

# The Data Engineering Cookbook

How to master the plumbing of data science

Andreas Kretz

January 6, 2019

v0.4

# Contents

<b>1</b>	<b>Introduction</b>	<b>6</b>
<b>2</b>	<b>What is Data Science?</b>	<b>7</b>
2.1	Data Scientist . . . . .	7
2.2	Data Engineer . . . . .	8
2.3	Data Analyst . . . . .	9
2.4	Who Companies Need . . . . .	9
<b>3</b>	<b>The Basic Skills</b>	<b>9</b>
<b>4</b>	<b>Learn to Write Code</b>	<b>9</b>
4.1	Coding Basics . . . . .	10
4.2	Learn To Use GitHub — available . . . . .	10
4.3	Agile Development – available . . . . .	11
4.3.1	Agile rules I learned over the years – available . . . . .	12
4.3.2	Scrum . . . . .	13
4.3.3	OKR . . . . .	13
<b>5</b>	<b>Computer Science Basics</b>	<b>14</b>
5.1	Learn how a Computer Works . . . . .	14
5.1.1	CPU,RAM,GPU,HDD . . . . .	14
5.1.2	Differences between PCs and Servers . . . . .	14
5.2	Computer Networking . . . . .	14
5.2.1	ISO/OSI Model . . . . .	14
5.2.2	IP Subnetting . . . . .	14
5.2.3	Switch, Level 3 Switch . . . . .	14
5.2.4	Router . . . . .	14
5.2.5	Firewalls . . . . .	14
5.3	Security and Privacy . . . . .	15
5.3.1	SSL Public & Private Key Certificates . . . . .	15
5.3.2	What is a certificate authority . . . . .	15
5.3.3	Java Web Tokens . . . . .	15
5.3.4	GDPR regulations . . . . .	15
5.3.5	Privacy by design . . . . .	15
5.4	Linux . . . . .	15
5.4.1	OS Basics . . . . .	15
5.4.2	Shell scripting . . . . .	15
5.4.3	Cron jobs . . . . .	15

5.4.4	Packet management . . . . .	15
5.5	The Cloud . . . . .	15
5.5.1	AWS,Azure, IBM, Google Cloud basics . . . . .	15
5.5.2	cloud vs on premise . . . . .	15
5.5.3	up & downsides . . . . .	15
5.5.4	Security . . . . .	15
<b>6</b>	<b>My Big Data Platform Blueprint – available</b>	<b>16</b>
6.1	Ingest – available . . . . .	17
6.2	Store – available . . . . .	17
6.3	Analyse / Process – available . . . . .	18
6.4	Display – available . . . . .	18
<b>7</b>	<b>Data Science Platform</b>	<b>19</b>
7.1	Security Zone Design . . . . .	19
7.1.1	How to secure a multi layered application . . . . .	19
7.1.2	Cluster security with Kerberos . . . . .	19
7.2	Lambda Architecture . . . . .	19
7.2.1	Stream and Batch processing – available . . . . .	19
7.3	Three methods of streaming — available . . . . .	21
7.4	Big Data . . . . .	23
7.4.1	What is big data and where is the difference to data science and data analytics? . . . . .	23
7.4.2	The 4Vs of Big Data — available . . . . .	23
7.4.3	Why Big Data? — available . . . . .	24
7.4.4	What are the tools associated? . . . . .	28
7.5	What is the difference between a Data Warehouse and a Data Lake . . . . .	28
7.6	Hadoop Platforms — available . . . . .	28
7.6.1	What makes Hadoop so popular? — available . . . . .	28
7.6.2	How does a Hadoop System architecture look like . . . . .	32
7.6.3	What tools are usually in a with Hadoop Cluster . . . . .	32
7.6.4	How to select Hadoop Cluster Hardware . . . . .	32
7.7	Is ETL still relevant for Analytics? . . . . .	32
7.8	Docker . . . . .	32
7.8.1	What is docker and what do you use it for — available . . . . .	32
7.8.2	How to create, start,stop a Container . . . . .	34
7.8.3	Docker micro services? . . . . .	34
7.8.4	Kubernetes . . . . .	34
7.8.5	Why and how to do Docker container orchestration . . . . .	34

<b>8</b>	<b>How to Ingest Data</b>	<b>34</b>
8.1	Application programming interfaces . . . . .	34
8.1.1	REST APIs . . . . .	35
8.1.2	HTTP Post/Get . . . . .	35
8.1.3	API Design . . . . .	35
8.1.4	Implementation . . . . .	35
8.1.5	OAuth security . . . . .	35
8.2	JSON . . . . .	35
8.2.1	Super important for REST APIs . . . . .	35
8.2.2	How is it used for logging and processing logs . . . . .	35
<b>9</b>	<b>Distributed Processing</b>	<b>35</b>
9.1	MapReduce – available . . . . .	35
9.1.1	How does MapReduce work – available . . . . .	36
9.1.2	What is the limitation of MapReduce? – available . . . . .	39
9.2	Apache Spark . . . . .	39
9.2.1	Spark Basics . . . . .	39
9.2.2	What is the difference to MapReduce? – available . . . . .	39
9.2.3	How does Spark fit to Hadoop? – available . . . . .	40
9.2.4	Available Languages – available . . . . .	42
9.2.5	How to do stream processing . . . . .	42
9.2.6	How to do batch processing . . . . .	42
9.2.7	How does Spark use data from Hadoop – available . . . . .	42
9.2.8	What is a RDD and what is a DataFrame? . . . . .	43
9.2.9	Spark coding with Scala . . . . .	43
9.2.10	Spark coding with Python . . . . .	43
9.2.11	How and why to use SparkSQL? . . . . .	43
9.2.12	Machine Learning on Spark? (Tensor Flow) . . . . .	43
9.2.13	Spark Setup – available . . . . .	43
9.2.14	Spark Resource Management – available . . . . .	44
9.3	Message queues with Apache Kafka . . . . .	45
9.3.1	Why a message queue tool? . . . . .	45
9.3.2	Kakfa architecture . . . . .	45
9.3.3	What are topics . . . . .	45
9.3.4	What does Zookeeper have to do with Kafka . . . . .	45
9.3.5	How to produce and consume messages . . . . .	45
9.4	Machine Learning . . . . .	45
9.4.1	Training and Applying models . . . . .	45
9.4.2	What is deep learning . . . . .	45

9.4.3	How to do Machine Learning in production — available . . . . .	45
9.4.4	Why machine learning in production is harder then you think – available . . . . .	46
9.4.5	How to convince people machine learning works — available . . .	47
<b>10</b>	<b>How to Store Data</b>	<b>49</b>
10.1	Data Modeling . . . . .	49
10.1.1	How to find out how you need to store data for the business case .	49
10.1.2	How to decide what kind of storage you need to use . . . . .	49
10.2	SQL . . . . .	49
10.2.1	Database Design . . . . .	49
10.2.2	SQL Queries . . . . .	49
10.2.3	Stored Procedures . . . . .	49
10.2.4	ODBC/JDBC Server Connections . . . . .	49
10.3	NoSQL . . . . .	49
10.3.1	KeyValue Stores (HBase) . . . . .	49
10.3.2	Document Store HDFS — available . . . . .	49
10.3.3	Document Store MongoDB . . . . .	51
10.3.4	Hive Warehouse . . . . .	51
10.3.5	Impala . . . . .	51
10.3.6	Time Series Databases (?) . . . . .	51
10.3.7	MPP Databases (Greenplum) . . . . .	51
<b>11</b>	<b>How to Visualize Data</b>	<b>51</b>
11.1	Mobile Apps . . . . .	51
11.1.1	Android & IOS basics . . . . .	51
11.1.2	How to design APIs for mobile apps . . . . .	51
11.2	How to use Webservers to display content . . . . .	51
11.2.1	Tomcat . . . . .	52
11.2.2	Jetty . . . . .	52
11.2.3	NodeRED . . . . .	52
11.2.4	React . . . . .	52
11.3	Business Intelligence Tools . . . . .	52
11.3.1	Tableau . . . . .	52
11.3.2	PowerBI . . . . .	52
11.3.3	QlikSense . . . . .	52
11.4	Identity & Device Management . . . . .	52
11.4.1	What is a digital twin? . . . . .	52
11.4.2	Active Directory . . . . .	52

<b>12 Case Studies</b>	<b>52</b>
12.1 7 Steps to Successfull Data Science Project . . . . .	52
12.2 Data Science @Airbnb . . . . .	52
12.3 Data Science @Netflix – available . . . . .	52
12.4 Data Science @Uber . . . . .	56
12.5 Data Sciecne @Zalando . . . . .	56
<b>13 DEV/OPS</b>	<b>56</b>
13.1 Hadoop . . . . .	56
13.1.1 Hadoop Cluster setup and management with Cloudera Manager (for example) . . . . .	56
13.1.2 Spark code from coding to production . . . . .	56
13.1.3 How to monitor and manage data processing pipelines . . . . .	56
13.1.4 Oozie . . . . .	56
13.1.5 Airflow Application management . . . . .	56
13.1.6 Creating Statistics with Spark and Kafka . . . . .	56

# 1 Introduction

What do you actually need to learn to become an awesome data engineer? Look no further, you find it here.

How to use this document: This is not a training! It's a collection of skills, that I value highly in my daily work as a data engineer. It's intended to be a starting point for you to find the topics to look into.

This project is a work in progress! Over the next weeks I am going to share with you my thoughts on why each topic is important. I also try to include links to useful resources.

How to find out what is new? You will always find the newest version on my Patreon <https://www.patreon.com/plumbersofds>

Help make this collection awesome! Join the discussion on Patreon or write me an email to [andreaskayy@gmail.com](mailto:andreaskayy@gmail.com). Tell me your thoughts, what you value, you think should be included, or where I am wrong.

– Andreas

## 2 What is Data Science?

### 2.1 Data Scientist

Data scientists aren't like every other scientist.

Data scientists do not wear white coats or work in high tech labs full of science fiction movie equipment. They work in offices just like you and me.

What differs them from most of us is that they are the math experts. They use linear algebra and multivariable calculus to create new insight from existing data.

How exactly does this insight look?

Here's an example:

An industrial company produces a lot of products that need to be tested before shipping.

Usually such tests take a lot of time because there are hundreds of things to be tested. All to make sure that your product is not broken.

Wouldn't it be great to know early if a test fails ten steps down the line? If you knew that you could skip the other tests and just trash the product or repair it.

That's exactly where a data scientist can help you, big-time. This field is called predictive analytics and the technique of choice is machine learning.

Machine what? Learning?

Yes, machine learning, it works like this:

You feed an algorithm with measurement data. It generates a model and optimises it based on the data you fed it with. That model basically represents a pattern of how your data is looking. You show that model new data and the model will tell you if the data still represents the data you have trained it with. This technique can also be used for predicting machine failure in advance with machine learning. Of course the whole process is not that simple.

The actual process of training and applying a model is not that hard. A lot of work for the data scientist is to figure out how to pre-process the data that gets fed to the algorithms.

Because to train a algorithm you need useful data. If you use any data for the training the produced model will be very unreliable.

A unreliable model for predicting machine failure would tell you that your machine is damaged even if it is not. Or even worse: It would tell you the machine is ok even when there is an malfunction.

Model outputs are very abstract. You also need to post-process the model outputs to receive health values from 0 to 100.

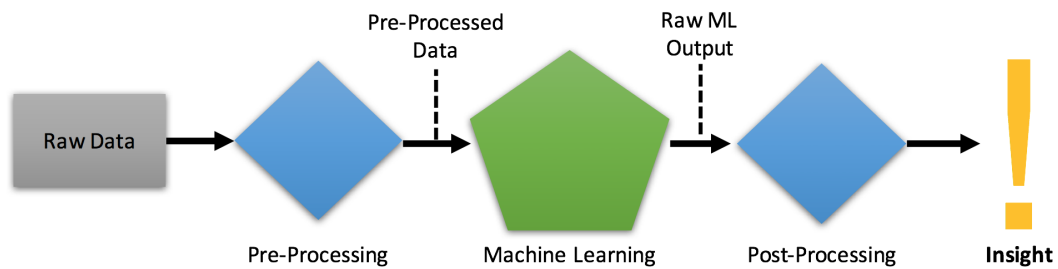


Figure 1: The Machine Learning Pipeline

## 2.2 Data Engineer

Data Engineers are the link between the management’s big data strategy and the data scientists that need to work with data.

What they do is building the platforms that enable data scientists to do their magic.

These platforms are usually used in four different ways:

- Data ingestion and storage of large amounts of data
- Algorithm creation by data scientists
- Automation of the data scientist’s machine learning models and algorithms for production use
- Data visualisation for employees and customers

– Most of the time these guys start as traditional solution architects for systems that involve SQL databases, web servers, SAP installations and other “standard” systems.

But to create big data platforms the engineer needs to be an expert in specifying, setting up and maintaining big data technologies like: Hadoop, Spark, HBase, Cassandra, MongoDB, Kafka, Redis and more.

What they also need is experience on how to deploy systems on cloud infrastructure like at Amazon or Google or on premise hardware.



## 2.3 Data Analyst

## 2.4 Who Companies Need

For a good company it is absolutely important to get well trained data engineers and data scientists.

Think of the data scientist as the professional race car driver. A fit athlete with talent and driving skills like you have never seen.

What he needs to win races is someone who will provide him the perfect race car to drive. That's what the solution architect is for.

Like the driver and his team the data scientist and the data engineer need to work closely together. They need to know the different big data tools Inside and out.

That's why companies are looking for people with Spark experience. It is a common ground between both that drives innovation.

Spark gives data scientists the tools to do analytics and helps engineers to bring the data scientist's algorithms into production.

After all, those two decide how good the data platform is, how good the analytics insight is and how fast the whole system gets into a production ready state.

## 3 The Basic Skills

## 4 Learn to Write Code

Why this is important: Without coding you cannot do much in data engineering. I cannot count the number of times I needed a quick Java hack.

The possibilities are endless:

- Writing or quickly getting some data out of a SQL DB
- Testing to produce messages to a Kafka topic
- Understanding Source code of a Java Webservice
- Reading counter statistics out of a HBase key value store

So, which language do I recommend then?

I highly recommend Java. It's everywhere!

When you are getting into data processing with Spark you should use Scala. But, after learning Java this is easy to do.

Also Python is a great choice. It is super versatile.

Personally however, I am not that big into Python. But I am going to look into it

Where to Learn? There's a Java Course on Udemy you could look at: <https://www.udemy.com/java-programming-tutorial-for-beginners>

## 4.1 Coding Basics

- OOP Object oriented programming
- What are Unit tests to make sure what you code is working
- Functional Programming
- How to use build management tools like Maven
- Resilient testing (?)

I talked about the importance of learning by doing in this podcast:

## 4.2 Learn To Use GitHub — available

Why this is important: One of the major problems with coding is to keep track of changes. It is also almost impossible to maintain a program you have multiple versions of.

Another is the topic of collaboration and documentation. Which is super Important.

Let's say you work on a Spark application and your colleagues need to make changes while you are on holiday. Without some code management they are in huge trouble:

Where is the code? What have you changed last? Where is the documentation? How do we mark what we have changed?

But if you put your code on GitHub your colleagues can find your code. They can understand it through your documentation (please also have in-line comments)

Developers can pull your code, make a new branch and do the changes. After your holiday you can inspect what they have done and merge it with your original code. and you end up having only one application

Where to learn: Check out the GitHub Guides page where you can learn all the basics: <https://guides.github.com/introduction/flow/>

This great GitHub commands cheat sheet saved my butt multiple times: <https://www.atlassian.com/git/git-cheatsheet>

Pull, Push, Branching, Forking

Also talked about it in this podcast:

## 4.3 Agile Development – available

Agility, the ability to adapt quickly to changing circumstances.

These days everyone wants to be agile. Big or small company people are looking for the “startup mentality”.

Many think it’s the corporate culture. Others think it’s the process how we create things that matters.

In this article I am going to talk about agility and self-reliance. About how you can incorporate agility in your professional career.

**Why is agile so important?** Historically development is practiced as a hard defined process. You think of something, specify it, have it developed and then built in mass production.

It’s a bit of an arrogant process. You assume that you already know exactly what a customer wants. Or how a product has to look and how everything works out.

The problem is that the world does not work this way!

Often times the circumstances change because of internal factors.

Sometimes things just do not work out as planned or stuff is harder than you think.

You need to adapt.

Other times you find out that you build something customers do not like and need to be changed.

You need to adapt.

That's why people jump on the Scrum train. Because Scrum is the definition of agile development, right?

#### **4.3.1 Agile rules I learned over the years – available**

**Is the method making a difference?** Yes, Scrum or Google's OKR can help to be more agile. The secret to being agile however, is not only how you create.

What makes me cringe is people try to tell you that being agile starts in your head. So, the problem is you?

No!

The biggest lesson I have learned over the past years is this: Agility goes down the drain when you outsource work.

**The problem with outsourcing** I know on paper outsourcing seems like a no brainer: Development costs against the fixed costs.

It is expensive to bind existing resources on a task. It is even more expensive if you need to hire new employees.

The problem with outsourcing is that you pay someone to build stuff for you.

It does not matter who you pay to do something for you. He needs to make money.

His agenda will be to spend as less time as possible on your work. That is why outsourcing requires contracts, detailed specifications, timetables and delivery dates.

He doesn't want to spend additional time on a project, only because you want changes in the middle. Every unplanned change costs him time and therefore money.

If so, you need to make another detailed specification and a contract change.

He is not going to put his mind into improving the product while developing. Firstly because he does not have the big picture. Secondly because he does not want to.

He is doing as he is told.

Who can blame him? If I was the subcontractor I would do exactly the same!

Does this sound agile to you?

**Knowledge is king: A lesson from Elon Musk** Doing everything in house, that's why startups are so productive. No time is wasted on waiting for someone else.

If something does not work, or needs to be changed, there is someone in the team who can do it right away. .

One very prominent example who follows this strategy is Elon Musk.

Tesla's Gigafactories are designed to get raw materials in on one side and spit out cars on the other. Why do you think Tesla is building Gigafactories who cost a lot of money?

Why is SpaceX building its one space engines? Clearly there are other, older, companies who could do that for them.

Why is Elon building tunnel boring machines at his new boring company?

At first glance this makes no sense!

**How you really can be agile** If you look closer it all comes down to control and knowledge. You, your team, your company, needs to as much as possible on your own. Self-reliance is king.

Build up your knowledge and therefore the teams knowledge. When you have the ability to do everything yourself, you are in full control.

You can build electric cars, rocket engines or bore tunnels.

Don't largely rely on others and be confident to just do stuff on your own.

Dream big and JUST DO IT!

PS. Don't get me wrong. You can still outsource work. Just do it in a smart way by outsourcing small independent parts.

#### **4.3.2 Scrum**

#### **4.3.3 OKR**

I talked about this in this Podcast: <https://anchor.fm/andreaskayy/embed/episodes/Agile-Development-Is-Important-But-Please-Dont-Do-Scrum-PoDS-041-e2e2j4>

## **5 Computer Science Basics**

### **5.1 Learn how a Computer Works**

#### **5.1.1 CPU, RAM, GPU, HDD**

#### **5.1.2 Differences between PCs and Servers**

I talked about computer hardware and GPU processing in this podcast: <https://anchor.fm/andreaskayy/e/the-hardware-and-the-GPU-is-super-important-PoDS-030-e23rig>

### **5.2 Computer Networking**

#### **5.2.1 ISO/OSI Model**

#### **5.2.2 IP Subnetting**

#### **5.2.3 Switch, Level 3 Switch**

#### **5.2.4 Router**

#### **5.2.5 Firewalls**

I talked about Network Infrastructure and Techniques in this podcast: <https://anchor.fm/andreaskayy/e/Networking-Infrastructure-and-Linux-031-PoDS-e242bh>

## **5.3 Security and Privacy**

### **5.3.1 SSL Public & Private Key Certificates**

### **5.3.2 What is a certificate authority**

### **5.3.3 JAva Web Tokens**

### **5.3.4 GDPR regulations**

### **5.3.5 Privacy by design**

## **5.4 Linux**

### **5.4.1 OS Basics**

### **5.4.2 Shell scripting**

### **5.4.3 Cron jobs**

### **5.4.4 Packet management**

Linux Tips are the second part of this podcast: <https://anchor.fm/andreaskayy/embed/episodes/IT-Networking-Infrastructure-and-Linux-031-PoDS-e242bh>

## **5.5 The Cloud**

### **5.5.1 AWS,Azure, IBM, Google Cloud basics**

### **5.5.2 cloud vs on premise**

### **5.5.3 up & downsides**

### **5.5.4 Security**

Listen to a few thoughts about the cloud in this podcast: <https://anchor.fm/andreaskayy/embed/episode/Be-Arrogant-The-Cloud-is-Safer-Then-Your-On-Premise-e16k9s>

## 6 My Big Data Platform Blueprint – available

Some time ago I have created a simple and modular big data platform blueprint for myself. It is based on what I have seen in the field and read in tech blogs all over the internet.

Today I am going to share it with you.

Why do I believe it will be super useful to you?

Because, unlike other blueprints it is not focused on technology. It is based on four common big data platform design patterns.

Following my blueprint will allow you to create the big data platform that fits exactly your needs. Building the perfect platform will allow data scientists to discover new insights.

It will enable you to perfectly handle big data and allow you to make data driven decisions.

**THE BLUEPRINT – available** The blueprint is focused on the four key areas: Ingest, store, analyse and display.

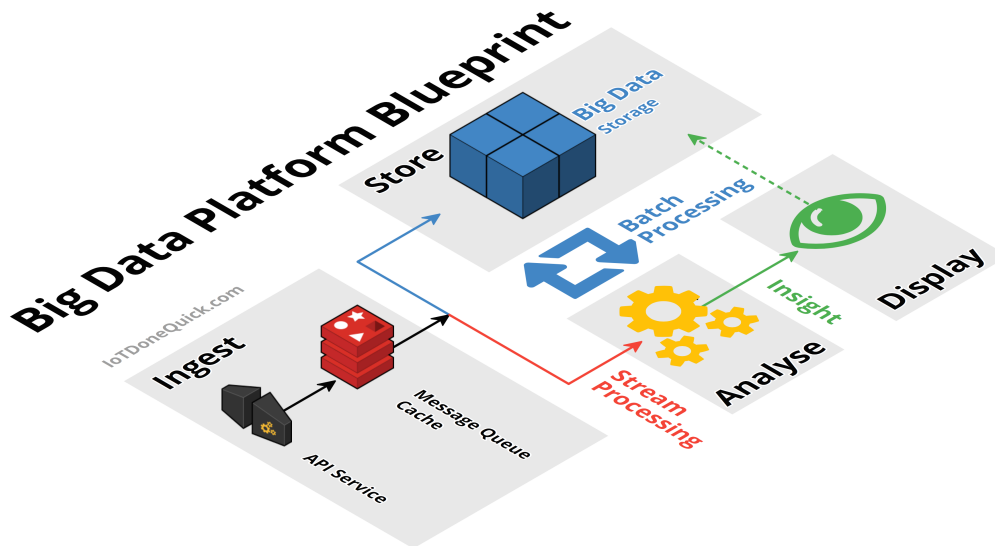


Figure 2: Platform Blueprint

Having the platform split like this turns it into a modular platform with loosely coupled interfaces.

Why is it so important to have a modular platform?

If you have a platform that is not modular you end up with something that is fixed or



hard to modify. This means you can not adjust the platform to changing requirements of the company.

Because of modularity it is possible to switch out every component, if you need it.

Now, lets talk more about each key area.

## **6.1 Ingest – available**

Ingestion is all about getting the data in from the source and making it available to later stages. Sources can be everything from tweets, server logs to IoT sensor data like from cars.

Sources send data to your API Services. The API is going to push the data into a temporary storage.

The temporary storage allows other stages simple and fast access to incoming data.

A great solution is to use messaging queue systems like Apache Kafka, RabbitMQ or AWS Kinesis. Sometimes people also use caches for specialised applications like Redis.

A good practice is that the temporary storage follows the publish, subscribe pattern. This way APIs can publish messages and Analytics can quickly consume them.

## **6.2 Store – available**

This is the typical big data storage where you just store everything. It enables you to analyse the big picture.

Most of the data might seem useless for now, but it is of utmost importance to keep it. Throwing data away is a big no no.

Why not throw something away when it is useless?

Although it seems useless for now, data scientists can work with the data. They might find new ways to analyse the data and generate valuable insight from it.

What kind of systems can be used to store big data?

Systems like Hadoop HDFS, Hbase, Amazon S3 or DynamoDB are a perfect fit to store big data.

## 6.3 Analyse / Process – available

The analyse stage is where the actual analytics is done. Analytics, in the form of stream and batch processing.

Streaming data is taken from ingest and fed into analytics. Streaming analyses the “live” data thus, so generates fast results.

As the central and most important stage, analytics also has access to the big data storage. Because of that connection, analytics can take a big chunk of data and analyse it.

This type of analysis is called batch processing. It will deliver you answers for the big questions.

To learn more about stream and batch processing read my blog post: [How to Create New and Exciting Big Data Aided Products](#)

The analytics process, batch or streaming, is not a one way process. Analytics also can write data back to the big data storage.

Often times writing data back to the storage makes sense. It allows you to combine previous analytics outputs with the raw data.

Analytics insight can give meaning to the raw data when you combine them. This combination will often times allow you to create even more useful insight.

A wide variety of analytics tools are available. Ranging from MapReduce or AWS Elastic MapReduce to Apache Spark and AWS lambda.

## 6.4 Display – available

Displaying data is as important as ingesting, storing and analysing it. People need to be able to make data driven decisions.

This is why it is important to have a good visual presentation of the data. Sometimes you have a lot of different use cases or projects using the platform.

It might not be possible for you to build the perfect UI that fits everyone. What you should do in this case is enable others to build the perfect UI themselves.

How to do that? By creating APIs to access the data and making them available to developers.

Either way, UI or API the trick is to give the display stage direct access to the data in the big data cluster. This kind of access will allow the developers to use analytics results as well as raw data to build the the perfect application.

## 7 Data Science Platform

### 7.1 Security Zone Design

#### 7.1.1 How to secure a multi layered application

(UI in different zone then SQL DB)

#### 7.1.2 Cluster security with Kerberos

I talked about security zone design and lambda architecture in this podcast: <https://anchor.fm/andreask-to-Design-Security-Zones-and-Lambda-Architecture-PoDS-032-e248q2>

### 7.2 Lambda Architecture

#### 7.2.1 Stream and Batch processing – available

**Batch Processing:** Ask the big questions. Remember your last yearly tax statement?

You break out the folders. You run around the house searching for the receipts.

All that fun stuff.

When you finally found everything you fill out the form and send it on its way.

Doing the tax statement is a prime example of a batch process.

Data comes in and gets stored, analytics loads the data from storage and creates an output (insight):

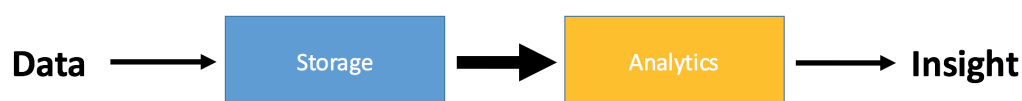


Figure 3: Batch Processing Pipeline

Batch processing is something you do either without a schedule or on a schedule (tax statement). It is used to ask the big questions and gain the insights by looking at the big picture.

To do so, batch processing jobs use large amounts of data. This data is provided by storage systems like Hadoop HDFS.

They can store lots of data (petabytes) without a problem.

Results from batch jobs are very useful, but the execution time is high. Because the amount of used data is high.

It can take minutes or sometimes hours until you get your results.

**Stream Processing:** Gain instant insight into your data.

Streaming allows users to make quick decisions and take actions based on “real-time” insight. Contrary to batch processing, streaming processes data on the fly, as it comes in.

With streaming you don’t have to wait minutes or hours to get results. You gain instant insight into your data.

In the batch processing pipeline, the analytics was after the data storage. It had access to all the available data.

Stream processing creates insight before the data storage. It has only access to fragments of data as it comes in.

As a result the scope of the produced insight is also limited. Because the big picture is missing.



Figure 4: Stream Processing Pipeline

Only with streaming analytics you are able to create advanced services for the customer. Netflix for instance incorporated stream processing into Chuckwa V2.0 and the new Keystone pipeline.

One example of advanced services through stream processing is the Netflix “Trending Now” feature. Check out the Netflix case study.

**Should you do stream or batch processing?** It is a good idea to start with batch processing. Batch processing is the foundation of every good big data platform.

A batch processing architecture is simple, and therefore quick to set up. Platform simplicity means, it will also be relatively cheap to run.

A batch processing platform will enable you to quickly ask the big questions. They will give you invaluable insight into your data and customers.

When the time comes and you also need to do analytics on the fly, then add a streaming pipeline to your batch processing big data platform.

### 7.3 Three methods of streaming — available

In stream processing sometimes it is ok to drop messages, other times it is not. Sometimes it is fine to process a message multiple times, other times that needs to be avoided like hell.

Today's topic are the different methods of streaming: At most once, at least once and exactly once.

What this means and why it is so important to keep them in mind when creating a solution. That is what you will find out in this article.

**At Least Once** At least once, means a message gets processed in the system once or multiple times. So with at least once it's not possible that a message gets into the system and is not getting processed.

It's not getting dropped or lost somewhere in the system.

One example where at least once processing can be used is when you think about a fleet management of cars. You get GPS data from cars and that GPS data is transmitted with a timestamp and the GPS coordinates.

It's important that you get the GPS data at least once, so you know where the car is. If you're processing this data multiple times, it always has the the timestamp with it.

Because of that it does not matter that it gets processed multiple times, because of the timestamp. Or that it would be stored multiple times, because it would just override the existing one.

**At Most Once** The second streaming method is at most once. At most once means that it's okay to drop some information, to drop some messages.

But it's important that a message is only only processed once as a maximum.

A example for this is event processing. Some event is happening and that event is not important enough, so it can be dropped. It doesn't have any consequences when it gets dropped.

But when that event happens it's important that it does not get processed multiple times. Then it would look as if the event happend five or six times instead of only one.

Think about engine misfires. If it happens once, no big deal. But if the system tells you it happens a lot you will think you have a problem with your engine.

**Exactly Once** Another thing is exactly once, this means it's not okay to drop data, it's not okay to lose data and it's also not okay to process one message what data said multiple times

A example for this is for instance banking. When you think about credit card transactions it's not okay to drop a transaction.

When dropped your payment is not going through. It's also not okay to have a transaction processed multiple times, because then you are paying multiple times.

**Check The Tools!** All of this sounds very simple and logical. What kind of processing is done has to be a requirement for your use case.

It needs to be thought about in the design process, because not every tool is supporting all three methods. Very often you need to code your application very differently based on the streaming method.

Especially exactly once is very hard to do.

So, the tool of data processing needs to be chosen based on if you need exactly once, at least once or if you need at most once.

## 7.4 Big Data

### 7.4.1 What is big data and where is the difference to data science and data analytics?

I talked about the difference in this podcast: <https://anchor.fm/andreaskayy/embed/episodes/BI-vs-Data-Science-vs-Big-Data-e199hq>

### 7.4.2 The 4Vs of Big Data — available

It is a complete misconception. Volume is only one part of the often called four V's of big data: Volume, velocity, variety and veracity.

**Volume** is about the size. How much data you have.

**Velocity** is about the speed of how fast the data is getting to you.

How much data is in a specific time needs to get processed or is coming into the system. This is where the whole concept of streaming data and real-time processing comes in to play.

**Variety** is the third one. It means, that the data is very different. That you have very different types of data structures.

Like CSV files, PDFs that you have stuff in XML. That you have JSON logfiles, or that you have data in some kind of a key value store.

It's about the variety of data types from different sources that you basically want to join together. All to make an analysis based on that data.

**Veracity** is fourth and this is a very very difficult one. The issue with big data is, that it is very unreliable.

You cannot really trust the data. Especially when you're coming from the IoT, the Internet of Things side. Devices use sensors for measurement of temperature, pressure, acceleration and so on.

You cannot always be hundred percent sure that the actual measurement is right.

When you have data that is from for instance SAP and it contains data that is created by hand you also have problems. As you know we humans are bad at inputting stuff.

Everybody articulates different. We make mistakes, down to the spelling and that can be a very difficult issue for analytics.

I talked about the 4Vs in this podcast: <https://anchor.fm/andreaskayy/embed/episodes/4-Vs-Of-Big-Data-Are-Enough-e1h2ra>

### 7.4.3 Why Big Data? — available

What I always emphasize is the four V's are quite nice. They give you a general direction.

There is a much more important issue: Catastrophic Success.

What I mean by catastrophic success is, that your project, your startup or your platform has more growth than you anticipated. Exponential growth is what everybody is looking for.

Because with exponential growth there is the money. It starts small and gets very big very fast. The classic hockey stick curve:

1,2,4,8,16,32,64,128,256,512,1024,2048,4096,8192,16384...BOOM!

Think about it. It starts small and quite slow, but gets very big very fast.

You get a lot of users or customers who are paying money to use your service, the platform or whatever. If you have a system that is not equipped to scale and process the data the whole system breaks down.

That's catastrophic success. You are so successful and grow so fast that you cannot fulfill the demand anymore. And so you fail and it's all over.

It's now like you just can make that up while you go. That you can foresee in a few months or weeks the current system doesn't work anymore.

**Planning is Everything** It's all happens very very fast and you cannot react anymore. There's a necessary type of planning and analyzing the potential of your business case necessary.

Then you need to decide if you actually have big data or not.

You need to decide if you use big data tools. This means when you conceptualize the whole infrastructure it might look ridiculous to actually focus on big data tools.

But in the long run it will help you a lot. Good planning will get a lot of problems out of the way, especially if you think about streaming data and real-time analytics.



**The Problem With ETL** A typical old-school platform deployment would look like the picture below. Devices use a data API to upload data that gets stored in a SQL database. An external analytics tool is querying data and uploading the results back to the SQL db. Users then use the user interface to display data stored in the database.

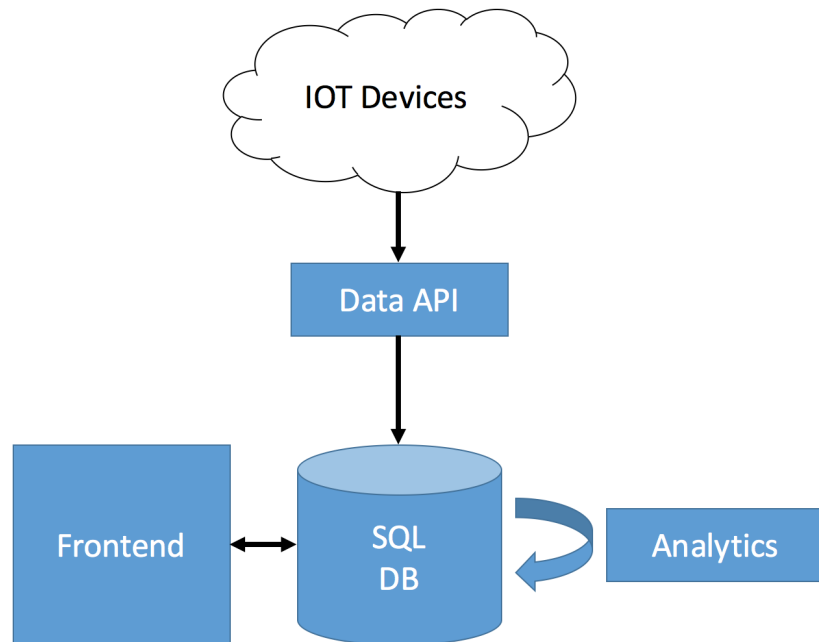


Figure 5: Common SQL Platform Architecture

Now, when the front end queries data from the SQL database the following three steps happen:

The database extracts all the needed rows from the storage. Extracted data gets transformed, for instance sorted by timestamp or something a lot more complex.

The extracted and transformed data gets loaded to the destination (the user interface) for chart creation. With exploding amounts of stored data the ETL process starts being a real problem.

Analytics is working with large data sets, for instance whole days, weeks, months or more. Data sets are very big like 100GB or Terabytes. That means Billions or Trillions of rows.

This has the result that the ETL process for large data sets takes longer and longer. Very quickly the ETL performance gets so bad it won't deliver results to analytics anymore.

A traditional solution to overcome these performance issues is trying to increase the performance of the database server. That's what's called scaling up.

**Scaling Up** To scale up the system and therefore increase ETL speeds administrators resort to more powerful hardware by:

Speeding up the extract performance by adding faster disks to physically read the data faster. Increasing RAM for row caching. What is already in memory does not have to be read by slow disk drives. Using more powerful CPU's for better transform performance (more RAM helps here as well) Increasing or optimising networking performance for faster data delivery to the front end and analytics Scaling up the system is fairly easy.

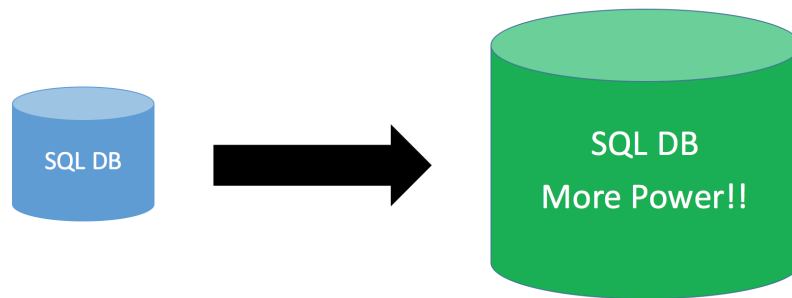


Figure 6: Scaling up a SQL Database

But with exponential growth it is obvious that sooner or later (more sooner than later) you will run into the same problems again. At some point you simply cannot scale up anymore because you already have a monster system, or you cannot afford to buy more expensive hardware.

The next step you could take would be scaling out.

**Scaling Out** Scaling out is the opposite of scaling up. Instead of building bigger systems the goal is to distribute the load between many smaller systems.

The simplest way of scaling out an SQL database is using a storage area network (SAN) to store the data. You can then use up to eight SQL servers, attach them to the SAN and let them handle queries. This way load gets distributed between those eight servers.

One major downside of this setup is that, because the storage is shared between the sql servers, it can only be used as an read only database. Updates have to be done periodically, for instance once a day. To do updates all SQL servers have to detach from the database. Then, one is attaching the db in read-write mode and refreshing the data. This procedure can take a while if a lot of data needs to be uploaded.

This [Link](#) to a Microsoft MSDN page has more options of scaling out an SQL database for you.

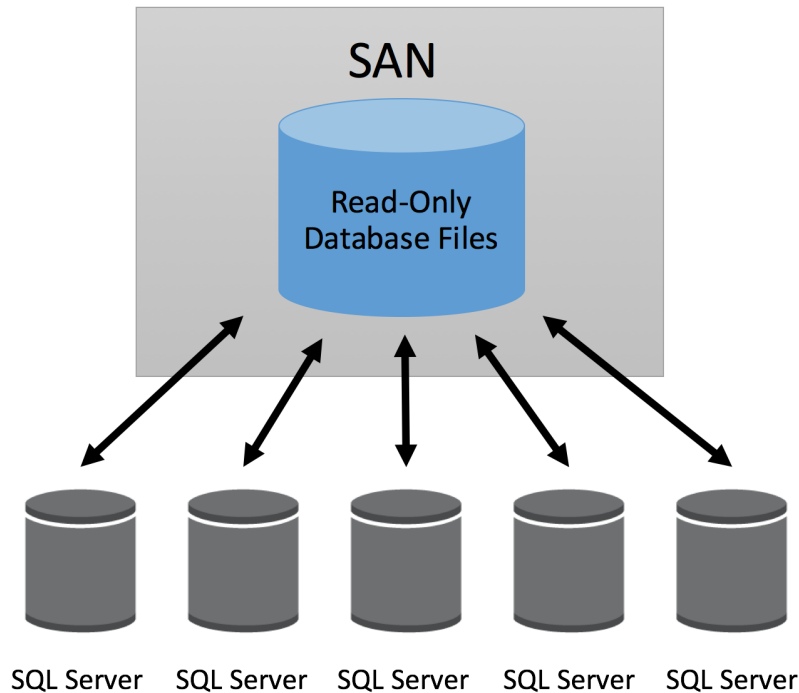


Figure 7: Scaling out a SQL Database

I deliberately don't want to get into details about possible scaling out solutions. The point I am trying to make is that while it is possible to scale out SQL databases it is very complicated.

There is no perfect solution. Every option has its up- and downsides. One common major issue is the administrative effort that you need to take to implement and maintain a scaled out solution.

**Please Don't go Big Data** If you don't run into scaling issues please, do not use big data tools!

Big data is a expensive thing. A Hadoop cluster for instance needs at least five servers to work properly. More is better.

Believe me this stuff costs a lot of money.

Especially when you are talking the maintenance and development on top big daat tools into account.

If you don't need it it's making absolutely no sense at all!

On the other side: If you really need big data tools they will save your ass :)

#### 7.4.4 What are the tools associated?

### 7.5 What is the difference between a Data Warehouse and a Data Lake

### 7.6 Hadoop Platforms — available

When people talk about big data, one of the first things come to mind is Hadoop. Google's search for Hadoop returns about 28 million results.

It seems like you need Hadoop to do big data. Today I am going to shed light onto why Hadoop is so trendy.

You will see that Hadoop has evolved from a platform into an ecosystem. It's design allows a lot of Apache projects and 3rd party tools to benefit from Hadoop.

I will conclude with my opinion on, if you need to learn Hadoop and if Hadoop is the right technology for everybody.

**WHAT IS HADOOP?** Hadoop is a platform for distributed storing and analyzing of very large data sets.

Hadoop has four main modules: Hadoop common, HDFS, MapReduce and YARN. The way these modules are woven together is what makes Hadoop so successful.

The Hadoop common libraries and functions are working in the background. That's why I will not go further into them. They are mainly there to support Hadoop's modules.

#### 7.6.1 What makes Hadoop so popular? — available

Storing and analyzing data as large as you want is nice. But what makes Hadoop so popular?

Hadoop's core functionality is the driver of Hadoop's adoption. Many Apache side projects use it's core functions.

Because of all those side projects Hadoop has turned more into an ecosystem. An ecosystem for storing and processing big data.

To better visualize this eco system I have drawn you the following graphic. It shows some projects of the Hadoop ecosystem who are closely connected with the Hadoop.

It is not a complete list. There are many more tools that even I don't know. Maybe I am drawing a complete map in the future.

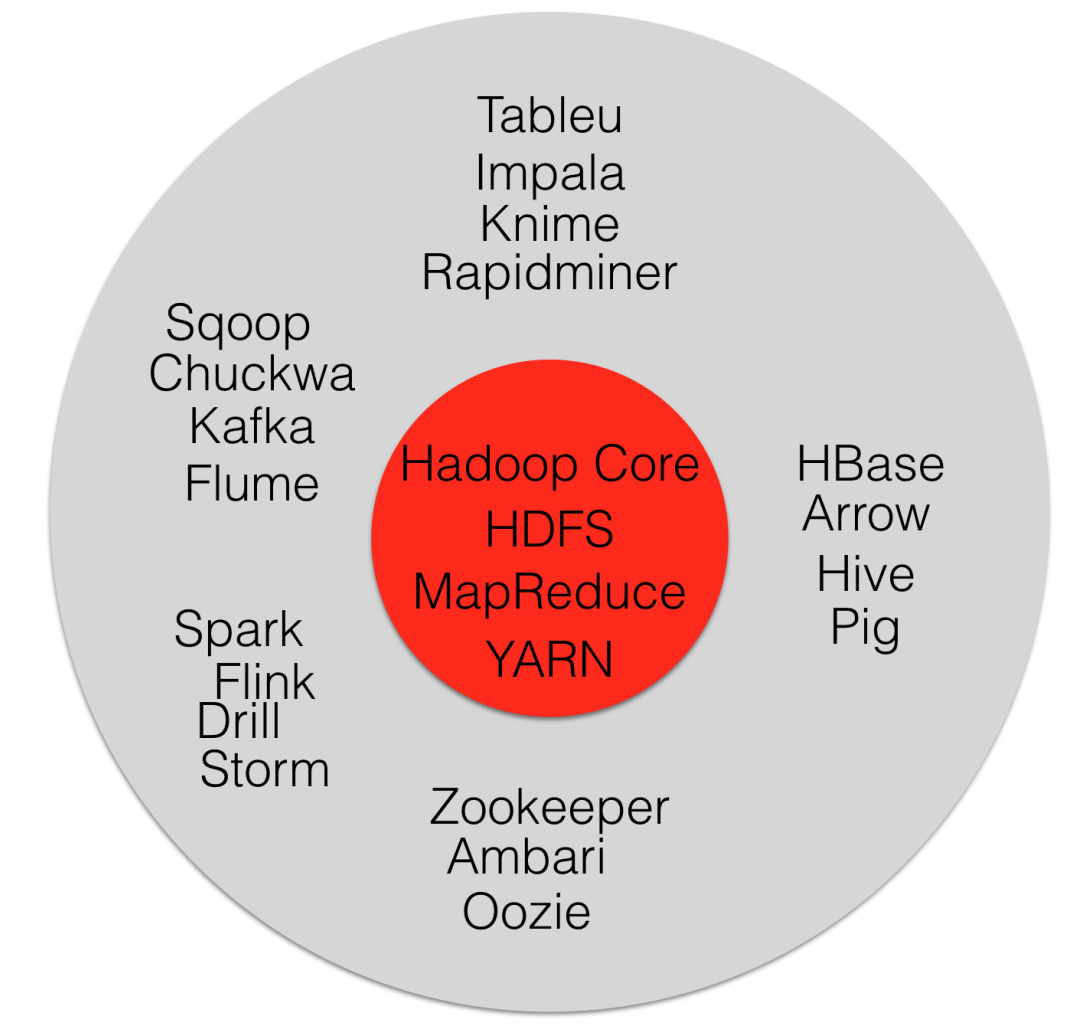


Figure 8: Hadoop Ecosystem Components

**HOW THE ECOSYSTEM'S COMPONENTS WORK TOGETHER** Remember my big data platform blueprint? The blueprint has four stages: Ingest, store, analyse and display.

Because of the Hadoop ecosystem” the different tools in these stages can work together perfectly.

Here's an example:

You use Apache Kafka to ingest data, and store the it in HDFS. You do the analytics with Apache Spark and as a backend for the display you store data in Apache HBase.

To have a working system you also need YARN for resource management. You also need

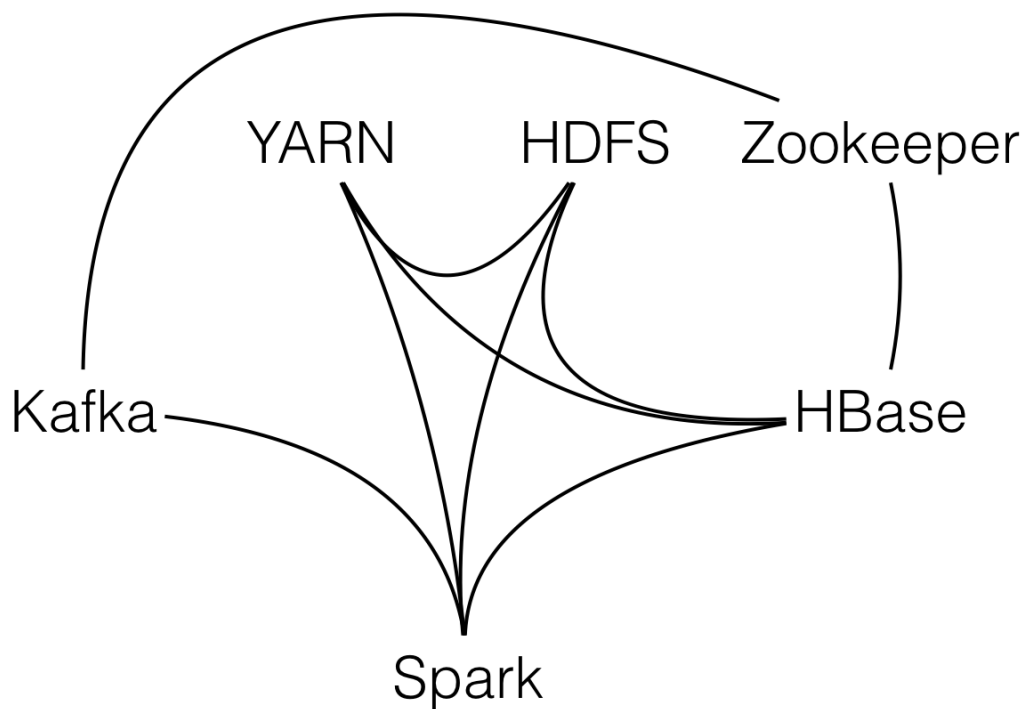


Figure 9: Connections between tools

Zookeeper, a configuration management service to use Kafka and HBase

As you can see in the picture below each project is closely connected to the other.

Spark for instance, can directly access Kafka to consume messages. It is able to access HDFS for storing or processing stored data.

It also can write into HBase to push analytics results to the front end.

The cool thing of such ecosystem is that it is easy to build in new functions.

Want to store data from Kafka directly into HDFS without using Spark?

No problem, there is a project for that. Apache Flume has interfaces for Kafka and HDFS.

It can act as an agent to consume messages from Kafka and store them into HDFS. You even do not have to worry about Flume resource management.

Flume can use Hadoop's YARN resource manager out of the box.

**IS HADOOP WORKING EVERYWHERE?** Although Hadoop is so popular it is not the silver bullet. It isn't the tool that you should use for everything.

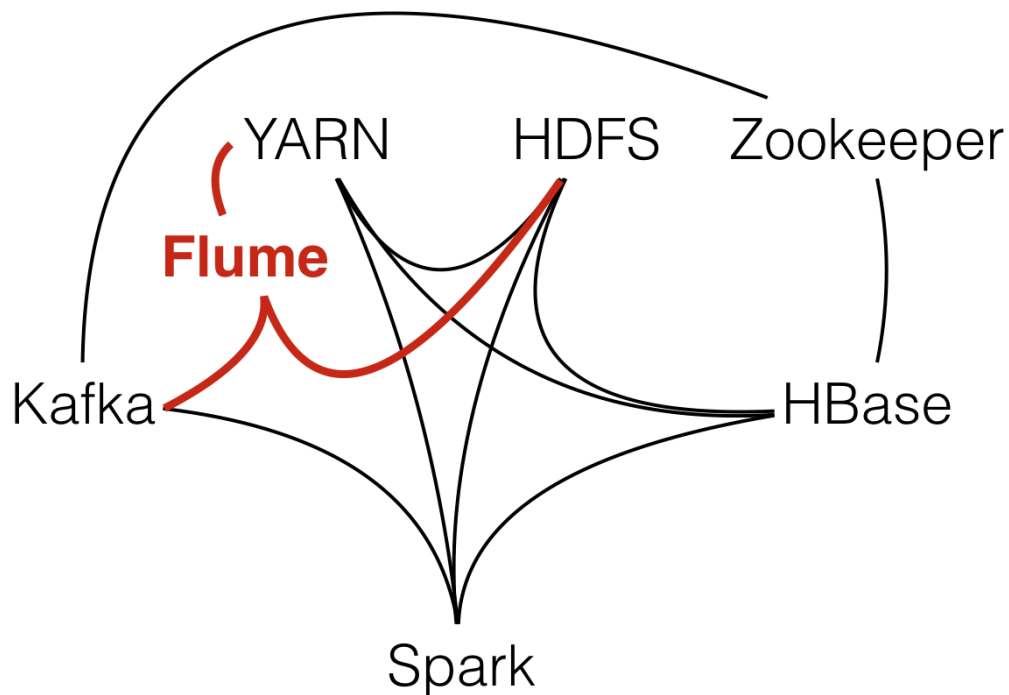


Figure 10: Flume Integration

Often times it does not make sense to deploy a Hadoop cluster, because it can be overkill. Hadoop does not run on a single server.

You basically need at least five servers, better six to run a small cluster. Because of that, the initial platform costs are quite high.

One option you have is to use a specialized systems like Cassandra, MongoDB or other NoSQL DB's for storage. Or you move to Amazon and use Amazon's Simple Storage Service, or S3.

Guess what the tech behind S3 is. Yes, HDFS. That's why AWS also has the equivalent to MapReduce named Elastic MapReduce.

The great thing about S3 is that you can start very small. When your system grows you don't have to worry about s3's server scaling.

**SHOULD YOU LEARN HADOOP?** Yes, I definitely recommend you to get to now how Hadoop works and how to use it. As I have shown you in this article, the ecosystem is quite large.

Many big data projects use Hadoop or can interface with it. That's why it is generally a good idea to know as many big data technologies as possible.

Not in depth, but to the point that you know how they work and how you can use them. Your main goal should be to be able to hit the ground running when you join a big data project.

Plus, most of the technologies are open source. You can try them out for free.

### **7.6.2 How does a Hadoop System architecture look like**

### **7.6.3 What tools are usually in a with Hadoop Cluster**

Yarn Zookeeper HDFS Oozie Flume Hive

### **7.6.4 How to select Hadoop Cluster Hardware**

## **7.7 Is ETL still relevant for Analytics?**

I talked about this in this podcast: <https://anchor.fm/andreaskayy/embed/episodes/Is-ETL-Dead-For-Data-Science-Big-Data—PoDS-039-e2b604>

## **7.8 Docker**

### **7.8.1 What is docker and what do you use it for — available**

Have you played around with Docker yet? If you're a data science learner or a data scientist you need to check it out!

It's awesome because it simplifies the way you can set up development environments for data science. If you want to set up a dev environment you usually have to install a lot of packages and tools.

**Don't Mess Up Your System** What this does is you basically mess up your operating system. If you're a starter you don't know which packages you need to install. You don't know which tools you need to install.

If you want to for instance start with Jupyter notebooks you need to install that on your PC somehow. Or you need to start installing tools like PyCharmor Anaconda.



All that gets added to your system and so you mess up your system more and more and more. What Docker brings you, especially if you're on a Mac or a Linux system is simplicity.

**Preconfigured Images** Because it is so easy to install on those systems. Another cool thing about docker images is you can just search them in the Docker store, download them and install them on your system.

Running them in a completely pre-configured environment. You don't need to think about stuff you go to the Docker library you search for deep learning GPU and Python.

You get a list of images you can download. You download one, start it up, you go to the browser hit up the URL and just start coding.

Start doing the work. The only other thing you need to do is bind some drives to that instance so you can exchange files. And then that's it!

There is no way that you can crash or mess up your system. It's all encapsulated into Docker. Why this works is because Docker has native access to your hardware.

**Take It With You** It's not a completely virtualized environment like a VirtualBox. An image has the upside that you can take it wherever you want. So if you're on your PC at home use that there.

Make a quick build, take the image and go somewhere else. Install the image which is usually quite fast and just use it like you're at home.

It's that awesome!

**Kubernetes Container Deployment** I am getting into Docker a lot more myself. For a bit different reasons.

What I'm looking for is using Docker with Kubernetes. Kubernetes you can automate the whole container deployment process.

The idea with is that you have a cluster of machines. Let's say you have 10 server cluster and you run Kubernetes on them.

Kubernetes lets you spin up Docker containers on-demand to execute tasks. You can set up how much resources like CPU, RAM, Network, Docker container can use.

You can basically spin up containers, on the cluster on demand. When ever you need to do a analytics task.

Perfect for Data Science.

### **7.8.2 How to create, start,stop a Container**

### **7.8.3 Docker micro services?**

### **7.8.4 Kubernetes**

### **7.8.5 Why and how to do Docker container orchestration**

Podcast about how data science learners use Docker (for data scientists): <https://anchor.fm/andreaskayy/Data-Science-Go-Docker-e10n7u>

## **8 How to Ingest Data**

### **8.1 Application programming interfaces**

Check out my podcast about how APIs rule the world: <https://anchor.fm/andreaskayy/embed/episodes/APIs-Rule-The-World-PoDS-033-e24ttq>

### **8.1.1 REST APIs**

### **8.1.2 HTTP Post/Get**

### **8.1.3 API Design**

### **8.1.4 Implementation**

### **8.1.5 OAuth security**

## **8.2 JSON**

### **8.2.1 Super important for REST APIs**

### **8.2.2 How is it used for logging and processing logs**

## **9 Distributed Processing**

I talked about why distributed processing is so super important in this Podcast:

### **9.1 MapReduce – available**

Since the early days of the Hadoop eco system, the MapReduce framework is one of the main components of Hadoop alongside the Hadoop file system HDFS.

Google for instance used MapReduce to analyse stored html content of websites through counting all the html tags and all the words and combinations of them (for instance headlines). The output was then used to create the page ranking for Google Search.

That was when everybody started to optimise his website for the google search. Serious search engine optimisation was borne. That was the year 2004.

How MapReduce is working is, that it processes data in two phases: The map phase and the reduce phase.

In the map phase, the framework is reading data from HDFS. Each dataset is called an input record.

Then there is the reduce phase. In the reduce phase, the actual computation is done and the results are stored. The storage target can either be a database or back HDFS or

something else.

After all it's Java – so you can implement what you like.

The magic of MapReduce is how the map and reduce phase are implemented and how both phases are working together.

The map and reduce phases are parallelised. What that means is, that you have multiple map phases (mappers) and reduce phases (reducers) that can run in parallel on your cluster machines.

Here's an example how such a map and reduce process works with data:

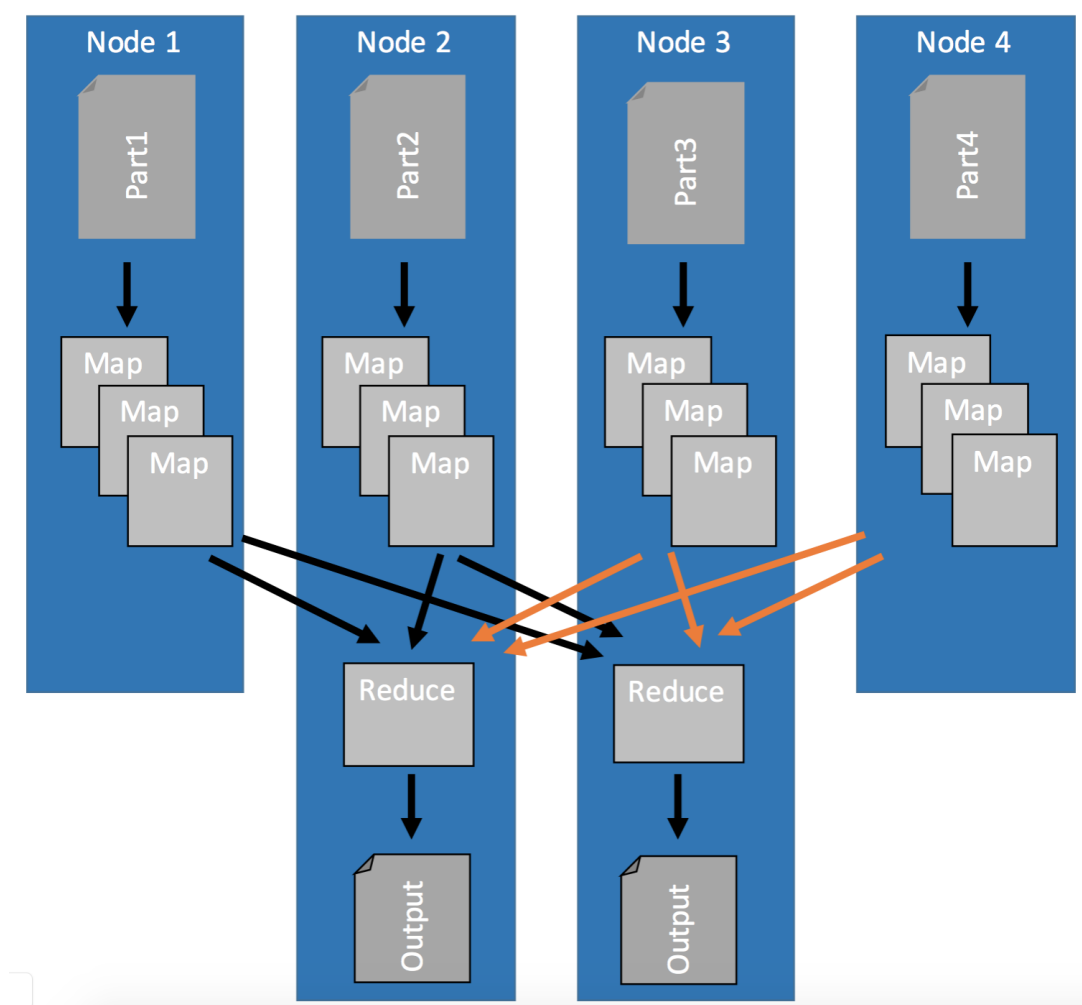


Figure 11: Mapping of input files and reducing of mapped records

### 9.1.1 How does MapReduce work – available

First of all, the whole map and reduce process relies heavily on using key/value pairs. That's what the mappers are for.

In the map phase input data, for instance a file, gets loaded and transformed into key/value pairs.

When each map phase is done it sends the created key/value pairs to the reducers where they are getting sorted by key. This means, that an input record for the reduce phase is a list of values from the mappers that all have the same key.

Then the reduce phase is doing the computation of that key and its values and outputting the results.

How many mappers and reducers can you use in parallel? The number of parallel map and reduce processes depends on how many CPU cores you have in your cluster. Every mapper and every reducer is using one core.

This means that the more CPU cores you actually have, the more mappers you can use, the faster the extraction process can be done. The more reducers you are using the faster the actual computation is being done.

To make this more clear, I have prepared an example:

**Example** As I said before, MapReduce works in two stages, map and reduce. Often these stages are explained with a word count task.

Personally, I hate this example because counting stuff is too trivial and does not really show you what you can do with MapReduce. Therefore, we are going to use a more real world use-case from the world of the internet of things (IoT).

IoT applications create an enormous amount of data that has to be processed. This data is generated by physical sensors who take measurements, like room temperature at 8.00 o’Clock.

Every measurement consists of a key (the timestamp when the measurement has been taken) and a value (the actual value measured by the sensor).

Because you usually have more than one sensor on your machine, or connected to your system, the key has to be a compound key. Compound keys contain additionally to the measurement time information about the source of the signal.

But, let’s forget about compound keys for now. Today we have only one sensor. Each measurement outputs key/value pairs like: Timestamp-Value.

The goal of this exercise is to create average daily values of that sensor’s data.

The image below shows how the map and reduce process works.

First, the map stage loads unsorted data (input records) from the source (e.g. HDFS) by key and value (key:2016-05-01 01:02:03, value:1).

Then, because the goal is to get daily averages, the hour:minute:second information is cut from the timestamp.

That is all that happens in the map phase, nothing more.

After all parallel map phases are done, each key/value pair gets sent to the one reducer who is handling all the values for this particular key.

Every reducer input record then has a list of values and you can calculate  $(1+5+9)/3$ ,  $(2+6+7)/3$  and  $(3+4+8)/3$ . That's all.

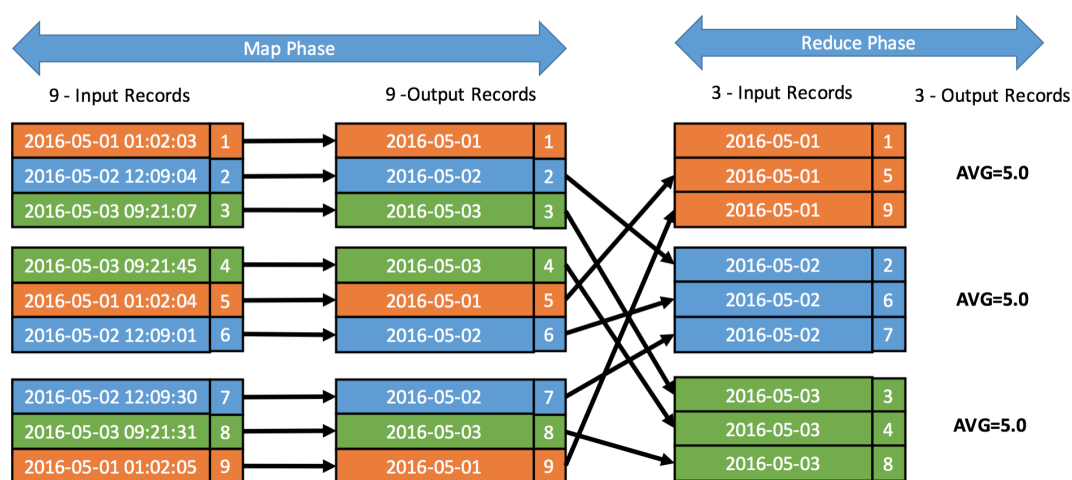


Figure 12: MapReduce Example of Time Series Data

What do you think you need to do to generate minute averages?

Yes, you need to cut the key differently. You then would need to cut it like this: “2016-05-01 01:02”. Keeping the Hour and minute information in the key.

What you can also see is, why map reduce is so great for doing parallel work. In this case, the map stage could be done by nine mappers in parallel because each map is independent from all the others.

The reduce stage could still be done by three tasks in parallel. One each for orange, blue and one for green.

That means, if your dataset would be 10 times as big and you'd have 10 times the machines, the time to do the calculation would be the same.

### 9.1.2 What is the limitation of MapReduce? – available

MapReduce is awesome for simpler analytics tasks, like counting stuff. It just has one flaw: It has only two stages Map and Reduce.



Figure 13: The Map Reduce Process

First MapReduce loads the data from HDFS into the mapping function. There you prepare the input data for the processing in the reducer. After the reduce is finished the results get written to the data store.

The problem with MapReduce is that there is no simple way to chain multiple map and reduce processes together. At the end of each reduce process the data must be stored somewhere.

This fact makes it very hard to do complicated analytics processes. You would need to chain MapReduce jobs together.

Chaining jobs with storing and loading intermediate results just makes no sense.

Another issue with MapReduce is that it is not capable of streaming analytics. Jobs take some time to spin up, do the analytics and shut down. Basically Minutes of wait time are totally normal.

This is a big negative point in a more and more real time data processing world.

## 9.2 Apache Spark

I talked about the three methods of data streaming in this podcast: <https://anchor.fm/andreaskayy/embeds/Methods-of-Streaming-Data-e15r6o>

### 9.2.1 Spark Basics

### 9.2.2 What is the difference to MapReduce? – available

Spark is a complete in memory framework. Data gets loaded from, for instance hdfs, into the memory of workers.

There is no longer a fixed map and reduce stage. Your code can be as complex as you want.

Once in memory, the input data and the intermediate results stay in memory (until the job finishes). They do not get written to a drive like with MapReduce.

This makes Spark the optimal choice for doing complex analytics. It allows you for instance to do Iterative processes. Modifying a dataset multiple times in order to create an output is totally easy.

Streaming analytics capability is also what makes spark so great. Spark has natively the option to schedule a job to run every X seconds or X milliseconds.

As a result, Spark can deliver you results from streaming data in “real time”.

### 9.2.3 How does Spark fit to Hadoop? – available

There are some very misleading articles out there titled Spark or Hadoop, Spark is better than Hadoop or even Spark is replacing Hadoop.

	Storage	Analytics	Resource Management
Hadoop	Hadoop Distributed File System HDFS	MapReduce	YARN (Yet Another Resource Negotiator)
Spark	--	Spark	Spark Resource Management

Figure 14: Hadoop vs Spark capabilities

So, it’s time to show you the differences between Spark and Hadoop. After this you will know when and for what you should use Spark and Hadoop.

You’ll also understand why Hadoop or Spark is the totally wrong question.

**Where’s the difference?** To make it clear how Hadoop differs from Spark I created this simple feature table:

Hadoop is used to store data in the Hadoop Distributed File System (HDFS). It can analyse the stored data with MapReduce and manage resources with YARN.

However, Hadoop is more than just storage, analytics and resource management. There’s a whole eco system of tools around the Hadoop core. I’ve written about tis eco system in this article: What is Hadoop and why is it so freakishly popular. You should check it out as well.

Compared to Hadoop, Spark is “just” an analytics framework. It has no storage capability. Although it has a standalone resource management, you usually don’t use that



feature.

**Spark and Hadoop is a perfect fit** So, if Hadoop and Spark are not the same things, can they work together?

Absolutely! Here's how the first picture will look if you combine Hadoop with Spark:

As Storage you use the Hadoop distributed file system. Analytics is done with Apache Spark and Yarn is taking care of the resource management.

Why does that work so well together?

From a platform architecture perspective, Hadoop and Spark are usually managed on the same cluster. This means on each server where a HDFS data node is running, a spark worker thread runs as well.

In distributed processing, network transfer between machines is a large bottle neck. Transferring data within a machine reduces this traffic significantly.

Spark is able to determine on which data node the needed data is stored. This allows a direct load of the data from the local storage into the memory of the machine.

This reduces network traffic a lot.

**As for YARN:** You need to make sure that your physical resources are distributed perfectly between the services. This is especially the case when you run Spark workers with other Hadoop services on the same machine.

It just would not make sense to have two resource managers managing the same server's resources. Sooner or later they will get in each others way.

That's why the Spark standalone resource manager is seldom used.

So, the question is not Spark or Hadoop. The question has to be: Should you use Spark or MapReduce alongside Hadoop's HDFS and YARN.

**My simple rule of thumb:** If you are doing simple batch jobs like counting values or doing calculating averages: Go with MapReduce.

If you need more complex analytics like machine learning or fast stream processing: Go with Apache Spark.

#### **9.2.4 Available Languages – available**

Spark jobs can be programmed in a variety of languages. That makes creating analytic processes very user-friendly for data scientists.

Spark supports Python, Scala and Java. With the help of SparkR you can even connect your R program to a Spark cluster.

If you are a data scientist who is very familiar with Python just use Python, its great. If you know how to code Java I suggest you start using Scala.

Spark jobs are easier to code in Scala than in Java. In Scala you can use anonymous functions to do processing.

This results in less overhead, it is a much cleaner, simpler code.

With Java 8 simplified function calls were introduced with lambda expressions. Still, a lot of people, including me prefer Scala over Java.

#### **9.2.5 How to do stream processing**

#### **9.2.6 How to do batch processing**

#### **9.2.7 How does Spark use data from Hadoop – available**

Another thing is data locality. I always make the point, that processing data locally where it is stored is the most efficient thing to do.

That's exactly what Spark is doing. You can and should run Spark workers directly on the data nodes of your Hadoop cluster.

Spark can then natively identify on what data node the needed data is stored. This enables Spark to use the worker running on the machine where the data is stored to load the data into the memory.

The downside of this setup is that you need more expensive servers. Because Spark processing needs stronger servers with more RAM and CPUs than a “pure” Hadoop setup.

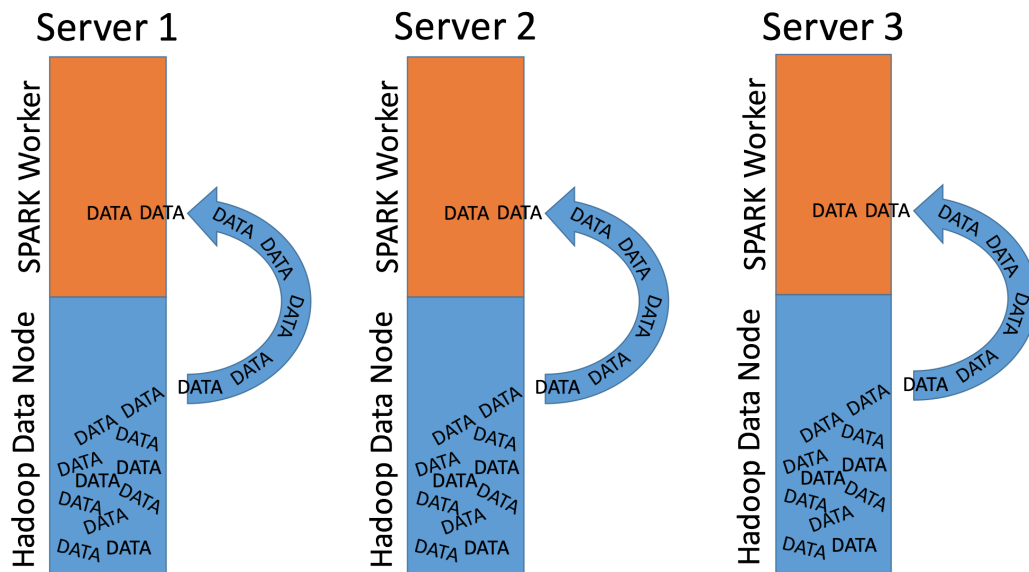


Figure 15: Spark Using Hadoop Data Locality

### 9.2.8 What is a RDD and what is a DataFrame?

### 9.2.9 Spark coding with Scala

### 9.2.10 Spark coding with Python

### 9.2.11 How and why to use SparkSQL?

### 9.2.12 Machine Learning on Spark? (Tensor Flow)

**MLlib:** The machine learning library MLlib is included in Spark so there is often no need to import another library.

I have to admit because I am not a data scientist I am not an expert in machine learning.

From what I have seen and read though the machine learning framework MLlib is a nice treat for data scientists wanting to train and apply models with Spark.

### Tensorflow

### 9.2.13 Spark Setup – available

From a solution architect's point of view Spark is a perfect fit for Hadoop big data platforms. This has a lot to do with cluster deployment and management.

Companies like Cloudera, MapR or Hortonworks include Spark into their Hadoop distributions. Because of that, Spark can be deployed and managed with the clusters Hadoop management web fronted.

This makes the process for deploying and configuring a Spark cluster very quick and admin friendly.

#### 9.2.14 Spark Resource Management – available

When running a computing framework you need resources to do computation: CPU time, RAM, I/O and so on. Out of the box Spark can manage resources with it's stand-alone resource manager.

If Spark is running in an Hadoop environment you don't have to use Spark's own stand-alone resource manager. You can configure Spark to use Hadoop's YARN resource management.

Why would you do that? It allows YARN to efficiently allocate resources to your Hadoop and Spark processes.

Having a single resource manager instead of two independent ones makes it a lot easier to configure the resource management.

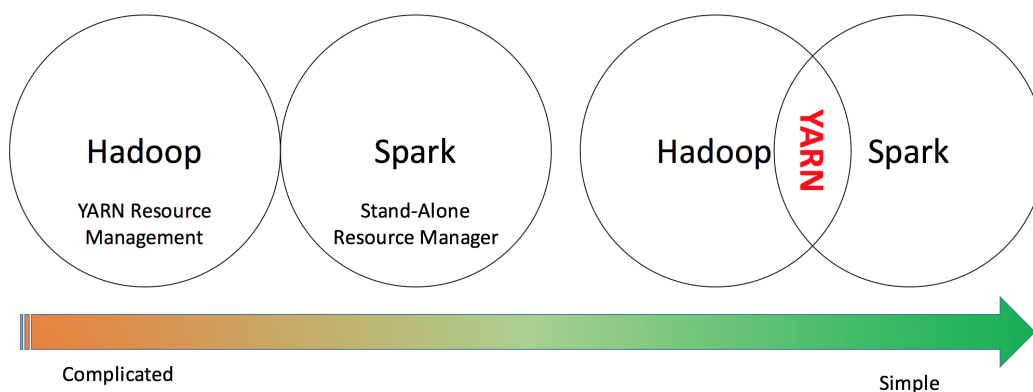


Figure 16: Spark Resource Management With YARN

## 9.3 Message queues with Apache Kafka

### 9.3.1 Why a message queue tool?

### 9.3.2 Kafka architecture

### 9.3.3 What are topics

### 9.3.4 What does Zookeeper have to do with Kafka

### 9.3.5 How to produce and consume messages

My YouTube video how to set up Kafka at home: <https://youtu.be/7F9tBwTUSEY>

My YouTube video how to write to Kafka: <https://youtu.be/RboQBZvZCh0>

## 9.4 Machine Learning

My podcast about how to do machine learning in production: <https://anchor.fm/andreaskayy/embed/episode-1-Learning-In-Production-e11bbk>

### 9.4.1 Training and Applying models

### 9.4.2 What is deep learning

### 9.4.3 How to do Machine Learning in production — available

Machine learning in production is using stream and batch processing. In the batch processing layer you are creating the models, because you have all the data available for training.

In the stream in processing layer you are using the created models, you are applying them to new data.

The idea that you need to incorporate is this is a constant a constant cycle. Training, applying, re-training, pushing into production and applying.

What you don't want to do is, you don't want to do this manually. You need to figure out a process of automatic retraining and automatic pushing to into production of models.

In the retraining phase the system automatically evaluates the training. If the model no longer fits it works as long as it needs to create a good model.

After the evaluation of the model is complete and it's good, the model gets pushed into production. Into the stream processing.

#### **9.4.4 Why machine learning in production is harder then you think – available**

How to automate machine learning is something that drives me day in and day out.

What you do in development or education is, that you create a model and fit it to the data. Then that model is basically done forever.

Where I'm coming from, the IoT world, the problem is that machines are very different. They behave very different and experience wear.

**Models Do Not Work Forever** Machines have certain processes that decrease the actual health of the machine. Machine wear is a huge issue. Models that that are built on top of a good machine don't work forever.

When the Machine wears out, the models need to be adjusted. They need to be maintained, retrained.

**Where The Platforms That Support This?** Automatic re-training and re-deploying is a very big issue, a very big problem for a lot of companies. Because most existing platforms don't have this capability (I actually haven't seen one until now).

Look at AWS machine learning for instance. The process is: build, train, tune deploy. Where's the loop of retraining?

You can create models and then use them in production. But this loop is almost nowhere to be seen.

It is a very big issue that needs to be solved. If you want to do machine learning in production you can start with manual interaction of the training, but at some point you need to automate everything.

**Training Parameter Management** To train a model you are manipulating input parameters of the models.

Take deep learning for instance. To train you are manipulating for instance:

How many layers do you use. The depth of the layers, which means how many neurons you have in a layer. What activation function you use, how long are you training and so on.

You also need to keep track of what data you used to train which model.

All those parameters need to be manipulated automatically, models trained and tested.

To do all that, you basically need a database that keeps track of those variables.

How to automate this, for me, is like the big secret. I am still working on figuring it out.

**What's Your Solution?** Did you already have the problem of automatic re-training and deploying of models as well?

Were you able to use a cloud platform like Google, AWS or Azure?

It would be really awesome if you share your experience :)

#### 9.4.5 How to convince people machine learning works — available

Many people still are not convinced that machine learning works reliably. But they want analytics insight and most of the time machine learning is the way to go.

This means, when you are working with customers you need to do a lot of convincing. Especially if they are not into machine learning themselves.

But it's actually quite easy.

**No Rules, No Physical Models** Many people are still under the impression that analytics only works when it's based on physics. When there are strict mathematical rules to a problem.

Especially in engineering heavy countries like Germany this is the norm:

“Sere has to be a Rule for Everyising!” (imagine a German accent) When you're engineering you are calculating stuff based on physics and not based on data. If you are constructing an airplane wing, you better make sure to use calculations so it doesn't fall off.

And that's totally fine.

Keep doing that!

Machine learning has been around for decades. It didn't quite work as good as people hoped. We have to admit that. But there is this preconception that it still doesn't work.

Which is not true: Machine learning works.

Somehow you need to convince people that it is a viable approach. That learning from data to make predictions is working perfectly.

**You Have The Data. USE IT!** As a data scientist you have one ace up your sleeve, it's the obvious one:

It's the data and it's statistics.

You can use that data and those statistics to counter peoples preconceptions. It's very powerful if someone says: "This doesn't work"

You bring the data. You show the statistics and you show that it works reliably.

A lot of discussions end there.

Data doesn't lie. You can't fight data. The data is always right.

**Data is Stronger Than Opinions** This is also why I believe that autonomous driving will come quicker than many of us think. Because a lot of people say, they are not safe. That you cannot rely on those cars.

The thing is: When you have the data you can do the statistics.

You can show people that autonomous driving really works reliably. You will see, the question of: Is this is this allowed or is this not allowed? Will be gone quicker than you think.

Because goverenment agencies can start testing the algorithms based on predefined scenarios. They can run benchmarks and score the cars performance.

All those opinions, if it works, or if it doesn't work, they will be gone.

The motor agency has the statistics. The stats show people how good cars work.

Companies like Tesla, they have it very easy. Because the data is already there.

**They just need to show us that the algorithms work. The end.**



## 10 How to Store Data

Check out my podcast how to decide between SQL and NoSQL: <https://anchor.fm/andreaskayy/embed/Vs-SQL-How-To-Choose-e12f1o>

### 10.1 Data Modeling

10.1.1 How to find out how you need to store data for the business case

10.1.2 How to decide what kind of storage you need to use

### 10.2 SQL

10.2.1 Database Design

10.2.2 SQL Queries

10.2.3 Stored Procedures

10.2.4 ODBC/JDBC Server Connections

### 10.3 NoSQL

10.3.1 Key/Value Stores (HBase)

10.3.2 Document Store HDFS — available

The Hadoop distributed file system, or HDFS, allows you to store files in Hadoop. The difference between HDFS and other file systems like NTFS or EXT is that it is a distributed one.

What does that mean exactly?

A typical file system stores your data on the actual hard drive. It is hardware dependent.

If you have two disks then you need to format every disk with its own file system. They are completely separate.

You then decide on which disk you physically store your data.

HDFS works different to a typical file system. HDFS is hardware independent.

Not only does it span over many disks in a server. It also spans over many servers.

HDFS will automatically place your files somewhere in the Hadoop server collective.

It will not only store your file, Hadoop will also replicate it two or three times (you can define that). Replication means replicas of the file will be distributed to different servers.

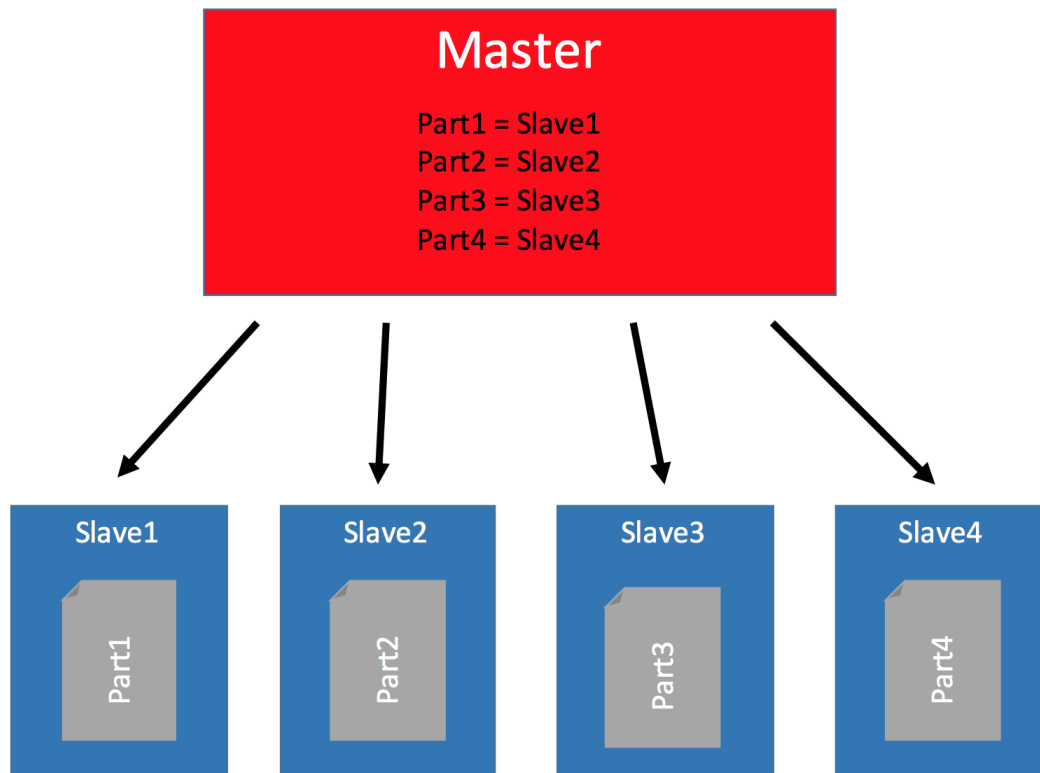


Figure 17: HDFS Master and Data Nodes

This gives you superior fault tolerance. If one server goes down, then your data stays available on a different server.

Another great thing about HDFS is, that there is no limit how big the files can be. You can have server log files that are terabytes big.

How can files get so big? HDFS allows you to append data to files. Therefore, you can continuously dump data into a single file without worries.

HDFS physically stores files different than a normal file system. It splits the file into blocks.

These blocks are then distributed and replicated on the Hadoop cluster. The splitting happens automatically.

In the configuration you can define how big the blocks should be. 128 megabyte or 1 gigabyte?

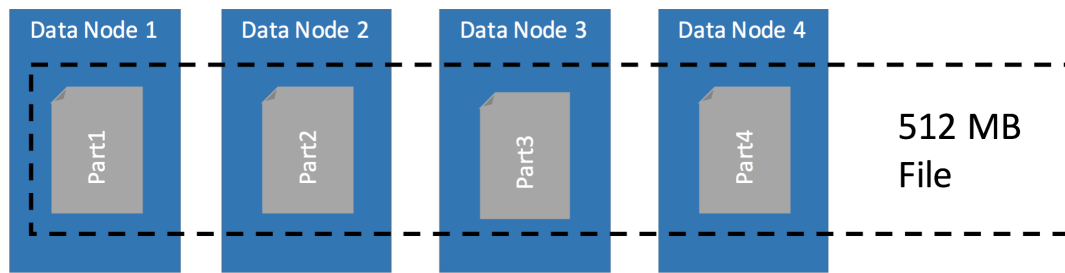


Figure 18: Distribution of Blocks for a 512MB File

No problem at all.

This mechanic of splitting a large file in blocks and distributing them over the servers is great for processing. See the MapReduce section for an example.

### 10.3.3 Document Store MongoDB

### 10.3.4 Hive Warehouse

### 10.3.5 Impala

### 10.3.6 Time Series Databases (?)

DalmatinerDB InfluxDB Prometheus Riak TS OpenTSDB KairosDB Elasticsearch Druid

### 10.3.7 MPP Databases (Greenplum)

## 11 How to Visualize Data

### 11.1 Mobile Apps

#### 11.1.1 Android & IOS basics

#### 11.1.2 How to design APIs for mobile apps

### 11.2 How to use Webservers to display content

This section does not contain any text that's why the page is messed up

### **11.2.1 Tomcat**

### **11.2.2 Jetty**

### **11.2.3 NodeRED**

### **11.2.4 React**

## **11.3 Business Intelligence Tools**

### **11.3.1 Tableau**

### **11.3.2 PowerBI**

### **11.3.3 Quliksense**

## **11.4 Identity & Device Management**

### **11.4.1 What is a digital twin?**

### **11.4.2 Active Directory**

## **12 Case Studies**

### **12.1 7 Steps to Successfull Data Science Project**

### **12.2 Data Science @Airbnb**

### **12.3 Data Science @Netflix – available**

Netflix revolutionized how we watch movies and tv. Currently over 75 million users watch 125 million hours of Netflix content every day!

Netflix's revenue comes from a monthly subscription service. So, the goal for Netflix is to keep you subscribed and to get new subscribers.

To achieve this, Netflix is licensing movies from studios as well as creating its own original movies and tv series.

But offering new content is not everything. What is also very important is, to keep you watching content that already exists.

To be able to recommend you content, Netflix is collecting data from users. And it is collecting a lot.

Currently, Netflix analyses about 500 billion user events per day. That results in a stunning 1.3 petabytes every day.

All this data allows Netflix to build recommender systems for you. The recommenders are showing you content that you might like, based on your viewing habits, or what is currently trending.

**The Netflix batch processing pipeline** When Netflix started out, they had a very simple batch processing system architecture.

The key components were Chuckwa, a scalable data collection system, Amazon S3 and Elastic MapReduce.

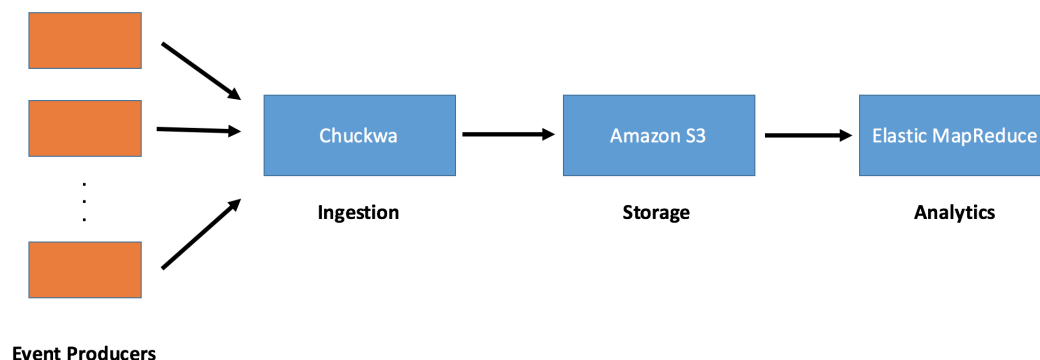


Figure 19: Old Netflix Batch Processing Pipeline

Chuckwa wrote incoming messages into Hadoop sequence files, stored in Amazon S3. These files then could be analysed by Elastic MapReduce jobs.

Netflix batch processing pipeline Jobs were executed regularly on a daily and hourly basis. As a result, Netflix could learn how people used the services every hour or once a day.

**Know what customers want:** Because you are looking at the big picture you can create new products. Netflix uses insight from big data to create new tv shows and movies.

They created House of Cards based on data. There is a very interesting Ted talk about this you should watch:

How to use data to make a hit TV show — Sebastian Wernicke:

Batch processing also helps Netflix to know the exact episode of a TV show that gets you hooked. Not only globally but for every country where Netflix is available.

Check out the article from TheVerge

They know exactly what show works in what country and what show does not.

It helps them create shows that work in everywhere or select the shows to license in different countries. Germany for instance does not have the full library that Americans have :(

We have to put up with only a small portion of tv shows and movies. If you have to select, why not select those that work best.

**Batch processing is not enough** As a data platform for generating insight the Cuckwa pipeline was a good start. It is very important to be able to create hourly and daily aggregated views for user behavior.

To this day Netflix is still doing a lot of batch processing jobs.

The only problem is: With batch processing you are basically looking into the past.

For Netflix, and data driven companies in general, looking into the past is not enough. They want a live view of what is happening.

**The trending now feature** One of the newer Netflix features is “Trending now”. To the average user it looks like that “Trending Now” means currently most watched.

This is what I get displayed as trending while I am writing this on a Saturday morning at 8:00 in Germany. But it is so much more.

What is currently being watched is only a part of the data that is used to generate “Trending Now”.

“Trending now” is created based on two types of data sources: Play events and Impression events.

What messages those two types actually include is not really communicated by Netflix. I did some research on the Netflix Techblog and this is what I found out:



Figure 20: Netflix Trending Now Feature

Play events include what title you have watched last, where you did stop watching, where you used the 30s rewind and others. Impression events are collected as you browse the Netflix Library like scroll up and down, scroll left or right, click on a movie and so on

Basically, play events log what you do while you are watching. Impression events are capturing what you do on Netflix, while you are not watching something.

**Netflix real-time streaming architecture** Netflix uses three internet facing services to exchange data with the client's browser or mobile app. These services are simple Apache Tomcat based web services.

The service for receiving play events is called "Viewing History". Impression events are collected with the "Beacon" service.

The "Recommender Service" makes recommendations based on trend data available for clients.

Messages from the Beacon and Viewing History services are put into Apache Kafka. It acts as a buffer between the data services and the analytics.

Beacon and Viewing History publish messages to Kafka topics. The analytics subscribes to the topics and gets the messages automatically delivered in a first in first out fashion.

After the analytics the workflow is straight forward. The trending data is stored in a Cassandra Key-Value store. The recommender service has access to Cassandra and is making the data available to the Netflix client.

The algorithms how the analytics system is processing all this data is not known to the public. It is a trade secret of Netflix.

What is known, is the analytics tool they use. Back in Feb. 2015 they wrote in the tech blog that they use a custom made tool.

They also stated, that Netflix is going to replace the custom made analytics tool with Apache Spark streaming in the future. My guess is, that they did the switch to Spark some time ago, because their post is more than a year old.

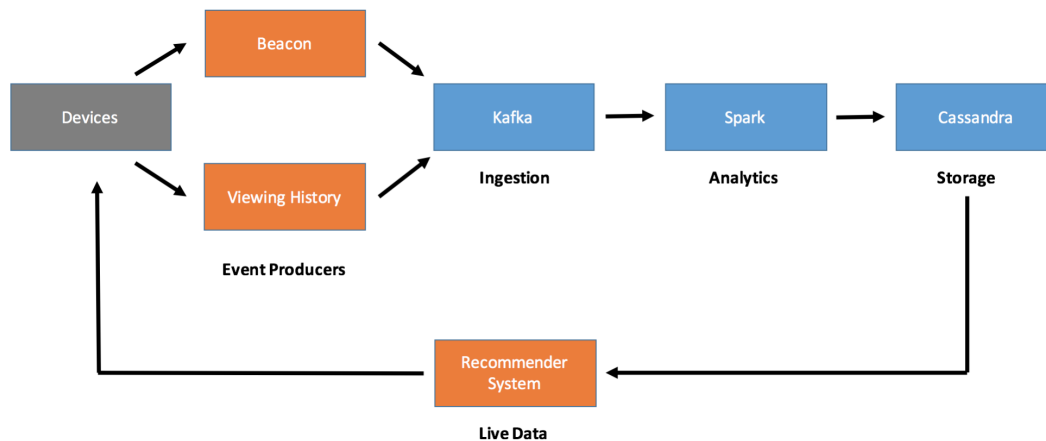


Figure 21: Netflix Streaming Pipeline

## 12.4 Data Science @Uber

## 12.5 Data Science @Zalando

# 13 DEV/OPS

## 13.1 Hadoop

Listen to an introduction to Hadoop for Data Scientists: [1]

### 13.1.1 Hadoop Cluster setup and management with Cloudera Manager (for example)

### 13.1.2 Spark code from coding to production

### 13.1.3 How to monitor and manage data processing pipelines

### 13.1.4 Oozie

### 13.1.5 Airflow Application management

### 13.1.6 Creating Statistics with Spark and Kafka



## References

- [1] J. Ely and I. Stavrov, *Analyzing chalk dust and writing speeds: computational and geometric approaches*, BoDine Journal of Mathematics **3** (2001), 14-159.

## List of Figures

1	The Machine Learning Pipeline . . . . .	8
2	Platform Blueprint . . . . .	16
3	Batch Processing Pipeline . . . . .	19
4	Stream Processing Pipeline . . . . .	20
5	Common SQL Platform Architecture . . . . .	25
6	Scaling up a SQL Database . . . . .	26
7	Scaling out a SQL Database . . . . .	27
8	Hadoop Ecosystem Components . . . . .	29
9	Connections between tools . . . . .	30
10	Flume Integration . . . . .	31
11	Mapping of input files and reducing of mapped records . . . . .	36
12	MapReduce Example of Time Series Data . . . . .	38
13	The Map Reduce Process . . . . .	39
14	Hadoop vs Spark capabilities . . . . .	40
15	Spark Using Hadoop Data Locality . . . . .	43
16	Spark Resource Management With YARN . . . . .	44
17	HDFS Master and Data Nodes . . . . .	50
18	Distribution of Blocks for a 512MB File . . . . .	51
19	Old Netflix Batch Processing Pipeline . . . . .	53
20	Netflix Trending Now Feature . . . . .	55
21	Netflix Streaming Pipeline . . . . .	56

## List of Tables