Computer-Assisted Structure Elucidation

Christoph Steinbeck

In a typical metabolome measurement, less than 40% of the features can be assigned to known compounds.

> Oliver Fiehn, UC Davis, USA Internal communication

Setting the scene

(i.e., what are we actually talking about?)



I have a molecular formula from HR MS.

My database has a hit (or two, or three).

Problem solved.

Or is it?

There are **known knowns**; there are things we know we know.

We also know there are **known unknowns**; that is to say, we know there are some things we do not know.



But there are also **unknown unknowns** – the ones we don't know we don't know.

-United States Secretary of Defense,

Donald Rumsfeld



Levels of Confidence

Table 1

Summary of levels of confidence in metabolite 'identification', as defined by the Chemical Analysis Working Group of the Metabolomics Standards Initiative [4]

Level of confidence	Description	Requisite analytical data
Level 1	'Identified metabolites'	Two orthogonal analytical techniques applied to the analysis of both the metabolite of interest and to a chemical reference standard of suspected structural equivalence, with all analyses performed under identical analytical conditions within the same laboratory Examples of appropriate orthogonal data: accurate mass via MS with retention time; accurate mass MS and fragmentation data or isotopic pattern; 2D NMR spectra; full ¹ H and/or ¹³ C NMR spectra and so on
Level 2	'Putatively annotated compounds'	As for levels 3 and 4, including spectral (NMR and/or MS) similarity with public or commercial libraries
Level 3	'Putatively characterised compound classes'	As for level 4, plus spectral and/or physicochemical properties consistent with a particular class of organic compounds
Level 4	'Unknown'	A discernible spectral signal (NMR, MS or other) that can be reproducibly detected and quantified

Reproduced from: Viant MR, Kurland IJ, Jones MR, Dunn WB (2017) How close are we to complete annotation of metabolomes? Current Opinion in Chemical Biology 36:64–69. doi: 10.1016/j.cbpa.2017.01.001

The More Complicated Cases



So to be clear:

For everything that follows, we need a **pure**, **isolated** compound in **sufficient amounts**.

Isomers

Isomers

share one molecular formula







 $\square \land \frown \emptyset \land \Box \emptyset \bowtie \Box \land \Box \emptyset \land \frown \Box 0 \land \Box 0 \land$ The 217 constitutional isomers of C₆H₆ **E**

 $\nabla \oslash \Box \nabla \bowtie \Box \oslash \lor \Box \lor \Box \oslash \Box \oslash$

The 217 constitutional isomers of C₆H₆

S ✓<>>>
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓
✓</p

Q



 $\succ \bigtriangledown \checkmark \checkmark \checkmark \checkmark$ $\mathbb{K}_{_} \quad \stackrel{}{\longrightarrow} \quad \quad \stackrel{}{\longrightarrow} \quad \quad \quad \rightarrow} \quad \quad \quad \quad \rightarrow} \quad \quad \quad \quad$ $\checkmark \Leftrightarrow \mathrel{\triangleleft} \checkmark \mathrel{\neg} \backsim \checkmark \checkmark \checkmark$

The 217 constitutional isomers of C₆H₆



Q



 \bigcirc A \Diamond Ĵ Q \square \succ V \mathbb{N} \rightarrow ~ I

The 217 constitutional isomers of C₆H₆

 $C_{10}H_{16}$ $C_{13}H_{16}O_3$ $C_{30}H_{48}O_2$

24938 constitutional isomers

 $C_{13}H_{16}O_{3}$

 $C_{10}H_{16}$

 $C_{30}H_{48}O_2$

 $C_{10}H_{16}$

 $C_{13}H_{16}O_{3}$

24938 constitutional isomers

> 2,000,000,000 constitutional isomers

 $C_{30}H_{48}O_2$

 $C_{10}H_{16}$ $C_{13}H_{16}O_3$ $C_{30}H_{48}O_2$

24938 constitutional isomers

> 2,000,000,000 constitutional isomers

>> 10¹² constitutional isomers

Assume you had a computational tool to make all isomers.

Structure Generators

are **computer programs** to generate **all constitutional isomers** which adhere to a given set of input constraints, such as a molecular formula Constitutional Isomers of $C_{10}H_{16}$



Computer-Assisted Structure Elucidation with **2D** NMR

The only viable way for slightly more complex problems upwards.



Cryomagnet Design



NMR Principle

Odd number of protons and/or neutrons -> intrinsic magnetic moment and angular momentum (nonzero spin)

Even numbers of both -> spin of zero.



NMR Principle

Odd number of protons and/or neutrons -> intrinsic magnetic moment and angular momentum (nonzero spin)

Even numbers of both -> spin of zero.





NMR Principle

Odd number of protons and/or neutrons -> intrinsic magnetic moment and angular momentum (nonzero spin)

Even numbers of both -> spin of zero.



1D Proton NMR

1D Proton NMR

1D Proton NMR

Exploded Pharmacy: HR-MS

yields information about elemental composition, such as C₁₀H₁₆

Experiments: J-Couplings: DEPT

Acronyms:

DEPT: Distortionless Enhancement by Polarization Transfer

APT: Attached Proton Test

• ¹³C-detekted 1D-Exp.

 Number of protons attached to each carbon is coded as signal phase

•Combining information from DEPT-135, DEPT-90 and bbdecoupled carbon nmr yields a complete list of carbon fragments in the molecule. After evaluation of DEPT experiment (or multiplicity edited HSQC) heavy atoms are labeled with chemical shifts and number of attached hydrogen atoms.

Experiments: J-Couplings: HSQC

alias HMQC, CH-COSY, HETCOR

Acronyms:

HSQC: Heteronuclear Single Quantum Coherence

HMQC: Heteronuclear Multiple Quantum Coherence

- Cross-Signals in HSQC-diagrams caused by CH-Couplings via one CH-bond(¹J_{CH}, ~140 Hz).
- Experiment yields list of pairs of directly bonded carbon and hydrogen atoms.



Experiments: J-Couplings: HSQC

alias HMQC, CH-COSY, HETCOR



After evaluation of HSQC, our carbon atoms are also labled with the proton shifts of their directly attached protons



Experiments: J-Couplings: HH-COSY

Acronyms:

DQF-COSY: Double Quantum Filtered COrrelation Spectroscopy

- ³J_{HH}-Couplings (and a few others, unfortunately)
- Proton-rich skeletons might be elucidated just by HHCOSY und HSQC
- Problem: Quarternary carbons, heteroatoms in the skeleton



Experiments: J-Couplings: HH-COSY



Experiments: J-Couplings: HH-COSY



Experiments: J-Couplings: HMBC

Acronyms:

HMBC: Heteronuclear Multiple Bond Coherence

COLOC: COrrelation via LOng range Kopplungen

- Cross-signals through scalar couplings between carbon and hydrogen via 2 or 3 bonds (${}^{2}J_{CH}/{}^{3}J_{CH}$, ~8 Hz).
- Problem: ²J_{CH}/³J_{CH}-couplings cannot be distinguished.
- Problem: ⁴J_{CH}/⁵J_{CH}-couplings, which cannot be easily distinguished from ²J_{CH}/ ³J_{CH}-couplings



Experiments: J-Couplings-> HMBC



After evaluation of HMBC, we are looking at a molecular puzzle of pairs of carbon atoms that are either 1 or 2 bonds apart (but we don't know which is the case).





Computer-Assisted Structure Elucidation (CASE): Step by Step



SOFTWARE





Egon L. Willighagen^{1*}, John W. Mayfield², Jonathan Alvarsson³, Arvid Berg³, Lars Carlsson⁴, Nina Jeliazkova⁵, Stefan Kuhn⁶, Tomáš Pluskal⁷, Miquel Rojas-Chertó⁸, Ola Spjuth³, Gilleain Torrance⁹, Chris T. Evelo¹, Rajarshi Guha¹⁰ and Christoph Steinbeck¹¹

Abstract

Background: The Chemistry Development Kit (CDK) is a widely used open source cheminformatics toolkit, providing data structures to represent chemical concepts along with methods to manipulate such structures and perform computations on them. The library implements a wide variety of cheminformatics algorithms ranging from chemical structure canonicalization to molecular descriptor calculations and pharmacophore perception. It is used in drug discovery, metabolomics, and toxicology. Over the last 10 years, the code base has grown significantly, however, resulting in many complex interdependencies among components and poor performance of many algorithms.

Results: We report improvements to the CDK v2.0 since the v1.2 release series, specifically addressing the increased functional complexity and poor performance. We first summarize the addition of new functionality, such atom typing and molecular formula handling, and improvement to existing functionality that has led to significantly better performance for substructure searching, molecular fingerprints, and rendering of molecules. Second, we outline how the CDK has evolved with respect to quality control and the approaches we have adopted to ensure stability, including a code review mechanism.

Conclusions: This paper highlights our continued efforts to provide a community driven, open source cheminformatics library, and shows that such collaborative projects can thrive over extended periods of time, resulting in a high-quality and performant library. By taking advantage of community support and contributions, we show that an open source cheminformatics project can act as a peer reviewed publishing platform for scientific computing software.

Keywords: Java, Cheminformatics, Bioinformatics, Metabolomics, Depiction

Structure Elucidation versus Identification/ Dereplication





40 years of CASE research

Applications of Artificial Intelligence for Chemical Inference. 22. Automatic Rule Formation in Mass Spectrometry by Means of the Meta-DENDRAL Program^{1a}

B. G. Buchanan,* D. H. Smith, W. C. White, R. J. Gritter,^{1b} E. A. Feigenbaum, J. Lederberg, and Carl Djerassi

Contribution from the Departments of Computer Science, Chemistry, and Genetics, Stanford University, Stanford, California 94305. Received January 27, 1976

Abstract: The DENDRAL computer program uses established rules of molecular fragmentation to help chemists solve complex structural problems from mass spectral data. This paper describes a computer program called Meta-DENDRAL, that can aid in the discovery of such rules from empirical data on known compounds. The program uses heuristic methods to search for common structural environments around those bonds that are found to fragment and abstracts plausible fragmentation rules. The program has been tested on the well-characterized, low-resolution mass spectra of aliphatic amines and the high-resolution mass spectra of estrogenic steroids. The program has also discovered new fragmentation rules for mono-, di-, and trike-toandrostanes.

The DENDRAL computer program is designed to aid chemists with complex structure elucidation problems. One main part uses established molecular fragmentation rules to help chemists interpret mass spectra;² another main part generates lists of isomers that satisfy constraints derived from a variety of spectroscopic techniques.³ Because the mass spectrometry rules used by the DENDRAL program have been culled from the literature, the program's growth depends upon manual examination of collections of spectra. But investigating the spectral data of new compound classes to determine frag-

Journal of the American Chemical Society / 98:20 / September 29, 1976

COMPUTER PROGRAM FOR STRUCTURE ELUCIDATION

J. Chem. Inf. Comput. Sci., Vol. 18, No. 4, 1978 211

CHEMICS-F: A Computer Program System for Structure Elucidation of Organic Compounds

SHIN-ICHI SASAKI,* HIDETSUGU ABE, YUJI HIROTA, YOSHIAKI ISHIDA, YOSHIHIRO KUDO, SHUKICHI OCHIAI, KEIJI SAITO, and TOHRU YAMASAKI

Miyagi University of Education, Aoba, Sendai 980 Japan

Received November 3, 1977

A computer program system CHEMICS-F for the structure elucidation of organic molecules containing C, H, and O is described in detail. CHEMICS-F involves the software used to analyze spectra, to convert the spectral information into "components" (defined substructures), and to construct a molecular structure based on the components by means of a newly developed method. Multiple structural formulas as final output are often constructed based on both desirable and undesirable components designated by the spectral information. To minimize the multiplicity of answers, a file handling procedure has been integrated in which the spectra of NMR, IR, and MS are stored in compressed and concentrated shapes capable of identification of compounds.

Fresenius Z Anal Chem (1982) 313:473-479

Computer-Assisted Structure Elucidation

Morton E. Munk¹, Craig A. Shelley², Hugh B. Woodruff³, and Mark O. Trulson⁴

¹ Department of Chemistry, Arizona State University, Tempe, Arizona 85287, USA

² Research Laboratories, B82, Kodak Park, Rochester, New York 14650, USA

³ Merck, Sharp & Dohme Research Laboratories, P.O. Box 2000, Rahway, New Jersey 07065, USA

⁴ Department of Chemistry, University of California, Berkeley, California 94720, USA

Strukturaufklärung mit Computer-Hilfe

Zusammenfassung. Ein in Entwicklung befindliches Computer-Modell ("program CASE") wird beschrieben, das hauptsächlich folgende drei Aufgaben erfüllt: Zuordnung der chemischen und spektralen Eigenschaften einer Verbindung unbekannter Struktur zu ihren strukturellen Bedeutungen; Aufstellung aller Strukturen, die mit den gefundenen Eigenschaften im Einklang stehen; Einordnung dieser Strukturen je nach Übereinstimmung mit den vorhergesagten und gefundenen Parametern. Der derzeitige Entwicklungsstand des Programms wird erläutert und die Anwendung diskutiert.

Aufstellung aller Strukturen, die mit den gelundenen Eigenschaften im Einklang stehen; Einordnung dieser Strukturen je nach Übereinstimmung mit den vorhergesagten und gefundenen Parametern. Der derzeitige Entwicklungsstand des Programms wird erläutert und die Anwendung diskutiert. Using intuition and experience! The solution may be a "procedure", perhaps resembling a set of rules.

Consider the problem of structure elucidation. A piochemist isolates an antipiotic from a fermentation peer, or an organic chemist detects a trace impurity in the commercial synthesis of a drug. The structure of the compound may be completely unknown or it is possible that it is related to a known compound and that a structure can be suggested. Of course, in the latter case confirmation is reduired. How does the chemist broceed? Generally, the chemical and spectral properties of the "unknown" are examined. Today there is a known compound and that a structure can be suggested. Of course, in the latter case confirmation is reduired. How does the chemist proceed? Generally, the chemical and spectral proberties of the "unknown" are examined. Today there is a heavy embhasis on spectral properties because these data – the infrared spectrum proton and carbon-13 nuclear masthe infrared spectrum broton and carbon-13 nuclear master infrared spectrum broton and carbon-13 nuclear master is a heavy emphasis on spectral problem broton and carbon-13 nuclear master is a heavy emphasis on spectral problem broton and carbon-13 nuclear master is a heavy emphasis on spectral problem broton and carbon-13







Munk, M.E. et al., 1982. Computer-assisted structure elucidation. Fresenius' Zeitschrift für analytische Chemie, 313(6), pp.473–479.

474

Application of Expert System CISOC-SES to the Structure Elucidation of Complex Natural Products

Chen Peng,[†] Shengang Yuan,^{•,†} Chongzhi Zheng,[†] Yongzheng Hui,[†] Houming Wu,[‡] and Kan Ma[‡]

Laboratory of Computer Chemistry and State Key Laboratory of Bioorganic and Natural Products Chemistry, Shanghai Institute of Organic Chemistry, Chinese Academy of Sciences, 354 Fenglin Lu, Shanghai 200032, People's Republic of China

Xiuwen Han

Dalian Institute of Chemical Physics, Chinese Academy of Sciences, P.O. Box 100, Dalian 116011, People's Republic of China

Received October 12, 1993*

This paper demonstrates the application of a newly developed computer-assisted structure elucidation system, CISOC-SES, to the structure determination of complex natural products. Discussed is the structure elucidation of four natural products, including one new compound, with up to 50 non-hydrogen atoms principally from their 2D NMR spectral data. The effectiveness of the novel approaches that exploit the direct and long-range distance constraints is analyzed in detail. The results show that the efficiency of structure elucidation can be substantially improved with the help of the system so that a manageable number of candidate structures, with the correct structure always included, can be obtained in a supportable period of CPU time.

LSD

(Logic for Structure Determination)



- LSD (Logic for Structure Determination
- Command-line driven
- Takes spectral constraints as input
- Generates lists of connection tables (molecules)
- Open Source
- Rocket-fast
- No early spectrum processing

Ab-Initio Structure Elucidation by 2D NMR

2D peak picking table

No.	¹³ C CPD	¹ J _{CH} (HMQC)	HMBC
1	148.24	-	4.0; 2.42; 2.28; 1.20
2	118.25	5.49	4.0; 2.28; 2.17
3	66.32	4.0	5.49
4	43.87	2.17	5.49; 4.0; 2.42; 1.32; 1.20; 0.86
5	41.40	2.14	5.49; 1.32; 0.86
6	38.32	-	2.42; 1.32; 1.20; 0.86
7	32.00	1.20/2.42	
8	31.52	2.28	2.42; 1.20
9	26.52	1.32	0.86; 2.17
10	21.46	0.86	1.20; 1.32

Table with heavy-atom relations Internally generated by CASE program



Steinbeck, C. Computer-Assisted Structure Elucidation. In Handbook on Chemoinformatics.; Gasteiger, J. Ed.; Wiley-VCH: Weinheim, 2003; Vol. 2; pp. 1378-1406.

• Deterministic methods suffer from combinatorial explosion

• Deterministic methods suffer from combinatorial explosion



No. of Heavy Atoms

• Deterministic methods suffer from combinatorial explosion

 $C_{10}H_{16}$ (10 Heavy Atoms)

24938 Constitutional Isomers

No. of Constitutional Isomer Calculation Time



No. of Heavy Atoms

• Deterministic methods suffer from combinatorial explosion

 $C_{10}H_{16}$ (10 Heavy Atoms)

24938 Constitutional Isomers

C₁₃H₁₆O₃ (16 Heavy Atoms)

> 2,000,000,000 Constitutional Isomers

No. of Constitutional Isomer Calculation Time f(t) = f(t) + f(t

No. of Heavy Atoms

• Deterministic methods suffer from combinatorial explosion

 $C_{10}H_{16}$ (10 Heavy Atoms)

24938 Constitutional Isomers

C₁₃H₁₆O₃ (16 Heavy Atoms)

> 2,000,000,000 Constitutional Isomers

 $C_{30}H_{48}O_2$ (32 Heavy Atoms)

>> 10¹² Constitutional Isomers

No. of Constitutional Isomer Calculation Time



No. of Heavy Atoms

• Deterministic methods suffer from combinatorial explosion

 $C_{10}H_{16}$ (10 Heavy Atoms)

24938 Constitutional Isomers

 $C_{13}H_{16}O_3$ (16 Heavy Atoms)

> 2,000,000,000 Constitutional Isomers

C₃₀H₄₈O₂ (32 Heavy Atoms)

>> 10¹² Constitutional Isomers

No. of Constitutional Isomer Calculation Time



No. of Heavy Atoms

• Prospective use of spectroscopic input information may make them error-intolerant

Deterministic Structure Generators: The LUCY Method

- Prospective use of spectral information for building isomers
- Needs 1D ¹³C, 2D HMQC, HMBC, HH COSY
- Example: Walk a decision tree while interpreting HMBC-Signals



Steinbeck, C.; Angewandte Chemie. International Ed. in English 1996, 35, 1984-1986.

LSD Input Syntax: Basics AcAlaOMe Example



Command prompt don't type this in



\$ outlsd 7 < acalaome.sol > acalaome.sdf

Command prompt don't type this in



will produce acalaome.sol
(LSD specific output file)

- \$ lsd AcAlaOMe
- \$ outlsd 7 < acalaome.sol > acalaome.sdf

Command prompt don't type this in



will produce acalaome.sol
(LSD specific output file)

\$ lsd AcAlaOMe

\$ outlsd 7 < acalaome.sol > acalaome.sdf

converts .sol into .sdf (standard molecular file format) Use e.g. MarvinView to view .sdf

MarvinView rendering of the one result in accordance with out input data from the previous slide


LSD CASE reveals one more possible solution

Liu et al., C₃₄H₅₂O₆, HRESI-MS (m/z 557.3833 [m+H])



No.	14					
	δ _C	$\delta_{\rm H}$ muti (J/Hz)	НМВС			
1	206.3					
2	74.5					
3	58.5					
4	39.2	1.70 m	C-2,C-3,C-5,C-6,			
			C-14,C-15, C-22			
5	35.6	(Ha)2.61dd(14.1,3.0)	C-1, C-3, C-4,			
		1.21 m				
		(Hb)1.21 m	C-6,C-7, C-21			
6	61.3					
7	195.8					
8	150.6					
0						
9	217.0					
10	42.7	3.20 dp(12.8, 6.4)	C-10 C-12 C-13			
12	42.7	1 28 d (6 5)	C-10, C-12, C-13			
12	10.8	1.36 U (0.3)	C 10 C 11 C 12			
13	17.0	0.00 c	C = 10, C = 11, C = 12			
14	25.2	0.90 S	(-2, (-3, (-4, (-15)))			
15	55.5	1.8/ddd(12.0,9.7,	C = 2, C = 3, C = 4,			
10	25.7	3.1), 1.42 III	C-14,C-10,C-17			
10	25.7	2.32 111	C-3,C-13,C-17,C-18			
17	55.3	2.88 t(10.0)	C-1.C-2.C-3.C-10.C-11,			
			C-16,C-18,C-19,C-20			
18	86.1					
19	24.9	1.19 s	C-17,C-18,C-20			
20	24.2	1.15 s	C-17,C-18,C-19			
21	29.4	2.02dd(12.8,5.5),	C-3,C-4,C-5,			
		1.66 m	C-22, C-23			
22	122.6	5.16 t(6.4)	C-4,C-21, C-23,			
			C-24,C-25			
23	133.4					
24	25.8	1.76 s	C-22, C-23, C-25			
25	18.0	1.60 s	C-22, C-23, C-24			
26	125.1	5.92 d(6.8)	C-7, C-8, C-18,			
			C-27,C-28			
27	74.0	4.19 d(6.8)	C-7,C-8,C-28,			
			C-26,C-29,C-30			
28	72.4					
29	26.6	1.19 s	C-27,C-28,C-30			
30	25.8	1.15 s	C-27,C-28,C-29			
31	34.1	2.35 m, 2.27 m	C-1,C-5,C-6,			
			67632633			
22	117.0	100 (7 2)	C-7,C-32,C-33			
32	117.0	4.90 L(7.2)	L-0,L-31,L-33, L-34			
33	135.2					
34	26.0	1.69 s	C-32,C-33,C-35			
35	18.1	1.55 s	C-32,C-33,C-34			

Liu, R., Su, Y., Yang, J., & Wang, A. (2017). Polyprenylated acylphloroglucinols from Hypericum scabrum. Phytochemistry, 142, 38–50.

LSD Input Syntax: Advanced

Liu et al., C₃₄H₅₂O₆, HRESI-MS (m/z 557.3833 [m+H])

ELIM 3 4

MULT 1 C 2 0 MULT 2 C 3 0 [... 35 more omitted ...]

; known carbonyls BOND 1 36 BOND 7 37 BOND 10 39

CARB L1 ; define list L1 containing all carbons HETE L2 LIST L3 2 3 6 8 17 18 27 28 LIST L4 9 38 40 PROP L1 1 L2 - ; Every carbon atom ; can carry one or less ; hetero-atoms, but not ; two
PROP L4 0 L3 ; Every oxygen which is ; not an sp2 O has a ; parter from L3 (based

; on a conservative

; chemical shift

; inspection)

COSY 4 5 [... 4 more omitted ...]

HMQC 4 4 [... more omitted ...]

HMBC 2 4 [... 98 more omitted ...]

http://eos.univ-reims.fr/LSD/index_ENG.html

Liu et al., C₃₄H₅₂O₆, HRESI-MS (m/z 557.3833 [m+H])



InChIKey=ANWFPAAUCGPEBV-MOHJPFBDNA-N



InChIKey=YAJAXOAXTCGOQA-HKOYGPOVNA-N

Liu et al., C₃₄H₅₂O₆, HRESI-MS (m/z 557.3833 [m+H])





InChIKey=ANWFPAAUCGPEBV-MOHJPFBDNA-N

InChIKey=YAJAXOAXTCGOQA-HKOYGPOVNA-N

InChIKey of published compound #14 = YAJAXOAXTCGOQA-HKOYGPOVNA-N

Quantum Chemistry Calculations for Validation



- Commercial Software (Gaussian 16, Spartan)
- Open Software (NWChem)
- Tricky setup, tricky analysis
- Potentially long calculation time (> day on your PC)

1. Buevich, A. V.; Elyashberg, M. E. Synergistic Combination of CASE Algorithms and DFT Chemical Shift Predictions: A Powerful Approach for Structure Elucidation, Verification, and Revision. J. Nat. Prod 2016, 79, 3105–3116.









Structure #6

	exptl	structure #1	structure #2 (4)	structure #3	structure #4	structure #5	structure #6
label	δC	$\delta C_{ m calc}$	δC_{calc}	$\delta C_{ m calc}$	δC_{calc}	δC_{calc}	$\delta C_{ m calc}$
C1	168.7	171.3	172.0	165.9	165.6	164.0	160.4
C2	79.8	81.6	79.6	69.6	85.1	84.4	102.9
C3	127.2	142.7	135.6	146.1	147.1	140.9	156.3
C4	143.1	125.4	144.9	134.1	135.8	136.1	125.7
C5	164.5	171.0	163.8	160.8	155.7	163.8	164.1
C7	155.7	154.6	154.7	158.8	160.3	159.6	154.7
C8	108.1	105.3	105.5	105.3	101.5	103.8	105.7
C9	135.8	134.8	134.3	125.3	126.4	126.7	134.1
C10	112.7	111.9	112.0	106.5	108.9	109.1	111.8
C11	160.6	160.5	160.5	145.3	150.2	148.9	160.8
C12	110.7	109.8	110.3	110.2	108.8	108.6	110.1
C13	176.0	173.9	172.9	159.1	173.1	180.5	172.3
C14	120.8	117.2	120.6	121.1	124.7	122.7	117.3
C15	52.9	54.1	53.7	53.8	53.2	51.7	52.0
RMSD, ppm		6.75	2.76	9.46	7.86	6.40	11.29
max_ <i>δ</i> , ppm		17.71	8.39	18.94	19.87	13.67	29.07

15

1. Buevich, A. V.; Elyashberg, M. E. Synergistic Combination of CASE Algorithms and DFT Chemical Shift Predictions: A Powerful Approach for Structure Elucidation, Verification, and Revision. J. Nat. Prod 2016, 79, 3105–3116.

Stochastic CASE

Stochastic Search Methods

Algorithms known to tackle large search spaces:

- Simulated Annealing
 - Traveling Salesman Problem
 - Finding the solution structure of large biomolecules
 - Integrated Circuits Layout
 - Robotic Path Planning
- Genetic Algorithms
 - Protein Folding
 - Immune System Simulation
 - Computer-Aided Design
- Quite a number of other options ...













Guided Walk in Constitution Space





 are in close chemical distance to each other





- are in close chemical distance to each other
- are likely to be similar in their spectroscopic properties





- are in close chemical distance to each other
- are likely to be similar in their spectroscopic properties





Evaluating a score for each point (constitution) in structure space



Score Function		
based on		
Spectroscopic Fitness		
S _{total} =		
$c_1 S_{HMBC}$	+	
c ₂ S _{HHCOSY}	+	
$c_3 S_{Shift}$	+	
	+	
c _n S _{Features}		



Acceptance criterion

$$p = \exp(-\frac{\delta f}{T})$$





Acceptance criterion

$$p = \exp(-\frac{\delta f}{T})$$





Acceptance criterion

$$p = \exp(-\frac{\delta f}{T})$$





Acceptance criterion

$$p = \exp(-\frac{\delta f}{T})$$



Pluggable Target Function

- Pluggable Target Function
 - If you can reliably calculate a measurable property for a given constitution, it can be part of your target function

- Pluggable Target Function
 - If you can reliably calculate a measurable property for a given constitution, it can be part of your target function
 - Spectroscopic information
 - IR
 - UV-VIS
 - Other types of NMR experiments
 - MS fragmentation (?)

Pluggable Target Function

- If you can reliably calculate a measurable property for a given constitution, it can be part of your target function
 - Spectroscopic information
 - IR
 - UV-VIS
 - Other types of NMR experiments
 - MS fragmentation (?)
 - Additional knowledge
 - Good-List/Bad-List fragments
 - Drug Likeness
 - Natural Product Likeness

- Pluggable Target Function
 - If you can reliably calculate a measurable property for a given constitution, it can be part of your target function
 - Spectroscopic information
 - IR
 - UV-VIS
 - Other types of NMR experiments
 - MS fragmentation (?)
 - Additional knowledge
 - Good-List/Bad-List fragments
 - Drug Likeness
 - Natural Product Likeness

General System for Optimization in Constitution Space

- Pluggable Target Function
 - If you can reliably calculate a measurable property for a given constitution, it can be part of your target function
 - Spectroscopic information
 - IR
 - UV-VIS
 - Other types of NMR experiments
 - MS fragmentation (?)
 - Additional knowledge
 - Good-List/Bad-List fragments
 - Drug Likeness
 - Natural Product Likeness

General System for Optimization in Constitution Space

• Artifacts only lead to slightly lower ranking of correct structure in hit list

Distributed computing at its cheapest.









Distributed Server

Qualitative assessment:

Computational complexity of deterministic and stochastic algorithms

Compound	LUCY	SENECA	Steps overall
α -pinene (C ₁₀ H ₁₆)	2 s	1 min	30,000
Eurabidiol(C ₁₅ H ₂₈ O ₂)	29 s	5 min	90 000
Polycarpol (C ₃₀ H ₅₀ O)	33 min	12 min	350,000









 α -Pinene (C₁₀H₁₆)

Eurabidiol (C₁₅H₂₈O₂)

Polycarpol ($C_{30}H_{48}O_2$).

C. Steinbeck, Journal of Chemical Information & Computer Sciences 2001, 41, 1500.

Ranking Solutions

- Heteroatom-rich/proton-poor skeletons can yield many solutions (hundreds, thousands)
- Possible ranking by
 - Spectrum Similarity
 - Natural Product Likeness

Natural Product-likeness Score and Its Application for Prioritization of Compound Libraries

Peter Ertl,* Silvio Roggo, and Ansgar Schuffenhauer

Novartis Institutes for BioMedical Research, CH-4002 Basel, Switzerland



Datasets



Components for Molecule Curation

Components for Molecule Curation


Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Jean-Loup Faulon,*,[†] Michael J. Collins,[‡] and Robert D. Carr[§]



Component for Signature Generation



$$Fragment_i = \log \left[\frac{NP_i}{SM_i} * \frac{SM_t}{NP_t} \right]$$

$$Fragment_i = \log \left[\frac{NP_i}{SM_i} * \frac{SM_t}{NP_t} \right]$$

NP - Natural product SM - Synthetic molecule

 $Fragment_i = \log \left| \frac{NP_i}{SM_i} * \frac{SM_t}{NP_t} \right|$

NP - Natural product SM - Synthetic molecule

In the fragment contribution (*Fragment*_i),

 \bigstar *NP_i* is the total number of molecules in the NP dataset in which the *Fragment*_i occurs,

 \bigstar *SM*^{*i*} is the total number of molecules in the SM dataset in which the *Fragment*^{*i*} occurs,

- \bigstar *SM*_t is the total number of molecules int he SM dataset
- \bigstar *NP*^{*t*} is the total number of molecules in the NP dataset.
- ◆N is the number of fragments in given molecule

 $Fragment_i = \log \left| \frac{NP_i}{SM_i} * \frac{SM_t}{NP_t} \right|$

NP - Natural product *SM* - Synthetic molecule

$$Score_N = \sum_{i=0}^{N} Fragment_i$$

In the fragment contribution (*Fragment*_i),

 \bigstar *NP_i* is the total number of molecules in the NP dataset in which the *Fragment_i* occurs,

 \bigstar *SM_i* is the total number of molecules in the SM dataset in which the *Fragment*_i occurs,

- A *SM*^{*t*} is the total number of molecules int he SM dataset
- \bigstar *NP*^{*t*} is the total number of molecules in the NP dataset.
- ◆N is the number of fragments in given molecule

$$Fragment_i = \log \left[\frac{NP_i}{SM_i} * \frac{SM_t}{NP_t} \right]$$

NP - Natural product SM - Synthetic molecule

$$Score_N = \sum_{i=0}^{N} Fragment_i$$

In the fragment contribution (*Fragment*_i),

 \bigstar *NP_i* is the total number of molecules in the NP dataset in which the *Fragment_i* occurs,

 \bigstar *SM*^{*i*} is the total number of molecules in the SM dataset in which the *Fragment*^{*i*} occurs,

- A *SM*^{*t*} is the total number of molecules int he SM dataset
- \bigstar *NP*^{*t*} is the total number of molecules in the NP dataset.
- ◆N is the number of fragments in given molecule

$$NormalisedScore = \frac{Score_N}{N}$$

Natural Product-likeness classification and integrated it into Taverna workflow tool

- (<u>http://sourceforge.net/projects/np-likeness/</u>).
- Included in second version of SENECA CASE



Jayaseelan KV, Moreno P, Truszkowski A, Ertl P & Steinbeck C (2012) Natural product-likeness score revisited: an open-source, open-data implementation. BMC Bioinformatics 13, 106.

Distribution of NP-likeness scores



Score

Availability





Publicationen

Publicationen

New developments on the cheminformatics open workflow environment CDK-Taverna

Andreas Truszkowski¹, Kalai Vanii Jayaseelan², Stefan Neumann³, Egon L Willighagen⁴, Achim Zielesny¹and Christoph Steinbeck^{*2}

Publicationen

New developments on the cheminformatics open workflow environment CDK-Taverna

Andreas Truszkowski¹, Kalai Vanii Jayaseelan², Stefan Neumann³, Egon L Willighagen⁴, Achim Zielesny¹and Christoph Steinbeck^{*2}

Natural product-likeness score revisited: an open-source, open-data implementation

Kalai Vanii Jayaseelan^{*1}, Christoph Steinbeck^{*1}, Pablo Moreno¹, Andreas Truszkowski², Peter Ertl³

