



Advisory [\(https://www.thedigitaltransformationpeople.com/services/\)](https://www.thedigitaltransformationpeople.com/services/).

[.https://www.thedigitaltransformationpeople.com](https://www.thedigitaltransformationpeople.com)

# ta Analytics Reference Architectures – Bi on Facebook, LinkedIn and Twitter

January 18,  
2016

---



**B**ig Data is becoming a new technology focus both in science and industry, and motivate technology shift to data centric architecture and operational models.

There is a vital need to define the basic information/semantic models, architecture components and operational models that together comprise a so-called Big Data Ecosystem.

[nsformationpeople.com/authors/birendra-](https://www.thedigitaltransformationpeople.com/authors/birendra-kumar-sahu/)

[kumar-sahu/](https://www.thedigitaltransformationpeople.com/authors/birendra-kumar-sahu/)).

Birendra Kumar Sahu  
FirstHive

[rmationpeople.com/supplier\\_directory/firsthive/](https://www.thedigitaltransformationpeople.com/supplier_directory/firsthive/)  
0590-484b-9e35-a81a31e59ad8).

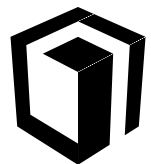
## Extended Relational Reference Architecture:

thive

[about Birendra Kumar Sahu](https://www.thedigitaltransformationpeople.com/authors/birendra-kumar-sahu/)

[rmationpeople.com/authors/birendra-kumar-sahu/](https://www.thedigitaltransformationpeople.com/authors/birendra-kumar-sahu/)).

---



[.https://www.thedigitaltransformationpeople.com](https://www.thedigitaltransformationpeople.com)

---

## Popular Now

The Case For Digital  
Transformation  
An Executive  
Summary:  
Leading Dig

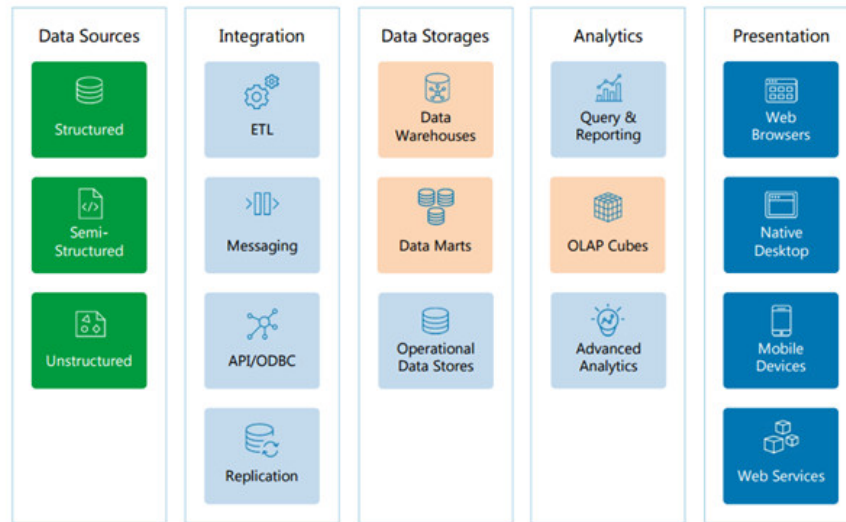


This is more about Non-Relational Reference Architecture but components with pink blocks cannot handle big data challenges completely.

by George Westernman, Didier Bonnet & Andrew McAfee (https://www.thedigitaltransformationpeople.com/case-for-digital-transformation/lead-digital-a-summary/)

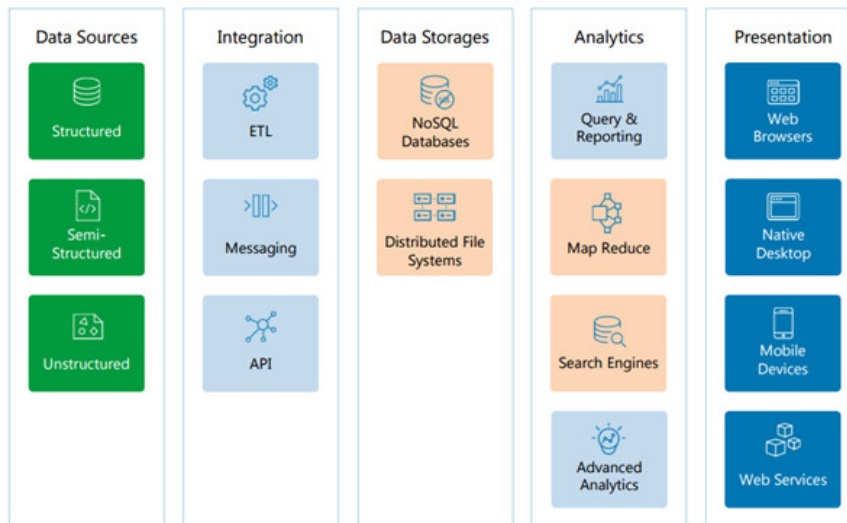
Join us for networking & quality resources to help you and your team succeed in digital transformation.

Join us



### Non-Relational Reference Architecture:

This is more about Non-Relational Reference Architecture but still components with pink blocks cannot handle big data challenges completely.



### Data Discovery: Big Data Architecture

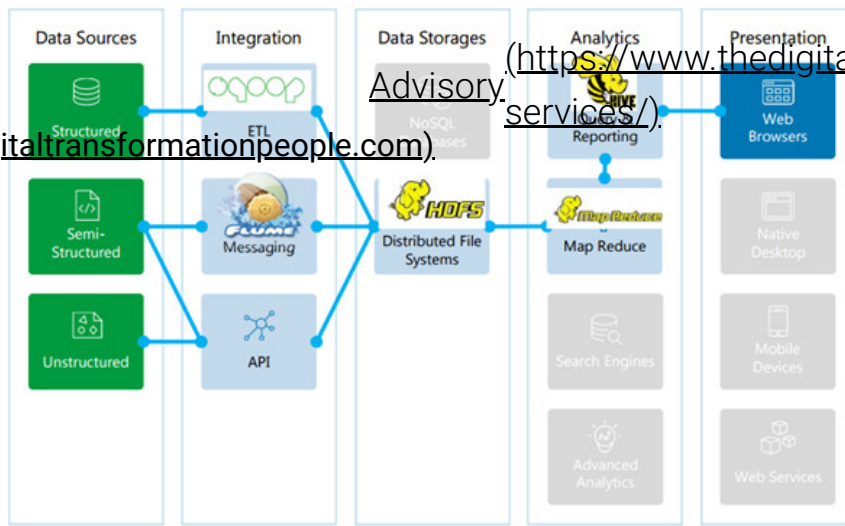
This is more about Hadoop based Big Data Architecture which can be handle few core components of big data challenges but not all (like Search Engine etc)

The Case For Digital Transformation  
 The Digital Transformation Pyramid: A Business-driven Approach for Corporate Initiatives (https://www.thedigitaltransformationpeople.com/case-for-digital-transformation-digital-transformation-pyramid-business-driven-approach-corporate-initiatives/)

Strategy & Innovation Target Operating Models & Roadmaps for Change (https://www.thedigitaltransformationpeople.com/innovation/target-operating-models-roadmaps-for-change/)



(<https://www.thedigitaltransformationpeople.com>).



Advisory

(<https://www.thedigitaltransformationpeople.com>) models- roadmaps-for-change/

### Data analytics Architecture adopted by Facebook:

Data analytics infrastructure at Facebook has been given below. Facebook collects data from two sources. Federated MySQL tier contains user data, and web servers generate event based log data. Data from the web servers is collected to Scribe servers, which are executed in Hadoop clusters.

The Scribe servers aggregate log data, which is written to Hadoop Distributed File System (HDFS). The HDFS data is compressed periodically, and transferred to Production Hive-Hadoop clusters for further processing. The Data from the Federated MySQL is dumped, compressed and transferred into the Production Hive-Hadoop cluster.

Digitisation alone will not transform your business, people will.



Talk to us to secure the people you need.

**Learn More**



(<https://cta-redirect.hubspot.com/cta/redirect/644390/8693db58-66ff-40e8-81af-8e6ca2658ecd>).

Delivery  
Data Asset  
Management  
(DAM)  
(<https://www.thedigitaltransformationpeople.com>)  
asset-  
management-  
dam/)

Strategy & Innovation  
The Innovation  
Management  
Theory  
Evolution Map  
(<https://www.thedigitaltransformationpeople.com>)  
and-  
innovation/the-  
innovation-  
management-  
theory-  
evolution-map/)

### Related Articles



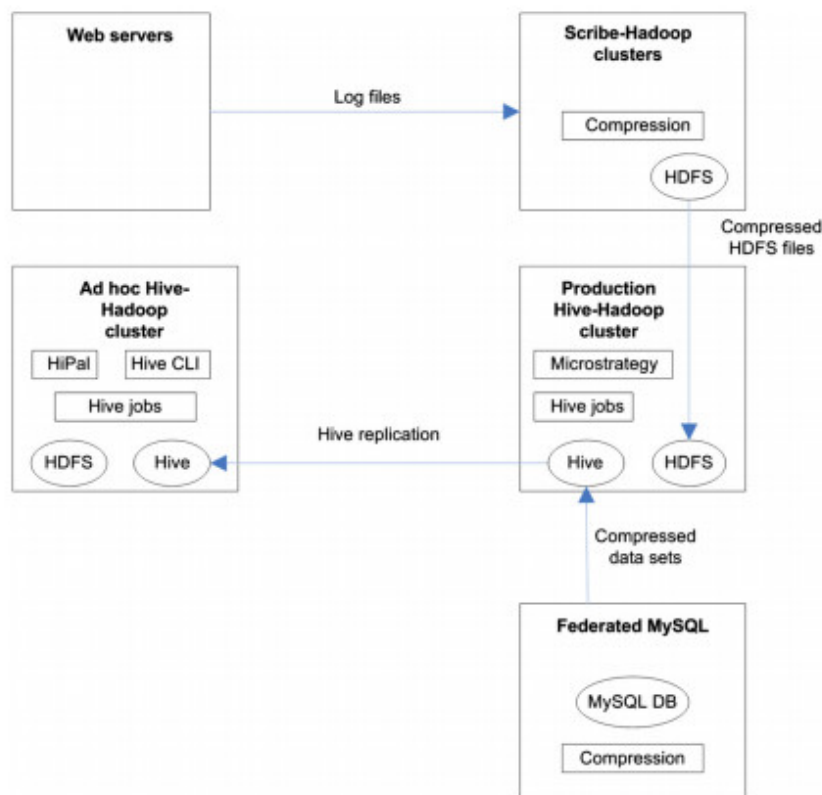
Data & Analytics



Facebook uses two different clusters for data analysis. Jobs with strict deadlines are executed in the Production Hive-Hadoop cluster. Lower priority jobs and ad hoc analysis jobs are executed in Ad hoc Hive-Hadoop cluster. Data is replicated from the Production cluster to the Ad hoc cluster.

The results of data analysis are saved back to Hive-Hadoop cluster or to the MySQL tier for Facebook users. Ad hoc analysis queries are specified with a graphical user interface (HiPal) or with a Hive command-line interface (Hive CLI).

Facebook uses a Python framework for execution (Databee) and scheduling of periodic batch jobs in the Production cluster. Facebook also uses Microstrategy Business Intelligence (BI) tools for dimensional analysis.



### Data analytics Architecture adopted by LinkedIn:

The data analytics infrastructure at LinkedIn has been given below. Data is collected from two sources: database snapshots and activity data from users of LinkedIn.

MR Realities – Conjoint Analysis: Making it Work for You”  
(<https://www.thedigitaltransformationpeople.com/mr-realities-conjoint-analysis-making-it-work-for-you/>).



People & Change  
Engaging employees in your safety program  
(<https://www.thedigitaltransformationpeople.com/people-change-engaging-employees-in-your-safety-program/>).





The activity data comprises streaming events, which is collected based on usage of LinkedIn's services. Kafka is a distributed messaging system, which is used for collection of the streaming events.

Kafka producers report events to topics at a Kafka broker, and Kafka consumers read data at their own pace. Kafka's event data is transferred to Hadoop ETL cluster for further processing (combining, de-duplication).

Data from the Hadoop ETL cluster is copied into production and development clusters. Azkaban is used as a workload scheduler, which supports a diverse set of jobs.

An instance of Azkaban is executed in each of the Hadoop environments. Scheduled Azkaban workloads are realised as MapReduce, Pig, shell script, or Hive jobs. Typically workloads are experimented in the development cluster, and are transferred to the production cluster after successful review and testing.

Results of the analysis in the production environment are transferred into an offline debugging database or to an online database.

Results may also be fed back to the Kafka cluster. Avatara is used for preparation of OLAP data. Analysed data is read from the Voldemort database, pre-processed, and aggregated/cubified for OLAP, and saved to another Voldemort read-only database.

Enabling Technologies  
Open-source  
bioinformatic  
solutions for  
'Big Data'  
analysis  
(<https://www.thedigitaltransformationpeople.com/technologies/open-source-bioinformatic-solutions-for-big-data-analysis/>).



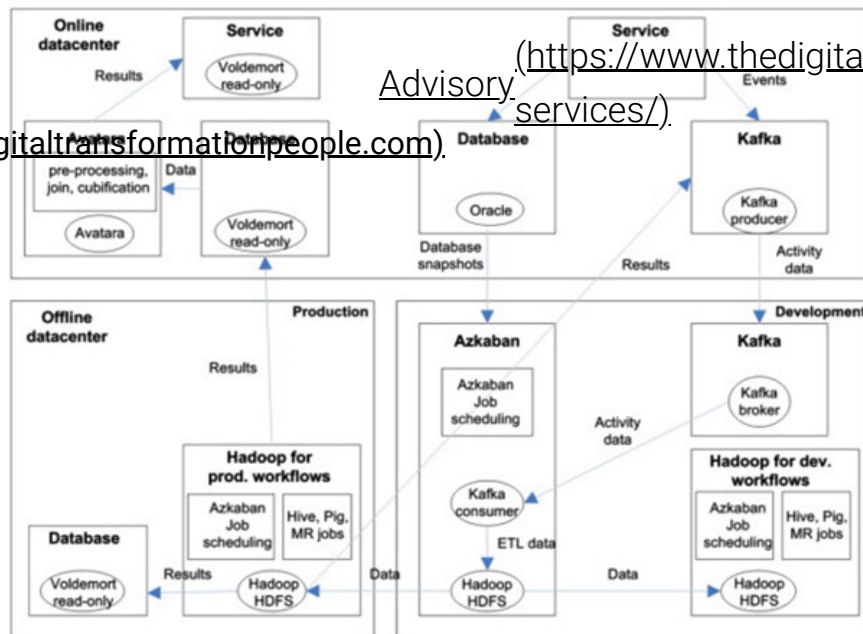
Learn from peers, thought leaders and expert practitioners for the best chance of success.

Digital transformation is a complex and difficult task. Join the community for quality resources, peer support and expert help.

**Join us** ➤







### Data analytics Architecture adopted by Twitter:

In the Twitter's infrastructure for real-time services, a Blender brokers all requests coming to Twitter. Requests include searching for tweets or user accounts via a QueryHose service. Tweets are input via a FireHose service to an ingestion pipeline for tokenization and annotation.

Subsequently, the processed tweets enter to EarlyBird servers for filtering, personalization, and inverted indexing .

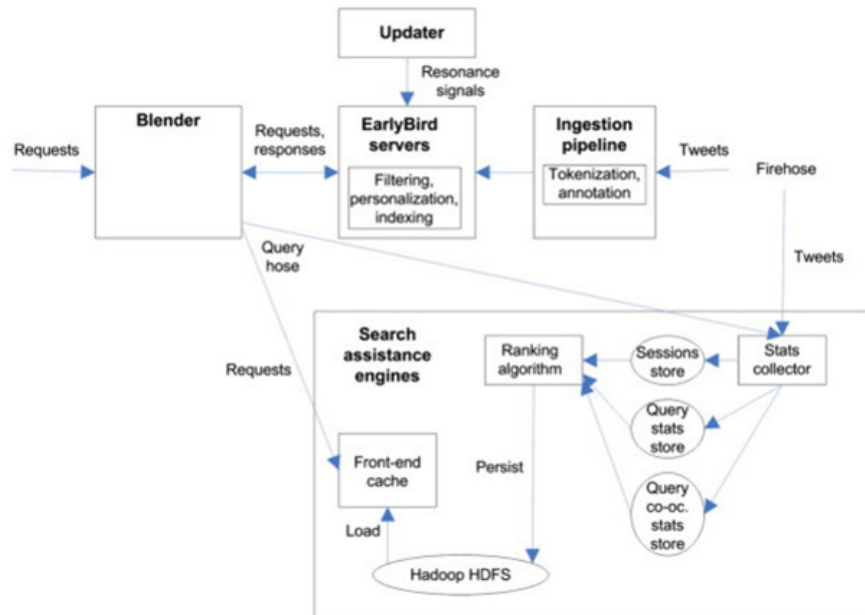
The EarlyBird servers also serve incoming requests from the QueryHose/Blender. The EarlyBird is a real-time retrieval engine, which was designed for providing low latency and high throughput for search queries.

Additionally, search assistance engines are deployed. Stats collector in the Search assistance engine saves statistics into three in-memory stores, when a query or tweet is served.

User sessions are saved into Sessions store, statistics about individual queries are saved into Query statistics store, and statistics about pairs of co-occurring queries are saved into Query co-occurrence store.



A ranking algorithm fetches data from the in-memory stores, and analyses the data. The results of analysis are persisted into Hadoop HDFS. Finally, Front-end cache polls results of analysis from the HDFS, and serves users of Twitter. (https://www.thedigitaltransformationpeople.com/)



Twitter has three streaming data sources (Tweets, Updater, queries), from which data is extracted. Tweets and queries are transmitted over REST API in JSON format.

Thus, they can be considered as streaming, semi-structured data. The format of data from Updater is not known (streaming data source). Ingestion pipeline and Blender can be considered as Stream temp data stores.

Tokenization, annotation, filtering, and personalization are modelled as stream processing. EarlyBird servers contain processed stream-based data (Stream data store). Stats collector is modelled as stream processing.

[Successful digital transformation is a matter of know how and access to the best talent. We connect you to both. Click for more.](https://cta-redirect.hubspot.com/cta/redirect/644390/07ba6b3c-83ee-4495-b6ec-b2524c14b3c5)

(https://cta-

redirect.hubspot.com/cta/redirect/644390/07ba6b3c-83ee-4495-b6ec-b2524c14b3c5).

The statistical stores may be considered as Stream data stores, which store structured information of processed data. The ranking algorithm performs Stream analysis functionality.





Hadoop HDFS storing the analysis results is modelled as a  
Stream analysis data store. Front end cache (Serving data  
store) serves the End user application (Twitter app).  
(<https://www.thedigitaltransformationpeople.com>)

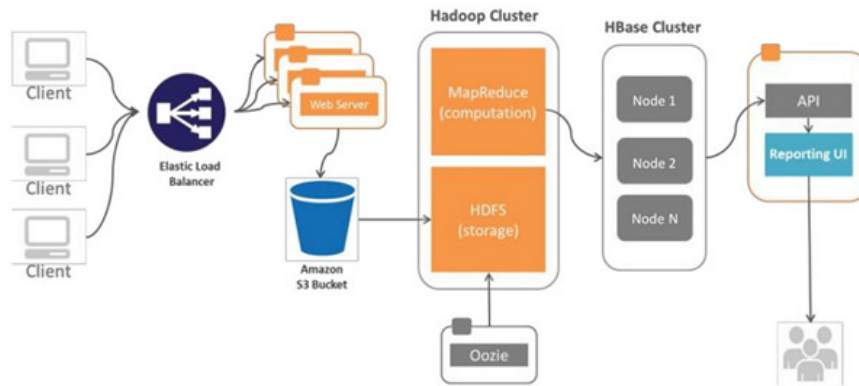
Reference: Reference Architecture and Classification of  
Technologies by Pekka Pääkkönen and Daniel Pakkala  
(facebook, twitter and linkedin Reference Architecture  
mentioned here are derived from this publication )

AWS cloud based Solution Architecture (ClickStream  
Analysis):

## Solution Architecture

### Technologies:

- Amazon S3
- Flume
- Hadoop/HDFS, MapReduce
- HBase
- Oozie
- Hive



(<https://www.thedigitaltransformationpeople.com>)

CTO & Vice President

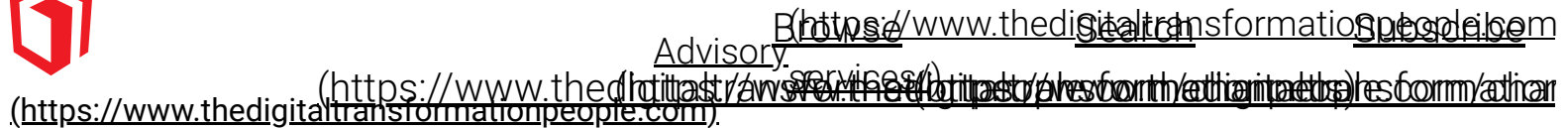
(<https://www.thedigitaltransformationpeople.com/authors/birendra-kumar-sahu/>)

(<https://www.thedigitaltransformationpeople.com>)

(<https://www.thedigitaltransformationpeople.com>)







The best  
articles,  
news and  
events  
direct to  
your inbox

(<https://www.thedigitaltransformationpeople.com/channels/tag/featured/>).



[Terms of Use](#)

[Medium](#)

[\(https://www.thedigitaltransformationpeople.com/interim-consultancy/\)](https://www.thedigitaltransformationpeople.com/interim-consultancy/)

[\(https://www.thedigitaltransformationpeople.com/terms-of-use/\)](https://www.thedigitaltransformationpeople.com/terms-of-use/)

[\(https://www.thedigitaltransformationpeople.com/terms-of-use/\)](https://www.thedigitaltransformationpeople.com/terms-of-use/)

[services/](#)

[\(https://www.thedigitaltransformationpeople.com/\)](https://www.thedigitaltransformationpeople.com/)

[. \(https://www.youtube.com/channel/UCPDv4OvDD4Z8uQIIQ78PFPg\)](https://www.youtube.com/channel/UCPDv4OvDD4Z8uQIIQ78PFPg)

[GooglePlus](#)

[\(https://plus.google.com/b/100832677067367192859/100832677067367192859/](https://plus.google.com/b/100832677067367192859/100832677067367192859/)

[gmbpt=true&pageld=100832677067367192859&hl=en-](#)

[GB\).](#)

## Talent Solutions

---

## Career Opportunities

---

[Executive Search](#)

[Candidate Registration](#)

[\(https://www.thedigitaltransformationpeople.com/candidate-registration/\)](https://www.thedigitaltransformationpeople.com/candidate-registration/)

[Interim Consultancy](#)

[\(https://www.thedigitaltransformationpeople.com/interim-consultancy/\)](https://www.thedigitaltransformationpeople.com/interim-consultancy/)

