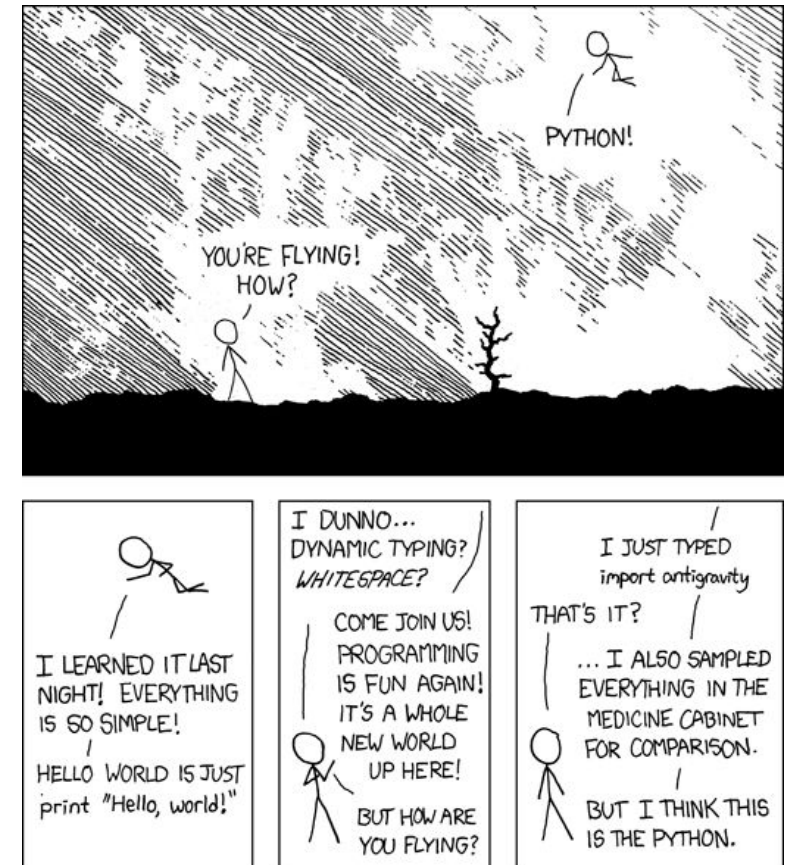# Python for Big Data Analysis

**Kim Hee**
**Graduate research assistant**
**at Heinrich-Lanz-Center (HLZ) for Digital Health**

# About us

- **Kim Hee  -  PhD candidate since 2016**
  - Computer engineering background
  - Research interest: microbiology image data analysis with consideration of scalability
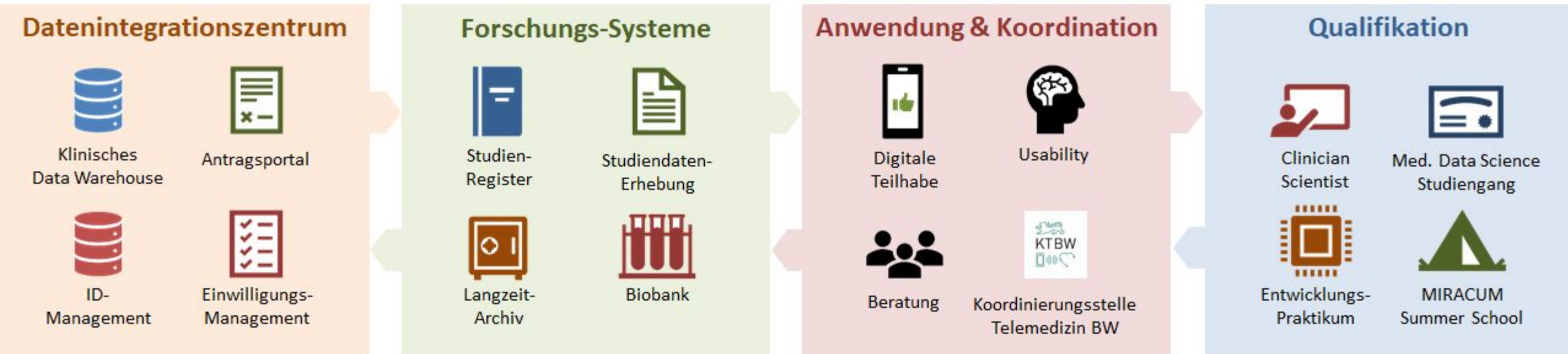  - Responsibility: ETL for laboratory data



Heinrich-Lanz-Zentrum für Digitale Gesundheit

# Heinrich-Lanz-Center (HLZ) for Digital Health

- **Our lab is currently active in the following areas:**
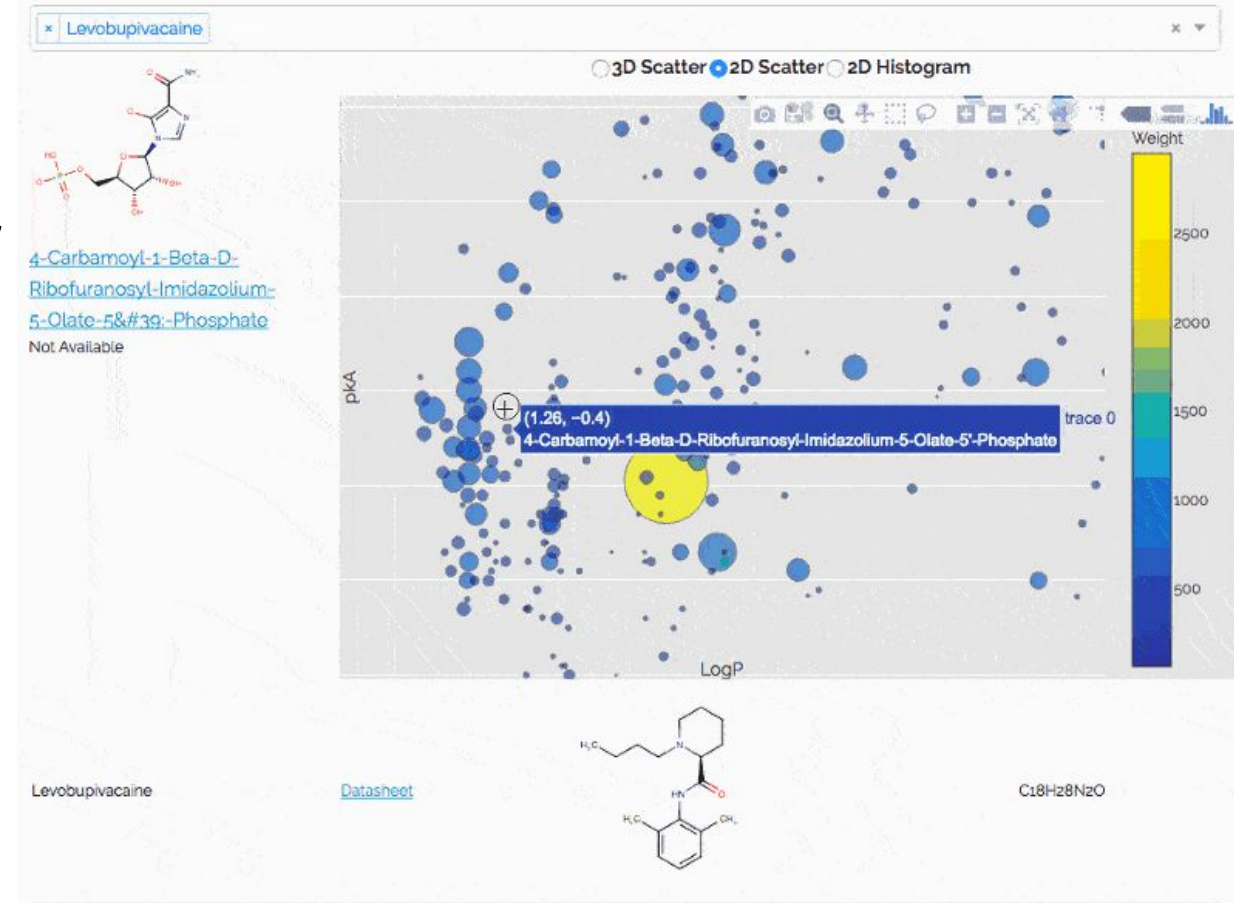
# Challenges in digital healthcare domain

- **Data silos across and within the healthcare providers**
- **Premature status of interoperability and data standards**
- **Data quality is low (ex: completeness, consistency, accuracy, timeliness and more)**
- **In-house database management system is not ready for data analysis**
- **Not many successful healthcare applications for patients**
- **Not enough study programs**
- **Legal issues**
- **and more**

# In-house database management system is for neither cyclic data analyse nor visualization

- **In-house database is built for the billing purpose**
- **Data analysis with a conventional database management system (CDBMS) is possible, but not optimal because too many I/Os occur and memory is not properly utilized for data analysis**
- **Data visualization with SQL is not possible**

# Tutorial organization

| Time | Topic |
| --- | --- |
| **09:00 - 09:30** | **Set the scene** |
| 09:30 - 10:45 | Data analysis and Data visualization |
| 10:45 - 11:15 | Break |
| 11:15 - 12:15 | Machine Learning |
| 12:15 - 13:00 | Python at Scale (PySpark) |

UMM
UNIVERSITÄTSMEDIZIN
MANNHEIM

Medizinische Fakultät Mannheim
der Universität Heidelberg
Universitätsklinikum Mannheim

# Agenda

- Big Data and Data Science
- Why Python for Data Science?
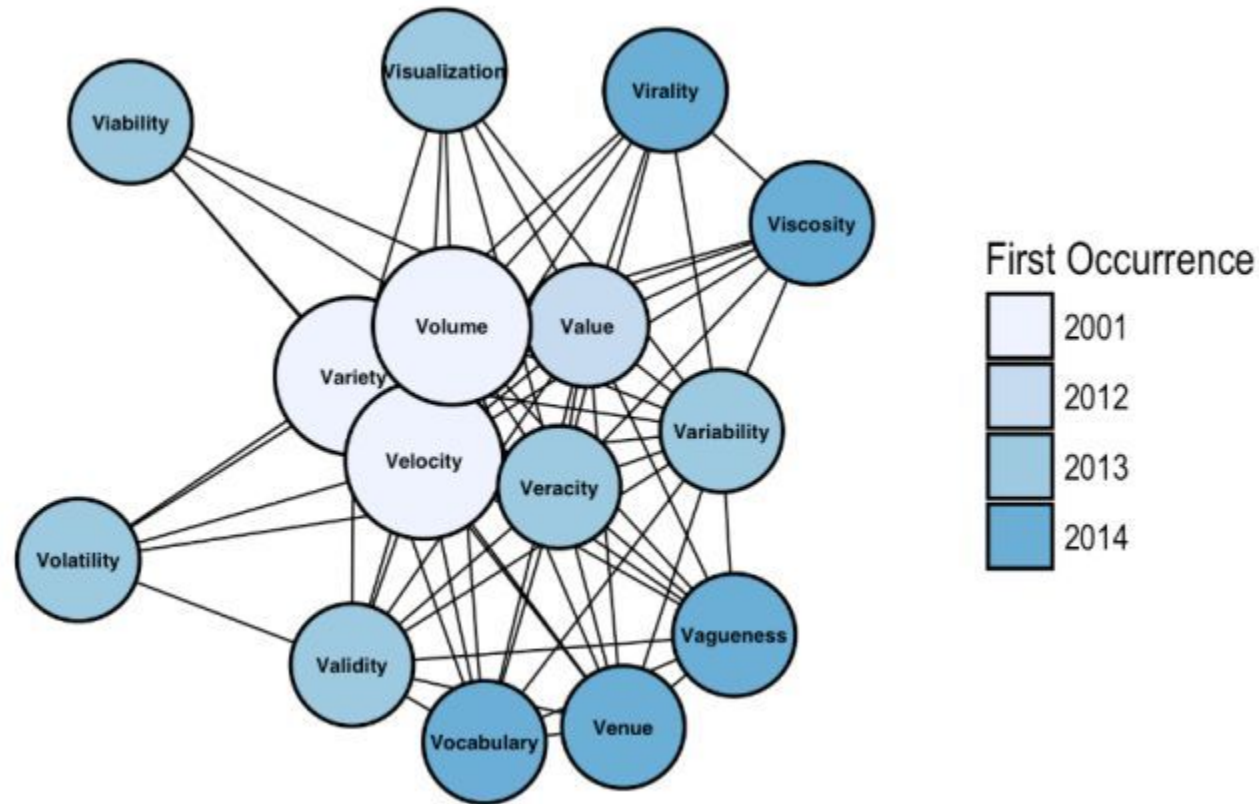- Hands-on: Jupyter Notebook Launch

# What is Big Data?

- Big data is NOT defined in terms of data size
- Big data refers to datasets that a conventional database software is not capable to capture, store, manage, process, analyze, and understand
- Those datasets are large, fast, complex, and not structured
- 3Vs (volume, variety and velocity) are three properties of big data introduced by D. Laney, "3D data management: Controlling data volume, velocity and variety," Appl. Deliv. Strateg. File, vol. 949, 2001.
  - **Volume** represents the ever - growing amount of data and challenges the current stage of storage systems.
  - **Velocity** describes how quickly the data is retrieved, stored and processed.
  - **Variety** describes the multitude of data sources like sensors, smart devices and social media often in unstructured data formats.

# The inflation of Big Data Vs (2017)

- These 42 Vs are orbiting the original three



42 V's of Big Data and Data Science visualized by Tom Shafer (2017)
https://www.elderresearch.com/blog/42-v-of-big-data

# What is Data Science?

- An informal definition of data scientist by Josh Wills:



**Josh Wills** @josh_wills

Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.

RETWEETS 1,046  FAVORITES 532

6:55 PM - 3 May 2012

# What is Data Science?

- Data science is a multidisciplinary domain
- Data science aims to extract knowledge from data with a set of tools, scientific methods and scalable systems

UMM
UNIVERSITÄTSMEDIZIN
MANNHEIM

Medizinische Fakultät Mannheim
der Universität Heidelberg
Universitätsklinikum Mannheim

# Data science educational programs in EU (2016)
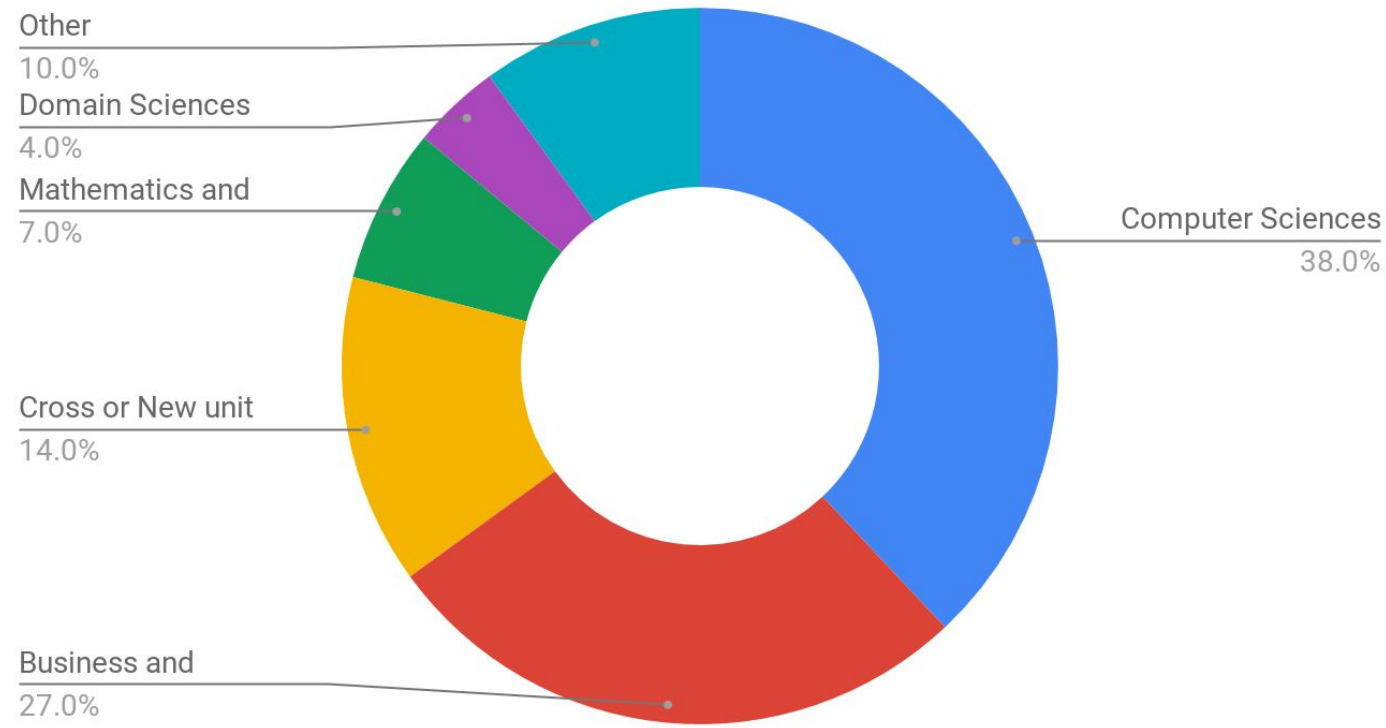
- Data Science is hard to teach due to the nature of the multidisciplinary
- ...and they are teaching different things...



| | |
|---|---|
| Other | 10.0% |
| Domain Sciences | 4.0% |
| Mathematics and | 7.0% |
| Computer Sciences | 38.0% |
| Cross or New unit | 14.0% |
| Business and | 27.0% |

EDISION survey result in 2016: Data Science educational programs providers
(https://cordis.europa.eu/project/rcn/198292/results/en)

# CRISP-DM: the most used methodology for a data science project

- CRISP-DM is an open standard process model that describes common approaches used by data mining experts.
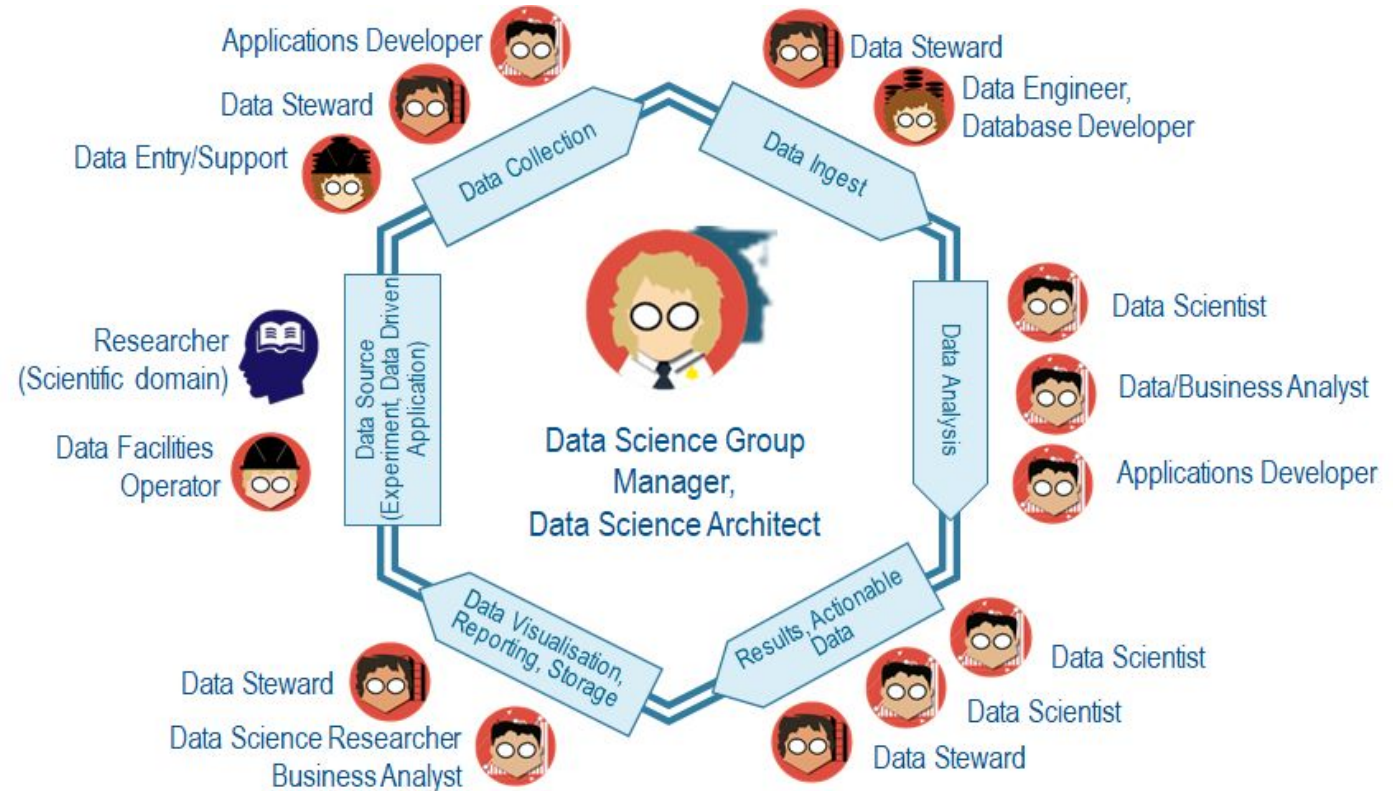- developed by DamilerChrysler, SPSS, NCR and OHRA and released in 1999

| Business Understanding | Data Understanding | Data Preparation | Modeling | Evaluation | Deployment |
|---|---|---|---|---|---|
| determine business objective | collect initial data | select data | select modeling technique | evaluate results | plan deployment |
| access situation | describe data | clean data | generate test design | approved models | plan monitoring and maintenance |
| determine data mining goal | explore data | construct data | build model parameter settings | determine next steps | produce final report |
| produce project plan | verify data quality | integrate data | access model | | review project |
| | | format data | | | |

which steps are unique compared to a research project?

# Data science team needs not only a data scientist

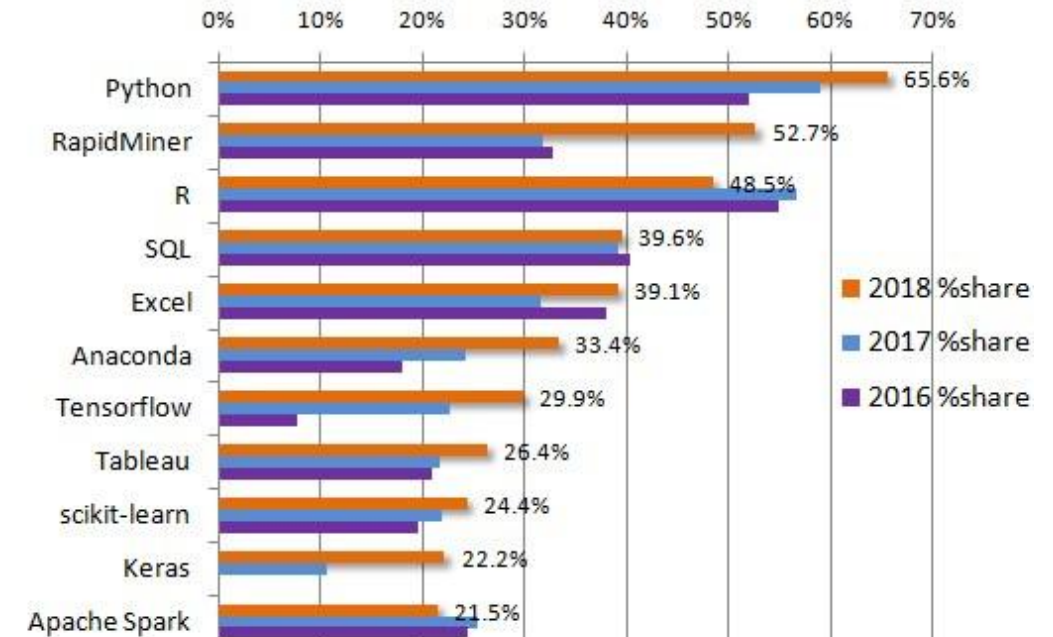| Profile Group | Profile Title |
|---|---|
| Data Science Managers | Data Science (group) Manager |
| | Data Science Infrastructure Manager |
| | Research Infrastructure Manager |
| Data Science Professionals | Data Scientist |
| | Data Science Researcher |
| | Data Science Architect |
| | Data Science Programmer/Engineer |
| | Data Analyst |
| | Business Analyst |
| Data handling Professionals | Data Stewards |
| | Digital data curator |
| | Digital Librarians |
| | Data Archivists |
| Database Professionals | Large scale database designer |
| | Large scale database admin |
| | Scientific database administrator |
| Technicians and associate profession | Big Data facilities Operator |
| | Large scale data storage operator |
| | Scientific database operator |

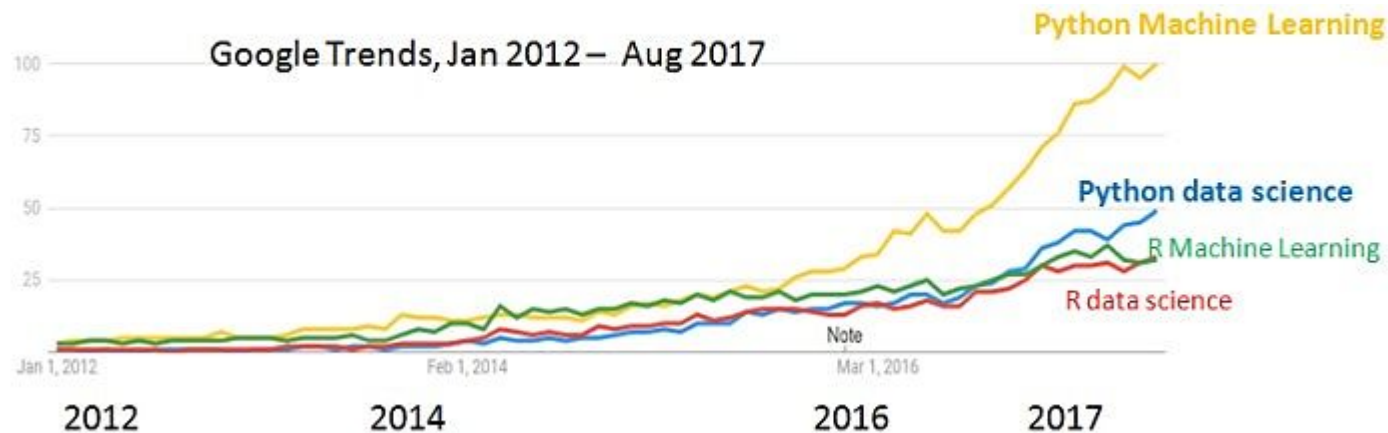# Agenda

- Big Data and Data Science
- Why Python for Data Science?
- Hands-on: Jupyter Notebook Launch

UMM
UNIVERSITÄTSMEDIZIN
MANNHEIM

Medizinische Fakultät Mannheim
der Universität Heidelberg

Universitätsklinikum Mannheim

# Python is the most popular language among data scientists

" Python remains the big kahuna,
but specialist languages hold their own "

Stephen Cass (2019), *The Top Programming Languages 2019*, IEEE Spectrum

# Python is highly compatible language in big data world



DATA & AI LANDSCAPE 2019

FIRSTMARK
EARLY STAGE VENTURE CAPITAL

# Python is highly compatible language in big data world



DATA & AI LANDSCAPE 2019

# Why is the compatibility matter?
## $1 million solution has never been used

- Netflix held the Netflix Prize open competition [1] for the best algorithm to predict user ratings for films.
- On September 21, 2009 Netflix awarded the $1M Grand Prize to team "BellKor's Pragmatic Chaos".
- The solution improved the recommendation algorithm by 10%. **But Netflix never implemented that solution itself.**
- According to the Netflix blog post [2]:

> *" We evaluated some of the new methods offline but* ***the additional accuracy gains that we measured did not seem to justify the engineering effort*** *needed to bring them into a production environment.* "

[1] https://www.kaggle.com/netflix-inc/netflix-prize-data
[2] https://medium.com/netflix-techblog/netflix-recommendations-beyond-the-5-stars-part-1-55838468f429

UMM
UNIVERSITÄTSMEDIZIN
MANNHEIM

Medizinische Fakultät Mannheim
der Universität Heidelberg

Universitätsklinikum Mannheim

# Compatibility Evaluation in Big Data stack (1/3)
## among data science layer

- Python is the first-class citizen in the deep learning/AI world



Ivanov, T., & Singhal, R.. (2018). *ABench: Big Data Architecture Stack Benchmark.*
Paper presented at the Companion of the 2018 ACM/SPEC International Conference
on Performance Engineering, ICPE 2018, Berlin, Germany, April 09 - 13, 2018.

|  | Python | R | Other APIs |
|---|---|---|---|
| TensorFlow | Yes | No | c++, Java, Go, Swift |
| Keras | Yes | Yes | No |
| MXNet | Yes | Yes | c++, Scala, Julia, Perl |
| PyTorch | Yes | No | No |
| CNTK | Yes | No | c++ |

compatibility table (https://goo.gl/na3DzY)

# Compatibility Evaluation in Big Data stack (1/3)
## among data science layer

## Java MapReduce

```
Source Code
                                                                              WordCount.java
1.   package org.myorg;
2.
3.   import java.io.IOException;
4.   import java.util.*;
5.
6.   import org.apache.hadoop.fs.Path;
7.   import org.apache.hadoop.conf.*;
8.   import org.apache.hadoop.io.*;
9.   import org.apache.hadoop.mapred.*;
10.  import org.apache.hadoop.util.*;
11.
12.  public class WordCount {
13.
14.    public static class Map extends MapReduceBase implements Mapper<LongWritable, Text, Text, IntWritable> {
15.      private final static IntWritable one = new IntWritable(1);
16.      private Text word = new Text();
17.
18.      public void map(LongWritable key, Text value, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
19.        String line = value.toString();
20.        StringTokenizer tokenizer = new StringTokenizer(line);
21.        while (tokenizer.hasMoreTokens()) {
22.          word.set(tokenizer.nextToken());
23.          output.collect(word, one);
24.        }
25.      }
26.    }
27.
28.    public static class Reduce extends MapReduceBase implements Reducer<Text, IntWritable, Text, IntWritable> {
29.      public void reduce(Text key, Iterator<IntWritable> values, OutputCollector<Text, IntWritable> output, Reporter reporter) throws IOException {
30.        int sum = 0;
31.        while (values.hasNext()) {
32.          sum += values.next().get();
33.        }
34.        output.collect(key, new IntWritable(sum));
35.      }
36.    }
37.
38.    public static void main(String[] args) throws Exception {
39.      JobConf conf = new JobConf(WordCount.class);
40.      conf.setJobName("wordcount");
41.
42.      conf.setOutputKeyClass(Text.class);
43.      conf.setOutputValueClass(IntWritable.class);
44.
45.      conf.setMapperClass(Map.class);
46.      conf.setCombinerClass(Reduce.class);
47.      conf.setReducerClass(Reduce.class);
48.
49.      conf.setInputFormat(TextInputFormat.class);
50.      conf.setOutputFormat(TextOutputFormat.class);
51.
52.      FileInputFormat.setInputPaths(conf, new Path(args[0]));
53.      FileOutputFormat.setOutputPath(conf, new Path(args[1]));
54.
55.      JobClient.runJob(conf);
56.    }
57.  }
58. }
59.
```

## Spark's Python API

```python
text_file = spark.textFile("hdfs://...")

text_file.flatMap(lambda line: line.split())
    .map(lambda word: (word, 1))
    .reduceByKey(lambda a, b: a+b)
```
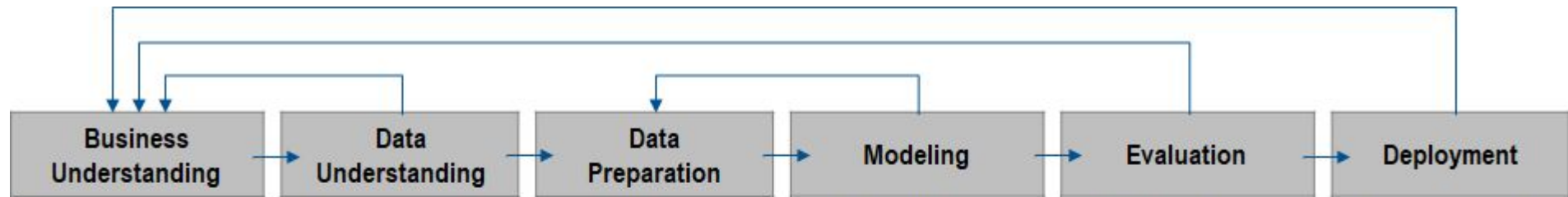
source: https://generalassemb.ly/data/data-science/spark

# Compatibility Evaluation in Big Data stack (1/3)
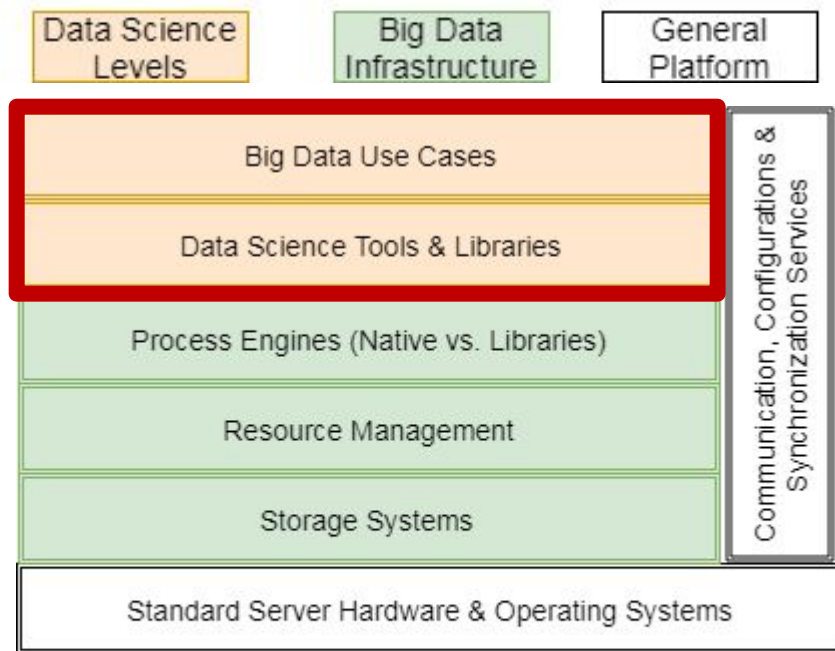## among data science layer

- **The requirements for a data science project cannot be done up-front.**
- **With data science, you learn as you go, not before you go.** You must be learned through experimentation, trial and error, and iteration.



CRISP-DM methodology

# Compatibility Evaluation in Big Data stack (2/3)
## application layer vs. data science layer



| Data Science Levels | Big Data Infrastructure | General Platform |
|---|---|---|

**Big Data Use Cases**

**Data Science Tools & Libraries**

Process Engines (Native vs. Libraries)

Resource Management

Storage Systems

Communication, Configurations & Synchronization Services

Standard Server Hardware & Operating Systems

Ivanov, T., & Singhal, R.. (2018). *ABench: Big Data Architecture Stack Benchmark*. Paper presented at the Companion of the 2018 ACM/SPEC International Conference on Performance Engineering, ICPE 2018, Berlin, Germany, April 09 - 13, 2018.

Popular web programming languages
https://t.ly/VAYIP

" costs scale super-linearly with the number of people involved. "

Eric Colson (2019). *Why Data Science Teams Need Generalists, Not Specialists.* Harvard Business Review

" The number of relationships (r) grows as a function number of members (n) per this equation: $r = (n^2-n) / 2$. And, each relationship bares coordination costs "

Hackman, J. R., & Hackman, R. J. (2002). *Leading teams: Setting the stage for great performances*. Harvard Business Press.

UMM
UNIVERSITÄTSMEDIZIN MANNHEIM

Medizinische Fakultät Mannheim der Universität Heidelberg
Universitätsklinikum Mannheim

# Compatibility Evaluation in Big Data stack (3/3)
## data science layer vs. process engine layer

- Two layers of data science tools and process engines have a high compatibility
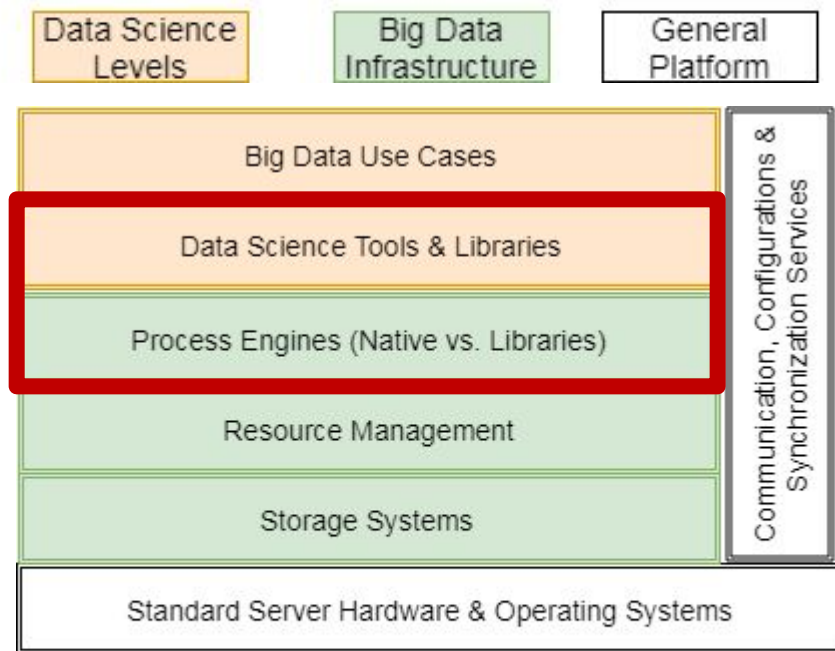


Ivanov, T., & Singhal, R.. (2018). *ABench: Big Data Architecture Stack Benchmark.*
Paper presented at the Companion of the 2018 ACM/SPEC International Conference
on Performance Engineering, ICPE 2018, Berlin, Germany, April 09 - 13, 2018.

|  | Java | Scala | Python | R | Other APIs |
|---|---|---|---|---|---|
| Hadoop | Yes | No | Yes | No | c/c++, ruby, groovy, Perl |
| Spark | Yes | Yes | Yes | Yes | |
| Flink | Yes | Yes | Yes | Yes | |
| Kafka | Yes | Yes | Yes | Yes | c/c++, ruby, groovy, Go, .NET, and more |

https://data-flair.training/blogs/hadoop-vs-spark-vs-flink/
https://www.confluent.io/blog/12-programming-languages-walk-into-a-kafka-cluster
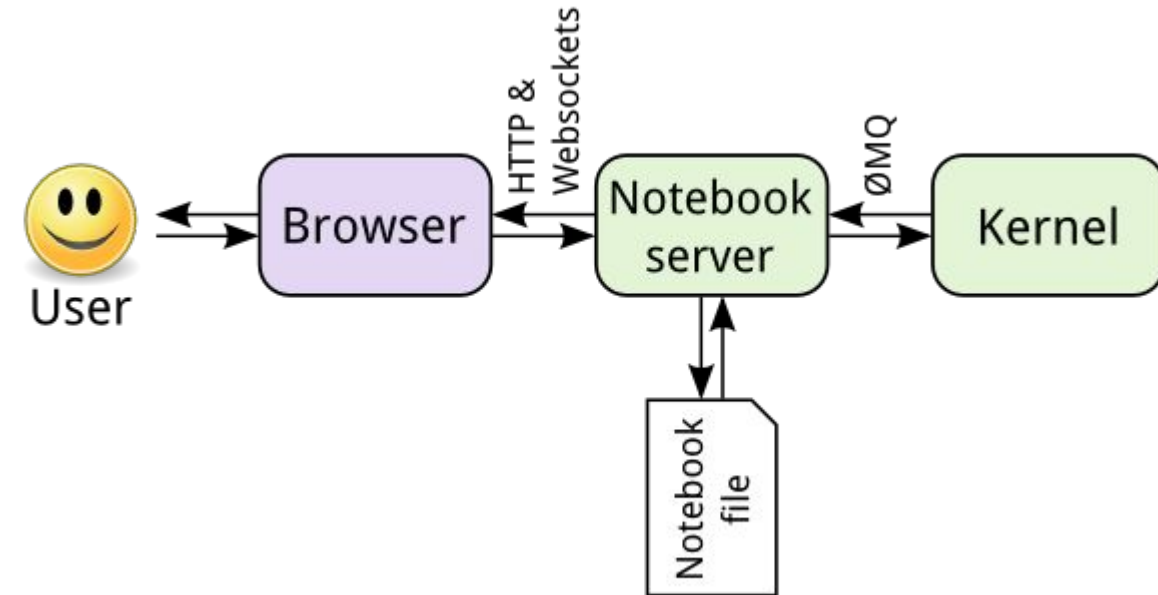
# Agenda

- Big Data and Data Science
- Why Python for Data Science?
- Data Mining Methodology
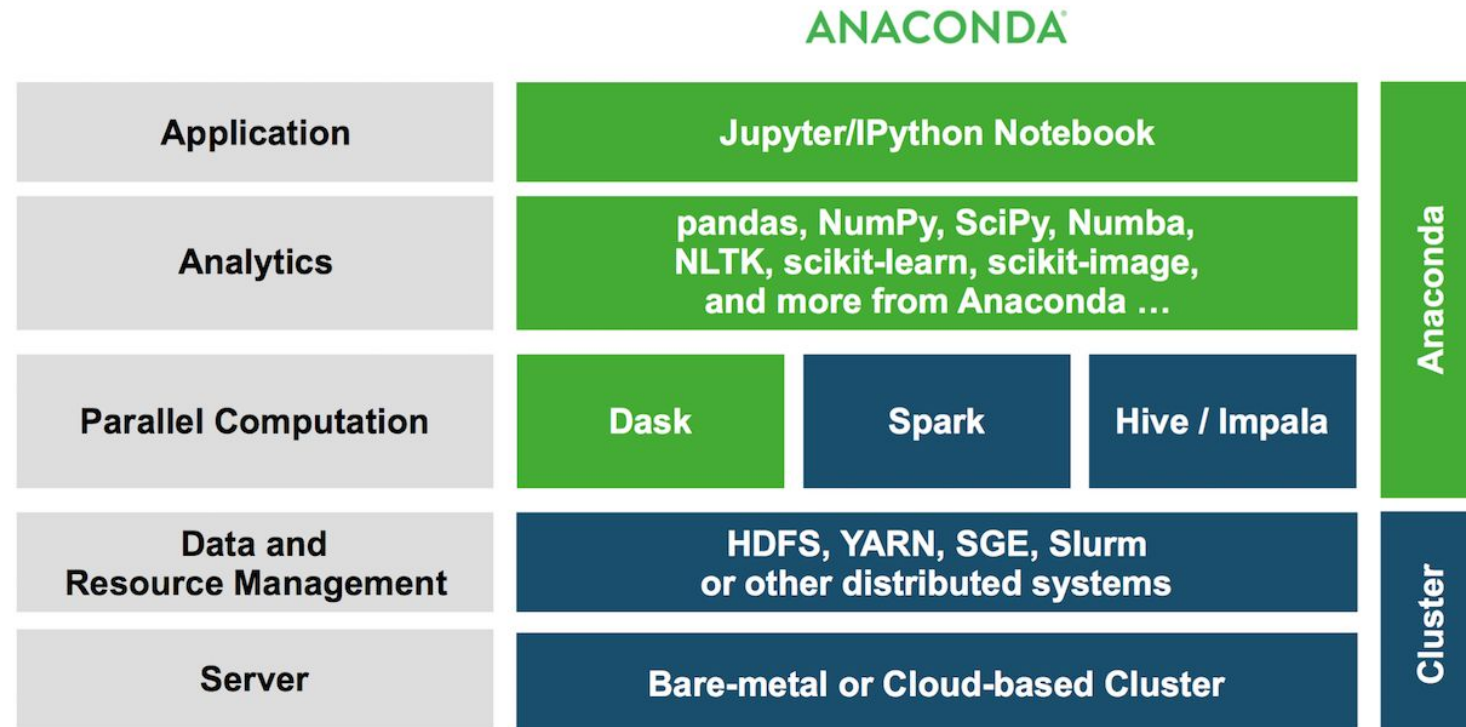- Hands-on: Jupyter Notebook Launch

# Jupyter Notebook

- The Jupyter Notebook is an open-source web application that allows you to create and share documents that contain live code, equations, visualizations and narrative text
- Jupyter supports over 40 programming languages, including **Ju**lia, **Pyt**hon, **R**, and Scala. ([https://github.com/jupyter/jupyter/wiki/Jupyter-kernels](https://github.com/jupyter/jupyter/wiki/Jupyter-kernels))



https://jupyter.readthedocs.io/en/latest/architecture/how_jupyter_ipython_work.html

# Install Jupyter Notebook using Anaconda

- Anaconda is a open-source distribution including Jupyter Notebook, Python, and popular libraries for data science project



ANACONDA®

| Application | Jupyter/IPython Notebook | Anaconda |
|---|---|---|
| Analytics | pandas, NumPy, SciPy, Numba, NLTK, scikit-learn, scikit-image, and more from Anaconda … | Anaconda |
| Parallel Computation | Dask   Spark   Hive / Impala | Anaconda |
| Data and Resource Management | HDFS, YARN, SGE, Slurm or other distributed systems | Cluster |
| Server | Bare-metal or Cloud-based Cluster | Cluster |

https://www.anaconda.com

UMM
UNIVERSITÄTSMEDIZIN
MANNHEIM

Medizinische Fakultät Mannheim
der Universität Heidelberg
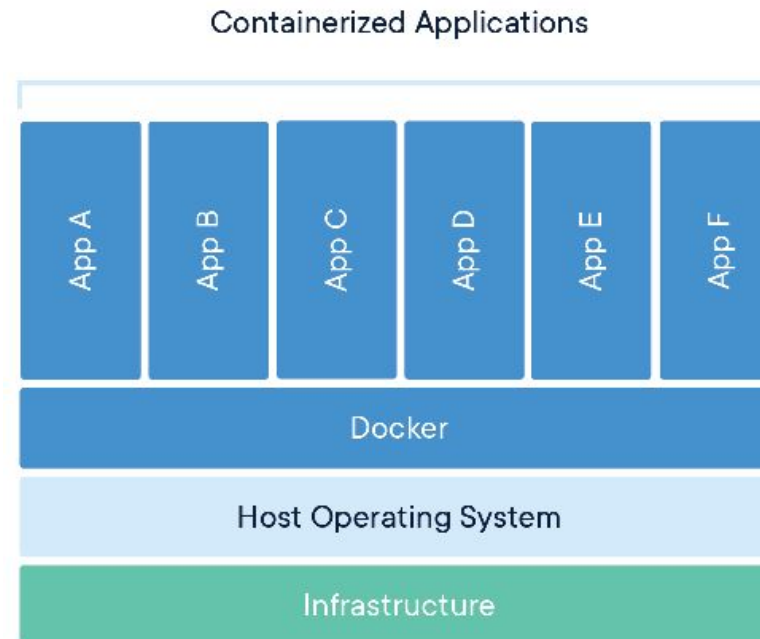Universitätsklinikum Mannheim

# docker pull Jupyter Notebook

- Docker allows you to decouple applications from infrastructure
  (i.e. it is no longer needed to match the development environment to the production environment)
- Docker
  - Docker is an open platform that performs OS-level virtualization
  - Unlike system level virtualization, kernel is the same for all users

Containerized Applications

| App A | App B | App C | App D | App E | App F |
|-------|-------|-------|-------|-------|-------|

Docker

Host Operating System

Infrastructure

UMM
UNIVERSITÄTSMEDIZIN
MANNHEIM

Medizinische Fakultät Mannheim
der Universität Heidelberg
Universitätsklinikum Mannheim

# Jupyter Notebook Launch

1. Open a browser and go to 129.206.5.27:PORT
2. Open a browser and go to 129.206.7.29:PORT

THANK YOU

Kim Hee
Graduate research assistant
at Heinrich-Lanz-Center (HLZ) for Digital Health