# Wikidata Statistics: What, Where, and How?

## Goran S. Milovanović, Phd

**Data Scientist for Wikidata and Wiktionary**
*Wikimedia Deutschland*

# Wikidata Statistics:
# What, Where, and How?

## All data products mentioned in this talk are available from

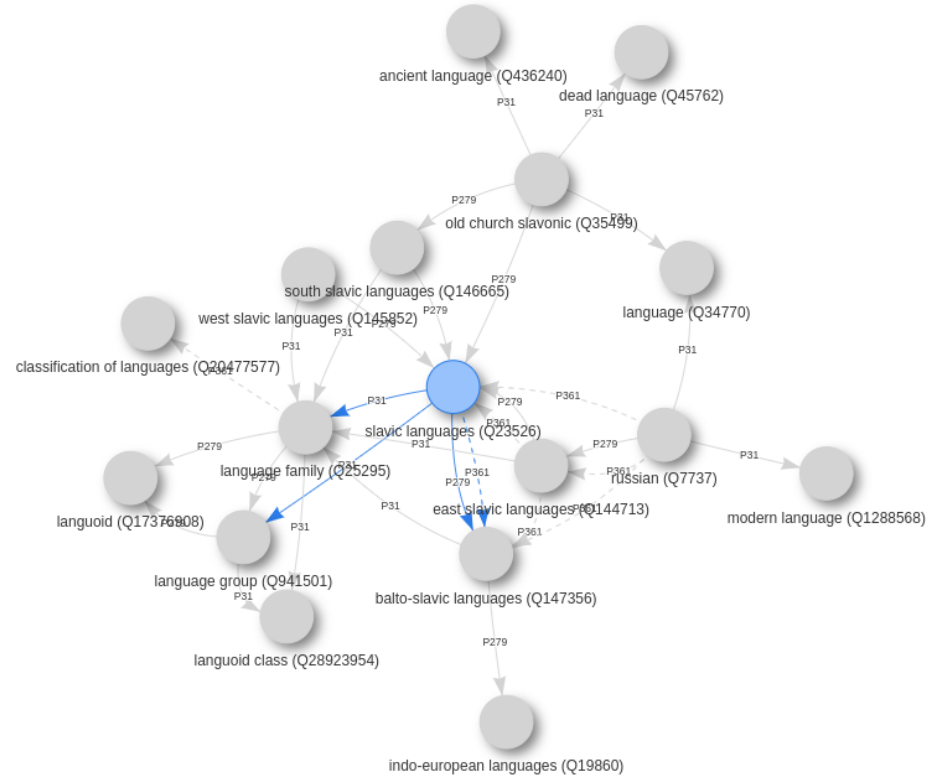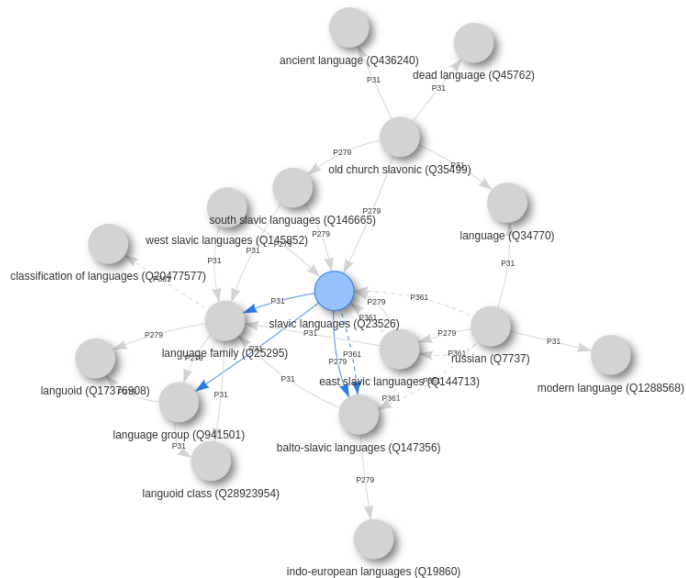# http://wmdeanalytics.wmflabs.org/

# Wikidata Statistics:
# What, Where, and How?

## Goals

- An Overview of **Wikidata Statistics & Analytics systems**
- Exemplify **the usage of our analytics** in several domains (Wikidata items, languages, external identifiers, item quality)
- Go just a bit under the hood to illustrate **how we are doing it**



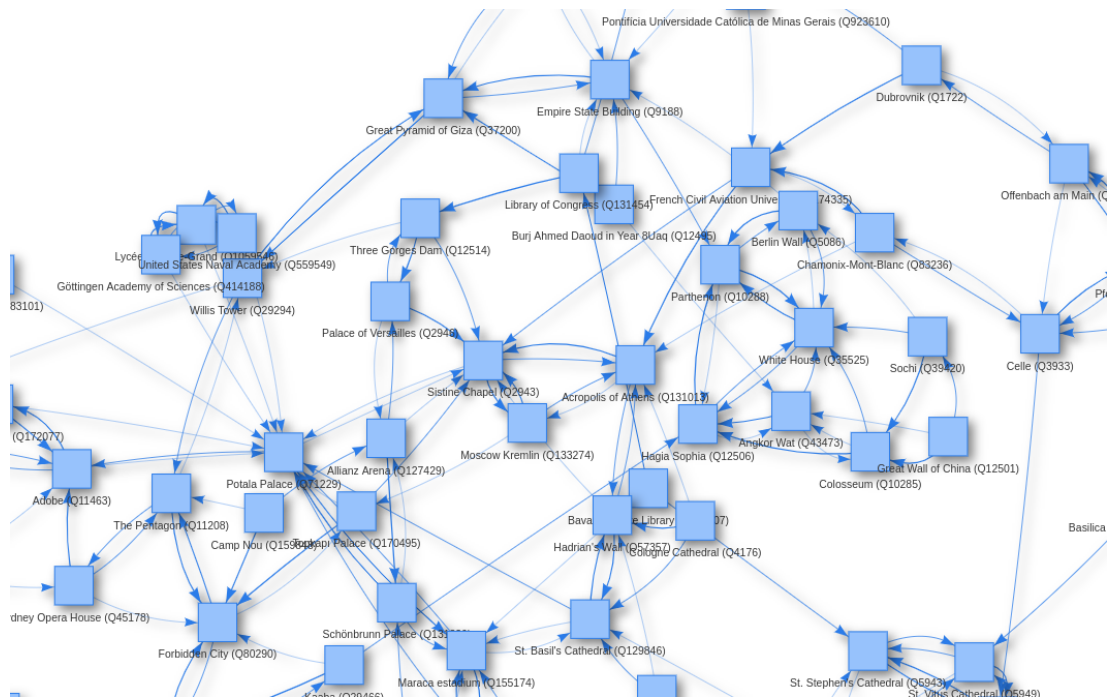http://wmdeanalytics.wmflabs.org/WD_LanguagesLandscape/

In 2017: we need an analytical system that will give us an insight into the ways the Wikidata items are reused across the Wikimedia projects (Wikipedia, Wikivoyage, Wikisource, etc)

**Wikidata Concepts Monitor**
→item reuse similarity structures
= items frequently used together across the Wikimedia projects are connected.
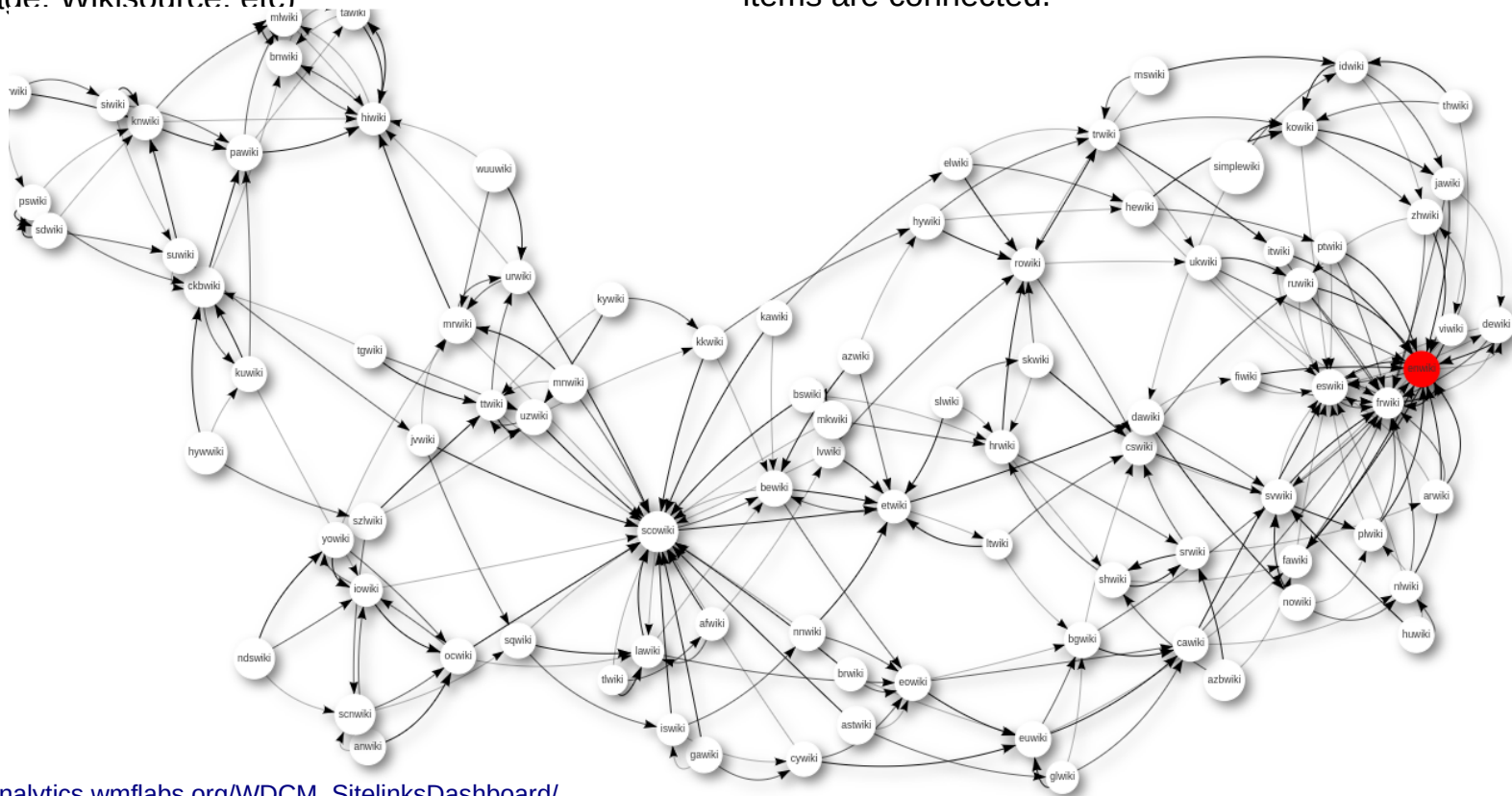


http://wmdeanalytics.wmflabs.org/WDCM_SitelinksDashboard/

In 2017: we need an analytical system that will give us an insight into the ways the Wikidata items are reused across the Wikimedia projects (Wikipedia, Wikivoyage, Wikisource, etc)

**Wikidata Concepts Monitor**
→Wikipedia similarity structures
= projects that reuse the similar Wikidata items are connected.



http://wmdeanalytics.wmflabs.org/WDCM_SitelinksDashboard/

In 2017: we need an analytical system that will give us an insight into the ways the Wikidata items are reused across the Wikimedia projects (Wikipedia, Wikivoyage, Wikisource, etc) →
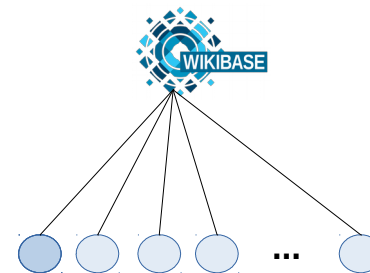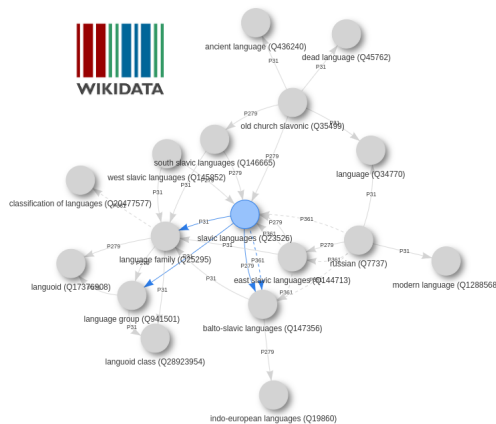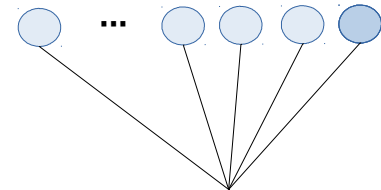**Wikidata Concepts Monitor**

**RDF Graph ← SPARQL**
**WDQS or WD Dump processing**



**Data Model**
**&**
**Structures**

**Motivation/Goals**
**→**
**Questions**

**Client Projects**

**RDBS ← MySQL**
**Big Data → Hadoop/Spark**

**...**

**Wikimedia Foundation**
**Data Lake**

**Production (stat1004)**

1 → m2 MariaDB replica
analytics.eqiad.wmnet:
wbc_entity_usage tables

WDCM_Sqoop_Clients.R
―――――
production (stat1004)

2 → Hadoop:
wdcm_clients_wb_entity_usage table
―――――
HDFS
filesystem

**Production (stat1007)**

wdcmConfig.xml
―――――
ETL, Apache Spark,
and Machine Learning
parameters

WDCM Taxonomy

3 →

4 →

SPARQL Endpoint

wdcmModule_CollectItems.R
―――――
Blazegraph GAS programs
to collect sets of Wikidata items
for analysis

5

wdcmModule_Orchestra.R
―――――
Orchestration of
analytical scripts

wdcmModule_ETL.py
―――――
Pyspark ETL procedures to ftech
re-use data for the selected Wikidata items
accross the Wikimedia projects

7a

6

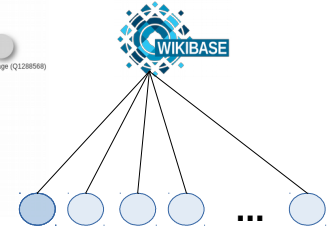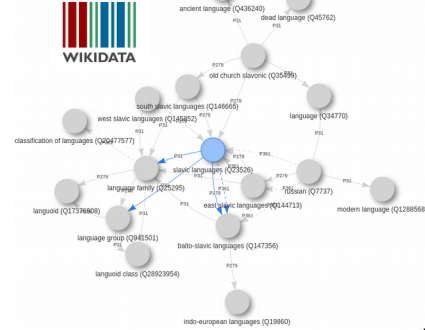WDCM Public Datasets
expose results to Cloud VPS
―――――
path:
https://analytics.wikimedia.org/datasets/wmde-analytics-engineering/wdcm/

7b

wdcmModule_ML.R
―――――
Machine Learning procedures:
Latent Dirichlet Allocation
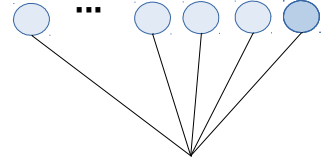t-SNE dimensionality reduction

**Many systems need to work together just in order for you to get your data and organize the data model appropriately...**
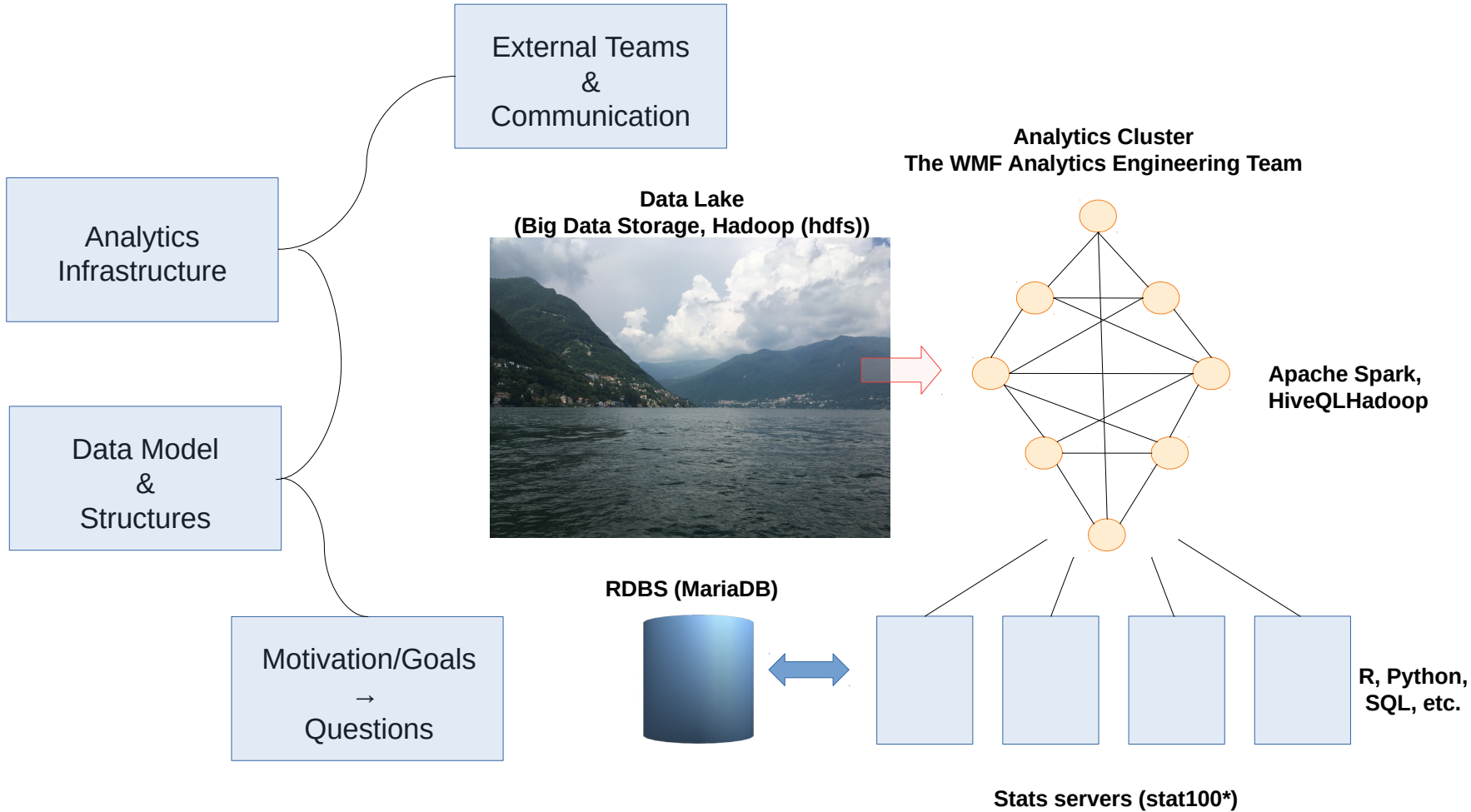
**RDF Graph ← SPARQL**
**WDQS or WD Dump processing**



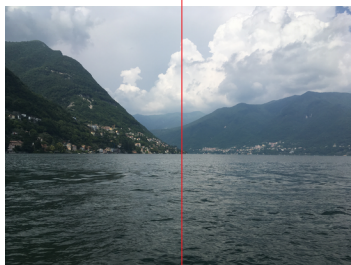**Client Projects**

**RDBS ← MySQL**
**Big Data → Hadoop/Spark**

**Wikimedia Foundation**
**Data Lake**

## Data + Systems Synchronization:
## it can get really nasty...



wmf.mediawiki_history table (huge)

Wikidata Concepts Monitor Reuse Statistics

WD JSON Dump in HDFS

ORES ML System Item Quality Predictions

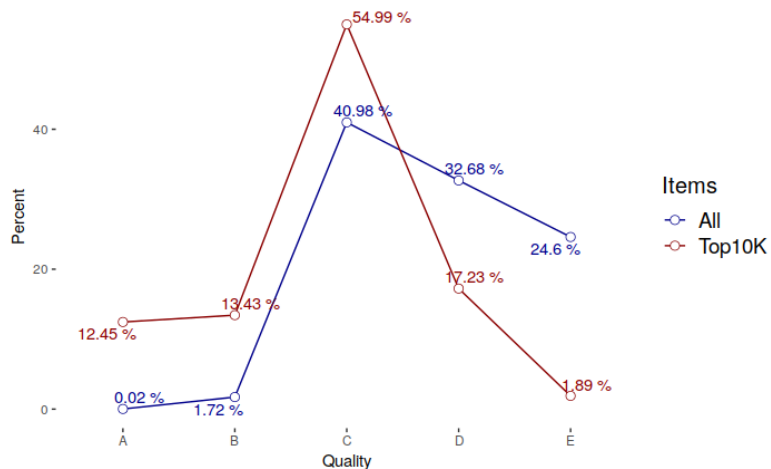http://wmdeanalytics.wmflabs.org/Wikidata%20Quality%20Report.nb.html

# Data + Systems Synchronization:
# it can get really nasty...
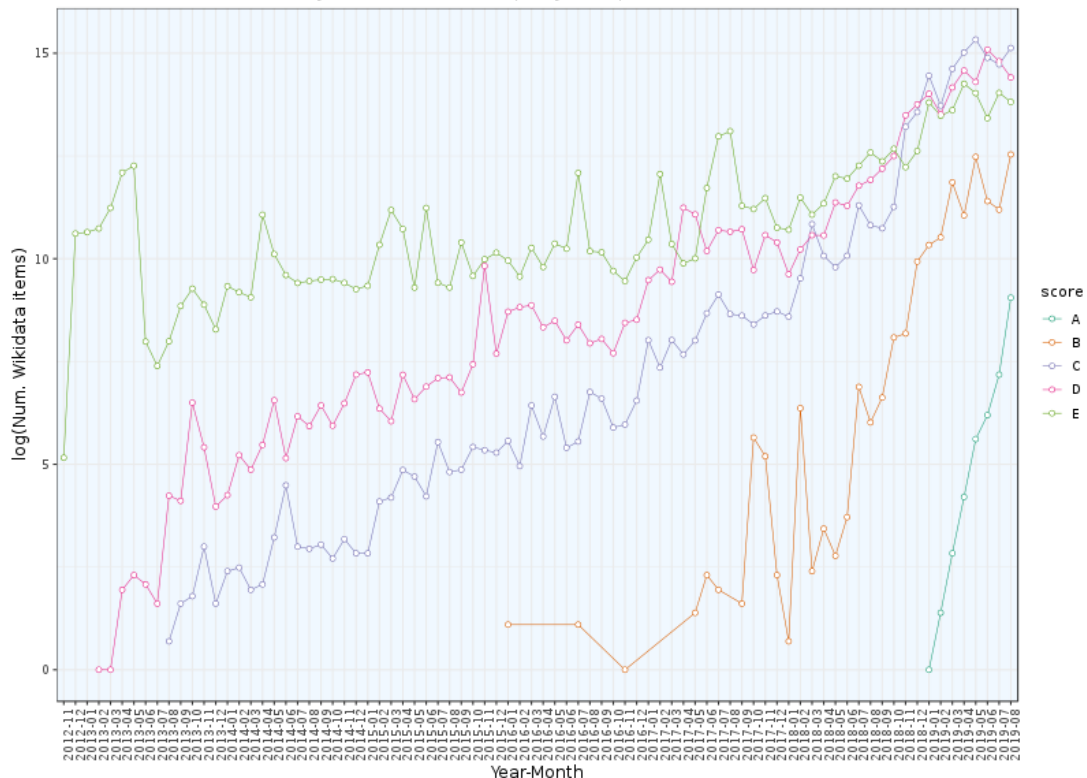
wmf.mediawiki_history table (huge)

Wikidata Concepts Monitor Reuse Statistics

WD JSON Dump in HDFS

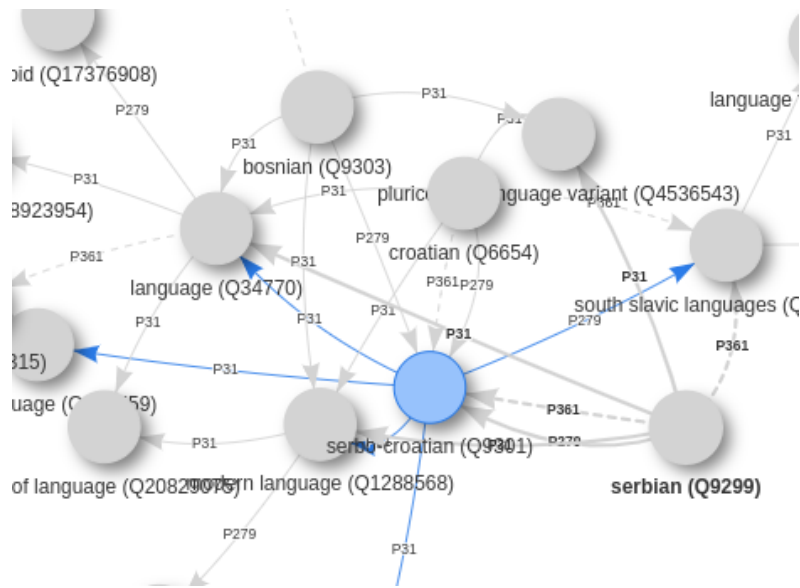ORES ML System Item Quality Predictions



When did the item receive its latest revision?
(NOTE. Only items with the ORES quality class prediction are considered)

**While we play with large and complex datasets, we try to make use of the byproducts of our work…**
**<u>And you should too!</u>**

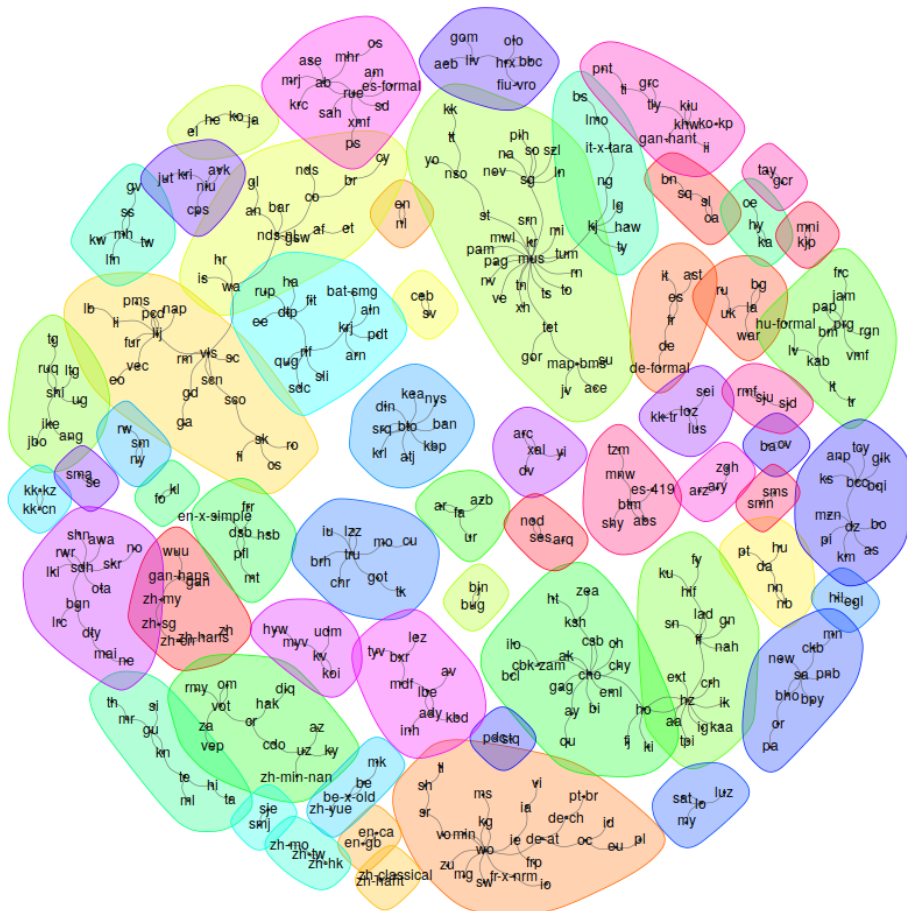**Fix the Ontology!**

Example: languages that are a part of (P361) and a subclass of (P279) something at the same time (e.g. both Serbian (Q9299) and Croatian (Q6654) are parts of and subclasses of Serbo-Croatian (Q9301) at the same time; so is Serbo-Croatian a language, or a set of languages (besides being a pluricentric language (Q250858))?

*Mereological and set-theoretic relations are not the same.*

# Speaking of languages...

**The Wikidata Languages Landscape**
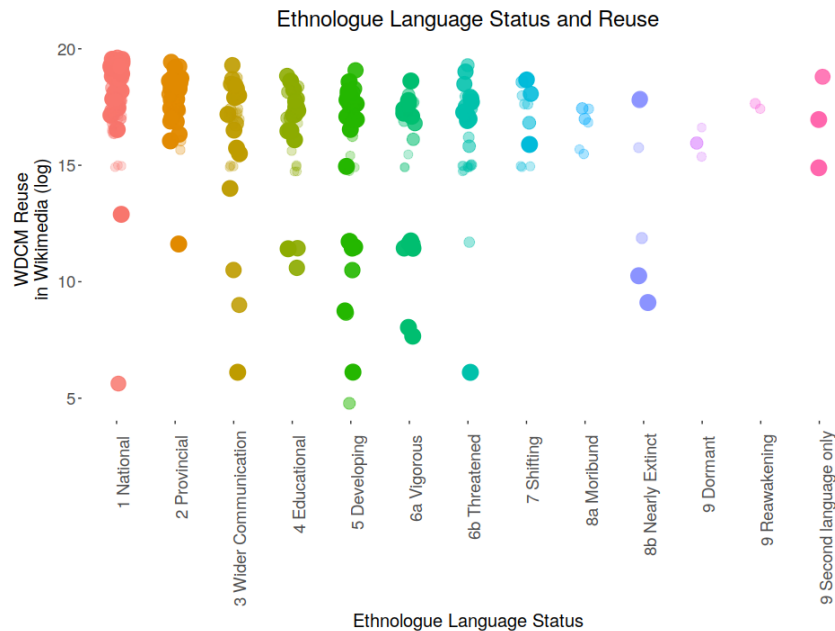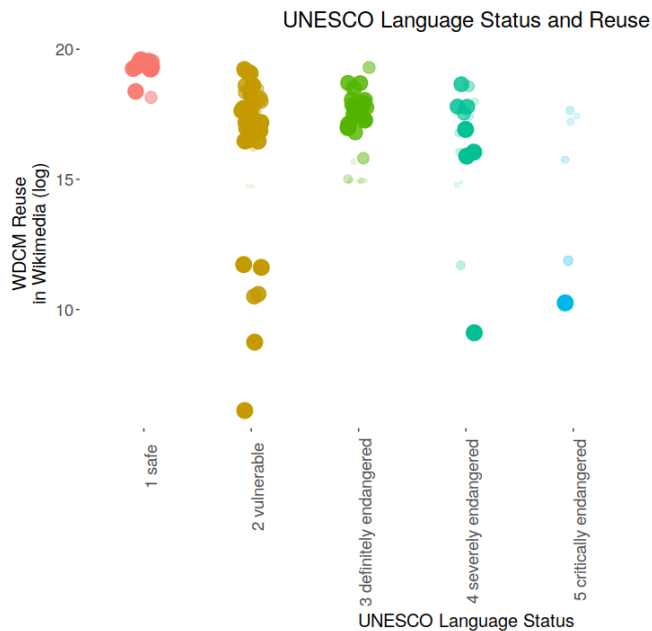
… relies on different data sources to provide a comprehensive picture of how different languages are used in Wikidata and - via the entities that they refer to - how they are mapped across the universe of Wikimedia project



http://wmdeanalytics.wmflabs.org/WD_LanguagesLandscape/

# Speaking of languages...

## The Wikidata Languages Landscape

… relies on external data found in Wikidata to make relevant assessments of the way languages are used...



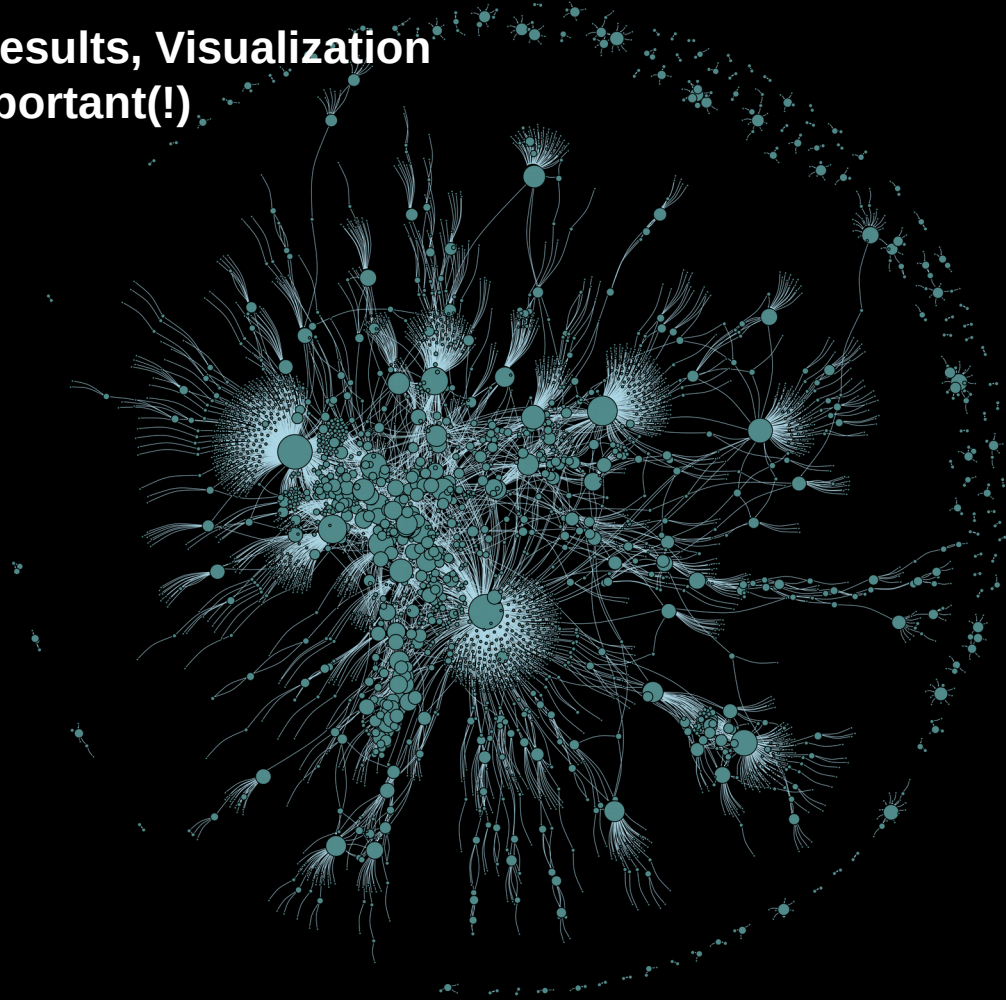http://wmdeanalytics.wmflabs.org/WD_LanguagesLandscape/

# Speaking of external resources...

**The Wikidata External
Identifiers Landscape**

… to provide insight into the structure of
the overlap in usage of various WD
external identifiers.



http://wmdeanalytics.wmflabs.org/WD_ExternalIdentifiersDashboard/

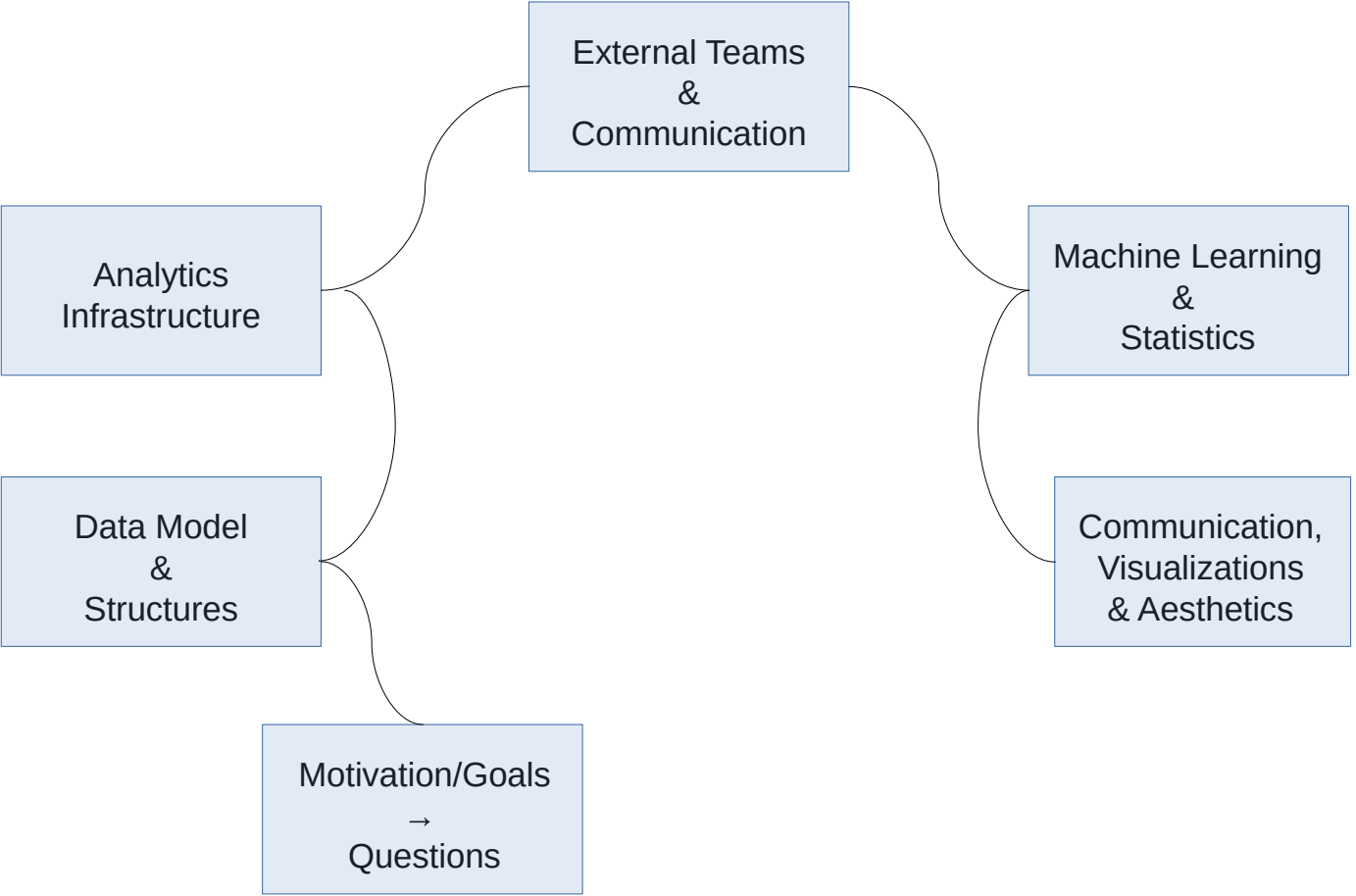# Communication of our results, Visualization
# & Aesthetics → Very Important(!)

Communication of our results, Visualization & Aesthetics → Very Important(!)

But one just does not get interesting results for free: **Machine Learning and Statistics**



← **Obtained by running a clustering algorithm against a Jaccard distance matrix derived from *408 languages* x ~60M items contigency matrix...**

# Machine Learning and Statistics

- t-distributed stochastic neighbor (t-SNE) embedding for dimensionality reduction

- Latent Dirichlet Allocation (LDA) for extracting semantic themes in projects

- Various Clustering algorithms

- Various distance metrics

- etc.

External Teams
&
Communication

Analytics
Infrastructure

Machine Learning
&
Statistics

Data Model
&
Structures

Communication,
Visualizations
& Aesthetics

Motivation/Goals
→
Questions

```mermaid
graph

External Teams & Communication

Analytics Infrastructure

Machine Learning & Statistics

Data Model & Structures

Communication, Visualizations & Aesthetics

Motivation/Goals → Questions

Deployment & Dashboards
```
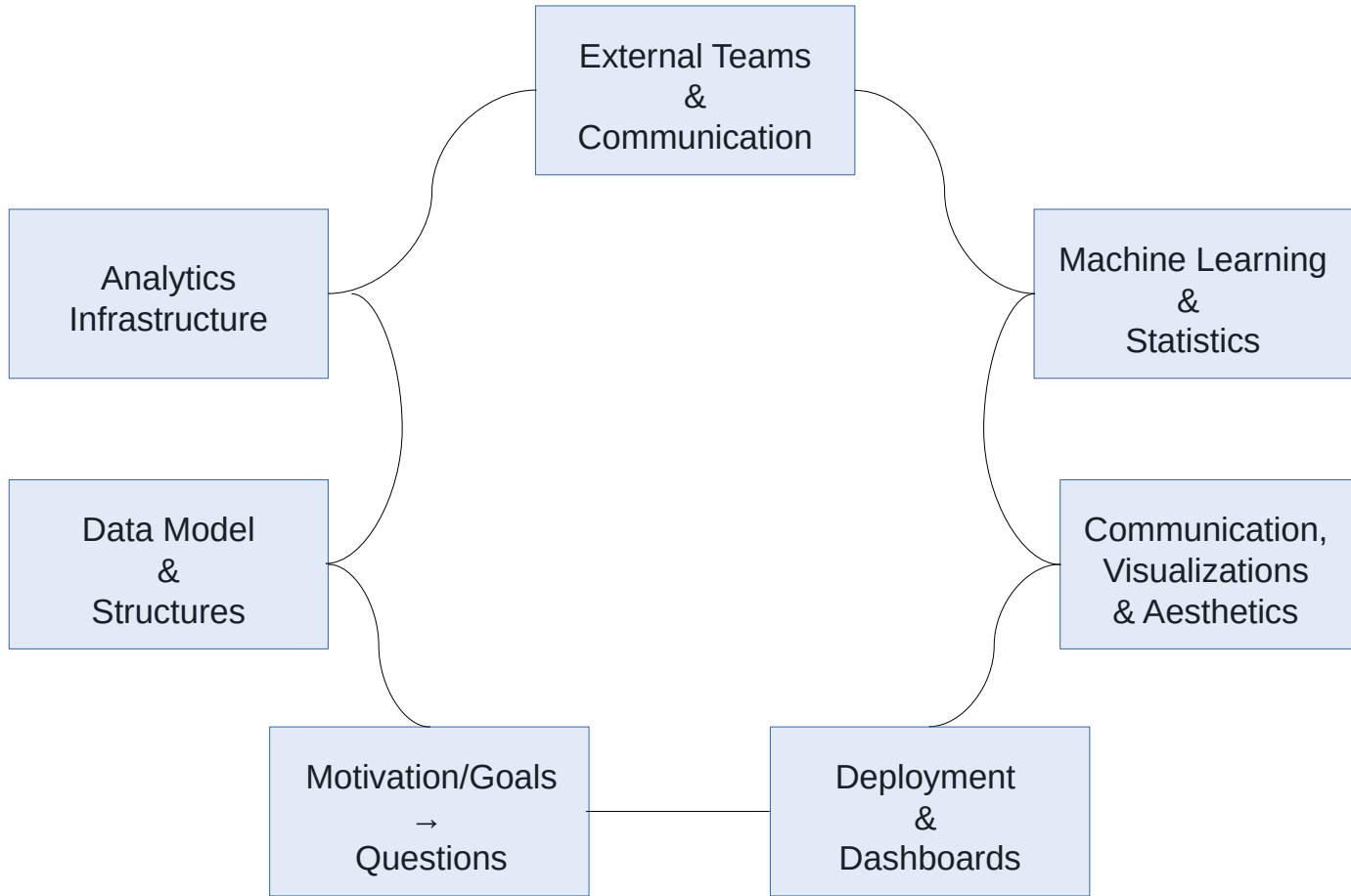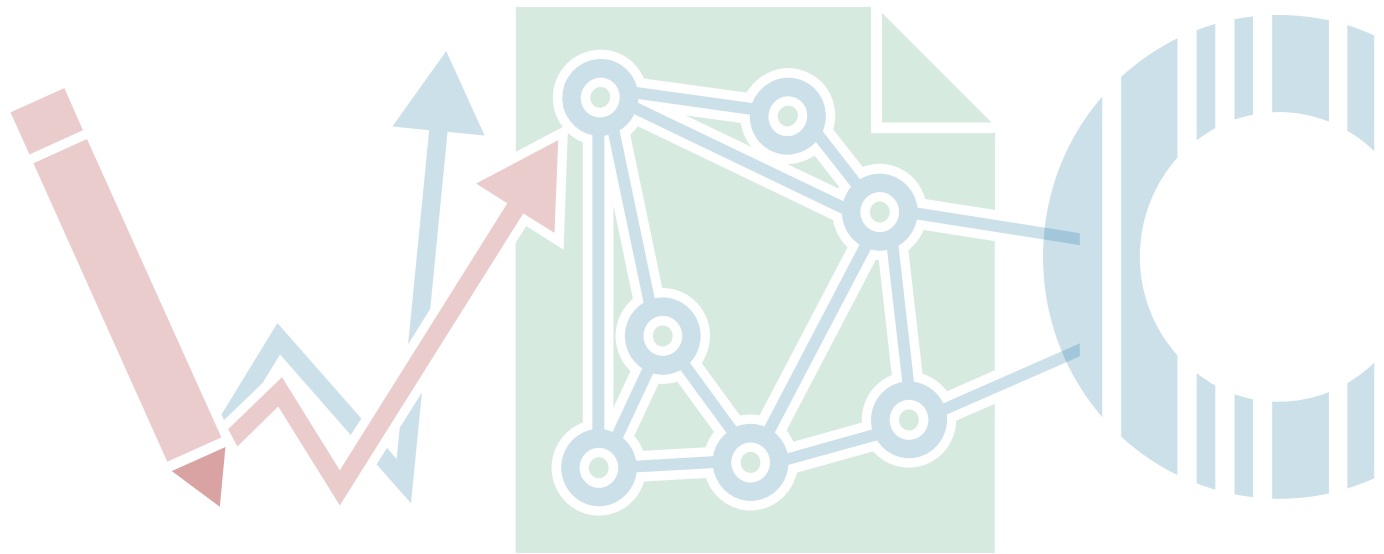
**The true picture:
all phases of the process are interdependent**

External Teams
&
Communication

Machine Learning
&
Statistics

Analytics
Infrastructure

Data Model
&
Structures

Communication,
Visualizations
& Aesthetics

Motivation/Goals
→
Questions

Deployment
&
Dashboards