# Laboratory of
# Big Data and Textual Analysis

Higher Institute of Information Technologies and Information Systems

28 April, 2017

# Outline

- Laboratory of Big Data and Textual Analysis
- Past projects
- Current projects

# Laboratory of Big Data and Textual Analysis

- Team:
  - Prof. V. Solovyev;
  - Postdocs: V. Ivanov, E. Tutubalina;
  - A. Sirotkin (HSE St. Petersburg), A. Kadurin (PDMI RAS, Mail.Ru Moscow);
  - 3 students (Master or PhD);
- Collaborations:
  - A. Tropsha, A. Kotov (USA);
  - I. Batyrshin (Mexico);
  - KFU Cheminformatics Lab, N. Loukachevitch, V. Polyakov, S. Nikolenko (Russia).

# Conferences

- International Conference on Web Search and Data Mining (WSDM 2016; USA)
- European Conference on Information Retrieval (ECIR 2015; Austria)
- International Conference on Intelligent Text Processing and Computational Linguistics (CICLing 2013; Greece)
- International Conference on Computational Linguistics and Intellectual Technologies "Dialogue" (Dialogue 2015, 2016, 2017; Russia)
- International Conference on Text, Speech, and Dialogue (TSD 2014, 2015, 2016; Czech Republic)
- Mexican International Conference on Artificial Intelligence (2013, 2014, 2015; Mexico)
- Women in Machine Learning Workshop (WiML 2015 in conj. with NIPS, Canada)
- The 8th IEEE International Conference on Social Computing and Networking (SocialCom 2015; China)
- International Conference on Knowledge Engineering and Semantic Web (KESW 2013, 2016; Russia, Czech Republic)

# Past projects

1. Development of Ontologies
   a. Ontology CIDOC CRM - RFBR / Kunstkamera / European Commission (europeana.eu)
   b. OntoMathPro ontology
   c. Nanomaterials Ontology - Ministry of Education and Science of the Russian Federation
2. Development of text mining methods for Russian
   a. NLP@CLOUD - Ministry of Education and Science of the Russian Federation
   b. Information Extraction - joint work with Hewlett-Packard (EvExRus-2012, InfEx-2013)
3. Development of linguistic resources for Russian
   a. Corpora: 18th Century Textual Documents in Russian National Corpus - joint work with Russian Language Institute of the Russian Academy of Sciences and Yandex
   b. Entities Dictionaries (e.g, Persons, Organizations)
   c. Thesaurus: Russian WordNet - joint work with Research Computing Center of Moscow State University

# Open Kunstkammer Data Project

**Goal:** Represent a dataset consisting of more than 40 000 digital images and their descriptions in English and Russian as Linked Data

**Constrains:** Follow ICOM-CIDOC recommendations; "SKOSify" MAE RAS's controlled bilingual vocabularies; interlink with **DBPedia and Geonames**

**The resulting data warehouse (more than 10 M triples):** interconnected descriptions of **museum objects, persons, places** and **events**

**SPARQL end-point: http://data.kunstkamera.ru/sparql**

# CIDOC Conceptual Reference Model and the Europeana

1) The **CIDOC CRM** ontology is a reference model in museum domain; it was translated into Russian
2) Then it was mapped to major thesauri in cultural heritage, including the Getty's Art and Architecture Thesaurus **(more than 30,000 concepts)**
3) Museum metadata for KFU museum was represented in the CIDOC CRM format
4) Then the metadata was transferred to the European CH project: www.europeana.eu



The CRM class hierarchy

A hierarchy of thesaurus descriptor blocks

New instances of the CRM E57.Material class

A P127F.has_broader_term property value

# 18th Century Texts in RNC

НАЦИОНАЛЬНЫЙ КОРПУС
РУССКОГО
ЯЗЫКА

Результаты поиска в основном корпусе

перейти на страницу поиска   выбрать подкорпус   версия без ударений   настройки   English

Объем всего корпуса: 115 645 документов, 23 803 881 предложение, 283 431 966 слов.

Найдено 1 838 документов общим объемом 285 063 предложения, 5 202 929 слов.

Сохранить подкорпус и перейти к странице поиска

Страницы: 1  2  3  4  5  6  7  8  9  10  11   следующая страница

1. И. Ф. Богданович. Письма (1767-1800)  [омонимия снята]

2. И. Ф. Богданович. Письма (1767-1800)  [омонимия снята]

3. А. Н. Радищев. Бова (1798-1799)  [омонимия снята]

4. А. Н. Радищев. Письмо М. Н. Радищеву (1797)  [омонимия снята]

5. А. Н. Радищев. Прошение на имя Павла I (1797)  [омонимия снята]

6. Н. М. Карамзин. Филалет к Мелодору (1795)  [омонимия снята]

7. Н. М. Карамзин. Что нужно автору? (1794)  [омонимия снята]

8. Григорий Сковорода. Толкование из Плутарха о тишине сердца (1766-1794)  [омонимия снята]

9. Н. М. Карамзин. Бедная Лиза (1792)  [омонимия снята]

10. Д. И. Фонвизин. Выбор гувернера (1790-1792)  [омонимия снята]

11. А. Н. Радищев. Записки путешествия в Сибирь (1790)  [омонимия снята]

# **Text mining toolkit**

Components for text processing:
- Text Segmenter
- Morphological analyzer
- PoS-tagging
- Lemmatizer
- Chunker
- Syntactic Analyzer
- Named Entity Recognizer
- Event Extractor

Tasks: information extraction (NER, EE)

Open-source toolkit: https://github.com/textocat/textokit-core

# Examples

GPE | Person
Турция лишилась одной из самых ярких фигур отечественной журналистики - Мехмета Али Биранда.

Person | Time | Time
Он скончался в четверг на 71-м году жизни вследствие сердечной недостаточности в

Facility | Time
стамбульской больнице, куда поступил в среду.

Person | Time | GPE
Биранд родился 9 декабря 1941 года в Стамбуле.

Time | Person | Organization
В 1964 году он начал свою журналистскую карьеру в газете "Миллиет".

Org | End-Position | Person | Org | Time
«Баффало» отправил в отставку главного тренера Линди Раффа, который работал с командой с июля 1997 года,

Person | Org | Org
За период, когда наставник стоял у руля «Клинков», в лиге произошло 170 изменений на тренерских мостиках друг

# Information Extraction

- Named Entity Recognition (NER)
- Event Extraction (EE)

# NER. Original Corpus

Contains:

    100 documents from different Russian online newsfeeds;

    ~35000 words.

Manually annotated:

    ~1300 Organization mentions;

    ~500 Person mentions.

Paper: *Gareev, Tkachenko, Solovyev, Simanovsky, Ivanov*.
**Introducing Baselines for Russian Named Entity Recognition**.
CICLING-2013.

# NER. Quality

Person F1                ~ **0.79**

Organization F1          ~ **0.55**

Future work:

    1) extend the corpus;

    2) annotate other NE types: GPE, Locations, Temporal Expressions;

    3) implement a Sequential Classification-based Tagger with morphological and chunker features.

# Event Extraction

Task: extract event mentions of predefined types and fill their arguments from Russian news texts.

Example event types:

- Merge & Acquisition

- Person Position Change

We follow Automatic Context Extraction

(ACE-2005) definitions mostly.

# EE. Experience

Finished the 1-year project in collaboration with HP Russia Labs (Saint-Petersburg)

Elaborated rule-based extraction

— use the NER component;

— match '*event triggers*' in a text;

— apply rules to match arguments around a trigger.

Prepared an evaluation corpus

— 100 docs, ~1500 sentences;

— ~130 event mentions (3 event types) have been annotated so far

# EE. Ongoing efforts

Moving from 'linear patterns' matching

   towards applying trigger-specific mappings

     fill an event argument by a nearby NP that satisfies certain
     morphological and semantic constraints


The experiments have shown that the latter approach significantly better, **but**
   requires more sophisticated preprocessing steps:

    - PoS-tagging;

    - NP recognition.

# Current Projects

1. Text mining models and methods for analysis of the needs, preferences and consumer behavior (RSCF grant, 15-11-10019)
2. Computational models and mathematical methods for big data analysis of trends and correlations in society (RFBR grant, 15-29-01173)
3. Mining Hypotheses concerning Drug Discovery and Drug Repurposing (joint work with Cheminformatics Lab of KFU)
4. Pharmacovigilance based on Social Media Posts written in Russian
5. Named Entity Recognition in Legal Documents (joint work with the Faculty of Law of KFU)

# Analysis of consumer behavior and user opinions

- Develop, implement and evaluate text mining methods for processing of consumers' texts:
  - Deep neural networks for prediction of demographic information from medical user reviews;
  - Deep neural networks for bilingual sentiment analysis of short texts;
  - Aspect-based sentiment lexicons constructed with topic modeling;
  - Methods for inferring sentiment-based priors in topic models;
  - Topic models for extracting failures from product reviews;
  - Knowledge-driven Event Extraction in Russian.
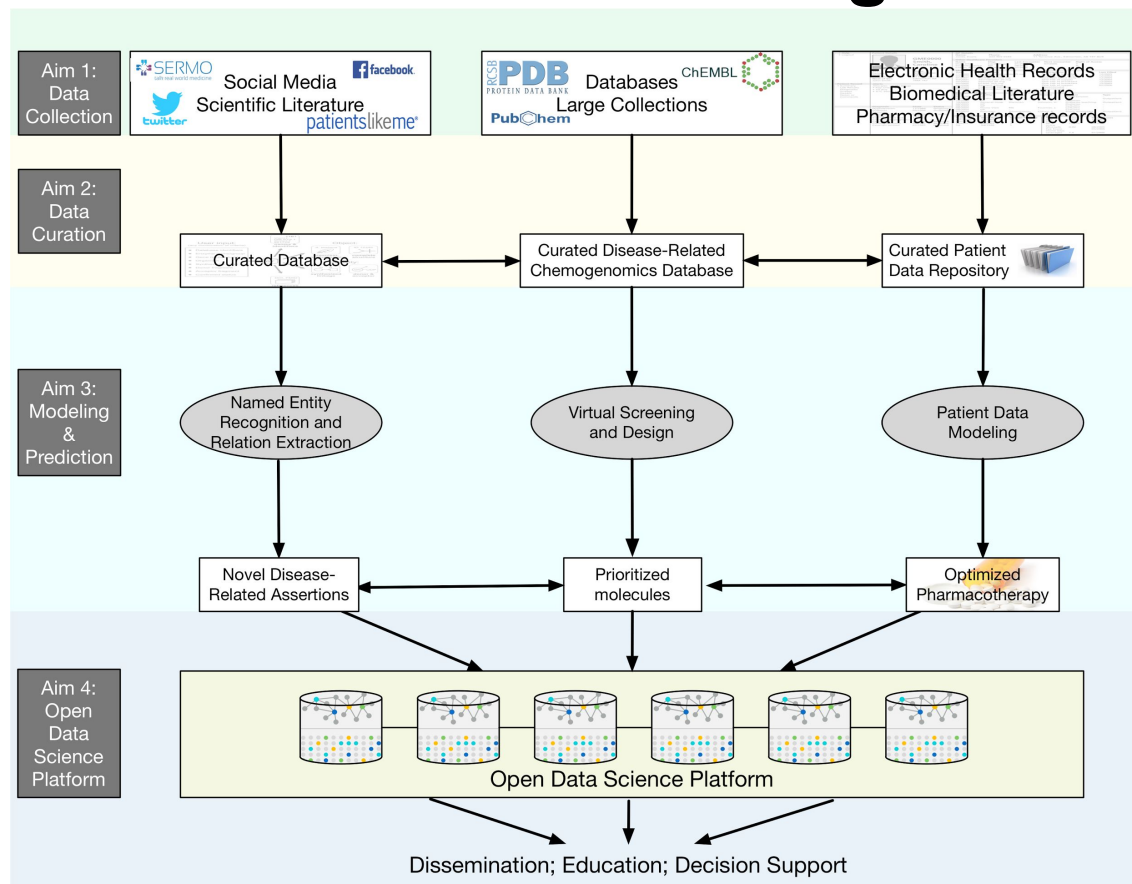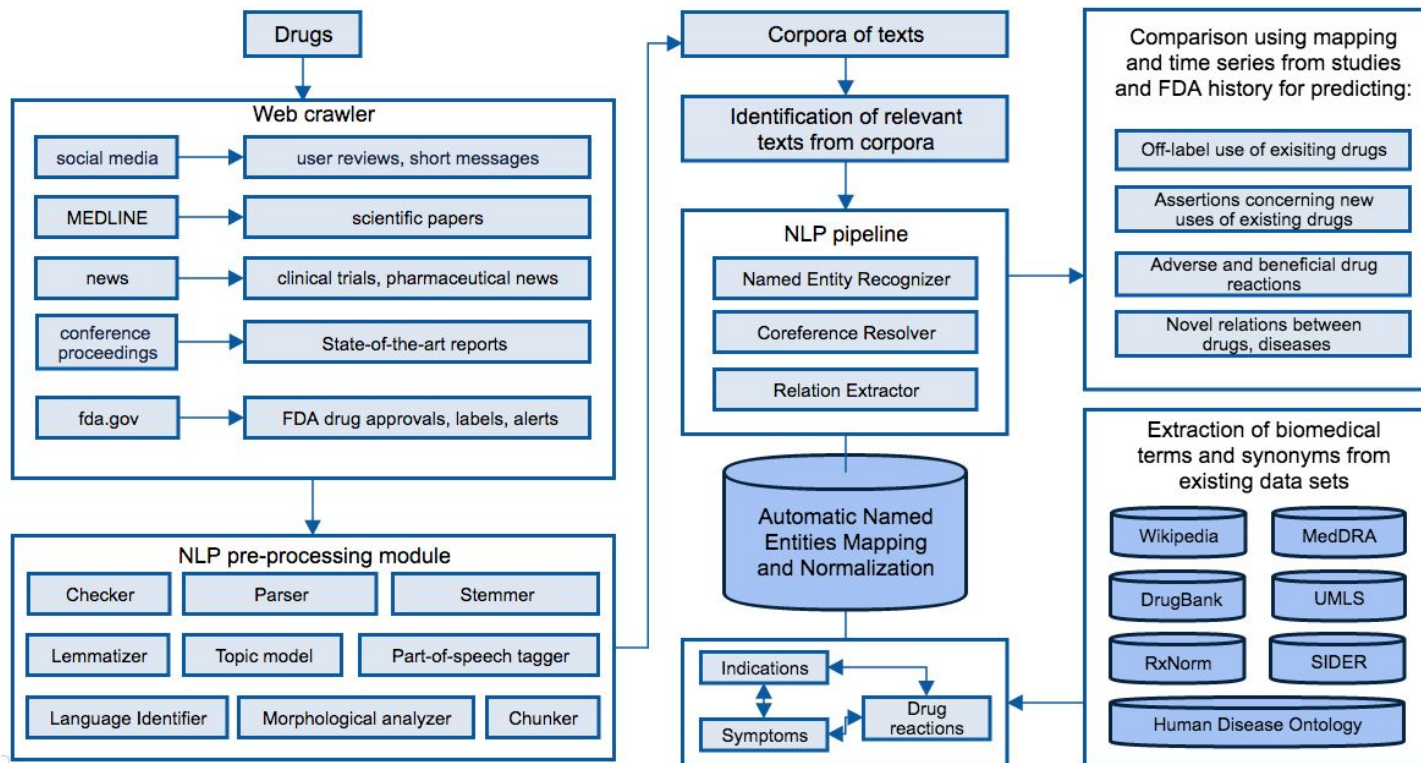
# Cheminformatics & Text Mining & Data Mining



Image credit to
A. Tropsha (UNC USA)

# NLP pipeline for Disease Extraction

# Examples



Disease [Indication]
Такому **гипертонику** со стажем как я, **Барбовал** посоветовала принимать на ночь мой лечащий врач.
Medication [Drugname]    Medication [Route]

Disease [Indication]                          Disease [Indication]                    Disease [Indication]
**Бессонница** стала меня мучить недавно, после перенесенного **стрессового состояния, появился** **страх,**

Medication [Drugname]        Disease [BNE-Pos]    Disease [BNE-Pos]    Disease [BNE-Pos]
**Барбовал** помогает не только **успокоится,** **заснуть,** но и **снижает давление.**

Medication [Drugname]            Disease [BNE-Pos]
Входящий в состав **фенобарбитал** оказывает **сосудорасширяющее действие.**

Medication [Drugname]    Disease [BNE-Pos]
**Валидол** **успокаивает центральную нервную систему.**

Medication [Drugname]
User Reviews & Ratings - **Mucinex DM** oral
TEXT

Disease [BNE-Pos][Indication]    Medication [Drugname]    Disease [ADE-Neg][Worse]    Disease [ADE-Neg][Worse]    Disease [ADE-Neg][Worse]
My coughing fits have lessened since taking **Mucinex DM,** however I **feel like I am high,** **drowsy,** and have **diarrhea.**

Disease [Indication][Unknown]
I'll take something else next time I have a **cough,** not worth it.