



Consulting

Beratung

Results,
no Excuses.

Lösungen

Grown from
Experience.

Ventum Consulting

SQL auf Hadoop

Oliver Gehlert

1 Ventum Consulting – Daten und Fakten

Results, no excuses



Fachwissen

- Strategisches Businessverständnis
- Fachliche und methodische Prozessexpertise
- Breites IT KnowHow

Branchenkenntnis



Umsetzungsstärke



Results, no excuses.

Unser Versprechen

1

Einführung

2

Hive, Impala, Tajo and Co

3

Performance

4

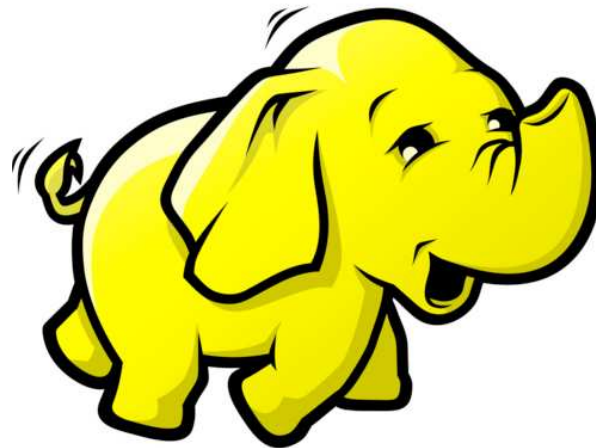
Ausblick

SQL auf Hadoop

Muss man für Hadoop Java programmieren?



Java Programmierung war die Basis für den Zugriff auf Hadoop



SQL auf Hadoop

Es gibt zahlreiche Ansätze für SQL auf Hadoop

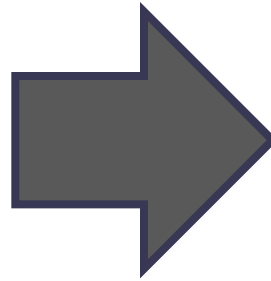


Ausgewählte SQL Frameworks



SQL auf Hadoop

SQL auf Hadoop verlangt strukturierte Daten



Images office.com

SQL Frameworks

Wie haben wir getestet



Testsetup

- **Sandboxserver mit**
 - Cloudera 5 Beta 2
 - 16 GB RAM
 - Testdatengröße 20 GB
- **Sandboxserver mit**
 - Hortonworks HDP 2.0
 - 16 GB RAM
 - Testdatengröße 20 GB
- **Cluster**
 - Amazon Elastic Map Reduce
 - 5 Nodes m2.4xl
 - Testdatengröße 100 GB

Testdaten

- **Hive-Testbench basierend auf TPC-DS**
 - Download unter <https://github.com/cartersha nklin/hive-testbench>
 - Starschema und ausgewählte Abfragen

1

Einführung

2

Hive, Impala, Tajo and Co

3

Performance

4

Ausblick

Überblick

- Ein Datawarehouse System um strukturierte Daten im HDFS zu speichern
- Unterstützt zahlreiche Dateiformate
- Strukturiert die Daten analog zu Datenbankkonzepten in
 - Tabellen
 - Partitionen
 - Spalten
 - Zeilen
- Ermöglicht einen einfachen SQL basierten Zugriff
- Hive Metadata Store wird von vielen Frameworks verwendet

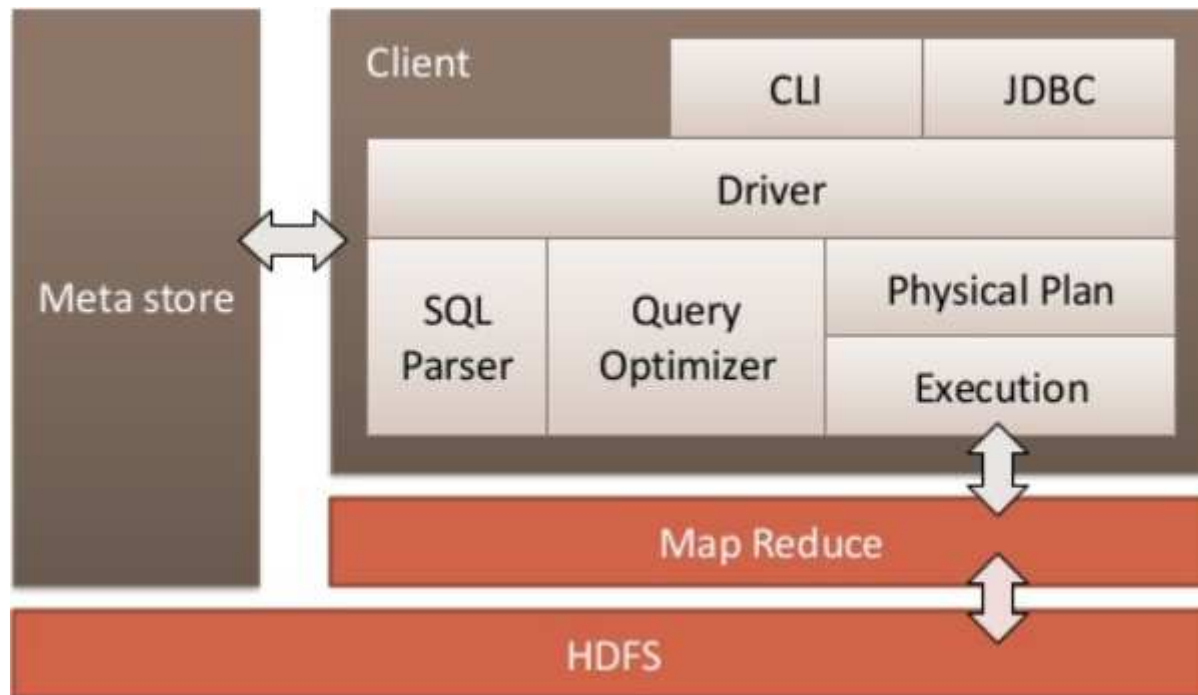


SQL Frameworks

Hive Architektur



Architektur



Images apache.org

Stärken

- Stabil
- Intensive Weiterentwicklung
- Memoryanforderungen
- Batchprocessing
- Umfangreiche SQL Unterstützung
- Hive Metastore
- UDF

Schwächen

- Interaktive Abfragen
- Keine Einzelsatzverarbeitung
- Map-Reduce

Unterstützung

- Analytische Funktionen
- Subselects
- Group by
- Rollup
- UNION ALL

Nicht unterstützt:

- UNION
- MINUS
- INTERSECT

SQL Frameworks

Stinger bündelt Performanceoptimierungen



Überblick

- Die Stinger Initiative vereint mehrere Performanceoptimierungen für Hive
- Optimiertes Fileformat orcFile
 - Komprimierung um Faktor 4 im Vergleich zu Text
 - Predicate Pushdown
- Analytische Funktionen
- YARN



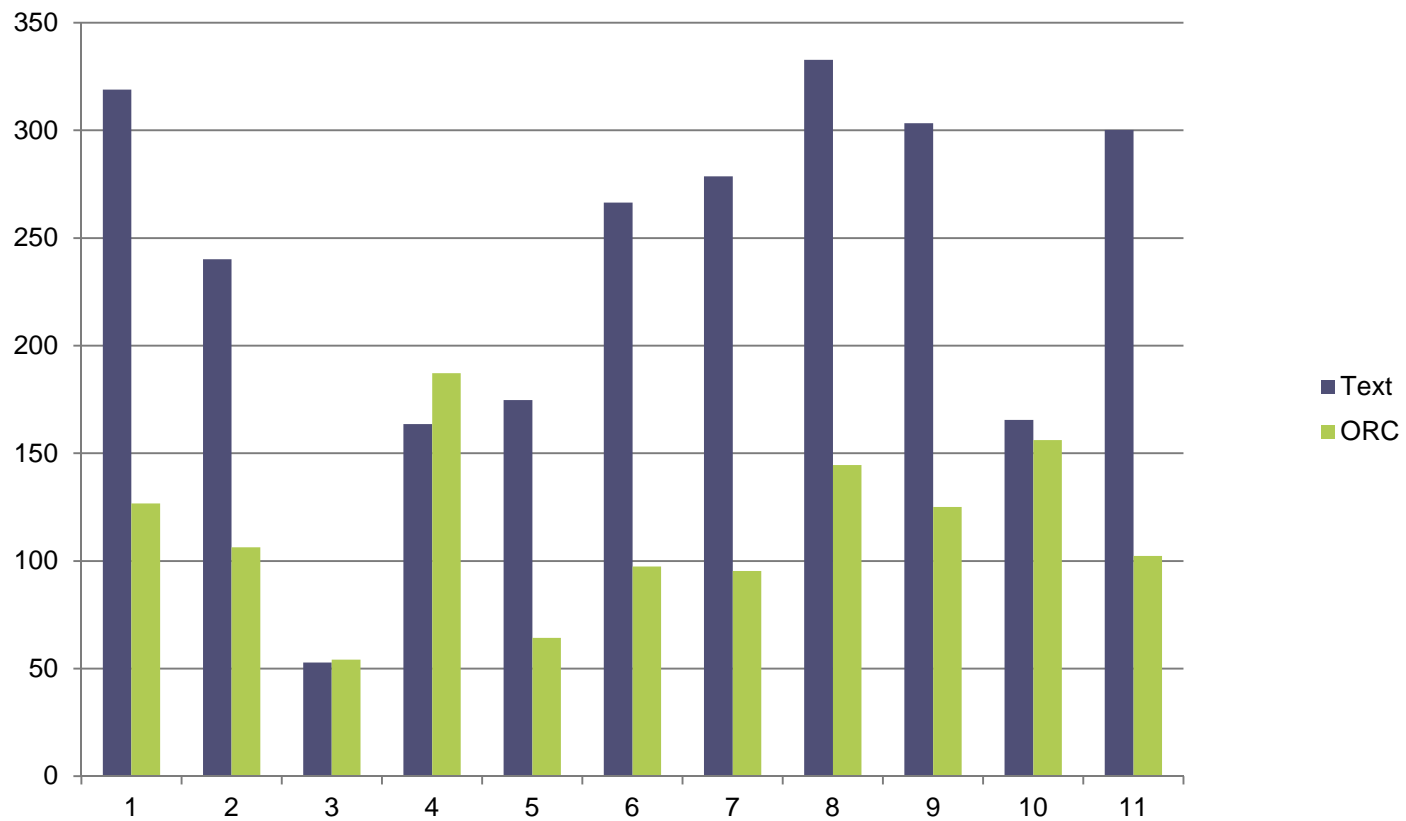
Stärken

- Stabil
- Intensive Weiterentwicklung
- Memoryanforderungen
- Batchprocessing
- Umfangreiche SQL Unterstützung
- Hive Metastore

Schwächen

- Interaktive Abfragen
- Map-Reduce
- Keine Einzelsatzverarbeitung

Performance Vergleich Speichertypen



Überblick

- Ein Datawarehouse System um strukturierte Daten im HDFS zu speichern
- Toplevel Projekt der Apache Foundation
- Unterstützt zahlreiche Dateiformate
- Ermöglicht einen einfachen SQL basierten Zugriff
- Kann Hive Metadata Store benutzen
- Fully distributed SQL processing



SQL Frameworks

Tajo Architecture



Architektur

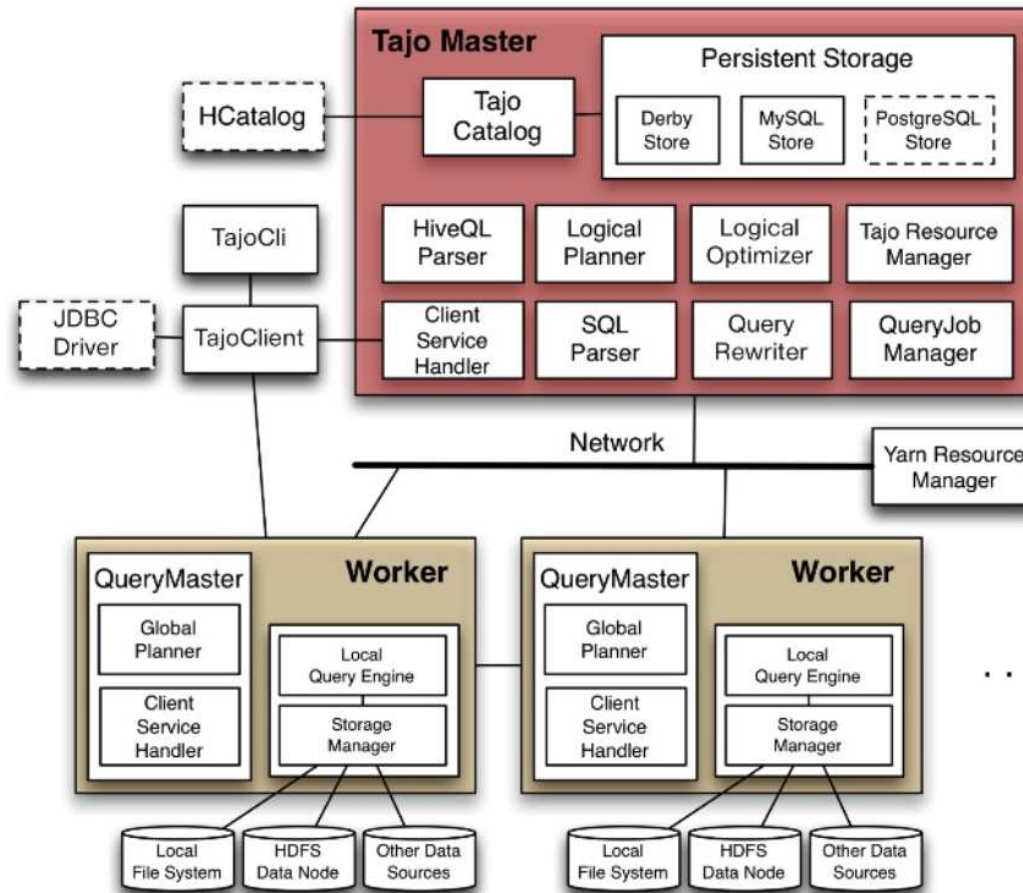


Image from apache.org

SQL Frameworks

Tajo – Stärken und Schwächen



Stärken

- **Hive Metastore möglich**
- **Performance**
- **Aufstieg zum Apache Toplevel Project**

Schwächen

- **Installation auf aktuellen Cloudera bzw. Hortonworks**
- **Dokumentation**

SQL Unterstützung

- Basis SQL
- Union (all)
- Group by

Nicht unterstützt

- Analytische Funktionen
- Subselect nur mit Alias
- substr
- Minus
- Intersect

Überblick

- MPP SQL engine
- Nicht Map-Reduce basiert
- Für interaktive SQL-Abfragen
- Unterstützt zahlreiche Dateiformate
- Verwendet Hive Metadata Store
- Daten im HDFS oder Hbase gespeichert



SQL Frameworks

Impala Architecture



Architektur

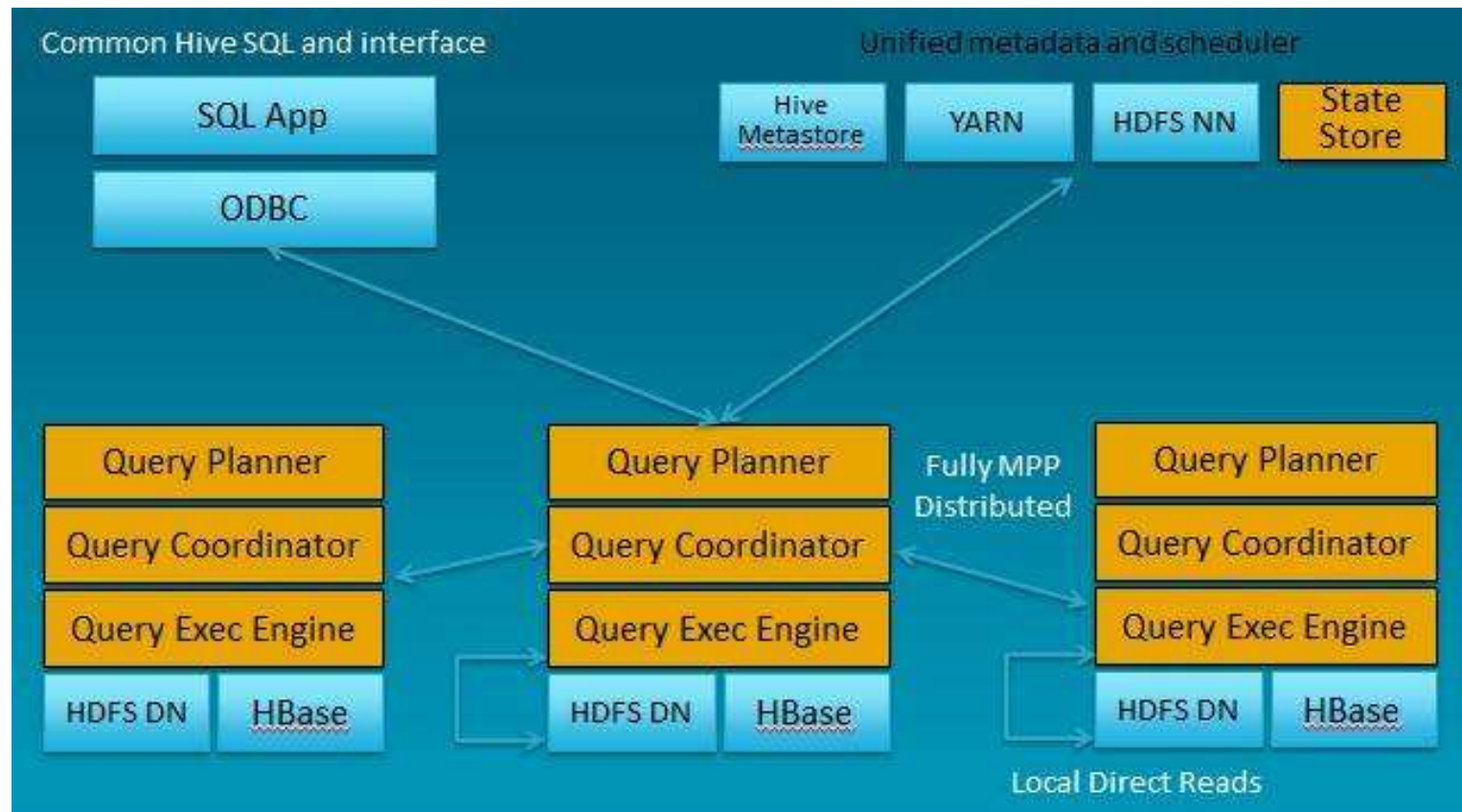


Image from Cloudera.com

Stärken

- Performance
- Nutzung von Hive-Komponenten
- Unterstützung zahlreicher Dateiformate, Hbase und Hive

Schwächen

- RAM-Anforderungen
128 GB empfohlen
- Abbrüche durch Out-of-Memory
- SQL Unterstützung geringer als bei Hive
- Keine User Defined Functions (UDF)

SQL Unterstützung

- Basis SQL
- UNION

Nicht unterstützt

- Analytische Funktionen
- Group by ohne Limit
- Rollup
- Subselects ohne Alias
- MINUS
- INTERSECT

Überblick

- Fully distributed SQL processing
- Nicht Map-Reduce basiert
- In-Memory oder disk basiert
- Unterstützt zahlreiche Dateiformate
- Verwendet Hive Metadata Store
- Verwendet HIVE Frontend
- Spaltenbasierte Tabellen in-Memory
- Setzt auf Apache Spark auf



SQL Frameworks

Apache Spark



Architektur

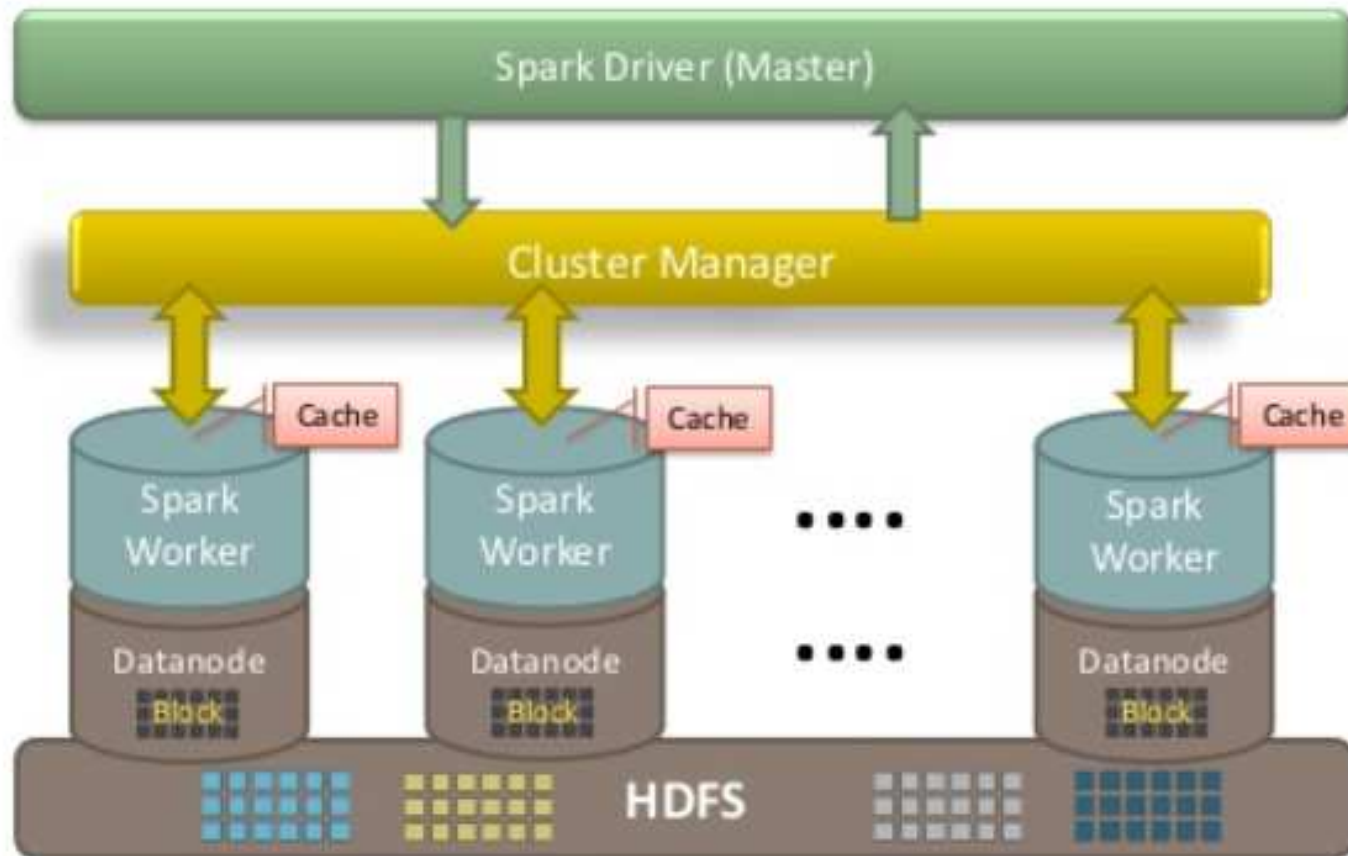


Image from apache.org

SQL Frameworks

Shark Architecture



Architektur

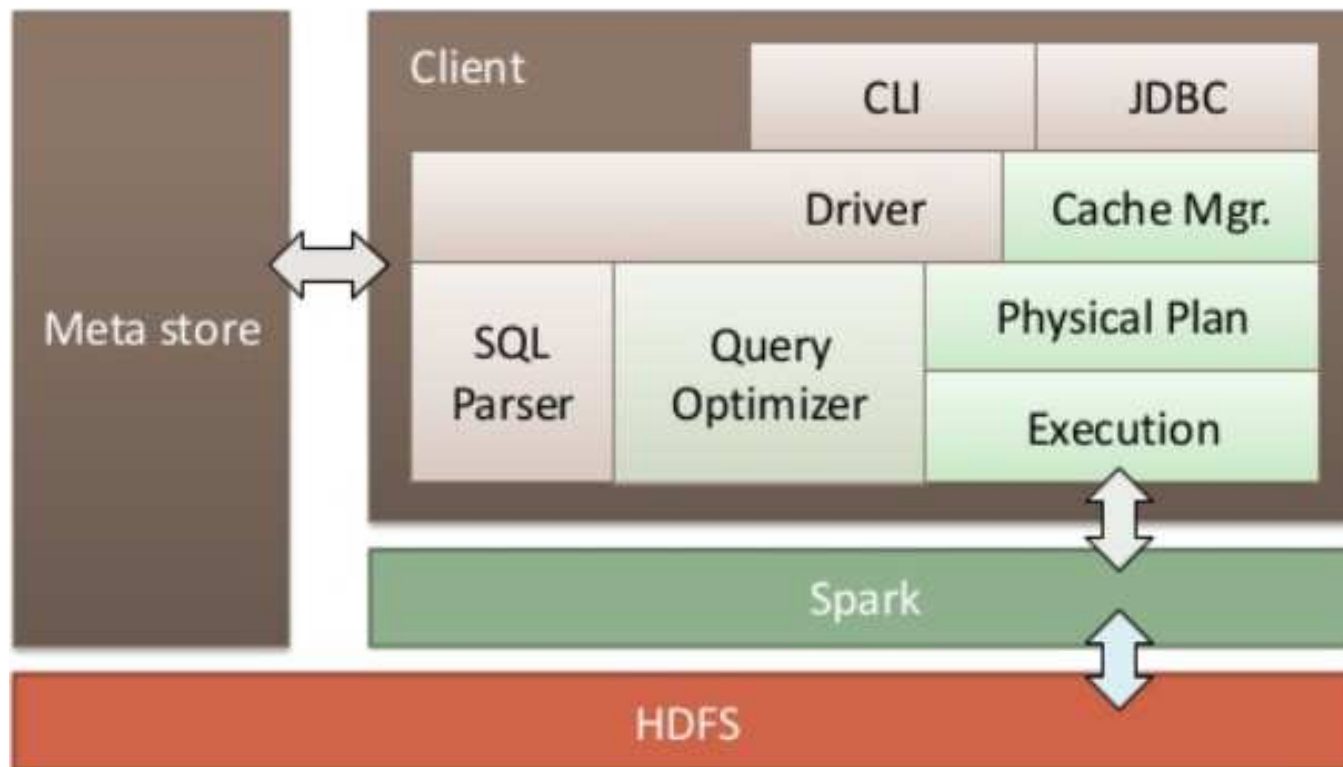


Image from apache.org

Stärken

- Performance
- Abbruchtolerant
- Nutzung von Hive-Komponenten
- Unterstützung zahlreicher Dateiformate und Hive

Schwächen

- Benötigt passende Hive Version
- Benötigt Spark
- Memorybedarf für hohe Performanceanforderungen

Unterstützung

- Analytische Funktionen
- Subselects
- Group by
- Rollup
- UNION ALL

Nicht unterstützt:

- UNION
- MINUS
- INTERSECT

1

Einführung

2

Hive, Impala, Tajo and Co

3

Performance

4

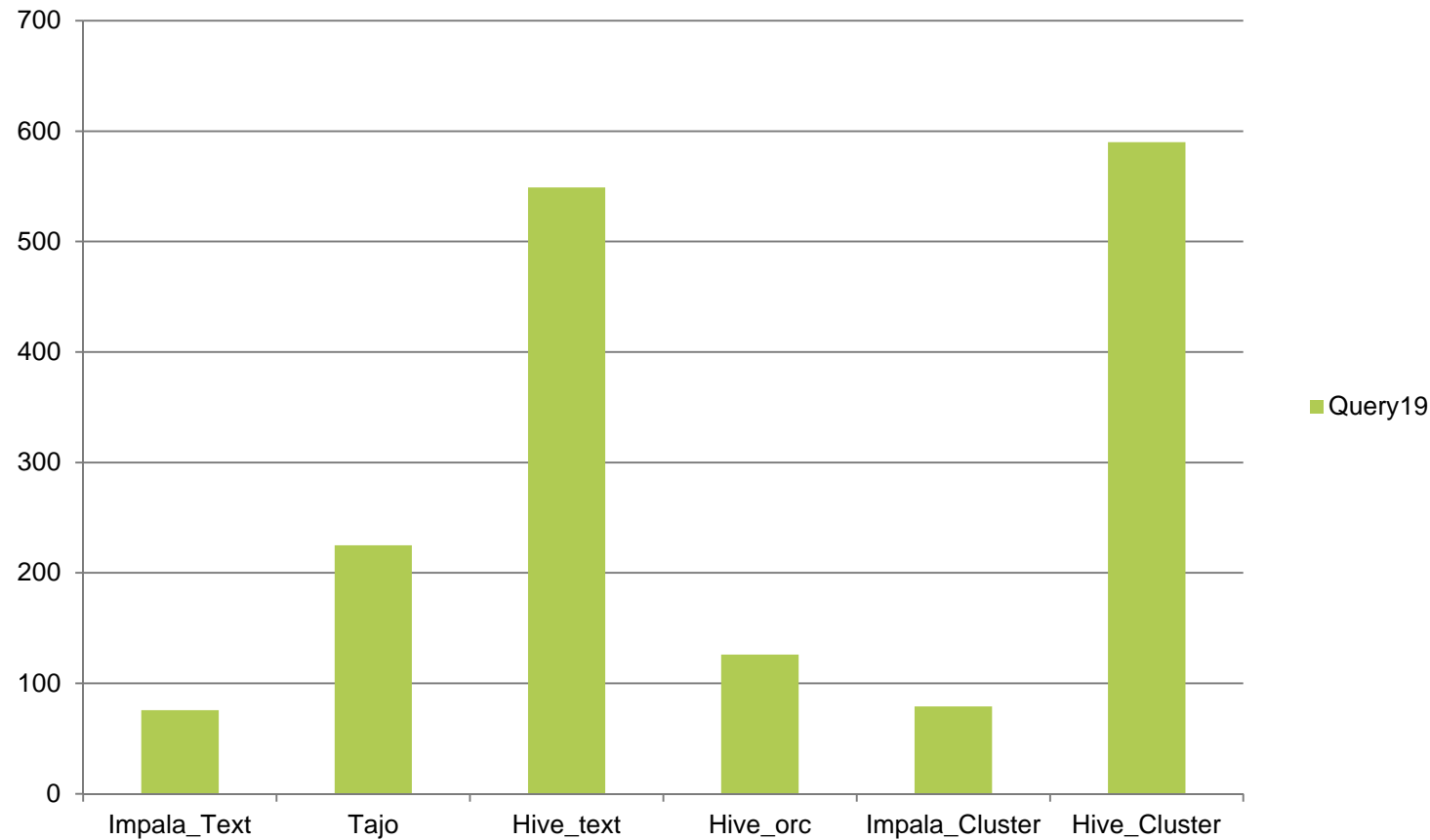
Ausblick

Performance

Deutliche Unterschiede zwischen den Frameworks



Query19



1

Einführung

2

Hive, Impala, Tajo and Co

3

Performance

4

Ausblick

Fazit

- Alle vorgestellten Frameworks erleichtern den Zugriff auf Big Data
- Hive Metastore als übergreifendes Element
- Tradeoff zwischen Hardwareanforderungen und Performance
- Die Entwicklung ist sehr dynamisch
- Unterschiedliche Reifegrade



Consulting

Beratung

Results,
no Excuses.

Lösungen

Grown from
Experience.

Ventum Consulting

Office München

Infanteriestr. 11a
D-80797 München

Telefon +49 89 1222 1964 2
Fax +49 89 1222 1964 25

Office Wien

Ernst Melchior Gasse 24
A-1020 Wien

Telefon +43 (1) 535 34 22 - 0
Fax +43 (1) 535 34 22 - 99

www.ventum-consulting.com