# From Mathematical Chemistry to Quantum and Medicinal Chemistry through Meta-Heuristics
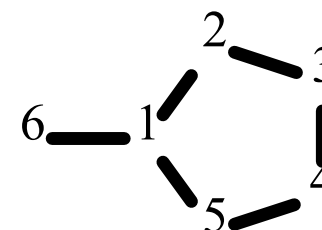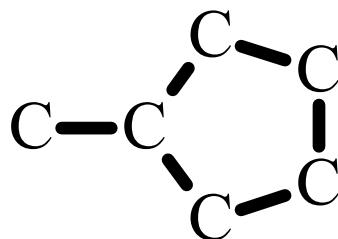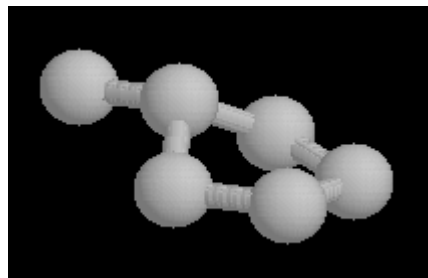
Lorentz JÄNTSCHI & Sorana D. BOLBOACĂ

*Technical & Medicine and Pharmacy Universities of Cluj-Napoca, Romania*

# Mathematical Chemistry Issues

- (Crum-Brown and Fraser, 1868): On the connection between chemical constitution and physiological action. Part 1. On the physiological action of the salts of the ammonium bases, derived from Strychnia, Brucia, Thebia, Codeia, Morphia, and Nicotia [Trans R Soc Edinb 25:151–203]
- (Sylvester, 1874): On an Application of the New Atomic Theory to the Graphical Representation of the Invariants and Covariants of Binary Quantics - With Three Appendices [Am J Math, 1, 64-90]
- (Harary, 1969): Graph Theory, Addison - Wesley, Reading, MA.
- (Kier and Hall, 1976): Molecular Connectivity in Chemistry and Drug Research, Acad Press, New York, NY.
- (Trinajstic, 1983): Chemical Graph Theory, CRC Press, Boca Raton, FL.
- (Diudea and others, 2001): Molecular Topology, Nova, Hutington, NY.

# Models of chemical structure



## 3D, 2D, and graph structure

- Molecular geometry – move to quantum chemistry
- Molecular topology:
  - Matrices (Adjacency, Laplacian, Distance, Detour, Combinatorial C(D,2), C(Δ,2), Wiener, Szeged, Path, Hosoya, Cluj, Distance-Extended, Detour-Extended, Reciprocal, Walk, Layer, Sequence);
  - Polynomials (Characteristic, Matching, Immanantal, Laplacian, Independence, Hosoya, Wiener, Z-counting);
  - Indices (half sum of elements from a symmetric matrix; so on - too long list).

# Moving to Quantum Chemistry

- (Schrödinger, 1926): An Undulatory Theory of the Mechanics of Atoms and Molecules [Phys Rev 28(6),1049–1070]
- $E\Psi = \hat{H}\Psi$ – Schrödinger time-independent
  - Where are the electrons and nuclei of a molecule in space?
    - configuration, conformation, size, shape, etc.
  - Under a given set of conditions, what are their energies?
    - heat of formation, conformational stability, chemical reactivity, spectral properties, etc.

# Molecular Modelling Software

- MPQC (Massively Parallel Quantum Chemistry) – computes properties of atoms and molecules from first principles using the time independent Schrödinger equation
- GAMESS (General Atomic and Molecular Electronic Structure System) – a general ab intio quantum chemistry package
- MOPAC (Molecular Orbital PACkage) - a semi-empirical quantum-mechanics code
- GAUSSIAN - predicts the energies, molecular structures, and vibrational frequencies of molecular systems, and other molecular properties derived from these
- HyperChem – molecular modeling environment that is known for its quality, flexibility, and ease of use
- Octopus - ab initio virtual experimentation; electrons are described quantum-mechanically in TD-DFT
- deMon2k - a software package for DFT calculations
- Many others (including ones attended here!)

# Medicinal Chemistry facts

- (Richet, 1893): first lipophilicity-activity relationship
  - "plus ils sont solubles, moins ils sont toxiques" [C R Seances Soc Biol Fil 45:775–776]
- (Hansch, 1962-1964) - foundations of QSAR:
  - Combination of several phycicochemical parameters in one regression equation;
  - Definition of the lipophilicity parameter π;
  - Formulation of the parabolic model for nonlinear lipophilicity-activity relationships.

# Property-Activity Measurements

- Quantitative
  - Absolute (two refs and a scale between)
    - Temperature (eg. boiling), Energy (eg. hidration)
  - Relative (one ref and a ratio)
    - Sweetness relative to fructose
- Semi-quantitative
  - Ordinal scale (precision, accuracy, confidence)
- Qualitative
  - Nominal (blood groups);
  - Binary (dead or alive; present or absent)

# Mathematical - Quantum - Medicinal

- These fields are somehow separated
- Mathematical:
  - Journal of Mathematical Chemistry (1987-)
  - MATCH Communications in Mathematical and in Computer Chemistry (1975-)
- Quantum:
  - International Journal of Quantum Chemistry (1967-)
  - Journal of Molecular Modeling (1995-)
- Medicinal:
  - Journal of Medicinal Chemistry (1959-)
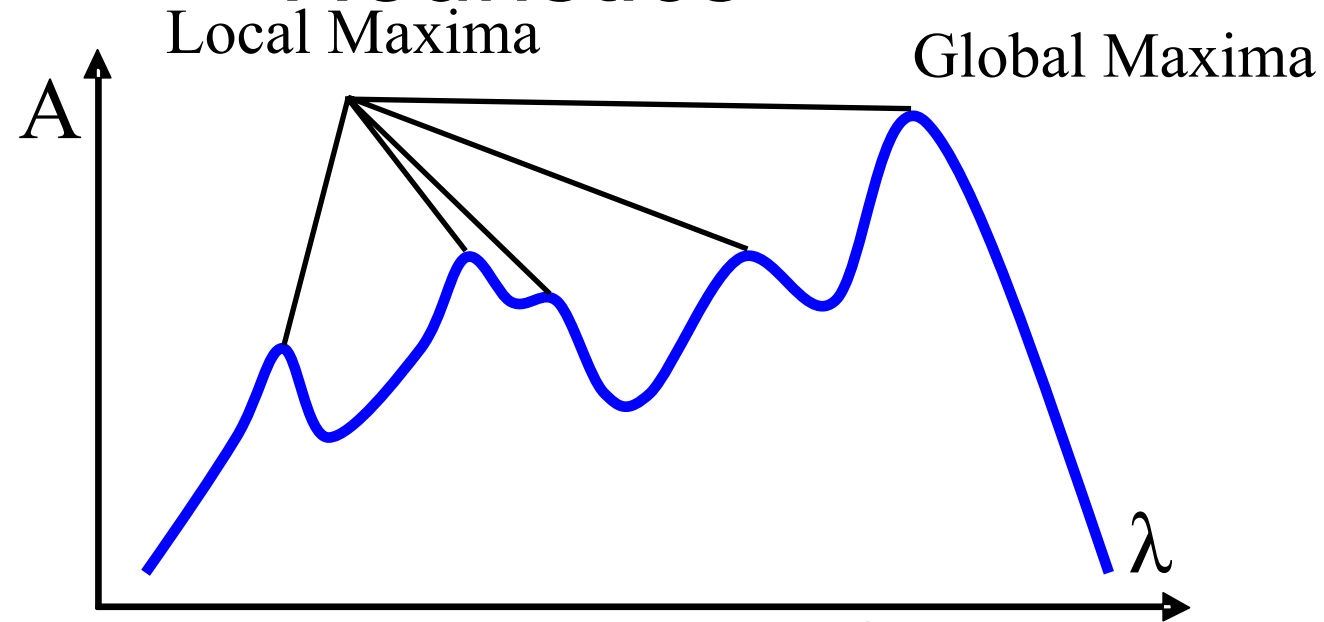  - Chemical Biology & Drug Design (1969-)

# Integrated approaches

- (Grassy and others, 1998): Computer Assisted Rational Design of Immunosuppressive Compounds [Nature Biotechnol 16:748-752] reports on a search for peptides possessing immunosuppressive activity. They used 27 structure descriptors (12 mathematical). From a combinatorial library of about 280000 compounds they selected 26 peptides for which high activity was predicted. Five of them were actually synthesized and tested experimentally. The most potent of these showed an immunosuppressive activity approximately 100 times higher than the lead compound.

# Hard problems and search for the answer

- **Hard problems**
  - Usually, we operate with problems. A problem has a precise meaning, very close to the meaning of the algorithm (eg. a recipe specifying what to do in certain conditions to obtain an objective). An algorithm uses two resources to solve a problem: time and space.
  - Not all problems (and their solving algorithms) has same complexity. A problem with exponential complexity (time increases exp. with the size of entry data) are called hard.
- **Search for the answer**
  - Hard problems: decision, classification, optimization, and simulation; theoretical studies prove that one type (eg. Decision) can be converted into another (eg. Optimization).
  - Hard problems optimum search solving algorithms for real applications often goes out of time. Not always a call for the optimum is required; frequently a good solution is enough.

# Heuristics



- Heuristics:
  - Sets of rules designed to solve a specific problem, usually based on common sense (regarding the expected solution) by avoiding of obvious mistakes.
  - These are still not designed to produce always an exact solution, or sometimes even may not produce an solution for any input data.
  - A lot of heuristics are ad-hoc and problem-dependent. Still, three heuristics are very general and has applicability to a lot of hard problems: meta-heuristics.
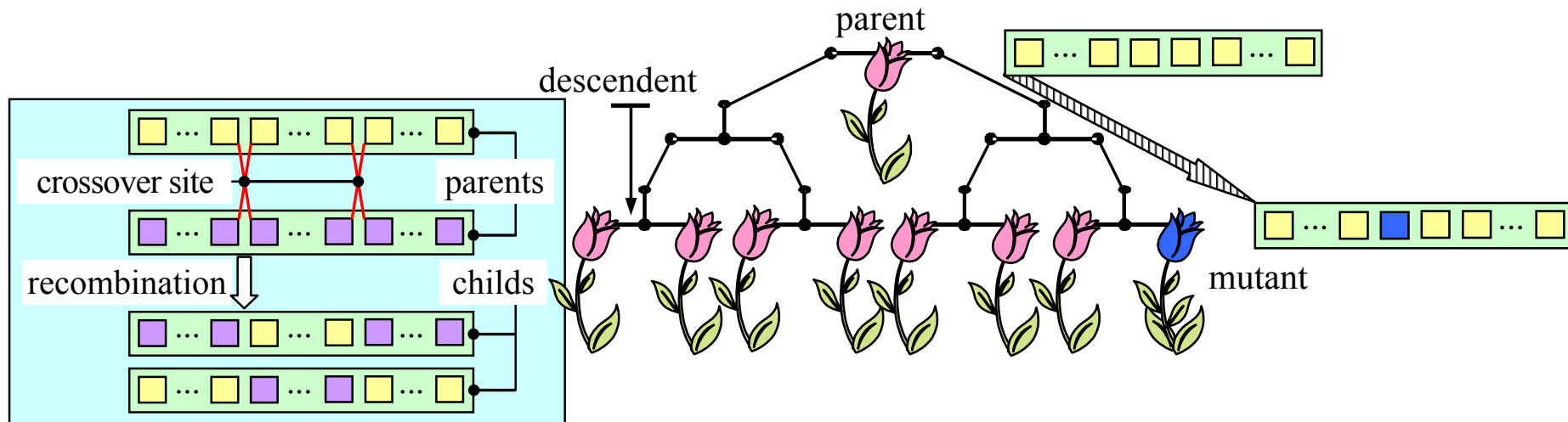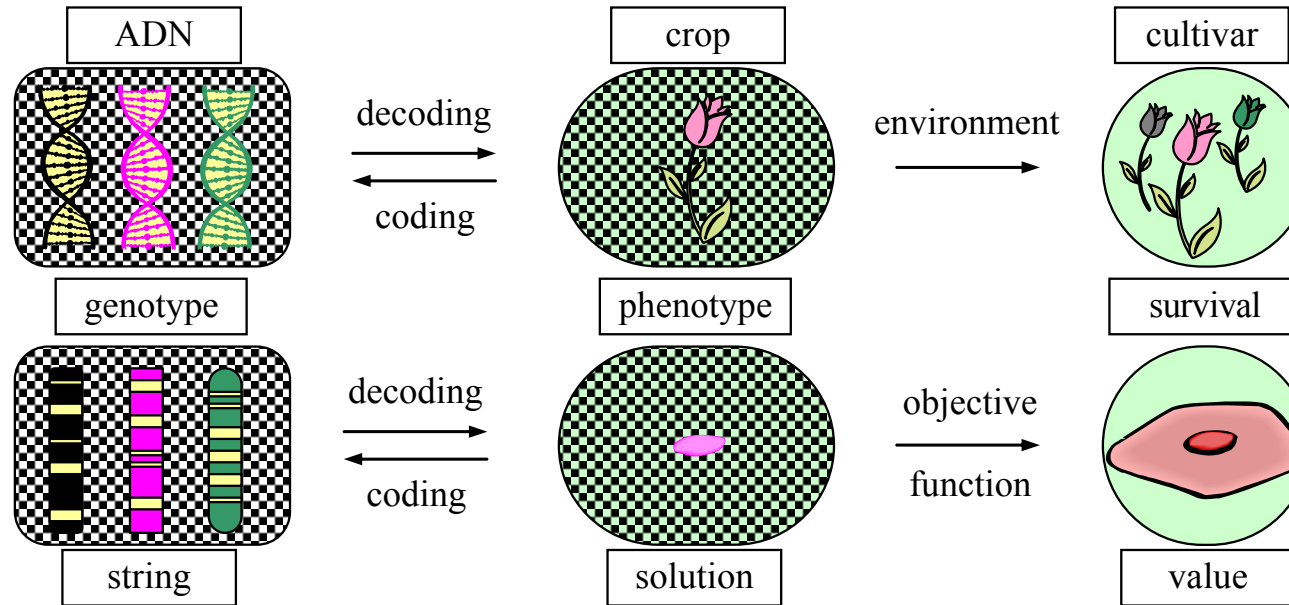
# Meta-heuristics

- All three are stochastic (implying the chance and probability):
  - Tabu Search (TS) – Glower (1977, 1986, 1992);
  - Simulated Annealing (SA) – van Laarhoven, Aarts, Davis (1987);
  - Genetic Algorithms (GA) | Evolutionary Algorithms (EA) – Fraser (1957-1970).
- Quality of an heuristic:
  - Speed (to solution);
  - Precision (how far is from global optima);
  - Scope (applicability domain) – how large is the subset of the subset of input data for which it performs well according to previous two criteria.

# Complexity

- NFLT – No Free Lunch Theorem (Wolpert and Macready, 1995 and 1997)
  - For A and B (algorithms) and input data Φ for which A better than B it exists input data Ψ for which B better than A

- $\Rightarrow$ Come back to the applicability domain
  - Of the algorithm
  - Of the QSAR equation (also an algorithm!)
  - So on

# Genetic Algorithms

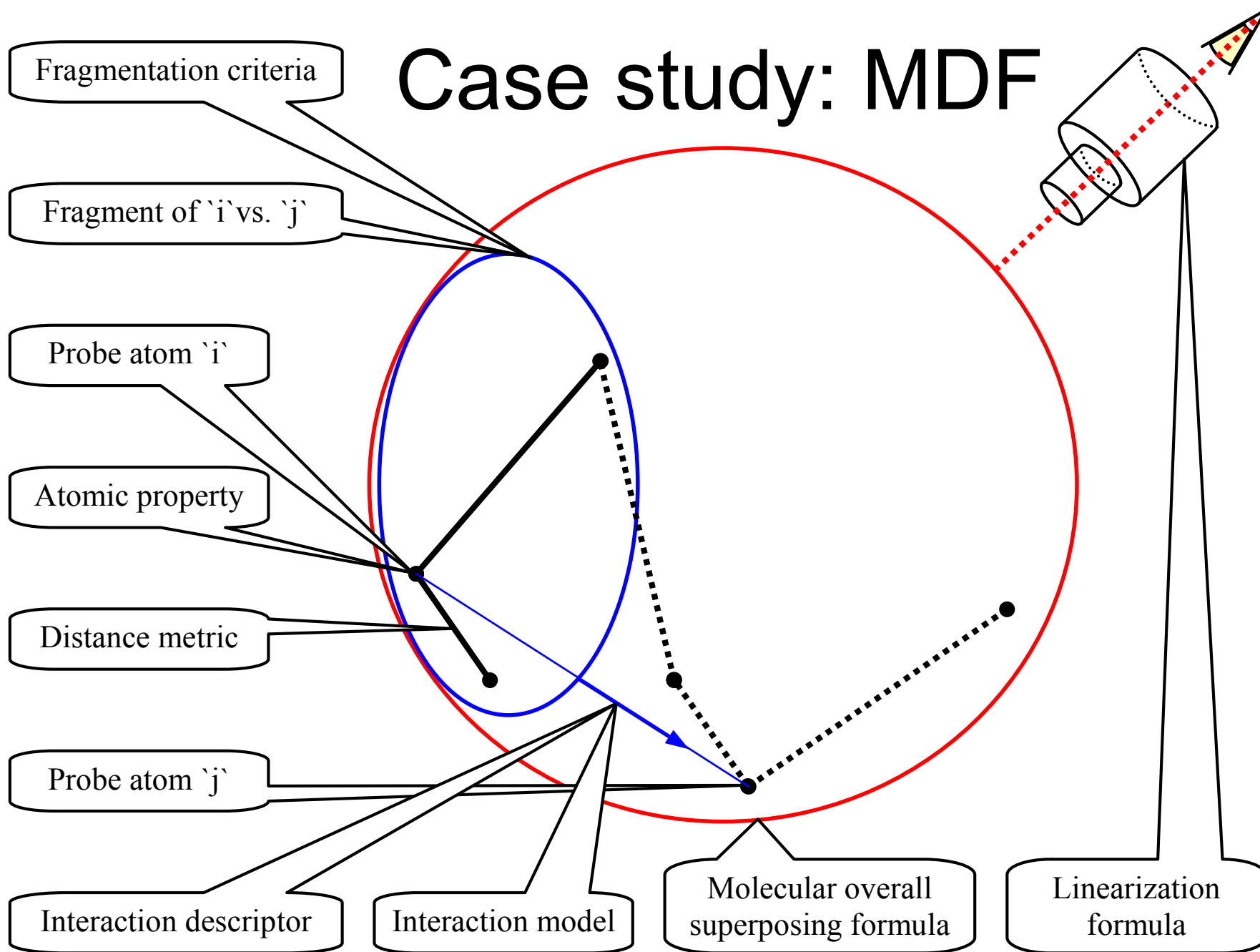gene    chromosome    population

ADN
decoding
coding
genotype

crop
decoding
coding
phenotype

environment
cultivar
survival

string
decoding
coding
solution

objective function
value

descendent
parent
mutant

crossover site    parents
recombination    childs

# Methodology

- Design methodology
  - Family of descriptors (FPIF[2000] or MDF[2005] or MDFV[2008])
- Linkage methodology
  - Genetic algorithms (for multivariate regressions)[2008]
- Assessment methodology
  - Descriptive (Jarque-Bera[2008]), Quantitative AND Qualitative stats measures (Pearson, Spearman, Kendall, Goodman–Kruskal)[2006]
- Prediction methodology
  - Combinatorial libs and Virtual screening[2009]

# Design: Families of descriptors

- FPIF – Fragmental Property Index Family
  - Match 40:151-188, 2000
  - SAR QSAR Environ Res 12:159-179, 2001
  - Chapt. 7 of Molecular Topology, 2001 & 2002
- MDF – Molecular Descriptors Family
  - Int J Quant Chem 107:1736-1744, 2007
  - Int J Mol Sci 8: 189-203 & 1125-1157, 2007
  - Chem Biol Drug Des 71:173-179, 2008
  - Env Chem Lett 6:175-181, 2008
  - Mar Drugs 6:372-388, 2008
  - Electronic J Biotechnol DOI: 10.2225/vol11-issue3-fulltext-9, 2008
- MDFV - ------''----- Vertex - In development

# Case study: MDF

Fragmentation criteria

Fragment of `i` vs. `j`

Probe atom `i`

Atomic property

Distance metric

Probe atom `j`

Interaction descriptor

Interaction model

Molecular overall superposing formula

Linearization formula

## Molecule filename:

000_PCB000.hin

## Distance operator:

Topological distance, t
Geometrical distance, g

## Atomic property:

Cardinality, C
Count of directly bounded hidrogen's, H
Relative atomic mass, M
Atomic electronegativity, E
Group electronegativity, G
Partial charge, Q

## Descriptor (of interaction) formula:

Distance, `D` = d
Inverted distance, `d` = 1/d
First atom's property, `O` = p1
Inverted O, `o` = 1/p1
Product of atomic properties, `P` = p1p2
Inverted P, `p` = 1/p1p2
Squared P, `Q` = p1p2^1/2
Inverted Q, `q` = 1/p1p2^1/2
First atom's Property multiplied by distance, `J` = p1d
Inverted J, `j` = 1/p1d
Product of atomic properties and distance, `K` = p1p2d
Inverted K, `k` = 1/p1p2d
Product of distance and squared atomic properties, `L` = d(p1p2)^1/2
Inverted L, `l` = 1/p1p2d
First atom's property potential, `V` = p1/d
First atom's property field, `E` = p1/d^2
First atom's property work, `W` = p1^2/d
Properties work, `w` = p1p2/d
First atom's property force, `F` = p1^2/d^2
Properties force, `f` = p1p2/d^2
First atom's property weak nuclear force, `S` = p1^2/d^3
Properties weak nuclear force, `s` = p1p2/d^3
First atom's property strong nuclear force, `T` = p1^2/d^4
Properties strong nuclear force, `t` = p1p2/d^4

## Interaction model:

Rare model and resultant relative to fragment's head, R
Rare model and resultant relative to conventional origin, r
Medium model and resultant relative to fragment's head, M
Medium model and resultant relative to conventional origin, m
Dense model and resultant relative to fragment's head, D
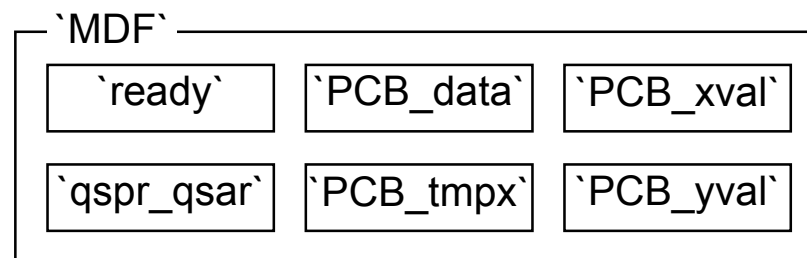Dense model and resultant relative to conventional origin, d

## Molecular overall superposing formula:

Cond., smallest, m
Cond., highest, M
Cond., smallest absolute, n
Cond., highest absolute, N
Avg., sum, S
Avg., average, A
Avg., S/count(fragments), a
Avg., Avg.(Avg./atom)/count(atoms), B
Avg., S/count(bonds), b
Geom., product, P
Geom., mean, G
Geom., P^1/count(fragments), g

## Fragmentation criteria:

Minimal fragments, m
Maximal fragments, M
Szeged distance based fragments, D
Cluj path based fragments, P

## Linearization operator:

Identity (no change), I
Inversed I, i
Absolute I, A
Inversed A, a
Logarithm of A, L
Logarithm of I, l

# Linkage: Genetic algorithms

| `MDF` | | |
|---|---|---|
| `ready` | `PCB_data` | `PCB_xval` |
| `qspr_qsar` | `PCB_tmpx` | `PCB_yval` |

- Genotype: 7 genes:
  - d (distance operator), p (atomic property), I (interaction descriptor), O (overlapping interactions operator), f (pair-based fragmentation criteria), M (fragments overlapping operator), L (linearizing operator)

- Phenotype: array of values for molecules from training set

- Population: a fixed size (eg. 100)

- Score: multiple regression: $r^2(Y, a_0 + \Sigma_i a_i \cdot Phenotype_i)$

- Objective: max.

- Mutation:
  - A low probability decides if applies;
  - A individual are selected (see selection)
  - A gene are random choused; a random replacing value replaces its value;

# Linkage using GA (2)

- Crossover:
  - Genotype$_1$=d$_1$p$_1$I$_1$O$_1$f$_1$M$_1$L$_1$ and Genotype$_2$=d$_2$p$_2$I$_2$O$_2$f$_2$M$_2$L$_2$
  - A double crossover site are randomly choused (let be 2,4)
  - Offspring$_1$=d$_1$p$_1$I$_2$O$_2$f$_2$M$_1$L$_1$ and Offspring$_2$= d$_2$p$_2$I$_1$O$_1$f$_1$M$_2$L$_2$
- Selection: based on fitness (score):
  - Proportional: $f_i$=Score(Genotype$_i$); $p_i = f_i/\Sigma_i f_i$
  - Deterministic (elitist): i | $f_i$ = max. OR min.
  - Tournament: one of {i,k}; max. OR min. ($f_i$,$f_k$)
  - Normalized: $g_i \leftarrow (f_i - F_0)(f_{max} - f_{min})/(F_1 - F_0)$; $p_i = g_i/\Sigma_i g_i$
  - Ranks: $h_i$=Rank($f_i$); $p_i = h_i/\Sigma_i h_i$

# Linkage using GA (3)

- Evolution:
  - Random generate N genotypes (seeds);
  - Construct phenotypes for the genotypes (cultivar);
  - Repeat
    - For every m-uple of phenotypes (for m-varied QSAR)
      - Obtain $r^2(Y, a_0 + \Sigma_i Phenotype_i)$
    - For every genotype (MyG)
      - Score(MyG) $\leftarrow$ min. $r^2(Y, a_0 + a_1 \cdot MyG + \ldots)$
    - Select two genotypes; do crossover; if alive then add to cultivar
    - Decide using probability of 1/7 if and then mutate one genotype; if alive then add to cultivar
    - Kill individuals as many as necessary to keep the size to N using Score($\cdot$)
  - Until max $r^2(Y, a_0 + \Sigma_i Phenotype_i) \geq$ Value (eg. 0.99)

# Assessment

- Reject a descriptor or an QSAR if:
  - Jarque-Bera: It has JB value larger than JB for measured property;
  - Pearson r: It has a not significant correlation;
  - Spearman ρ: ---''---
  - Kendall т-a, т -b, т-c: ---''---
  - Goodman–Kruskal Γ: ---''---
- http://l.academicdirect.org/Statistics/linear_dependence/

# Assessment example on MDFV

Descriptive Correlation Analysis on 31aa set.

| Id | Prop | Mols | Vars | r2Pearson | r2Spearman | r2Ken_a | r2Ken_b | r2Ken_c | r2Gamma | r2Geometry | Equation |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | MP | 23 | 1 | 0.6287 | 0.6012 | 0.2890 | 0.3692 | 0.2644 | 0.4718 | 0.4139 | 7.944e+1+GLsIPAdR*2.036e+2 |
| 2 | MP | 23 | 2 | 0.6803 | 0.6580 | 0.3468 | 0.4219 | 0.3173 | 0.5181 | 0.4699 | 1.992e+2+GLsIPAdR*4.522e+1+GLPICFdR*4.145e+1 |
| 3 | MP | 23 | 2 | 0.8058 | 0.7856 | 0.5347 | 0.5347 | 0.4892 | 0.5347 | 0.6015 | 2.268e+2+GLsIPAdR*5.395e+1+GQfAIcDR*-1.000e+2 |
| 4 | MP | 23 | 2 | 0.8248 | 0.7306 | 0.4462 | 0.4462 | 0.4082 | 0.4462 | 0.5288 | 2.265e+2+GLsIPAdR*6.037e+1+TQfAPpDR*-2.562e+2 |
| 5 | MP | 23 | 2 | 0.8513 | 0.8738 | 0.6438 | 0.6438 | 0.5890 | 0.6438 | 0.6993 | 2.260e+2+GLPICFdR*5.663e+1+GQfAIcDR*-1.085e+2 |
| 6 | MP | 23 | 2 | 0.8859 | 0.8500 | 0.5819 | 0.5819 | 0.5324 | 0.5819 | 0.6551 | 2.260e+2+GLPICFdR*6.390e+1+TQfAPpDR*-2.871e+2 |
| 7 | MP | 23 | 3 | 0.8959 | 0.8409 | 0.5699 | 0.5699 | 0.5215 | 0.5699 | 0.6461 | 2.227e+2+GLsIPAdR*1.903e+1+GLPICFdR*4.886e+1+TQfAPpDR*-2.697e+2 |
| 8 | MP | 23 | 3 | 0.9045 | 0.8573 | 0.6063 | 0.6063 | 0.5547 | 0.6063 | 0.6765 | 3.754e+2+GLsIPAdR*4.097e+1+GQfAIcDR*-1.073e+2+GLhIacdI*-1.583e+2 |
| 9 | MP | 23 | 3 | 0.9308 | 0.8609 | 0.6063 | 0.6063 | 0.5547 | 0.6063 | 0.6802 | 3.813e+2+GLsIPAdR*4.715e+1+TQfAPpDR*-2.784e+2+GLhIacdI*-1.645e+2 |
| 10 | MP | 23 | 3 | 0.9321 | 0.8943 | 0.6694 | 0.6694 | 0.6125 | 0.6694 | 0.7314 | 2.440e+2+GLPICFdR*5.367e+1+TQfAPpDR*-2.770e+2+GL5IPIdI*-8.048e+0 |
| 11 | MP | 23 | 3 | 0.9369 | 0.8924 | 0.6694 | 0.6694 | 0.6125 | 0.6694 | 0.7318 | 2.344e+2+GLPICFdR*5.534e+1+TQfAPpDR*-2.841e+2+TA3PIpDL*9.303e+0 |
| 12 | MP | 23 | 3 | 0.9443 | 0.9244 | 0.7222 | 0.7222 | 0.6607 | 0.7222 | 0.7753 | 2.276e+2+GLPICFdR*6.116e+1+TQfAPpDR*-2.685e+2+GApaaCDR*-1.400e-2 |
| 13 | MP | 23 | 3 | 0.9505 | 0.9187 | 0.7222 | 0.7222 | 0.6607 | 0.7222 | 0.7754 | 2.184e+2+GLPICFdR*6.911e+1+GESACFdI*6.700e+0+TQoAPidI*-4.621e+0 |
| 14 | MP | 23 | 3 | 0.9528 | 0.9050 | 0.6890 | 0.6890 | 0.6304 | 0.6944 | 0.7508 | 1.849e+2+GLPICFdR*6.916e+1+TQ1IFfDL*9.745e+0+GQ1ICPdL*1.677e+1 |
| 15 | MP | 23 | 3 | 0.9580 | 0.8780 | 0.6502 | 0.6502 | 0.5949 | 0.6553 | 0.7194 | 2.099e+2+GLPICFdR*7.675e+1+TQOAAidI*-1.301e+1+TQbFiFdL*3.992e+0 |
| 36 | MP | 23 | 4 | 0.9580 | 0.9187 | 0.7088 | 0.7088 | 0.6485 | 0.7088 | 0.7668 | 1.850e+2+GLsIPAdR*-1.545e+1+GLPICFdR*8.133e+1+TQ1IFfDL*1.051e+1+GQ1ICPdL*1.813e+1 |
| 37 | MP | 23 | 4 | 0.9589 | 0.8887 | 0.6824 | 0.6824 | 0.6244 | 0.6824 | 0.7436 | 2.113e+2+GLsIPAdR*-6.273e+0+GLPICFdR*8.195e+1+TQOAAidI*-1.350e+1+TQbFiFdL*4.060e+0 |
| 38 | MP | 23 | 4 | 0.9595 | 0.9093 | 0.6955 | 0.6955 | 0.6364 | 0.6955 | 0.7561 | 3.898e+2+GLsIPAdR*4.039e+1+TQfAPpDR*-2.588e+2+GQZaaiDL*-8.680e+0+GLhIacdI*-1.673e+2 |
| 39 | MP | 23 | 4 | 0.9603 | 0.9339 | 0.7222 | 0.7222 | 0.6607 | 0.7222 | 0.7788 | 4.687e+2+GLsIPAdR*3.586e+1+TQfAPpDR*-2.823e+2+GLhIacdI*-1.965e+2+GA3AaPdI*-7.081e+1 |
| 40 | MP | 23 | 4 | 0.9644 | 0.9093 | 0.6955 | 0.6955 | 0.6364 | 0.6955 | 0.7567 | 4.785e+2+GLsIPAdR*3.559e+1+TQfAPpDR*-2.830e+2+GLhIacdI*-1.994e+2+GAbAaPdI*-8.212e+1 |
| 41 | MP | 23 | 4 | 0.9653 | 0.8789 | 0.6630 | 0.6630 | 0.6066 | 0.6682 | 0.7298 | 2.671e+2+GLsIPAdR*4.763e+1+TQfAPpDR*-2.905e+2+GLMIiPdI*-6.891e+1+GAuAIcDR*2.000e-3 |
| 42 | MP | 23 | 4 | 0.9756 | 0.9282 | 0.7630 | 0.7630 | 0.6981 | 0.7630 | 0.8092 | 1.557e+2+GLsIPAdR*9.125e+1+GL7aCFDR*-5.000e-3+GAgAiCdL*7.958e+0+GLfICFdI*6.171e+0 |
| 43 | MP | 23 | 4 | 0.9794 | 0.9667 | 0.8481 | 0.8481 | 0.7760 | 0.8481 | 0.8748 | 4.195e+2+GLsIPAdR*5.413e+1+GQSIPIdI*5.528e+0+GL7IacDL*2.876e+1+GQYaFiDL*-1.049e+1 |
| 44 | MP | 23 | 4 | 0.9804 | 0.9706 | 0.8628 | 0.8628 | 0.7894 | 0.8628 | 0.8856 | 2.254e+2+GLPICFdR*6.150e+1+GQaFCPdR*-3.025e+2+GESACFdI*6.054e+0+TQoAPidI*-4.036e+0 |
| 45 | MP | 23 | 4 | 0.9815 | 0.9589 | 0.7909 | 0.7909 | 0.7236 | 0.7909 | 0.8342 | 2.348e+2+GLPICFdR*6.624e+1+GQaFCPdR*-4.292e+2+TQHAPpdI*-3.987e+0+GQ1ICPdL*1.710e+1 |
| 46 | MP | 23 | 4 | 0.9824 | 0.9488 | 0.7561 | 0.7561 | 0.6918 | 0.7622 | 0.8093 | 2.368e+2+GLPICFdR*6.699e+1+GQ1FCCdR*-1.209e+2+TQHAPpdI*-4.565e+0+GQ1ICPdL*1.605e+1 |
| 47 | MP | 23 | 4 | 0.9454 | 0.9017 | 0.6955 | 0.6955 | 0.6364 | 0.6955 | 0.7532 | 2.495e+2+GLPICFdR*7.222e+1+GLvIFPdR*0.000e-1+TQOAAidI*-1.073e+1+TQsFPIdR*-9.437e+3 |
| 48 | MP | 23 | 4 | 0.9841 | 0.9339 | 0.7630 | 0.7630 | 0.6981 | 0.7630 | 0.8112 | 2.396e+2+GLPICFdR*6.638e+1+GQqFICDR*-1.565e+2+TQHAPpdI*-4.661e+0+GQ1ICPdL*1.767e+1 |
| 49 | MP | 23 | 4 | 0.9847 | 0.9411 | 0.7700 | 0.7700 | 0.7045 | 0.7761 | 0.8184 | 2.337e+2+GLPICFdR*6.816e+1+GQmFIFdR*-5.391e+3+TQHAPpdI*-4.365e+0+GQ1ICPdL*1.737e+1 |
| 50 | MP | 23 | 4 | 0.9871 | 0.9474 | 0.8050 | 0.8050 | 0.7365 | 0.8050 | 0.8432 | 2.042e+2+GLPICFdR*6.818e+1+GL0IPadI*-1.262e+1+TQbFiFdL*3.311e+0+GQ1ICPdL*1.775e+1 |

# Interpretation – Monovariate

- Based on descriptor formula (from Chem Biol Drug Des 2008; 71:173-9)

**Table 1:** SAR models for amino acids

| Amino acid property | Hyd(20) |
|---|---|
| MDF SAR equation | $\hat{Y} = -160X - 0.065$ |
| SAR determination (%) | 65 |
| MDF descriptor (X) | AbmrEQg |
| Dominant atomic property | Charge (Q) |
| Interaction via | Space (geometry) |
| Interaction model | $Qd^2$ |
| Structure on activity scale | Proportional |

- A: absolute value
- b: avg. by bonds
- m: min. fragments
- r: rare interactions
- E: charge field
- Q: charge
- g: geometry

# Interpretation - Multivariate

- $\hat{Y}_{2v}$ = -2.261 + 0.037·*ASMmVQt* -0.216·*IfDdOQg (from* Electron J Biotechnol DOI 10.2225/vol11-issue3-fulltext-9)

- # Monovariate
  - Should characterize a global property/measure
    - Free energy, …

- # Multivariate (2-, 3-, …)
  - Should characterize a partial property/measure
    - Enthalpy
    - Entropy
    - Environment (pH, $H_2O$)
    - ...

# Prediction

- Internal validation
  - Training vs. Test experiment (split the data into Training and Test sets)
  - Leave-one-out
- External validation
  - External set
- Use combinatorics to generate new compounds
- Use software to construct 3D-model
- Use obtained QSARs to do prediction

# Comparison

- Steiger's Z test - Correlated correlations analysis
  - Y (measured), Y1 (predicted), Y2 (predicted)
    - A Z-value based on r(Y,Y1), r(Y,Y2), and r(Y1,Y2)
  - Overlapping of predictors; Compute (/check if):
    - Probability that Y1 and Y2 express (comes) from same reasoning (i.e. more than 95% confidence)
    - Probability that Y1 and Y2 express (comes) from different reasoning (i.e. less than 5% confidence)

# Conclusion: tools

- Property
  - Design: Gauss - Fisher
  - Characterization: JB (skewness and kurtosis)
- Structure
  - Design – software (Quantum)
  - Characterization – descriptors (Mathematical)
  - Population Diversity – Families of descriptors (Physical)
  - Assessment: JB
- Relationship
  - Design: Evolutionary (genetic) algorithms
  - Assessment (1): JB, r, $\rho,\tau$-a,b,c, $\Gamma$
  - Assessment (2): TvT, cv-loo, Steiger, External data sets
- Usage
  - Virtual screening (Combinatorial)
  - Design (medicinal)

# Acknowledgements