



# Big Data techniques with Python

Project and overview of the relevant  
technologies

Supervisor: Dr. hab. Inż. Dariusz Król  
Presented by: Oleksii Kyrylchuk 223224  
17.11.2017



HR EXCELLENCE IN RESEARCH



Wrocław University  
of Science and Technology

# Project Overview



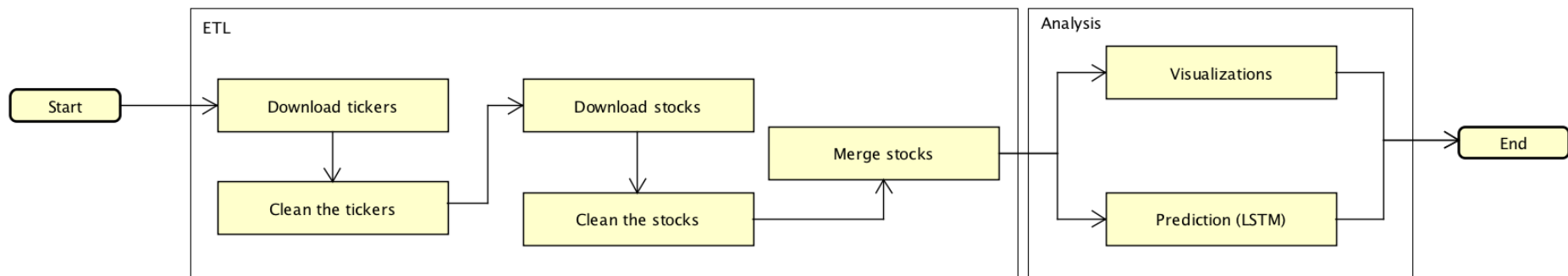
HR EXCELLENCE IN RESEARCH



Wrocław University  
of Science and Technology

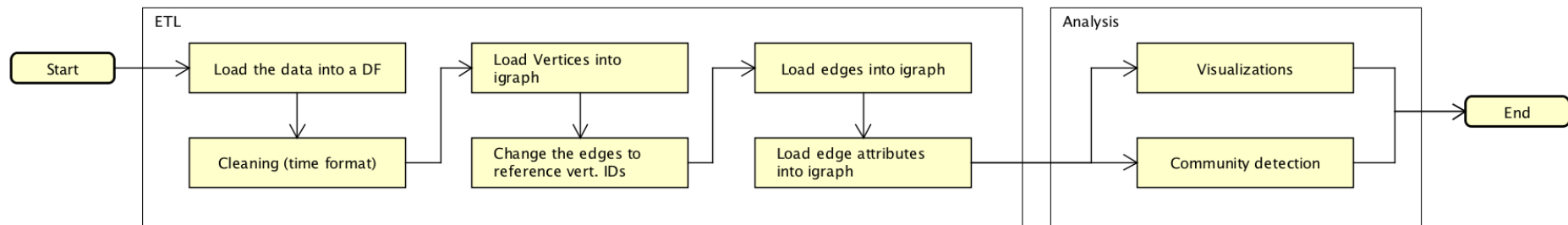
# Stock data analysis (NYSE)

Goal: Demonstrate how to handle time series data in Python



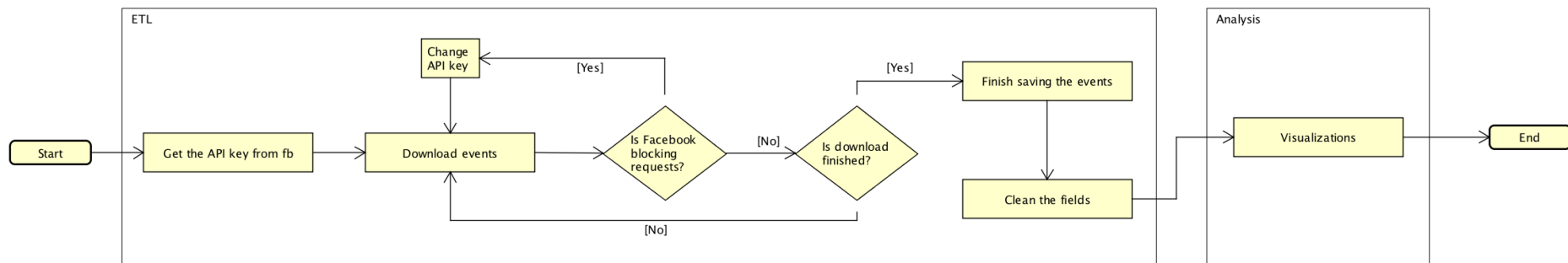
# YouTube network analysis

Goal: Demonstrate how to handle network data in Python



# Facebook geodata analysis

Goal: Demonstrate how to handle geodata in Python



# Requirements and more



HR EXCELLENCE IN RESEARCH



Wrocław University  
of Science and Technology

# FURPS

- **Functional**
  - Code should be able to produce results similar to those found in other papers
- **Usability**
  - The experiments' code must be well-explained in the thesis
- **Reliability**
  - Code shouldn't crash on its own
  - (I am not responsible for the website uptime)
- **Performance**
  - Experiments should be able to finish in less than a day on a suitable server
- **Supportability**
  - The code should be written with the future (python3) in mind

# Verification

- NYSE Stock market (time series data)
  - Charts look good and are understandable
  - Prediction is better than random
- YouTube network (network data)
  - Graphs look similar to those in other papers
- Facebook locations (geolocation data)
  - Charts look good and are understandable



# Methodology

- Following the PEP 8 code standard
- Trying to make code as simple as possible so the reader can easily understand it

# Technology overview



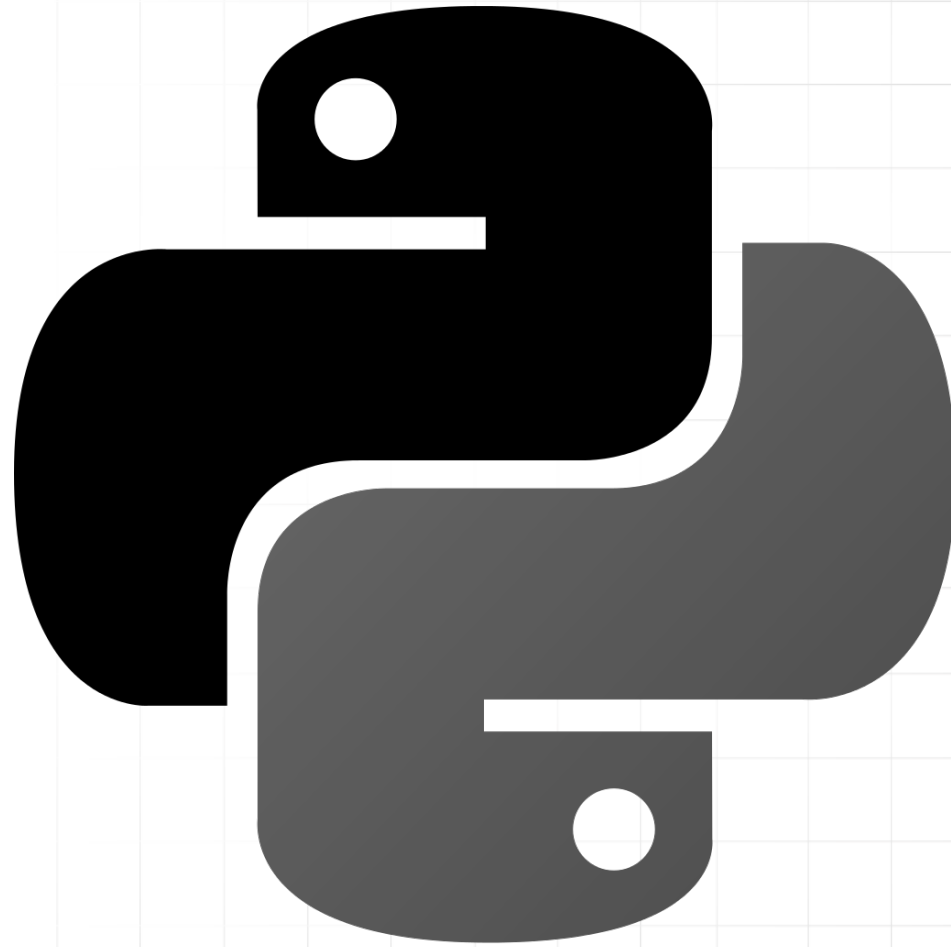
HR EXCELLENCE IN RESEARCH



Wrocław University  
of Science and Technology

# Language

- Python 3
- Async processing is only available in Python3, as well as some other libraries
- Python2 will not be maintained past 2020



# Processing server

- FloydHub<sup>1</sup>
  - Made for data science programming in Python
  - Provides useful CLI tools to upload data to their servers
  - Can run Jupyter notebooks on GPU clusters
  - Easier to get into than AWS<sup>2</sup> or Google cloud<sup>3</sup>



<sup>1</sup> <https://www.floydhub.com/>

<sup>2</sup> <https://aws.amazon.com/>

<sup>3</sup> <https://cloud.google.com/>

# Data retrieval

- Stocks
  - requests<sup>1</sup> – for getting stock tickers
  - BeautifulSoup<sup>2</sup> – for parsing web pages
  - Pandas DataReader<sup>3</sup> – for getting the stocks themselves
- YouTube
  - Pre-generated dataset from KONECT<sup>4</sup>
- Facebook
  - aiohttp<sup>5</sup> (asyncio) – to get locations asynchronously

<sup>1</sup> <http://docs.python-requests.org/>

<sup>2</sup> <https://www.crummy.com/software/BeautifulSoup/>

<sup>3</sup> <https://pypi.python.org/pypi/pandas-datareader>

<sup>4</sup> <http://konect.uni-koblenz.de/networks/youtube-u-growth>

<sup>5</sup> <https://aiohttp.readthedocs.io/en/stable/>



# Data transformation & cleaning

- Numpy<sup>1</sup> – performant arrays, math operations
- Pandas<sup>2</sup> – access to statistic functions and easy transformations
- Built-ins – loop comprehensions, itertools<sup>3</sup>, ...

<sup>1</sup> <http://www.numpy.org/>

<sup>2</sup> <http://pandas.pydata.org/>

<sup>3</sup> <https://docs.python.org/3/library/itertools.html>



# Data storage

- Stocks
  - Pandas (CSV files) – DataReader provides you with CSV files
- YouTube
  - igraph<sup>1</sup> files
- Facebook
  - SQLAlchemy ORM<sup>2</sup> – Load geodata from facebook into a DB

<sup>1</sup> <http://igraph.org/python/>

<sup>2</sup> <https://www.sqlalchemy.org/>

# Data visualization (current)

- Matplotlib<sup>1</sup> – The base for many python plotting packages
  - Seaborn<sup>2</sup> – High-level interface for matplotlib
  - Pandas – Basic plots that don't require importing an extra library
- Datashader<sup>3</sup> – Geodata visualization
- Bokeh<sup>4</sup> – interactive plots that run in browser
- More?

<sup>1</sup> <https://matplotlib.org/>

<sup>2</sup> <http://seaborn.pydata.org/>

<sup>3</sup> <http://datashader.readthedocs.io/>

<sup>4</sup> <https://bokeh.pydata.org/>





# Data Analysis

- igraph – network analysis
- Keras<sup>2</sup> – predictive analysis, machine learning

<sup>1</sup> <https://keras.io/>

# Thank you for your attention!



HR EXCELLENCE IN RESEARCH



Wrocław University  
of Science and Technology