# Using cutting-edge open-source technologies to build a World-class Industrial Data Lake
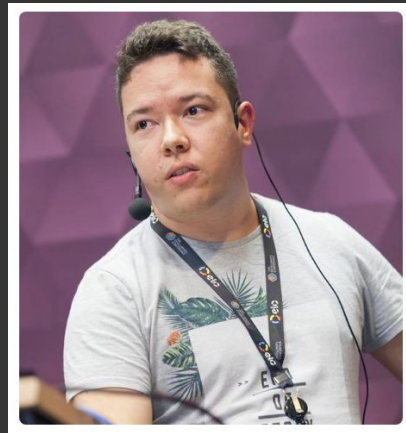
# Hello!
## My name is Allan Sene

Co-Founder & CTO @**Data Sprints**

Co-Founder & Podcaster & Instructor @**Data Hackers**

+10 years in **Data & Software**, +4 years as **Data Engineer**

# Agenda

🤔 The Challenge

🤔 What is a Data Lake?

🤔 Awesome cutting-edge data tools

🤔 Putting everything together do build a IDL

🤔 Results and Next Steps

🤔 Q&A

# THE CHALLENGE

# The Challenge

- **Multinational Steel Industry, with plants worldwide**

- **They need to give to global managers the capability of track the production line**

- **Data columns have BLOBs, Arrays and complex data types**

- **Migrating from on-prem to cloud**

- **Data sets with 70 Gb (compressed) and 8 million lines, more than 1600 columns, very complex queries (+200 joins), 15 minutes delay**

- **Maximum Query Response time of 30s**

- **Very limited cloud budget - USD 15.000,00/year**
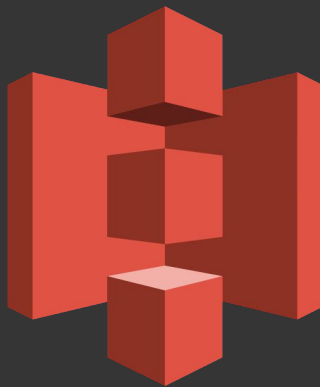
# Stakeholder's name, can you guess it?

# What is a Data Lake?

"Is a Data Repository that holds a huge quantity of data on raw state, structured or not. The schema is only defined when is necessary for consumption"

(Anne Buff - Best Practices Leader at SAS)

# Not really...

# A Data Lake must have

- **Security & Audit**
- **Catalog & Access**
- **Data Pipelines & Orchestration**

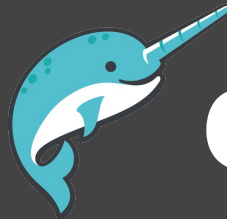Basically, a **Platform** that ensure safe and efficient data consumption

www.datasprints.com

# Awesome cutting-edge data tools

# Dremio's UI

# dbt

- Dbt = "Data Build Tool"
- Built by Fisthtown Analytics
- Data Pipeline Orquestration
- "Airflow for Data Analysts"
- Open-source
- SaaS Version

**" dbt**

- ETL over a MPP
- Code Versioning
- Data Lineage & Docs
- Data Validation
- Seed Data

# dbt's UI

Putting **everything together** do build a **DL**

# Bringing all together

- **Pipelines & Orchestration => dbt**

- **Catalog & Access => Dremio**

- Processing => Spark

- Storage => Amazon S3

# Bringing all together

Repartitioned Data

Coalesced Data

Micro-batch Data

Repartition & Schema Evolution

Optimization

dremio

JDBC

Data Port

Data Explorer

# Results and Next Steps

# Results, until now

- **Platform built on 6 months by a team of 4 engineers full-time**

# Results, until now



# of Different Queries

Query Time (in seconds)

# Results, until now

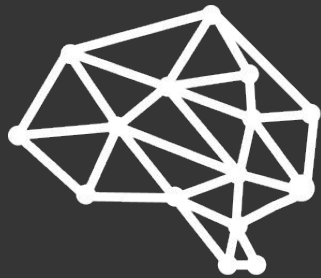| concurrent_queries | rowCount sum | meet_the_requirements sum | difference_in_seconds min | max | median | mean |
|---|---|---|---|---|---|---|
| 45 | 3262636 | 34 | 0.393 | 44.031 | 7.3050 | 12.090267 |
| 70 | 1478566 | 61 | 0.258 | 32.023 | 2.4400 | 7.007314 |
| 95 | 5961920 | 70 | 0.234 | 76.075 | 5.5820 | 14.498263 |
| 120 | 3658127 | 113 | 0.196 | 31.908 | 1.3755 | 4.608908 |
| 145 | 9041435 | 104 | 0.209 | 57.661 | 5.5160 | 12.167186 |
| 170 | 7531672 | 88 | 0.189 | 108.799 | 18.5740 | 30.950112 |
| 195 | 6629865 | 181 | 0.178 | 26.454 | 1.4020 | 4.658010 |
| 220 | 7134027 | 166 | 0.192 | 55.385 | 5.8280 | 11.232900 |
| 245 | 8909619 | 214 | 0.201 | 61.665 | 3.4950 | 8.003478 |
| 270 | 8387467 | 231 | 0.202 | 55.081 | 4.2275 | 8.923437 |
| 295 | 14860819 | 226 | 0.216 | 89.322 | 5.0350 | 12.182417 |
| 320 | 18289194 | 206 | 0.281 | 148.988 | 9.6635 | 22.119000 |
| 345 | 17224813 | 187 | 0.415 | 186.446 | 15.6190 | 33.190449 |
| 370 | 21406306 | 268 | 0.193 | 119.080 | 5.8045 | 14.771841 |

# Next steps

- **Integrating with very common Analytics tools (Power BI, Tableau...), deployed worldwide**
- **Data Science Models consuming data through the Data Bus**
- **Data Quality Monitoring**

**Data Hackers Podcast** Spotify

linkedin.com/company/data-hackers/

www.datahackers.com.br

Questions?

www.datasprints.com

# Thanks!



allan@datasprints.com

https://www.linkedin.com/in/allansene/

www.datasprints.com