

Büyük Veri & Alet Çantam





Spotify “Discover Weekly” Listesi

The screenshot shows a Spotify interface for a "Discover Weekly" playlist. At the top, there's a small thumbnail labeled "Your Discover Weekly". Below it, the title "Discover Weekly" is displayed in large white letters. To the right are buttons for "PAUSE", a heart icon, and three dots. The main area is a table with four columns: "TITLE", "ARTIST", "ALBUM", and a timestamp. The first row shows "C Section" by Dot Hacker from "N°3" added "5 days ago". The second row shows "Animal" by Badflower from "Temper" added "5 days ago". The third row shows "Bottom of the Lake" by The Builders and The Butchers from "The Builders and the Butchers" added "5 days ago". The fourth row shows "Tired Old Dog" by The Devil and the Almighty Bl... from "The Devil and the Almighty Bl..." added "5 days ago". The fifth row shows "Impossible Dreams" by Versus Them from "Six A.M. Salvation" added "5 days ago". The sixth row shows "The Rains Of Castomere - From The "Game Of Thron..." by The National from "Game Of Thrones: Season 2 (...)" added "5 days ago". The seventh row, which is highlighted with a blue border, shows "Death in Greenpoint" by Mishka Shubaly from "When We Were Animals" added "5 days ago". The eighth row shows "Sugar Boats" by Modest Mouse from "Strangers to Ourselves" added "5 days ago". The ninth row shows "Love" by Colour Haze from "Colour Haze" added "5 days ago".

TITLE	ARTIST	ALBUM	
C Section	Dot Hacker	N°3	5 days ago
Animal	Badflower	Temper	5 days ago
Bottom of the Lake	The Builders and The Butchers	The Builders and the Butchers	5 days ago
Tired Old Dog	The Devil and the Almighty Bl...	The Devil and the Almighty Bl...	5 days ago
Impossible Dreams	Versus Them	Six A.M. Salvation	5 days ago
The Rains Of Castomere - From The "Game Of Thron..."	The National	Game Of Thrones: Season 2 (...)	5 days ago
Death in Greenpoint	Mishka Shubaly	When We Were Animals	5 days ago
Sugar Boats	Modest Mouse	Strangers to Ourselves	5 days ago
Love	Colour Haze	Colour Haze	5 days ago

 **dave horwitz**
@Dave_Horwitz

It's scary how well @Spotify Discover Weekly playlists know me. Like former-lover-who-lived-through-a-near-death experience-with-me well.

190 people are talking about this



 **Amanda Whitbred**
@amandawhitbred

At this point @Spotify's discover weekly knows me so well that if it proposed I'd say yes

12:36 AM - Aug 19, 2016



Instagram'ın aynı zamanda sunduğu farklı profil sayfaları



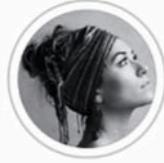
521 posts 141K followers 714 following

Message  

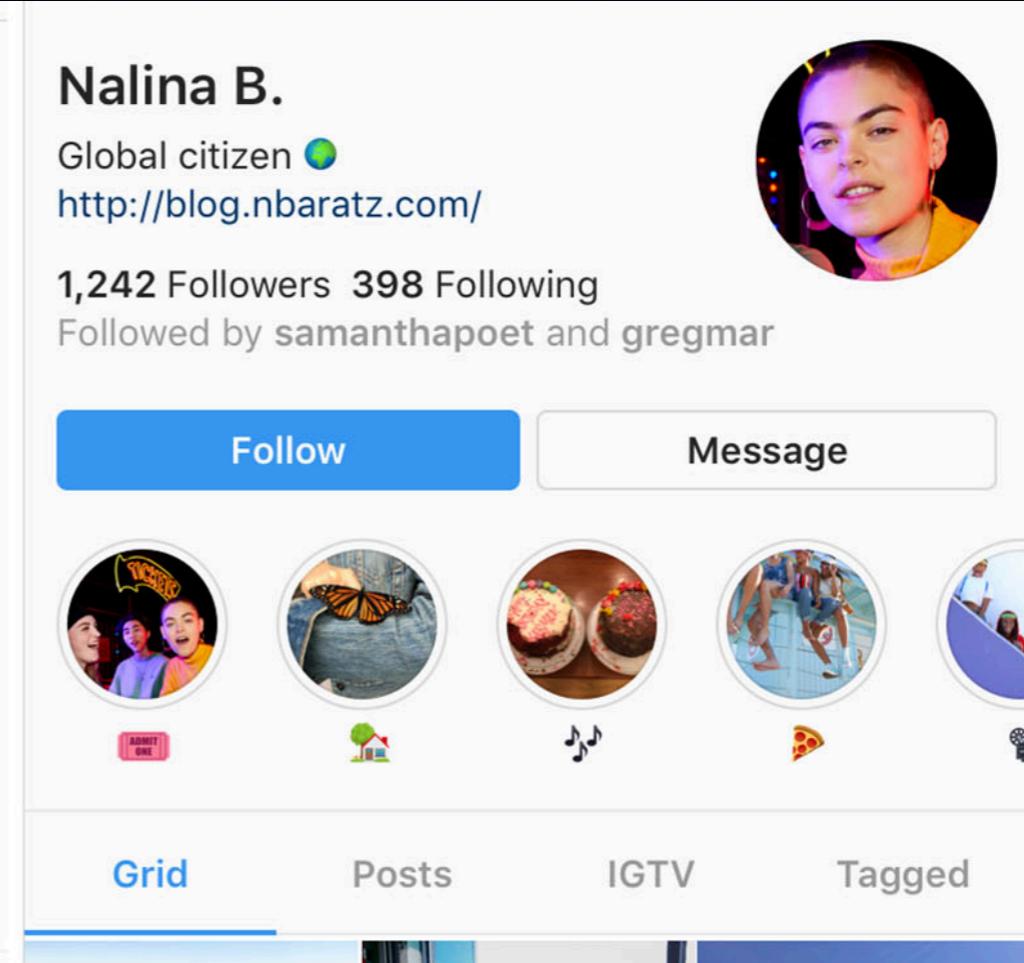
Jeremy Cowart 

Artist
Photographer / Founder of @ThePurposeHotel.
My new book "I'm Possible" is available for pre-orders now!
PossibleBook.com/

Followed by [peter_hurley](#), [srlounge](#), [resourcемаг](#) + 4 more

 **BTS**

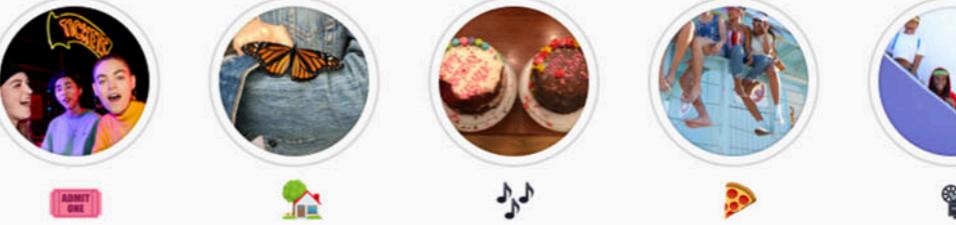
[Message](#) [Email](#)



Nalina B.
Global citizen 
<http://blog.nbaratz.com/>

1,242 Followers 398 Following
Followed by [samanthapoet](#) and [gregmar](#)

Follow  **Message**



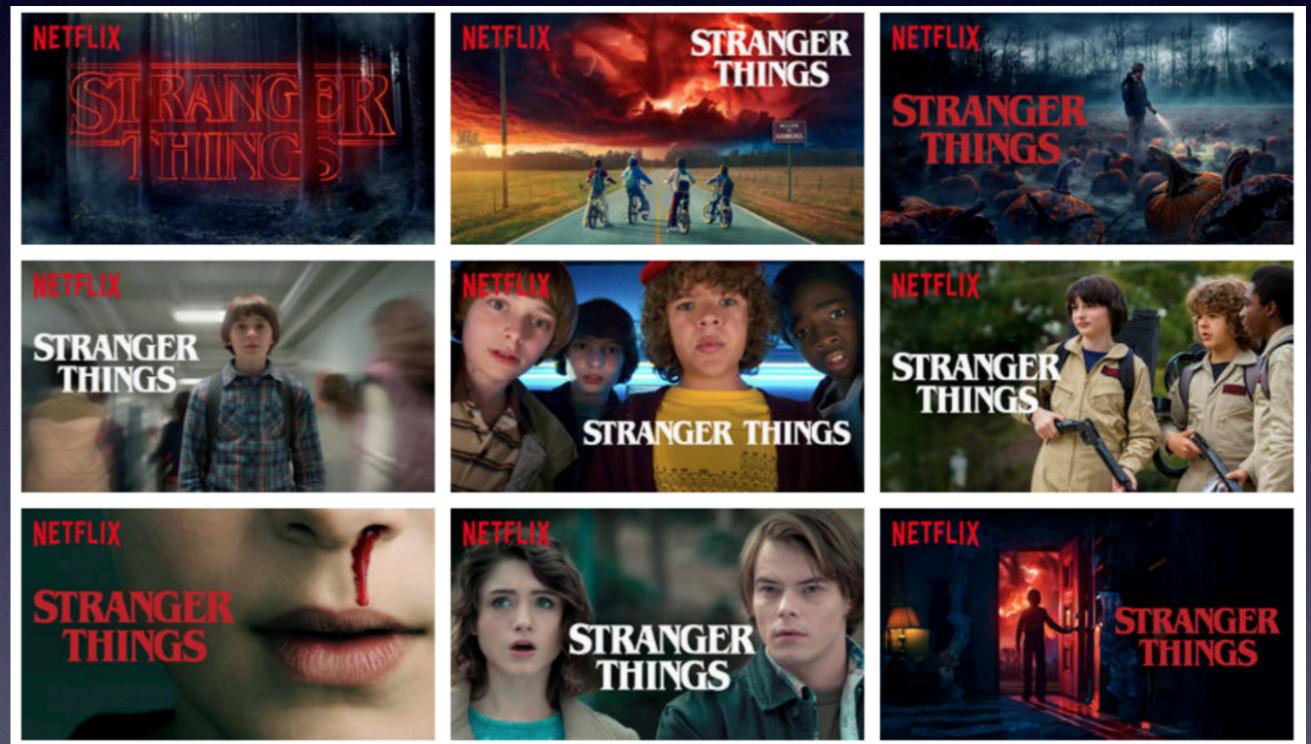
Grid **Posts** **IGTV** **Tagged**

Netflix'in Kişileştirilmiş Özellikleri

Where was I at?

14m

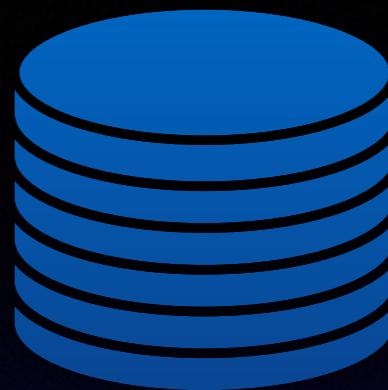
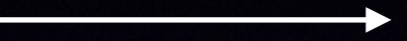
49m



Alışveriş Sitelerinin Bizi Çevremizden iyi Tanıması



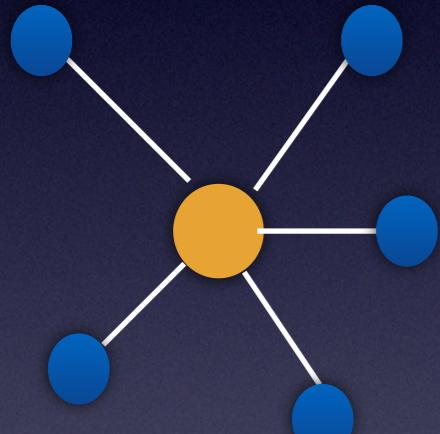
How Target Figured Out A Teen
Girl Was Pregnant Before Her
Father Did



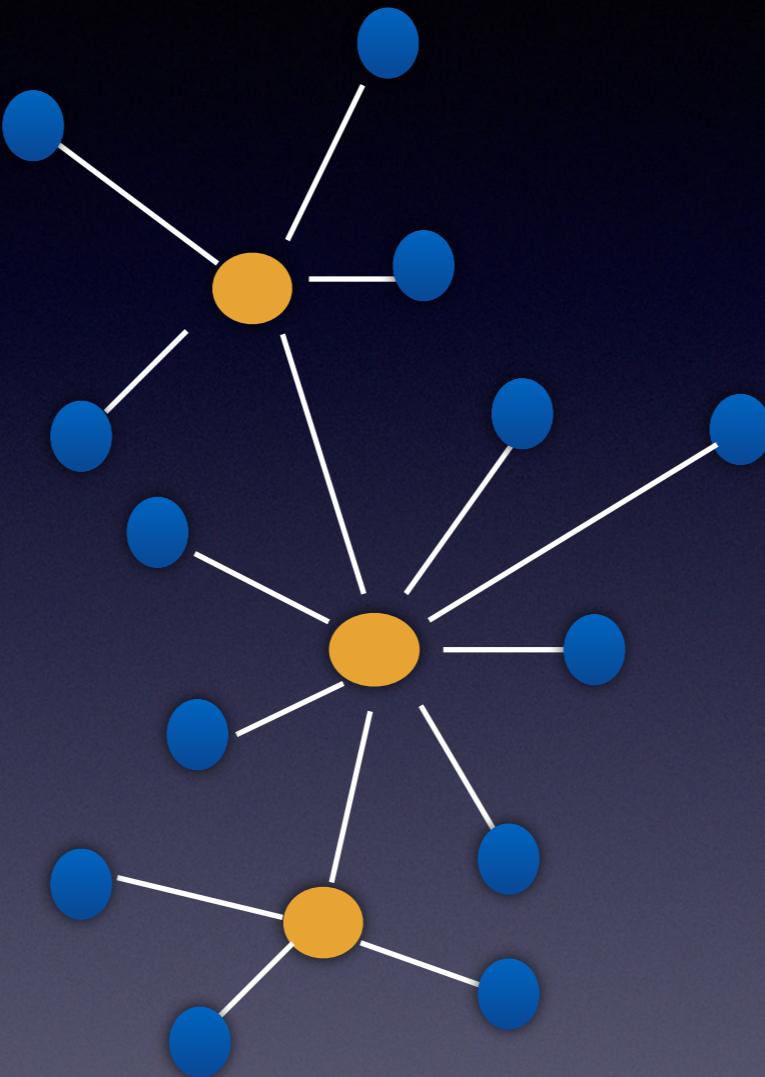




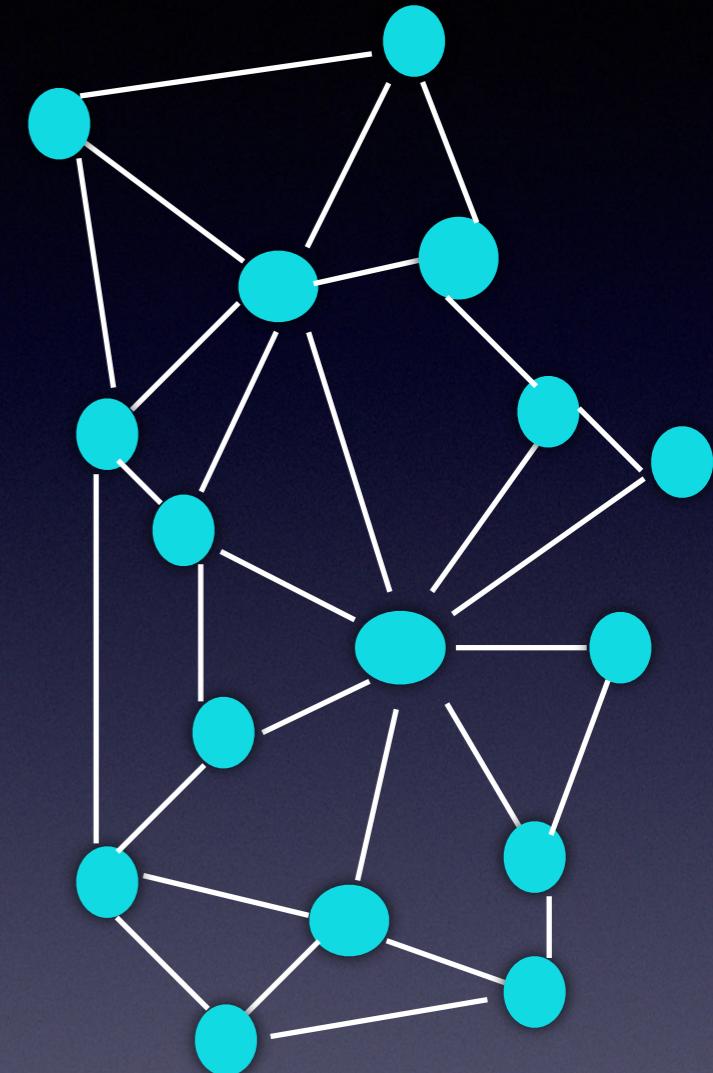
Veriyi Saklama ve Erişim



Centralized

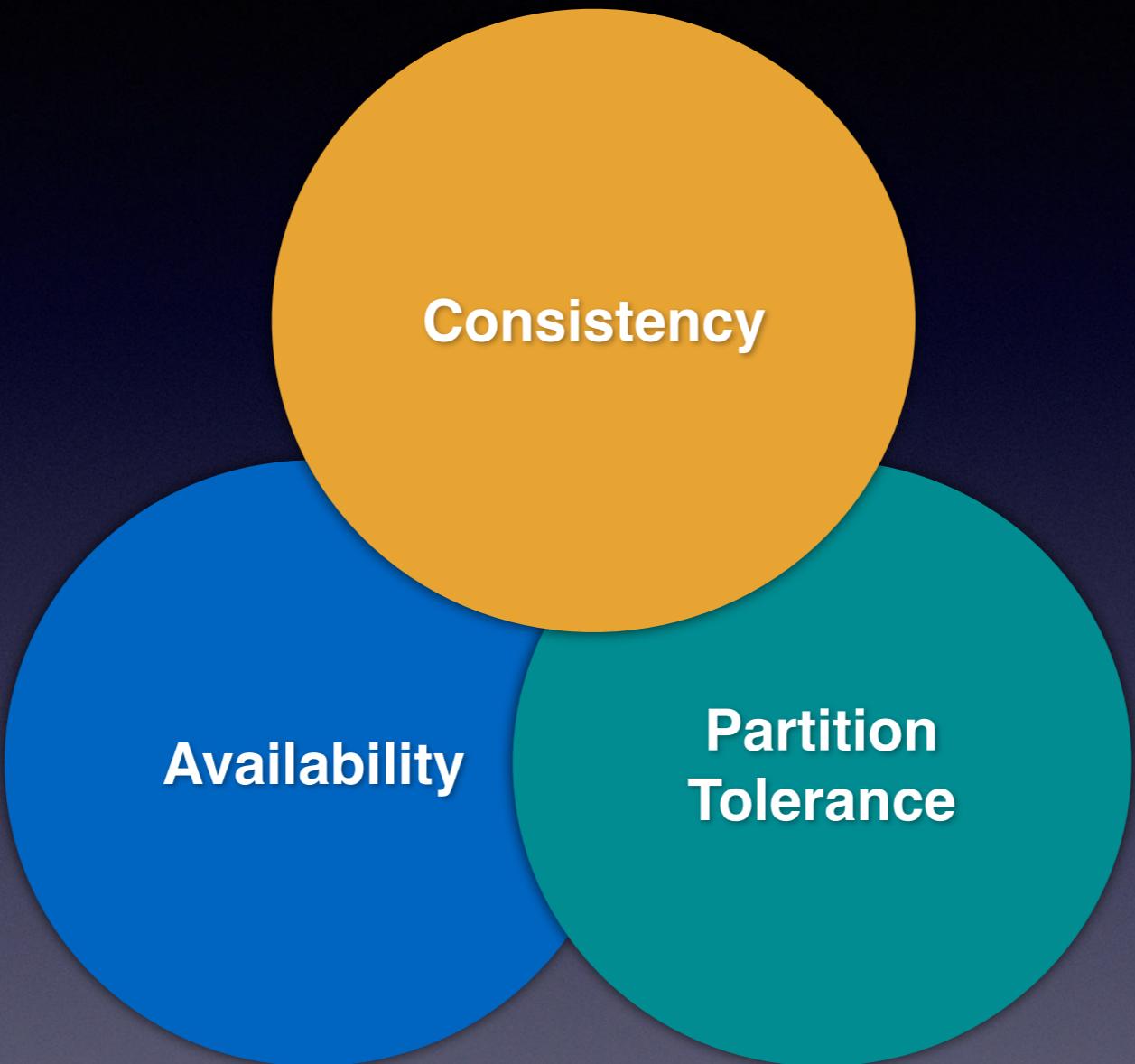


Decentralized

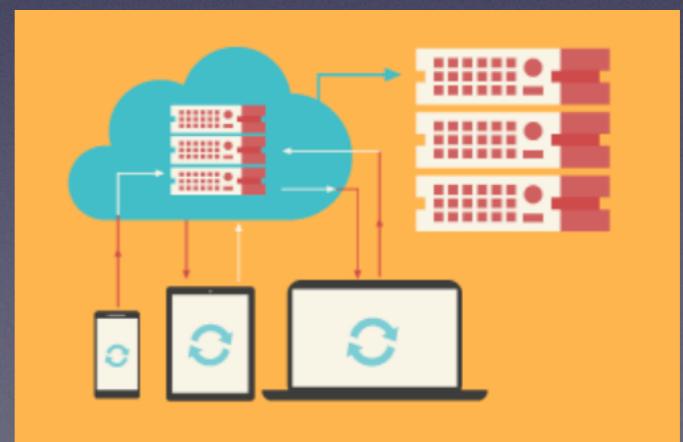
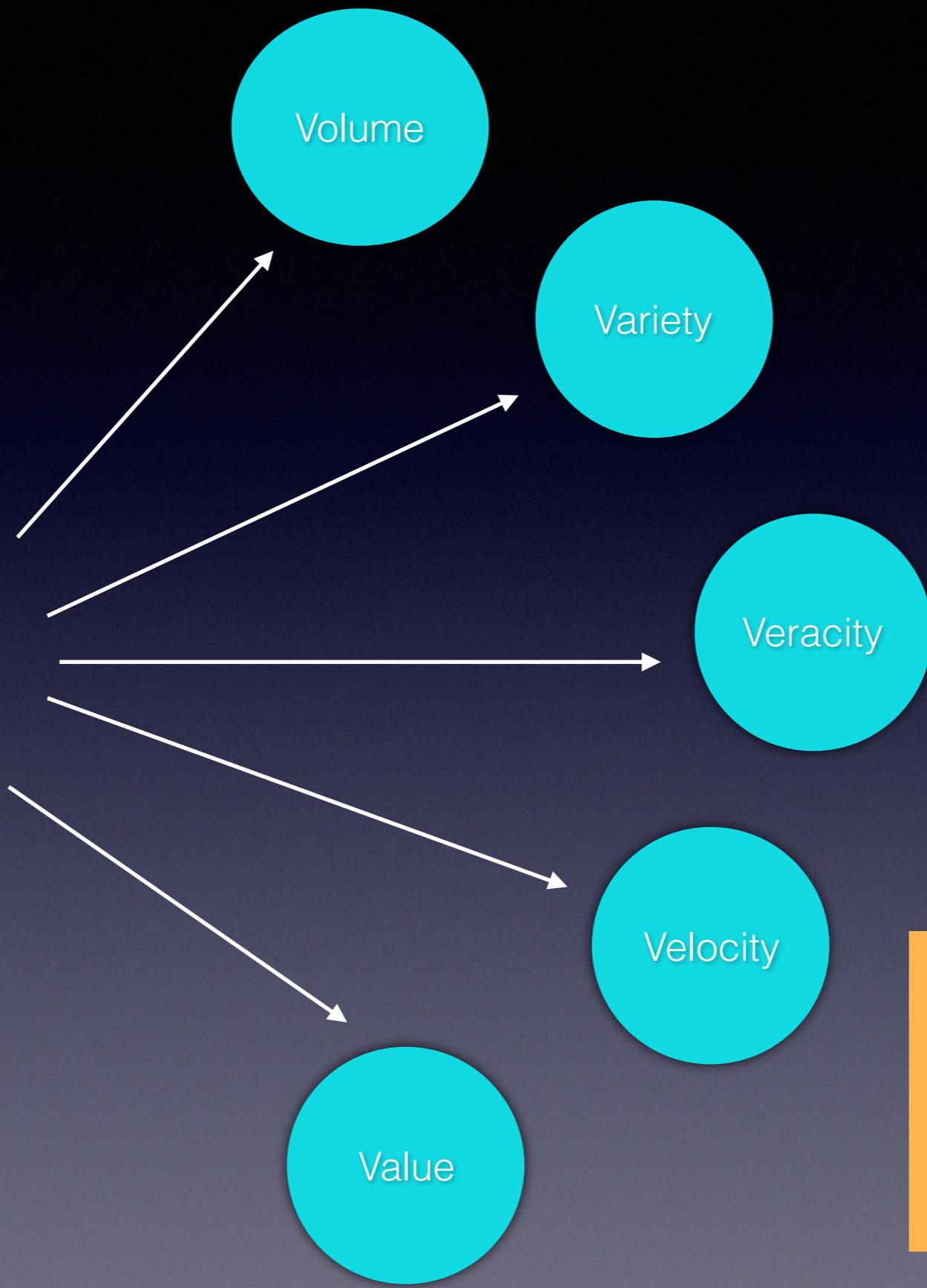


Distributed

CAP Teoremi



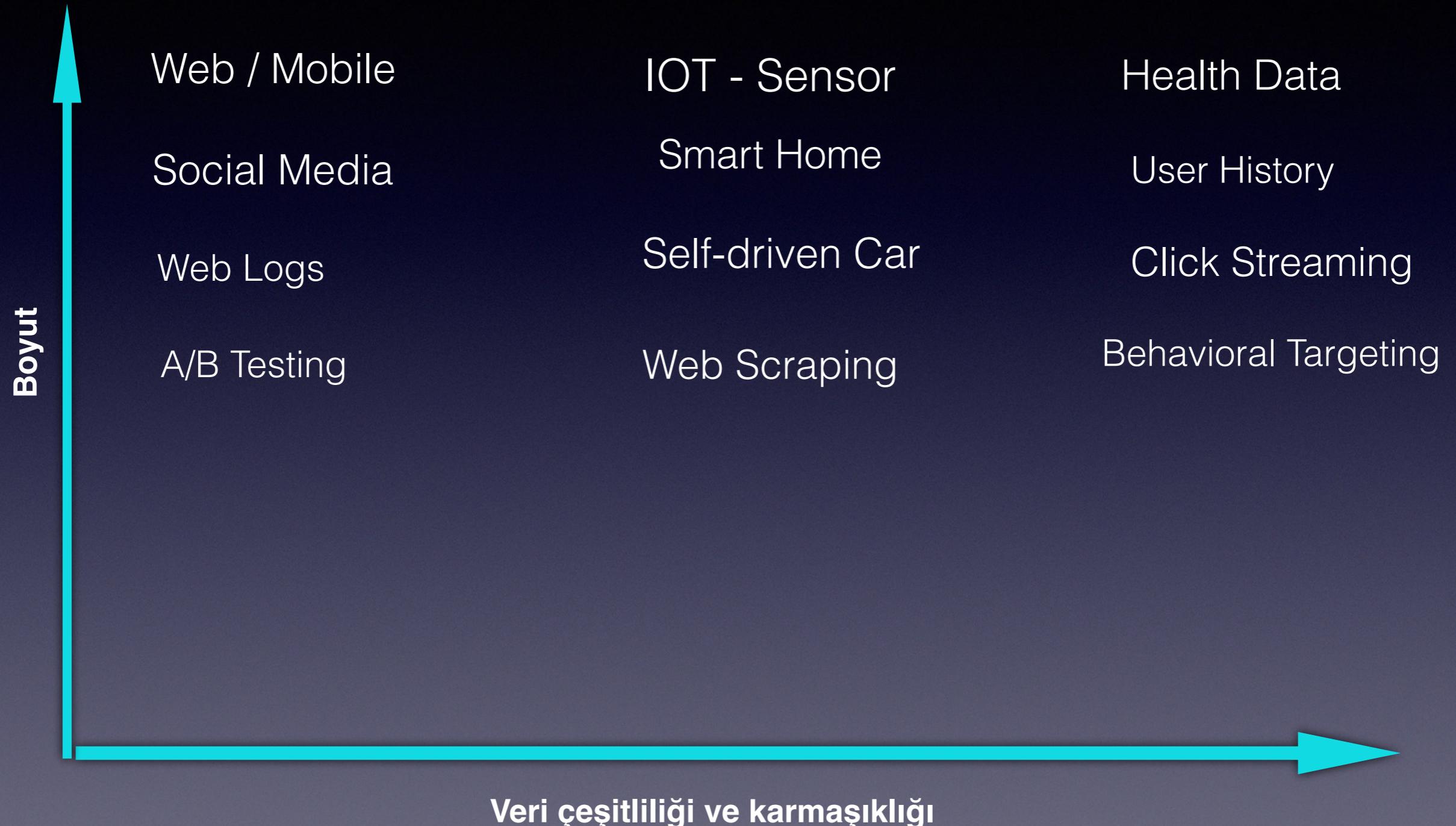
Big Data 5V



Veri Modelleri

- Structured Data
- Semi-Structured Data
- Unstructured Data





Aynı zamanda ...



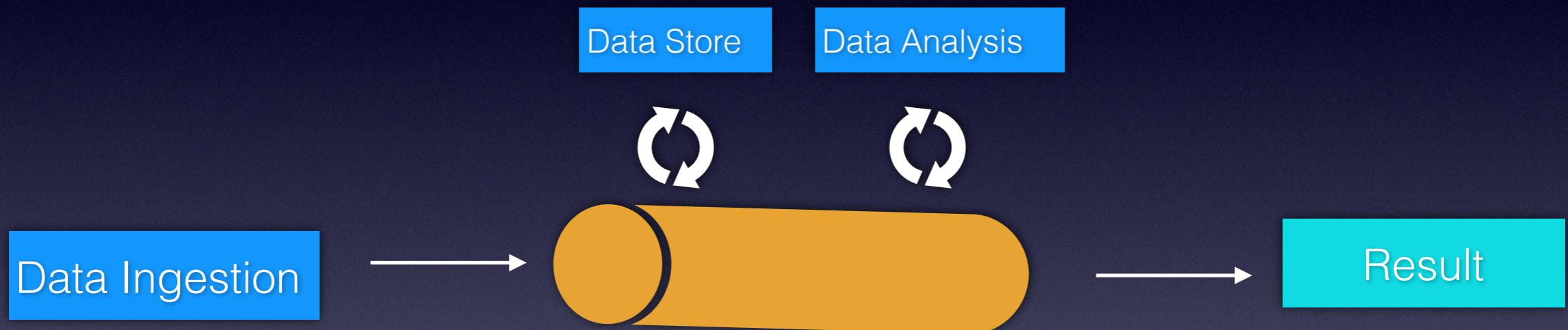
- Saklama maliyeti
- CPU maliyeti



- Network Erişimi

Big Data = İşlem + Etkileşim + Gözlem

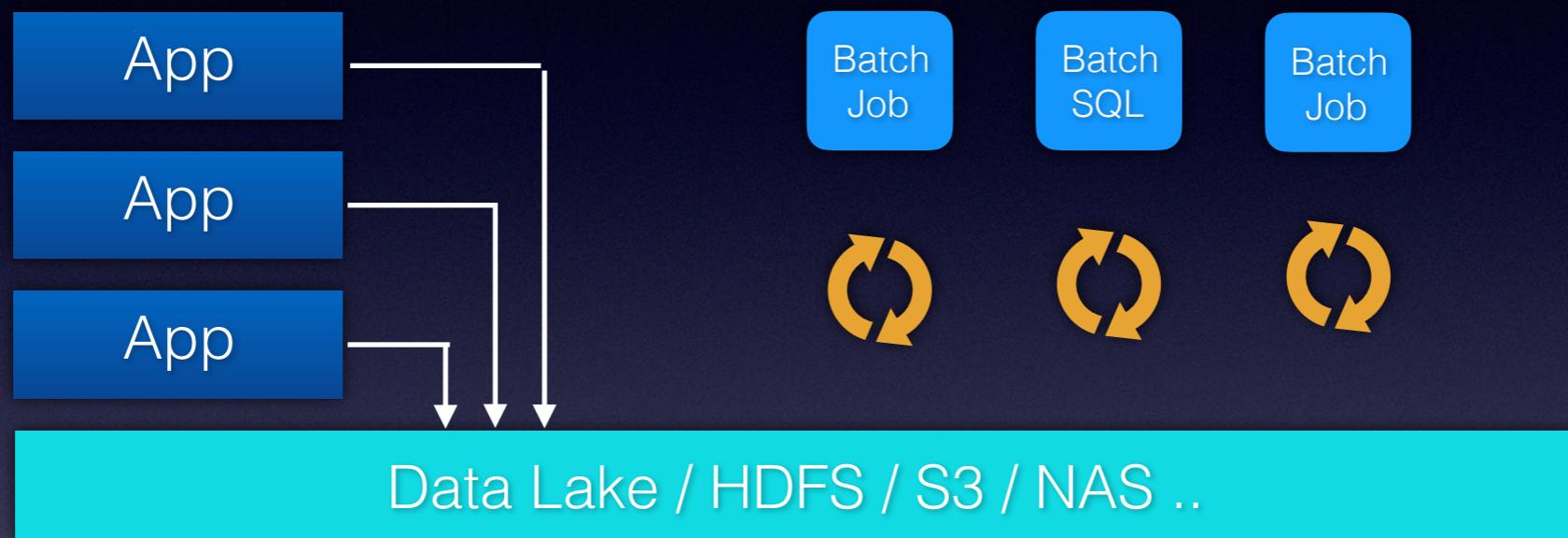
Data Pipeline



VERİ İŞLEME MODELLERİ

- Batch Processing
- Stream Processing / Real-time Processing

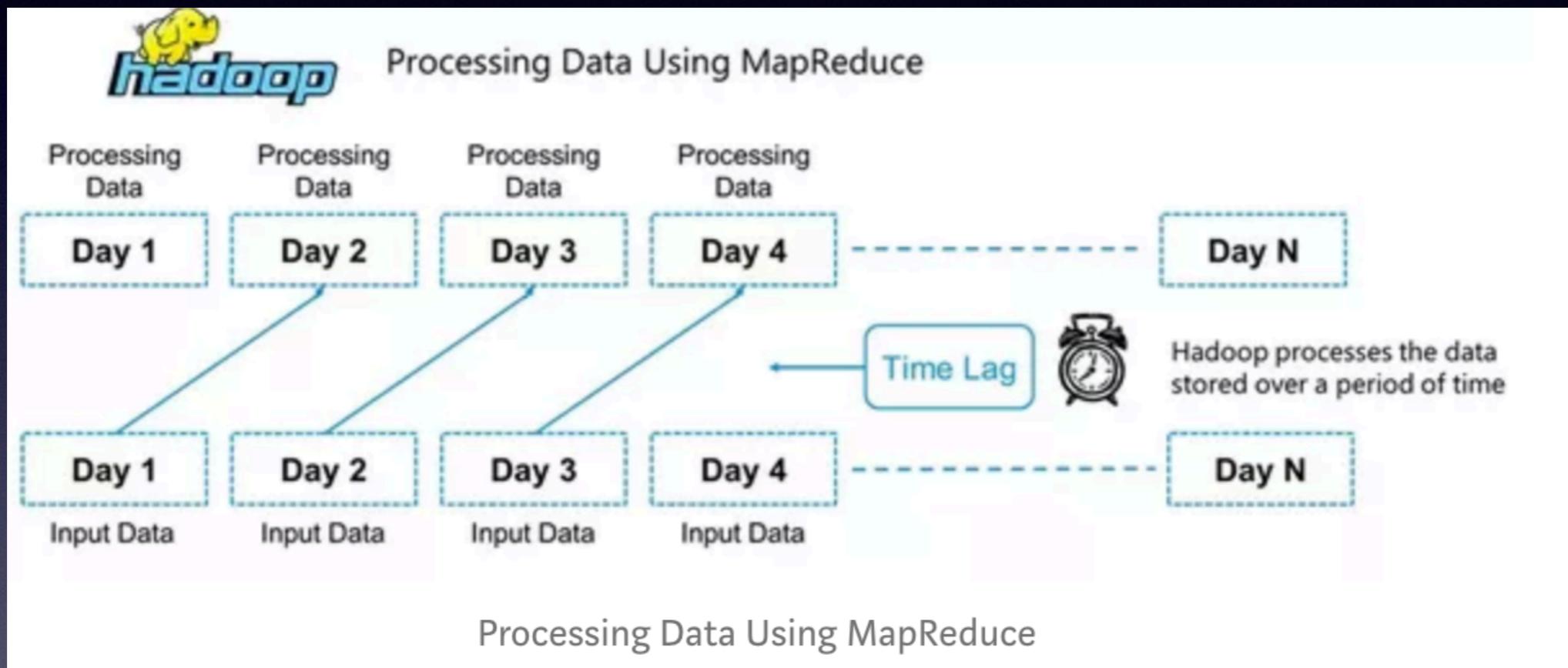
BATCH PROCESSING



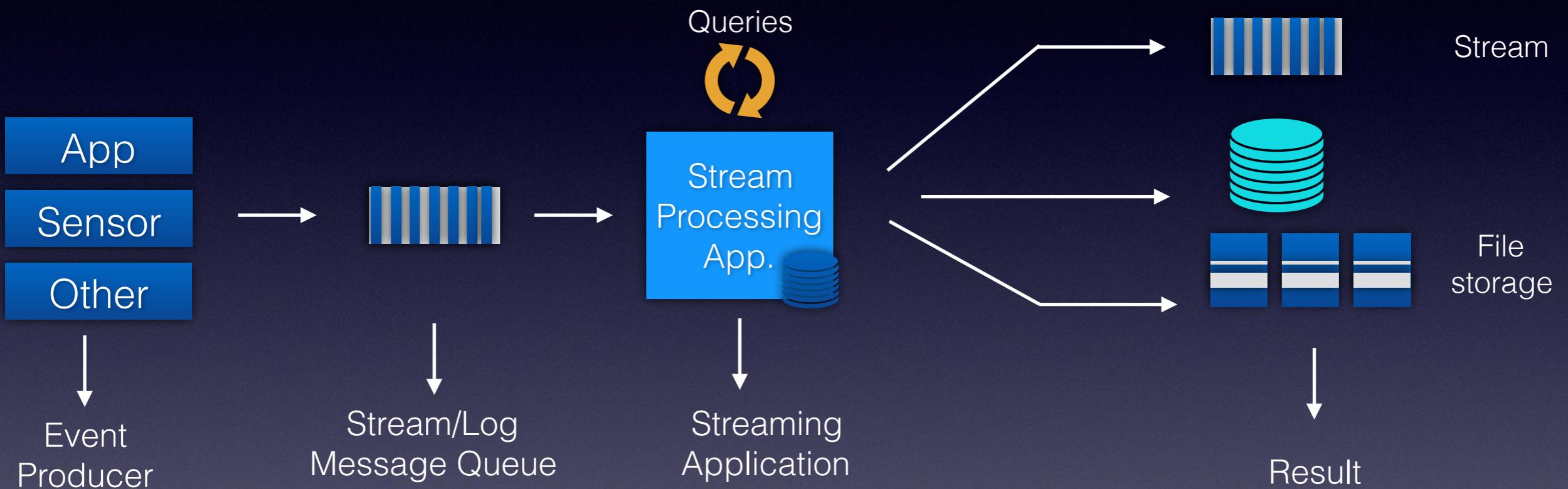
Batch Processing

- Önceden depolanmış veriler üzerinde yapılan analizlerdir
- Veri limitlidir
- Belirtilmiş zaman aralığındaki analizlerde kullanılır.
- Anlık işlemeye ihtiyaç duyulmadığı durumlarda kullanılır.
- Bazı gerçek zamanlı teknolojilerde Micro-batch olarak da adlandırılır.

BATCH PROCESSING MAPREDUCE



STREAMING PROCESSING



STREAMING PROCESSING

- Verinin limiti yoktur.
- Veriyi elde ettiği anda çok kısa bir sürede kurgulanan şekilde işler.
- Hızlıdır.
- Paraleldir.



NEDEN GEREKLİ ?

- Bazı durumlarda verinin sonsuz olması
- Anlık takibini yapmak
- Birden fazla akış üzerinde veriyi eşzamanlı erişim
- Toplu şekilde işlemenin yüksek gecikmeyle sonucu vermesi ve depolama gereksinimi
- Veri depolanamayacak kadar büyük olabilir.

KULLANILAN DURUMLAR

- Sensör Verileri
- Web Site Aktiviteleri
- Finansal Veriler
- Zaman Serisi Verileri

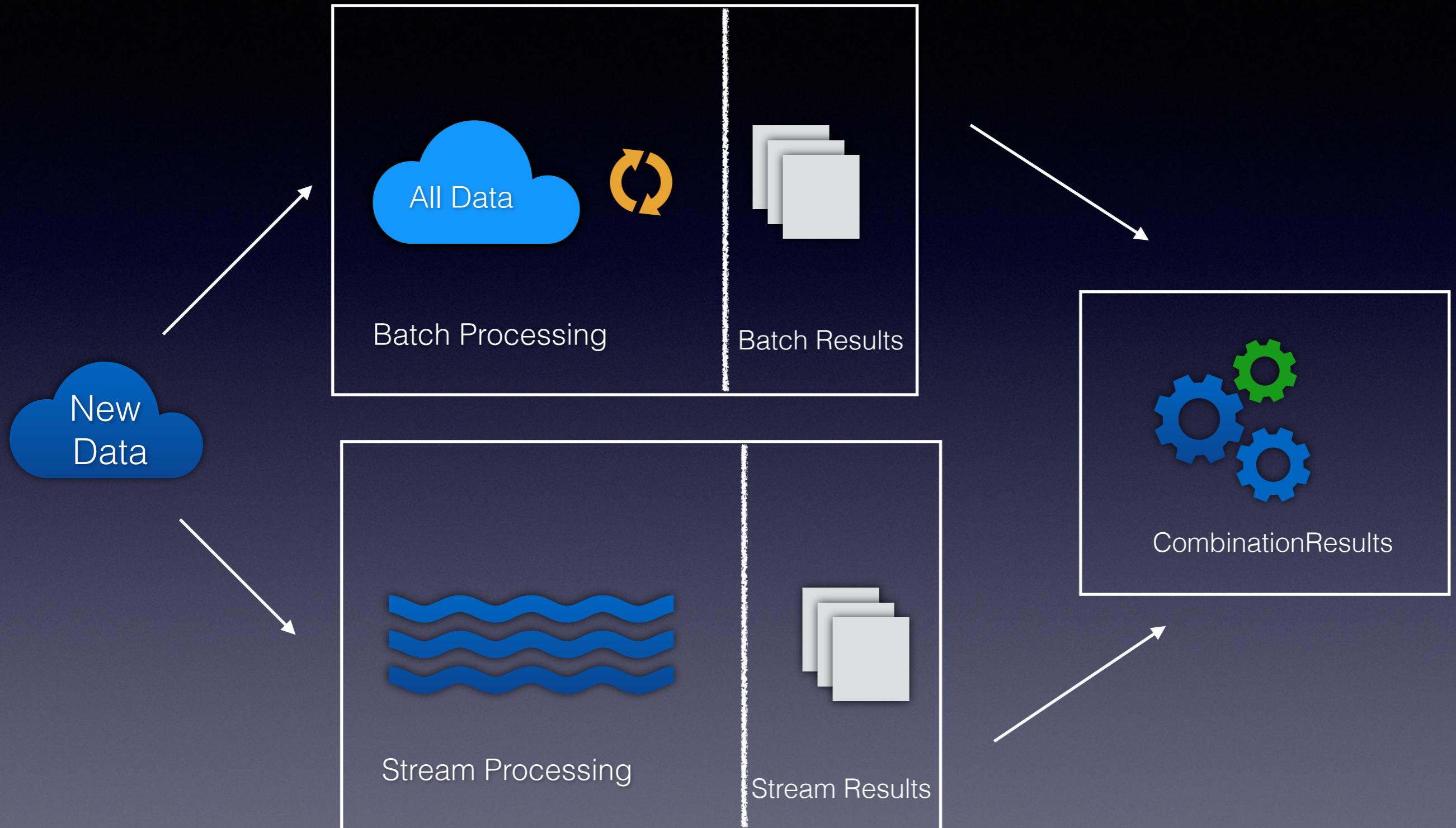
Sonuçlarında Neler Oluyor ?

- Uyarı Sistemleri
- Fraud Detection
- Akıllı Ev, Tarım sistemleri
- Sürücüsüz araçlar
- Makine öğrenmesine eğitim verisi

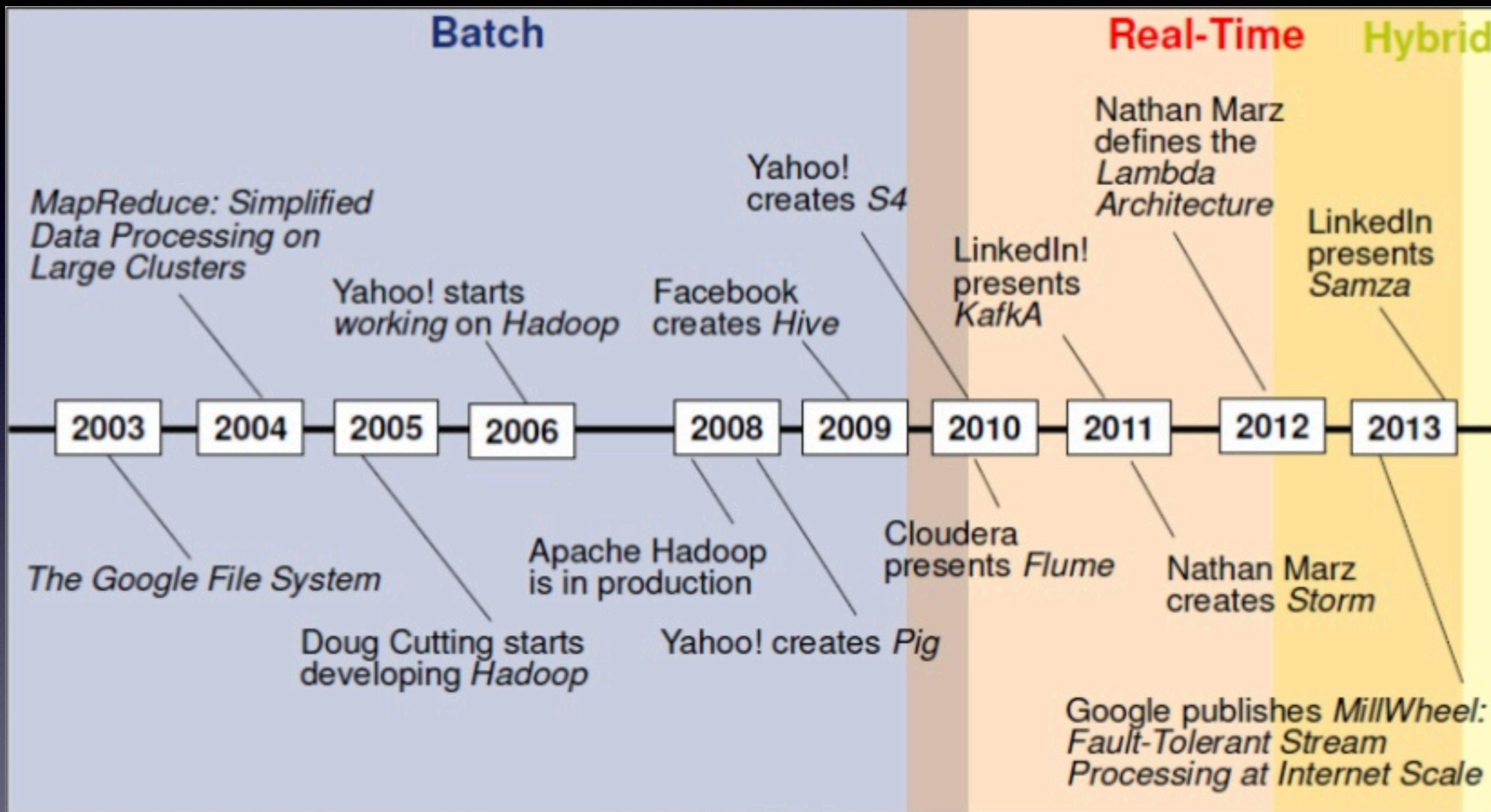
Veriyi İşleme Tarzı - Stream Processing

En fazla bir defa (At Most Once)	En az bir defa (At Least Once)	Sadece bir defa (Exactly Once)
Mesaj bir defa gönderilir	Mesaj bir veya daha fazla gönderilir, her defa işlenir	Mesaj bir veya daha fazla gönderilir ama bir defa işlenir
Mesaj karşılanabilir.	Mesajın tüketileceği garantidir.	Mesajın tüketileceği garantidir.
Veri çöktürülmez	Veri çöktürülür	Veri çöktürülmez
Veri kaybı olabilir.	Veri kaybı olmaz	Veri kaybı olmaz

Hibrit Model



[Source](#)



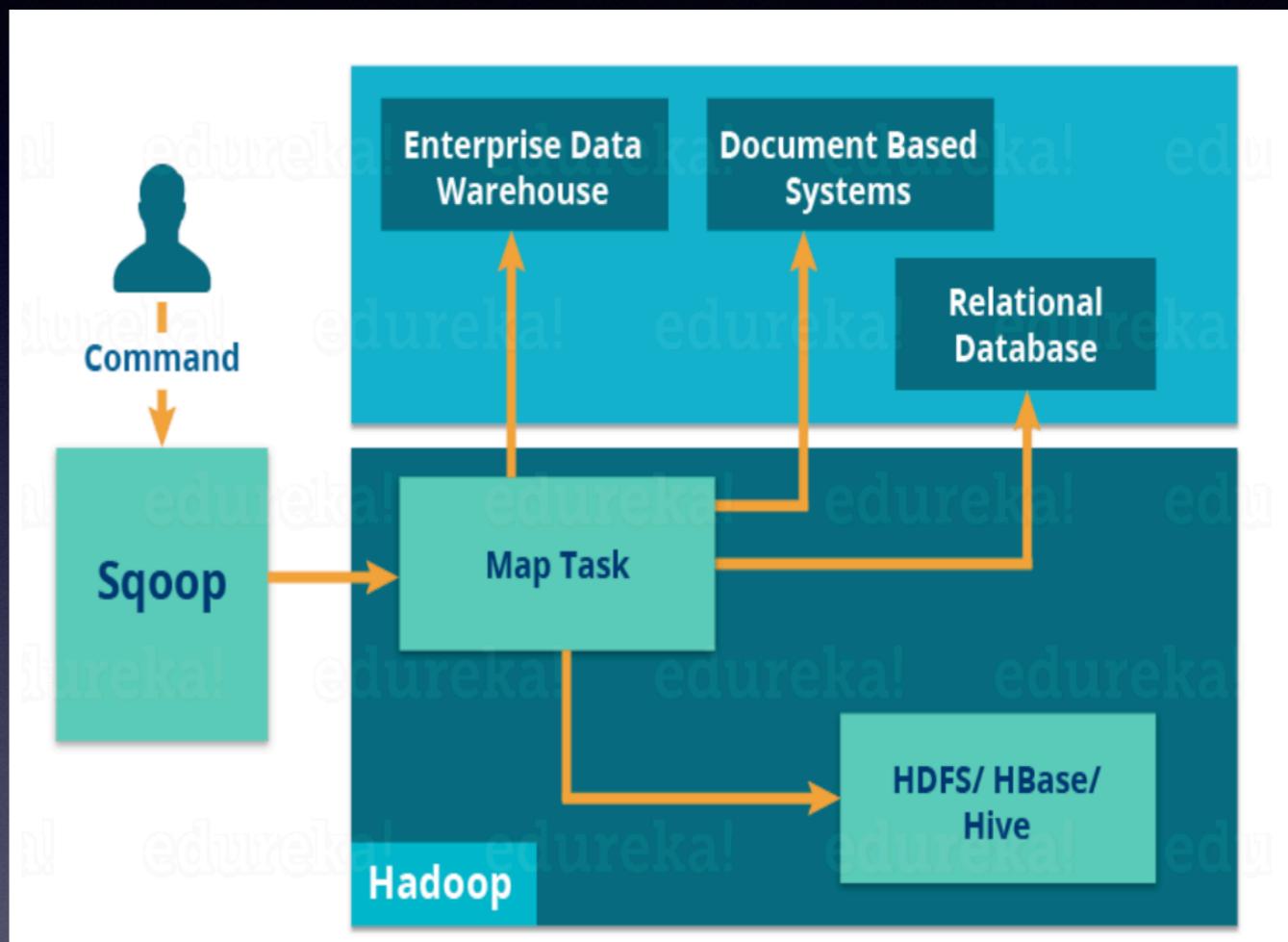
Büyük Veri Teknolojileri

Veri Toplama Teknolojileri



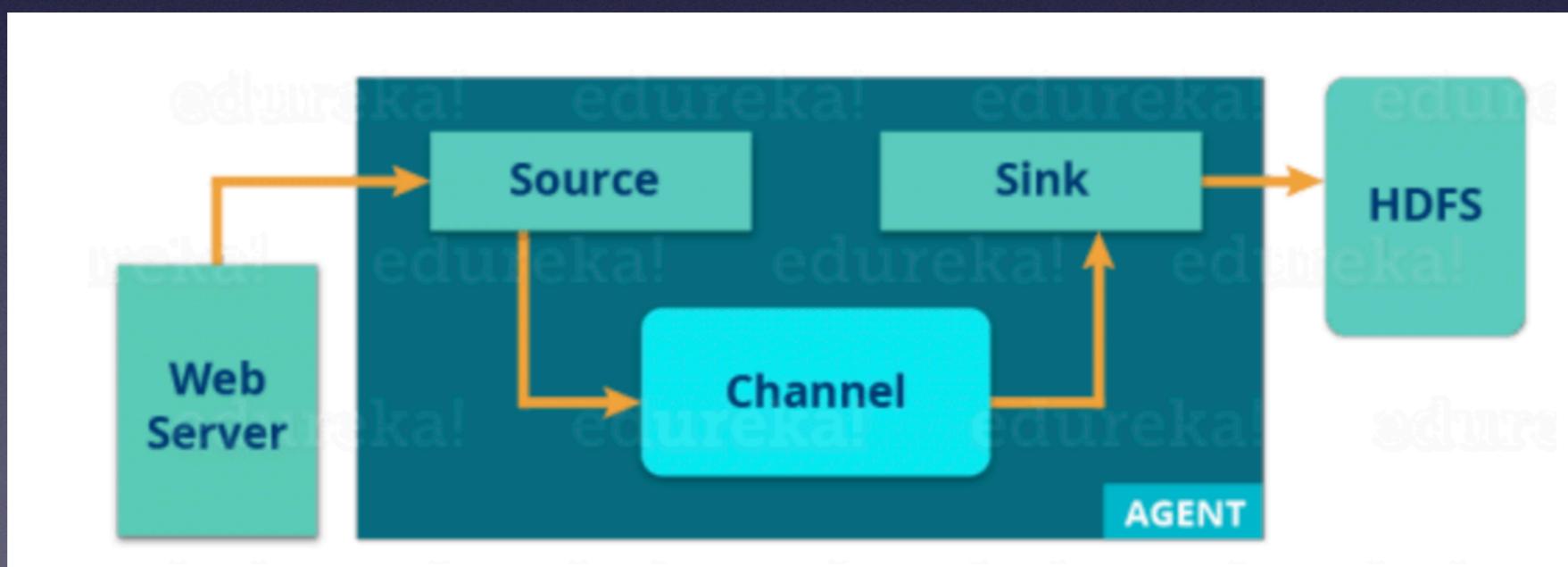
Apache Sqoop

- İlişkisel veritabanları ve Hadoop arasında veri aktarımı yapar
- Mapreduce Job'ları oluşturur.
- Yarn kullanıldığından paralel olarak import ve export işlemi yapabilir.
- Komut satırı ile istenilen boyuttaki veriyi istenilen yere aktarır.



Apache Flume

- Birden fazla kaynaktan veriyi toplayıp HDFS'e yazar
- Dağıtıktır.
- Channel bazlı transaction api'yi vardır. (Sender & receiver)



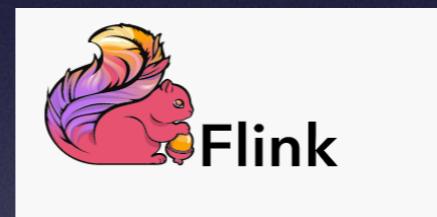
Apache Nifi

- Farklı kaynaklardan veri toplanmasına, işlenmesine ve farklı ortamlara aktarılmasını sağlar
- Web arayüzü sunar.
- Dağıtıktır.
- Kafka, HDFS, NoSQL gibi ortamlara veri akışını yönlendiren hazır işlemciler (processor) vardır.
- Kendi custom processor'lerimizi yazabiliriz.
- Kendi içersinde kuyruk yapısı vardır.

Fluentd

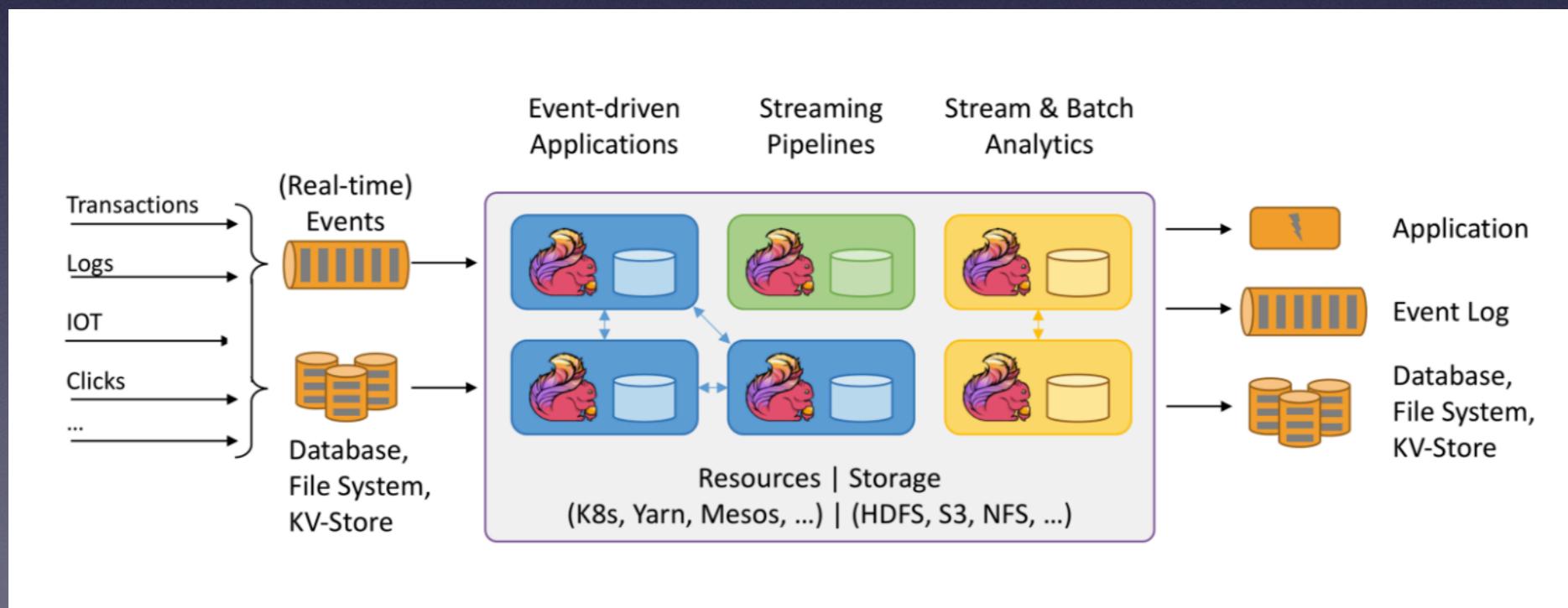
- Log yapılarını birleştirmek için tasarlanmıştır
- Mümkin olduğunda JSON veri formatında çalışır
- Çoklu kaynaklara rağmen log toplamayı, filtrelemeyi ve çıktı vermeyi kolaylaştırır.
- Kaynak tüketimi ve kurulumu oldukça kolaydır.

Realtime / Near-Realtime Streaming Teknolojileri



Apache Flink

- Hadoop ile uyumlu veri akış platformudur.
- Batch ve stream programlar çalıştırabilir.
- FlinkML kütüphanesi ile makine öğrenmesi ortamı sunar.
- Stateful Streaming ile “exactly once” garantisı verir.



Source

Apache Beam

- Hem batch hem de streaming processing API'sıdır.
- Platform bağımsızdır.
 - Oluşturulan data pipeline Apache flink, spark, Google dataflow gibi ortamlarda çalıştırılabilir.

Saklama Teknolojileri

NoSQL



KeyValue



etcd



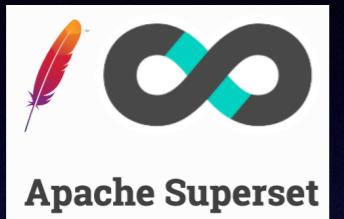
GraphQL



Time Series



Görselleştirme Teknolojileri



Apache Superset



matplotlib



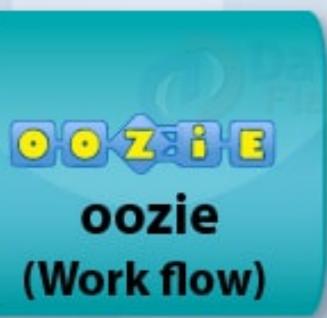
Hadoop Ekosistemi

Hadoop Nedir ?

- Framework
- Büyük veri üzerinde çalışan uygulamalara ortam sunar.
- Açık kaynak kodludur.
- Dağıtıktır.
- Ölçeklenebilir yapıdadır.

- İki ana katmanı vardır
 1. Dağıtık Dosya Sistemi (HDFS)
 2. Dağıtık veri işleme - MapReduce

Hadoop Ekosistemi



oozie
(Work flow)



HCatalog
Table & schema Management



Pig (Scripting) **Hive (Sql Query)**



mahout
(Machine Learning) **Drill**
(Interactive Analysis)



AVRO
(JSON) **Thrift**
(Cross Language Service)



APACHE HBASE
HBASE
(Columnar Store)



Sqoop
(Data Collection)



Zookeeper
(Coordination)



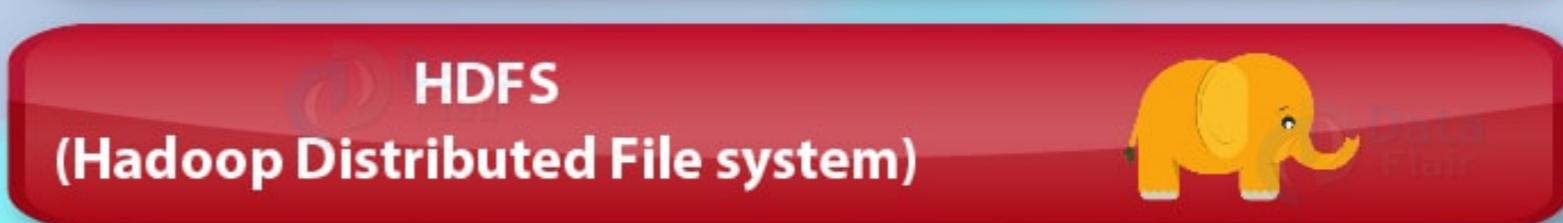
FLUME
Flume
(Data Collection)



Mapreduce
(Data Processing)

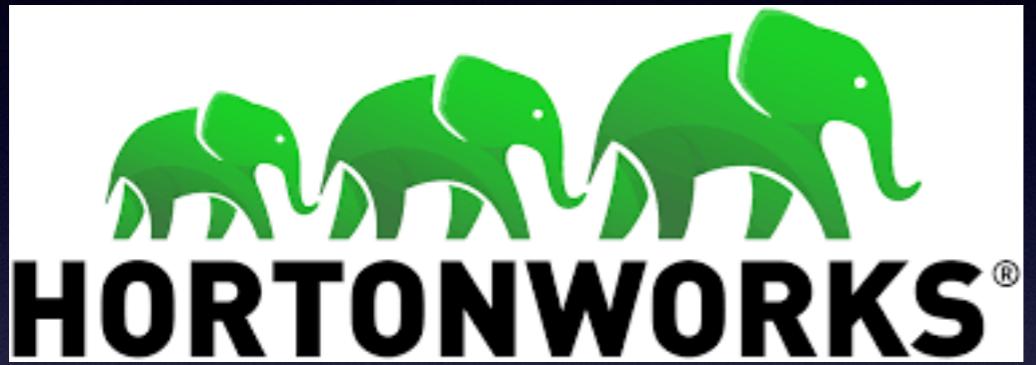


Yarn
(Cluster Resource Management)



HDFS
(Hadoop Distributed File system)

HADOOP YONETİM PLATROMLARI

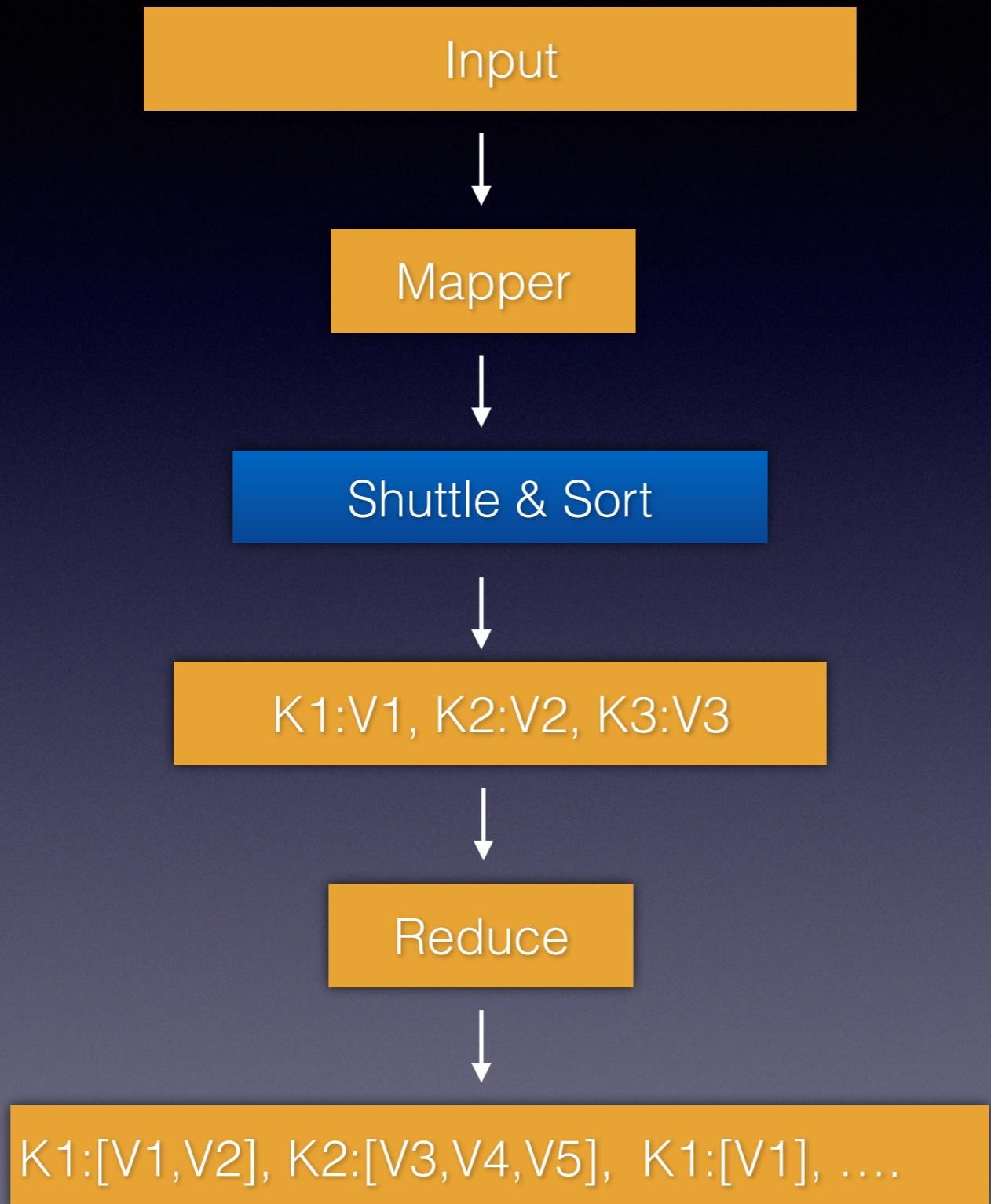


	Apache Hadoop	Hortonworks	Cloudera	MapR
Open Source	Yes	Yes	Partially	No
Support	Community	Enterprise Support	Enterprise Support	Enterprise Support
Frontend	Apache Ambari	Apache Ambari	Cloudera Manager	MapR Control System
Price	Free	\$\$	\$\$	\$\$\$
Focus	Open Source, reliable, scalable, distributed computing	Enterprise capabilities	Enterprise capabilities	Enterprise & Performance

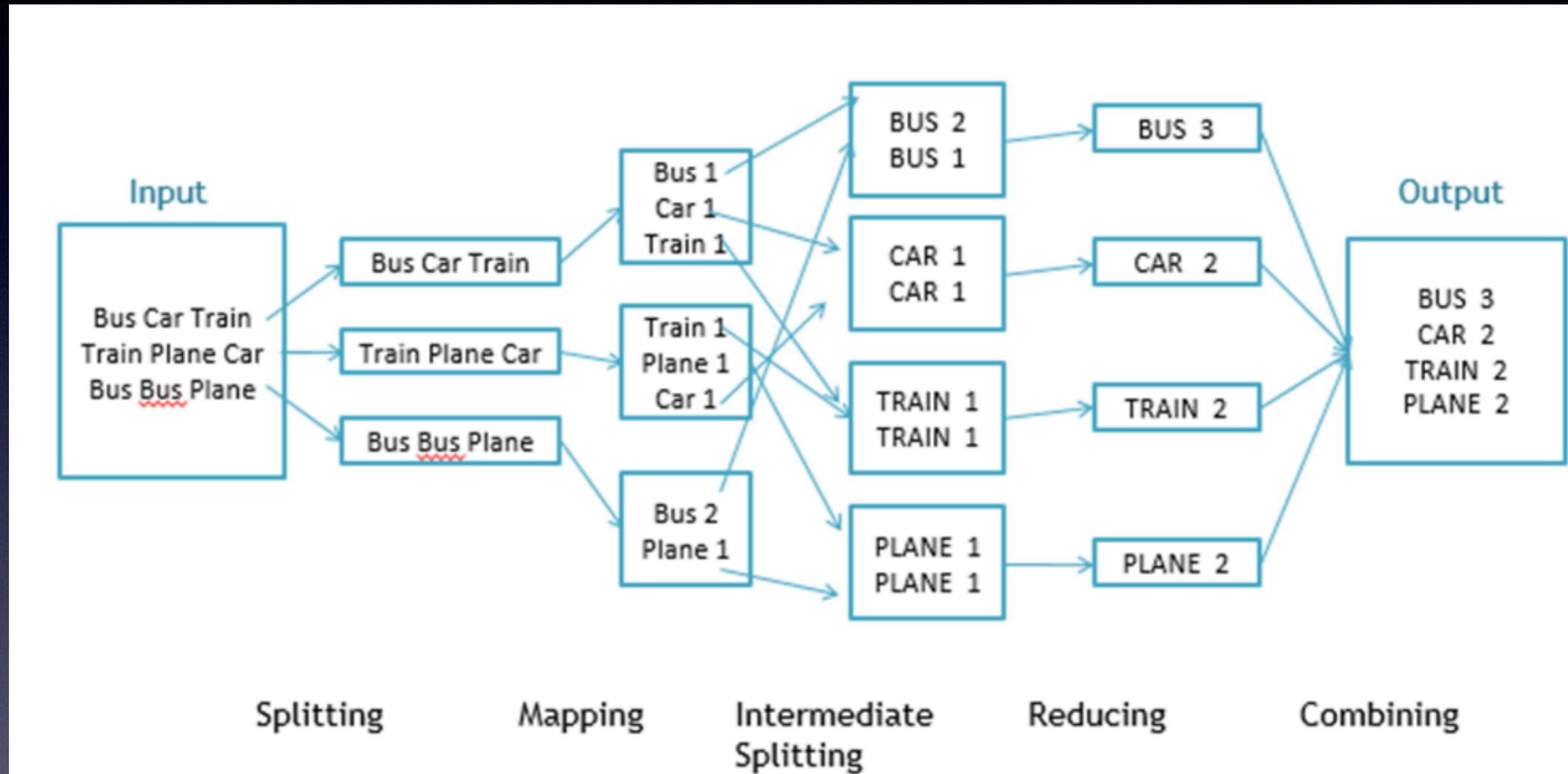
Mapreduce

- **Dağıtık** veri analiz framework'ü
- **Hadoop V1 ile hem kaynak yönetimini hem de veri analizi**
- Birçok dil ile geliştirme yapılır.
- **Pig** ile daha kolay geliştirme sunar.
- **Batch** işlemler gerçekleştirir.
- **Yüksek gecikme**

Mapreduce

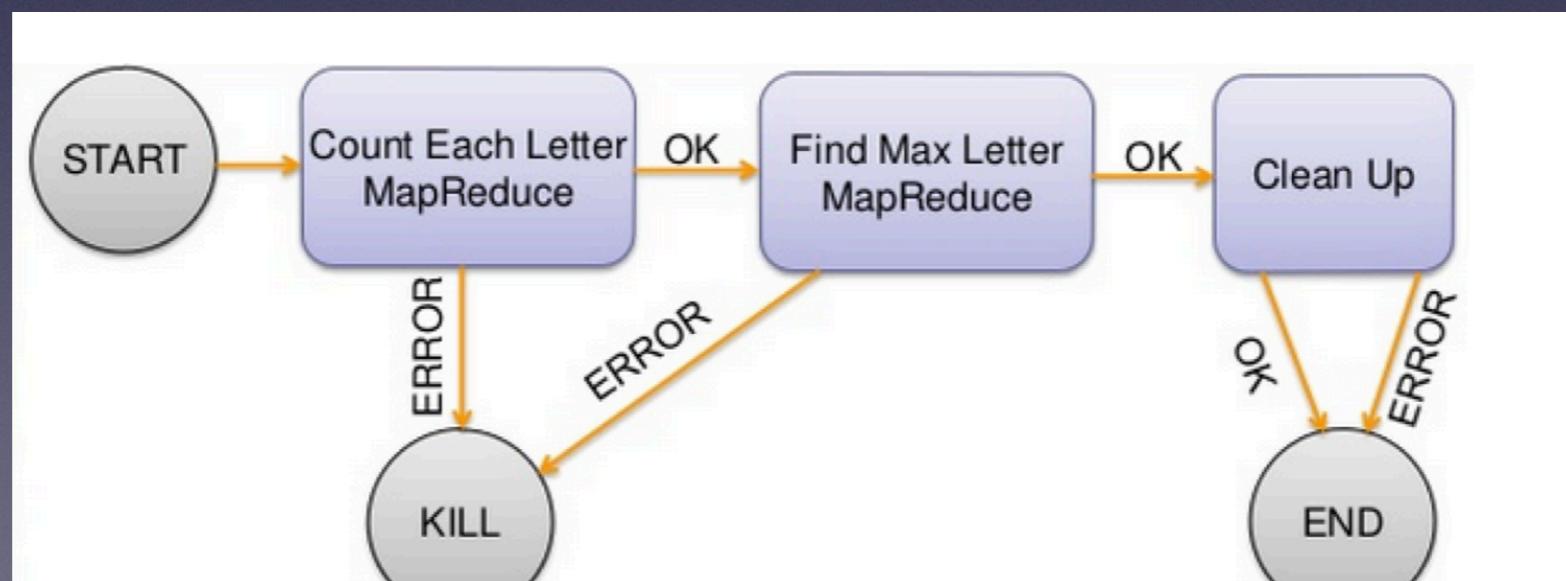


Mapreduce Word Count



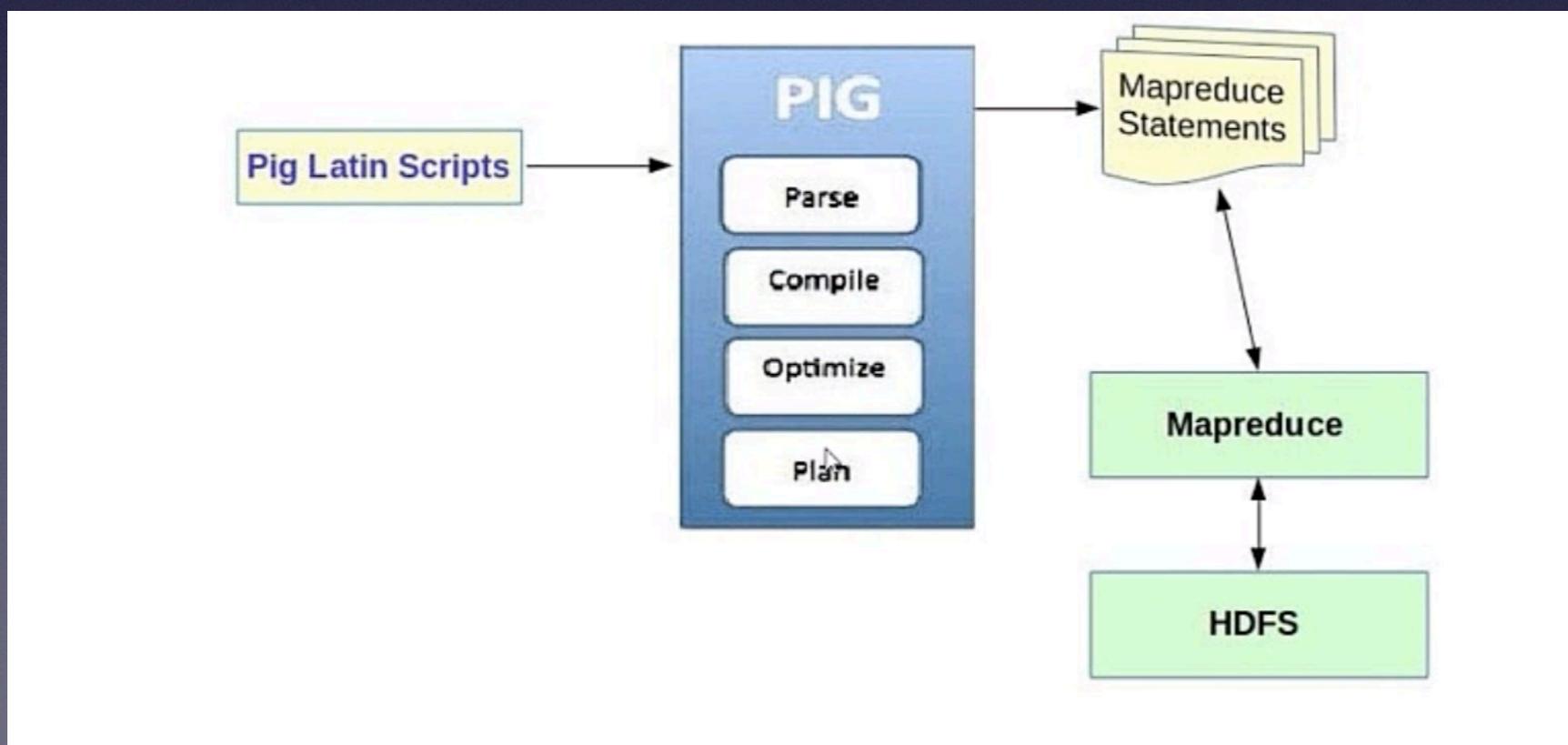
Oozie

- Workflow yönetimi & bunların koordinasyonu
- Akış, “action node” lardan oluşur.
 - Start, End, Kill
- Akışları belli zamanlarda otomatik çalıştırır.
- Pig script, bash script, Hive Sorguları



Pig

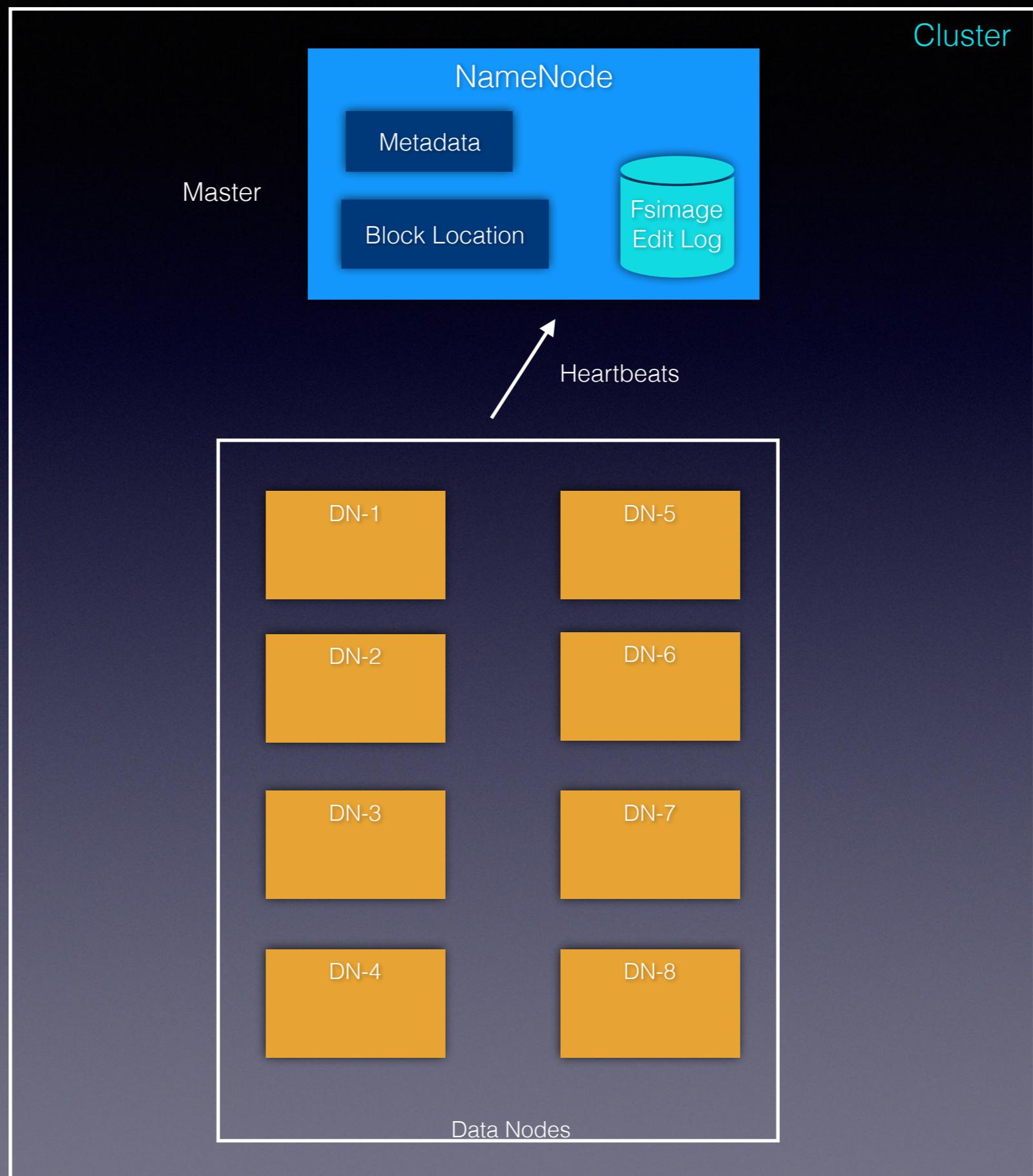
- Mapreduce ile veri analizi yapmamızı sağlar.
- Optimizasyon mekanizması vardır.
- Kendine özgü Pig Latin dili ile geliştirme yapılır.



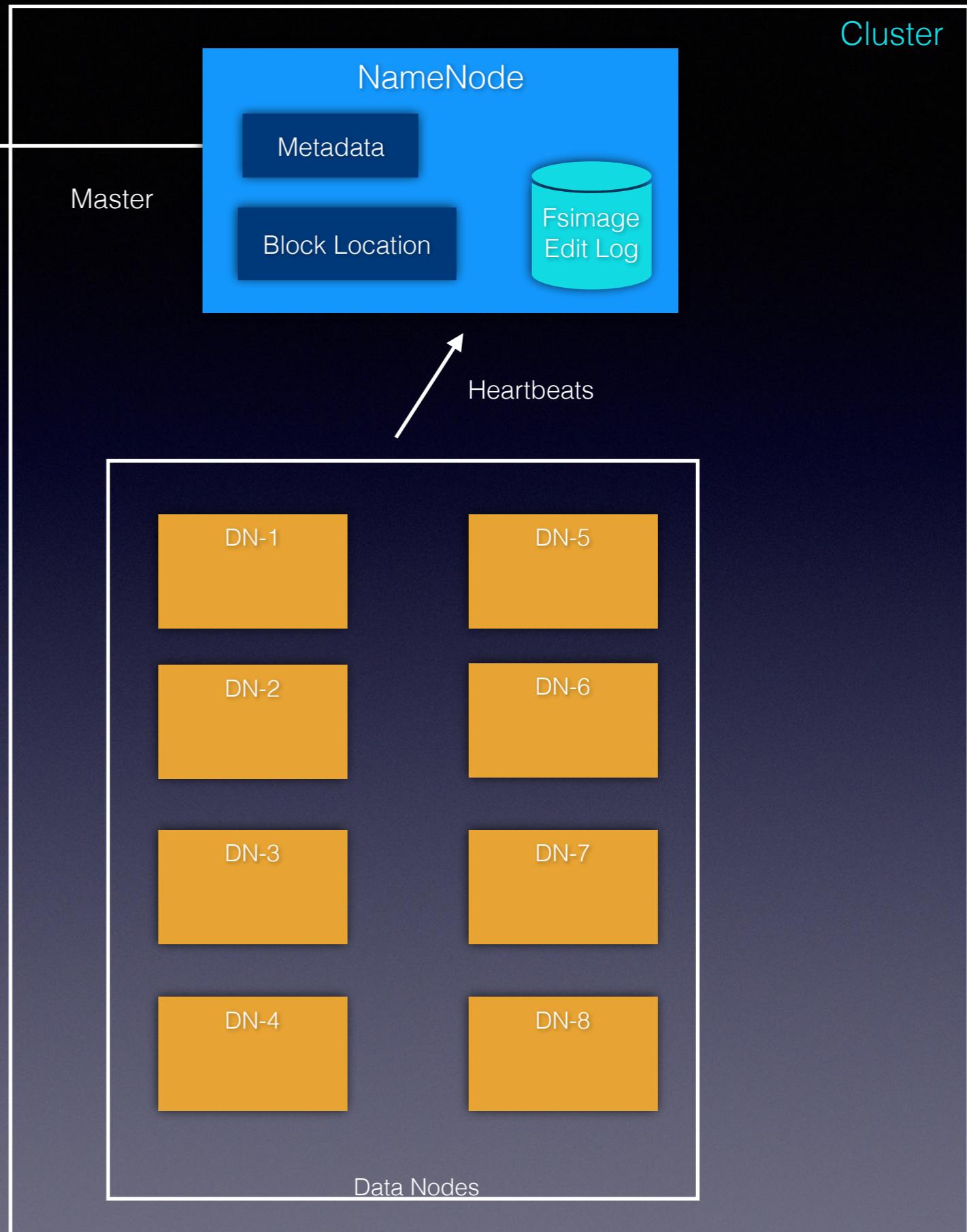
HDFS - Hadoop Distributed File System

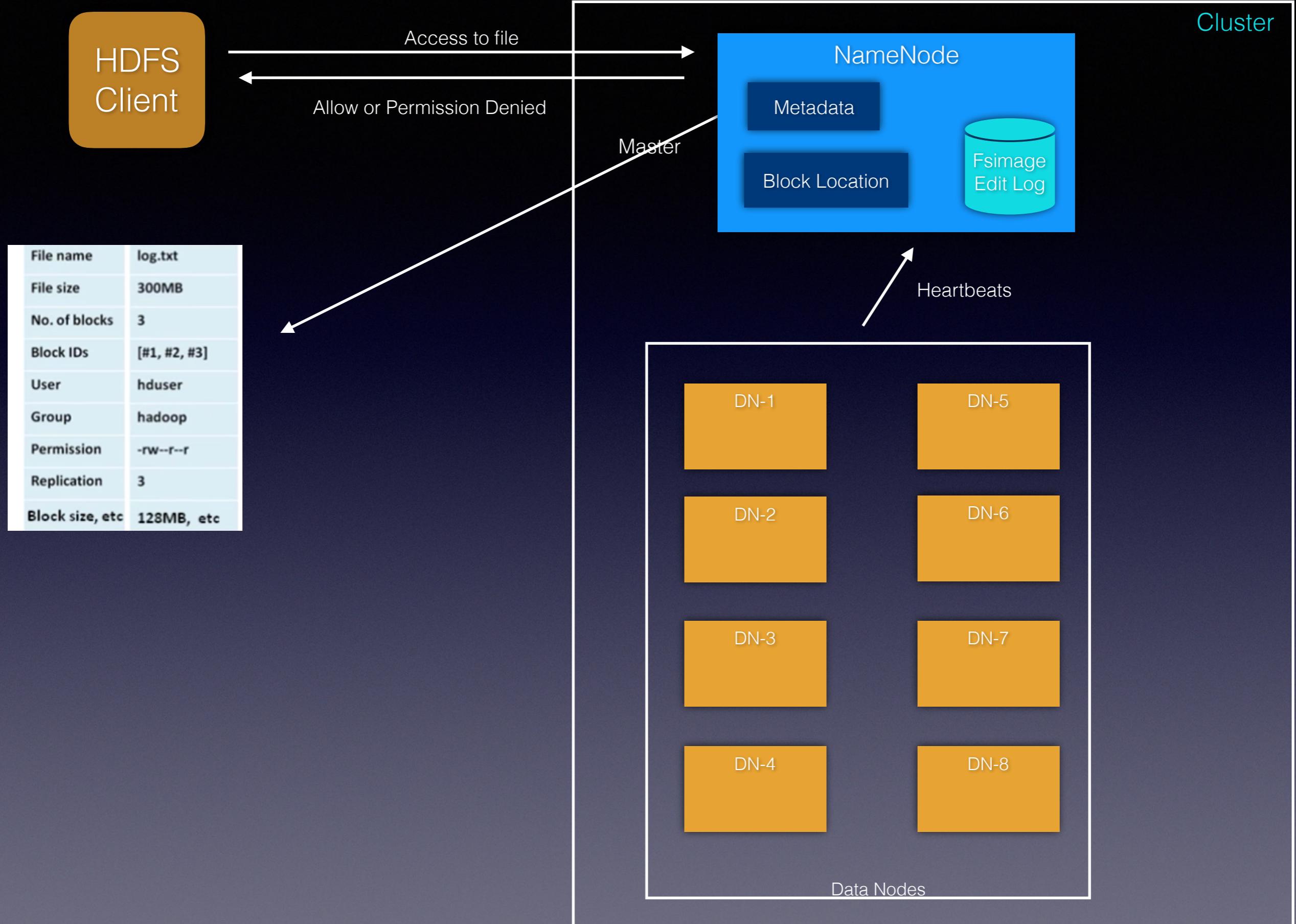
- Veriyi **bloklar** halinde **dağıtık** olarak saklar.
- Master - Slave mimarisi ile çalışır.
 - Data Nodes
 - Name Nodes

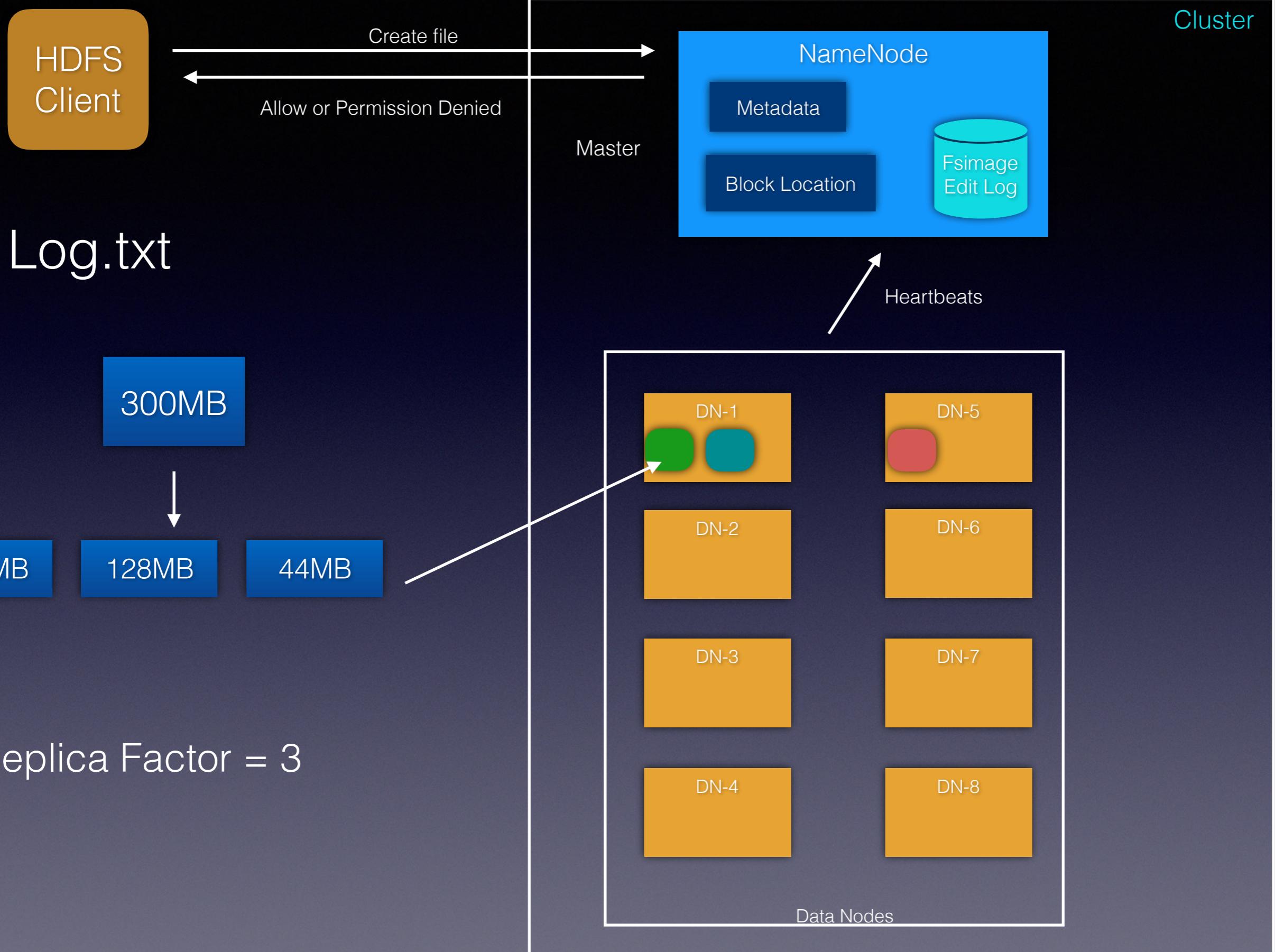
HDFS Mimarisi



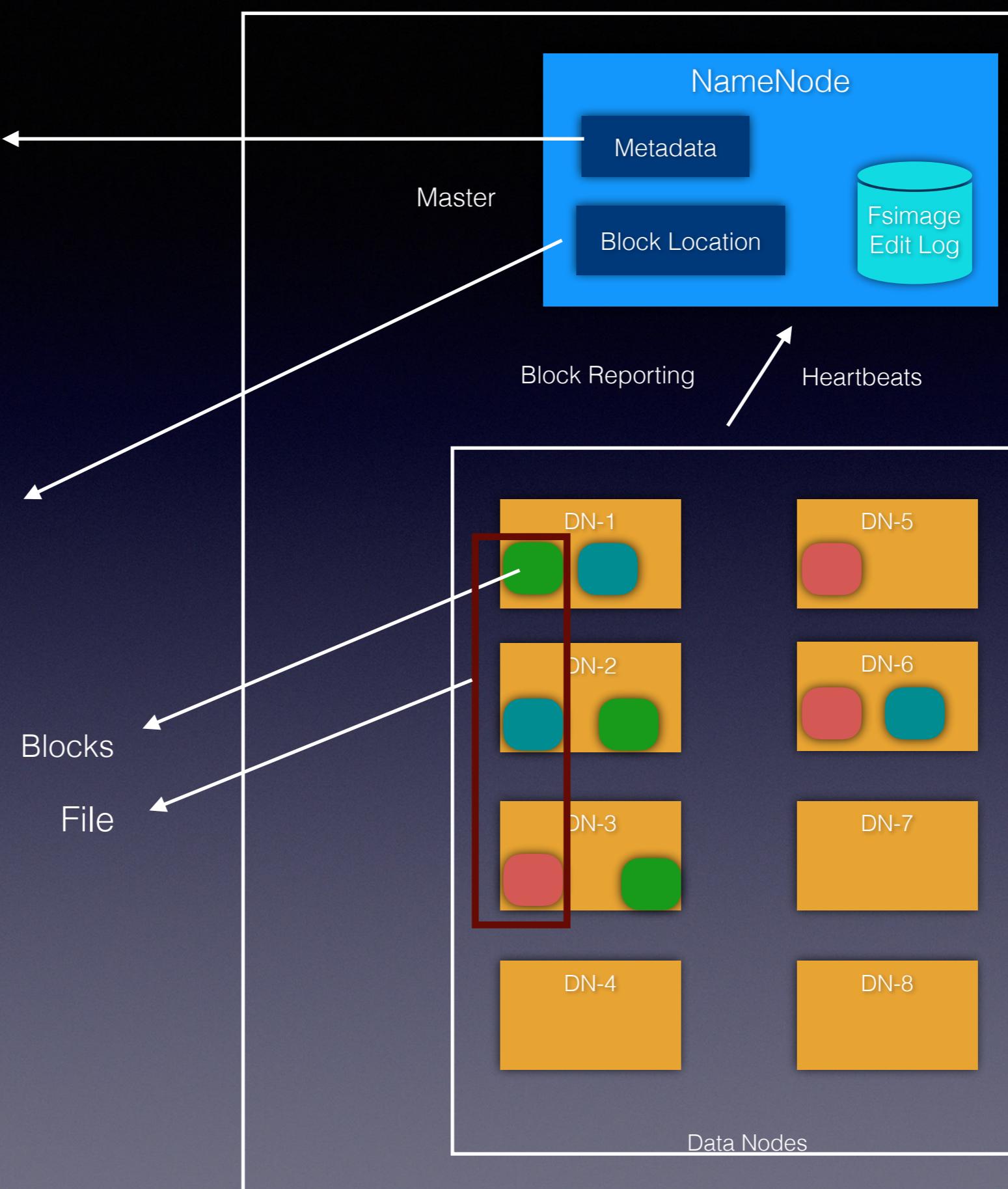
File name	log.txt
File size	300MB
No. of blocks	3
Block IDs	[#1, #2, #3]
User	hduser
Group	hadoop
Permission	-rw--r--r
Replication	3
Block size, etc	128MB, etc







File name	log.txt
File size	300MB
No. of blocks	3
Block IDs	[#1, #2, #3]
User	hduser
Group	hadoop
Permission	-rw--r--r
Replication	3
Block size, etc	128MB, etc

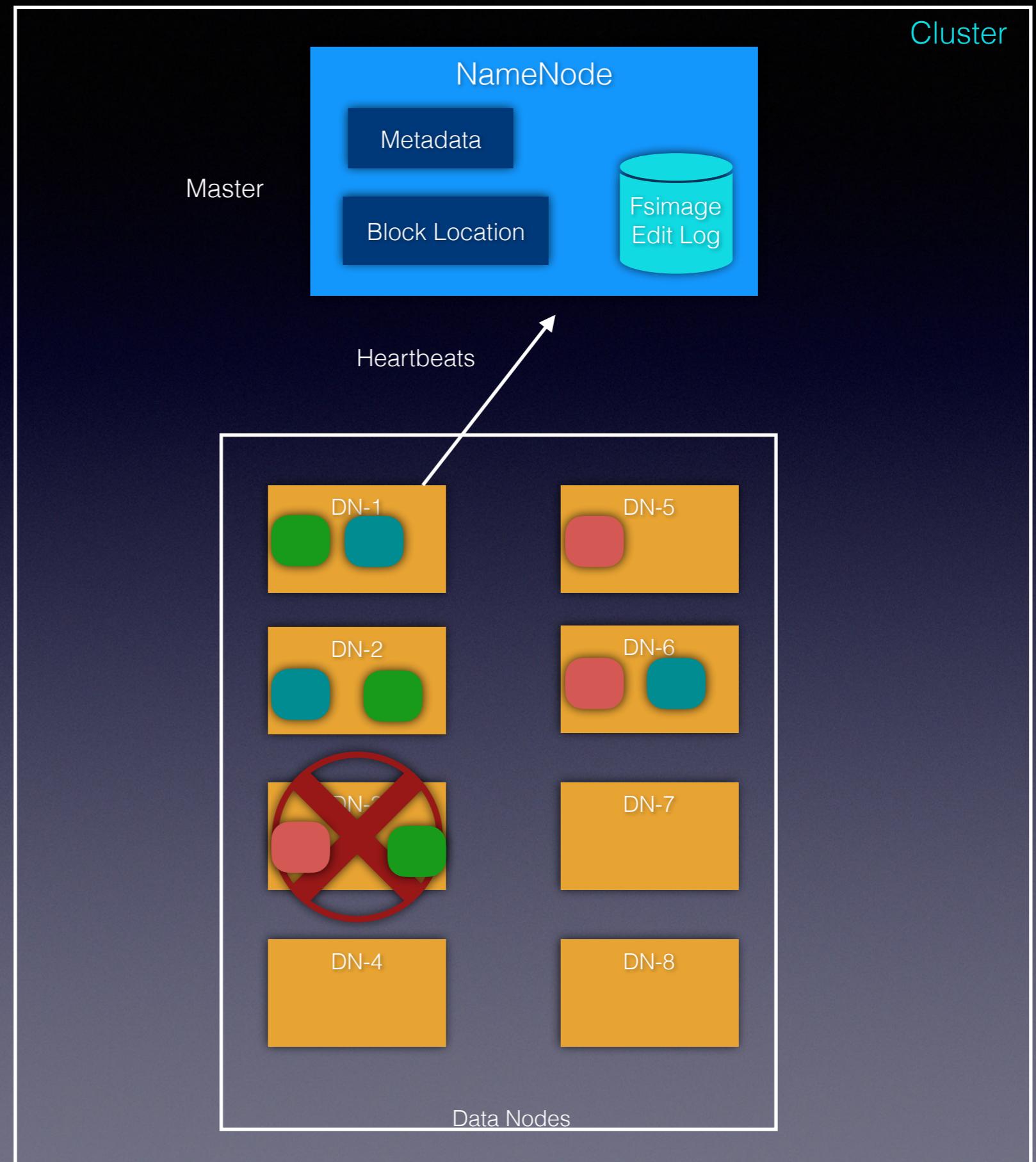


Heartbeats

Total Space
Used Space
Free Space
Data Transfer in Progress

LoadBalancing
Block Allocation

DN servis dışı kalırsa ?

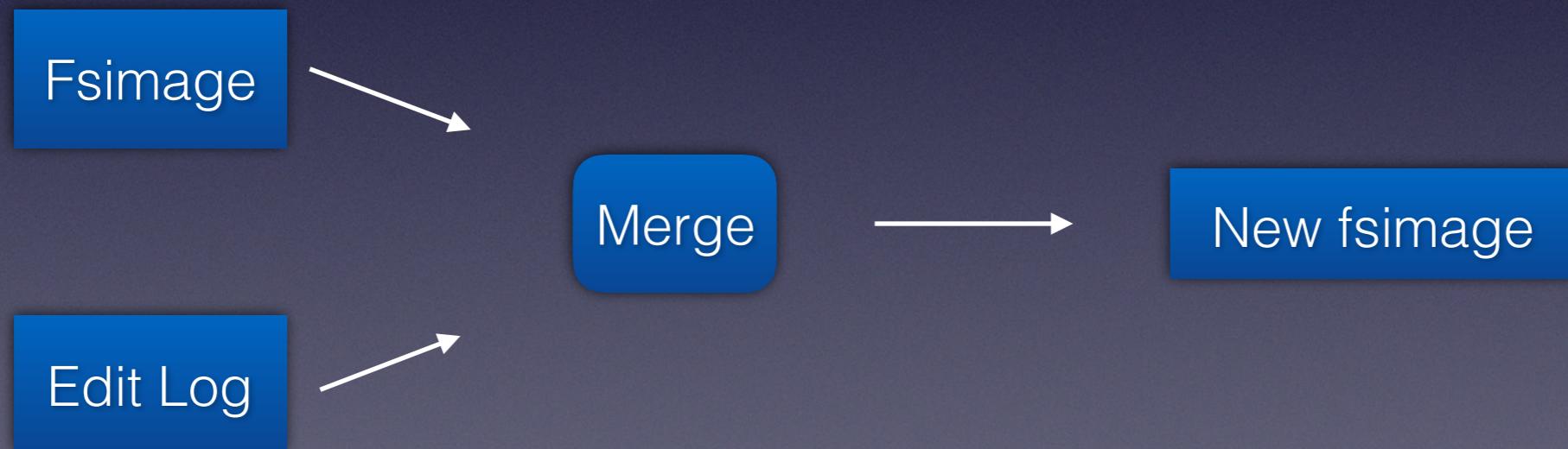


Replica Factor

- Replica faktör arttıkça;
 - Erişilebilirlik ve güvenilirlik artar.
 - Yazma ve network performansı azalır.

HDFS Check Pointing

- Secondary Name Node,
 - Name Node'un yerini **tutmaz**.
 - Aynı görevi yapmaz.



HDFS'i Kullanırken Nelere Dikkat Etmeliyiz

- Küçük boyutlu dosyalar yerine tek ve büyük boyutlu dosyalar halinde saklamaya
- Ham veriyi silmemeye
- Dosyaları kategorize edebileceğimiz şekilde tutmaya
 - /datawarehouse/tweets/2019-12-20

En önemlisi “MONITORING”



Nelerden Bahsettik

- Günlük hayatımızdaki sonuçlarından
- Büyük veri işleme modellerinden
- Teknolojilerinden
- Stream & Batch processing
- Hadoop Ekosistemi
- Kafka ve Kafka Streams

Büşra UMAN

Trendyol

Yazılım Geliştiricisi

Github: busrauman

E-mail: busrauman93@gmail.com