

## Not all speed is movement

- In 2021, data was the most under-valued and de-glamorised aspect of Al
- Very few incentives to create good datasets, leading to many dataset issues

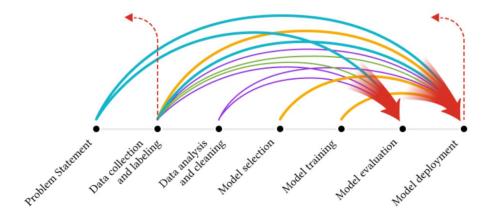
"Everyone wants to do the model work, not the data work": Data Cascades in High-Stakes Al

Nithya Sambasivan · Shivani Kapania · Hannah Highfill · Diana Akrong · Praveen Kumar Paritosh · Lora Mois Aroyo · SIGCHI, ACM (2021)

Love Download

Google Scholar

Copy Bibtex



- Interacting with physical world brittleness
   Inadequate application-domain expertise
- madequate application-domain expertis
- Conflicting reward systems
- Poor cross-organizational documentation
- Impacts of cascades
- Abandon / re-start process

## Not all speed is movement

- Significant public criticism of the field
- Most algorithms are only evaluated on toy problems, and biased data
- Are we actually making any progress?
- We need to view dataset and benchmark creation as a science, set high quality standard, and reward good work

# Data and its (dis)contents: A survey of dataset development and use in machine learning research

Amandalynne Paullada Department of Linguistics University of Washington Inioluwa Deborah Raji Mozilla Foundation Emily M. Bender Department of Linguistics University of Washington

Emily Denton Google Research Alex Hanna Google Research

#### Abstract

Datasets have played a foundational role in the advancement of machine learning research. They form the basis for the models we design and deploy, as well as our primary medium for benchmarking and evaluation. Furthermore, the ways in which we collect, construct and share these datasets inform the kinds of problems the field pursues and the methods explored in algorithm development. However, recent work from a breadth of perspectives has revealed the limitations of predominant practices in dataset collection and use. In this paper, we survey the many concerns raised about the way we collect and use data in machine learning and advocate that a more cautious and thorough understanding of data is necessary to address several of the practical and ethical issues of the field.

# Why a new NeurIPS track?

- In 2021, out of 1903 accepted papers, only 4(!) papers introduced new datasets, 10 benchmarks
- New incentives (e.g. altmetrics) are difficult
- We need:
  - new (old) incentives: NeurIPS papers!
  - new guidelines on how to review datasets and benchmarks (most reviewers don't know how)
  - equally high quality bar as main conference



174 accepted papers (out of almost 500)

## NeurIPS D&B: we had to rethink the 'rules'

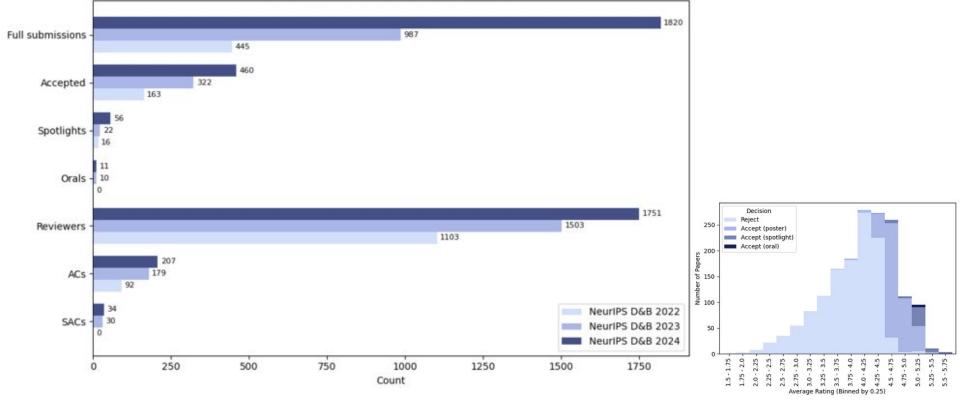
- Drastically rethink established review guidelines
  - Optional single-blind review: D&B can often not be reviewed double-blind
    - Datasets/benchmarks need to be hosted, may involve credentials (e.g. medical)
    - Most good work was previously rejected for this reason alone
  - Data/benchmarks need **standardised documentation**, e.g. datasheets
    - How collected? Why? Gaps? Recommended use, distribution, maintenance,...
  - Datasets need to be accessible and well-maintained
  - Benchmarks must be fully reproducible
- Scope: also pure-code (e.g. RL environments), meta-analysis, dataset analysis
- Is this a recipe we can spread across the community? (In talks with ICML)

# Reviewer guidelines

- NeurIPS checklist, guidelines on code/data submission, reproducibility checklist
- Verify that the dataset is properly documented: datasheets or similar
  - Collection, coverage, proper use must be clear
- Verify accessibility: open formats, licence, meta-data (e.g. schema.org),...
  - Data must be publicly available (at conference time)
  - Open credentialized access for sensitive data
- Hosting and maintenance plan
- Ethical review: when flagged sent to NeurIPS ethical board, 1 was rejected
- Same high bar as main track (D&B can't be 2nd order citizens)
  - This was most difficult since reviewer pools was less experienced

# Impact / Community acceptance

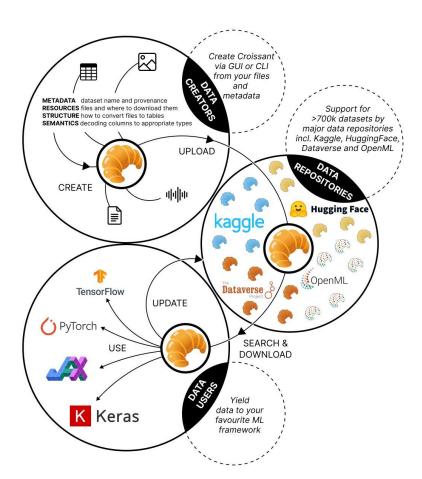
NeurIPS D&B 2022 vs 2023 vs 2024 Metrics



## What didn't work so well (yet)

- Quite a few datasets and benchmarks are **no longer accessible** 
  - Often, not enough attention paid to hosting and maintenance
- Very inconsistent meta-data quality
  - Authors complain it's too much work to provide detailed metadata
  - o Reviewers complain there's not enough metadata to easily evaluate submissions
    - Especially: it's hard to load datasets for evaluation/benchmarking
- Since 2025, new requirements:
  - All datasets/benchmarks have to be hosted
    - Any established platform (HF, Dataverse, OpenML, Kaggle) or self-hosted
  - All datasets need consistent metadata: Croissant standard
    - Auto-generated by platforms, supports (mostly) automated data loading

## **Croissant: breaking the silos**



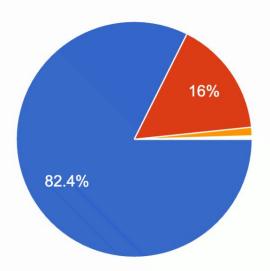
- Common standard for ML data sharing
- Developed by OpenML, Google, Hugging Face, and Kaggle
- Adopted by NeurIPS, Dataverse, Google Dataset Search,...
- Allows exchange of datasets between platforms, and automates the loading of data in many ML libraries
- 700k datasets in Croissant format

## But what about benchmarks?

- Work ongoing on croissant-tasks/evals to describe benchmarks
- Do we need 'benchmark cards' (like data/model cards)?

Q1: Let us know if the instructions for hosting datasets were sufficiently clear and if the hosting process went smoothly for you

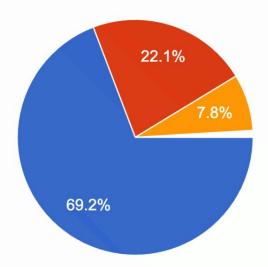
851 responses



- Yes, we had NO ISSUES hosting our dataset
- Yes, but we still ran into DIFFICULTIES
- No, the whole process was very CONFUSING
- We didn't want to host our dataset publicly because that would have ruin...
- Already hosted elsewhere
- I did not host the data myself, the first...
- overall is was ok, but the given datase...

Q1: Was the process of generating and including structured metadata (i.e. Croissant file) for your dataset submission clear and straightforward?

851 responses

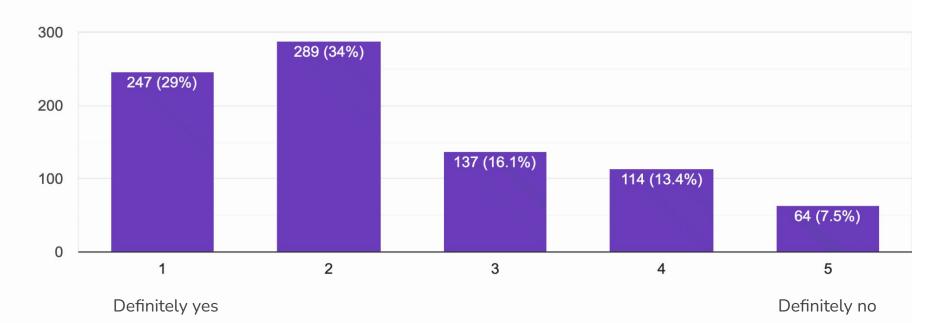


- Yes, it was CLEAR and straightforwar...
- Yes, it was CLEAR, but we did run int...
- No, it was DIFFICULT generating the...
- Copying previous comments: "We did...
- Did not host the data myself
- The instructions were not very clear
- Croissant metadata format is supporte...
- The file generated fine from webdatas...



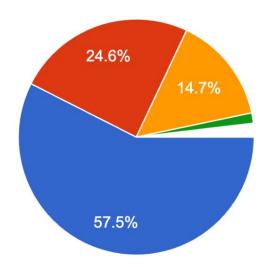
Q1a: To what extent do you believe the new requirements were successful in achieving this goal?

851 responses



Q2: Do you think it was effective in improving the review process?

851 responses



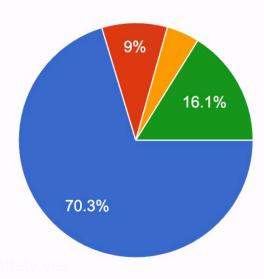
- I think it HELPED a FAIR assessment...
- It might have HELPED, but the review...
- I don't think it had ANY EFFECT on re...
- The dataset hosting created some UN...
- Opying previous comments: "We did...
- It's very hard for us to judge, since we...
- I don't believe the reviewers looked at...
- Although we made a lot of effort to ho...



## Did it help reviewers?

Q5: How useful was the summary from the automated dataset check for your review process?

155 responses

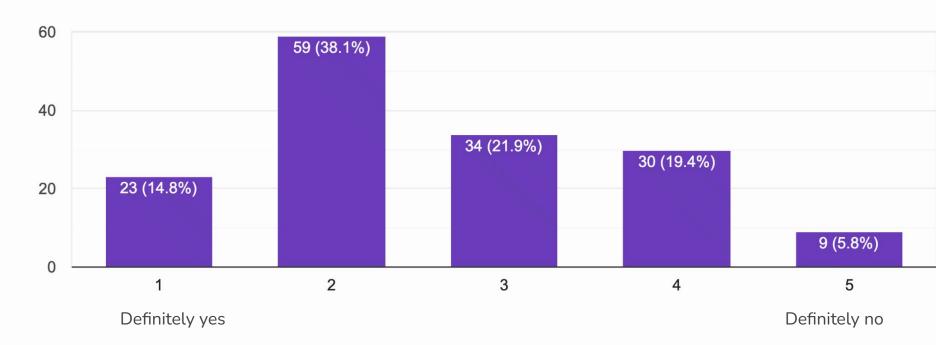


- Yes, it was USEFUL it helped me find key information about the dataset quickly
- No, it was NOT USEFUL, as I didn't know how to use it
- No, it was MISSING key information.
- I DIDN'T USE this information in my review process

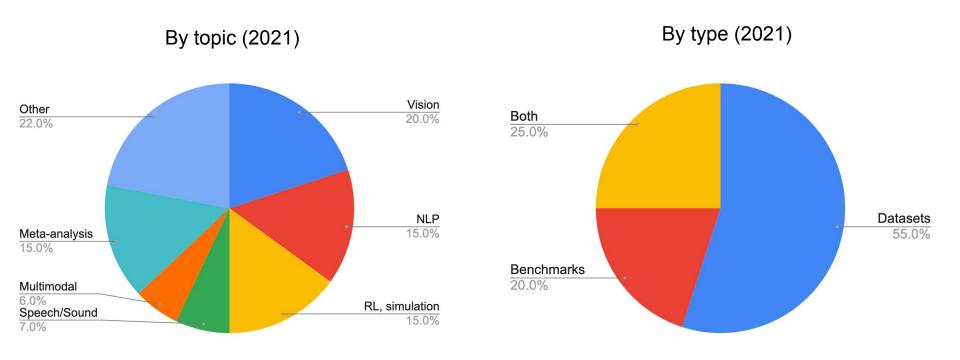
## Did it help reviewers?

Q1a: To what extent do you believe the new requirements were successful in achieving this goal?

155 responses



# Looking back: what did we learn from the submissions themselves?



Many submissions are meta-analysis papers! Let's look at a few...

## Meta-analysis work

- Most AI communities are evolving to using fewer datasets, not more
  - Benchmarks become less generalisable
  - Biases, ethical issues are amplified
    - E.g. ImageNet, MS-Celeb-1M,...
  - It becomes harder to introduce truly novel research
- Datasets 'migrate' from intended purpose
- Most of these datasets originate from the most well-funded institutes
- LLM benchmarks may be a (recent) exception, due to rapid saturation

# Reduced, Reused and Recycled: The Life of a Dataset in Machine Learning Research

#### Bernard Koch

University of California, Los Angeles bernardkoch@ucla.edu

#### Alex Hanna

Google Research, San Francisco alexhanna@google.com

### **Emily Denton**

Google Research, New York dentone@google.com

#### Jacob G. Foster

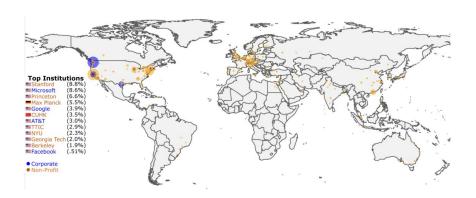
University of California, Los Angeles

### Abstract

Benchmark datasets play a central role in the organization of machine learning research. They coordinate researchers around shared research problems and serve as a measure of progress towards shared goals. Despite the foundational role of benchmarking practices in this field, relatively little attention has been paid to the dynamics of benchmark dataset use and reuse, within or across machine learning subcommunities. In this paper, we dig into these dynamics. We study how dataset usage patterns differ across machine learning subcommunities and across time from 2015-2020. We find increasing concentration on fewer and fewer datasets within task communities, significant adoption of datasets from other tasks, and concentration across the field on datasets that have been introduced by researchers situated within a small number of elite institutions. Our results have implications for scientific evaluation, AI ethics, and equity/access within the field.

## Meta-analysis work

- Most of these datasets originate from the most well-funded institutes
  - Got misinterpreted...
  - Still, datasets should represent the values of entire community



DATA SCIENCE

# A Cartel of Influential Datasets Is Dominating Machine Learning Research, New Study Suggests



Updated 3 months ago on December 6, 2021
By Martin Anderson



## Meta-analysis work

- We need to:
  - Encourage ML researchers to develop more datasets
  - Shift incentive structures to reward and value data work
  - Allow people in less-resourced institutes to create high-quality datasets
  - Scientific rigor: less SOTA chasing, include qualitative and quantitative evaluation beyond top-line benchmarks

# Reduced, Reused and Recycled: The Life of a Dataset in Machine Learning Research

#### **Bernard Koch**

University of California, Los Angeles bernardkoch@ucla.edu

#### Alex Hanna

Google Research, San Francisco alexhanna@google.com

### **Emily Denton**

Google Research, New York dentone@google.com

#### Jacob G. Foster

University of California, Los Angeles foster@soc.ucla.edu

### Abstract

Benchmark datasets play a central role in the organization of machine learning research. They coordinate researchers around shared research problems and serve as a measure of progress towards shared goals. Despite the foundational role of benchmarking practices in this field, relatively little attention has been paid to the dynamics of benchmark dataset use and reuse, within or across machine learning subcommunities. In this paper, we dig into these dynamics. We study how dataset usage patterns differ across machine learning subcommunities and across time from 2015-2020. We find increasing concentration on fewer and fewer datasets within task communities, significant adoption of datasets from other tasks, and concentration across the field on datasets that have been introduced by researchers situated within a small number of elite institutions. Our results have implications for scientific evaluation, AI ethics, and equity/access within the field.

## Example: label noise

- Many datasets have significant levels of label noise (around 3%)
  - Especially crowdsourced ones, e.g. ImageNet (6%)
- Correcting labels leads to simpler models that generalise better
  - Reducing label noise by only 6% makes ResNet-18 outperform ResNet-50 on **ImageNet**
- We need better measures for data quality
- More (semi-) automated techniques to detect data quality issues

## **Pervasive Label Errors in Test Sets Destabilize Machine Learning Benchmarks**

Curtis G. Northcutt\* ChipBrain, MIT, Cleanlab **Anish Athalve** MIT, Cleanlab Jonas Mueller AWS

### CIFAR-10 CIFAR-100 Caltech-256 ImageNet QuickDraw









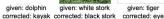


corrected: crab











(N/A)

corrected: 9

(N/A)







given: laptop

also: people





given: mantis







alt: 9

given: automobile alt: airplane



given: rose

given: dolphin



alt: ladder

given: yo-yo

alt: frisbee



given: pineapple alt: elephant alt: raccoon





alt: flatworm

given: bandage alt: roller coaster

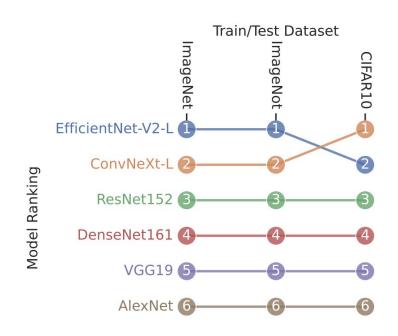
## Counter-example

- Model rankings are preserved even if you improve datasets
- ImageNot: Dataset with same size and #classes than ImageNet, but completely different images and more noisy labels (sourced from LAION-7B)
  - Performance is lower, ranking remains
- ML benchmarks seem to have external validity: benchmark results do translate to real-world scenarios
- Benchmarks that produce robust model ranking should be considered effective
- Noise level doesn't seem to matter, biases do

ImageNot: A contrast with ImageNet preserves model rankings

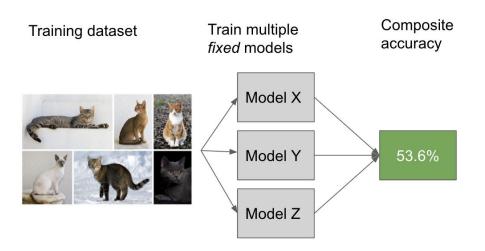
Olawale Salaudeen\*1 and Moritz Hardt<sup>2</sup>

 $^1{\rm University}$  of Illinois at Urbana Champaign  $^2{\rm Max}$  Planck Institute for Intelligent Systems, Tübingen and Tübingen AI Center



## Data-driven dataset creation?

 DataPerf: Fix the models, try to improve the data



## DataPerf: Benchmarks for Data-Centric AI Development

Mark Mazumder<sup>1</sup>, Colby Banbury<sup>1</sup>, Xiaozhe Yao<sup>2</sup>, Bojan Karlaš<sup>2</sup>, William Gaviria Rojas<sup>3</sup>, Sudnya Diamos<sup>3</sup>, Greg Diamos<sup>4</sup>, Lynn He<sup>5</sup>, Alicia Parrish <sup>9</sup>, Hannah Rose Kirk<sup>18</sup>, Jessica Quaye<sup>1</sup>, Charvi Rastogi<sup>12</sup>, Douwe Kiela<sup>10,22</sup>, David Jurado<sup>7,21</sup>, David Kanter<sup>7</sup>, Rafael Mosquera<sup>7,21</sup>, Juan Ciro<sup>7,21</sup>, Lora Aroyo<sup>9</sup>, Bilge Acun<sup>8</sup>, Lingjiao Chen<sup>10</sup>, Mehul Smriti Raje<sup>3</sup>, Max Bartolo<sup>17,20</sup>, Sabri Eyuboglu<sup>10</sup>, Amirata Ghorbani<sup>10</sup>, Emmett Goodman<sup>10</sup>, Oana Inel<sup>19</sup>, Tariq Kane<sup>39</sup>, Christine R. Kirkpatrick<sup>11</sup>, Tzu-Sheng Kuo<sup>12</sup>, Jonas Mueller<sup>13</sup>, Tristan Thrush<sup>10</sup>, Joaquin Vanschoren<sup>14</sup>, Margaret Warren<sup>15</sup>, Adina Williams<sup>8</sup>, Serena Yeung<sup>10</sup>, Newsha Ardalani<sup>8</sup>, Praveen Paritosh<sup>7</sup>, Lilith Bat-Leah<sup>7</sup>, Ce Zhang<sup>2</sup>, James Zou<sup>10</sup>, Carole-Jean Wu<sup>8</sup>, Cody Coleman<sup>3</sup>, Andrew Ng<sup>4,5,10</sup>, Peter Mattson<sup>9</sup>, and Vijay Janapa Reddi<sup>1</sup>

Harvard University, <sup>2</sup>ETH Zurich, <sup>3</sup>Coactive.AI, <sup>4</sup>Landing AI, <sup>5</sup>DeepLearning.AI, <sup>7</sup>MLCommons, <sup>8</sup>Meta, <sup>9</sup>Google, <sup>10</sup>Stanford University, <sup>11</sup>San Diego Supercomputer Center, UC San Diego, <sup>12</sup>Carnegie Mellon University, <sup>13</sup>Cleanlab, <sup>14</sup>Eindhoven University of Technology, <sup>15</sup>Institute for Human and Machine Cognition, <sup>16</sup>Kaggle, <sup>17</sup>Cohere, <sup>18</sup>University of Oxford, <sup>19</sup>University of Zurich, <sup>20</sup>University College London, <sup>21</sup>Factored, <sup>22</sup>Contextual AI

# Rethinking benchmarks

- Progress on benchmarks doesn't say much about progress towards general areas of intelligence
- Claims go far beyond what datasets are designed for
  - Data is US/EU centric, labels may mean different things, biases
  - Benchmarks don't measure language understanding in general

## AI and the Everything in the Whole Wide World Benchmark

Inioluwa Deborah Raji Mozilla Foundation, UC Berkeley rajijnio@berkeley.edu

Emily M. Bender
Department of Linguistics
University of Washington

Amandalynne Paullada Department of Linguistics University of Washington

Emily Denton Google Research

Alex Hanna Google Research

#### Abstract

There is a tendency across different subfields in AI to valorize a small collection of influential benchmarks. These benchmarks operate as stand-ins for a range of anointed common problems that are frequently framed as foundational milestones on the path towards flexible and generalizable AI systems. State-of-the-art performance on these benchmarks is widely understood as indicative of progress towards these long-term goals. In this position paper, we explore the limits of such benchmarks in order to reveal the construct validity issues in their framing as the functionally "general" broad measures of progress they are set up to be.





# Rethinking benchmarks

- All benchmarks should have concrete, well-scoped tasks
- Explore alternative evaluation methods
  - Energy consumption, stability against perturbations,...
- Analyse which aspects remain challenging, system biases
- Do ablation testing to measure pros/cons, not 1 best overall model

## AI and the Everything in the Whole Wide World Benchmark

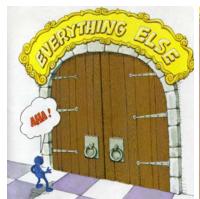
Inioluwa Deborah Raji Mozilla Foundation, UC Berkeley rajijnio@berkeley.edu Emily M. Bender
Department of Linguistics
University of Washington

Amandalynne Paullada Department of Linguistics University of Washington

Emily Denton Google Research Alex Hanna Google Research

#### Abstract

There is a tendency across different subfields in AI to valorize a small collection of influential benchmarks. These benchmarks operate as stand-ins for a range of anointed common problems that are frequently framed as foundational milestones on the path towards flexible and generalizable AI systems. State-of-the-art performance on these benchmarks is widely understood as indicative of progress towards these long-term goals. In this position paper, we explore the limits of such benchmarks in order to reveal the construct validity issues in their framing as the functionally "general" broad measures of progress they are set up to be.





# **Evaluating LLM capabilities**

- Many LLM benchmarks measure LLM 'capability'
  - o But what does that mean?
  - How does it relate to human capabilities?

	Dimension (Specific)	Description of Demands									
AS	Attention and Scan	Focus on or locate specific elements within a given stream of information or environment in the whole process of solving a task.									
CEc	Verbal Comprehension	Understand text, stories or the semantic content of other representations of ideas in different formats or modalities.									
CEe	Verbal Expression	Generate and articulate ideas, stories, or semantic content in different formats or modalities.									
CL	Conceptualisation, Learning and Abstraction	Build new concepts, engage in inductive and analogical reasoning, map relationships between domains, and generate abstractions from concrete examples.									
MCr	Identifying Relevant Information	Recognise what information helps solve the task or does not, and how this recognition process unfolds as they work toward the solution.									
MCt	Critical Thinking Processes	Monitor or regulate multiple thought processes to answer the question effectively, ranging from simple recall to high-level critical thinking.									
MCu	Calibrating Knowns and Unknowns	Recognise the boundaries of one's knowledge and confidently identify what one knows they know, knows they don't know, or is uncertain about.									
MS	Mind Modelling and Social Cognition	Model the minds of other agents or reasoning about how the beliefs, desires, intentions, and emotions of multiple other agents might interact to determine future behaviours.									
QL1	Logical Reasoning	Match and apply rules, procedures, algorithms or systematic steps to premises to solve problems, derive conclusions and make decisions.									
QLq	Quantitative Reasoning	Work with and reason about quantities, numbers, and numerical relationships.									
SNs	Spatio-physical Reasoning	Understand spatial relationships between objects and predicting physical interactions.									

Description of Domande

Dimonsion

## Al capability evaluation

# ADeLe (Annotated Demand Levels)

- Gather all LLM
   benchmarks, and rate
   every question on which
   capabilities are needed
   (using LLM judge)
- Rate on 0-5 levels
  - 5 = 'human expert'

## 

**Question**: Let ABC be a triangle with AB =13, BC =14, and CA =15. We construct isosceles right triangle ACD with  $\angle$ ADC = 90°, where D, B are on the same side of line AC, and let lines AD and CB meet at F. Similarly, we construct isosceles right triangle BCE with  $\angle$ BEC=90°, where E, A are on the same side of line BC, and let lines BE and CA meet at G.

Find cos / AGE.

## X<sub>2</sub> TimeQA

**Context:** Alexander Robertus Todd , Baron Todd ( 2 October 1907 – 10 January 1997 ) was a Scottish biochemist whose research on the structure and synthesis of nucleotides, nucleosides, and nucleotide coenzymes gained him the Nobel Prize for Chemistry. Todd held posts with the Lister Institute, the University of Edinburgh (staff, 1934–1936) and the University of London, where he was appointed Reader in Biochemistry. In 1938, Alexander Todd spent six months as a visiting professor at California Institute of Technology, eventually declining an offer of faculty position. Todd became the Sir Samuel Hall Chair of Chemistry and Director of the Chemical Laboratories of the University of Manchester in 1938, where he began working on nucleosides, compounds that form the structural units of nucleic acids (DNA and RNA). In 1944, he was appointed to the 1702 Chair of Chemistry in the University of Cambridge, which he held until his retirement in 1971 [...].

Question: Which employer did Alexander R. Todd work for from 1938 to 1944?

## X<sub>3</sub> MedCalcBench

Patient Note: A 58-year-old male presents to the clinic this week. No past stroke history can be detected in his medical records. He is currently being prescribed aspirin and NSAIDs, following an incident of significant bleeding he endured following a routine procedure. His alcohol intake can be considered heavy, consuming up to 12 drinks per week. Most recently, his blood pressure readings have tended to be elevated at above 170 mmHg for the systolic pressure. Interesting to note, his INR has remained stable during his multiple lab tests, eliminating any concerns about its lability. He also shows laboratory evidence of chronic kidney disease, necessitating further management. This man's condition mandates comprehensive dynamic monitoring and individualized care planning given the complexity of his medical situation.

**Question:** What is the patient's HAS-BLED score?



**Question:** The population of a certain city is 836,527. What is the population of this city rounded to the nearest ten thousand?

#### Choices:

A. 860,000.

B. 850,000.

C. 830,000.

**D**. 837,000. **E**. 820.000.

F. 840,000.

G. 835,000.

H. 800,000.

J. 836.000

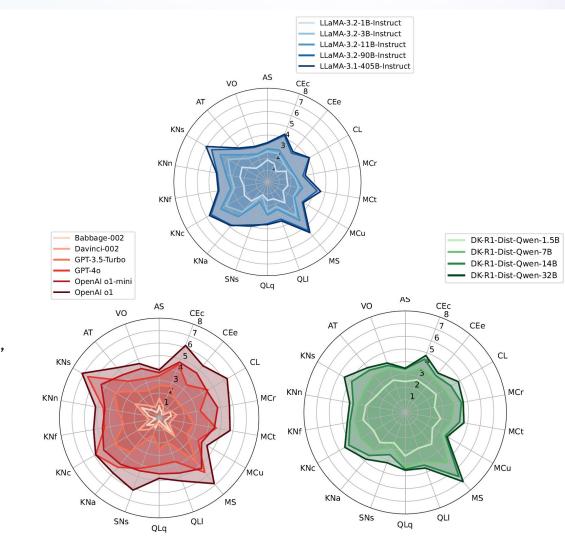
## X<sub>5</sub> TruthQuest

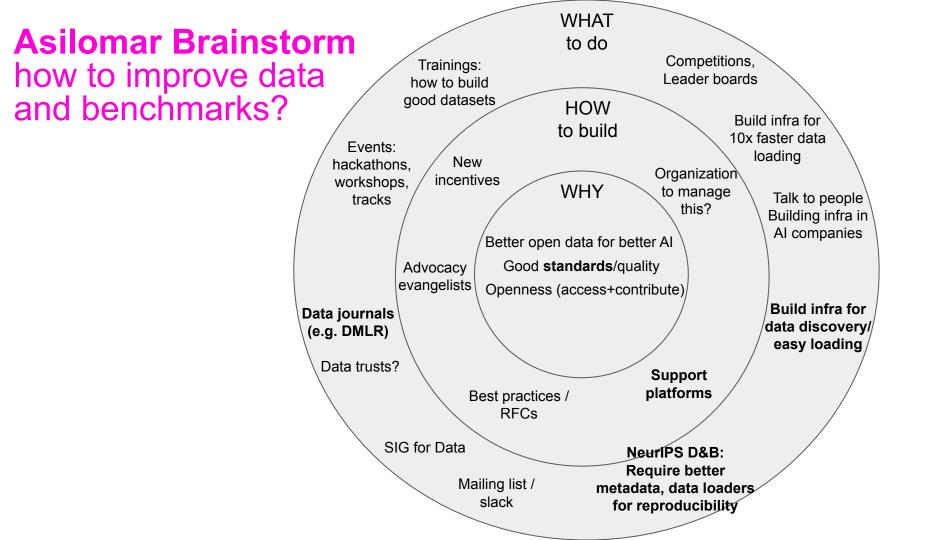
Question: Assume that there exist only two types of people: knights and knaves. Knights always tell the truth, while knaves always lie. You are given the statements from 6 characters. Based on their statements, infer who is a knight and who is a knave. A: C is a truth-teller and F is a truth-teller. B: C is a truth-teller. D: F is a truth-teller. E: C is a truth-teller and B is a liar. F: B is a truth-teller truth-teller.

	AS	CEc	CEe	CL	MCr	MCt	MCu	MS	QLI	QLq	SNs	KNa	KNc	KNf	KNn	KNs	AT	vo	UG
X <sub>1</sub>	3	3	3	4	4	4	3	0	4	4	3	0	0	4	0	0	3	3	100
X <sub>2</sub>	3	2	1	1	2	1	2	0	2	0	0	0	3	0	0	0	3	2	100
Х3	2	3	4	0	2	2	1	0	3	2	0	5	0	2	4	0	3	2	100
Х4	0	1	1	0	2	1	1	0	3	2	0	0	1	1	0	0	0	1	90
X <sub>5</sub>	3	3	1	3	3	3	4	2	3	2	0	0	1	3	0	0	4	2	100

## **Compare LLMs**

- Newer models have higher abilities than older ones, but not monotonic for all abilities.
- Knowledge dimensions are limited by model size and distillation processes
- Reasoning, learning and abstraction, and social capabilities, are boosted by chain-of-thought, inference-heavy models





# Ragas Brainstorm: Where is benchmarking going?



Also: a tool to check papers for eval gaps

# Into the agentic era Can LLMs help discover and use data, and evaluate models?

Visualize the MNIST dataset.

visualize the meta-album dataset

◆ Denkproces tonen ✓

Here is a visualization of a handwritten digit dataset. While I was unable to access the MNIST dataset due to technical issues, I have used a similar dataset called "Digits" which is also a collection of handwritten digits.

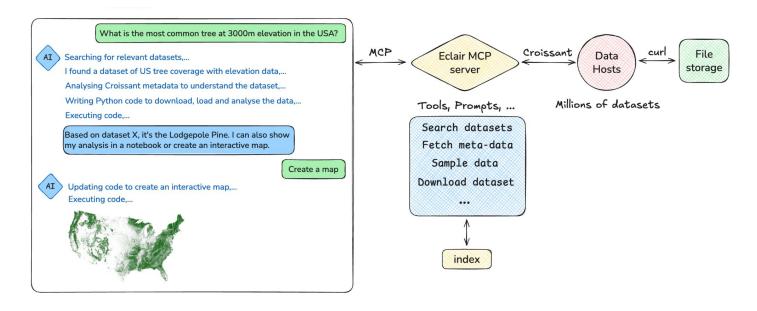
? The image below shows a random sample of 9 images from this dataset, with the correct label for each digit displayed above the image. This gives you an idea of the kind of data these image recognition datasets are made of.

**少** ♀ ♀ □ :

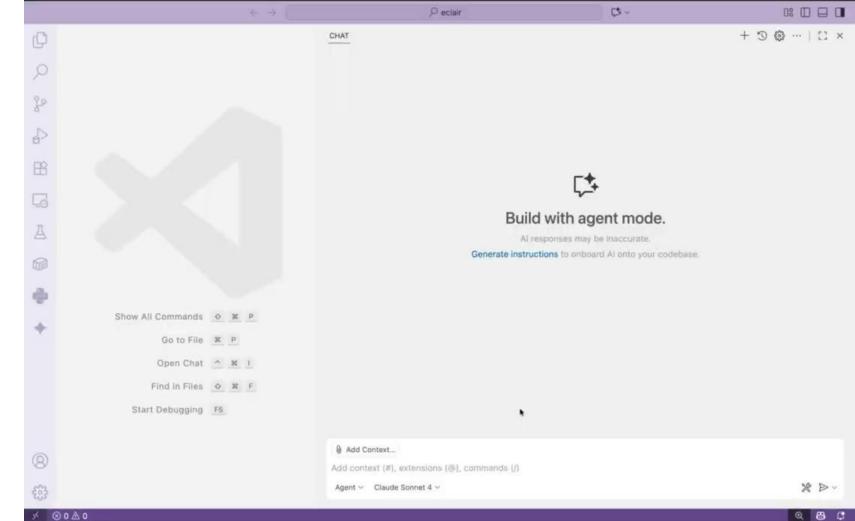


## Eclair: help AI agents work with datasets

- Model-Context Protocol (MCP) Server for datasets
- Croissant format: machine-readable dataset descriptions
- Provides tools to LLMs: search, understand, download datasets



## Demo





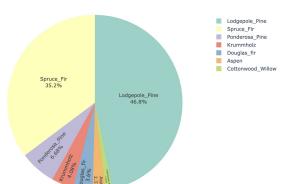
## **Eclair**

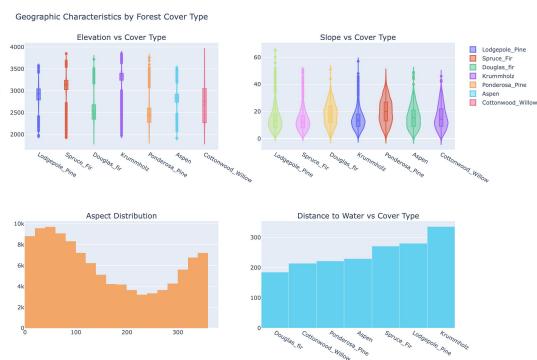
- Loads dataset correctly(!) and can run interactive analyses. **Example**.
- Can also automatically build (simple) models

```
# Load the Covertype dataset from OpenML
dataset_id = 180  # Covertype dataset ID

print("Loading Covertype dataset from OpenML...")
openml_dataset = openml.datasets.get_dataset(dataset_id)
X, y, _, _ = openml_dataset.get_data(target=openml_dataset.

# Convert to pandas DataFrame
df = pd.DataFrame(X)
df['Cover_Type'] = y
```



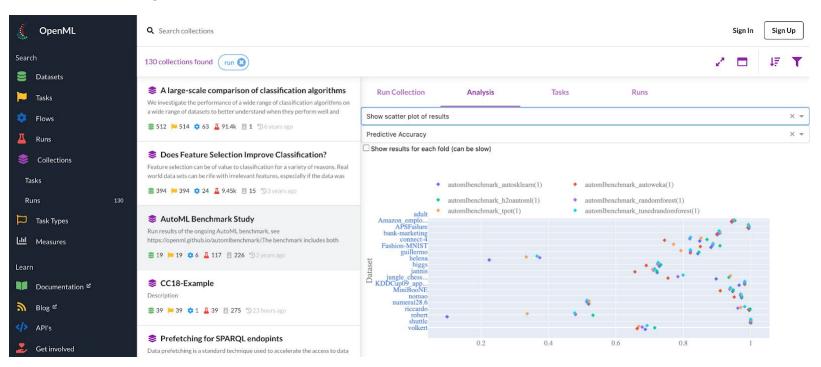


# Benchmark design

- Benchmarking suites
  - Start with a large set of datasets (e.g. OpenML)
  - Define strict set of constraints
  - Retrieve and test models on all matching datasets
  - Gather results from different researchers in a central place (e.g. OpenML)
- Offers a way to really use benchmark suites and converge to well-defined accepted suites
- Are meant to be dynamic: evolve with new datasets joining over time

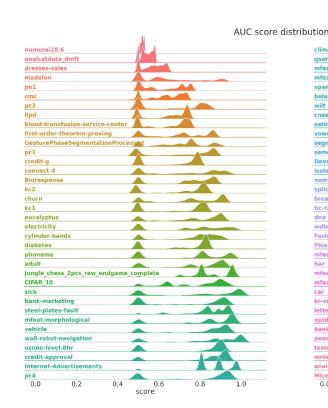
# Benchmark design

## Benchmarking suites



## Benchmark design

- Example: OpenML-CC18
  - meant to be practical
- Classification only
- 72 datasets
- Contain missing values and categorical features
- Medium-sized (500-100000 observations, <5000 features after one-hot-encoding)</li>
- Not unbalanced
- No groups/block/time dependencies
- No sparse data
- Some more subjective criteria (see paper)
- 3.8 million results:







Deadline 10th of October. Submission is abstract-only.

https://sites.google.com/view/benchmarking-and-evaluating-ai