

2 – Delta Lake vs Iceberg

Delta Lake to rozwiązanie stworzone przez Databricks i mocno wspierane przez Microsoft. Jest ono naturalnie zintegrowane z platformą Azure i środowiskiem Spark. Zapewnia pełne wsparcie transakcji ACID, łatwe wersjonowanie danych oraz automatyczna optymalizacja danych – bez potrzeby dodatkowej konfiguracji. Schema evolution działa bezproblemowo, co oznacza że zmiany w strukturze danych nie wymagają dużej ingerencji w pipeline.

Apache Iceberg natomiast to projekt rozwijany w ramach Apache Software Foundation – powstał w Netflixie i stawia na otwartość i elastyczność. Iceberg lepiej sprawdza się w środowiskach z wieloma silnikami zapytań np. Trino, Presto, Hive czy Flink. Ma bardziej zaawansowane mechanizmy wersjonowania i zarządzania metadanymi, ale ich wdrożenie wymaga więcej pracy i znajomości konfiguracji. Nie jest tak dobrze zintegrowany z Azure jak Delta.

Jeśli chodzi o optymalizację – Delta oferuje mechanizmy typu Z-Ordering i Auto Optimize, co znacznie upraszcza zarządzanie danymi. Iceberg natomiast bazuje na filtracji po metadanych, co daje większą kontrolę, ale wymaga zrozumienia więcej niż tylko samego kodu SQL.

Delta sprawdzi się najlepiej gdy:

- korzystamy z Databricks lub ekosystemu Azure
- zależy nam na szybkim wdrożeniu z automatyczną optymalizacją
- chcemy mieć wsparcie dla machine learningu, streamingu
- potrzebujemy prostych upsertów i transformacji typu batch/stream

To dobre rozwiązanie do budowy lakehouse'ów w środowisku silnie opartym o Spark.

Iceberg to lepszy wybór gdy:

- pracujemy w środowisku z różnymi silnikami (Presto, Hive, Flink itd.)
- działamy on-prem lub hybrid
- zależy nam na pełnej kontroli wersjonowania i pracy na snapshotach danych
- mamy częste zmiany schematu i potrzebujesz zaawansowanej ewolucji danych

Iceberg świetnie radzi sobie w dużych organizacjach z rozproszoną architekturą danych.

3 – Krytyka medalion

Architektura medalionowa czyli podział danych na warstwy bronze silver i gold ma swoje zalety, ale też sporo ograniczeń. W wielu przypadkach jej stosowanie jest przesadne lub wręcz przeszkadza w efektywnej pracy z danymi.

Po pierwsze jest po prostu zbyt złożona – trzy warstwy oznaczają więcej kodu, więcej testów i więcej punktów awarii. Dane są kopiowane między warstwami co podnosi koszty i utrudnia śledzenie ich drogi. Dla czystych źródeł danych warstwa bronze jest często niepotrzebna, a mimo to ją się tworzy na siłę.

Kolejna sprawa to opóźnienia – dane muszą przejść przez wszystkie warstwy, zanim dotrą do użytkownika, co jest problematyczne np. w analizach bliskich rzeczywistego czasu. Architektura nie wspiera też dobrze real-time streamingu. Dodatkowo logika transformacji potrafi się dublować w warstwach silver i gold, przez co trudniej to utrzymać.

Dla analityków biznesowych system bywa nieprzejrzyisty – trzeba wiedzieć z której warstwy korzystać, nie ma jednej wersji prawdy. Jeśli chodzi o uprawnienia, to każda warstwa może wymagać osobnych konfiguracji, a bez Unity Catalog robi się to szybko problematyczne. Nie ma też dobrego mechanizmu do śledzenia historii zmian – dokumentacja często nie nadąża za transformacjami.

Schematy danych trzeba aktualizować oddzielnie dla każdej warstwy, co bywa upierdliwe. Każda dodatkowa warstwa zużywa więcej zasobów: CPU, IO, storage. Debugowanie problemów w takich pipeline'ach jest jak szukanie igły w stogu siana – trudne, wolne i kosztowne.

Medalion wymusza też tworzenie dużej liczby tabel, monitoring, pipeline'ów – dla każdej warstwy osobno. Nowoczesne narzędzia jak DuckDB albo dbt pozwalają pominąć część warstw, zachowując przejrzystość. W efekcie architektura medalionowa potrafi doprowadzić do izolacji danych – dane zamknięte w jednej warstwie mogą być niewidoczne w innej.

Na koniec trzeba wspomnieć, że ta architektura nie jest łatwa do wdrożenia bez doświadczonego zespołu data engineeringowego. I że jest mocno promowana przez Databricks – ciężko ją przenieść do innych platform 1:1 bez utraty funkcjonalności.