

# Trabajando con **APACHE ATLAS**

BIG DATA – BUSINESS INTELLIGENCE – MACHINE LEARNING

strate**bi**  
open business intelligence



## CONTENIDO

1. INTRODUCCIÓN .....	3
PRICING.....	3
ARQUITECTURA .....	4
INSTALACIÓN Y DESPLIEGUE .....	5
2. DATA DISCOVERY: ALMACÉN DE METADATOS .....	6
INTERFAZ DE USUARIO .....	7
GLOSARIOS.....	8
3. DATA LINEAGE: SEGUIMIENTO DE TRANSFORMACIONES .....	9
4. DATA SECURITY: AUTENTIFICACIÓN Y PERMISOS .....	10
5. CONCLUSIONES.....	11
6. SOBRE STRATEBI.....	12
7. TECNOLOGÍAS.....	13
8. REFERENCIAS.....	16
9. EJEMPLOS DE DESARROLLOS ANALYTICS .....	17

# 1. INTRODUCCIÓN

[Apache Atlas](#) es una herramienta *open-source*, con licencia Apache 2.0, para la gobernanza del dato la cual permite la integración con todo el ecosistema de datos de las empresas.

Atlas permite crear un almacén de metadatos centralizado para saber dónde encontrar un conjunto de datos dentro de la empresa (**Data Discovery**), permite saber qué cambios ha sufrido y que transformaciones se le han realizado a los datos originales a lo largo del tiempo (**Data Lineage**) y centraliza la seguridad para saber quién puede acceder a esos datos y quién puede modificarlos (**Data Security**) apoyándose en [Apache Ranger](#).

## Pricing

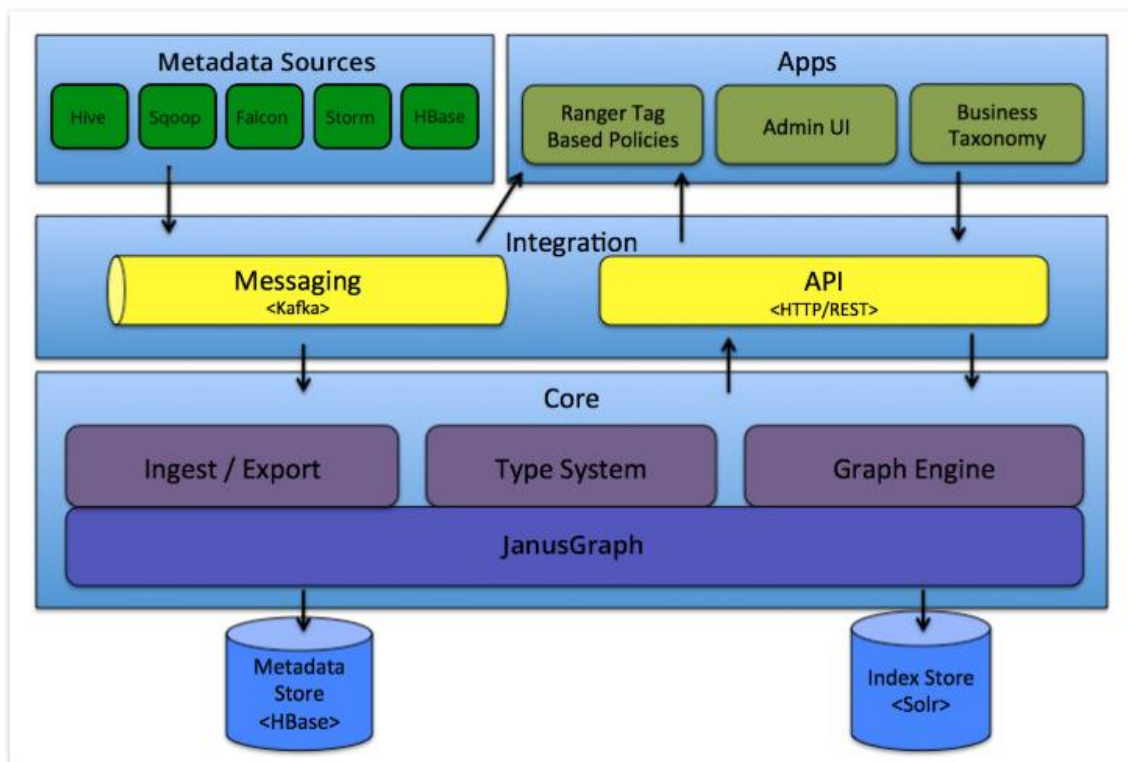
Actualmente está integrada dentro de la plataforma de Cloudera, la cual ofrece además una versión Enterprise que añade soporte por parte de la empresa. Su [pricing](#) es el siguiente:

Platform pricing			
Cloudera Platforms Annual subscription pricing	CDP Data Center	Enterprise Data Hub	HDP Enterprise Plus
<b>Base price</b>			
Per node	\$10,000 <sup>1</sup>	\$10,000	\$10,000
<b>Variable price - compute<sup>2</sup></b>			
Per CCU over 16 core, 128 GB node cap	\$75 <sup>3</sup>	\$75 <sup>4</sup>	\$75 <sup>4</sup>
<b>Variable price - storage<sup>2</sup></b>			
Per TB over 48 TB node cap	\$25 <sup>3</sup>	\$25 <sup>4</sup>	\$25 <sup>4</sup>
<b>Conversion entitlement<sup>5</sup></b>			
For CDP Data Center	-	Yes	Yes

Se puede obtener una comparativa entre las distintas versiones a través de este [enlace](#).

## Arquitectura

La [arquitectura](#) de Apache Atlas y las herramientas en las que se apoya son las siguientes:



Como fuente de metadatos puede utilizar: [HBase](#), [Hive](#), [Sqoop](#), [Storm](#), [Kafka](#) y Falcon (en desuso).



## Instalación y despliegue

Apache Atlas se puede instalar a través de la plataforma de [Apache Ambari](#) que se puede instalar en los siguientes sistemas operativos: [RHEL/CentOS/Oracle Linux 7](#), [SLES 12](#), [Ubuntu 16](#), [Ubuntu 18](#) y [Debian 9](#).

Los **prerrequisitos** para su instalación son: HBase, Kafka y SolrCloud/Elasticsearch

Atlas se ha diseñado para desplegarse como **servidor**, pudiendo conectarse a ella desde un navegador.

## 2. DATA DISCOVERY: ALMACÉN DE METADATOS

Atlas permite a los usuarios definir un modelo para los objetos de metadatos que desean administrar. Todos los objetos de metadatos administrados por Atlas que están listos para ser utilizados se modelan usando "Tipos" y se representan como "Entidades". Antes de empezar a utilizar Atlas es importante entender estos conceptos.

Para almacenar y acceder a un tipo particular de objetos de metadatos, Atlas, permite definir "Tipos". Similar a las "clases" de Java o al "esquema de una tabla" en base de datos relacionales.

An example of a type that comes natively defined with Atlas is a Hive table. A Hive table is defined with these attributes:

```
Name:      hive_table
TypeCategory: Entity
SuperTypes: DataSet
Attributes:
  name:      string
  db:        hive_db
  owner:     string
  createTime: date
  lastAccessTime: date
  comment:   string
  retention: int
  sd:        hive_storagedesc
  partitionKeys: array<hive_column>
  aliases:   array<string>
  columns:   array<hive_column>
  parameters: map<string,string>
  viewOriginalText: string
  viewExpandedText: string
  tableType: string
  temporary: boolean
```

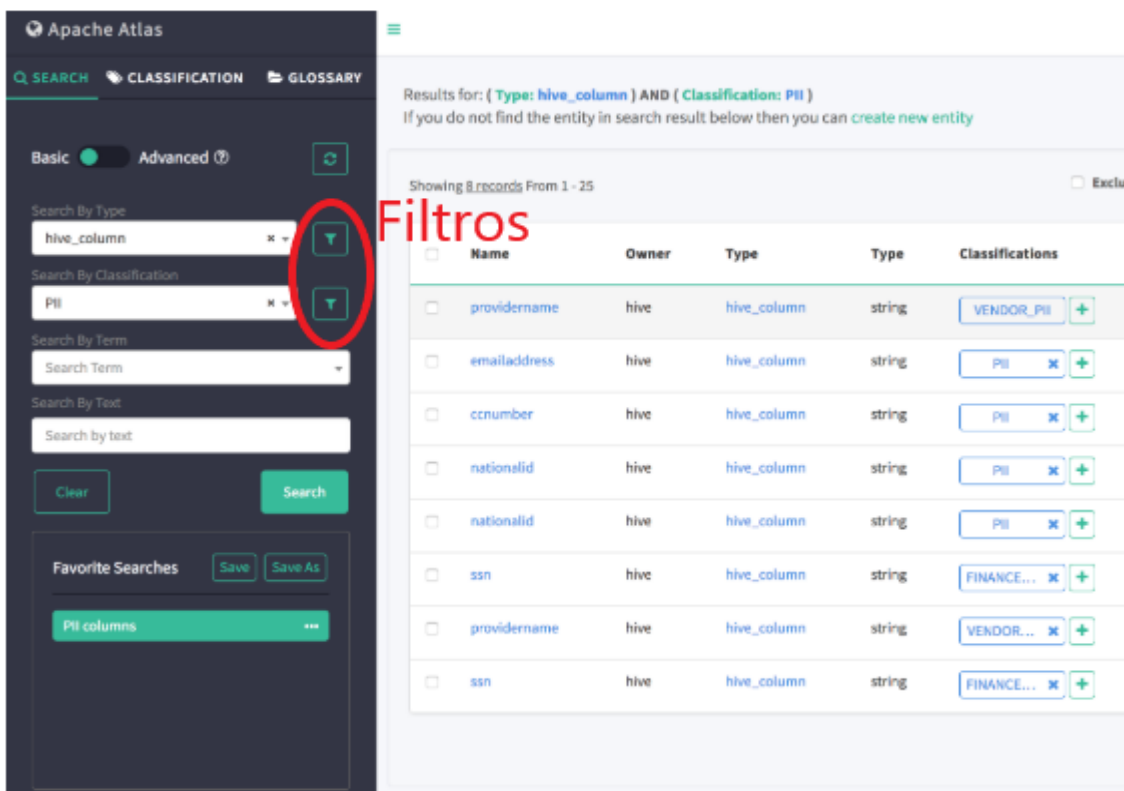
Por cada "Tipo" creado en Atlas se pueden crear diferentes instancias o "Entidades". Similar a la "instanciación de objetos de una clase" en Java.

An example of an entity will be a specific Hive Table. Say Hive has a table called 'customers' in the 'default' database. This table will be an 'entity' in Atlas of type hive\_table. By virtue of being an instance of an entity type, it will have values for every attribute that are a part of the Hive table 'type', such as:

```
guid: "9ba387dd-fa76-429c-b791-ffc338d3c91f"
typeName: "hive_table"
status: "ACTIVE"
values:
  name: "customers"
  db: { "guid": "b42c6cfc-c1e7-42fd-a9e6-890e0adf33bc", "typeName": "hive_db" }
  owner: "admin"
  createTime: 1490761686029
  updateTime: 1516298102877
  comment: null
  retention: 0
  sd: { "guid": "ff58025f-6854-4195-9f75-3a3058dd8dcf", "typeName": "hive_storagedesc" }
  partitionKeys: null
  aliases: null
  columns: [ { "guid": "65e2204f-6a23-4130-934a-9679af6a211f", "typeName": "hive_column" }, { "guid": "d726de70-faca-46fb-9c99-cf04f6b579a6", "typeName": "hive_column" }, ... ]
  parameters: { "transient_lastDdlTime": "1466403208" }
  viewOriginalText: null
  viewExpandedText: null
  tableType: "MANAGED_TABLE"
  temporary: false
```

## Interfaz de usuario

Para empezar a utilizar datos se pueden buscar y encontrar fácilmente con su buscador ya sea de manera [básica](#) utilizando texto y apoyándose en filtros, o de forma más [avanzada](#) mediante DSL. En la siguiente captura se puede observar como sería el buscador básico:



## Glosarios

Para facilitar la utilización de la herramienta a los usuarios de negocios Atlas permite crear **Glosarios** donde se puede relacionar términos y categorías con los datos para que sean más fácil de encontrar. En la siguiente captura se puede ver una clasificación de distintos datos mediante categorías y términos:

The screenshot shows the Apache Atlas GLOSSARY interface. On the left, a sidebar lists various terms under categories like Automotive, Finance, and Insurance. The 'Claim' term is highlighted. The main panel shows the details for 'Claim', including short and long descriptions, a classification of 'INSURANCE', and a table of related entities. A red circle highlights the 'Insurance' category in the sidebar, and another red circle highlights the 'Claim' term. Red text labels 'Categorías' and 'Términos' are overlaid on the image.

Name	Owner	Description	Type	Classifications
/hive_data/cost_savings			hdfs_path	INSURANCE

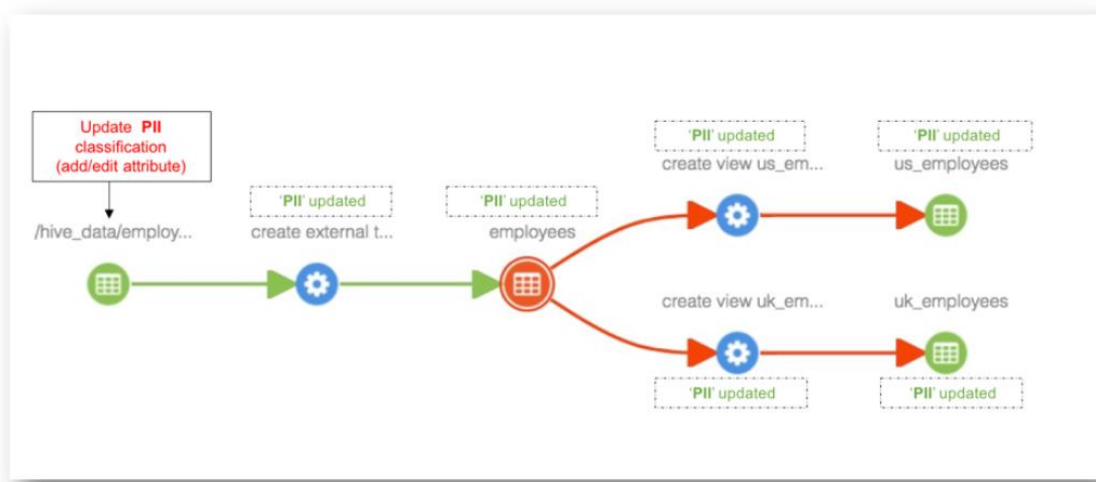


### 3. DATA LINEAGE: SEGUIMIENTO DE TRANSFORMACIONES

Apache Atlas permite hacer un **seguimiento de las transformaciones** que han sufrido los datos originales y en que datos derivados han terminado para poder tener un control total sobre ellos. En esta captura se puede observar un caso de uso:

#### Update classification associated with an entity

Any updates to classifications associated with an entity will be seen in all entities the classification is propagated to as well.



Para ver más ejemplos se puede acceder a [este enlace](#).

## 4. DATA SECURITY: AUTENTIFICACIÓN Y PERMISOS

Atlas permite la [autenticación de personal](#) con los siguientes métodos: LDAP, File y Kerberos. Se pueden utilizar varios a la vez y en caso de que falle el que se ha establecido como método principal pasará a utilizar otro como escenario alternativo.

Además, permite la creación de usuarios con distintos **permisos** de [forma básica](#) mediante Java reg-ex y a estos usuarios se les puede modificar la política de autorizaciones de [forma avanzada](#) apoyándose en Apache Ranger. En la siguiente captura se puede ver un ejemplo del editor avanzado:

Following authorization policy allows user 'admin' perform all operations on metadata entities of Hive database named "my\_db".

**Policy Details :**

Policy Type **Access** ⊕ Add V...

Policy Name \* Hive entiles for my\_db **enabled** **normal**

Policy Label Policy Label

entity-type \* **hive\*** **include**

Entity Classification \* **\*** **include**

Entity ID \* **my\_db.\*** **include**

Description Access to metadata entiles for Hive database my\_db

Audit Logging **YES**

**add/edit permissions**

- Read Entity
- Create Entity
- Update Entity
- Delete Entity
- Read Classification
- Add Classification
- Update Classification
- Remove Classification
- Select/Deselect All

**Allow Conditions :**

Select Group	Select User	gate min
Select Group	<b>admin</b>	<input type="checkbox"/>

**Add Permissions**

## 5. CONCLUSIONES

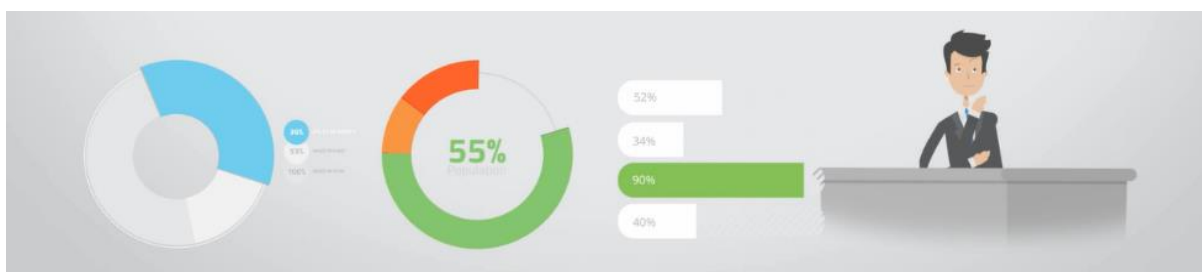
Apache Atlas es una herramienta eficaz para implementar el Gobierno del Dato ya que abarca los múltiples campos que se requieren para una buena gobernanza de dato, la herramienta es *open-source* y está diseñada para ser extensible y añadir funcionalidades a través de una API REST para así poder ofrecer una solución global. La empresa Comcast, por ejemplo, ha añadido a Atlas la posibilidad de [añadir esquemas de Apache Avro](#), lo que le permite manejar de manera más eficiente la calidad del dato.

Aunque eficaz, desplegarla y utilizarla puede ser una tarea compleja ya que requiere de gran número de herramientas para que funcione correctamente (Apache Ambari, Apache Ranger, Kafka, HBase...), por lo que puede terminar en no ser una solución muy eficiente.

## 6. SOBRE STRATEBI

En Stratebi ofrecemos **gran cantidad de soluciones analíticas** por una compañía de **rápido crecimiento**, innovando en las áreas tecnológicas de mayor desarrollo en la actualidad: **Business Intelligence, Big Data y Social Intelligence**, muchas de ellas, basadas en soluciones **Open Source**.

Además, somos **Partners Certificados en Microsoft PowerBI, Talend, Kylin, Snowflake, Pentaho y Vertica**, con gran número de proyectos con ambas tecnologías



**Stratebi** es una empresa española, con sede en Madrid y oficinas en Barcelona, Alicante y Sevilla, creada por un grupo de profesionales con amplia experiencia en sistemas de información, soluciones tecnológicas y procesos relacionados con soluciones de Open Source y de inteligencia de Negocio.

Esta experiencia, adquirida durante la participación en proyectos estratégicos en compañías de reconocido prestigio a nivel internacional, se ha puesto a disposición de nuestros clientes a través de Stratebi.

**Stratebi es la única empresa española que ha estado presente todos los Pentaho Developers celebrados en Europa** (Mainz-Alemania, Barcelona, Lisboa, Roma, Amsterdam, Sintra, Amberes (2) Londres...), habiendo organizado el de Barcelona.

En Stratebi nos planteamos como **objetivo** dotar a las compañías e instituciones, de herramientas escalables y adaptadas a sus necesidades, que conformen una estrategia Business Intelligence capaz de rentabilizar la información disponible. Para ello, nos basamos en el desarrollo de soluciones de Inteligencia de Negocio, mediante tecnología Open Source.

Stratebi son **profesores y responsables de proyectos** del Master en Business Intelligence de la Universidad UOC, UCAM, EOI...

Los profesionales de Stratebi son los creadores y autores del primer weblog en español sobre el mundo del Business Intelligence, Data Warehouse, CRM, Dashboards, Scorecard y Open Source.

Stratebi es partner de las principales soluciones Analytics: Microsoft PowerBI, Talend, Pentaho, Vertica, Snowflake, Kylogence, Cloudera...

**Todo Bi**, se ha convertido en una referencia para el conocimiento y divulgación del Business Intelligence en español.

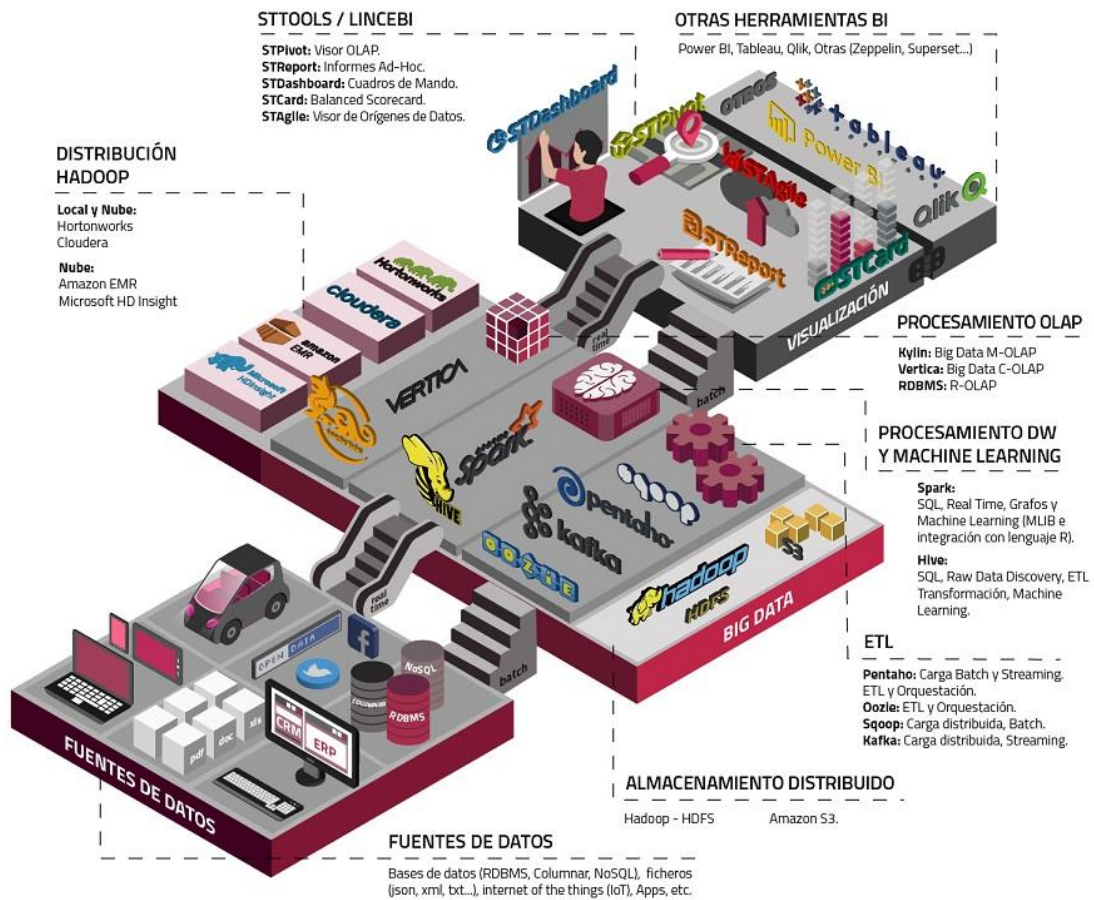
Desarrollamos nuevas soluciones analíticas basadas en Open Source, para la generación de Cuadros de Mando en tiempo real, con tecnologías IoT para SmartCities, machine learning, etc...



## 7. TECNOLOGÍAS

Recientemente, hemos sido nombrados Partners Certificados de Vertica, Talend, Microsoft, Snowflake, Kylligence, Pentaho, etc...





## 8. REFERENCIAS

Trabajamos en los principales sectores y con algunas de las compañías y organizaciones más importantes de España.

### SECTOR PRIVADO



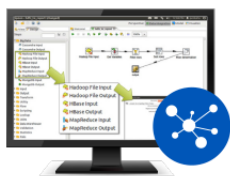
### SECTOR PÚBLICO





## 9. EJEMPLOS DE DESARROLLOS ANALYTICS

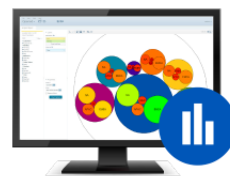
A continuación se presentan **ejemplos de algunos screenshots** de cuadros de mando diseñados por Stratebi, con el fin de dar a conocer lo que se puede llegar a obtener, así como Demos Online en la web de Stratebi:



Data Ingestion  
Manipulation  
Integration



Enterprise and  
Ad Hoc Reporting



Data Discovery  
Visualization



Predictive  
Analytics

Pentaho Analytics Platform

Hadoop

NoSQL

Analytic  
Databases

Relational



