# Trino & OPA @ Stackable

Sönke Liebau
& Sebastian Bernauer
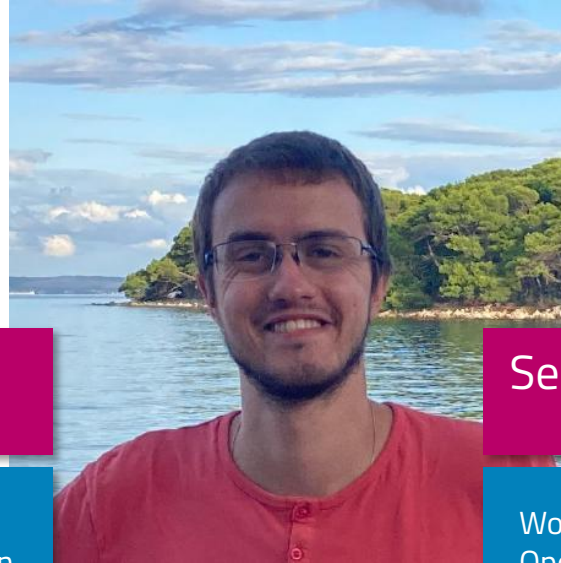
Stackable

# Agenda

1. What is Stackable?
2. Open Policy Agent (OPA) authorization plugin
   - History
   - Recent development
   - Compatibility layer to Trino's File-based access control
   - Quick demo on row filtering and column masking
3. Auto-scale Trino clusters using trino-lb
   - Differences between trino-gateway and trino-lb
4. Lessons learned running Trino on Kubernetes
   - What our trino-operator is doing
   - Potential next steps

Stackable

# About us



## Sönke Liebau
### CPO Stackable

Co-Founder of Stackable, working with Big Data Open Source Software since 2012, speaker, contributor, husband & father…

## Sebastian Bernauer
### Software developer

Working with Big Data Open Source Software since 2019
Big Open Source and Trino fan :)

# Stackable in a Nutshell

## Founded

2020

OpenCore

b.telligent

IONOS

## Stackable Data Platform

> Open Source

> Infrastructure as Code

> Cloud-native (Kubernetes)

> On-Premises, Cloud, Hybrid

## Our Team: 20 People

International
in Germany & Europe

## Our Services

> Product Support

> Big Data Consulting

> Trainings

## Network - Collaborations

OSB Open Source Business ALLIANCE

KI BUNDESVERBAND

gaia-x

bitkom

eco

Stackable

# Popular Data Apps. Kubernetes-native. Easy to deploy and operate.

# OPA plugin - History



https://www.youtube.com/watch?v=fbqqapQbAv0

# OPA plugin - History

History

1. 2021/10: Stackable creates the stackabletech/trino-opa-authorize repo
2. 2023/02: After Bloomberg reached out license was changed to ASL2
3. 2023/05: Bloomberg created Trino PR upstream with much improved version
4. 2024/01: OPA plugin was merged into Trino and released with version 438 🚀

Recent development

5. 2024/07: Bloomberg improve the performance of column masks by batching requests send to OPA, released in 453 🚀
   https://github.com/trinodb/trino/pull/21997

🦓 Stackable

# Compatibility layer to Trino's File-based access control

- Trino already offers a great and flexible access control

```
access-control.name=file
security.config-file=etc/rules.json
```

- We want users to be able to migrate to OPA as easy as possible

  → Compatibility layer written in rego, which takes the same JSON definition as input and emulates the Trino behaviour

  → Can server as a starting point

https://github.com/stackabletech/trino-operator/tree/main/tests/templates/kuttl/opa-authorization/trino_rules

Stackable

# Compatibility layer

```json
{
 "tables": [
   {
     "user": "admin",
     "privileges": ["SELECT", "INSERT", "DELETE", "UPDATE", "OWNERSHIP"]
   },
   {
     "schema": "hr",
     "table": "employee",
     "privileges": ["SELECT"],
     "filter": "user = current_user"
   }
 ]
}
```

https://trino.io/docs/current/security/file-system-access-control.html

🗄 Stackable

# Compatibility layer

```json
{
 "schemas": [
    {
      "user": "admin",
      "schema": ".*",
      "owner": true
    },
    {
      "group": "finance|human_resources",
      "schema": "employees",
      "owner": true
    }
  ]
}
```

https://trino.io/docs/current/security/file-system-access-control.html

# Userinfo Fetcher
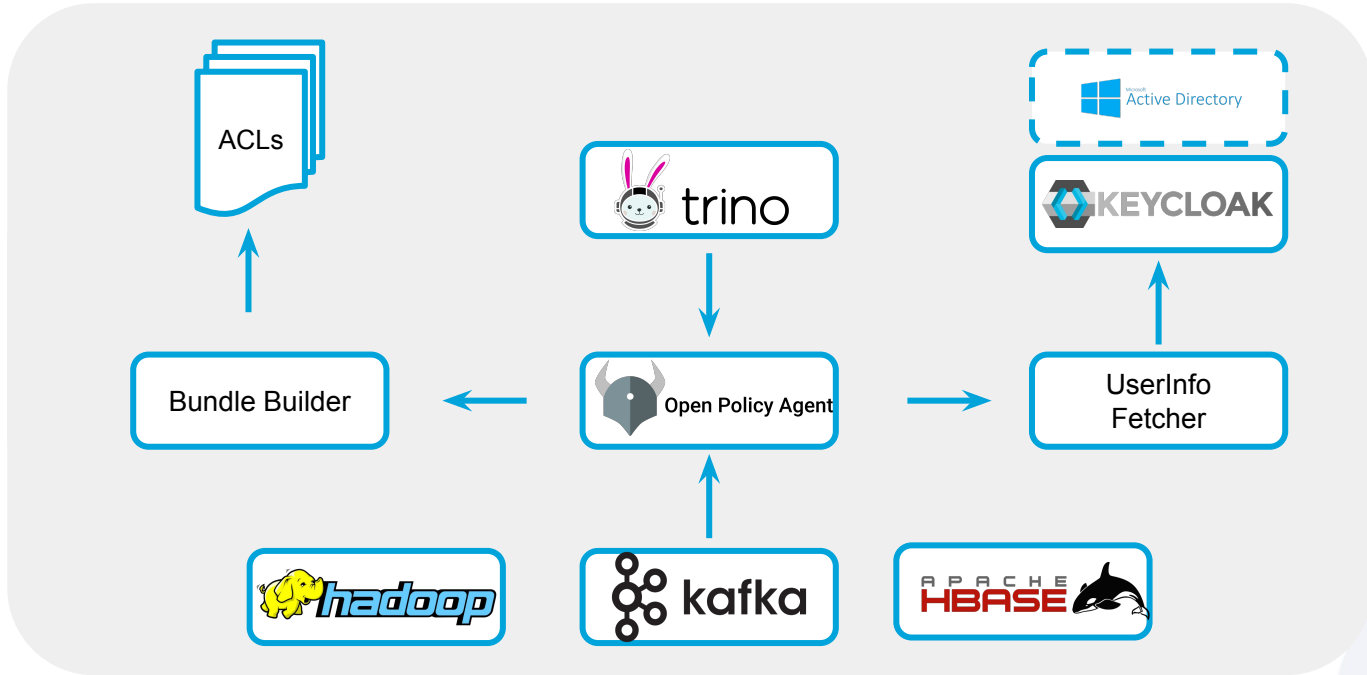
## User info fetcher

**WARNING**

This feature is experimental, and subject to change.

The *User info fetcher* allows for additional information to be obtained from the configured backend (for example, Keycloak). You can then write Rego rules for OpenPolicyAgent which make an HTTP request to the User info fetcher and make use of the additional information returned for the username or user id.

Stackable

# Userinfo Fetcher

```
{
  "id": "af07f12c-a2db-40a7-93e0-874537bdf3f5",
  "username": "alice",
  "groups": [
    "/admin"
  ],
  "customAttributes": {}
}
```

Stackable

# The Big Picture

# Quick demo on row & column level security



https://www.youtube.com/watch?v=ATlq_l3WNiA

# Column Level Security



**Marketing**

Mark Ketting

**Customer Service**

Customer Analytics

Justin Martin

**Compliance**

Compliance Analytics

Sophia Clarke

Column-level security:

Read from customer, but the following rules should apply:

1. Prohibit reading first and last name, birth_month and birth_day
2. Instead of seeing the customer_id they only see the sha256 hash of it
3. Email-addresses are masked to abcXXXX@domain.com

**customer_analytics schema**

customers table
- customer_id
- first_name
- last name
- birth_year
- birth_month
- birth_day
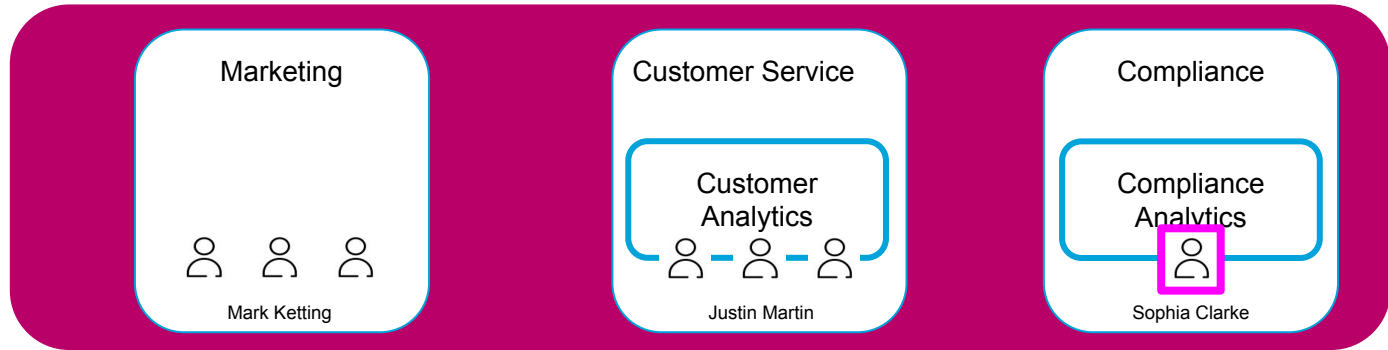- login
- email_address…

**compliance_analytics schema**

customer_enriched view
- customer_id
- birth_year
- email_address…

# How does it look in code?

```json
{
  "group": "/Compliance and Regulation/Analytics",
  "catalog": "lakehouse",
  "schema": "customer_analytics",
  "table": "customer",
  "privileges": ["SELECT"],
  "columns" : [
    {"name": "c_first_name", "allow": false},
    {"name": "c_last_name", "allow": false},
    {"name": "c_birth_day", "allow": false},
    {"name": "c_birth_month", "allow": false},
    {
      "name": "c_customer_id",
      "mask": "'sha256:' || to_hex(sha256(to_utf8(c_customer_id)))",
    },
    {
      "name": "c_email_address",
      "mask": "regexp_replace(c_email_address, '([^@]{1,3})([^@]+)@', '$1---@')",
    },
  ]
},
```
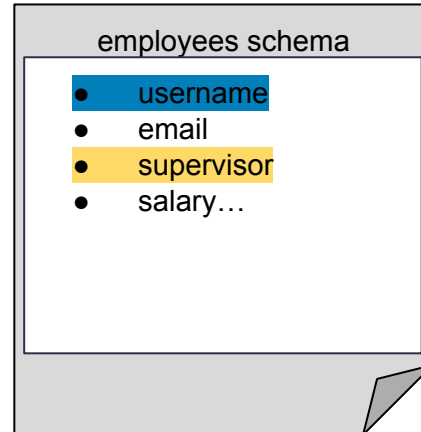
Stackable

# Row Level Security



Row-level security:

Read from employees, but the following rules should apply:

1. Everyone can only see themselves
2. Supervisor additionally see their reports

employees schema
- username
- email
- supervisor
- salary…

Stackable

# How does it look in code?

```json
{
  "catalog": "lakehouse",
  "schema": "employees",
  "table": "employees",
  "privileges": ["SELECT"],
  "filter": "username = current_user or supervisor = current_user",
},
```

Stackable

# Demo time

Stackable

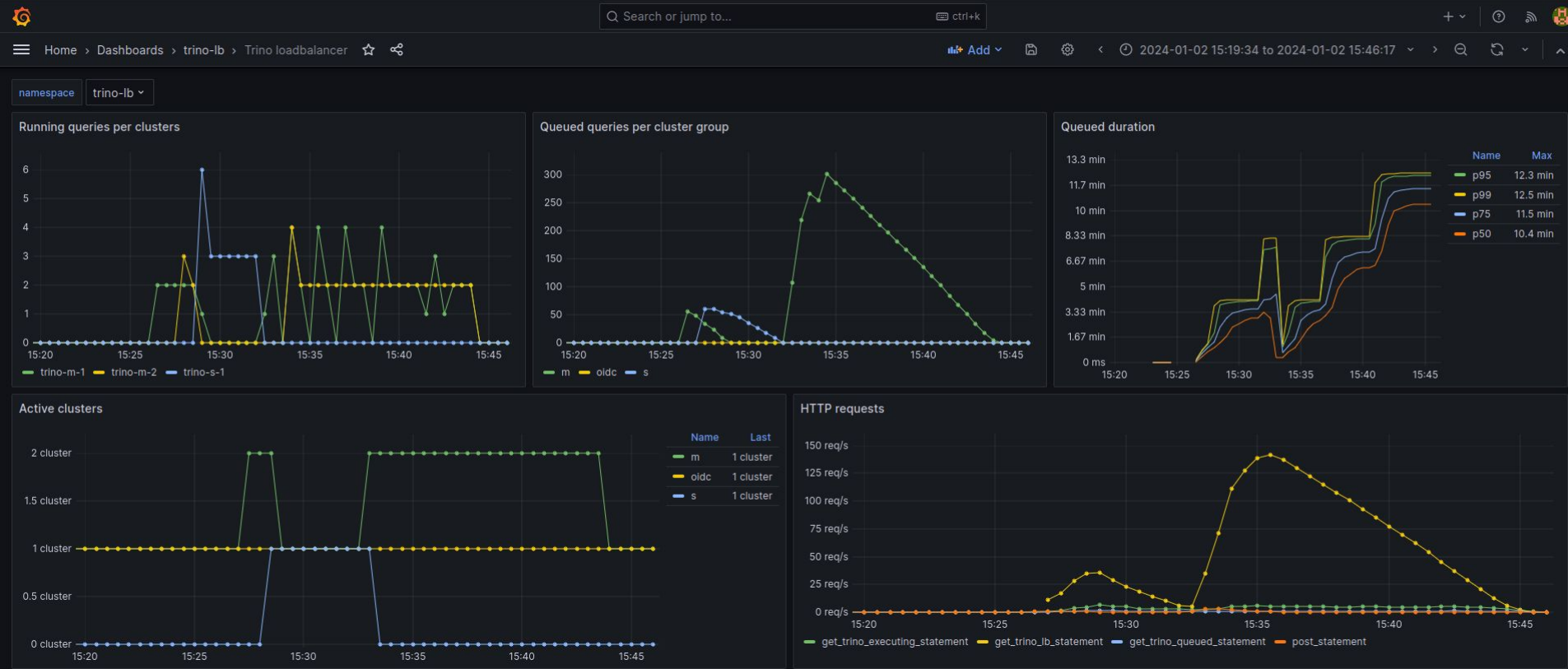# Quick demo on row & column level security

# Auto-scale Trino clusters using trino-lb

- trino-lb development started around 2023/10, just before trino-gateway was first released

- The primary goals are
  - Queuing of queries in case all available Trino clusters are already full
  - Auto-scaling of entire Trino clusters (load and time based)
  - Performance (trino-lb is horizontally scalable)
  - High availability (trino-lb is stateless)
  - Very flexible routing strategies (e.g. Python script)
  - Modularity to supported different persistence, routing and scaling implementations

https://github.com/stackabletech/trino-lb



Stackable

# Auto-scale Trino clusters using trino-lb

# Auto-scale Trino clusters using trino-lb

- OpenTelemetry tracing
  - Trace propagation to Trino

# Lessons learned running Trino on Kubernetes

- First off: We don't run any production Trino on Kubernetes

- But our customers do :)

- We have written an operator to manage Trino on Kubernetes:
  https://github.com/stackabletech/trino-operator
- Documentation: https://docs.stackable.tech/home/stable/trino/

Stackable

# Try to avoid coordinators restarts

- A coordinator restart kills all running queries
- Mitigation:
  - We have a flag: Don't touch this cluster at all costs!
- Potential future work:
  - Trino HA [#391]? :)
  - Maintenance windows
  - Graceful shutdown of coordinator
    i. Remove coordinator from trino-lb/trino-gateway
    ii. Wait till all queries finished
    iii. Restart
    iv. Add coordinator to trino-lb/trino-gateway again
    v. Requires Kubernetes nodes to wait long enough while draining!

Stackable

# Graceful shutdown of workers

- A worker restart kills all running queries (without fault tolerant execution)
- Mitigation:
  - Graceful shutdown of workers
    i. Requires Kubernetes nodes to wait long enough while draining!
    ii. We also set query.max-execution-time
  - Fault tolerant execution :)

Stackable

# Pod placement

- Avoid too many workers being down at the same time
- Mitigation:
  - By default we spread all workers across as many nodes as possible
    i. Can be customized by customer based on their topology
    ii. Avoid impact of node/rack/room/datacenter failures
    iii. Assumption: Big worker nodes to reduce internal Trino traffic
  - PodDisruptionBudets: Only X nodes can be down simultaneously

Stackable