# Who's using your face? The ugly truth about facial recognition

Researchers are scraping our images from social media and CCTV. We may not like the consequences

Madhumita Murgia SEPTEMBER 18 2019

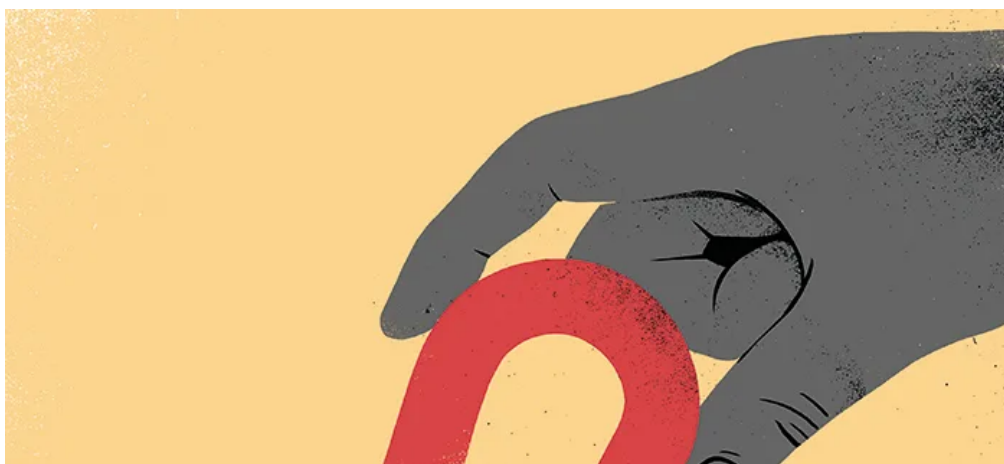This article was originally published on 19 April 2019.

When Jillian York, a 36-year-old American activist, was on vacation in February, she received an unexpected text. Her friend Adam Harvey, another activist and researcher, had discovered photos of her in a US government database used to train facial-recognition algorithms, and wondered whether she knew about it.

York, who works in Berlin for the Electronic Frontier Foundation, a digital rights non-profit group, did not. She was stunned to discover that the database contained nearly a dozen images of her, a mixture of photos and YouTube video stills, taken over a period of almost a decade.

When she dug into what the database was used for, it dawned on her that her face had helped to build systems used by the federal government to recognise faces of interest, including suspected criminals, terrorists and illegal aliens.

"What struck me immediately was the range of times they cover," York says. "The first images were from 2008, all the way through to 2015." Two of the photos, by a photographer friend, had been scraped from Google. "They were taken at closed meetings. They were definitely private in the sense that it was me goofing around with friends, rather than me on stage," she adds.

Another half-dozen photos had been clipped from YouTube videos of York speaking at events, on topics including freedom of expression, digital privacy and security. "It troubles me that someone was watching videos of me and clipping stills for this purpose," she says.

To teach a machine to recognise a human face it has to be trained using hundreds of thousands of faces. The more natural, varied and unposed the faces are, the better they simulate real-life scenarios in which surveillance might take place © Sébastien Thibault

York is one of 3,500 subjects in this database, which is known as Iarpa Janus Benchmark-C (IJB-C). Iarpa is a US government body that funds innovative research aimed at giving the US intelligence community a competitive advantage; Janus — named after the two-faced Roman god — is its facial-recognition initiative.

The dataset, which was compiled by a government subcontractor called Noblis, includes a total of 21,294 images of faces (there are other body parts too), averaging six pictures and three videos per person, and is available on application to researchers in the field. By their own admission, its creators picked "subjects with diverse occupations, avoiding one pitfall of 'celebrity-only' media [which] may be less representative of the global population."

Other subjects in the dataset include three EFF board members, an Al-Jazeera journalist, a technology futurist and writer, and at least three Middle Eastern political activists, including an Egyptian scientist who participated in the Tahrir Square protests in 2011, the FT can confirm.

None of the people described above was aware of their inclusion in the database. Their images were obtained without their explicit consent, as they had been uploaded under the terms of Creative Commons licences, an online copyright agreement that allows images to be copied and reused for academic and commercial purposes by anyone.

The primary use of facial-recognition technology is in security and surveillance, whether by private companies such as retailers and events venues, or by public bodies such as police forces to track criminals. Governments increasingly use it to identify people for national and border security.

The biggest technical obstacle to achieving accurate facial recognition thus far has been the inability of machines to identify human faces when they are only partially visible, shrouded in shadow or covered by clothing, as opposed to the high-resolution, front-facing portrait photos the computers were trained on.

To teach a machine how to better read and recognise

# Now somebody's face is used as a tracking number to watch them as they move across locations on video, which is a huge shift

**Dave Maass, senior investigative researcher at the Electronic Frontier Foundation**

a human face in these conditions, it has to be trained using hundreds of thousands of faces of all shapes, sizes, colours, ages and genders. The more natural, varied and unposed the faces are, the better they simulate real-life scenarios in which surveillance might take place, and the more accurate the resulting models for facial recognition.

In order to feed this hungry system, a plethora of face repositories — such as IJB-C — have sprung up, containing images manually culled and bound together from sources as varied as university campuses, town squares, markets, cafés, mugshots and social-media sites such as Flickr, Instagram and YouTube.

To understand what these faces have been helping to build, the FT worked with Adam Harvey, the researcher who first spotted Jillian York's face in IJB-C. An American based in Berlin, he has spent years amassing more than 300 face datasets and has identified some 5,000 academic papers that cite them.

The images, we found, are used to train and benchmark algorithms that serve a variety of biometric-related purposes — recognising faces at passport control, crowd surveillance, automated driving, robotics, even emotion analysis for advertising. They have been cited in papers by commercial companies including Facebook, Microsoft, Baidu, SenseTime and IBM, as well as by academics around the world, from Japan to the United Arab Emirates and Israel.

"We've seen facial recognition shifting in purpose," says Dave Maass, a senior investigative researcher at the EFF, who was shocked to discover that his own colleagues' faces were in the Iarpa database. "It was originally being used for identification purposes . . . Now somebody's face is used as a tracking number to watch them as they move across locations on video, which is a huge shift. [Researchers] don't have to pay people for consent, they don't have to find models, no firm has to pay to collect it, everyone gets it for free."

The dataset containing Jillian York's face is one of a series compiled on behalf of Iarpa (earlier iterations are IJB-A and -B), which have been cited by academics in 21 different countries, including China, Russia, Israel, Turkey and Australia.

They have been used by companies such as the Chinese AI firm SenseTime, which sells facial-recognition products to the Chinese police, and the Japanese IT company NEC, which supplies software to law enforcement agencies in the US, UK and India.

The images in them have even been scraped by the National University of Defense Technology in China, which is controlled by China's top military body, the Central Military Commission. One of its academics collaborated last year in a project that used IJB-A, among other sets, to build a system that would, its architects wrote, "[enable] more detailed understanding of humans in crowded scenes", with applications including "group behaviour analysis" and "person re-identification".

In China, facial scanning software has played a significant role in the government's mass surveillance and detention of Muslim Uighurs in the far-western region of Xinjiang. Cameras made by Hikvision, one of the world's biggest CCTV companies, and Leon, a former partner of SenseTime, have been used to track Muslims all over Xinjiang, playing a part in what human-rights campaigners describe as the systematic repression of millions of people.
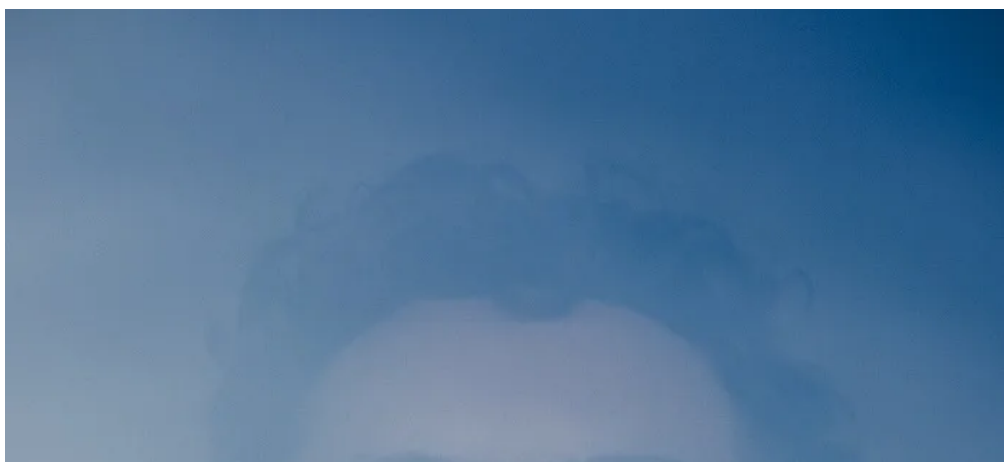
Earlier this week, it emerged that SenseTime had sold its 51 per cent stake in a security joint venture with Leon in Xinjiang after the growing international outcry over the treatment of the Uighurs.

"That was the shocking part," York says, as she considers the ways multiple companies and agencies have used the database. "It's not that my image is being used, it's about how it's being used."

———

**Harvey has been investigating** face datasets since 2010. The collection he has built up in that time comprises datasets that are readily accessible to researchers for academic and commercial purposes. The 37-year-old has been analysing where these faces come from, and where they've ended up.

By mapping out these biometric trade routes, he has started to slowly piece together the scale of distribution of faces, which may have contributed to commercial products and surveillance technologies without any explicit permission from the individuals in question.

"There's an academic network of data-sharing, because it's considered publicly beneficial to collaborate. But researchers are ignoring the stark reality that once your face is in a dataset, it's impossible to get out of it because it's already been downloaded and re-used all over the world," he says over coffee in Berlin's Mitte neighbourhood.

'People didn't have the expectation that a photo of them at a party would end up all over the world and would never be able to be removed,' Adam Harvey, a technology researcher and privacy activist based in Berlin © Mustafah Abdulaziz

Harvey likes being off the map; he prefers to pay in cash, uses the anonymous Tor browser and communicates through the encrypted app Signal, where messages disappear within a few hours. After studying engineering and photojournalism at Pennsylvania State University, he worked as a photographer at private events in New York.

"At the time, more and more people were putting photos online, because digital cameras were relatively new. What bothered me was you couldn't get those photos offline once someone put them up there," he says. "People didn't have the expectation that a photo of them at a party would end up all over the world and would never be able to be removed. You lose control of your data, you lose control of your narrative."

His fascination with surveillance technology resulted in CV Dazzle, his first biometrics project in 2011, where he created a series of make-up and hairstyling designs to enable ordinary people to hide from automated facial-recognition systems. His latest project, MegaPixels, for which he has received funding from the Mozilla Foundation, a non-profit, is also aimed at empowering citizens to

understand the new world they are inhabiting: it launched this week as a searchable database of all the papers citing every dataset he has unearthed.

Over the following months, he plans to develop the search tool to allow people to type in their names and see if their faces have been used to train an artificial intelligence system in any part of the world. "When everyone talks about facial recognition and surveillance, they're usually talking about the implementation of it. But if you take a few steps back, none of that would exist without the faces, and this project looks at where that data comes from," he explains.

One of the first large-scale face databases was created by the US defence department in the 1990s,

by paying military personnel to pose for photographs in studios. By 2007, researchers had started to realise that studio photos were too unrealistic to use as training sets, and that facial-recognition systems built on them would fail in the real world. What was needed was more images in "the wild" — natural, blurred and unposed.



The market for facial recognition has grown 20 per cent annually over the past three years, and will be worth $9bn by 2022. The speed and accuracy of the software has advanced thanks to recent strides in machine learning © Sébastien Thibault

That year, a new set was released by the University of Massachusetts, Amherst, called Labeled Faces In The Wild (LFW), a set of images scraped from news stories on the internet, with each face named. These "wild" images showed subjects with different poses and expressions, under different lighting conditions, and researchers realised they had tapped a gold mine: the internet. Here was a fount of face data that was varied, unrestricted and easy to come by. LFW has now become the most widely used facial recognition benchmark globally.

After 9/11, when defence and intelligence agencies were spurred on to invest further in facial recognition, millions of dollars were injected into labs with the explicit aim of collecting more photos in the wild, for more robust biometric profiles. Today, the default method if you need a training set is to go to search engines such as Google or Bing, or social networks such as Flickr and YouTube, where multimedia are often uploaded with a Creative Commons licence, and take what you need.

# 37,000

The number of photos taken from Flickr albums to create Facebook's People in Photo Albums

The market for facial recognition has grown 20 per cent annually over the past three years, and will be worth $9bn by 2022, according to estimates by Market Research Future. The speed and accuracy of the software has advanced thanks to recent strides in machine learning, the technology by which computers can learn to recognise specific objects — such as faces — by training on large datasets such as IJB-C.

"In 2019, there are dozens of datasets created this way that have been passed all over the world, some of them funded by defence departments, others funded directly by commercial facial recognition companies, and some of those working with controversial cyber-authoritarian regimes like in China," Harvey says. "But no one ever stopped to think if it was ethical to collect images of people's weddings [and] family photo albums with children, or if people who uploaded the photos even knew what they were doing when they got the licence."

---

**Today, as law enforcement bodies** and governments are keen to push the limits of facial recognition to lower the costs and time required to identify suspects, the consensus is that even internet faces aren't truly "wild" because people tend to put up edited photos. The appetite for more diverse face data has led down a rabbit hole of capturing people's images as naturally as possible, often without their knowledge.

Take the [UnConstrained College Students Dataset](https://www.ft.com). About 1,700 students at the University of Colorado, Colorado Springs, were photographed on 20 different days, between February 2012 and September 2013, using a "long-range high-resolution surveillance camera without their knowledge," according to the set's creator, Professor Terry Boult, a computer scientist at the university.

"Even [LFW] photos aren't that wild because people know they are being photographed and uploaded on the internet. But these are students walking on a sidewalk on campus, who are unaware they are part of a data collection," Boult tells me. "When you're watching students on a sidewalk, there's an awful lot of facing down looking at your phone. In Colorado, where it's cold and

snowy, they cover up in a natural way with scarves and hats. Our goal is to make it the most realistic unconstrained video surveillance facial recognition dataset in the world."

Boult's lab has been funded by the US Navy, as well as by Iarpa, which has disseminated his facial recognition research to the Department of Homeland Security and other government bodies. Because Boult was filming students in a public place (although it is within the university), he says he did not need their permission as long as he does not know their identities, according to Colorado state law.

"We didn't make it publicly available for anyone to download. They have to contact us through a

website — we check if they're researchers," Boult says. "It's not for commercial use, but if corporate researchers are trying to make facial recognition better for their company's products, we are OK with them doing that as long as they publish."

This apparently porous relationship between an "academic" use of data that has been collected without consent, and the commercial exploitation of that same data, highlights the complex ethical questions surrounding face sets. The paper in which Boult introduced his dataset has been cited by six universities in Chile, Italy and the US and downloaded by at least four private companies based in Europe, China and the US, none of which seems to have published its work.

Another campus dataset, Duke-MTMC, was collected at Duke University in North Carolina. Funded partly by the US Army Research Office, it has become one of the most popular pedestrian recognition training sets, and has been cited by 96 institutions globally.

## China dominates: a breakdown of the countries where academics have used the Duke MTMC database

Number of citations

| China | US | UK |
|-------|-----|-----|
| 80 | 37 | 17 |
| | Australia 21 | Germany 4 · Italy 4 · Jap. 3 |
| | | Switzerland 3 |
| | | Others 10 |

Source: Megapixels
© FT

Researchers used eight cameras to film students walking on campus, and notified them through posters placed around the perimeter of the surveilled area. Ergys Ristani, one of the authors, said the work had been approved by an institutional review board and, despite the signs, not a single student had asked to be excluded. From the video footage, it is unclear whether the students have seen the posters or are aware they are being filmed.

In a third case, footage of customers in a café called Brainwash in San Francisco's Lower Haight district, taken through a livestreaming camera, was turned into a pedestrian dataset — Brainwash — that has been cited by companies including Huawei and Qualcomm.

"When you're in urban space you have a reasonable expectation of anonymity; this is recognised in [US] jurisprudence, this is something so deeply a part of common sense that it interferes with our ability to understand [that] our tech companies are impacting our privacy, even in nominally private spaces," says Adam Greenfield, a technology writer and urbanist. Greenfield has recently discovered himself in a database of one million faces created by Microsoft called MSCeleb.

———

**Researchers point out** that facial analysis technologies aren't just for surveillance — they could be used for health monitoring; for instance, scanning faces to see if someone is developing dementia or type 2 diabetes, or to check for drowsiness or inebriation in drivers.

If datasets weren't shareable, corporations such as Facebook and Google, which have billions of user photos and videos uploaded to their sites each day, would be the only organisations with access to an ocean of high-quality face data and so would have the best face recognition algorithms, some researchers argue.

"I'm not worried about government, I'm worried about Google and Facebook," says Karl Ricanek, a professor at the University of North Carolina Wilmington who has built two publicly accessible face datasets. "In my opinion, they have more info on citizens than governments themselves, and we can't affect leadership at these companies. I think our government at least has a good mission. From an academic perspective we are trying to solve problems that we think will make life better in our world. Most of us aren't attempting to make money."

## I'm worried about Google and Facebook. In my opinion, they have more info on citizens than governments themselves

**Karl Ricanek, professor at the University of North Carolina Wilmington**

Despite commercial companies often having their own extensive data pools, they too have increasingly turned to the internet to cull larger and more natural datasets used to benchmark and train algorithms. For instance, Facebook created a dataset called People in Photo Albums, consisting of more than 37,000 photos of 2,000 individuals, including children, from personal Flickr photo albums. "While a lot of progress has been made recently in recognition from a frontal face, non-frontal views are a lot more common in photo albums than people might suspect," Facebook researchers wrote in their paper.

Their dataset specifically picks out photos with "large variations in pose, clothing, camera viewpoint, image resolution and illumination", and the paper describes a new algorithm that can recognise even these partially hidden faces with high accuracy.
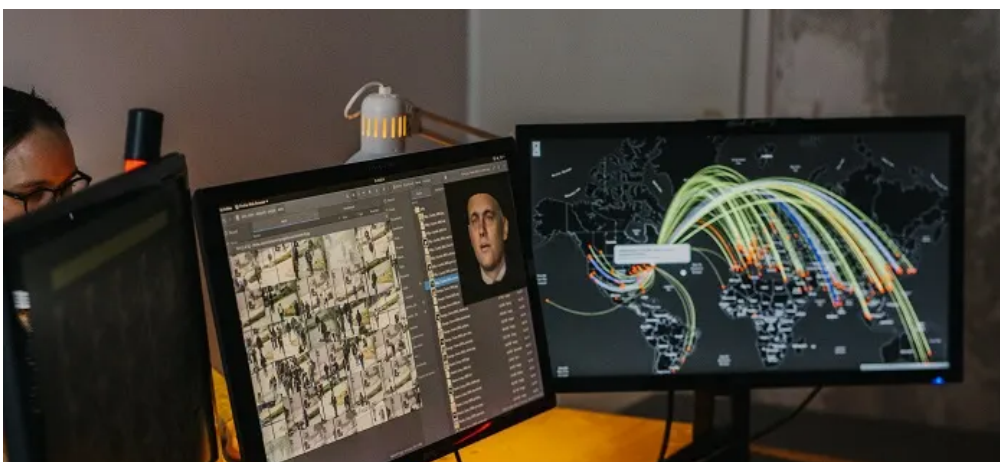
"We hope our dataset will steer the vision community towards the very important and largely unsolved problem of person recognition in the wild," the Facebook researchers conclude. This dataset has now been reused all over the world, including by the National University of Defense Technology in China to improve video surveillance technology.

"To be clear, we are not collaborating with the Chinese government on facial recognition and never have," a Facebook spokesperson said. "However, there will always be a question if advancements in technology should be shared widely or closely held. Facebook and other leading technology companies believe that the scientific community can share learnings to advance technology, while also working to prevent abuse."

Maass, the Electronic Frontier Foundation researcher, says: "This is not a question of legality but of morals and ethics. I'm not sure where a research project like this would fall, but I wonder if creating a dataset to train surveillance tech is the same as conducting surveillance yourself."

Where images have been scraped off the internet, researchers — even at large companies such as Microsoft, IBM and Facebook — have relied on the Creative Commons licence to stand in for user consent. But the application of these photos as training data for surveillance and facial analysis is so far removed from what the licence was originally intended to cover that Creative Commons itself, a non-profit, recently put out a statement to clarify the change.

"CC licences were designed to address a specific constraint, which they do very well: unlocking restrictive copyright. But copyright is not a good tool to protect individual privacy, to address research ethics in AI development, or to regulate the use of surveillance tools employed online," chief executive Ryan Merkley wrote in March. "Those issues rightly belong in the public policy space . . . [we hope to] engage in discussion with those using our content in objectionable ways, and to speak out on . . . the important issues of privacy, surveillance, and AI that impact the sharing of works on the web."

Adam Harvey's workspace. He likes being off the map; prefers to pay in cash, uses the anonymous Tor browser and communicates through the encrypted app Signal © Mustafah Abdulaziz

Ultimately, experts believe it is too late to put the face data back in a box, or to restrict its movement across geographic borders. "We may trust the entity that gathers or captures that information, and we may have consented to an initial capture, but custody over that dataset leaks," says Greenfield, the tech writer who features in Microsoft's celebrity dataset. "It can leak through hacking, corporate acquisition, simple clumsiness, it can leak through regime change. Even if [creators] attempted to control access, there's no way they could stop it coming into the hands of the Israeli, American or Chinese state, or anyone who wants to train up facial-recognition algorithms."

For Harvey, who has spent almost a decade trying to illustrate the scale of the issue, there doesn't seem to be an end in sight. "There are so many egregious examples in these datasets that are just brazen abuses of privacy," he says. "Some of them come from public cameras pointed at the street, and there are even a few that came from cameras in cafés. After looking at these, you never know when you could walk in front of a camera that may one day be part of a training dataset."

In fact, recognising a face is only the first step of biometric surveillance, he suggests. "It's really like an entry-level term to much broader, deeper analysis of people's biometrics. There's jaw recognition — the width of your jaw can be used to infer success as CEO, for example. Companies such as Boston-based Affectiva are doing research that analyses faces in real time, to determine from a webcam or in-store camera if someone is going to buy something in your store."

Other analyses, he adds, can be used to determine people's tiredness, skin quality and heart rate, or even to lip-read what they are saying. "Face recognition is a very deceiving term, technically, because there's no limit," he concludes. "It ends ultimately only with your DNA."

*Madhumita Murgia is the FT's European technology correspondent. Additional reporting by Max Harlow*

*Follow @FTMag on Twitter to find out about our latest stories first. Subscribe to FT Life on YouTube for the latest FT Weekend videos*

## Letter in response to this article:

*Sunning, strolling — and captured on film / From Amanda Nicholls, London E17, UK*
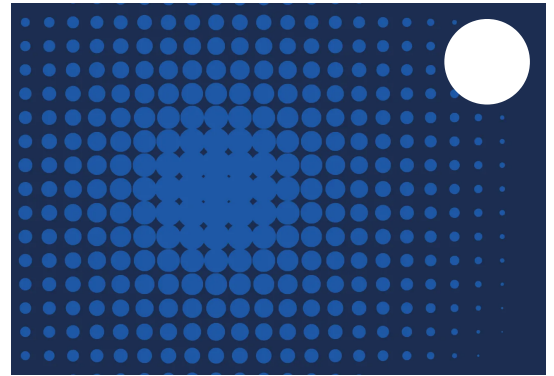
# FT-OPTIV CISO BRIEFING

London
13 November 2019

Cyber is broken. Is the damage irreparable?

**Register now**

Presented by