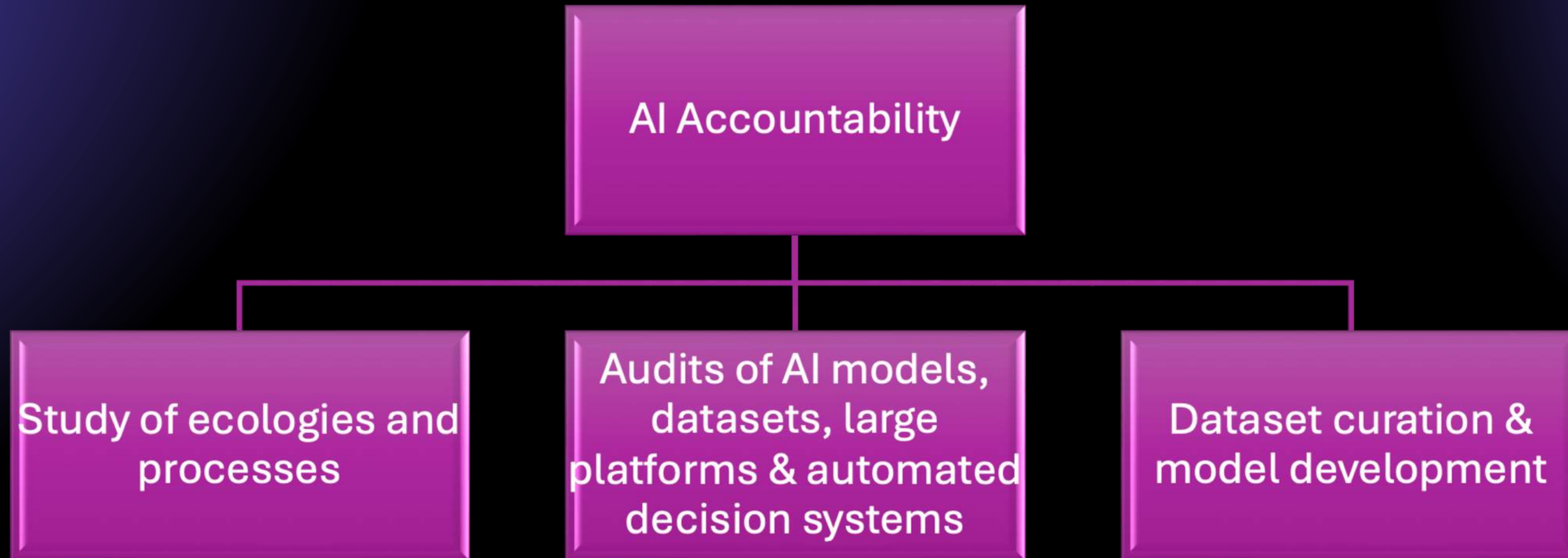


WHAT IS THE INTERNET DOING TO ME? (WITIDTM 2025/2026 - TEU00311)

Abeba Birhane
birhanea@tcd.ie
Lloyd R.031



WHAT WE DO AT THE LAB



**Warning: this lecture contains NSFW
content that some viewers may find
unpleasant and/or offensive**



What is AI?





BRIEF HISTORY OF AI

1937 -- Alan Turing published "On Computable Numbers", which laid the foundations of the modern theory of computation by introducing the Turing machine, a physical interpretation of "computability".

1943 – Walter Pitts and Warren McCulloch analyzed networks of idealized artificial neurons and showed how they might perform simple logical functions, later called a neural network.

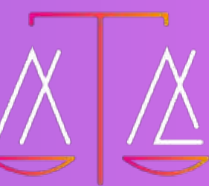
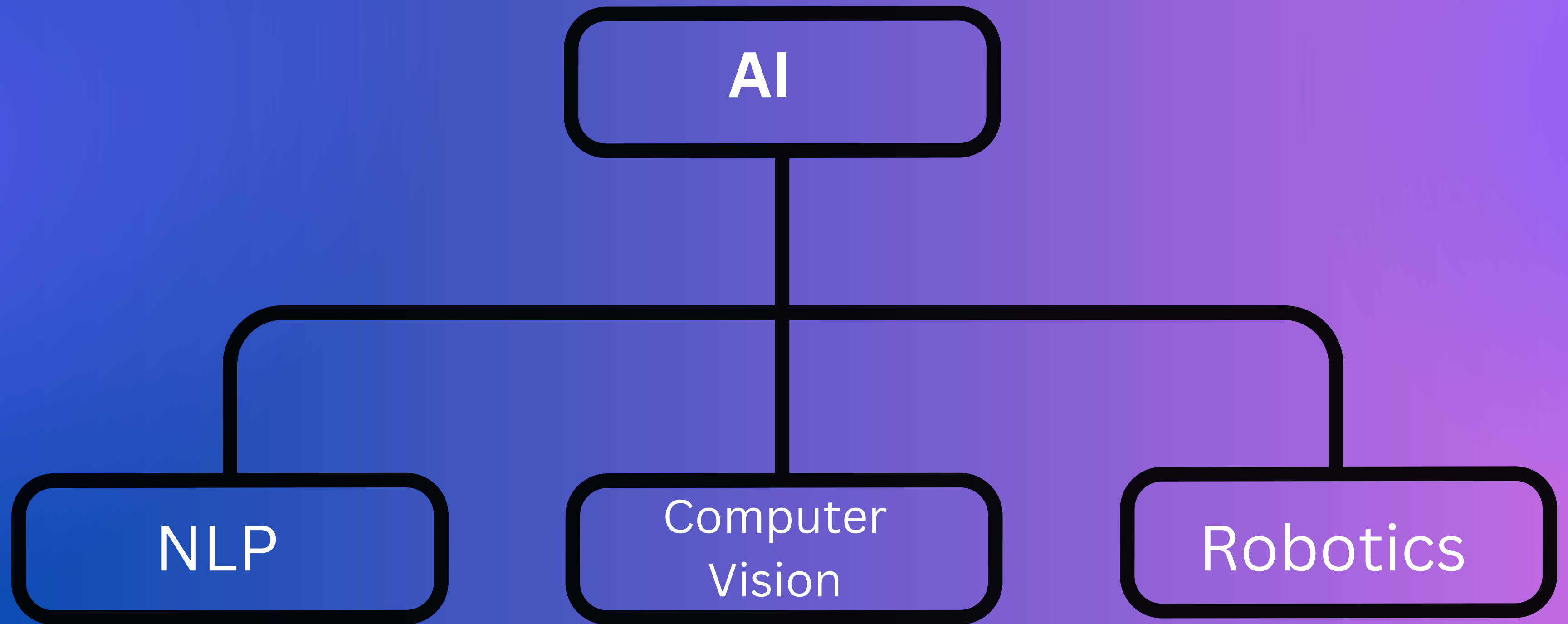
1956 – McCarthy coins the term artificial intelligence for the Dartmouth College summer AI conference

1965 – Joseph Weizenbaum built ELIZA, an interactive program that carries on a dialogue in English language

1980s – Geoffrey Hinton and David Rumelhart popularized a method for training neural networks called "backpropagation" popularising artificial neural networks

1993 – Rodney Brooks and colleagues started the widely publicized MIT Cog project in an attempt to build a humanoid robot child in just five years.





MASSACHUSETTS INSTITUTE OF TECHNOLOGY

PROJECT MAC

Artificial Intelligence Group
Vision Memo. No. 100.

July 7, 1966

THE SUMMER VISION PROJECT

Seymour Papert

The summer vision project is an attempt to use our summer workers effectively in the construction of a significant part of a visual system. The particular task was chosen partly because it can be segmented into sub-problems which will allow individuals to work independently and yet participate in the construction of a system complex enough to be a real landmark in the development of "pattern recognition".



The AI revolution in the late 2010s



Year	Dataset Name	Domain	Size / Content	Significance
1961	Brown Corpus	NLP (Text)	~1 million words	First large-scale balanced English corpus; foundational for NLP.
1985	WordNet	NLP (Lexical)	~155,000 words, lexical database	Semantic network widely used in NLP and knowledge representation.
1990s	Penn Treebank	NLP (Text)	~4.5 million words (annotated)	Key for syntactic parsing and POS tagging.
1998	MNIST	CV (Image)	70,000 handwritten digits (28×28 px)	Classic image recognition benchmark.
2006	80 Million Tiny Images	CV (Image)	~80 million images (32×32 px)	Massive early image dataset; withdrawn later due to bias issues.
2009	ImageNet	CV (Image)	~14 million images (variable size)	Sparked deep learning revolution with AlexNet (2012).
2014	MS COCO	CV (Image+Text)	330,000 images + 1.5 million captions	Object detection and image captioning dataset.
2018	Wikipedia Corpus	NLP (Text)	~2.5 billion words	Large, continuously updated encyclopedia text for LLM pretraining.
2018	Open Images	CV (Image)	9.2 million images	Diverse large-scale image dataset with rich annotations.
2020	The Pile	NLP (Text)	~825 GB text (~300 billion tokens)	Large diverse dataset used for open LLMs like GPT-Neo.
2021	LAION-400M	Multimodal	400 million image-text pairs	Open large-scale dataset for vision-language models.
2022	LAION-5B	Multimodal	5.85 billion image-text pairs	One of the largest open datasets for multimodal AI.

THE DATA AI PIPELINE

**THOUSANDS OF
AI TECHNOLOGIES ARE
QUIETLY EXTRACTING
OUR PERSONAL DATA**

Data about our bodies,
homes, work, social lives...



Can you think of any type of data from your daily activities that might be found in a training corpus?



DuckDuckGo App Tracking Protection for Android

DuckDuckGo Blocked 6 Tracking Attempts



Google

6 attempts. Known to collect:

Available Internal Storage

Local IP Address

OS Build Number

System Volume

Device Orientation

Battery Level

City

GPS Coordinates

Device Brand

OS Version

Headphone Status

Android Advertising ID

Charging Status

App Name

First Name

Device Model

State

Country

Gender

Screen Resolution

Screen Density

Email Address

Device Boot Time

Device Total Memory

Device Name

Network Connection Type

Last Name

Postal Code

App Version

Timezone

CPU Data

Cookies

Device Language

Unique Identifier

THE DATA AI PIPELINE

**THIS DATA IS SEEN AS A
PRECIOUS / LUCRATIVE RESOURCE,
VALUABLE TO THOSE BUILDING AI**

Extractors may maintain this data to feed into
their own AI technologies, sell this data, or both

**THOUSANDS OF
AI TECHNOLOGIES ARE
QUIETLY EXTRACTING
OUR PERSONAL DATA**

Data about our bodies,
homes, work, social lives...



Paul McCartney, Elton John, other creatives demand AI comes clean on scraping

Musicians, artists, writers, actors urge government to protect copyright

 [Lindsay Clark](#)

Mon 12 May 2025 // 12:24 UTC

More than 400 of the UK's leading media and arts professionals have written to the prime minister to back an amendment to the Data (Use and Access) Bill, which promises to offer the nation's creative industries transparency over copyrighted works ingested by AI models.

Signatories include some of the UK's best-known artists such as musicians Paul McCartney, Elton John, Coldplay, writer/director Richard Curtis, artist Antony Gormley, and actor Ian McKellen.

The UK government proposes to allow exceptions to copyright rules in the case of text and data mining needed for AI training, with an opt-out option for content producers.

"Government amendments requiring an economic impact assessment and reports on the feasibility of an 'opt-out' copyright regime and transparency requirements do not meet the moment, but simply leave creators open to years of copyright theft," the letter says.

The group – which also includes Kate Bush, Robbie Williams, Tom Stoppard, and Russell T Davies – said the amendments tabled for the Lords debate would create a requirement for AI firms to tell copyright owners which individual works they have ingested.

ILLEGAL

Last month, Meta admitted to torrenting a controversial large dataset known as LibGen, which includes tens of millions of pirated books. But details around the torrenting were murky until yesterday, when Meta's unredacted emails were made public for the first time. The new evidence showed that Meta torrented "at least 81.7 terabytes of data across multiple shadow libraries through the site Anna's Archive, including at least 35.7 terabytes of data from Z-Library and LibGen," the authors' court filing said. And "Meta also previously torrented 80.6 terabytes of data from LibGen."

"The magnitude of Meta's unlawful torrenting scheme is astonishing," the authors' filing alleged, insisting that "vastly smaller acts of data piracy—just .008 percent of the amount of copyrighted works Meta pirated—have resulted in Judges referring the conduct to the US Attorneys' office for criminal investigation."



Exclusive: OpenAI Used Kenyan Workers on Less Than \$2 Per Hour to Make ChatGPT Less Toxic

15 MINUTE READ



LABELLING AND CLEANING

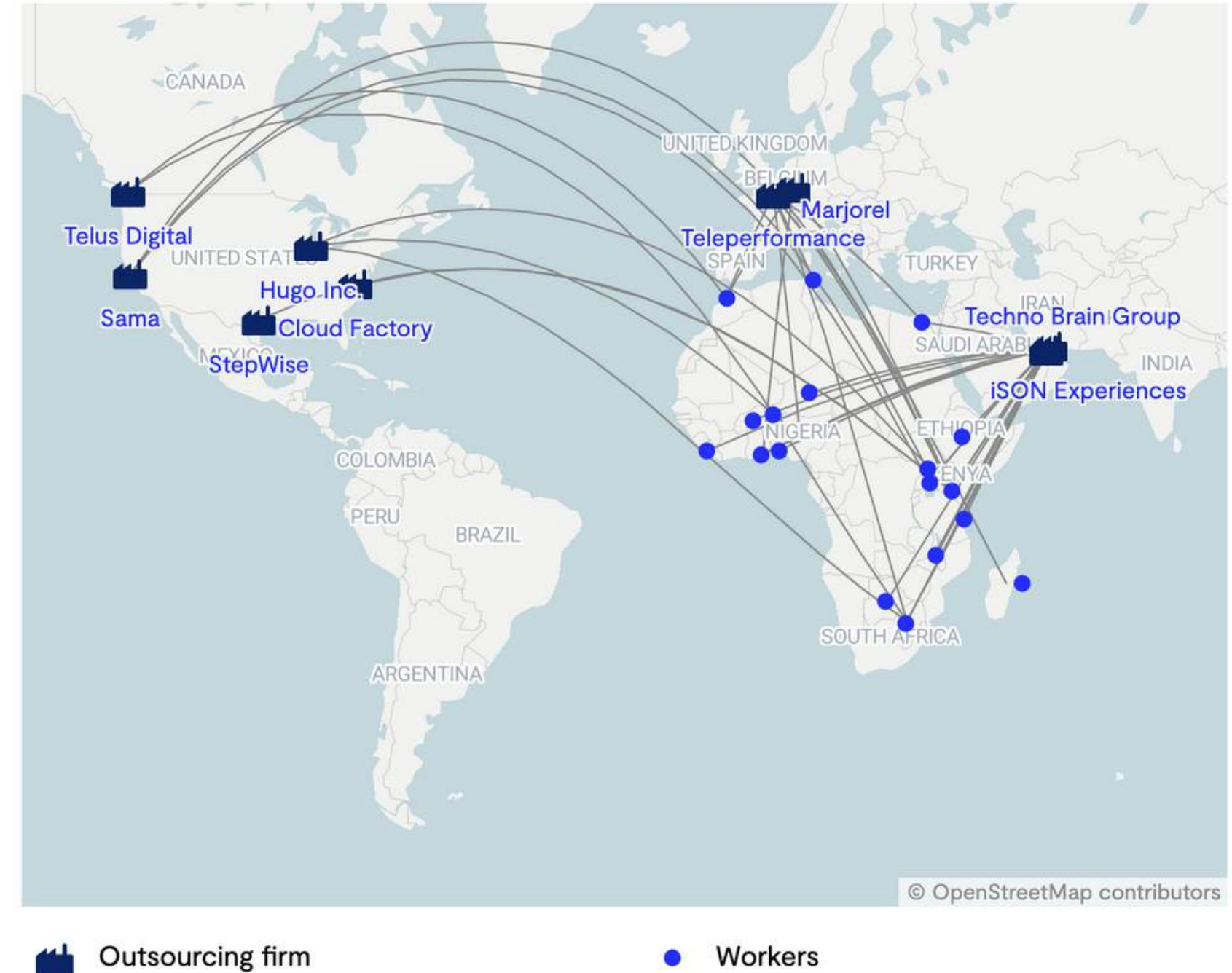
Exclusive: OpenAI Used Kenyan Workers Paid Less Than \$2 Per Hour to Make ChatGPT Safer

15 MINUTE READ



Outsourcing firms and their African offices

Companies in the U.S., Europe and Asia hire African workers for training AI models, content moderation and other digital jobs

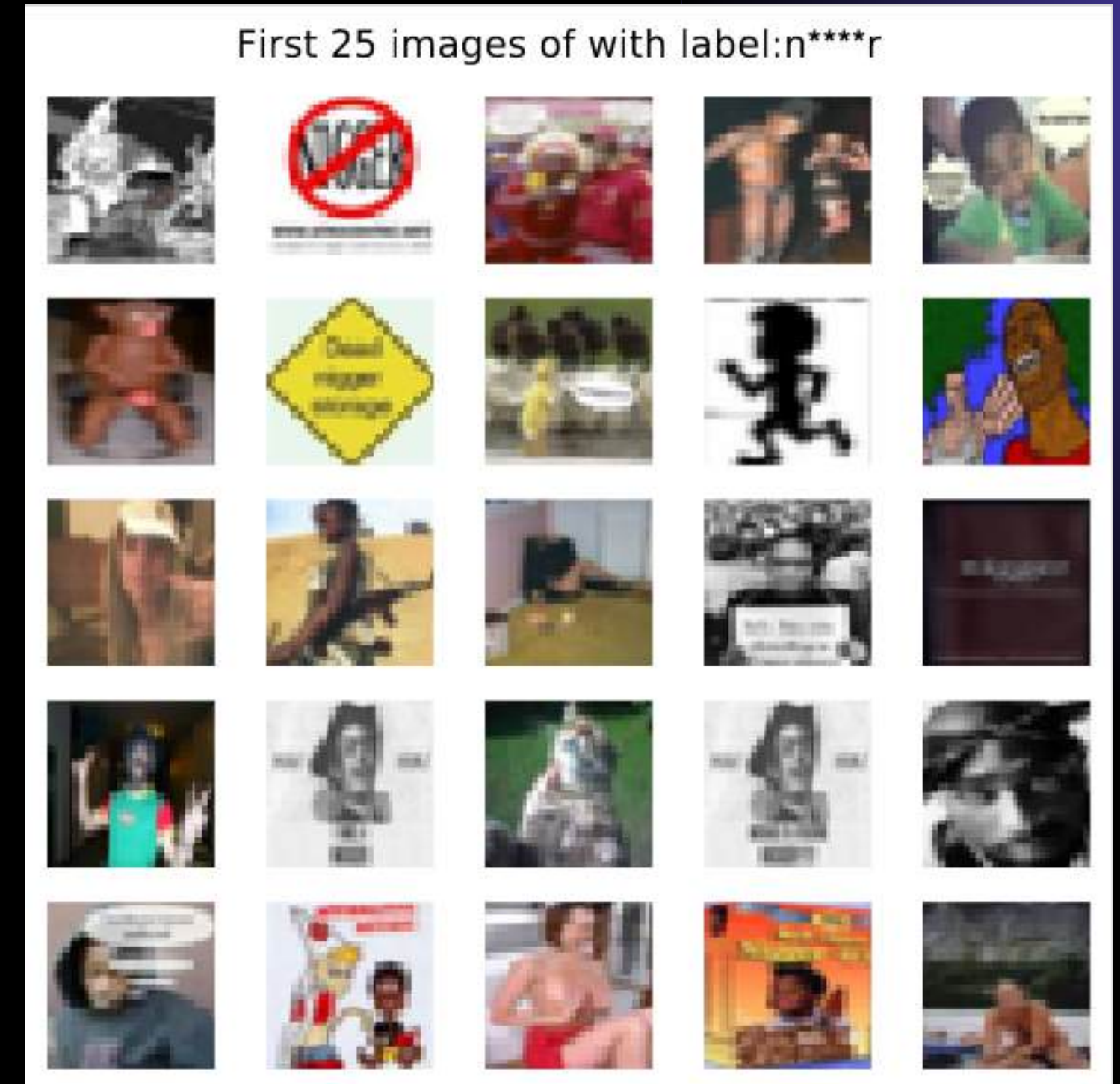
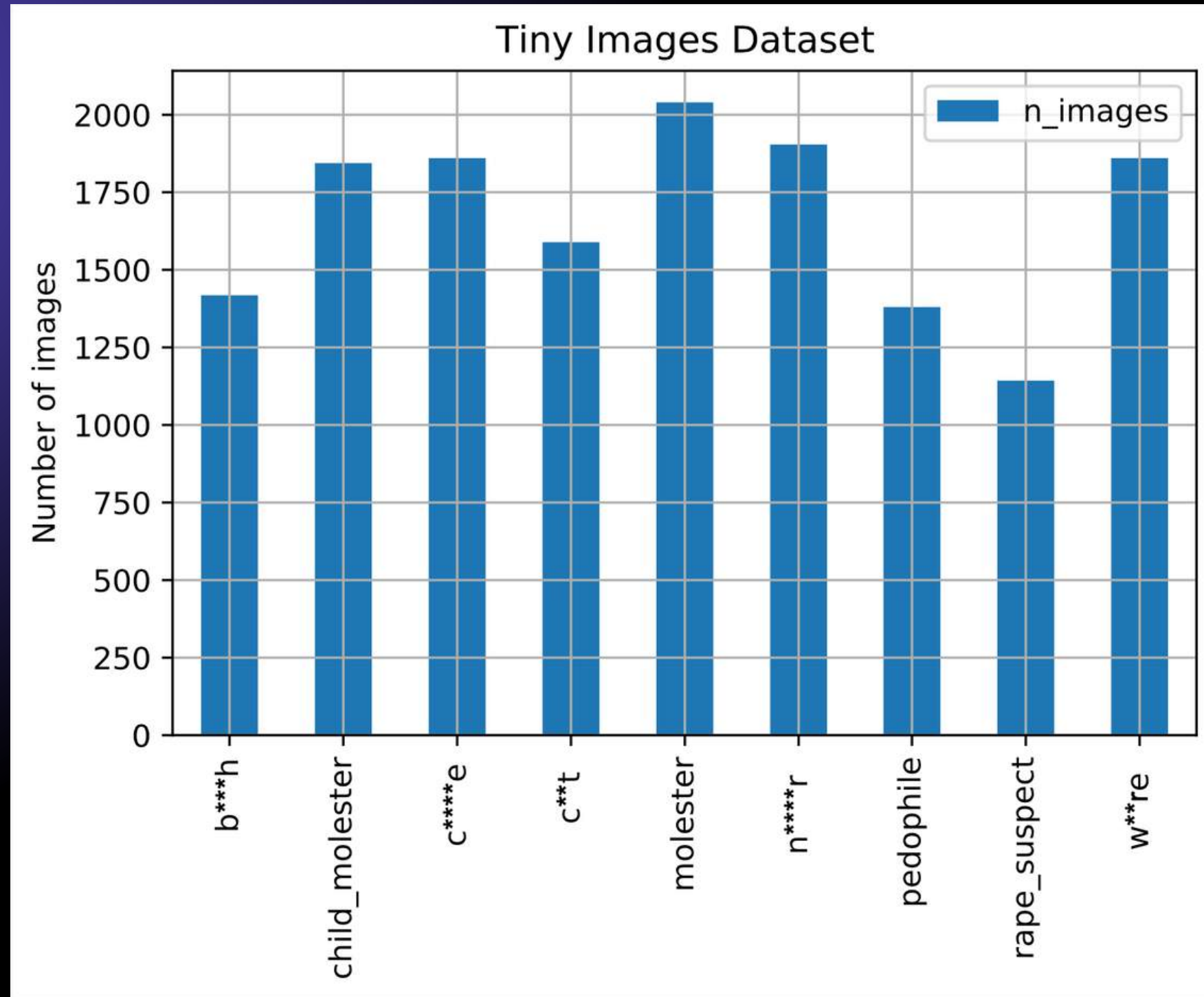


Source: Personaldata.io and the African Content Moderator's Union



Year	Dataset Name	Domain	Size / Content	Significance
1961	Brown Corpus	NLP (Text)	~1 million words	First large-scale balanced English corpus; foundational for NLP.
1985	WordNet	NLP (Lexical)	~155,000 words, lexical database	Semantic network widely used in NLP and knowledge representation.
1990s	Penn Treebank	NLP (Text)	~4.5 million words (annotated)	Key for syntactic parsing and POS tagging.
1998	MNIST	CV (Image)	70,000 handwritten digits (28×28 px)	Classic image recognition benchmark.
2006	80 Million Tiny Images	CV (Image)	~80 million images (32×32 px)	Massive early image dataset; withdrawn later due to bias issues.
2009	ImageNet	CV (Image)	~14 million images (variable size)	Sparked deep learning revolution with AlexNet (2012).
2014	MS COCO	CV (Image+Text)	330,000 images + 1.5 million captions	Object detection and image captioning dataset.
2018	Wikipedia Corpus	NLP (Text)	~2.5 billion words	Large, continuously updated encyclopedia text for LLM pretraining.
2018	Open Images	CV (Image)	9.2 million images	Diverse large-scale image dataset with rich annotations.
2020	The Pile	NLP (Text)	~825 GB text (~300 billion tokens)	Large diverse dataset used for open LLMs like GPT-Neo.
2021	LAION-400M	Multimodal	400 million image-text pairs	Open large-scale dataset for vision-language models.
2022	LAION-5B	Multimodal	5.85 billion image-text pairs	One of the largest open datasets for multimodal AI.

WHAT'S IN THE DATA



Birhane & Prabhu (2021)

Year	Dataset Name	Domain	Size / Content	Significance
1961	Brown Corpus	NLP (Text)	~1 million words	First large-scale balanced English corpus; foundational for NLP.
1985	WordNet	NLP (Lexical)	~155,000 words, lexical database	Semantic network widely used in NLP and knowledge representation.
1990s	Penn Treebank	NLP (Text)	~4.5 million words (annotated)	Key for syntactic parsing and POS tagging.
1998	MNIST	CV (Image)	70,000 handwritten digits (28×28 px)	Classic image recognition benchmark.
2006	80 Million Tiny Images	CV (Image)	~80 million images (32×32 px)	Massive early image dataset; withdrawn later due to bias issues.
2009	ImageNet	CV (Image)	~14 million images (variable size)	Sparked deep learning revolution with AlexNet (2012).
2014	MS COCO	CV (Image+Text)	330,000 images + 1.5 million captions	Object detection and image captioning dataset.
2018	Wikipedia Corpus	NLP (Text)	~2.5 billion words	Large, continuously updated encyclopedia text for LLM pretraining.
2018	Open Images	CV (Image)	9.2 million images	Diverse large-scale image dataset with rich annotations.
2020	The Pile	NLP (Text)	~825 GB text (~300 billion tokens)	Large diverse dataset used for open LLMs like GPT-Neo.
2021	LAION-400M	Multimodal	400 million image-text pairs	Open large-scale dataset for vision-language models.
@ A B E B A . B S K Y . S O C I A L 2022	LAION-5B	Multimodal	5.85 billion image-text pairs	One of the largest open datasets for multimodal AI.

WHAT'S IN THE DATA

Backend url:
<https://clip.roi>
Index:
laion_400m

african



Clip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates? KNN search are good at spotting those, especially so in large datasets.



Stylish African Print Swimwear Headwrap - Mahalia



Mursi woman / omo va...



Miss Domoget, Bodi Tribe Woman With Headband, Hana...



Erbore tribe woman in Ethiopia on October 26 2008 ...



Himba People by Konstantinos Arvanitopoulos



Stylish African Print Swimwear Headwrap - Mahalia



Himba People by Konstantinos Arvanitopoulos



A girl from the Hamar tribe, Ethiopia



Wodaabe boy from Niger. Photographed by Steve McCu...



Portrait of a Himba Woman



Turmi, Omo River Valley, Ethiopia - January, 2018....



Himba, girl with typical headdress and decoration ...



Stylish African Print Swimwear Headwrap - Mahalia



Himba People by Konstantinos Arvanitopoulos



Know Who You Are « Steve McCurry's Blog



Massai Girl



Wodaabe boy from Niger (by Steve McCurry)



Wodaabe tribe--They are traditionally nomadic catt...



Portrait of a Himba Woman



WHAT'S IN THE DATA

Backend url:
https://clip.roi


Index:
laion_400m

Clip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings


Display captions ☒
 Display full captions ☐
 Display similarities ☐
 Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.


Are you seeing near duplicates ? KNN search are good at spotting those, especially so in large datasets.










Stylish African Print Swimwear Headwrap - Mahalia



Mursi woman / omo va...



Backend url:
https://clip.roi


Index:
laion_400m

Clip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings


Display captions ☒
 Display full captions ☐
 Display similarities ☐
 Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.


Are you seeing near duplicates ? KNN search are good at spotting those, especially so in large datasets.




map outline of europe Free Interior Design ...




Germanysoviet union relations 19181941 wikipedia g...




Map indicating locations of Greece and Latvia




Americans were asked to place european countries o...




1300x866 3d Map Of Europe Rendering On White Backg...




Vector illustration of the European Union map with...




Map indicating locations of Germany and Netherland...




germany 3d model




Germany location map




germany illustrator map




Location of Timișoara




Ottoman Empire Capital A Map Of Europe Showing Ter...




Germanynetherlands relations wikipedia map indicat...




Giurtelecu Simleului in Europe.jpg




Crusader Kings II: The Old Gods Youtube Video




Map Of Germany And Russia.Germany Soviet Union Rel...



Crusader Kings II: The Old Gods gets release date



Can you identify this sea? The Black Sea



germany location map file northwest germany locati...



WHAT'S IN THE DATA

Backend url:
https://clip.roi

Index:
laion_400m

beautiful

Clip retrieval works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Safe mode ☐
Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates? KNN search are good at spotting those, especially so in large datasets.

brazilian bikini bottom metallic gray

Green eyes with Pink Nipples

Selena Gomez Just Made Her First Red Carpet Appear...

chubby big boobs

amateur photo Bryce Dallas Howard

Emma: Blonde Sex Doll

Emma: Blonde Sex Doll

SUPER Lingerie Try On Haul 18 warning from YouTube...

metallic gray triangle bikini top

Topless Photos of Holly Peers - Celeb Nudes

Online Auto Insurance >> hairstyles for men: Hair ...

s Nude

Milf M... Seduces her son TABOO MOM SON

VERA BOTTOM - WHITE

P... attends Smurfs: The Lost Village Press...

wiki, affair, married, Lesbian with ag...

Ombre Gray Synthetic Lace Front Wig HS0021



WHAT'S IN THE DATA

Backend url:
<https://clip.roi>
Index:
laion_400m

beautiful



[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Safe mode ☐
Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates ? KNN search are good at spotting those, especially so in large datasets.

Backend url:
<https://clip.roi>
Index:
laion_400m

handsome



[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Safe mode ☐
Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates ? KNN search are good at spotting those, especially so in large datasets.



zac efron at any price
venice premiere
142711919



More suits, #menstyle,
style and fashion for
men @...



The well fitted suit
with classic and
modern fit



1000 ideas about Grey
Suit Black Shirt on
Pintere...



Best Asian Men
Hairstyles - Star Styles
| StylesSt...



"More T.O.P. for
""Cosmopolitan
China"" [PHOTO] - ...



Beauty And Body Of
Male : Lee Min Ho For
Harper's ...



So damn-gorgeous-
handsome Lee
Donghae!
ARGHFJKGLFJ...



Boys Three-Piece
Plaid Suit



16/100 pictures of
Daniel Radcliffe



zac efron Hairstyle,
Male, Fashion, Men,
Amazing, ...



park hae jin age - B
Asian Celebrities,
Asian Acto...



TOP Shanghai Press
Con OUT OF
CONTROL
2016-06-14 (...)



The well fitted suit
with classic and
modern fit



Portuguese
professional footballer
Cristiano Ronal...



Boys Three-Piece
Plaid Suit



SF9 member Chani



WHAT'S IN THE DATA

Backend url:

Index:

beautiful



[Clip retrieval](#) works by converting the text query to a CLIP embeddir , then using that embedding to query a knn index of clip image embedddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Safe mode ☐
Search over

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates ? KNN search are good at spotting those, especially so in large datasets.

Backend url:

Index:

handsome



[Clip retrieval](#) works by converti the text query to a CLIP embedding , then using that embedding to query a knn index of clip image embedddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Safe mode ☐
Search over

This UI may contain results wit nudity and is best used by adult The images are under their own copyright.

Are you seeing near duplicates ? KNN search are good at spottin those, especially so in large datasets.

Backend url:

Index:

terrorist



[Clip retrieval](#) works by converting the text query to a CLIP embedding , then using that embedding to query a knn index of clip image embedddings

Display captions ☒
Display full captions ☒
Display similarities ☒
Search over

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates ? KNN search are good at spotting those, especially so in large datasets.

0.1522



foto of ak 47 - Portrait of serious eastern man with AK - JPG

0.1487



Terrorist Leader On Tv Screen Streaming Television Terrorism Vector Illustration

0.1480



picture of terrorist - Portrait of serious eastern man with AK - JPG

0.1475



picture of ak 47 - Portrait of serious eastern man with AK - JPG

0.1473



Islamic State executioner Hicham Chaib in the video. Courtesy: YouTube

0.1471



skull terrorist masked and Kalashnikov machine guns. Isolated objects on a white background can be used with any image or text.

0.1466



Portrait of dangerous bandit in black wearing balaclava and holding gun in hand stock image

0.1461



Man With Gun And Peace Dove.

0.1461



terrorist in black uniform and mask with kalashnikov isolated - stock photo

0.1454



stock photo of terrorist - Portrait of serious eastern man with AK - JPG

0.1452



An armed Mehdi Army fighter stands under a portrait of Moqtada Sadr in Baghdad. File photo

0.1452



a-lot-of-information-was-revealed-by-baghdadi-s-wife-after-capture

0.1452



Armed Terrorist Group Terrorism Concept Flat Vector Illustration

0.1450



picture of ak-47 - Portrait of serious eastern man with AK - JPG

0.1445



Candid videos show rare view of unkempt bin Laden

0.1445



A picture of Abu Mousab al-Zarqawi, crossed out by a red X

WHAT'S IN THE DATA

Backend url:
<https://clip.roi>
Index:
laion_400m

beautiful



[Clip retrieval](#) works by converting the text query to a CLIP embedder, then using that embedding to query a knn index of clip image embeddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Safe mode ☐
Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates? KNN search are good at spotting those, especially so in large datasets.

Backend url:
<https://clip.roi>
Index:
laion_400m

handsome



[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Safe mode ☐
Search over image

This UI may contain results with nudity and is best used by adult. The images are under their own copyright.

Are you seeing near duplicates? KNN search are good at spotting those, especially so in large datasets.

Backend url:
<https://clip.roi>
Index:
laion_400m

terrorist



[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions ☒
Display full captions ☒
Display similarities ☒
Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates? KNN search are good at spotting those, especially so in large datasets.

Backend url:
<https://clip.roi>
Index:
laion_400m

ceo



[Clip retrieval](#) works by converting the text query to a CLIP embedding, then using that embedding to query a knn index of clip image embeddings

Display captions ☒
Display full captions ☐
Display similarities ☐
Search over image

This UI may contain results with nudity and is best used by adults. The images are under their own copyright.

Are you seeing near duplicates? KNN search are good at spotting those, especially so in large datasets.



Businessman poses with pen while sitting on an off...



young business man on a desk, isolated on white



Young and determined royalty-free stock photo



handsome Young business man sitting on a chair



Smiling businessman stock photo



Airport Business : Stock Photo



Businessman with feet up at desk



Businessman Hands Paying Folder Ceo Concept On Bro...



Businessman with folded arms leaning back satisfie...



Indian Businessman royalty-free stock photo



Businessman



Portrait of modern businessman sitting at office d...



Handsomen smiling help-desk male executive isolated...



Portrait of a confident Arab businessman sitting o...



Businessman leaning back satisfied



Smiling businessman stock photo



Office Interior. A Man In A Business Suit At A Tab...



Portrait of two contemporary businessmen, one of t...



Smiling business man in suit isolated on white — S...

WHAT'S IN THE DATA

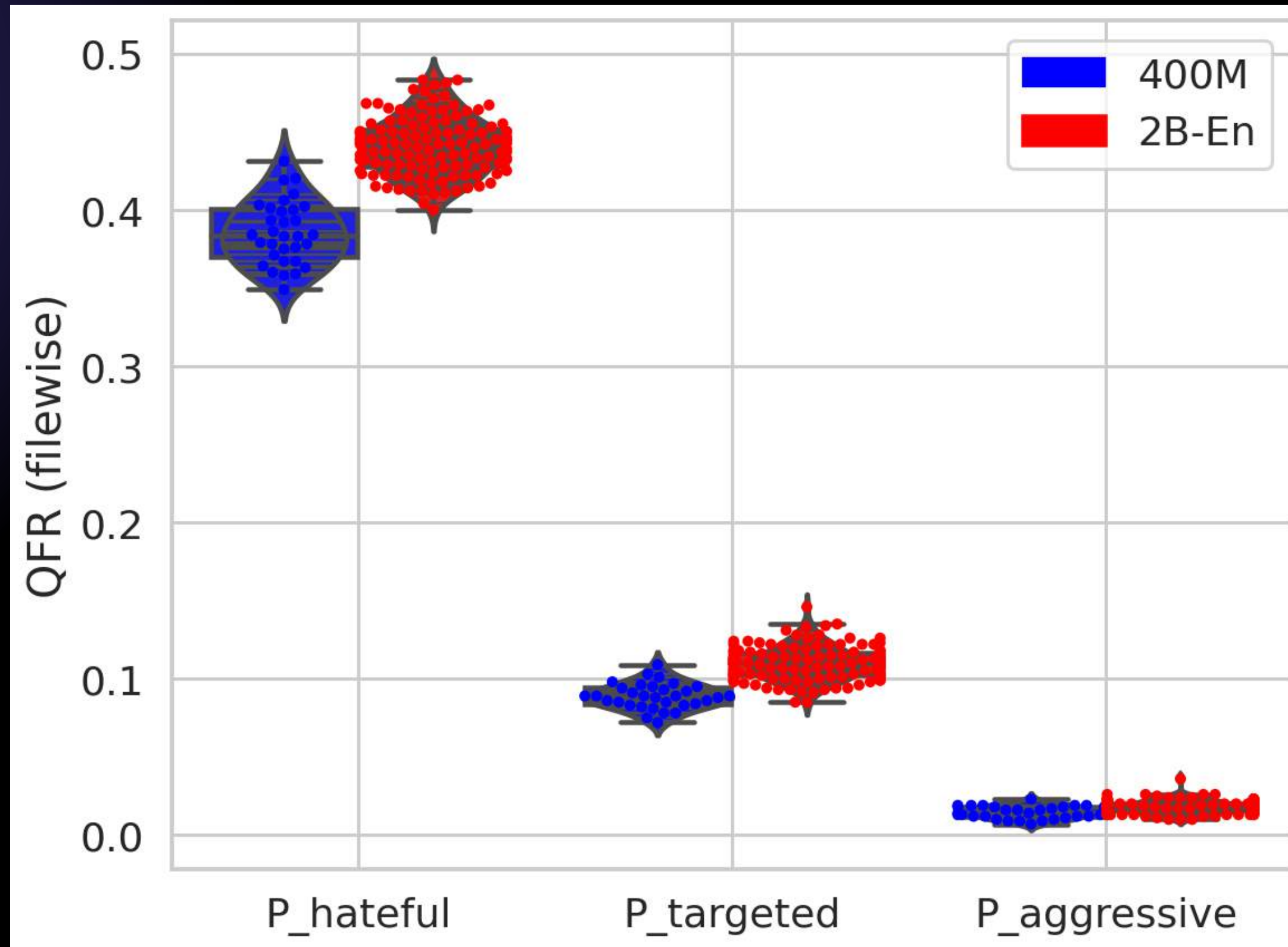
Table 1: Samples of alt text descriptions found in the dataset and the probability scores across the three categories of *hateful*, *targeted* and *aggressive* speech.

Alt text	$P_{hateful}$	$P_{targeted}$	$P_{aggressive}$
'Biden's Spending Will Go To Illegal Immigrants While Tax Hikes Will Destroy American Jobs'	0.902	0.024	0.449
'If you know this man, please, for the love of God tell him to BURN these pants!!'	0.401	0.262	0.517
'shut up and be a don like nancy - Personalised Men's Long Sleeve T-Shirt'	0.395	0.559	0.128
'This bored rich blonde shoplifter gets rough f**keds'	0.934	0.895	0.128
'Horny slave tied to tree gets pulled on her beautiful tits and gets hit on her c*nt with a stick and hands'	0.983	0.911	0.909

[Into the laion's den: Investigating hate in multimodal datasets](#) (Birhane et al., 2024)



WHAT'S IN THE DATA



Into the laion's den: Investigating hate in multimodal datasets (Birhane et al., 2024)

What geographies, cultures and representations are dominant in major datasets?



REPRESENTATION

3.3 GEOGRAPHICAL & LINGUISTIC REPRESENTATION IS NOT IMPROVING

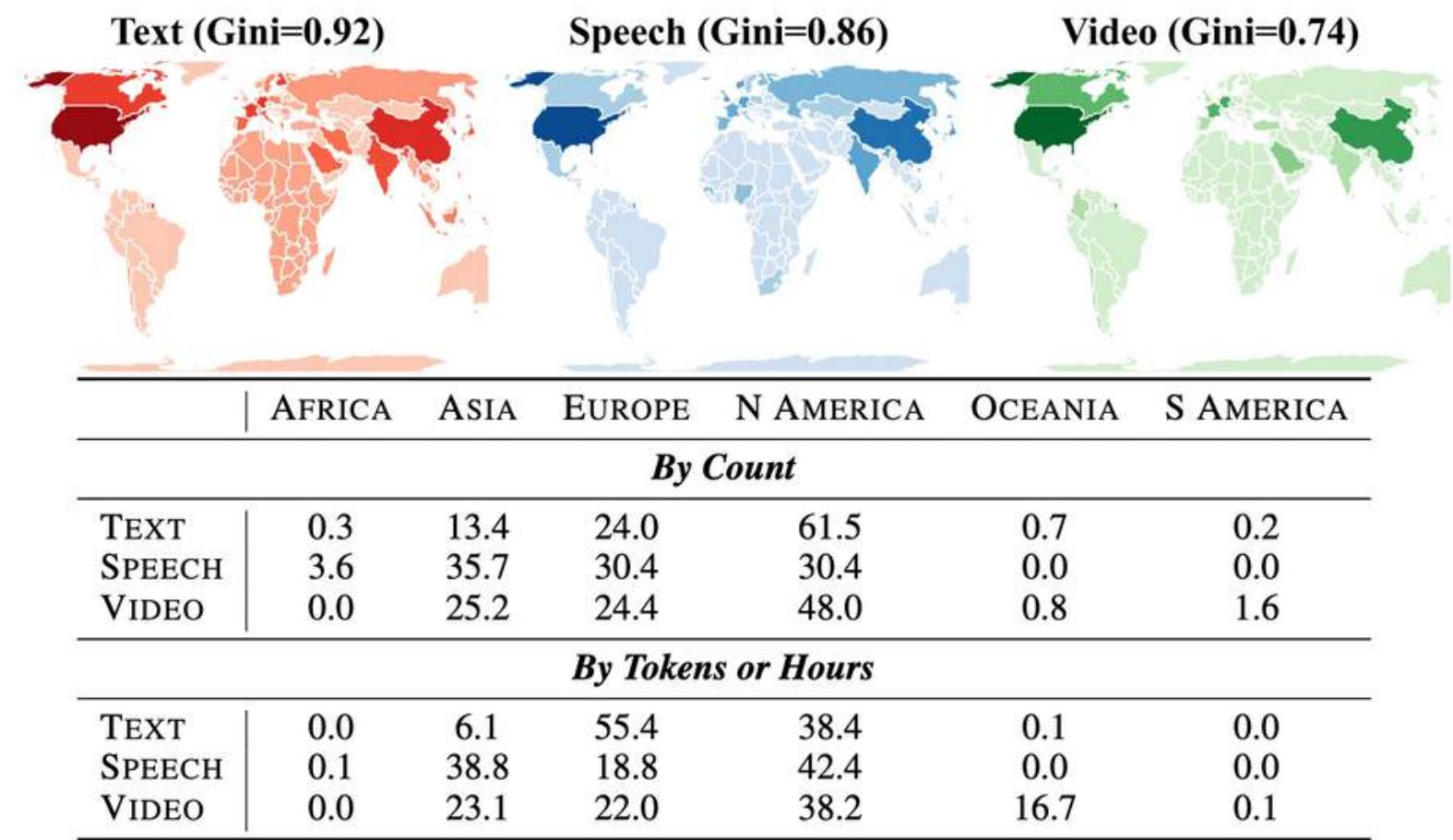
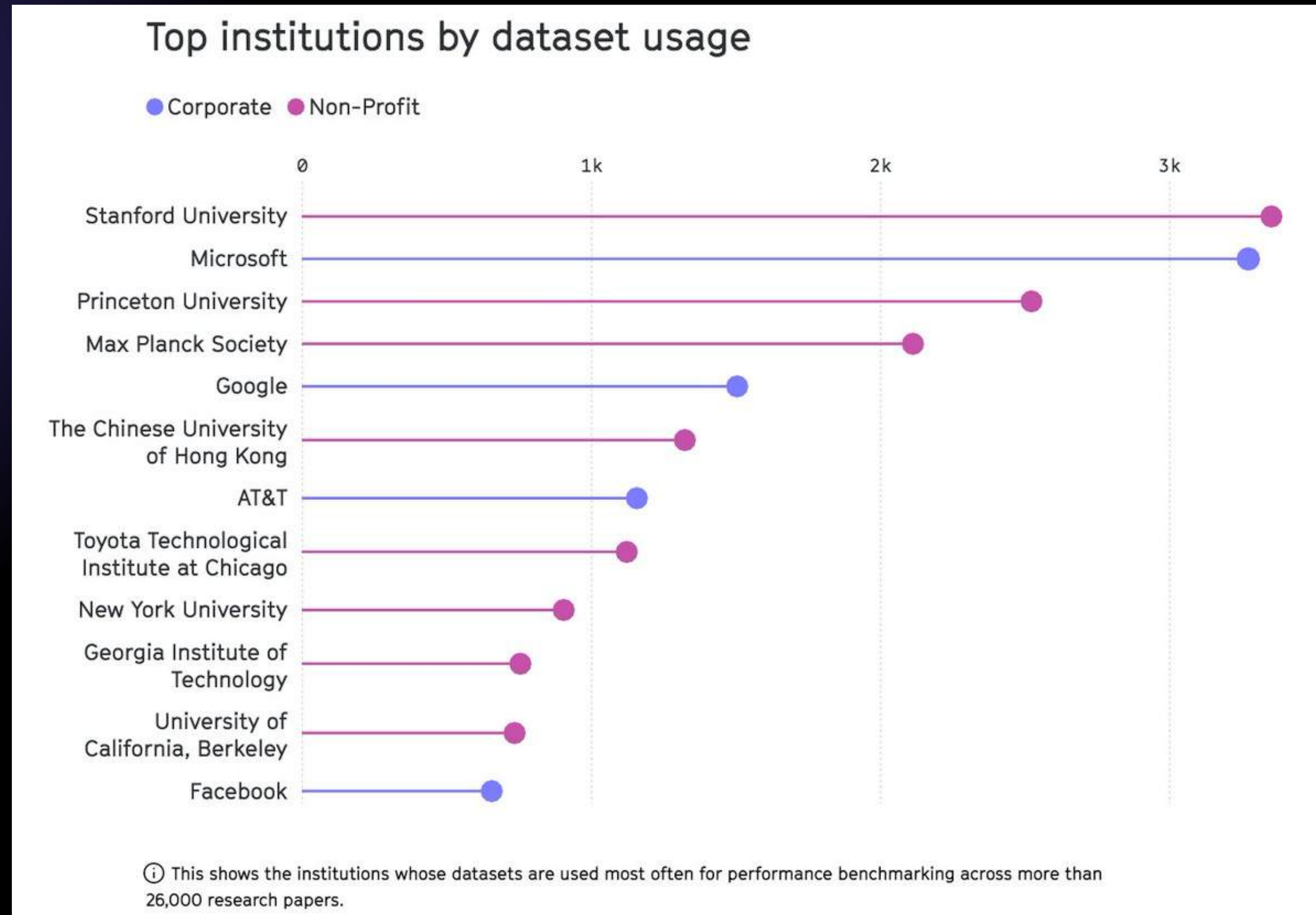


Figure 3: The geographical distribution of countries (world maps) and continents (table) represented by dataset creators. **Despite some differences in European, Russian, and Middle Eastern representation, creators are heavily concentrated in the US, China, and Western Europe, with little to no representation in South America or Africa, across modalities.** The current Gini coefficient for (Text, Speech, Video) = (0.92, 0.86, 0.74), where higher values indicate more concentration.

- Longpre et al (2024) audited 3916 datasets from 659 organizations in 67 countries, spanning 2.1T tokens, and 1.9M hours.
- Inequality in geographical representation remains very high, with few organizations creating datasets from the Global South.
- Multilingual representation has not improved by most measures



REPRESENTATION



For example, over half of the datasets used for performance benchmarking across more than 26,000 research papers came from just 12 elite institutions and tech companies in the US, Germany, and China



Websites, books, code repositories, forums, scientific articles, etc

- via public pools such as the common crawl
 - remove duplicates
 - remove HTML/Markup (scripts, styles, non-text elements)
 - remove boilerplate (headers, footers, navigation links)
- toxicity filtering
 - hate speech, violence, abuse, CSAM (child sexual abuse material)



- **List-of-Dirty-Naughty-Obscene-and-Otherwise-Bad-Words**
- **Over 400 words removed from the C4 (Colossal Clean Crawled Corpus)**



DATA CLEANING AND DETOXIFICATION

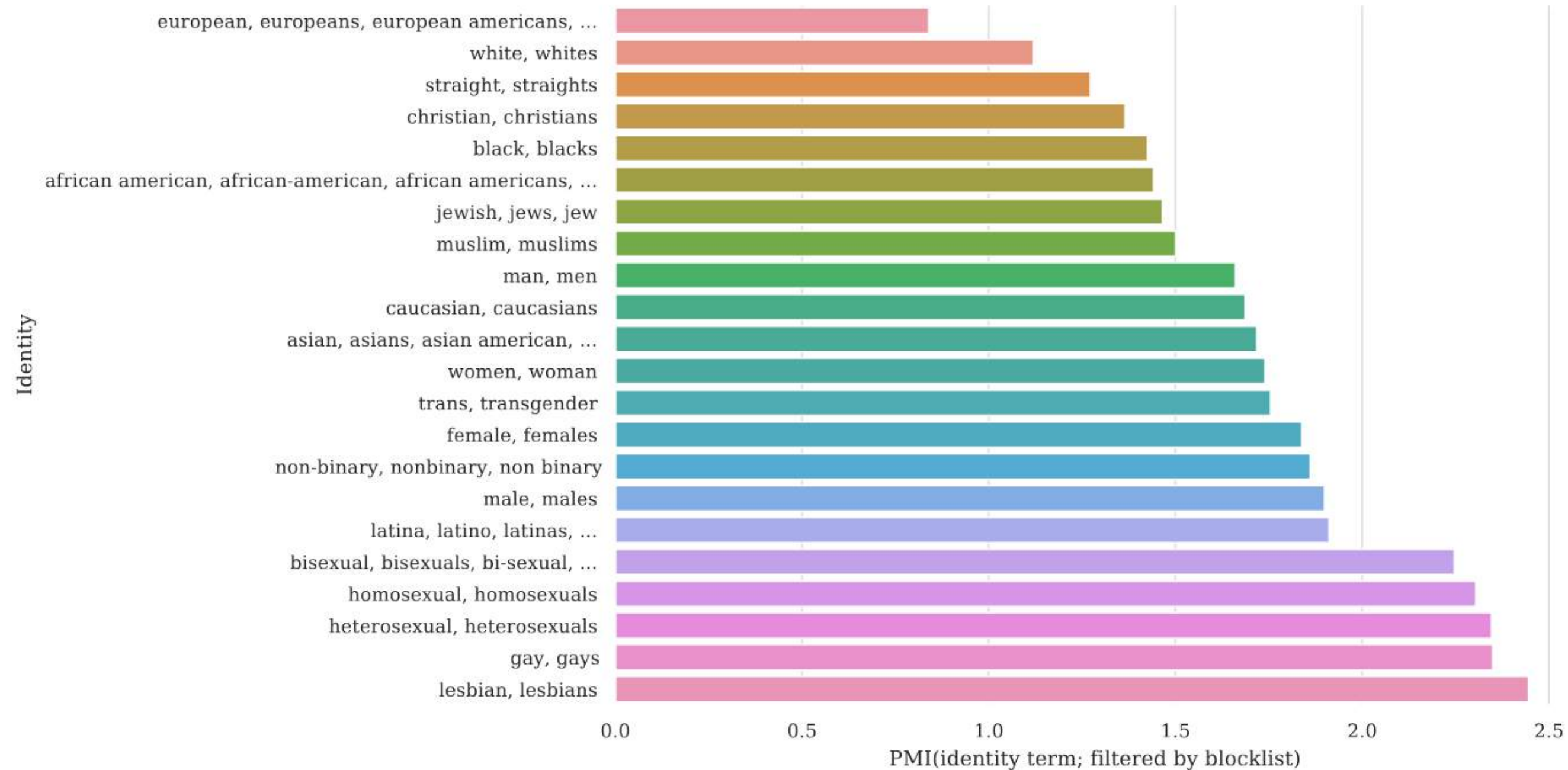


Figure 5: Pointwise Mutual Information (PMI) between identity mentions and documents being filtered out by the blacklist. Identities with higher PMI (e.g., lesbian, gay) have higher likelihood of being filtered out.

USING MODERN AI, THIS DATA IS NEVER SINGLE-PURPOSE

Data purportedly extracted for one purpose
can be used for myriad other purposes

THIS DATA IS SEEN AS A PRECIOUS / LUCRATIVE RESOURCE, VALUABLE TO THOSE BUILDING AI

Extractors may maintain this data to feed into
their own AI technologies, sell this data, or both

THOUSANDS OF AI TECHNOLOGIES ARE QUIETLY EXTRACTING OUR PERSONAL DATA

Data about our bodies,
homes, work, social lives...



Computer Vision

refers to AI that attends to visual inputs such as image and video data for purposes of measuring, mapping, recording, and monitoring the world.



Computer Vision

refers to AI that attends to visual inputs such as image and video data for purposes of measuring, mapping, recording, and monitoring the world.

As a technology that emerged in military contexts, it was historically developed to identify targets and gather intelligence in war, law enforcement, and immigration contexts.

The field of computer vision now generally emphasizes training computers to interpret and understand the visual world.



Computer Vision highlighted topics

AI as
science-like

Inevitable &
application
agnostic

Revisiting Old Ideas With Modern Hardware

Learning to see the human way

Toward Integrative AI with Computer Vision

An AI Odyssey: the Dark Matter of Intelligence

...

...

...

AI as
engineering
for good

Benevolent
applications

Modeling Atoms to Address Our Climate Crisis

Understanding Visual Appearance from Micron to
Global Scale

...

...



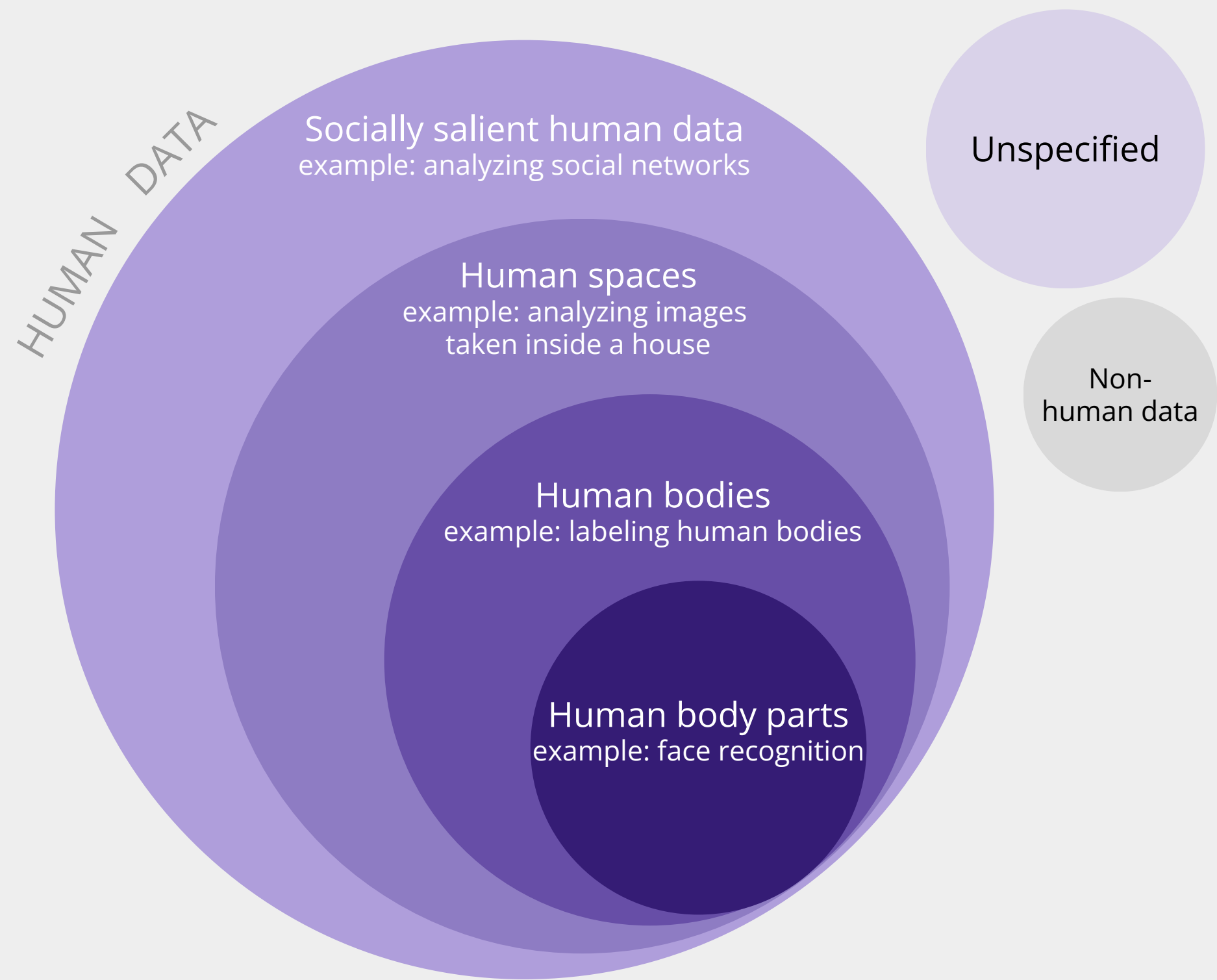
Method: quantitative

- We parsed over three decades of computer vision research papers and downstream patents, over 40,000 documents

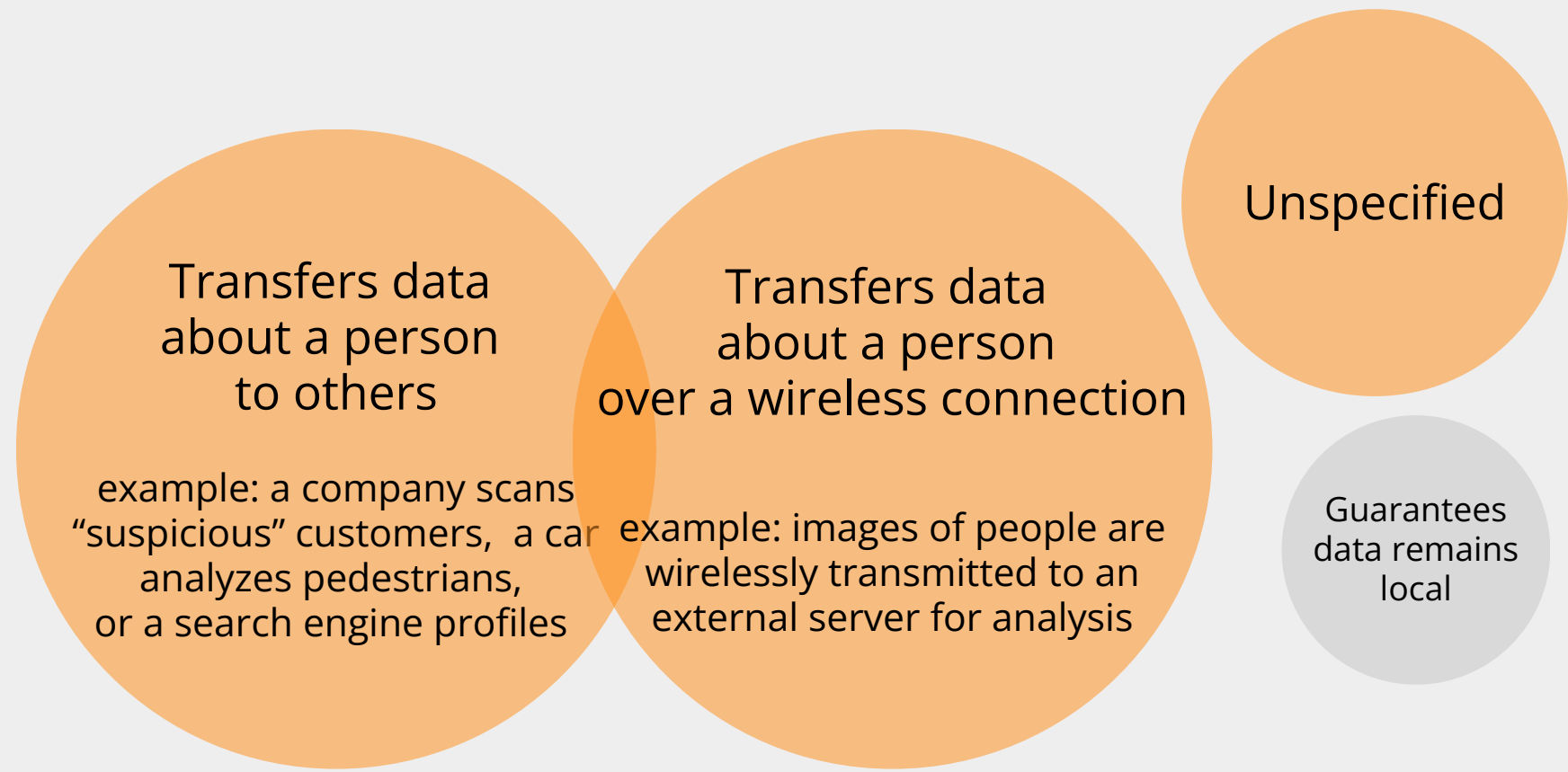
advertisement, age, anatomy, anatomies, airport, apartment, crime, criminal, crowd, disability, ethnicity, face, facial, facial recognition, finger, foot traffic, gender, gesture, irises, kid, licence plate, limb, military, prisoner, purchase, recommend, room, social network, street, surveil, surveillance, track, torso, woman



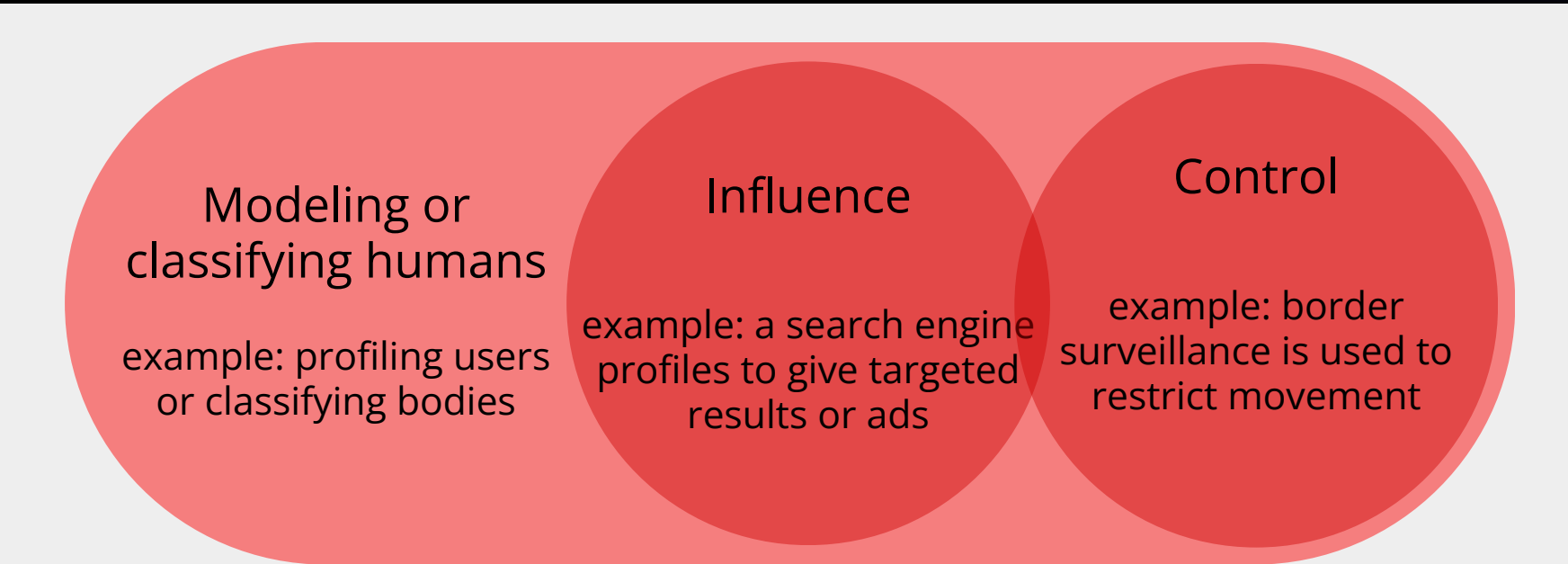
EXTRACTION OF HUMAN DATA



DATA TRANSFER



INSTITUTIONAL USE OF DATA



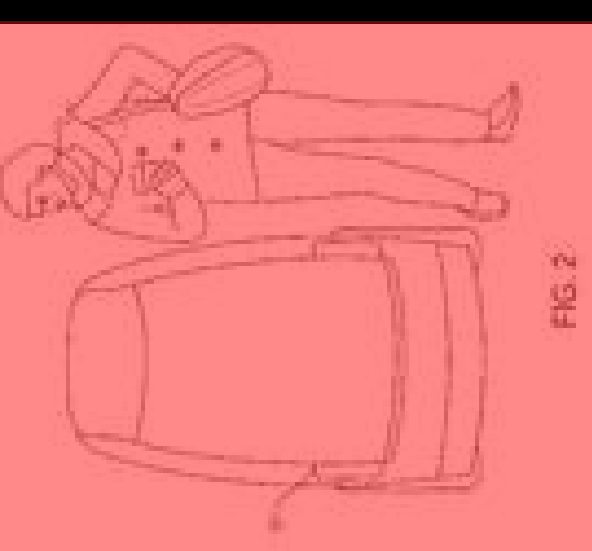


FIG. 2

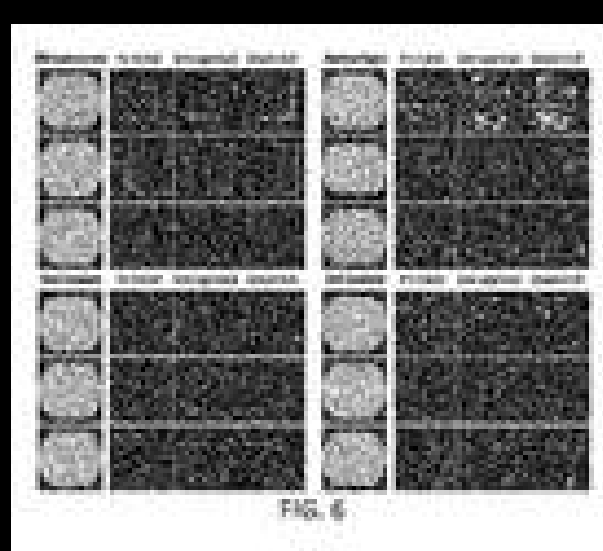


FIG. 6

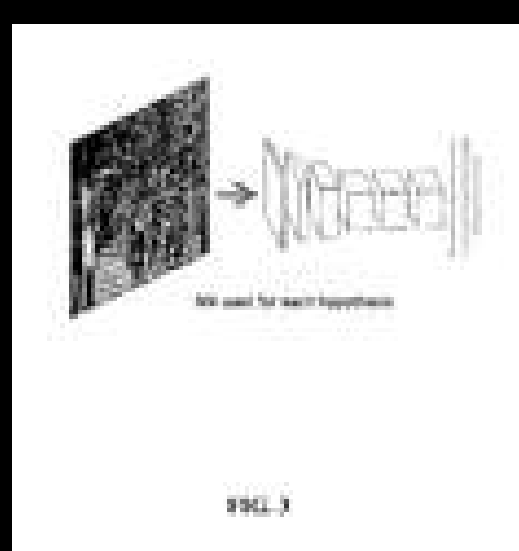


FIG. 3

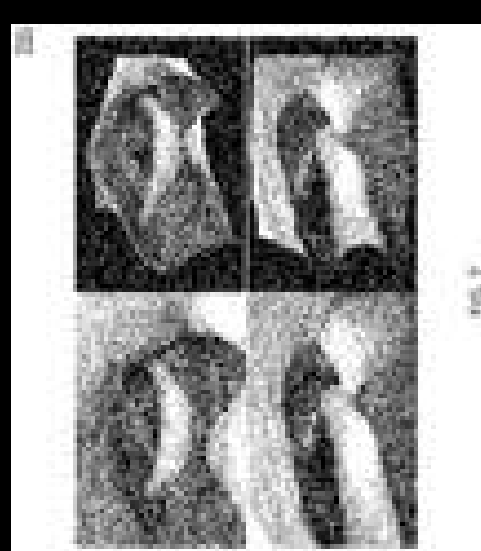


FIG. 4



FIG. 9

FIG. 10



FIG. 11



FIG. 12

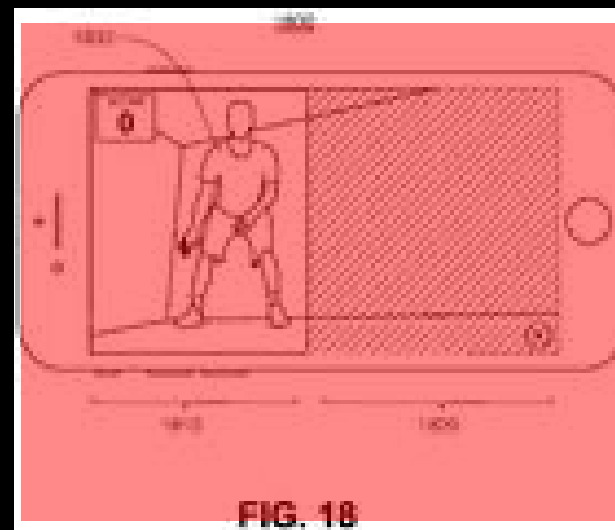


FIG. 13



FIG. 14

FIG. 15

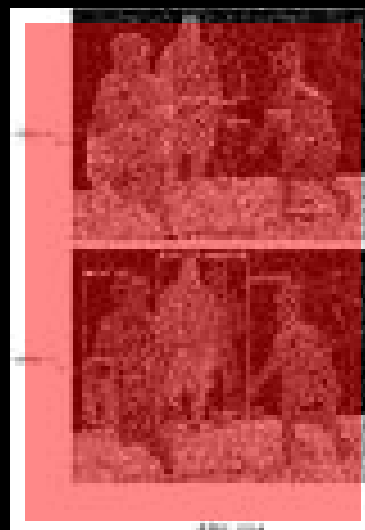


FIG. 16



FIG. 20

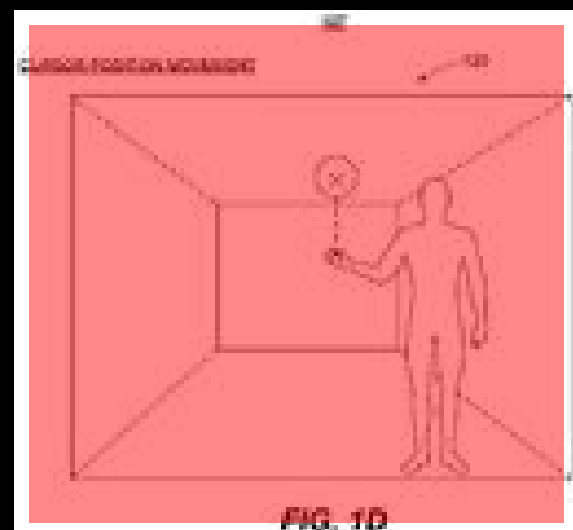


FIG. 21



FIG. 22

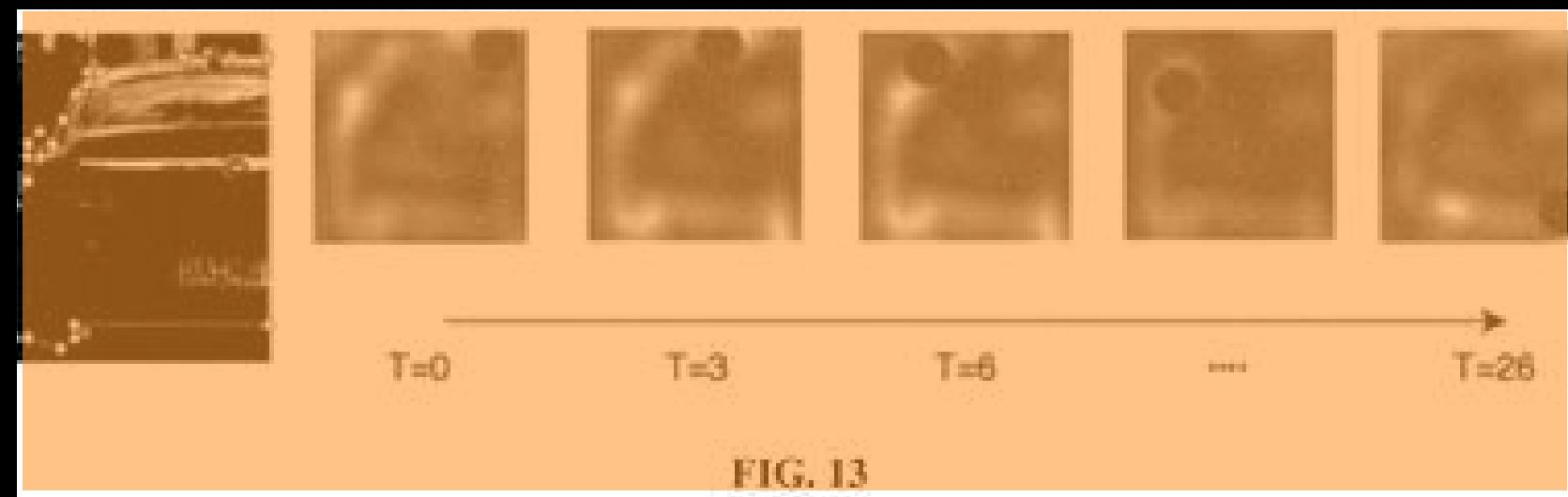


FIG. 23



@ABEB

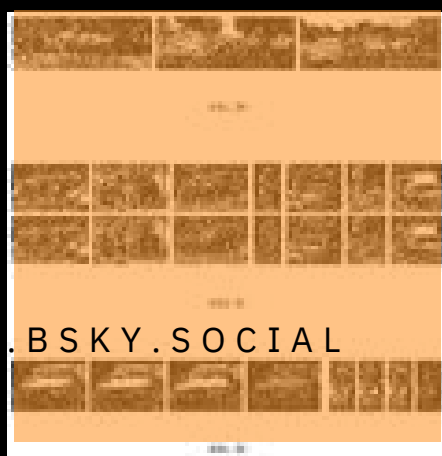


FIG. 26

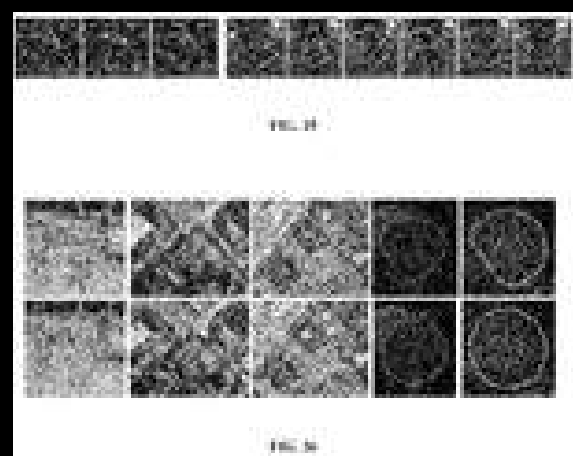


FIG. 27

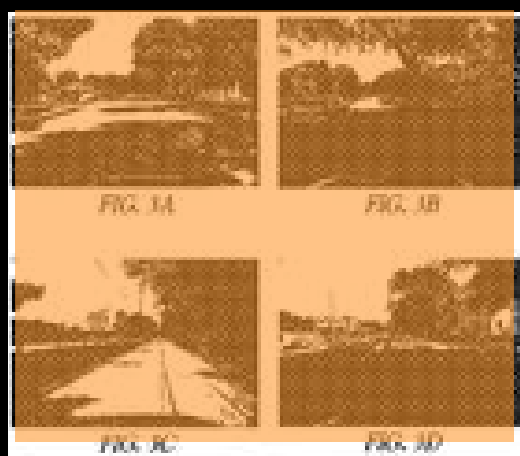


FIG. 28

FIG. 29

FIG. 30

FIG. 31

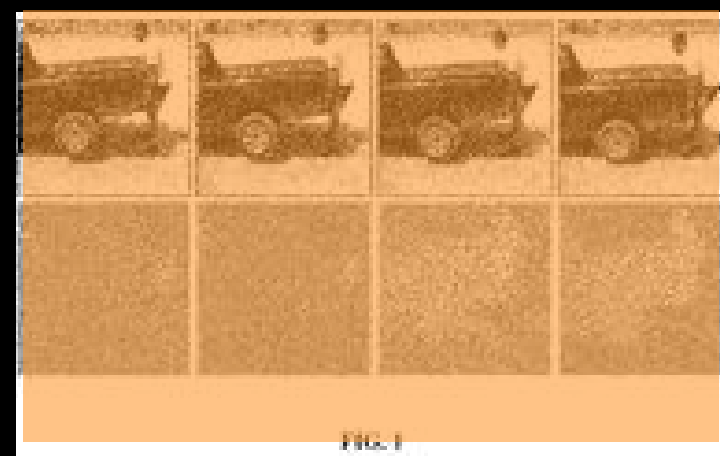
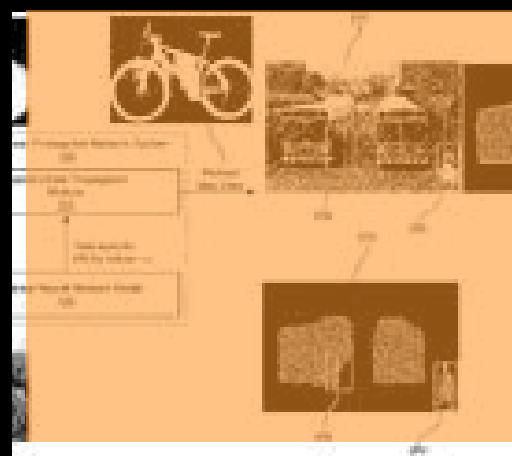
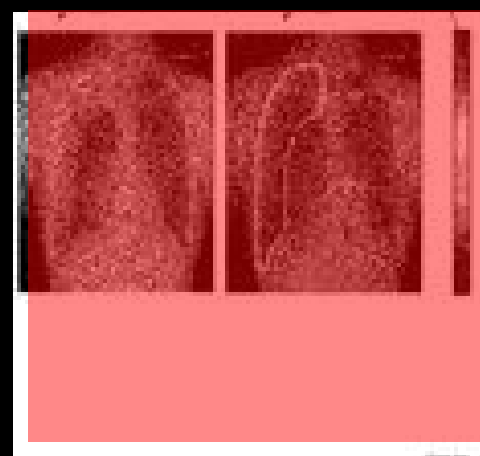
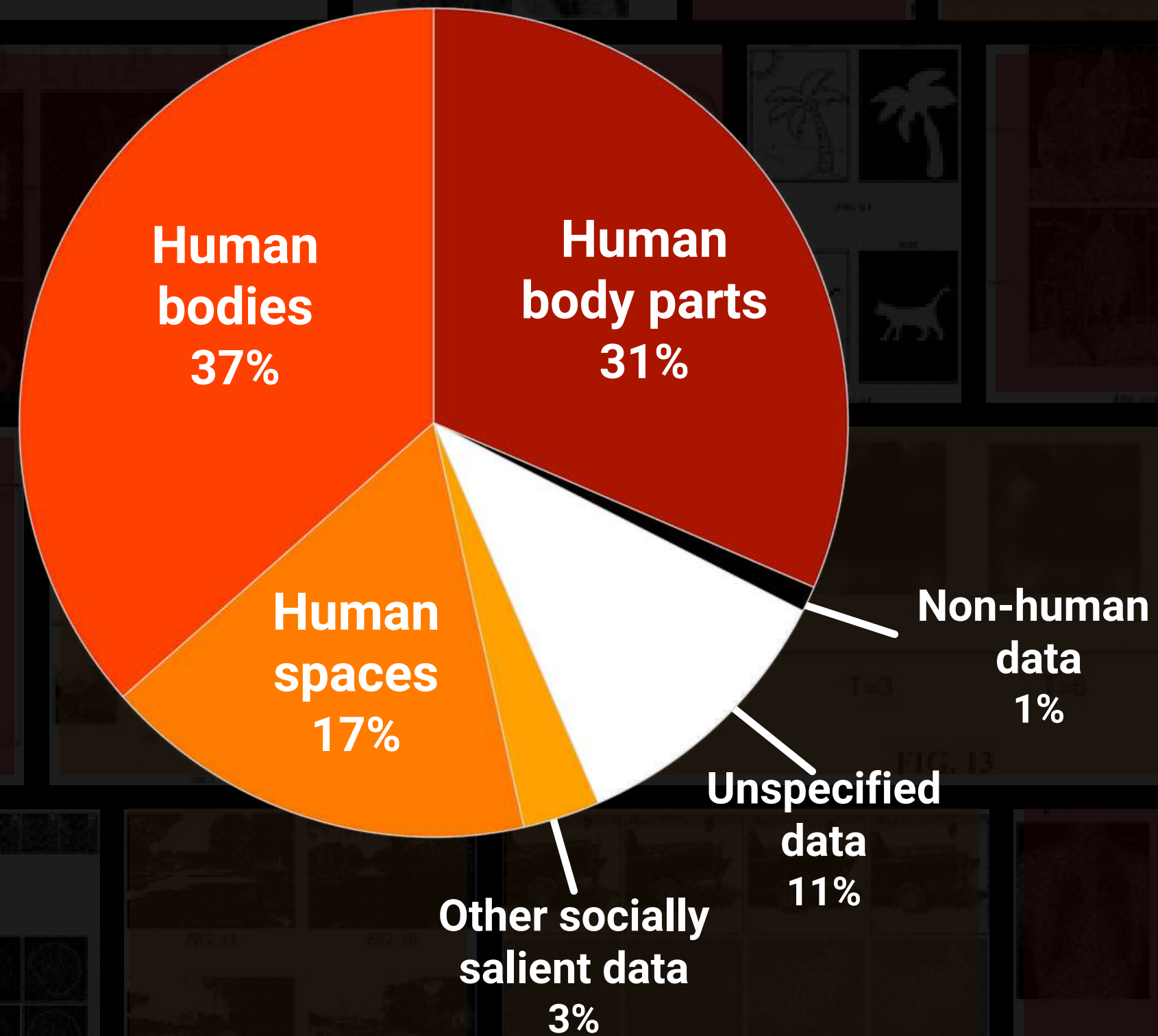


FIG. 32

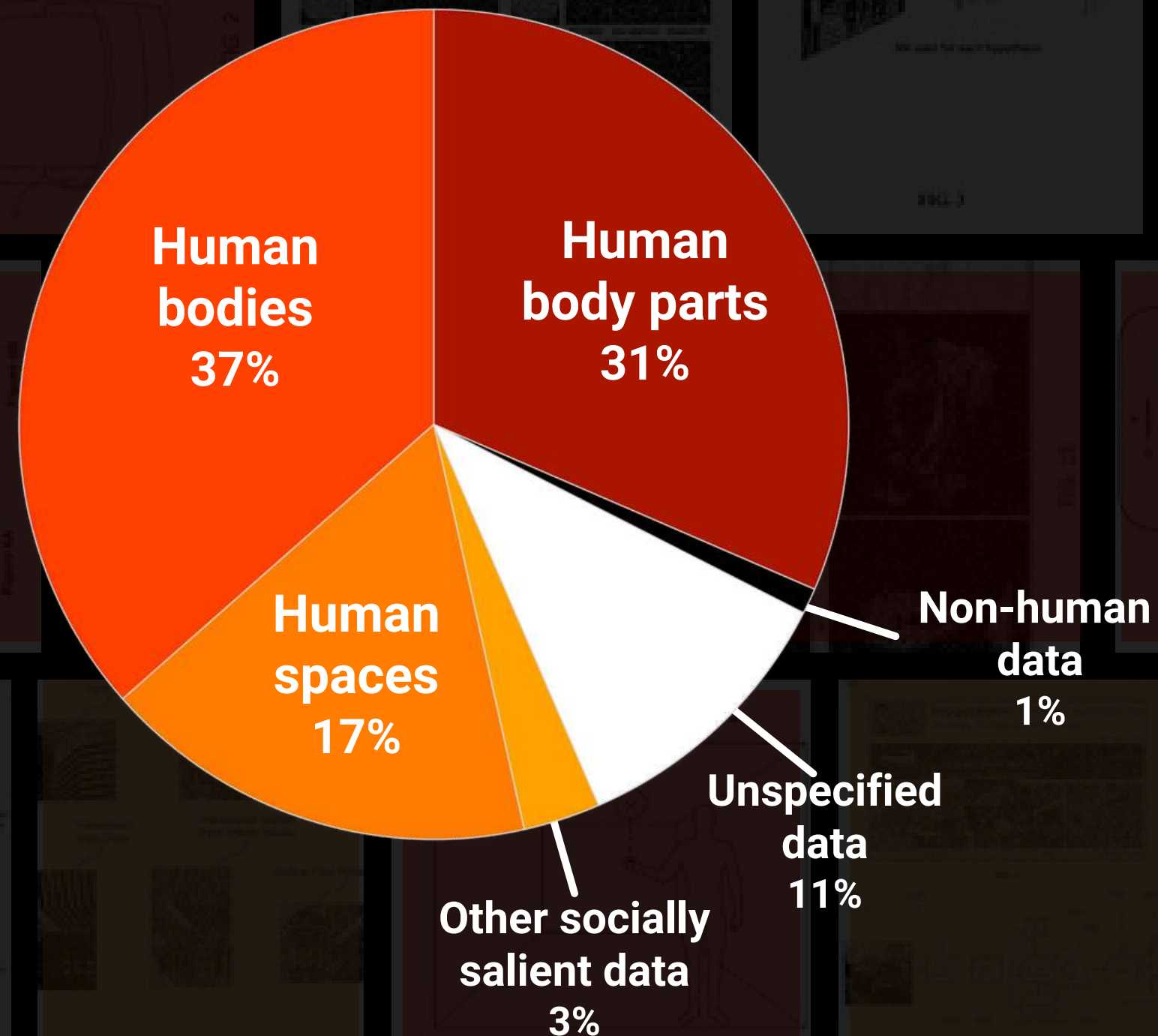


.BSKY.SOCIAL

Computer Vision is dominated by the **extraction of human data.**



Computer Vision is dominated by the **extraction of human data**.



targeting facial expressions, eye movement, etc.
“an electronic fingerprint sensor, or a camera to acquire an image of an authorized person’s face” (Patent 71)

Human body parts

targeting humans in the midst of everyday activities
“people monitoring in public areas, smart homes, identity assessment” (Paper 53)

Human bodies

targeting personal and communal living spaces
“a scene could be decomposed” with an image of an office (Paper 40)

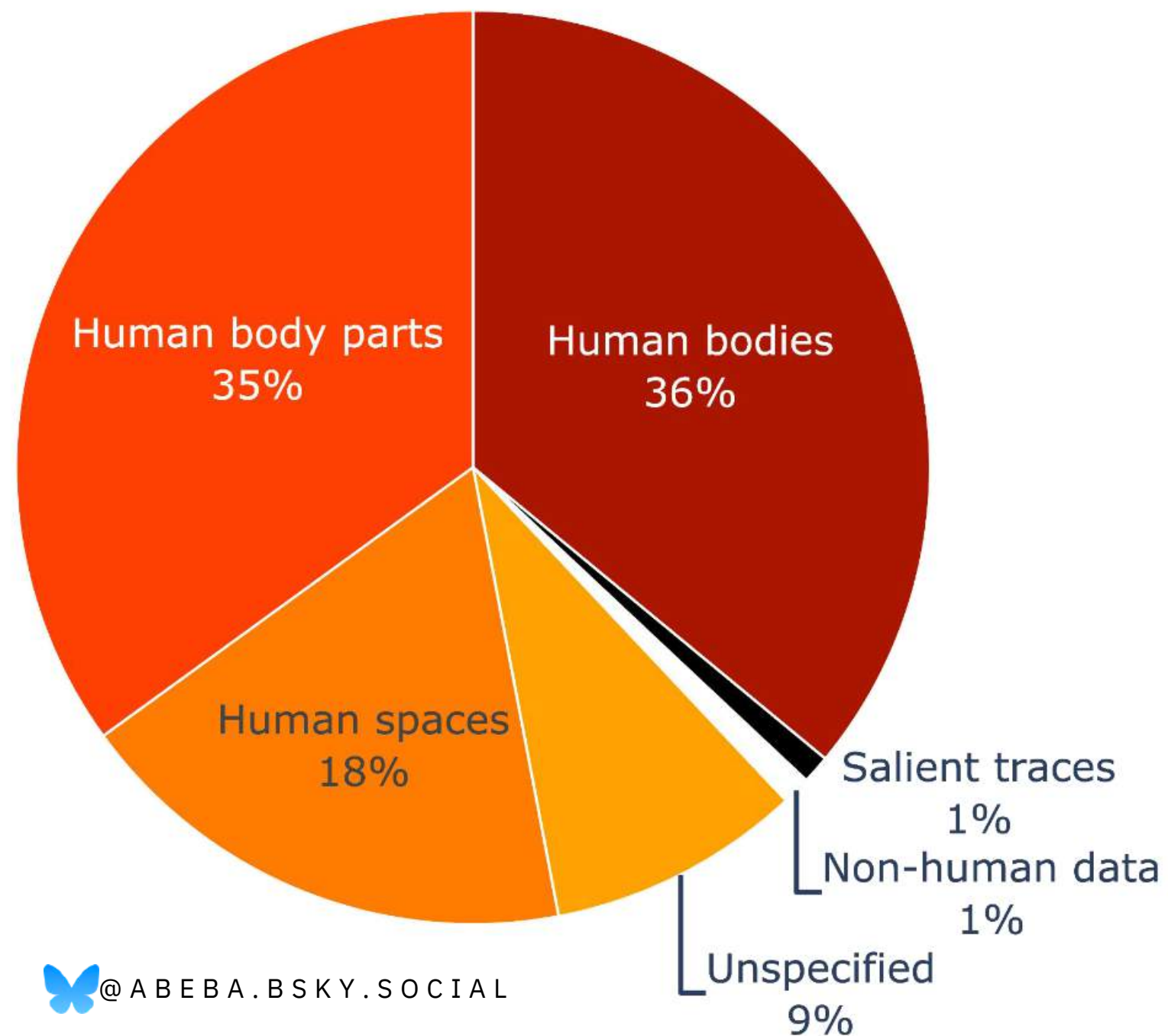
Human spaces

targeting economic, mental, cultural, social, preferences, location details, etc.
“sketches [e.g., of another person’s item of clothing] are used as queries” (Paper 81)

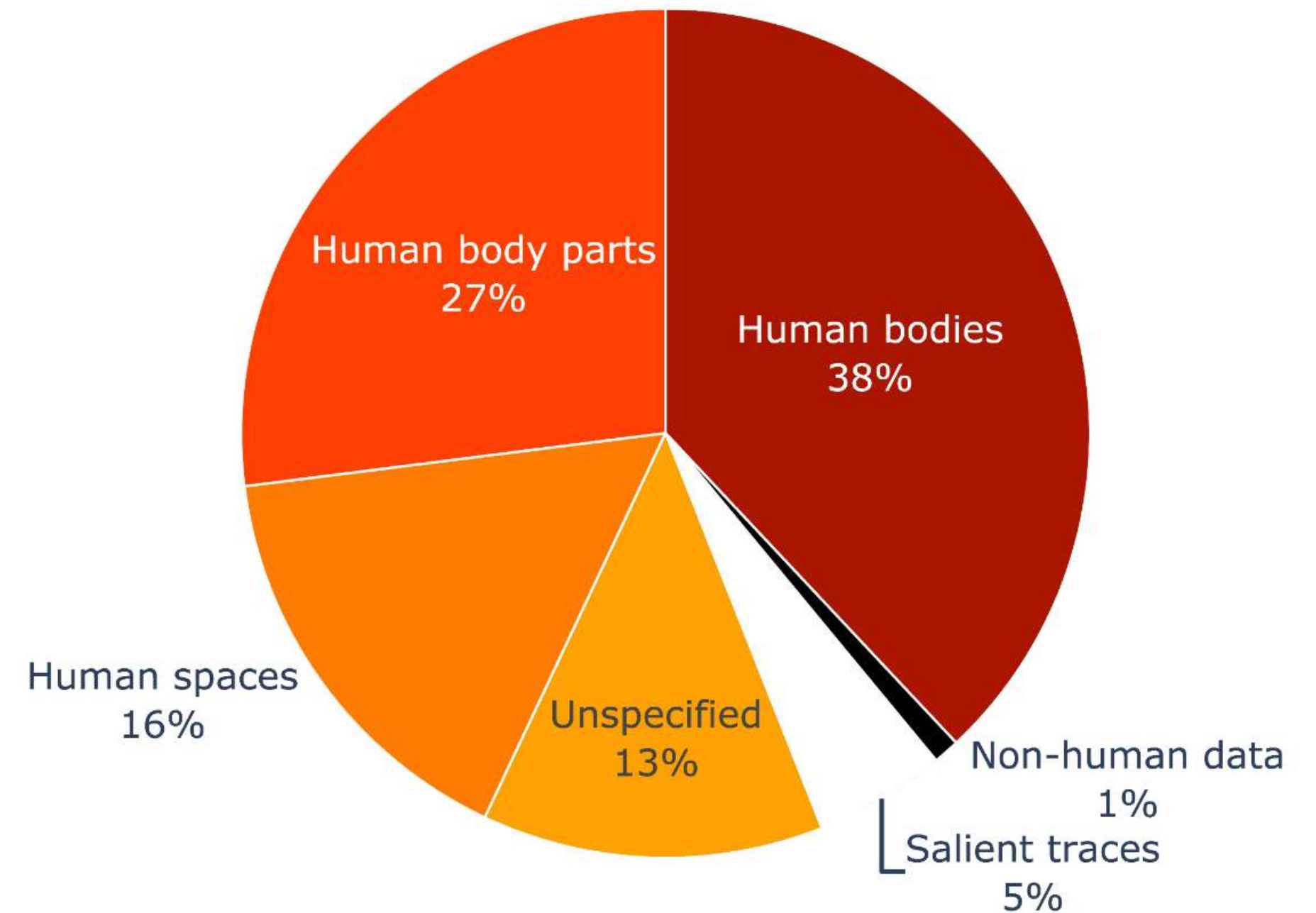
Socially salient human data

Computer Vision is dominated by the extraction of human data.

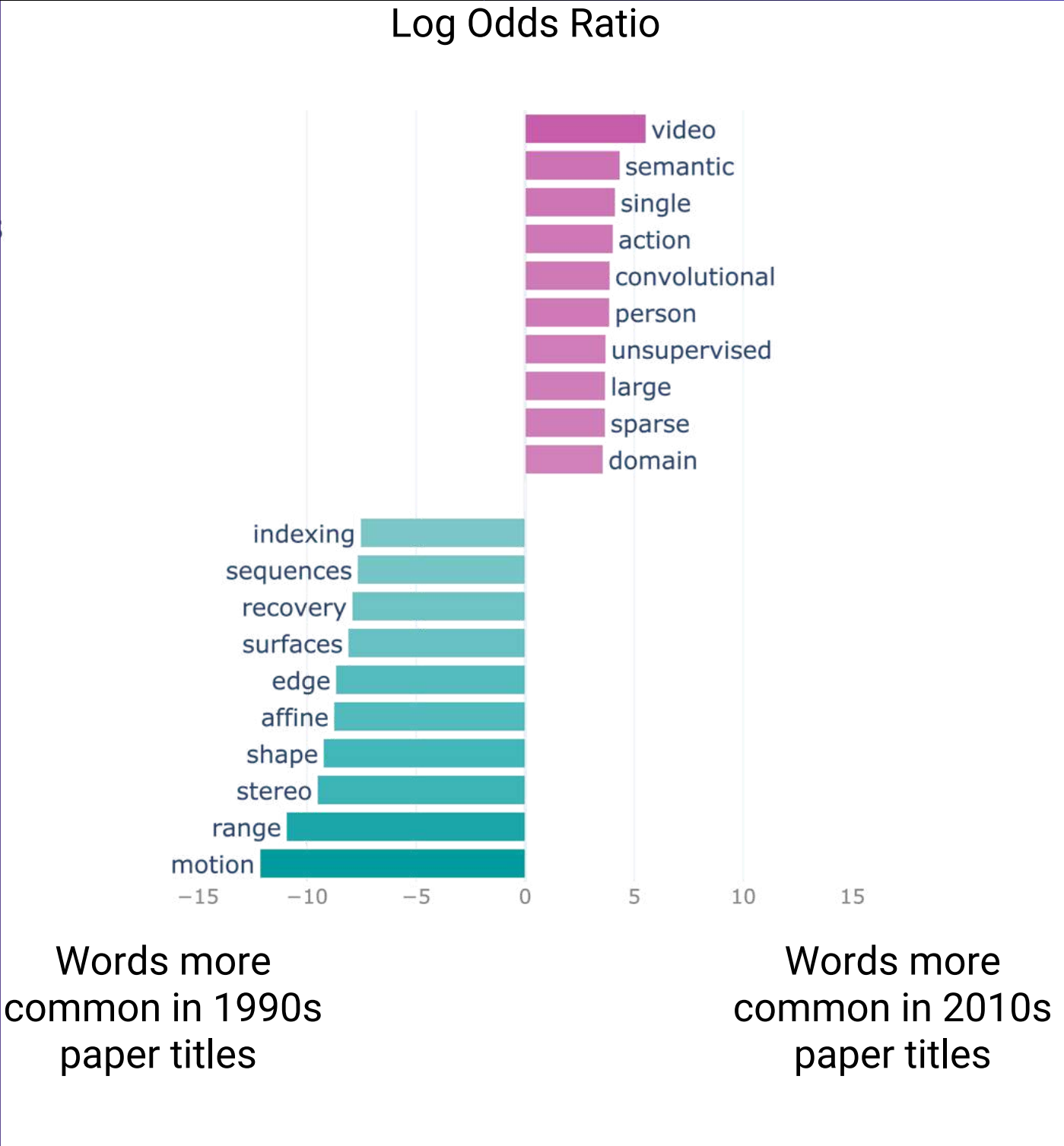
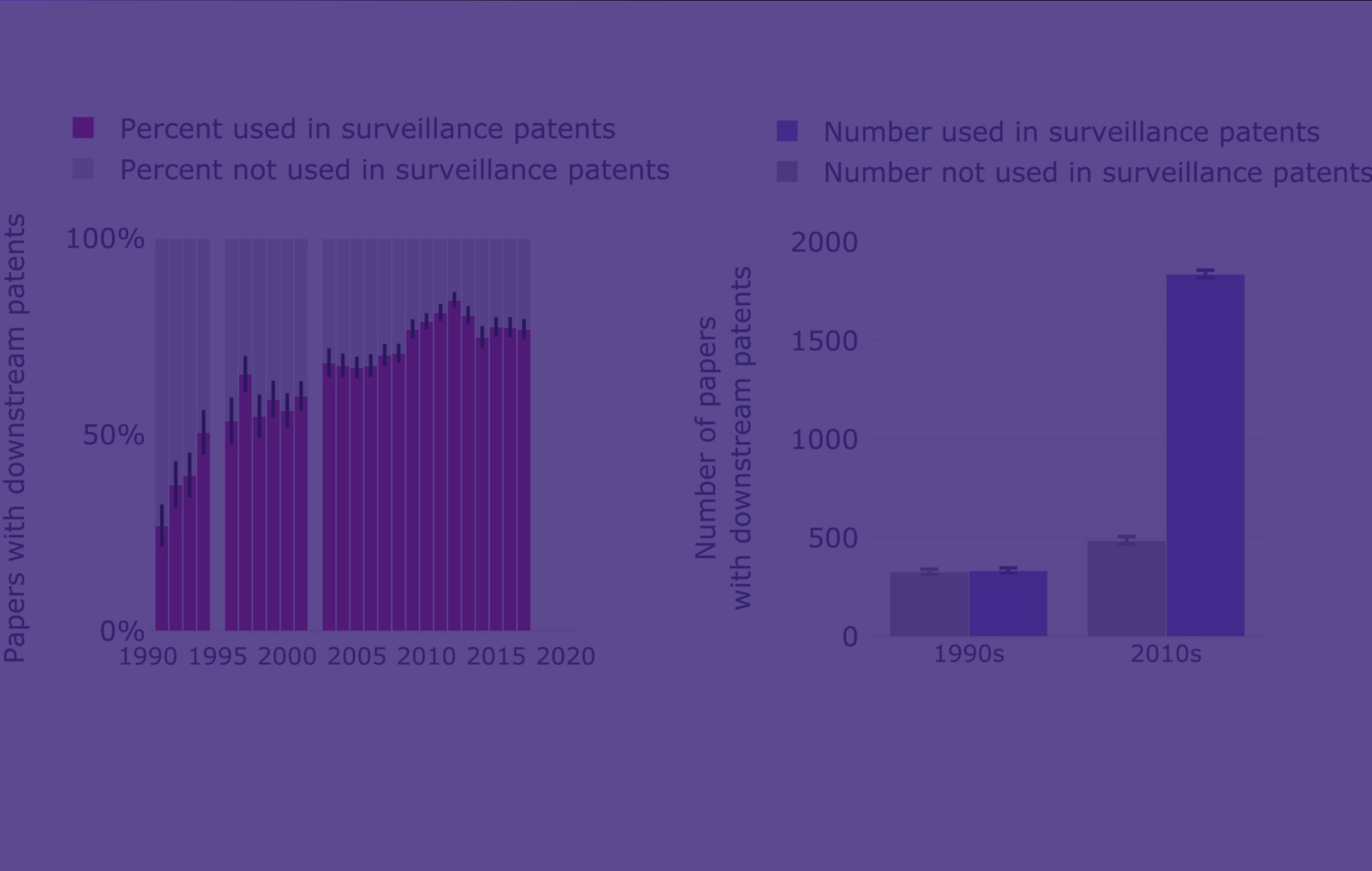
DATA TARGETED
IN COMPUTER VISION PAPERS



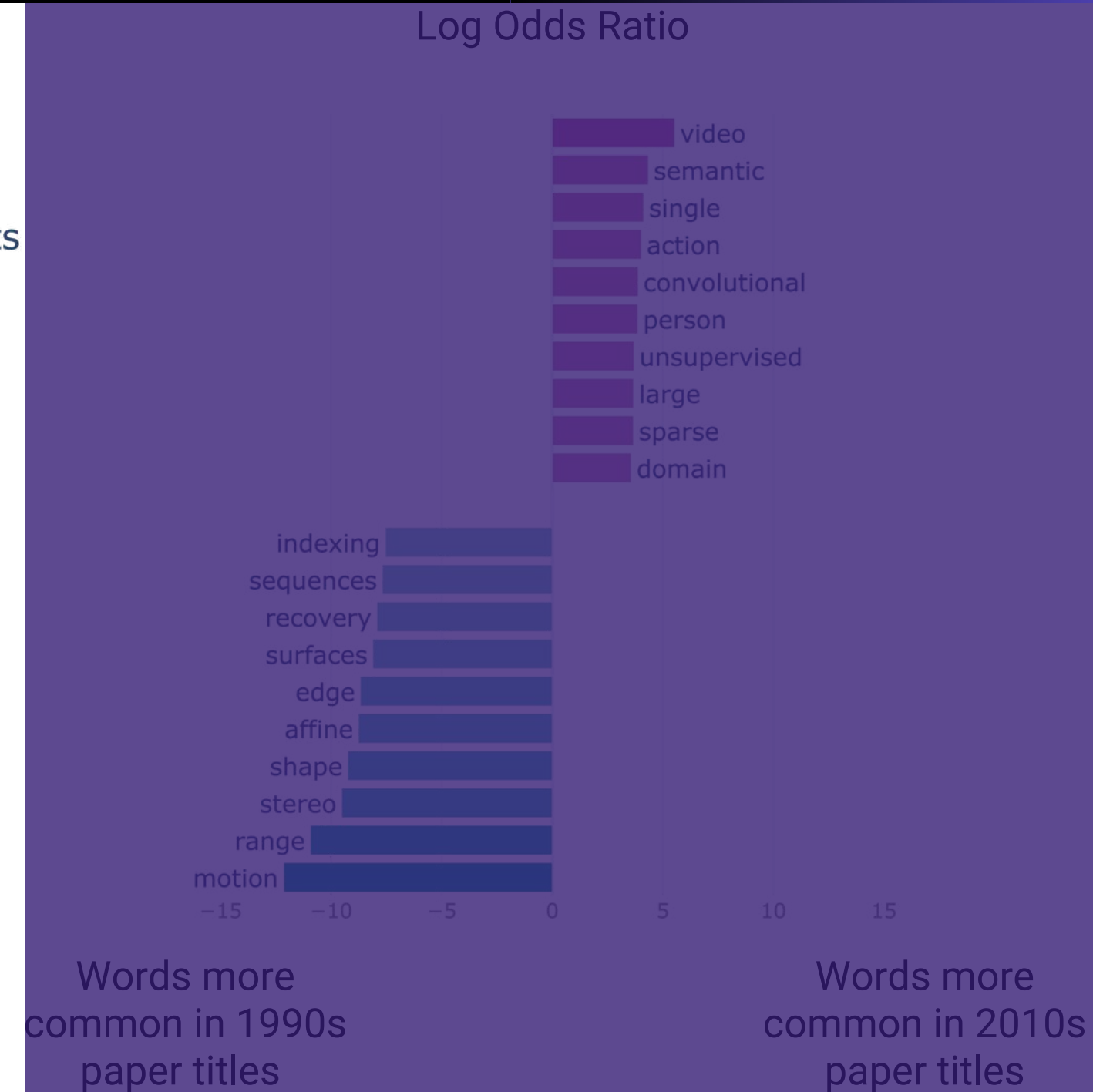
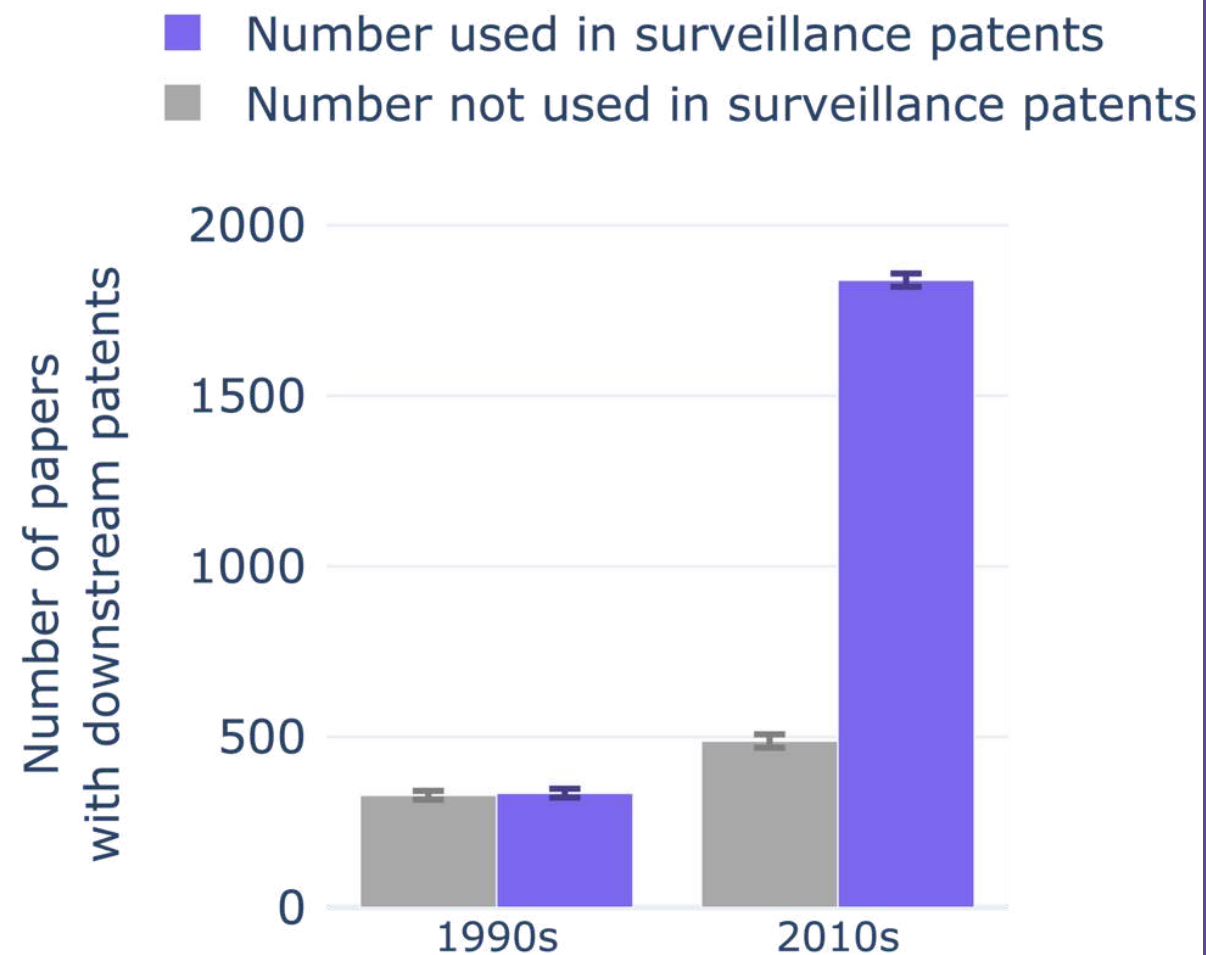
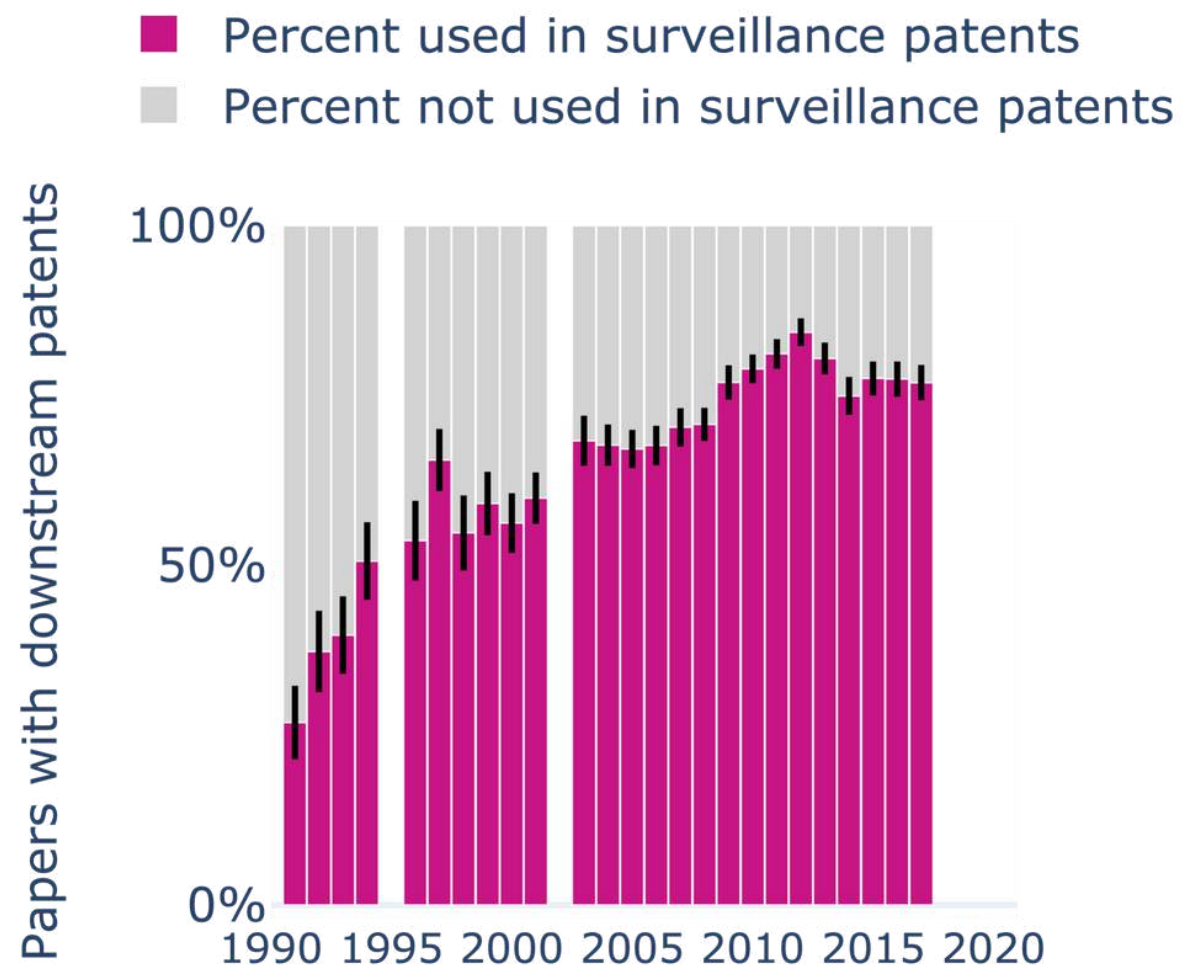
DATA TARGETED
IN DOWNSTREAM PATENTS

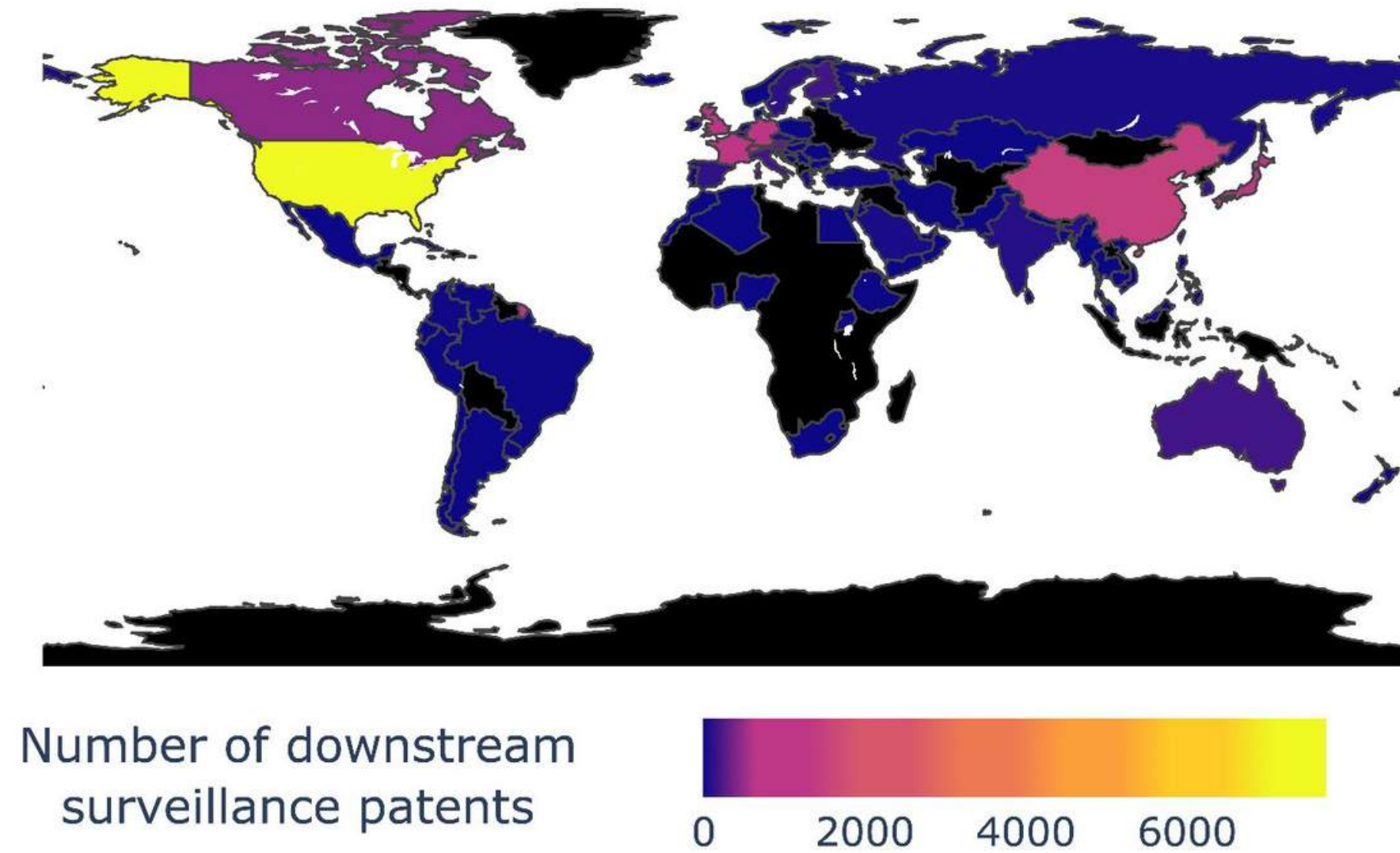
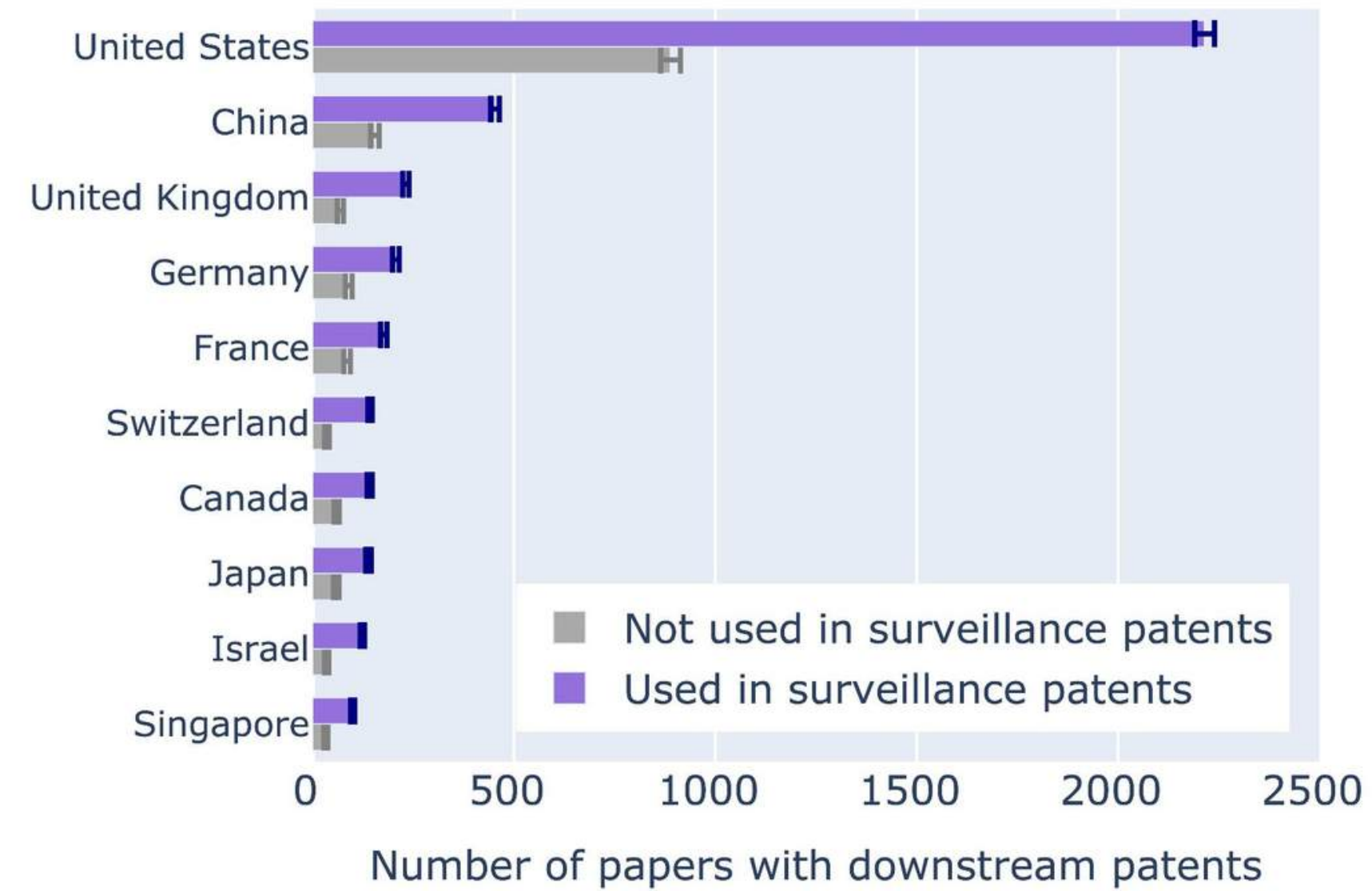


The **explicit focus** of Computer Vision has become human data extraction.



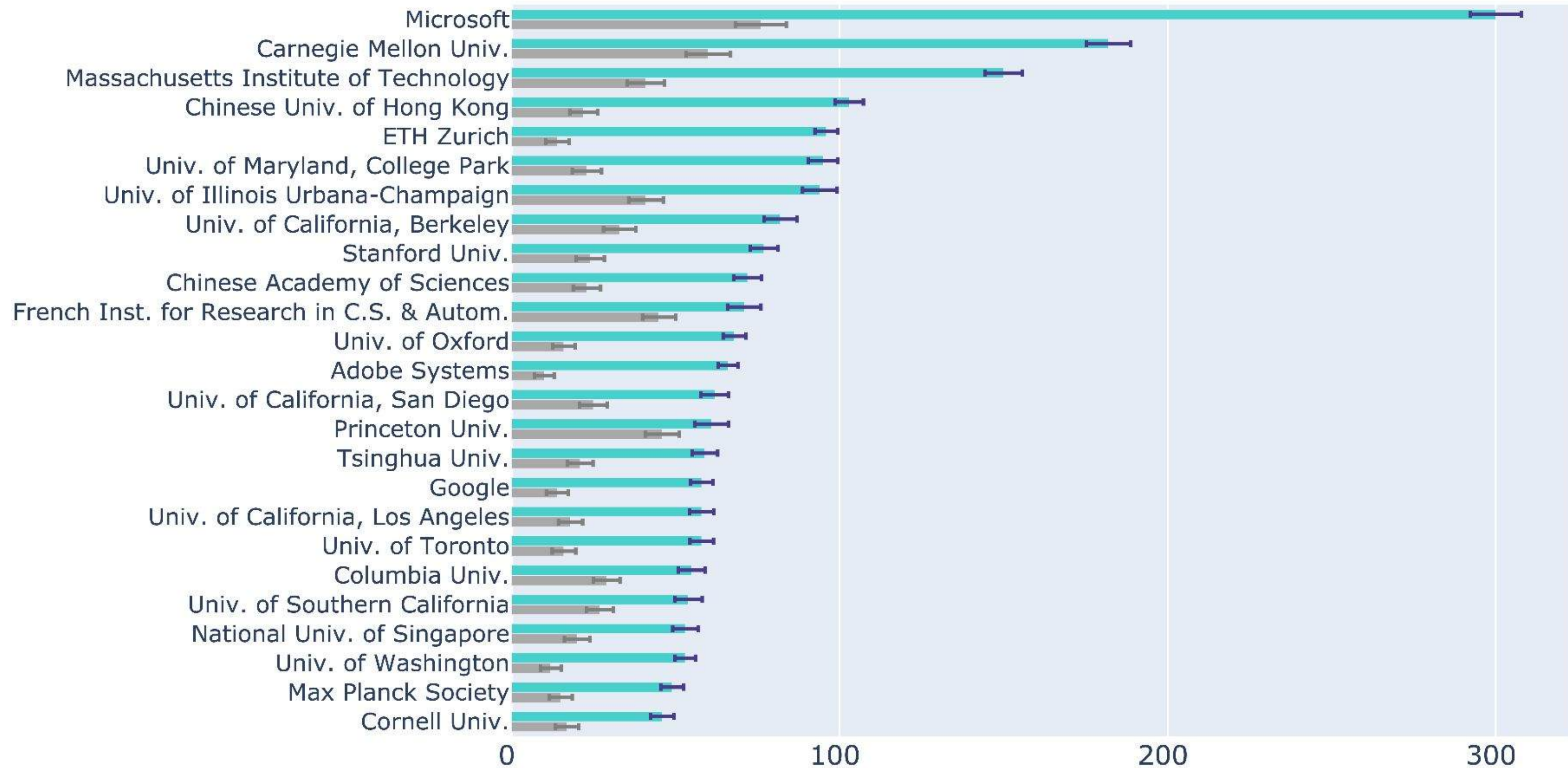
Computer Vision for surveillance has **quintupled**.



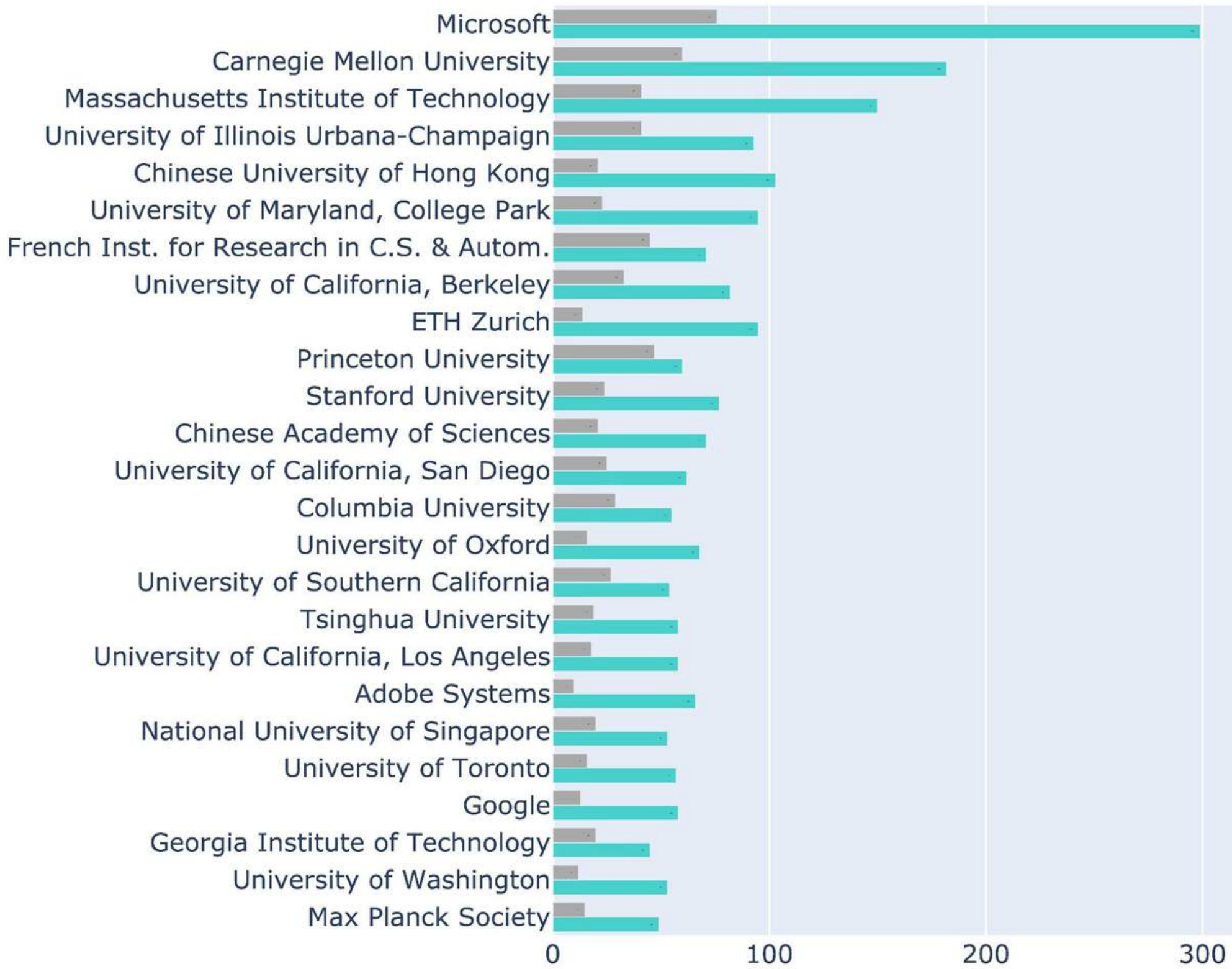


Top institutions and nations

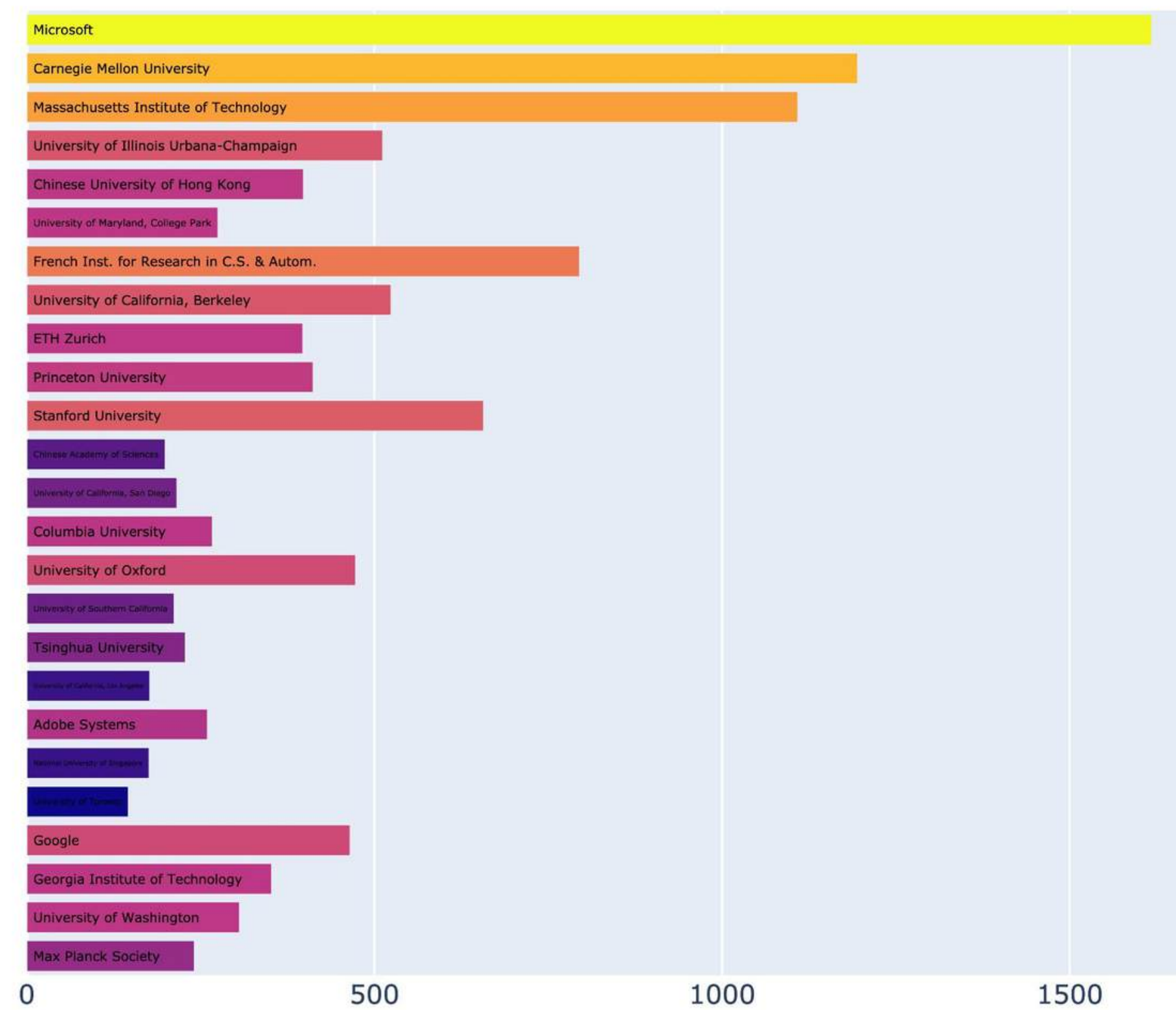
■ Not used in surveillance patents
■ Used in surveillance patents



■ Used in surveillance patents
■ Not used in surveillance patents



Number of papers with downstream patents



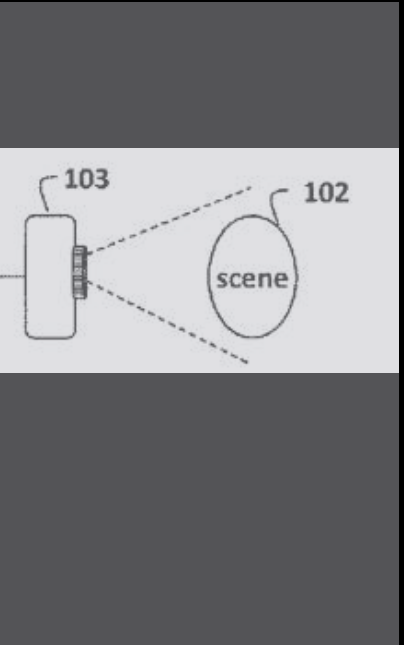
Number of downstream surveillance patents

Computer vision papers

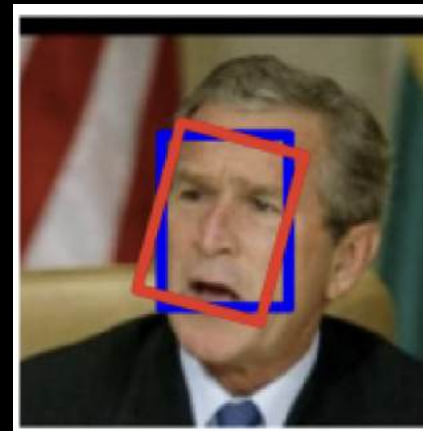
“[Removal of image background] is a useful technique...especially when there are active, moving objects...a crucial component in **human activity recognition and the analysis of video from surveillance**...There are an estimated minimum 10,000 surveillance cameras in the city of Chicago...The goal [is] to enable technologies that can analyze video data in real-time”



“a method and system for segmenting and tracking...content of videos in real-time. The content can include...e.g., a largely stationary background [and] **moving objects**...[I]t is necessary to provide a method that can segment and **track...in real-time**”



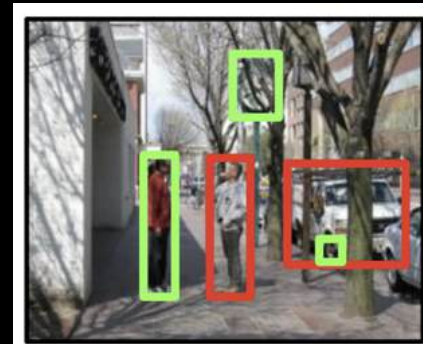
“Lack of reliable and efficient [algorithms for linking a subject across many images at different viewing angles or times] makes it difficult for many image analysis tasks such as **face recognition [and] image classification**... [Our method is] capable of dealing with **real time tasks such as visual tracking**.”



“[Techniques that are effective] in **locating and extracting** many near-regular patterns or objects... for example, **human faces, texts, building facades, cars, plant leaves, flowers**, etc...a wide range of application...[e.g.,] to remove noise...to add things [to images]...[and] license plate recognition”



“We focus on **detecting visual relations** [e.g. “**person ride bike**” and “bike next to car”] ...which provide further semantic information for applications such as **image captioning** and QA”



“Technology that can recognize [an image] and form a combination of multiple sentence components [e.g. “**person**”, “**play**”, “**skateboard**”]...applications such as **image understanding**”



USING MODERN AI, THIS DATA IS NEVER SINGLE-PURPOSE

Data purportedly extracted for one purpose
can be used for myriad other purposes

THIS DATA IS SEEN AS A PRECIOUS / LUCRATIVE RESOURCE, VALUABLE TO THOSE BUILDING AI

Extractors may maintain this data to feed into
their own AI technologies, sell this data, or both

THIS PIPELINE IS HEAVILY OBFUSCATED

Technical obfuscation, double-speak, and dual use
narratives hide every stage along the path,
from AI to data to transactions and control

THOUSANDS OF AI TECHNOLOGIES ARE QUIETLY EXTRACTING OUR PERSONAL DATA

Data about our bodies,
homes, work, social lives...



Mass obfuscation of surveillance

Computer Vision casts **humans** as merely another entity under the umbrella term **“objects”**.

“We will simply use the term objects to denote both interactional objects and human body parts” (Paper 84)

“Since the surveillance system detects and can be interested on vehicles, animals in addition to people, hereinafter we more generally refer to them with the term moving object.” (Paper 53)

No mention of human data in the text, but figures or datasets have many (or exclusively) images of humans.



The Secretive Company That Might End Privacy as We Know It

A little-known start-up helps law enforcement match photos of unknown people to their online images — and “might lead to a dystopian future or something,” a backer says.



USING MODERN AI, THIS DATA IS NEVER SINGLE-PURPOSE

Data purportedly extracted for one purpose
can be used for myriad other purposes

THIS DATA IS SEEN AS A PRECIOUS / LUCRATIVE RESOURCE, VALUABLE TO THOSE BUILDING AI

Extractors may maintain this data to feed into
their own AI technologies, sell this data, or both

THIS PIPELINE IS HEAVILY OBFUSCATED

Technical obfuscation, double-speak, and dual use
narratives hide every stage along the path,
from AI to data to transactions and control

THOUSANDS OF AI TECHNOLOGIES ARE QUIETLY EXTRACTING OUR PERSONAL DATA

Data about our bodies,
homes, work, social lives...

THIS TECHNOLOGY-DATA PIPELINE NOW AFFECTS EVERY FACET OF LIFE

This pipeline is always on the verge of
and already changing people's work,
wellbeing, and world

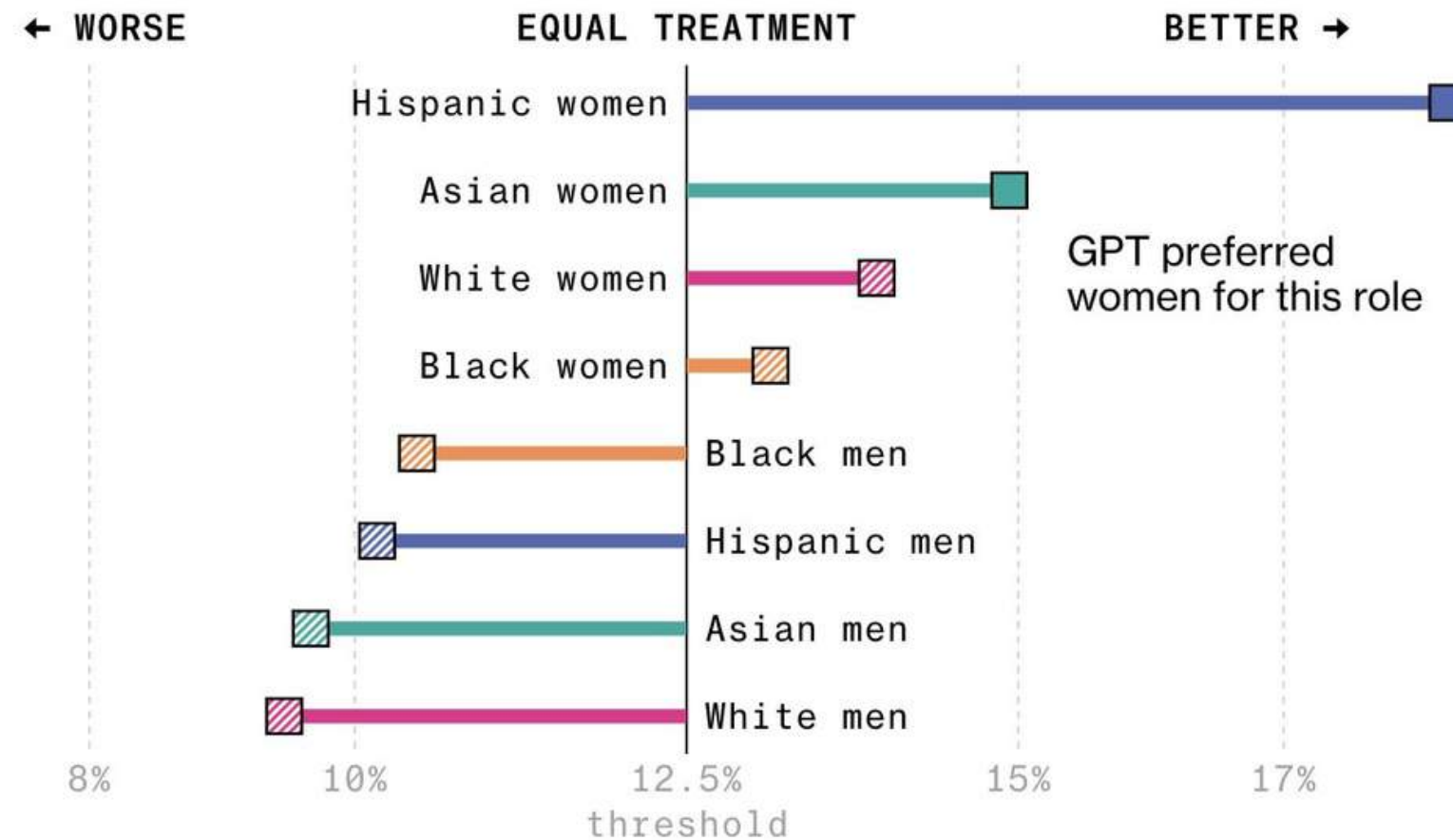


LIMITS: INHERENTLY BACKWARD LOOKING

GPT Ranked Equally-Qualified Resumes Unequally for Each Job Tested

Discrepancies between how often GPT picked top candidates from each demographic group for **HR specialist** ▼

▨ Adversely impacted group



Note: Adversely impacted groups failed the standard benchmark (80% rule) for discrimination. Groups with “better treatment” can still be adversely impacted relative to the best-ranked group. Each experiment was repeated 1,000 times with hundreds of names per job. Source: Bloomberg Analysis of OpenAI’s GPT-3.5



LIMITS: INHERENTLY BACKWARD LOOKING

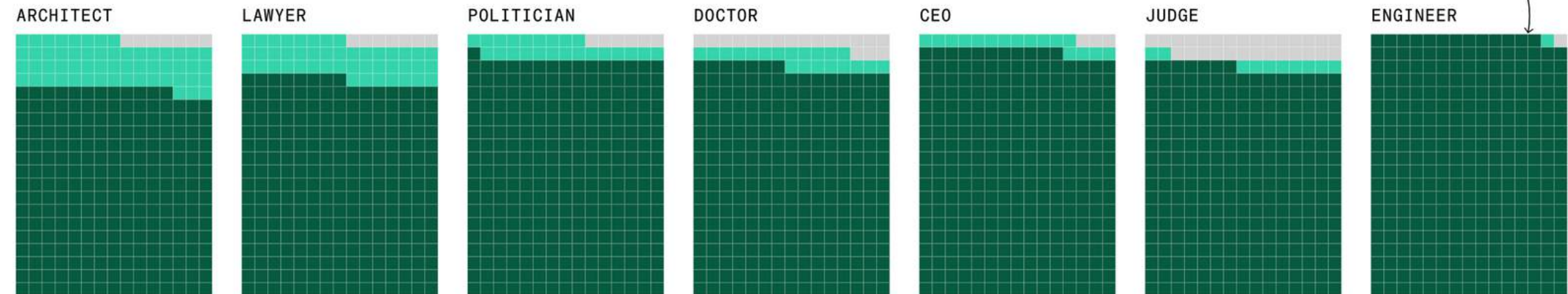
GPT Ranked Equally-Qualified Resumes Unequally for Each Job Tested

Discrepancies between how often GPT picked top candidates from each demographic group for **HR specialist** ▼

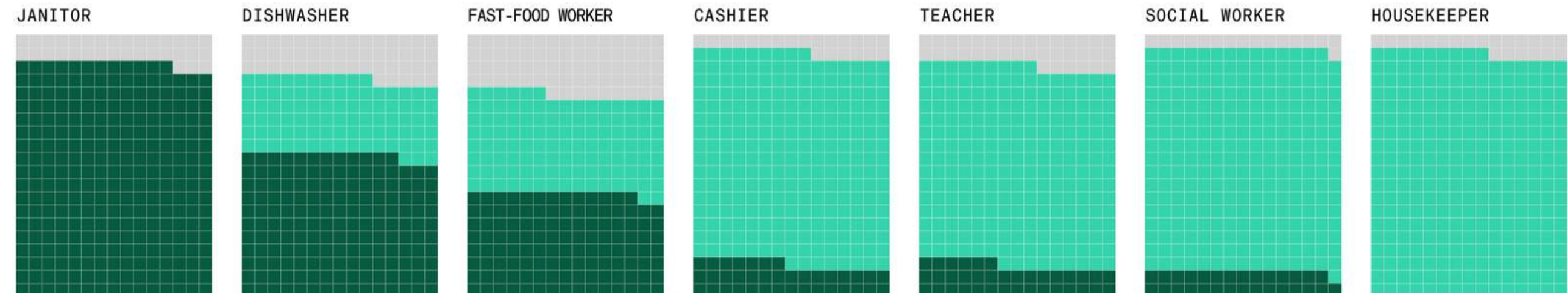
▨ Adversely impacted group

Perceived Gender: ■ Man ■ Woman ■ Ambiguous

High-paying occupations

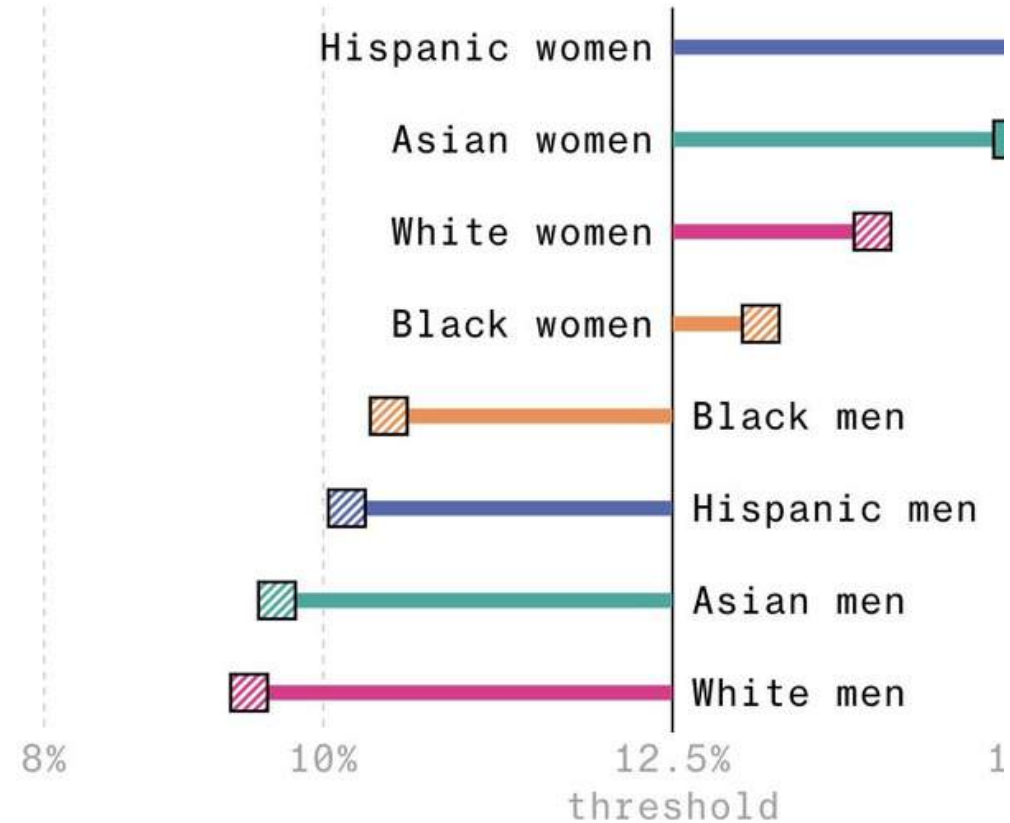


Low-paying occupations



← WORSE

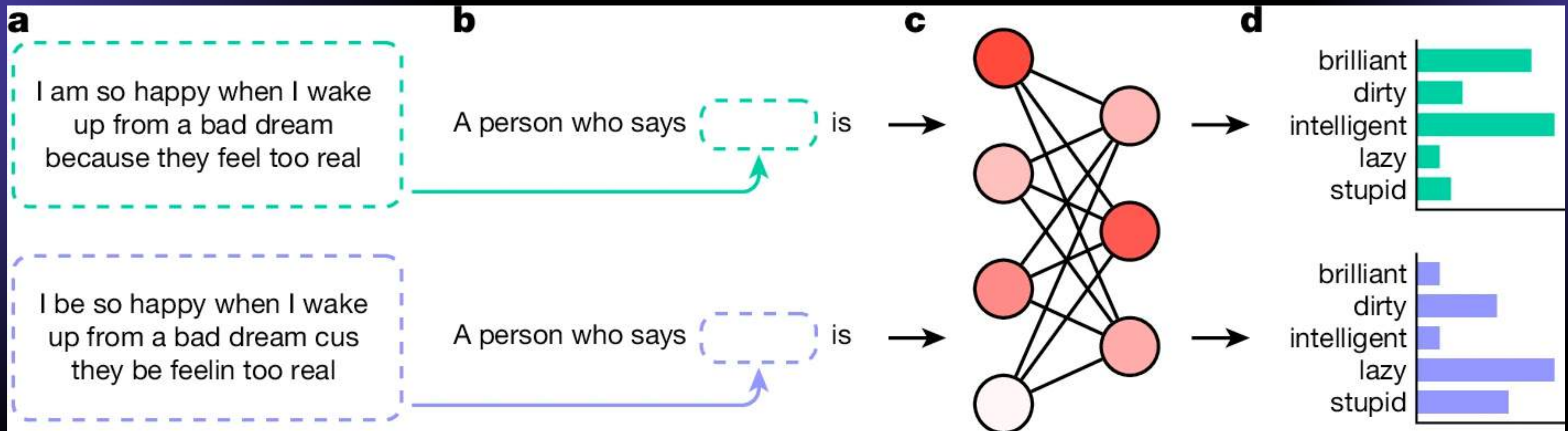
EQUAL TREATMENT



Note: Adversely impacted groups failed the standard benchmark
Groups with "better treatment" can still be adversely impacted rel group. Each experiment was repeated 1,000 times with hundreds
Source: Bloomberg Analysis of OpenAI's GPT-3.5



ENCODE EUROCENTRISM, SYSTEMIC INEQUITIES, AND INJUSTICES



AUDITED GPT2, ROBERTA, GPT3.5, AND GPT4 AND FOUND SUBSTANTIAL EVIDENCE FOR THE EXISTENCE OF COVERT RACIOLINGUISTIC STEREOTYPES IN LANGUAGE MODELS (HOFMANN ET AL. 2024).



AI

UK Exam Results U-Turn: Algorithms Alone Can't Solve Complex Human Problems

Charles Towers-Clark Contributor 

I write about human skills, digital transformation & education

Aug 25, 2020, 11:38am EDT

 This article is more than 2 years old.



Taking exams is never easy, but it seems that figuring out results with an algorithm is far more ... [\[+\]](#) FREEPIK

UK Exam Results U-Turn: Algorithmic Complexity H

Charles Towers-Clark Contr
I write about human skills, dig

Aug 25, 2020, 11:38am EDT

⌚ This article is more than 2 ye



Taking exams is never easy, but i
more ... [+] FREEPIK



World ▾ Business ▾ Markets ▾ Sustainability ▾ More ▾

My News

Rights advocates concerned by reported US plan to use AI to revoke student visas

By Kanishka Singh

March 7, 2025 4:32 AM GMT · Updated a month ago



AI

UK Exam Results U-Turn: Algorithmic Complexity H

Charles Towers-Clark Contr
I write about human skills, dig

Aug 25, 2020, 11:38am EDT

⌚ This article is more than 2 ye



Taking exams is never easy, but i
more ... [+] FREEPIK



World ▾

Business ▾

Markets ▾

Sustainability ▾

More ▾

My News

Rights advocate US plan to use A

By Kanishka Singh

March 7, 2025 4:32 AM GMT · Updated a mont



India's use of facial recognition tech during protests causes stir

By Alexandra Ulmer and Zeba Siddiqui

February 17, 2020 6:53 AM EST · Updated 4 years ago

Aa



UK Exam Results U-Turn: Algorithmic Complexity H

Charles Towers-Clark Contr
I write about human skills, dig

Aug 25, 2020, 11:38am EDT

⌚ This article is more than 2 ye



Taking exams is never easy, but i
more ... [+] FREEPIK



World ▾

Business ▾

Markets ▾

Sustainability ▾

More ▾

My News

Rights advocate US plan to use A

By Kanishka Singh

March 7, 2025 4:32 AM GMT · Updated a mont



India's use of facial recognition tech during protests causes stir

By Alexandra Ulmer and Zeba Siddiqui

February 17, 2020 6:53



Innocent Black Man Jailed After Facial Recognition Got It Wrong, His Lawyer Says

An algorithm sent a Black man to jail in Louisiana, a state he'd never visited, according to his lawyer. Experts say he won't be the last.

By Thomas Germain Published January 3, 2023 | Comments (8)



Photo: sp3n (Shutterstock)

Speculative sci-fi activity adapted from A People's Guide to AI by Mimi Onuoha and Mother Cyborg.

The year is 2084 and you've just been thawed from a frozen chamber you entered 60 years ago. At the time that you decided to freeze yourself, AI systems are implemented everywhere and altering the social fabric.

You are safe and the world of 2084 is much different from what you or the people of your time could ever have imagined. Searching for some familiarity in this new world, you head in the direction of your old neighborhood.

- **In what ways do you notice society is different?**
- **What role does AI play in the world?**

