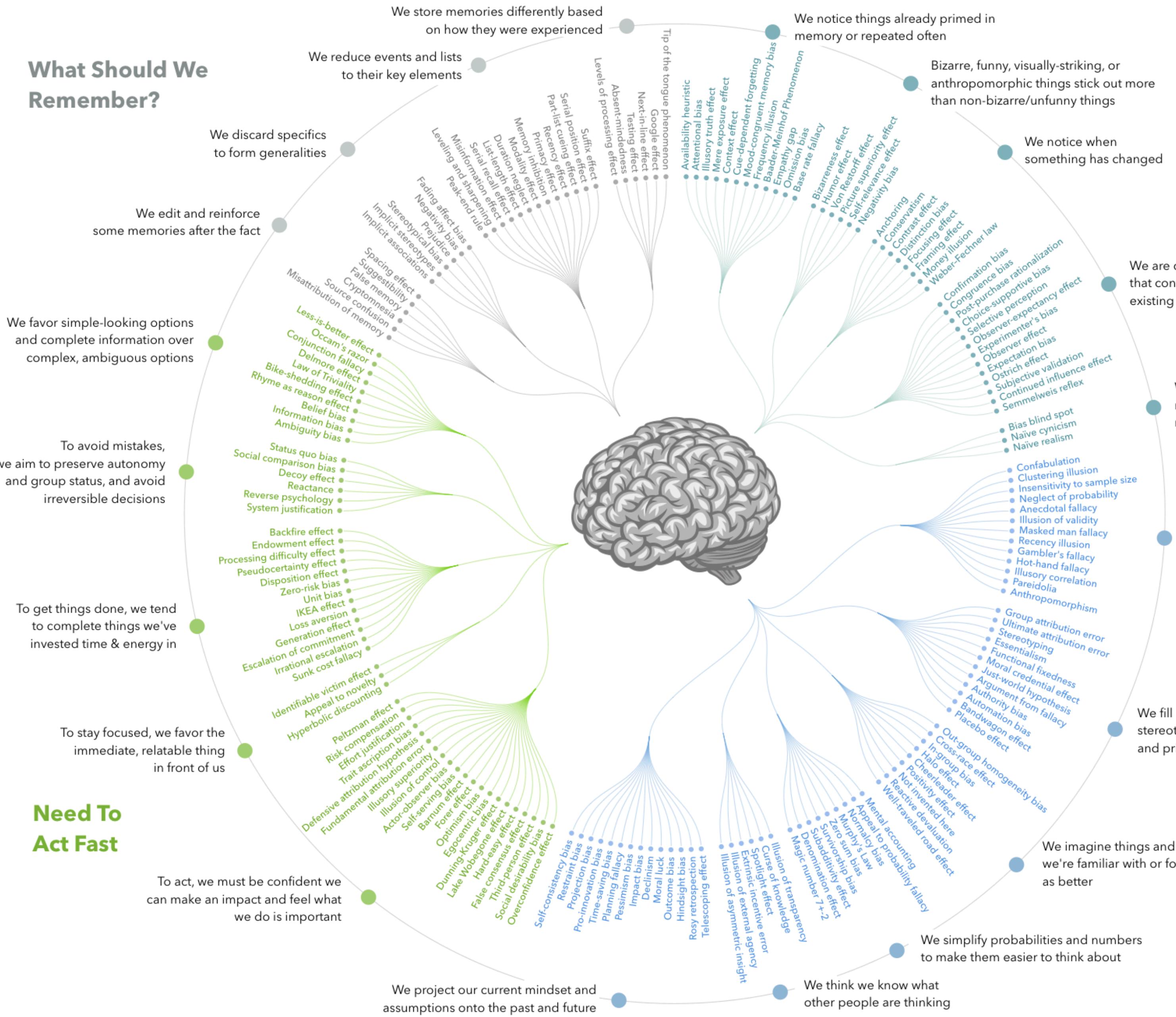


# **Errors of human cognition**

# COGNITIVE BIAS CODEX

## What Should We Remember?



## Too Much Information

We notice flaws in others more easily than we notice flaws in ourselves

We tend to find stories and patterns even when looking at sparse data

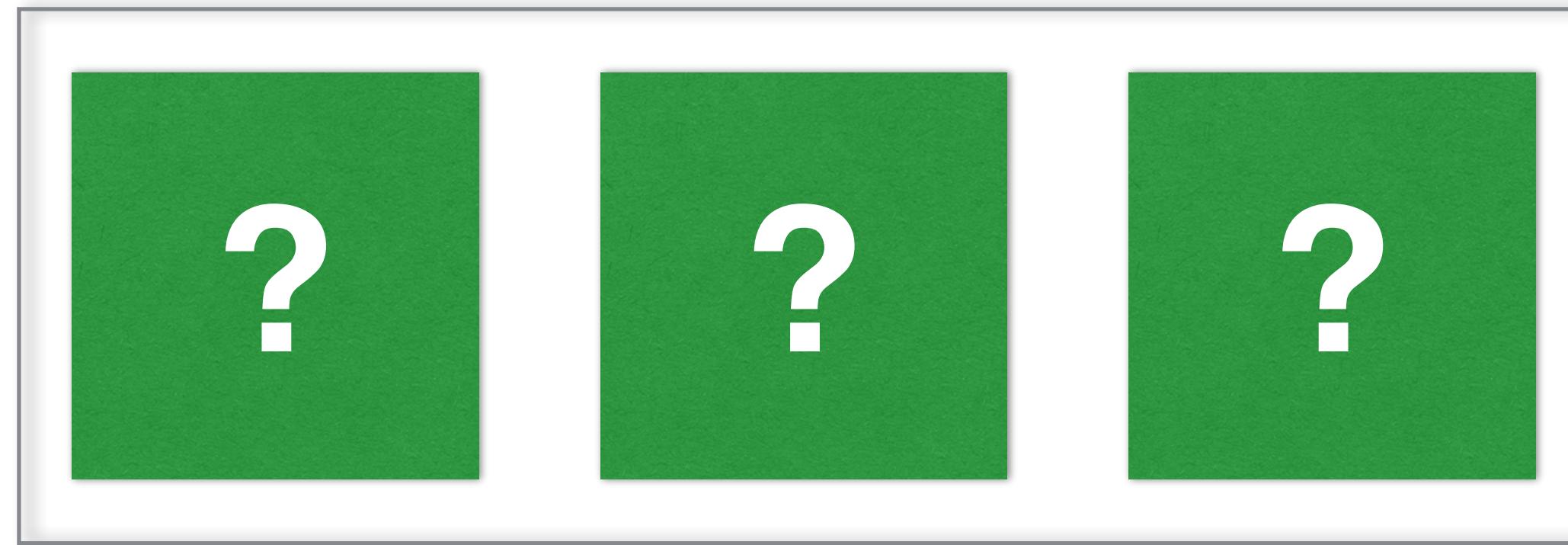
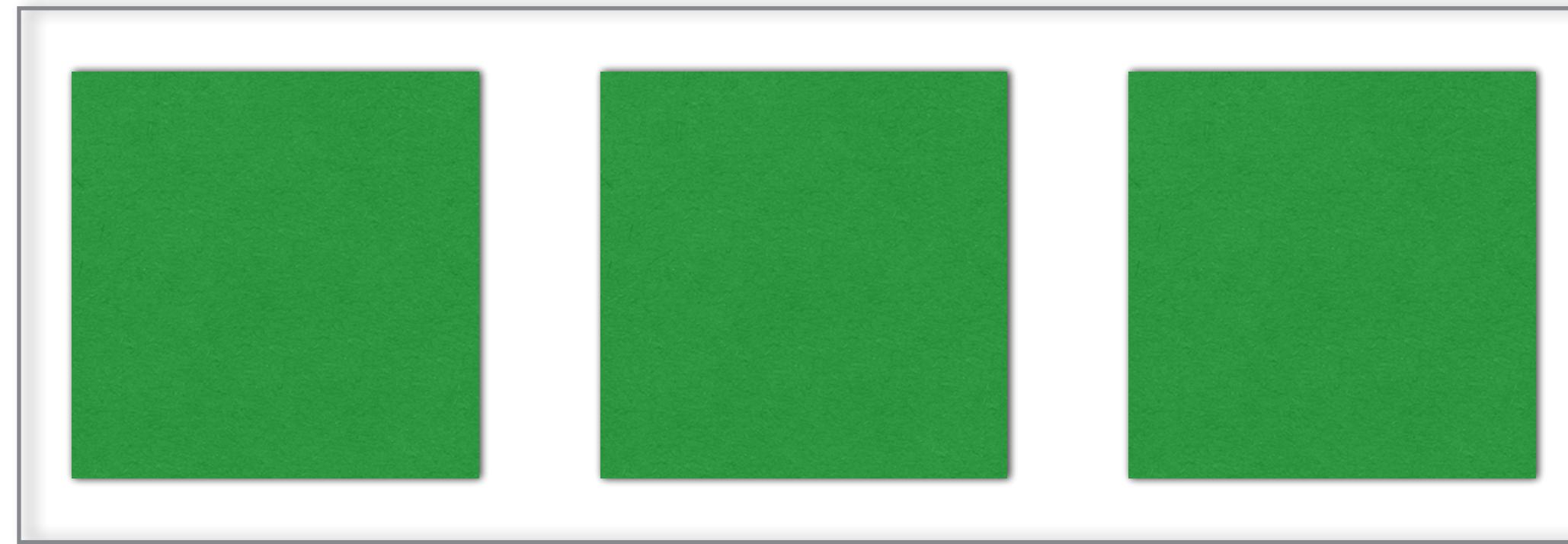
We fill in characteristics from stereotypes, generalities, and prior histories

We imagine things and people we're familiar with or fond of as better

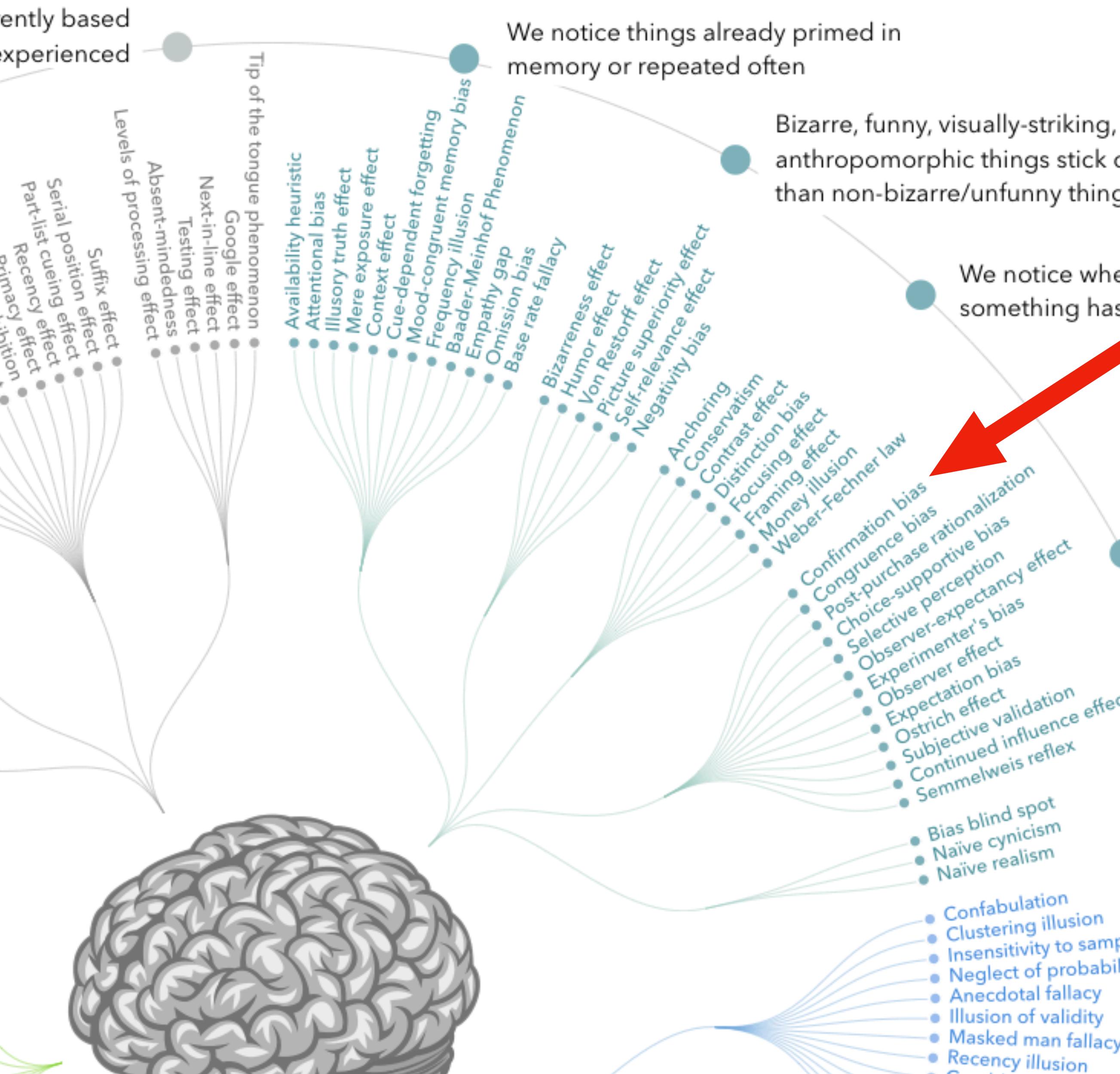
## Not Enough Meaning

We are drawn to details that confirm our own existing beliefs

# Figure out the rule



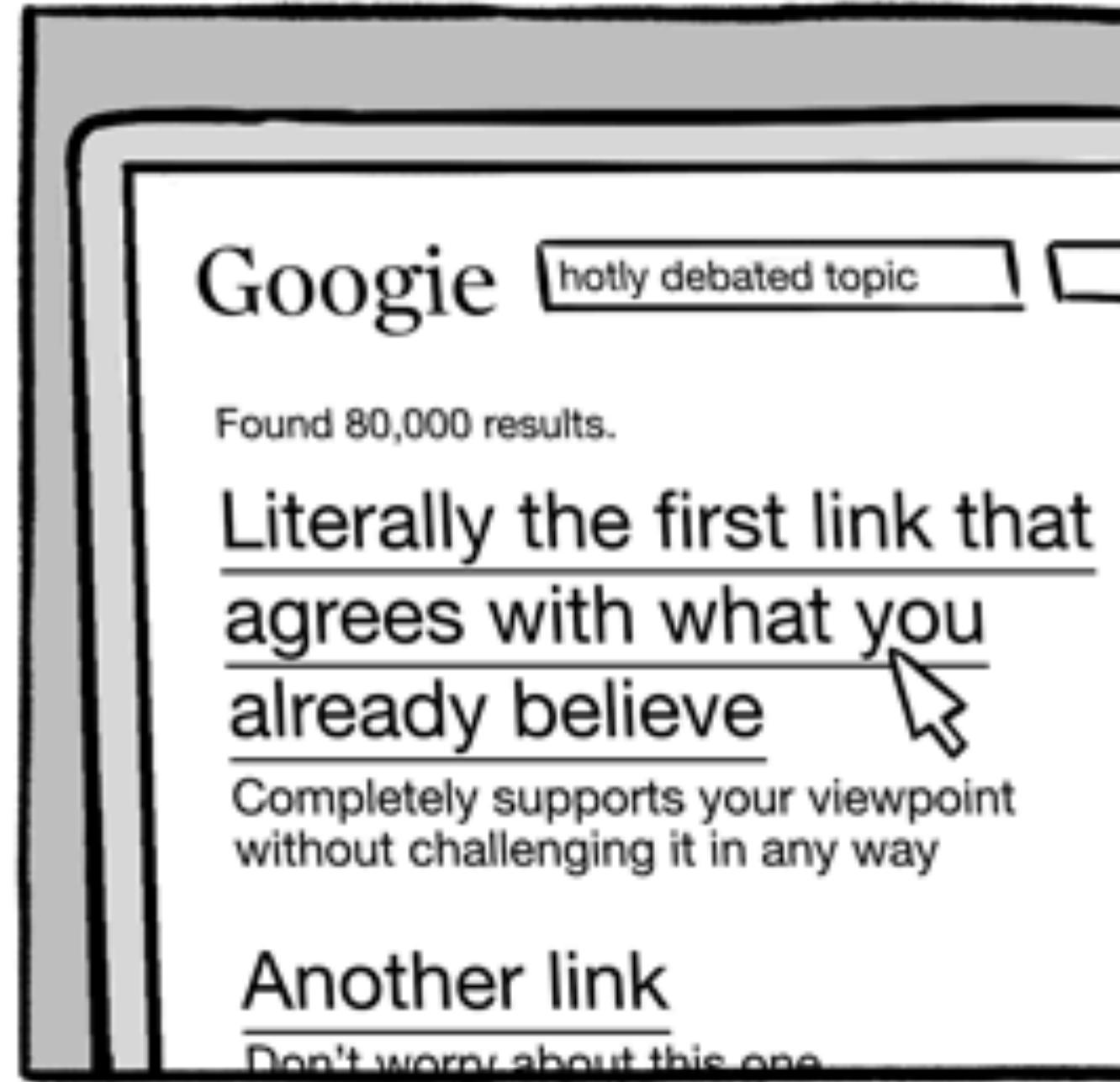
# ITIVE BIAS CODEX



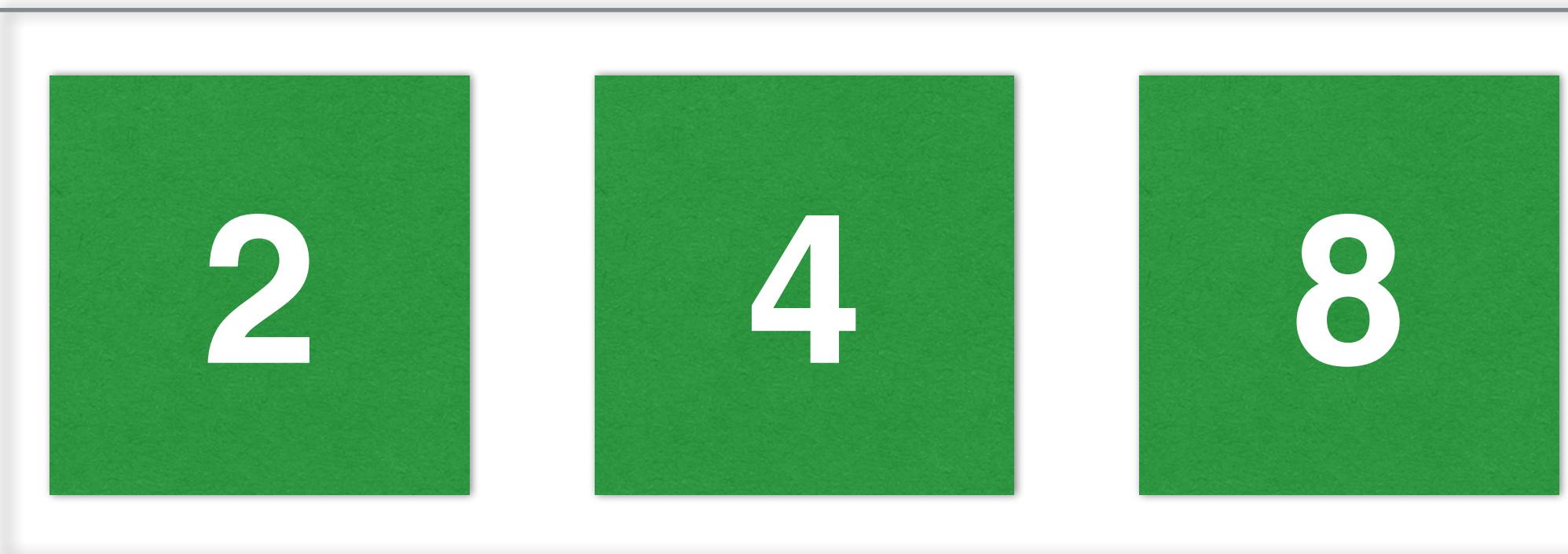
## Too Much Information

# Confirmation bias

CHAINSAWSUIT.COM



**The only way to show you're right  
is to try to show you're wrong!**



**Confirmation bias...  
and programming?**

We favor simple-looking options and complete information over complex, ambiguous options

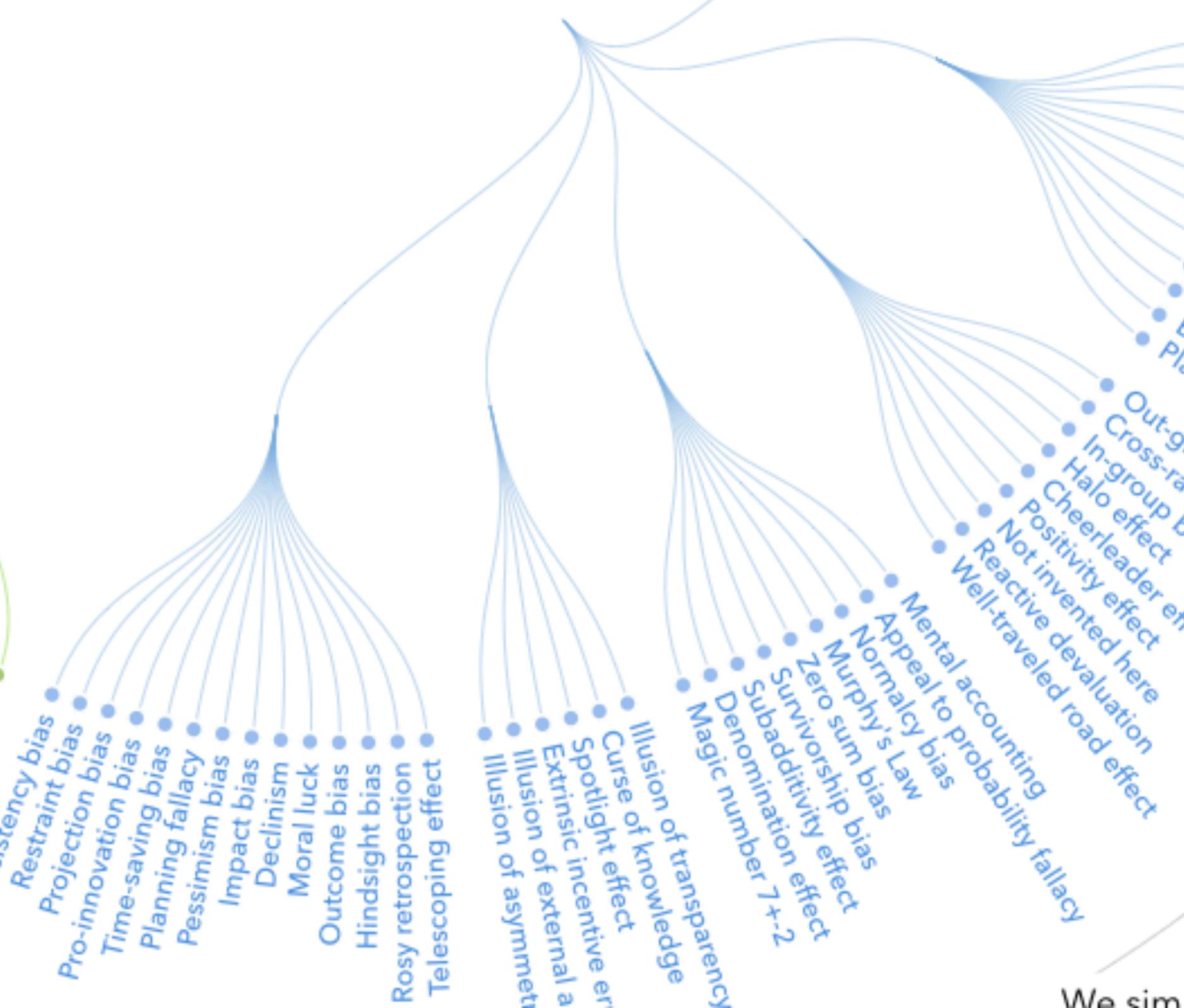
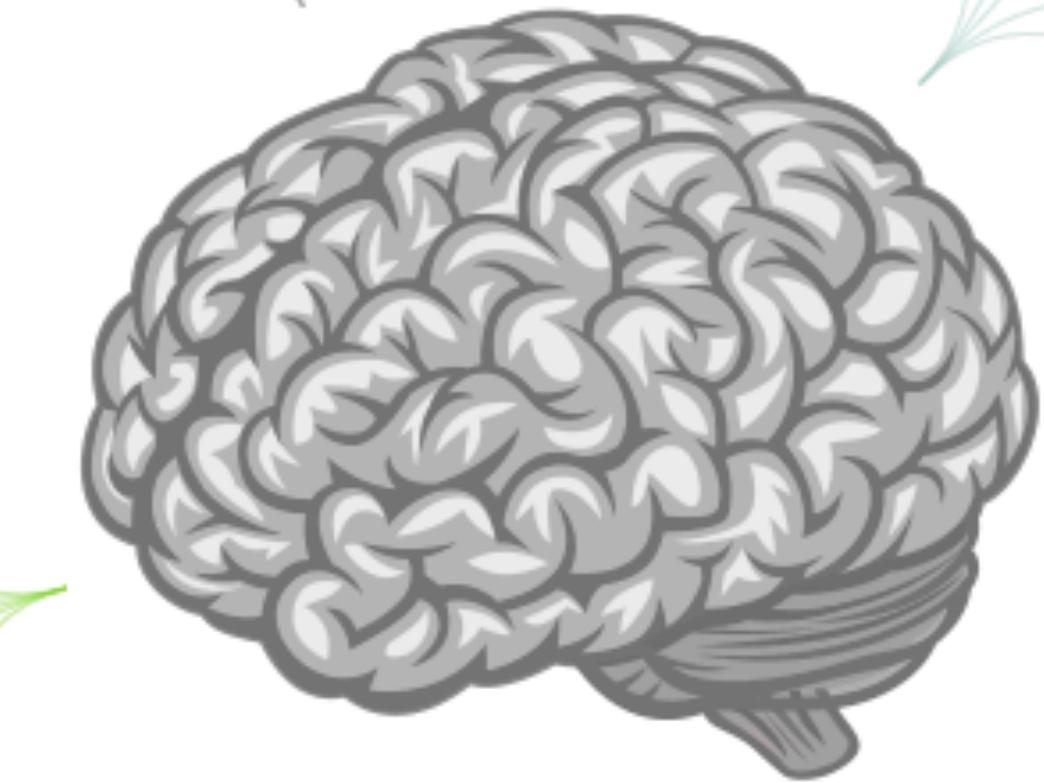
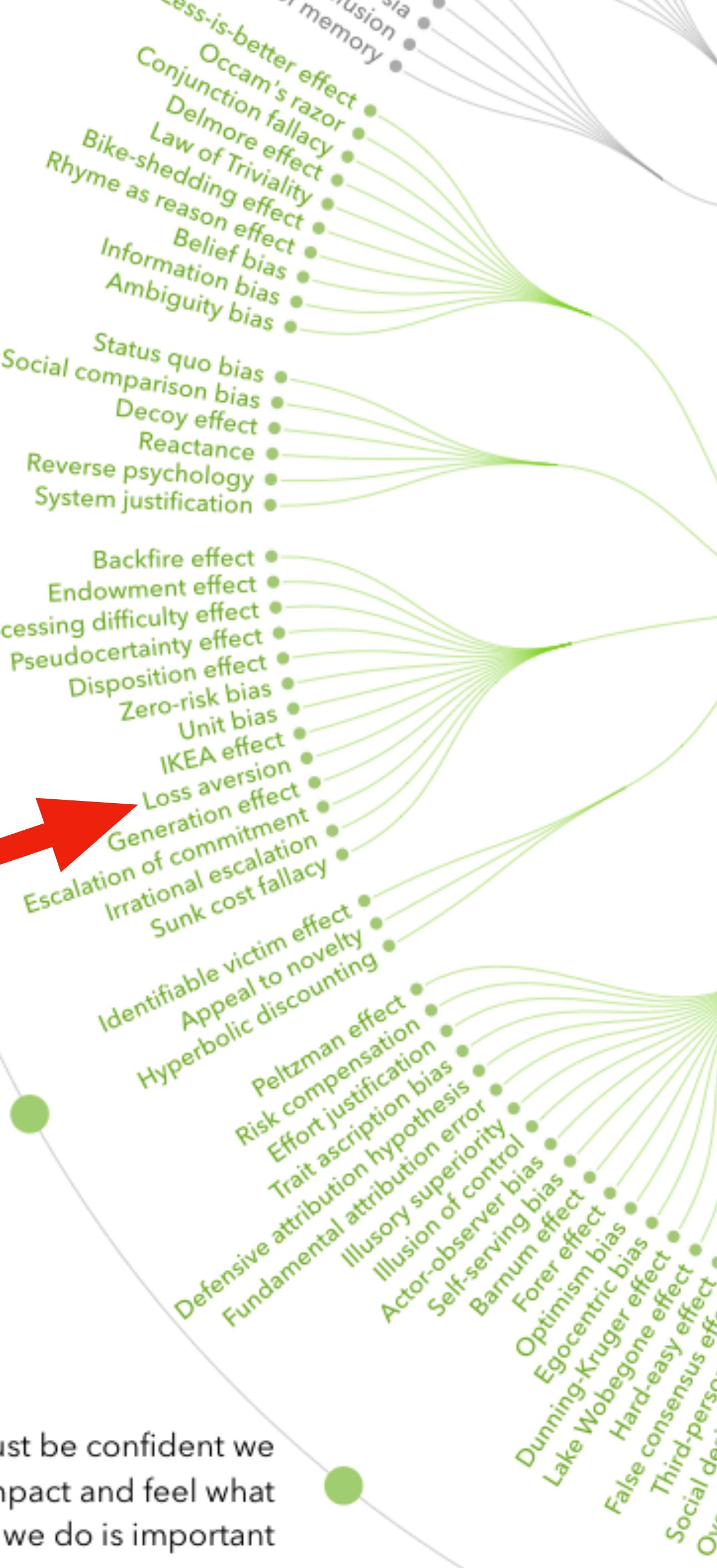
To avoid mistakes, we aim to preserve autonomy and group status, and avoid irreversible decisions

To get things done, we tend to complete things we've invested time & energy in

To stay focused, we favor the immediate, relatable thing in front of us

## Need To Act Fast

To act, we must be confident we can make an impact and feel what we do is important



We sim

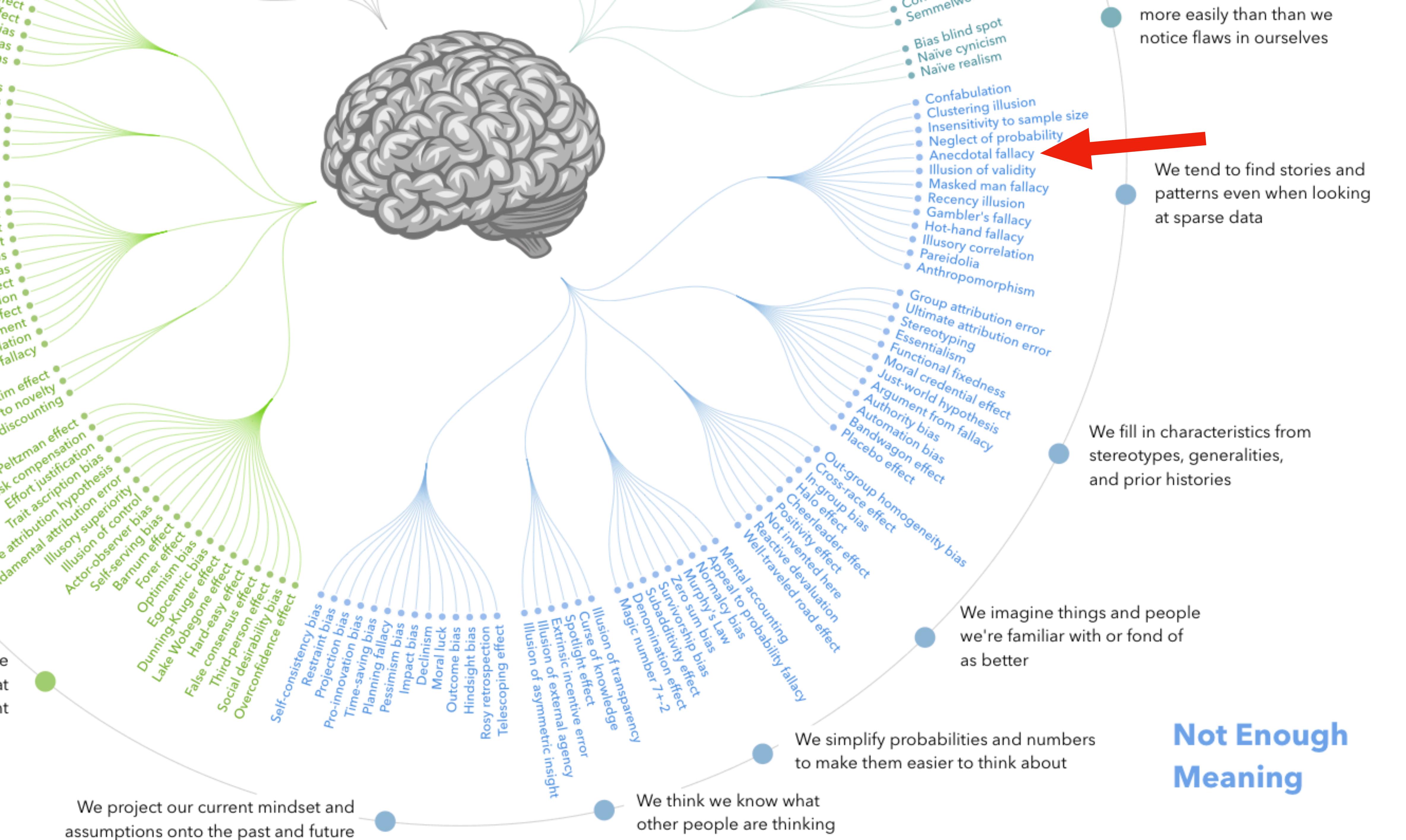
# Should I try base jumping? It could be fun!

## Loss aversion

If you tell someone they have a 99.5% of surviving they think it's fine

...but tell them 1 in 200 people die and they think it's too scary







## On AI, Analytics, and the New Machine Age

If you read nothing else on how intelligent machines are revolutionizing business, read these definitive articles from Harvard Business Review.

**HBR'S  
10  
MUST  
READS**

BONUS ARTICLE  
"Why Every Organization Needs an Augmented Reality Strategy"  
By Michael E. Porter and James E. Heppelmann



**Anecdata**  
**Rigorous analysis supported by anecdote**

# **Errors of communication**

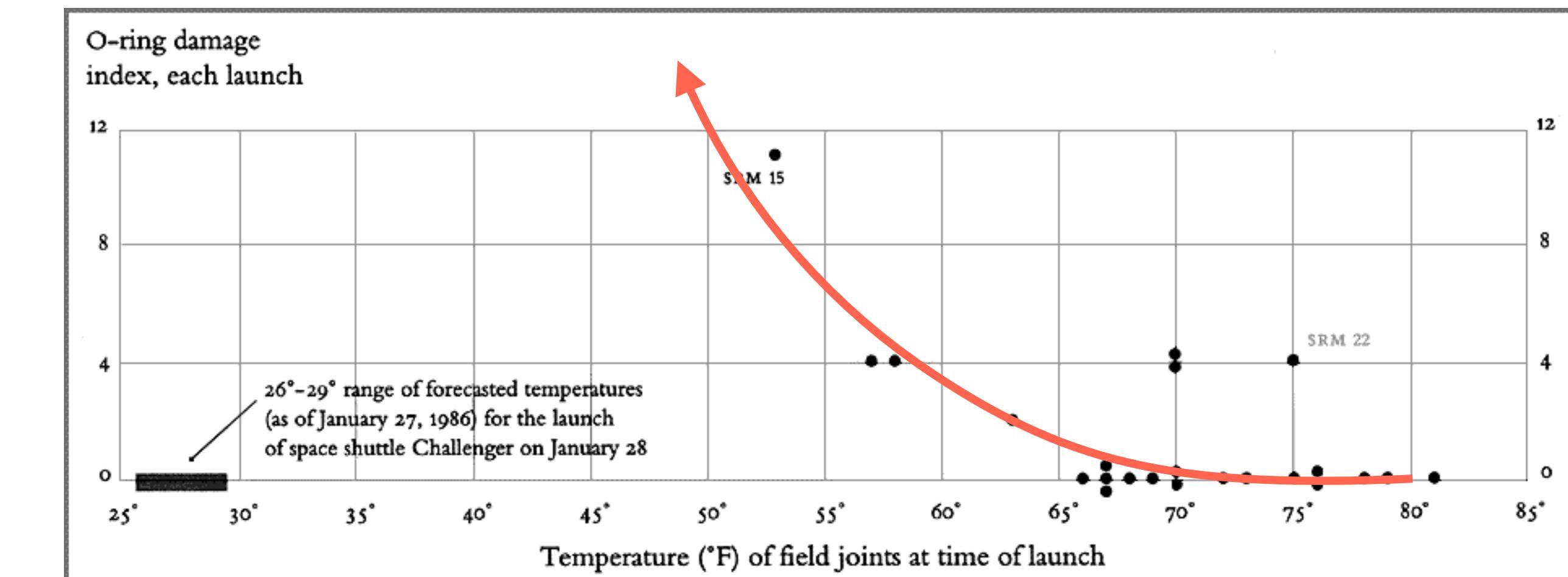
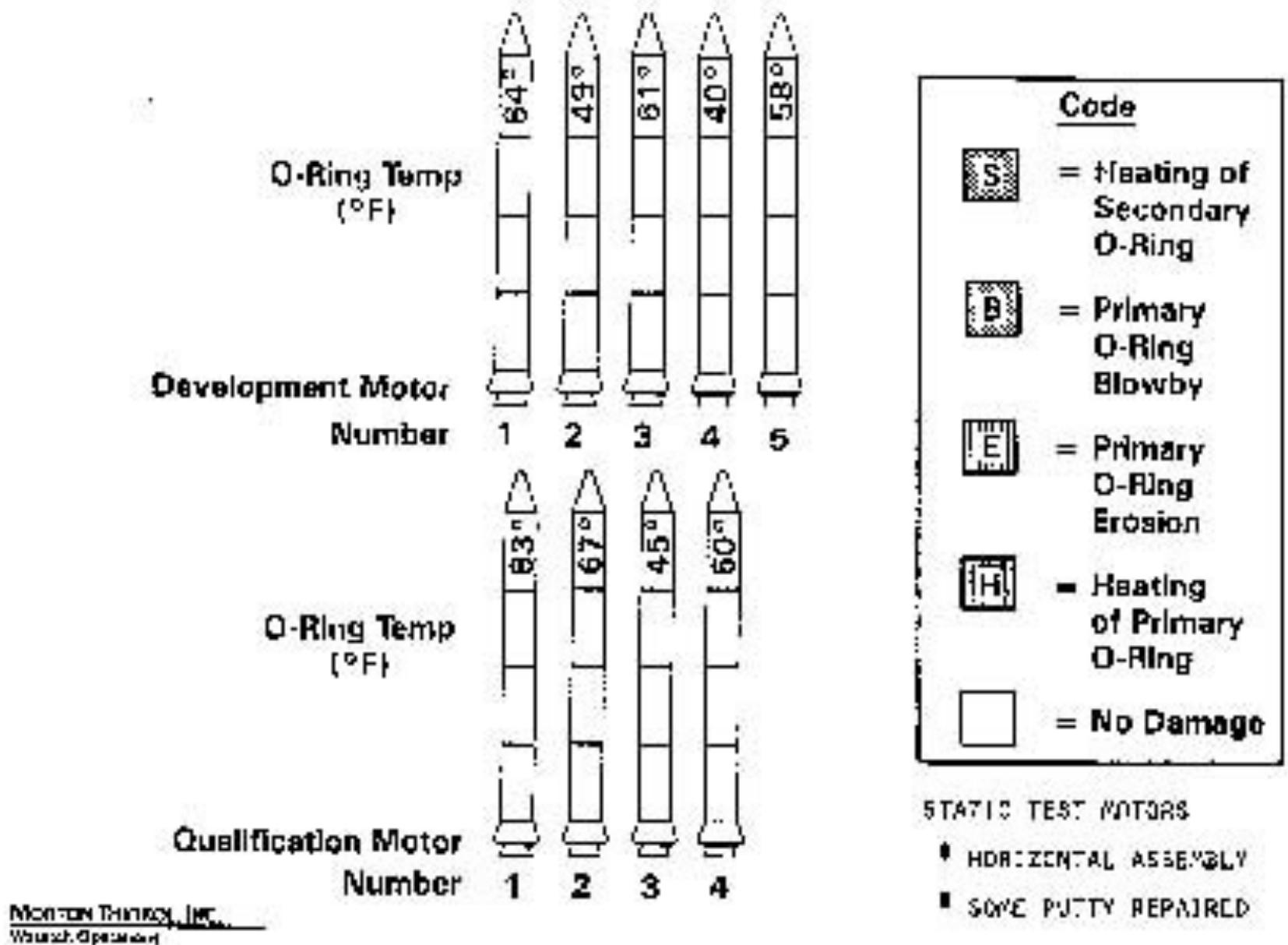


Jan 28, 1986

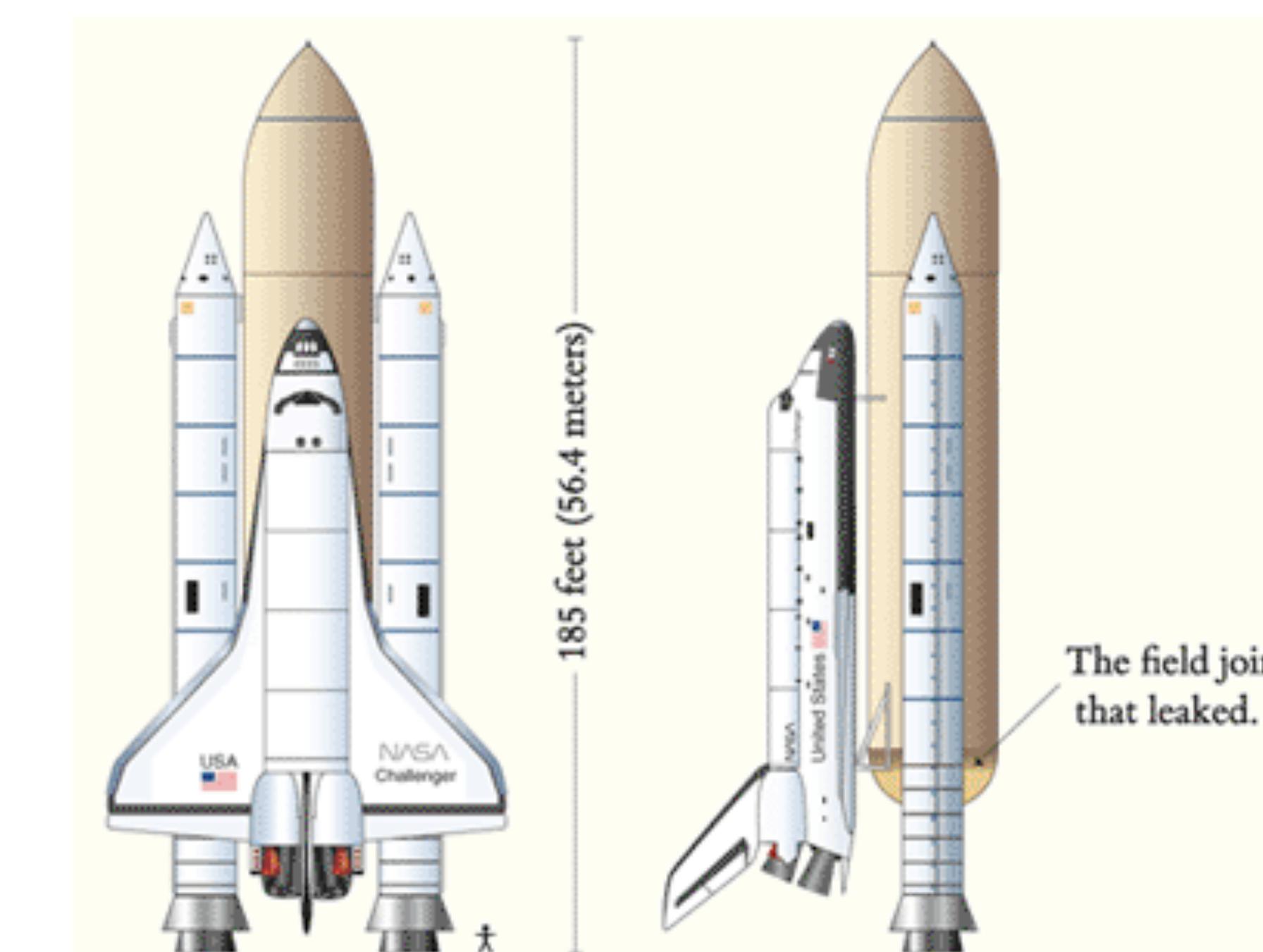
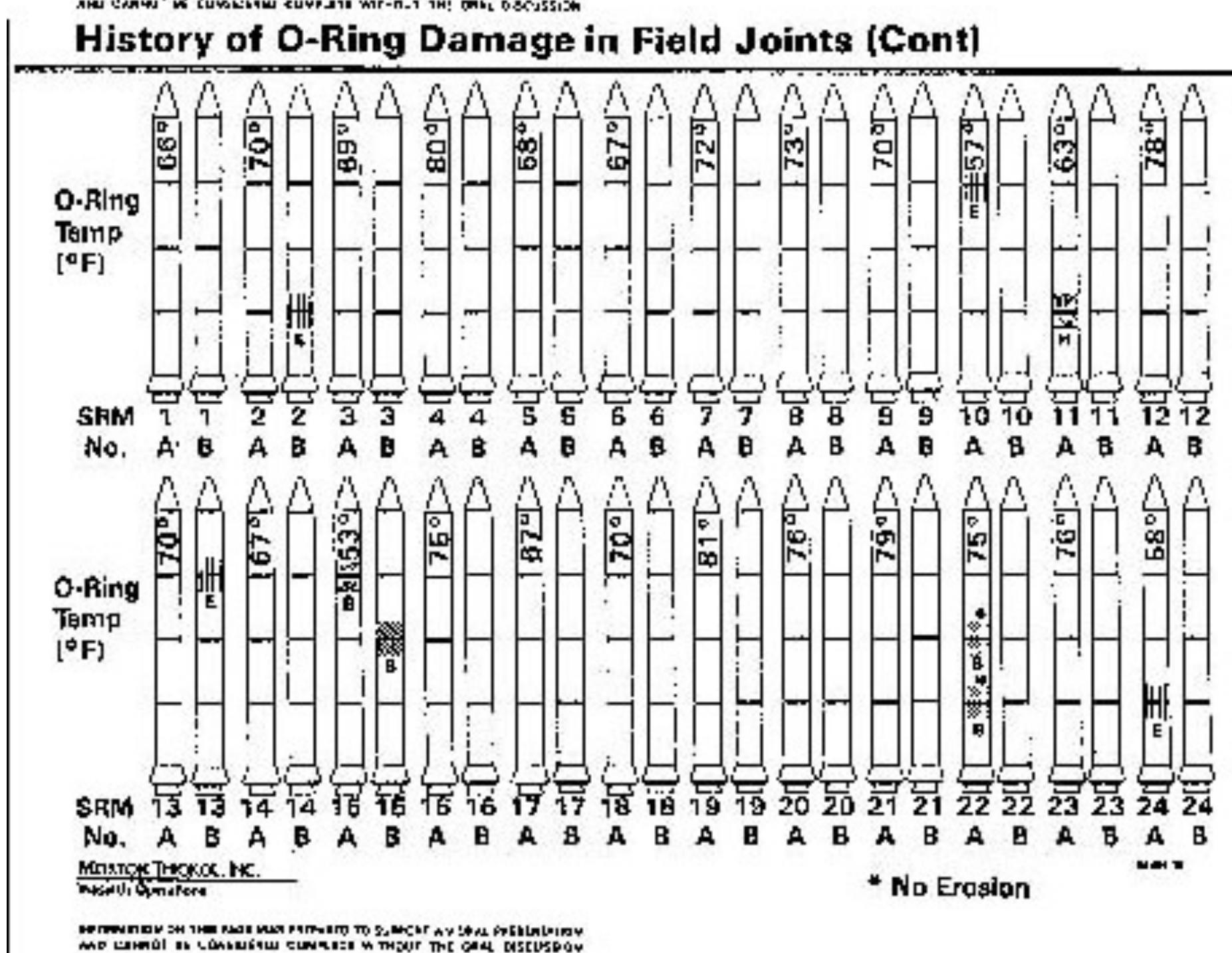


Feb 1, 2003

## **History of O-Ring Damage in Field Joints**



# Images and graphs from Edward T



On this one Columbia slide, a PowerPoint festival of bureaucratic hyper-rationalism, 6 different levels of hierarchy are used to display, classify, and arrange 11 phrases:

- Level 1 Title of Slide
- Level 2 ● Very Big Bullet
- Level 3 — big dash
- Level 4 • medium-small diamond
- Level 5 • tiny square bullet
- Level 6 ( ) parentheses ending level 5

The analysis begins with the dreaded Executive Summary, with a conclusion presented as a headline: "Test Data Indicates Conservatism for Tile Penetration." This turns out to be unmerited reassurance. Executives, at least those who don't want to get fooled, had better read far beyond the title.

The "conservatism" concerns the *choice of models* used to predict damage. But why, after 112 flights, are foam-debris models being calibrated during a crisis? How can "conservatism" be inferred from a loose comparison of a spreadsheet model and some thin data? Divergent evidence means divergent evidence, not inferential security. Claims of analytic "conservatism" should be viewed with skepticism by presentation consumers. Such claims are often a rhetorical tactic that substitutes verbal fudge factors for quantitative assessments.

As the bullet points march on, the seemingly reassuring headline fades away. Lower-level bullets at the end of the slide undermine the executive summary. This third-level point notes that "Flight condition [that is, the debris hit on the Columbia] is significantly outside of test database." How far outside? The final bullet will tell us.

This fourth-level bullet concluding the slide reports that the debris hitting the Columbia is estimated to be  $1920/3 = 640$  times larger than data used in the tests of the model! The correct headline should be "Review of Test Data Indicates Irrelevance of Two Models." This is a powerful conclusion, indicating that pre-launch safety standards no longer hold. The original optimistic headline has been eviscerated by the lower-level bullets.

Note how close readings can help consumers of presentations evaluate the presenter's reasoning and credibility.

The Very-Big-Bullet phrase fragment does not seem to make sense. No other VBB's appear in the rest of the slide, so this VBB is not necessary.

Spray On Foam Insulation, a fragment of which caused the hole in the wing

A model to estimate damage to the tiles protecting flat surfaces of the wing

## Review of Test Data Indicates Conservatism for Tile Penetration

- The existing SOFI on tile test data used to create Crater was reviewed along with STS-87 Southwest Research data
  - Crater overpredicted penetration of tile coating significantly
    - Initial penetration is described by normal velocity
      - Varies with volume/mass of projectile (e.g., 200ft/sec for 3cu. In)
    - Significant energy is required for the softer SOFI particle to penetrate the relatively hard tile coating
      - Test results do show that it is possible at sufficient mass and velocity
    - Conversely, once tile is penetrated SOFI can cause significant damage
      - Minor variations in total energy (above penetration level) can cause significant tile damage
  - Flight condition is significantly outside of test database
    - Volume of ramp is 1920cu in vs 3 cu in for test

*BOEING*

Here "ramp" refers to foam debris (from the bipod ramp) that hit Columbia. Instead of the cryptic "Volume of ramp," say "estimated volume of foam debris that hit the wing." Such clarifying phrases, which may help upper level executives understand what is going on, are too long to fit on low-resolution bullet outline formats. PP demands the shorthand of acronyms, phrase fragments, and clipped jargon in order to get at least some information into the tight format.

Edward Tufte

# Our models are irrelevant

Debris hitting the wing was **640x** larger than the experimental data used to build these models

We have **no clue** what will happen on re-entry

# **Communication is key**

**Identify audience & setting**

**Identify key insight, main points of evidence, and assumptions**

**Organize into a story focussed on** 

**Create supporting visualizations**

**Revise to be as precise and concise as possible**

**Errors of measurement - are we measuring what we think we are?**

**Errors of analysis - did we use the right methods to address the question?**

**Errors of borked tools - choosing the wrong tools or using them poorly leads to bad results**

**Errors of human cognition - data science is a human endeavor with all the usual frailties and foibles**

**Errors of communication - sometimes you get everything right, but the group and the decision makers never understand properly**

# Your future in DS

**Jason G. Fleischer, Ph.D.**

**Asst. Teaching Professor**

**Department of Cognitive Science, UC San Diego**

**jfleischer@ucsd.edu**



**@jasongfleischer**

**<https://jgfleischer.com>**

Slides in this presentation are from material kindly provided by  
Shannon Ellis and Brad Voytek

# Courses in DS and ML at UCSD

- DS
- CSE
- CS
- ECE
- COGS
- But also many other departments like ECON, MATH, LING, BENG, etc

My list of '20-21 ML (and ML adjacent) courses

# Some job titles and what they do

- Analytics or statistician: data handling, analysis
- Data scientist: programming, data handling, analysis
- Data engineer: programming, databases, management
- Data architect: programming, databases, design
- Data manager: databases, design, management
- \*OPs (eg, devOPs, dataOPs, full stack): programming, tool development, mangagement concentrating on end to end process
- ML Engineer: programming, tool development, management of infrastructure
- ML researcher: programming, algorithm design and testing

# Glut of new data scientists

First, let's talk about the oversupply of junior data scientists. The [continuing media hype cycle around data science](#) has enormously exploded the amount of junior talent available on the market over the past five years.

This is purely anecdotal evidence, so take it with a large grain of salt. But, based on my own participation as a resume screener, mentor to data scientists leaving boot camps, interviewer, interviewee, and from conversations with friends and colleagues in similar positions, I've developed an intuition that the number of candidates per any given data science position, particularly at the entry level, has grown from 20 or so per slot, to 100 or more. I was talking to a friend recently who had to go through 500 resumes for a single opening.

This is not abnormal. More anecdotal evidence comes from job openings [like this one](#), from machine learning's godfather, Andrew Ng, whose AI startup demanded 70-80 hours a week. He was flooded with applications, after blithely noting that previously many people had tried to volunteer for free. As of this latest writing, they [ran out of space](#) in their current office.

It's very, very hard to estimate the true gap between market demand and supply, but [here's a starting point](#).

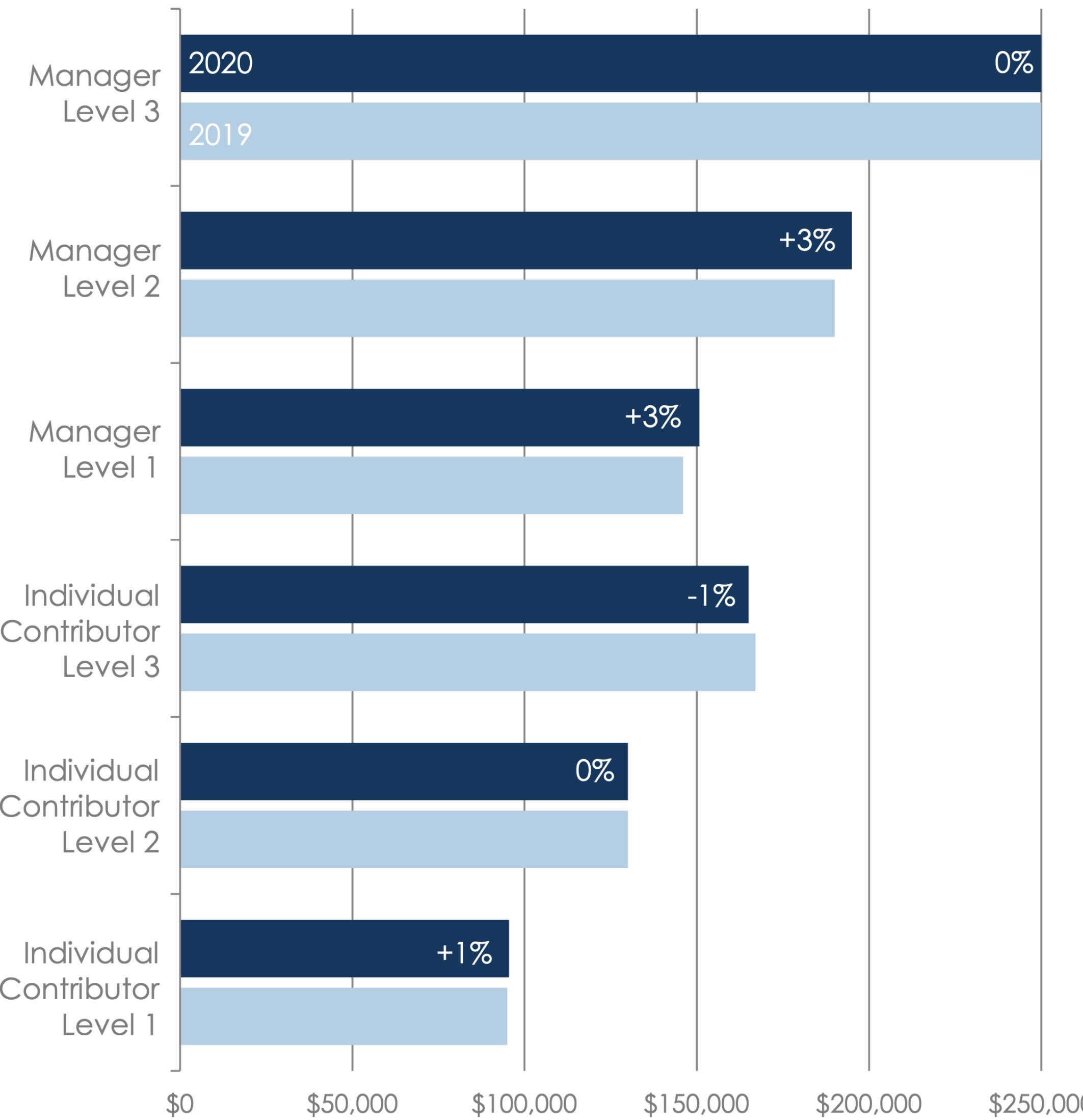
# Advice from Vicki Boykis

Sr. Manager, Data Science + Engineering at CapTech Ventures, Inc

1. Learn SQL
2. Learn a programming language extremely well and learn programming concepts.
3. Learn how to work in the cloud.
4. This stuff is really hard **for everyone**, and there are a million things it seems like you have to know. Don't get discouraged.

# Burtsworks annual predictions and report on DS hiring

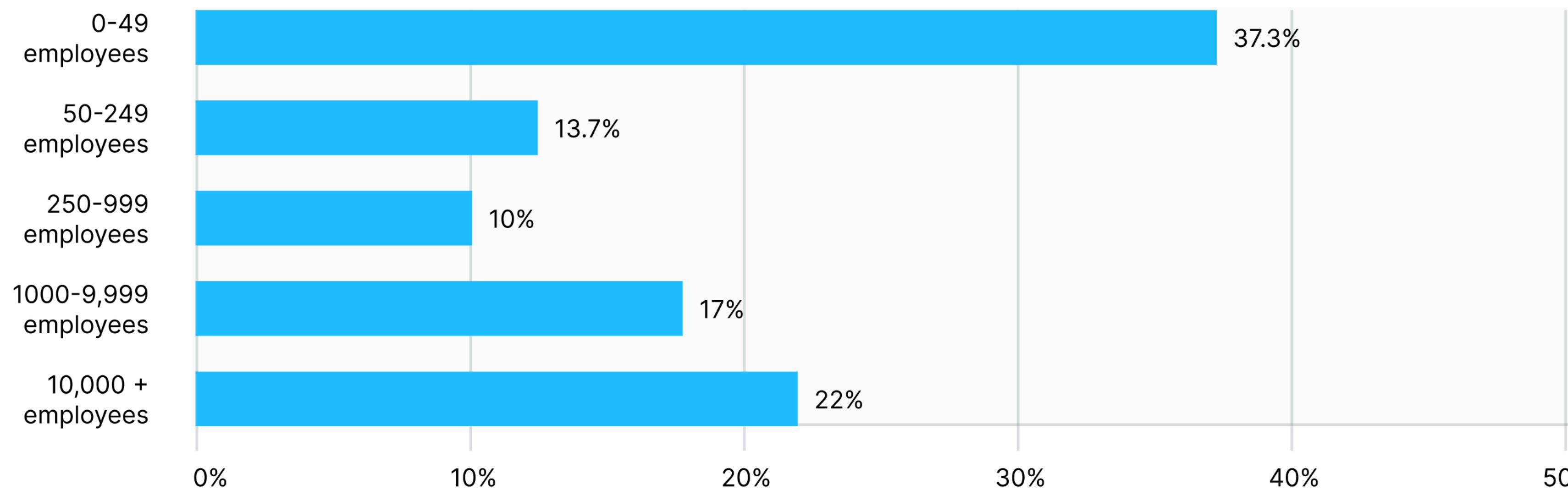
**Figure 2** Comparison of Data Scientists' Median Base Salaries by Job Category



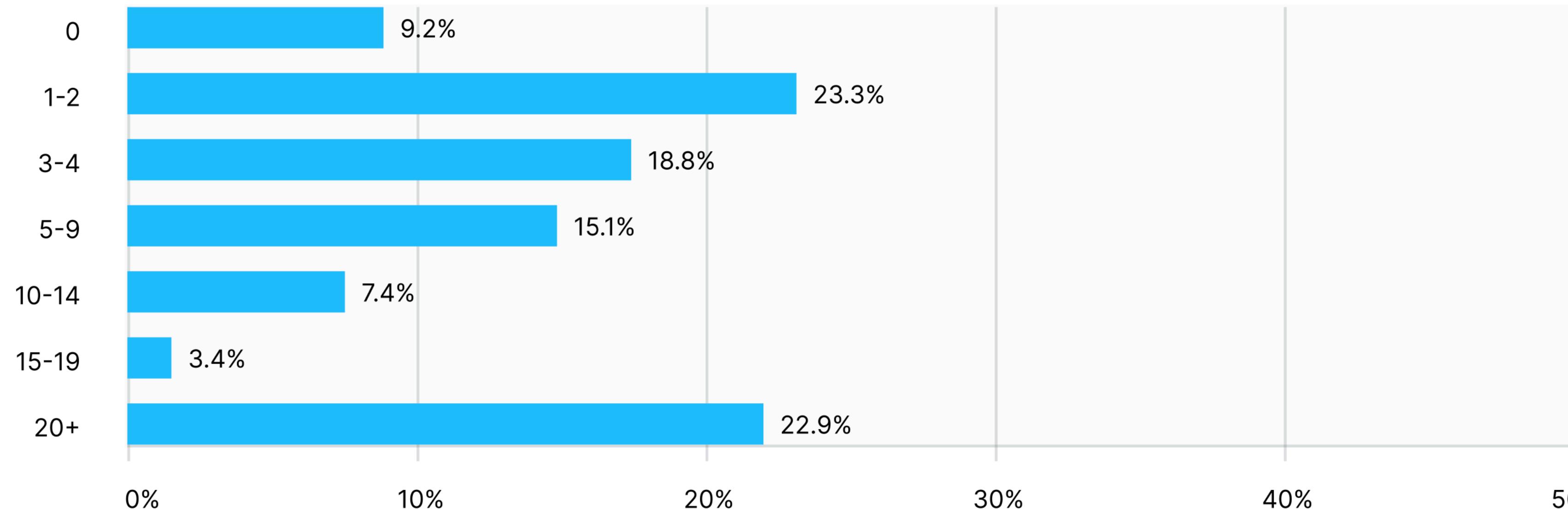
- Salary remained strong, pandemic may change that
- Concerns that supply + demand for DS may be narrowing at entry level
- WFH becomes normal, people move out of hot cities
- Current hot industries: Health, Supply chain

# Kaggle 2020 State of ML & DS

#### COMPANY SIZE (# OF EMPLOYEES)



#### DATA SCIENCE TEAMS (# OF EMPLOYEES)

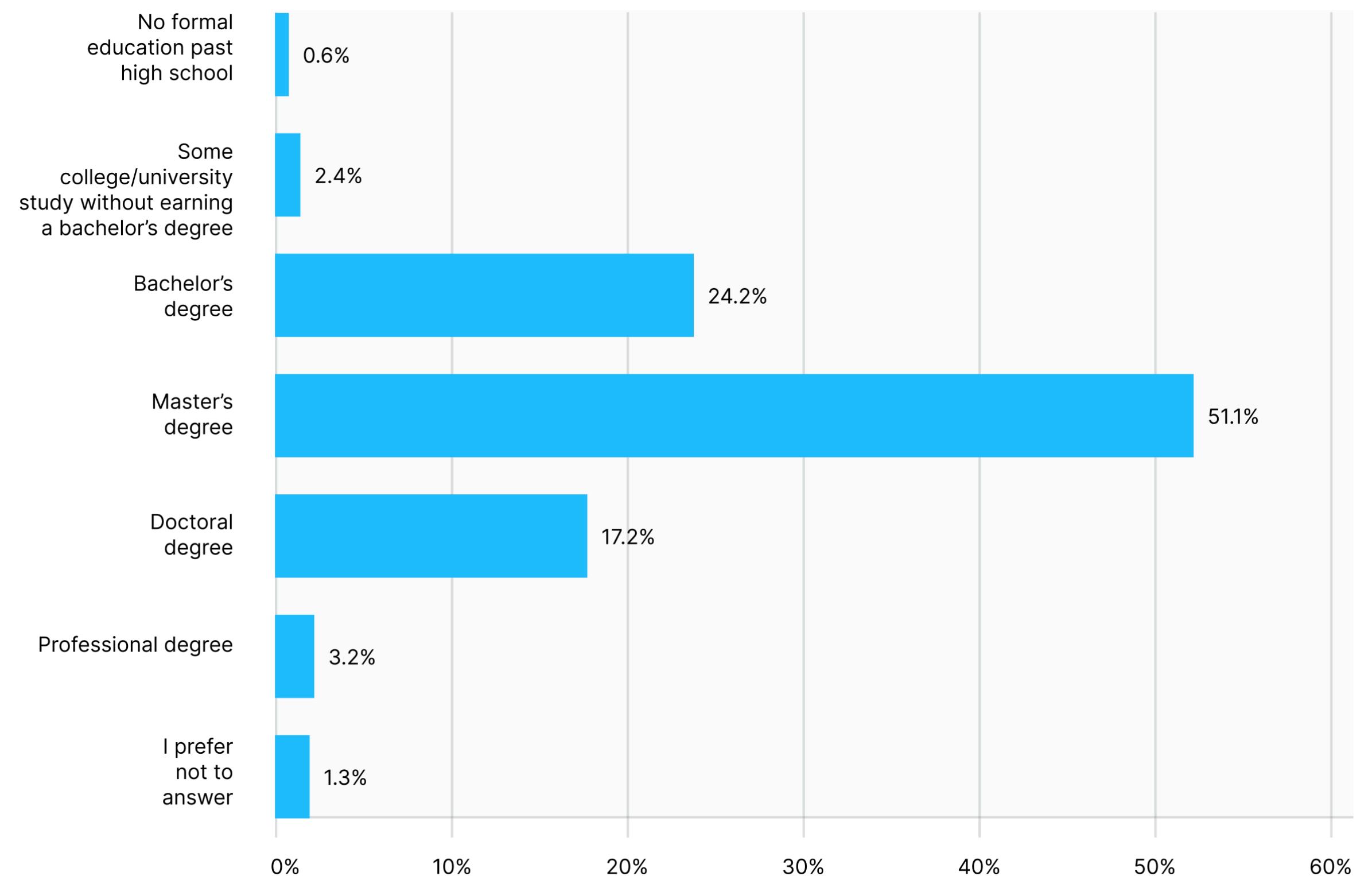


# Ongoing Learning

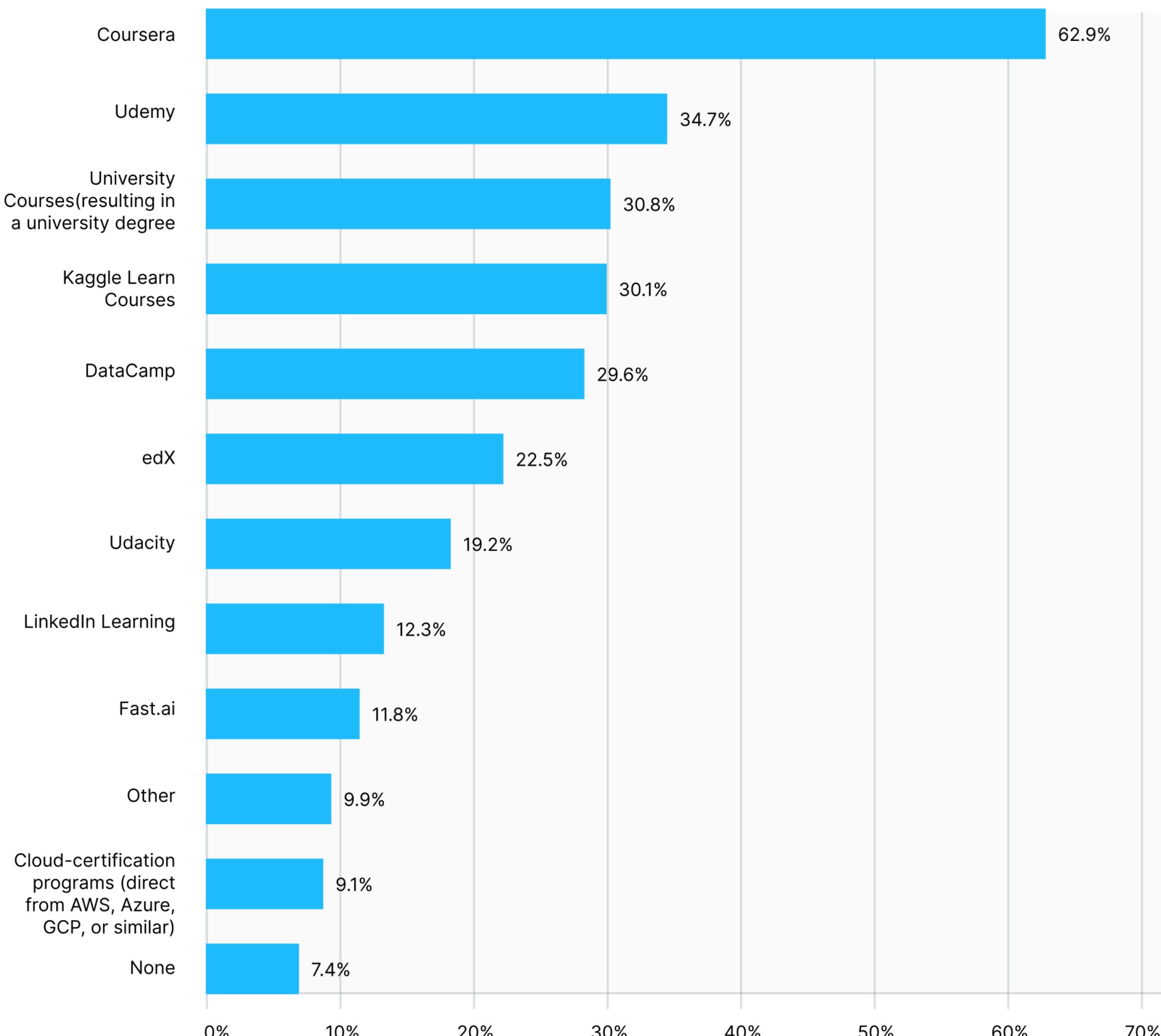
Data science and machine learning are quickly changing, so it's no surprise over 90% of Kaggle data scientists maintain ongoing education. While about 30% take traditional higher education courses, many more learn through online materials.

Coursera, Udemy, and Kaggle Learn top the most common mediums in our survey. Unsurprisingly, many Kaggle data scientists chose multiple resources in the survey, with an average of 2.8 mediums selected.

EDUCATION LEVEL OF KAGGLE DATA SCIENTISTS

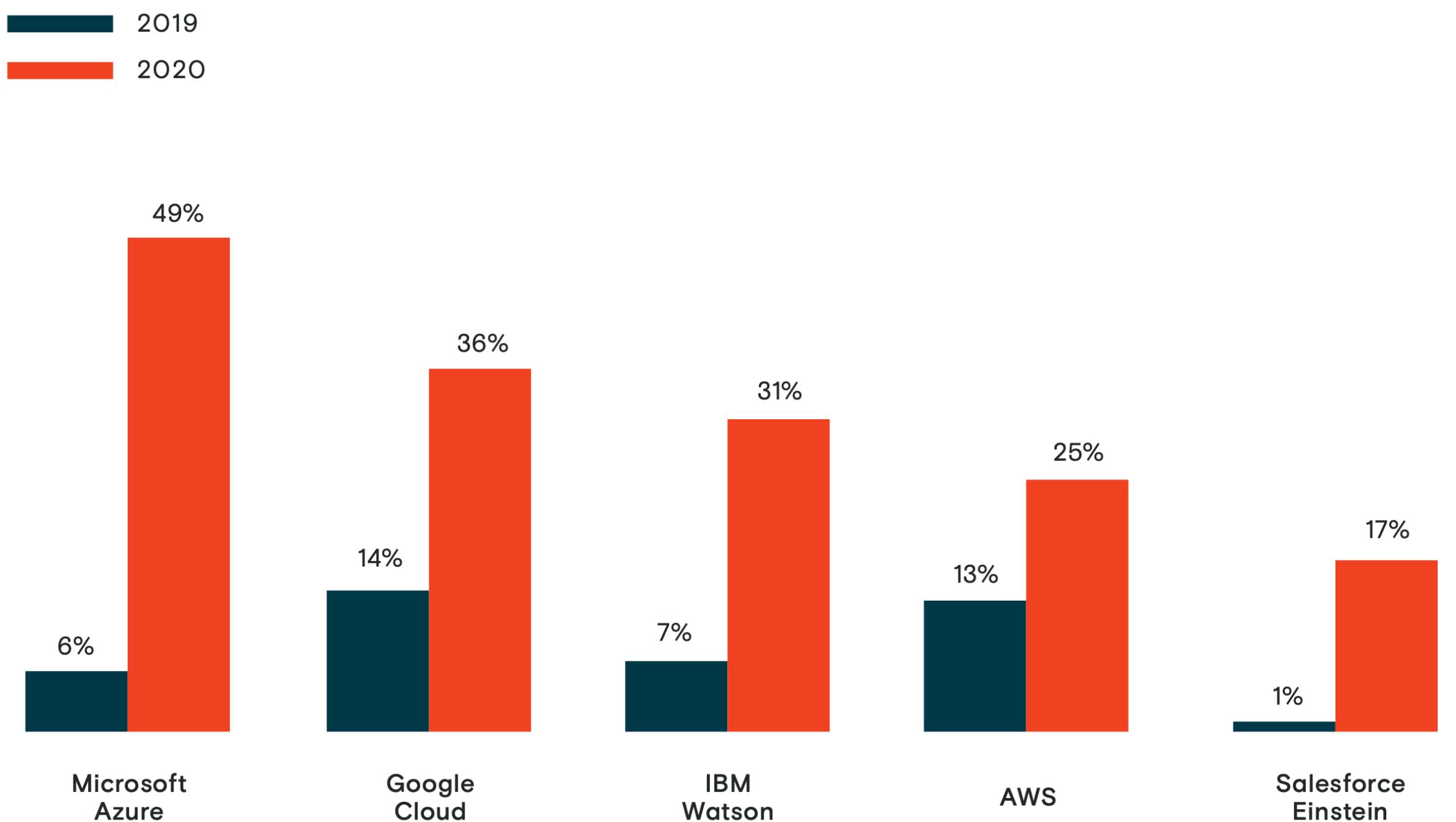


POPULAR ONGOING LEARNING RESOURCES



Appen (aka Figure-Eight aka  
Crowdflower) State of AI report

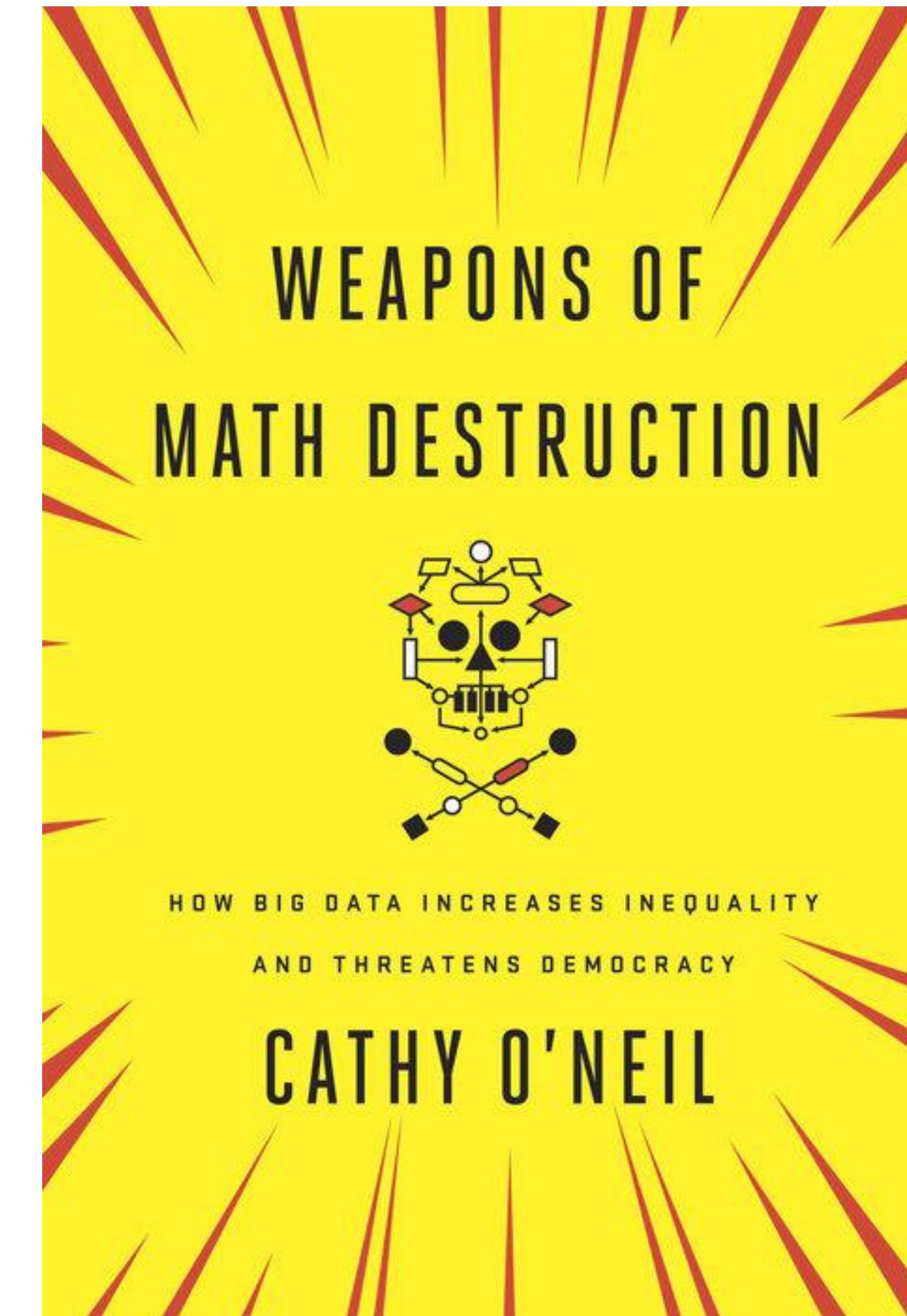
• Figure 7: What data science and machine learning tools/frameworks do you use?



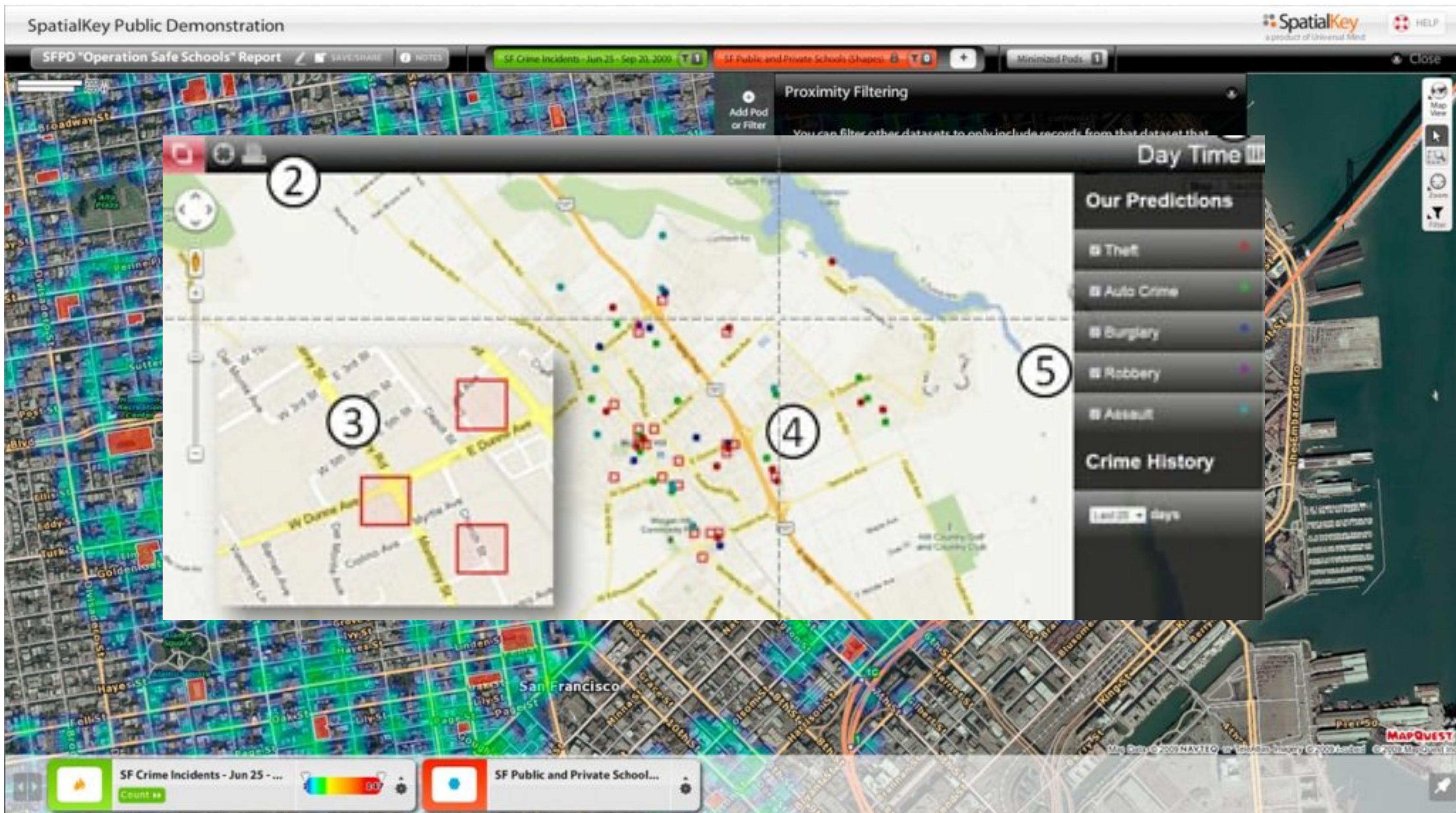
- AI/ML is a big deal at corporate level
- Everything is moving to the cloud (accelerating from WFH?)
- Lots of love for Azure, slowing growth for AWS

# Don't be a tool for creating WMDs

- Algorithms (and DS!) implement our biases, yet look objective
- Can implement our biases at scale
- Can have huge impacts on people's lives
- Are not transparent or accountable to the people being impacted

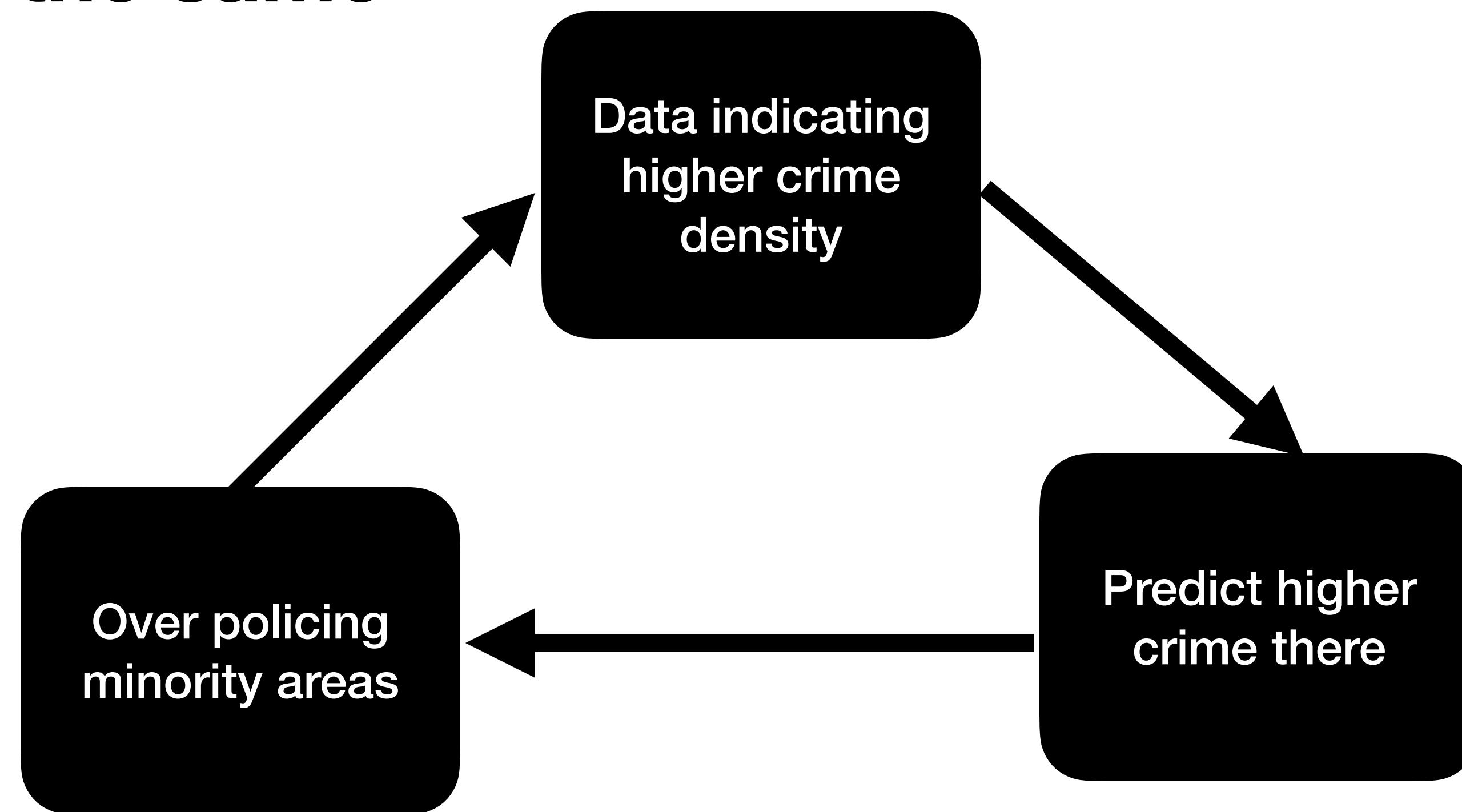


# Predictive policing & sentencing



# Predictive policing & sentencing

**Blacks arrested for possession at 4x the rate of whites  
Usage rates the same**





**“A lot of times, people are talking about bias in the sense of equalizing performance across groups. They’re not thinking about the underlying foundation, whether a task should exist in the first place, who creates it, who will deploy it on which population, who owns the data, and how is it used?”**

**-Timnit Gebru**

# **Thank you!**

## **Teaching Assistants:**

**Tyler Chang**

**Mia Lucio**

**Pooja Pathak**

**Areeb Syed**

## **Instructional Assistants:**

**Scott Yang**

**Harrison Ma**

**Richard Duong**

**Viki Zhao**

**And thanks to YOU!**