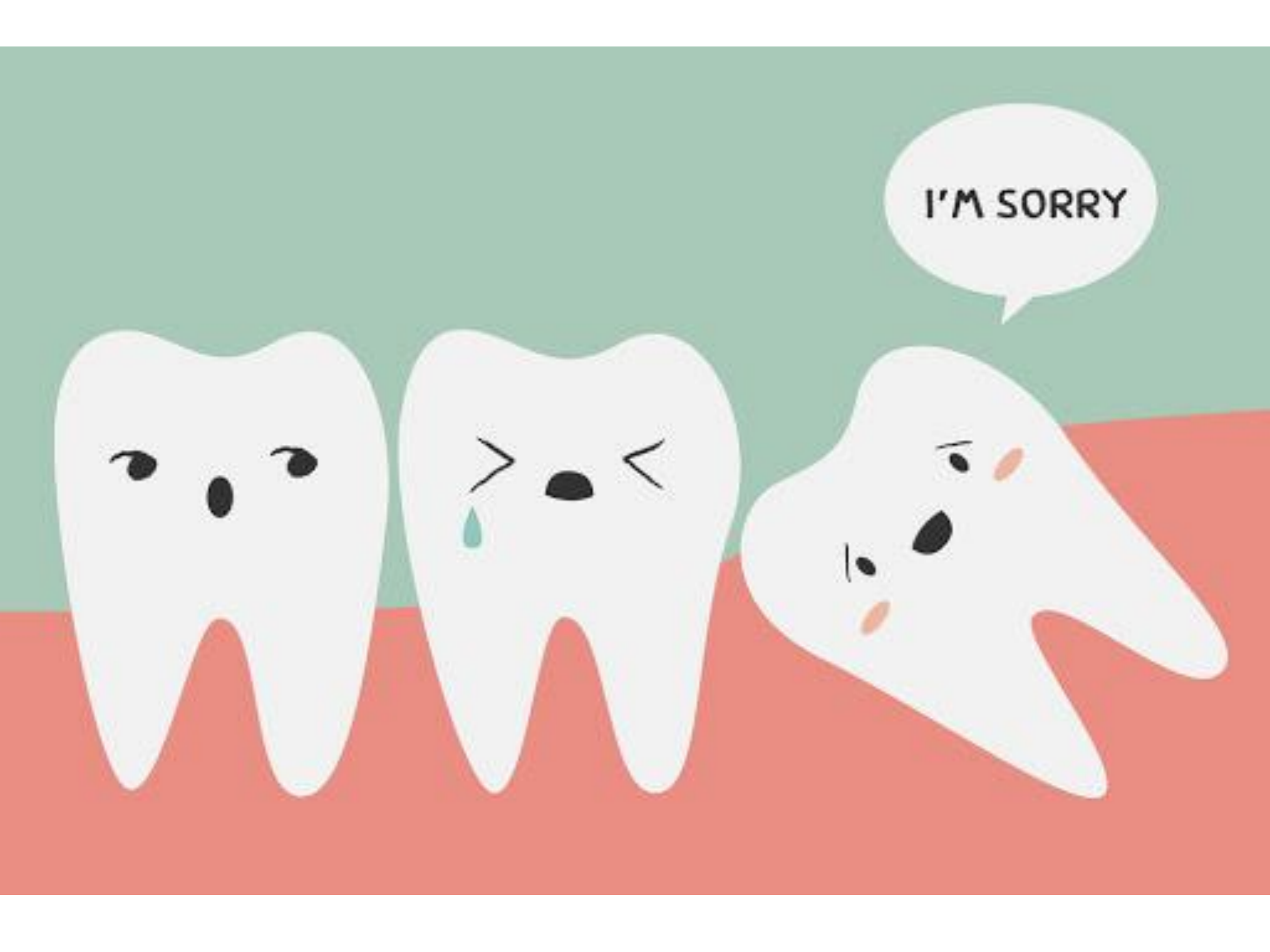


Delta Lake

ACID tranzakciók Big Data rendszerekben



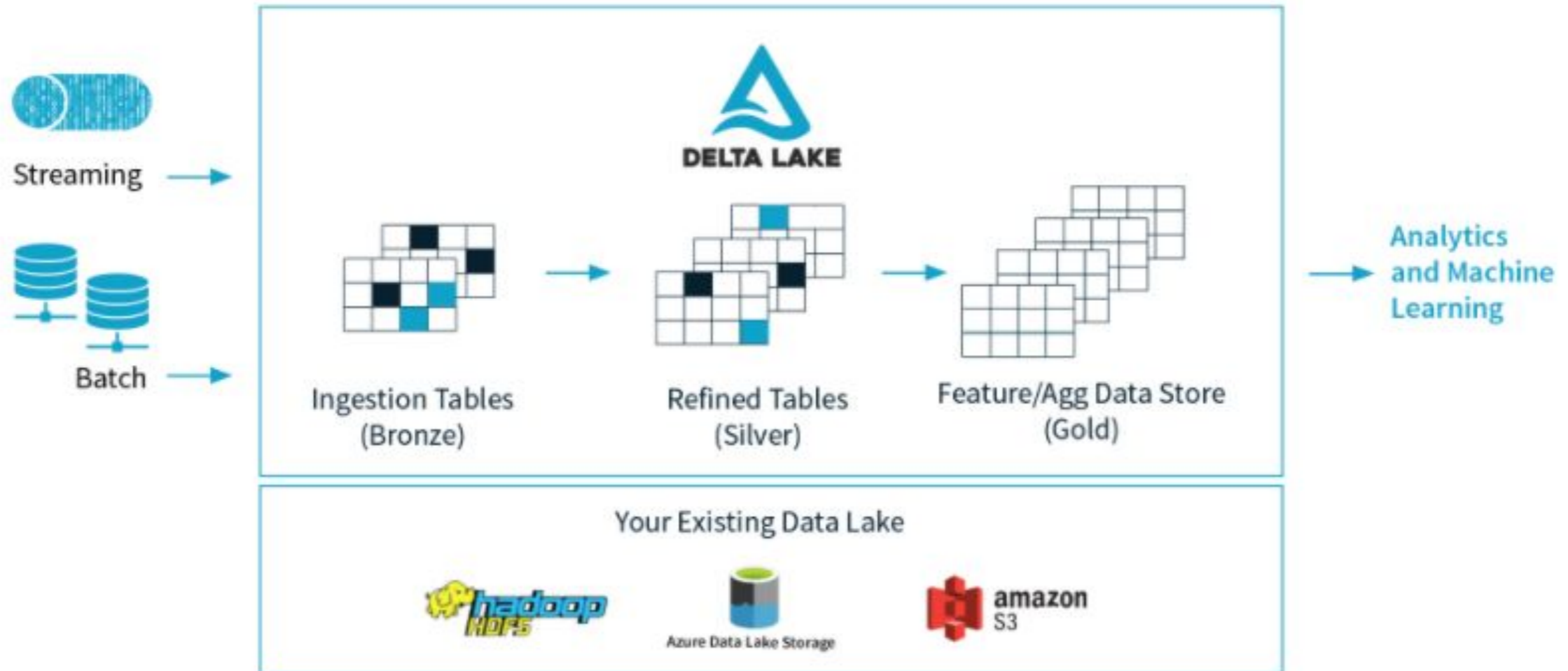
Tóth Zoltán
CTO Datapao



I'M SORRY



Delta Lake is an open-source storage layer that brings ACID transactions to Apache Spark™ and big data workloads.



Atomicity	<ul style="list-style-type: none">• File alapú transaction log• Cloud storage/HDFS file írás atomikus
Consistency	<ul style="list-style-type: none">• Az adatbázis tartha a definiált invariánsokat• <i>C in ACID was “tossed in to make the acronym work”*</i>
Isolation	<ul style="list-style-type: none">• Konkurens írások nem interferálnak• Megoldás: File alapú, szekvenciális transaction log (atomic) olvasása/írása
Durability	<ul style="list-style-type: none">• Nem veszik el az adat• A Cloud storage/HDFS garantálja



DELTA LAKE

Schema Evolution

- NOT NULL → NULLABLE
- Új oszlop hozzáadása
- Típusok “upcast”-olása (pl Int -> Long)



DELTA LAKE

Adat-változtatás

- DELETE
- UPDATE
- UPSERT/MERGE



DELTA LAKE

Time travel

A delta tábla korábbi snapshotjain is könnyen lehet dolgozni.



DEMO

Mikor használd?

- Spark 3-at használasz, és ...:
 - Konkurens írásra van szükséged
 - UPDATE/DELETE/MERGE műveleteket végzel

Mikor ne használd?

- Nem Sparkot használasz (= nincs Delta támogatás)
- Nincs sok adatod és egy fájl/SQLite épp elég
- Egy “writer” ír egyszerre appenddel/overwrite-tal mindig
(ilyenkor használj parquet-t, ami sokkal gyorsabb)



DELTA LAKE

Köszönöm!



Contact

Tóth Zoltán

zoltan@datapao.com

+36 30 291 3599

LinkedIn: zoltanctoth