

Integracja źródeł NoSQL

Arkadiusz Kasprzak

GFT

Krzysztof Pupiec



Nadia Sikorska



Marcin Siudziński



Różnorodne źródła danych



Narzędzia Business Intelligence



fppt.com

Wybrane technologie

ORACLE



mongoDB



 pentaho®
A Hitachi Data Systems Company

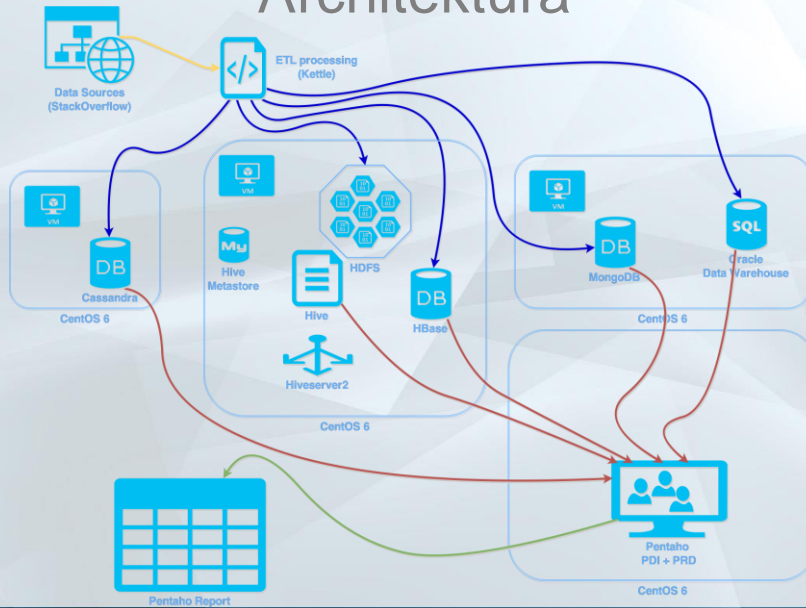


APACHE
HBASE

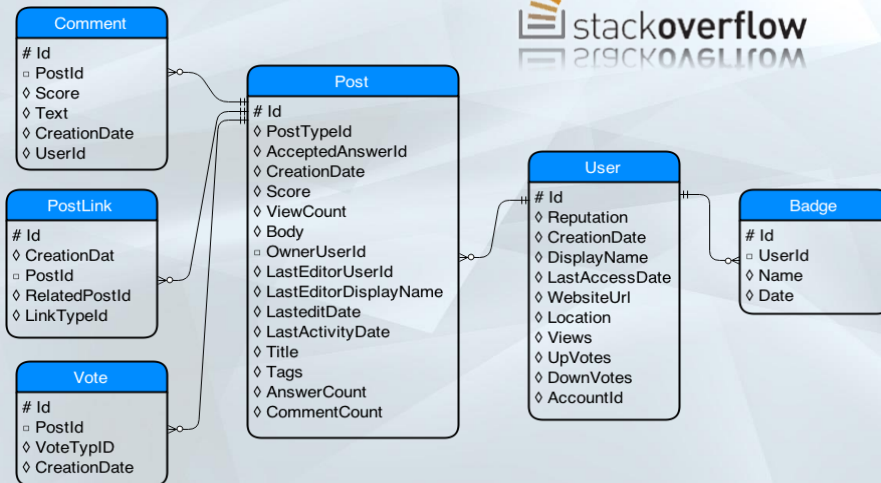
cassandra

fppt.com

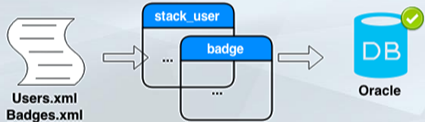
Architektura



Import danych



Import danych



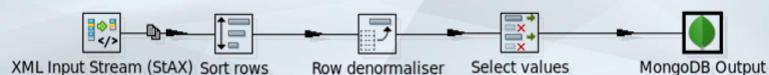
Proces ETL



Proces ETL



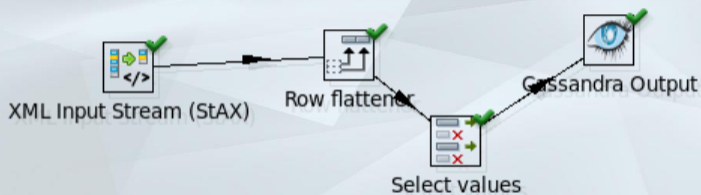
✗ brak wsparcia języka SQL



Proces ETL



✗ tabele muszą być kompatybilne wstecz
- WITH COMPACT STORAGE
✗ brak konektora PRD



Proces ETL

APACHE
HBASE

- ✗ brak konektora PRD
- ✗ błędnie działający komponent w PDI
- ✗ dostęp do bazy tylko przez Hive
- ✗ wymagana dodatkowa konfiguracja sterownika „Generic Database”

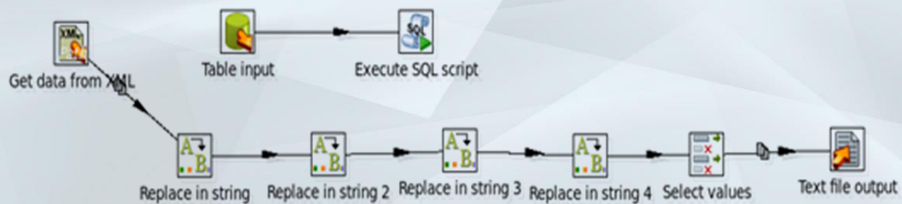


fppt.com

Proces ETL

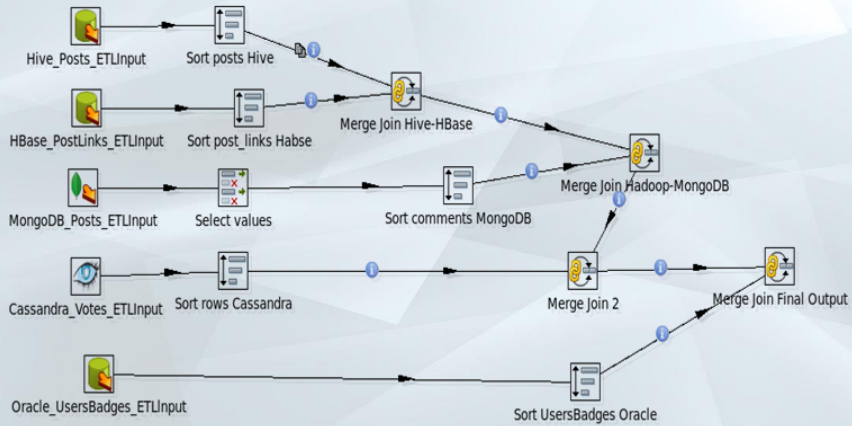


- ✗ niezgodność Pentaho Shims i Hadoop
- ✗ problemy z obsługą distributed cache
- ✗ brak wsparcia dla nowych wersji
- ✗ błędnie działające komponenty w PDI
- ✗ procesy ETL wymagają obsługi HDFS z poziomu OS

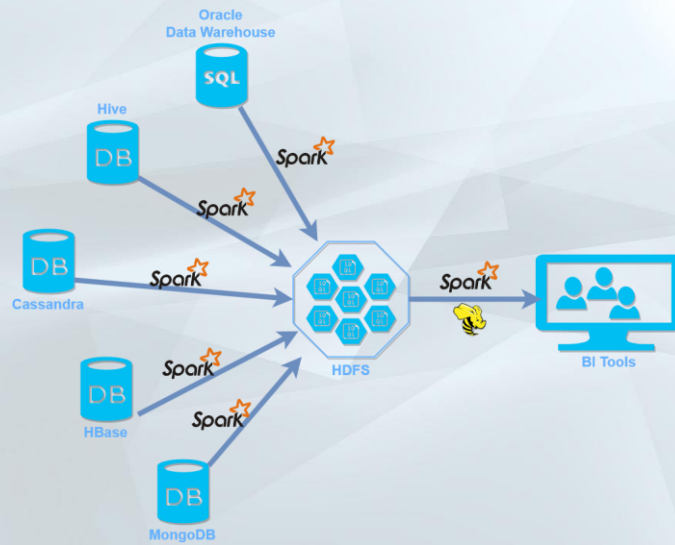


fppt.com

Konsolidacja źródeł



Rozwiązanie alternatywne



Cechy Apache Spark

- API w Java, Scala, Python, R
- Szybkość
- In-memory
- Równoległe operacje
- RDD (Resilient Distributed Datasets)
- Kompatybilność
- Odporność na awarie
- Dostępność bibliotek
- Strumieniowanie
- Algorytmy uczenia maszynowego

Apache Spark

Open Source
Ecosystem



Spark + Zeppelin

Apache Zeppelin (incubating) Community Docs Download GitHub Apache

Apache Zeppelin (incubating)

A web-based notebook that enables interactive data analytics.
You can make beautiful data-driven, interactive and collaborative documents with SQL, Scala and more.

Watch the video Get Zeppelin

Multi-purpose Notebook

The Notebook is the place for all your needs

- Data Ingestion
- Data Discovery
- Data Analytics
- Data Visualization & Collaboration

Bank

The screenshot shows a Zeppelin notebook titled 'Bank'. It contains two visualizations: a line chart on the left showing data points over time, and a pie chart on the right showing the distribution of data. The line chart has a legend with 'Bank', 'Other', and 'Expense'. The pie chart has a legend with 'Bank', 'Other', and 'Expense'.

fppt.com

Dziękujemy za uwagę.

Arkadiusz Kasprzak

GFT

Krzysztof Pupiec

NASPERS

Nadia Sikorska

naviexpert

Marcin Siudziński

apollologic
it
KONTRAKT

fppt.com