# "From (Text) Mining to Models: Applying Large-Scale Text Mining on Patents and Electronic Patient Records"

## Fraunhofer SCAI

**Martin Hofmann-Apitius**
Head of the Department of Bioinformatics
Fraunhofer Institute for Algorithms and Scientific
Computing (SCAI)

# Fraunhofer: Applied Research for Industrial Applications

Fraunhofer stands for:

- sustainable (applied) research
- focus on contract research and innovation
- bridging between excellent academic research and industrial application
- clear mission towards improving and fostering innovation
- research done with the idea in mind to generate added value in a commercial sense

# SCAI Department of Bioinformatics: R&D in a nutshell

**Fraunhofer SCAI Department of Bioinformatics R&D activities:**

1. Information extraction in the **life sciences**:

   I. Text Mining - Recognition of named entities & relationships in text

   II. Image Mining - Reconstruction of chemical information from chemical structure depictions

2. Disease modelling (focus on neurodegenerative diseases)

3. eScience, Grid-/Cloud- Computing and HPC (Clus~

*Making Scientific Content available for Computing*

**Fraunhofer**

SCAI

# Productive Use of Large Compute Infrastructures

**High Throughput Extraction of Scientific**

**Information from Full Text Sources:**

The UIMA-HPC Project

Fraunhofer

SCAI

# Vision

## Scientific Challenge:

The knowledge in Chemistry, Biology and Pharmaceutical Sciences growths with impressive speed. As a result, the number of publications in these areas is reaching unparalleled dimensions. However, knowledge is being communicated in non-standardised ways.
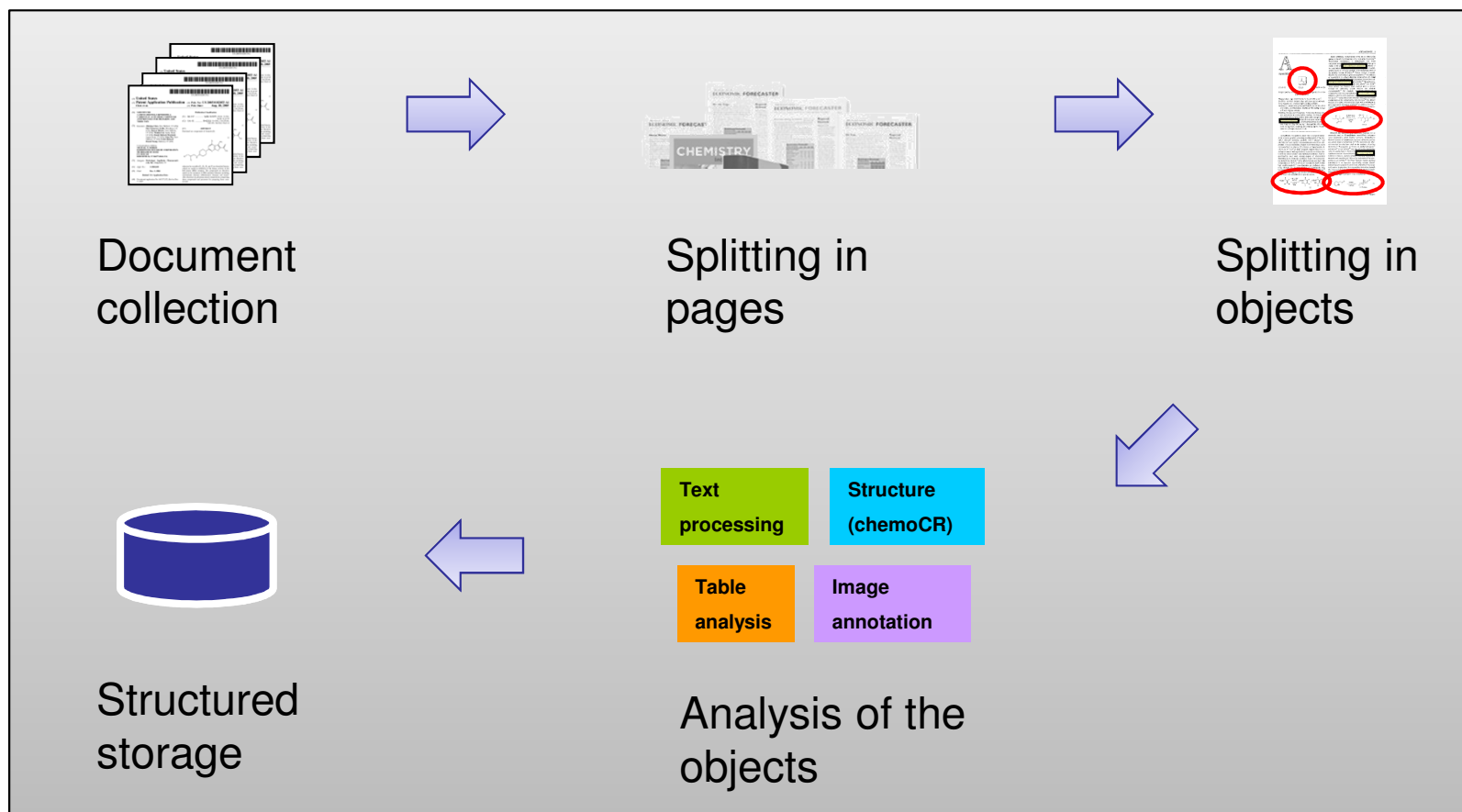
**Relevant knowledge sources are not well standardized, let alone is the knowledge structured.  This limits the ability to query knowledge sources.**

## Problem-solving approach:

Development of technology that – based on HPC – allows for high throughput extraction of structured information from unstructured knowledge sources
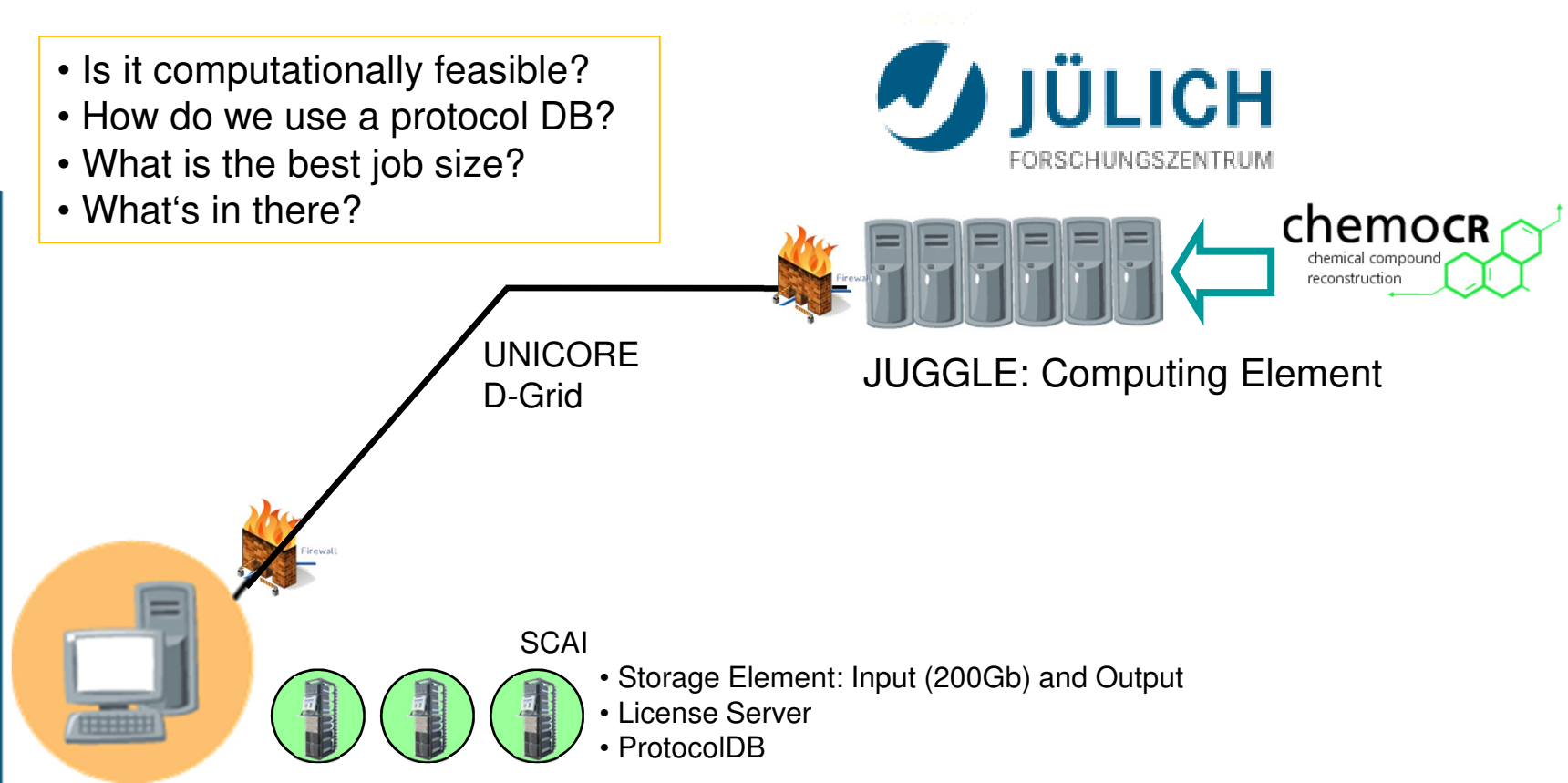
**Structured Knowledge**

# UIMA-HPC

## Use case scenario:
### automatic patent structuring

**Document collection**

**Splitting in pages**

**Splitting in objects**

| | |
|---|---|
| Text processing | Structure (chemoCR) |
| Table analysis | Image annotation |

**Structured storage**

**Analysis of the objects**

# Previous study: The grand patent challenge 2009

- Is it computationally feasible?
- How do we use a protocol DB?
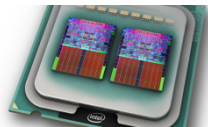- What is the best job size?
- What's in there?

**UIMA-HPC**

JÜLICH
FORSCHUNGSZENTRUM

chemoCR
chemical compound
reconstruction

JUGGLE: Computing Element

UNICORE
D-Grid

SCAI

- Storage Element: Input (200Gb) and Output
- License Server
- ProtocolDB

# UIMA-HPC

## Technical Issues and Pitfalls

- User accession rights (files, scheduler, installed tools and libs, ...)
- Firewall (ports: MySQL, denial of service attack, time outs, ...)
- Missing files (NFS down, package lost, not installed, ...)
- Too many requests on license server
- Too many connections in database
- Ressources (reservation, priorities, ...)

# UIMA AS in the context of HPC

Support of many-core architecture
- several instances of a service
- eff. usage of shared memory (JVM)
- asynchronous execution

Support of clusters
- several remote services (eg SOAP)
- communication via JMX and http

Control via pre-configured parameters
- CAS pool size
- casMultiplier poolSize
- …

Manual tuning

# UIMA-HPC

## Problem-Solving Approach



UNICORE Rich Client
- Grid Workflow Experts-

Web Portal
- Researchers -

Gateway

UNICORE

UNICORE services

Storage and Resource Management

UIMA

Software

**Third parties**

Hardware

HPC

http://www.unicore.eu/

Multi Core

http://uima.apache.org/

# System Architecture

# First Prototype: PDF Annotation of Patents

IMA-HPC

Lesezeichen

- ChemicalStructure
- SourceDocumentInformation
- Synthesis
- DISEASE
- IUPAC
  - 2-benzyl hexanoic acid 1
  - 2-benzylidene hexanoic acid 2
  - thioether 3
  - carbon 4
  - 2-(3,5-diphenethoxybenzyl)hexanoic acid 5
  - thioether 6
  - phosphonate 7
  - phosphonate 8
  - glucose 9
  - pyrimidine 10
  - dihydroxybenzaldehyde 11
  - dimethylthiocarbamoyl chloride 12
  - Dihydroxybenzaldehyde 13
  - 2-[(3,5-diphenethoxyphenyl)thio]hexanoic acid 14
  - Dihydroxybenzaldehyde 15
  - Dihydroxybenzaldehyde 16
  - ethanol 17
  - phenylethyl bromide 18
  - Phloroglucinol 19

Bioorganic & Medicinal Chemistry 19 (2011) 5372–5382

Contents lists available at SciVerse ScienceDirect

## Bioorganic & Medicinal Chemistry

journal homepage: www.elsevier.com/locate/bmc

ELSEVIER

## SAR studies of acidic dual γ-secretase/PPARγ modulators

Martina Hieke [a,1], Julia Ness [b], Ramona Steri [a], Christine Greiner [c], Oliver Werz [d], Manfred Schubert-Zsilavecz [a], Sascha Weggen [b,*], Heiko Zettl [a,*,1]

[a] Institute of Pharmaceutical Chemistry, Goethe-University Frankfurt, Max-von-Laue-Str. 9, D-60438 Frankfurt am Main, Germany
[b] Department of Neuropathology, Heinrich-Heine-University Duesseldorf, Moorenstrasse 5, D-40225 Duesseldorf, Germany
[c] Department of Pharmaceutical Analytics, Eberhard-Karls-University Tuebingen, Auf der Morgenstelle 8, D-72076 Tuebingen, Germany
[d] Institute of Pharmacy, Friedrich-Schiller-University Jena, Philosophenweg 14, D-07743 Jena, Germany

ABSTRACT

A novel set of dual γ-secretase/PPARγ modulators characterized by a ... presented. Synthetic efforts were focused on the variation of the substitution pattern of the central ben-zene. Finally, we obtained a new class of 2,5-disubstituted 2-benzylidene hexanoic acid derivatives, which act as dual γ-secretase/PPARγ modulators in the low micromolar range. We have explored broad SAR and successfully improved the dual pharmacological activity and the selectivity profile a ... targets such as NOTCH and COX. Compound 17 showed an IC$_{50}$ Aβ42 = 2.4 μM and an I ... and could be a valuable tool to further evaluate ... concept of dual γ-secretase/P ... animal models of Alzheimer's disease.

© 2011 Elsevier Ltd.

2-benzyl hexanoic acid

**1. Introduction**

... proliferation and B-cell differentiation in pre ... mice.[6,7] As an appealing alternative to GSIs, so- ...

Chemical Popups

Statistics of highlighted entities of class IUPAC (page 1)

| Entity | Frequency |
|---|---|
| thioether | 5 |
| benzene | 5 |
| phloroglucinol | 3 |
| pyrimidine | 3 |
| Dihydroxybenzaldehyde | 3 |
| Phloroglucinol | 2 |
| phenylethanol | 2 |
| carbon | 2 |
| sodium hydride | 2 |
| 2-(3,5-phenethoxyphenoxy)hexanoic acid | 2 |
| phosphonate | 2 |
| 2-benzylidene hexanoic acid | 1 |
| phenol | 1 |
| dimethylthio-carbamoylchloride | 1 |
| 3-hydroxy-4-(2-cyclohexylethoxy)benzaldehyde | 1 |
| glutathione | 1 |
| 2-[(3,5-diphenethoxyphenyl)thio]hexanoic acid | 1 |
| 2-benzyl hexanoic acid | 1 |
| 3-(2-cyclohexylethoxy)-4-phenethoxybenzal | 1 |
| tetramethylsilane | 1 |
| dimethylthiocarbamoyl chloride | 1 |

Chemical Index

Annotations + Linkouts

Overview

18. April 2012

# Performance of Fraunhofer SCAI in international Benchmarking Competitions

| | BiTeM10 PAx | SCAI10C IENTP | SCAI10C ITENT | SCAI10C ITNP | SCAI10C ITTOK | SCAI10N RMENT | SCAI10N RMNP | SCAI10N RMTOK | York10C aPA01 |
|---|---|---|---|---|---|---|---|---|---|
| ■ map | 0.2657 | 0.4121 | 0.2336 | 0.2065 | 0.0947 | 0.0665 | 0.0551 | 0.0172 | 0.0136 |
| ■ ndcg | 0.4975 | 0.5834 | 0.4119 | 0.3764 | 0.1888 | 0.2547 | 0.2224 | 0.0868 | 0.0885 |
| ■ p30 | 0.3485 | 0.4554 | 0.2794 | 0.2485 | 0.1126 | 0.1169 | 0.1 | 0.0366 | 0.0309 |
| ■ r100 | 0.4724 | 0.5491 | 0.3596 | 0.3265 | 0.1511 | 0.1974 | 0.1743 | 0.0629 | 0.0583 |
| ■ mrr | 0.6121 | 0.7153 | 0.5324 | 0.4769 | 0.1956 | 0.3456 | 0.3088 | 0.1133 | 0.1022 |
| ■ bpref | 0.6592 | 0.7075 | 0.5468 | 0.511 | 0.2804 | 0.4171 | 0.3702 | 0.1536 | 0.1681 |

**TRECCHEM 2009 → better not talk about it ...**

**TRECCHEM 2010 → winner of the Prior Art Search Task**

**TRECCHEM 2011 → winner of the Technology Survey Task**

**I2B2 challenge 2010 → rank 3 (out of 22) in the concept id task**

**TRECMED 2011 → rank 4 (amongst 24 participants)**

Fraunhofer
SCAI

# Direct Usage of Unstructured Information Sources for Disease Modelling

**From Medline Mining**

**to**

**Modelling Neurodegenerative Diseases**

# Why Modelling of Neurodegeneration?

In 2009 the Federal Government of Germany decided to start a new research centre that focuses on translational research on neurodegenerative diseases. In fact, neurodegenerative diseases (Alzheimer, Parkinson, Multiple Sclerosis; Epilepsy; „rare" NDDs)

The total costs of Alzheimer is estimated to exceed 20 trillion US$ in the US in the years between 2020 - 2050. (source: Alzheimer.org). Current costs / year in the US (according to Alzheimer.org): **183 billion US$**

The incidence rate of Alzheimer and other dementias is almost 50% in the population older than 85 years.  Next generation will regularly have a life span of >100 years.

# Diseases specific mortality rate

**Changes in selected causes of death in USA , 2000-2010[1]**



[1] www.alz.org

# The Starting Conditions

What we have:

- An ontology capturing relevant knowledge on Alzheimer´s Disease (ADO)

- An ontology representing and integrating brain regions and cell types (BRCO)

- A method for the automated identification of hypotheses in text based on regular expressions

- An excellent machinery for biomedical text mining (ProMiner) with top performing gene and protein name recognition

Fraunhofer
SCAI

# Alzheimer´s Disease Ontology (ADO)

Alzheimer's ontology:

- Captures more than 700

  classes/concepts

- BFO already implemented

# Brain Region and Cell-type Ontology (BRCO)



**Current state: more than 3000 concepts; more than 5000 synonyms**

# Expression of Speculative Statements in Scientific Text

**Kallikrein-related peptidase 6** in Alzheimer's disease and vascular dementia.

PubMed 20846516 **Authors:** Ashby, Emma L; Kehoe, Patrick G; Love, Seth **Date:** 2010-12- **Journal:** Brain research **Affiliation:** Dementia research Group, Institute of Clinical Neurosciences, Clinical Science at North Bristol, University of Bristol, UK.

Statistics ☐ Select ID with comment:

Human kallikrein-related peptidase 6 (KLK6) is highly expressed in the central nervous system. Although the physiological roles of this serine protease are unknown, in vitro substrates include amyloid precursor protein and components of the extracellular matrix, which are altered in neurological disease, particularly Alzheimer's disease (AD). We have compared KLK6 expression in post-mortem brain tissue in AD, vascular dementia (VaD) and controls. We studied the distribution of KLK6 in the temporal cortex and white matter by immunohistochemistry, and measured KLK6 mRNA and protein levels in the frontal and temporal cortex from 15 AD, 15 VaD and 15 control brains. Immunohistochemistry showed KLK6 to be restricted to endothelial cells. After adjustment for variations in vessel density by measurement of factor VIII-related antigen, we found KLK6 protein and mRNA levels to be significantly decreased in the frontal but not the temporal cortex in AD. In VaD, KLK6 protein level was significantly increased in the frontal cortex. Our findings suggest that an altered KLK6 expression may contribute to vascular abnormalities in AD and VaD.

**Hypothesis =   KLK6   +   may contribute   +   AD**

Fraunhofer
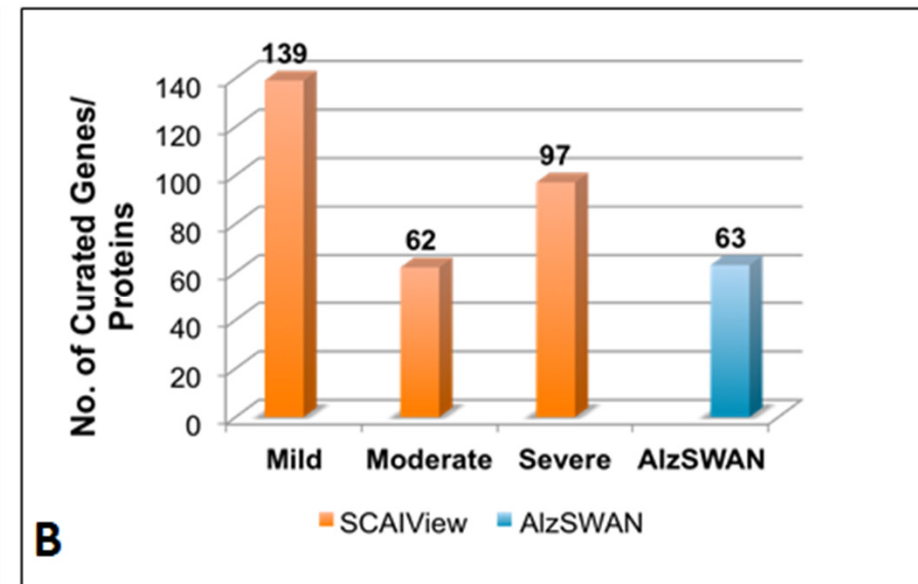SCAI

# Hypotheses finder∩ AD ontology ∩Human genes and proteins

# Performance of Hypotheses finder

| S.No | Data type | Source | Precision | Recall | F score |
|------|-----------|--------|-----------|--------|---------|
| 1 | 200 abstracts related to Alzheimer's | PubMed | 0.84 | 0.86 | 0.85 |
| 2 | 2 full text articles related to Alzheimer's | Journal of Medical Hypotheses | 0.85 | 0.88 | 0.86 |
| 3 | 143 abstracts related to Alzheimer's | Alzswan/PubMed | 0.90 | 0.97 | 0.93 |
| 4 | 100 abstracts related to Epilepsy | PubMed | 0.96 | 0.91 | 0.94 |
| 5 | 100 abstracts related to Parkinson's | PubMed | 0.90 | 0.93 | 0.92 |

# Performance of Hypotheses finder

# The Knowledge – Discovery Strategy

**Alzheimer ontology**

**Hypotheses finder**

SCAI**VIEW** *pharma*

**Knowledge retrieval and extraction**

**Genes and proteins associated to hypotheses**

**Hypothetical sentences on disease mechanism**

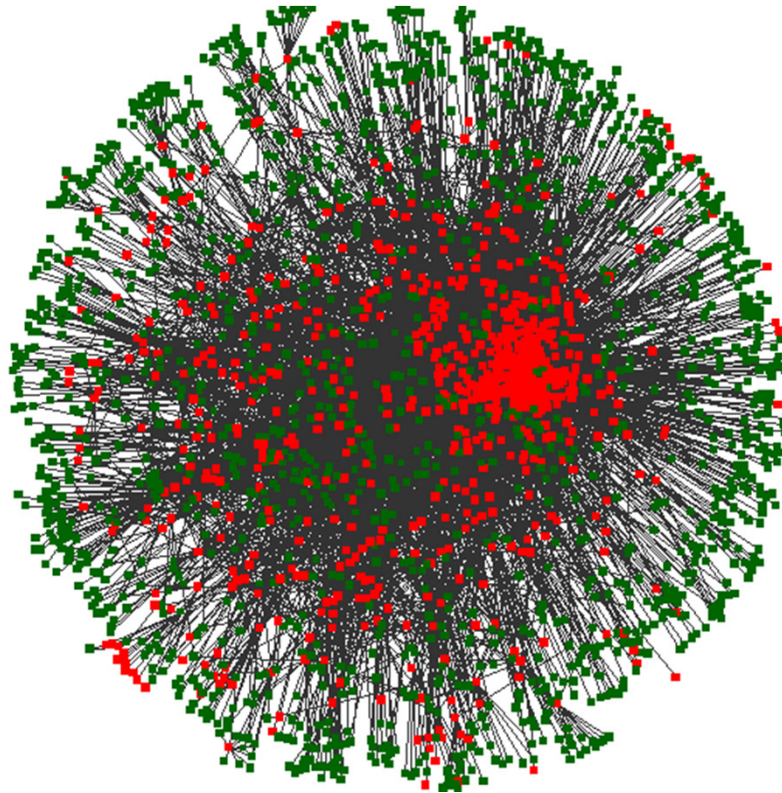**Molecular representation of disease hypotheses**

Fraunhofer
SCAI

# Analysis of hypotheses patterns across disease stages
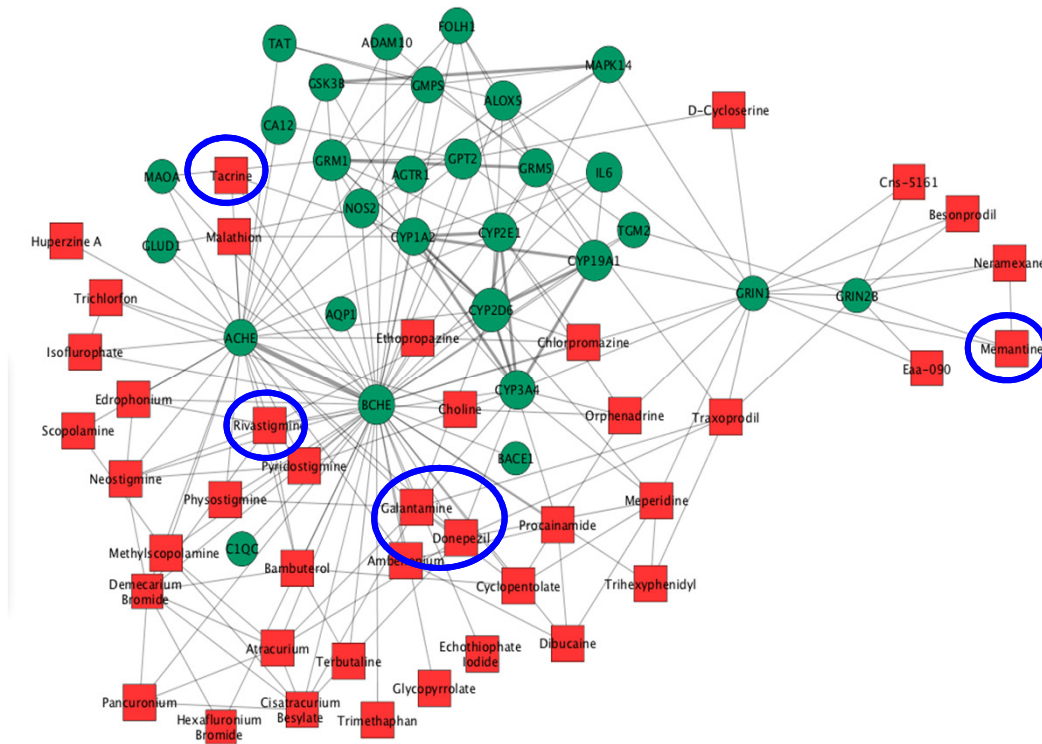
# Current Work at SCAI: NDD Pharmacome



NDD Pharmacome Network:
Total nodes: 2,153
Total edges: 11,630
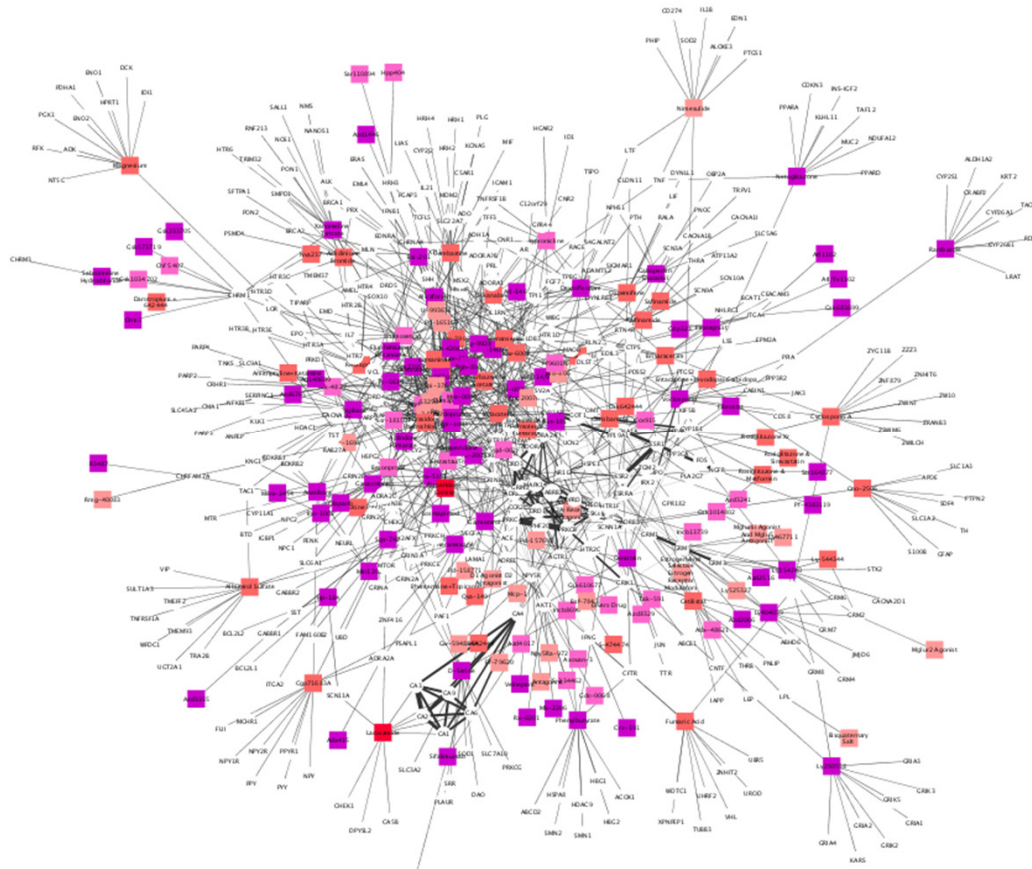
- Drugs
- Targets

Fraunhofer
SCAI

# Using the Pharmacome for VS



1. Galantamine
2. Memantine
3. Rivastigmine
4. Tacrine
5. Donepezil

Drugs having structural similarity and targets having binding site similarities are connected by edges.
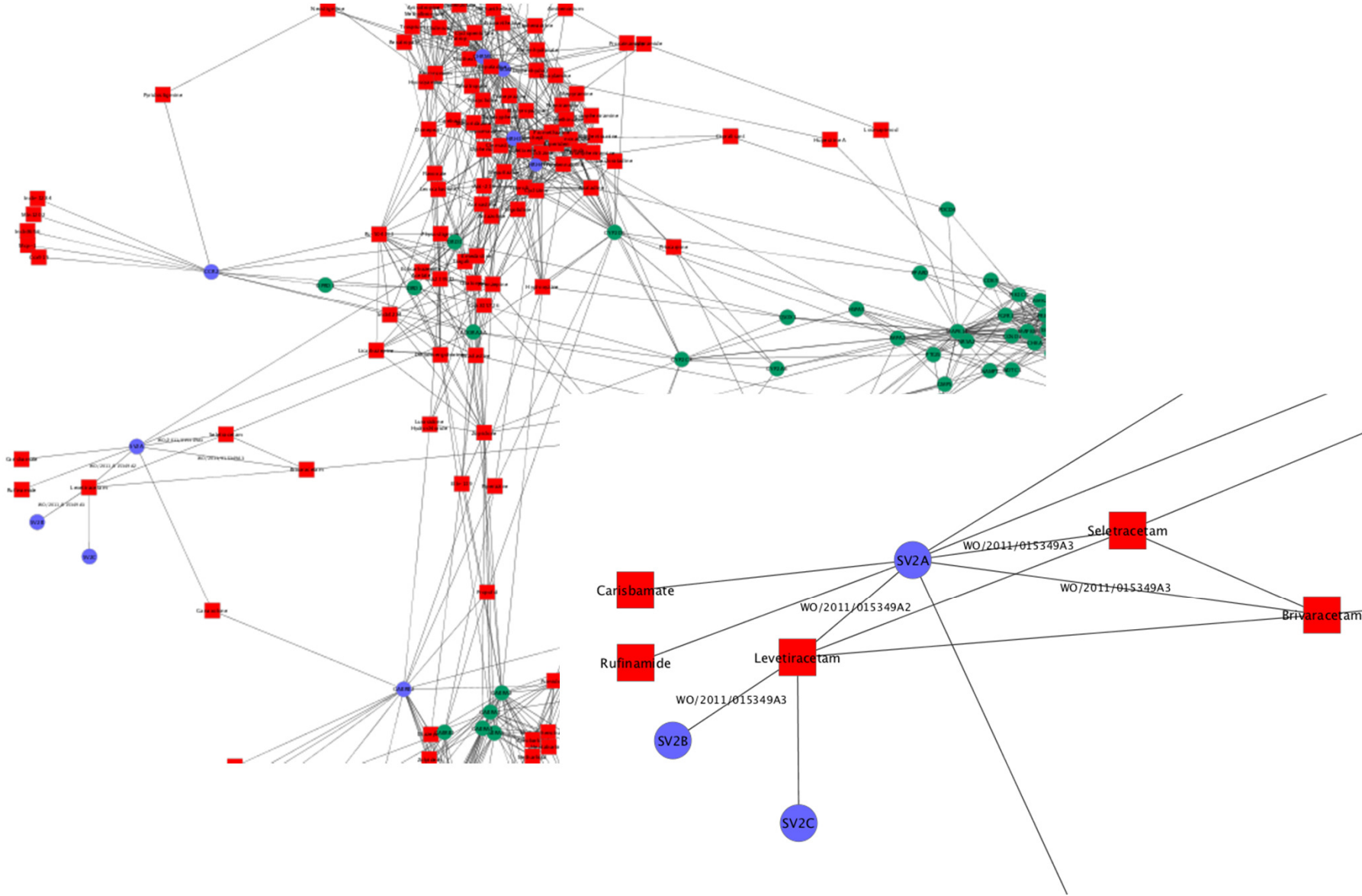
# What is Competition doing? Competitive intelligence and strategic partnering



Drugs in various phases

- Phase I
- Phase II
- Phase III
- Phase IV

Fraunhofer

SCAI

# Brain Drug Portfolios in the Pharmacome

# Summary

**Fraunhofer SCAI Department of Bioinformatics stands for:**

- Advanced technologies in text- and data mining, disease modelling in the area of neurodegeneration and high performance computing

- We are using our competence in HPC to enable large-scale information extraction from full text documents (focus on patents)

- Internal usage of technologies (information extraction; distributed and high performance computing) for biomedical application: modelling of neurodegenerative diseases

- Information integration and aggregation in models: the brain pharmacome

Fraunhofer

SCAI