



UNLOCKING THE TREASURE CHEST OF DATA

A QUEST WITH CLICKHOUSE[©] & DBT[™]

EMIR KARAMEHMETOGLU

DATA ENGINEER & DUNGEON MASTER



AXIS COMMUNICATIONS

DIAGNOSTIC DATA MANAGEMENT

CHARACTER BACKGROUND: ASTRONOMER



Translation: **Self-service analytics**
data-products with **ClickHouse** and
Dagster orchestrated **dbt**

OUR ADVENTURE MAP



: Agenda

- **TREASURE CHEST:** SELF-SERVICE ANALYTICS DATA PRODUCTS
- **ARCANE LORE:** DATA PRODUCT? PLATFORM?
 - **THE DRAGON'S LAIR:** THE CHALLENGES
- **SWORD AND SHIELD:** OUR SOLUTION
 - **BATTLE PLANS:** TECHNICAL ARCHITECTURE
- **SLAYING THE DRAGON:** PROBLEMS OVERCOME
 - **A DEAD END?:** CLICKHOUSE-DBT PLUGIN AND DISTRIBUTED DATABASE ARCHITECTURE
 - **TREASURES AND TRAPS:** LESSONS LEARNED
- **NEW HORIZONS:** FUTURE ROADMAP

ClickHouse and dbt - A Gift from the Community

The ClickHouse Team
Feb 1, 2023

Take this!

+



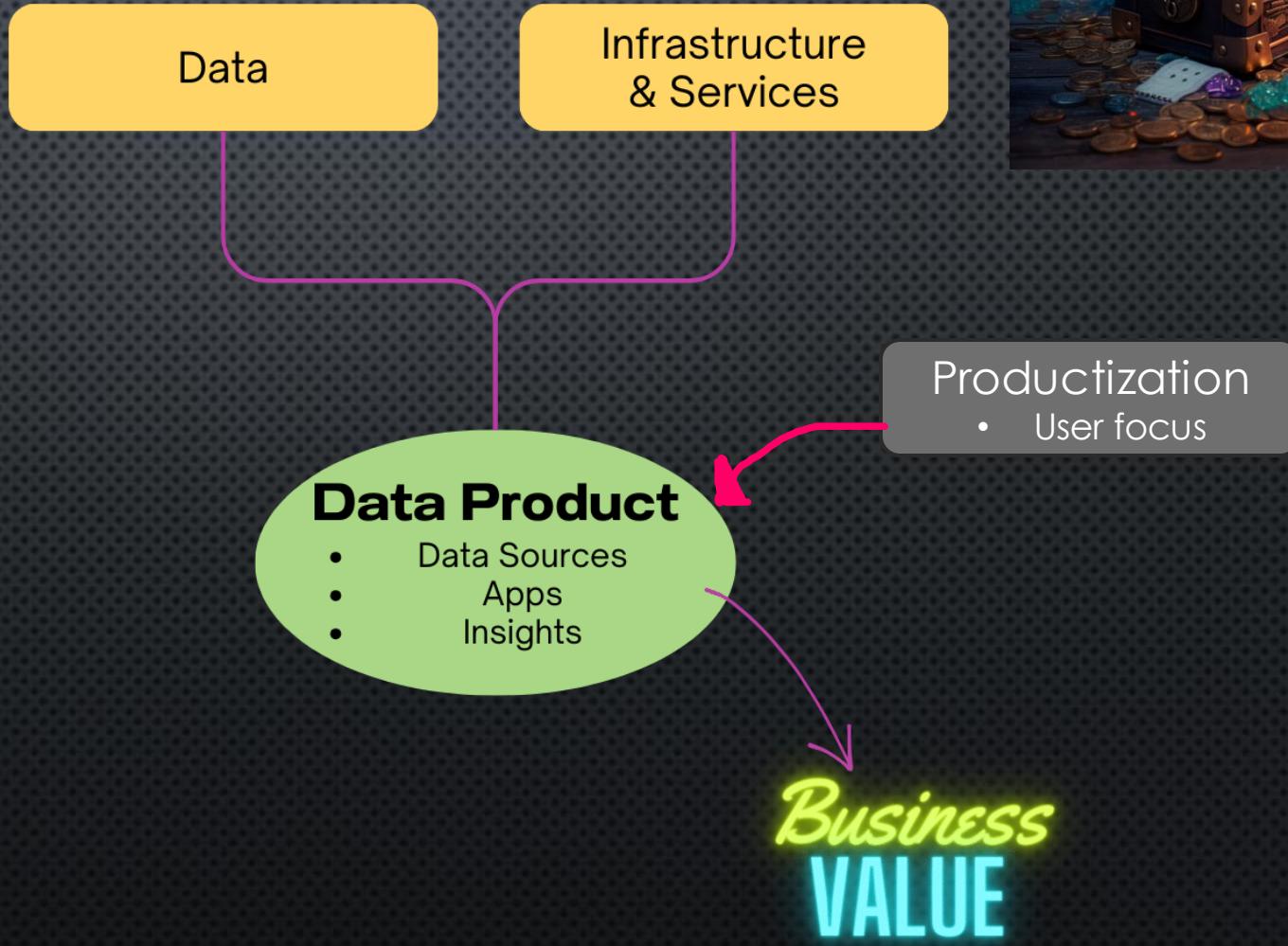
THE TREASURE CHEST

- UNLOCK VALUE OF DATA
 - HALT SOUL-CRUSHING DATA DIGGING EXPEDITIONS
- DEMOCRATIZE (AND MODERNIZE) DATA DRIVEN DECISION MAKING
- ELEVATE DATA MATURITY

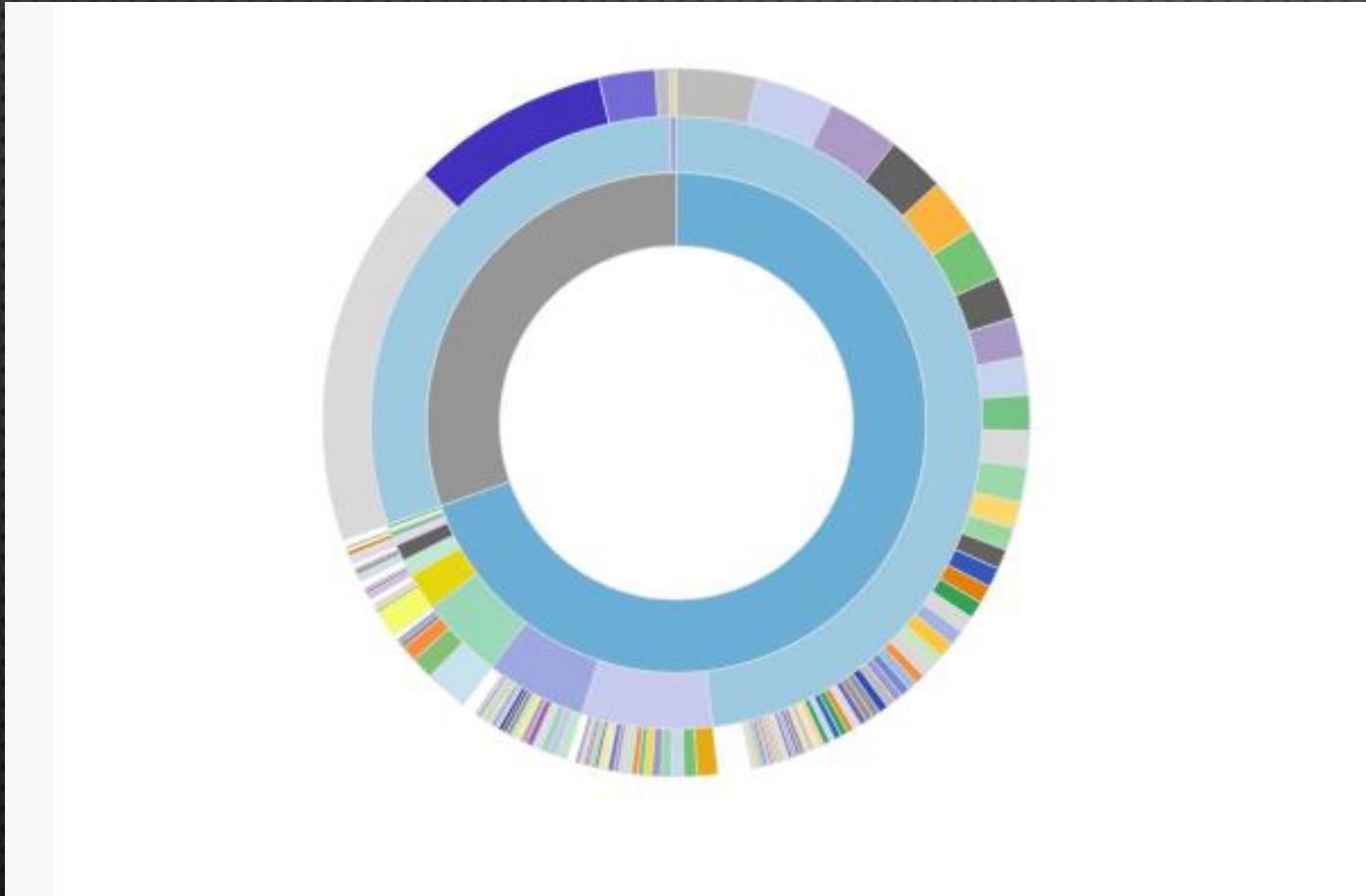


ARCANE LORE: DATA PRODUCTS

- BI DASHBOARD
- DATA API
- ML MODEL
- DATA VISUALIZATION
- ANALYTICS APP
- "QUERYABLE" DATASET



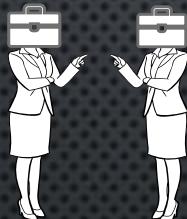
BI LAYER VISUALIZATION



WHAT IS IN A GOOD DATA PRODUCT?

DATA PRODUCT IS..

- DISCOVERABLE
- TRUSTWORTHY
- SELF-DESCRIBING
- INTEROPERABLE
- ACCOUNTABLE
- CONSUMED
- ADDRESSABLE & SECURE



MEANING..

- SEARCH, FIND, GET DOCS.
- DATA QUALITY. **OBSERVABILITY**. VALIDATION. METADATA.
- DESCRIPTION. SCHEMA. NO SHERPAS NEEDED.
- STANDARDS (OPEN). UNIFIED NAMING.
- CLEAR OWNERSHIP. RESPONSIBILITY. CONTACT.
- COMMUNICATION WITH STAKEHOLDERS. USER-FOCUS.
- ACCESS & USE, SECURELY.

VALUEABLE

~~our team~~

~~DDM~~ CHOSE CLICKHOUSE FOR ITS ANALYTICS DATA PRODUCTS

- OPEN SOURCE
- ACTIVE COMMUNITY & DEVELOPMENT
- OPEN STANDARDS (SQL)
- LAUGHS AT SCALE & SPEED



MASTER'S THESIS 2022

Evaluating ClickHouse as a Big Data Processing Solution for IoT-Telemetry

Adrian Göransson, Oskar Wändesjö

Compression Selection for Columnar Data using Machine-Learning and Feature Engineering

*Master's thesis submitted in partial fulfilment
of the requirements for the degree
of*

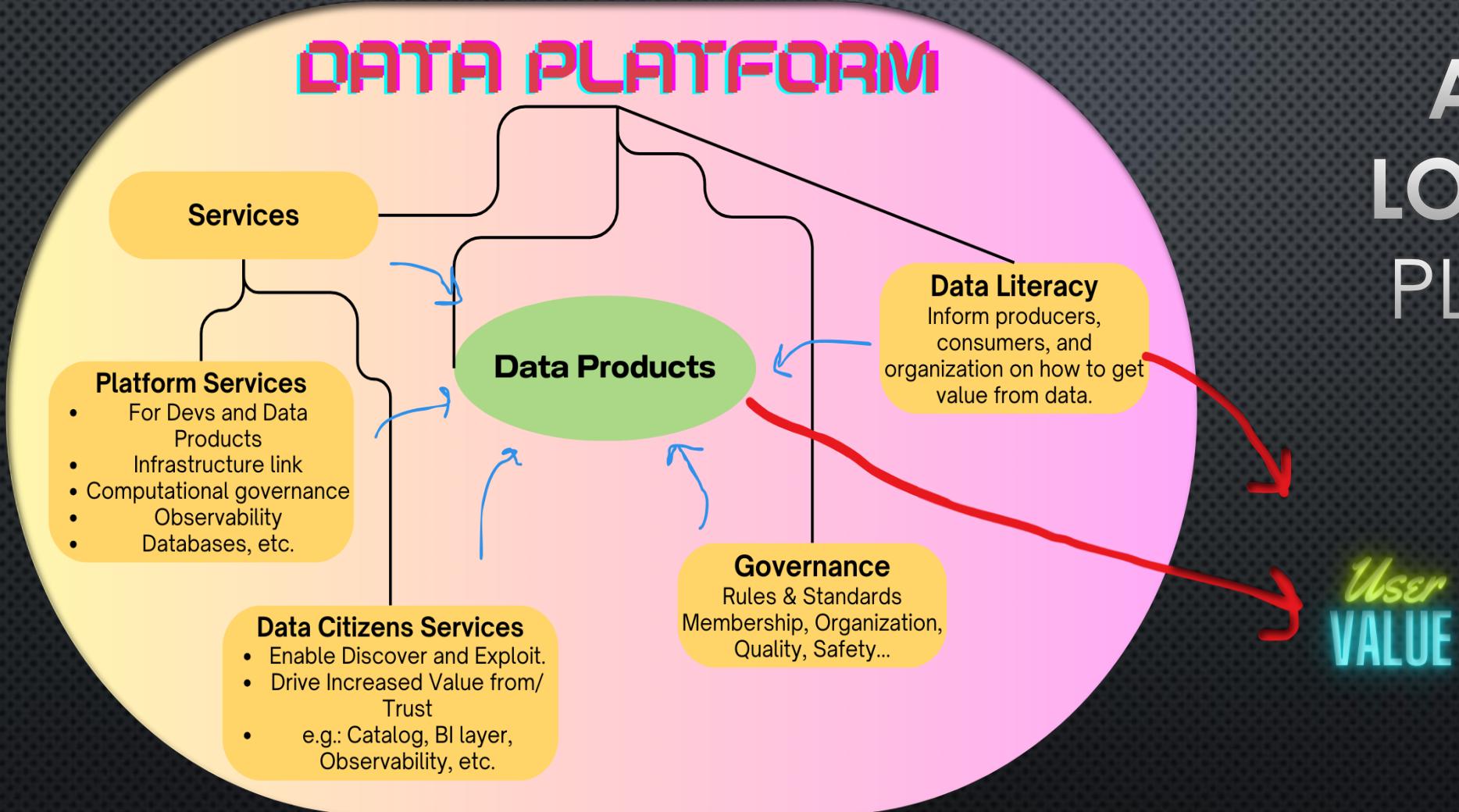
Master's in Applied Data Science
by

Ludvig Juelsson Larsen | Douglas Persson

Under the supervision of
Emir Karamehmetoglu | Oskar Wändesjö | Bengt J. Nilsson

And examined by
Johan Holmgren

ARCANE LORE: DATA PLATFORM



THE DRAGON'S LAIR

- PRODUCTION CLICKHOUSE CLUSTER WITH HUNDREDS OF TRILLIONS OF ENTRIES
 - SHARDED & REPLICATED.
 - TABLES WITH TENS OF BILLIONS OF ROWS.
- NEED FOR DATA TRANSFORMATION AND MODELING.
- DATA ORCHESTRATION, OBSERVABILITY, AND REPORTING.
- BI-LAYER NEEDS RAW & TRANSFORMED DATA ACCESS.
- REQUIREMENTS:
 - SPEED.
 - DATA QUALITY & GOVERNANCE.
 - MINIMIZE DATA DOWNTIME.



THE DRAGON'S LAIR

- SHARDED REPLICATED*MERGETREE TABLES
 - EXISTING DATA SOURCE PRODUCT HAVE INTELLIGENT SHARDING KEY.
 - HARD TO REPLICATE FOR USERS
- SERVICES & USERS CAN ACCESS ANY NODE IN A LOAD BALANCED WAY
- HOW TO SMOOTH OVER THE DIFFICULT UNSOLVED PROBLEM OF DISTRIBUTED DATABASE DESIGN FOR SELF-SERVICE?



QUEST REQUIREMENTS

STAKEHOLDERS

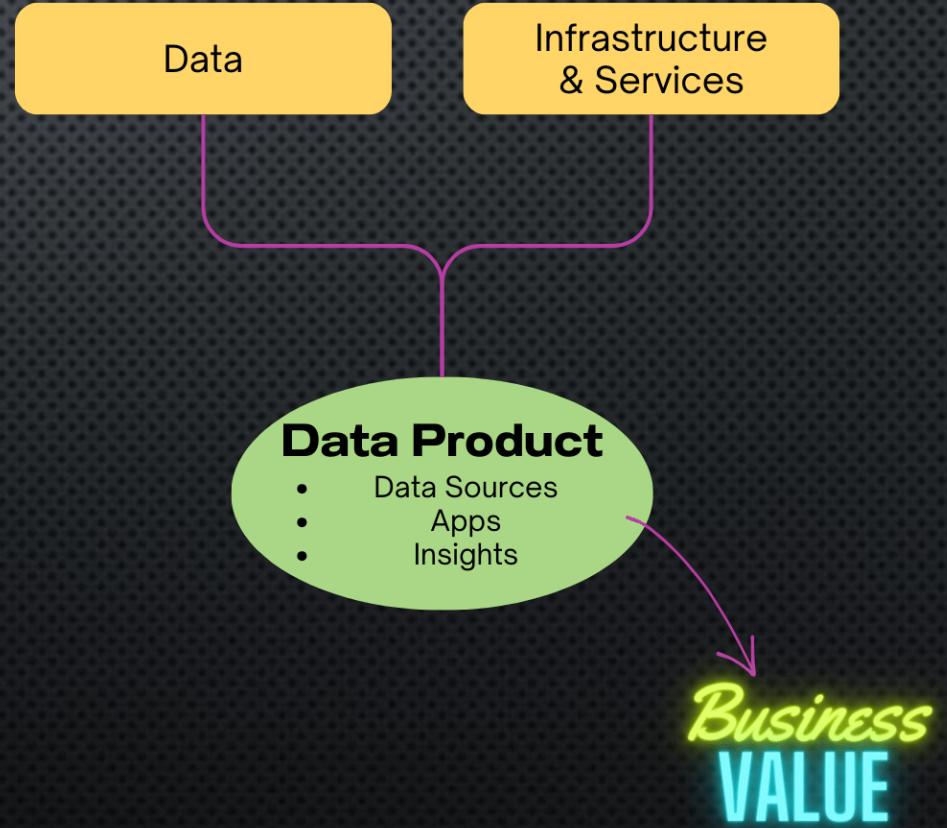
- INTERNAL FACING
- DATA ANALYSTS/SCIENTISTS
- (FUTURE) DATA PLATFORM CITIZENS
- (FUTURE) DATA STEWARDS AND PRODUCT OWNERS.

REQUIREMENTS

- **DATA TRANSFORMATIONS & DATA MODELING**
- DATA OPS WAY OF WORKING WITH ANALYSIS
- X AS CODE. CI/CD.
- **SERVE PRODUCTION MODELS AND DEV DATA EXPLORATION**
- **SELF-SERVICE**
- **DATA QUALITY**
- OBSERVABILITY
- COMPOSABILITY & **INTEROPERABLE**
- **FASTER DASHBOARDS**
- **SELF SCHEDULED**
- **ALERTS** AND MONITORING
- DOCUMENTED AND EASY TO CREATE DOCUMENTATION
- **PYTHON AND SQL.**
- HOOK UP TO DATA SOURCES OUTSIDE THE PLATFORM
- SERVES FIKA AND WHATEVER THE ANALYSTS NEED.

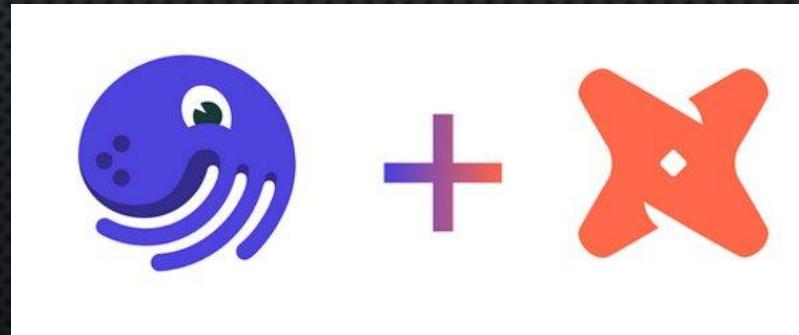
SELF-SERVICE ANALYSIS DATA PRODUCTS – JUST ESSENTIALS

- DEMOCRATIZE ACCESS
- TRANSFORMATION & MODELING
- MAKE DBT WORK WITH (OUR) CLICKHOUSE
- SQL & PYTHON
- AUTOMATE & SCHEDULE



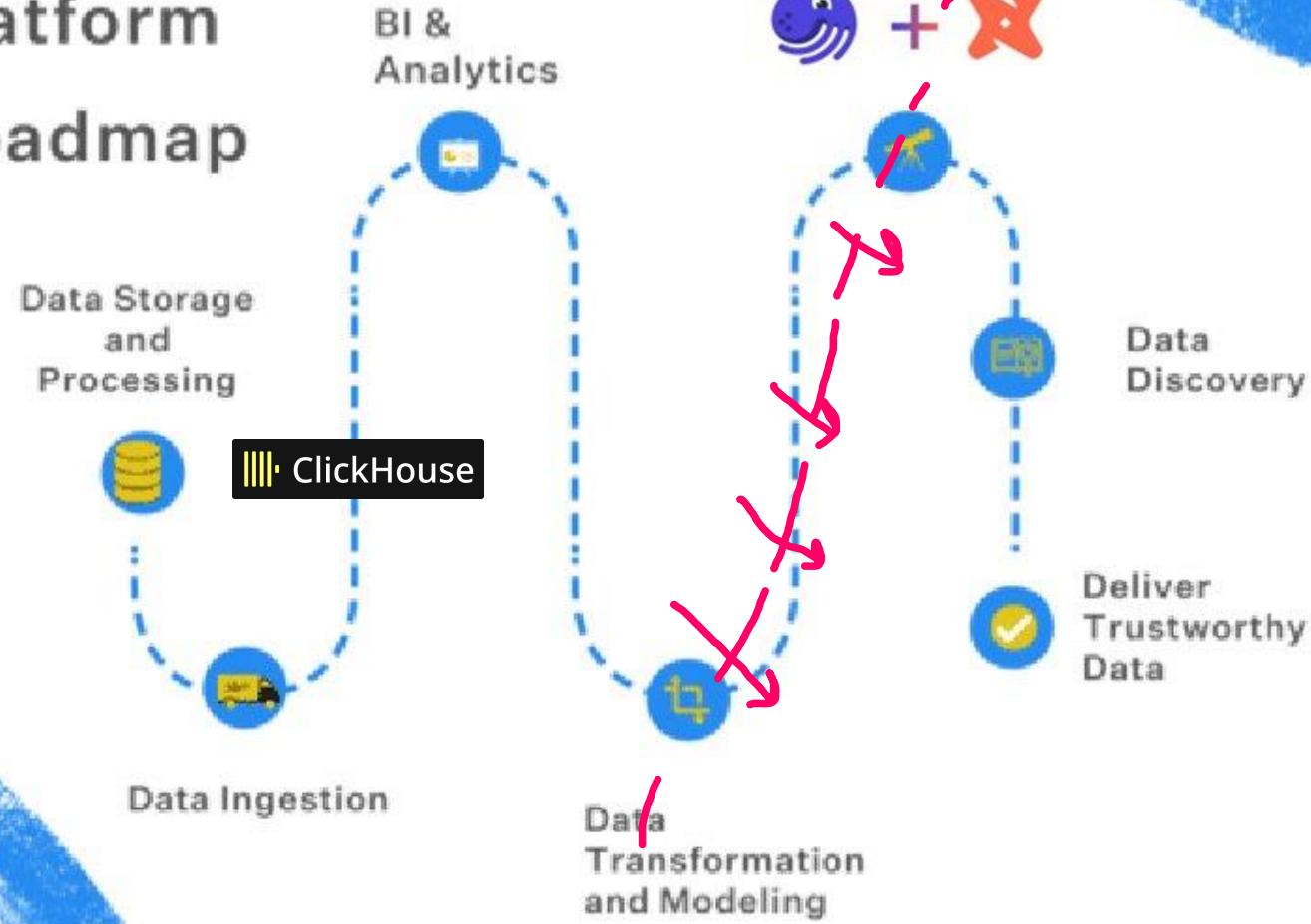
DON'T GO ALONE, TAKE THIS!

- DAGSTER ORCHESTRATED DATA BUILD TOOL (DBT) PLATFORM
- INTEGRATION WITH DATA ORCHESTRATOR DAGSTER
- SELF-SERVICE DBT + PYTHON DATA PLATFORM FOR 'ANALYTICS ENGINEERING'

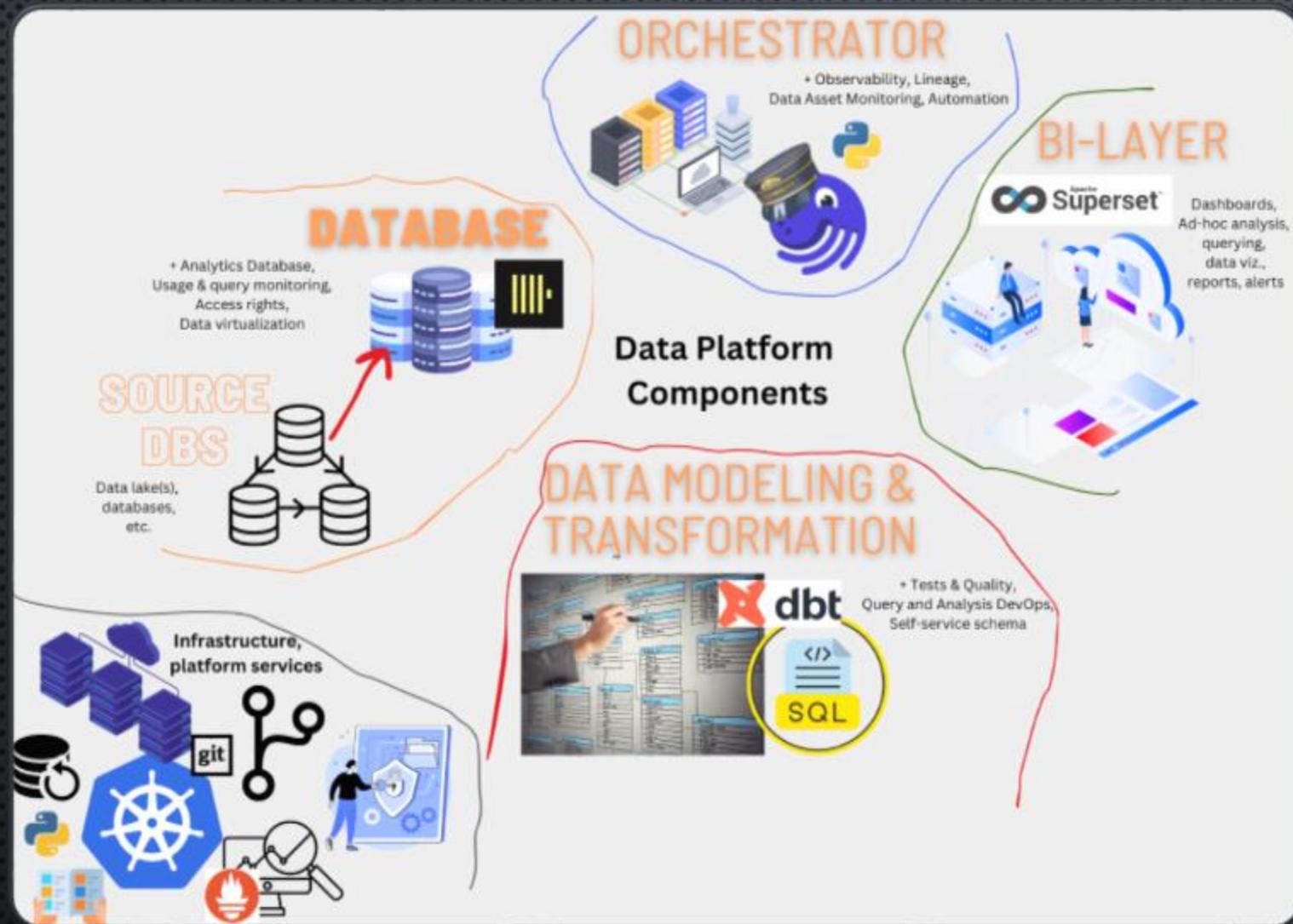


BATTLE PLANS

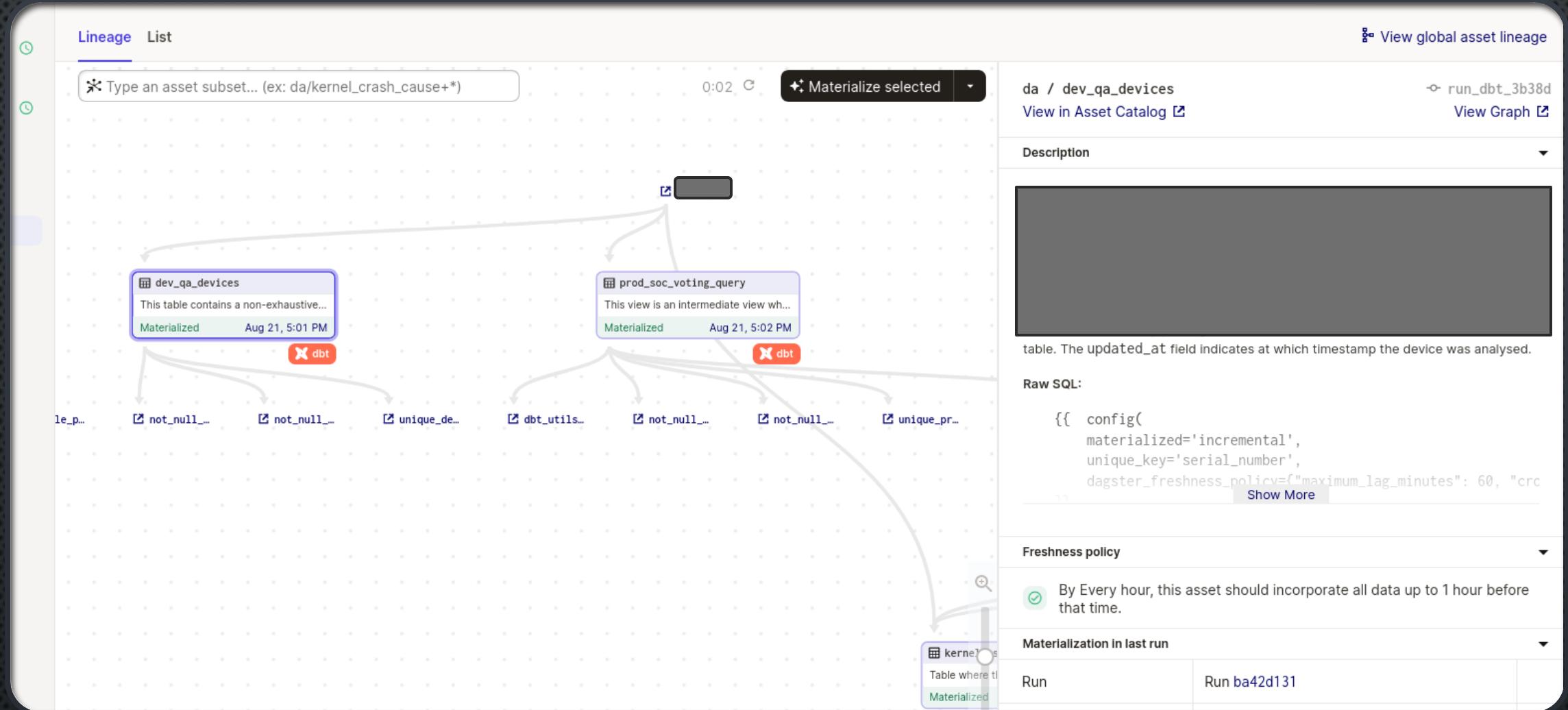
Data Platform Roadmap



BATTLEPLANS



- CLICKHOUSE: FAST ANALYTICAL DATABASE
- DBT: DATA TRANSFORMATION AND MODELING
- DAGSTER: DATA ORCHESTRATOR
- SUPERSET: BI DASHBOARD



DBT + DAGSTER DATA TRANSFORMATION AND MODELING PLATFORM

BY ORCHESTRATING DBT WITH DAGSTER: AUTOMATED, OBSERVABLE, DOCUMENTED, AND EVEN ALLOWS CONSUMERS OF THE ANALYSIS DATA PRODUCTS TO SELF-SERVICE (MONITOR OR TRIGGER DATA REFRESH, SEE ISSUES, STATUS, TESTS/QUALITY, DOCUMENTATION)

SLAYING THE DRAGON

- CLICKHOUSE-DBT PLUGIN LIMITATIONS
- SHARDED AND DISTRIBUTED SETUP
- NEW 'CLUSTER' WORKAROUND
- TRADE-OFFS



genzgd commented on Jan 24

Contributor ...

Hi @simpl1g -- I believe you are correct. I've been thinking too much in ClickHouse Cloud terms which doesn't require sharding, but even with a Replicated Database it looks like you need a distributed table for multiple shards. We'll take another look at this issue and see where it fits on our roadmap.

1 3 1

TREASURES AND TRAPS

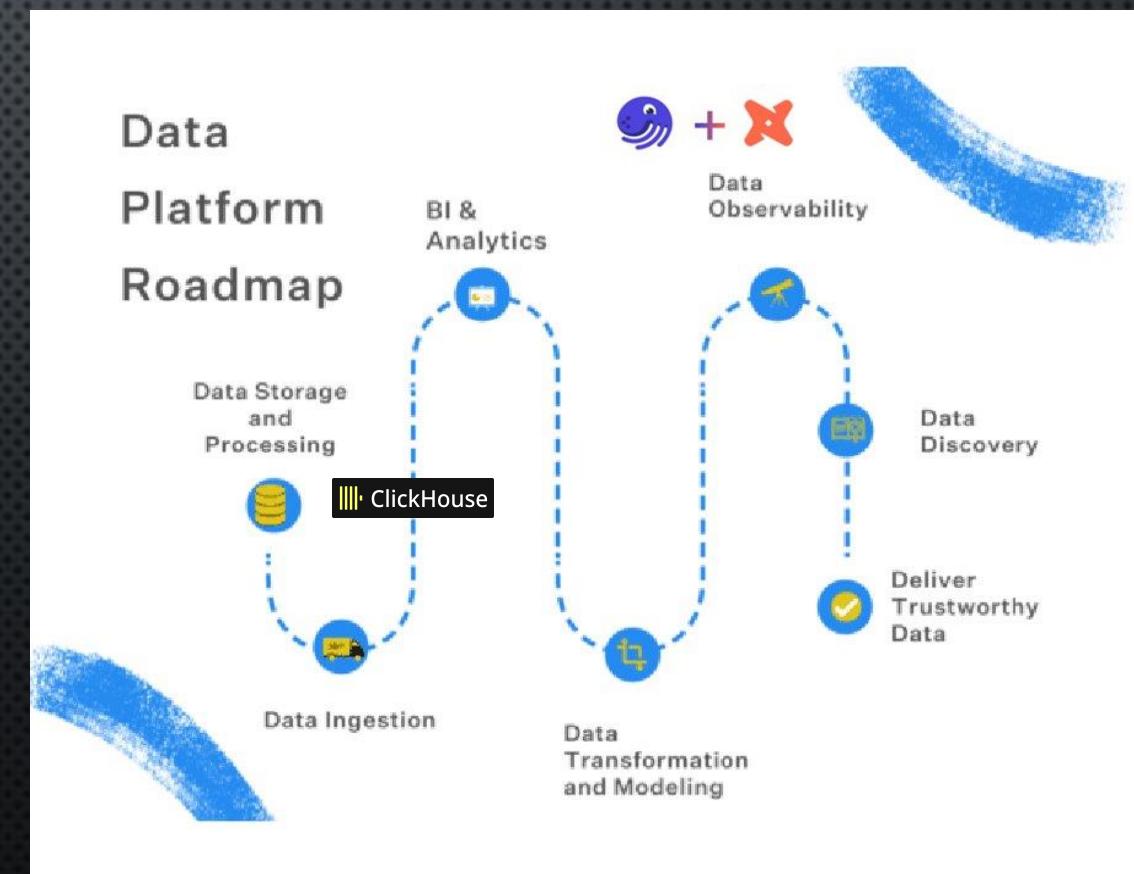
- IMPORTANCE OF COMMUNITY SUPPORT, OPEN SOURCE.
 - WHAT WE WANT DOES NOT EXIST AS A PRODUCT, BUT WE CAN BUILD IT.
- DATA AND CODE VERSIONING
- OBSERVABILITY AND REPORTING
- GOVERNANCE IS HARD

WHAT WAS NOT AN ISSUE

- DATA DUPLICATION
 - COMPRESSION IS HIGH LEVEL MAGIC.
 - (IN FACT, IT PROVIDED SOME USER VALUE).

FUTURE ROADMAP

- RAPID PROGRESS IN THE CLICKHOUSE-COMMUNITY PLUGIN
 - DISTRIBUTED TABLES
- AREAS FOR FURTHER SUPPORT
 - DISTRIBUTED TABLES
 - MATERIALIZED VIEWS.
 - DBT AND CLICKHOUSE DRIVER.
 - PYTHON INTEGRATION.
 - TIME TO LIVE (TTL)
- NEXT STEPS IN OUR DATA PLATFORM JOURNEY
 - INTELLIGENT USER LIMITS.



Have we tamed the data dragon, or has our adventure just begun?

CLICKHOUSE® IS A REGISTERED TRADEMARK OF CLICKHOUSE, INC. [HTTPS://CLICKHOUSE.COM](https://clickhouse.com)

DAGSTER IS AN OPEN-SOURCE PROJECT MAINTAINED BY DAGSTER LABS. COPYRIGHT © 2023 ELEMENTL, INC.
D.B.A. DAGSTER LABS. ALL RIGHTS RESERVED.

DBT™, DBT CORE, AND DBT CLOUD
ARE TRADEMARKS OF DBT LABS, INC.

YOUR TURN: ANY QUESTIONS?