# Collaboration on data with Lakehouses and Iceberg

Christian Thiel, Vakamo

+ Demo

Vakamo

# whoami
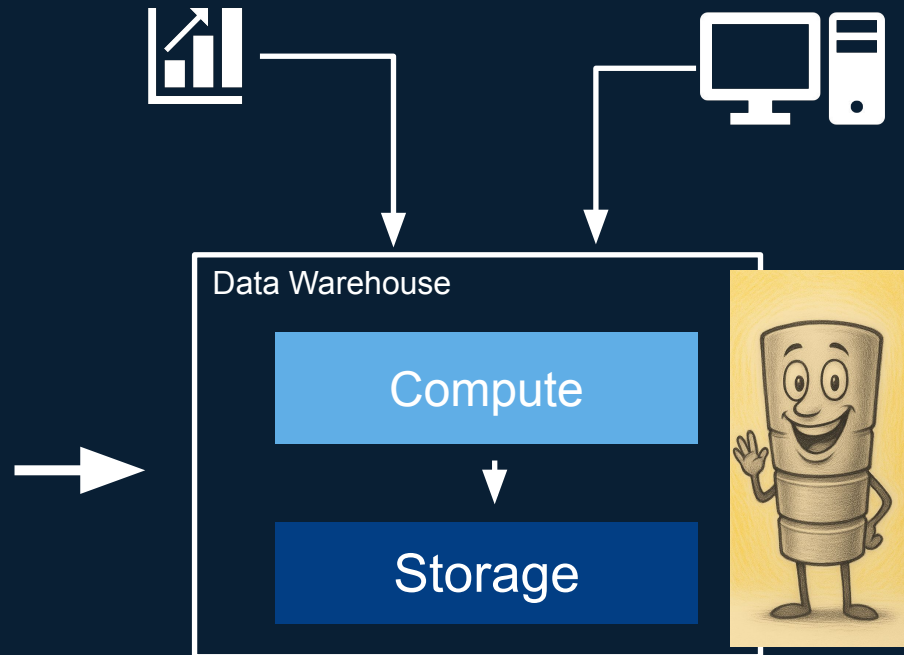


> Christian, OSS Enthusiast

📍 Hamburg

> Co-Founder & Software Engineer
   @ Vakamo

> Developing **Lakekeeper**
A Rust-Native, Apache Licensed
Iceberg Rest Catalog

> Apache Iceberg Contributor

Vakamo

# A not too lengthy and reasonably accurate history of Data & Analytics systems
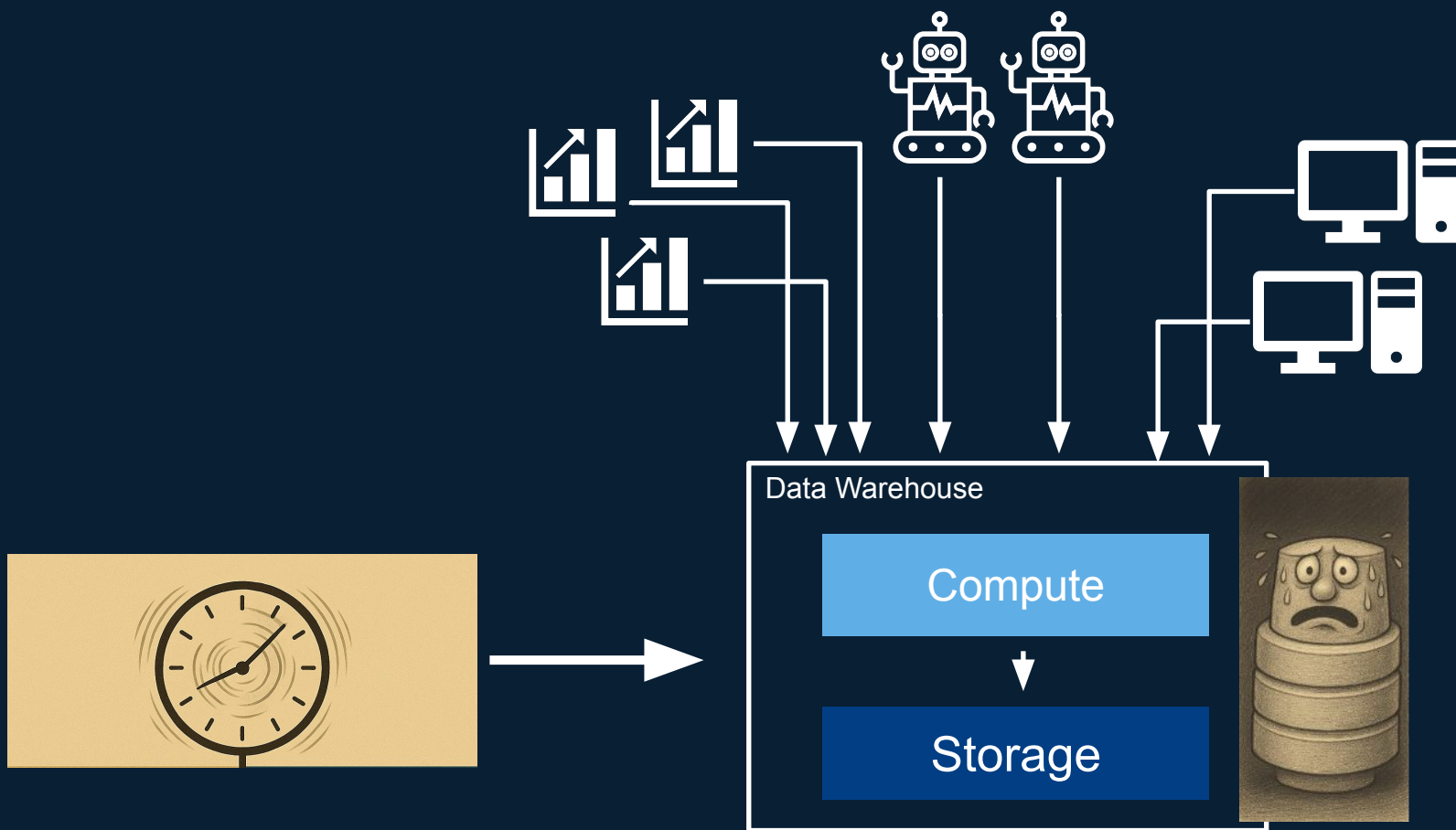
Vakamo

# A long time ago …



Big Bang

## Data Warehouse

Compute

Storage

🔒 Transaction Safe (ACID)

🔑 Schema Enforcement & Evolution

🔄 Governance & Access Control

Vakamo

# … more recently

Data Warehouse

Compute

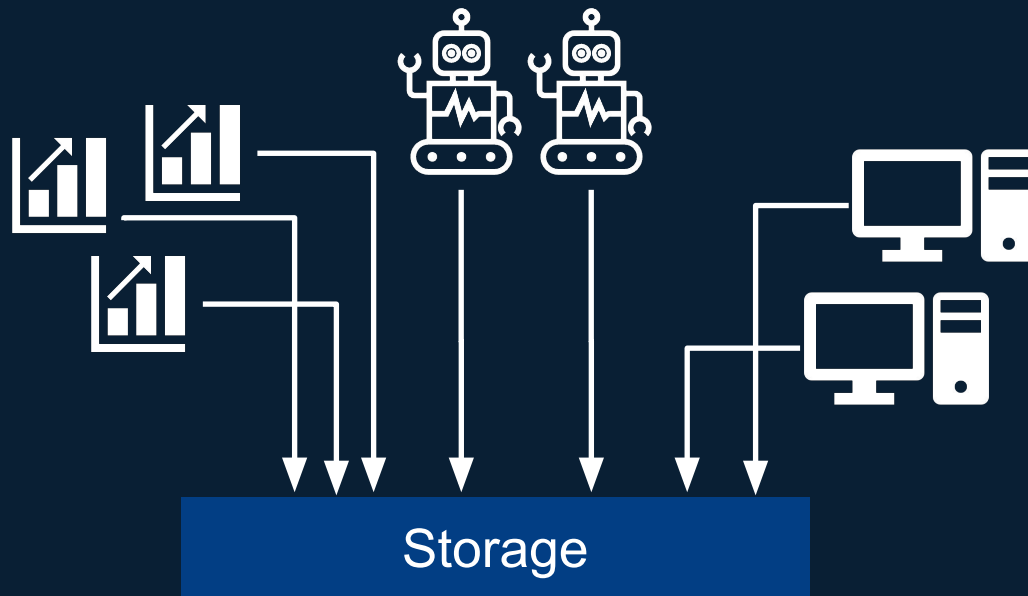Storage

- Central Compute is bottleneck
- Doesn't scale well
- Expensive
- Very Slow Data Extraction
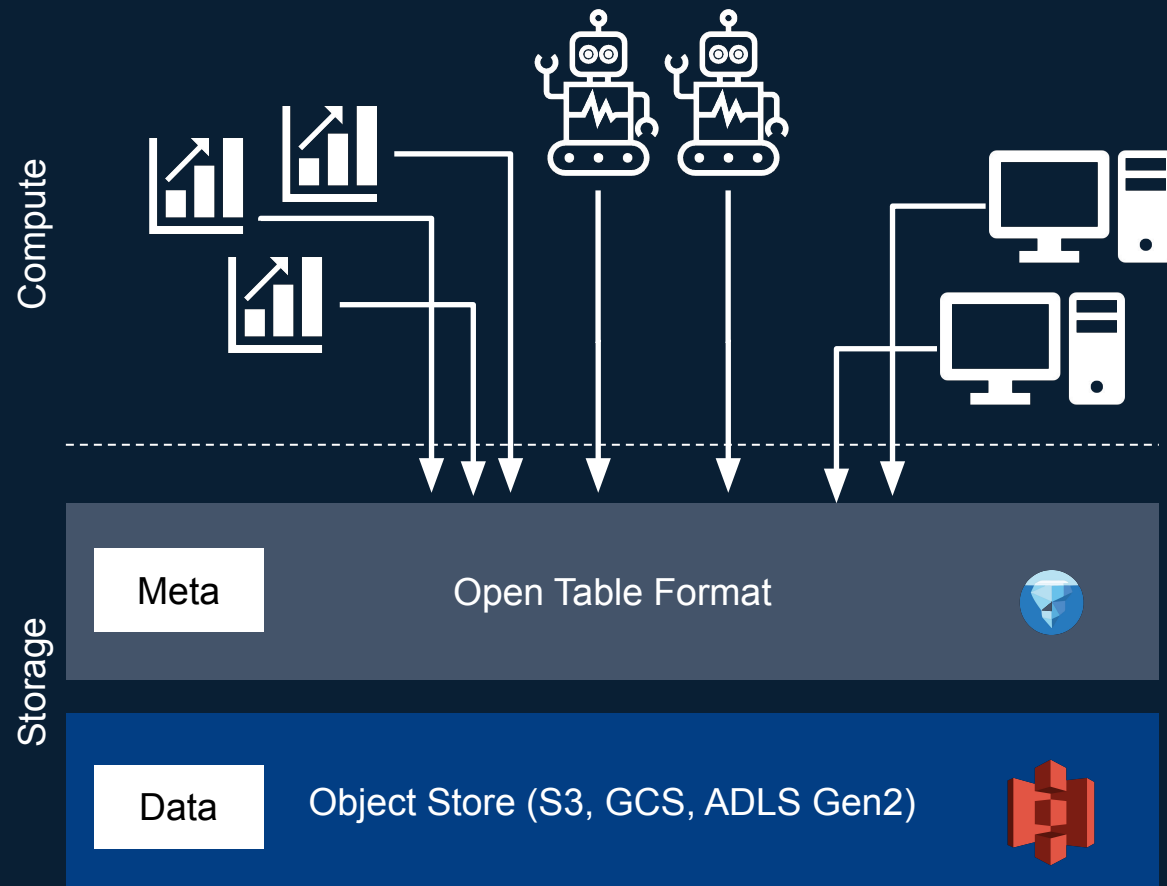- Not open to all use cases and external Compute

Vakamo

# Data Lakes



Storage

Scalable & Performant

Heterogeneous Compute Engines

Cost Reduction

Vakamo

# Lakehouse / The Marketing Slide

Storage

**Meta** — Open Table Format

**Data** — Object Store (S3, GCS, ADLS Gen2)

## Best of both Worlds

### Warehouse

- Transaction Safe (ACID)
- Schema Evolution
- Governance

### Data Lake

- Scalable & Performant
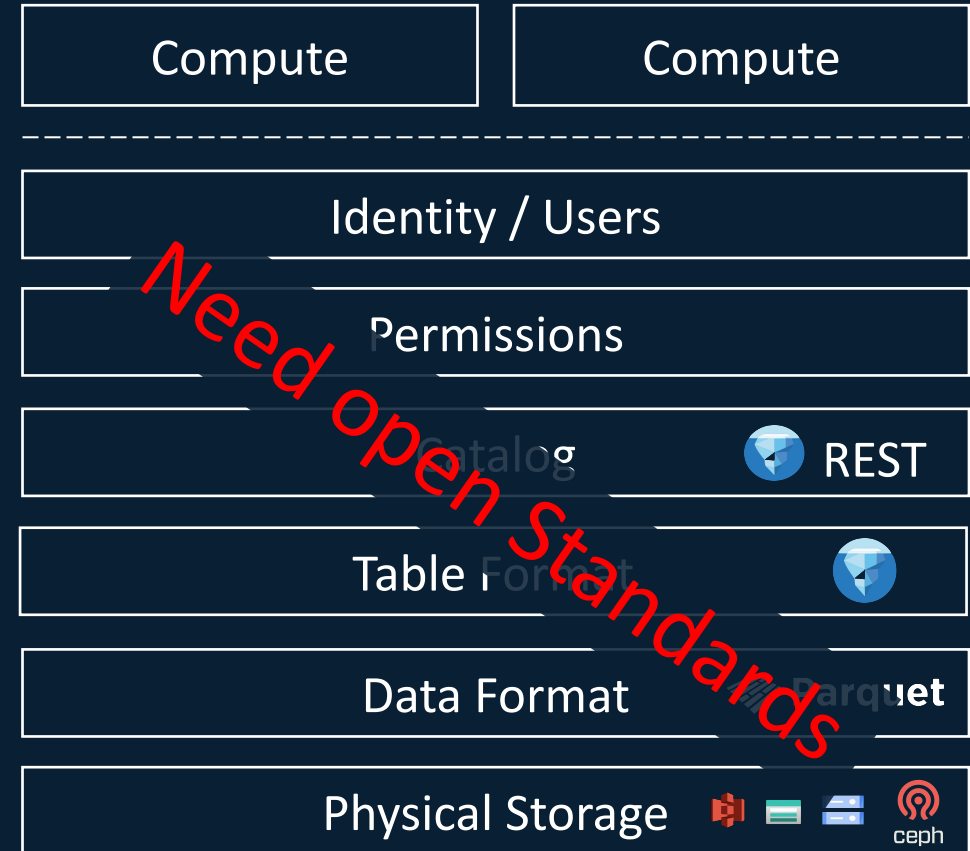- Open, Heterogeneous Compute
- Cost Reduction

Vakamo

# Collaborating on Data is Hard

| Compute | Compute |
|---|---|

Shared
- Identity / Users
- Permissions
- Catalog — 🌐 REST
- Table Format — 🔷
- Data Format — Parquet
- Physical Storage — ceph

Need open Standards

Vakamo

# What is Apache Iceberg?

**"The open table format for analytic datasets"**

- Developed at Netflix 2017

- Apache Project since 2018

- High-Performant format for huge analytic tables

- Brings reliability and simplicity of SQL tables to analytical data

Contributors from all Major D&A Vendors

Vakamo

# REST Catalog

OpenAPI Specification as part of the Iceberg Project
-> Multiple Implementations, i.e. **Lakekeeper**

- Catalog is Source-of-Truth for Table Metadata

- Transactions over multiple Tables (ACID)

- Table Discovery

- Change-based Commits

- Performant load of Table Metadata

- Hand out temporary Credentials for Data Access

- Language agnostic REST Interface, Replaces Legacy Catalogs:
  Hadoop, Hive, Glue, JDBC, Dynamo, ECS, Nessie, …

---

**Supported by all Major Query Engines:**

PyIceberg    Apache Spark    kafka

trino    DuckDB    APACHE DATAFUSION

aws

---

Vakamo

# Collaborating on Data is Hard

| Compute | Compute |
|---------|---------|

**Shared**

| Identity / Users | |
| Permissions | |
| Catalog | 🌐 REST |
| Table Format | |
| Data Format | 📊 **Parquet** |
| Physical Storage | |

Vakamo

# The Goal

> **I need access controls for my Iceberg Tables!**
>
> - Naive Business Users

Shared Query Engine

Catalog

Permissions

Revenue Table

Vakamo

# Where should permissions be Stored?

**Compute**     **Compute**     ❌ ∘ ∘ ○ ☁ Not shared

**Share**d

**Identity / Users**

**Permissions**     ✔ … A dedicated Permission System sounds open

**Catalog** 🌐 REST     (✔)

**Table Format** 🌐

**Data Format** ▱ **Parquet**

**Physical Storage** 🔶 🟩 🟦 🔴     ❌ ∘ ∘ ○ ☁ Ask me in the break

Vakamo

# Permissions in Dedicated Systems

Open Policy Agent    OpenFGA™

Permissions

Compute    Compute    Catalog

IdP
(OAuth2)

Physical Storage

🔒 Permission
s
Identity

o **Single place for shared permissions**

o **Semantic understanding with permission hierarchy**

o **Consistent across engines**

o **Share permissions across Catalogs and Table Formats and Applications**

⚠ System that create PATs or Client Credentials or PWs are IdPs and a security risk!

# How Permissions can be Enforced?

Vakamo

# How can Security be Enforced?

**1** Catalog enforces Permissions
(Iceberg REST Vended Credentials)

Vakamo

# Catalog enforces Security: Vended Credentials

① GET /table/xyz
Authorization XXXX

② Access Storage

Current TableMetadata
Temporary Credentials 🔒

**Lakekeeper**

**Object Store**

Vended Credentials & Remote Signing work on File-Level
-> Table-Level Security
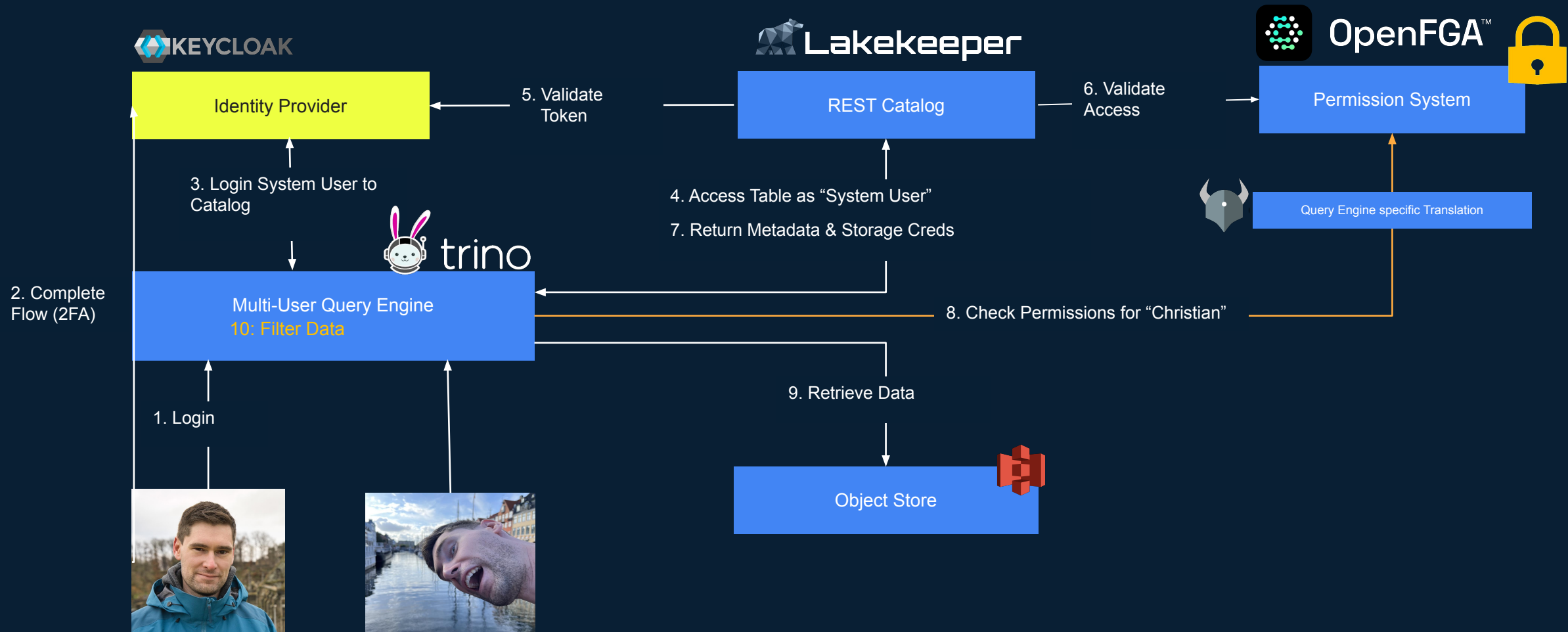
Vakamo

# How can Security be Enforced?

**1**    Catalog enforces Permissions
(Iceberg REST Vended Credentials)

**2**    Compute enforces Permissions

Vakamo

# Query engine enforced Security

Vakamo

# How can security be enforced?

## Catalog enforced Security

+ Table-Level Security
+ Doesn't rely on filtering in client: Works for untrusted Clients
+ Standardized: Works with all Iceberg REST Clients

− No row or column level security

## Query engine enforced Security

+ Row / Column Level Security
+ Performance optimizations in Query Engine

○ Requires secure, isolated User Sessions

− Requires trusted compute

**We need Both!**

Vakamo

# How can Security be Enforced?

**1**    Catalog enforces Permissions
(Iceberg REST Vended Credentials, Table Level)

**2**    Compute enforces Permissions
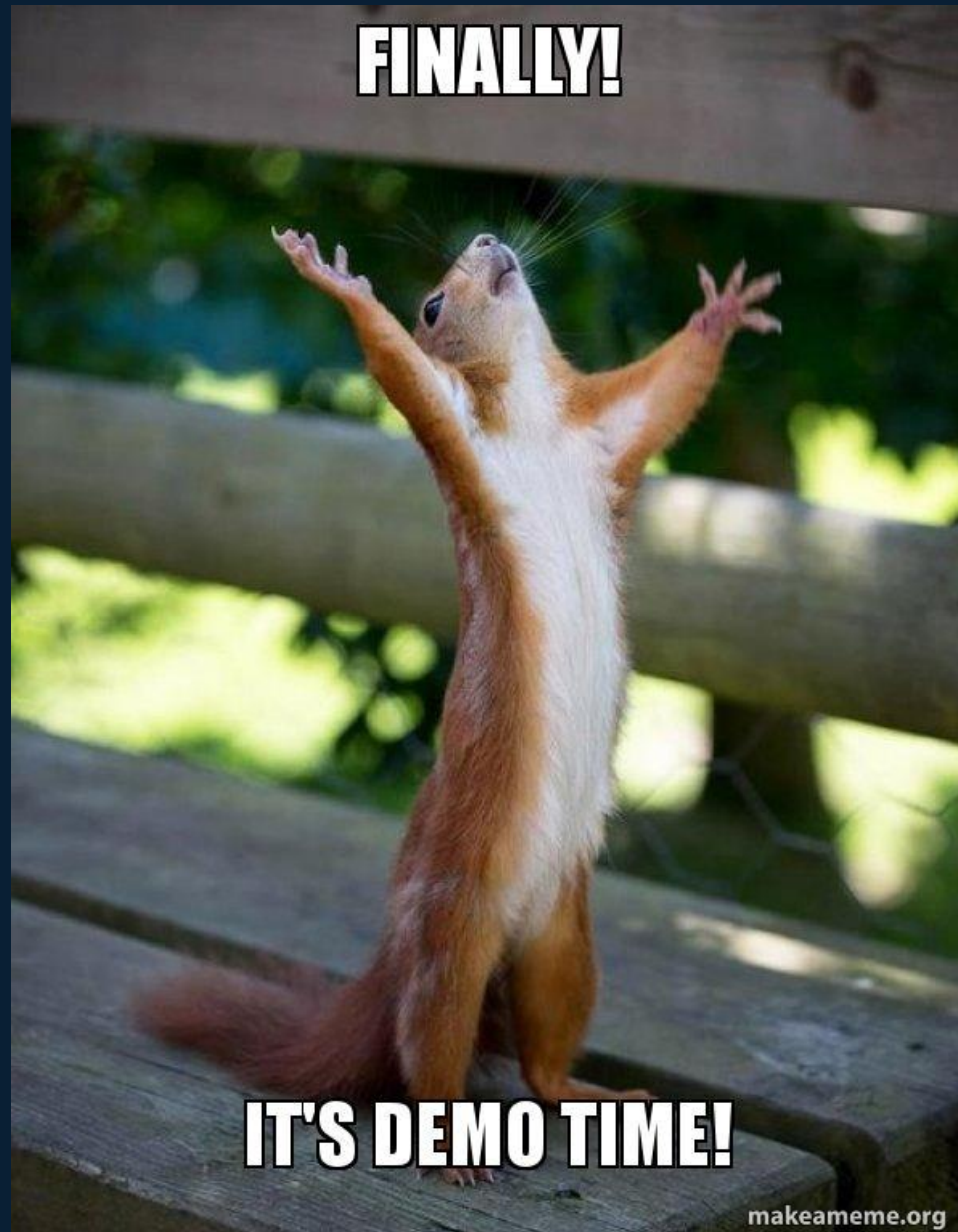(Requires trusted Compute)

**3**    Storage enforces Permissions
(Requires Iceberg-Aware Storage)

**4**    Encryption of Files or Columns
(Application level)

Vakamo

Vakamo

# Lakekeeper

o Iceberg REST Catalog implementation

o Apache Licensed

o Written in Rust

o Kubernetes Integration

o Secure: OAuth2 and Permissions

o Production Ready & Easy to Use

o Extendible (Rust Traits!)

o Multi-IdP Support (Kubernetes + OIDC)

Vakamo

# Open?!



Your Datawarehouse will never need another migration after moving to an open Iceberg Lakehouse.

CHANGE MY MIND

Permissions stored in open, independent System

Authentication using YOUR Identity Provider

Vakamo

# Open Lakehouse with Open Governance is possible today!
## Try it now!

```
git clone https://github.com/lakekeeper/lakekeeper.git
cd /examples/access-control-advanced
docker compose up

-> localhost:8888
```

1 Identity Provider
1 Source of truth for Permissions
Different Computes
Some of them shared
100% OSS

## Lakekeeper

Lakekeeper is a Rust-Native, Apache Licensed Iceberg Rest Catalog

Vakamo