# RIGA**CODING**SCHOOL

# Data Analysis Platform

# Par mums

- Viena no pirmajām programmēšanas skolām Lietuvā un Latvijā;

- Skolas filiāles Viļņā, Klaipēdā, Kauņā un Rīgā;

- Vairāk nekā 1500 absolventi;

- 80% absolventu, kas vēlas strādāt IT jomā;

- Vairāk nekā 50 sadarbības uzņēmumu Rīgā;

- Skolas lektori – ar pieredzi IT sfērā.

# VALDIS

- Izglītība: Maģistra grāds datorzinātnēs
- Pieredze programmēšanā: 20+ gadi
- Specialitāte: grafu teorija sociālo tīklu analizēšanā
- Hobiji: prāta spēles, riteņbraukšana, šahs
- RCS Pasniedzu: Python iesācējiem, Datu Analīzes kursus

# Data Lake

# Brief History of Data Analysis

- ~ 18,000BC – Uganda, Ishango Bone
- ~ 2400BC – Babylon abacus, libraries
- 300BC – 48AD – Library of Alexandria
- ~ 100-200AD Antikythera Mechanism
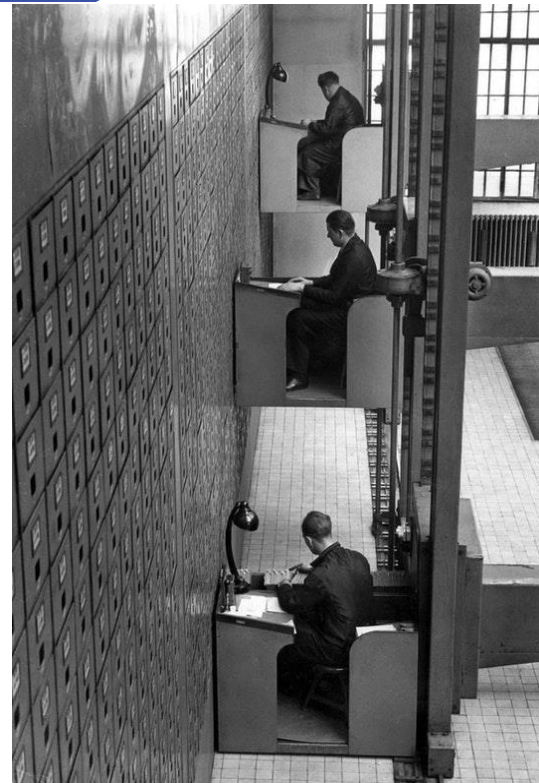
# Brief History of Data Analysis II

- 1663 – London, J.Graunt mortality analysis
- 1865 – banker H. Furnese business intelligence
- 1880-90 US Census Hollerith Machine -> IBM
- 1928 – F. Pfleumer magnetic tape invention

# Brief History of Data Analysis III

- 1950s - Flat Files
- 1958 – IBM's Luhn defines Business Intelligence
- 1960s - CODASYL
- 1970s – Codd's relational DBs -> SQL
- 1980s – Data Werehouses / Marts
- 2000s – Big Data / noSQL DBs
- 2010s – Rise of accessible ML/DL libraries

# BIG DATA LANDSCAPE 2017

# Buzzword bingo

- Big Data
- Data Mining (datizrace)
- Machine Learning – subset of AI
- Data Science – statistics

- Big Data or Pokemon
- https://pixelastic.github.io/pokemonorbigdata/

| Random Forests | Neural Network | Reinforcement Learning | Supervised Learning | Cognitive Computing |
|---|---|---|---|---|
| Caffe | Support Vector Machine | Artificial Intelligence | Python | Cloud |
| Unstructured Data | Bot | DATA SCIENCE BUZZWORD BINGO (free square) | K-means | GPU |
| Spark | Data Wrangling | Deep Learning | Ensemble | Machine Learning |
| Keras | Tensorflow | Big Data | Algorithm | Feature Engineering |

# Data Mining

- Anomalies
- Classification
- Clusters
- Dimension Reduction
- Regression
- Relationship finding
- Summarization / Visualization

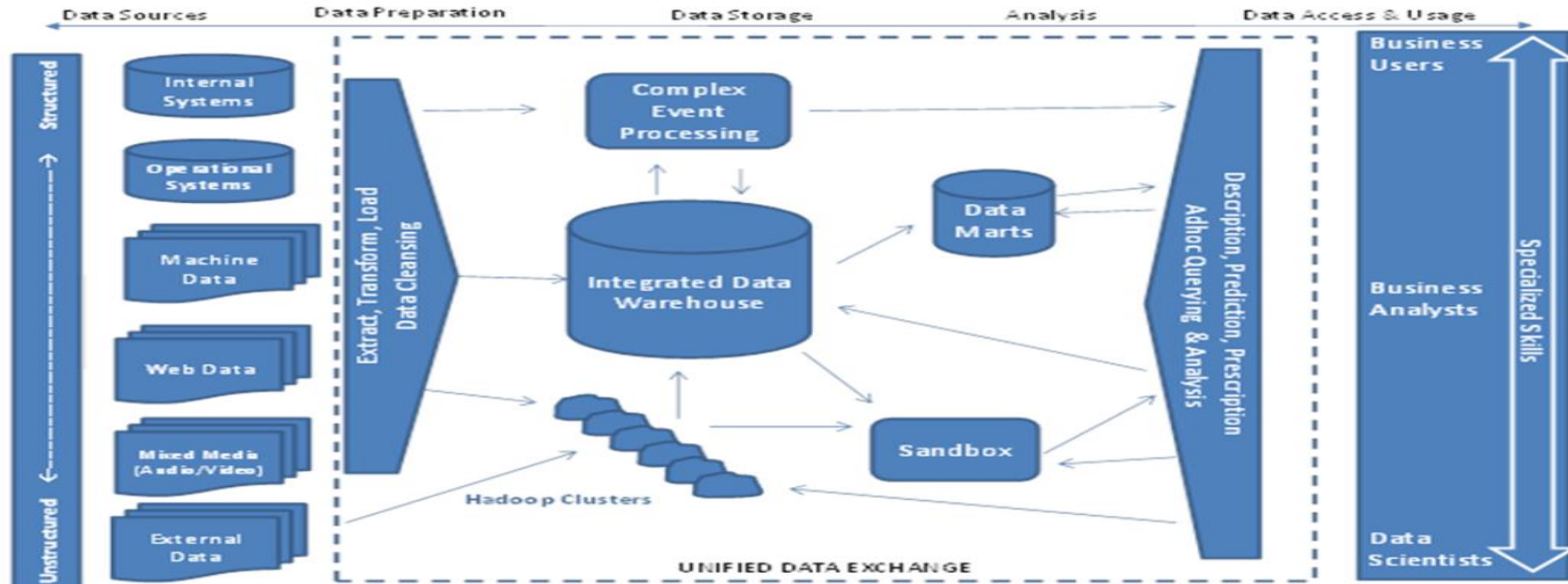# Building a Pipeline

- Data Security
- Cleanup
- Organising Database
- Analysis
- Visualization – Dashboard

- Emphasis on Analysis less on Infrastructure

FIND

GET

VERIFY

CLEAN

ANALYSE

PRESENT
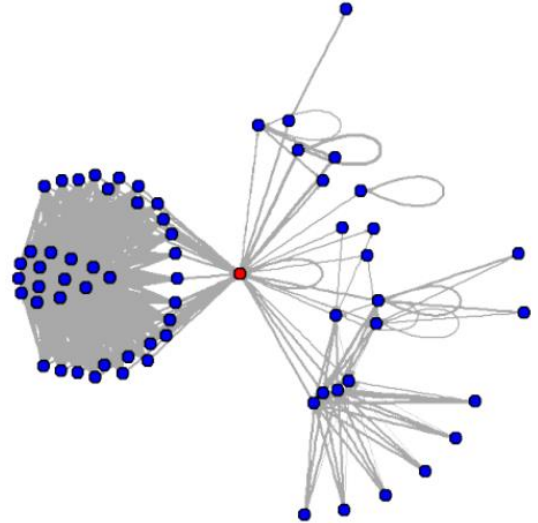
# Full Analysis Framework

# Sandbox for solutions

# Common Projects

- Recommendation System / Churn Prediction
- Web Comments Sentiment Analysis
- Network Analysis (such as transactions on blockchain)
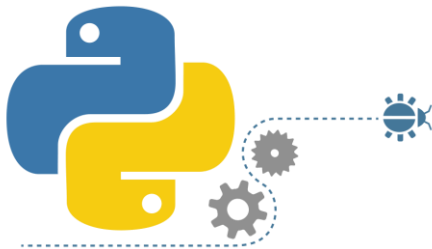
# Avoiding Pitfalls

- Overfitting
- Data Dredging / p-hacking
- https://xkcd.com/882/

# Python Programming Language

- #3 in TIOBE language index
- #1 in Popularity of Programming Languages index
- General purpose / readable
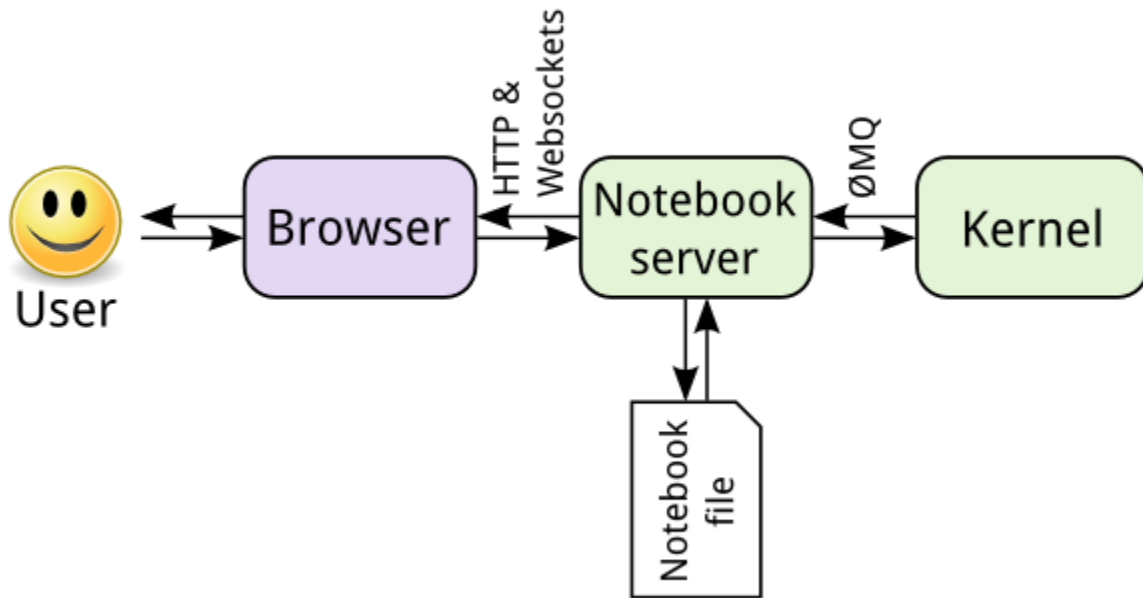- 2018 Nobel Economist Paul Romer recommended!

# Resources

- Anaconda - **The Most Popular Python Data Science Platform**

- (includes Jupyter) : https://www.anaconda.com/download/
- Github: https://github.com/ValRCS/RigaComm_DataAnalysis

# How Jupyter Works

# Jupyter Hosted in the Cloud

- Microsoft - https://notebooks.azure.com/
- Google - https://colab.research.google.com/notebooks/welcome.ipynb
- IBM - https://dataplatform.cloud.ibm.com/
- Anaconda Cloud  : https://anaconda.org/
- MyBinder – https://mybinder.org/
- Self Hosted or just run locally by installing Anaconda