# Developing, using, and evaluating the next generation of DESTs to produce AI-Powered LES for Climate & Health

**The DESTINY project**

**Dr. Maria-Inti Metzendorf**

**3 March 2026**

**Period of IPCC AR**

| AR1 | AR2 | AR3 | AR4 | AR5 |
|-----|-----|-----|-----|-----|
| 1,099 | 7,778 | 18,766 | 34,007 | 115,42 |

Number of publications relating to climate change (thousands)

80

60

40

20

0

1990    1995    2000    2005    2010

| | Measures to reduce infection risk at individual level | Measures to identify and isolate those who are infectious or may become infectious | Measures to reduce the number of contacts | Measures to protect the most vulnerable | Travel and border restrictions |
|---|---|---|---|---|---|
| Randomised controlled trial | | | | | |
| Longitudinal | | | | | |
| Cross-sectional | | | | | |
| Ecological | | | | | |
| Modelling | | | | | |

# DESTINY
## Digital Evidence Synthesis Tool Innovation Yielding Improvements in Climate & Health

DESTINY is co-developing a **new generation of digital evidence synthesis tools (DESTs)**

by showcasing the delivery of rigorous **living evidence** in climate and health

**that matters** to policymakers & other evidence users.
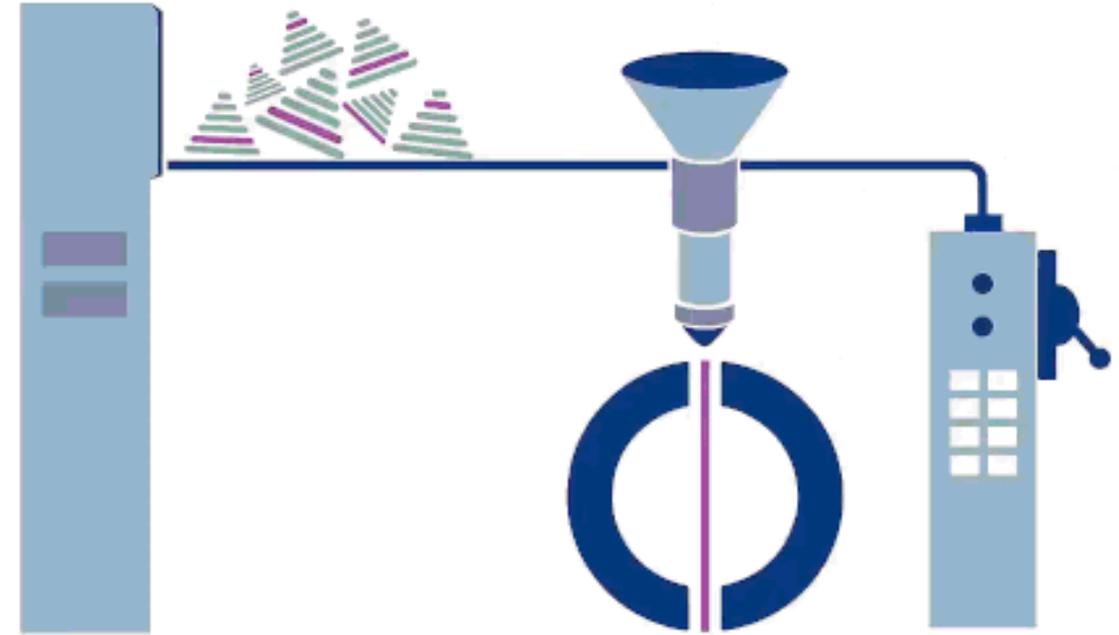
Potsdam Institute for Climate Impact Research

PIK

Leibniz Association

# DESTINY consortium

# DESTINY
## Living evidence for climate & health

- **New DESTs** – explore AI to create faster, cheaper, and more useful evidence synthesis tools (WP2)

- **Responsible use** – ensure safe and responsible DEST applications without compromising standards (WP3)

- **Impact through co-production**
  – work with decision-makers to apply DESTs in key impact cases (WP1, WP4)

- **Mainstreaming DESTs**
  – help users, producers, and funders establish best practices (WP5)

# Six DESTINY Impact Cases

## Why this case selection?

Ensure that impact cases are representative of real-world problems

**Evidence users** Different government levels & organisation types

**Geographies** Different availability of resources & evidence

**Evidence gaps** Different evidence needs for pressing decisions

**Methods** Different types of evidence & synthesis methods

## Case studies

| 1 GLOBAL EVIDENCE | 2 LOCAL EVIDENCE | 3 MORTALITY & MORBIDITY | 4 FOOD SYSTEMS | 5 LOCAL ADAPTATION | 6 GLOBAL SDGs |
|---|---|---|---|---|---|

### Who needs the evidence?

| 1 GLOBAL EVIDENCE | 2 LOCAL EVIDENCE | 3 MORTALITY & MORBIDITY | 4 FOOD SYSTEMS | 5 LOCAL ADAPTATION | 6 GLOBAL SDGs |
|---|---|---|---|---|---|
| International organisations; Evidence curators; National governments | City networks; Local governments | International organisations; NGOs; Local governments; National governments | International organisations; NGOs; Local governments; National governments | City networks; Local governments | International organisations |
| **Current & planned partnerships\*:** WHO, IDRC, OECD, Lancet Countdown, Campbell, Cochrane | **Current & planned partnerships\*:** ICLEI, C40, CDP, GLA | **Current & planned partnerships\*:** WHO member states, EH!WOZA, NWRA, SECTION27 | **Current & planned partnerships\*:** WHO member states, UK-EA, EH!WOZA, SECTION27, WRC | **Current & planned partnerships\*:** ICLEI, C40, CDP, GLA | **Current & planned partnerships\*:** The Global SDG Synthesis Coalition, Sustainable Development Solutions Network |
| *Abbreviations in the annex | | | | | |

### Which evidence gap is addressed?

| 1 GLOBAL EVIDENCE | 2 LOCAL EVIDENCE | 3 MORTALITY & MORBIDITY | 4 FOOD SYSTEMS | 5 LOCAL ADAPTATION | 6 GLOBAL SDGs |
|---|---|---|---|---|---|
| Comprehensive global evidence base on climate & health | Comprehensive local evidence base of climate & **health impacts** and **benefits of climate actions** | Effective mitigation and adaptation **actions to reduce mortality and morbidity** | Interventions to advance **uptake of sustainable diets** including barriers and enablers | **Adaptation** strategies to climate-related **heat in cities** | Interventions to **accelerate** progress towards **climate & health related SDGs** |
| | EVIDENCE SCARCITY | | | EVIDENCE SCARCITY | |

### What method is used to address the gap?

| 1 GLOBAL EVIDENCE | 2 LOCAL EVIDENCE | 3 MORTALITY & MORBIDITY | 4 FOOD SYSTEMS | 5 LOCAL ADAPTATION | 6 GLOBAL SDGs |
|---|---|---|---|---|---|
| Living evidence map | Living evidence map; Evidence transfer | Living systematic review | Living systematic review | Living systematic review; Evidence transfer | Living science assessment |
| **Combining bibliometrics with evidence mapping** methodologies to deal with vast amounts of evidence | Automating traditional evidence gap mapping and qualitative synthesis; **focus on evidence transfer** | Automatic **quantitative synthesis** | Automating **mixed methods synthesis** and integration of empirical and modelling data | Automating mixed methods synthesis with **focus on evidence transfer** | Radical automation strategies for **science assessments** focusing on UN evaluation |

### What is the potential impact and use?

| Agenda setting, horizon scanning | Policy design, policy advice, policy advocacy | Policy design, policy learning |
|---|---|---|
| More efficient evidence ecosystems and improved priority and agenda setting in climate & health across in low-, middle-, and high-income countries | Comprehensive, timely and relevant evidence that informs more effective policies to protect people's health and reduce emissions | Accelerated progress towards climate and health related SDGs |

# Living evidence map of the global climate and health literature



**WP4 IMPACT CASES** Showcasing the transformational power of Digital Evidence Synthesis Tools in six communities of practice for the delivery of rigorous and living evidence that matters to evidence users

## Why this case selection?
Ensure that impact cases are representative of real-world problems

| Evidence users | Geographies | Evidence gaps | Methods |
|---|---|---|---|
| Different government levels & organisation types | Different availability of resources & evidence | Different evidence needs for pressing decisions | Different types of evidence & synthesis methods |

## Case studies

| 1 GLOBAL EVIDENCE | 2 LOCAL EVIDENCE | 3 MORTALITY & MORBIDITY | 4 FOOD SYSTEMS | 5 LOCAL ADAPTATION | 6 GLOBAL SDGs |
|---|---|---|---|---|---|

### Who needs the evidence?

| 1 GLOBAL EVIDENCE | 2 LOCAL EVIDENCE | 3 MORTALITY & MORBIDITY | 4 FOOD SYSTEMS | 5 LOCAL ADAPTATION | 6 GLOBAL SDGs |
|---|---|---|---|---|---|
| ■ International organisations<br>■ Evidence curators<br>■ National governments | ■ City networks<br>■ Local governments | ■ International organisations<br>■ NGOs<br>■ Local governments<br>■ National governments | ■ International organisations<br>■ NGOs<br>■ Local governments<br>■ National governments | ■ City networks<br>■ Local governments | ■ International organisations |
| **Current & planned partnerships\*:** WHO, IDRC, OECD, Lancet Countdown, Campbell, Cochrane<br><br>*Abbreviations in the annex | **Current & planned partnerships\*:** ICLEI, C40, CDP, GLA | **Current & planned partnerships\*:** WHO member states, EH!WOZA, NWRA, SECTION27 | **Current & planned partnerships\*:** WHO member states, UK-EA, EH!WOZA, SECTION27, WRC | **Current & planned partnerships\*:** ICLEI, C40, CDP, GLA | **Current & planned partnerships\*:** The Global SDG Synthesis Coalition, Sustainable Development Solutions Network |

### Which evidence gap is addressed?

| 1 GLOBAL EVIDENCE | 2 LOCAL EVIDENCE | 3 MORTALITY & MORBIDITY | 4 FOOD SYSTEMS | 5 LOCAL ADAPTATION | 6 GLOBAL SDGs |
|---|---|---|---|---|---|
| Comprehensive global evidence base on climate & health | Comprehensive local evidence base of climate & **health impacts** and **benefits of climate actions**<br><br>EVIDENCE SCARCITY | Effective mitigation and adaptation **actions to reduce mortality and morbidity** | Interventions to advance **uptake of sustainable diets** including barriers and enablers | **Adaptation** strategies to climate-related **heat in cities**<br><br>EVIDENCE SCARCITY | Interventions to **accelerate** progress towards **climate & health related SDGs** |

### What method is used to address the gap?

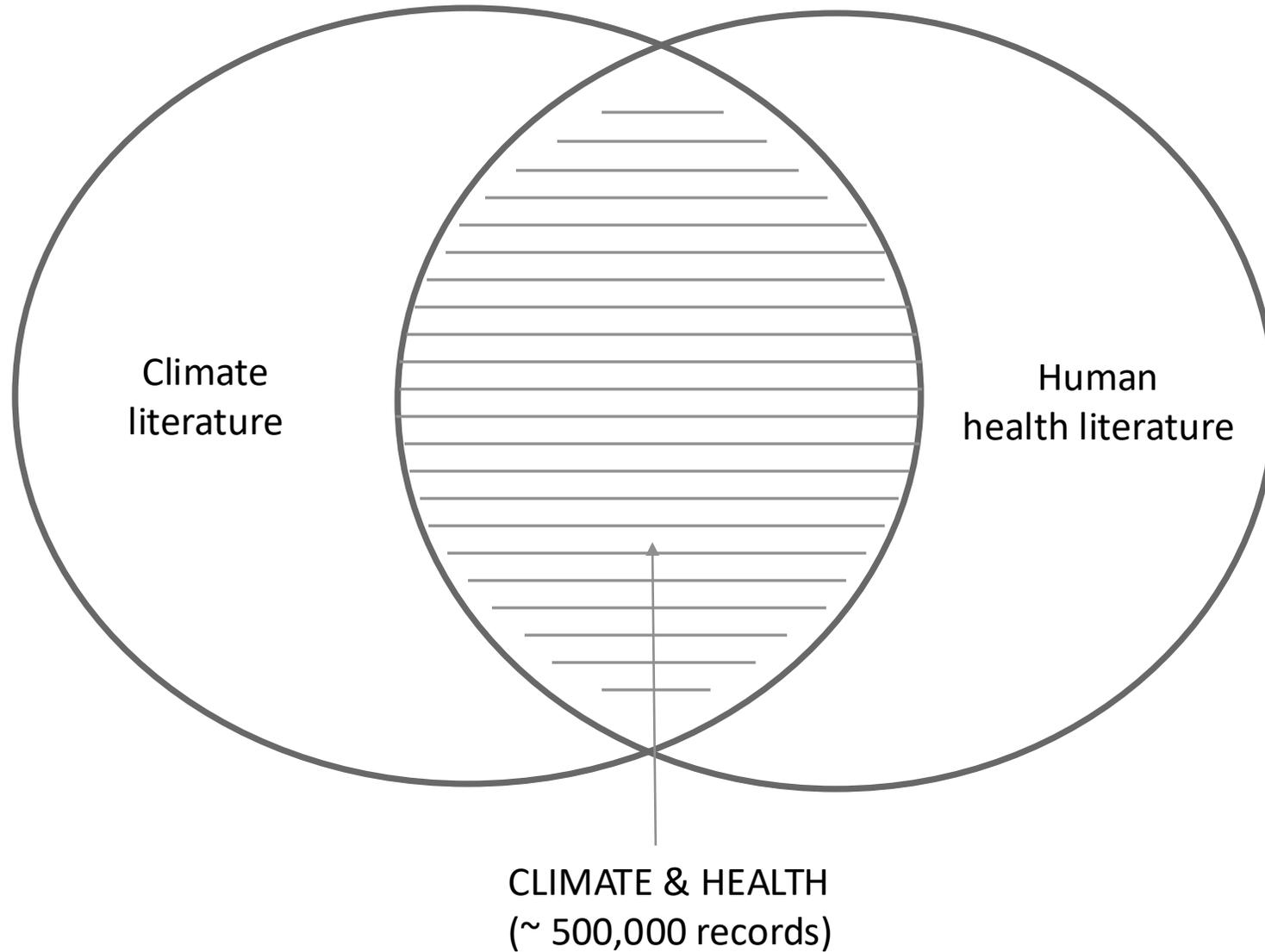| 1 GLOBAL EVIDENCE | 2 LOCAL EVIDENCE | 3 MORTALITY & MORBIDITY | 4 FOOD SYSTEMS | 5 LOCAL ADAPTATION | 6 GLOBAL SDGs |
|---|---|---|---|---|---|
| ■ Living evidence map | ■ Living evidence map<br>■ Evidence transfer | ■ Living systematic review | ■ Living systematic review | ■ Living systematic review<br>■ Evidence transfer | ■ Living science assessment |
| **Combining bibliometrics with evidence mapping** methodologies to deal with vast amounts of evidence | Automating traditional evidence gap mapping and qualitative synthesis; **focus on evidence transfer** | Automatic **quantitative synthesis** | Automating **mixed methods synthesis** and integration of empirical and modelling data | Automating mixed methods synthesis with **focus on evidence transfer** | Radical automation strategies for **science assessments** focusing on UN evaluation |

### What is the potential impact and use?

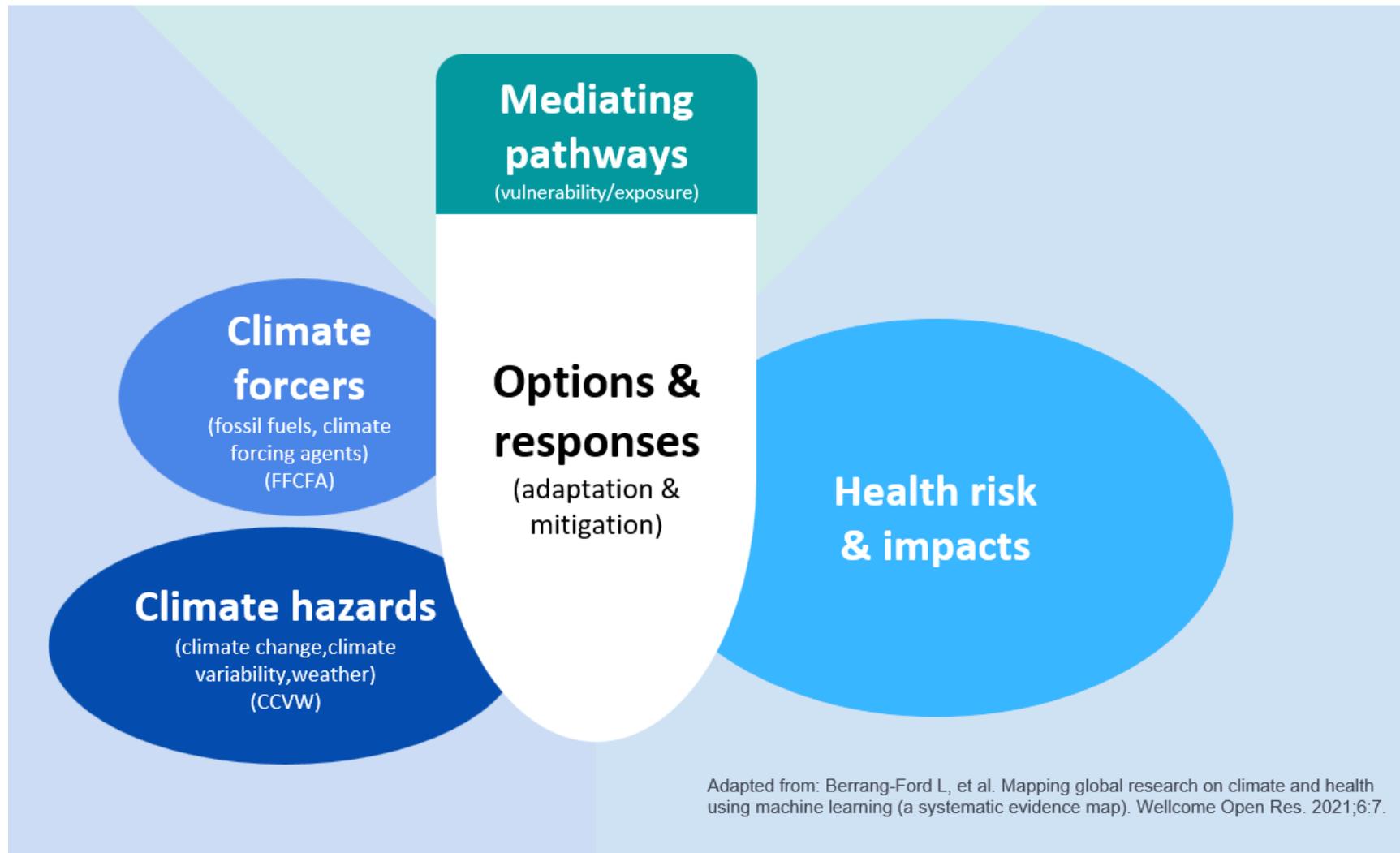| **Agenda setting, horizon scanning** | **Policy design, policy advice, policy advocacy** | **Policy design, policy learning** |
|---|---|---|
| More efficient evidence ecosystems and improved priority and agenda setting in climate & health across in low-, middle-, and high-income countries | Comprehensive, timely and relevant evidence that informs more effective policies to protect people's health and reduce emissions | Accelerated progress towards climate and health related SDGs |

# The scope in simple terms
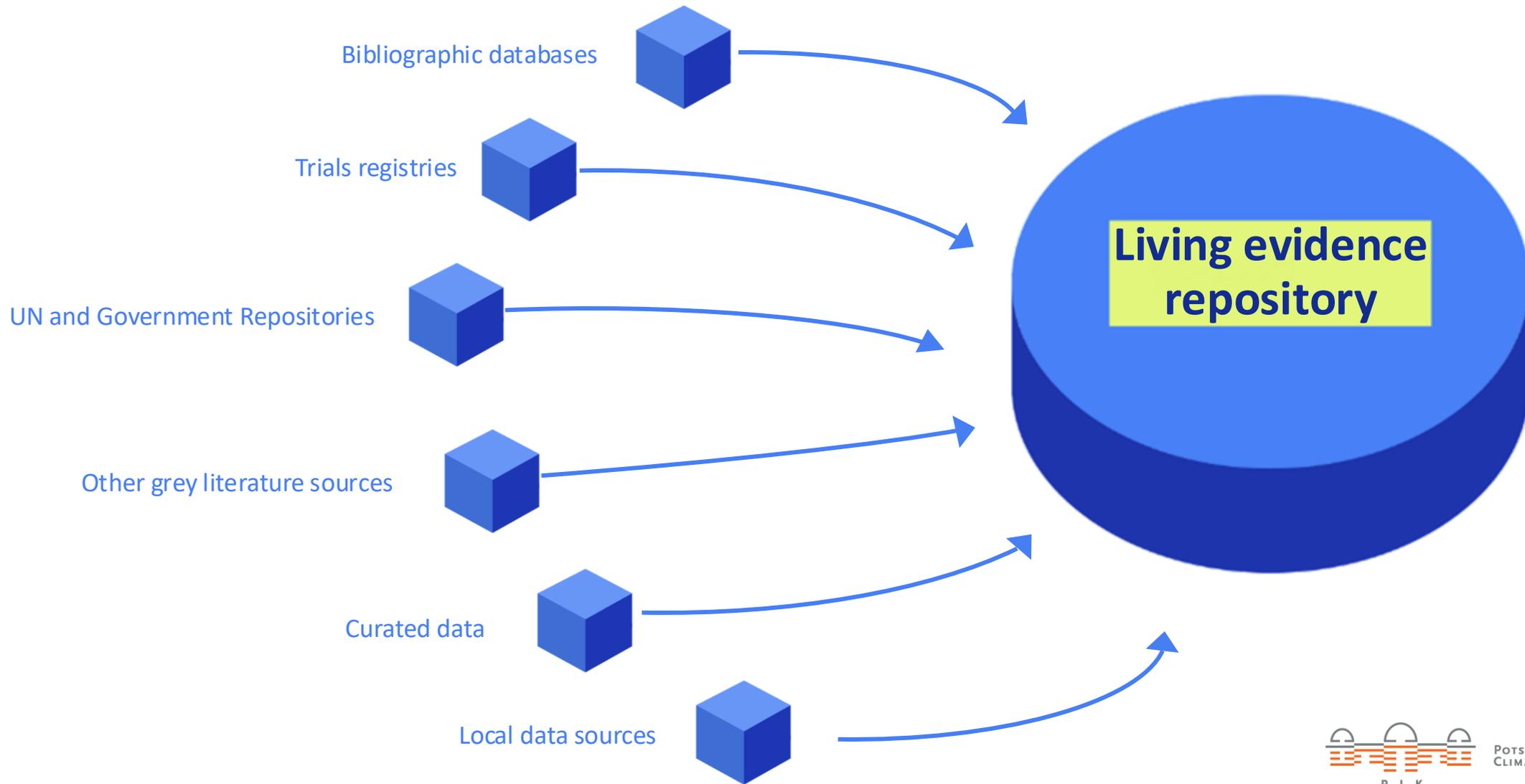
Climate
literature

Human
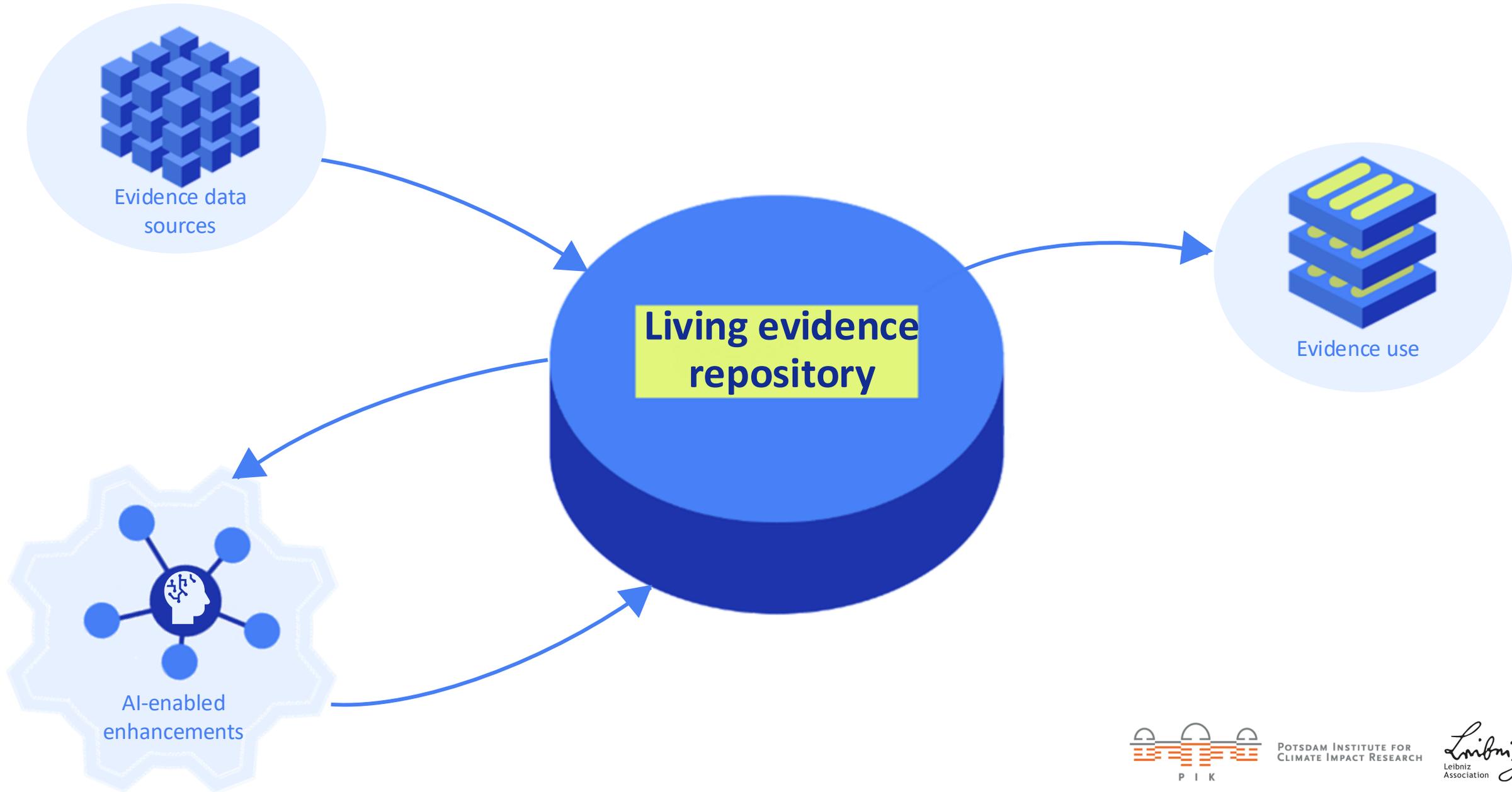health literature

CLIMATE & HEALTH
(~ 500,000 records)

# Key concepts



Mediating pathways (vulnerability/exposure)

Climate forcers (fossil fuels, climate forcing agents) (FFCFA)

Options & responses (adaptation & mitigation)

Climate hazards (climate change, climate variability, weather) (CCVW)

Health risk & impacts

Adapted from: Berrang-Ford L, et al. Mapping global research on climate and health using machine learning (a systematic evidence map). Wellcome Open Res. 2021;6:7.

# Evidence sources/types of studies

# A living data repository for CC & Health

# Taxonomy

**Population** (sub-categories: vulnerability, gender/sex, age)

**Topic** (sub-categories: mitigation, adaptation, impact);

**Climate factors** (sub-categories: climate drivers, extreme weather events, climate-forcing agents, fossil fuels)

**Interventions/measures/responses** (sub-categories: policy, technology/infrastructure, behavior/culture, ecosystem-based, institutional);

**Sectors** (sub-categories: agriculture, fishing and forestry; buildings and housing; education; energy and extractives; finance; healthcare; social protection; industry, trade and services; information and communications technologies; public administration; transportation; water, sanitation and waste management);

**Governance scale** (sub-categories: global, international, national, sub-national, household and individual);

**Geographic location** of affected area, and derive location characteristics, e.g. income category (sub-categories: low, lower middle, upper middle, high);

**Geographic feature** (sub-categories: urban, rural, desert, forest, freshwater, grassland, island, mountain, ocean/coastal, polar, rainforest, valley, temperate, tropical, wetland);

**Exposure** (sub-categories: heat stress, air quality, food supply and safety, water quality and quantity, extreme weather events, vector distribution and ecology, social factors (forced displacement, violence, conflict));

**Actors** (sub-categories: international organizations, government, private sector, civil society, households and individuals);

**Methods** (sub-categories: primary research, evidence synthesis);

**Outcomes** (sub-categories: death/mortality, outpatient visits, hospitalization, and the first three levels of the ICD-11).

POTSDAM INSTITUTE FOR
CLIMATE IMPACT RESEARCH

PIK

Leibniz
Association

# Evidence products / outputs



Living evidence repository

Open APIs

Discovery interface

Mapping interfaces

Evidence synthesis tools

Data provider

Bulk export

# DESTINY Dashboard

## Whole database overview

| Total Records | Included Records ? | Recently Published Records ? | Most Recent Ingested Publication ? |
|---|---|---|---|
| 17,916,584 | 492,648 | 3,513,245 | 2026 |

## Taxonomy

Choose Type ?

○ ANTD
● ST_tree

> ☐ Intervention
> ☐ Context
> ☐ General exposure
> ☐ Health outcomes

💡 ANTD allows selecting any node, while ST_tree only allows selecting leaf nodes for more precise filtering.

## Yearly Record Distribution

**Records per Year (showing latest 10 years)**



### Sidebar

DESTINY

- 🏠 Home

**REPOSITORY OVERVIEWS**
- Query Search
- Advanced Search
- Dashboard
- DESTINY Taxonomy

**DEMO VISUALISATIONS**
- Demo Page 1: Bar Chart
- Demo Page 2: Pie Chart
- Demo Page 3: Hierarchical V...
- Demo Page 4: Attribute Cros...

⚙️ **Settings Menu**

Select a color theme

viridis

Maximum Number of Items Displayed in Plots

10

Maximum Length of Labels in Plots

Current label length: 50

# Thank you and over to Max...

maria-inti.metzendorf@pik-potsdam.de

# Tools and Evaluation

Enabling systematic reviewers and evidence users to conduct and access trustworthy and responsible AI-assisted evidence synthesis

# DEET

# DEET
## Data Extraction Evaluation Toolkit

- There are **many** platforms which allow you to use LLMs to screen, categorise, or extract data from studies

- There are **none** that are **open source**, that give you complete **control of models and prompts**, and that encourage **good evaluation practice**

# DEET Workflow

# Chunked evaluation workflow



- We formalise an evaluation protocol that
  o strictly separates development/training from testing
  o gives a realistic estimation of available work savings

# Future of DEET

- Initial release in coming weeks
  - Can be used via CLI, or as a python package
- **Self-hosted UI**
- **DESTINY-hosted UI**

# Open operational questions for evaluation

- How many documents do I need in my development sets?

- How do I deal with the fact that "gold-standard" annotations are not solid gold?

- How can I communicate the implications of any evaluation I do on the answer to my review's research question?

# A wider research agenda on the evaluation for AI-assisted systematic reviews

- **Dealing with unreliable human annotations**

- Can we correct disagreements between human-annotated and AI-annotated data without biasing our evaluations of the AI? Can we *not* correct disagreements without biasing our evaluations? How should we approach this?

- Under what conditions do more human annotators improve the fidelity of gold-standard data? Does accuracy always approach 1 with more annotators?

- Can we infer the number of annotators we would need to achieve a given level of fidelity to ground truth data from measures of inter-rater reliability?

- Does a ground-truth exist? How do we deal with this epistemologically? What would superhuman performance mean?

- **Managing uncertainty**

- Are binomial confidence intervals for recall and precisions good baselines for estimating confidence intervals around performance metrics?

- Can confidence intervals be narrowed?
    - By using Bayesian statistics and priors based on performance on similar tasks?
    - By estimating jointly across predicted categories / tasks?
    - Automation and the results of systematic reviews

- How do errors and uncertainties around errors compound across tasks, when automation is used for multiple stages of a systematic review?

- How do errors affect the results of systematic reviews?

- How can we incorporate uncertainty around the accuracy of specific tasks (screening, data extraction, criticial appraisal) into the overall uncertainty in our results?

- Optimal stopping criteria / stopping criteria for living reviews / LLMs for screening / alternatives to prioritised screening

- **Optimising the distribution of labour between humans and machines**

- How can we quantify the costs and benefits of conducting evidence synthesis tasks by hand and using different AI approaches?

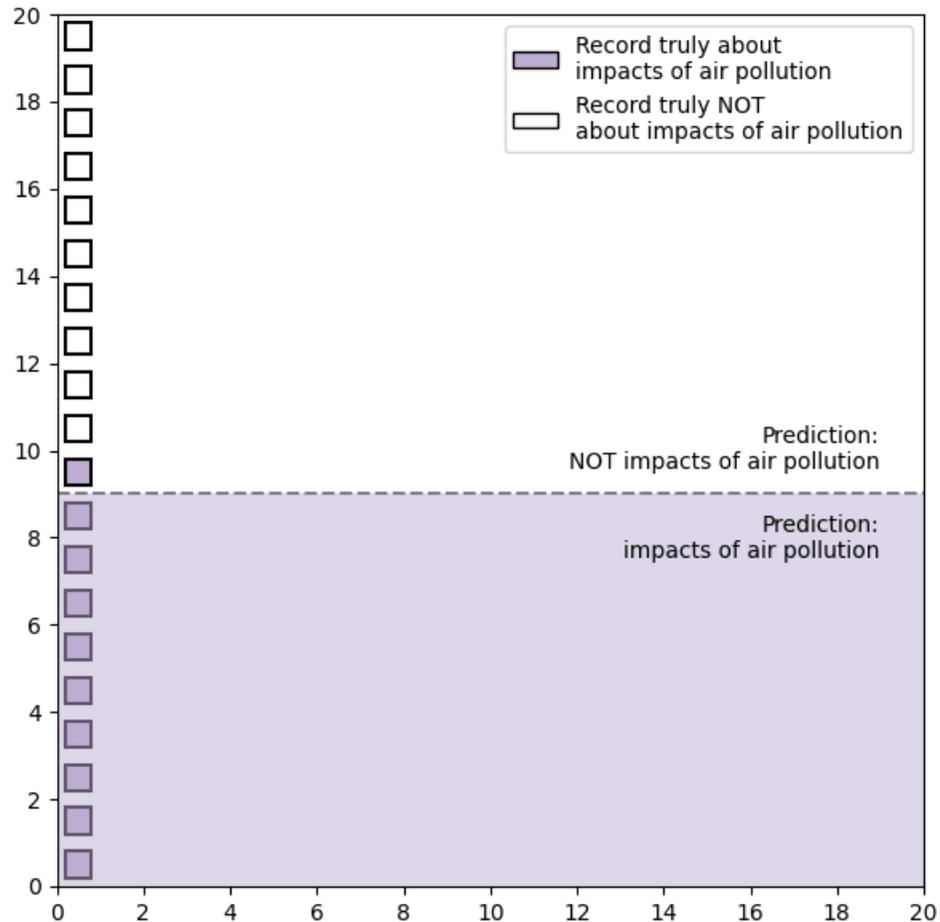- How can we manage trade-offs and allocate resources efficiently?

https://destiny-evidence.github.io/evaluation-book/research-questions/

POTSDAM INSTITUTE FOR
CLIMATE IMPACT RESEARCH

P I K

Leibniz
Association
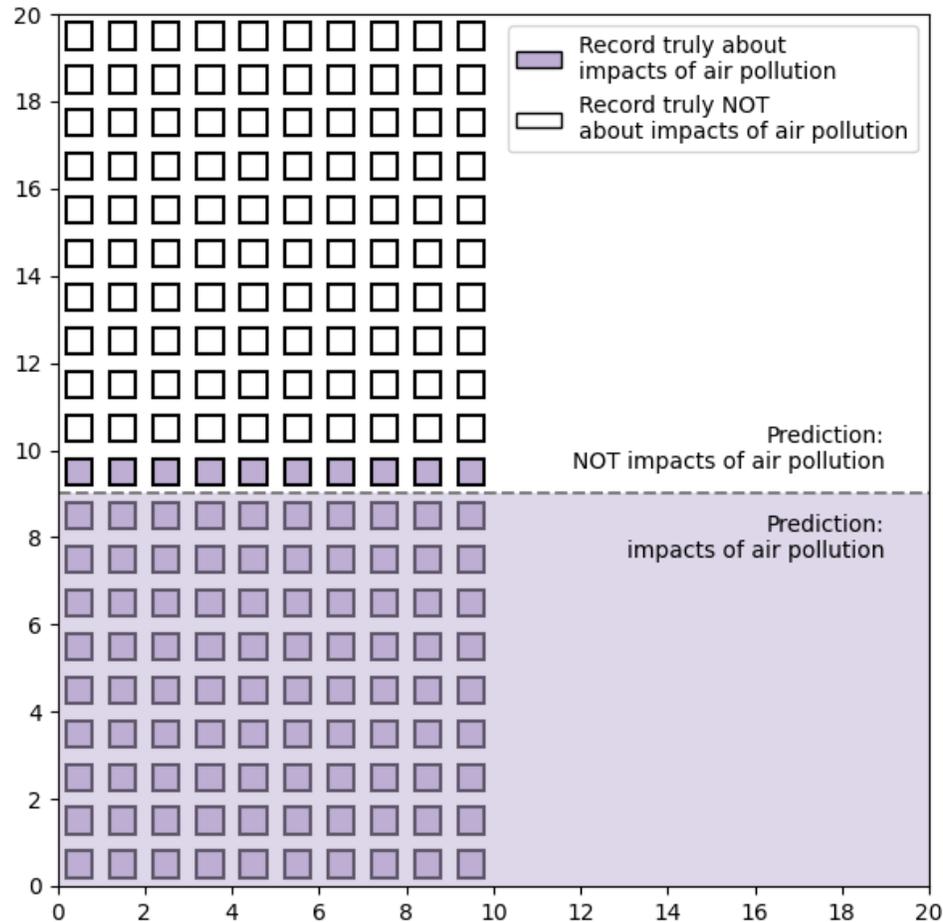
# Uncertainty-aware evaluation

# Uncertainty in evaluation metrics is often unaccounted for



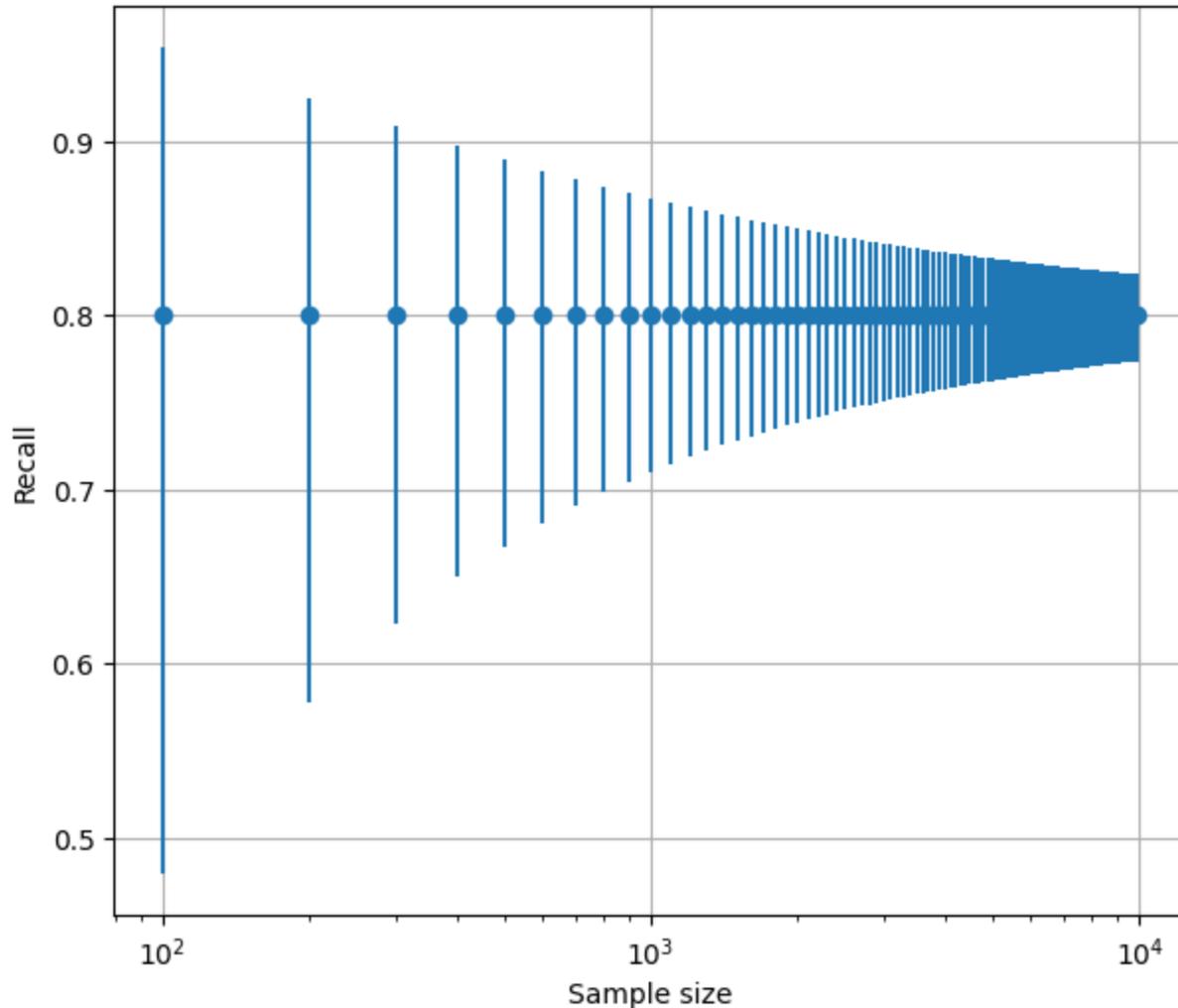- **Given this evaluation data, we estimate our recall to be 90%**

# Uncertainty in evaluation metrics is often unaccounted for



- **Given this evaluation data, we estimate our recall to be 90%**
- **With more evaluation data, we still estimate 90%, but our confidence interval should shrink**
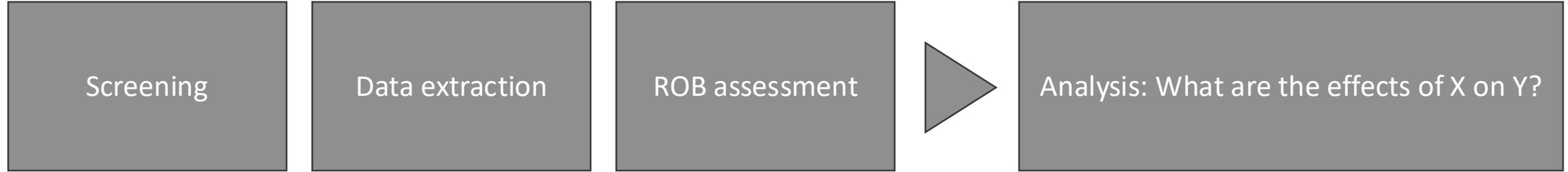
# Binomial confidence intervals for evaluation metrics
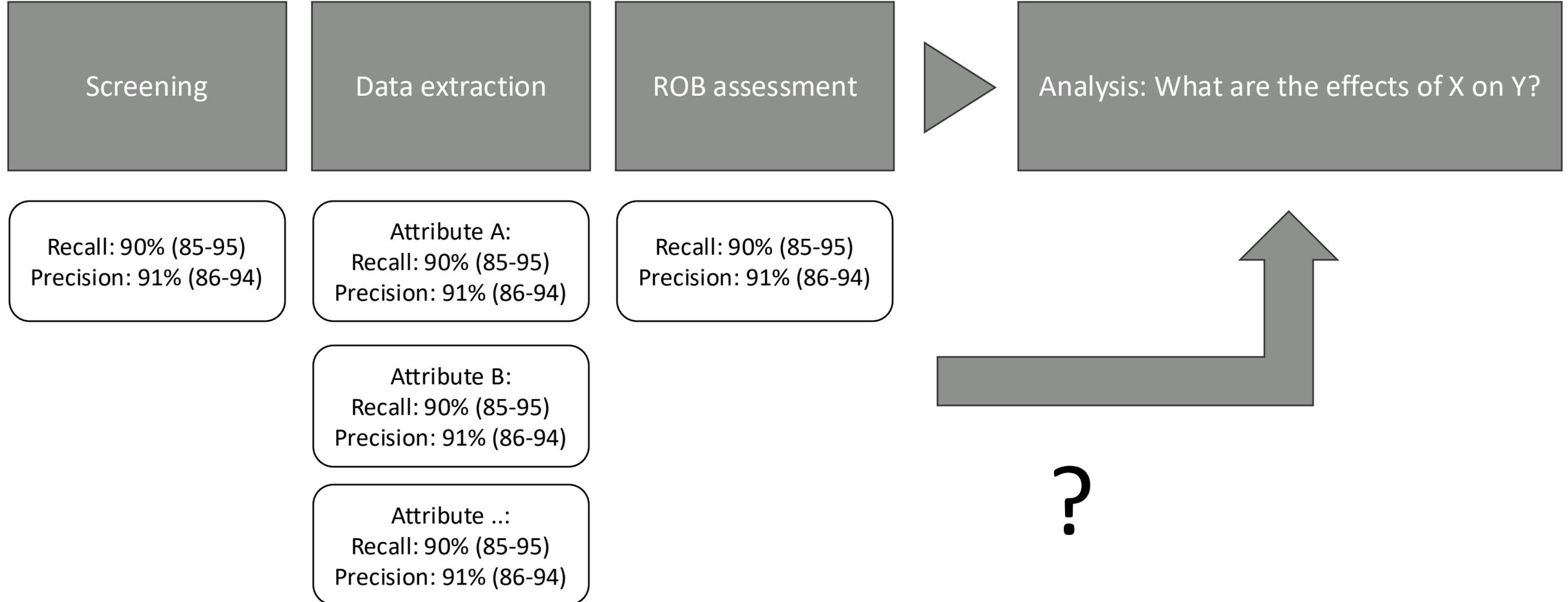


True recall = 0.8, prevalence = 0.1

- We can use binomial confidence intervals (as in diagnostic test accuracy) to estimate proportions

- These are large!

- Can they be narrowed with
  - Bayesian statistics
  - Bootstrapping
  - Joint estimation across classes

POTSDAM INSTITUTE FOR CLIMATE IMPACT RESEARCH

Leibniz Association

# When we use AI across the systematic review process, how do errors and uncertainties compound?

| Screening | Data extraction | ROB assessment | Analysis: What are the effects of X on Y? |
|:---:|:---:|:---:|:---:|

# When we use AI across the systematic review process, how do errors and uncertainties compound?

**Screening**

**Data extraction**

**ROB assessment**

**Analysis: What are the effects of X on Y?**

Recall: 90% (85-95)
Precision: 91% (86-94)

Attribute A:
Recall: 90% (85-95)
Precision: 91% (86-94)

Recall: 90% (85-95)
Precision: 91% (86-94)

Attribute B:
Recall: 90% (85-95)
Precision: 91% (86-94)

Attribute ..:
Recall: 90% (85-95)
Precision: 91% (86-94)

**?**

**Meaningful communication to evidence users**

Ultimately we do not care if we are 95% confident we achieve 95% recall in screening, 97% recall on average across attributes in data extraction, …

We want to know how potential errors introduced through automation (OR through human error) affect our findings.

How can we propagate these errors into the confidence interval of our findings?
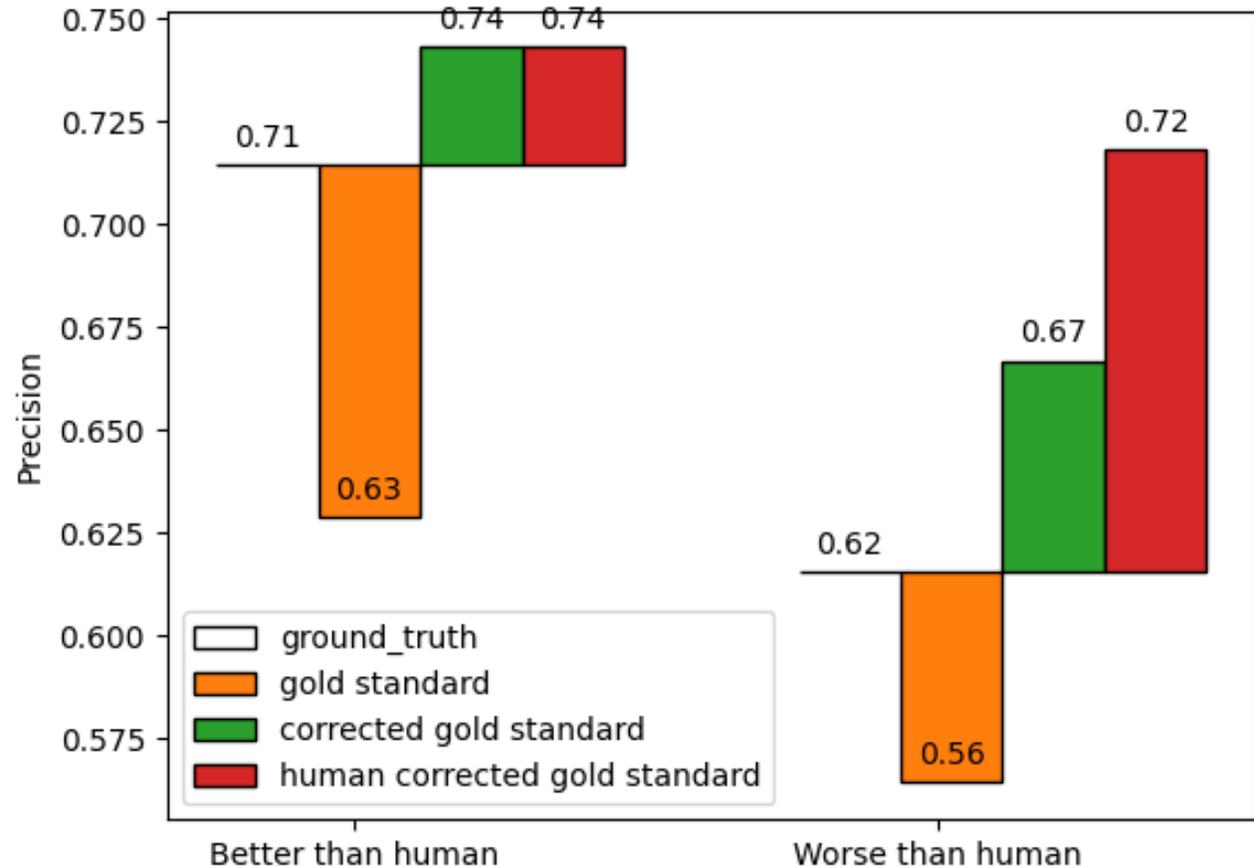
# Dealing with imperfect human annotations

## What happens if we "correct" disagreements between human annotations and AI predictions?

When human and AI annotations disagree, the evaluation metrics for our AI pipeline go down.
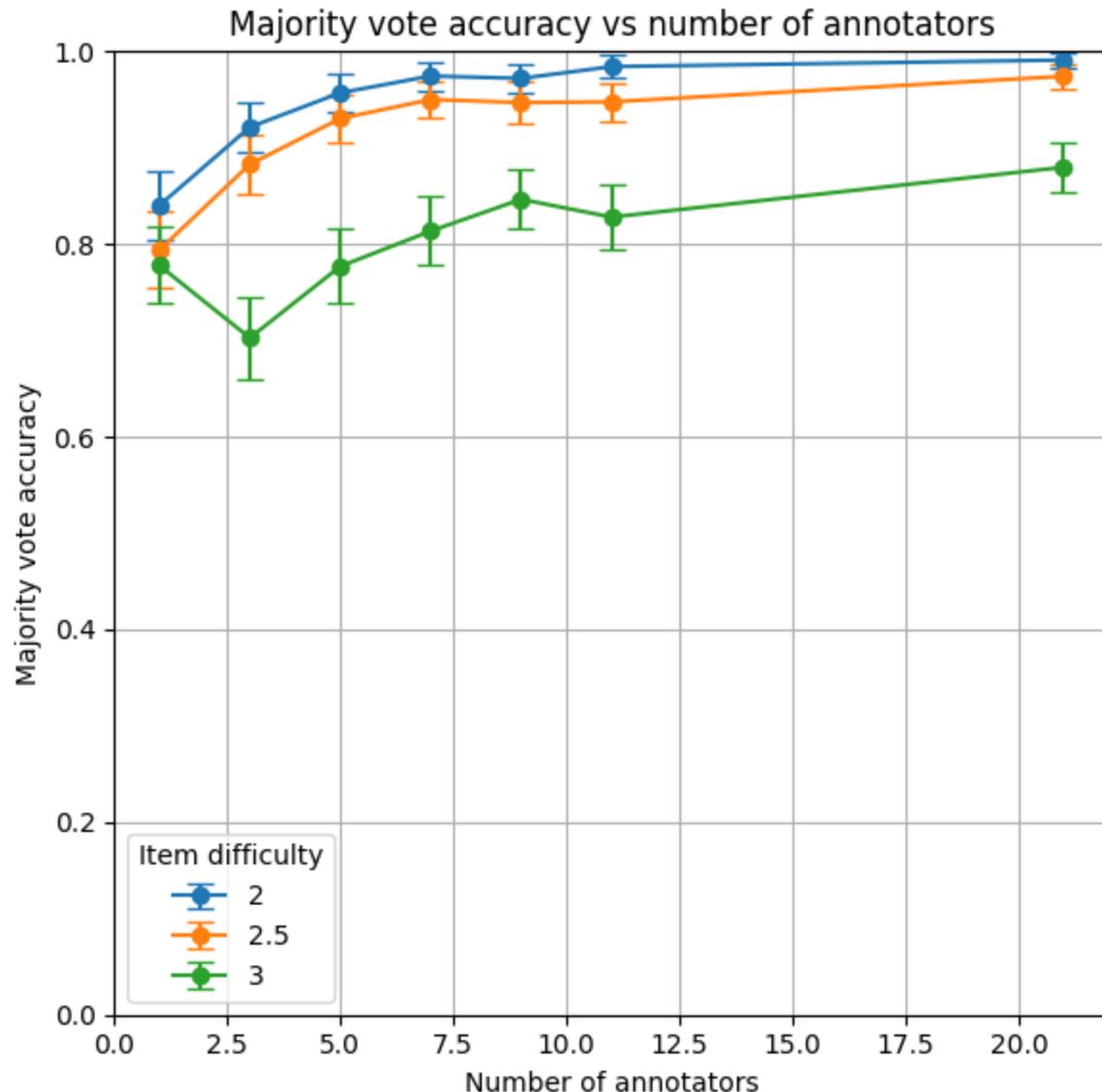
But are these true disagreements? Or were they coding errors?

# What happens if we "correct" disagreements between human annotations and AI predictions?



- **Not correcting errors means we underestimate AI performance**
- **Correcting them means we over-estimate performance (especially if AI is working less well than average human)**
  - **This is because we do not catch those cases where the AI and human erred**

# How many human annotations do we need?



Majority vote accuracy vs number of annotators

- **Two human annotations, with disagreements resolved by a 3rd is a convenient heuristic**
- **AI gives us the chance to shift how we allocate human labour**
- **Under non-correlation of errors, accuracy approaches 1 as we add annotators (Condorcet's jury theorem)**
- **But can we estimate the utility of an additional annotator, given observed inter-rater reliability**

# An optimal distribution of labour between humans and AI

## An optimal distribution of labour between humans and AI

If we make progress on the questions above, operational questions around organising work and the distribution of work between humans and AI start to become answerable

We can move towards an optimal distribution of labour between humans and AI that is likely to improve, rather than degrade, the quality and integrity of systematic reviews

# Thank you

max.callaghan@pik-potsdam.de