# Direct Use of Information Extraction from Scientific Text for Modeling and Simulation in the Life Sciences

**SCAI**

**Fraunhofer** Institute
Algorithms and
Scientific Computing

Martin Hofmann-Apitius

**Department of Bioinformatics**
**Fraunhofer Institute for Algorithms and Scientific Computing (SCAI)**

# Fraunhofer-Campus Schloss Birlinghoven



## Institutes

- Algorithms and Scientific Computing **SCAI**

- Intelligent Analysis and Information Systems **IAIS**

- Applied Information Technology **FIT**

600 Scientists, 200 Students

Linked to Universities Bonn, Aachen and Cologne

**Fraunhofer** Institute Algorithms and Scientific Computing

b-it

universität**bonn**

# Direct Use of Information Extraction from Scientific Text

# for Modeling and Simulation in the Life Sciences

SCAI

**Fraunhofer** Institute
Algorithms and
Scientific Computing

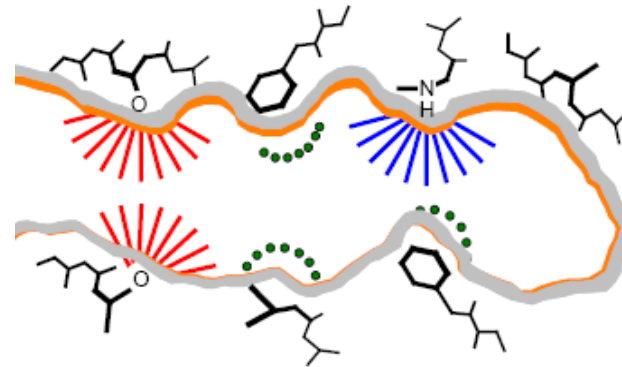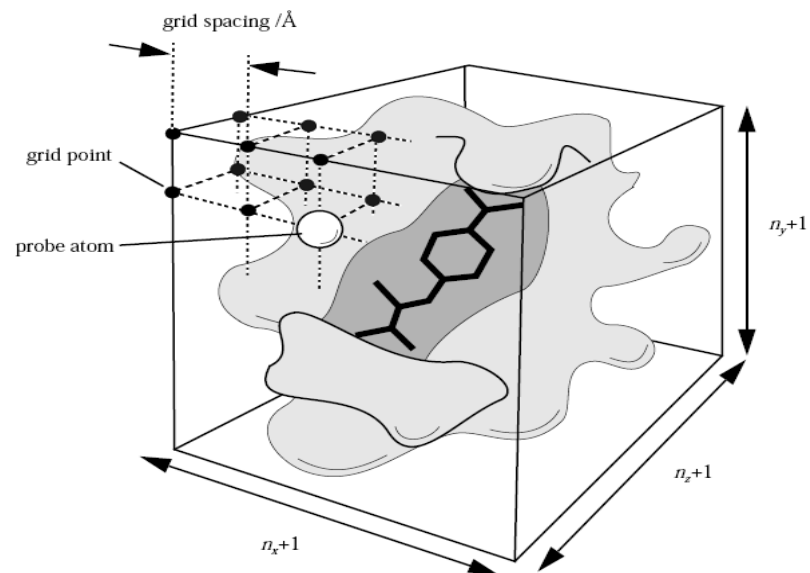# Paradigm Changes in the Life Sciences

**In the Life Sciences, we currently observe significant paradigm changes**

- the "omics" paradigm has lead to a flood of data and a flood of publications

- a single researcher cannot keep track with all the relevant (and related) literature any more

- everything is connected; genetics, molecular biology, biochemistry, pharmaceutical chemistry and organic chemistry are "networked"

- Biology and Chemistry and Medicine are more and more turning into quantitative sciences, described by mathematical models and with the option of using simulation (*in silico* experimentation)

SCAI

Fraunhofer Institute
Algorithms and
Scientific Computing

# Virtual Screening as an Example

SCAI

Fraunhofer Institute
Algorithms and
Scientific Computing

AutoDock - affinity grid maps for each atom type

grid spacing /Å

grid point

probe atom

$n_y+1$

$n_z+1$

$n_x+1$

FlexX – different types of interactions (interaction points)

# Large-Scale in silico Experimentation: eSciences

**Large-scale *in silico* experimentation & eSciences: the WISDOM project**

- Large-scale virtual screening for novel drugs against Malaria (*plasmodium falciparum*)

- International collaboration based on the EGEE grid computing infrastructure (with thousands of CPUs connected worldwide)

- Millions of protein – drug interaction simulations; equivalent to more than 80 years of permanent computing on a single CPU

- However, WISDOM was based on a rather "physical" scenario: a virtual representation of a chemical compound is positioned into the binding geometry of a protein, of which the 3D-structure is known.

Fraunhofer Institute
Algorithms and
Scientific Computing

# Simulation and Knowledge – Driven Approaches

**First principle – based sciences**

- Physics

- Engineering

- Physico-Chemistry


→ Based on mathematical models

→ Simulation approaches can be easily applied

**Descriptive, empirical sciences**

- Biology

- Pharmaceutical Chemistry

- Medicine


→ Based on knowledge represented in the literature

→ Very complex, difficult to simulate

SCAI

Fraunhofer Institute Algorithms and Scientific Computing

# Making Use of The Wealth of Knowledge that is Out There

**SCAI**

**Fraunhofer** Institute
Algorithms and
Scientific Computing

# Literature as a Main Source for Knowledge in the Life Sciences

- Biology and Medicine are still to a large extend *empirical* sciences

- Complex: very high number of entities and relationships

- Lots of data on genes and proteins in databases

- However, biodatabases do only comprise data and not necessarily knowledge (data + models)

- Expressiveness of natural language in text is much higher; therefore scientific text is a much better source for biomedical knowledge

→ How do we get access to the knowledge and how can we model it?

SCAI

**Fraunhofer** Institute
Algorithms and
Scientific Computing

# Biological, Medical and Chemical Objects in Text

## Abstract

**Background:** Mast cell-derived prostaglandin $D_2$ ($PGD_2$), may contribute to eosinophilic inflammation and mucus production in allergic asthma. Chemoattractant receptor homologous molecule expressed on $TH_2$ cells (CRTH2), a high affinity receptor for prostaglandin $D_2$, mediates trafficking of $TH_2$-cells, mast cells, and eosinophils to inflammatory sites, and has recently attracted interest as target for treatment of allergic airway diseases. The present study involving mice explores the specificity of CRTH2 antagonism of TM30089, which is structurally closely related to the dual TP/CRTH2 antagonist ramatroban, and compares the ability of ramatroban and TM30089 to inhibit asthma-like pathology.

**Methods:** Affinity for and antagonistic potency of TM30089 on many mouse receptors including thromboxane $A_2$ receptor mTP, CRTH2 receptor, and selected anaphylatoxin and chemokines receptors were determined in recombinant expression systems *in vitro*. *In vivo* effects of TM30089 and ramatroban on tissue eosinophilia and mucus cell histopathology were examined in a mouse asthma model.

**Results:** TM30089, displayed high selectivity for and antagonistic potency on mouse CRTH2 but lacked affinity to TP and many other receptors including the related anaphylatoxin C3a and C5a receptors, selected chemokine receptors and the cyclooxygenase isoforms 1 and 2 which are all recognized players in allergic diseases. Furthermore, TM30089 and ramatroban, the latter used as a reference herein, similarly inhibited asthma pathology *in vivo* by reducing peribronchial eosinophilia and mucus cell hyperplasia.

**Conclusion:** This is the first report to demonstrate anti-allergic efficacy *in vivo* of a highly selective small molecule CRTH2 antagonist. Our data suggest that CRTH2 antagonism alone is effective in mouse allergic airway inflammation even to the extent that this mechanism can explain the efficacy of ramatroban.

Fraunhofer Institute
Algorithms and
Scientific Computing

# Technologies for Information Extraction from Literature

During the last five years, substantial progress has been made in the area of automated text analysis. In particular in life science informatics there is a strong community developing new methods and tools for the automated recognition and extraction of information from scientific literature.

Our group at Fraunhofer SCAI has developed three tools that enable mining in literature:

- **ProMiner,** a solution for named entity recognition based on rules and dictionaries (a "reading machine" for biomedical text)

- **ChemoCR**, a software that identifies chemical structure depictions in full text (a "reading machine" for chemical structure depictions)

- **SCAIView**, a text mining environment that supports end-users

Fraunhofer Institute
Algorithms and
Scientific Computing

# ProMiner & SCAIView



**SCAIVIEW**

Entity Tree View, select Entity Class to view and search

Documents | Entity | Analysis

Alzheimer

+ Human Genes / Proteins
- Chromosomal Location
- STS Marker
- non Normalized SNP
- Normalized SNP
- Normalized CRF SNP
- Drug Names
- IUPAC-like
- OMIM Reference
- HuGeNet Genetic Associations
- Epigenetics
- Arabidopsis Genes
- Mouse Genes
- Interaction Verbs
- MeSH Disease
- Relations
- @neurIST Ontology

Select Confidence:

1 ☐  2 ☐  3 ☐  4 ☐  5 ☑

## Your Search:
(once the color changed from red to green the query is ready)

- Use **Medline** as the Document Base.
- Limit Corpus using Full Text Search **'Alzheimer'**
- Entities of the class **Human Genes / Proteins** must be in the document
- Display entities of type **Human Genes / Proteins** in Entity View.

## Help

Reset Search

Show (this) Information Screen

Start Search

Filter Results

Expand / Collapse Tree Viewing

Show results of this entity class

Show Saved Search Queries

A manual can be obtained here

A demo video as **mpeg** or **quick time**

## Steps to pose a Query (Use Firefox version >2.0)

1. Enter a Full Text Search into the grey field located below the icons on the top-left. (click on the blue arrow to access standard searches)

   > Select 'Intracranial AND Aneurysm*'

2. Click once on the name of an item in the tree to include it in the entity tree (click on it again to not include it and again to disregard it). Use the «» button to increase the size of the tree's viewing area.

# ProMiner & SCAIView

# ProMiner & SCAIView

# ProMiner & SCAIView

**Making Use of The Wealth of Knowledge that is Out There:**

**<u>Example 1:</u>**

**Using Text-based Information for the Prediction of Pharmaceutical Activities of Drugs**

(Master Thesis of Harsha Gurulingappa, B-IT and Fraunhofer SCAI)

SCAI

**Fraunhofer** Institute
Algorithms and
Scientific Computing

# Task: ATC Classification of yet Unclassified Drugs

Goal of this study:

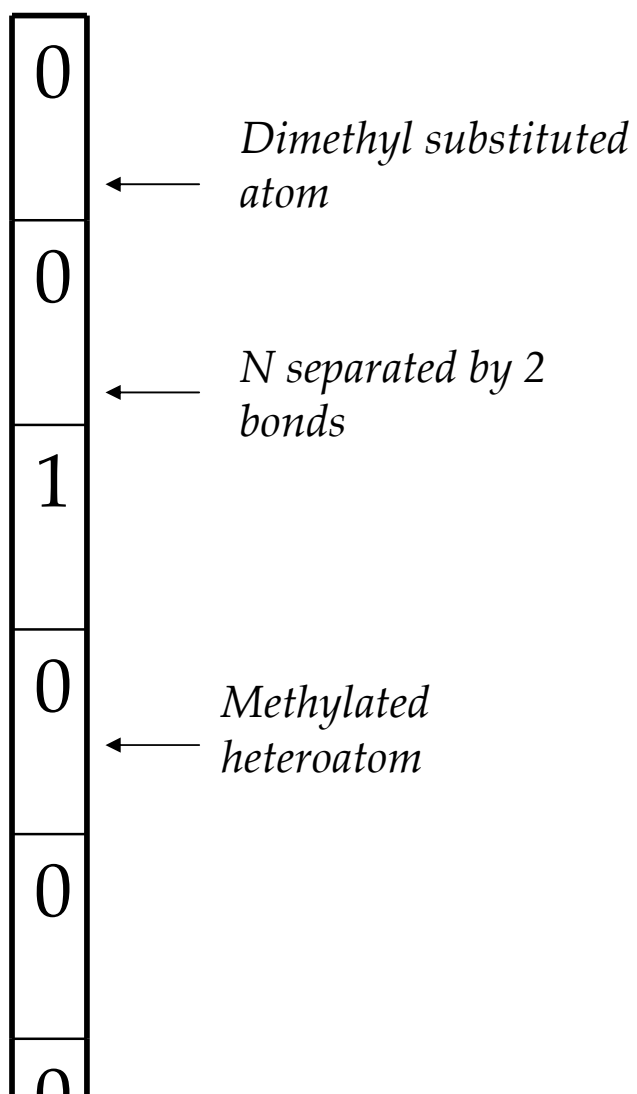- developed of a method for predicting putative ATC classes for unclassified drug terms.

- develop a new paradigm for strategies aiming at identifying potential secondary applications for existing drugs.

- Use of textual features/evidences for characterization of drugs

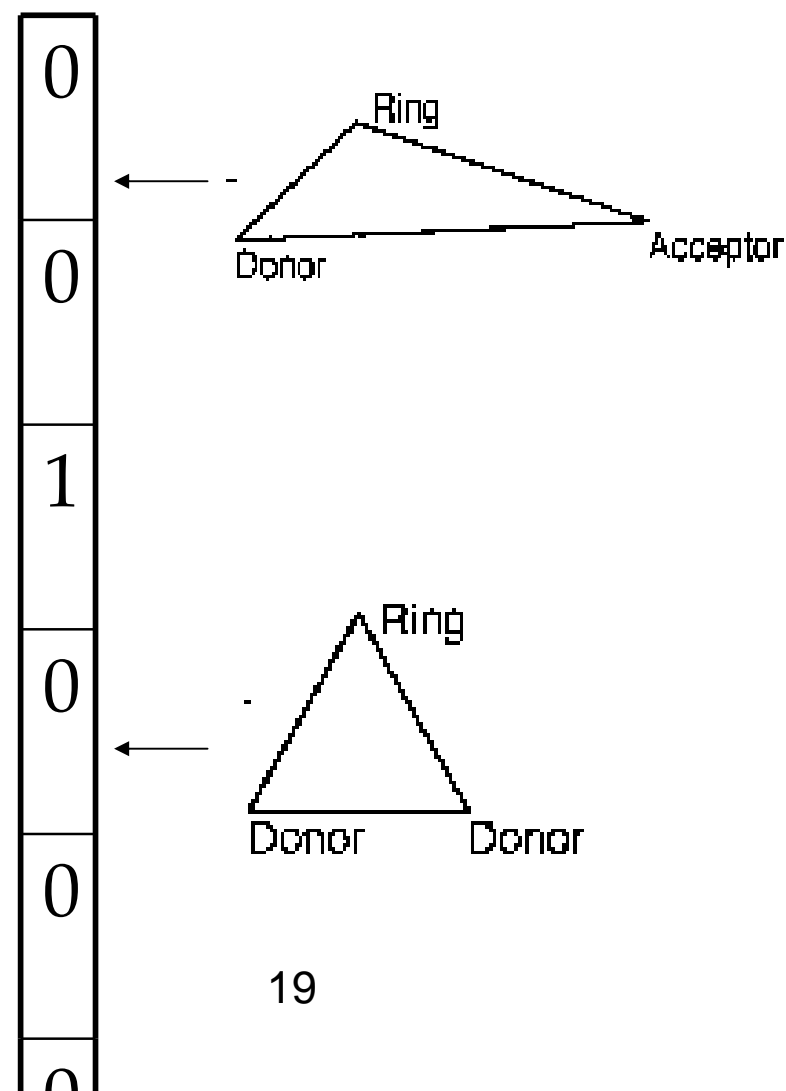→ Disadvantage: Highly dependent on Information Availability.

SCAI
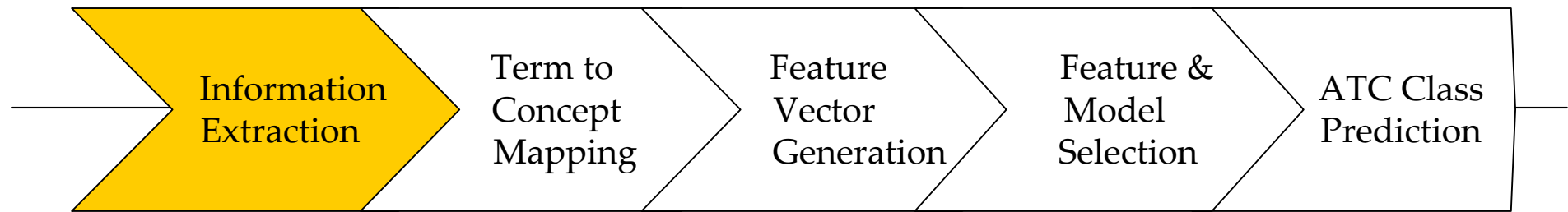
Fraunhofer Institute
Algorithms and
Scientific Computing

# Fingerprints: Method for the Prediction of Chemical Properties

MACCS keys: Structural Keys

TGT keys: 3 Point Pharmacophore based fingerprints

| 0 |
|---|
| 0 |
| 1 |
| 0 |
| 0 |

*Dimethyl substituted atom* ←

*N separated by 2 bonds* ←

*Methylated heteroatom* ←

| 0 |
|---|
| 0 |
| 1 |
| 0 |
| 0 |

Information Extraction → Term to Concept Mapping → Feature Vector Generation → Feature & Model Selection → ATC Class Prediction

Timolol is a beta adrenoceptor blocker

A vasodilator like propatyl nitrate, can open the ...
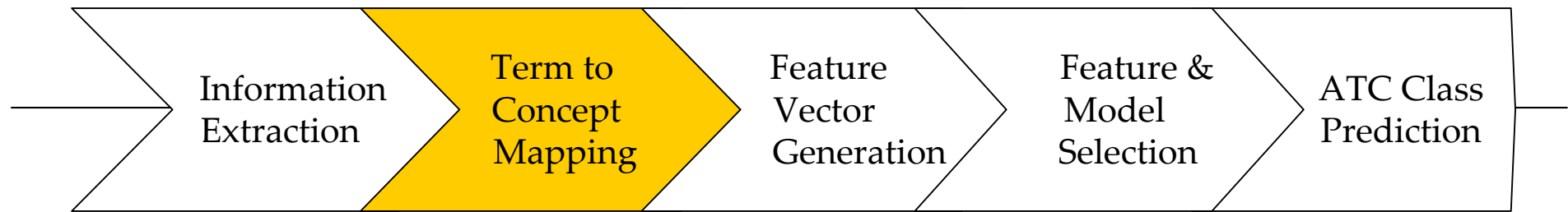
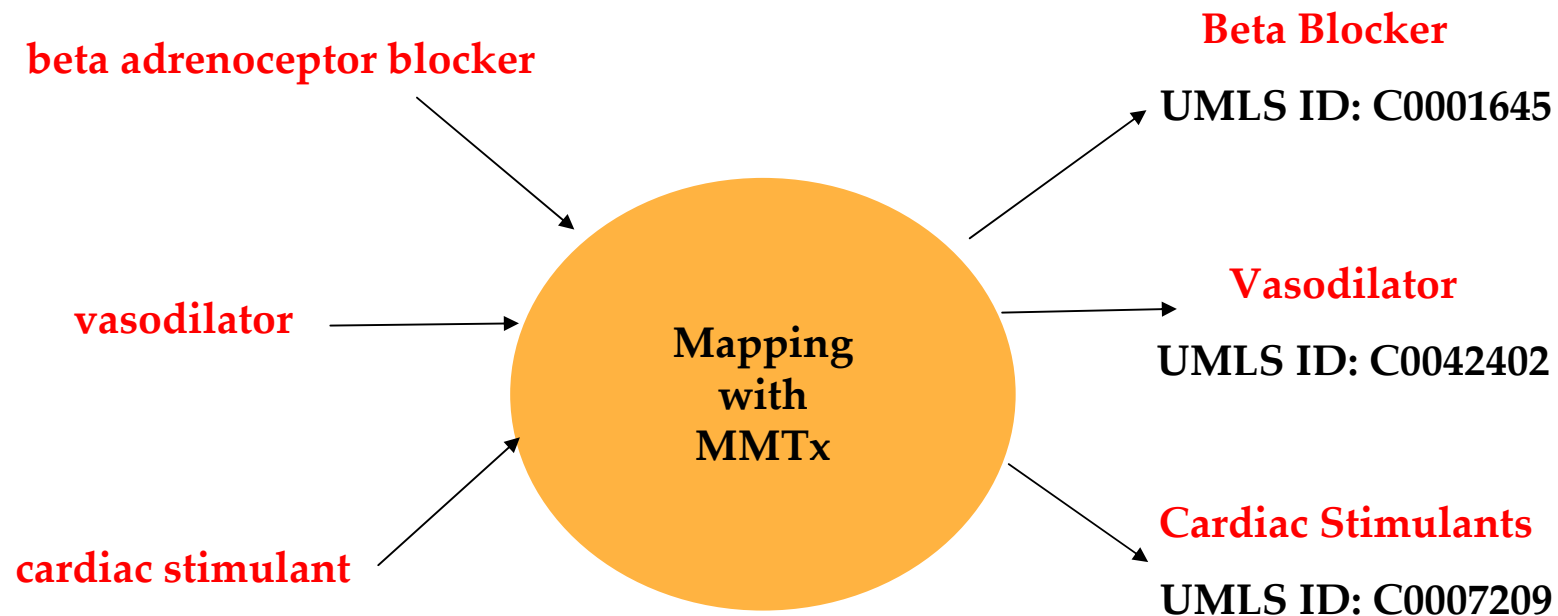Digitoxin, a cardiac stimulant is responsible for ...
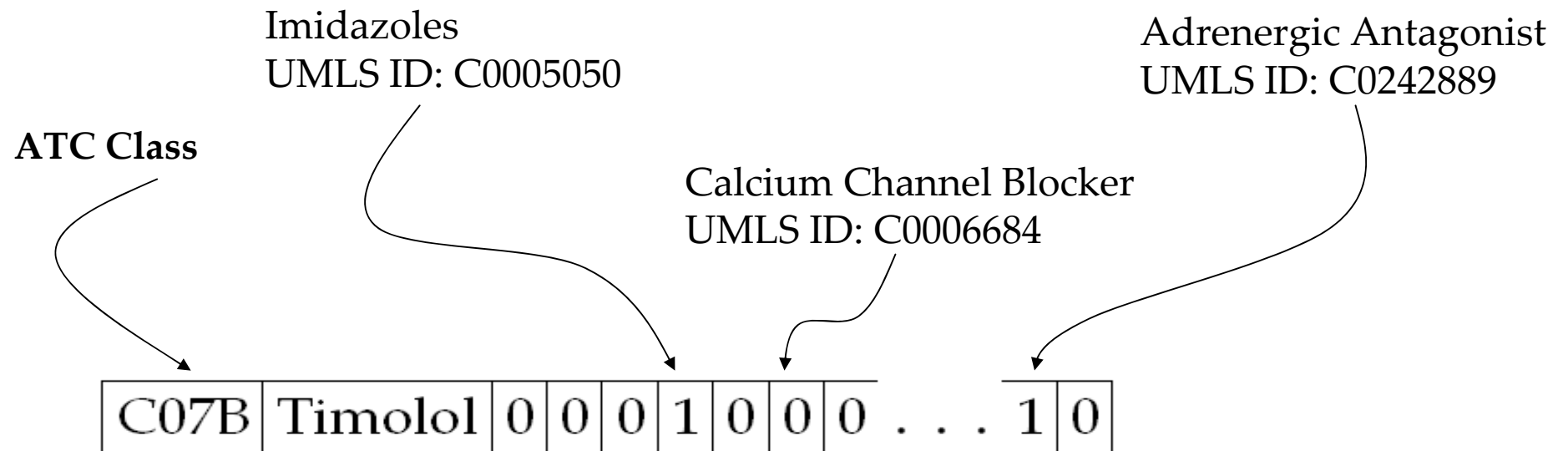
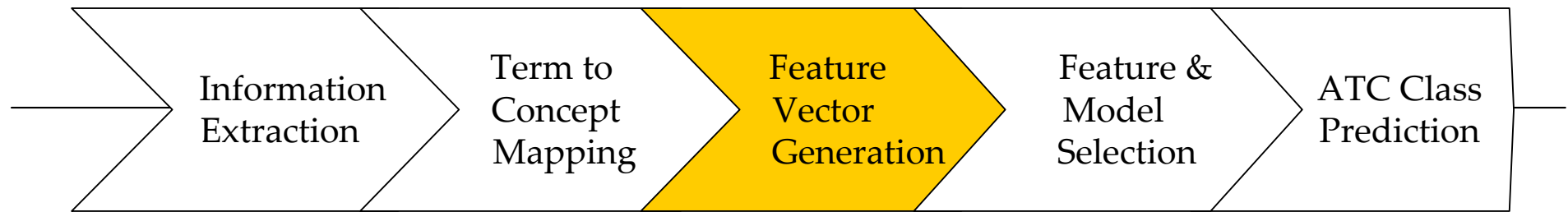Drug Term    Property Term    Free Text

Hearst Patterns: Express direct relationship between drugs and their properties.

[Kolářik et al., 2007]

Fraunhofer Institute Algorithms and Scientific Computing
SCAI

Map extracted property terms to concepts in UMLS*

**beta adrenoceptor blocker** →

**vasodilator** →

**cardiac stimulant** →

**Mapping with MMTx**

→ **Beta Blocker**
**UMLS ID: C0001645**

→ **Vasodilator**
**UMLS ID: C0042402**

→ **Cardiac Stimulants**
**UMLS ID: C0007209**

SCAI

**Fraunhofer** Institute
Algorithms and
Scientific Computing

*UMLS: Unified Medical Language System

Information Extraction → Term to Concept Mapping → **Feature Vector Generation** → Feature & Model Selection → ATC Class Prediction

**ATC Class**

Imidazoles
UMLS ID: C0005050

Calcium Channel Blocker
UMLS ID: C0006684

Adrenergic Antagonist
UMLS ID: C0242889

| C07B | Timolol | 0 | 0 | 0 | 1 | 0 | 0 | 0 | . . . | 1 | 0 |

**Binary Feature Vector:** '0': Feature Absent & '1': Feature Present

| C07B | Timolol | 0 | 0 | 0 | 7 | 0 | 0 | 0 | . . . | 4 | 0 |

**Weighted Feature Vector:** '0': Feature Absent & '≥1': Corpus Frequency of the Feature

22

| | Information Extraction | | Term to Concept Mapping | | Feature Vector Generation | | Feature & Model Selection | | ATC Class Prediction | |

**Feature Selection**
> Mutual Information
> Chi-square criterion

**Models/Classifiers**
> Naïve Bayes
> Nearest Neighbor
> Decision Tree
> Support Vector Machine

| Rank | Feature/Concept | Concept ID | Chi-square score |
|---|---|---|---|
| 1 | Diuretic | C0012798 | 390.0000 |
| 2 | Anti-arrhythmic | C0003195 | 390.0000 |
| 3 | Dihydroxyphenylalanine | C0012315 | 378.3875 |
| 4 | Steroids | C0338671 | 345.7591 |
| 5 | Cardenolide | C0007143 | 345.7591 |
| 6 | Loop diuretic | C0354100 | 345.7591 |
| 7 | AT1 receptor blocker | C1449680 | 328.5642 |
| 8 | Vasoconstrictor | C0042397 | 321.2956 |
| 9 | Coronary dilator | C0596385 | 317.6199 |
| 10 | Potassium channel agonist | S10000044 | 316.9350 |

Fraunhofer Institute
Algorithms and
Scientific Computing

SCAI

Information Extraction → Term to Concept Mapping → Feature Vector Generation → Feature & Model Selection → ATC Class Prediction

**Concept Based Vs Structure Based Approaches**
**Test Set = 114 Drugs**

Concepts: 77.2%
SuperPred: 53.5%
MACCS Keys: 36.8%
TGT Keys: 35.8%

Number of Drugs

True Predictions: 88, 61, 42, 41
False Predictions: 6, 40, 65, 64
No Predictions: 20, 13, 7, 9

Fraunhofer Institute Algorithms and Scientific Computing

**Making Use of The Wealth of Knowledge that is Out There:**

**Example 2:**

**Using Text-based Information for the Identification of Genes likely to mediate Susceptibility to Breast Cancer**

(Master Thesis of Erfan Younesi, B-IT and Fraunhofer SCAI)

Fraunhofer Institute
Algorithms and
Scientific Computing

# Task: Predicting Networks of Genes / Proteins that are Linked to the Clinical Progression and Outcome of the Disease

Goal of this study:

- Identification of networks of proteins functionally linked to tumorigenesis of breast cancer

- Identification of combinations of nodes in a network that can serve as markers for susceptibility to clinical treatment

→ Vision: using text mining to extract evidences for best clinical practice

SCAI

Fraunhofer Institute
Algorithms and
Scientific Computing

# Decision on Clinical Treatment Strategies for Breast Cancer

Current situation:

- Decisions made up on very few factors (e.g. Lymph Node status; ER+/-)

- Cooperativity of genes and proteins not taken into account

# Definition of a Network of Interacting Proteins in Breast Cancer

Approach:

- Definition of a network of molecular entities strongly associated with breast cancer

- Network based on simple co-occurences in text

28

# Reduction of Complexity: Definition of a Minimum Network

Approach:

- Selection of subgraph

  based on

  - Network topology

  - Functional characterization

  - association with clinical

    outcome

# Identification of Novel Breast Cancer Susceptibility Associations

23 novel associations between the minimum gene set associated with breast cancer susceptibility and clinical outcome could be identified

| Associations | Cocitation frequencies | Novel association (not exists in PIANA) | Evidence of general relation between two genes from Literature (PMID) | Shared GO process | Shared KEGG pathway |
|---|---|---|---|---|---|
| AKT1 - EGFR | 93 | Y | 14981538 -17686159 - 18351692-16774943 - 16419029-16546981- 16288304-15800944 | Nitric oxide anabolism, protein amino acid phosphorylation | MAPK signaling pathway, Focal adhesion, Colorectal cancer, Pancreatic cancer, Glioma |
| TP53 - EGFR | 71 | Y | 18311481 | Cell cycle, Response to stress, regulation of cell proliferation, | MAPK signaling pathway, Colorectal cancer, Pancreatic cancer, Glioma |
| PGR - EGFR | 23 | Y | 1616857-1911227 | regulation of epithelial cell proliferation | -- |
| STAT3 - AKT1 | 16 | Y | 10853013-16288304- 16728588 | -- | Jak-STAT signaling pathway, Adipocytokine signaling pathway, Pancreatic cancer |
| TNF - EGFR | 10 | Y | 9829842-11221831 | regulation of protein amino acid phosphorylation, cell-cell adhesion, regulation of cell proliferation | MAPK signaling pathway |
| CDH1 - PGR | 6 | Y | 16512896 | -- | -- |
| VDR - EGFR | 4 | Y | 16087726-17377416 | skeletal development | -- |

**Making Use of The Wealth of Knowledge that is Out There:**

**Example 3:**

**A Look into the Future: A Computational Grand Challenge in the Area of Patent Mining**

(ongoing collaboration between Fraunhofer SCAI and FZ Jülich)

# Task: Annotation of All Chemical Structure Depictions in All Pharmaceutical Patents from EPO

Goal of this study:

- Feasibility study for large-scale annotation of patents

- Grand computing challenge in the area of knowledge computing

- Demonstration of enhancement of retrieval in the area of chemistry by intelligent software (ChemoCR – chemical structure reconstruction)

→ Vision: using image mining to mine chemical IP at large scale

**SCAI**

**Fraunhofer** Institute
Algorithms and
Scientific Computing

# Reconstruction of Chemical Information

## Step1: PDF Conversion



or

Normalization of image: 250 DPI, grey scale

# Reconstruction of Chemical Information

## Step 2: Page Segmentation



1. Classification of interesting regions
2. Grouping of chemical reaction schemata
3. Transfer of chemical reaction schemata segment to reconstruction module

# Reconstruction of Chemical Information

## Step 3: Reconstruction of Chemical Information



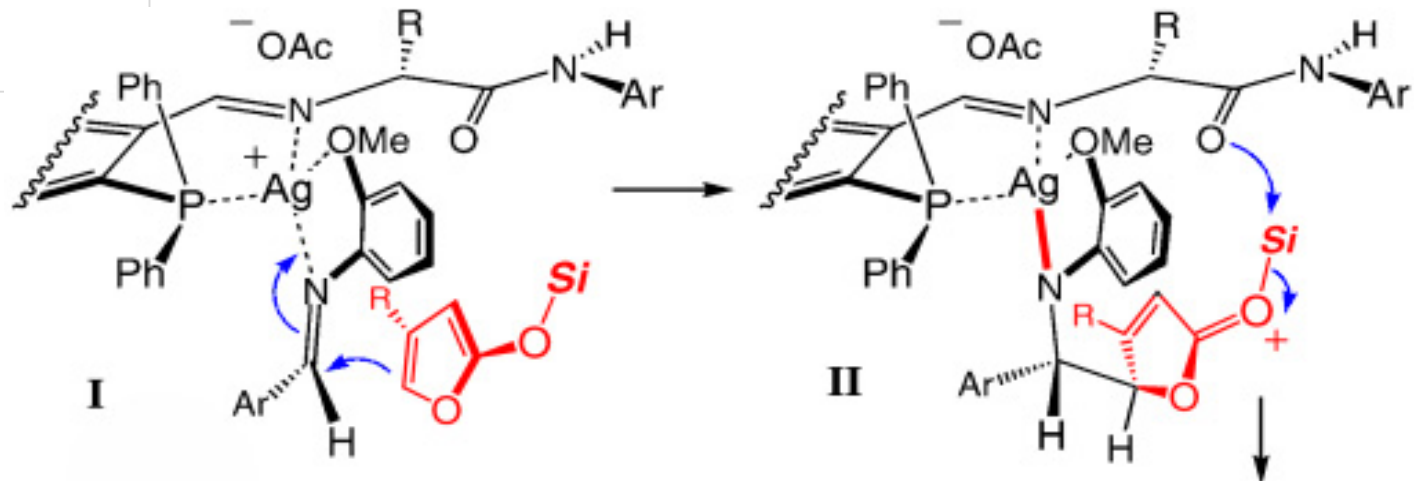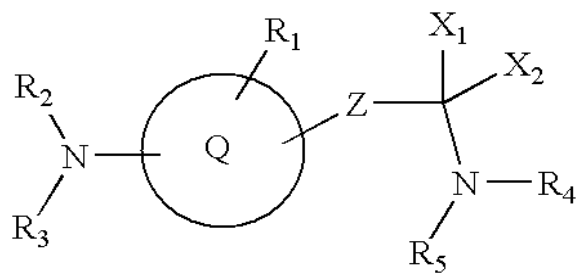1. Classification of characters and symbols
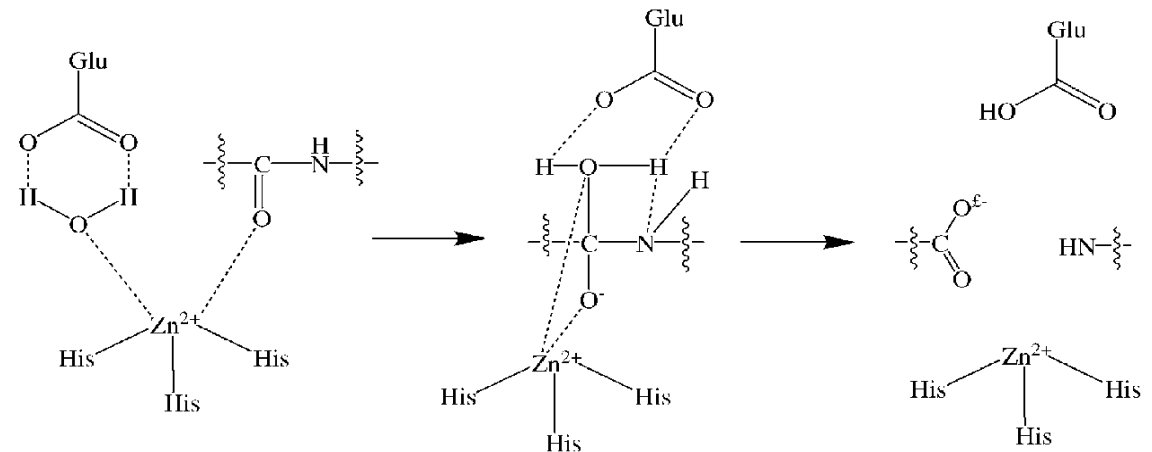2. OCR
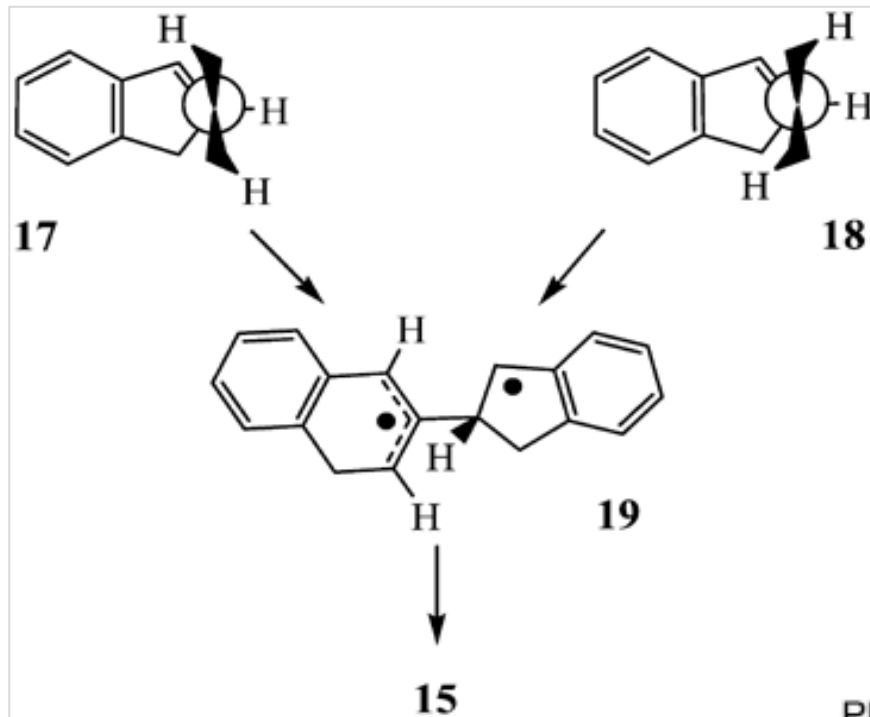3. vectorization
4. Chemical rule set / expert system
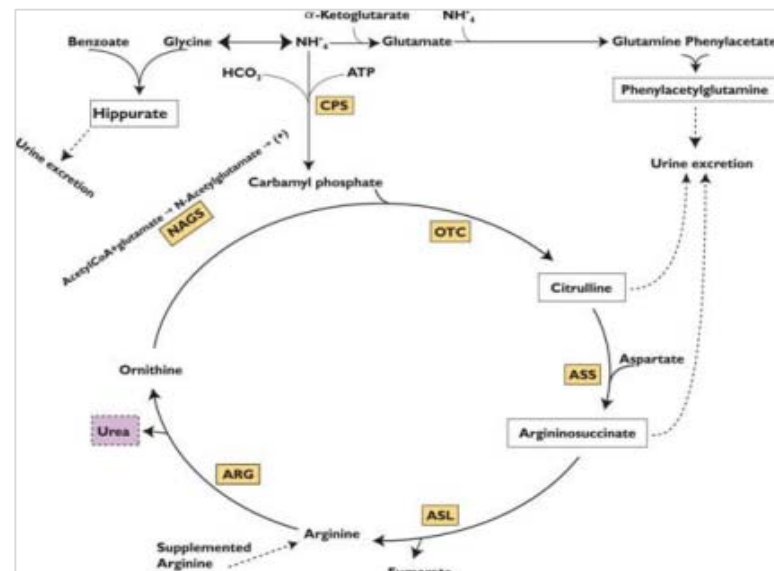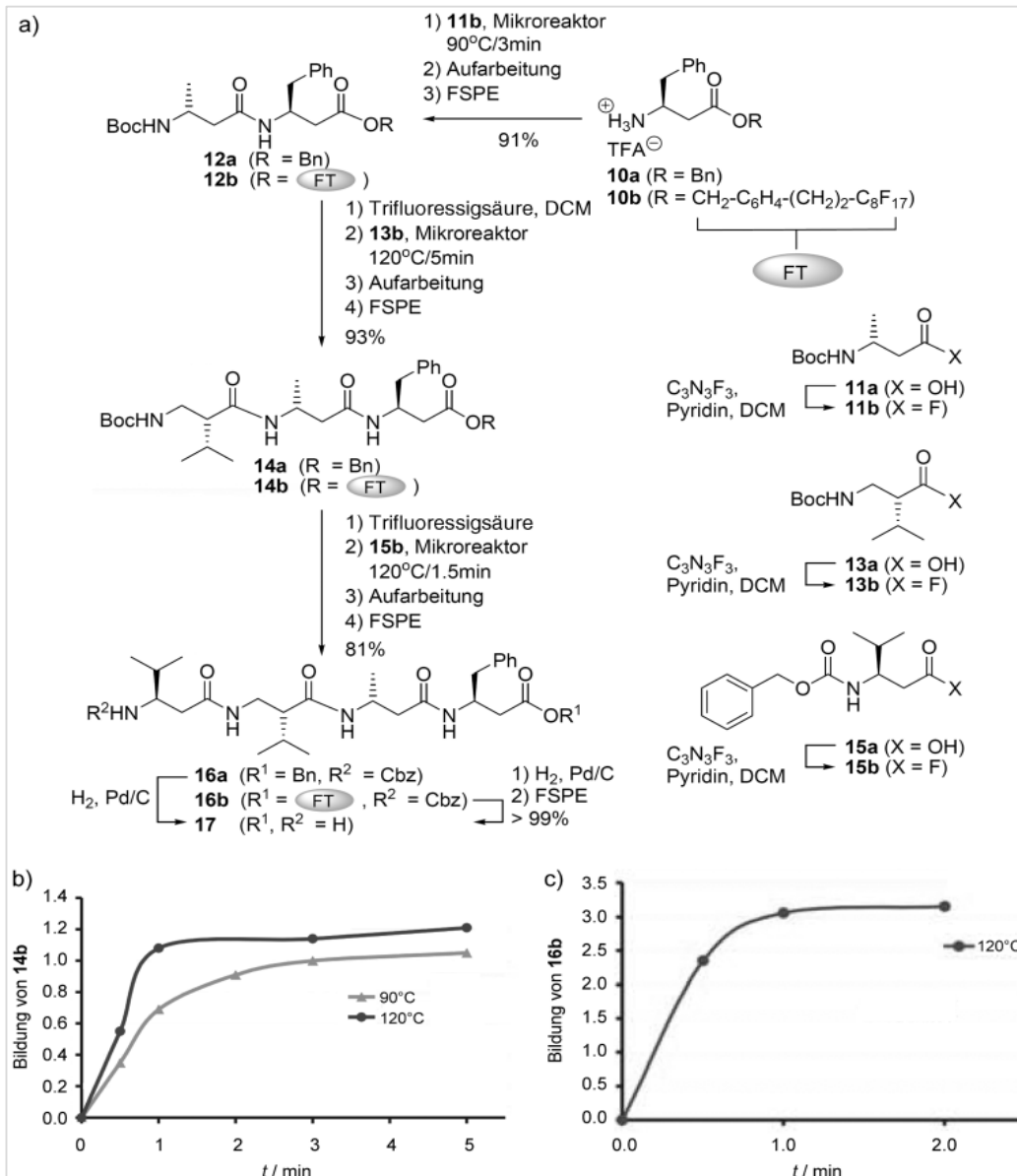5. Assembly of reconstructed chemical reaction

# The Challenges

The following images give an idea how our current way of communicating

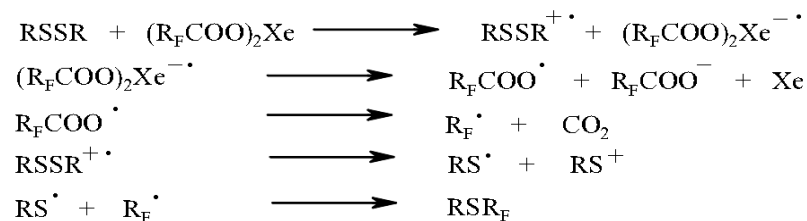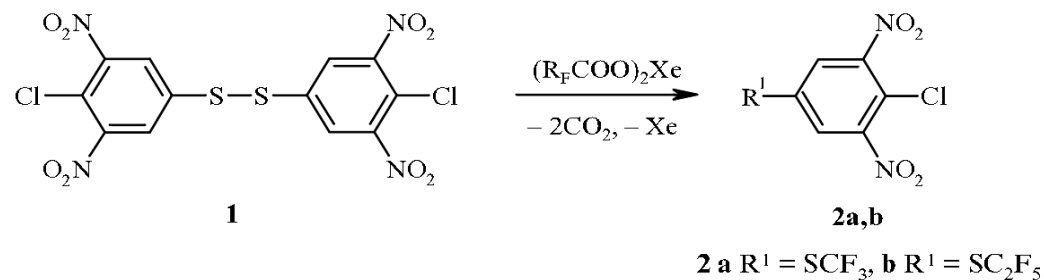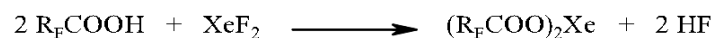chemical knowledge makes the life of computer scientists interesting

Fraunhofer Institute
Algorithms and
Scientific Computing

SCAI

# Semantik unklar

# Komplex Zusammengesetzt

# Komplex Zusammengesetzt

$$2\ R_FCOOH\ +\ XeF_2\ \longrightarrow\ (R_FCOO)_2Xe\ +\ 2\ HF$$



**1**  →  **2a,b**

**2 a** $R^1 = SCF_3$, **b** $R^1 = SC_2F_5$

$$RSSR\ +\ (R_FCOO)_2Xe\ \longrightarrow\ RSSR^{+\cdot}\ +\ (R_FCOO)_2Xe^{-\cdot}$$

$$(R_FCOO)_2Xe^{-\cdot}\ \longrightarrow\ R_FCOO^{\cdot}\ +\ R_FCOO^{-}\ +\ Xe$$

$$R_FCOO^{\cdot}\ \longrightarrow\ R_F^{\cdot}\ +\ CO_2$$

$$RSSR^{+\cdot}\ \longrightarrow\ RS^{\cdot}\ +\ RS^{+}$$

$$RS^{\cdot}\ +\ R_F^{\cdot}\ \longrightarrow\ RSR_F$$

R = 1-chloro-2,6-dinitrophenyl;  $R_F = C_nF_{2n+1}$ (when $n = 1$–2)

# Zeichnerisch schwierig
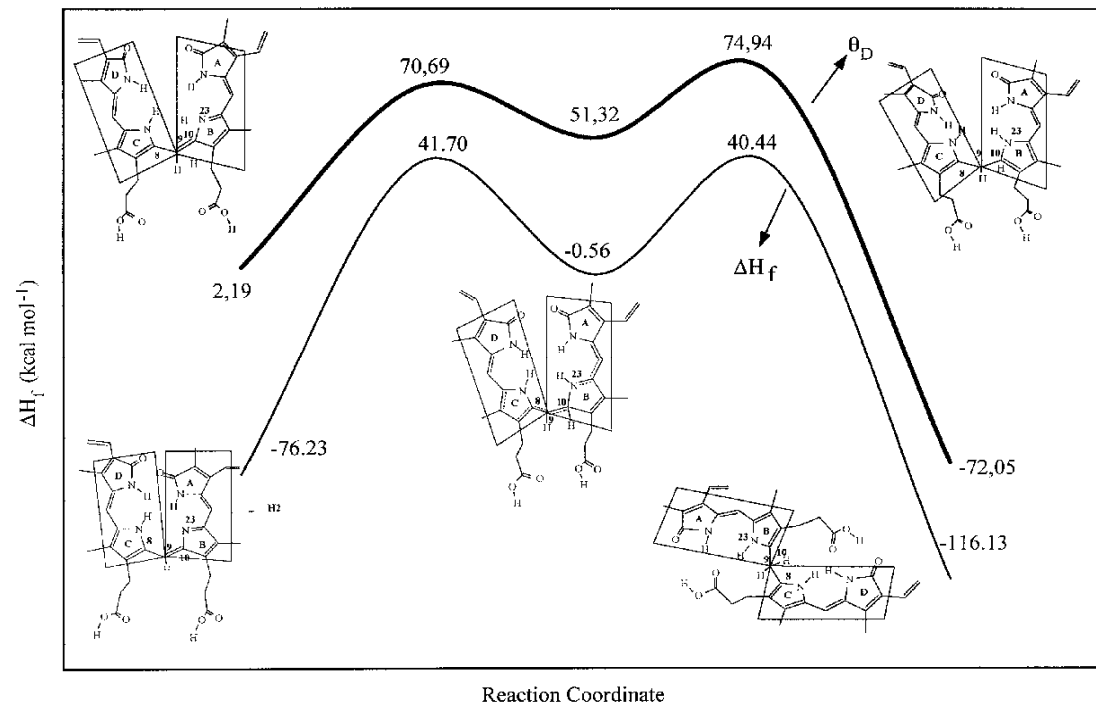


SCHEME 3. Dehydroabietic acid-based surfactants.

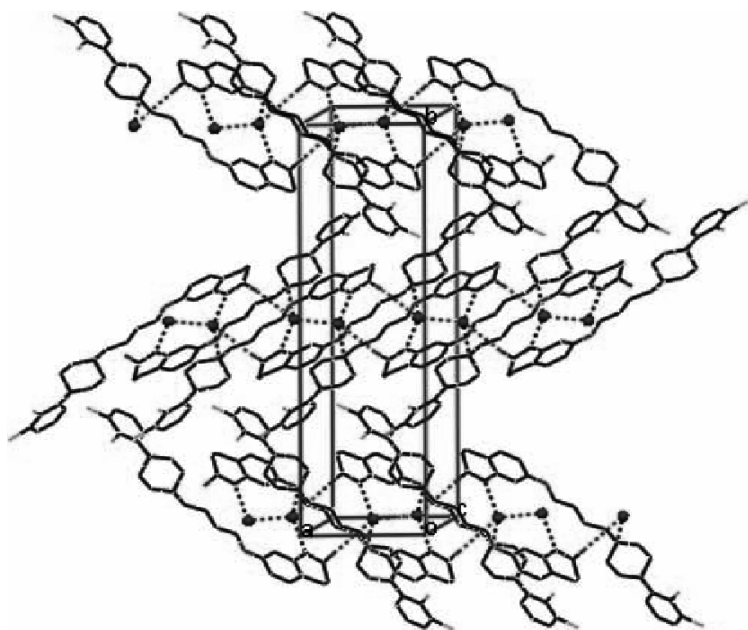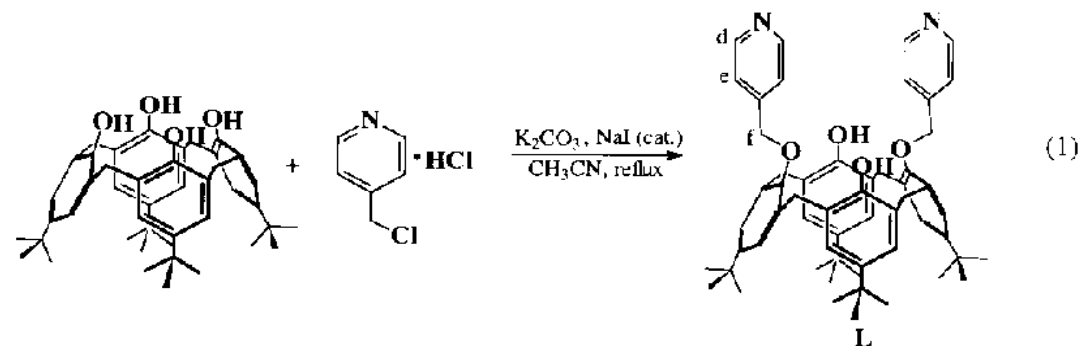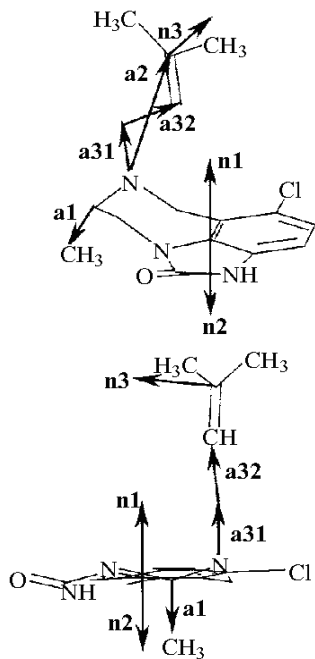# Zeichnerisch schwierig



Scheme 5 Biosynthesis of manzamine A proposed by Baldwin and Whitehead

# Völliger Wahnsinn



$$CH_3\text{-}N\text{-}CH_3 \quad \text{(molecular diagram with labels n3, a2, a32, a31, a1, n1, n2, Cl, NH, O)}$$



$$(1)$$



Reaction Coordinate

$\Delta H_f$ (kcal mol$^{-1}$)

Dihedral Angle $\theta_D$ (C8-CC9-C10-N23)

74,94

70,69

51,32

41.70

40.44

2,19

-0.56

$\theta_D$

$\Delta H_f$

-76.23

-72,05

-116.13

# Acknowledgement

**SCAI**

**Fraunhofer** Institute
Algorithms and
Scientific Computing