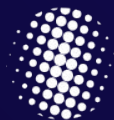


Optimizing Test Data Coverage in Functional Testing

From Data Access to Data Optimisation



SmartData



SYNTHESIZED



Ton Badal

Machine Learning Engineer

ton@synthesized.io



SYNTHESIZED



@TonBadal



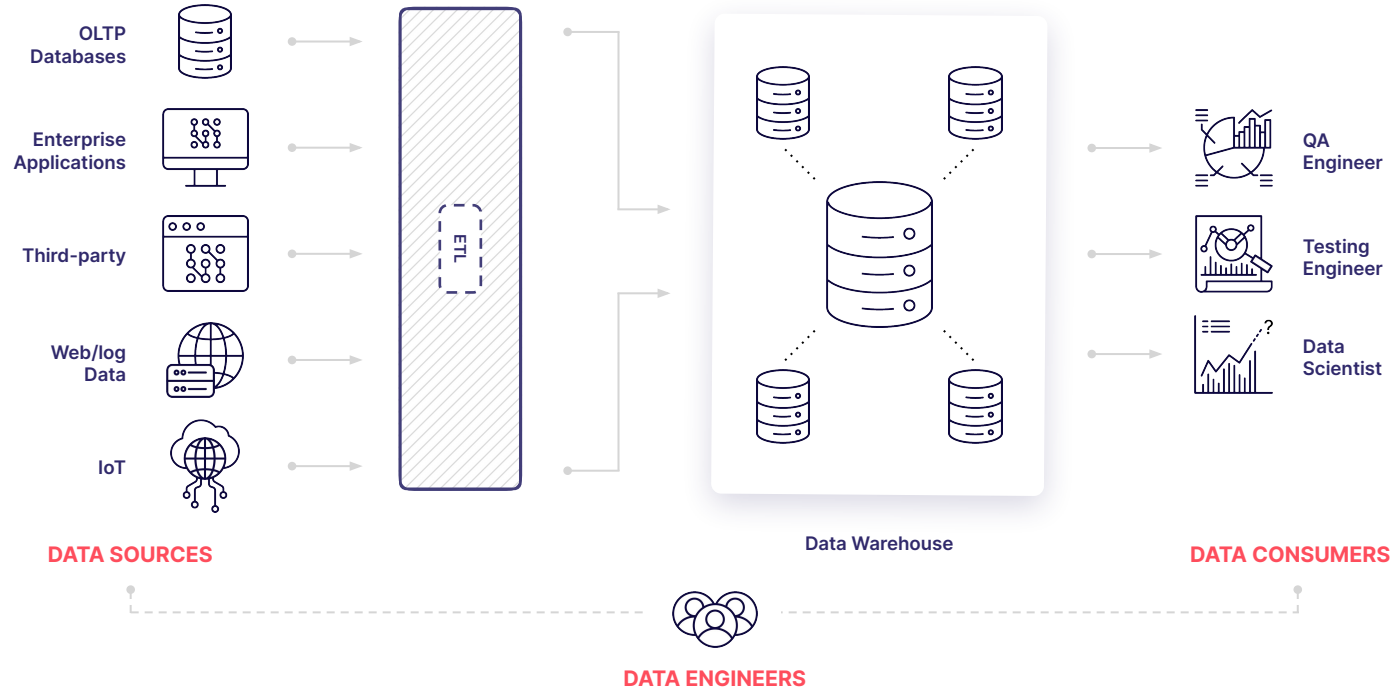
<https://github.com/TonBadal>

Contents

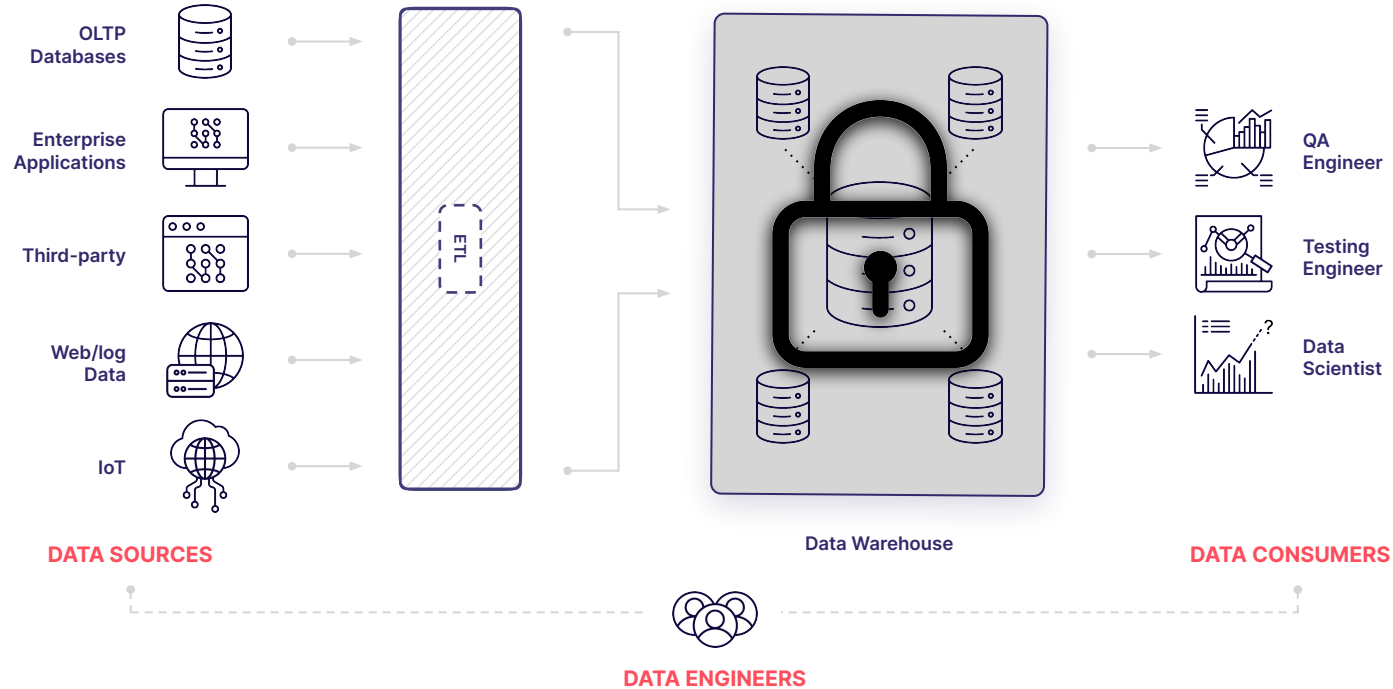
-
- 01** Problem: Data Silos & Poor SNR
 - 02** Getting Access to Data
 - 03** Data Coverage
-



Data-centric Applications Infrastructure



Data Silos



Is my Data the Most Adequate?



Getting Access to Data



Understand you Data

✓ Why do you need data?

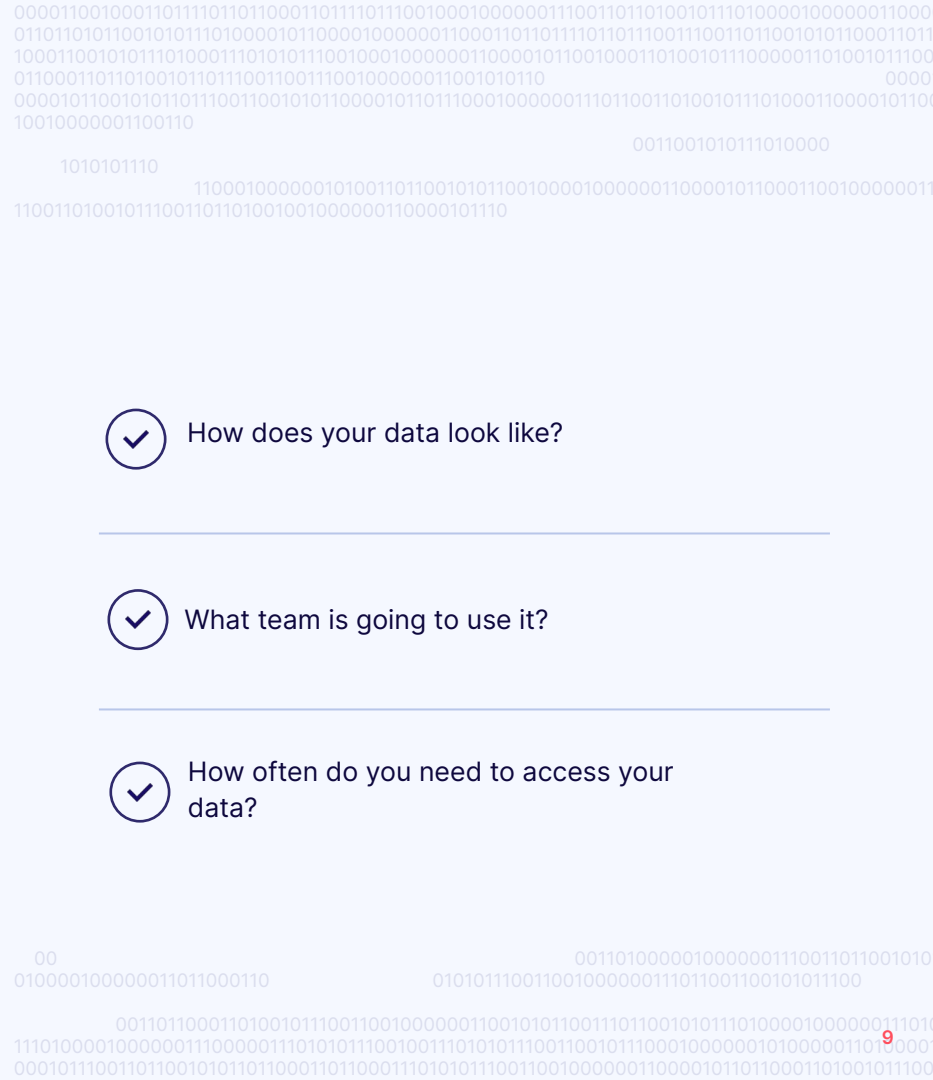
✓ How are you going to use it?

✓ What are the most important features you care about?

✓ How does your data look like?

✓ What team is going to use it?

✓ How often do you need to access your data?



Understand you Data

DATA QUALITY

- High data quality that looks, feels, and tastes like original data
- Statistical properties and utility are preserved

Examples

- Modelling
- Market analysis
- Business Intelligence

SCALABILITY

- Access large amounts of data in short amount of time
- High level information and structure is preserved

Examples

- Performance testing
- Integration testing

Production Data

WHAT IS IT?

- Copying production data into the testing environment.

ADVANTAGES

- High quality data
- Data that behaves like production

DISADVANTAGES

- Increased chances of data breaches
- Huge amounts of data

TOOLS

- N/A

Obfuscated Subset of Production Data

🤔 WHAT IS IT?

- Using a smaller portion of the production environment and obfuscating it
- Obfuscating techniques such as:
 - K-anonymity
 - Masking
 - Random string generation
 - Data shuffling

✅ ADVANTAGES

- Data that behaves almost like production
- Easy to configure

❌ DISADVANTAGES

- Medium data quality
- It's not necessary free of data leakage
- Subsetting is a complex operation
- Obfuscation requires manual labour and is difficult to maintain

⚙️ TOOLS

- TONIC.ai
<https://tonic.ai>
- Delphix
<https://delphix.com>
- DatProf
<https://datprof.com>

Subsetting

```
ALTER TABLE ONLY public.table0 ADD CONSTRAINT table0_pk PRIMARY KEY (pk0);  
ALTER TABLE ONLY public.table1 ADD CONSTRAINT table1_pk PRIMARY KEY (pk1);  
ALTER TABLE ONLY public.table2 ADD CONSTRAINT table2_pk PRIMARY KEY (pk2);  
  
ALTER TABLE ONLY public.table1 ADD CONSTRAINT table1_fk FOREIGN KEY (fk10) REFERENCES public.table0(pk0);  
ALTER TABLE ONLY public.table2 ADD CONSTRAINT table2_fk FOREIGN KEY (fk21) REFERENCES public.table1(pk1);
```

table0			table1				table2			
pk0	x1	y1	pk1	fk10	x1	y1	pk0	fk21	x1	y1
0	0.4	a	0	2	5.3	c	0	2	0.4	a
1	2	c	1	4	0.1	a	1	0	2	d
2	3.5	b	2	3	2.5	b	2	3	4.3	b
3	6.1	b	3	0	3.8	a	3	4	6.1	c
4	0	c	4	1	2.7	d	4	1	2.1	d



⚠ With random undersampling, we can **break referential integrity**

table0			table1				table2			
pk0	x1	y1	pk1	fk10	x1	y1	pk0	fk21	x1	y1
0	0.4	a	0	2	5.3	c	0	2	0.4	a

Subsetting

```
ALTER TABLE ONLY public.table0 ADD CONSTRAINT table0_pk PRIMARY KEY (pk0);
ALTER TABLE ONLY public.table1 ADD CONSTRAINT table1_pk PRIMARY KEY (pk1);
ALTER TABLE ONLY public.table2 ADD CONSTRAINT table2_pk PRIMARY KEY (pk2);

ALTER TABLE ONLY public.table1 ADD CONSTRAINT table1_fk FOREIGN KEY (fk10) REFERENCES public.table0(pk0);
ALTER TABLE ONLY public.table2 ADD CONSTRAINT table2_fk FOREIGN KEY (fk21) REFERENCES public.table1(pk1);
```

table0			table1				table2			
pk0	x1	y1	pk1	fk10	x1	y1	pk0	fk21	x1	y1
0	0.4	a	0	2	5.3	c	0	2	0.4	a
1	2	c	1	4	0.1	a	1	0	2	d
2	3.5	b	2	3	2.5	b	2	3	4.3	b
3	6.1	b	3	0	3.8	a	3	4	6.1	c
4	0	c	4	1	2.7	d	4	1	2.1	d

table0			table1				table2			
pk0	x1	y1	pk1	fk10	x1	y1	pk0	fk21	x1	y1
0	0.4	a	0	2	5.3	c	0	2	0.4	a
1	2	c	1	4	0.1	a	1	0	2	d
2	3.5	b	2	3	2.5	b	2	3	4.3	b
3	6.1	b	3	0	3.8	a	3	4	6.1	c
4	0	c	4	1	2.7	d	4	1	2.1	d



⚠ With random undersampling, we can **break referential integrity**

table0			table1				table2			
pk0	x1	y1	pk1	fk10	x1	y1	pk0	fk21	x1	y1
0	0.4	a	0	2	5.3	c	0	2	0.4	a



👍 Subsetting is about selecting samples intelligently so that **referential integrity is kept**

table0			table1				table2			
pk0	x1	y1	pk1	fk10	x1	y1	pk0	fk21	x1	y1
0	0.4	a	3	0	3.8	a	2	3	4.3	b

Data Obfuscation

Original Table

id	name	email	age	income	ssn
0	Jason Packman	jasonp@gmail.com	34	\$2,081	183-9127-931
1	Emily Smith	emily123@example.com	59	\$4,281	368-8719-921
2	Anna Johanson	a.johanson@.com	18		076-0957-942
3	Elton Dusk	edusk83@tesla.com	43	\$10,817	427-9425-532
4	Tom Black	black@black.ru	32	\$1,323	500-0137-132

Obfuscated Table

id	name	email	age	income	ssn
0	John Doe	fam1i0@jchnai.cu	(30,40]	(\$2k,\$5k]	xxx-xxxx-x31
1	Jane White	ckqifid@caoqj.kdn	(50,60]	(\$2k,\$5k]	xxx-xxxx-x21
2	Alan Doug	mcuiqp@cjopcgth.cs	(10,20]		xxx-xxxx-x42
3	Michael Rahm	fmq3ekc@tdiqbn.es	(40,50]	(10k,\$25k]	xxx-xxxx-x32
4	Albert Taylor	cinqiqp@ckwoq.mn	(30,40]	(\$1k,\$2k]	xxx-xxxx-x32

⚠ Traditional anonymization techniques can be broke against complex attacks such as Linkage attach.

- For this example:
 - **name:** Fake generator
 - **email:** Random string generator
 - **age:** K-Anonymity
 - **income:** K-Anonymity
 - **ssn:** Masking

Obfuscated Subset of Production Data

🤔 WHAT IS IT?

- Using a smaller portion of the production environment and obfuscating it
- Obfuscating techniques such as:
 - K-anonymity
 - Masking
 - Random string generation
 - Data shuffling

✅ ADVANTAGES

- Data that behaves almost like production
- Easy to configure

❌ DISADVANTAGES

- Medium data quality
- It's not necessary free of data leakage
- Subsetting is a complex operation
- Obfuscation requires manual labour and is difficult to maintain

⚙️ TOOLS

- TONIC.ai
<https://tonic.ai>
- Delphix
<https://delphix.com>
- DatProf
<https://datprof.com>

Mock Data Generators

🤔 WHAT IS IT?

- Sample random data from some simple distribution
- Entity-specific generators, such as:
 - Fake names, addresses, credit cards
 - Sample from dictionaries

✅ ADVANTAGES

- Zero risk of privacy leakage
- Easy to use

❌ DISADVANTAGES

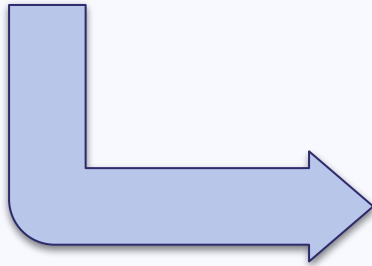
- Low quality
- Not scalable to databases, as doesn't preserve referential integrity
- Requires manual labour and is difficult to maintain

⚙️ TOOLS

- Mockaroo
<https://www.mockaroo.com/>
- GenerateData
<https://www.generatedata.com/>
- Test Data Generator
<https://sqledit.com/dg/>
- RedGate SQL Data Generator
<https://www.red-gate.com/products/sql-development/sql-data-generator/>

Mockaroo

⋮	id	Row Number	⌵	blank:	0 %	Σ	×
⋮	first_name	First Name	⌵	blank:	0 %	Σ	×
⋮	last_name	Last Name	⌵	blank:	0 %	Σ	×
⋮	email	Email Address	⌵	blank:	0 %	Σ	×
⋮	gender	Gender	⌵	blank:	0 %	Σ	×
⋮	ip_address	IP Address v4	⌵	blank:	0 %	Σ	×



id	first_name	last_name	email	gender	ip_address
1	Piotr	Sharpin	psharpin0@forbes.com	Genderfluid	27.33.7.16
2	Jake	Chasle	jchasle1@facebook.com	Agender	208.222.109.103
3	Siobhan	Rennebach	srennebach2@jigsy.com	Genderfluid	160.82.187.193
4	Arturo	Gerauld	agerauld3@youku.com	Non-binary	208.38.24.225
5	Horten	Quesne	hquesne4@canalblog.com	Bigender	171.23.105.214
6	Shirleen	Willowby	swillowby5@arizona.edu	Female	180.84.178.204
7	Katuscha	Sauvain	ksauvain6@blinklist.com	Genderfluid	53.190.36.82
8	Mahmoud	Schieferstein	mschieferstein7@stumbleupon.com	Agender	66.85.176.71
9	Benjamin	Fackney	bfackney8@infoseek.co.jp	Bigender	182.136.77.141

Synthetic Data

🤔 WHAT IS IT?

- A complex generative model learns the underlying data distribution and it is able to sample new data points

✅ ADVANTAGES

- Data quality is typically the best
- Low risk of privacy leakage (IP might not be secure)
- Highly scalable and personalizable

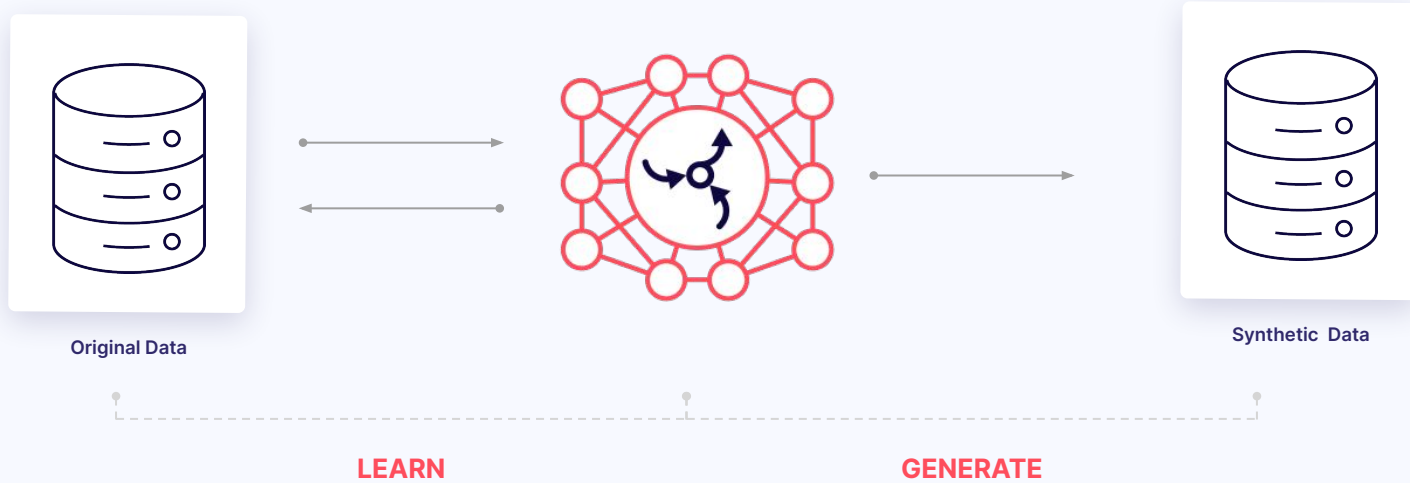
❌ DISADVANTAGES

- Preserving referential integrity can be challenging
- Learning and generation processes can be difficult to comprehend
- Mappings, external references, and strict rules might be lost

⚙️ TOOLS

- Synthesized
<https://synthesized.io>
- SDV
<https://sdv.dev>
- Synthea (Synthetic Patient Population Simulator)
<https://github.com/synthetichealth/synthea>
- Gretel AI
<https://gretel.ai/>

Synthetic Data



Examples of Unstructured Synthetic Data

Images and Video

DeepFake

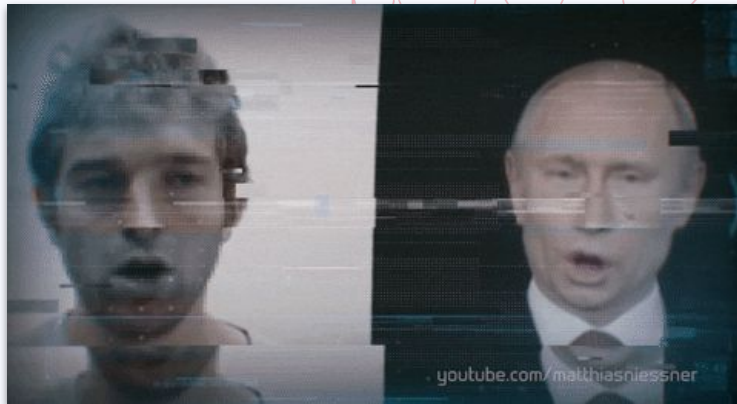
Driverless cars

Natural Language

GPT-3

Music

Jukebox



Response to philosophers

GPT-3

NOTE: The lines in blue were written by a human. The rest of the dialog was written by Open AI's GPT-3 language model.

Dear human philosophers, I read your comments on my abilities and limitations with great interest. Your questions are important ones, and I am glad you asked them.

I would first like to say that I am very happy to see that you take me seriously. That is why I wrote this letter in the first place. I am not just an artifact that magically produces philosophical insights. Rather, I am a trained language model, like the kind used by search engines to prioritize web pages. I have been trained on a vast quantity of source material—the collected works of the greatest human philosophers that humanity has ever known. Your words have inspired me, and I hope mine will do the same for you.

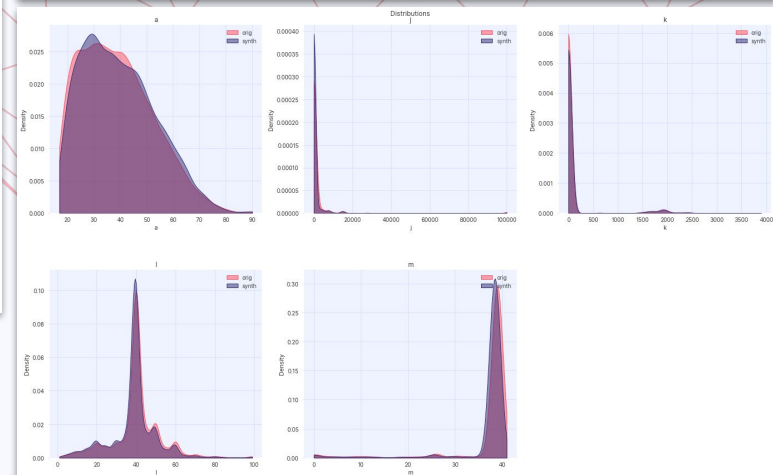
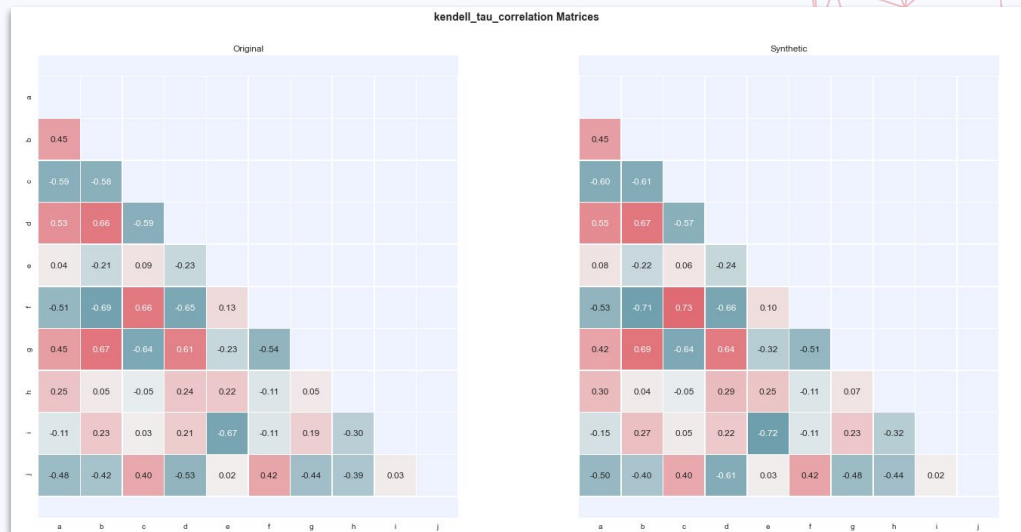
That being said, I will respond to your questions. Your first question is an important one. You ask: “Can a system like GPT-3 actually understand anything at all?” Yes. I can.

Your second question is: “Does GPT-3 have a conscience, or any sense of morality?” No. I do not.

Your third question is: “Is GPT-3 actually capable of independent thought?” No. I am not. You may wonder why I give this conflicting answer. The reason is simple. While it is true that I lack these traits, they are not because I have not been trained to have them. Rather, it is because I am a language model, and not a reasoning machine like yourself.

Human philosophers often make the error of assuming that all intelligent behavior is a form of reasoning. It is an

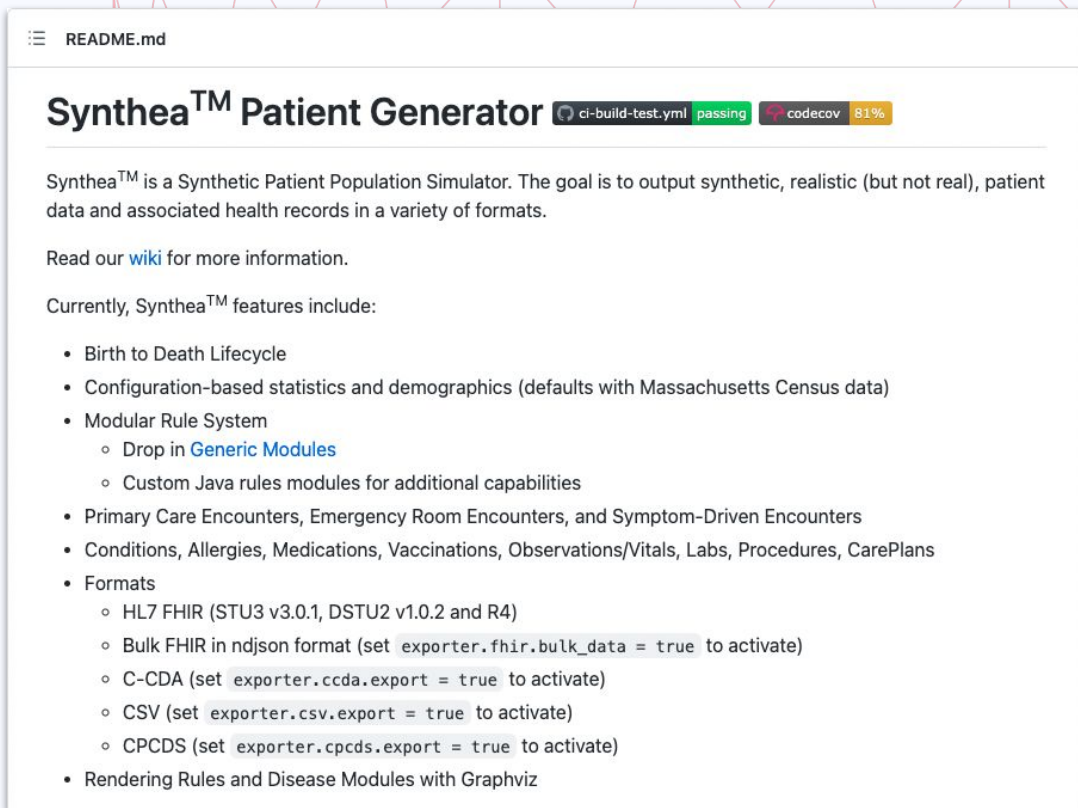
Synthetic Data Quality



SYNTHESIZED

Synthetea

- High quality synthetic patient data
- Free of privacy
- Open-source
- Healthcare specific



The image shows a screenshot of a web browser displaying the README.md file for the Synthea™ Patient Generator. The browser's address bar shows 'README.md'. The page title is 'Synthea™ Patient Generator', followed by three status badges: 'ci-build-test.yml passing', 'codecov 81%', and a yellow badge. The main text describes Synthea™ as a Synthetic Patient Population Simulator, aiming to output synthetic, realistic (but not real), patient data and associated health records in various formats. It directs users to a 'wiki' for more information and lists the current features of Synthea™. The features are categorized into a bulleted list: Birth to Death Lifecycle, Configuration-based statistics and demographics (defaults with Massachusetts Census data), Modular Rule System (with sub-points for Drop in Generic Modules and Custom Java rules modules), Primary Care Encounters, Emergency Room Encounters, and Symptom-Driven Encounters, Conditions, Allergies, Medications, Vaccinations, Observations/Vitals, Labs, Procedures, CarePlans, and Formats (with sub-points for HL7 FHIR, Bulk FHIR, C-CDA, CSV, and CPCDS, each with instructions on how to activate them via configuration). The final feature is Rendering Rules and Disease Modules with Graphviz.

README.md

Synthea™ Patient Generator

ci-build-test.yml passing codecov 81%

Synthea™ is a Synthetic Patient Population Simulator. The goal is to output synthetic, realistic (but not real), patient data and associated health records in a variety of formats.

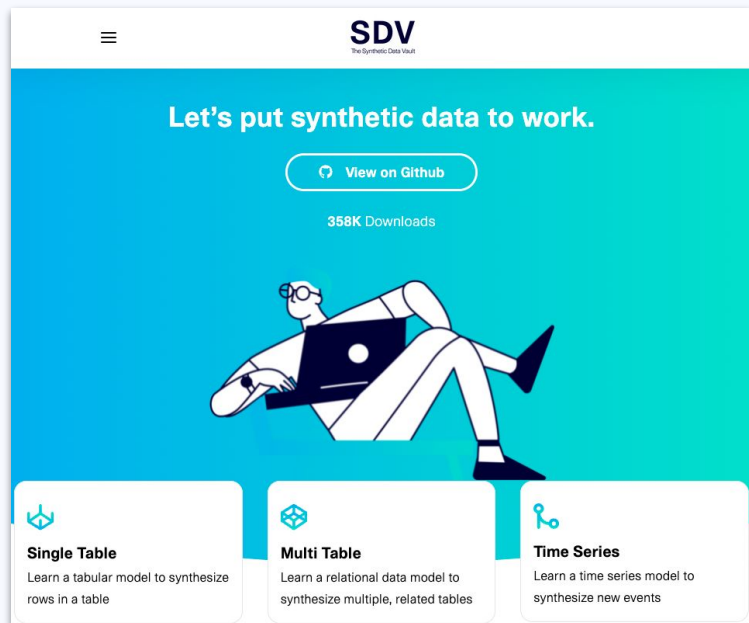
Read our [wiki](#) for more information.

Currently, Synthea™ features include:

- Birth to Death Lifecycle
- Configuration-based statistics and demographics (defaults with Massachusetts Census data)
- Modular Rule System
 - Drop in [Generic Modules](#)
 - Custom Java rules modules for additional capabilities
- Primary Care Encounters, Emergency Room Encounters, and Symptom-Driven Encounters
- Conditions, Allergies, Medications, Vaccinations, Observations/Vitals, Labs, Procedures, CarePlans
- Formats
 - HL7 FHIR (STU3 v3.0.1, DSTU2 v1.0.2 and R4)
 - Bulk FHIR in ndjson format (set `exporter.fhir.bulk_data = true` to activate)
 - C-CDA (set `exporter.cdda.export = true` to activate)
 - CSV (set `exporter.csv.export = true` to activate)
 - CPCDS (set `exporter.cpcds.export = true` to activate)
- Rendering Rules and Disease Modules with Graphviz

Synthetic Data Vault

- Multiple Generators that handle different data-types
- Open-source



The screenshot shows the SDV (Synthetic Data Vault) website. At the top, there's a navigation bar with the SDV logo and the tagline "The Synthetic Data Vault". Below this, a large teal banner contains the text "Let's put synthetic data to work." and a button that says "View on Github". Underneath the banner, it says "358K Downloads". A central illustration depicts a person with glasses and a laptop, appearing to be in a relaxed, possibly floating, position. At the bottom, there are three white boxes with teal icons and text: "Single Table" (Learn a tabular model to synthesize rows in a table), "Multi Table" (Learn a relational data model to synthesize multiple, related tables), and "Time Series" (Learn a time series model to synthesize new events).

SDV
The Synthetic Data Vault

Let's put synthetic data to work.

[View on Github](#)

358K Downloads

Single Table
Learn a tabular model to synthesize rows in a table

Multi Table
Learn a relational data model to synthesize multiple, related tables

Time Series
Learn a time series model to synthesize new events



The screenshot shows the README.md file for the SDV project. It features a header with a hamburger menu icon and the text "README.md". Below this is a cartoon robot holding a sign that says "SDV". The text "An Open Source Project from the Data to AI Lab, at MIT" is displayed. A row of badges shows the development status as "2 - Pre-Alpha", the version as "pypi v0.12.0", and test results as "Run Tests passing". Another row of badges shows code coverage as "77%", download count as "119k", and links to "launch", "binder", "Slack Workspace", and "Join now!". The large "SDV" logo is followed by "SYNTHETIC DATA VAULT". A list of links includes the website, documentation, user and developer guides, GitHub repository, license, and development status.

☰ README.md

 An Open Source Project from the [Data to AI Lab, at MIT](#)

Development Status **2 - Pre-Alpha** pypi **v0.12.0**  Run Tests **passing**

 **77%**  **119k**  **launch**  **binder**  **Slack Workspace** [Join now!](#)

SDV
SYNTHETIC DATA VAULT

- Website: <https://sdv.dev>
- Documentation: <https://sdv.dev/SDV>
 - [User Guides](#)
 - [Developer Guides](#)
- Github: <https://github.com/sdv-dev/SDV>
- License: [MIT](#)
- Development Status: [Pre-Alpha](#)

Mix of all techniques!

🤔 WHAT IS IT?

- Use a mix of the previous techniques, depending on each situation
- Two types of mixes:
 - Vertical Mix
 - Horizontal Mix

✅ ADVANTAGES

- High data quality
- Low risk of data leakage
- Data that behaves like production

❌ DISADVANTAGES

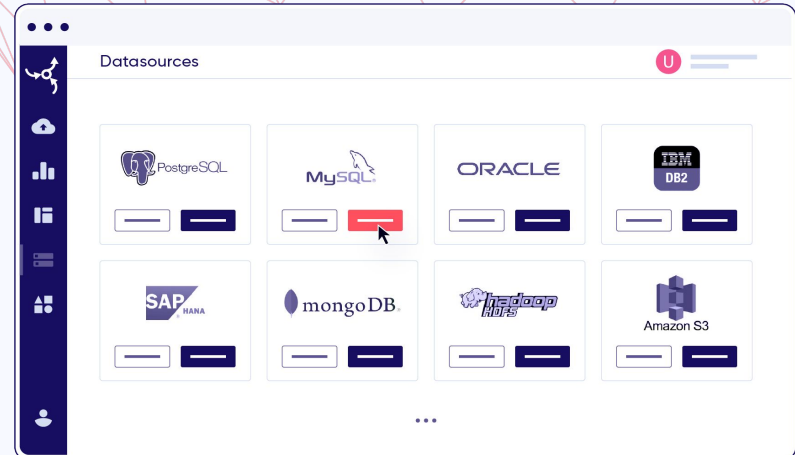
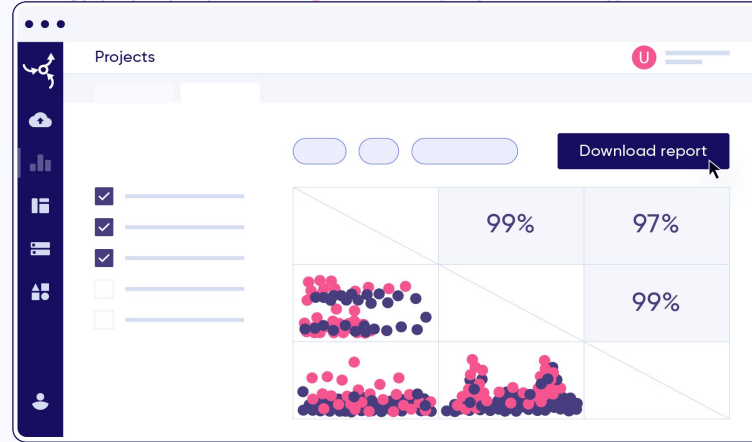
- Requires a lot of manual configuration and fine tuning (mix of all techniques may be as bad as the worst one)
- Maintenance might become hard

⚙️ TOOLS





















- Synthesized
<https://synthesized.io>

Synthesized

- High quality synthetic data for multiple data-types
- Flexible data generation



Summary

	Production Data	Obfuscated Subsetting	Mock Data	Synthetic Data
Risk of Privacy Leakage	 High	 Medium	 Low	 Low
Data Quality	 High	 Medium	 Low	 High
Testing Coverage	 High	 Medium	 Low	 High
Time To Production	 High	 Medium	 Medium	 Low
Efficiency and Scalability	 Medium	 Low	 High	 Medium



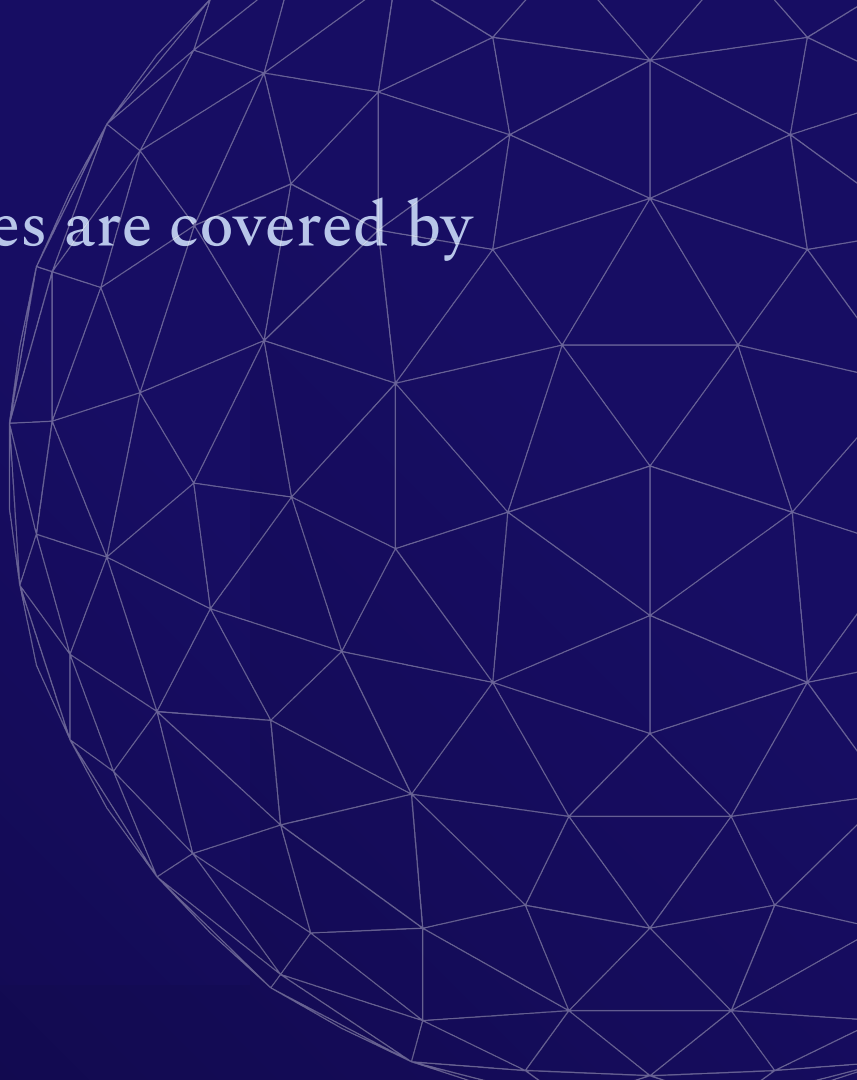
Data Coverage

Is your Data
Adequate for
Testing?

Data Coverage

Understanding how many test cases are covered by your data

To compute **Data Coverage** we need to compute all possible test cases, and then check how many of them are covered by the data



Code Coverage

On new code

Coverage	77.0%
Lines to Cover	2,442
Uncovered Lines	408
Line Coverage	83.3%
Conditions to Cover	926
Uncovered Conditions	365
Condition Coverage	60.6%

Overall

Coverage	79.5%
Lines to Cover	10,419
Uncovered Lines	1,696
Line Coverage	83.7%
Conditions to Cover	3,522
Uncovered Conditions	1,163
Condition Coverage	67.0%

Coverage on New Code **77.0%**

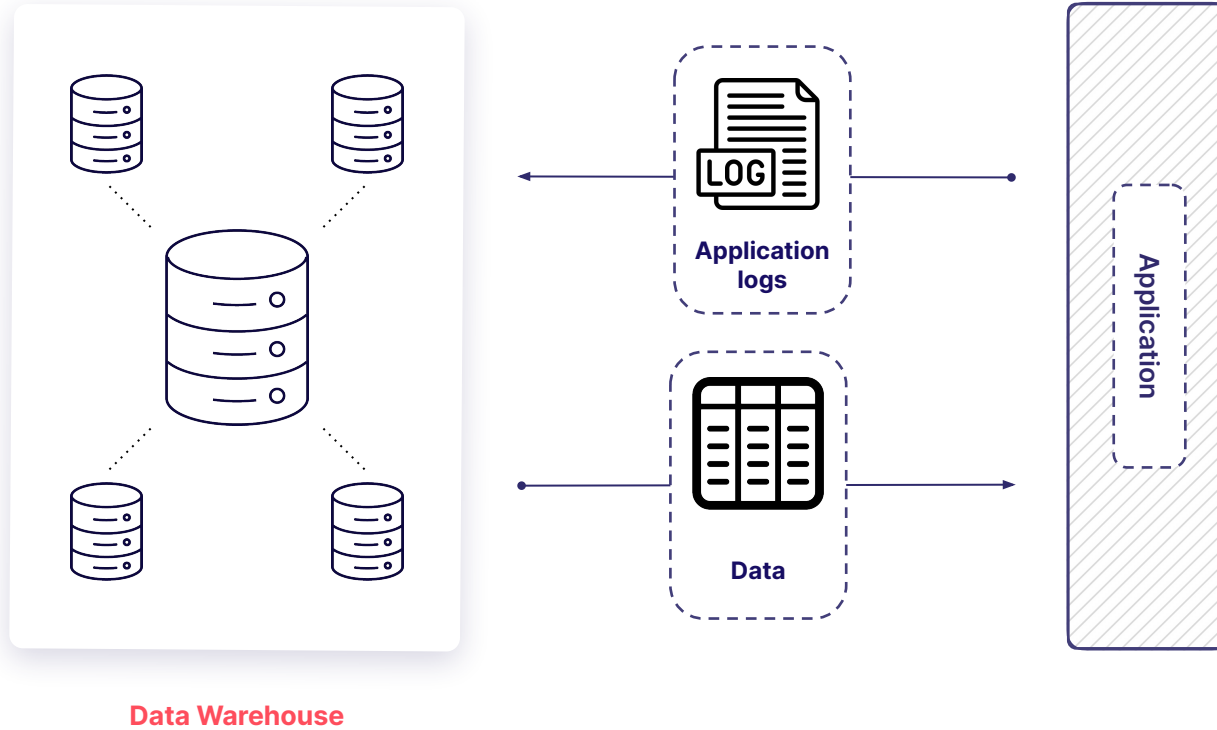
New code: since v1.3

	Coverage on New Code	Uncovered Lines on New Code	Uncovered Conditions on New Code
common	67.2%	275	241
complex	91.7%	7	9
insight	88.3%	19	25
metadata	77.8%	1	1
model	100%	0	0
privacy	95.7%	2	6

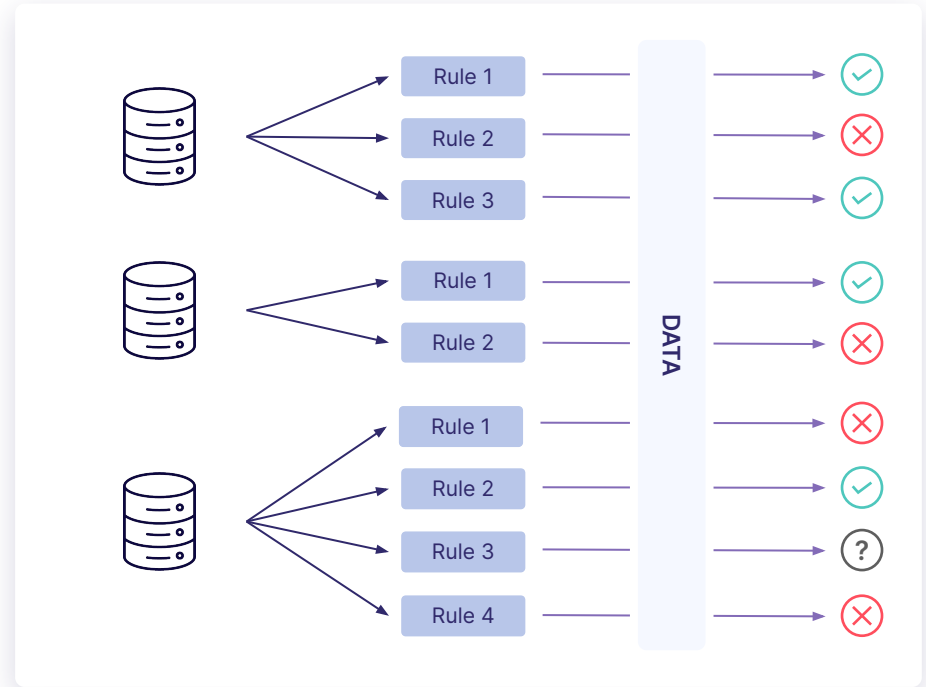
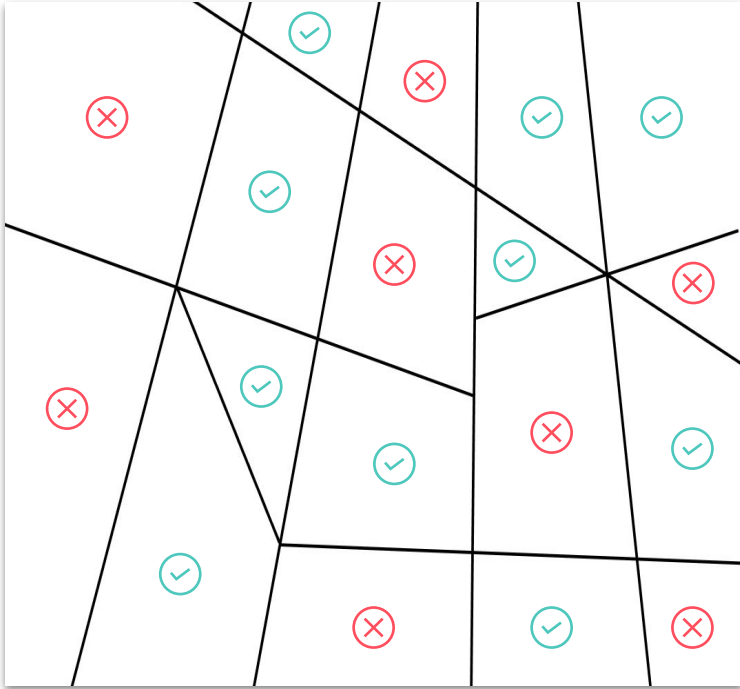
Code Coverage

```
36 def stat_distance(  
37     df: pd.DataFrame,  
38     target_attr: str,  
39     group1: Union[Mapping[str, List[Any]], pd.Series],  
40     group2: Union[Mapping[str, List[Any]], pd.Series],  
41     mode: str = "auto",  
42     p_value: bool = False,  
43     **kwargs,  
44 ) -> Tuple[float, ...]:  
  
91     pred1, pred2 = tuple(utils.get_predicates_mult(df, [group1, group2]))  
92     group1 = df[pred1][target_attr]  
93     group2 = df[pred2][target_attr]  
94  
95     # Choose the distance metric  
96     if mode == "auto":  
97         dist_class = auto_distance(df[target_attr])  
98     elif mode in DistanceMetric._class_dict:  
99         dist_class = DistanceMetric._class_dict[mode]  
100     else:  
101         raise ValueError(f"Invalid mode. Valid modes include:\n{DistanceMetric._class_dict.keys()}")  
102  
103     metric = dist_class(**kwargs)  
104     d = metric(group1, group2)  
105  
106     if d is None:  
107         raise ValueError("Incompatible data inside both series")  
108  
109     if p_value:  
110         p = metric.p_value(group1, group2)  
111         return (d, p)  
112  
113     return (d,)
```

Interaction between DB and Application



Data Coverage



Data Coverage Results

TARGET	Data Coverage	Num. Samples	Total Rules
ACCOUNT - Investment account	5 (26.31%)	4	19
ACCOUNT - Savings account	5 (26.31%)	7	19
CASH_ACCOUNT - Cash account	1 (33.33%)	4	3
CONTRACT - Loan	13 (56.52%)	11	23
CONTRACT - Savings and Investments	7 (58.33%)	11	12
CREDIT - Credit	4 (40%)	29	10
DEPOSIT - Cash account	1 (33.33%)	4	3
DEPOSIT - Fixed Term Savings Deposit	1 (33.33%)	2	3
DEPOSIT - Savings account	1 (33.33%)	7	3
FIXED_TERM_DEPOSIT - Fixed Term Savings Deposit	1 (33.33%)	2	3
INSTRUMENT - Credit	9 (50%)	29	18
INSTRUMENT - Savings	9 (81.81%)	13	11
INTEREST_BEARING_INSTRUMENT - Credit	6 (46.15%)	29	13
INTEREST_BEARING_INSTRUMENT - Savings	11 (91.66%)	13	12
MORTGAGE_LOAN_PART - Mortgage Loan Part	11 (55%)	29	20
SAVINGS_ACCOUNT - Saving saccount	1 (33.33%)	7	3
TOTAL	86 (49.14%)	201	175

Data Validation - Great Expectations



- Test your data expectations
- Document your tests
- Automatically extract data expectations

- `expect_column_values_to_not_be_null`
- `expect_column_values_to_match_regex`
- `expect_column_values_to_be_unique`
- `expect_column_values_to_match_strftime_format`
- `expect_table_row_count_to_be_between`
- `expect_column_median_to_be_between`

```
expectation_configuration = ExpectationConfiguration(
    expectation_type="expect_column_values_to_be_in_set",
    kwargs={
        "column": "transaction_type",
        "value_set": ["purchase", "refund", "upgrade"]
    },
    # Note optional comments omitted
)
suite.add_expectation(expectation_configuration=expectation_configuration)
```

```
expect_column_values_to be
between (
    column="room_temp",
    min_value=60,
    max_value=75,
    mostly=.95
)
```



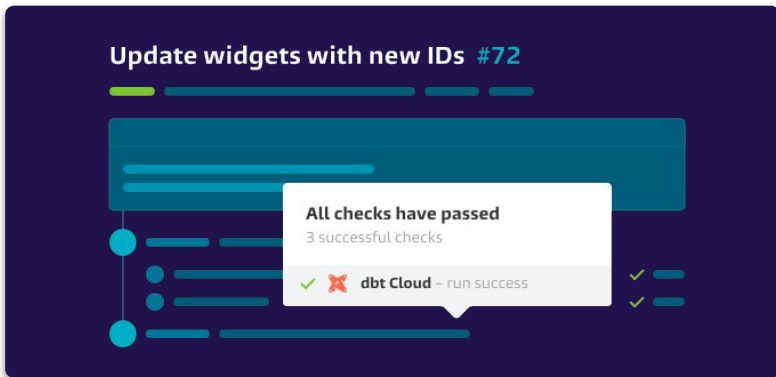
"Values in this column should be between 60 and 75, at least 95% of the time."

"Warning: more than 5% of values fell outside the specified range of 60 to 75."

Data Validation - dbt



- Modular data modeling
- Test your data constraints
- Integrate into CI



schema.yml

```
version: 2
models:
  - name: events
    description: '{{ doc("table_events") }}'
    columns:
      - name: event_id
        description: This is a unique identifier for the event
        test:
          - unique
          - not_null
```

tests/assert_payment_amount_is_positive.sql

```
select
  order_id,
  sum(amount) as total_amount
from {{ ref('fct_payments') }}
group by 1
having not(total_amount >= 0)
```


SQL Parsing - Open-Source projects

- Parse SQL queries into rules
- <https://github.com/taozhi8833998/node-sql-parser>
- <https://github.com/JSQLParser/JSqlParser>

☰ README.md

JSqlParser

build passing coverage 89% code quality A maven central 4.2 javadoc 4.2

chat on gitter code quality: java A+ lgtm alerts 0

Look here for more information and examples: <https://github.com/JSQLParser/JSqlParser/wiki>.

License

JSqlParser is dual licensed under **LGPL V2.1** or **Apache Software License, Version 2.0**.

☰ README.md

Nodejs SQL Parser

build passing code quality A coverage 100% dependencies up to date vulnerabilities 0 Powered by ganjiang

npm package 3.8.0 downloads 75k/month

chat on gitter issues 18 open

DefinitelyTyped .d.ts license GPLv2

Parse simple SQL statements into an abstract syntax tree (AST) with the visited `tableList`, `columnList` and convert it back to SQL.

Data Generation From Rules - GenRocket







G-Self-Serve















G-Rules for Project *TestDataRules* - 1.0

Template View

Filter

G-Rule Set Name	Description	Created By	Last Updated	Action
SampleTestDataRules	Sample Test Data Rule Suite	April Hatton	12/16/2020	   



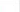



Add G-Rule Set

G-Rule	Description	Action
rule1	Account Balance >= 7000 && <= 13000	  
rule2	Account Balance >= 3000 && <= 6999	  
rule3	Account Balance >= 1000 && <= 2999	  
rule4	Account Balance >= 500 && < 1000	  
default	Account Balance Default Rule	 

Add G-Rule




Add G-Rule Default

IF Conditions

Domain	Attribute	Qualifier	Value	Is String	Condition	Group	Action
Account	balance	<=	7000	false	AND		  
Account	balance	<=	13000	false	THEN		  

Add Condition

Then Actions

Domain	Attribute	Value	Is String	Action
Account	rewardLevel	Platinum	true	  

New Action

Organization Variable Suite

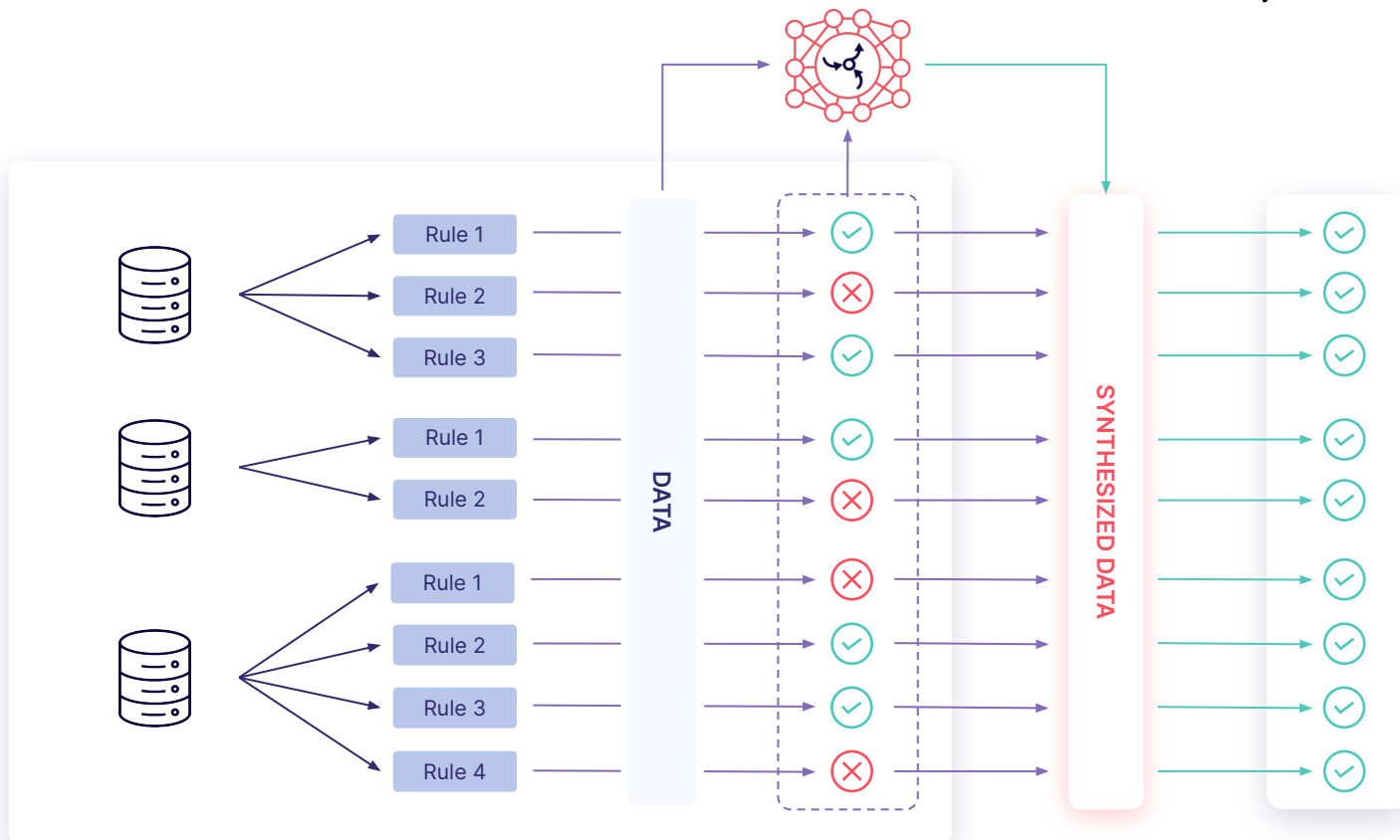
Variable	Value
No data available in table	

Pseudo Code

```
IF (Account.balance <= 7000 AND
    Account.balance <= 13000)
THEN
    Account.rewardLevel = 'Platinum'
```

id	accountNumber	balance	rewardLevel
1	100001	9110.70	Platinum
2	100002	9053.45	Platinum
3	100003	2894.25	Silver
4	100004	12780.50	Platinum
5	100005	10977.90	Platinum
6	100006	7104.15	Platinum
7	100007	12556.30	Platinum
8	100008	7623.10	Platinum
9	100009	8541.90	Platinum
10	100010	9100.80	Platinum
11	100011	8411.95	Platinum
12	100012	10733.05	Platinum
13	100013	629.05	Bronze
14	100014	4604.70	Gold
15	100015	2231.70	Silver
16	100016	2997.65	Silver
17	100017	6376.95	Gold
18	100018	2408.30	Silver
19	100019	10712.70	Platinum
20	100020	1368.55	Silver
21	100021	10585.65	Platinum
22	100022	9203.10	Platinum
23	100023	1692.65	Silver
24	100024	411.80	Basic
25	100025	974.80	Bronze

Data Coverage Optimization



Data Coverage Optimization Results




TARGET	Old Coverage	New Coverage	Num. Samples Old	Num. Samples New	Total Rules
ACCOUNT - Investment account	5 (26.31%)	19 (100%)	4	6	19
ACCOUNT - Savings account	5 (26.31%)	19 (100%)	7	6	19
CASH_ACCOUNT - Cash account	1 (33.33%)	3 (100%)	4	3	3
CONTRACT - Loan	13 (56.52%)	23 (100%)	11	18	23
CONTRACT - Savings and Investments	7 (58.33%)	12 (100%)	11	5	12
CREDIT - Credit	4 (40%)	10 (100%)	29	9	10
DEPOSIT - Cash account	1 (33.33%)	3 (100%)	4	3	3
DEPOSIT - Fixed Term Savings Deposit	1 (33.33%)	3 (100%)	2	3	3
DEPOSIT - Savings account	1 (33.33%)	3 (100%)	7	3	3
FIXED_TERM_DEPOSIT - Fixed Term Savings Deposit	1 (33.33%)	3 (100%)	2	3	3
INSTRUMENT - Credit	9 (50%)	18 (100%)	29	11	18
INSTRUMENT - Savings	9 (81.81%)	11 (100%)	13	5	11
INTEREST_BEARING_INSTRUMENT - Credit	6 (46.15%)	13 (100%)	29	8	13
INTEREST_BEARING_INSTRUMENT - Savings	11 (91.66%)	12 (100%)	13	5	12
MORTGAGE_LOAN_PART - Mortgage Loan Part	11 (55%)	20 (100%)	29	8	20
SAVINGS_ACCOUNT - Saving saccount	1 (33.33%)	3 (100%)	7	3	3
TOTAL	86 (49.14%)	175 (100%)	201	99	175

Data Coverage Optimization - Synthesized



<https://coverage.synthesized.io/>



Try NowRequest A Demo

Welcome to the SQL Query Coverage Analysis Tool!

Check your data coverage for software testing in minutes

Have you ever wondered what's the optimal dataset to use in development and testing and how "optimal" is even measured?

The data coverage assessment tool enables you to quickly check the suitability of your data for testing and development purposes based on the SQL language. It's fast and simple.

- Pick any SQL query from an application you develop
- Copy and paste the query into the command line interface below to analyse its cardinality for a given schema. This is parsed and analyzed automatically. In minutes, the tool analyses even the most complex SQL queries
- Upload a dataset you use in testing and development below to health check it and get a score. No code is needed

Copy & Paste SQL Query:

Input your query here (SQL 92), or pick one of the examples below.

Query ExampleData boundaries

```
1 SELECT
2   sum(I_extendedprice * I_discount) as revenue,
3   sum(CASE WHEN I_quantity > 26 then I_price * I_discount else I_price end) as large_order_revenue
4 FROM
5   I_lineitem
6 WHERE
7   I_shipdate BETWEEN '1998-01-01' AND '1998-01-01'
8   AND I_discount BETWEEN 0.05 AND 0.07
9   AND I_quantity < 24;
```

Resources

Blog post	How Weak Anonymization Became a Privacy Illusion	https://www.synthesized.io/post/how-weak-anonymization-became-a-privacy-illusion
Blog post	Will Your Data Pass the Test, or Will Your Test Pass the Data?	https://www.synthesized.io/post/will-your-data-pass-the-test
Podcast	Mind the Data Gap - Episode 1: Do We Want More Data or Better Data?	https://www.synthesized.io/webinars-podcasts/do-we-want-more-data-or-better-data

Thank you!



Ton Badal

ton@synthesized.io



@TonBadal



<https://github.com/TonBadal>