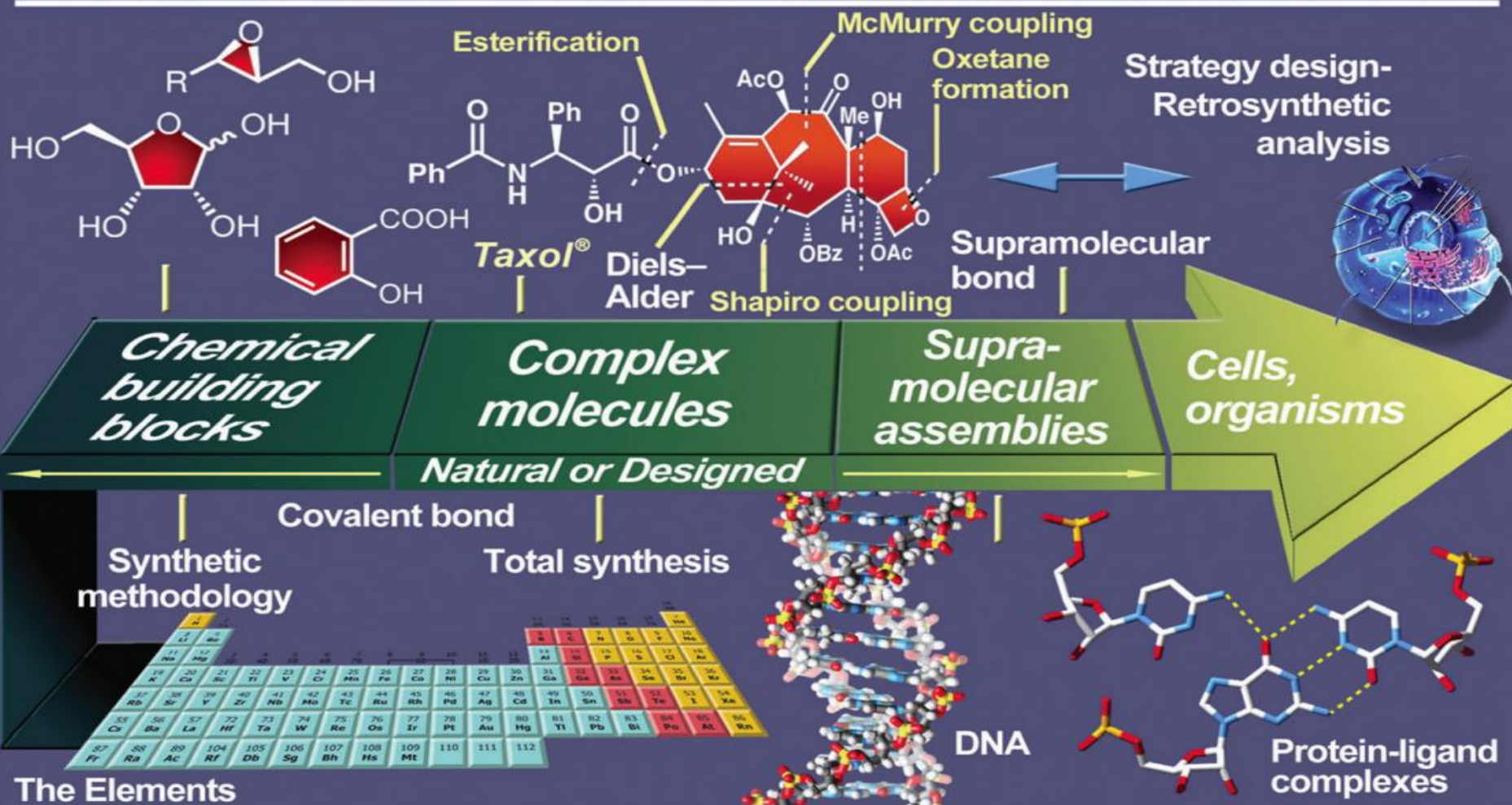


Literature Seminar

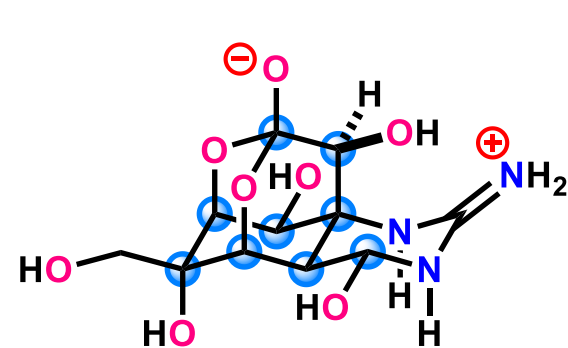
2020/11/21 Masanori Nagatomo

Molecular Complexity

Molecular Complexity and Chemical Synthesis

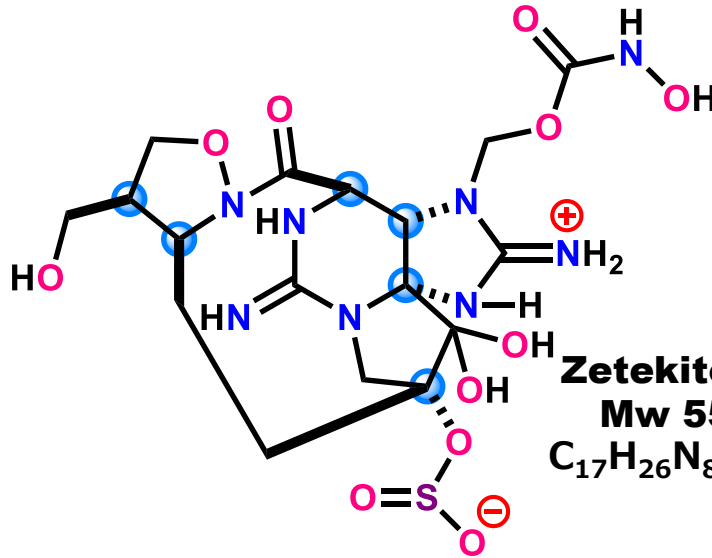
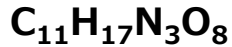


What Is Complexity?



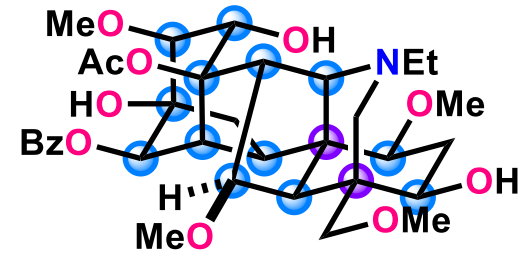
Tetrodotoxin

Mw 319



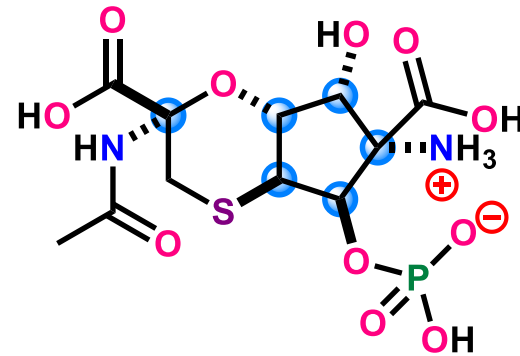
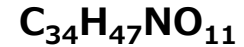
Zetekitoxin

Mw 550



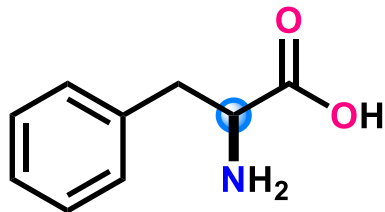
Aconitine

Mw 645



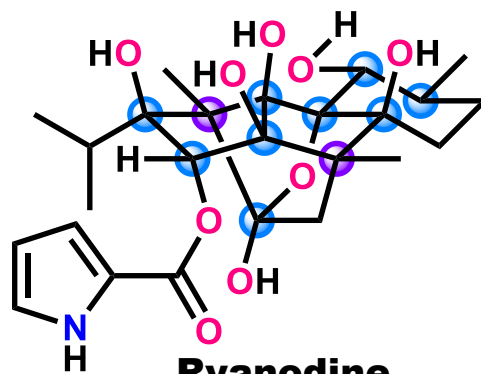
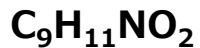
Tagetitoxin

Mw 416



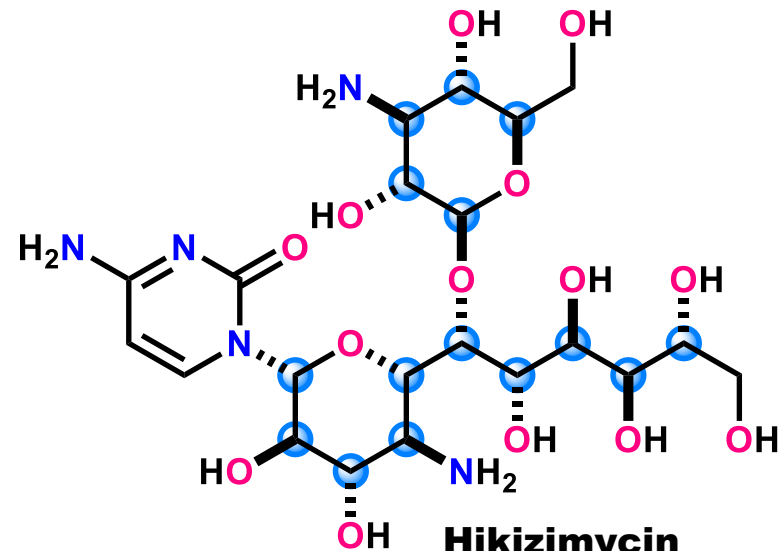
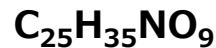
L-Phenylalanine

Mw 116



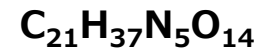
Ryanodine

Mw 493

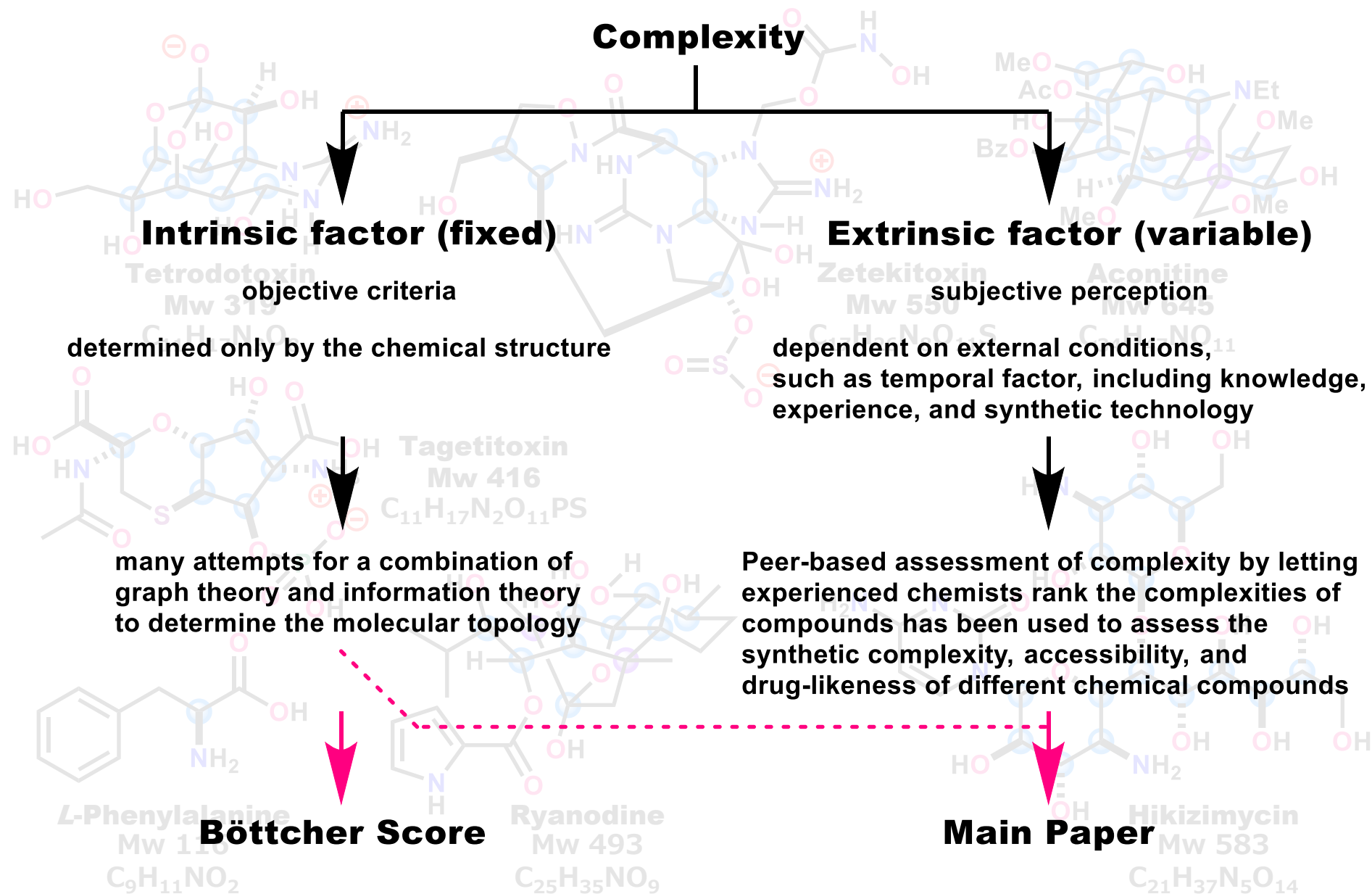


Hikizimycin

Mw 583

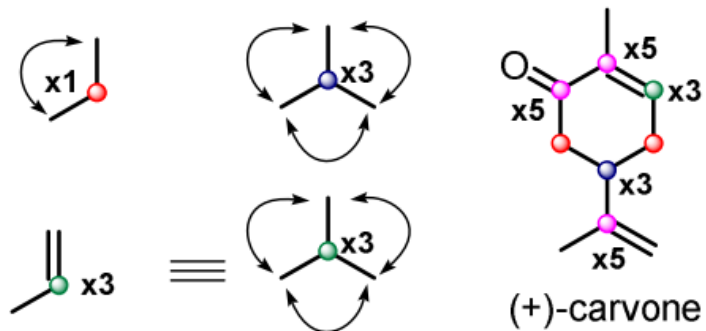


How Do We Chemist Define Molecular Complexity? ³



Bertz-(Hendrickson) Index

A Bertz's principle

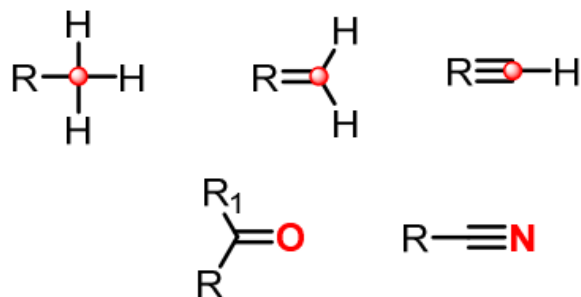


$$C_T = C_\eta + C_\epsilon$$

C_η : bond connectivities - skeletal complexity

C_ϵ : diversity of elements - kinds of atoms

C Zero complexity

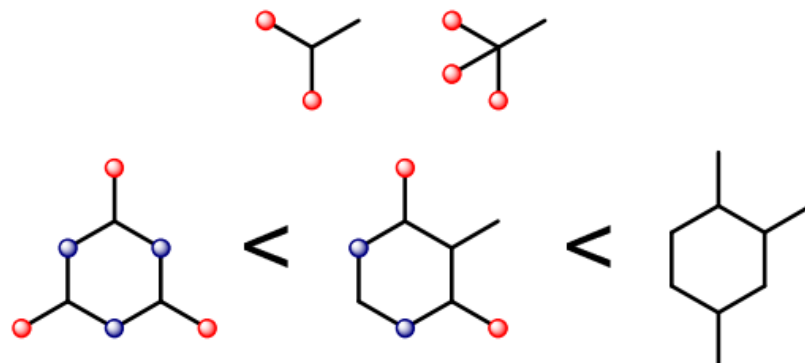


terminal atoms does not increment the complexity calculation

B complexity gain



D decrease in complexity (symmetry)



1. Bertz, S. H. *J. Am. Chem. Soc.* **1981**, *103*, 3599-3601.
2. Siguara Bastos de Lemos E Silva. Chemistry and biosynthesis of highly complex marine alkaloids from Mediterranean biodiversity. Organic chemistry. Université Paris Saclay (COMUE), 2017.

Applications and Shortcomings

C_T of many chemical compounds is publically available on PubChem.¹⁾ Molecular complexity also has important implications for organic synthesis planning, **in-silico drug design**, and pharmaceutical development, including **QSPR (Quantitative Structure-Property Relationship)** and **QSAR (Quantitative Structure-Affinity Relationship)** approaches.



Frequently criticized shortcomings are the **failure to address chirality in graph-theoretical approaches** and **missing sensitivity to skeletal structure**, branching, and symmetry of other indices.



Proposing a framework for molecular complexity that relies on both **mathematical rigor** and **chemically consistent inherent logic**.

Böttcher Score

1. Kim, S.; Thiessen, P. A.; Bolton, E. E.; Chen, J.; Fu, G.; Gindulyte, A.; Han, L.; He, J.; He, S.; Shoemaker, B. A.; Wang, J.; Yu, B.; Zhang, J.; Bryant, S. H. Pubchem Substance and Compound Databases. *Nucleic Acids Res.* **2016**, *44*, D1202–D1213.

Thomas Böttcher



2009: Ph.D @ LMU, Munich (Prof. Stephan A. Sieber)
2010: Postdoctoral research@ TMU, Munich (Prof. Stephan A. Sieber)

2011-2014: Postdoctoral research@ Harvard Medical School, Department of Biological Chemistry & Molecular Pharmacology (Prof. Jon Clardy)

2014-2020: Independent group leader @ University of Konstanz

2020: An Emmy Noether research group leader @ University of Konstanz, Professor of Microbial Biochemistry, Faculty of Chemistry, Department of Biological Chemistry and Department of Microbiology and Ecosystem Science (DOME), University of Vienna

Research interest: the isolation and identification of natural products that modify and manipulate bacterial population behavior. Aim to inhibit bacterial virulence and discover the chemistry of ecological interactions of microorganisms.

BIOSPHERE COMPLEXITY:
<https://www.youtube.com/watch?v=xFNWVSEuuxc>



Prof. Stephan A. Sieber



Prof. Jon Clardy

Böttcher Score

The information content is defined by the entropy H (Unit is bit).

$$H = - \sum_i p_i \log_x p_i$$

One variable is needed to identify the nature of the element by its **valence shell**, and four variables are required as descriptors of the bonding environment: **the number of bonds**, **the number of chemically different bonds**, **the element diversity**, and **the stereochemistry**.

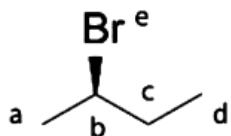
C_m : Molecular Complexity (Unit is mcbit), V_i : valence electrons, B_i : total number of bonds, d_i : introduced to characterize the number of chemically nonequivalent bonds to atoms with $V_{ibi} > 1$ at the i th position, e_i : giving the number of different non-hydrogen elements or isotopes involved in the bonding situation, including atom i and its direct neighbors, to include heteroatoms.

$$C_m = \sum_i d_i e_i s_i \log_2 (V_i b_i)$$

To account for symmetries of a molecule, the corresponding atom positions of chemically equivalent sets of atoms for each symmetric position j are subtracted.

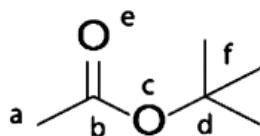
$$C_m = \sum_i d_i e_i s_i \log_2 (V_i b_i) - \frac{1}{2} \sum_j d_j e_j s_j \log_2 (V_j b_j)$$

Examples: How to Calculate C_m



i	a	b	c	d	e
d_i	1	3	2	1	1
e_i	1	2	1	1	2
s_i	1	2	1	1	1
V_i	4	4	4	4	7
b_i	1	3	2	1	1

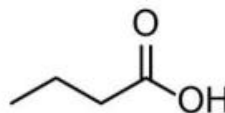
$$C_m = \log(4) + 3 \cdot 2 \cdot 2 \cdot \log(4 \cdot 3) + 2 \cdot \log(4 \cdot 2) + \log(4) + 2 \cdot \log(7) = 58.6 \text{ mcbits}$$



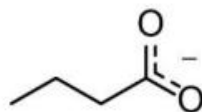
i	a	b	c	d	e	f
d_i	1	3	2	2	1	1
e_i	1	2	2	2	2	1
s_i	1	1	1	1	1	1
V_i	4	4	6	4	6	4
b_i	1	4	2	4	2	1

$$C_m = \log(4) + 3 \cdot 2 \cdot \log(4 \cdot 4) + 2 \cdot 2 \cdot \log(6 \cdot 2) + 2 \cdot 2 \cdot \log(4 \cdot 4) + 2 \cdot \log(6 \cdot 2) + 3 \cdot \log(4) - 0.5 \cdot 3 \cdot \log(4) = 66.5 \text{ mcbits}$$

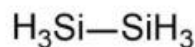
Examples illustrating how to calculate C_m on a per-atom basis for a molecule. Different atom positions i are labeled by letters (a–f), and for tert-butyl acetate the symmetry correction for the positions f is included. In the equations, “log” stands for \log_2 .



$$C_m = \log(4) + 2 \cdot \log(4 \cdot 2) + 2 \cdot \log(4 \cdot 2) + 3 \cdot 2 \cdot \log(4 \cdot 4) + 2 \cdot \log(6 \cdot 2) + 2 \cdot \log(6) = 50.34 \text{ mcbits}$$

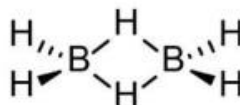


$$C_m = \log(4) + 2 \cdot \log(4 \cdot 2) + 2 \cdot \log(4 \cdot 2) + 2 \cdot 2 \cdot \log(4 \cdot 4) + 2 \cdot [2 \cdot \log(6 \cdot 1.5)] - 0.5 \cdot 2 \cdot [2 \cdot \log(6 \cdot 1.5)] = 36.34 \text{ mcbits}$$



$$C_m = 2 \cdot [\log(4)] - 0.5 \cdot 2 \cdot [\log(4)] = \log(4) = 2.00 \text{ mcbits}$$

Special bond situations and non-carbon based compounds. Log stands for \log_2 .



$$C_m = 2 \cdot [\log(3 \cdot 2) + \log(1 \cdot 2)] - 0.5 \cdot 2 \cdot [\log(3 \cdot 2) + \log(1 \cdot 2)] = 3.58 \text{ mcbits}$$

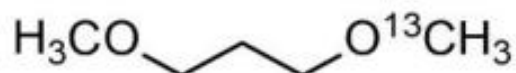
Examples: Influence of Isotopes on C_m

A



symmetric

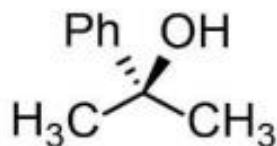
$$C_m = 2 \cdot \log(4) + 2 \cdot 2 \cdot \log(6 \cdot 2) + 2 \cdot 2 \cdot \log(4 \cdot 2) + \log(4 \cdot 2) \\ = 33.34 \text{ mcbits}$$



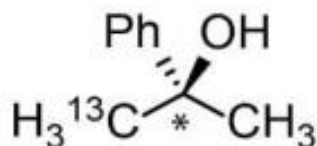
asymmetric

$$C_m = 2 \cdot \log(4) + 2 \cdot 2 \cdot \log(6 \cdot 2) + 2 \cdot 2 \cdot \log(4 \cdot 2) + \textcolor{red}{2} \cdot \log(4 \cdot 2) \\ + \textcolor{red}{2} \cdot 2 \cdot \log(4 \cdot 2) + \textcolor{red}{2} \cdot 3 \cdot \log(6 \cdot 2) + \textcolor{red}{2} \cdot \log(4) = 73.85 \text{ mcbits}$$

B

symmetric
achiral

$$C_m = \log(4 \cdot 3) + 2 \cdot \log(4 \cdot 3) + 2 \cdot \log(4 \cdot 3) + 2 \cdot \log(4 \cdot 4) + 3 \cdot 2 \cdot \log(4 \cdot 4) \\ + 2 \cdot \log(6 \cdot 2) + \log(4) = 57.09 \text{ mcbits}$$

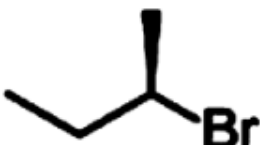

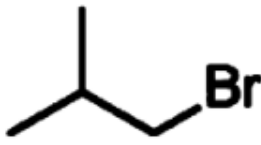
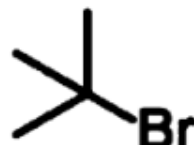
asymmetric
chiral

$$C_m = \log(4 \cdot 3) + 2 \cdot \log(4 \cdot 3) + 2 \cdot \log(4 \cdot 3) + 2 \cdot \log(4 \cdot 4) + \textcolor{red}{4} \cdot \textcolor{red}{3} \cdot 2 \cdot \log(4 \cdot 4) \\ + 2 \cdot \log(6 \cdot 2) + \log(4) + \textcolor{red}{\log(4)} = 131.09 \text{ mcbits}$$

Influence of isotopes on molecular complexity.

- A) The ^{13}C isotope of the methyl group disrupts the compound's symmetry, manifesting in changes in d_i for the central CH_2 group and the omission of symmetry correction. Additionally, the element diversity term e_i at one of the oxygens increases with the ^{13}C isotope.
- B) In addition to symmetry breaking at the quaternary carbon, stereochemical information (isotope chirality) is introduced, altering s_i and d_i and eliminating the symmetry correction term. Changes are highlighted in red. An asterisk labels the stereocenter. Log stands for \log_2 .

Comparison of C_m with Other Complexity Indices 10

				
non-equivalent protons (NMR)	5	4	3	1
C_m (mcbit)	58.6	31.6	26.8	24.6
N_S	10	9	9	9
N_T	17	15	9	9
$C(\eta, \epsilon)$	19.6	13.1	17.6	25.1
S (Whitlock)	3	1	1	1
Barone	45	42	45	54

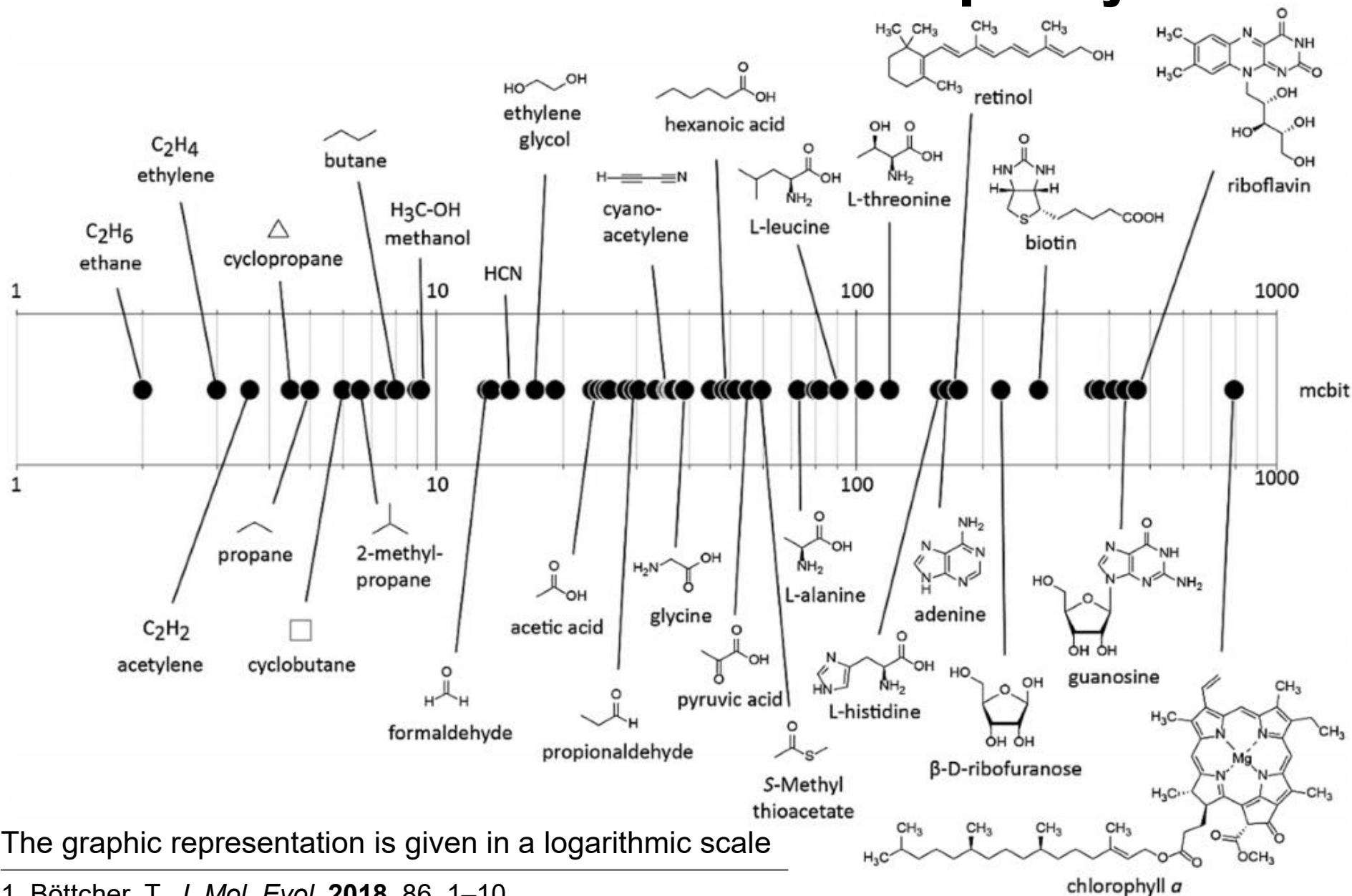
Molecular complexity of bromobutane isomers. The order of the isomers is given from the most complex (left) to the least complex (right) according to the spectral complexity of the ^1H NMR spectra as a proxy of molecular complexity. The numbers of nonequivalent protons are in line with the relative complexities given by distinct signals and multiplicities of the corresponding ^1H NMR spectra. Complexity values are given for C_m and other indices.

Changes in Molecular Complexity ΔC_m

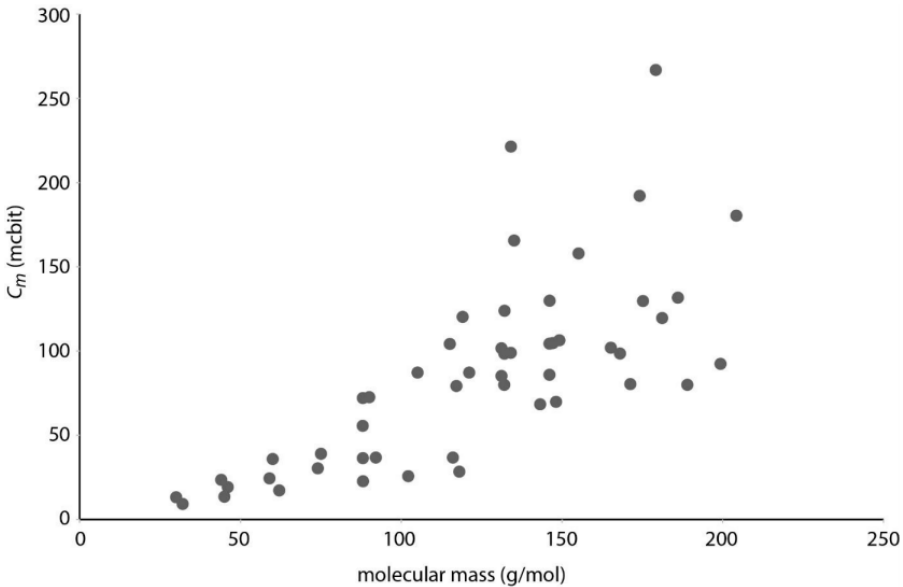
Name	Reaction	ΔC_m
Aldol condensation	$R-\overset{\text{O}}{\parallel}{\text{C}}-\text{H} + \text{CH}_3-\overset{\text{O}}{\parallel}{\text{C}}-\text{CH}_3 \xrightarrow[\text{-H}_2\text{O}]{[\text{OH}^-]} R-\text{CH}=\text{CH}-\overset{\text{O}}{\parallel}{\text{C}}-\text{CH}_3$	8.00
Aldol reaction	$R-\overset{\text{O}}{\parallel}{\text{C}}-\text{H} + \text{CH}_3-\overset{\text{O}}{\parallel}{\text{C}}-\text{CH}_3 \xrightarrow{(\text{S})\text{-proline}} R-\text{CH}(\text{OH})-\text{CH}_2-\overset{\text{O}}{\parallel}{\text{C}}-\text{CH}_3$	40.68
Clemmensen reduction*	$R-\overset{\text{O}}{\parallel}{\text{C}}-R' \xrightarrow[\text{HCl}]{\text{Zn(Hg)}} R-\text{CH}_2-R'$	-25.17
1,3-dipolar cycloaddition	$R-\text{N}_3 + \text{CH}_2=\text{CH}-R' \xrightarrow{\text{Cu}^{\text{I}}, \text{cat.}} R-\text{N}=\text{N}-\text{CH}(\text{R}')-\text{CH}_2-R'$	47.80
Jones oxidation	$R-\text{CH}_2-\text{OH} \xrightarrow[\text{H}_2\text{SO}_4, \text{H}_2\text{O}]{\text{CrO}_3} R-\overset{\text{O}}{\parallel}{\text{C}}-\text{OH}$	5.17
Rosenmund reduction	$R-\overset{\text{O}}{\parallel}{\text{C}}-\text{Cl} \xrightarrow{\text{H}_2/\text{Pd}} R-\overset{\text{O}}{\parallel}{\text{C}}-\text{H}$	-27.27
Schotten-Baumann reaction	$R-\overset{\text{O}}{\parallel}{\text{C}}-\text{Cl} + \text{H}_2\text{N}-R' \xrightarrow{\text{NaOH}} R-\overset{\text{O}}{\parallel}{\text{C}}-\text{NH}-R'$	3.03

C_m with Calculated for the Structures of Various Small Molecules in a Universal Complexity Scale

12



Molecular Complexity Per Molecular Mass Unit

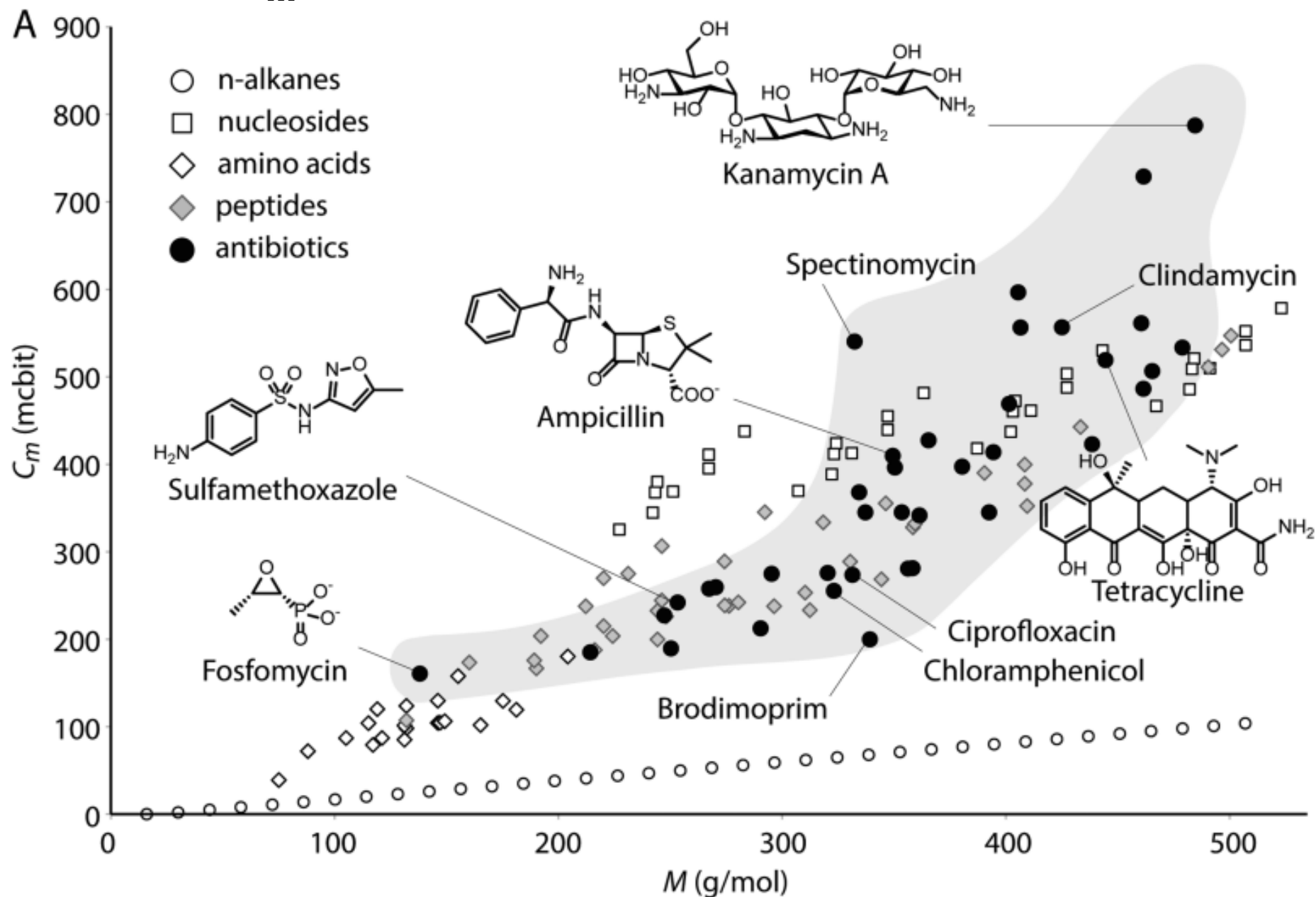


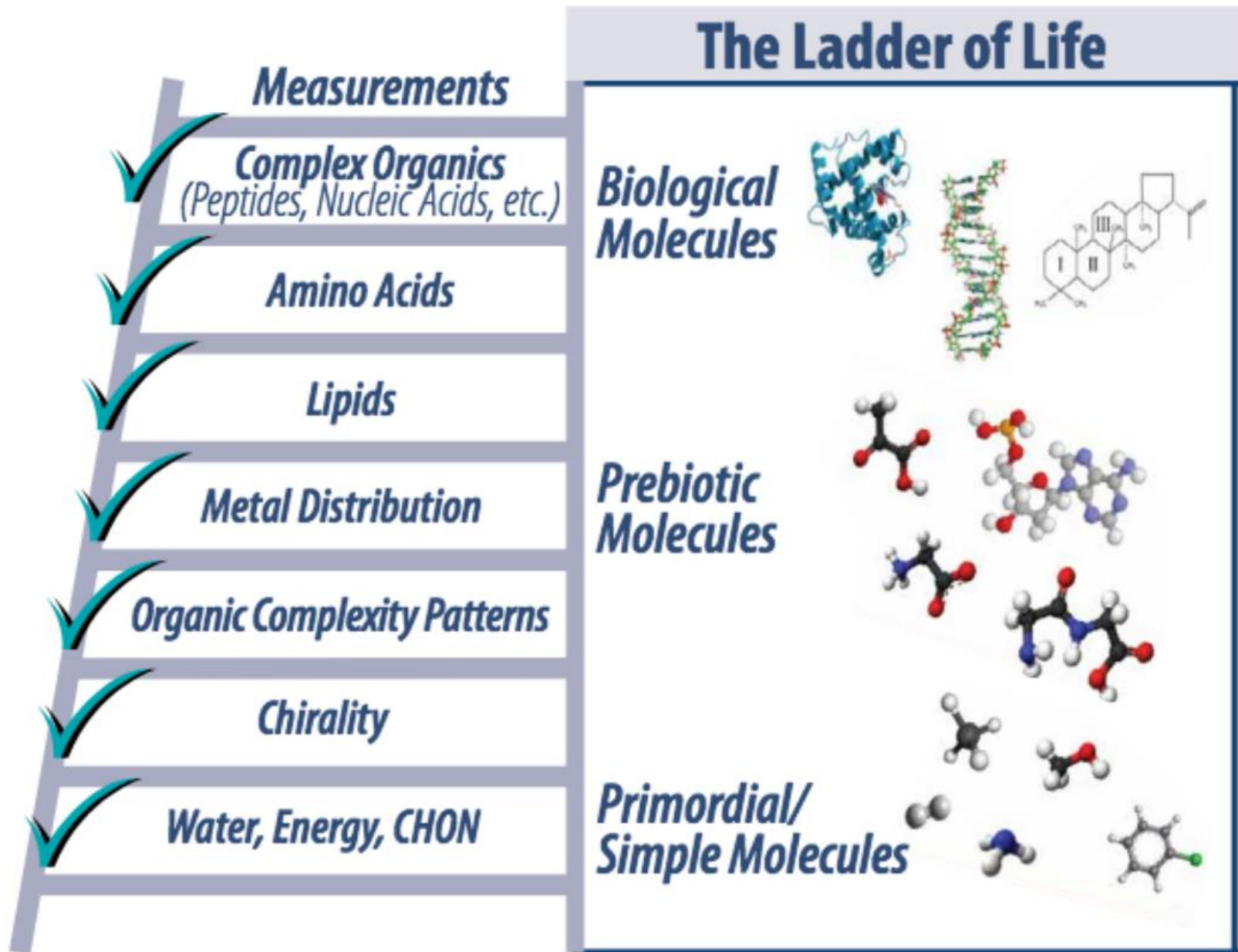
For 51 common metabolites ranging from acetate via amino acids and ribose to phosphoenolpyruvate, the molecular complexity C_m was calculated and plotted against the molecular mass of the compounds, demonstrating that molecular complexity is not directly dependent on molecular mass.

In fact, even small molecules with a comparable mass may differ by 100-200 mcbits.

Compound	C_m/M (mcbits•mol/g)	C_m (mcbits)
amino acids		
Alanine	0.817	72.00
Arginine	0.740	129.65
Asparagine	0.937	123.82
Aspartate	0.744	98.34
Cysteine	0.720	87.17
Glutamine	0.888	129.82
Glutamate	0.714	104.34
Glycine	0.519	38.98
Histidine	1.017	157.83
Isoleucine	0.774	101.51
Leucine	0.692	90.76
Lysine	0.711	104.65
Methionine	0.713	106.34
Phenylalanine	0.617	101.93
Proline	0.905	104.16
Serine	0.829	87.17
Threonine	1.009	120.19
Tryptophane	0.883	180.31
Tyrosine	0.660	119.51
Valine	0.724	84.76
fatty acids		
Hexanoic acid	0.416	48.34
Octanoic acid	0.418	60.34
Decanoic acid	0.420	72.34
Dodecanoic acid	0.421	84.34
Sugars		
Ribose	1.648	221.10
Galactose	1.484	267.29
Glucosamine	1.489	266.76

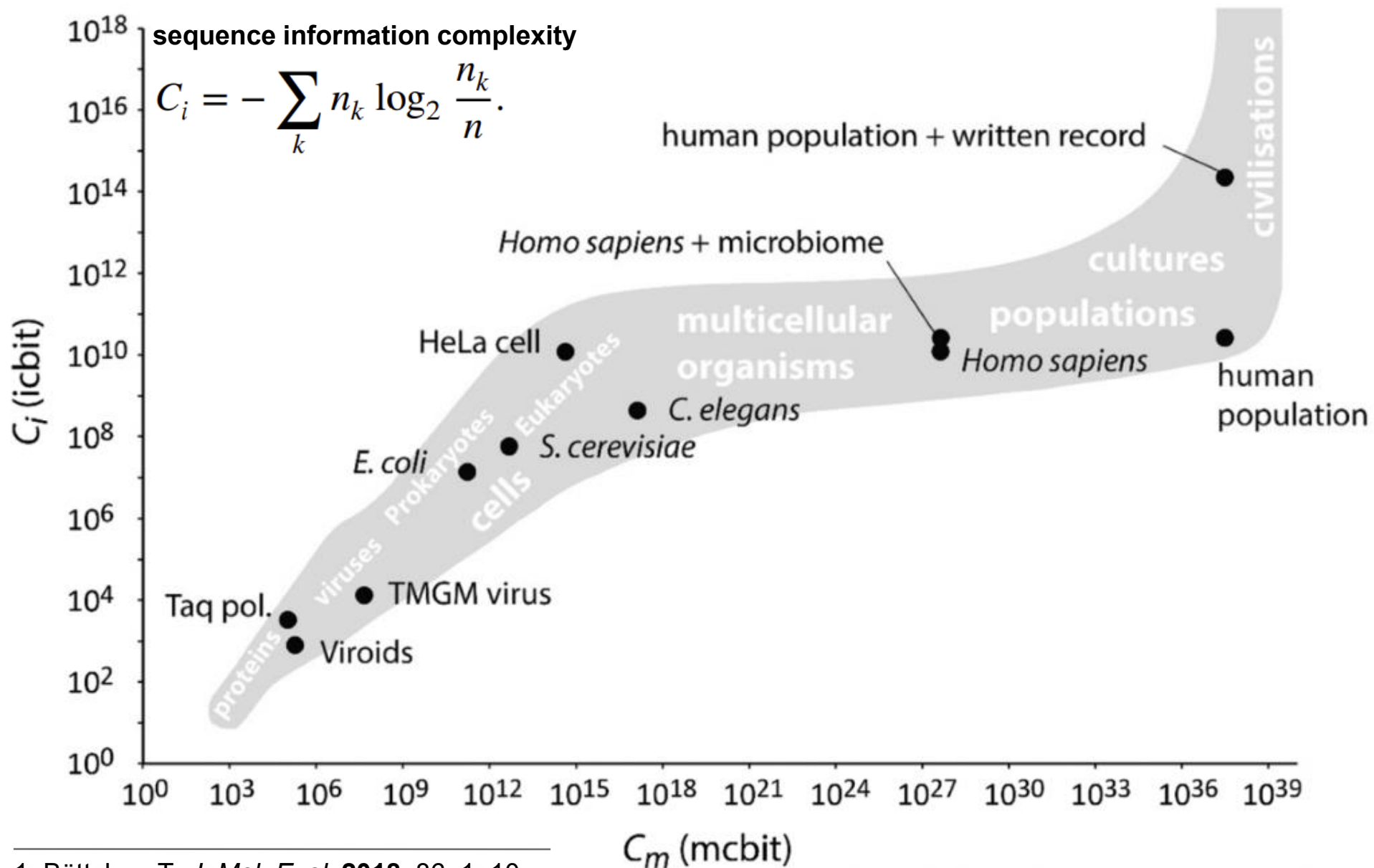
C_m Plotted against Molecular Mass





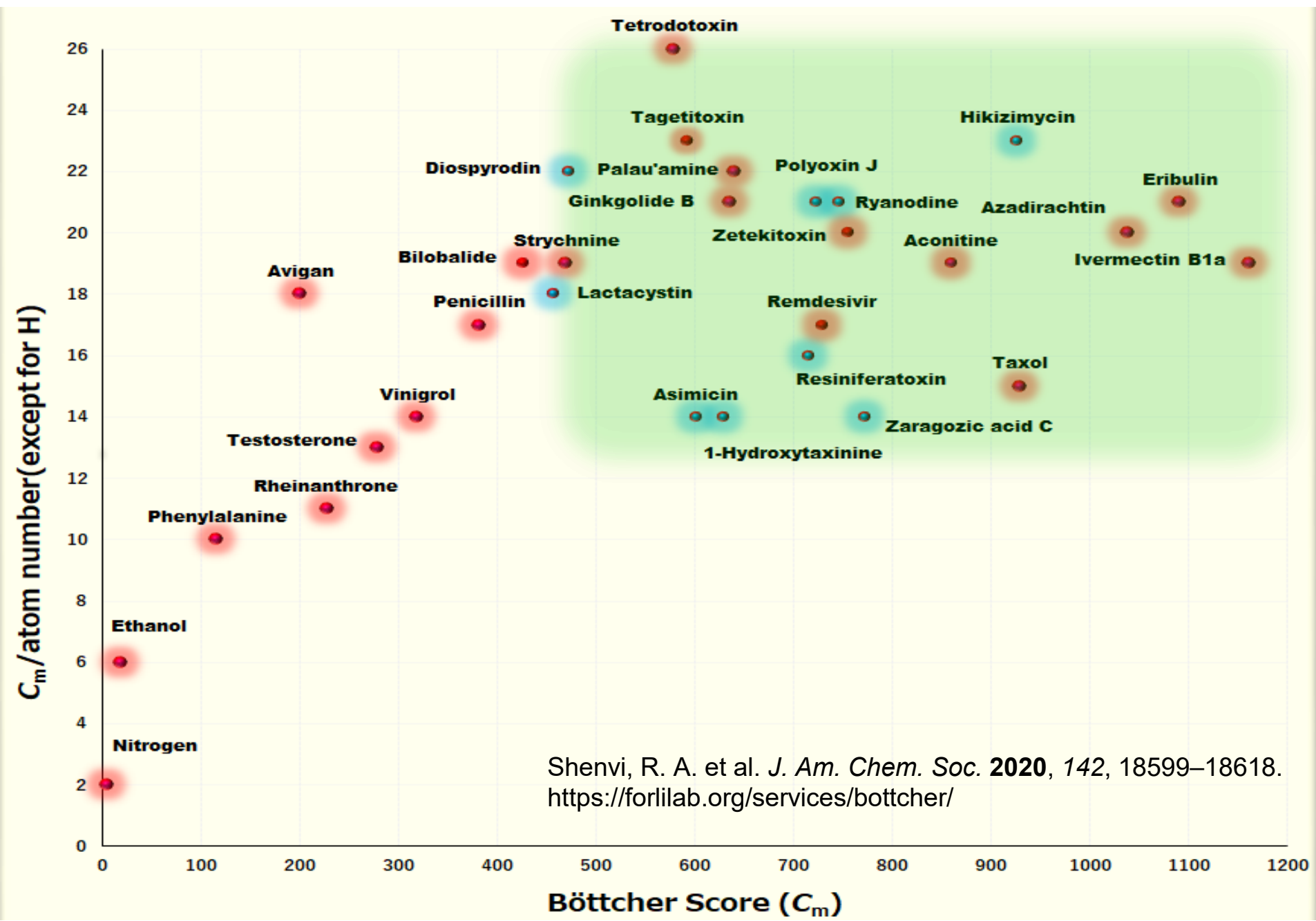
Universal Complexity Scale Plot with Representative Biogenic Units of Earth's Biosphere

16



1. Böttcher, T. *J. Mol. Evol.* **2018**, 86, 1–10.

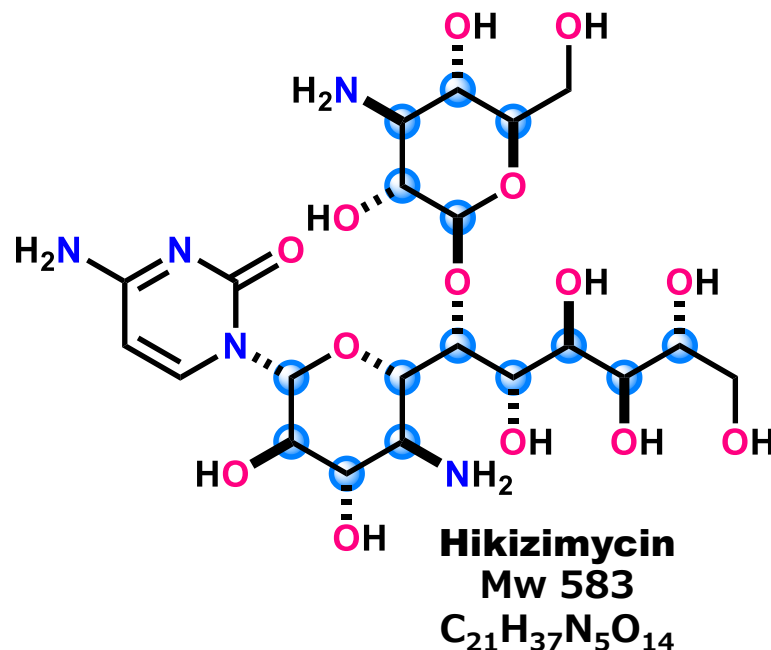
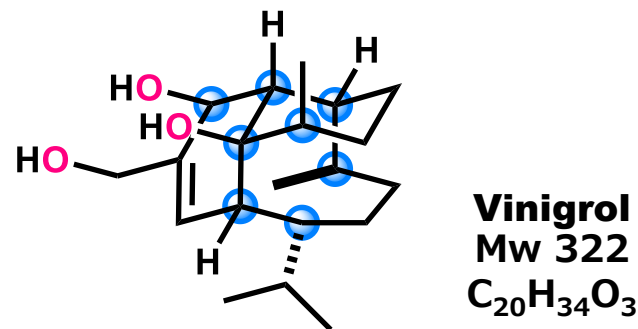
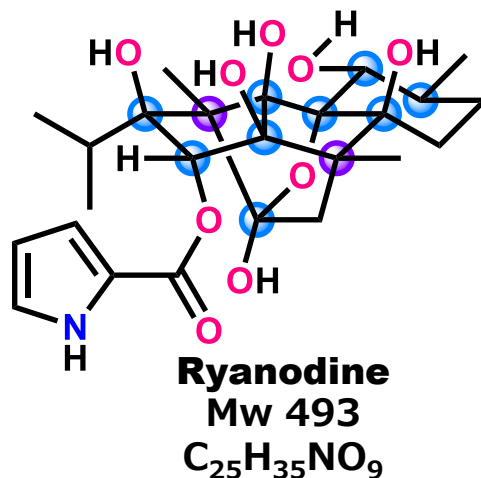
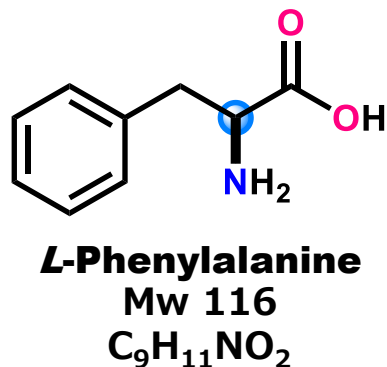
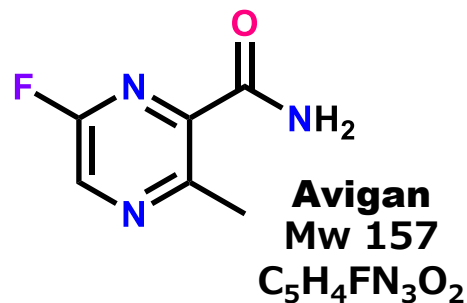
$C_m \propto$ Synthesizability?



How Do We Chemist Define Molecular Complexity? ¹⁸

In general, the following features of a target structure increase the challenge of chemical synthesis:

- (i) the number and **density of functional groups**
- (ii) the number of **stereocenters**
- (iii) the number and **types of rings**
- (iv) the overall **size of the molecule**.



1. a) Peterson, E. A.; Overman, L. E. *Proc. Natl. Acad. Sci. USA*. **2004**, *101*, 11943-11948.

b) Urabe, D.; Asaba, T.; Inoue, M. *Chem. Rev.* **2015**, *115*, 9207-9231.

Martin D. Eastgate

1977: Born in England

1999: B.S in Chemistry @ the University of Surrey, UK
Graduating with first class honors.

2002: Ph.D. in Organic Chemistry @ the University of Cambridge, UK (Prof. Stuart Warren)

sulfur participation chemistry, specifically the generation of thiiranium ions under basic conditions and their use in pyrrolidine synthesis.

2002-2005: Post-doctoral fellow @ the University of Illinois Urbana-Champaign (Prof. Scott E. Denmark)

the Lewis-base activation of Lewis-acids and understanding ligand-field theory in hyper-valent silyl cations.

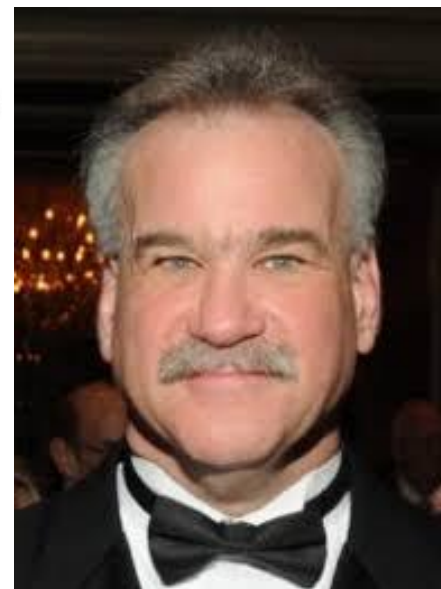
2005: Bristol-Myers Squibb

Currently: a Director in Chemical and Synthetic Development.



The late Dr. Stuart Warren

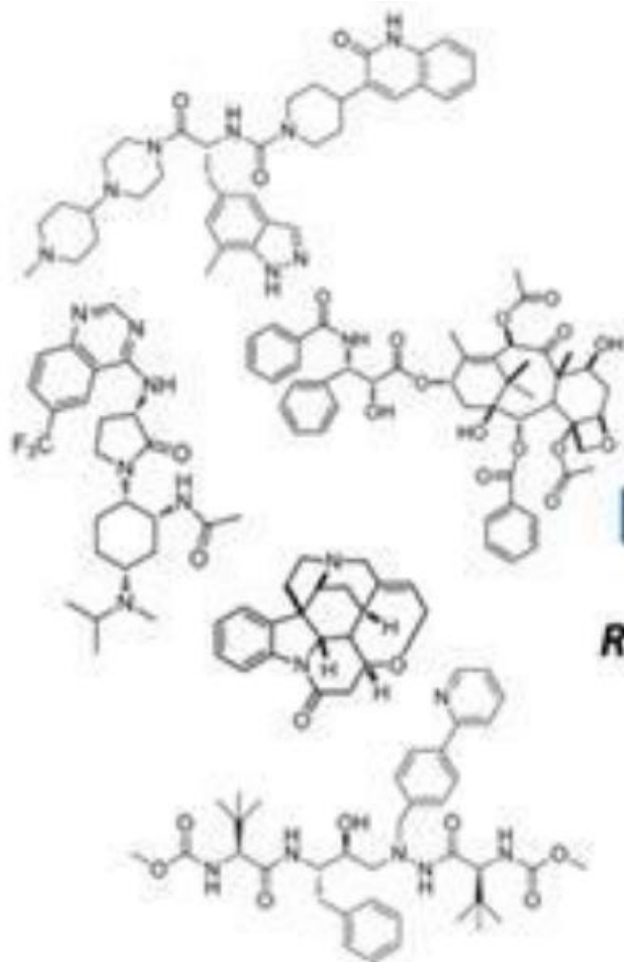
Prof. Scott E. Denmark



“Current” Complexity

Complexity postulate

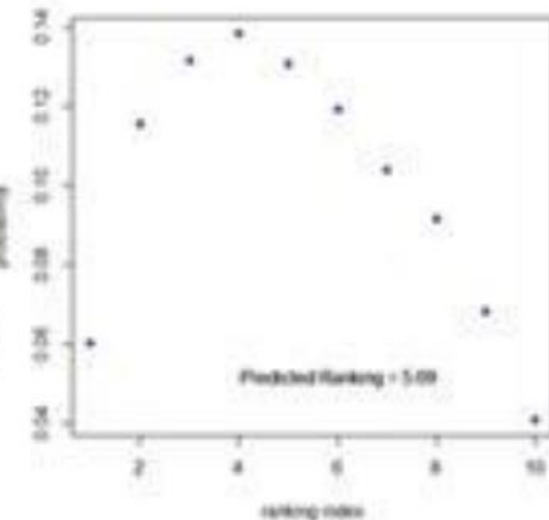
$$\text{Current Complexity} \propto \int_{\text{Fixed}} \text{Intrinsic Complexity} + \int_{\text{Today}} \text{Extrinsic Complexity}$$



Ranking

Collective Intelligence

Prediction



Current Complexity Index

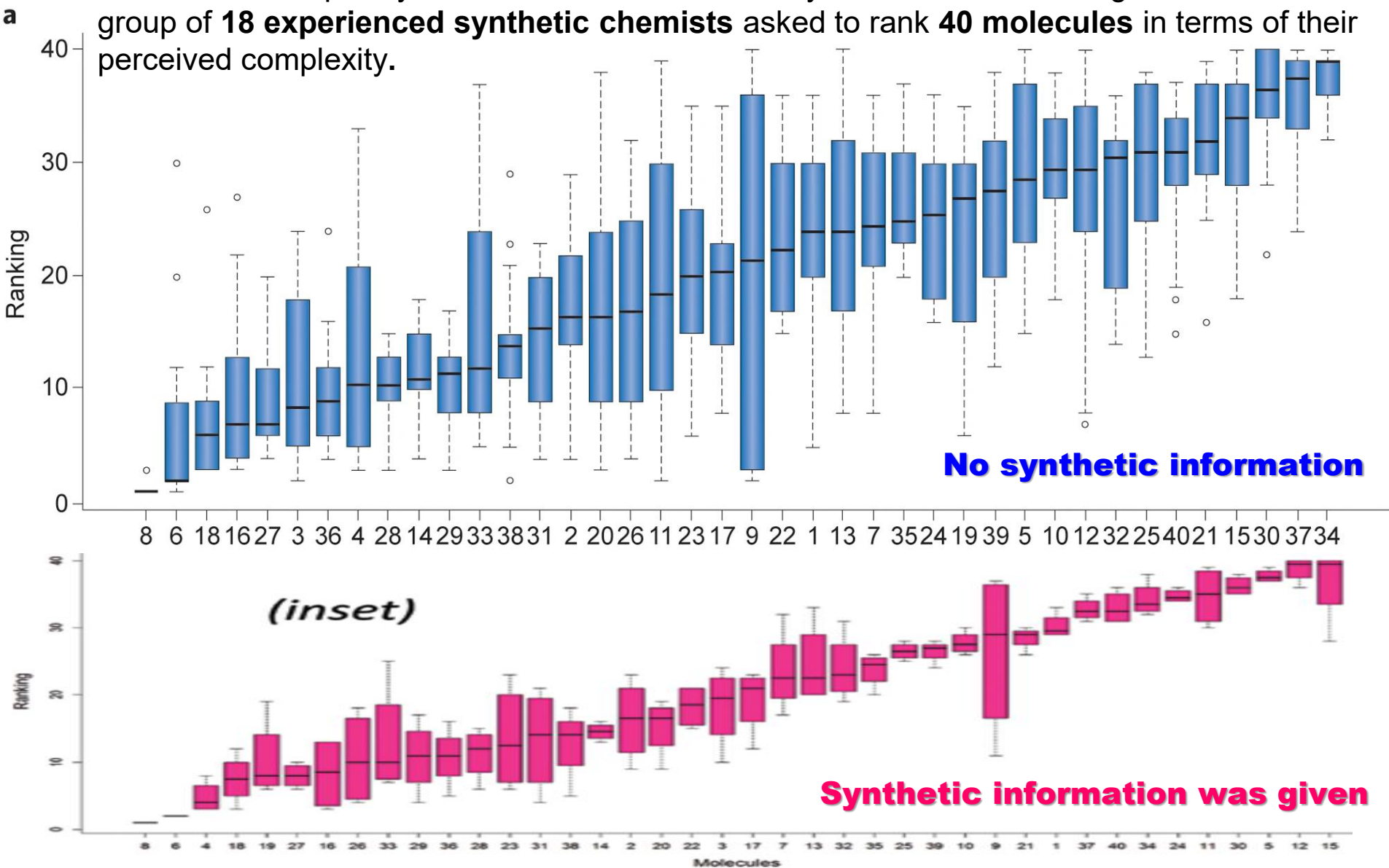
The index is based on a community's perception of complexity, within the context of current technology

1. a) Li, J.; Eastgate, M. D. *Org. Biomol. Chem.* **2015**, 13, 7164-7176.

b) Woolford, J. *Chemistry World* 22, May, 2015.

Analysis of Collective Intelligence

The current complexity index was based on an analysis of collective intelligence from a group of **18 experienced synthetic chemists** asked to rank **40 molecules** in terms of their perceived complexity.



Refinement by Intrinsic and Extrinsic Factors

The data obtained from the chemist's intuition was then refined by considering a large series of **intrinsic** and **extrinsic** factors and applying a **Bayesian regression model** to determine the five major factors that impacted the complexity of a structure the most. These are as follows:

- (i) the structure's molecular topological (Randic) index²
- (ii) the number of stereogenic centres established in the synthesis
- (iii) the number of heteroatoms on and in aromatic rings
- (iv) the number of steps
- (v) ideality of the route (as defined by P. S. Baran)³

$$\text{ideality} = \frac{(\text{numbers_of_construction_rxns}) + (\text{numbers_of_strategic_redox_rxns})}{\text{total_numbers_of_steps}}$$

(i) and (iii) are **intrinsic** and unchangeable, whereas the others are **extrinsic** variables reflecting advances that occur over time. From this was established an easily comprehensible 1–10 rating scale, with 1 being the most complex and 10 being least complex.

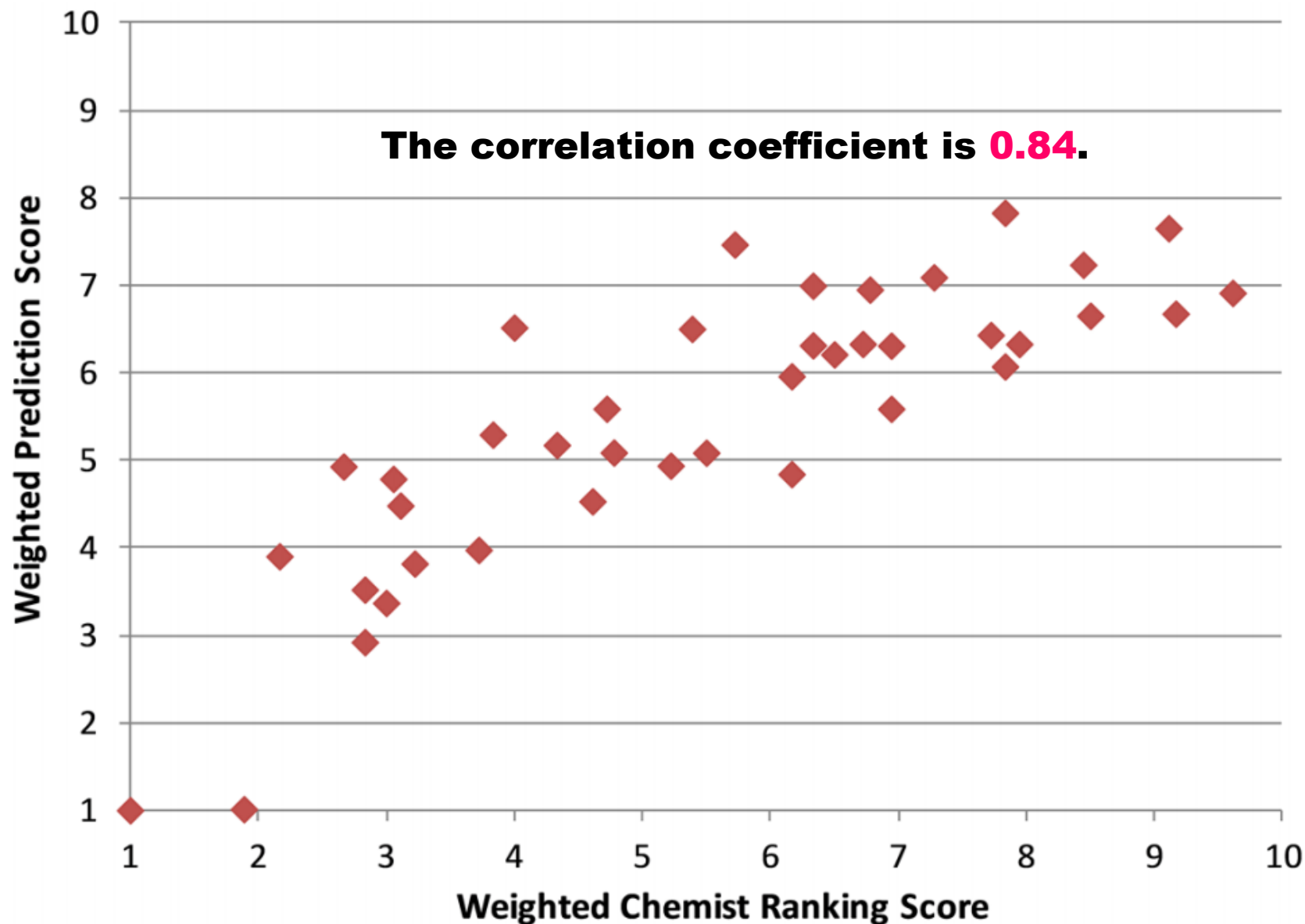
Regression model used in the current complexity index

$$\mu = \beta_0 + \beta_{\text{Randic}}^i x_{\text{Randic}} + \beta_{\text{SS}}^e x_{\text{SS}} + \beta_{\text{HAA}}^i x_{\text{HAA}} + \beta_{\text{Steps}}^e x_{\text{Steps}} + \beta_{\text{Ideality}}^e x_{\text{Ideality}} + \varepsilon$$

Latent response factor (μ) proportional to five weighted factor coefficients (β). Randic = Randic topology index; SS = number of stereocenters made; HAA = heteroatoms in or on aromatic rings; Steps = longest-linear + 50% of the branching steps; Ideality = ideality score. **Intrinsic (i)** and **extrinsic (e)** factors.

1. Li, J.; Eastgate, M. D. *Org. Biomol. Chem.* **2015**, 13, 7164-7176.
2. Randic, M. *J. Am. Chem. Soc.* **1975**, 97, 6609-6615.
3. Gaich, T.; Baran, P. S. *J. Org. Chem.* **2010**, 75, 4657-4673.

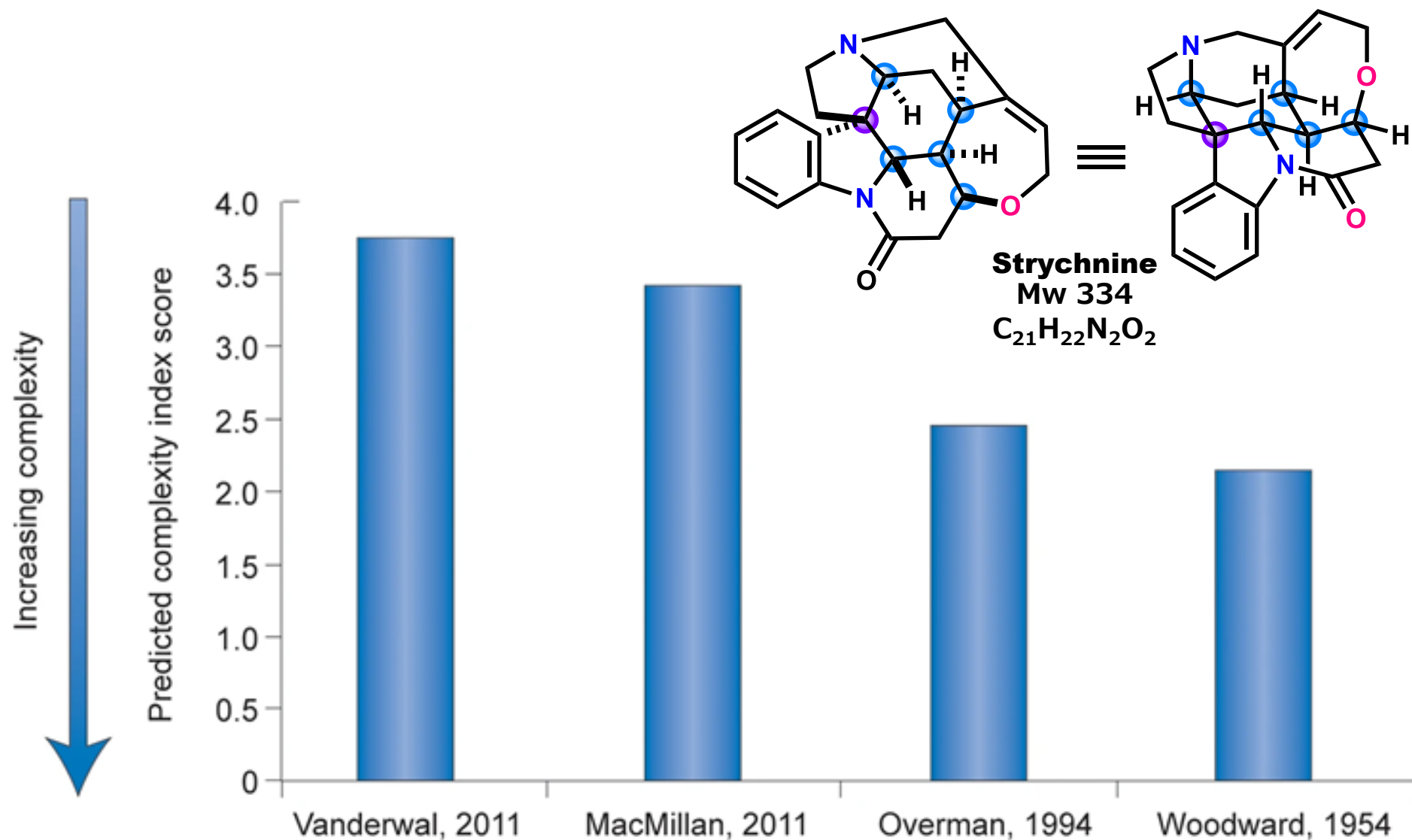
Correlation of Weighted Predict Complexity Scores²³ with Chemist Ranking Scores



Input Parameters for Selected Molecules and Weighted Predicted Complexity Score

	Randic	Steps	Ideality	Chiral_made	HAA	predicted complexity score	
Topiramate	10.99264	3	0.67	0	0	7.83	
Sildenafil	15.75559	8	0.89	0	7	6.42	
Atazanavir	25.40878	8.5	0.6	0	1	6.27	
BMS-2	23.16366	11	0.56	1	4	5.16	
BMS-1	18.35124	16	0.56	3	3	4.04	
BILN2061	27.91184	19.5	0.54	3	6	2.88	
EpoA	18.26984	27	0.47	5	2	2.38	Danishefsky
Strychnine	13.72358	8.5	0.60	6	1	3.75	Vanderwal
Strychnine	13.72358	13	0.69	6	1	3.43	MacMillan
Strychnine	13.72358	25	0.56	6	1	2.45	Overman
Strychnine	13.72358	30	0.53	6	1	2.14	Woodward
Discodermolide	23.87016	30	0.59	10	0	1.25	Novartis-Smith-Paterson
Taxol	31.71272	37	0.48	11	0	1.06	Nicolaou
Halaven	31.3614	56	0.38	13	0	1.00	Eisai
Welwit C isonitrile	13.3244	21	0.33	4	1	3.42	Rawal
Welwit C isonitrile	13.3244	23	0.39	4	1	3.30	Garg
Welwit A isonitrile	12.45351	23	0.43	4	1	3.38	Wood
Welwit A isonitrile	12.45351	9	0.78	4	1	4.98	Baran

Predictive Complexity Index Scores for Some of the Strychnine Syntheses



1. a) Li, J.; Eastgate, M. D. *Org. Biomol. Chem.* **2015**, 13, 7164-7176.

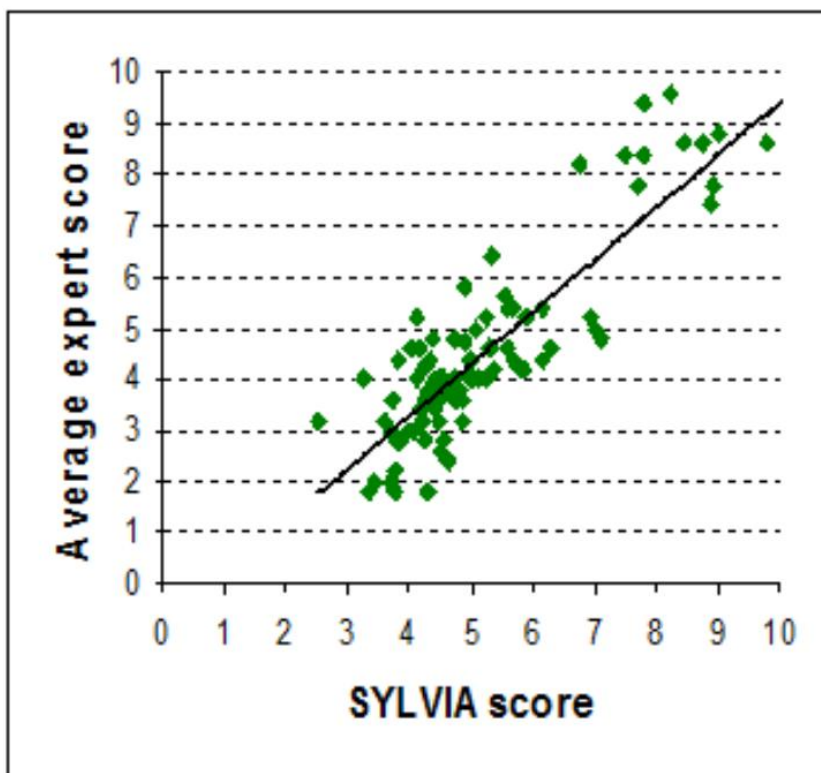
b) Woolford, J. *Chemistry World* 22, May, 2015.

SYLVIA

Similarly, **Gasteiger** has defined the '**synthetic accessibility**' of a compound; an estimate to reflect how easily a molecule can be synthesized, based on an analysis of the molecular structure and a comparison with the contents of an organic reaction database.

To adjust the synthetic accessibility estimates, **five chemists representing three different pharmaceutical companies** were asked to rank 100 molecules on a ten-point scale.

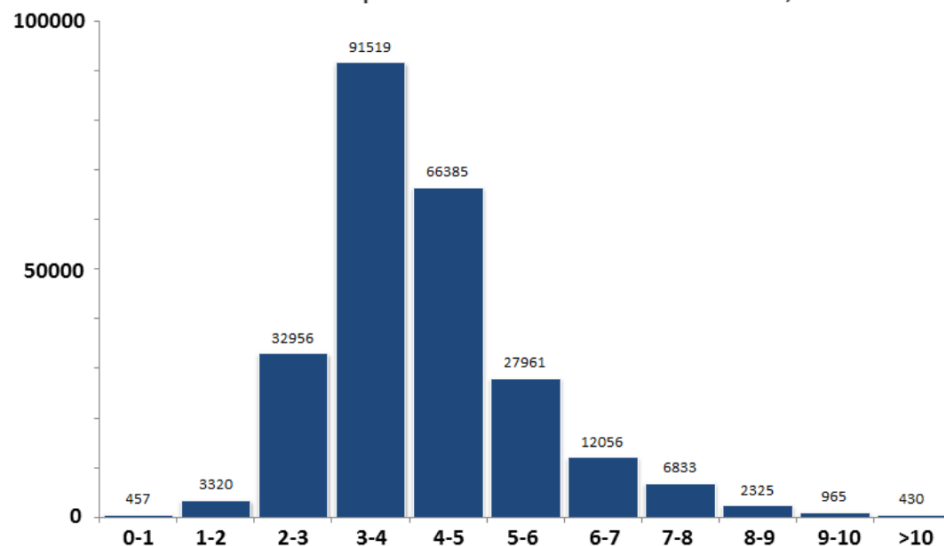
An accessibility tool called **SYLVIA** that was developed based on these studies is freely available (<http://www.molecular-networks.com/products/sylvia>).



Correlation: 0.89



► Distribution of scores in Open NCI Database with 245,208 records



Summary

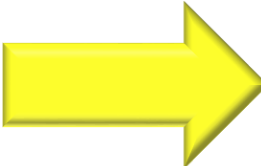


Johann Gasteiger, an expert in cheminformatics at the University of Erlangen-Nürnberg, Germany, says:

The idea of quantifying or assigning the complexity of a molecular structure a number has been around for some time without a suitable solution. In other words, **no system has found broad acceptance among the organic community yet.**

The reasons for this are manifold, but not least because of community resistance. Many organic chemists consider synthesis design as an “art” where computers should not have a place.

It seems that organic and process chemists have finally started to recognize the value of computer tools in their work on organic synthesis.



In my opinion, complexity should be defined by the intrinsic factor. The Böttcher Score is considerably better than previous indices but seems to overestimate the degree of complexity reduction due to molecular symmetry. It also seems to lack consideration of fused ring structures (e.g., transannular interaction).

Appendix

Examples: Various Simple Aliphatic Hydrocarbons ³¹

Compound	C_m (mcbit)	No. distinct carbons
Ethane	2.00	1
<i>n</i> -Propane	5.00	2
<i>n</i> -Butane	8.00	2
2-Methylpropane	6.58	2
2-Methylbutane	17.17	4
<i>n</i> -Pentane	11.00	3
2,2-Dimethylpropane	8.00	2
2-Methylpentane	23.17	5
2,2-Dimethylbutane	19.00	4
3-Methylpentane	17.17	4
<i>n</i> -Hexane	14.00	3
2,3-Dimethylbutane	11.17	2
2,3-Dimethylpentane	46.26	7*
3-Methylhexane	45.51	7*
2-Methylhexane	29.17	6
2,2-Dimethylpentane	25.00	5
2,2,3-Trimethylpentane	20.17	4
3,3-Dimethylpentane	18.00	4
<i>n</i> -Heptane	17.00	4
3-Ethylpentane	15.58	3
2,4-Dimethylpentane	14.17	3

*Chiral molecule

Molecular complexity C_m for various simple aliphatic hydrocarbons from ethane to heptane and their isomers correlates with the number of chemically nonequivalent (distinct) carbon atoms.

Comparison of C_m with Other Complexity Indices

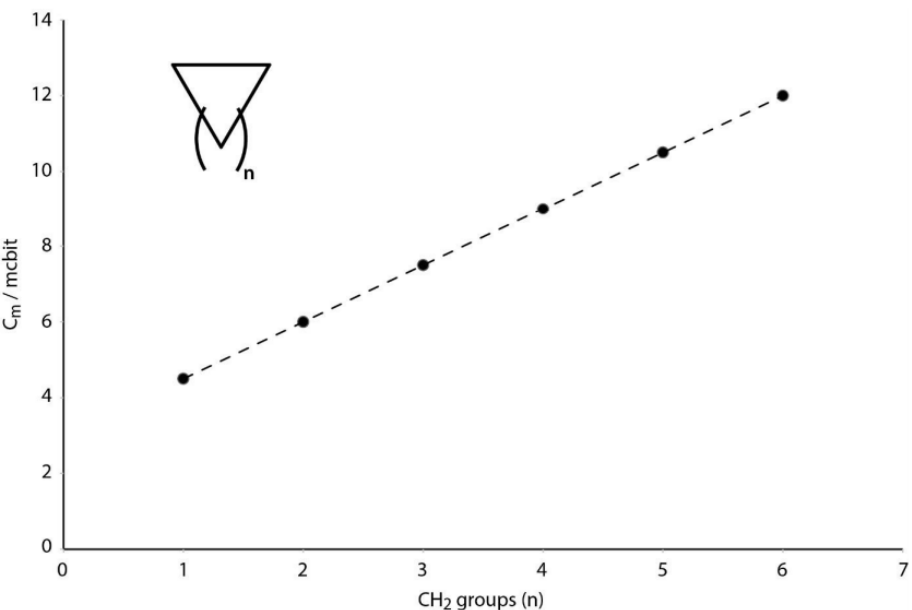
32

Compound	C_m (mcbrit)	$C(\eta, \epsilon)$	N_T	N_S
Methane	_*	_*	1	1
Ethane	2.00	_*	3	2
n-Propane	5.00	0.0	6	3
n-Butane	8.00	2.0	10	4
n-Pentane	11.00	7.5	15	5
n-Hexane	14.00	12.0	21	6
Ethylene	3.00	0.0	5	3
Propene	12.17	7.5	10	5
1-Butene	18.17	14.0	16	7
Methanol	9.17	2.0	3	3
Ethanol	19.17	2.8	6	5
1-Propanol	25.17	7.2	10	7
2-Propanol	21.51	10.8	11	7
Acetone	25.17	26.3	19	10
Acetaldehyde	23.51	10.3	10	7
Propanal	29.51	17.2	16	10
Dimethyl ether	11.17	2.8	6	4
Acetylene	3.58	4.8	9	4

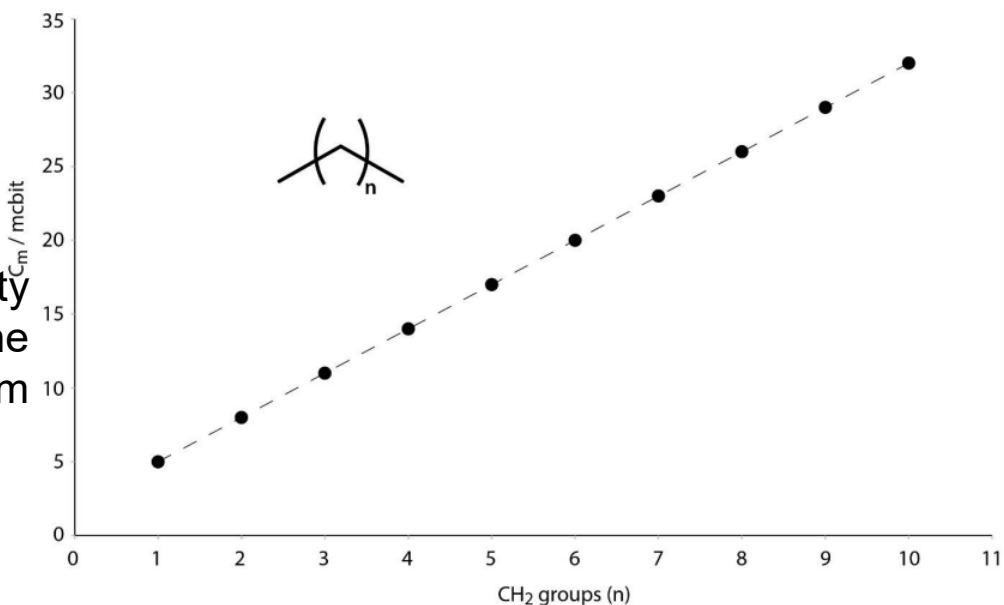
*not defined for these molecules

1. Böttcher, T. *J. Chem. Inf. Model.* **2016**, 56, 462–470.

Example: Linear Dependence of C_m

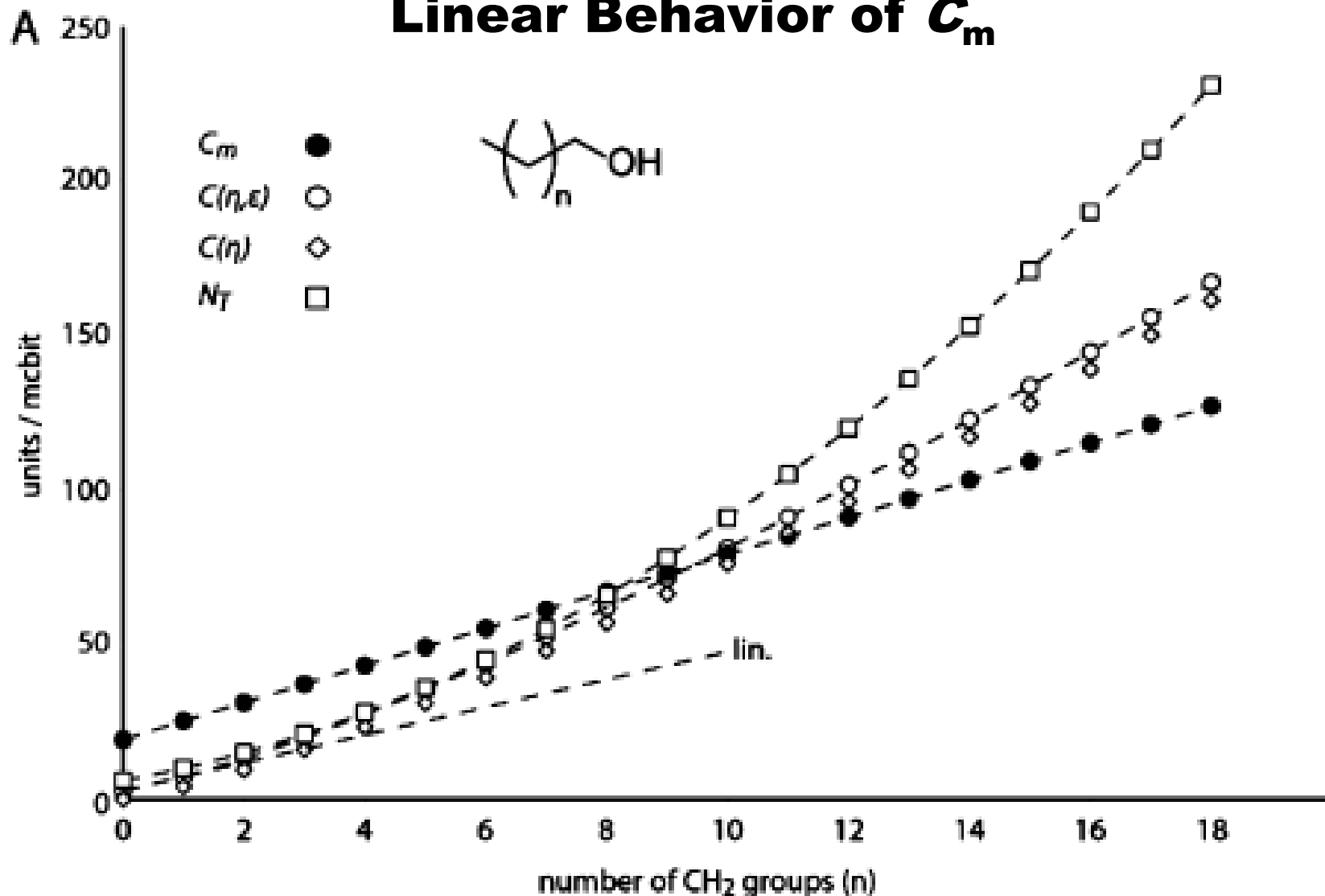


The linear dependence of molecular complexity C_m from the number of methylene groups in the homologous series of cycloalkanes from cyclopropane ($n = 1$) to cyclooctane ($n = 6$).



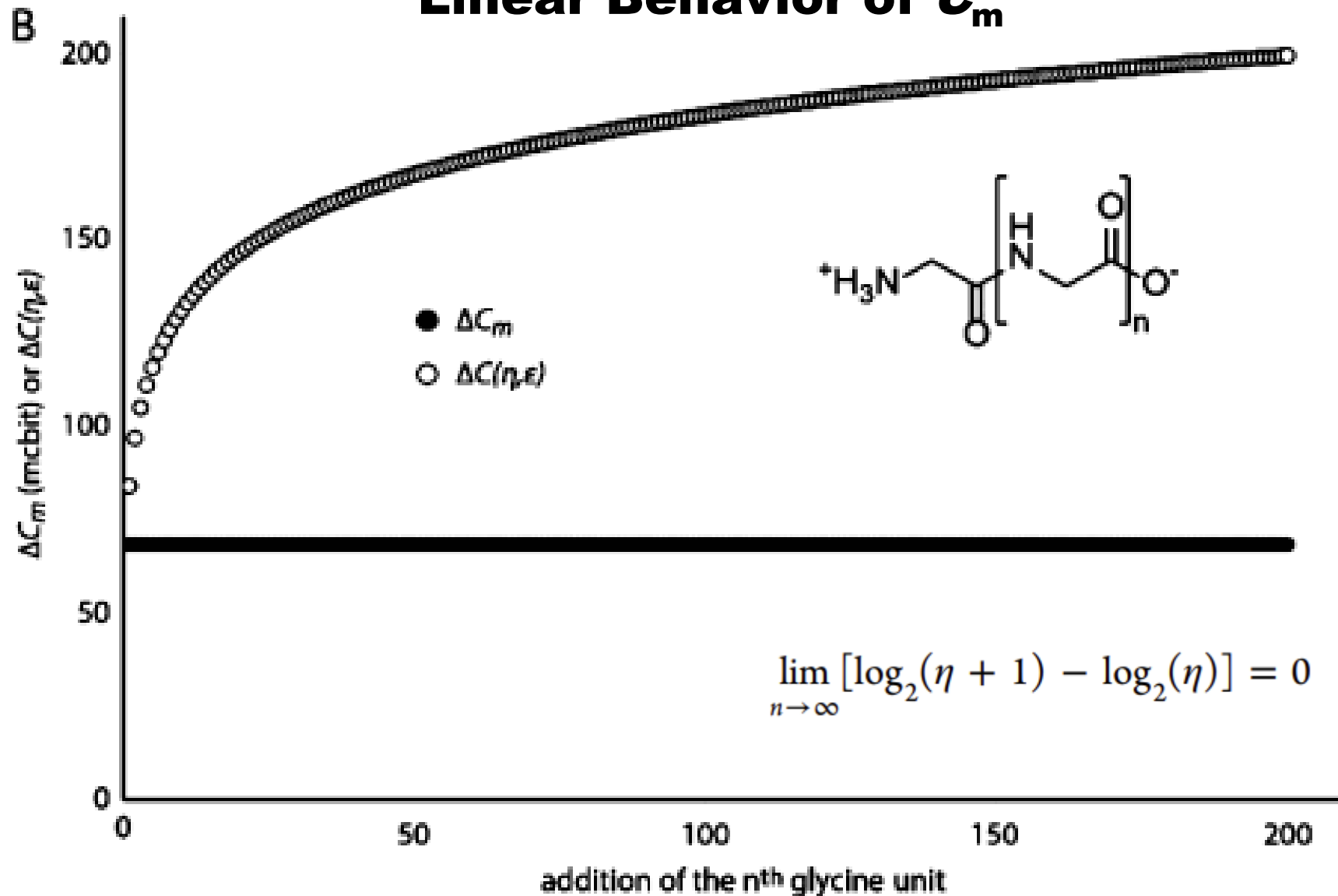
The linear dependence of molecular complexity C_m from the number of methylene groups in the homologous series of alkanes from propane ($n = 1$) to dodecane ($n = 10$).

Nonlinear Behavior of $C_{(\eta,\epsilon)}$ in Comparison to the Linear Behavior of C_m



Homologous series from ethanol to $C_{20}H_{41}OH$, where C_m gives a linear increase and $C(\eta,\epsilon)$, $C(\eta)$, and N_T give nonlinear increases. A linear extrapolation (lin.) from the first two values of $C(\eta,\epsilon)$ is given for better visualization.

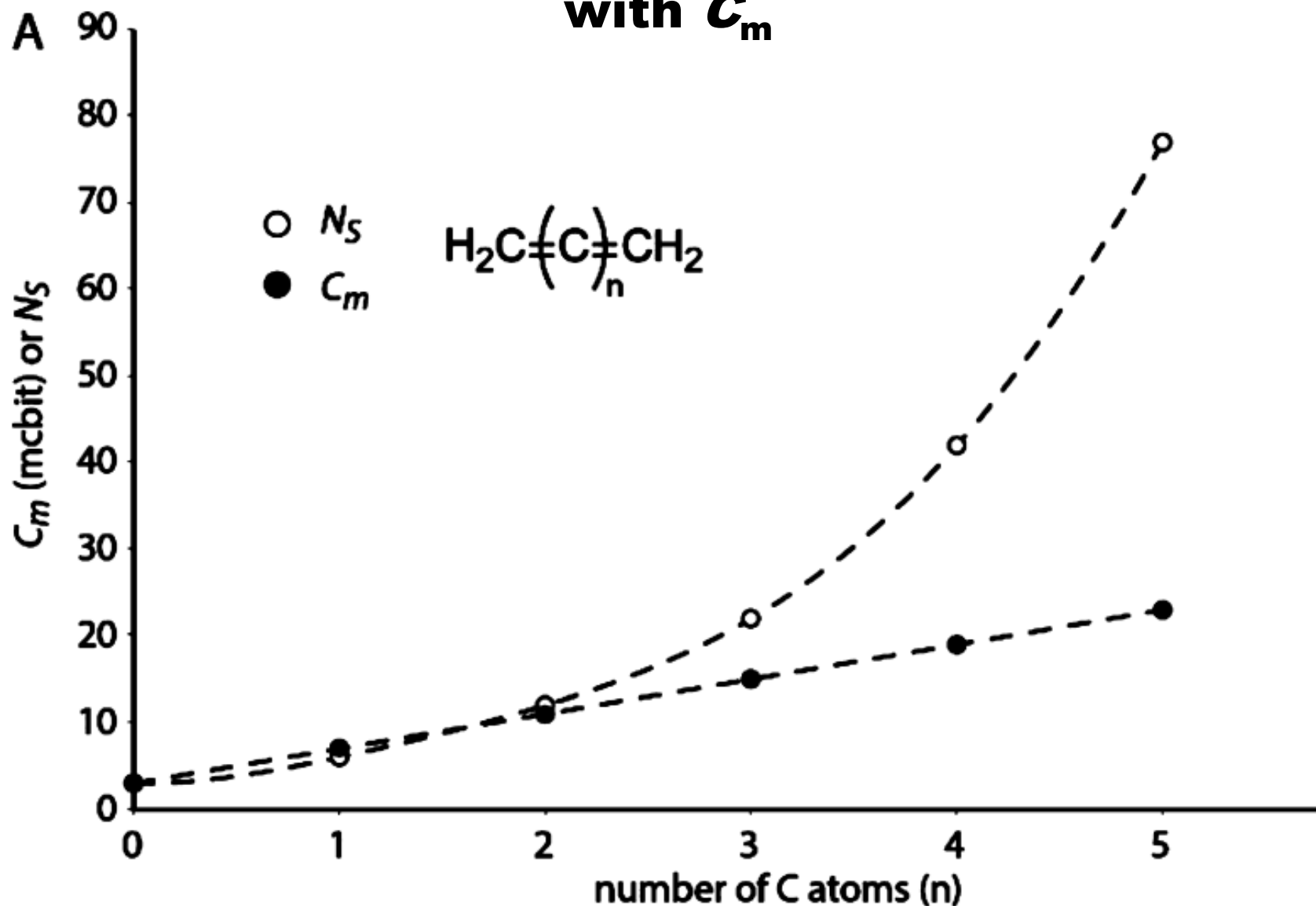
Nonlinear Behavior of $C_{(\eta,\epsilon)}$ in Comparison to the Linear Behavior of C_m



Changes in complexity $\Delta C_{(\eta,\epsilon)}$ and ΔC_m for the addition of the n^{th} glycine residue to a polyglycine chain, showing the nonlinear behavior of $C_{(\eta,\epsilon)}$.

Comparison of the Graph-Theory-Based Index N_S with C_m

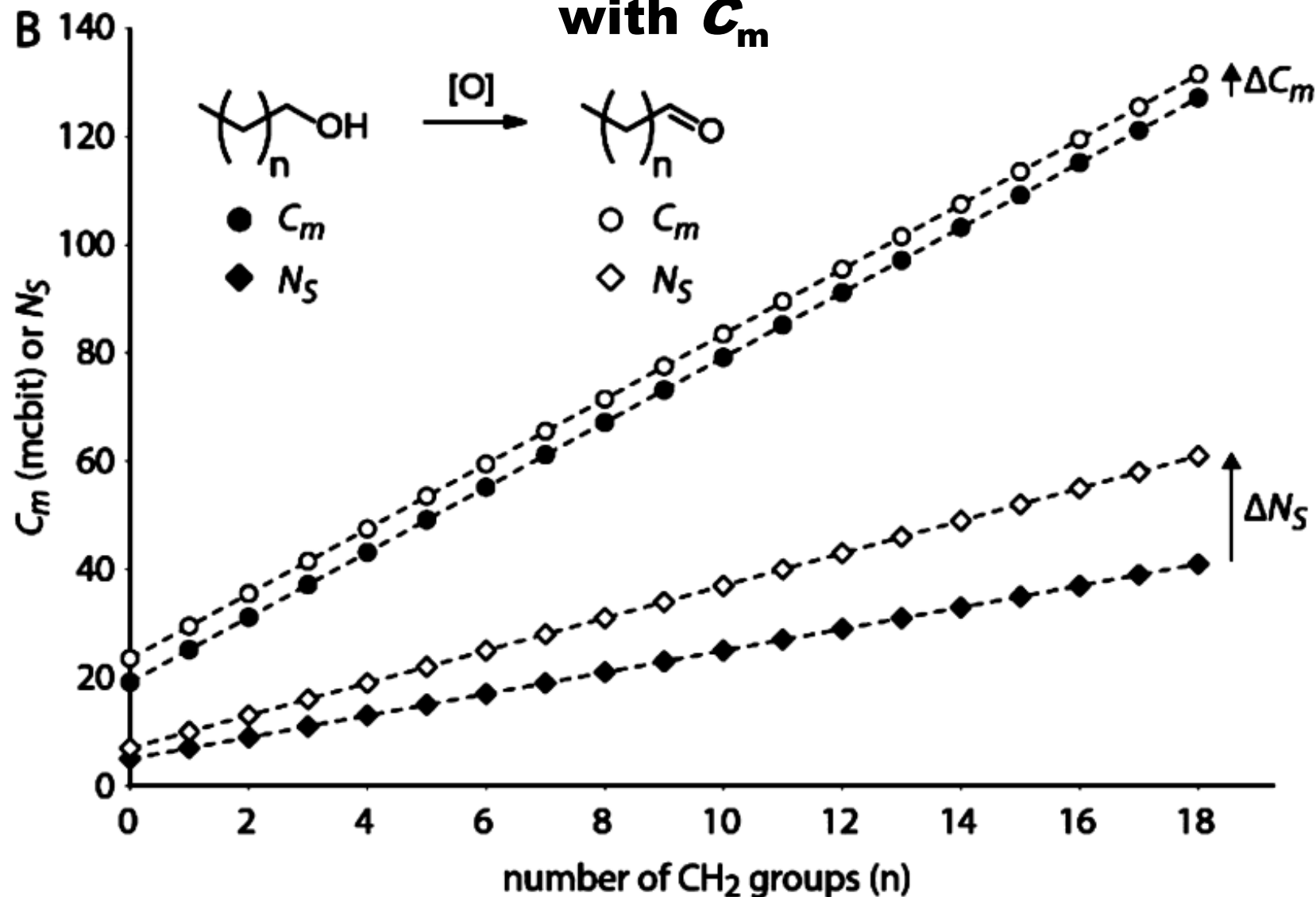
36



Linear and nonlinear increases within the homologous series of cumulenes for C_m and N_S , respectively.

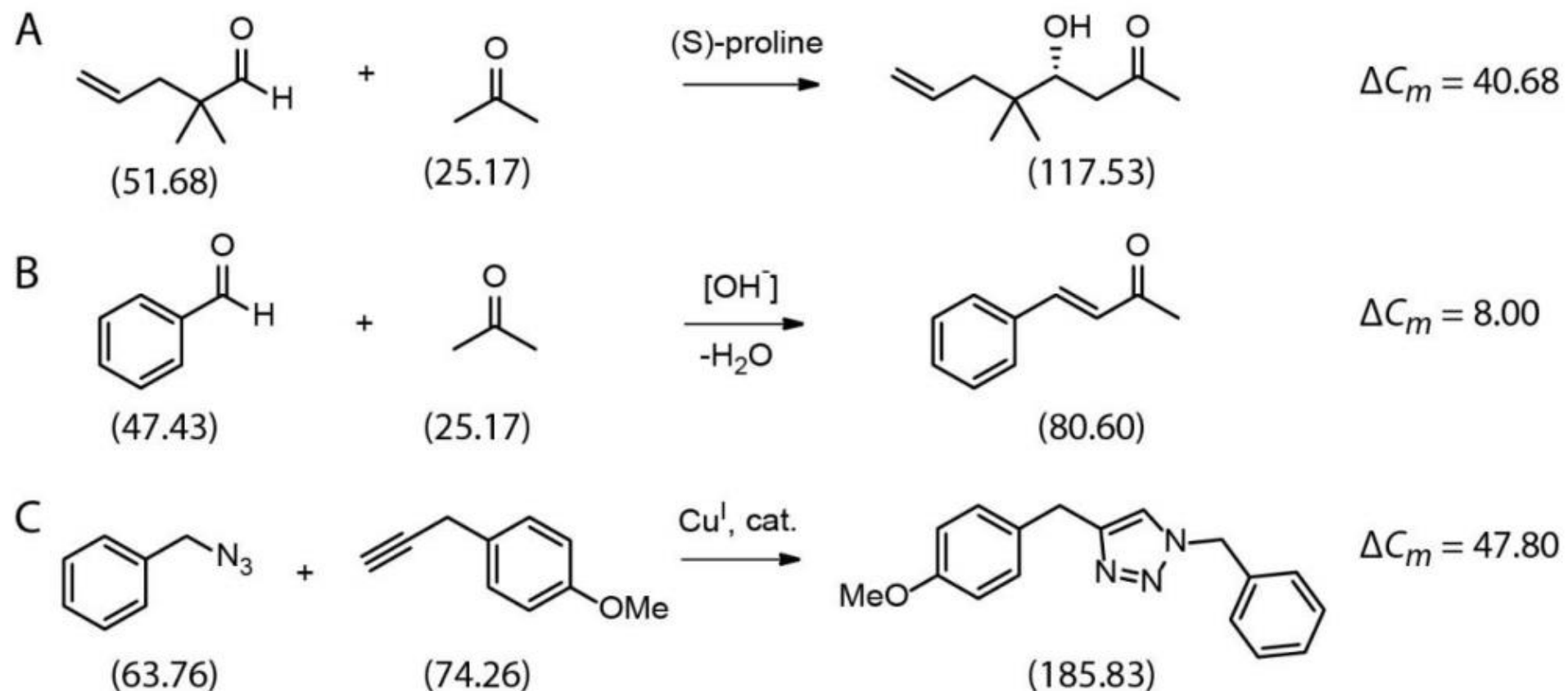
Comparison of the Graph-Theory-Based Index N_S with C_m

37



Invariance of ΔC_m to the chain length of aliphatic alcohols in the oxidation reaction to give the corresponding aldehydes, in contrast to an increase in ΔN_S for the same reaction with increasing hydrocarbon chain length.

Changes in Molecular Complexity ΔC_m



Changes of molecular complexity ΔC_m for different chemical reactions.

A) Stereoselective aldol reaction

B) aldol condensation

C) 1,3-dipolar azide-alkyne cycloaddition.

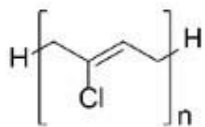
C_m values for all reactants are given in brackets.

Molecular Complexity for Various Artificial and Biological Polymers without Sequence Information

39

Artificial polymers

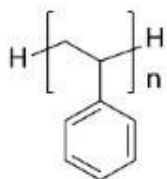
Neoprene



$$C_M = (56 \cdot n) - 8$$

$$C_I = 0$$

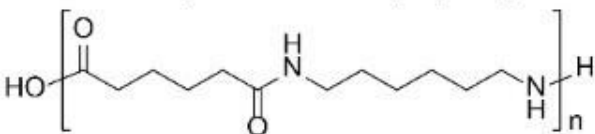
Polystyrene



$$C_M = (61 \cdot n) - 19.5$$

$$C_I = 0$$

Polyamide 6.6 (Nylon)

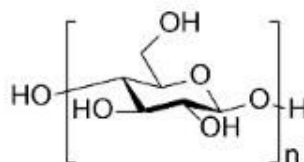


$$C_M = (172.9 \cdot n) - 15.5$$

$$C_I = 0$$

Biological polymers

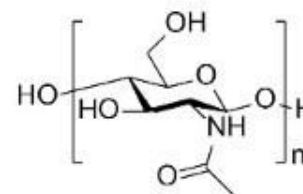
Cellulose



$$C_M = (271.3 \cdot n) - 4$$

$$C_I = 0$$

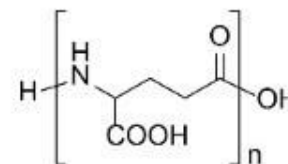
Chitin



$$C_M = (324.6 \cdot n) - 4$$

$$C_I = 0$$

Poly-γ-glutamic acid

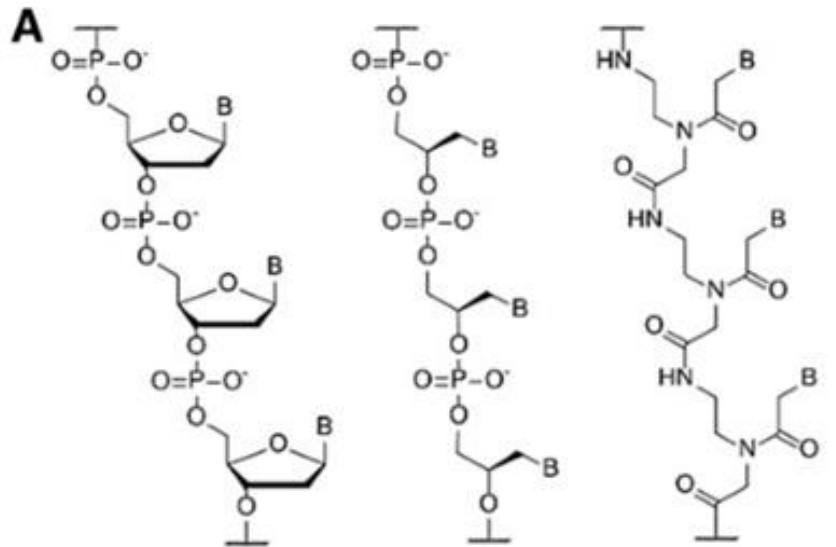


$$C_M = (138.1 \cdot n) - 21.5$$

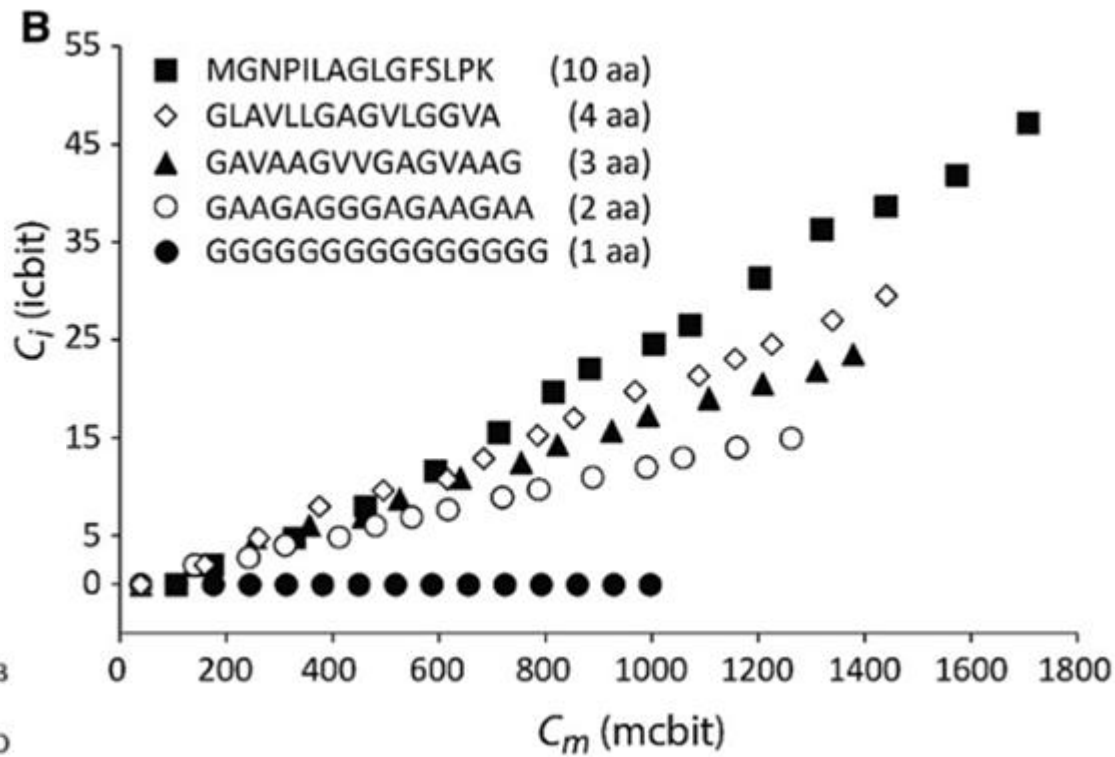
$$C_I = 0$$

Quantification of Molecular (C_m) and Information Complexity (C_i) for Various Types of Biogenic Units

A) Chemical structures of DNA, (R)-GNA, and PNA and calculated complexity values for an arbitrary model sequence (ATGTGA).

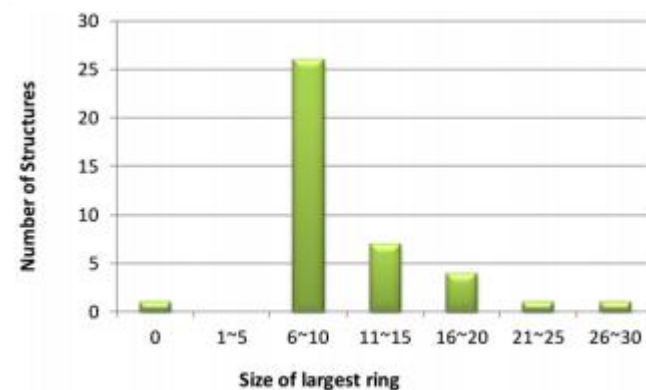
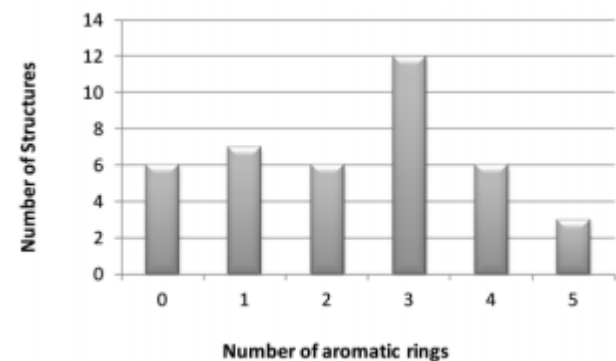
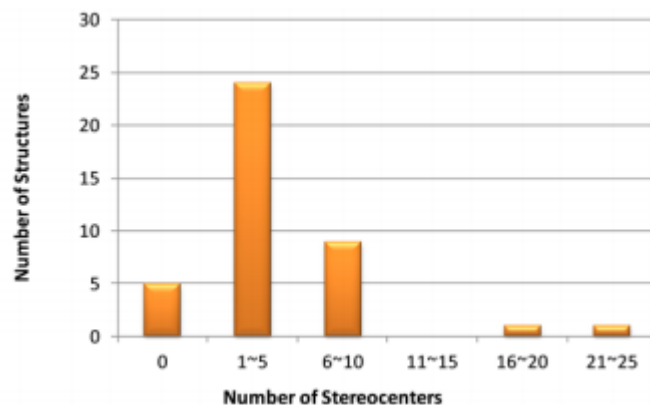
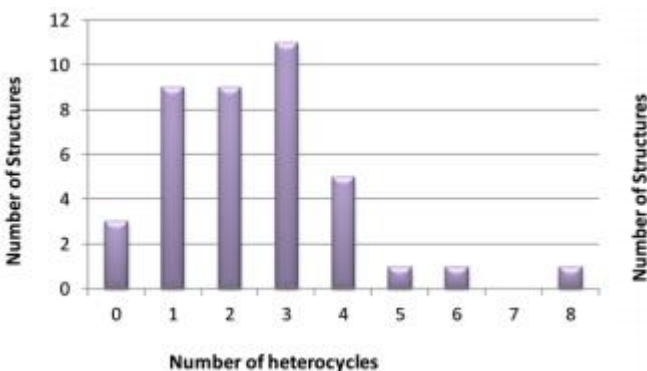
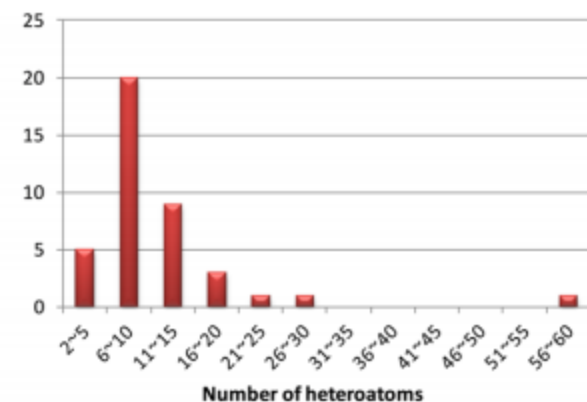
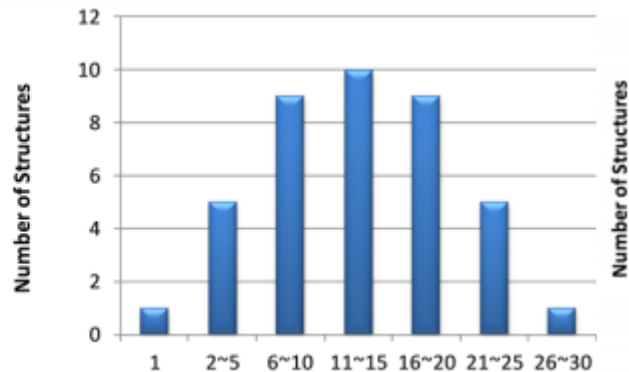
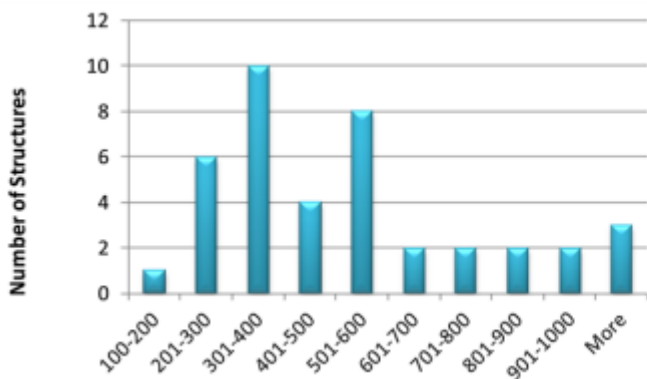


ATGTGA	DNA	(R)-GNA	PNA
C_m	2596.2	1831.1	1985.6
C_i	11.5	11.5	11.5



B) Information complexity plotted against molecular complexity for arbitrary protein sequences of different alphabet sizes (1 aa to 10 aa) as a function of length. B: nucleobase, aa: amino acid

Training Dataset



Distributions of molecular properties from MW, double bond equivalent (DBE), heteroatoms, heterocycles, aromatic rings, stereocenters, and the largest ring size in the training dataset.