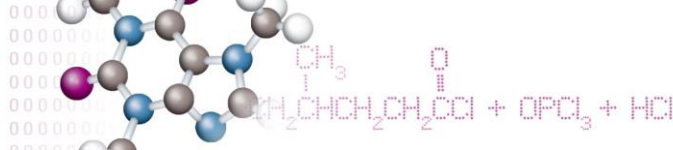# Extraction of structural information from ChemDraw CDX files: easy, or an underestimated, difficult challenge?

**Josef Eiblmaier, Hans Kraut, Sascha Hausberg, Peter Loew**

*ICIC 2013 Vienna, October 13 – 16*

# Outline

» **ChemDraw files:**

  **Relevance and the Challenge**

» **Approach**

» **Projects**

    » **InfoChem *ChemProspector***

    » **Wiley *Smart Article***

    » **Thieme *Science of Synthesis Update / Pharmaceutical Substances***
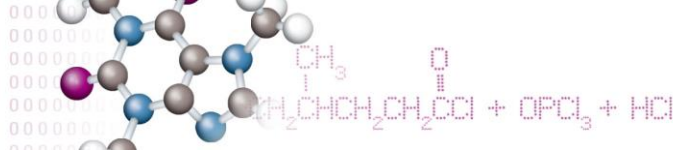
» **Conclusion / Outlook**



**© cora / PIXELIO, www.pixelio.de**

# Patents, Journal Articles and MRW's: a Buried Treasure?

Chemical structures (images)

Chemical names/fragments (text)

Markush structures (text, images, CDX)

Chemical structures (CDX files)

Reactions (CDX files)

Dr. Josef Eiblmaier

# Manuscript → Article → Database …

Publishing

Manuscript submission

Manual Indexing

Database production e.g. SciFinder, Reaxys, SPRESI eEROS, ...

# CDX Scheme vs. Database Record

| ChemDraw file | Database |
|---|---|
| Purpose: <u>presentation</u> / publishing <u>no</u> search | Purpose: search / retrieval |
| Unstructured | Structured |
| Structures: no strict rules | Structures: strict rules |
| General rules: none | Database rules: strict |



| Reactant | Product | Reagent | Solvent | Catalyst |
|---|---|---|---|---|
| | | $SOCl_2$ | | |
| | | LiOH | $H_2O$, THF | $Pd(OAc)_2$ |
| | | Cl-Co2Et, $Et_3N$ | Acetone, $H_2O$ | |

Source: Thieme *Pharmaceutical Substances*, Ticagrelor (in production)

*ICIC 2013 Vienna, October 13 – 16*

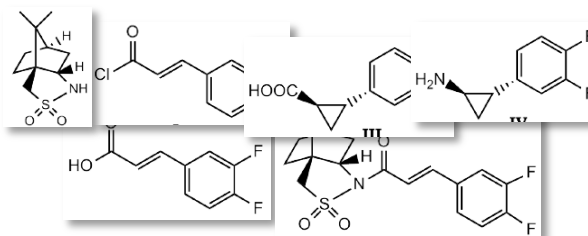# CDX Scheme Processing, what does that mean?
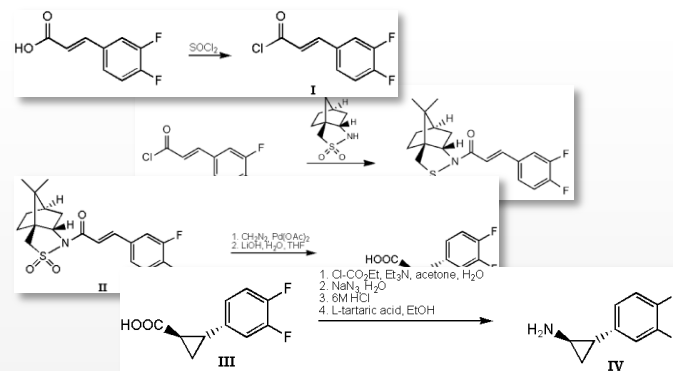


**ICSchemeProcessor**

Source: Thieme *Pharmaceutical Substances*, Ticagrelor (in production)
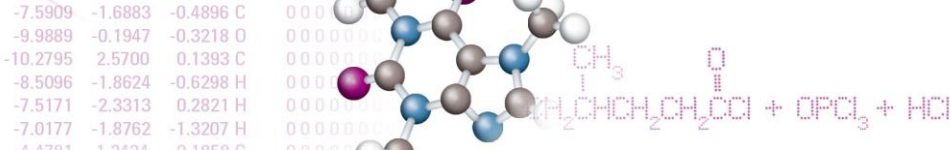
## Chemical structures (SD files)



## Reactions (RD files)



## Conditions (RD files)

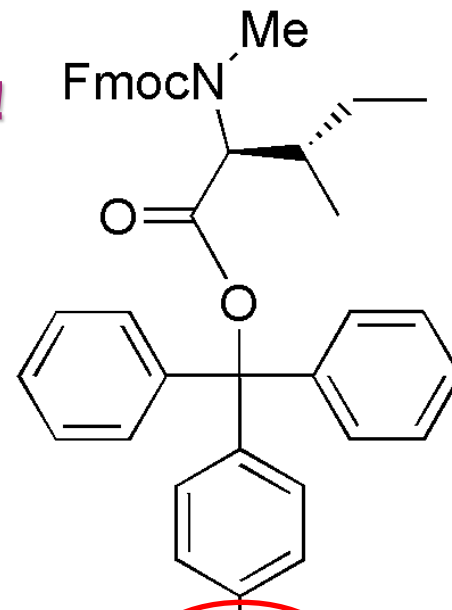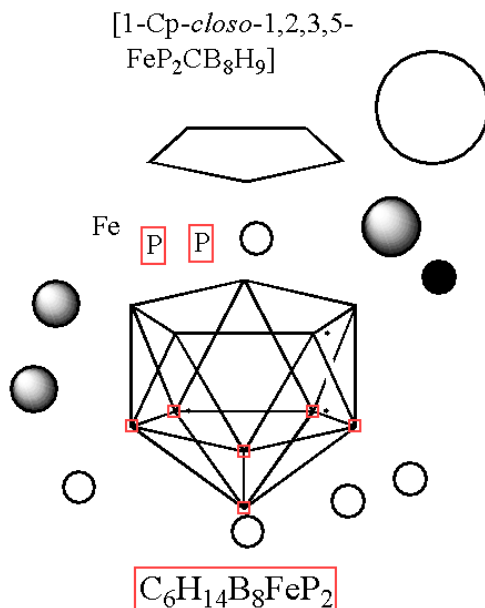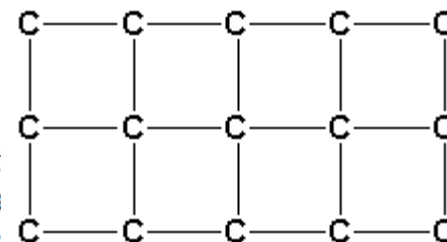| Reagent | Solvent | Catalyst |
|---|---|---|
| SOCl$_2$ | | |
| LiOH | H$_2$O, THF | Pd(OAc)$_2$ |
| Cl-Co2Et, Et$_3$N | Acetone, H$_2$O | |

# But: CDX files, often an optical illusion!
# Authors are very inventive for a ,perfect' layout!
# Appearences are deceiving!

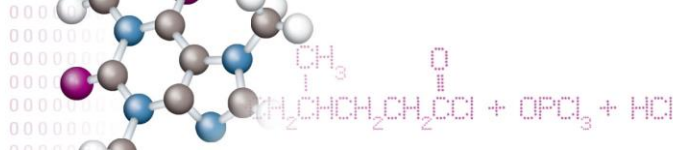» **Usage of graphical symbols**

  • Polymer supports

  • Heteroatoms

$[1\text{-}Cp\text{-}closo\text{-}1,2,3,5\text{-}FeP_2CB_8H_9]$

Fe

$C_6H_{14}B_8FeP_2$

Me

FmocN

C Grid:

I chloride SynPhase
lly available trityl
ol/lantern) with Ac
lethylamine (DIEA)/CH₂Cl₂ (Scheme 2).[12]
of **13** from the polymer-supported **19** with
Cl₂, the loading amount of **13** was deter-
33 mmol per lantern. After removal of the
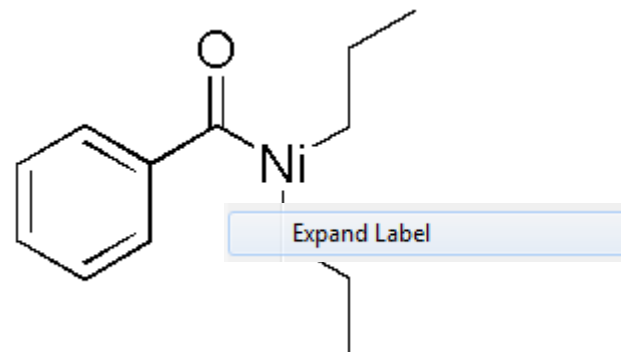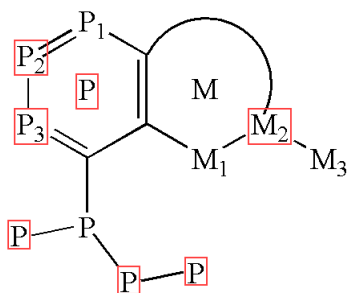ith 20% piperidine/DME, acylation with

# Optical illusions 2

» **Unresolvable labels**
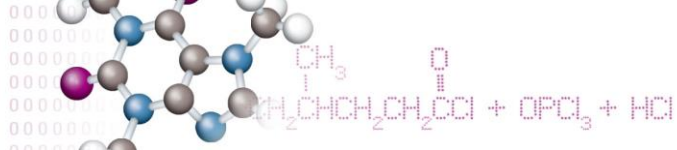
- Labels not defined

- Element symbols used as R-group labels
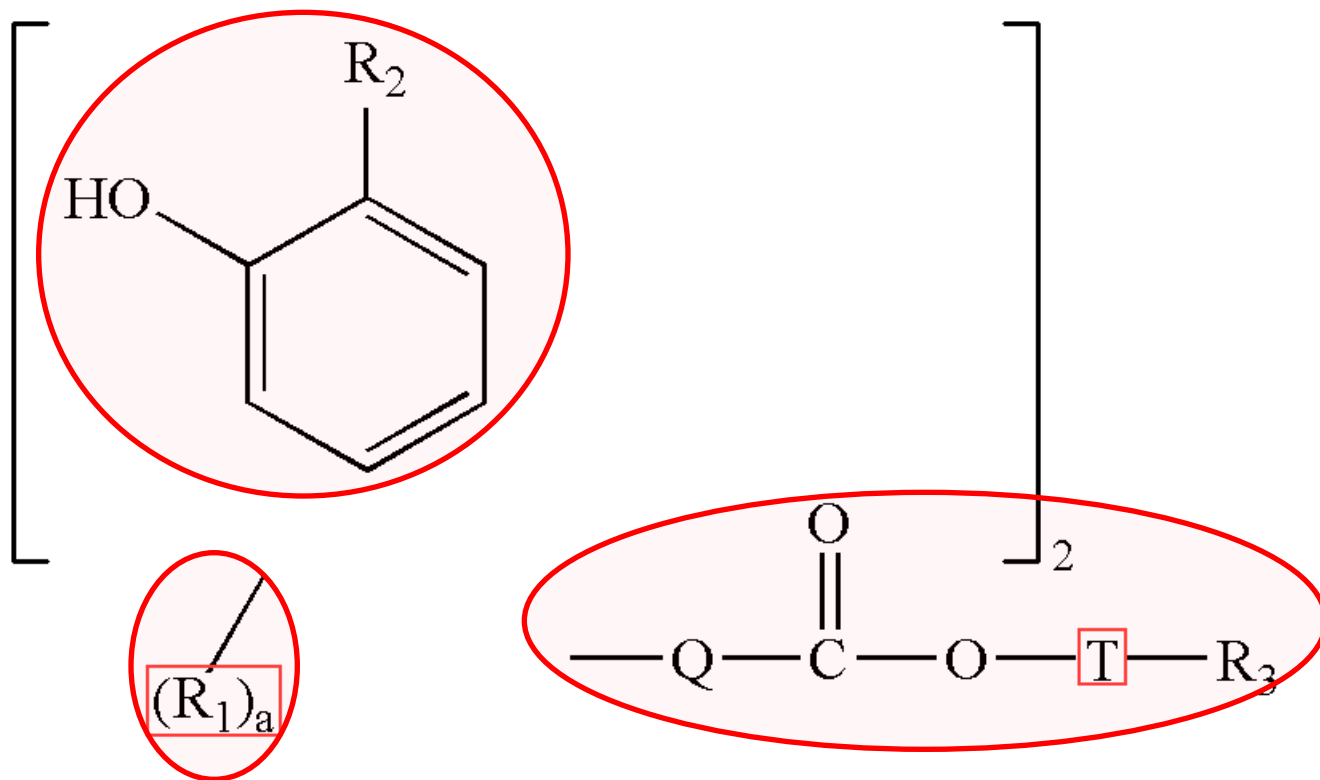
- Ambiguous fragment labels (e.g. molecular formula)

# Optical illusions 3

» **Variable points of attachment**

# Optical illusions 4

» **Reaction arrows / forked arrows / brackets**

# Approach



© Gerd Altmann / PIXELIO, www.pixelio.de

# Approach

» **The algorithmic approach:**

- Application of a set of rules in the software (generic, project unspecific). Software should recognize all cases that might occur!

- project (title-) specific rules (drawing conventions must not change), otherwise further development necessary

- manual post correction required (cost/time intensive)

- problem is infinite, unprecedented issues can not be handled
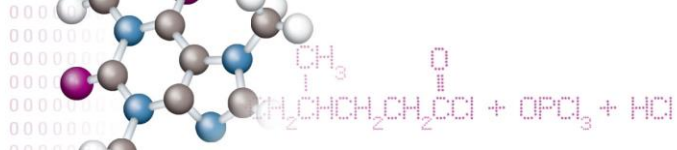
» **The templating approach:**

- software is developed to recognize a defined set of problems (PS)

- all content must be manually pre-templated (cost intensive) according to the capabilities of the software

» **The hybrid approach:**

- depending on the source the focus can be laid on either approach

# Templating

» *Templating*: **Guidelines for authors and typesetters**

- Syntax definitions for tables, R-groups etc.

- Syntax rules for captions

- Reaction arrangement, forked arrows

- Rules for reaction conditions

  (reactants, catalysts, solvents, yields, temperature)

# Examples:

» **Algorithmic detection of features**

  » Resolution of repeating groups

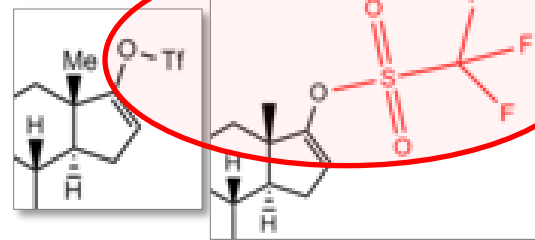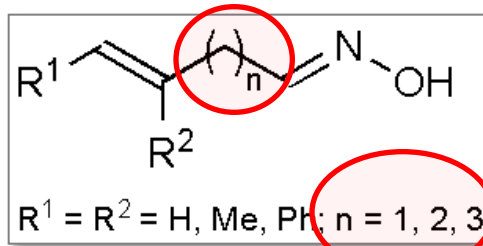  » Enumeration of R-groups

  » Resolution of aliases/labels

  • source specific alias databases

  • continuously extended

  » Table Enumeration
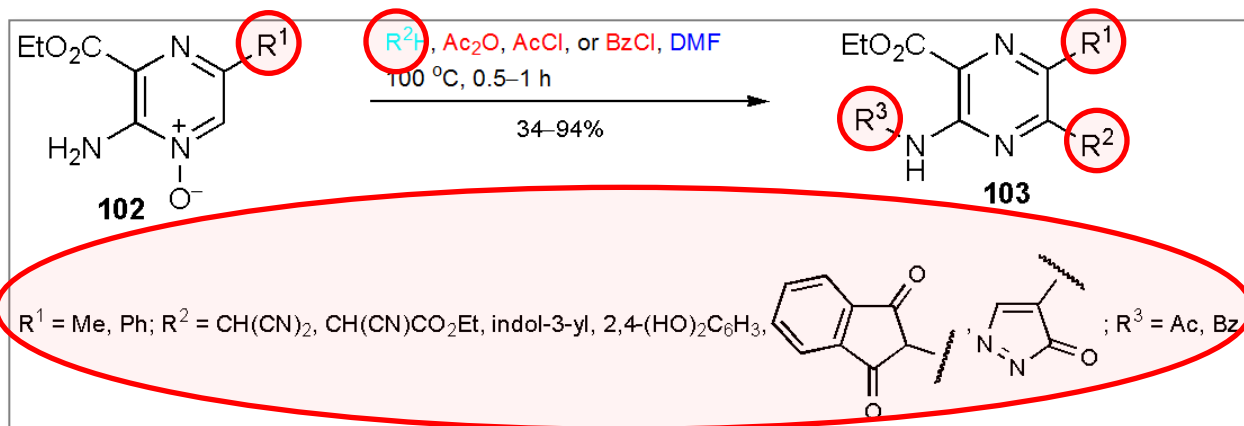
  • compound enumeration

  • reaction factual data: Caption/Yield

  » Variable points of attachment

  » Forked arrows
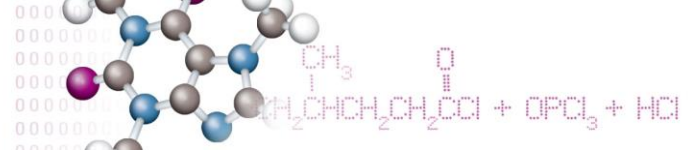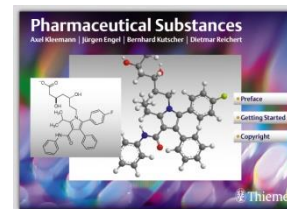


| Caption | X | R |
|---|---|---|
| 3i | H | Ph |
| 3j | (pyridin-2-yl) | Ph |
| 3k | (isoquinolin-3-yl) | Ph |

$R^1 = R^2 = H$, Me, Ph; $n = 1, 2, 3$

$R^1 = $ Me, Ph; $R^2 = CH(CN)_2$, $CH(CN)CO_2Et$, indol-3-yl, $2,4-(HO)_2C_6H_3$, ... ; $R^3 = $ Ac, Bz
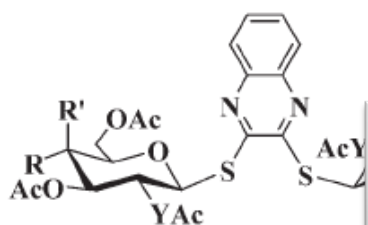
**Projects**

# Sucessful Application of CDX Processing: Chemistry Enrichment Workflow*, (Wiley *Smart Article*)



*Reinhard Neudert: *Enhancing the User Experience for Wiley Chemistry Content, ICIC 2012 14. – 17. October, Berlin*
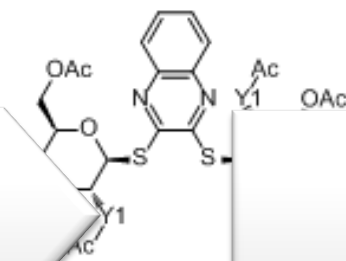
# Templating*

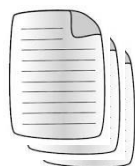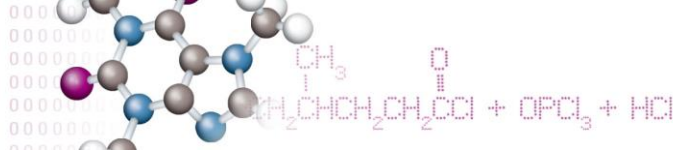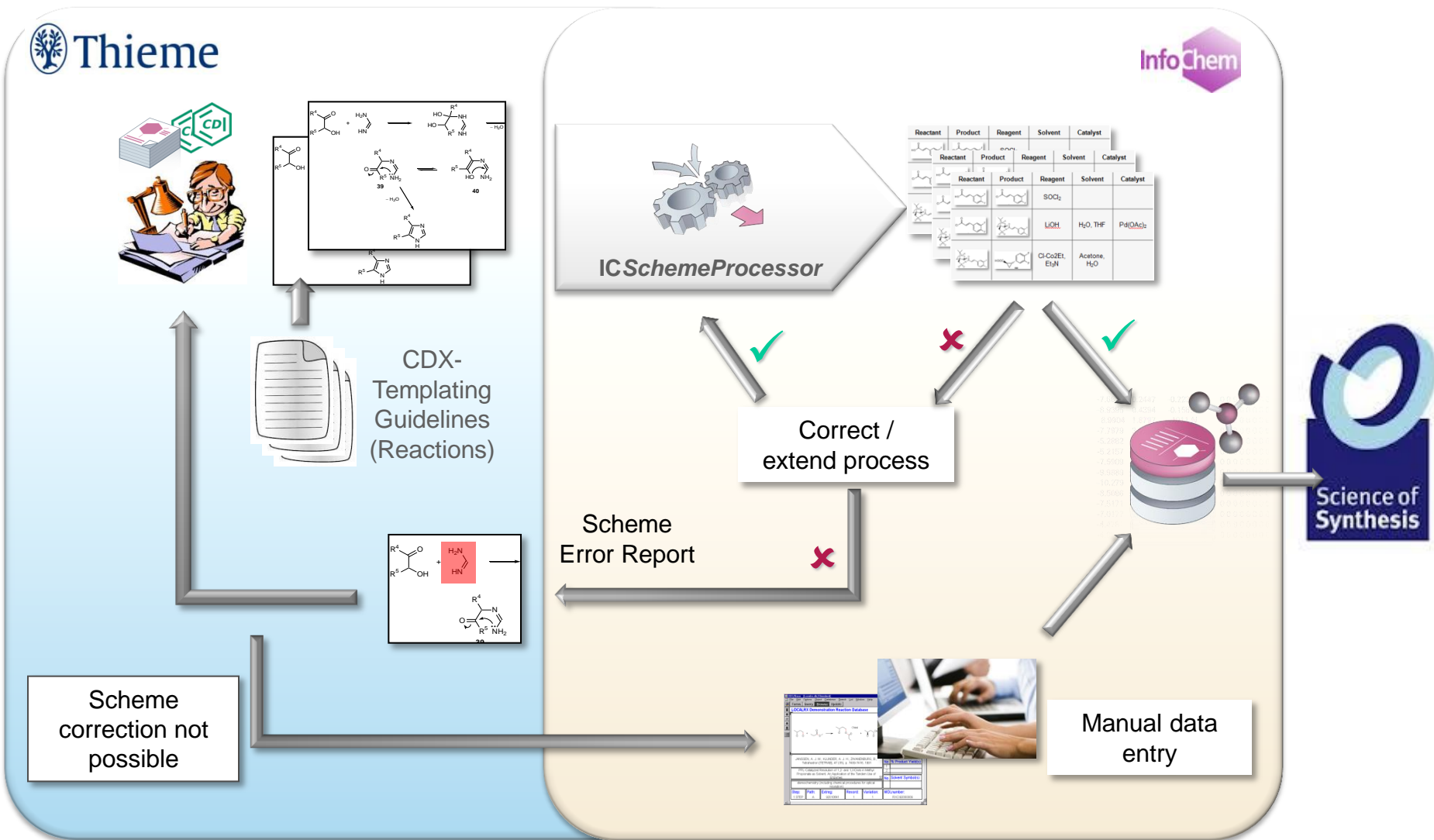| Author's CDX File | CDX Template | Enumerated structures |
|---|---|---|



Templating

ICSchemeProcessor

CDX-Templating
Guidelines (Structures)

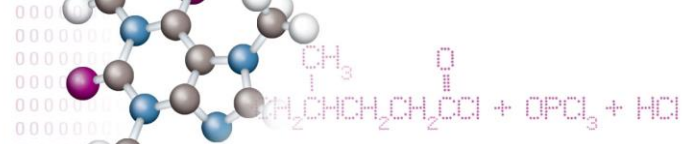*Reinhard Neudert: *Enhancing the User Experience for Wiley Chemistry Content, ICIC 2012 14. – 17. October, Berlin*

# Workflow *Science of Synthesis Update*

*ICIC 2013 Vienna, October 13 – 16*

# Sample *Pharmaceutical Substances Update*



**Source: Thieme *Pharmaceutical Substances*, Abiraterone**

# Conclusion

» **As much as possible algorithmic processing desirable**

- generic: can be applied to other contents as well

- cheaper (humans cost!)

» **100% conversion (without human interaction) never possible**

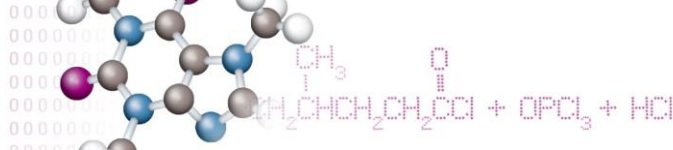» **Solutions are project / source specific**

» **Relevance of automatic extraction will continuously increase**

» **Authors / Publishers play an essential role in a successful conversion**

# Acknowledgements

» **Wiley**

- Michael Forster
- Reinhard Neudert

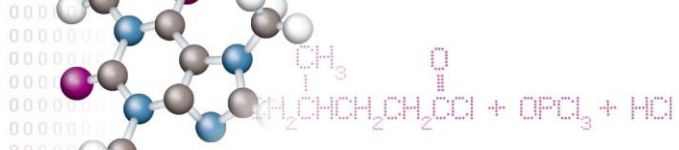» **Thieme**

- Guido Herrmann
- Rolf Hoppe
- Klaus Köberlein

» **InfoChem**

- Hans Kraut, Sascha Hausberg, Thomas Menke, Manuela Rauh
  Fanny Irlinger, Huyen Ngyen, Dagmar Kunzmann

© Thomas Link / Flickr

# Thank you!

*ICIC 2013 Vienna, October 13 – 16*

Dr. Josef Eiblmaier

# Questions?

*ICIC 2013 Vienna, October 13 – 16*

Dr. Josef Eiblmaier