



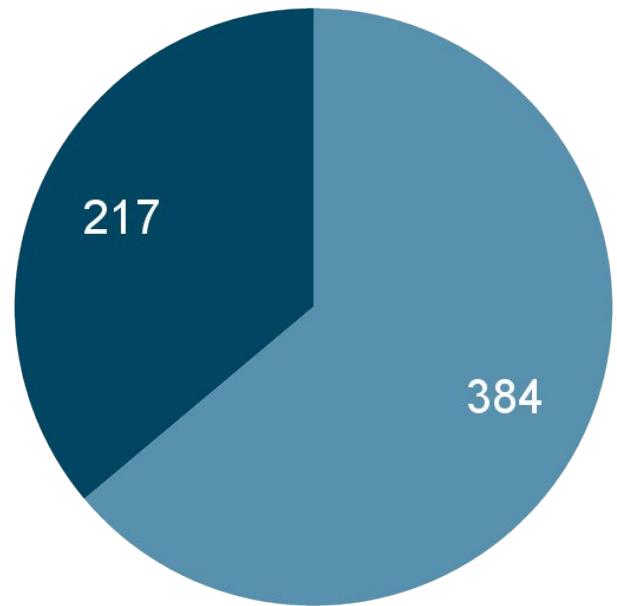
Reproduzierbarkeit(skrise) in den Daten-Wissenschaften

Forschungssymposium des Profilbereichs „Digitale Arbeit und Lebenswelten“
5. Mai 2023

Timo Breuer
Institut für Informationsmanagement,
Fakultät für Informations- und Kommunikationswissenschaften

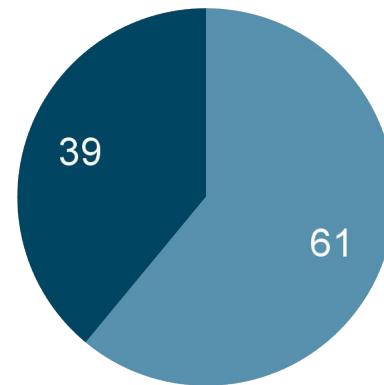
Reproducibility of Published Research Articles

Computer Science

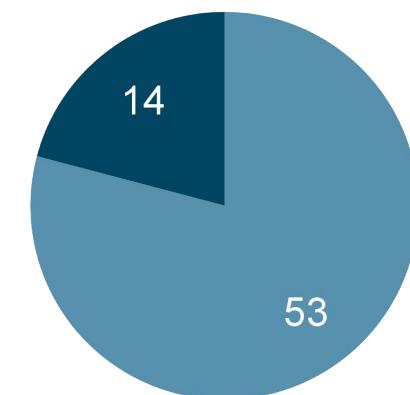


● Failed ● Reproducible

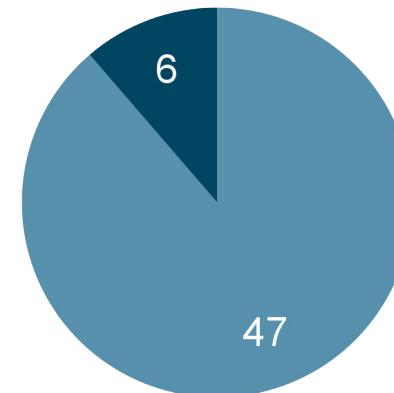
Psychology



Pharmacology



Cancer Studies



Repeatability in Computer Systems Research, Collberg and Proebsting; Commun. ACM; 2015

Estimating the Reproducibility of Psychological Science, Open Science Collaboration; Science; 2015

Believe it or not: How Much can we Rely on Published Data on Potential Drug Targets?; Prinz et al.; Nature Reviews on Drug Discovery, 2011

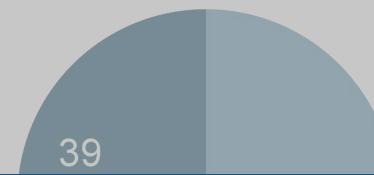
Raise Standards for Preclinical Cancer Research, Begley and Ellis; Nature; 2012

Reproducibility of Published Research Articles

Computer Science

Psychology

Pharmacology



Is there a **reproducibility crisis**?

- 70% out of 1,500 scientists failed to reproduce another researcher's experiment
- 50% out of 1,500 scientists failed to reproduce their own experiment

1,500 Scientists Lift the Lid on Reproducibility, Baker, Nature, 2016

● Failed ● Reproducible

Repeatability in Computer Systems Research, Collberg and Proebsting; Commun. ACM; 2015

Estimating the Reproducibility of Psychological Science, Open Science Collaboration; Science; 2015

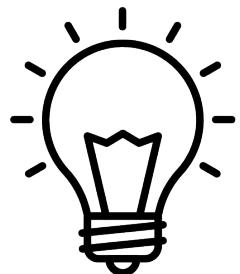
Believe it or not: How Much can we Rely on Published Data on Potential Drug Targets?; Prinz et al.; Nature Reviews on Drug Discovery, 2011

Raise Standards for Preclinical Cancer Research, Begley and Ellis; Nature; 2012

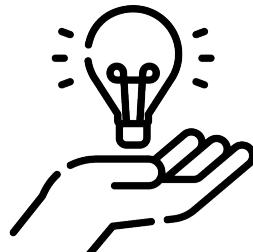
Overview



Common Issues



Countermeasures

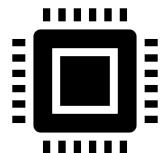


Our Contributions

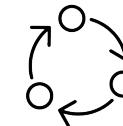
Common Issues and Countermeasures



Platform



Implementation



Method



Research goal



Actor



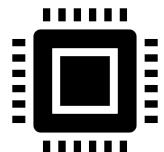
Data



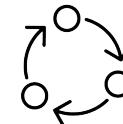
Common Issues and Countermeasures



Platform



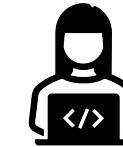
Implementation



Method



Research goal



Actor



Data

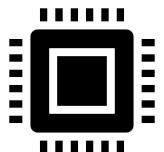


- Missing hardware, kernel, or operating system
- No source code, dependency management

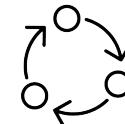
Common Issues and Countermeasures



Platform



Implementation



Method



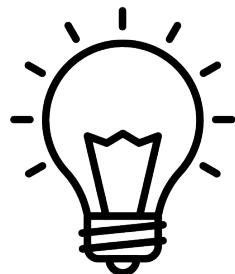
Research goal



Actor



Data



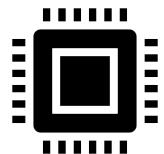
- Missing hardware, kernel, or operating system
- No source code, dependency management
- Virtualization or containerization
- Open-source code



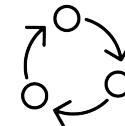
Common Issues and Countermeasures



Platform



Implementation



Method



Research goal



Actor



Data

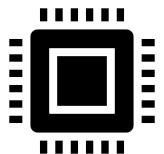


- Missing statistical evidence / low statistical power
- Leaderboard chasing (data dredging / p-hacking / fraud)
- Cognitive biases

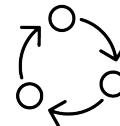
Common Issues and Countermeasures



Platform



Implementation



Method



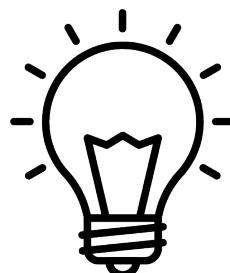
Research goal



Actor



Data

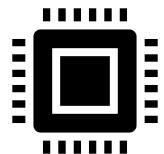
- 
- Missing statistical evidence / low statistical power
 - Leaderboard chasing (data dredging / p-hacking / fraud)
 - Cognitive biases

- 
- Research standards and guidelines, journal policies
 - Common task framework / shared tasks
 - Open reviews

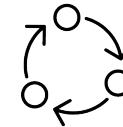
Common Issues and Countermeasures



Platform



Implementation



Method



Research goal



Actor



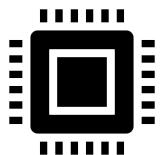
Data

- 
- Separation between code and data
 - Private / closed / paywalled datasets
 - Leakage
 - Data bias
- 

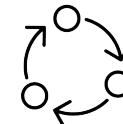
Common Issues and Countermeasures



Platform



Implementation



Method



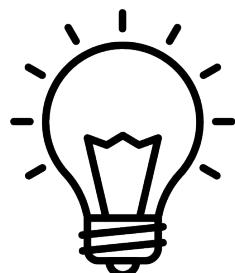
Research goal



Actor



Data



- Separation between code and data
- Private / closed / paywalled datasets
- Leakage
- Data bias

- Open data
- Data archival
- Data citation principles
- Datasheets



F

indable

A

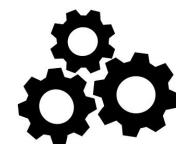
ccessible

I

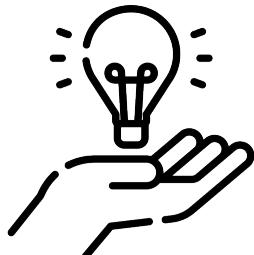
nteroperable

R

eusable



What About Information Retrieval?



A screenshot of a Google search results page for "restaurant südstadt". The top navigation bar shows "Google" and "Sign in". Below the search bar, there are links for "All", "Maps", "Images", "News", "Videos", "Shopping", and "B". The main content area is titled "Places" and shows a map of the Südstadt area in Cologne, Germany, with several restaurant pins. Below the map, three restaurants are listed with their names, ratings, prices, and descriptions. Each listing includes a small thumbnail image of the restaurant's interior.

A screenshot of a Google search results page for "restaurant südstadt". The top navigation bar shows "Google" and a search bar with "restaurant südstadt". Below the search bar, there are links for "X", "Camera", and "Search". The main content area is titled "Orte :". It shows a map of the Südstadt area with restaurant locations marked. Below the map, there are three cards for restaurants: "frau maher", "Jakob Gyros Köln Südstadt", and "HARDY KUGEL". Each card includes a thumbnail image, the restaurant's name, its rating (from 4.5 to 4.6 stars), price range (€€), cuisine type (Restaurant, Gyros, French), address, opening status (closed/open at 12:00/17:00), and delivery information. A "Weitere Orte →" button is at the bottom right.

What About Information Retrieval?



System A

q1	q2	q3
1	11	21
2	12	22
3	13	23
4	14	24
5	15	25
6	16	26
7	17	27
8	18	28
9	19	29
10	20	30

System B

q1	q2	q3
25	11	21
23	12	326
1	13	23
24	14	24
105	115	327
106	116	329
107	17	27
108	18	28
9	19	29
10	20	30

System C

q1	q2	q3
1	220	21
402	219	326
3	218	23
404	217	524
405	16	625
406	15	426
497	14	627
438	13	728
429	12	929
410	11	430

P@10 0.4 0.6 0.5

Avg. P@10 0.5000

0.5 0.4 0.6

0.5000

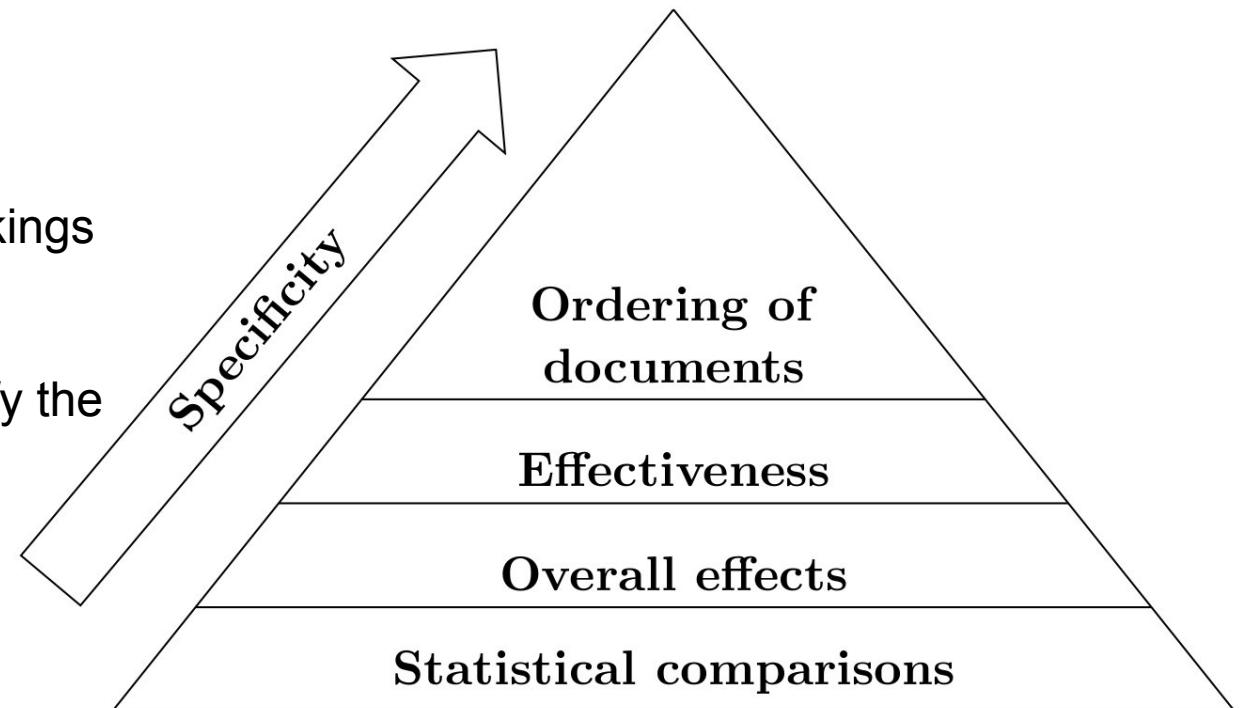
0.2 1.0 0.3

0.5000

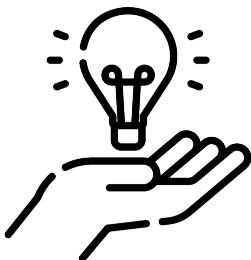
Our Contributions

Reproducibility measures

- Analysis of **reimplementations**
- **Comparison** of reimplemented rankings to the original results
- **Framework of measures** to quantify the degree of reproducibility
- Different **levels of fidelity**



Our Contributions



System A

q1	q2	q3
1	11	21
2	12	22
3	13	23
4	14	24
5	15	25
6	16	26
7	17	27
8	18	28
9	19	29
10	20	30

System B

q1	q2	q3
25	11	21
23	12	326
1	13	23
24	14	24
105	115	327
106	116	329
107	17	27
108	18	28
9	19	29
10	20	30

System C

q1	q2	q3
1	220	21
402	219	326
3	218	23
404	217	524
405	16	625
406	15	426
497	14	627
438	13	728
429	12	929
410	11	430

P@10 0.4 0.6 0.5

Avg. P@10 0.5000

RMSE 0.0000

RBO 1.0000

KTU 1.000

0.5 0.4 0.6

0.5000

0.1414

0.5663

0.4370

0.2 1.0 0.3

0.5000

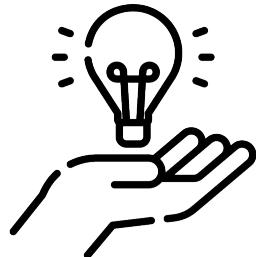
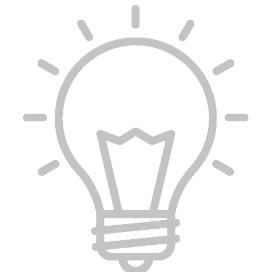
0.2828

0.3984

0.0815

Our Contributions

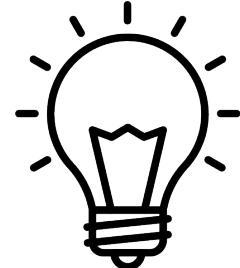
Metadata scheme



Wrap Up



- **Reproducibility issues in data science are manifold!**
- **Countermeasures** and solutions are **diverse**
 - Technological, organizational, social, ethical
- Our **contributions** to reproducible information retrieval
 - **Reactive support** by reproducibility measures
 - **Proactive support** by metadata scheme



Wrap Up

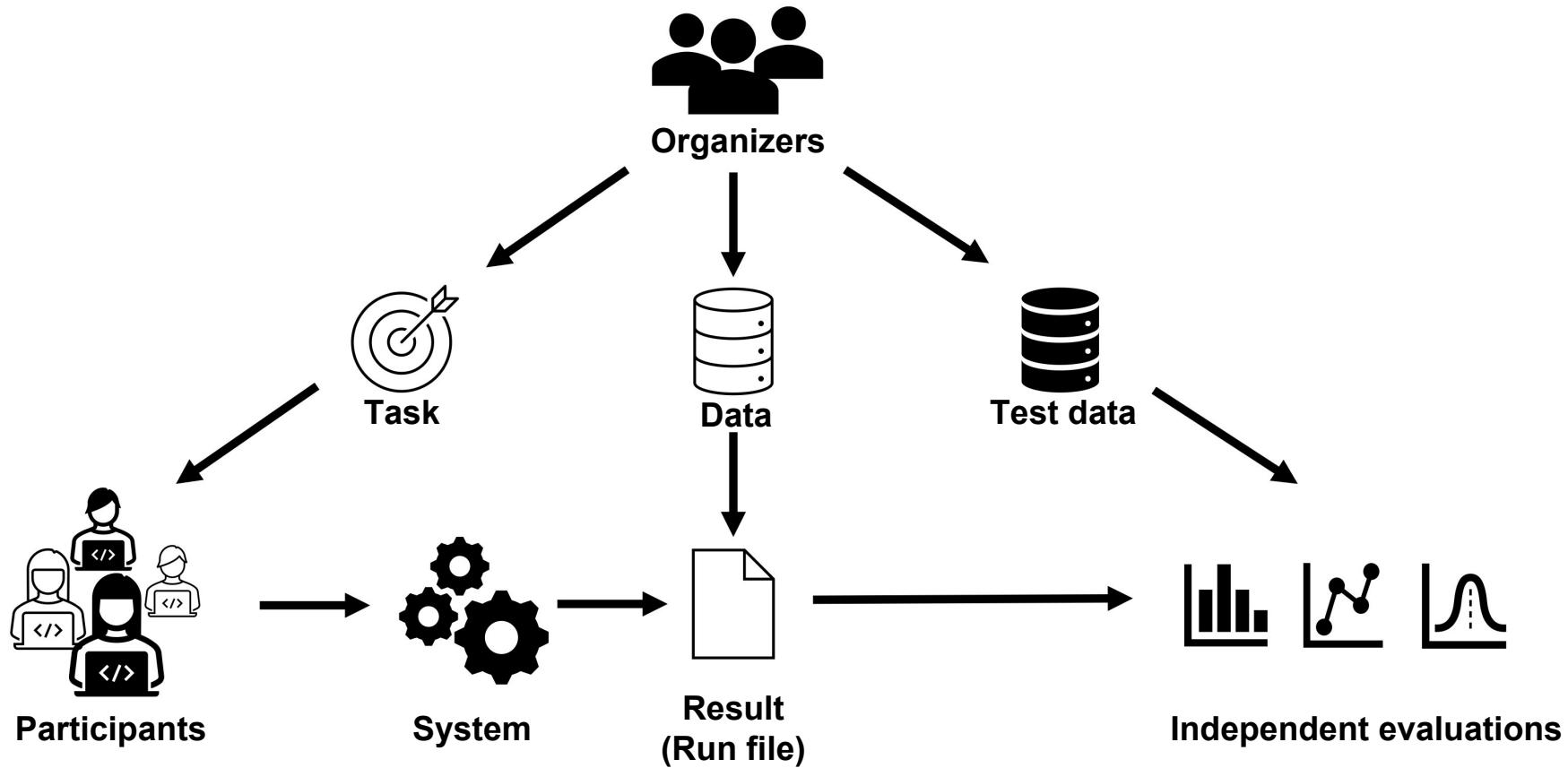
- Re
- Co
- Ou
- :

Thank you!

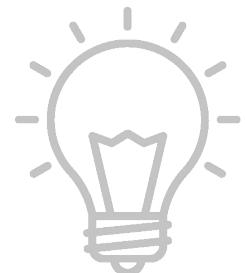
retrieval



Common task framework



Common task framework



IR is deeply rooted in experimentation

- Common evaluation protocols (evaluation measures of relevance)
- Common test collections (reusable data)
- Common data outputs (run files as de-facto standard)

IR conferences host dedicated reproducibility tracks

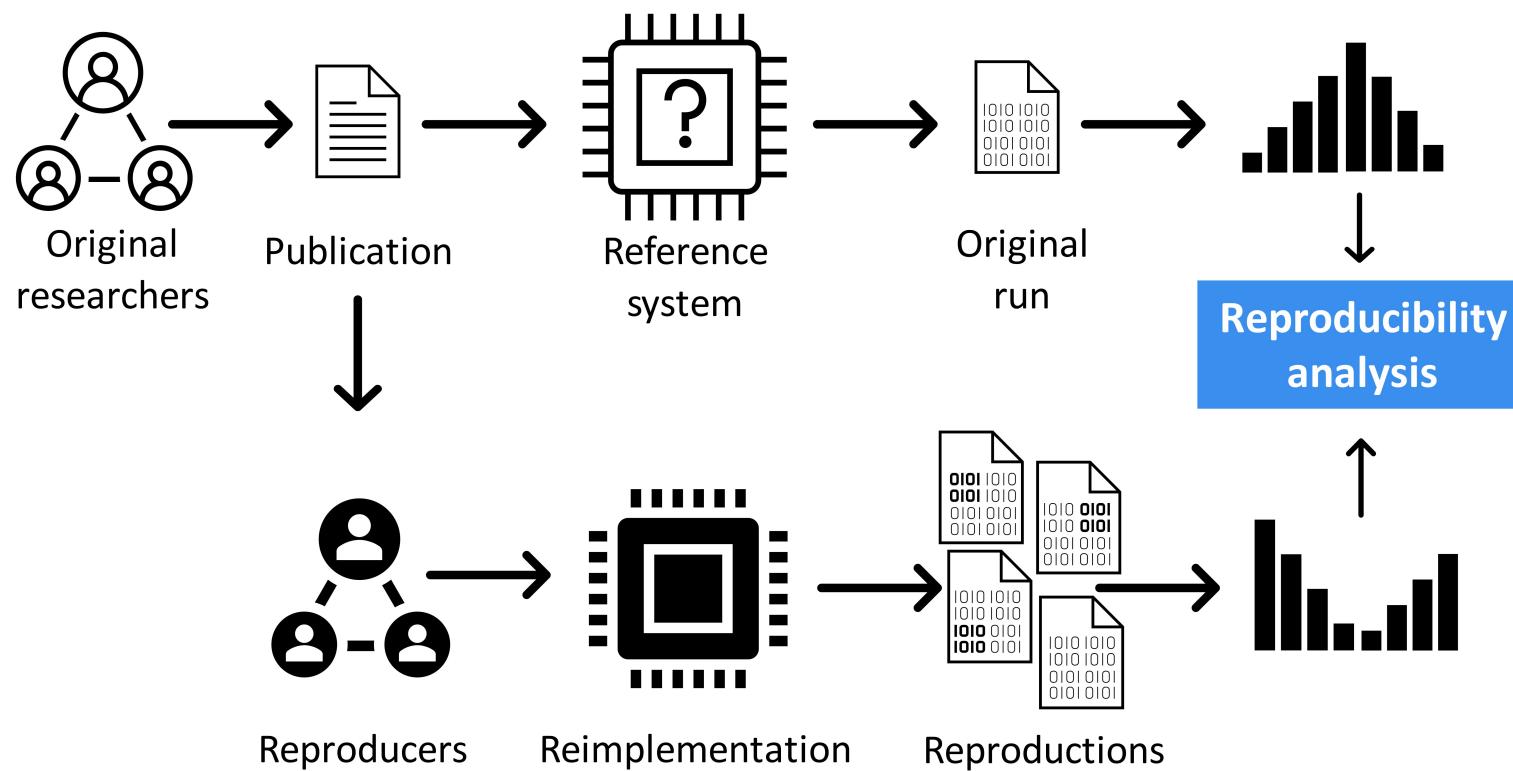
- ECIR since 2015
- SIGIR since 2022

IR community engages in several reproducibility initiatives

- Workshops (OSIRRC@SIGIR, CENTRE@CLEF)
- ACM SIGIR Artifact Badging
<https://sigir.org/general-information/acm-sigir-artifact-badging/>

Our Contributions

Reproducibility analysis



Our Contributions

Metadata scheme



Supporting reproducibility by metadata annotations

- PRIMAD is the logical plan
- Annotation of TREC run files
- YAML formatted header as comment
- Focus on extensibility
- Public resource hosted on <https://www.ir-metadata.org/>

Actor?

