# DATA MINING BY ONTOLOGIES CONSTRUCTION IN BIOMEDICAL DOMAIN

Ingrid Petrič

Tanja Urbančič

Bojan Cestnik

DEPARTMENT OF
KNOWLEDGE
TECHNOLOGIES
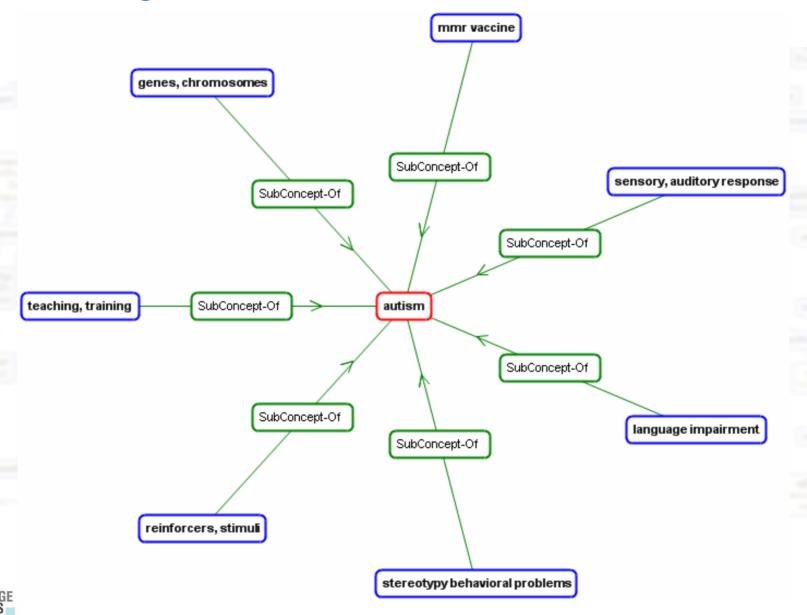Jožef Stefan Institute

# Overview

- Challenges:
  - To obtain systematic insight into domain of interest
  - To gain implicit knowledge from biomedical literature

- Focus domain: autism, a spectrum of pervasive developmental disorders

- Data mining support system: OntoGen

- Data source: PubMed

# Input data source

# Data mining: Identification of domain structure I.

# Data mining: Identification of domain structure II.

# Hypothesis generation: Swanson's model

Literature on agent A that causes phenomenon B

**+**

Literature on agent B that influences phenomenon C

---

Hypothesis: agent A may influence phenomenon C

# Investigation for rare data relations

Dataset of articles from PubMed Central

↓

Data mining by OntoGen

↓

*.txt.stat files

↓

Rare terms

```
'TEXT_BOX_CLICKS'013264.        1        'TEXT_BOX'013265.    1        'TEXTBOOKS'013266.        1
'TETRAHYDROPTERIN'013267.       1        'TETRAHYDROBIOPTERIN_BIOSYNTHESIS_REGENERATING'013268.    1
'TETRAHYDROBIOPTERIN_BIOSYNTHESIS'013269.    1        'TETRAHYDROBIOPTERIN'013270.        1
'TEST_SERIES'013271.    1        'TEST_SAM'013272.    1        'TEST_PICTURE'013273.    1
'TEST_DERIVED_TRANSFER'013274.    1        'TENSION'013275.    1        'TENNIS_BALL'013276.    1
'TELEPHONED_RADIO'013277.    1        'TEACHING_VOCAS'013278. 1        'TEACHING_TRIAL'013279. 1
'TEACHING_RECEPTION'013280.    1        'TEACHERS_BEHAVIOR'013281.        1        'TDI'013282.    1
'TBS'013283.    1        'TASTE_PREFER'013284.    1        'TASK_SEQUENCE'013285.    1        'TASK_INTERSPERSED'0
13286.    1        'TASK_INSTRUCTOR'013287.        1        'TARGETED_TOPOGRAPHY'013288.    1
'TARGETED_SOCIAL_SKILLS'013289. 1        'TARGETED_SOCIAL'013290.        1        'TARGETED_PROMPTED'013291.
1        'TARGETED_COMMUNICATION'013292. 1        'TARGETED_APPEARED'013293.    1        'TAQMAN_RT'013294.
1        'TALK_FRIENDS'013295.    1        'TACTS_TEST'013296.        1        'TACTS_INTRAVERBALS'013297.    1
'T203M'013298.    1        'SYSTEM_HYPOTHESIS'013299.        1        'SYNTHETIC_CHEMICAL'013300.    1
'SYNONYMOUS_SNPS'013301.        1        'SYNAPTOPHYSIN'013302.    1        'SYMPTOMS_COUNTED'013303.        1
'SYMMETRICAL_REQUESTS'013304.    1        'SWISS_PROT'013305.    1        'SWISS'013306.    1
'SWEDISH_MÖBIUS'013307. 1        'SUZIE_MOTHERS'013308.    1        'SUZIE'013309.    1        'SURVIVORS'013310.
1        'SURVEYS_REPORTS'013311.        1        'SURVEYS_POPULATION'013312.    1        'SURVEYS_BASED'0
13313.    1        'SUPERIMPOSITION_PROCEDURE'013314.        1        'SUPERIMPOSED_EDIBLE'013315.
'SUNDBERG_PARTINGTON'013316.    1        'SUCROSE_SOLUTION'013317.    1        'SUCROSE_CITRIC_ACID'013318.
1        'SUCROSE_CITRIC'013319. 1        'SUCCESSFUL_LOW'013320. 1        'SUBTYPES_GROUPS'013321.        1
'SUBSTITUTE_REINFORCEMENT'013322.        1        'SUBJECT_CONTROL'013323.        1        'SUBJECT_ASD'013324.
1        'SUBITIZING'013325.    1        'STUDY_PUNISHED'013326. 1        'STUDENT_DIRECT_INSTRUCT'013327.
```
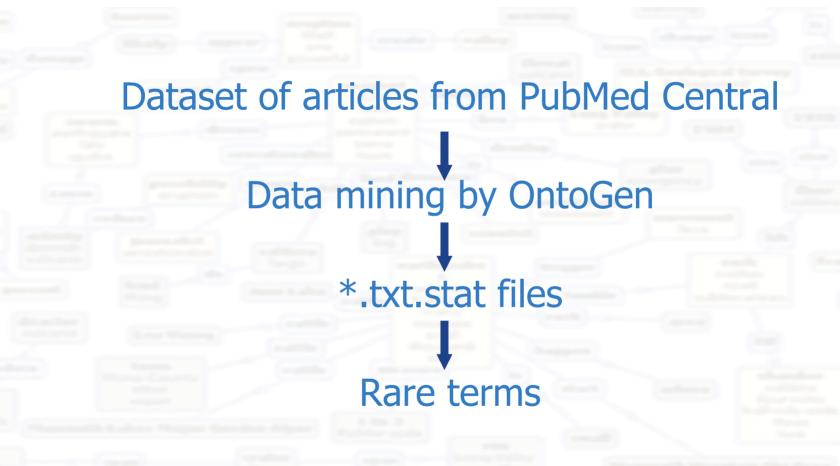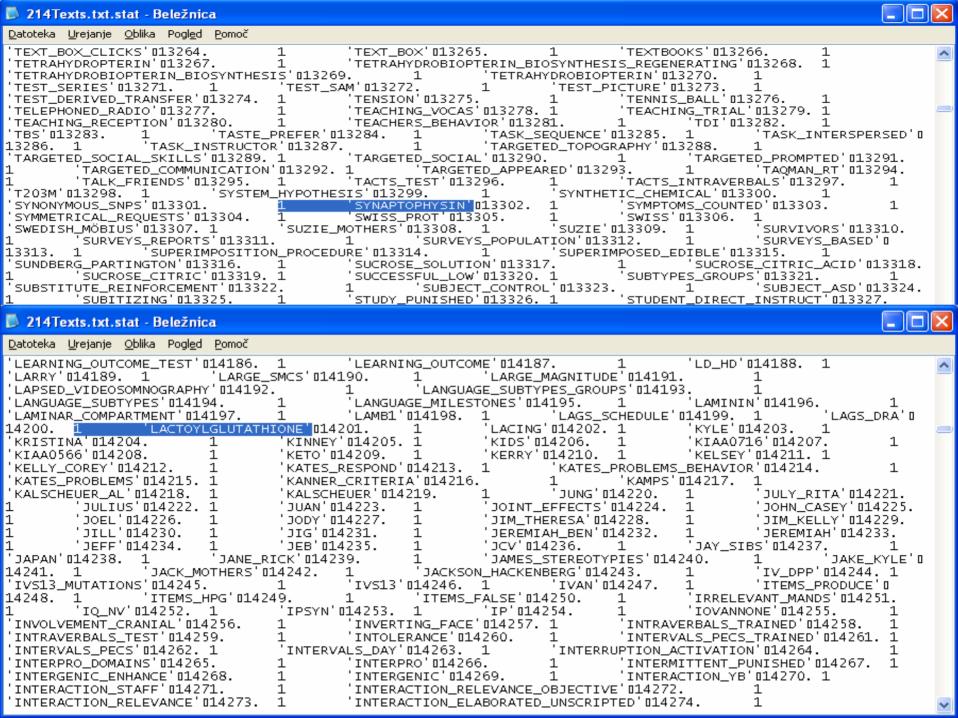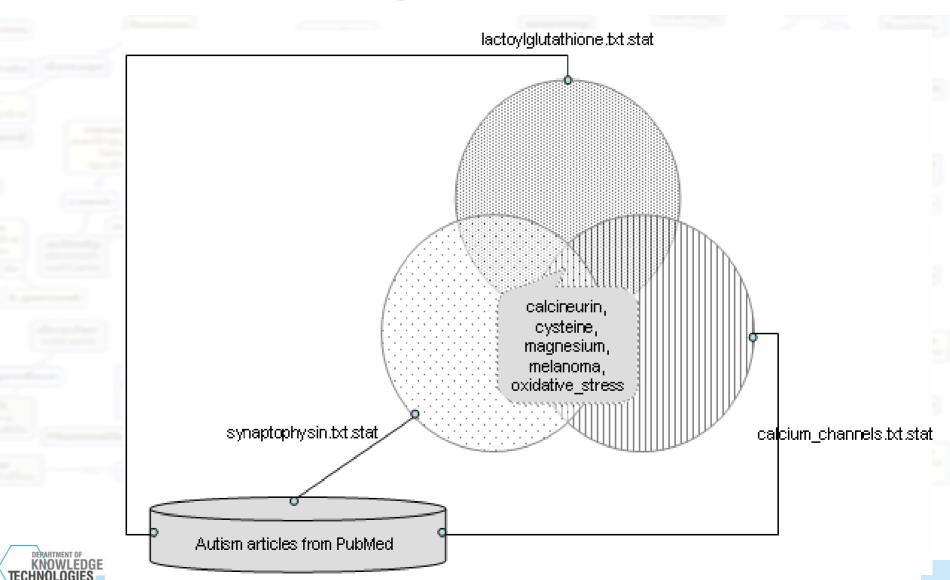
```
'LEARNING_OUTCOME_TEST'014186.    1        'LEARNING_OUTCOME'014187.        1        'LD_HD'014188.    1
'LARRY'014189.    1        'LARGE_SMCS'014190.        1        'LARGE_MAGNITUDE'014191.        1
'LAPSED_VIDEOSOMNOGRAPHY'014192.        1        'LANGUAGE_SUBTYPES_GROUPS'014193.        1
'LANGUAGE_SUBTYPES'014194.        1        'LANGUAGE_MILESTONES'014195.        1        'LAMININ'014196.        1
'LAMINAR_COMPARTMENT'014197.    1        'LAMB1'014198. 1        'LAGS_SCHEDULE'014199.    1        'LAGS_DRA'0
14200.    1        'LACTOYLGLUTATHIONE'014201.        1        'LACING'014202. 1        'KYLE'014203.    1
'KRISTINA'014204.        1        'KINNEY'014205. 1        'KIDS'014206.    1        'KIAA0716'014207.        1
'KIAA0566'014208.        1        'KETO'014209.    1        'KERRY'014210.    1        'KELSEY'014211. 1
'KELLY_COREY'014212.        1        'KATES_RESPOND'014213. 1        'KATES_PROBLEMS_BEHAVIOR'014214.        1
'KATES_PROBLEMS'014215. 1        'KANNER_CRITERIA'014216.        1        'KAMPS'014217. 1
'KALSCHEUER_AL'014218.        1        'KALSCHEUER'014219.        1        'JUNG'014220.        1        'JULY_RITA'014221.
1        'JULIUS'014222. 1        'JUAN'014223.        1        'JOINT_EFFECTS'014224.    1        'JOHN_CASEY'014225.
1        'JOEL'014226.    1        'JODY'014227.    1        'JIM_THERESA'014228.        1        'JIM_KELLY'014229.
1        'JILL'014230.    1        'JIG'014231.    1        'JEREMIAH_BEN'014232.        1        'JEREMIAH'014233.
1        'JEFF'014234.    1        'JEB'014235.    1        'JCV'014236.    1        'JAY_SIBS'014237.        1
'JAPAN'014238. 1        'JANE_RICK'014239.        1        'JAMES_STEREOTYPIES'014240.        1        'JAKE_KYLE'0
14241. 1        'JACK_MOTHERS'014242.    1        'JACKSON_HACKENBERG'014243.        1        'IV_DPP'014244. 1
'IVS13_MUTATIONS'014245.        1        'IVS13'014246. 1        'IVAN'014247.    1        'ITEMS_PRODUCE'0
14248. 1        'ITEMS_HPG'014249.    1        'ITEMS_FALSE'014250.        1        'IRRELEVANT_MANDS'014251.
1        'IQ_NV'014252. 1        'IPSYN'014253. 1        'IP'014254.    1        'IOVANNONE'014255.        1
'INVOLVEMENT_CRANIAL'014256.    1        'INVERTING_FACE'014257. 1        'INTRAVERBALS_TRAINED'014258.    1
'INTRAVERBALS_TEST'014259.    1        'INTOLERANCE'014260.    1        'INTERVALS_PECS_TRAINED'014261. 1
'INTERVALS_PECS'014262. 1        'INTERVALS_DAY'014263. 1        'INTERRUPTION_ACTIVATION'014264.        1
'INTERPRO_DOMAINS'014265.    1        'INTERPRO'014266.    1        'INTERMITTENT_PUNISHED'014267. 1
'INTERGENIC_ENHANCE'014268.    1        'INTERGENIC'014269.    1        'INTERACTION_YB'014270. 1
'INTERACTION_STAFF'014271.    1        'INTERACTION_RELEVANCE_OBJECTIVE'014272.        1
'INTERACTION_RELEVANCE'014273. 1        'INTERACTION_ELABORATED_UNSCRIPTED'014274.    1
```

# Connecting terms

| Word | Total | calcium_channels | lactoylglutathione | synaptophysin |
|------|-------|------------------|--------------------|---------------|
| 'CYPRUS_INSTITUTE' | 1 | 1 | | |
| 'CYS' | 2 | 1 | | 1 |
| 'CYS_CYS' | 1 | | | 1 |
| 'CYS_CYS_CYS' | 1 | | | 1 |
| 'CYSTEINE' | 3 | 1 | 1 | 1 |
| 'CYSTEINE_MOTIF' | 1 | | | 1 |
| 'CYSTEINE_RESIDUE' | 1 | 1 | | |
| 'CYSTEINE_RESIDUES' | 2 | | 1 | 1 |
| 'CYSTEINE_RICH' | 1 | 1 | | |
| 'CYSTEINE_RICH_DOMAINS | 1 | 1 | | |
| 'CYSTEINE_STRING' | 1 | | | 1 |
| 'CYSTEINE_STRING_PROTE | 1 | | | 1 |
| 'CYSTEINYL' | 1 | | 1 | |
| 'CYSTEINYL_GLYCINE' | 1 | | 1 | |

Record: |◄ ◄ 5654 ► ►| ►* of 30071

# Connecting sets of literature

# Conclusions

- Ontology construction is useful for systematical datasets exploration.
- OntoGen statistical data can lead to discovery of potentially meaningful information.
- Expert's involvement is crucial for speeding up the selections and the evaluations of candidate hypotheses.