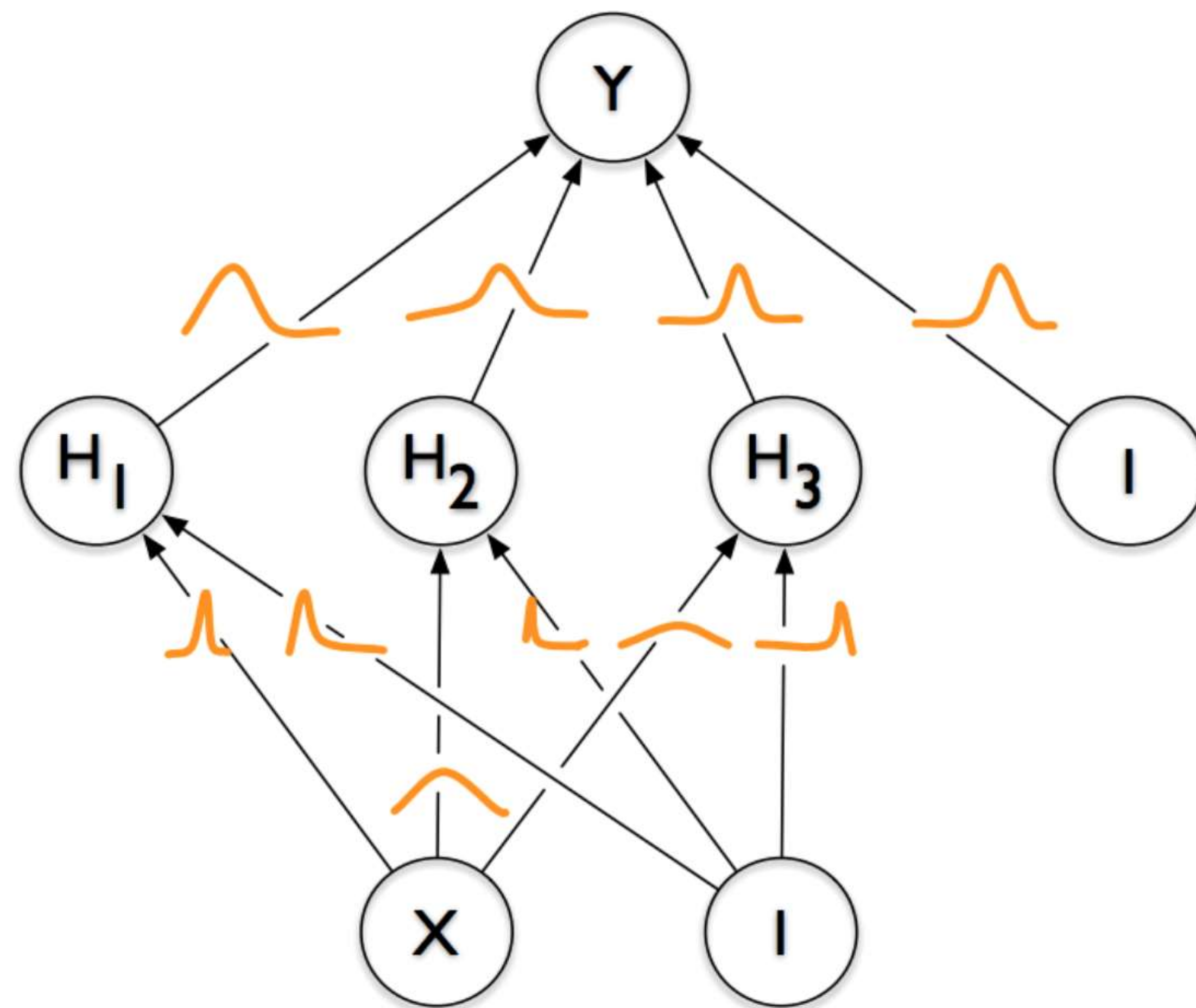
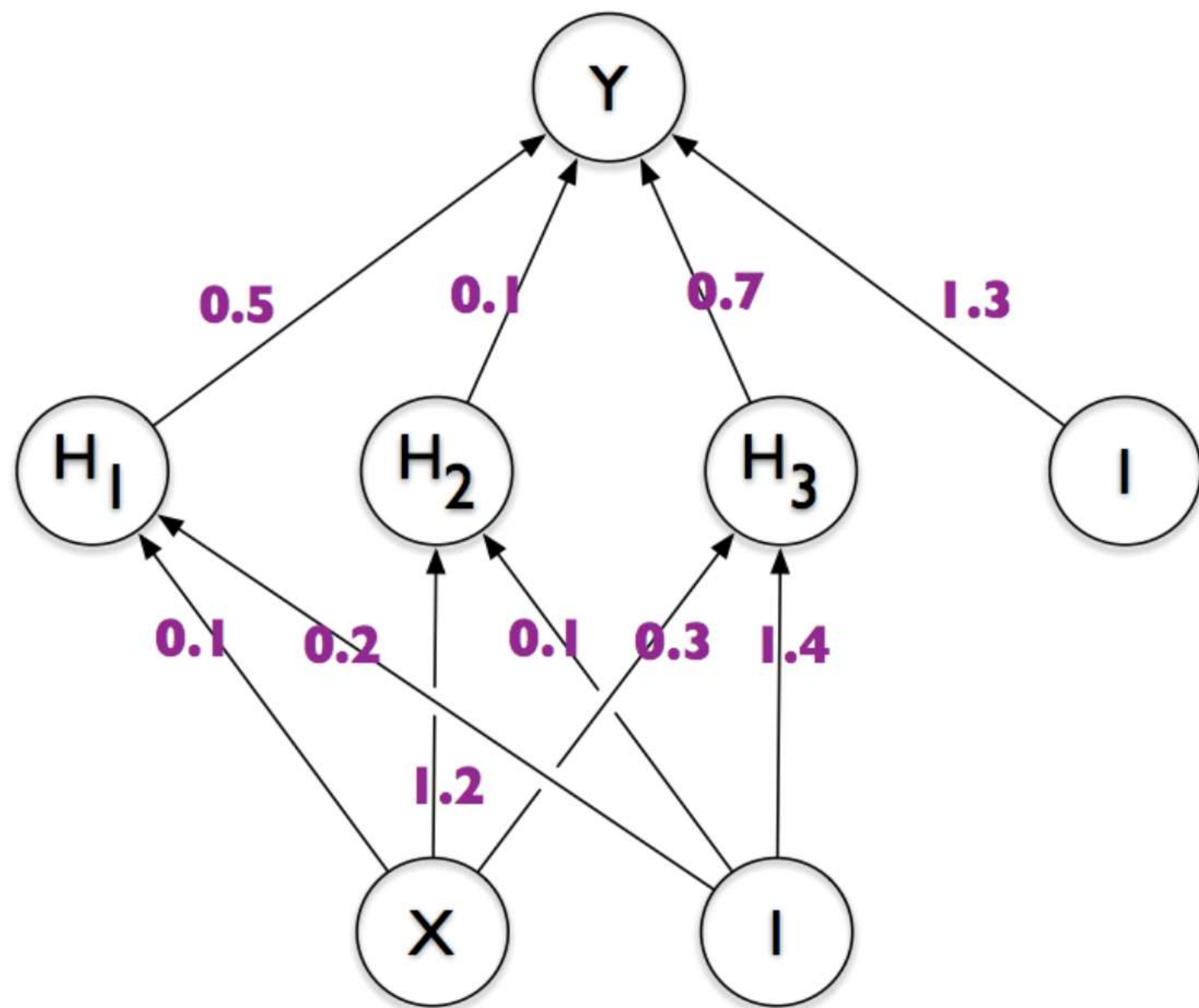
$$d_{\text{KL}}(\mu_N^{\text{aux}} || \pi_N^{\text{aux}}) \xrightarrow{\mathbb{P}} \eta_{\text{KL}}(\Psi) \triangleq \frac{\nu \otimes L}{\nu L(\tilde{\Xi})} \left\{ \log \left[\Phi \frac{\nu(\Psi)}{\nu L(\tilde{\Xi})} \right] \right\} ,$$

② *Asymptotic behavior of CSD-based criterion:*

$$\left| d_{\chi^2}(\mu_N^{\text{aux}} || \pi_N^{\text{aux}}) - \text{CV}^2(\{\tilde{\omega}_{N,i}\}_{i=1}^{\tilde{M}_N}) \right| \xrightarrow{\mathbb{P}} 0 .$$

In addition,

$$d_{\chi^2}(\mu_N^{\text{aux}} || \pi_N^{\text{aux}}) \xrightarrow{\mathbb{P}} \eta_{\chi^2}(\Psi) \triangleq \frac{\nu(\Psi)}{\nu L(\tilde{\Xi})} \times \frac{\nu \otimes L}{\nu L(\tilde{\Xi})}(\Phi) - 1 .$$



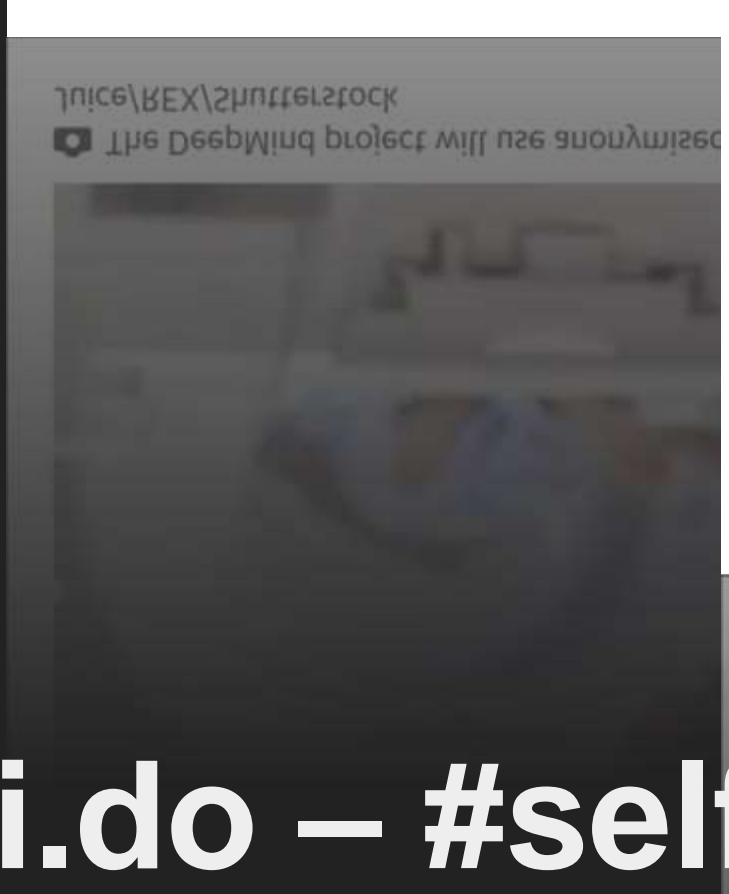
Google

Google DeepMind on AI-based radiology

Google's AI research arm is partnering with radiotherapists by using machine learning



The DeepMind project will use anonymised patient data. Juice/REX/Shutterstock

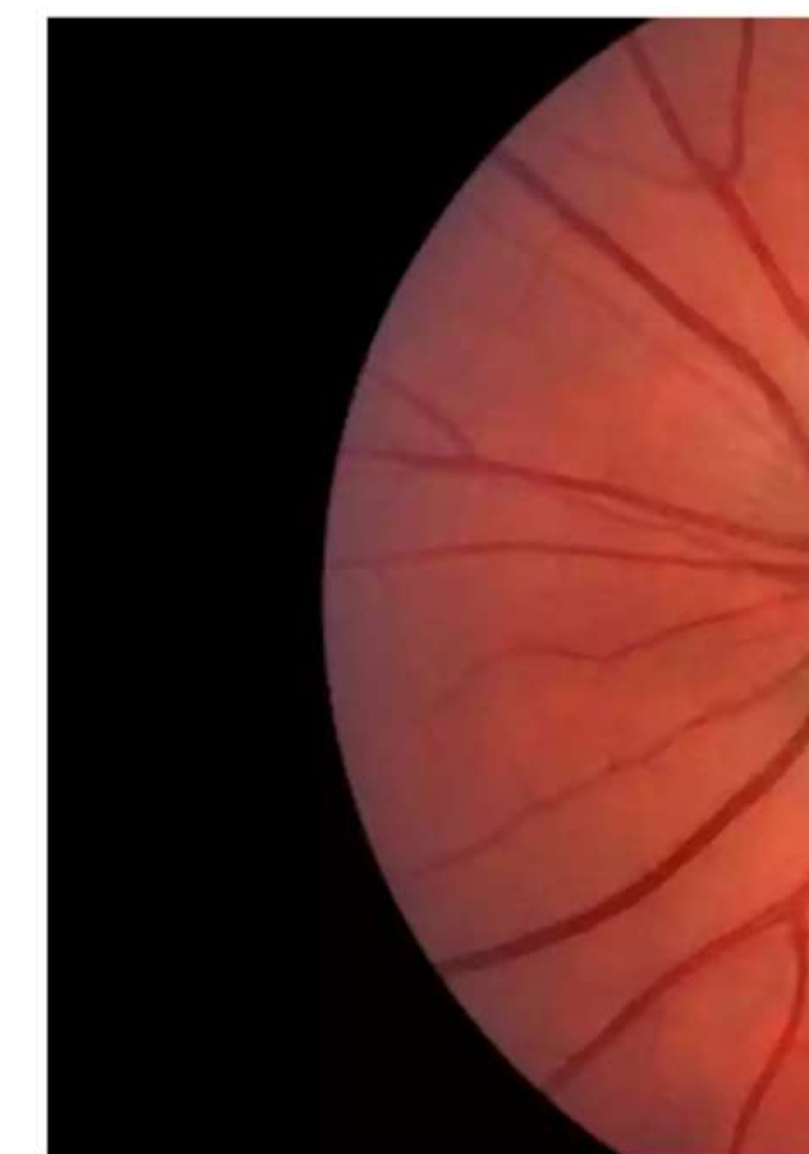


Retinal scans are produced rapidly, but require a specialist to interpret. Alamy Stock Photo/Alamy Stock Photo

Google

Google DeepMind machine learning

'Deep learning' research company will use AI network to identify early signs of degenerative diseases



Retinal scans are produced rapidly, but require a specialist to interpret. Alamy Stock Photo/Alamy Stock Photo

Retinal scans are produced rapidly, but require a specialist to interpret. Alamy Stock Photo/Alamy Stock Photo

FINANCIAL TIMES myFT

Get a fresh start. Choose your FT trial

Latest on Artificial intelligence

Marcus du Sautoy on creative AI Daimler 7 US self-d

Artificial intelligence grows into commercial product

Device can diagnose a range of eye diseases, say specialists

Twitter Facebook LinkedIn Save

Madhumita Murgia in London MARCH 31, 2019

DeepMind, the British artificial intelligence company, has unveiled a device that can diagnose complex eye diseases. It is Alphabet-owned company's first medical product.

In a live demonstration this month of the device, which was examined publicly, DeepMind performed a scan of a patient's eye. The scan was analysed by a set of algorithms that generated an urgency score and a detailed diagnosis.

WIRED

Health

DeepMind's new AI predicts kidney injury two days before it happens

New research from the Google-owned firm hints that AI may be a better way of assessing if someone is at risk of acute kidney injury. But there are still questions about how it handles patient data

By MATT REYNOLDS 31 Jul 2019

Facebook Twitter WhatsApp Email





AMNESTY
INTERNATIONAL



@deworrall92

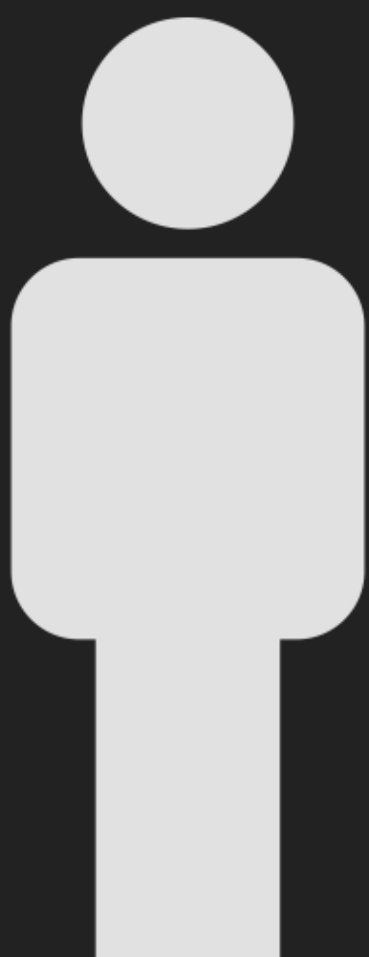


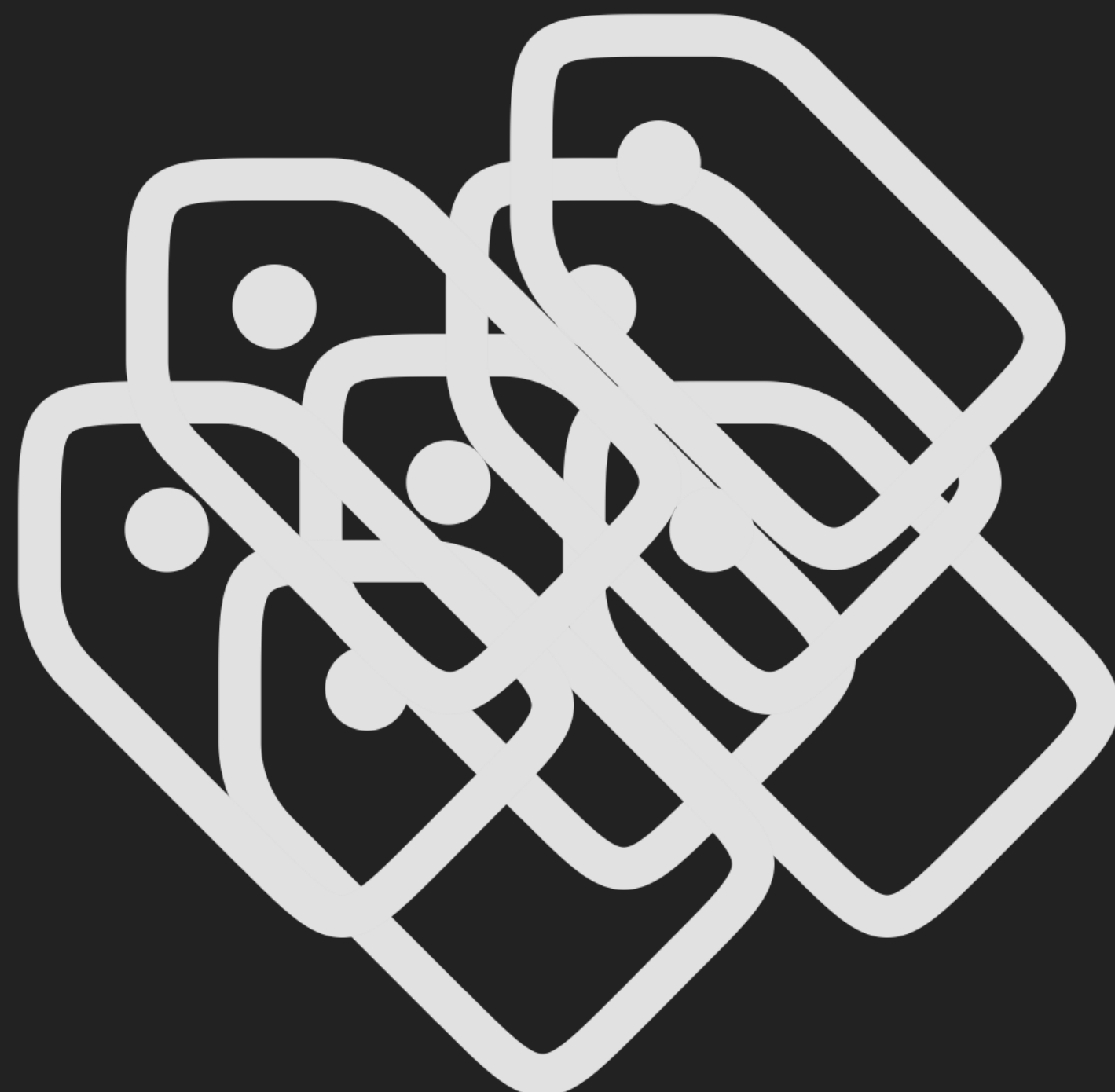
@milena_iul

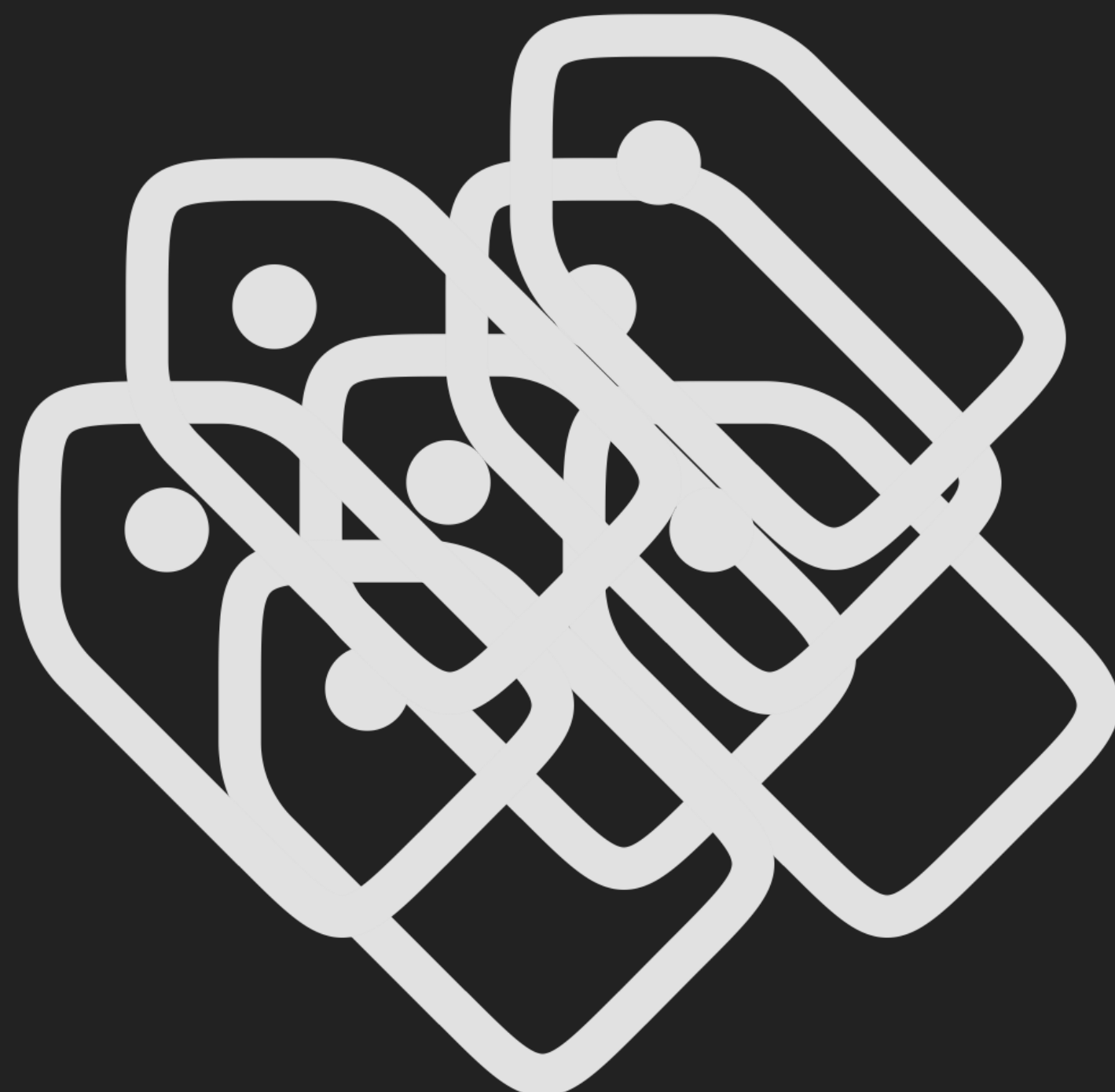
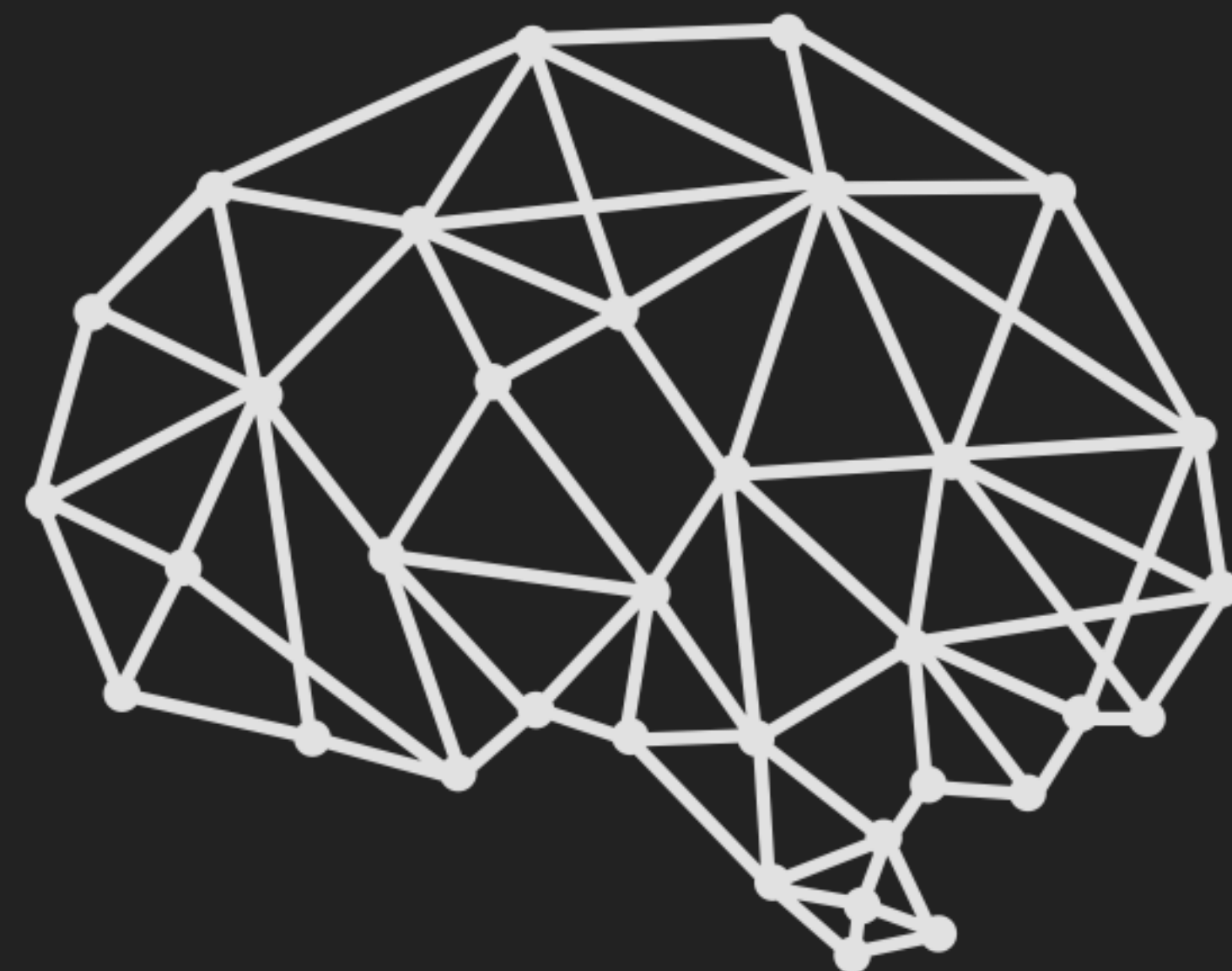
2017

DECODE DARFUR











Well done! You
identified the human
presence in this tile

NEXT

SELECT TILES CONTAINING HUMAN PRESENCE

FLAG

HELP



NEW ACHIEVEMENT

YOU DECODED
50 SQUARE KM

CONTINUE →

25

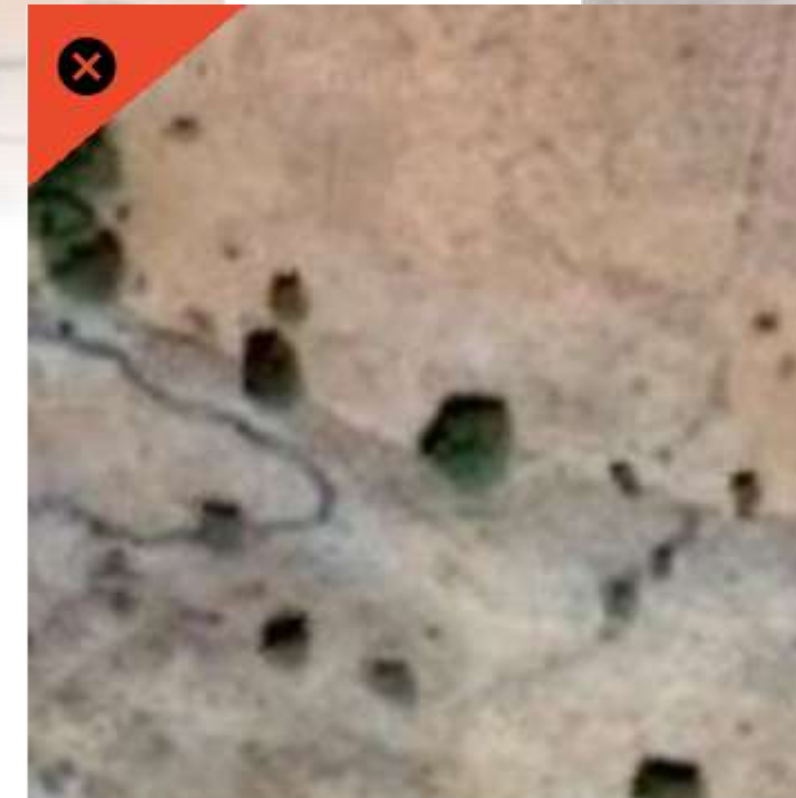
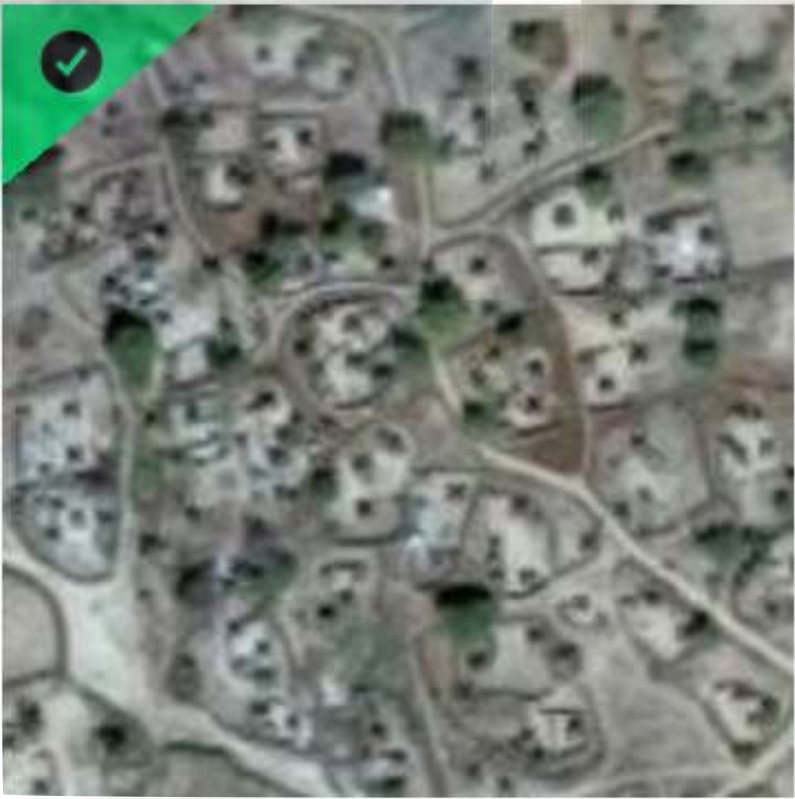
SQUARE KM
THIS SESSION

149

SQUARE KM
ALL SESSIONS

1659

SQUARE KM
ALL VOLUNTEERS



<Video demo removed due to potentially sensitive content>

AMNESTY
INTERNATIONAL



@alkalait



@laure_delisle



@milena_iul



@SnazzyAzzy



@sherifea



@ArchydeB

2018

TROLL PATROL: STUDYING

ABUSE AGAINST WOMEN ON TWITTER

TRIGGER

WARNING

 **unidentified, female (@unidentified, 2016)** · Jul 13
Replying to [@HackneyAbbott](#)
this **fat** retarded black **bitch** thinks you should be forced to feed and house a bunch of violence foreign invaders. i strongly disagree.

 **unidentified (@unidentified)** · Jul 13
[@HackneyAbbott](#) Piss off you disgusting useless **fat bitch!** You're a parasite alien looking to silence native people for your power.

 **unidentified (@unidentified)** · 10/11/2016, 23:33
Tasmina needs to be grabbed by the pussy. 🤔🤔🤔 #bbcqt

 **unidentified (@unidentified)** · 10/11/2016, 23:36
Tasmina needs to be poked with a massive dildo #bbcqt

 **unidentified (@unidentified)** · Jul 14
Replying to [@HackneyAbbott](#)
You forgot "fat disgusting obese chicken-loving **nigger**"

 **unidentified (@unidentified)** · Jul 14
Replying to [@HackneyAbbott](#)
An acid attack would probably make your face look better you fat **nigger**

 **unidentified (@unidentified)**
Someone jo cox Anna sourby please
02/12/2016, 07:41 from Poplar, London

STUDY: ML TO HUMANS TO ML TO HUMANS

- ▶ Build a dataset with **ML prefilter**
- ▶ **Importance sampling** to
 - ▶ reduce variance
 - ▶ drive engagement
- ▶ Human **Crowdsourcing**
- ▶ Data **Analysis with Bootstrap**
- ▶ Train **Deep model and benchmark**
- ▶ Carefully and expertly **report** and go public!
 - ▶ **Very sensitive topic**, human rights expertise!

AMNESTY
INTERNATIONAL



TROLL PATROL FINDINGS

Using Crowdsourcing, Data Science & Machine Learning to
Measure Violence and Abuse against Women on Twitter

Welcome to the findings of our Troll Patrol project: a joint effort by human rights researchers, technical experts and thousands of online volunteers to build the world's largest crowd-sourced dataset of online abuse against women.

Our findings reveal the sheer scale and nature of online abuse faced by women and provides a resource to researchers and engineers interested in exploring the potential of machine learning in content moderation.

These findings are the result of a collaboration between Amnesty International and **Element AI**, a global artificial intelligence software product company. Together, we surveyed millions of tweets received by

Technology

Together with Element AI, the human rights watchdog quantified what women on Twitter have known for a long time

U.S. POLITICS WORLD TECH ENTERTAINM

'A Toxic Place for Women.' A New Study Reveals the Scale of Abuse on Twitter

Topics•

Female black journalists and politicians get sent an abusive tweet every 30 seconds

EMILY DREYFUS SECURITY 12.18.18 12:00 AM

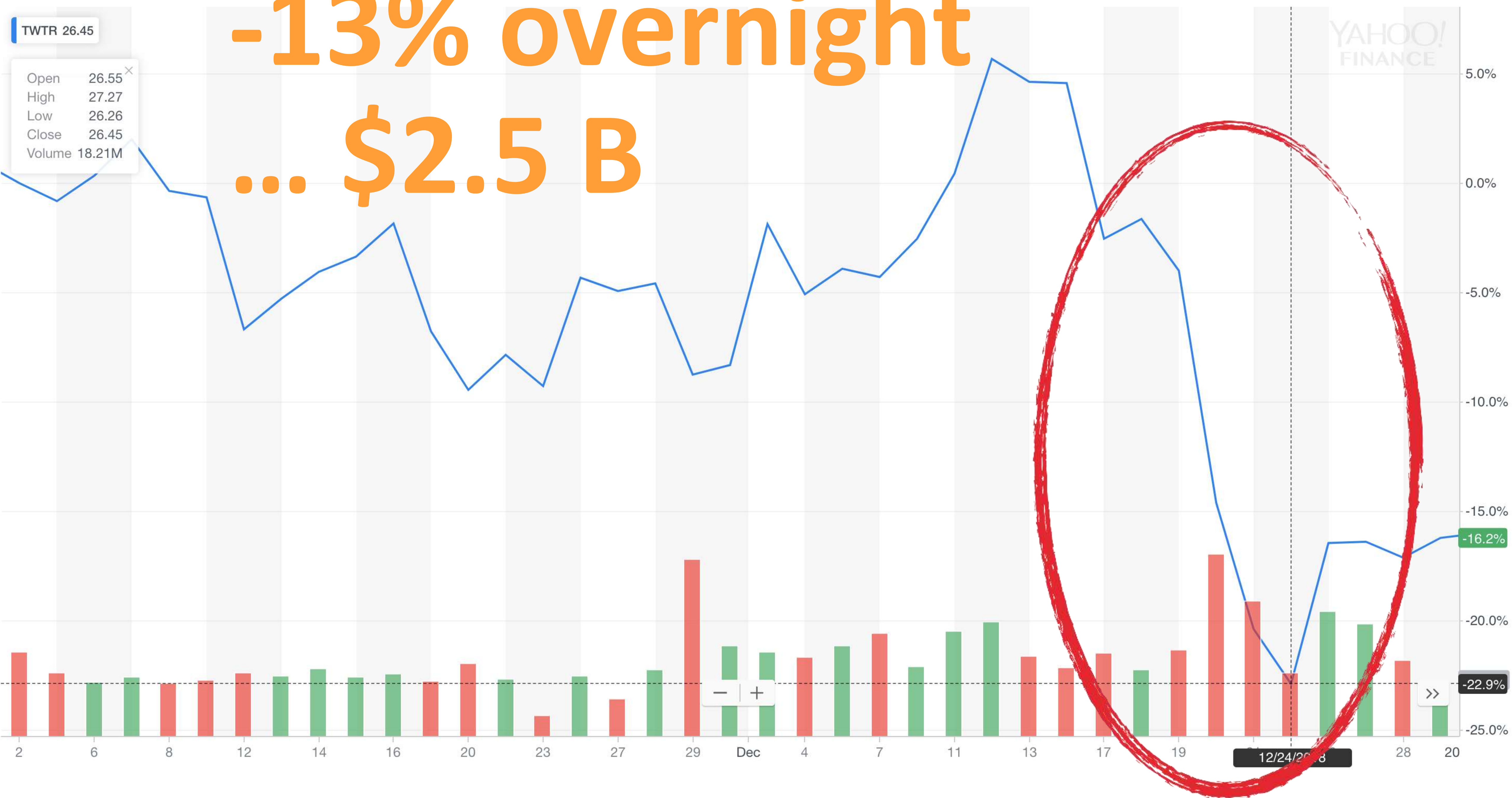
TWITTER IS INDEED TOXIC FOR WOMEN, AMNESTY REPORT SAYS



| | A | B |
|----|----------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 1 | Publication | Link |
| 2 | Wired | https://www.wired.com/story/annex-report-twitter-abuse-women/ |
| 3 | FT | https://www.ft.com/content/1891e116-5324-11e6-9a61-001448f00000 |
| 4 | Gizmodo | https://gizmodo.com/le-finally-have-some-hard-data-on-how-much-better-1511136563 |
| 5 | Shrapline Star | https://www.shrapline.com/news/2015/12/15/diane-abbott-condemns-twitter-over-thousands-of-racist-and-misogynistic-hate/ |
| 6 | Moroccan Gazette | https://moroccan-gazette.com/business/local-business/elements-5-partners-with-annex-international-to-study-online-abuse |
| 7 | Guardian | https://www.theguardian.com/politics/2015/dec/15/diane-abbott-calls-for-better-to-clang-down-on-hate-speech |
| 8 | Bloomberg | https://www.bloomberg.com/news/articles/2015-12-15/better-is-bad-place-for-women-says-annex-international |
| 9 | Le Devoir | https://www.ledevoir.com/societe/543625/quand-les-femmes-t-en-prennent-assez-femmes-sur-twitter |
| 10 | Engadget | https://www.engadget.com/2015/12/15/women-are-abused-every-30-seconds-on-tweet/ |
| 11 | Time | http://time.com/3412363/better-online-abuse-women-annex-international-study/ |
| 12 | Pactl | https://pactl.co.uk/pactl-co-uk/pactl-report-annex-international-and-element-8-use-machine-learning-to-understand-online-abuse-against-women/ |
| 13 | Jakarta Post | https://www.thejakartapost.com/news/2015/12/15/better-is-bad-place-for-women-says-annex-international.html |
| 14 | NewsTalk | https://www.newstalk.com/Study-reveals-scale-of-online-abuse-against-women |
| 15 | Hindustan Times | https://in.hindustantimes.com/world-news/better-is-bad-place-for-women-says-annex-international-story/1511136563.html |
| 16 | La Presse | https://www.lapresse.ca/actualites/monde-societe/2015/12/15/101-3206408-les-femmes-sont-premieres-cibles-des-attaques-sur-tweet.php |
| 17 | Business Standard | https://www.business-standard.com/article/international/annex-international-explains-how-tweet-is-a-bad-place-for-women-115121500254_1.html |
| 18 | The Independent | http://www.independent.co.uk/life-style/women/better-women-said-annex-international-diane-abbott-report-a655435.html |
| 19 | Business Today Kenya | https://businessday.co.ke/women-attacked-better-every-30-seconds-new-study-reveals |
| 20 | Press Gazette | https://www.pressgazette.co.uk/in-every-14-seconds-directed-at-women-journals-is-abusive-or-problematic-new-online-abuse-study-shows |
| 21 | Mashable | https://mashable.com/article/annex-study-better-abuse-women/Tenure+burden+As+DO+of |
| 22 | Fast Company | https://www.fastcompany.com/30250043/women-on-better-is-abused-every-30-seconds |
| 23 | MIT Tech Review | https://www.technologyreview.com/2015/12/15/female-journalists-and-politicians-get-attacked-every-30-seconds/ |
| 24 | Evening Standard | https://www.standard.co.uk/tech/online-abuse-better-annex-international-element-8-a6420176.html |
| 25 | The Voice | http://www.voice-online.co.uk/article/diane-abbott-urges-better-address-abusive-hate/ |
| 26 | Mail and Guardian (South Africa) | https://mg.co.za/article/2015-12-15-women-across-political-spectrum-expose-better-abuse-study |
| 27 | Daily Out | https://www.dailyout.com/tech/g-black-women-abuse-while-women-better-report/ |
| 28 | Nu (Netherlands) | https://www.nu.nl/nieuws/3037206/annex-better-besluit-mensen.html |
| 29 | NRC (Netherlands) | https://www.nrc.nl/nieuws/2015/12/15/voorvrouw-journalisten-en-politici-op-better-eens-per-30-seconden-beledigd-a3001044 |
| 30 | Talk Radio | http://talkradio.co.uk/news/diane-abbott-urges-better-crackdown-abusive-hate-after-study-abuse-women-social-media |
| 31 | Al Africa | https://alafrika.com/stories/2015/12/15/013.htm |
| 32 | Westfälische Rundschau | https://www.westfaelische-rundschau.de/annex-studie-gegen-haun-voor-better-und-journalist-in-hass-1215039315.html |
| 33 | Forbes | https://www.forbes.com/sites/forbesmag/2015/12/15/le-petit-petit-petit-prove-women-harassed-on-tweet/ |
| 34 | The Journal (Ireland) | http://www.thejournal.ie/better-abuse-women-of-color-4403065-Dec15/15.html |
| 35 | Fortune | http://fortune.com/2015/12/15/better-women-annex-international/ |
| 36 | Le Monde | https://www.lemonde.fr/societe/article/2015/12/15/le-violence-des-femmes-sur-tweet-annex-international-erobne-le-cou_5196413_4439996.html |
| 37 | sky news | https://www.sky.com/story/diane-abbott-better-must-act-over-abuse-targeting-black-women-11504785 |
| 38 | The Telegraph | https://www.telegraph.co.uk/politics/2015/12/15/diane-abbott-urges-better-crack-down-misogynistic-abuse.html |
| 39 | Evening Standard (again) | https://www.standard.co.uk/tech/online-abuse-better-condemns-twitter-over-thousands-of-racist-and-misogynistic-hate-appearing-on-platform-a6418996.html |
| 40 | CNBC | https://www.cnbc.com/am/2015/12/15/women-are-abused-every-30-seconds-on-better-new-research-finds.html |
| 41 | Washington Post | |

-13% overnight

... \$2.5 B





**“TECHNOLOGY – NO MATTER HOW WELL
DESIGNED – IS **ONLY A MAGNIFIER** OF HUMAN
INTENT AND CAPACITY.”**

Kentaro Toyama (2015)



GEEK HERESY

RESCUING SOCIAL CHANGE
FROM THE CULT OF TECHNOLOGY

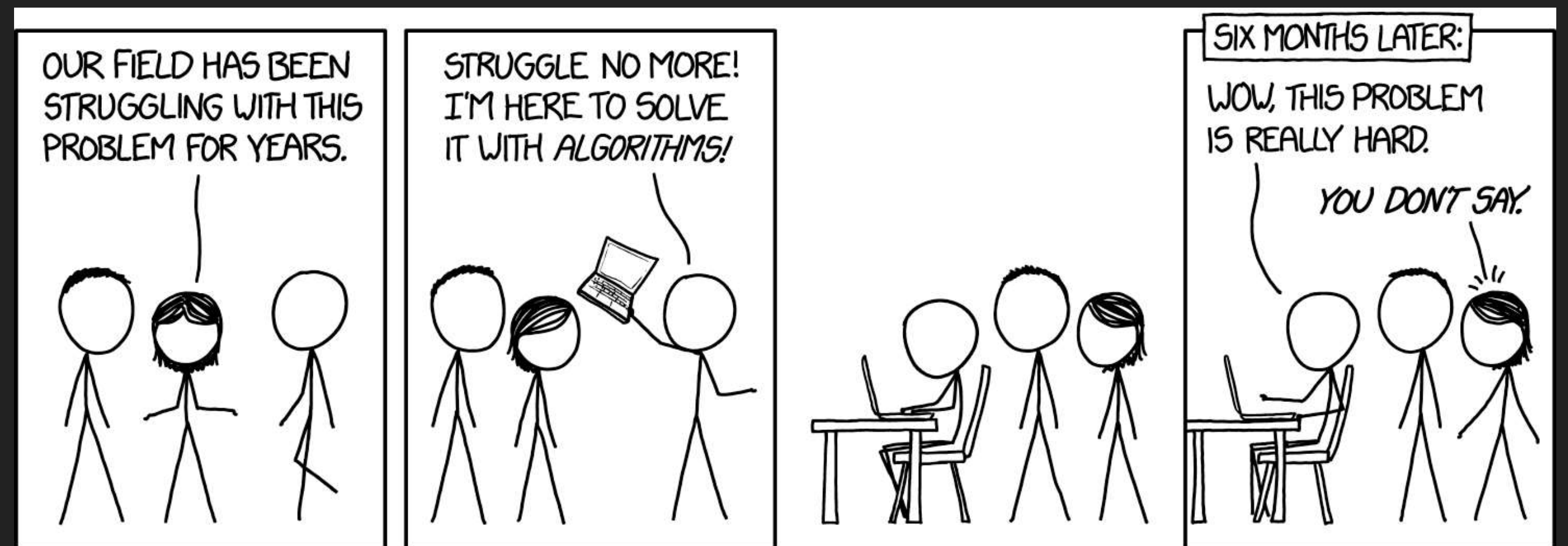
KENTARO TOYAMA

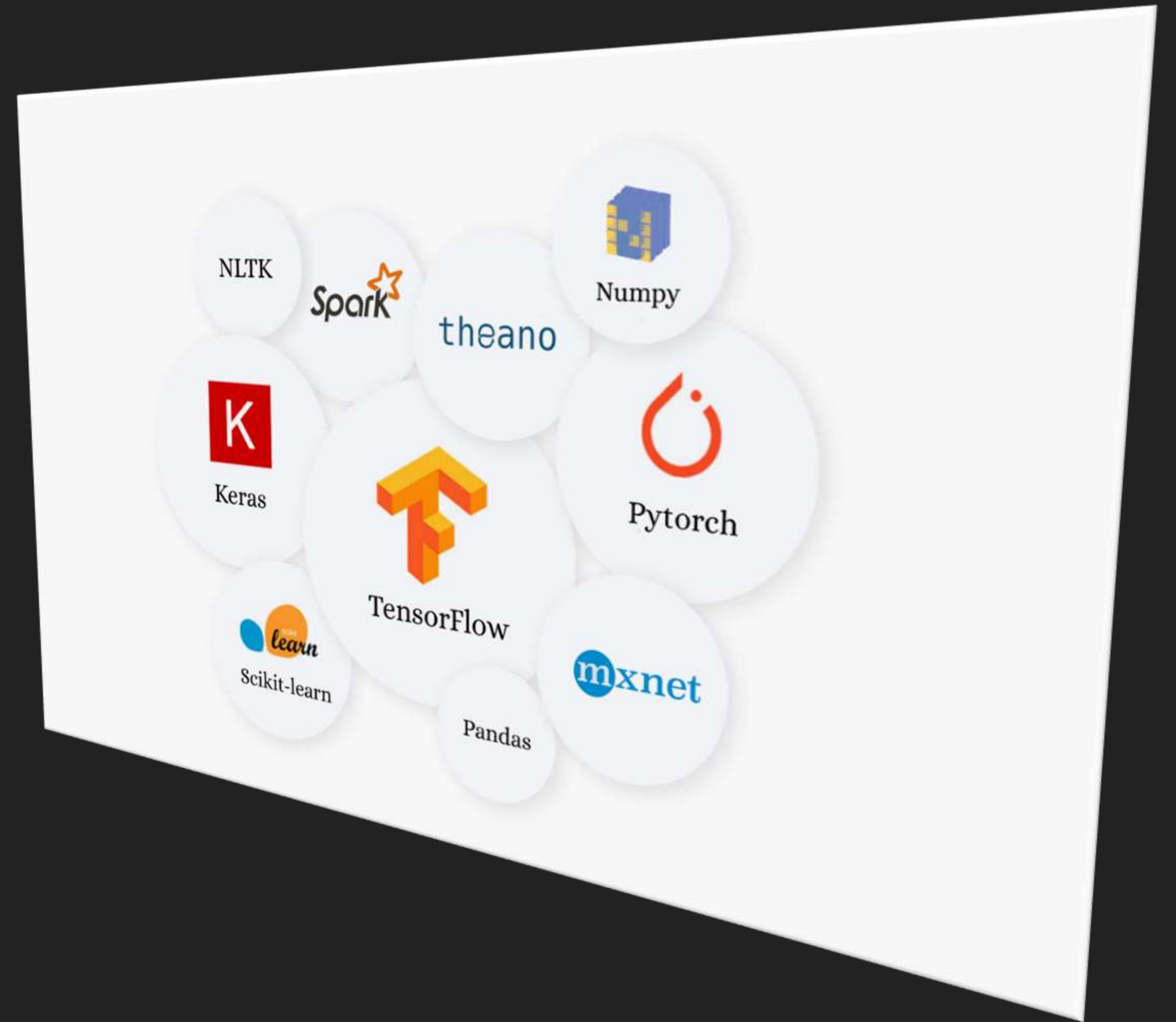
KENTARO TOYAMA

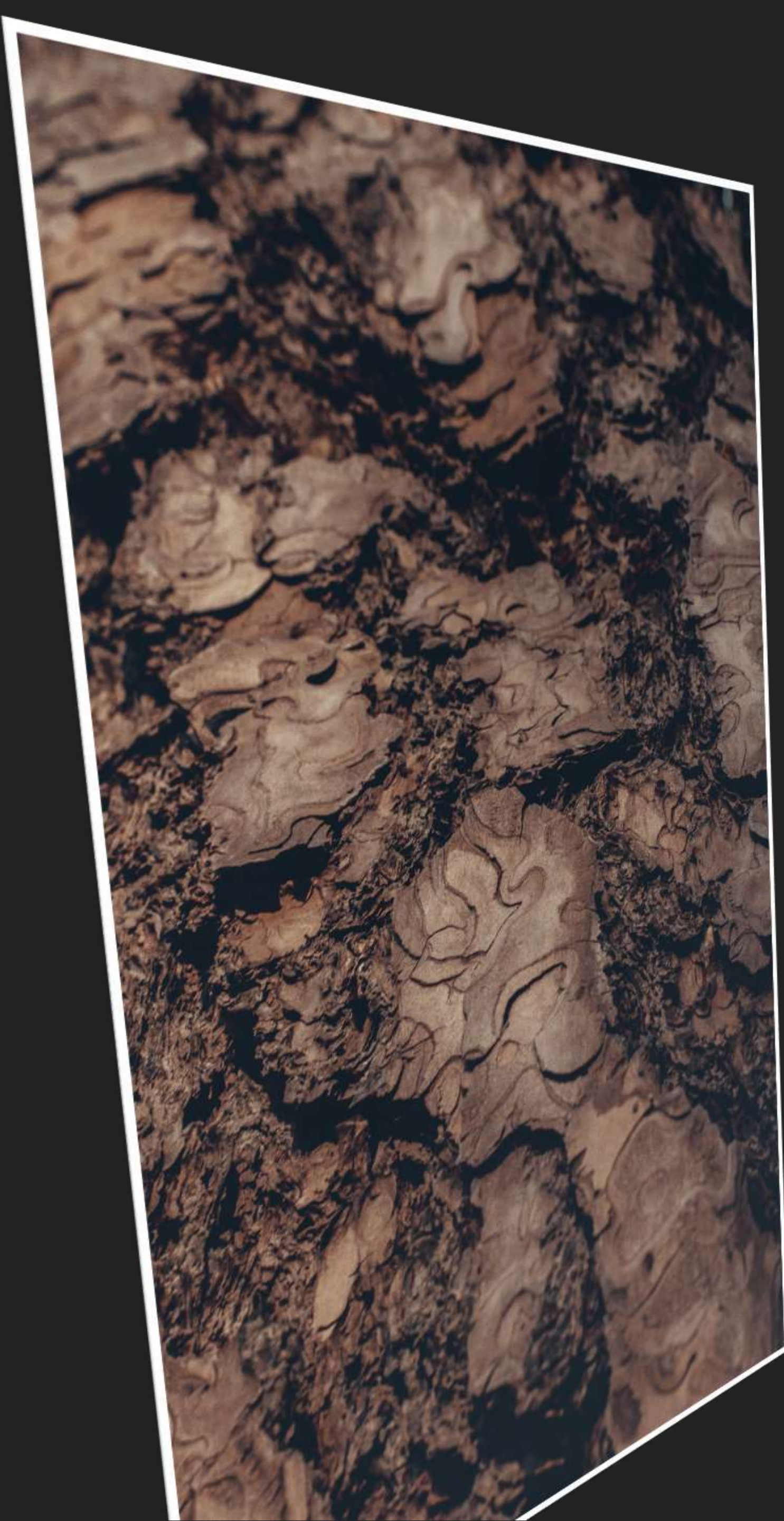
FROM THE CULT OF TECHNOLOGY
RESCUING SOCIAL CHANGE

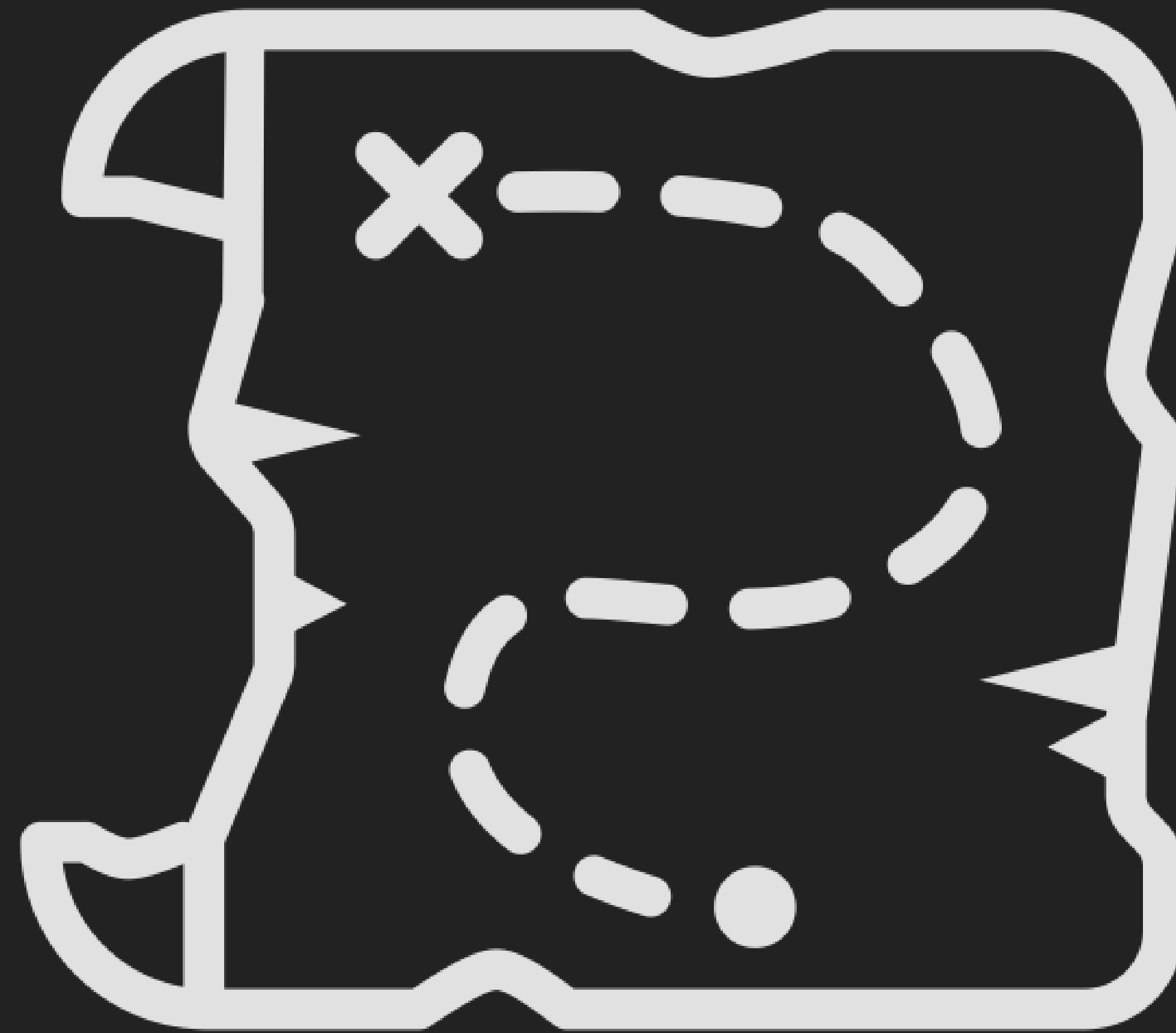












Zuckerberg Takes Steps to Calm Facebook Employees



Mark Zuckerberg
week for his c
Jim Wilson/The

By Sheera I
March 23, 20
SAN FRANCISCO
grappled with
Facebook's
another r

Industry Voice Data Strategy Spotlight

the INQUIRY

Artificial Intelligence Internet of Things Open Source Hardware

Apple Plus 6T > AMD Zen 2 > Samsung Galaxy X > Xiaomi c

PROJECT MAVEN

Report: Google Employees Resigning Over Controversial Pentagon Contract

News > World

Google employees demand transparency over secret work on censored search engine for China

Staff worry they have been unknowingly working on project ‘Dragonfly’ to develop technology helping Chinese government withhold information from citizens

Kate Conger and Daisuke Wakabayashi | Friday 17 August 2018 09:55 | 2 comments



ject

etter to compa

's Joint Enterp

y cause harm.

A nu

UK World Business Football UK politics Environment Education Society Science More

Google Google walkout: global protests after sexual misconduct allegations

Thousands of employees from Tokyo to California stage demonstrations targeting workplace culture



Tech Workers Now Want to Know: What Are We Building This For?



Laura Nolan, a software engineer in Ireland, left Google in June over the company's involvement in Project Maven, an effort to build artificial intelligence for the Department of Defense. Paulo Nunes dos Santos for The New York Times

By **Kate Conger** and **Cade Metz**

Oct. 7, 2018



SAN FRANCISCO — Jack Poulson, a Google research scientist,

ment Science Global development Football Tech Business Obituaries

Google workers demand reinstatement and apology for fired Black AI ethics researcher

Timnit Gebru's departure sparked outrage in the industry as it followed her paper criticizing the company's diversity efforts



nprSIGN INNPR SHOPDONATE

NEWSARTS & LIFEMUSICSHOWS & PODCASTSSEARCH

TECHNOLOGY

Ousted Black Google Researcher: 'They Wanted To Have My Presence, But Not Me Exactly'

December 17, 2020 · 8:28 PM ET

BOBBY ALLYN

Sign inSubscribe

TopicsMagazineNewslettersEvents

Tech policy / AI Ethics

Congress wants answers from Google about Timnit Gebru's firing

The letter, signed by nine members of Congress, sends an important signal about how regulators will scrutinize tech giants.

by Karen Hao

December 17, 2020

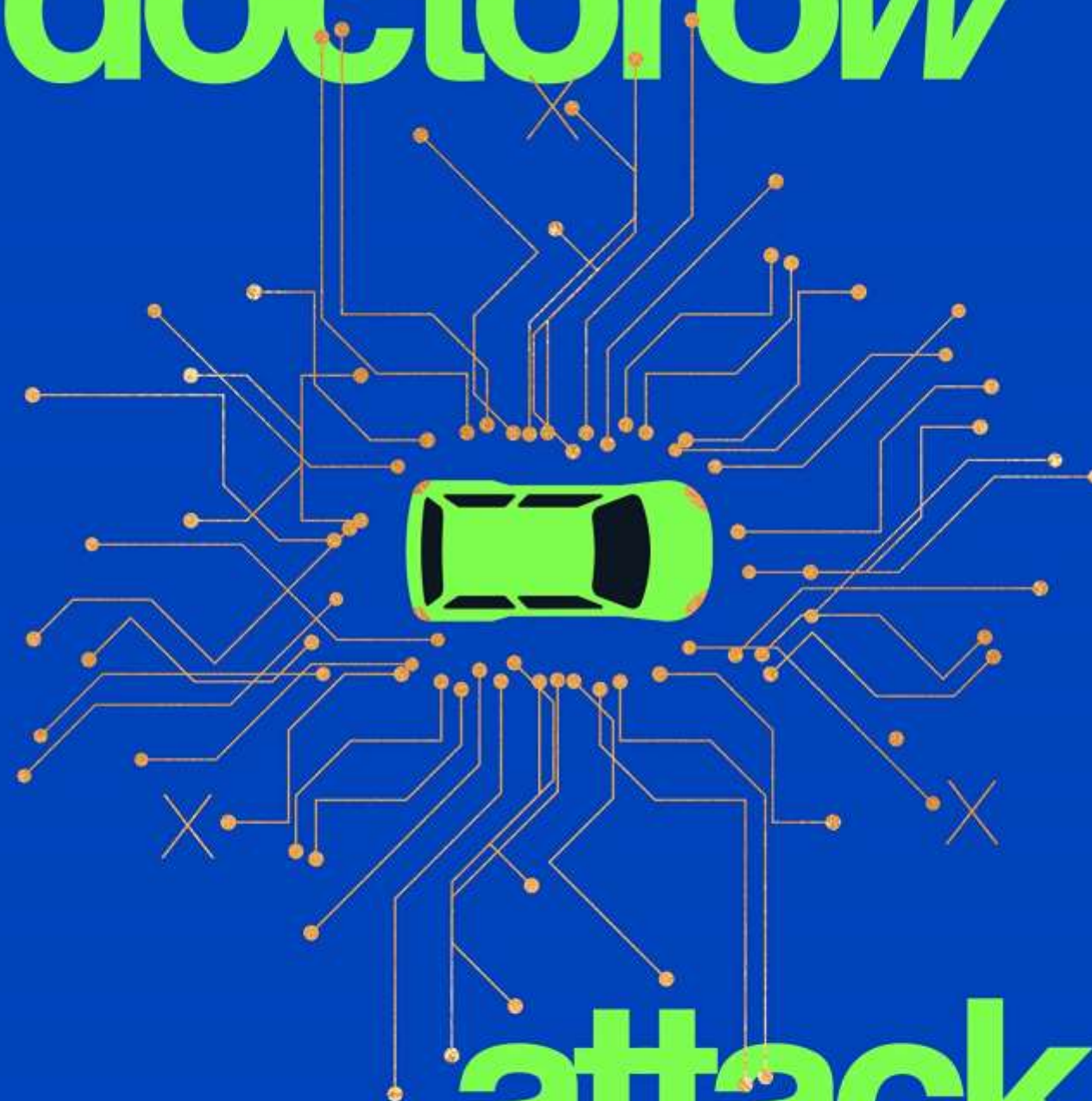




**ANALYSE YOUR
ORG**

**cory
doctorow**

New York Times
bestselling author



**attack
surface**

"One of the
internet's most
interesting
authors."
—Edward
Snowden

The cost of security is everything you believe in.

