

## **1. Walidacja Danych na Etapie Ekstrakcji**

Walidacja danych na etapie ekstrakcji pomaga w wykrywaniu i naprawianiu błędów we wczesnej fazie procesu ETL. Należy upewnić się, że dane pasują do oczekiwanego schematu (np. sprawdzenie typów danych, wymaganych pól i zakresów wartości) oraz mają prawidłowy format (np. poprawność dat).

## **2. Automatyczne Przekształcanie i Czyszczenie Danych**

Automatyczne przekształcanie i czyszczenie danych pozwala na eliminację błędów wynikających z niepoprawnych lub niekompletnych danych. Użycie narzędzi do skalowalnego przetwarzania danych (takich jak Apache Spark) umożliwia m.in. usuwanie wierszy z brakującymi wartościami oraz zastępowanie niepoprawnych wartości.

## **3. Mechanizmy Retry i Idempotency**

Zastosowanie mechanizmów retry (ponawianie operacji) i idempotency (bezpieczne powtarzanie operacji) pozwala na bezpieczne ponowne wykonanie operacji w przypadku nieoczekiwanych błędów, takich jak awarie sieci czy problemy z połączeniem do baz danych.

## **4. Monitorowanie i Alertowanie**

Implementacja systemów monitorowania i alertowania pozwala na szybkie wykrycie i reakcję na błędy oraz anomalie w pipeline. Dzięki temu można na bieżąco śledzić wydajność systemu, identyfikować problemy oraz podejmować natychmiastowe działania naprawcze.

## **5. Dokumentacja i Szkolenie**

Zapewnienie dokładnej dokumentacji oraz szkolenie zespołu na temat najlepszych praktyk i typowych pułapek pozwala na uniknięcie błędów użytkownika. Dokumentacja powinna obejmować wszystkie kroki pipeline, oczekiwane formaty danych, oraz procedury postępowania w przypadku błędów.