# chemoCR
## chemical compound reconstruction

# From Bitmap Images to Connection Tables

**chemoCR™**, a system for the reconstruction of chemical information from chemical structure depictions. This tool has been developed to make one of the largest sources of chemical information accessible: the information communicated through images containing chemical structure depictions. Although communication of chemical information through images is quite common amongst chemists, the information contained in images could not be used by machines. chemoCR™ solves this problem by reconstructing chemical information from images.

**SCAI**

**Fraunhofer** Institute Algorithms and Scientific Computing

## chemoCR Workflow:

I Check Configuration
II Load Input
III Start Reconstruction Workflow
IV Post-process Molecule

**I**

Known organic scaffolds

Known super atoms and aliases

Trained characters of the OCR

Loaded Expert Rules

**Machine learning:**
- instance based learning
- knowledge bases
- chemical expert rules

Check the validity of the licenses for a compute cluster

## Input and Output

**II**

All bitmap formats (GIF, PNG, BMP, ...)

Connection tables (MOL, SDF)

Select 1 Parameter Set for 1 Image Source

Parameters for the Extractors

## Reconstruction

**III**

Connected Components

Character Extraction

Bitmap Textures

Chemical Expert System

Choose Rule

Thick Chirals

OCR

Vectorization

Chemical Graph

Create Molecule

Color Coding

Select the appropriate Modules into a Workflow

## Validation

**IV**

Validate and search for possible conversion errors

Fix Errors using an structure editor

Format Conversions into 3D, SMILES, InChI, ...

Search the WWW

## Technology

In order to solve the problem of recognizing and translating chemical structures in image documents, our chemoCR™ system combines pattern recognition techniques with supervised machine-learning concepts. The method is based on the idea of identifying from depictions the most significant semantic entities (e.g. chiral bonds, super atoms, reaction arrows…). The workflow consists of three phases: image pre-processing, semantic entity recognition, and molecule reconstruction plus validation of the result. All steps of the reconstruction process make use of chemical knowledge. The following validation process is for detection and fixing of possible reconstruction errors. The system can be adapted to different sets of input images.

## Technical Specification

The chemoCR™ core functionality is based on platform independent JAVA libraries. It has been tested on UNIX™ operating systems (Fedora Linux, Sun Solaris) and on Windows XP. External tools can be easily included in the workflow. Our software can be used interactively by a graphical user interface or it can be run distributed in batch processing mode. Our benchmark test set consisting of 8000 images of natural products is processed in about 2 hours.

For contact and further details: http://www.scai.fraunhofer.de/chemocr.html