

Data Observability Monte Carlo-val – Mit, miért, hogyan?

Budapest Data Forum
2022

Nagy András István
Senior Data Solution Architect



<epam>

Nagy András István
Senior Data Solution Architect

Data Observability

Mi az a "Data observability"?

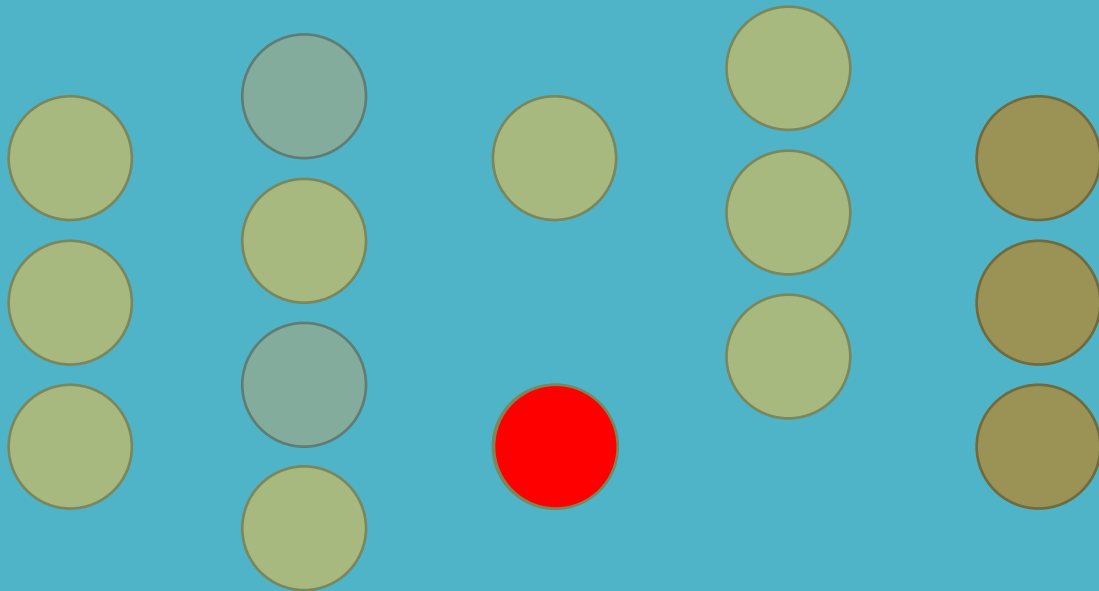
Miért kell nekünk ilyen?

Hogyan tudjuk elérni?

Ügyfelünknel hogyan vezettük be:

- Eszköz kiválasztása
- Bevezetés
- Beillesztés a folyamatainkba

Downtime



Data Observability

Egy rendszer képessége arra, hogy megfigyelhető legyen a működése

- 1) Egyáltalán mit kell megfigyelni?
- 2) Mi a teendő, ha gond van?
- 3) Észrevenni előre nem definiált eseményekből is az esetleges hibát

Eszközök:

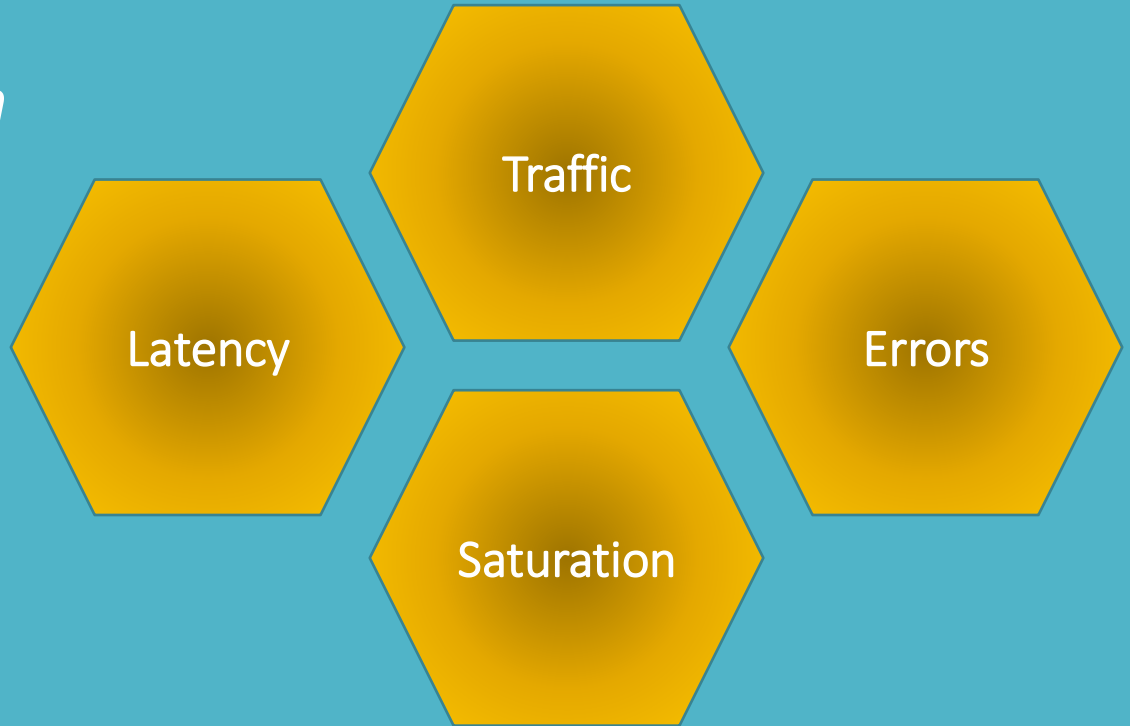
- Metrikák
- Logok
- Trace-ek
- Health check

Data Observability

4 Golden Signals –
*Site Reliability Engineering:
How Google Runs Production
Systems. Sebastopol, CA:
O'Reilly, 2016.*

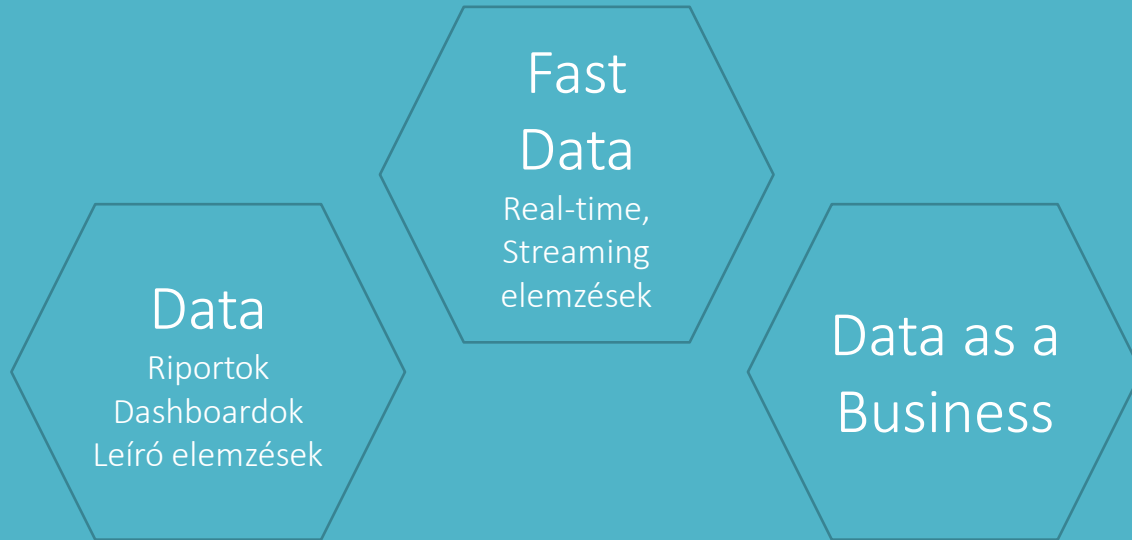
Eszközök:

- White-box monitoring
- Black-box monitoring



Data Observability

Data Downtime

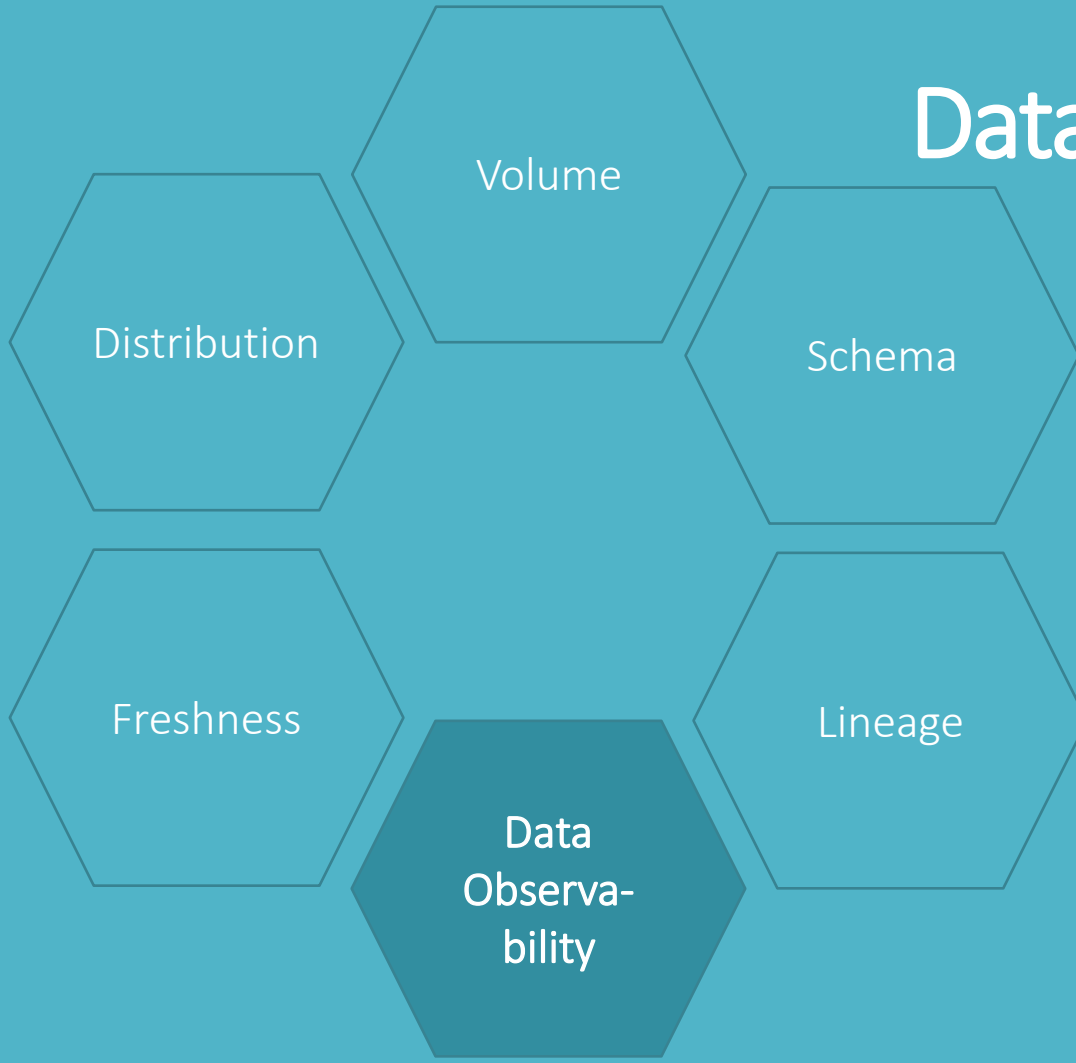


Data Observability

Data Downtime

- Megfigyelhető (Observable)
- Kontextus
- Reagálás

Data Observability



Sources



amazon
DynamoDB



amazon
DynamoDB

SaaS

SaaS

Manufacturing

Integration



kafka



API GW

Data Lake



amazon
REDSHIFT



S3



Glue

Mit válasszunk?



Apps

Forrásrendszerek:

- Egyedi alkalmazások (DynamoDB, D. Streams)
- SaaS (ServiceNow, Coupa, Jira,...)
- Egyéb rendszerek (API, CDC)

DWH / Data Lake:

- S3 – Raw,
Transformed: ML
training
- Redshift: staging,
ODS, Reporting layer

Analitikus alkalmazások:

- Tableau
- Egyedi data app-ok

Mit válasszunk?

Fő szempontok:

- Data Observability képességek
- A bevezetés ráfordításigénye
- Illeszkedés a környezetünkbe

Mit válasszunk?

- Data Lineage feltérképezése, lineage alapján impact analízis
- Sémaváltozás detektálása, processzek
- Data Security and privacy
- Szabályok egyedi definiálása (DQ expectations)
- Szabályok automatikus felállítására, automatizált anomália detekció
- Adatok frissességének monitorozása
- Csatolás a pipeline-ok és az eszköz között - hozzá kell-e nyúlni a pipeline-jainkhoz?
- Bevezetés erőforrásigénye
- Üzemeltetés erőforrásigénye
- Teljesítmény, skálázhatóság

Mit válasszunk?

Interfészek:

- Interaktív, vizuális felhasználói felület
- Incidensek esetében értesítések kiküldése, triaging
- API az integrációhoz
- Access Control
- Integráció: Redshift, Tableau, S3, Glue Data Catalog, Kafka
- Illeszkedés AWS-re épülő data infrastruktúránkhoz

Mit válasszunk?

- Célzott, custom megoldás Apache Deequ / AWS DQAF alapon
- BigEye (DO) + Alation (DC)
- Collibra DQ (DO/DQ) + Collibra (DC)
- Monte Carlo Data (DO)

Data Lake DEV



Data Lake QA



Data Lake PROD



Query Logs

System tables

S3 Events

Glue Data Catalog



AWS Transit Gateway

Data Collector DEV



Data Collector QA



Data Collector PROD



CloudFormation

Monte Carlo SaaS



Sync alert rules

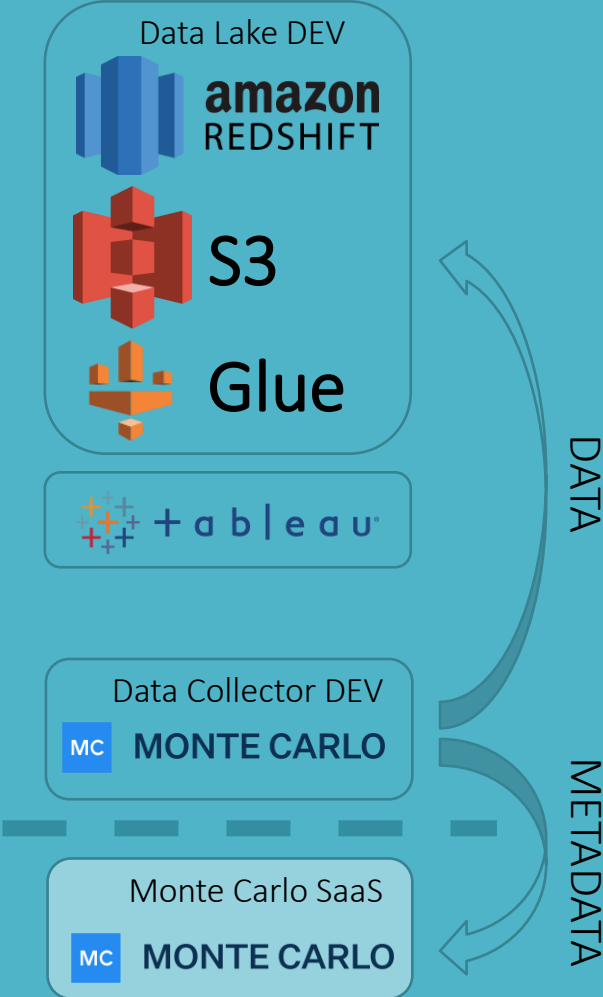


Table Owners



SLACK

Privacy



Következő lépések

- Lineage kiterjesztése
- Circuit Breakers: SQL Monitor integráció Airflow DAG-ból

Tapasztalatok

Közös slack csatorna, gyors reakció

Dedikált Customer Success Manager

Része a folyamatainknak, incidenseket mi vesszük észre, nem az ügyfeleink

Házon belüli elterjedtség nő