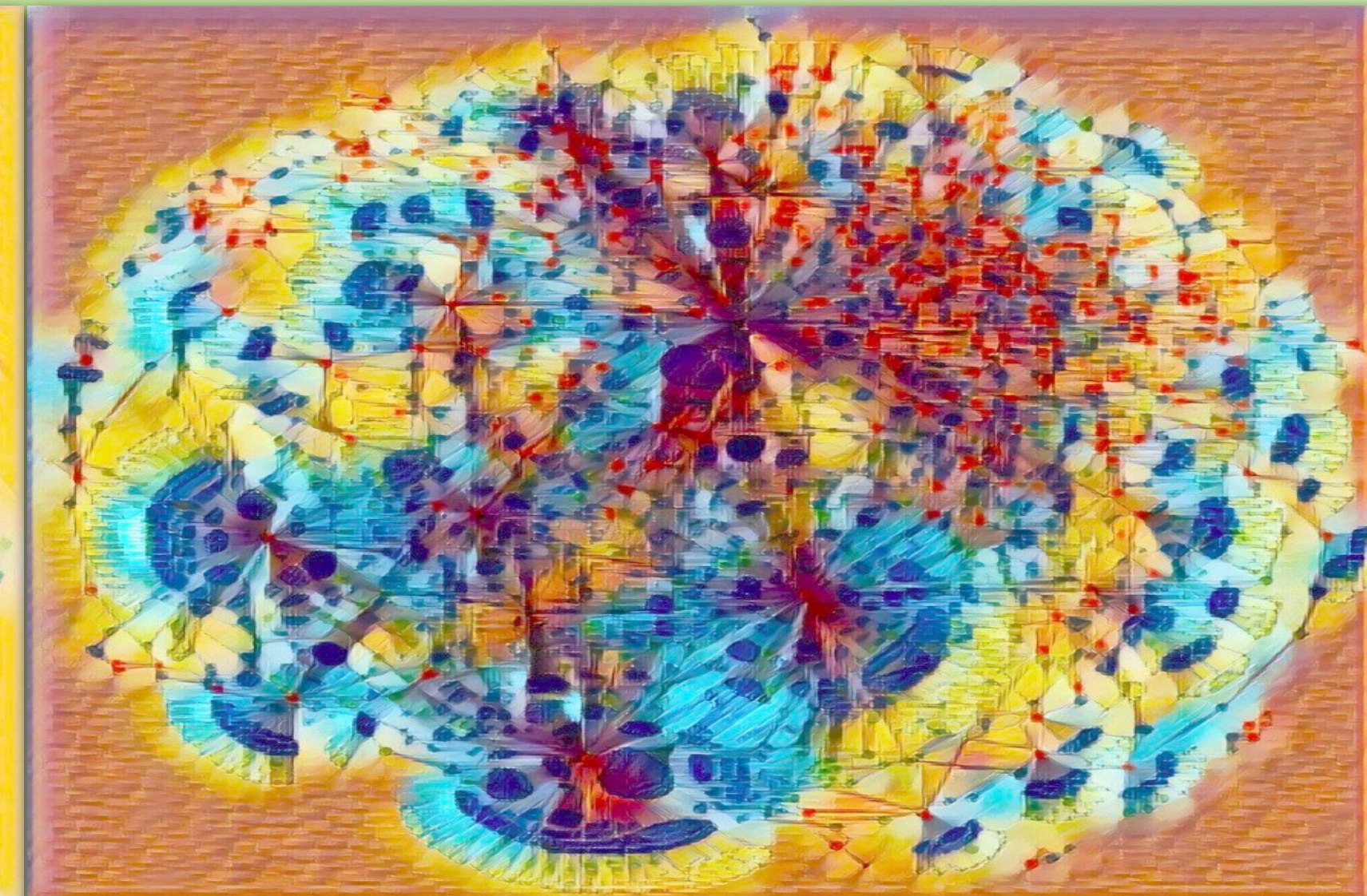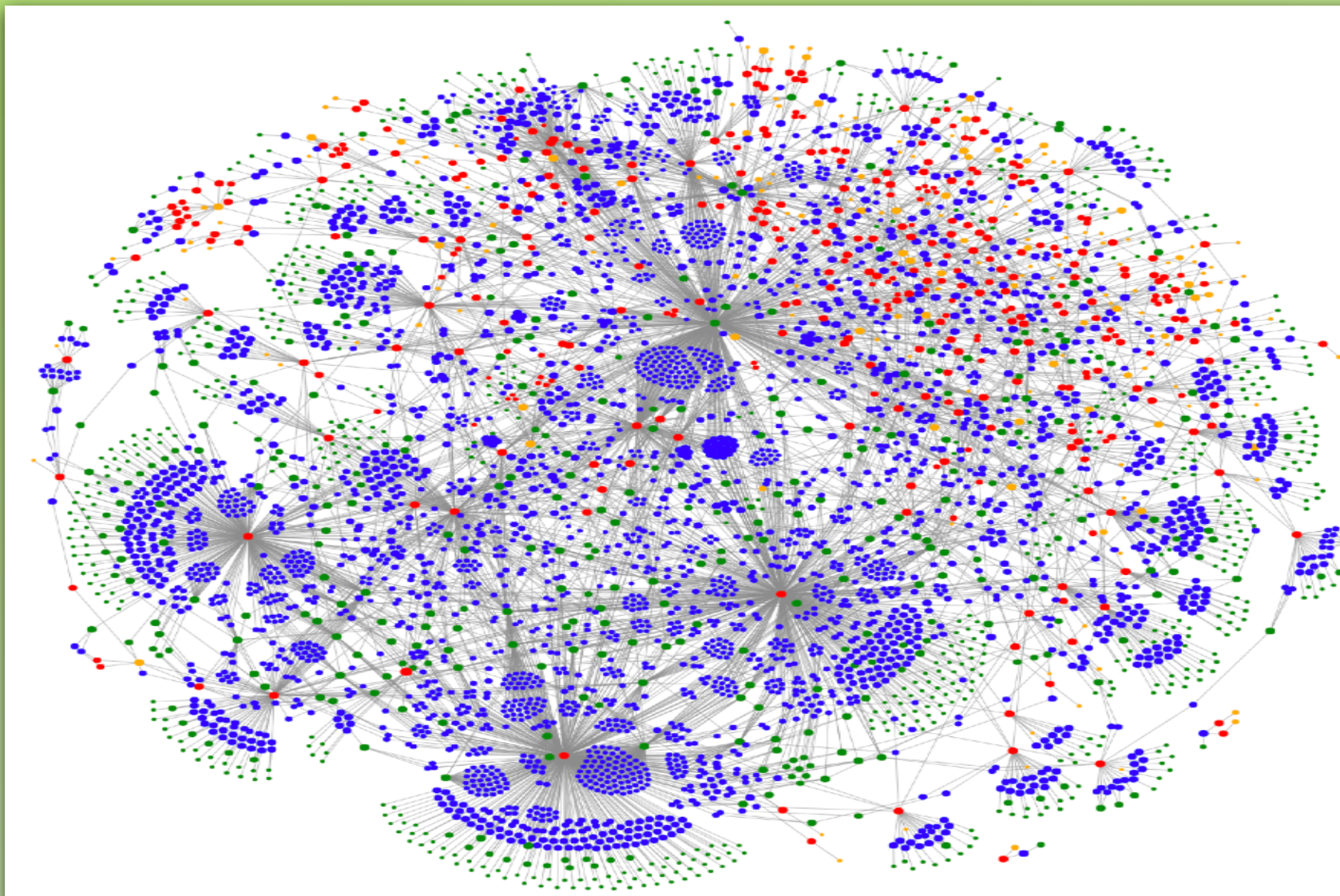# Rich Context:
## support for cross-agency data stewardship, measuring dataset impact on public policy

**Paco Nathan** **@pacoid**
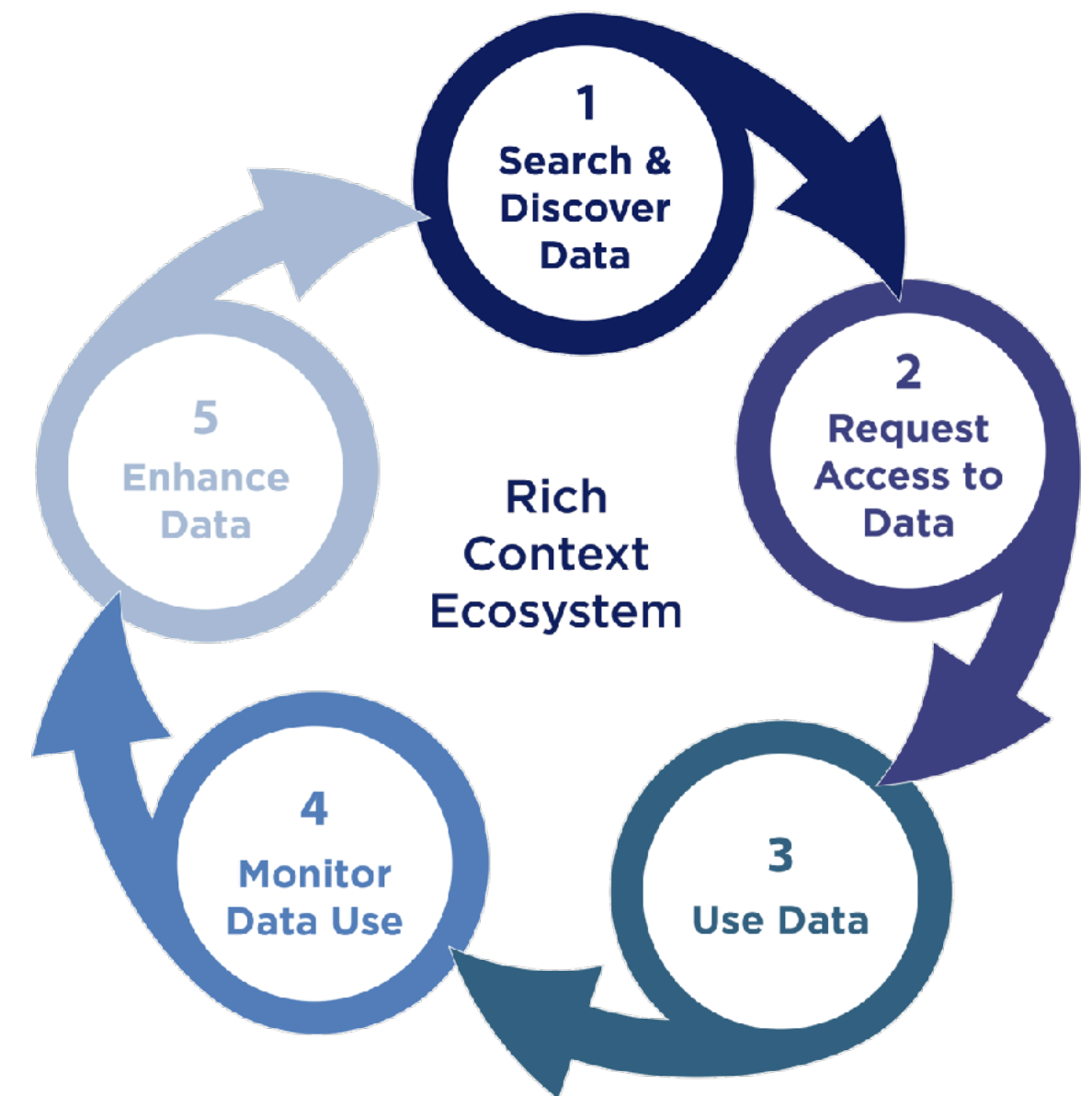
**derwen.ai**

**part 1:
public private partnership –
a case study**

# Administrative Data Research Facility

## Coleridge Initiative
**Julia Lane**, et al.  NYU Wagner

- FedRAMP-compliant **ADRF framework** on AWS GovCloud: "public agency capacity to accelerate the effective use of new datasets"

- for research projects using cross-agency sensitive data, in US and EU – **now in use by 30+ agencies**

- cited as the first federal example of Secure Access to Confidential Data in the final report of the Commission on **Evidence-Based Policymaking**

- augments **Data Stewardship** practices; collaboration with Project Jupyter on the related **data gov features**
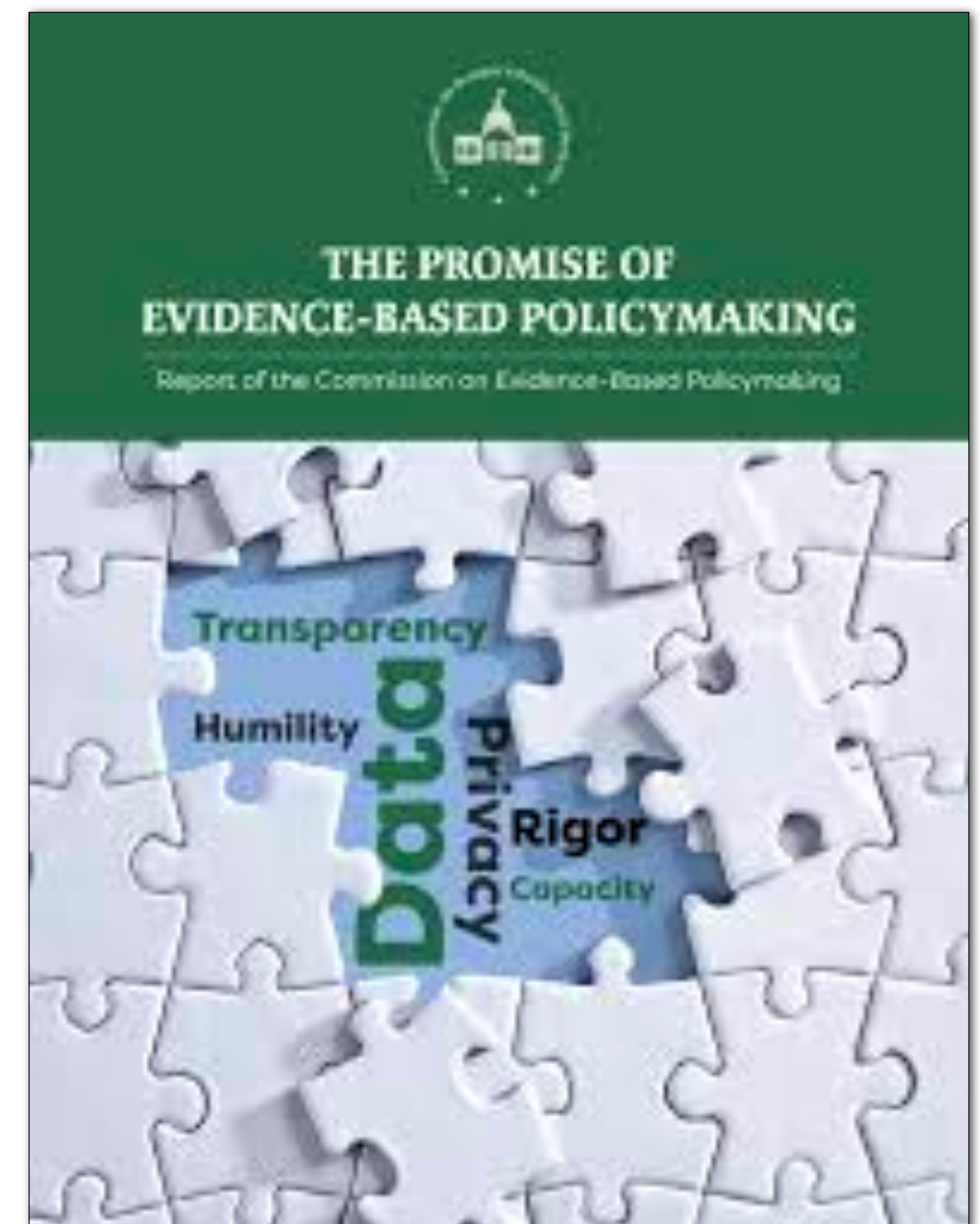
# Evidence-Based Policymaking

- **Foundations for Evidence-based Policymaking Act** (2018)

- **Information Quality Act** (2001)

- **NIH Strategic Plan for Data Science** (2018)

- **US federal data strategy**

- **Year-1 Action Plan** (2019)

See also:
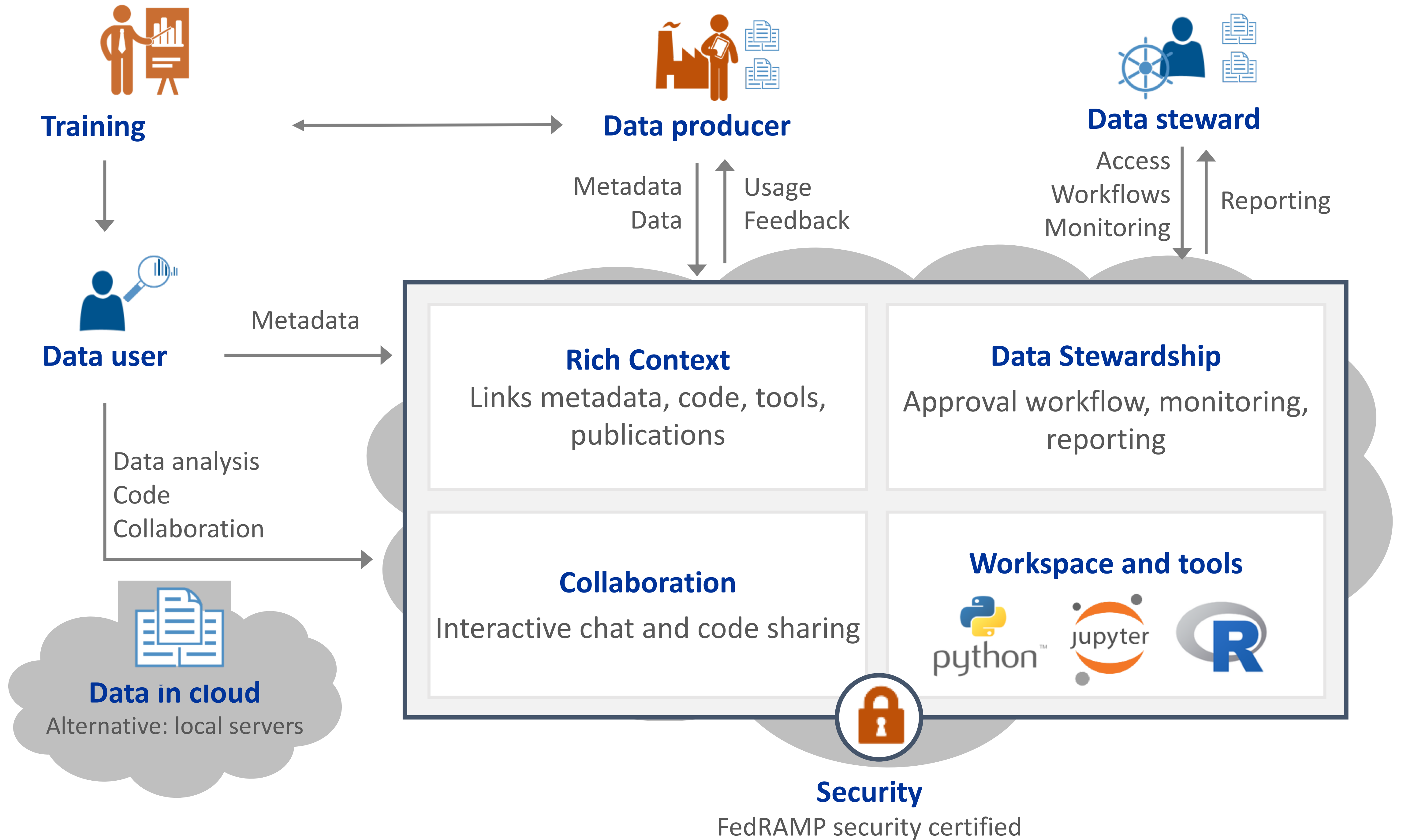**"Evidence-based decision making: What DOE, USDA and others are learning"**
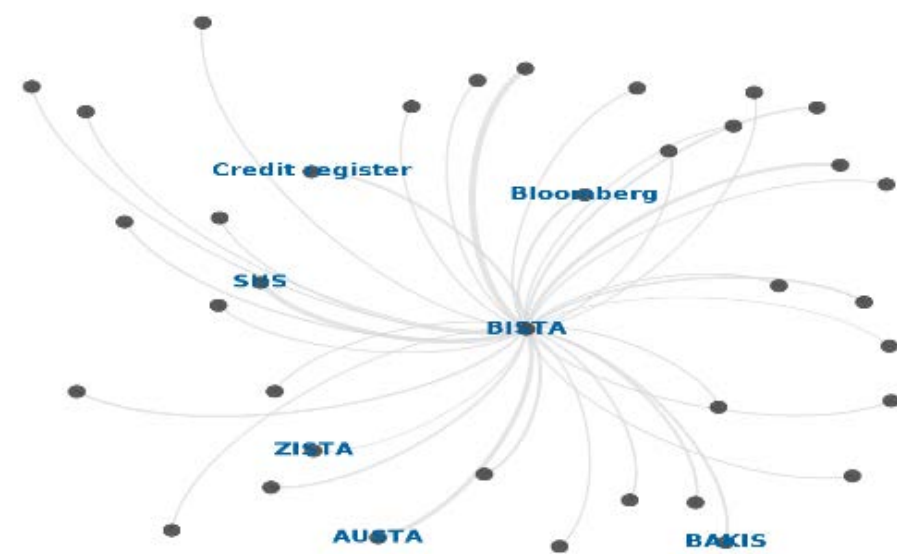Wyatt Kash
*FedScoop* (2019-06-28)



THE PROMISE OF
EVIDENCE-BASED POLICYMAKING
Report of the Commission on Evidence-Based Policymaking

Transparency
Humility
Data
Privacy
Rigor
Capacity

**Training**

**Data producer**

**Data steward**

Metadata
Data

Usage
Feedback

Access
Workflows
Monitoring

Reporting

**Data user**

Metadata

Data analysis
Code
Collaboration

**Rich Context**
Links metadata, code, tools, publications

**Data Stewardship**
Approval workflow, monitoring, reporting

**Collaboration**
Interactive chat and code sharing

**Workspace and tools**

python™        jupyter        R

**Data in cloud**
Alternative: local servers

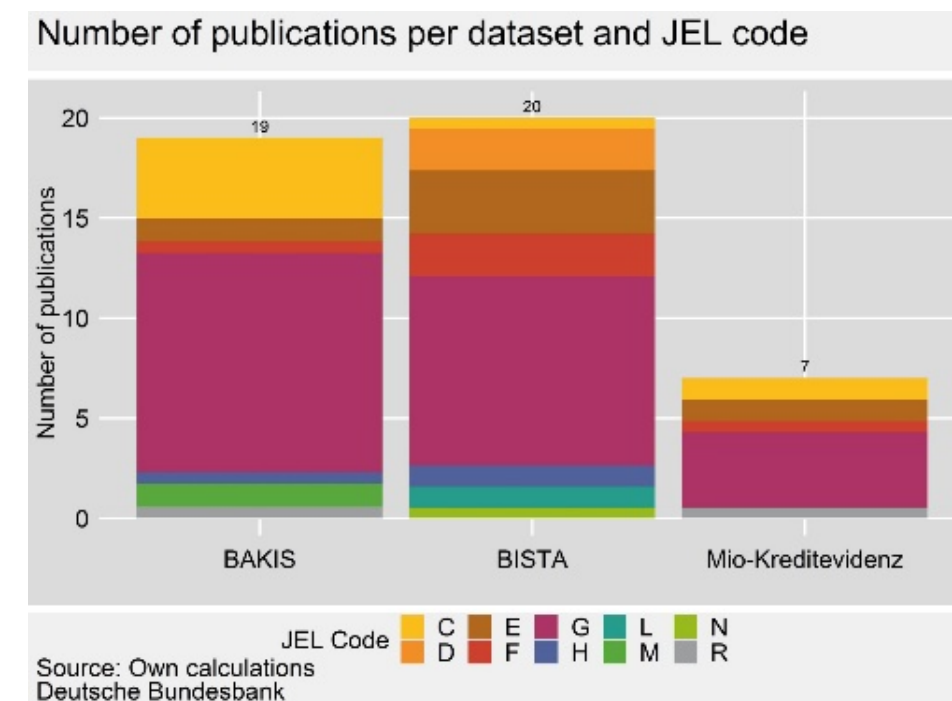**Security**
FedRAMP security certified

# Research Data Centers

To date, we successfully launched a machine learning competition and obtained algorithms to extract used datasets from publications. Based on the results, we implemented two prototypes that will serve as building blocks for a unified system.



**Recommend data** to researchers ("*based on your interest, you might also like this data*").

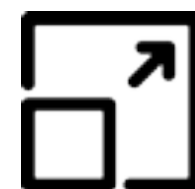**Data impact factor** ("*This dataset generates most research/ societal / policy value*")

The two prototypes already create value by enabling more optimal data usage and by supporting effective resource allocation into valuable data that generates results.



*– Hendrik Doll, Deutsche Bundesbank*

**The bigger picture**

Our vision is a system, where researchers and analysts apply for confidential data access, get data recommendations, have a secure remote digital workspace, and provide feedback on data. Such a system relies on three non-traditional data sources: Data usage information, implicit knowledge in researchers' heads (incentives for sharing), and structured administer metadata (*annodata*) to automatically govern data access.

# Research Data Centers

**Factsheet: Building an integrated data access system for empirical research**

Hendrik Doll, Stefan Bender, Jannick Blaschke, Christian Hirsch, Christian Resch[1]

*Federal Reserve Board, Washington, D.C., October 1-2, 2019*
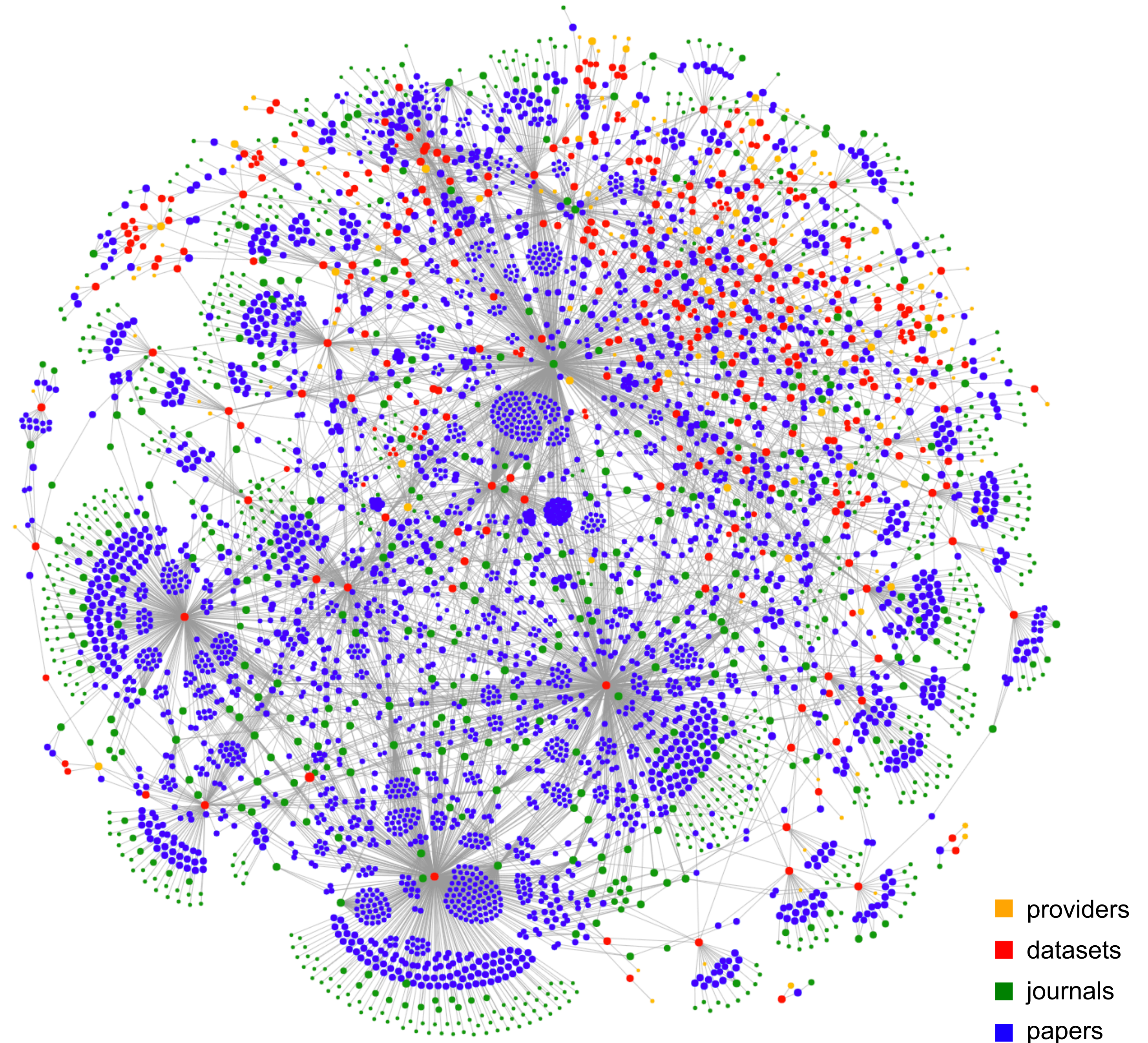
**Our idea**

Empirical research and evidence-based policy decisions increasingly rely on microdata. Research publications are well referenced and findable, however structured information on data usage is less available. The current project aspires to change this by building a data-centric ecosystem with rich context and a community around microdata.

# Need: recommenders for researchers/analysts

- Objective: provide better means of *search and discovery* for social science researchers and agency analysts.

- Collect workflow telemetry and query logs to augment the graph.

- Currently developing recommender systems based on the graph.

- This accelerates research and also assists training (e.g., onboarding agency analysts).

- Near-term goal: identify people with specific expertise.

- Long-term goal: learn workflow configurations to support AutoML meta-learning.

# Rich Context

- Focus on *socioeconomic impact*

- Funded by Schmidt Futures, Sloan, Overdeck

- Partnering with Bundesbank, USDA, etc.

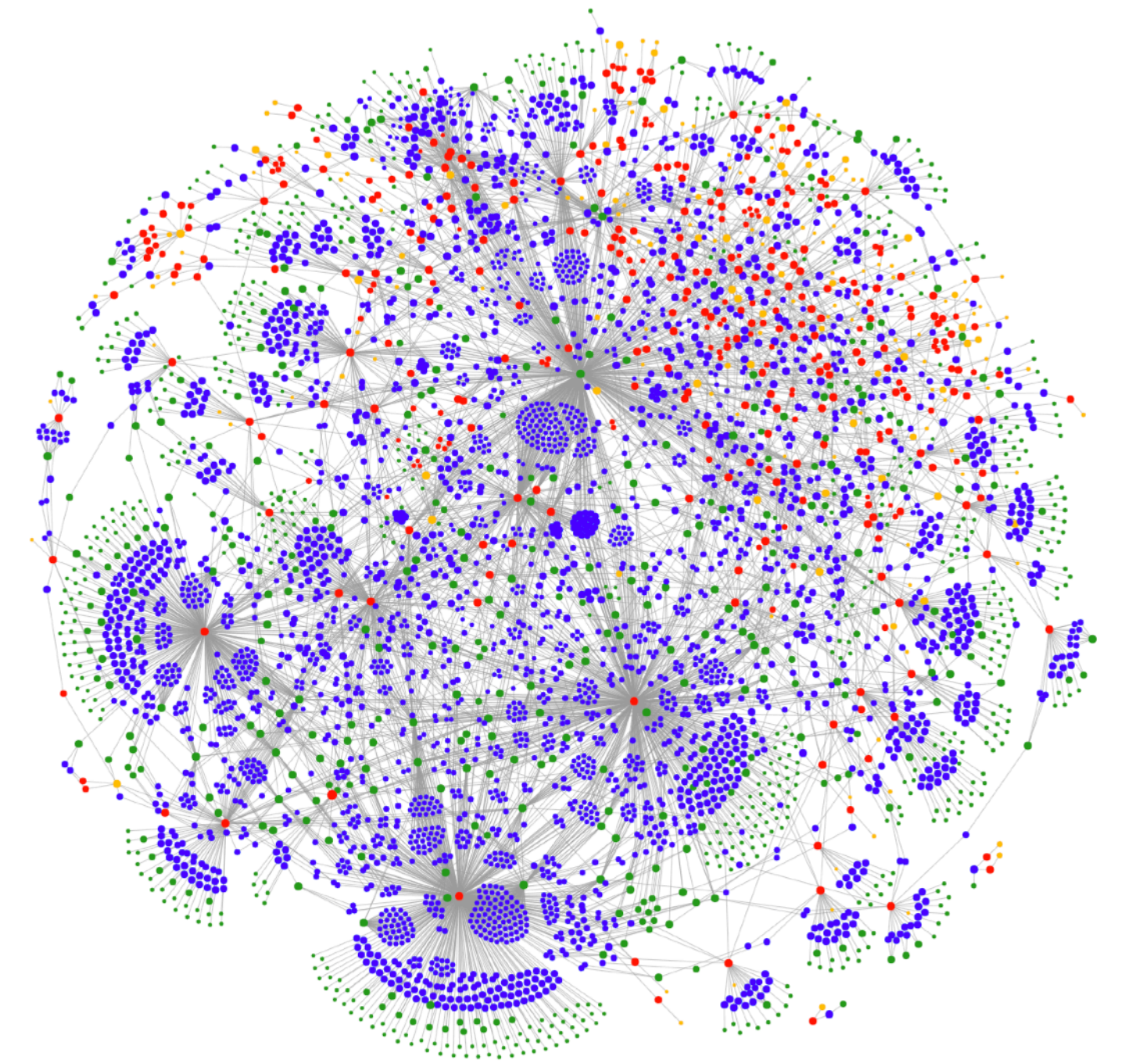- Collaboration with SAGE Pub, RePEc, ResearchGate, Digital Science, etc.



- providers
- datasets
- journals
- papers

# Rich Context

- Focus on *socioeconomic impact*

- Funded by Schmidt Futures, Sloan, Overdec...

- Partnering with... USDA, etc.

- Collaboration w... Pub, RePEc, R... Digital Science, etc.

Challenges that empirical researchers face: for a given dataset, find out **who** has worked with the data before, **what methods and code** were used, and **what results** were produced.

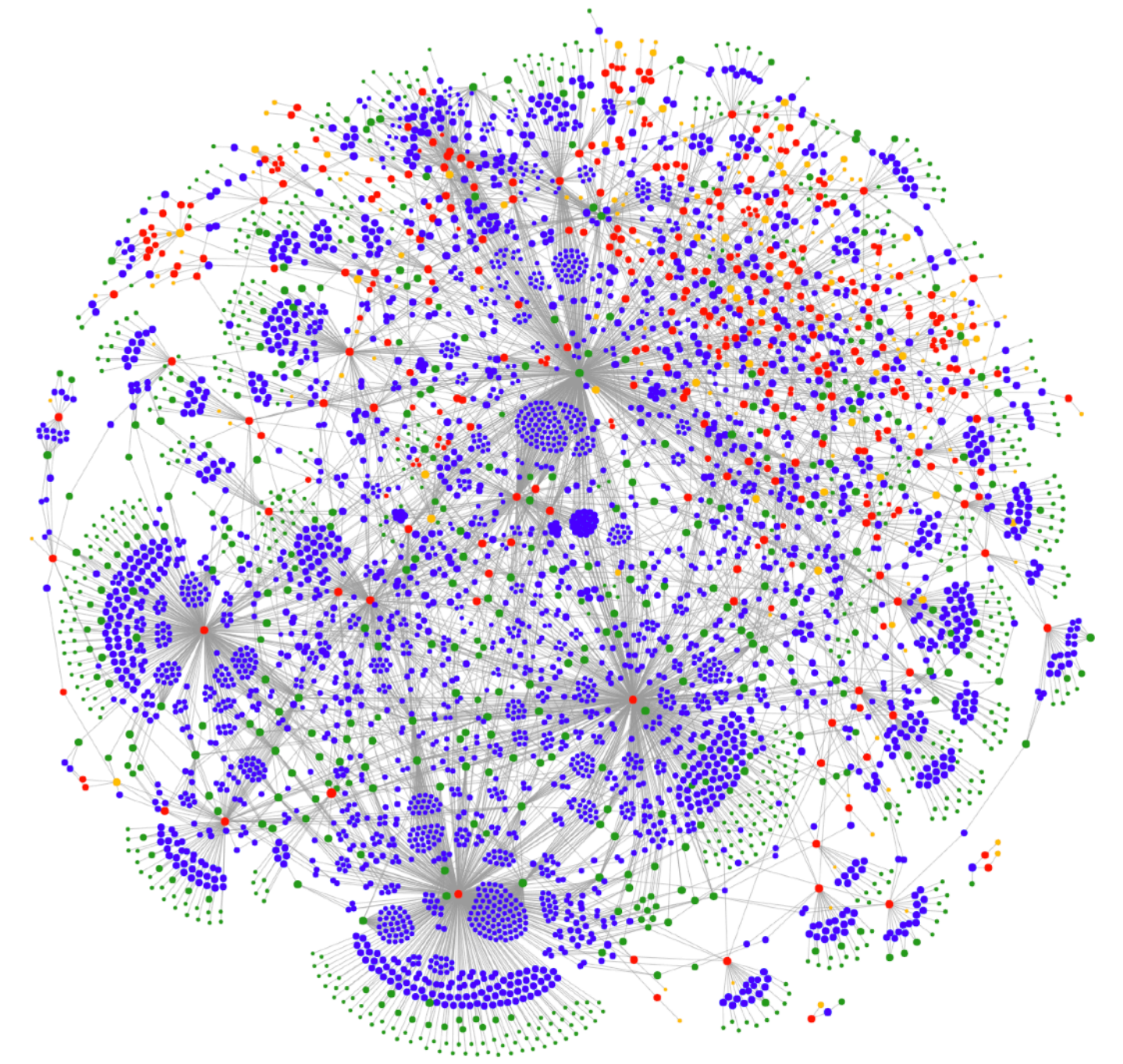providers

datasets

journals

papers

# Knowledge Graph – why?

- Allow flexibility for metadata representation

- Measure metadata quality

- Prepare features for ML models

- Build recommenders for *experts*, *topics*, *tools*, etc.

- Engage the public with automated data inventories

- Recommend configurations to new analysts

- Identify which datasets get used with others

- Quantify impact of datasets on policy

# Knowledge Graph – how?

- Manual data entry and curation of linked data

- Use persistent identifiers whenever possible:
  DOI, ISSN, ROR, ORCID, etc.

- Leverage ML models to infer missing metadata

- Federate queries of discovery services APIs

- Suggest corrections for metadata errors

- Use HITL to build feedback loops that engage experts,
  and provide convenient means for manual override

- Identify errors by using unit tests, ontology axioms,
  graph analytics, etc.

- **Collaborate with agency libraries!**

# KG process

- who are the expert people?
- which topics are emerging?
- how can methods be shared?

**activities** → **outputs** → **outcomes** → **impact**

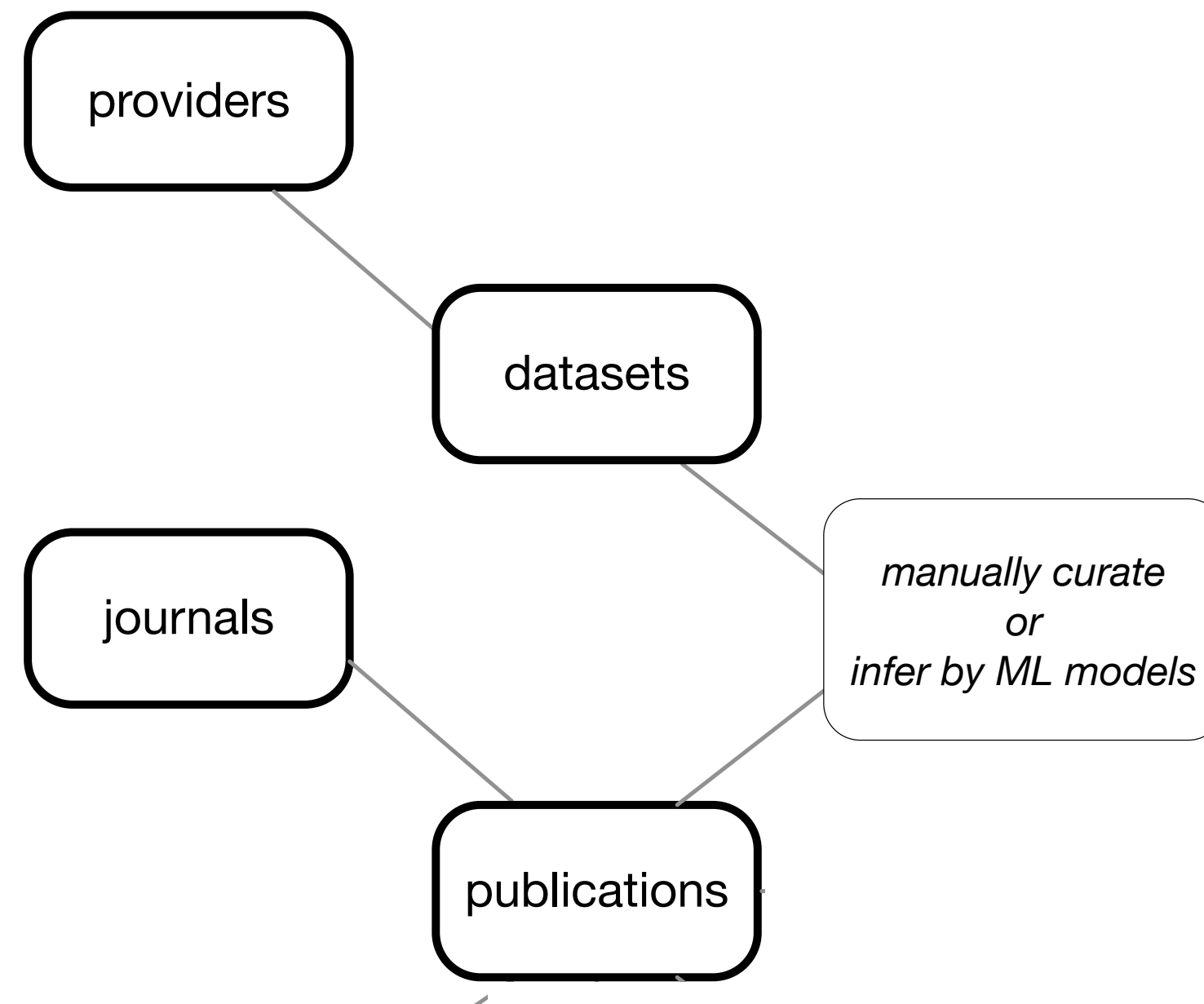curated datasets     research projects     published research     better science, government, education

*ML models infer new metadata links*
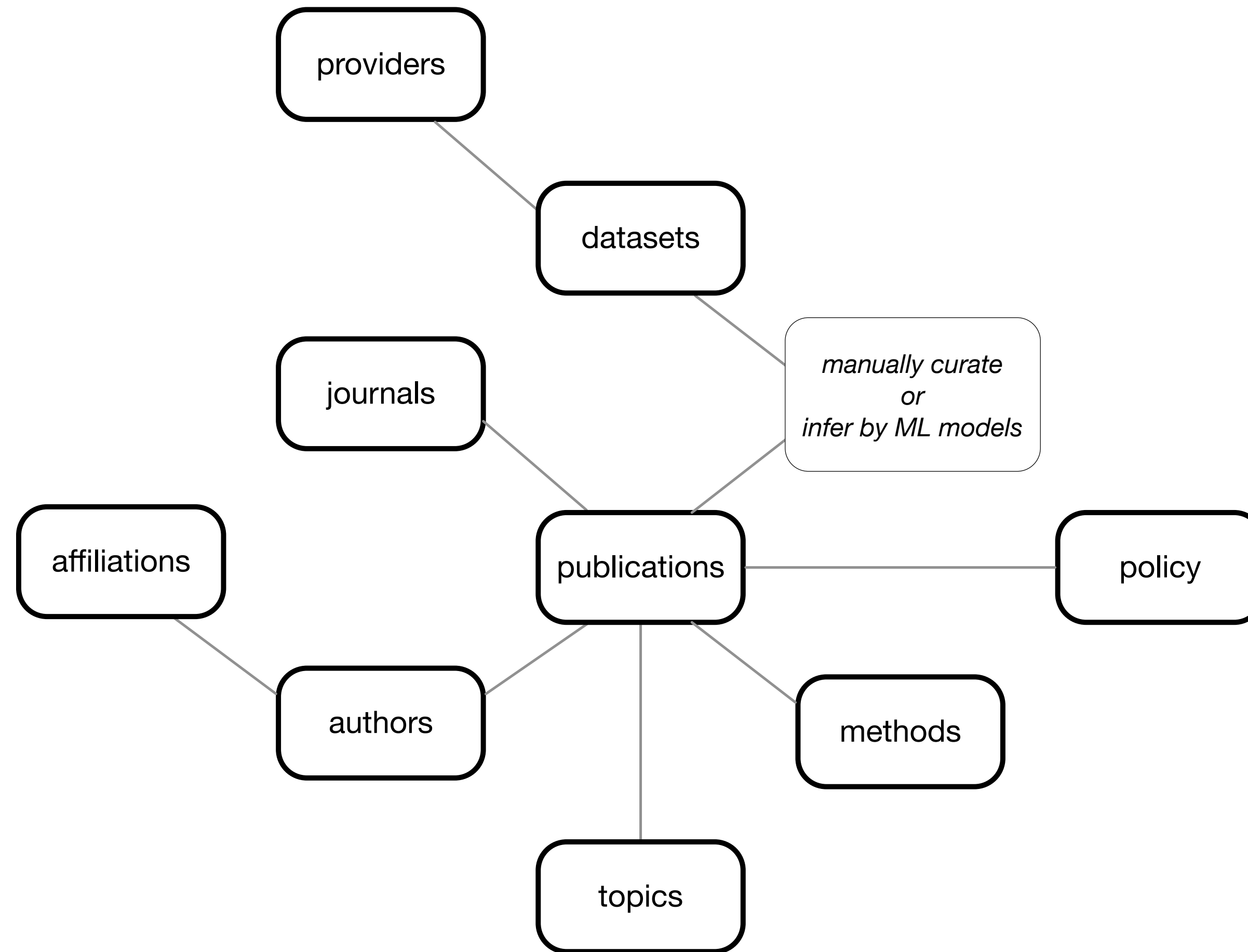
*how do we track the linkage??*
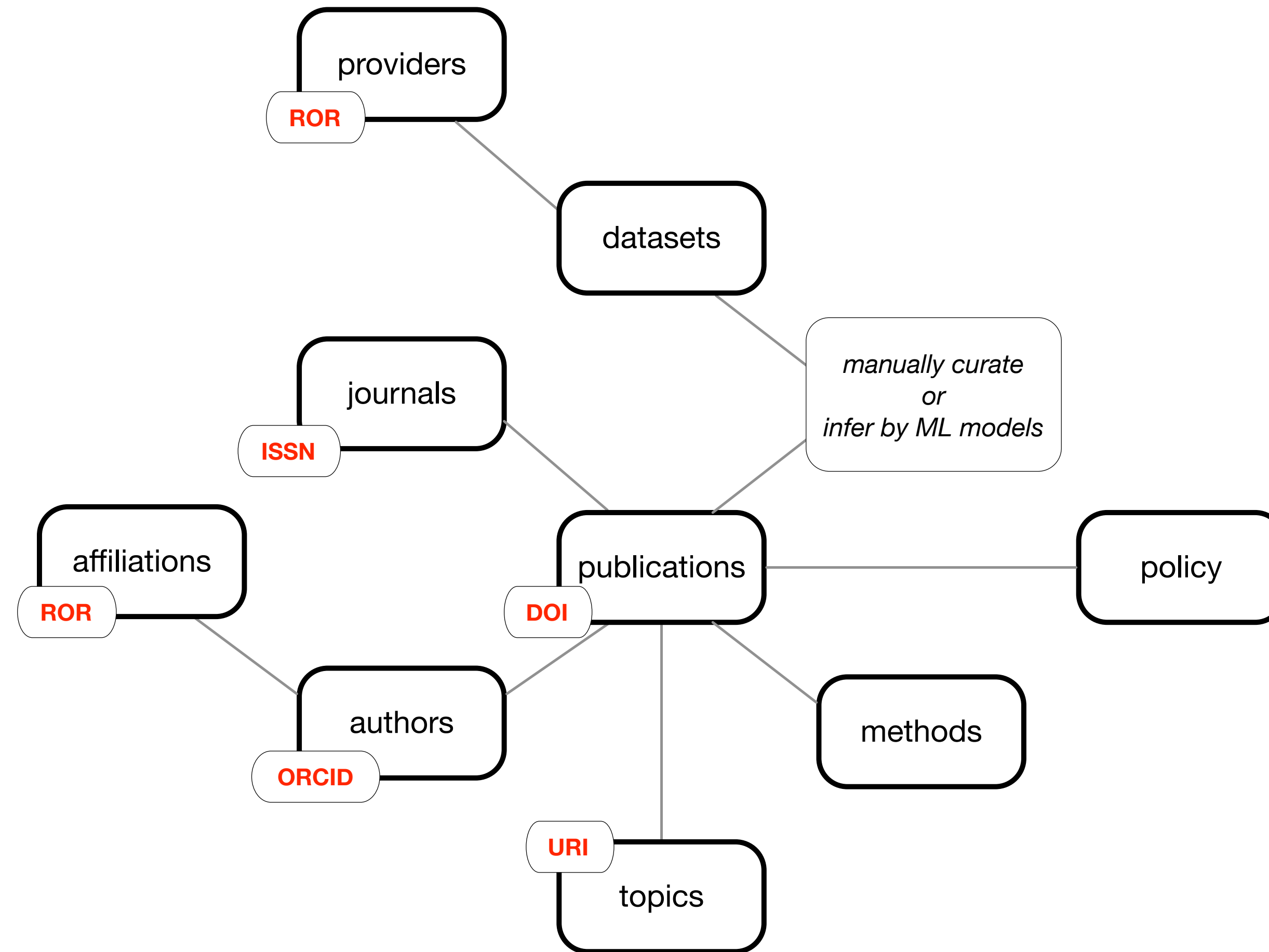
*how do we measure these behaviors??*

# KG construction and representation

providers

datasets

journals

publications

*manually curate*
*or*
*infer by ML models*

# KG construction and representation

# KG construction and representation



providers — ROR

datasets

*manually curate
or
infer by ML models*

journals — ISSN

affiliations — ROR

publications — DOI

policy

authors — ORCID

methods

topics — URI

See **github.com/Coleridge-Initiative/rclc/wiki/Corpus-Description**

# KG construction and representation



data catalog

providers
ROR

datasets

federated
queries

manually curate
or
infer by ML models

**Discovery
Services**
*Unpaywall
Dimensions
RePEc
ResearchGate
Crossref
DataCite
ORCID
OpenAIRE
PubMed
EuropePMC
Semantic Scholar
dissemin
Elsvier
SSRN
etc.*

scholarly
infra

journals
ISSN

affiliations
ROR

publications
DOI

policy

authors
ORCID

methods

metadata
updates

URI

topics

subject headings

*Library of Congress LCSH
PubMed MeSH
Wikidata + DBPedia
EuroVoc*

# KG construction and representation



data catalog

providers

ROR

datasets

*federated queries*

scholarly infra

journals

ISSN

projects

**Discovery Services**

*Unpaywall*
*Dimensions*
*RePEc*
*ResearchGate*
*Crossref*
*DataCite*
*ORCID*
*OpenAIRE*
*PubMed*
*EuropePMC*
*Semantic Scholar*
*dissemin*
*Elsvier*
*SSRN*
*etc.*

affiliations

ROR

publications

DOI

policy

authors

ORCID

methods

URI

topics

*metadata updates*

**HITL**
*author/expert feedback*

subject headings

*Library of Congress LCSH*
*PubMed MeSH*
*Wikidata + DBPedia*
*EuroVoc*

# KG construction and representation

# Open Source Projects

- **RCGraph** – Rich Context knowledge graph management
  **github.com/Coleridge-Initiative/RCGraph**

- **richcontext.scholapi** – federated discovery services and metadata exchange across scholarly infrastructure APIs
  **pypi.org/project/richcontext-scholapi**

- **adrf-onto** – controlled vocabulary for ADRF and Rich Context using OWL, SKOS, DCAT, PAV, CITO, FaBiO, etc.
  **github.com/Coleridge-Initiative/adrf-onto**

- **RCLC** – ML leaderboard competition
  **github.com/Coleridge-Initiative/rclc**

See also:

**"Machine Learning Highlights for Rich Context"**

# Funded additions to Project Jupyter

Make datasets and projects top-level constructs, support metadata exchange and privacy-preserving telemetry from notebook usage:

- JupyterLab **Commenting** and real-time collab similar to Google Docs

- JupyterLab **Data Explorer**: register datasets within research projects

- JupyterLab **Metadata Explorer**: browse metadata descriptions, get recommendations through knowledge graph inference (via extension)

- **Data Registry** (original proposal)

- **Telemetry** (privacy-preserving, reports usage)

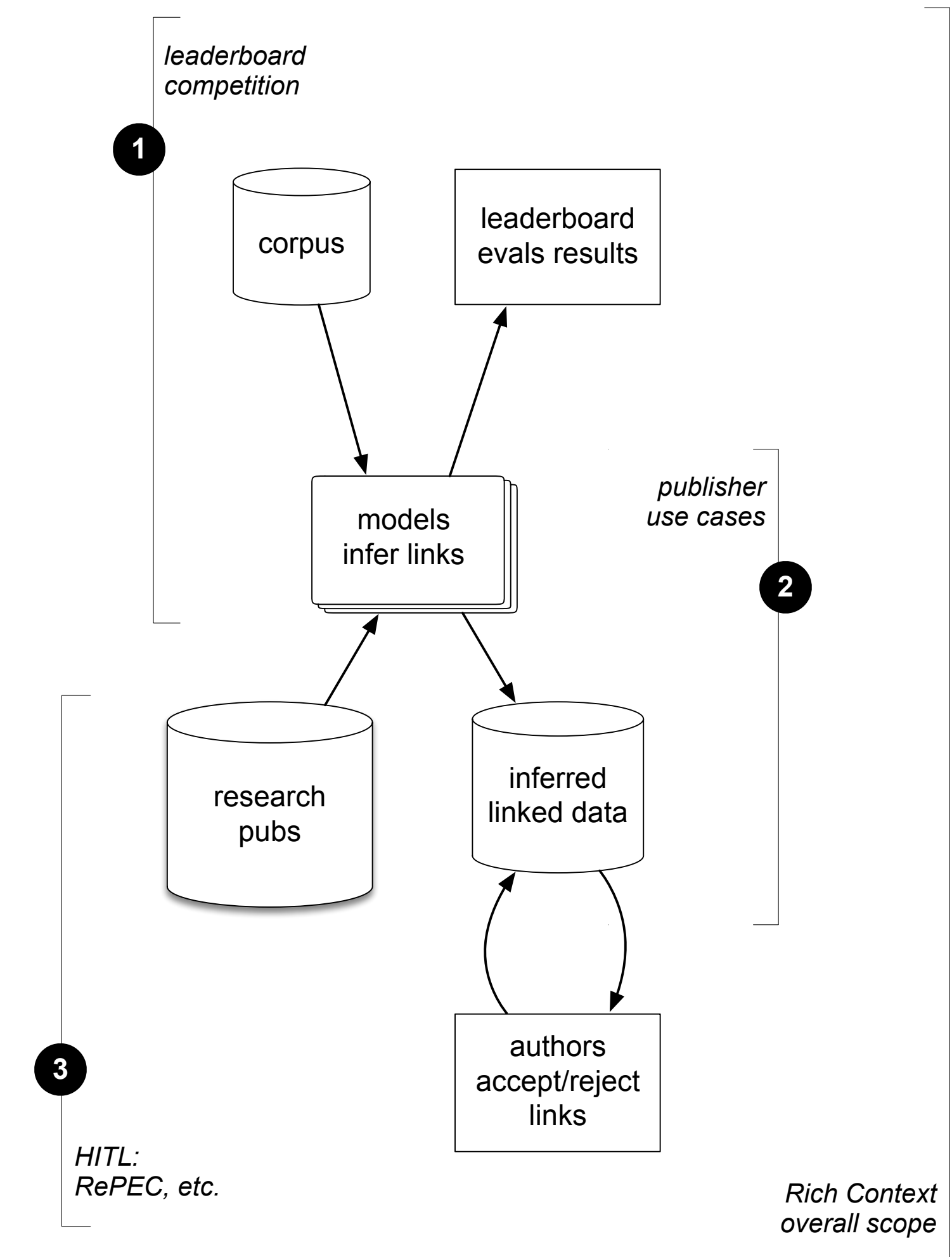# ML Leaderboard Competition

**github.com/Coleridge-Initiative/rclc**

- update from RCC competition in 2018-2019
- ongoing ML leaderboards (similar to **NLP-progress**)
- open source, hosted on GitHub
- highly curated test sets, all open-access publications
- teams collaborate via GH issues on corpus data quality, etc.
- focus on *precision* for ML model evaluation

**Current SOTA**

| source | precision | entry | code | paper | corpus | submitted | notes |
|---|---|---|---|---|---|---|---|
| LARC @philipskokoh | 0.7836 | ipynb | repo | RCC_1 | v0.1.5 | 2019-09-26 | RCLC baseline experiment using RCC_1 approach |
| KAIST @HaritzPuerto | 0.6319 | ipynb | repo | RCC_1 | v0.1.5 | 2019-11-01 | model trained a different dataset using DocumentQA and Ultra-Fine Entity Typing -- NB: this approach is able to identify new datasets |

# Human-in-the-loop

- semi-supervised learning, aka "human-in-the-loop" – in progress via RePEc

- interact with authors to confirm metadata inferred by ML models

- feedback from experts improves the corpus metadata and the ML modeling

See also:

**"Human-in-the-loop AI for scholarly infrastructure"**

**"New initiative to help with discovery of dataset use in scholarly work"**, Christian Zimmerman

# Collaboration with NOAA

- Initial focus on coastal inundation and community resilience, working with NOS

- Develop reusable dataset discovery services, so that the public and researchers can find trustworthy, high-impact data

- Identify experts who have used the data and the associated research topics, associated analytical methods and tools, and related datasets (e.g., Zillow, EPA, NASA, FEMA, etc.)

- Generalize for other federal agencies, such as USDA and NSF, as well as to international organizations, such as Deutsche Bundesbank

- Bring in AI expertise from industry and academia: KAIST, LARC, Recognai, DLA, Primer AI, GESIS, AllenAI, etc.

# Agency Benefits

- Understand more about NOAA's user community, to help outreach and get feedback, especially for researchers or commercial entities using NOAA data in novel ways.

- Help quantify the value and impact of data and research, especially in the context of NOAA's Blue Economy initiative.

- Relevance given the Federal Data Strategy and the OPEN Government Data Act: comprehensive data inventory, tools to help data users find trustworthy and relevant data.

- Explore a novel solution to dataset search.

# Additional Information

Rich Context @ NYU Coleridge Initiative
**coleridgeinitiative.org/richcontext**

- **white paper**

- **upcoming book** (Jan 2020)

- **feedback/propose collaboration**



**"Empty rhetoric over data sharing slows science"**
*Nature* (2017-06-12)

**"Experiences of the Deutsche Bundesbank"**
Stefan Bender
*CEMLA* (2019-05-28)

**"Where's Waldo: Finding datasets in empirical research publications"**
Julia Lane
*AKBC* (2019-05-22)

**"Google data set search"**
Ian Mulvany
*ScholCommsBlog* (2019-11-19)

**"Impact for social science researchers"**
Ian Mulvany
*FORCE11* (2019-11-17)

**part 2:**
**AI practices circa 2020 –**
**perspectives from industry**

# "Two Cultures" for AI

> **Mat Velloso**
> @matvelloso
>
> Difference between machine learning and AI:
>
> If it is written in Python, it's probably machine learning
>
> If it is written in PowerPoint, it's probably AI
>
> 5:25 PM - 22 Nov 2018
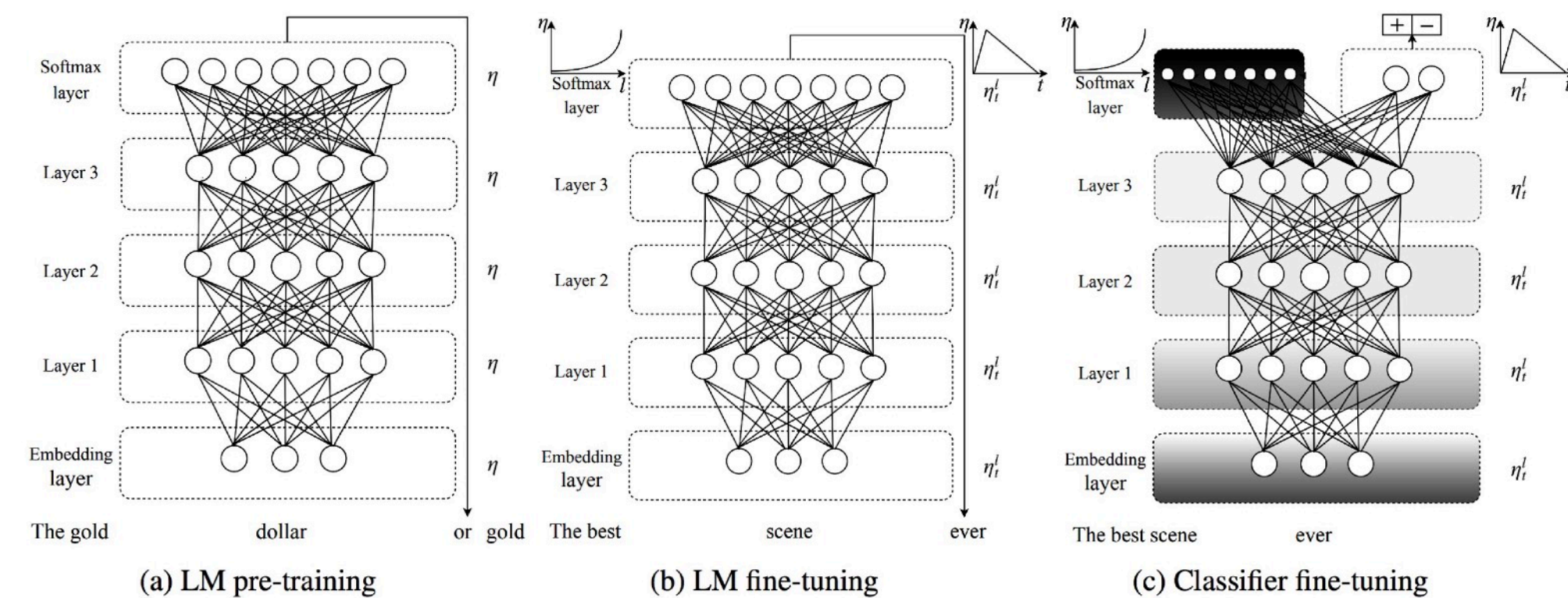
A more useful distinction:

- **ML** is about the tools and technologies
- **AI** is about use case impact on social systems

# So many ways to use optimized gradients...

AI finds pervasive use throughout industry: surpassing
human benchmarks – although it's probably better to think
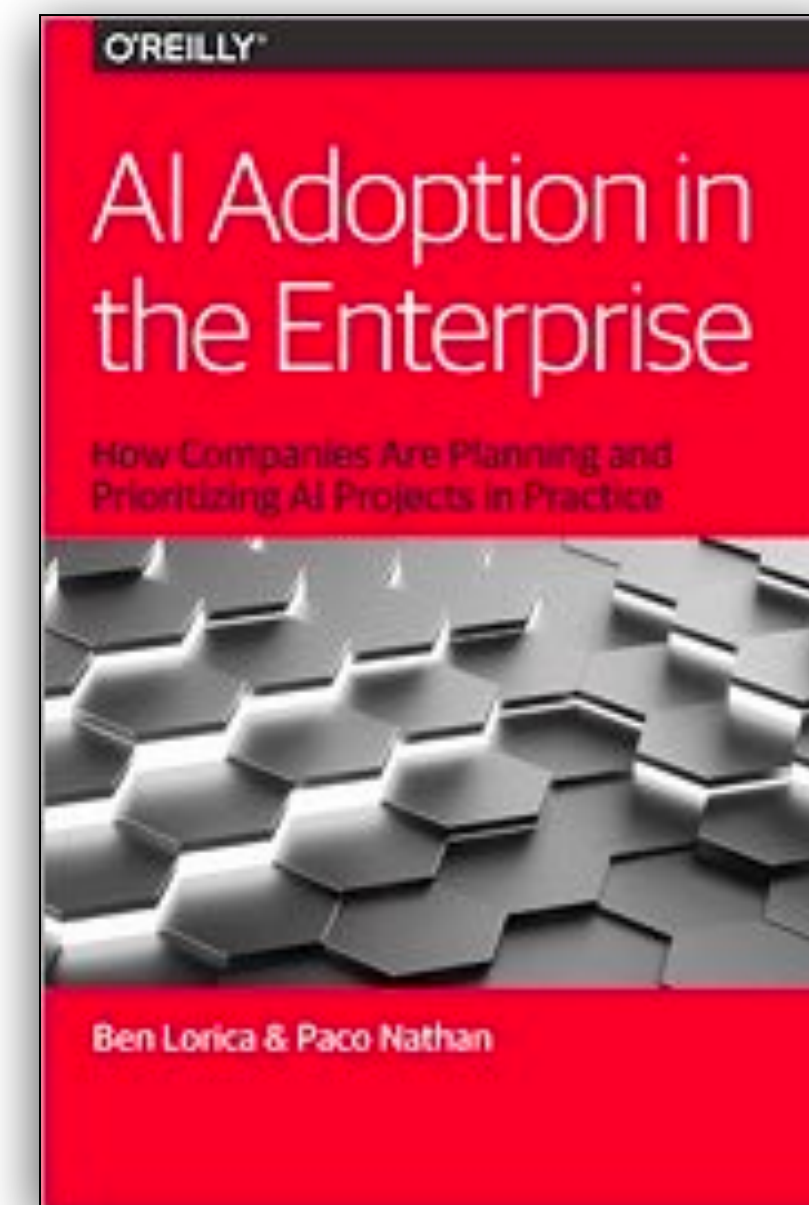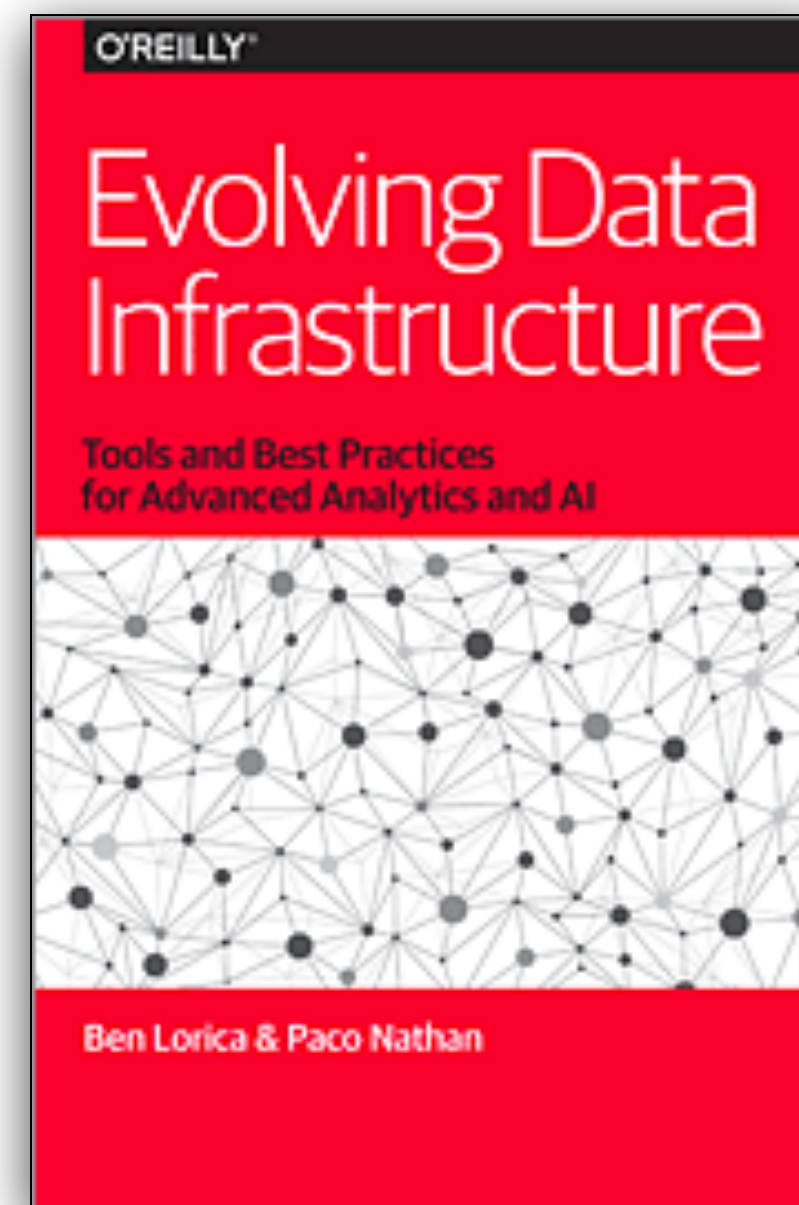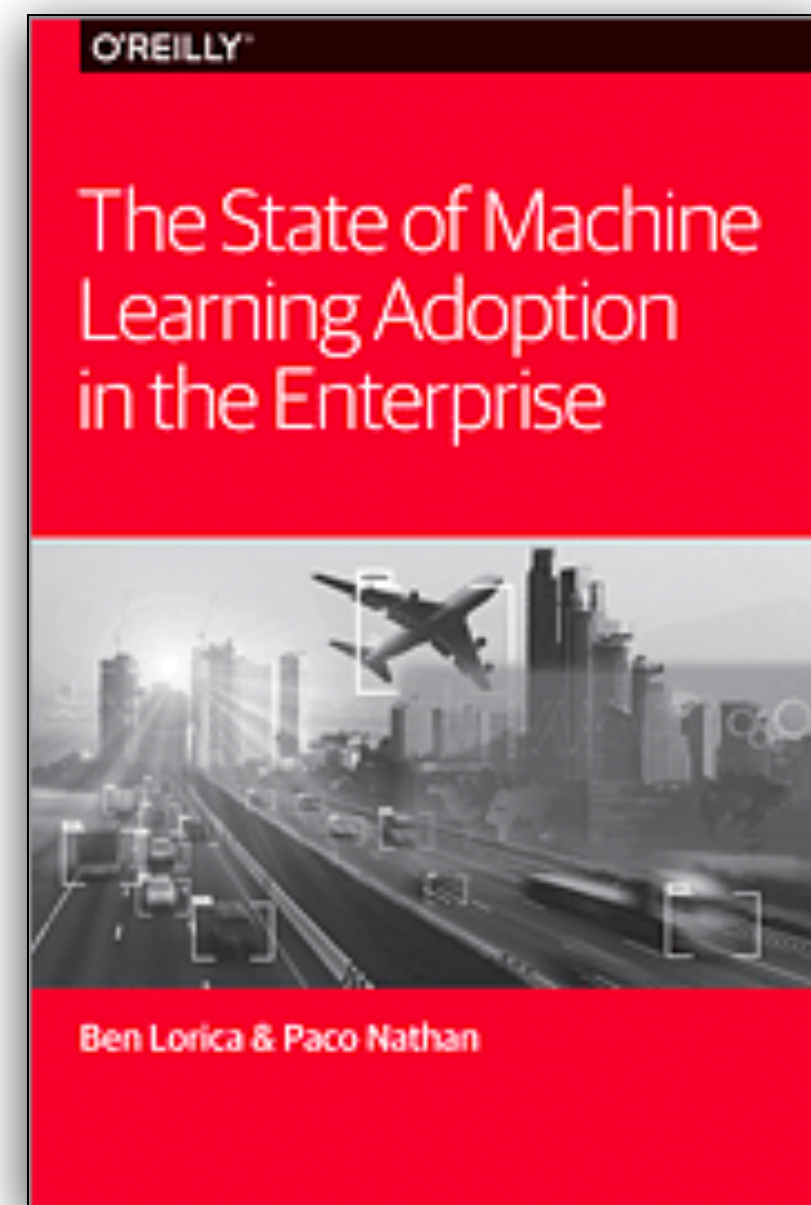in terms of **how to augment social systems**:

- *deep learning*; **Krizhevsky 2012**

- *reinforcement learning*; **Silver 2017**

- *transformers* (NLP); **Radford 2019**

- *weak supervision*; **Ratner 2017**

As well as: *self-supervision, knowledge graph,
transfer learning, active learning* ("human-in-the-loop"),
and so on…



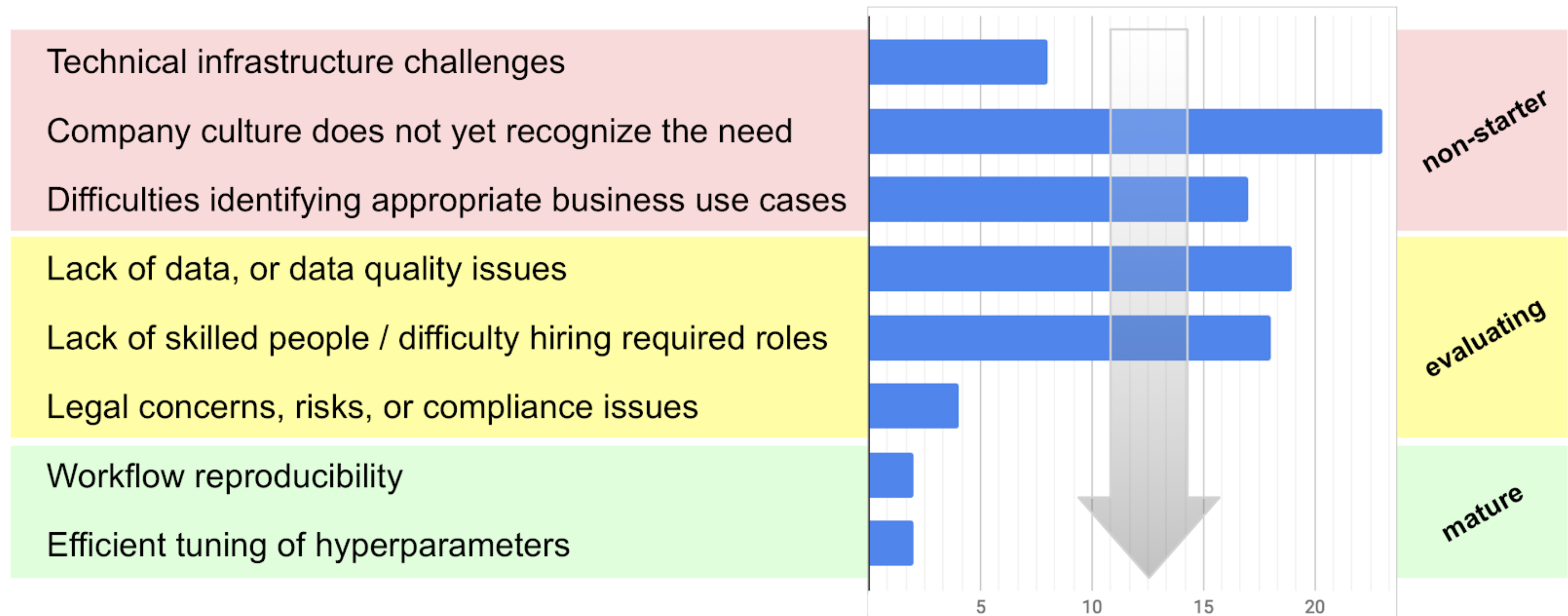(a) LM pre-training    (b) LM fine-tuning    (c) Classifier fine-tuning

# Industry surveys for AI and Cloud adoption

- **"Three surveys of AI adoption reveal key advice from more mature practices"** **Ben Lorica**, **Paco Nathan** *O'Reilly Media* (2019-02-20)

- **Episode 7**, *Domino*: surveying "ABC" adoption in enterprise (2019-03-03)
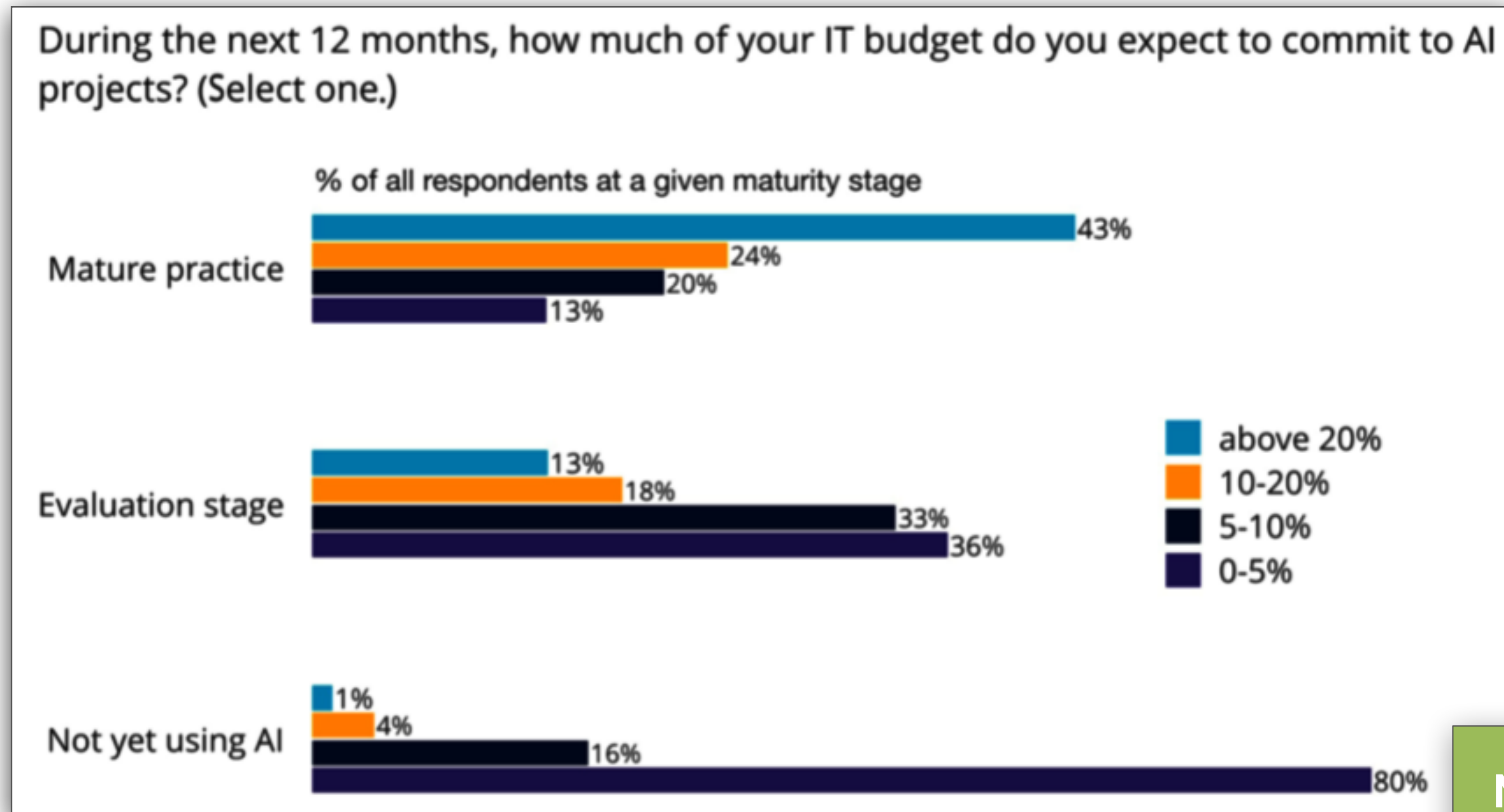
# Challenges in "ABC" adoption in enterprise

Leading risks for "ABC" (AI, Big Data, Cloud) adoption in enterprise, based on ORM surveys – viewed as a "survival analysis":

# An accelerating gap in AI funding

During the next 12 months, how much of your IT budget do you expect to commit to AI projects? (Select one.)

% of all respondents at a given maturity stage

**Mature practice**
- 43%
- 24%
- 20%
- 13%

**Evaluation stage**
- 13%
- 18%
- 33%
- 36%

**Not yet using AI**
- 1%
- 4%
- 16%
- 80%

Legend:
- above 20%
- 10-20%
- 5-10%
- 0-5%

**Note: firms with early advantage are investing more, moving still further away from the pack.**
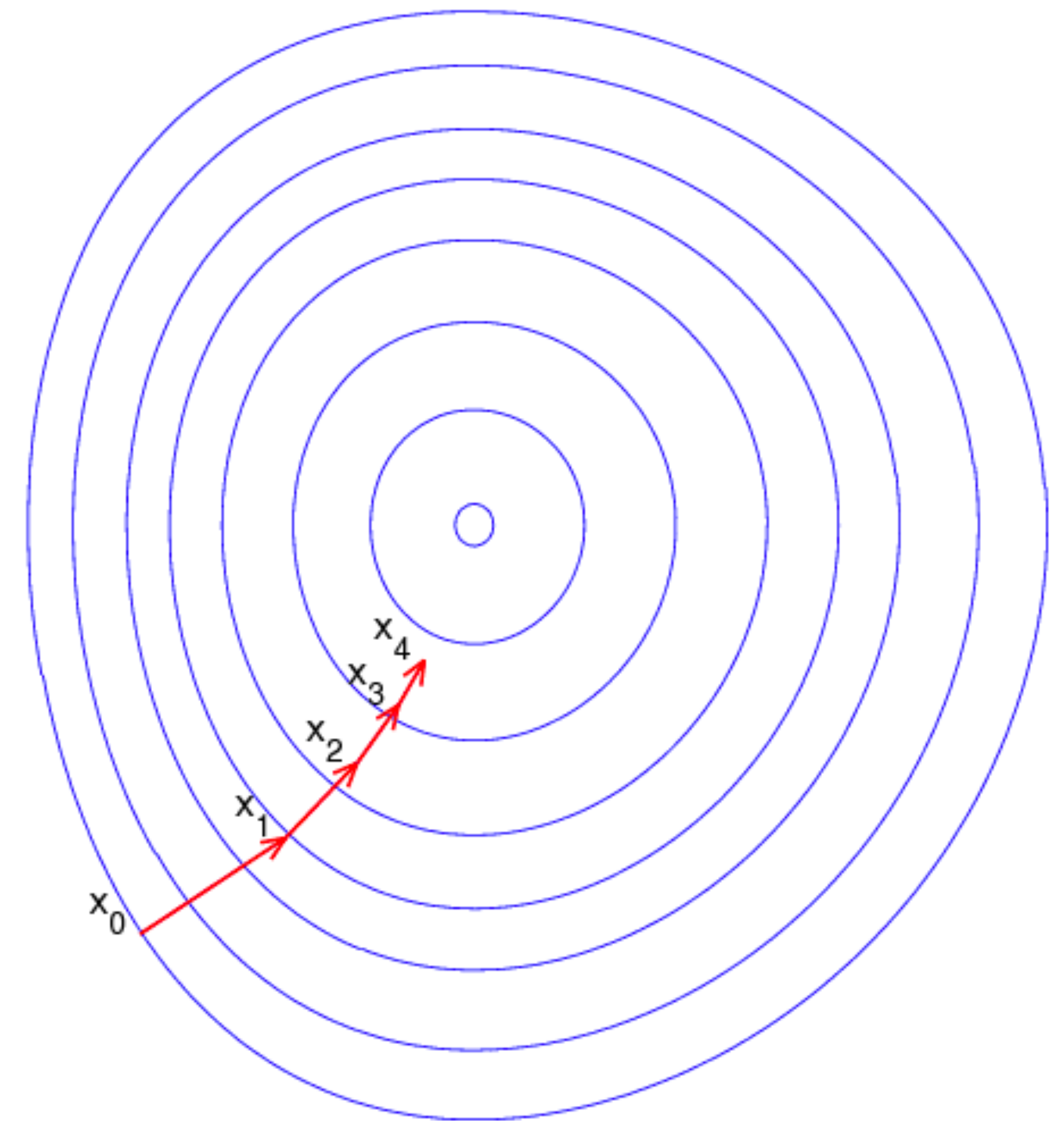
**2020 spending expected to increase 31%** – *WSJ*

# Concerns…

Large investments by **hyperscalers** which own public clouds, have access to enormous amounts of training data, and hired large teams of AI talent … are now sold back through a range of AutoML services.

**TensorFlow Dev Summit 2018**
Jeff Dean: "Can we replace scarce ML expertise with 100x computation?"

# Concerns...

Potential *attack surfaces* are getting exposed much more rapidly than they're becoming understood.

Contemporary work in data privacy technologies **injects noise** into data streams – mirrored by generative AI efforts creating *security exploits* that inject carefully designed noise into data streams:

**github.com/kenny-co/procedural-advml**

The math of DL is not clearly understood, e.g., why don't we encounter local minima more frequently? Advanced work in numerical analysis is exploring these phenomena:

**"Understanding deep neural networks"**
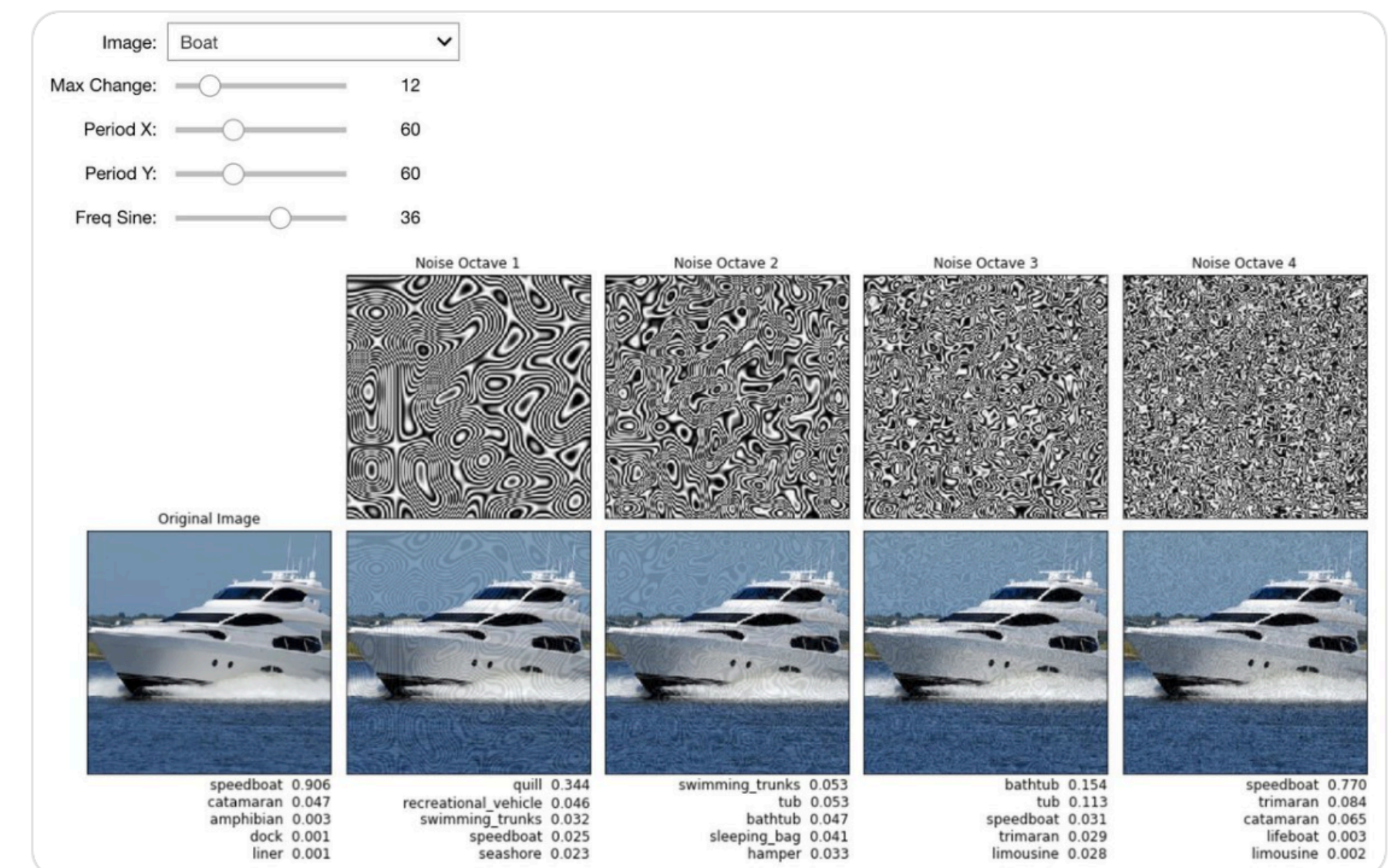**Michael Mahoney**

**Kenneth Co**
@kennyco26

Following

Released code for our black-box attack on DNNs. You can now generate your own procedural adversarial examples!
github.com/kenny-co/proce...

Image: Boat
Max Change:        12
Period X:          60
Period Y:          60
Freq Sine:         36

Noise Octave 1    Noise Octave 2    Noise Octave 3    Noise Octave 4

Original Image

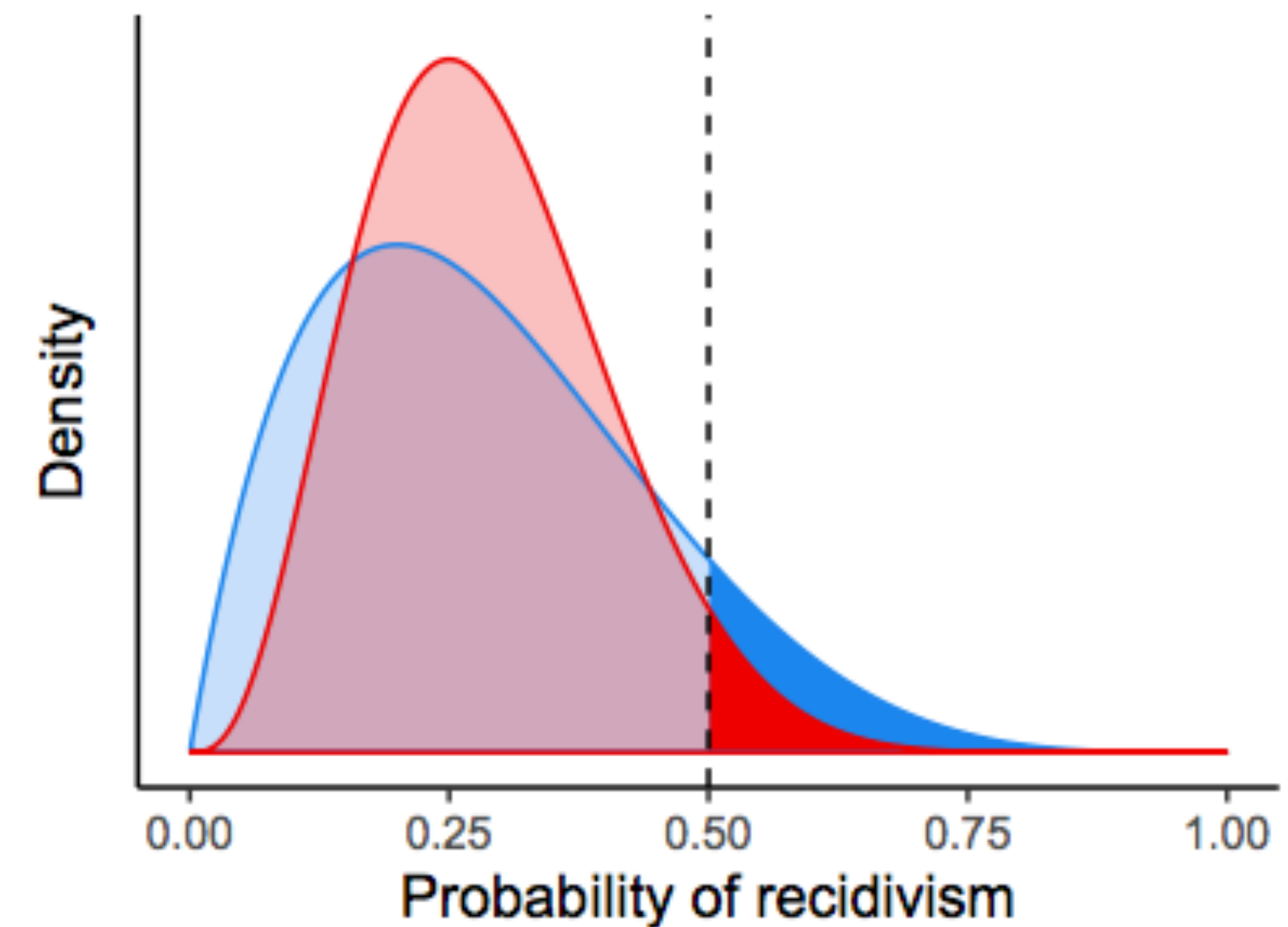| speedboat 0.906 | quill 0.344 | swimming_trunks 0.053 | bathtub 0.154 | speedboat 0.770 |
| catamaran 0.047 | recreational_vehicle 0.046 | tub 0.053 | tub 0.113 | trimaran 0.084 |
| amphibian 0.003 | swimming_trunks 0.032 | bathtub 0.047 | speedboat 0.031 | catamaran 0.065 |
| dock 0.001 | speedboat 0.025 | sleeping_bag 0.041 | trimaran 0.029 | lifeboat 0.003 |
| liner 0.001 | seashore 0.023 | hamper 0.033 | limousine 0.028 | limousine 0.002 |

12:20 PM - 25 Feb 2019

# Concerns...

Understanding fairness and bias in use of data is a provably hard problem:

**"Why it's hard to design fair machine learning models"**
**Sharad Goel**, **Sam Corbett-Davies**

recommended: open source **AIF360** toolkit
**http://aif360.mybluemix.net/**

# Concerns...

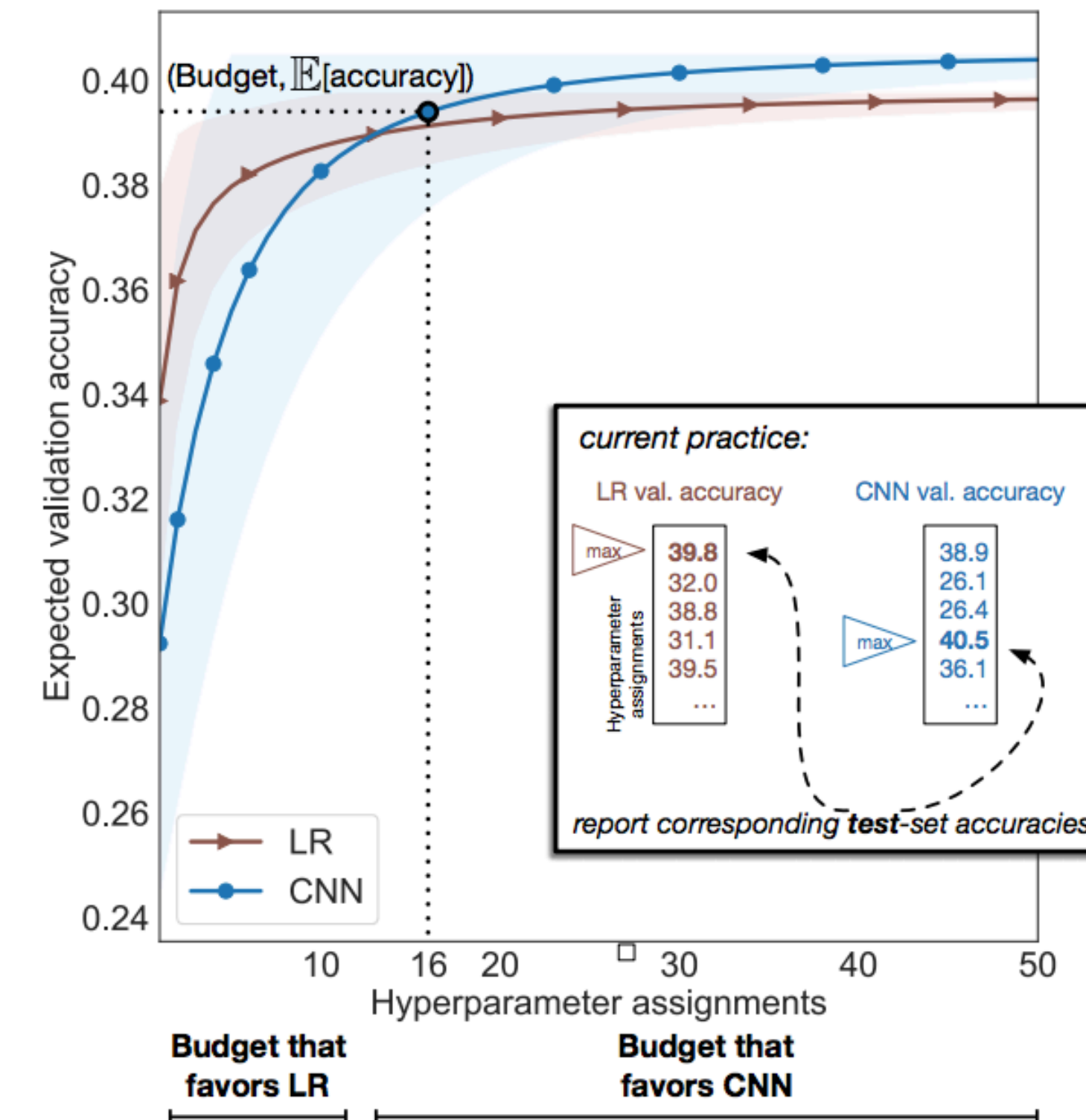Results in widely cited papers are too expensive for most organizations (except **hyperscalers)** to reproduce…

**"Show Your Work: Improved Reporting of Experimental Results"**
**Jesse Dodge**, **Suchin Gururangan**, **Dallas Card**, **Roy Schwartz**, **Noah A. Smith**

Policy: training a *transformer* model (GPT-2 w/ NAS) requires 5x the carbon footprint of operating an automobile over its lifetime…

**"Energy and Policy Considerations for Deep Learning in NLP"**
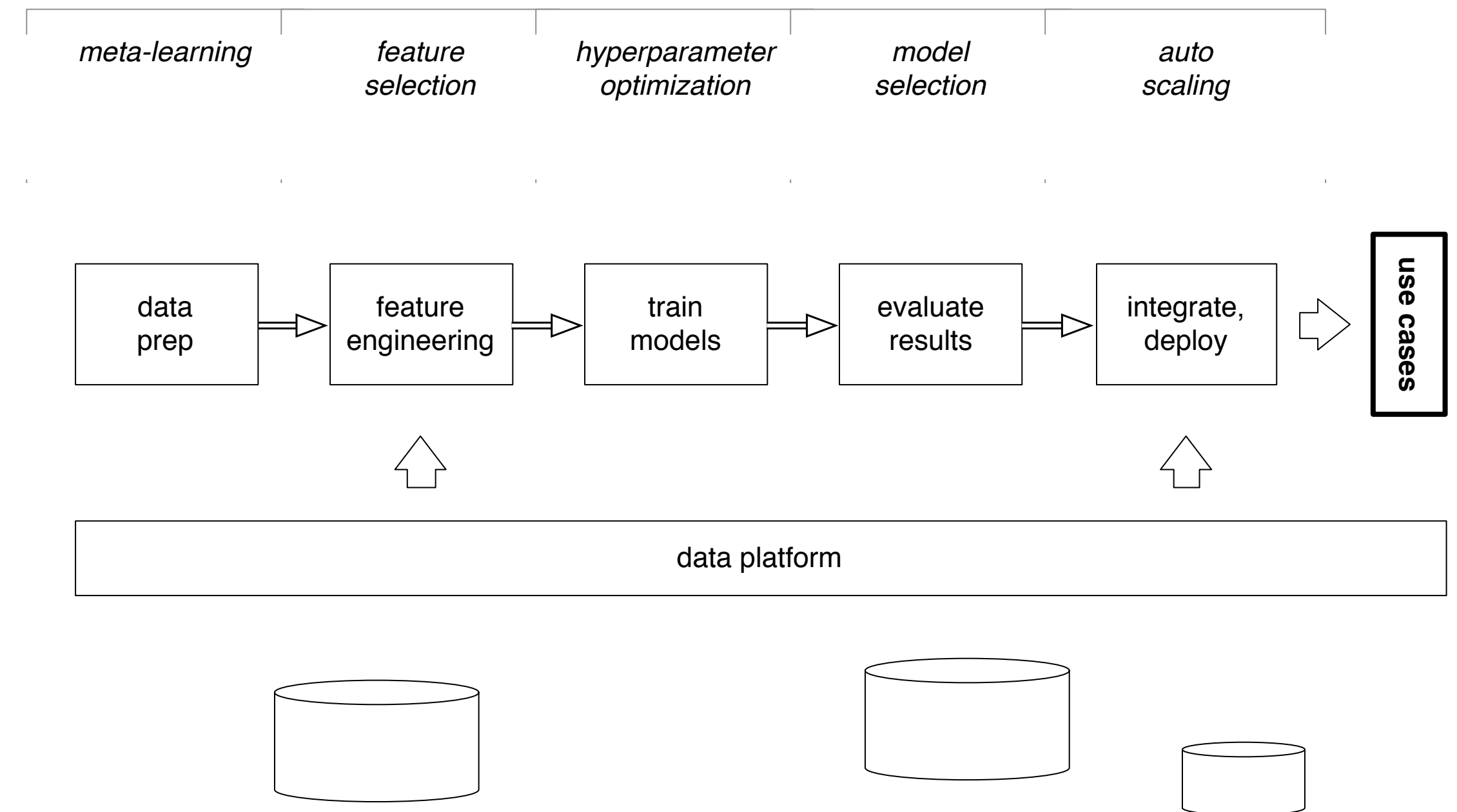**Emma Strubell**, **Ananya Ganesh**,
**Andrew McCallum**

# Unpacking AutoML



**derwen.ai/s/yvkg**

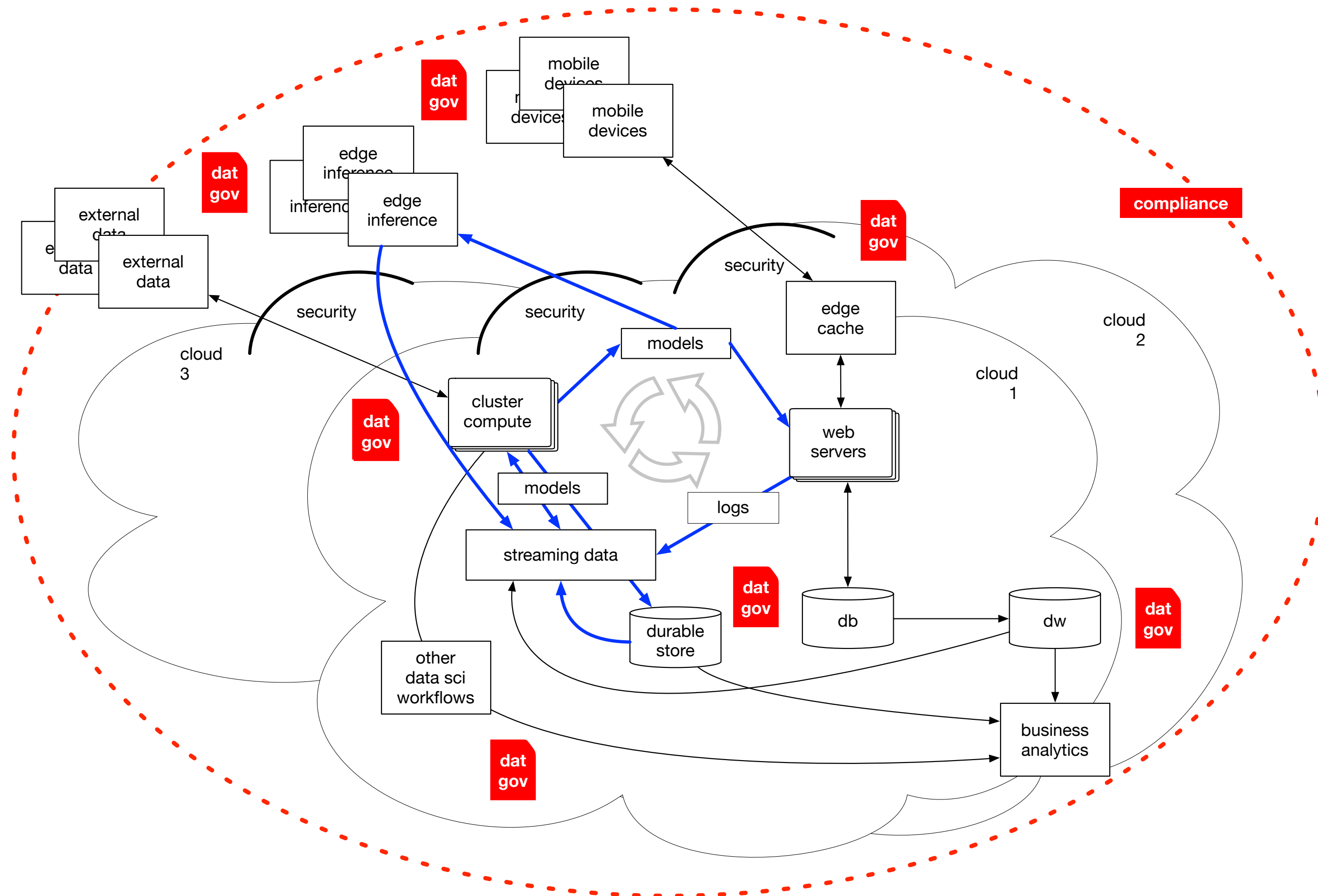*we noted an uptick in adoption for a third aspect, co-evolving along with DG and MLOps*



| meta-learning | feature selection | hyperparameter optimization | model selection | auto scaling |
|---|---|---|---|---|

data prep → feature engineering → train models → evaluate results → integrate, deploy → use cases

data platform

# Overview of Data Governance



**derwen.ai/s/6fqt**

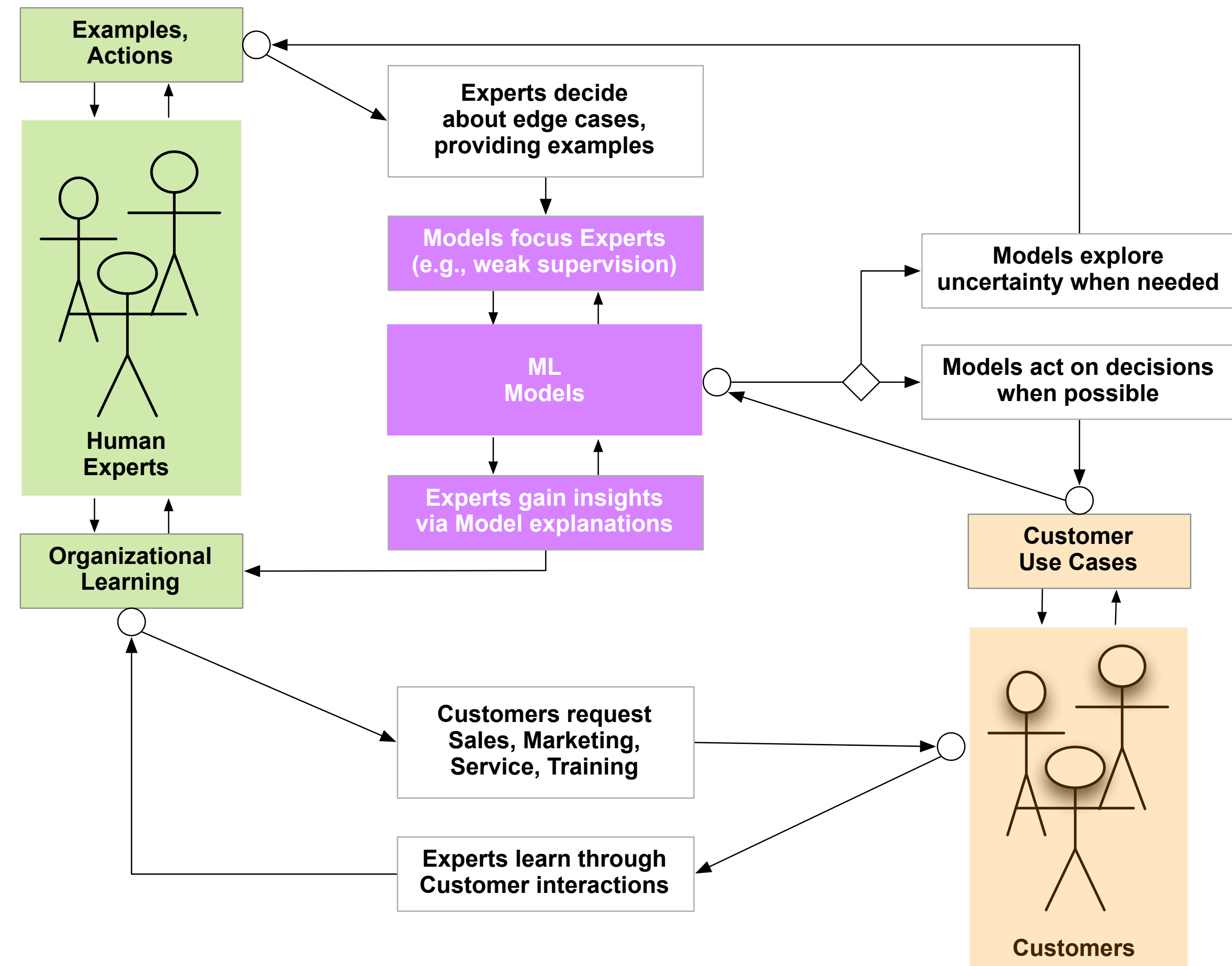*we noted a resurgence in data governance – this report examines key themes, vendors, issues, etc.*
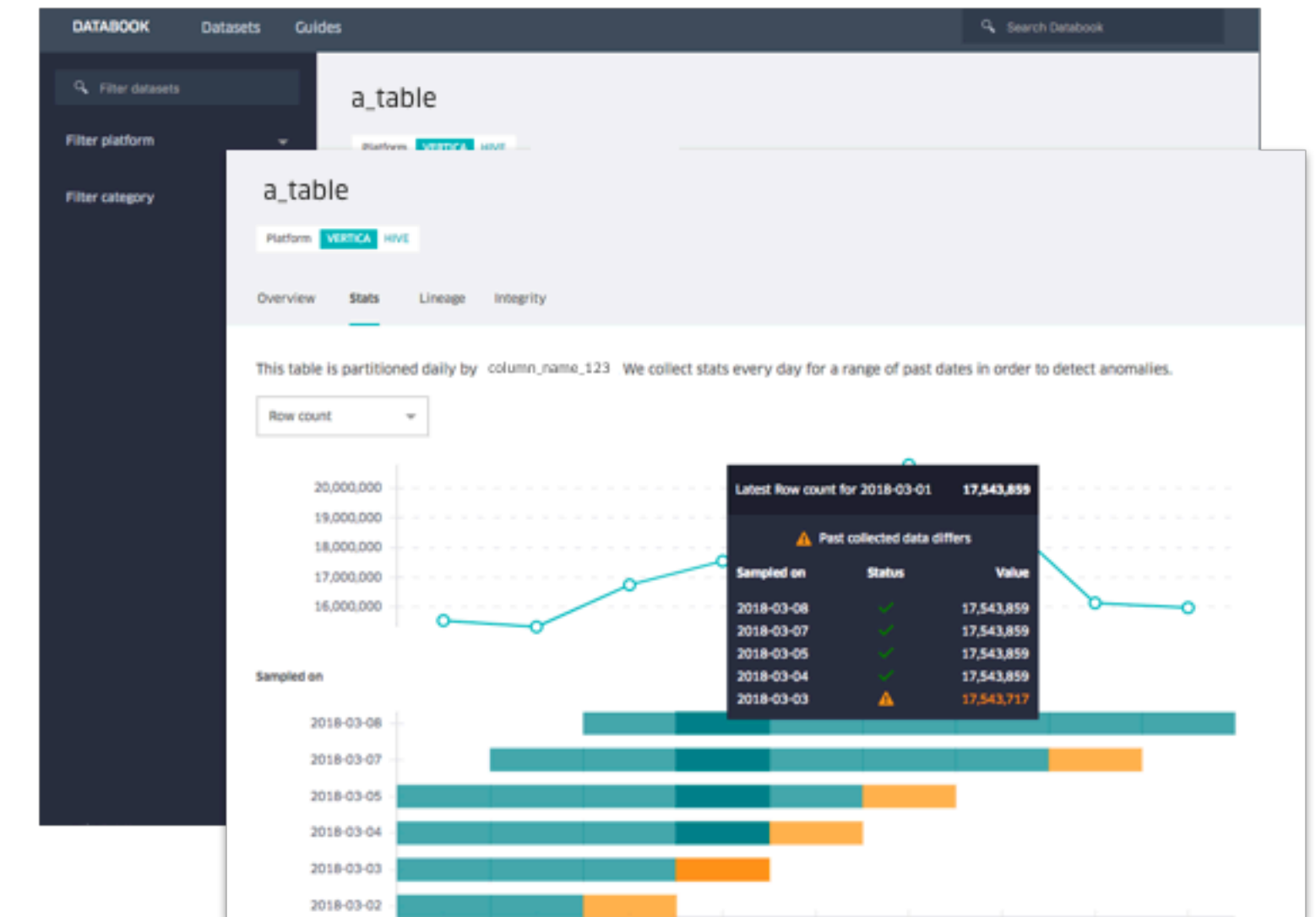
# Active Learning as a Data Strategy



Active Learning Use Cases

Paco Nathan  @pacoid  derwen.ai

**derwen.ai/s/d8b7**

*teams of people + machines, leveraging the complementary strengths of both*



Examples, Actions

Experts decide about edge cases, providing examples

Models focus Experts (e.g., weak supervision)

ML Models

Models explore uncertainty when needed

Models act on decisions when possible

Experts gain insights via Model explanations

Human Experts

Organizational Learning

Customer Use Cases

Customers request Sales, Marketing, Service, Training

Experts learn through Customer interactions

Customers

# Emerging category: watch the "AI Natives"

Projects (mostly OSS) that leverage **knowledge graph**
of metadata about datasets and their usage:

- **Amundsen** @ Lyft
  *data discovery and metadata*

- **Data Hub** @ LinkedIn
  *data discovery and lineage*

- **Marquez** @ Stitch Fix
  *collect, aggregate, visualize metadata*

- **UMS** @ Uber
  *manage metadata about datasets*

- **Metcat** @ Netflix
  *data discovery, metadata service*

- **Dataportal** @ Airbnb
  *integrated data-space (not OSS)*

**part 3:**
**watch this closely –**
**rapid evolution of hardware**

# Hardware in perspective

An **emerging trend** disrupts the past 15-20 years of software engineering practice:

**hardware > software > process**

Hardware is now evolving more rapidly than software, which is evolving more rapidly than effective process

**Moore's Law** is all but dead, although ironically many inefficiencies grew to be based on it

**Project Jupyter**, **Apache Arrow**, **NumPyWren** and the related **Ray** are emblematic for data infrastructure transformation in enterprise









GitHub search hits predicted until 2020-02-24 (95% confidence interval)

# ASICs the size of iPads, for ML



Cerebras WSE
1.2 Trillion transistors
46,225 mm² silicon

Largest GPU
21.1 Billion transistors
815 mm² silicon

# Impact of edge devices and low-power ML

- **"EIE: Efficient Inference Engine on Compressed Deep Neural Network"**
  **Song Han**, et al., (2016-05-03)

- **Alasdair Allan**, Babilim Light Industries:
  **The intelligent edge and the demise of big data?**
  (2019-05-02)

- **Pete Warden**, Google:
  **TensorFlow.js on low-power devices**
  (2018-10-10)

- **tinyML Summit**, 2020 Feb 12-13, San Jose
  **https://www.tinymlsummit.org/**

# Implications of rapid hardware evolution

- ~6 firms (hyperscaler cartel) now drive the **Intel** roadmap … and extract most of the value?

- innovation cost curves have gone exponential, with side-effects becoming difficult to anticipate

- $4-10B per fab now, with potentially < 2 yrs utility?

- geopolitically orchestrated IP theft + cyberattacks:

    - undercutting R&D investments in the US

    - competitive pressure on US public cloud vendors

# Evolution of cloud patterns

UC Berkeley published a **2009 report** about early use cases for cloud computing, which foresaw the shape of industry deployments over much of the next decade, and led directly to **Apache Mesos** and **Apache Spark**

It's fascinating to study the **contrasts** between that 2009 report and its 2019 follow-up.

(minor footnote: **vimeo.com/3616394**)

**Above the Clouds: A Berkeley View of Cloud Computing**

Michael Armbrust
Armando Fox
Rean Griffith
Anthony D. Joseph
Randy H. Katz
Andrew Konwinski
Gunho Lee
David A. Patterson
Ariel Rabkin
Ion Stoica
Matei Zaharia

Electrical Engineering and Computer Sciences
University of California at Berkeley

Technical Report No. UCB/EECS-2009-28
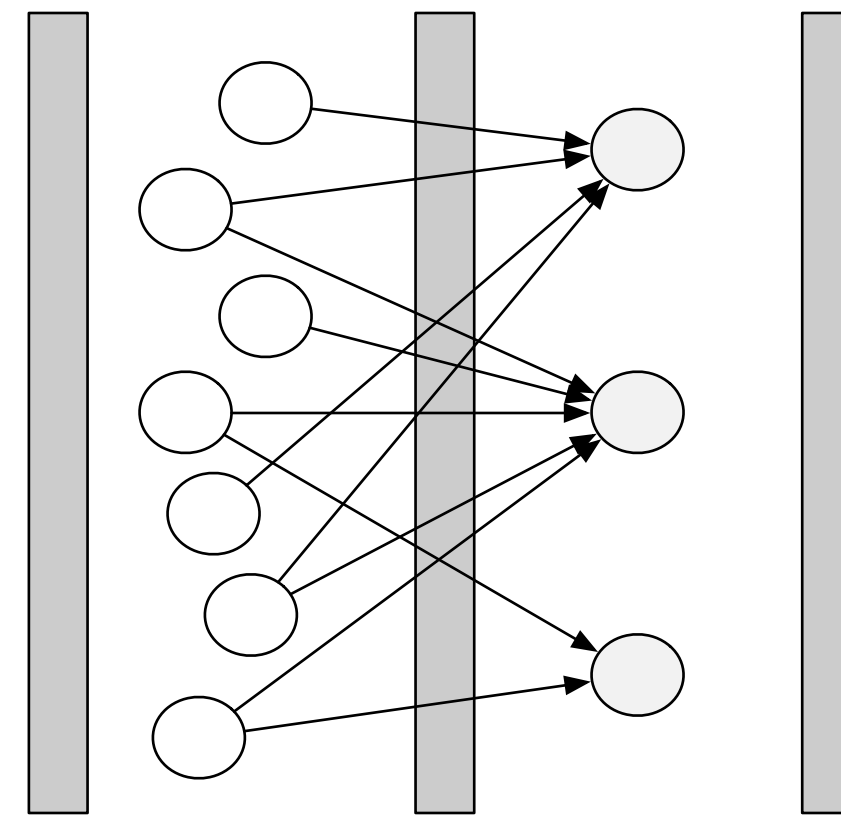http://www.eecs.berkeley.edu/Pubs/TechRpts/2009/EECS-2009-28.html

February 10, 2009

**2009**

# Evolution of cloud patterns

Early patterns of cloud use mirrored the **virtualization** familiar to enterprise firms in their existing on-prem architectures, mostly as a convenience – *for migration*

More **contemporary patterns** force restructuring – *for efficiency and security –* **decoupling computation and storage**

Cloud Programming Simplified: A Berkeley View on Serverless Computing

Eric Jonas
Johann Schleier-Smith
Vikram Sreekanti
Chia-Che Tsai
Anurag Khandelwal
Qifan Pu
Vaishaal Shankar
Joao Menezes Carreira
Karl Krauth
Neeraja Yadwadkar
Joseph Gonzalez
Raluca Ada Popa
Ion Stoica
David A. Patterson

Electrical Engineering and Computer Sciences
University of California at Berkeley

**2019**

# Key takeaways from "A Berkeley View"

- physics + cloud economics imply less "framework" layers

- service offerings **migrate up** the tech stack

- issues with **data gravity**, increased **vendor lock-in**


- also, we cannot overstate the importance of **data integrity**

  - the risks tend to **undermine science**

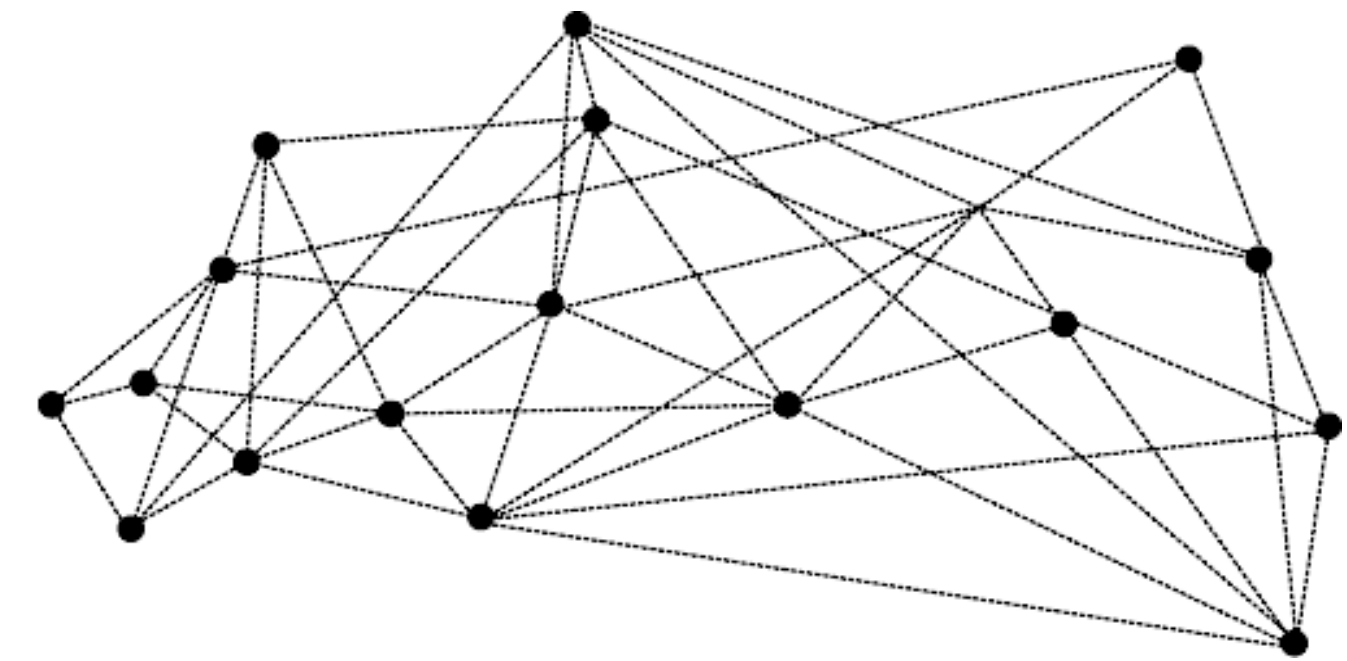  - questions of **global competitiveness** for training ML models; do we prioritize having volumes of training data or providing carefully **curated test sets**? (compare: China, US, EU

# Cluster topologies, by generation
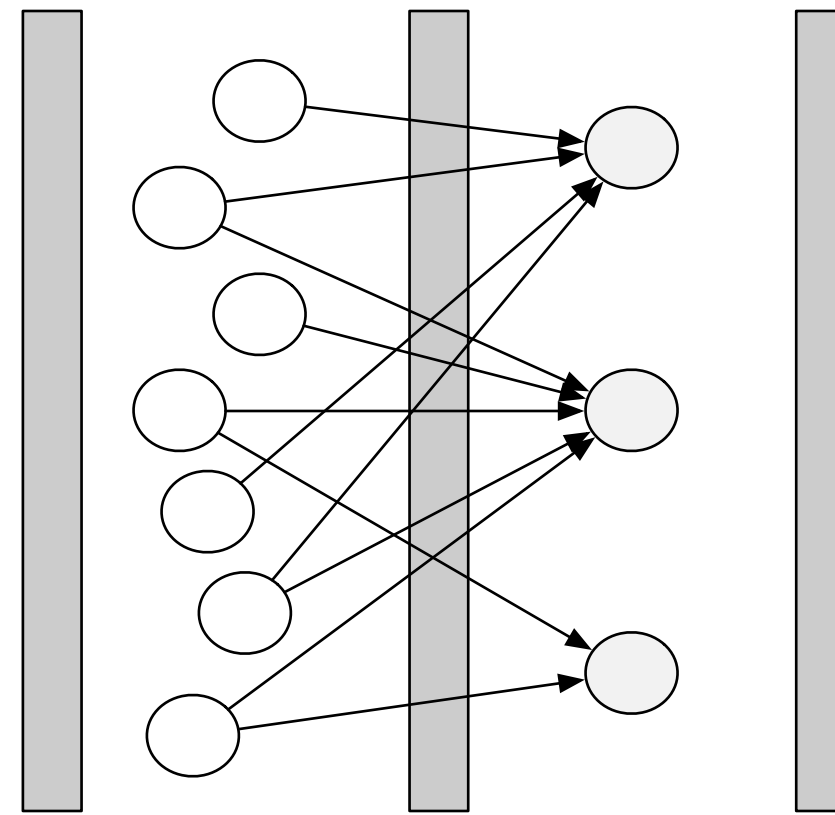


**1990s**

**mid-2000s**

**current**

Opinion: one problem with software/hardware interface for distributed systems is that it's taken *decades* to prioritize the need for handling **graphs/tensors** directly within popular, accessible open source libraries, without having some **commercial database vendor** intermediate.
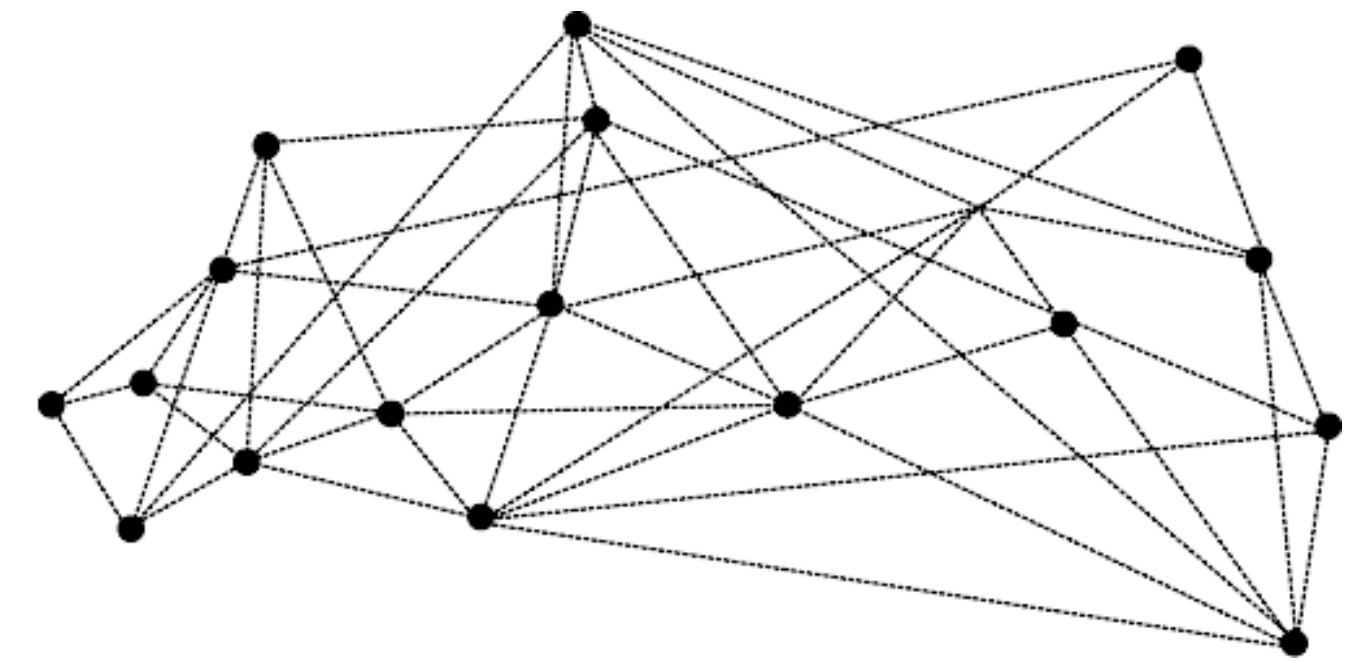
# Cluster topologies, by generation



**1990s**

**mid-2000s**

**current**

see also: **Jeff Dean** (2013)
**youtu.be/S9twUcX1Zp0**

# Next Generation Frameworks

Watch carefully about **Ray**, the next generation of distributed systems, i.e., what comes after Apache Spark:
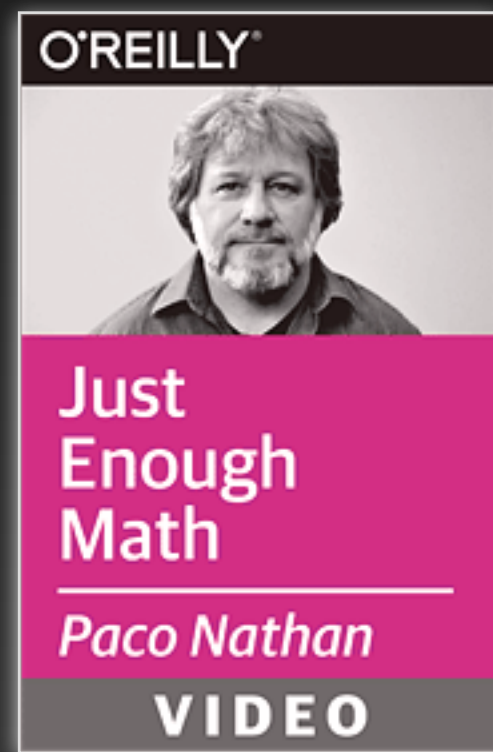
- focus on open source machine learning in Python

- accelerates work on laptops, scales to cloud without changing code; ~4 lines of code needed to parallelize

- supports ad-hoc queries, deep learning (TensorFlow, PyTorch), and reinforcement learning

- developed in response to contemporary ML research, cloud economics (after +10 years), taking advantage of rapid hardware evolution

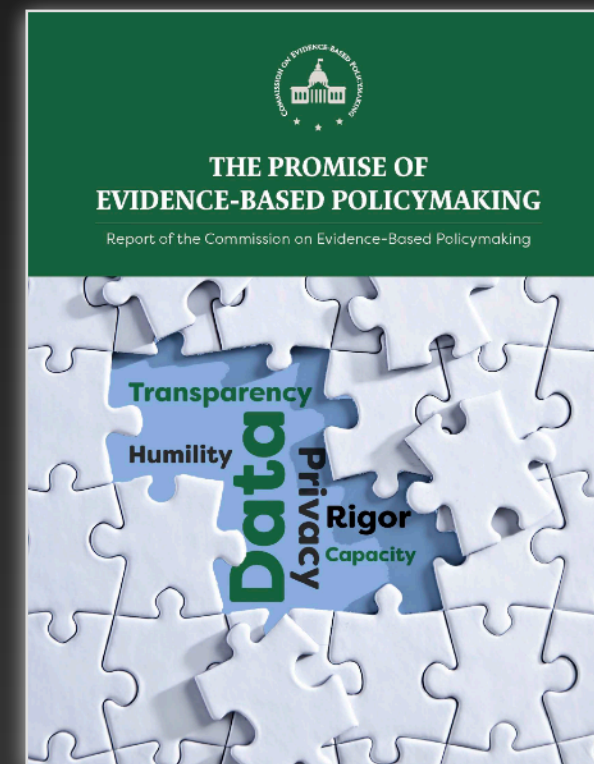- unifies *actor abstraction* with the *task-parallel abstraction*

See: **"Ray for the Curious"** by **Dean Wampler**

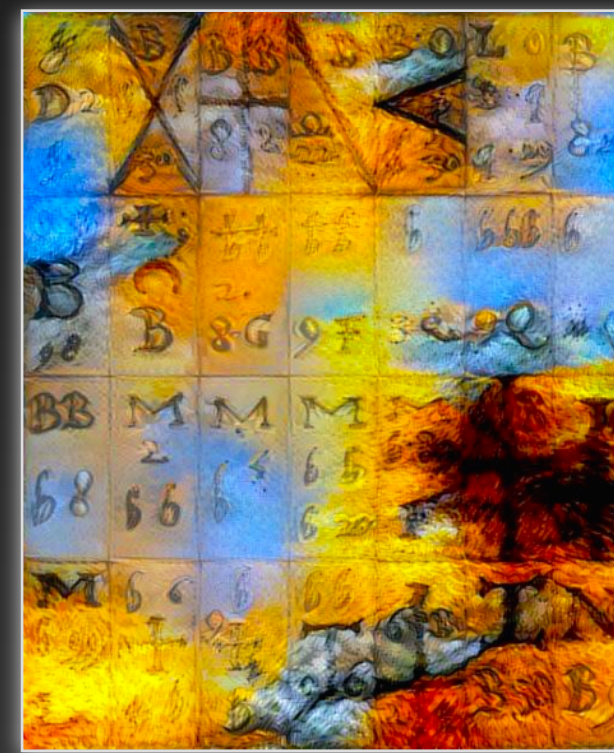# publications, interviews, conference summaries...

## https://derwen.ai/paco
## @pacoid



Just Enough Math

Rich Context

Hylbert-Speys

Rev conf

Themes + Confs
per Pacoid

**derwen.ai**