Siriraj Informatics and Data Innovation Center



## OHDSI Thailand Meetup #1



Thailand

OMOP
ETL Pipeline
by Siriraj

วันพฤหัสบดีที่ 26 มิถุนายน 2568 เวลา 12:30 - 13:30 น.



Join Zoom

https://bit.ly/omop-siriraj Meeting ID: 991 2848 3445

Passcode: omop



# Today's Agenda



26 June 2025 @ 12:30 - 13:30

12:30 – 12:40 ประชาสัมพันธ์ข่าวสาร OHDSI

**12:40 – 13:10 OMOP ETL Pipeline** โดย ณัฐวุฒิ อดุลยานุโกศล (Max) รองหัวหน้าศูนย์นวัตกรรมข้อมูลศ**ิ**ริราช SiData+

13:10+ ถาม-ตอบ



มีการบันทึก Zoom

และจะอัพโหลดบน https://www.youtube.com/@OHDSIThailand





## Submission Deadline Tuesday 1 July 2025

- Oct. 7-9 at New Brunswick, New Jersey
- https://www.ohdsi.org/ohdsi2025/

- Past submissions:
  - https://www.ohdsi.org/2024-collaborator-showcase/
  - https://www.ohdsi.org/2023-collaborator-showcase/





## **OHDSI APAC**

## https://www.ohdsi.org/apac/

- APAC Scientific Forums are held monthly every first Thursday at 10 am Thailand time
  - July 3: Generative Al-based OMOP mapping with HDR UK's Lettuce tool (TBC)
- APAC Community Calls are held monthly every third Thursday at 10 am Thailand time
  - July 17: European Symposium Recap, RWD Studies by European Medicines Agency (EMA)
- APAC Symposium 2025 December 6-7 in Shang Hai, China
  - Abstract submissions due September 7





## **OHDSI Thailand**

Next meetup: Thu 24 July "Overview of OHDSI"

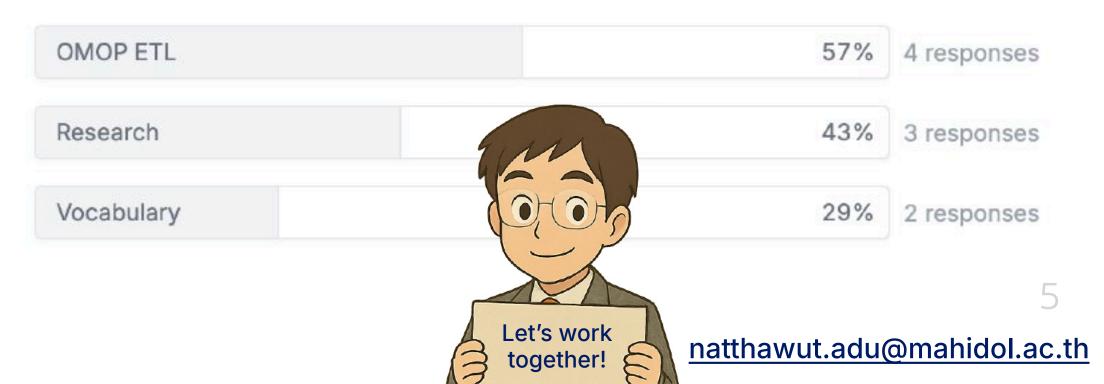


## **APAC Country Chapters**



Activities Community Support 2025
Research Funding ?

หากท่านสนใจเข้าร่วมเป็นผู้จัดกิจกรรมใน OHDSI Thailand Chapter กรุณาเลือกด้านที่ท่านถนัด 7 of 18 answered









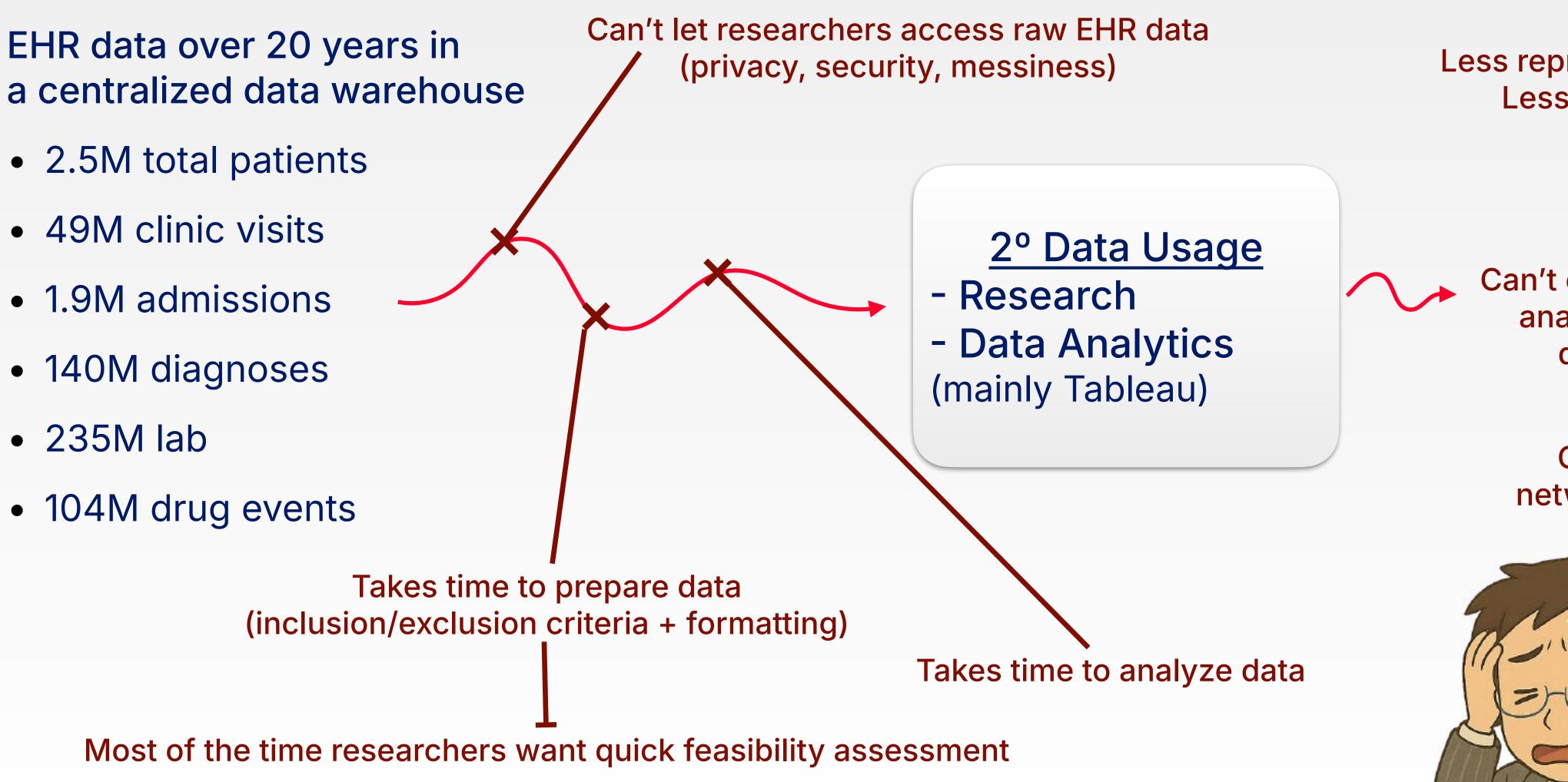
# OMOP ETL Pipeline

using on-premise SQLMesh



# Siriraj's Data and Challenges

pre-OMOP



Less reproducible research, Less trust in results Can't effectively share analytics scripts/ dashboards Challenging network research studies



Reproducible research,

More trust in results

# Siriraj believes OMOP CDM can help

Can't let researchers access raw EHR data (privacy, security, messiness)



ATLAS helps cohort generation
i2b2 helps feasibility assessment
GenAl helps SQL code generation

Takes time to prepare data (inclusion/exclusion criteria + formatting)

2º Data Usage

- Research
- Data Analytics
- Al

Use publicly available codes:
Official HADES R libraries,
Community contribution, Al libs

Takes time to analyze data

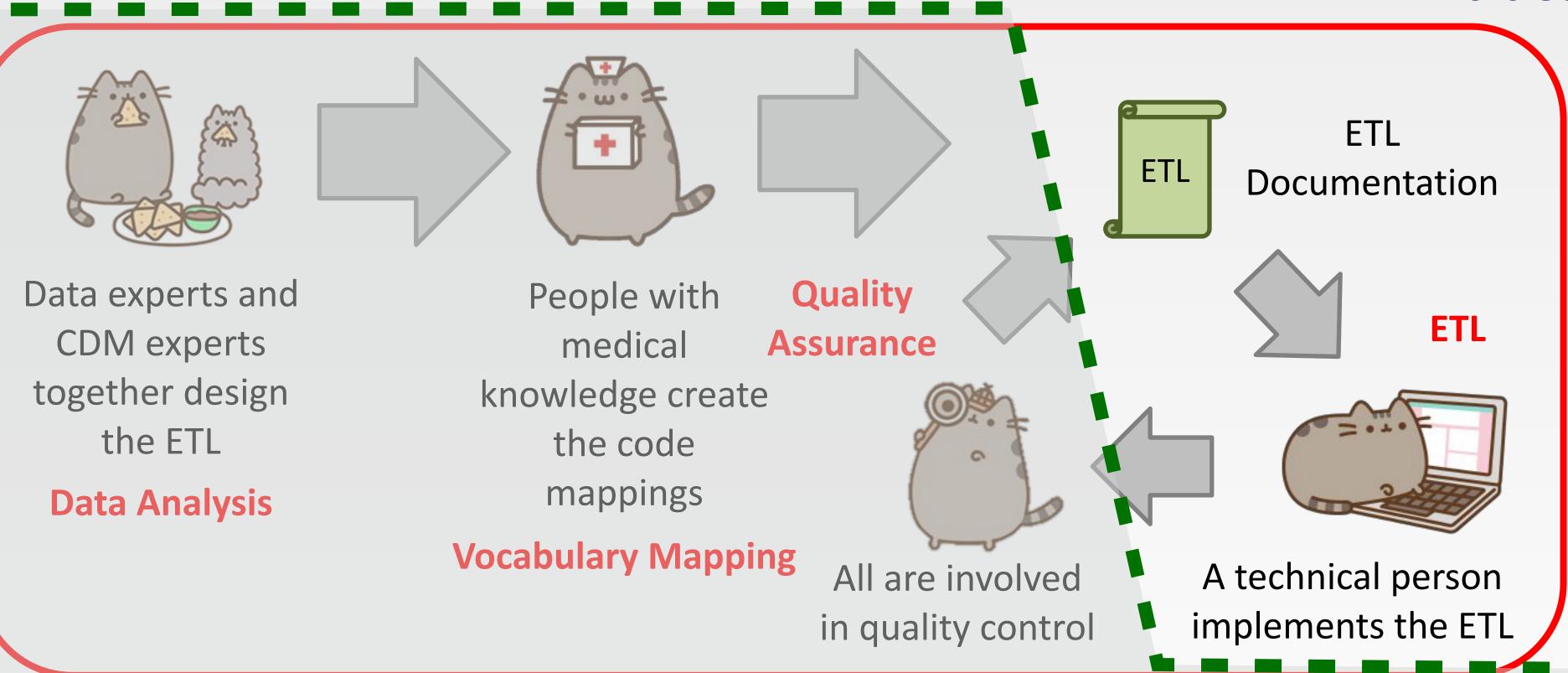
Sharable analytics scripts/dashboards **Easier network** research studies

Most of the time researchers want quick feasibility assessment



## **OMOP Data Conversion**

## ETL = Extract Transform Load

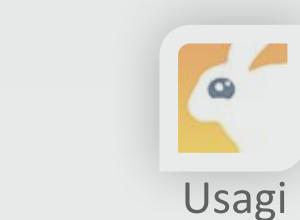


"เนื่องจากการแปลง ข้อมูลมีความซับซ้อนเซิง เทคนิคสูง (ประเภท DB, logics, security, 4a4) OHDSI จึงไม่มีเครื่องมือ มาตรฐานมาให้ครับ"









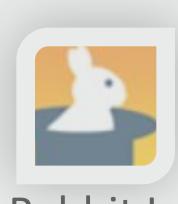


Rabbit









Rabbit In a Hat

# Siriraj ELT Pipeline

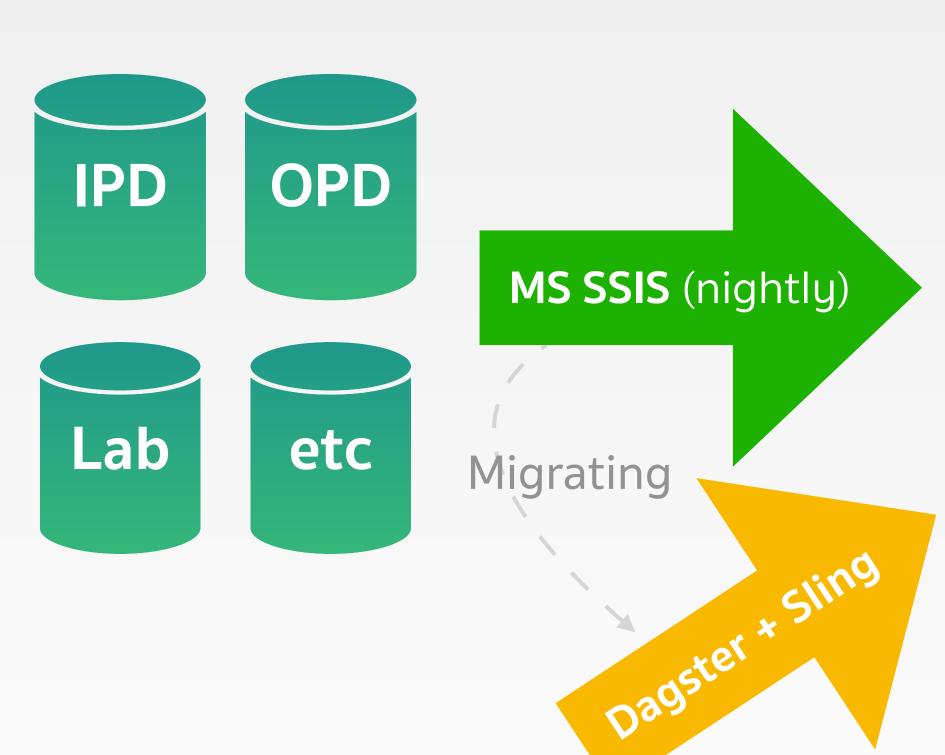
Extract - Load - Transform

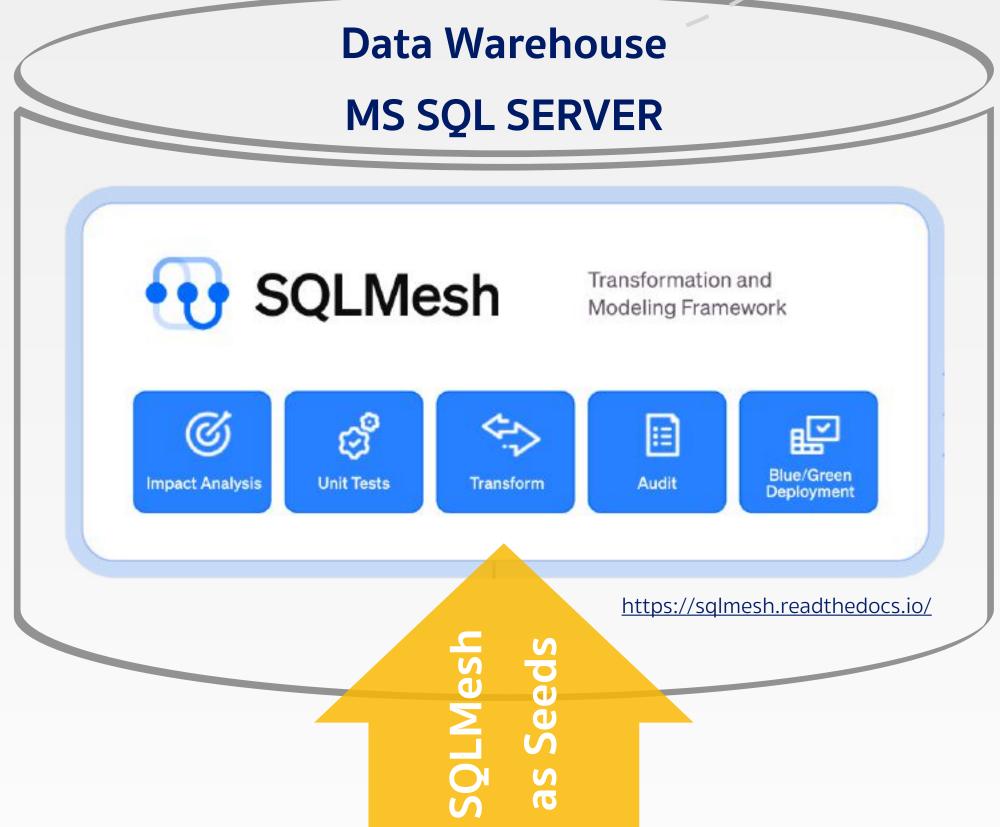
Data Lakehouse

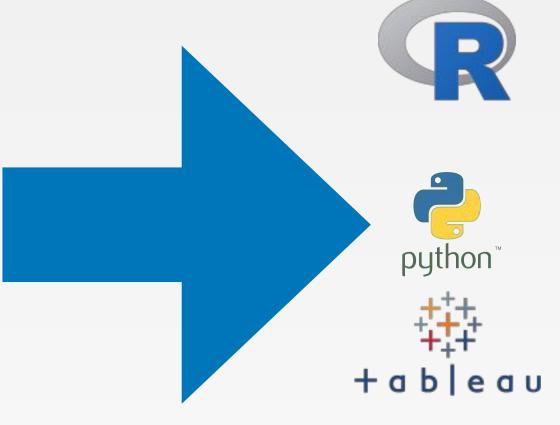
DuckLake

Migrating









OMOP Vocab from
Athena (.csv)

Vocab Mapping from Usagi (.csv)



## Who are you?

- 1. Data/IT/Dev
- 2. Researcher
- 3. Clinician
- 4. Student
- 5. Other

# Does your org have OMOP dataset?

- 1. Yes, on premise
- 2. Yes, on cloud
- 3. Developing
- 4. No
- 5. I don't know



Data Tool Landscape 1000



great\_expectations

SODA:

pandera V

lakeFS

Project Nessie

X dbt Core

"ไม่มีเครื่องมือที่ดีที่สุด ต้องพิจารณาเลือก ตามความเหมาะสม และข้อจำกัดครับ"



### Commercial



https://portable.io/learn/best-etl-tools



- OrientDB

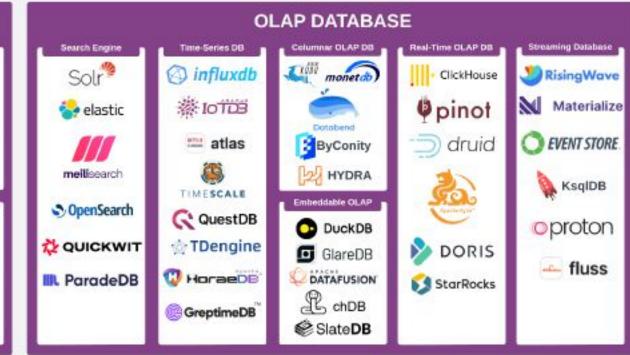
EDGE DB

ArangoDB

SurrealDB

**S**lanite

ReadySet





( influxdb

Mimir

KNOX

Apache Ranger

DATA INFRASTRUCTURE & MONITORING

ELK

graphite

WETRICS VICTORIA

ZABBIX

<sup>™</sup> Grafana

kibana

R console

logstash

VECTOR

**OCEANBASE** 

ShardingSpher

NEON

CrateDB

SQLite

MESOS

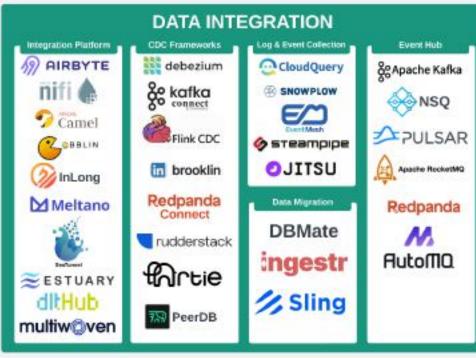
kubernetes

YUNIKORN

Apache Ambari

HELIX

in alirezasadeghi



💋 Dgraph

AG=

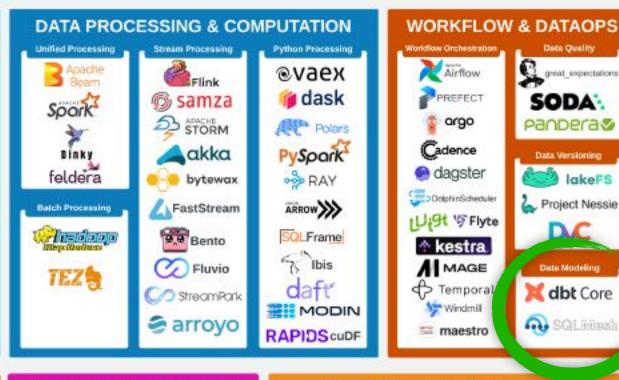
Falkor DB

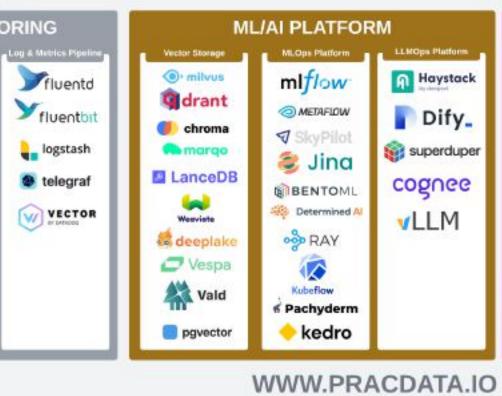
M Kvrocks

RocksDB

BadgerDB

MyRocks









# Siriraj: from dbt to SQLMesh

Using dbt—a free and open-source software framework-to transform data into OMOP CDM in the ETL process

RESENTER: Thanapat 'Thane' Pitchayarat thanapat.pit@mahidol.edu

### INTRO:

- The conversion of medical data into the OMOP CDM format requires a managed data engineering pipeline commonly referred to as the extract, transform, and load
- The main transformation tasks in a typical OMOP CDM conversion include combining data from multiple sources, changing the original data models to match the OMOP CDM, retrieving the concept IDs of source values, and mapping the source concept IDs to the standard IDs.
- . The complexity of the data transformation SQL scripts may grow rapidly beyond manageable. To keep the ETL pipeline maintainable, Siriraj Hospital uses dbt™ to transfrom its data to the OMOP CDM.
- dbt<sup>tM</sup> (shortened from data build tool) is a free and opensource software (FOSS) framework available at https:// www.getdbt.com. It could be applied to data transformation at other institutions.

The data conversion pipeline at Siriraj Hospital can be summarized as:

- 1. Extraction of data from hospital sources with Apache
- 2. Load the data into data lake and Development environment with Apache Spark
- 3. Transform the data to match the OMOP CDM specifications with dbt
- 4. Load the OMOP CDM-ed data into QA and Production

Each step is containerized with Docker. All steps are ochestrated and scheduled by Apache Airflow. Codes are version controlled with GitHub.

- · dbt comes with a command-line interface with commands that compile SQL scripts and execute the code on the connected database engines, as well as a graphical user
- . The core library of dbt is a Python package that supplements traditional SQL scripts with Pythonic Jinja templating.
- With the Jinja templating,
- any frequently used SQL command can be packaged as a modular macro that can take parameters similar to a Python function, and;
- . the Jinja tags enable data lineage tracking that can be visualized on an interactive web application generated by dbt command. The web application referred to as dbt documentation also presents metadata, such as upstream and downstream tables. The metadata are partly generated automatically and can be added manually as YAML files.
- To verify data quality, dbt can run automated tests during transformation execution or on demand.
- . Given the popularity of dbt in the enterprise analytics space, there are many tools that can be integrated with dbt, namely Airflow for data pipeline orchestration, GreatExpectations for data quality, and DataHub for data

"An organized approach to build a maintainable ETL pipeline for the OMOP CDM with minimal cost while keeping our data engineers sane "

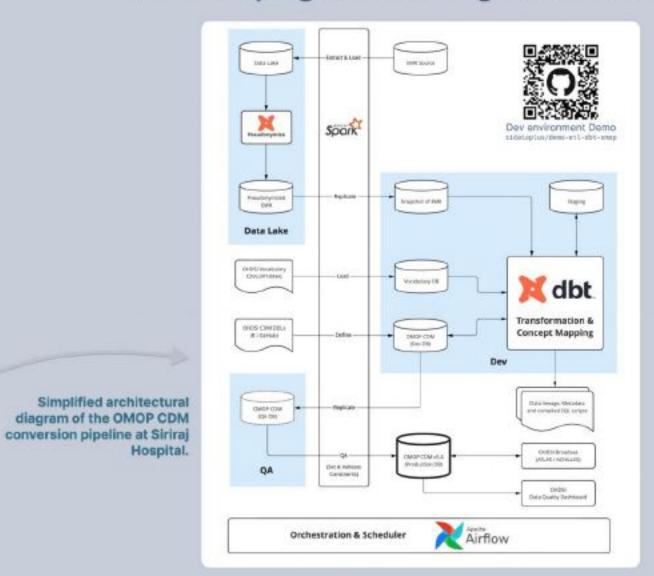




Table data lineage automatically generated by dbt. Each node represents a table or a view of data. Each linking edge represents a data flow from the source(s) to its destination(s), with data transformation in between. Each of the data transformation step is programmed as an SQL SELECT script.



Sirirai Informatics and Data Innovation Center







Simplified SQL snippets (a) to create the CDM PERSON table with data from a staging table joined with the vocabulary concept tables via macros (b) to set a macro template for concept mapping. These SQL snippets with Jinja tags are to be compiled and submitted to the database engine by dbt.

### CONCLUSION:

- . dbt is a promising free and open-source software framework that massively facilitates the data conversion process into OMOP CDM.
- . dbt programmatically manages the SQL transformation scripts in the ETL process, and consequently enhances the maintainability of the data pipeline
- . Data engineers with proficiency in SQL and Python could learn dbt in a few days and probably take a few weeks to implement dbt in the pipeline.

- 1. dbt Labs, Inc., dbt-core [Internet], 2022. Available from https://github.com/
- 2. OHDSL WhiteRabbit [Internet], 2022, Available from: https://github.com/
- OHDSL Rabbit-in-a-Hat [Internet], 2022, Available from http://ohdsi.ghtub.io/WhitaRabbit/RabbittnAHat.html
- A. OHDSI, Usagi Enternel]. 2022. Available from: https://github.com/OHDSI/
- 5. The Pallets Projects, Jinja (Internet), 2022, Available from: https://
- 6. Superconductive Health, Inc., Welcome to great expectations [internet], 2022. Available from: https://greatexpectations.in
- 7. Calogica, dbt\_expectations [Internet], 2022, Available from: https://
- bub.getdbt.com/calogica/dbt\_expectations/0.1.2 8. dbt Labs, Inc., Success Stories [Internet], 2022. Available from: https://
- Apache Software Foundation, dot Cloud Operators (Internet), 2022, Assisible
- rom: https://eirflow.epache.org/docs/apache-airflow-providers-dist-cloud/
- DataHub Project. dot [Internet], 2022. Available from: https://

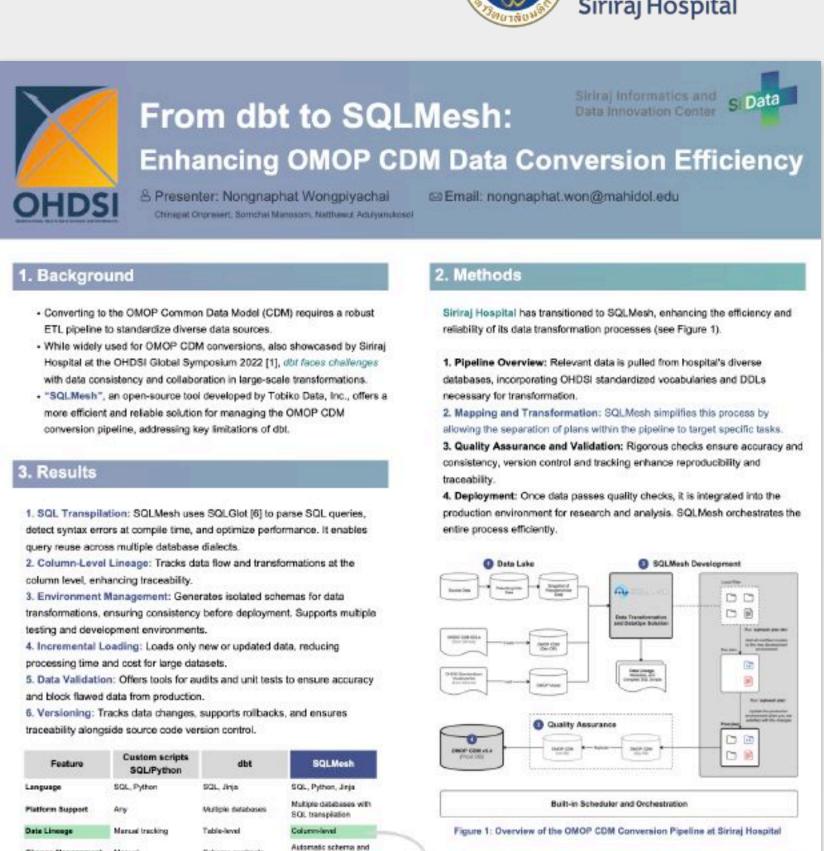
Thanapat Pitchayarat, Gun Pinyo, Watcharaporn Tanchotsrinon, Somkid Khamsrimuang, Chalita Issarasittiphap, Chaiyanun Bootnumpech, Noppon Siangchin, Kanphitcha Promma, Nattachai Bovornmongkolsak, Prapat Suriyaphol,

Natthawut Adulyanukosol Siring Informatics and Data Innovation Center (SiData+), Faculty of Medicine Sirinaj Hospital, Mahidol University, Bangkok,









Can be costly for large Efficient for large Table 1: Side-by-side comparison of data transformation tools Figure 2: Visualization of data lineage automatically generated by SQLMesh SQLMesh: Streamlining OMOP CDM Conversion

clata contracts.

Doponds on

SQLMesh has optimized data transformation at Siriraj Hospital,

boosting efficiency and reliability. This transition standardizes data

for research while offering developers enhanced tools.

Wongpiyachai\_From-dbt-to-SQLMesh.pdf

Built-in ánduding unit



© 2024 Siriraj Informatics and Data Innovation Center, Siriraj Hospital, Mahidol University, Thailand

# Data Modeling Tools

# dbt & SQLMesh ทำด้านปั้นข้อมูลโดยเฉพาะ





Tools	ความเหมาะสม	ข้อจำกัด
1. Custom scripts, e.g., SQL, Python, R, Java	ต้องการอะไรก็เขียนเอาเองได้เลย flexible ที่สุด	เหนื่อยทั้งเขียน ดูแล ซ่อม technical debt บาน
2. Commercial ETL tools, e.g., MS SSIS, Talend, Pentaho	ถ้ามีคนใช้เครื่องมือเหล่านี้เป็นอยู่แล้ว ในองค์กรก็ง่าย	ถ้าไม่มี ก็เสียตังค์ เสียเวลา ส่วนใหญ่เป็น GUI เยอะ maintain ยาก
3. Cloud ETL tools, e.g., Azure Data Factory, AWS Glue, GCP Cloud Dataflow, Dataform	ถ้าใช้ Cloud เจ้านั้น ๆ อยู่แล้ว ก็ integrate ได้ดี	ไม่ได้เกิดมาเพื่อปั้นโมเดลข้อมูลโดยเฉพาะ ขาด features ช่วยเหลือ
4. <b>abt</b>	เป็น free open source library ตัวแรก ๆ ที่มี features data modeling ครบครัน	ไม่เข้าใจ SQL โดยตรง เป็นแค่ Jinja templating engine หลาย ๆ features ต้องเสียตังค์เพิ่ม
5. SQLMesh	พัฒนาแก้ไขข้อจำกัดของ dbt เข้าใจ SQL semantic ແລະເປิດฟรีเกือบทุก features (slides 17)	ยังไม่ stable version (v0.196) มี learning curve (แต่ไม่ยากกว่า dbt)

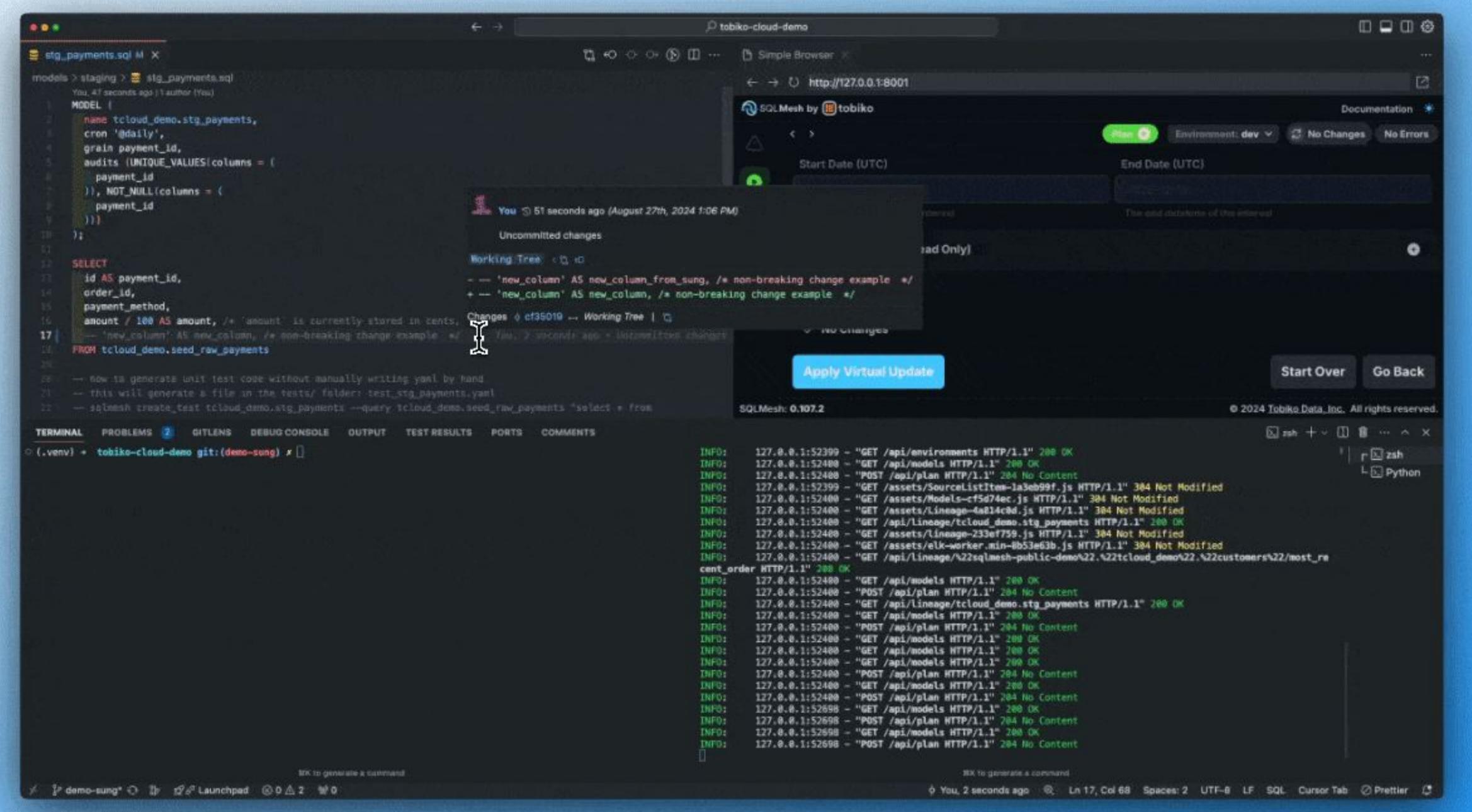


# แล้ว SQLMesh ใช้งานยังไง (แบบย่อ ๆ)

แบบเต็ม ๆ (1 ชั่วโมง): <a href="https://www.youtube.com/watch?v=J0\_2hkXz-HY">https://www.youtube.com/watch?v=J0\_2hkXz-HY</a>

## https://github.com/sidataplus/demo-etl-sqlmesh-omop-synthea/

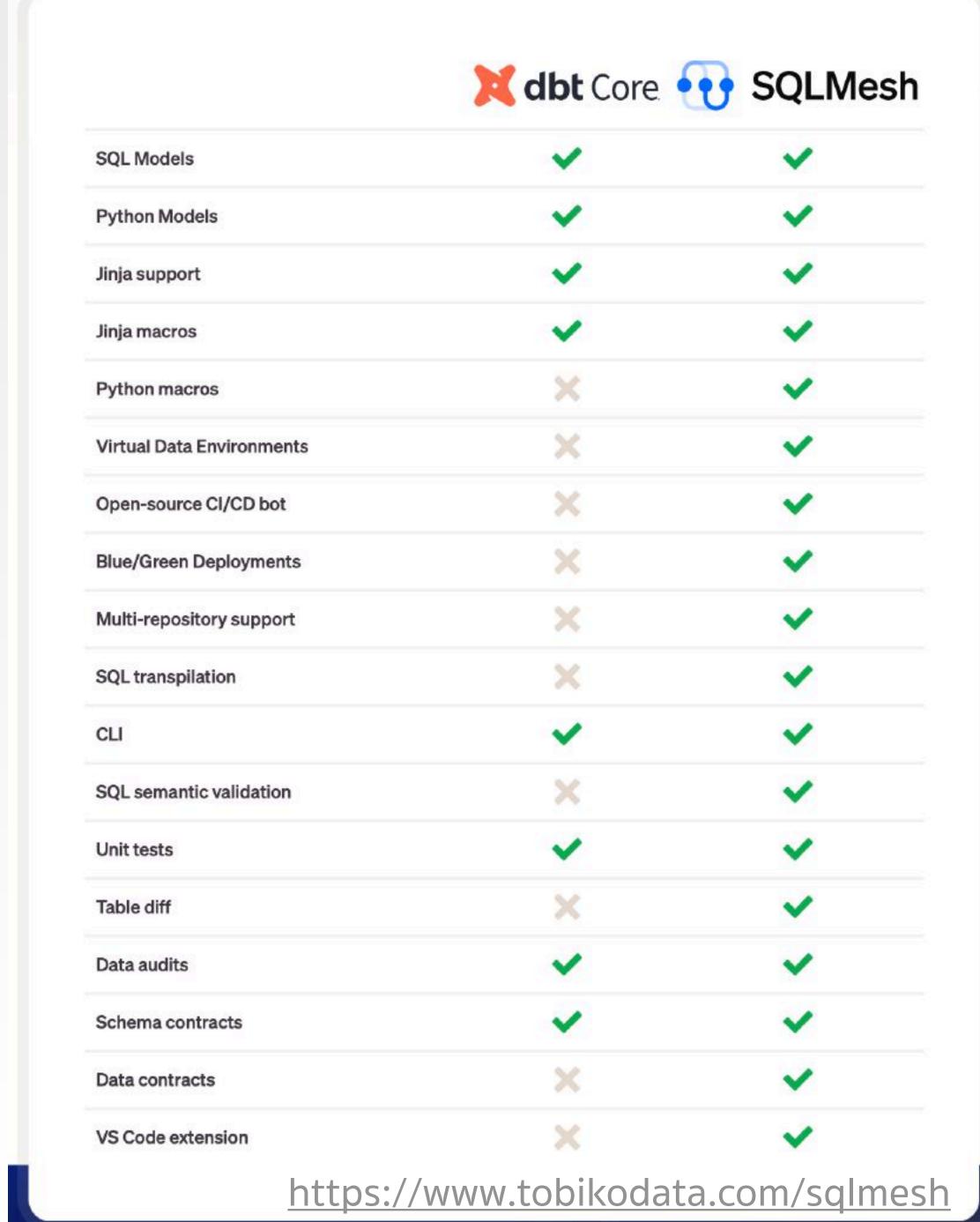
- 1. สร้าง Python environment แล้ว pip install sqlmesh (แต่ตอนแปลงข้อมูลใช้แค่ SQL ก็ได้)
- 2. ตั้งค่าการเชื่อมต่อกับ database ใน config.yaml
- 3. เขียน SQL ในการแปลงข้อมูลตามที่ต้องการแต่ละ table
- 4. เขียน Model Definition (metadata) ในไฟล์ .sql เดียวกัน บอกให้ SQLMesh จัดการข้อมูลอย่างไร เช่น incremental load, data type, audit
- 5. เขียน SQL table อื่นต่อไปเรื่อย ๆ แล้ว SQLMesh จะนำไปต่อกันให้เป็น DAG Lineage
- 6. สั่ง `sqlmesh plan dev` เพื่อให้ SQLMesh render SQL scripts แล้วรันบน DB ตาม DAG (ไม่มีข้อมูลเข้าเครื่อง dev) บน dev environment (schema omop\_\_dev.[table\_name]) ระหว่างที่รันสามารถให้เช็ค audit ตรวจสอบคุณภาพข้อมูลไปพร้อมกัน
- 7. เช็คข้อมูลใน dev env จนพอใจ แล้วรัน `sqlmesh plan` เพื่อให้ promote เข้า prod env (schema omop.[table\_name])
- 8. ตั้ง schedule อัพเดทข้อมูลด้วย `sqlmesh run`



# Features ช่วยการแปลงข้อมูล

## ที่ tools อื่นไม่มี

- 1. Stateful Virtual Data Environment ช่วยแยก dev, prod ให้ โดย ยังอยู่ใน DB เดียวกัน หรือจะแยก DB (multiple gateway) ก็ทำได้ เก็บ state history ทำให้ dev หลายคนทำงานร่วมกันได้อย่าง ปลอดภัย
- 2. SQL Semantic Understanding เข้าใจ SQL ช่วยให้เจอ bug ก่อน ไปรันบน DB ช่วยให้ track change by column สร้าง column lineage ได้
- 3. SQL Transpilation ใช้ SQLGlot แปล syntax ข้าม DBMS ได้
- **4. Data contracts** เช็คว่าการแปลงข้อมูลสมบูรณ์ทั้ง pipeline ผ่าน audits ถึงจะเปิดให้ access ได้ (dbt แปลงได้แค่ไหนก็ปล่อยเลย)



# Demo Siriraj





- 1. SQLMesh Docs Get Started: <a href="https://sqlmesh.readthedocs.io/en/stable/quick\_start/">https://sqlmesh.readthedocs.io/en/stable/quick\_start/</a>
- 2. **Demo Synthea -> OMOP:** <a href="https://github.com/sidataplus/demo-etl-sqlmesh-omop-synthea/">https://github.com/sidataplus/demo-etl-sqlmesh-omop-synthea/</a> ใช้ synthetic data บน DuckDB ทำตาม instructions ใน README ได้
- 3. เริ่มแปลงข้อมูลจาก table ง่าย ๆ / หลัก ๆ ก่อน



8.death	
9.condition_occurrence	
10.observation	
11.procedure_occurrence	
12.drug_exposure	
13.condition_era	
14.drug_era	
15.measurement	
16.cost	
17.payer_plan_period	