

# Chemical information presentation in the Crystallography Open Database

Andrius Merkys<sup>a</sup>, Agnė Matusevičiūtė<sup>b</sup>, Antanas Vaitkus<sup>b,c</sup>, Armel Le Bail<sup>d</sup>, Daniel Chateigner<sup>e</sup>, Luca Lutterotti<sup>f</sup>, Miguel Quirós-Olozabal<sup>g</sup>, Mykolas Okulič-Kazarinas<sup>a</sup>, Peter Moeck<sup>h</sup>, Peter Murray-Rust<sup>i</sup>, Nicholas E. Day<sup>j</sup>, Robert T. Downs<sup>j</sup>, Saulė Girdzijauskaitė<sup>c</sup> and Saulius Gražulis<sup>a,b</sup>

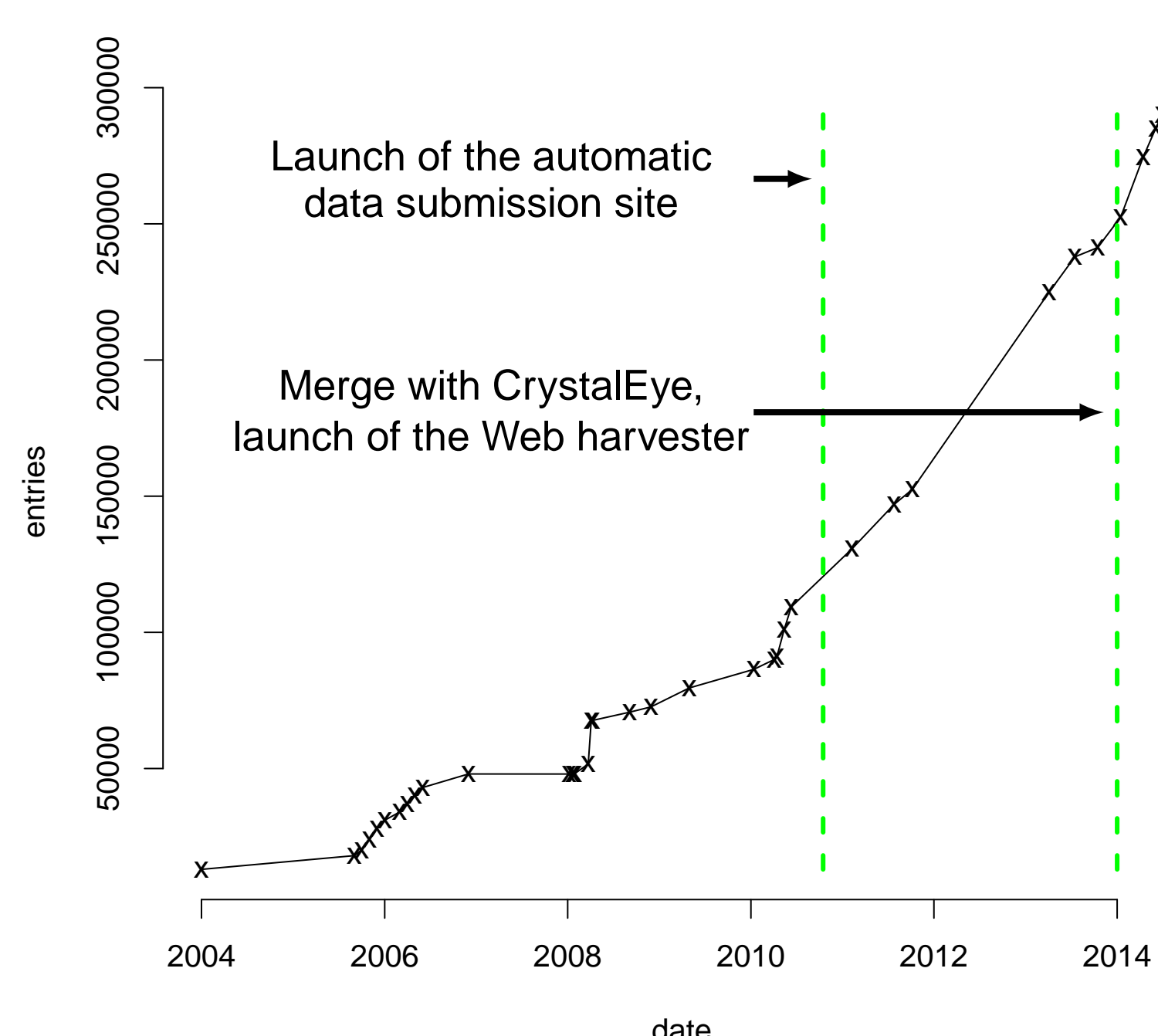
<sup>a</sup>Vilnius University, Institute of Biotechnology, Department of Protein-DNA Interactions, Vilnius, Lithuania; <sup>b</sup>Vilnius University, Faculty of Mathematics and Informatics, Department of Mathematical Computer Science, Vilnius, Lithuania; <sup>c</sup>Vilnius University, Faculty of Mathematics and Informatics, Department of Software Engineering, Vilnius, Lithuania; <sup>d</sup>Universite du Maine, Laboratoire des Oxydes et Fluorures, Université du Maine, Faculté des Sciences, Le Mans, France; <sup>e</sup>Universite de Caen-Basse Normandie, CRISMAT-ENSICAEN, Caen, F-14050 Caen, France; <sup>f</sup>University of Trento, Department of Materials Engineering, Trento, Italy; <sup>g</sup>Universidad de Granada, Facultad de Ciencias, Departamento de Química Inorganica, Granada, Spain; <sup>h</sup>Portland State University, Department of Physics, Portland, USA; <sup>i</sup>University of Cambridge, Department of Chemistry, Cambridge, United Kingdom; <sup>j</sup>University of Arizona, Department of Geosciences, Tucson, USA

## Abstract

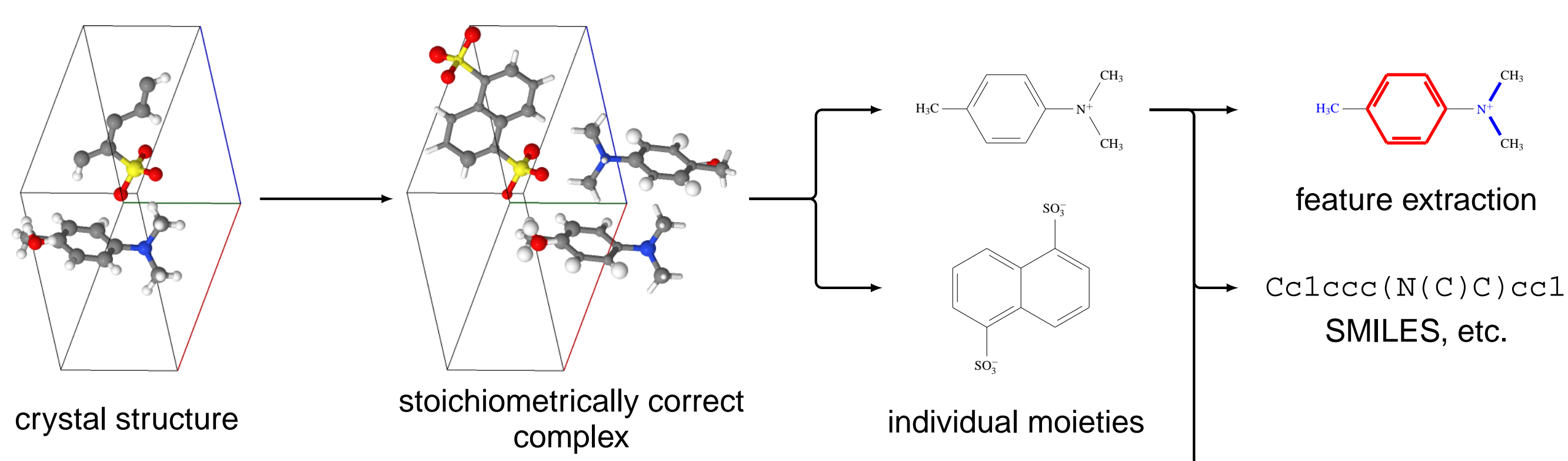
Crystallography Open Database (COD, <http://www.crystallography.net>) is the largest to date curated open-access collection of small to medium sized unit cell crystal structures [4, 3]. Over 11 years of development, COD has accumulated over 1/4 million structures. COD has an automated data submission Web site, performs routine automatic quality checks on all incoming structures and is now recommended as a database for crystallographic deposition by several scientific journals. To facilitate automatic use and discoverability of COD data, and to increase usefulness of our database for chemists, two steps were undertaken. COD was supplemented with software and data from the CrystalEye data aggregator [2]. The new software permits extracting chemical data and presenting them as structural formula, unique moieties and chemically significant fragments. We have also implemented search of crystal structures by the structural chemical formulae of the target compounds. To facilitate data curation, a new software platform for data review is being developed to automatically detect unusual structures and collect expert opinions from qualified human reviewers concerning credibility of such structures.

## Growth of the COD

- ▶ Data sources:
  - ▶ donators (IUCr, AMCS and others);
  - ▶ Web harvesters of open journals;
  - ▶ depositions via automatic data submission site, including personal communications.
- ▶ Journals recommending COD for data deposition:
  - ▶ Inorganic Chemistry;
  - ▶ Mineralogical Magazine;
  - ▶ Nature Data.

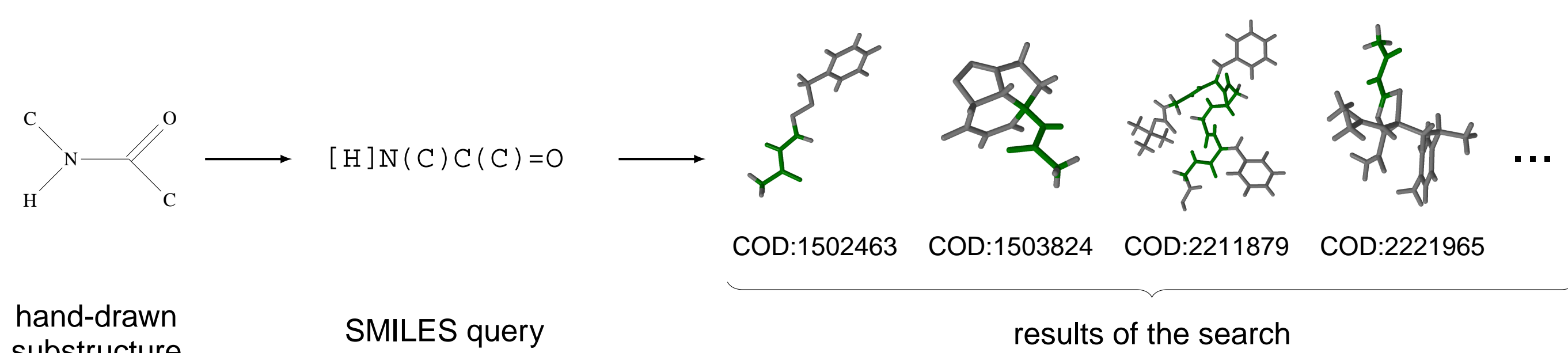


## Extraction of the chemical information



- ▶ Fully automatic and unsupervised pipeline is devised;
- ▶ Software from CrystalEye is employed:
  - ▶ heuristics for calculation of partial charges;
  - ▶ heuristics for determination of bond orders;
  - ▶ algorithm to isolate individual moieties;
  - ▶ algorithms to extract ring (*red*) and chain (*blue*) nuclei.
- ▶ Input and output use common file formats (CIF, CML and SDF).

## Search by substructure formulae

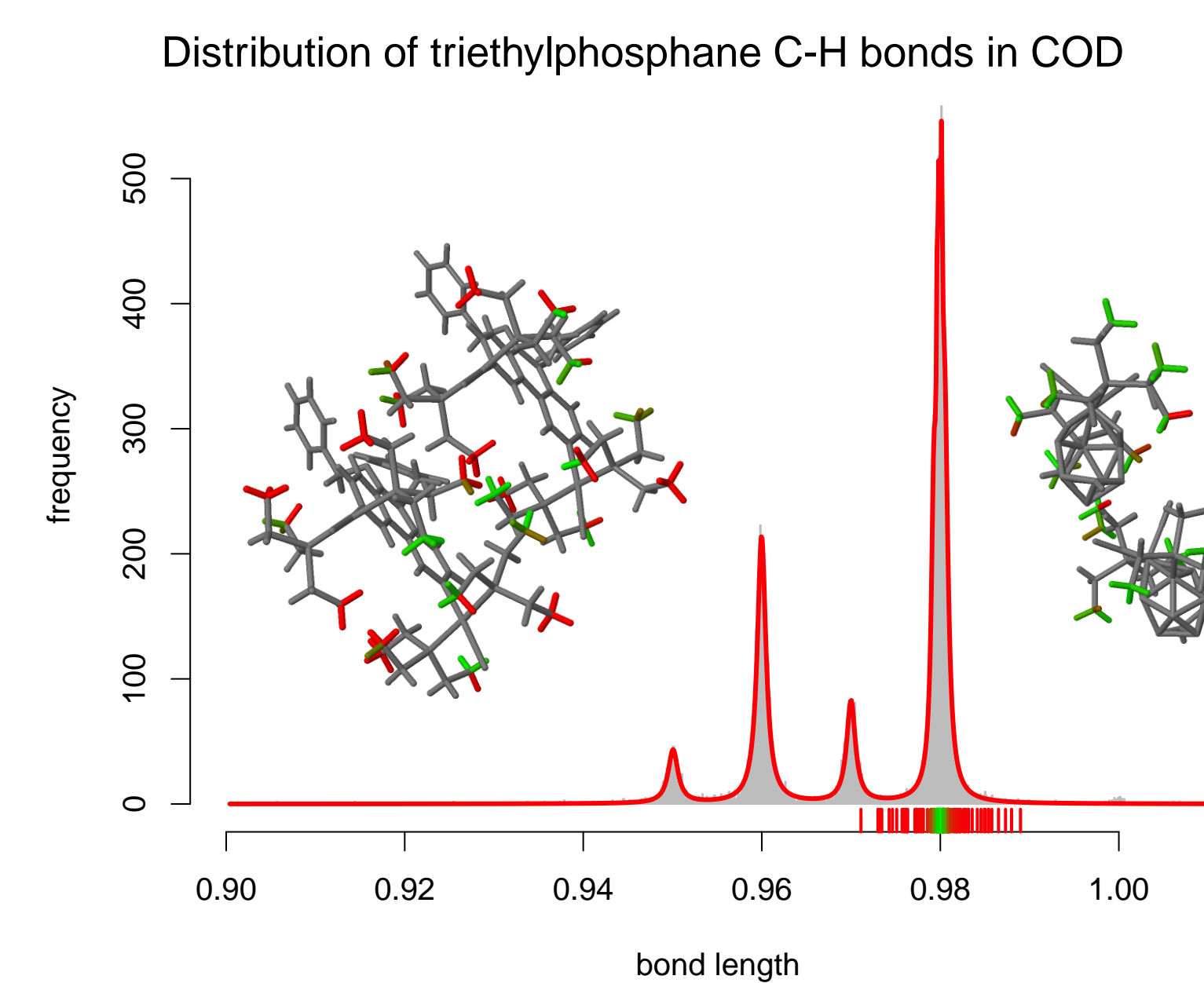


- ▶ Queries can be submitted by drawing substructures with Web browser applet or entering SMILES [1] manually;
- ▶ Currently the search is performed on a set of 70 000 hand-curated SMILES descriptors and can be extended to automatically generated descriptors.

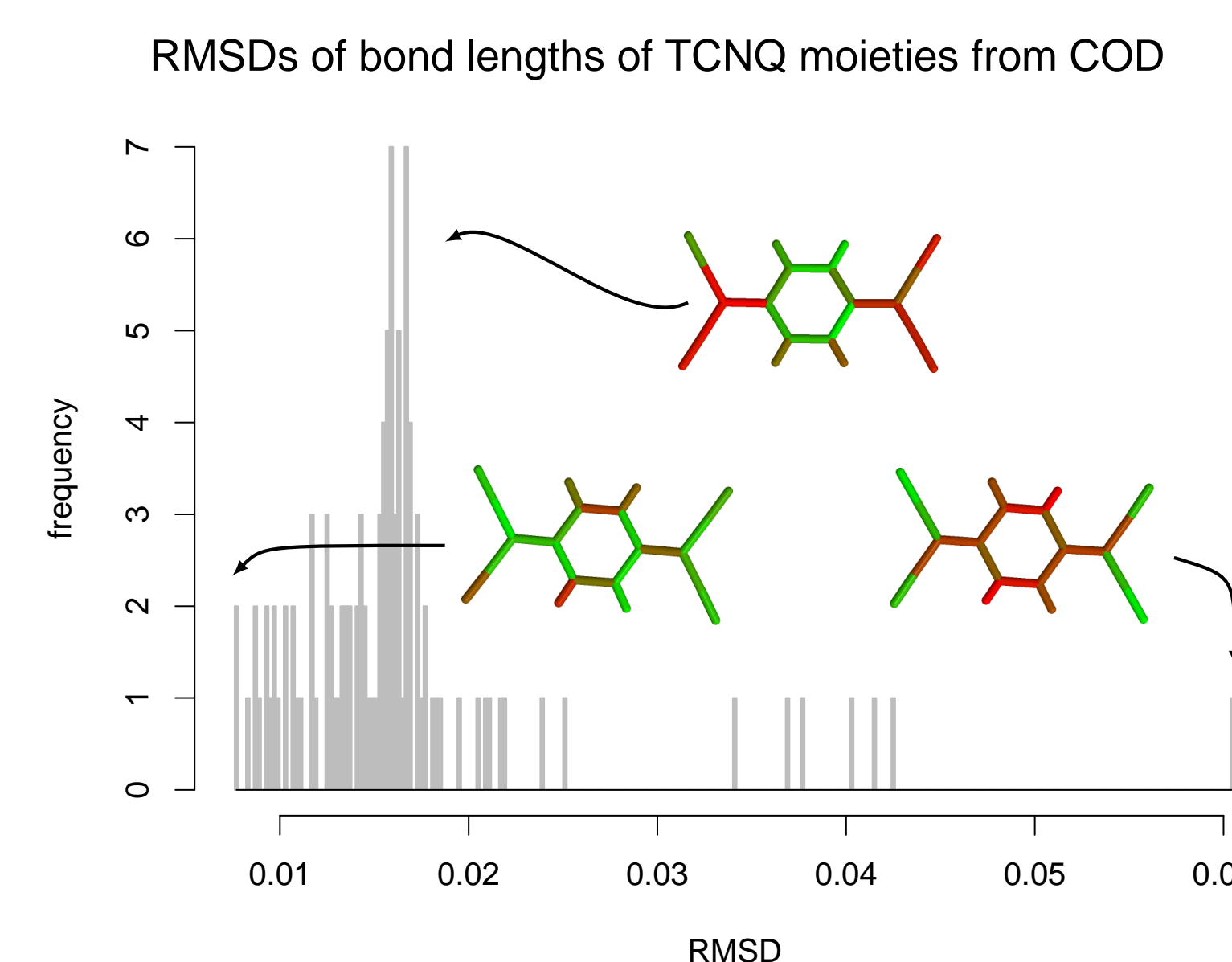
## Evaluating the geometry

- ▶ Bond lengths, valence and dihedral angle sizes are compared to the statistical distributions of the bulk data;
- ▶ Statistical models to describe the distributions are generated and updated automatically;
- ▶ Models depend on the overall structural data quality.

- ▶ COD:4027109 (*left*)
  - ▶ Most of the C-H bond lengths deviate from database's average;
  - ▶ Hydrogens are attached as riding atoms.
- ▶ COD:1101162 (*right*)
  - ▶ Most of the C-H bond lengths are close to the mode;
  - ▶ Hydrogens are also attached as riding atoms.

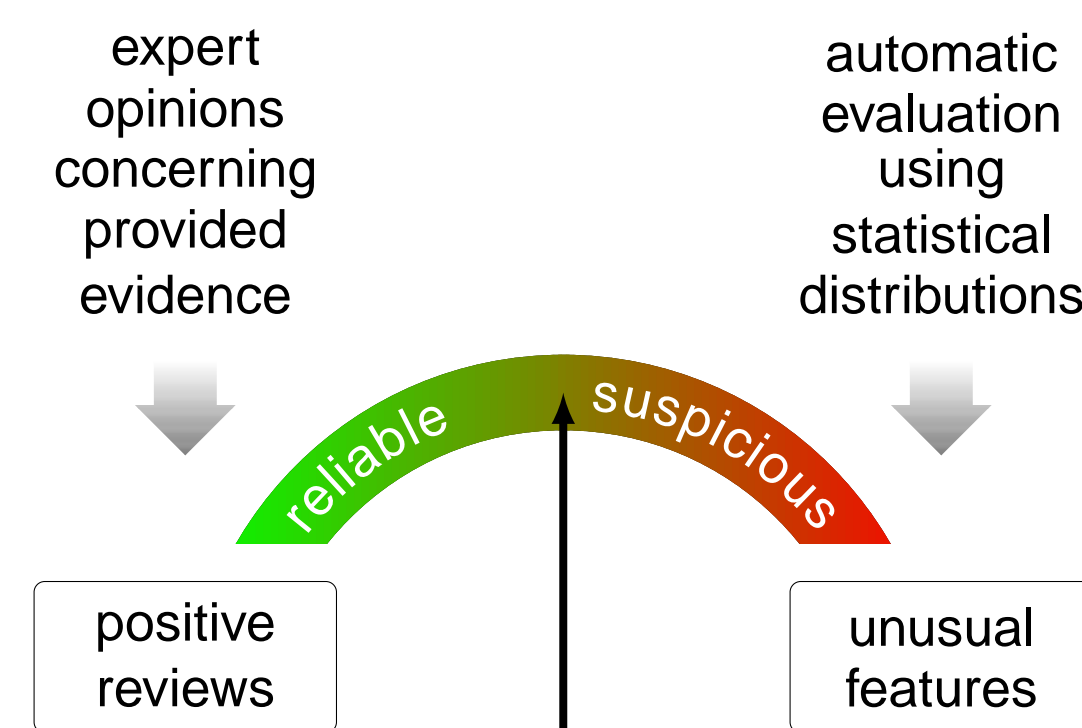


- ▶ A criterion to evaluate whole moieties can be derived
  - ▶ Allows selection of moieties with the "most usual" geometry;
  - ▶ Allows fast detection of outliers.



## Platform for data reviews

- ▶ "Unusual" does not necessarily mean "wrong"
  - ▶ Extraordinary claims require extraordinary evidence;
  - ▶ The most statistically unusual structures will be forwarded to a COD reviewer Web forum for verification;
  - ▶ Convincing evidence confirms validity of unusual structures.
- ▶ The set of usual and verified unusual structures should be used for reliable scientific inferences;
- ▶ Unusual structures require special attention.



## Cross-linking the COD with Open Data

- ▶ RDF (*Resource Description Framework*) descriptions are provided for each database entry
  - ▶ example: <http://www.crystallography.net/1516168.rdf>
- ▶ Cross-links are made with ChemSpider, PubChem, AMCS and MPOD;
- ▶ Links to Wikipedia and DrugBank are provided.

## Acknowledgements

This research was funded by a grant (No. MIP-025/2013) from the Research Council of Lithuania.

## Bibliography

- [1] Anderson et al. SMILES: A line notation and computerized interpreter for chemical structures. Technical report, Environmental Research Laboratory-Duluth, 1987.
- [2] Day. Automated Analysis and Validation of Open Chemical Data. PhD thesis, University of Cambridge, nov 2008.
- [3] Gražulis et al. Crystallography open database – an open-access collection of crystal structures. *Journal of Applied Crystallography*, 42(4):726–729, Aug 2009.
- [4] Gražulis et al. Crystallography Open Database (COD): an open-access collection of crystal structures and platform for world-wide collaboration. *Nucleic Acids Research*, 40(D1):D420–D427, Jan 2012.

Created with Ubuntu → MySQL → R Sweave → L<sup>A</sup>T<sub>E</sub>X beamerposter