**Large Scale Evaluation of Chemical Structure Recognition**
**4th Text Mining Symposium in Life Sciences**
**October 10, 2006**

SCAI

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

Dr. Marc Zimmermann

# Overview

- Brief introduction Chemical Structure Recognition (chemOCR)

- Manual conversion of images

- Up scaling and automatisation

- Protocol database and parameter evaluation

- 2 methods of validation

- Test and benchmark data sets

- Examples, results and  lessons learned

# Chemical Structure Recognition – an Overview

**1** Document

**2** Depiction

**3** Reconstruction

**4** SDF file

**5** *in silico* Chemistry



created from
/home/marc/workspace/CSR/results/CSR/examples/US20051820
53/US2005182053_result.pnm
MZCSRv0.5010050621162D  0.00000   0.00000   0

26 28  0   1 0 0 0 0 0999 V2000
 204.0000 102.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
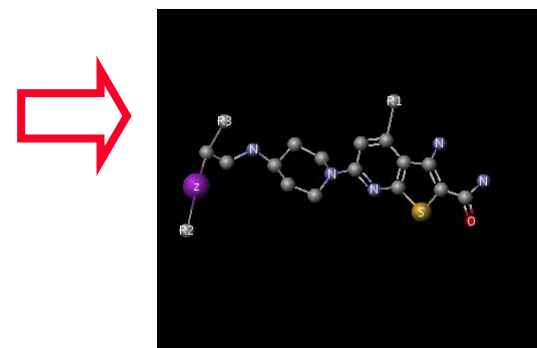 275.0000  61.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 201.0000  59.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 422.0000 178.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 311.0000 164.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 384.0000 165.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 447.0000 144.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 383.0000 123.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 131.0000  60.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 239.0000 123.0000   0.0000 C  0 0 0 0 0 0 0 0 0 0 0 0
 349.0000 218.0000   0.0000 R# 0 0 0 0 0 0 0 0 0 0 0
 447.0000 207.0000   0.0000 R# 0 0 0 0 0 0 0 0 0 0 0

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

SCAI

# The chemOCR Process

- is a multi step process:

    1. image preprocessing

    2. image conversion

    3. semantic entity recognition

    4. chemical structure assembly

    5. reconstruction validation

    6. post processing

- for each step a specific module has been implemented

- modules can be assembled into workflows

# Look And Feel Of CSR

Fraunhofer Institut
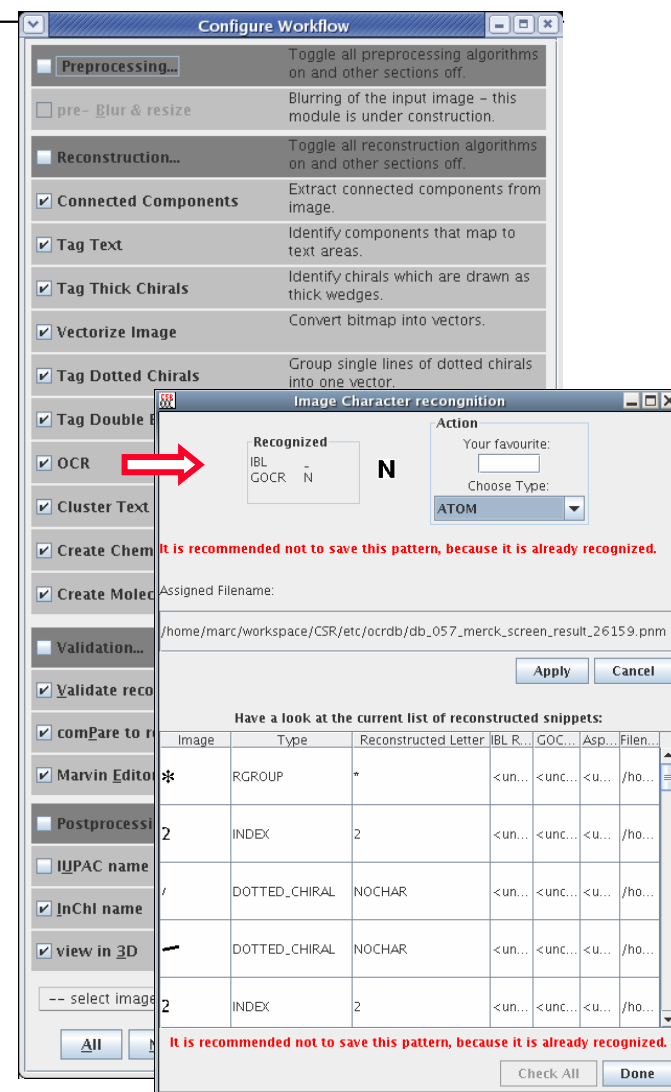Algorithmen und Wissen-
schaftliches Rechnen

# The Interactive User Mode

- a graphical user interface has been developed

- the user can trigger each module separately

- there are curators and editors to interfere with the process

the main advantages are:

- full control of the process

- easier than redrawing of the image

- teaching and learning capabilities of the system

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

# Adding New Modules – Using JAVA APIs and RPCs

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

# The Distributed Batch Mode: Scaling Up the Process

setting up the batch mode:

- a specific workflow is predefined

- a suitable parameter set is chosen

- each image becomes one job which is send to one computer

- all results are assembled

*advantages*:

- large speed up

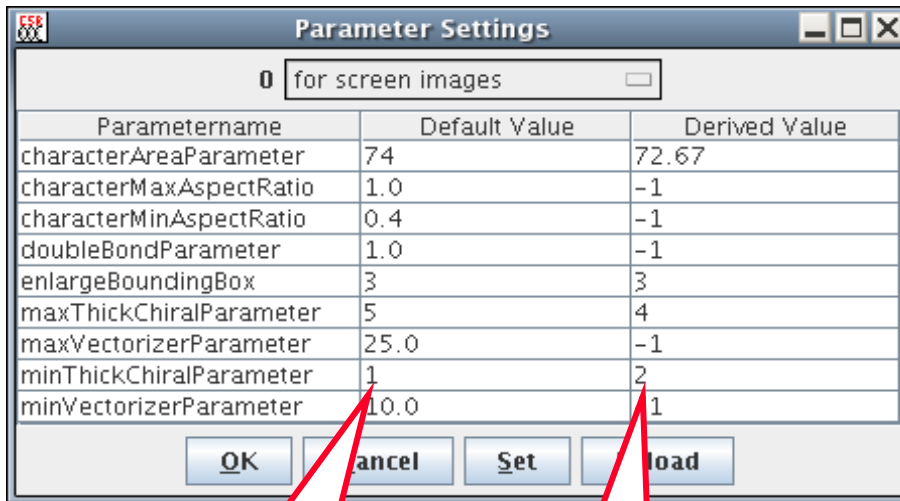- less human resources

- vast number of results

*disadvantages*:

- no control

- errors occur

- checking the results is time consuming

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

SCAI

# Many Images ⇨ Many Parameters?

**Parameter Settings**

| Parametername | Default Value | Derived Value |
|---|---|---|
| characterAreaParameter | 74 | 72.67 |
| characterMaxAspectRatio | 1.0 | -1 |
| characterMinAspectRatio | 0.4 | -1 |
| doubleBondParameter | 1.0 | -1 |
| enlargeBoundingBox | 3 | 3 |
| maxThickChiralParameter | 5 | 4 |
| maxVectorizerParameter | 25.0 | -1 |
| minThickChiralParameter | 1 | 2 |
| minVectorizerParameter | 10.0 | 1 |

**0** for screen images

OK    Cancel    Set    Reload

*predefined* for image sets – currently 4

*estimated* from the image itself

geometric constraints

*chiral*:
- # segments
- orientation

*bond*:
- min length
- max lenght

*character*:
- area size
- shape

File  Workflow  View  Help

Reconstruction... (1)   Image... (1)   Local Direct

Loaded Image (1)        Connec  mponents (1)

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

SCAI

# Technical Solution For Up Scaling

job description
sh-script

Sun Grid Engine

Cluster

license server

workstation

query / analysis

MySQL database

CSR

DB server

# Protocol Database for the Reconstruction Process
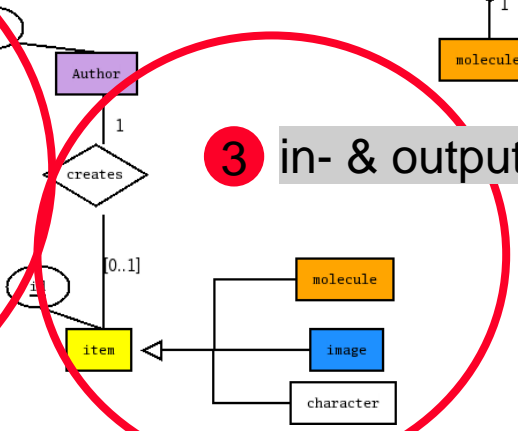


**1** session information

**2** workflow configuration

**3** in- & outputs

ER diagram

*query*: show me all molecules which have an atom error after changing the text_area parameter

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

SCAI

# Result Validation Using Training and Test Data



image corpora

CSR reconstruction

result validation

reconstruction          test molecules

result analysis

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

SCAI

# Validation Classes – A Closer Look



reconstruction / test molecule

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

# Reconstruction Error Prediction As an Alert System

*result validation* can only be used if the molecule is already known

or the expert is checking the result:

- good for bug fixing and training of the process
- can't be used for the data generation process

⇨ need a different strategy for the *batch mode*:

- identify and predict reconstruction errors
- alert the user only if interaction is needed
- choose a threshold for the precision

# Error Prediction – the Theory

*prediction* and *recognition* can be based on

- the use of chemical knowledge bases

- image properties, i.e. measure the complexity of the problem

- instance based machine learning, i.e. teach the system

the main goal is to assemble a *reconstruction score* without

knowing the correct solution

$$R_{score} = w_1 \cdot \text{complexity} + w_2 \cdot \text{chemical likelihood} + w_3 \cdot \text{known errors} < T_{alert} \ ?$$

weights *w* can be set by regression analysis

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

SCAI

# Established Error Classes

chemical knowledge bases

- OCR errors and unknown super atoms
- valence checking
- known scaffolds



image properties

- strange bond drawings (size, angles, ...)
- pixel density, size of connected components
- complexity



instance based machine learning (IBL)

- atom and bond distributions
- Lipinski score (i.e. drug like)

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

SCAI

# The Results – Current Status

test sets:

| | 2005 | 2006 | 2007 |
|---|---|---|---|
| 100x | 100x | 7600x | 750,000x |



measurement:
*correct*: input identical to output
*incorrect*: more than one error

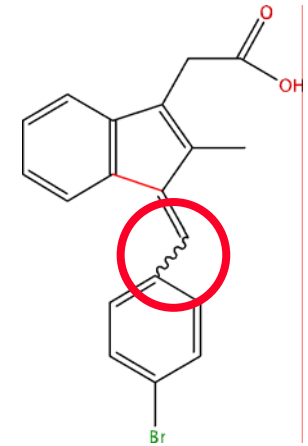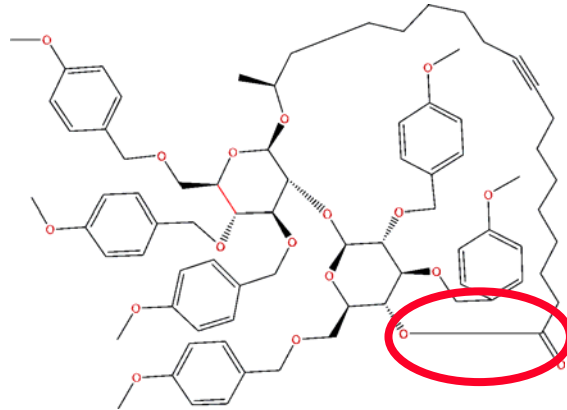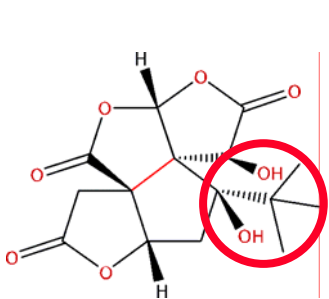| thin & clean > 95% | thick & clean > 75% | challenging mixture > 55% | patents > ??% |
|---|---|---|---|

# Not Too Bad...
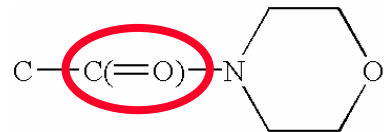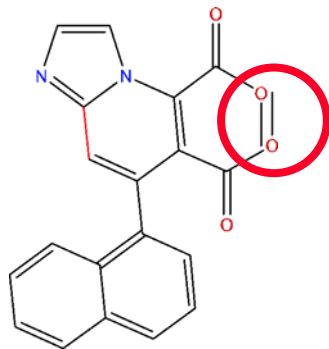


simple bridges
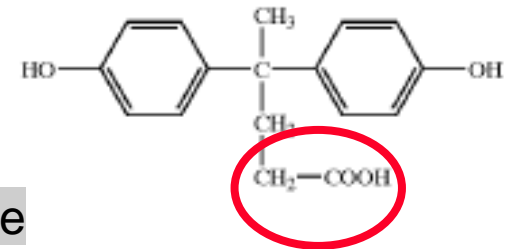
large molecules

salts and chirals

# Questionable…



extreme long bonds

new bond types

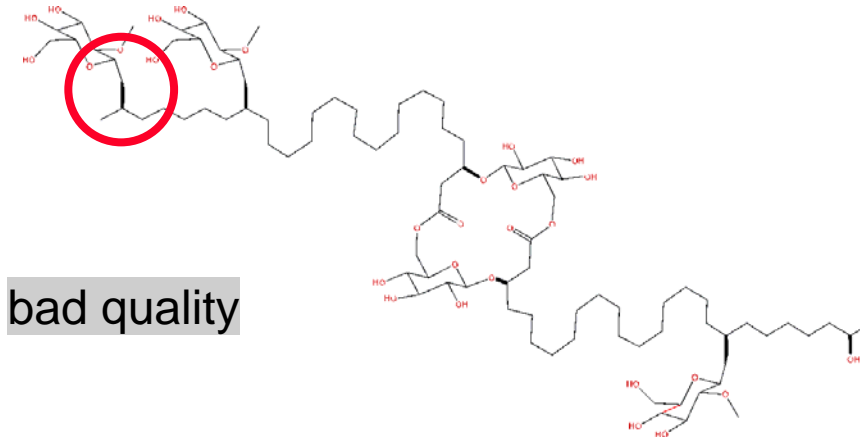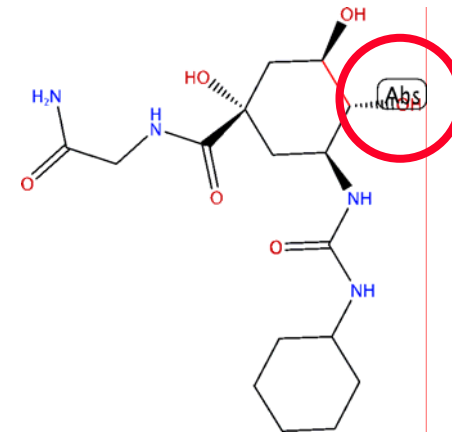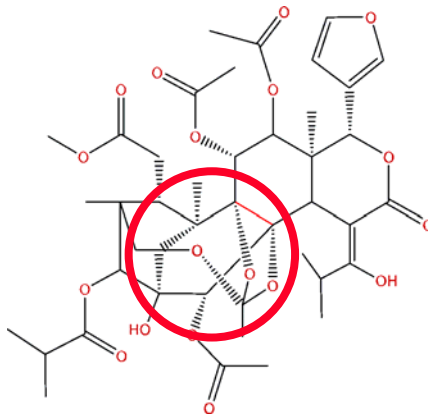close connections

SMILES and alike

Fraunhofer Institut
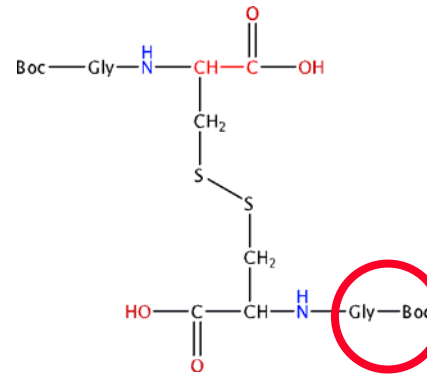Algorithmen und Wissen-
schaftliches Rechnen

SCAI

# Really Bad...



overlaps

bad quality
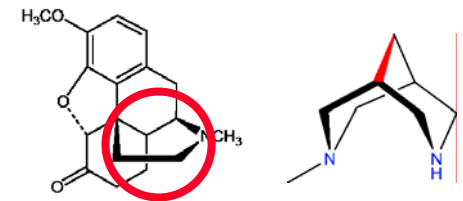
complicated bridged rings

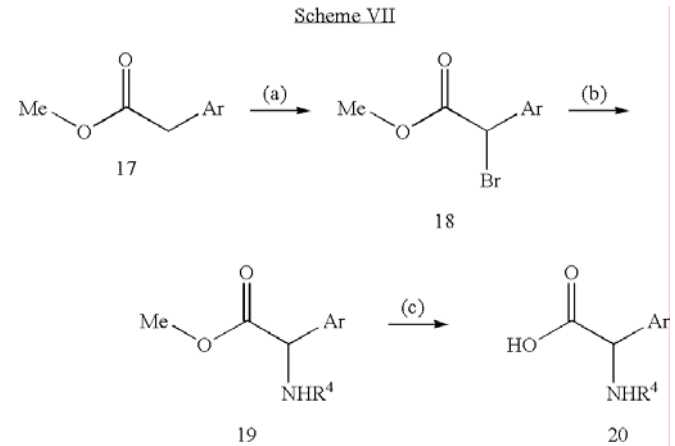amino acid orientation

thick bridges

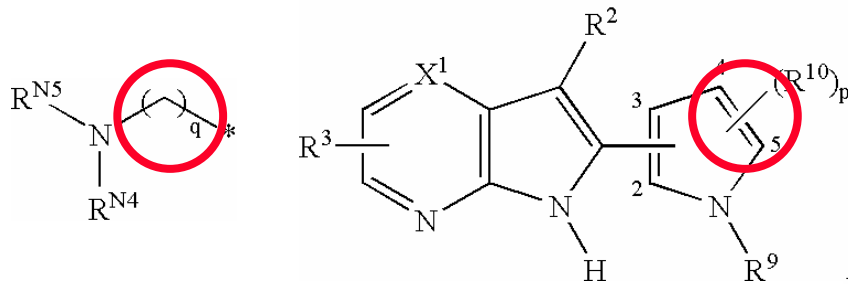Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen
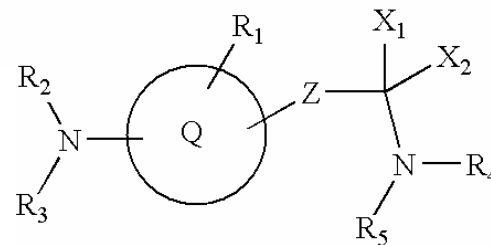
SCAI

# Patent Images...

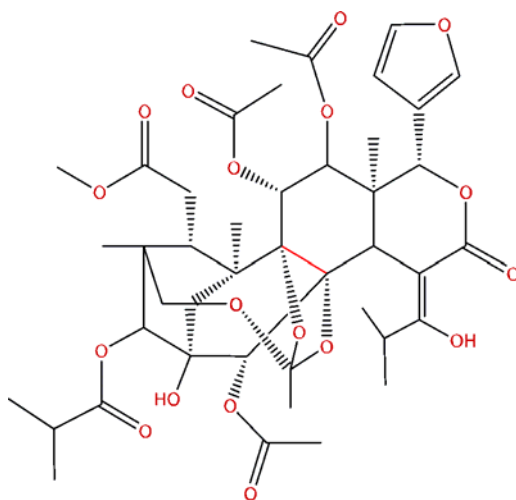R-groups and captions  [✓] SDF
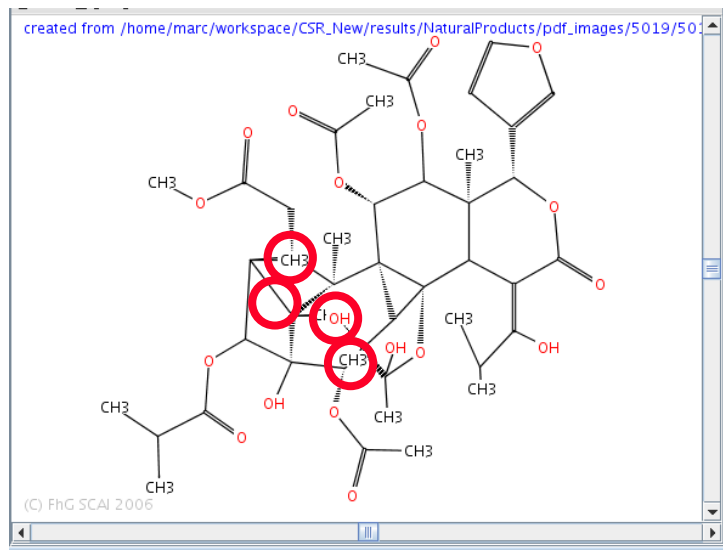
Scheme VII

reaction schema    [?] RDF

Markush structures    [×] file format

# So We Got an Error Reported

- need perfect reconstruction ⇨ start molecule editor

- need for indexing and retrieval

  ⇨ use similarity and substructure searches

  ⇨ specify reporting threshold



to be tolerated?

# A Glimpse at the Future: Multi Modal Extraction From Patents



**ProMiner (NER)**

**PatentMine by FhG**

**CSR**

**PatentSpace**

**graphical mining**

proteins
protein families
protein complex
compound
process
drug class
disease
pathways

# Lessons Learned

- a generic chemOCR framework has been established

- there is and there will not be a "one-fits-all" solution

- CSR can be adapted and optimized (parameters, error models, image preprocessing, …)

- although we have looked into many examples, we have not seen so far all sorts of image sources (e.g. legacy of old documents)

- we will continuously improve our methods as new challenges come along

  - *You* can get hands on experience on CSR in an *evaluation project*

  - *SCAI* provides: training, installation support, bug fixing, fitting CSR to the data, *long term research agenda*

Fraunhofer Institut
Algorithmen und Wissen-
schaftliches Rechnen

SCAI
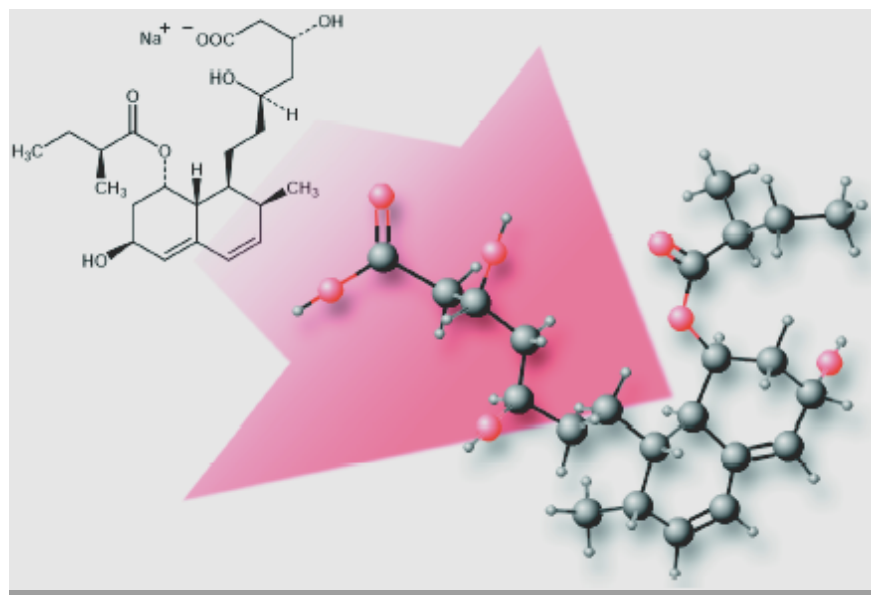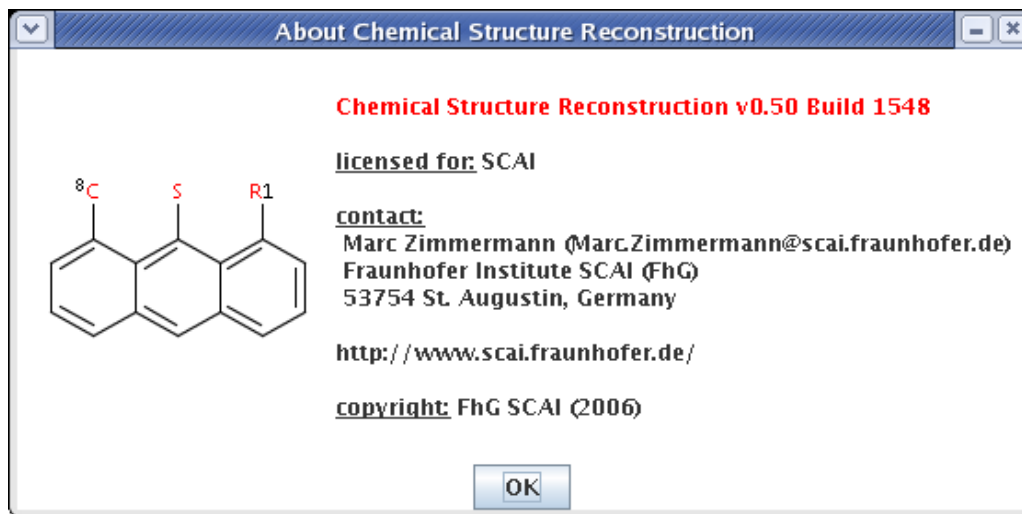
# The Team (in the order of appearance)

- Marc Zimmermann*

- Tanja Fey

- Le Thuy Bui Thi

- Christoph Friedrich*

- Yuan Wang

- Maria-Elena Algorri*

- Miguel Alvarez

- Angelika Weihermüller*

- Wei Wang

- Peter Kral*

- Carina Haupt*



*) currently improving CSR

# **CSR** Online Demo Available During The Break



CSR can extract chemical depictions from various image sources and convert them into SMILES and SD files, which can be further used in nearly all chemical software; it allows for the modification of reconstructed molecules by a structure editor; it maintains the superatom and bond (single, double, triple, or chiral) information; and it accepts user curation in each stage and scoring schema to improve its performance.