



UNIVERSITÄT
LEIPZIG

Medizinische Fakultät



Universitätsklinikum
Leipzig

Medizin ist unsere Berufung.

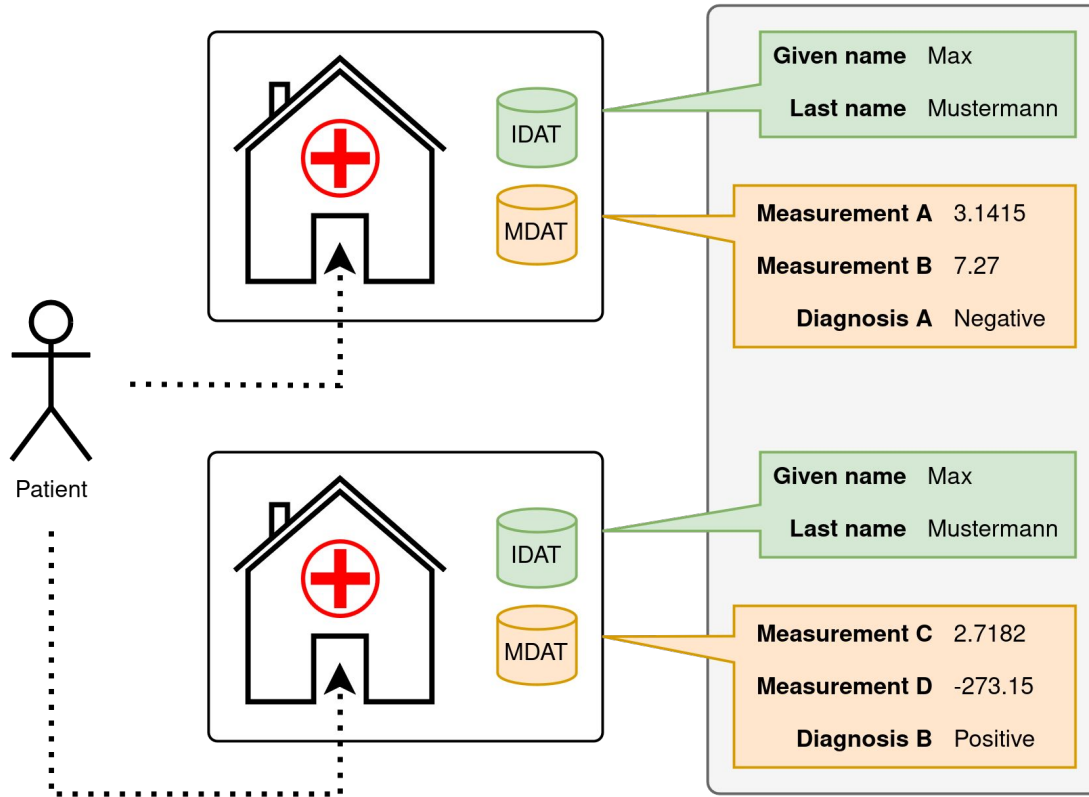
8th Freiburger PhD Conference

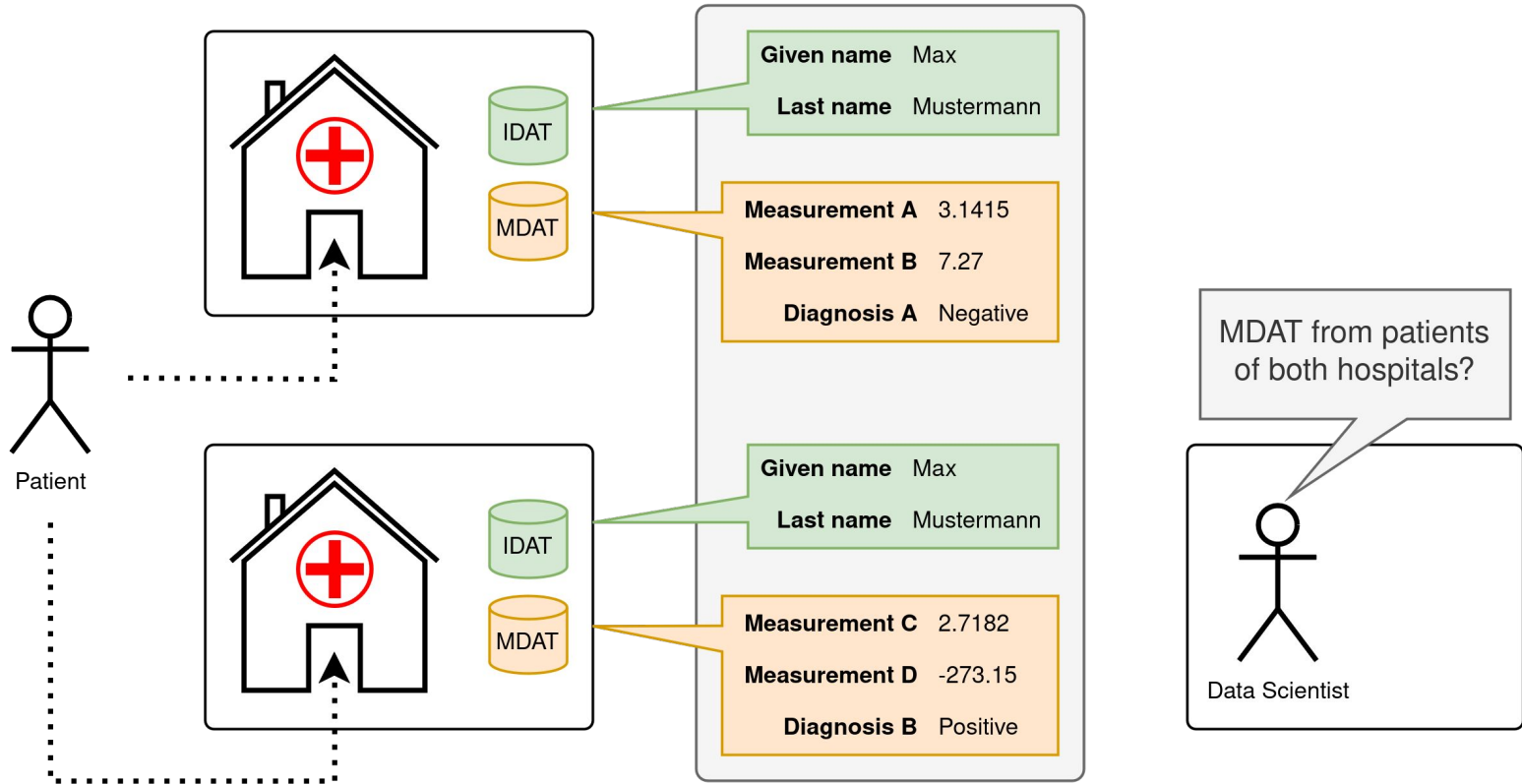
Gecko: Generation and Mutation of Realistic Identification Data at Scale for Record Linkage Evaluation

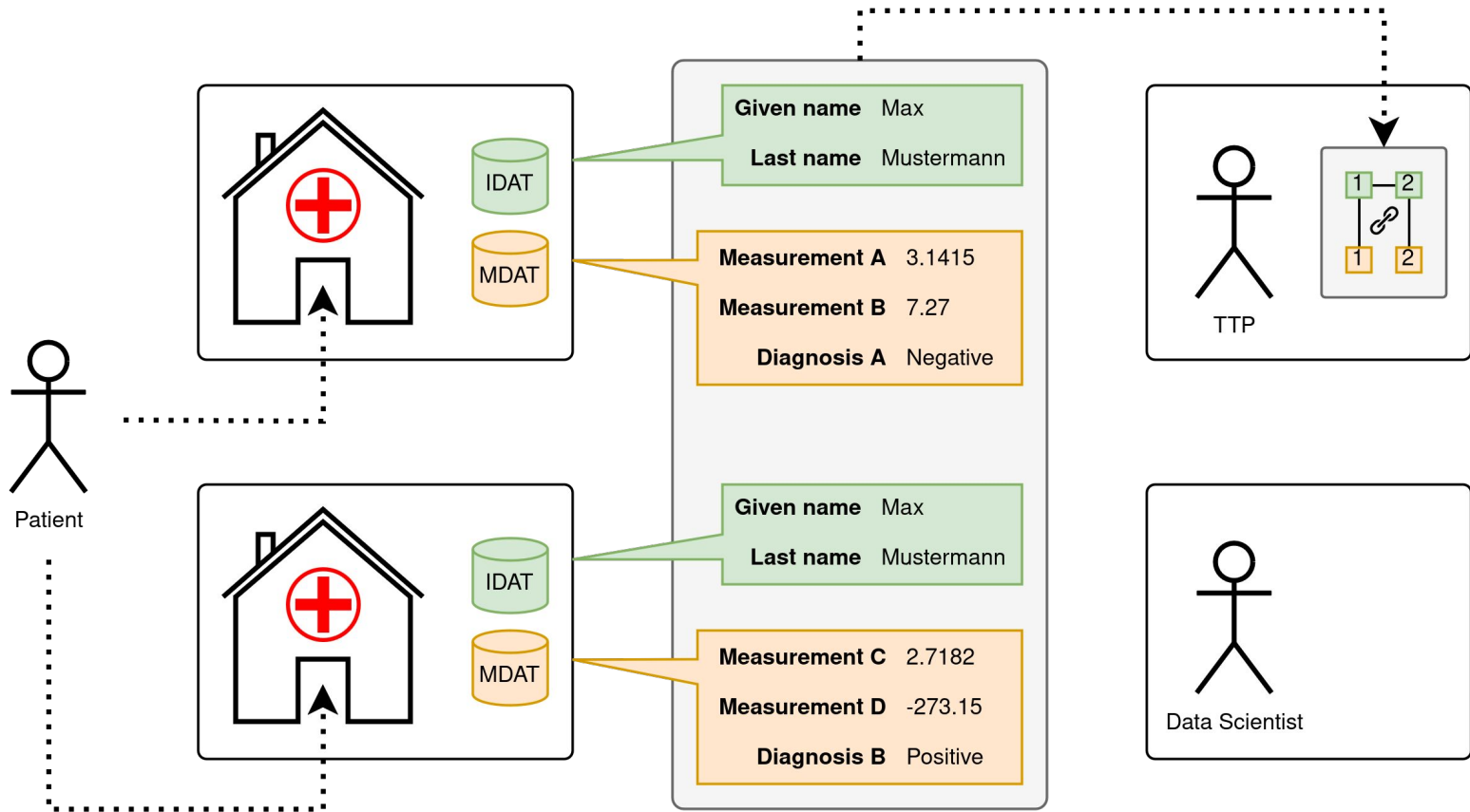
Freiberg, 07.06.2024

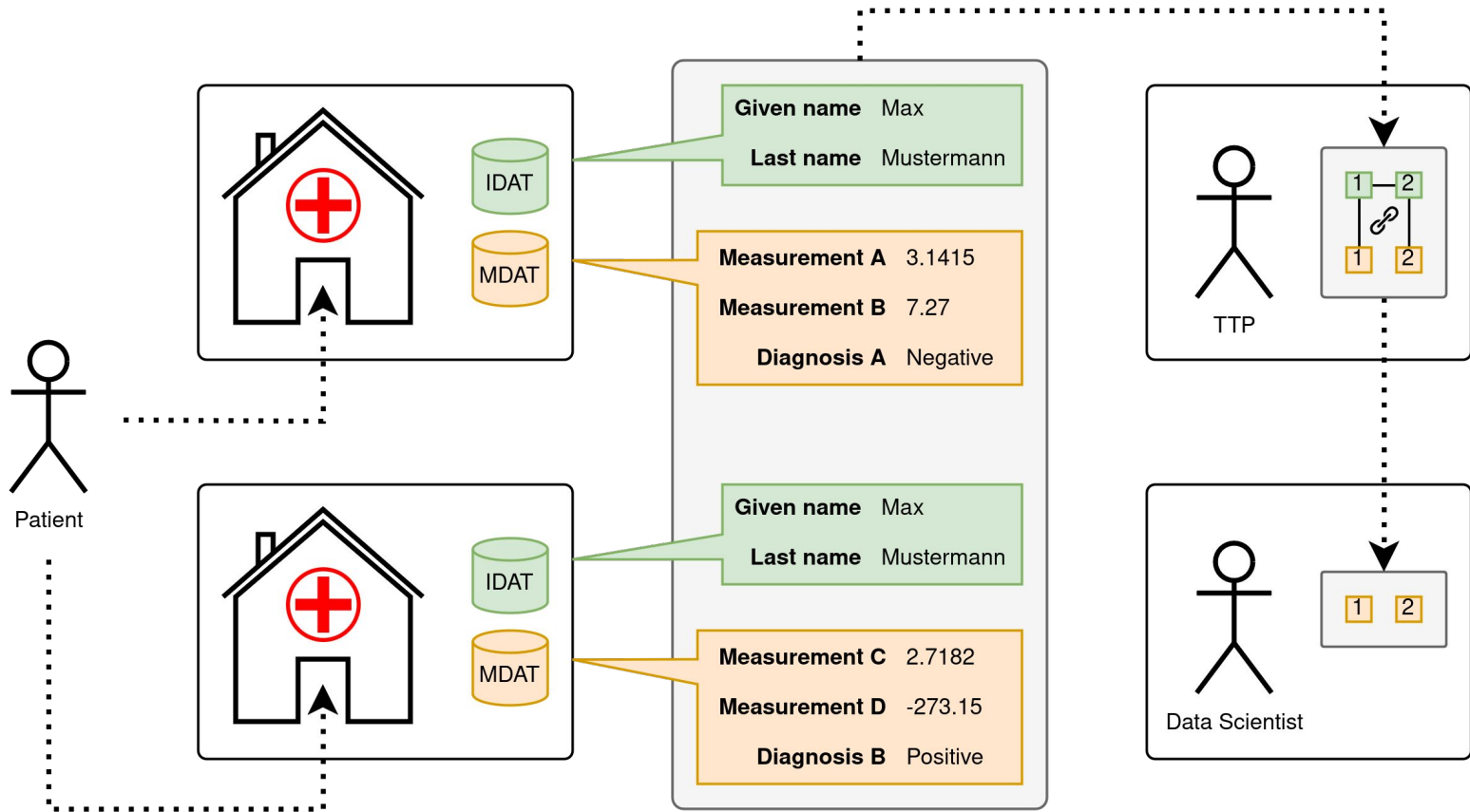
Maximilian Jugl



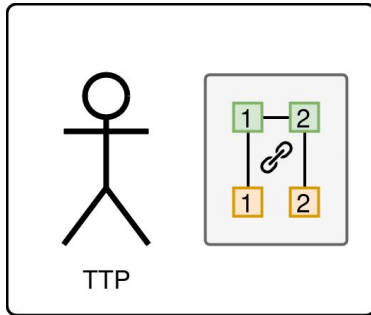








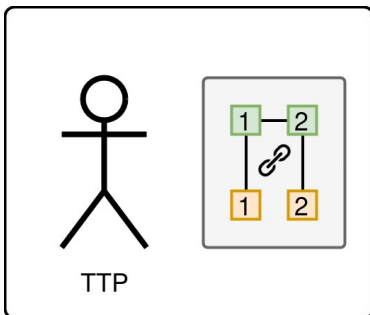
A PRIMER ON RECORD LINKAGE



- Testing of various record linkage strategies
- Testing against error sources and varying data schemas
- Limited access to real-world data

⇒ **Generation of realistic-looking test data**

A PRIMER ON RECORD LINKAGE

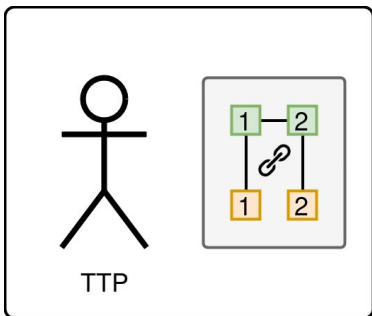


- Testing of various record linkage strategies
- Testing against error sources and varying data schemas
- Limited access to real-world data

⇒ **Generation of realistic-looking test data**

Given name	Last name	Gender	Date of birth
Axel	Schweiss	Male	1981-01-10
Grube	Claire		1970-10-01
Anna	Kond4	Female	1991-08-02

A PRIMER ON RECORD LINKAGE



- Testing of various record linkage strategies
- Testing against error sources and varying data schemas
- Limited access to real-world data

⇒ **Generation of realistic-looking test data**

	Given name	Last name	Gender	Date of birth
Typos?	Axel	Schweiss	Male	1981-01-10
Flipped values?	Grube	Claire		1970-10-01
OCR errors?	Anna	Kond4	Female	1991-08-02

Ambiguous format? (points to 1970-10-01)

Missing values? (points to empty gender cell)

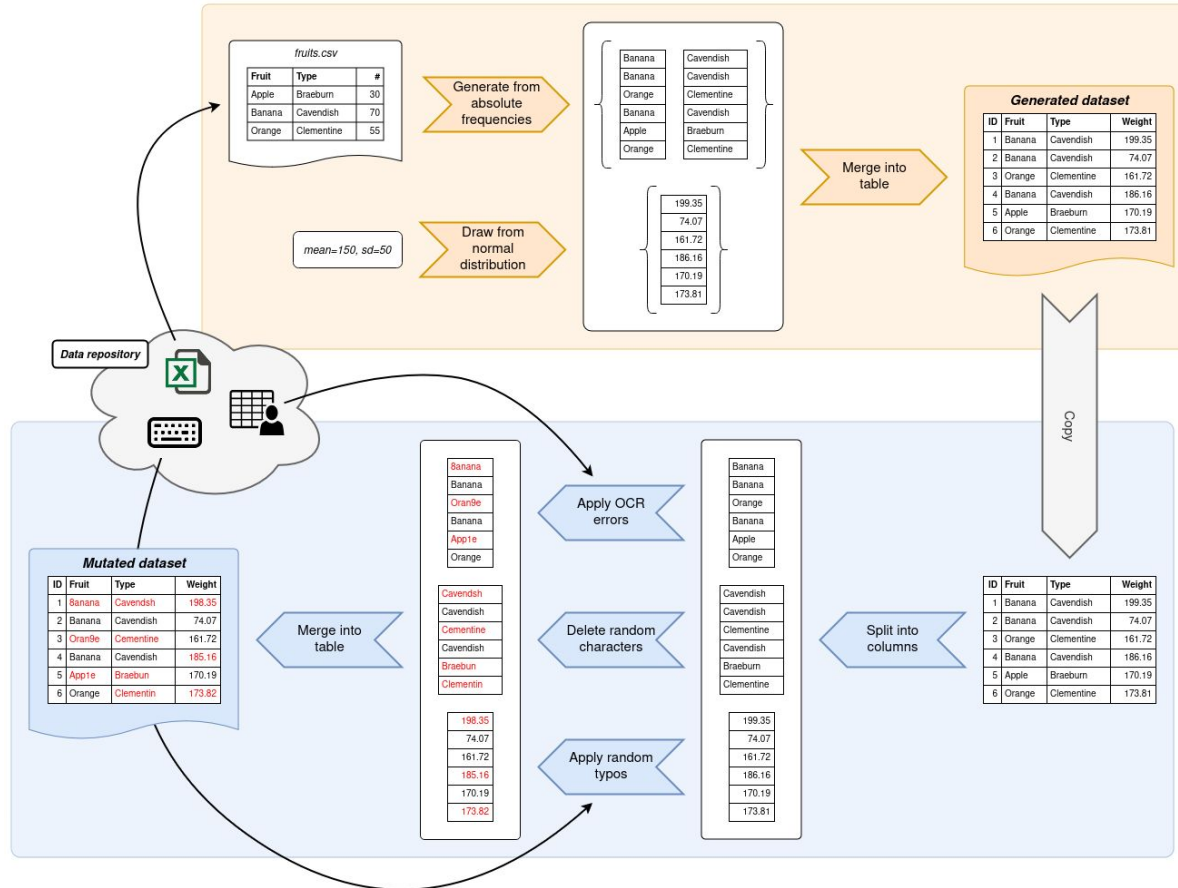
PRESENTING GECKO

- Modern approach to the ideas put forward by GeCo (Tran et al. 2013)
 - Completely reworked from the ground up for modern Python
 - Based on NumPy and Pandas to integrate into data science applications
 - Domain-independent, highly configurable and scalable
- Source code: <https://github.com/ul-mds/gecko>
- Documentation: <https://ul-mds.github.io/gecko/>
- Python Package Index: <https://pypi.org/project/gecko-syndata/>

The screenshot shows the GitHub repository page for 'gecko'. At the top, it indicates the repository is 'Public' and has 1 star. Below this, there are navigation options for 'main', '2 Branches', and '4 Tags'. A search bar and 'Add file' button are visible. The 'About' section describes it as a 'Python library for the generation and mutation of realistic personal identification data at scale' and provides the URL 'ul-mds.github.io/gecko/'. Below the description, there are tags for 'python', 'data-science', 'numpy', 'record-linkage', and 'pandas'. The commit history shows a recent merge pull request by 'mjugi' and three recent commits: adding cache to CI steps, removing a benchmark pipeline, and adding release notes to docs.



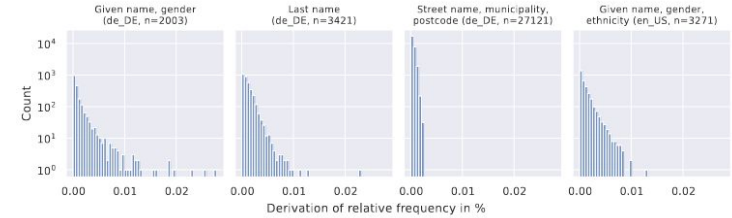
GECKO: REALISTIC DATA GENERATION FOR RECORD LINKAGE | Implementation



ID	Given name	Last name	Gender	Street name	Municipality	Postcode
254	Helmut	Jahn	m	Peenestraße	Stolpe	17391
M-254	Jahn	Helmut	m	Peenestraße	Stolpe	17391
1226	Rudolf	Franzen	m	Birkenweg	Suthfeld	31555
M-1226	Rudolf	Franzen	m	Birkenweg	Suthfeld	31565
2397	Erna	Eickhoff	f	Schulweg	Krautheim	74238
M-2397	Erna	Eickhoff	(empty)	Schulweg	Krautheim	74238
9960	Ingrid	Reinhold	f	Hochstraße	Mogendorf	56424
M-9960	Ingrid	Reinhold	m	Hochstraße	Mogendorf	56424

CORRECTNESS AND PERFORMANCE

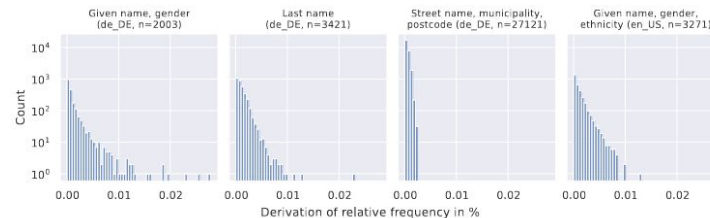
- Generate data from four distinct datasets with varying column and row counts



⇒ **Minor deviations from original distributions** (<0.005% in 99.5% of all cases)

CORRECTNESS AND PERFORMANCE

- Generate data from four distinct datasets with varying column and row counts



⇒ **Minor deviations from original distributions (<0.005% in 99.5% of all cases)**

- Benchmark with generation and mutation of 100k to 10m records
- Evaluation of single-core performance

Dataset	Records	CPU time in s				
		Min	Q5	Q50	Q95	Max
American	100 000	0.30	0.30	0.31	0.32	0.33
	1 000 000	2.82	2.83	2.87	3.07	3.11
	10 000 000	28.00	28.05	28.28	30.18	30.62
German	100 000	0.80	0.80	0.81	0.85	0.87
	1 000 000	6.63	6.63	6.74	6.84	6.86
	10 000 000	65.12	65.26	66.09	66.86	67.09

⇒ **Gecko is 15~100x faster than its modern alternatives**

CORRECTNESS AND PERFORMANCE

Dataset	Records	CPU time in s				
		Min	Q5	Q50	Q95	Max
American	100 000	0.30	0.30	0.31	0.32	0.33
	1 000 000	2.82	2.83	2.87	3.07	3.11
	10 000 000	28.00	28.05	28.28	30.18	30.62
German	100 000	0.80	0.80	0.81	0.85	0.87
	1 000 000	6.63	6.63	6.74	6.84	6.86
	10 000 000	65.12	65.26	66.09	66.86	67.09



Tool	year of release	interfaces	configurability	conf. effort (minimal)	schema/domain independence	degree of pollution	horizontal scalability	runtime (100k)	runtime (1mio)	runtime (10mio)
DBGen	1997	GUI/CLI	low	low	✗	✗	✗	5 s	7 s	25 s
FebrIDG	2008	CLI	low	low	✗	✗	✗	2 min	4.2 min	✗
GeCo	2012	Lib/Web-UI	high	high	✓	✗	✗	80 min	13.4 h	✗
TDGen	2012	GUI/(CLI)	high	medium	✓	✗	✗	4.5 min	✗	✗
ProbGee	2012	GUI	high	medium	✓	✓	✗	5.9 h	✗	✗
DaPo	2015	CLI	high	low	✓	✓	✓	31 s ¹ 17 s ²	2.1 min ¹ 56 s ²	26 min ¹ 14 min ²

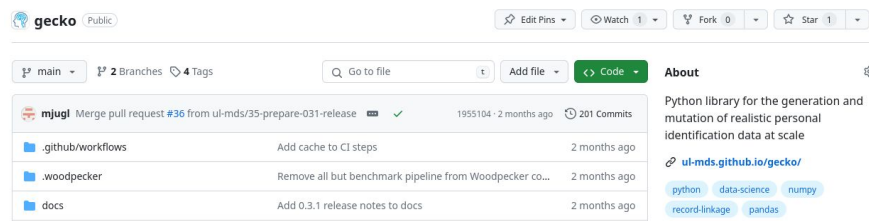
- Not bad, but...
 - Source code and test data for DaPo not available
 - DaPo performs much better when workloads are distributed
 - Other tools have more complex configuration options than Gecko

- Take these results with a large (but tasty) grain of salt :)

WHERE DO WE GO FROM HERE?

- Continuous testing of old and new record linkage algorithms
- Stress-testing of input forms that validate user-generated data
- PoC training data for machine learning models
- *This line could summarize your use case!*

⇒ Reach out! Maximilian.Jugl@medizin.uni-leipzig.de





UNIVERSITÄT
LEIPZIG

Medizinische Fakultät



Universitätsklinikum
Leipzig

Medizin ist unsere Berufung.

Thank you!

Freiberg, 07.06.2024

Maximilian Jugl



MDS
Medical Data Science