The
University
Of
Sheffield.

# Chemoinformatics: the first half century

## Peter Willett

Presented at the Third Strasbourg Summer School on Chemoinformatics, 25th June 2012

# Overview

- Introduction to chemoinformatics
  - What it is
  - How it has developed
- Historically important papers
  - Roughly chronological ordering
  - A personal choice
  - Many, many omissions

# Chemoinformatics' role

- The pharmaceutical industry has been one of the great success stories of scientific research in the latter half of the twentieth century
  - Range of novel drugs for important therapeutic areas
- The computer has revolutionised how the industry uses chemical (and increasingly biological) information
  - Many of these developments are within the discipline we now know as chemoinformatics
  - Focus on lead discovery and lead optimisation phases of drug discovery (also applicable to other types of specialty chemicals)

# Definitions

- F.K. Brown (1998) Chemoinformatics: what it is and how does it impact drug discovery? *Annual Reports in Medicinal Chemistry*, **33**, 375-384
  - "The use of information technology and management has become a critical part of the drug discovery process. Chemoinformatics is the mixing of those information resources to transform data into information and information into knowledge for the intended purpose of making better decisions faster in the area of drug lead identification and optimization"

- G. Paris (August 1999 ACS meeting), quoted by W.A. Warr at http://www.warr.com/warrzone.htm
  - "Chem(o)informatics is a generic term that encompasses the design, creation, organization, management, retrieval, analysis, dissemination, visualization and use of chemical information"

- J. Gasteiger and T. Engels (editors) (2003). *Chemoinformatics: a textbook*. Wiley-VCH.
  - "Chemoinformatics is the application of informatics methods to solve chemical problems."

# Emergence of chemoinformatics: I

- The principal driving force for current interest was the emergence during the Nineties of combinatorial synthesis and high-throughput screening

- Aims to assist more directly in the discovery of novel bioactive molecules than heretofore

  - Integration of chemical information (archival functions) and molecular modelling (small-scale correlation)

  - Development of effective and efficient tools that can exploit the chemical and biological data explosion

- Other types of *-informatics* becoming common

- But chemoionformatics is by no means new...

# Emergence of chemoinformatics: II

- First appearance of the core, journal, *Journal of Chemical Documentation*, in 1961
- First book on the subject appeared in 1971
  - M.F. Lynch et al., *Computer Handling of Chemical Structure Information*
- The first two textbooks with "chemoinformatics" in the title appeared in 2003
  - A.R. Leach and V.J. Gillet, *An Introduction to Chemoinformatics*
  - J. Gasteiger and T. Engel (eds.) *Chemoinformatics*
- The first international conference on the subject was held at Noordwijkerhout in 1973, and every three years since 1987

# Emergence of chemoinformatics: III

- Introduction of first full university courses in 2001
  - D.J. Wild and G. Wiggins (2006) Challenges for chemoinformatics education in drug discovery. *Drug Drug Discovery Today*, **11**, 436-439

- Nomenclature
  - Chemical informatics, chemical information (management/science), chemiinformatics, cheminformatics
    P. Willett (2008) A bibliometric analysis of chemoinformatics. *Aslib Proceedings*, **60**, 4-17
  - M. Hann and R. Green (1999) Chemoinformatics - a new name for an old problem? *Current Opinion in Chemical Biology*, **3**, 379-383

# Basis in the literature

The chemical literature has been established for many years

- *Chemisches Journal* first appeared in 1778

- *Chemical Abstracts* first appeared in 1907

- Computerised systems first appeared in the early Sixties, shortly after the start of the core journal, the *Journal of Chemical Documentation*



The very first paper: N. Lozac'h (1961) *Journal of Chemical Documentation*, **1**, 1-4

# Development of the journal

- Early issues of the *Journal of Chemical Documentation* largely comprised papers presented at meetings of the ACS Division of Chemical Literature

- First issue of *Journal of Chemical Information and Computer Sciences* (1975)
  - Largely comprised papers given at a National Academy of Sciences conference on databases (both textual and chemical)

- First issue of *Journal of Chemical Information and Modeling* (2005) with sections on
  - Chemical information, Computational chemistry, Computational biochemistry, Pharmaceutical modeling, and Bioinformatics

# Many other journals now cover the subject, most obviously...

- *Molecular Informatics*
  - Started in 1982 as *Quantitative Structure-Activity Relationships*
- *Journal of Molecular Graphics and Modelling*
  - Started in 1983 as the *Journal of Molecular Graphics*
- *Journal of Computer-Aided Molecular Design*
  - Started in 1987
- *Journal of Cheminformatics*
  - Open access journal started in 2009

L.C. Ray and R.A. Kirsch (1957) Finding chemical records by digital computers, *Science*, **126**, 814-819

Introduced the use of graphs to represent 2D chemical structure diagrams

Applied a graph matching algorithm to a file of such representations to enable substructure searching



**Finding Chemical Records by Digital Computers**

Louis C. Ray and Russell A. Kirsch

# Graph theory describes sets of objects (*nodes*) and the relationships between pairs of them (*edges*)

# Representation of molecules by graphs

- Graph theory is applicable to any context that can be described by nodes and edges

- Can hence be used to represent and search both 2D and 3D chemical structures

- 2D chemical structure
  - Connection tables
  - Nodes correspond to atoms
  - Edges correspond to bonds
  - 2D graph describes topology

- 3D chemical structure (see later)
  - Edges correspond to distances
  - 3D graph describes geometry

- Throughout the early Sixties, Chemical Abstracts Service received very substantial funding to develop textual and chemical processing

- Heart of the new processing was the CAS Registry (now contains ca. 68M organic and inorganic molecules – see http://www.cas.org/)

- Adopted a graph-based approach: at the heart is a systematic naming procedure for chemical graphs:  H. L. Morgan (1965)  The generation of a unique machine description for chemical structures - a technique developed at Chemical Abstracts Service, *Journal of Chemical Documentation* **5**, 107-113.

# Wisweser Line Notation

- CAS work of long-term importance but most early industrial systems were based on Wiswesser Line Notation (WLN)

- Use of WLN till late Seventies
  - Cf SMILES and now InChI

- Need for conversion to connection tables
  - ICI CROSSBOW system
  - L.H. Thomson *et al.*, (1967) Organic search and display using a connectivity matrix derived from Wiswesser notation, *Journal of Chemical Documentation* **7**, 204-209

ZR CVR CE

# 2D substructure search output

# Substructure searching

- Ability to retrieve all molecules in a database containing a user-defined substructure
  - Use of a *subgraph isomorphism algorithm*
  - Completely *effective*, but *efficiency* very low
- Standard methods such as set reduction (Sussenguth, 1965) and relaxation (Ullmann, 1976) underlie all operational substructure searching systems (both 2D and 3D)
  - Still not sufficiently fast so need for initial filter to eliminate molecules from graph processing
  - Encoding fragment screens describing query substructures and database structures in a *bit-string* or *fingerprint*

# Binary vector

- Each bit in the bit-string (binary vector) records the presence ("1") or absence ("0") of a particular fragment in the molecule.
  - Typical length is a few hundred or few thousand bits
- A database structure is passed on for subgraph matching only if its bit-string contains all of the bits that have been set in the query's bit-string
- How to select the fragments?
  - J.E. Crowe *et al.* (1970) Analysis of structural characteristics of chemical compounds in a large computer-based file. *Journal of the Chemical Society (C)* 990-996.

# Example fragments

**a. Augmented Atom**
C rs C rd C rs C

**b. Atom Sequence**
C rs C rs C rd C

**c. Bond Sequence**
AA rs AA rs AA rd AA

**d. Ring Composition**
N rs C rd C rs C rs C rs

**e. Ring Fusion**
XX3 XX3 XX3 XX2 XX2

**f. Atom Pair**
N 0;3 - 2 - C 0;3

# Reaction databases

- How to search for structural changes occurring in a reaction?

- G.E. Vleduts (1963) Concerning one system of classification and codification of organic reactions, *Information Storage and Retrieval* **1**, 17-146

  - Index a reaction by just those parts that have changed, the *reaction centre*, to allow searches for both changed and unchanged substructures

  - Practical realisation of his ideas not till early Eighties

# Computer-aided synthesis design

- Vleduts' paper was also the first to suggest computer-aided synthesis design

- Potential syntheses of a target molecule using a reactions database plus appropriate inference mechanisms ("retrosynthesis")

- An early example of an expert system

- First implemented in OCSS (subsequently LHASA) in 1969

- CASD programs can also work in the synthetic direction

# Computer-aided structure elucidation

- Identification of an unknown substance based on spectral information

- DENDRAL project started in 1965: first paper J. Lederberg *et al.* (1969) Applications of artificial intelligence for chemical inference. 1. *Journal of the American Chemical Society*, **91**, 2973-2976

- First of a whole series of AI projects at Stanford (e.g., MYCIN, Prospector, XCON)

# Hansch analysis

- C. Hansch *et al.* (1962) Correlation of biological activity of phenoxyacetic acids with Hammett substituent constants and partition coefficients, *Nature*, **194**, 178-180.

- Use of linear regression analysis to correlate physicochemical parameters with bioactivity

- The standard technique for Quantitative Structure-Activity Relationship (QSAR) studies for over two decades

# Free-Wilson analysis

Use of structural, rather than physicochemical, variables in the regression , these denoting the  presence of substituents on a common scaffold

$R_X$ —⬤— $R_Y$

$R_X = X1, X2,....$
$R_Y = Y1, Y2, ....$

| Molecule ID | Indicator variables | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | X1 | X2 | ….. | Y1 | Y2 | ..... |
| 1 | 1 | 0 | ….. | 1 | 0 | ….. |
| 2 | 0 | 1 | ….. | 1 | 0 | ….. |
| 3 | 1 | 0 | ….. | 0 | 1 | ….. |
| 4 | 0 | 1 | ….. | 0 | 1 | ….. |
| ….. | ….. | ….. | ….. | ….. | ….. | ….. |

$$y = C + a_1 X1 + a_2 X2 + a_3 X3 + ....b_1 Y1 + b_2 Y2 + b_3 Y3 + ....$$

$$y = C + \sum_{i=1}^{n} a_i Xi + \sum_{j=1}^{m} b_j Yj$$

# Substructural analysis

- R.D. Cramer *et al.* (1974) Substructural analysis. A novel approach to the problem of drug design, *Journal of Medicinal Chemistry,* **17**, 533-535

- Extension of Free-Wilson ideas to encompass
  - Structurally diverse molecules
  - Qualitative activity data

- Calculation of fragment weights
  - Probability that a molecule containing that fragment will be active, e.g., $N_{act}/(N_{act}+N_{inanct})$

- Used in US government anti-cancer programme in Eighties, but then in abeyance for many years till "re-discovery" as naive Bayesian classifier
  - First application of machine learning in virtual screening

# Moving on

- Throughout the Seventies, chemical search systems (mainly WLN-based) became widely available across the pharmaceutical industry

- Extensive use of QSAR methods in lead optimisation

- Computer hardware/software limitations meant processing slow

- Things did not change much till the late-Seventies/early-Eighties, with the advent of MDL (now Accelrys) and CAS Online

# Similarity searching

- Substructure searching very powerful but requires a clear view of the types of structures of interest

- Given a *target* (or *reference*) structure find molecules in a database that are most similar to it ("give me ten more like this")

- The *similar property principle* states that structurally similar molecules tend to have similar properties (cf *neighbourhood principle*)



Morphine     Codeine     Heroin

# How to define chemical similarity?

- Most obvious way is use of a maximum common subgraph isomorphism procedure but far too time-consuming for database-scale applications

- Use of fingerprint comparisons
  - G.W. Adamson and J.A. Bush (1973) A method for the automatic classification of chemical structures, *Information Storage and Retrieval*, **9**, 561-568

- How to use this idea?
  - Operational implementations from mid-Eighties with systems at Lederle, Pfizer/Sheffield and Upjohn

- Still the most widely used approach, despite inherent simplicity

# Tanimoto-based 2D similarity searching



**Query**

**Chemical patents are an important source of chemical information**



R = 2-chlorophenyl or 2,3-dichlorophenyl

R1 = $CH_3$

R2 = $C_2H_5$

N = 2

R3 = H or $CH_3$

R4 = C-O-R5 or C-S-R6 or S-O-R7

R5 = H or $NHCH_3$ or $NHCH_2CONH_2$ or 2-pyridon-5-yl

R6 = $NH_2$ or C(=NHCN)$NHCH_3$

R7 = $NH_2$ or $NHCH_3$ or NH-cyclopentyl or 2-thienyl

or 8-quinolyl or 2-(4-methypiperazin-1-yl)pyrid-5-yl

# Markush structures: II

- This example encodes 192 specific molecules; for many patents, the number is not defined

- M.F. Lynch *et al*. (1981) Computer storage and retrieval of generic chemical structures in patents, Part 1. *Journal of Chemical Information and Computer Sciences,* **21**, 148-150.

- Extension of fingerprint and graph matching methods for specifics

- Work in collaboration with Derwent and CAS, resulting in the operational systems Markush DARC (now Merged Markush Service MMS) and MARPAT

# 3D substructure searching: I

- P. Gund (1977) Three-dimensional pharmacophoric pattern searching, *Progress in Molecular and Subcellular Biology,* **5**,117-143

- Recognition that the nodes and edges of a graph could represent the atoms and inter-atomic distances (where 'atom' may include pharmacophore points, e.g., lone pairs) of a 3D molecule

- But ideas not taken up for a decade:

  - Lack of structural data (except for the Cambridge Structural Database)

  - There was no obvious way of carrying out a search efficiently

# 3D substructure searching: II

- Intense interest from mid/late Eighties as both problems addressed

- Approximate 3D coordinates from structure-generation programs
  - CONCORD (Pearlman group at Austin, Texas)
  - CORINA (Gasteiger group at Erlangen)

- Fingerprint- and graph-based searching methods
  - S.E. Jakes and P. Willett (1986) Pharmacophoric pattern-matching in files of 3-D chemical structures - selection of interatomic distance screens, *Journal of Molecular Graphics,* **4**, 12-20
  - Basis of first systems at Pfizer and Lederle. Later extensions to encompass conformational flexibility, with industrial systems widely available from the mid-Nineties.

a = 8.62± 0.58 Angstroms

b = 7.08± 0.56 Angstroms

c = 3.35± 0.65 Angstroms

# Pharmacophore mapping

- Given a set of bioactive molecules: what is the common pharmacophore?

- Specify the points

  - G.R. Marshall *et al*. (1979) The conformational parameter in drug design: the active analogue approach in computer-assisted drug design, *ACS Symposium Series*, **112**, 205-226.

- No constraints

  - Y.C. Martin *et al*. (1993) A fast new approach to pharmacophore mapping and its application to dopaminergic and benzodiazepine agonists, *Journal of Computer-Aided Molecular Design*, **7**, 83-102.

- Inclusion of flexibility

  - G. Jones *et al*. (1995) A genetic algorithm for flexible molecular overlay and pharmacophore detection, *Journal of Computer-Aided Molecular Design*, **9**, 532-549.

# Ligand docking: I

- Fitting a molecule into a binding site
  - "Lock and key" model
- Two-part problem
  - Search algorithm to investigate possible poses
  - Scoring function to prioritise poses/molecules
- I.D. Kuntz *et al.* (1982) A geometric approach to macromolecule-ligand interactions, *Journal of Molecular Biology*, **161**, 269-288
- The DOCK program for fitting an individual molecule into an active site

# Ligand docking: II

- Extensions for

  - Scanning an entire database, taking each molecule in turn

  - Including ligand flexibility: G. Jones et al. (1995) Molecular recognition of receptor sites using a genetic algorithm with a description of desolvation". *Journal of Molecular Biology*, **245**, 43-53.

- Now a standard technique for virtual screening



4PHV docked (red) into HIV protease

# Comparative Molecular Field Analysis (CoMFA)

- R.D. Cramer *et al.* (1988) Comparative Molecular-Field Analysis (CoMFA). 1. *Journal of the American Chemical Society,* **110**, 5959-5967.

- Extension of Free-Wilson to 3D

  - Align a set of molecules, and place them in a 3D grid

  - Treat the computed interactions at each grid point as a variable for PLS analysis

- Now the standard QSAR tool

# Molecular diversity analysis: I

- Technological developments in the early Nineties led to a data explosion in the volumes of chemical and biological data

- Many more compounds **could** be made: which **should** be made?

- Need for tools to:
    - Quantify diversity
    - Select molecules so as to maximise diversity (NB not a new question, e.g., early Pfizer/Upjohn clustering work)

# Molecular diversity analysis: II

- Huge range of papers, focussing on fingerprint-based similarity approaches
    - E.J. Martin *et al*. (1995) Measuring diversity - experimental-design of combinatorial libraries for drug discovery, *Journal of Medicinal Chemistry,* **38**, 1431-1436.
    - R.D. Brown and Y.C. Martin (1996) Use of structure-activity data to compare structure-based clustering methods and descriptors for use in compound selection, *Journal of Chemical Information and Computer Sciences,* **36**, 572-584.

# Diversity alone is not enough

- It soon became clear that many of the molecules being generated had poor ADME characteristics

- ADME traditionally studied during optimisation

  - "Fail fast" paradigm implies that such molecules should be filtered out as early as possible

- C.A. Lipinski *et al*. (1997) Experimental and computational approaches to estimate solubility and permeability in drug discovery and development settings. *Advanced Drug Delivery Reviews,* **23**, 3-25

  - Criteria for oral activity: ideally, not more than 5 donors or 10 acceptors, MW under 500 and logP under 5

- Idea of drugability or drug-likeness

# Conclusions

- Chemoinformatics integrates long-established research into structure searching and bioactivity prediction

- Ever-increasing demands on the pharmaceutical industry will make it even more important in the future, e.g.
  - ADMETox prediction, Chemogenomics, and Virtual screening

- Histories
  - W.L. Chen (2006) "Chemoinformatics: past, present and future" *Journal of Chemical Information and Modeling*, **46**, 2230-2255
  - J. Gasteiger (2006) "Chemoinformatics: a new field with a long tradition" *Analytical and Bioanalytical Chemistry*, **384**, 57-64
  - A.G. Maldonado *et al*. (2006) "Molecular similarity and diversity in chemoinformatics: from theory to applications" *Molecular Diversity*, **10**, 39-79
  - P. Willett (2008) "From chemical documentation to chemoinformatics: fifty years of chemical information science." *Journal of Information Science*, **34**, 477-499