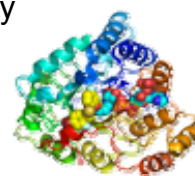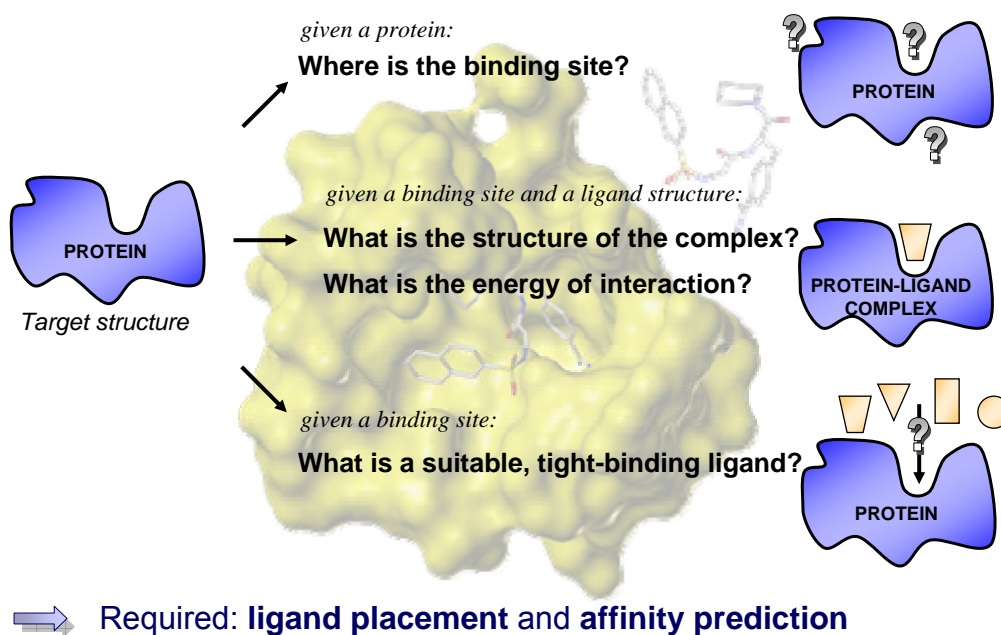# Docking and scoring

# of protein-ligand complexes:

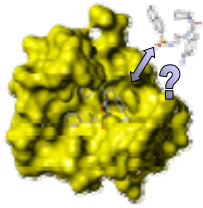## *What is possible and what is not?*

Christoph Sotriffer

Institute of Pharmacy and Food Chemistry
University of Würzburg
Am Hubland
D – 97074 Würzburg

---

## Key questions in structure-based drug design

*given a protein:*
**Where is the binding site?**

PROTEIN

PROTEIN

*Target structure*

*given a binding site and a ligand structure:*
**What is the structure of the complex?**
**What is the energy of interaction?**

PROTEIN-LIGAND
COMPLEX

*given a binding site:*
**What is a suitable, tight-binding ligand?**

PROTEIN

Required: **ligand placement** and **affinity prediction**

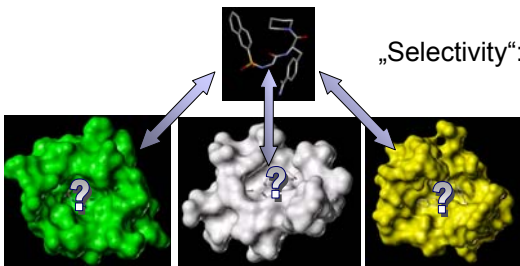## Docking problems & scoring tasks
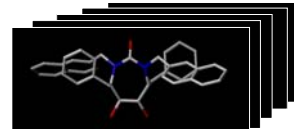


„Single Docking":   *1 protein  –  1 ligand*

looking for:   binding mode (and affinity) of the ligand

„Virtual Screening":   *1 protein  –  many (potential) ligands*

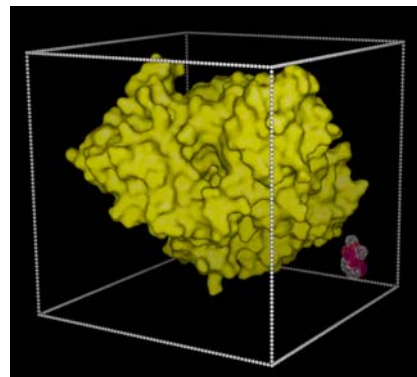looking for:   ligands with high affinity for target protein



„Selectivity":   *many proteins  –  one or more ligands*

looking for:  ligands with high
                    selectivity for one target
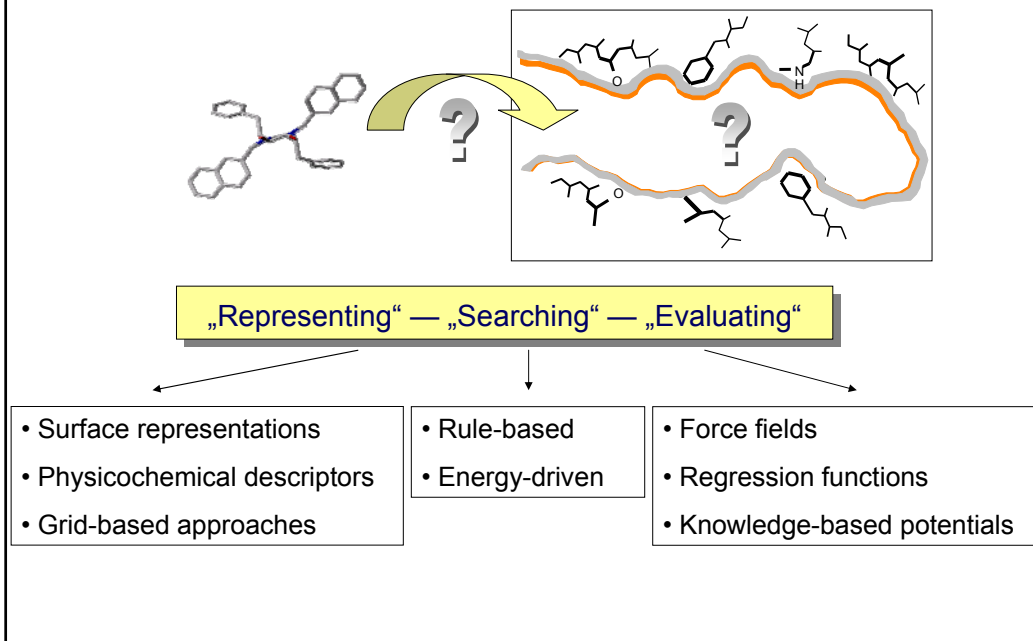
---

## Why is docking a „problem"?

• complex 3D jigsaw puzzle

• conformational flexibility

• mutual adaptations („induced fit")

• solvation in aqueous media

• complexity of thermodynamic contributions

• no easy route to $\Delta G$ evaluation for scoring



⟹ Simplifications und heuristic approaches necessary

*Modelling and computer-aided drug design are frequently
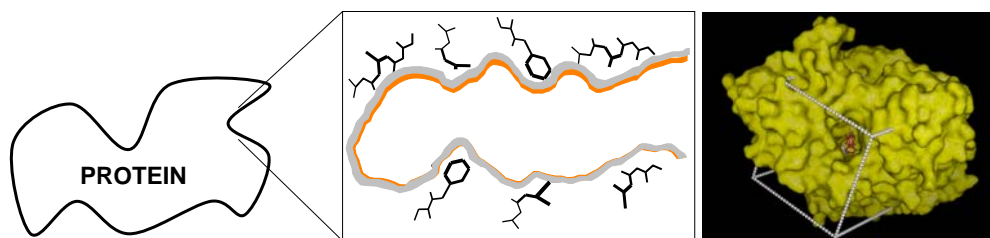a quest for suitable simplifications*

## Approaches to solve the docking problem



„Representing" — „Searching" — „Evaluating"

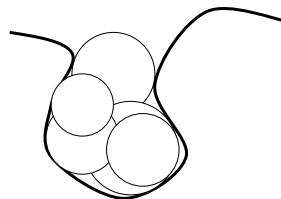| | | |
|---|---|---|
| • Surface representations | • Rule-based | • Force fields |
| • Physicochemical descriptors | • Energy-driven | • Regression functions |
| • Grid-based approaches | | • Knowledge-based potentials |

## „Representing": Molecular representations for docking

A.) Protein

0. Restricting the search space to the binding pocket
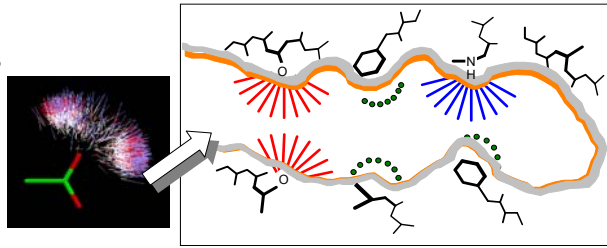


1. Geometric surface descriptors

   e.g., sphere representation of binding pockets

   (→ program DOCK)

2. Underline{Physicochemical descriptors}
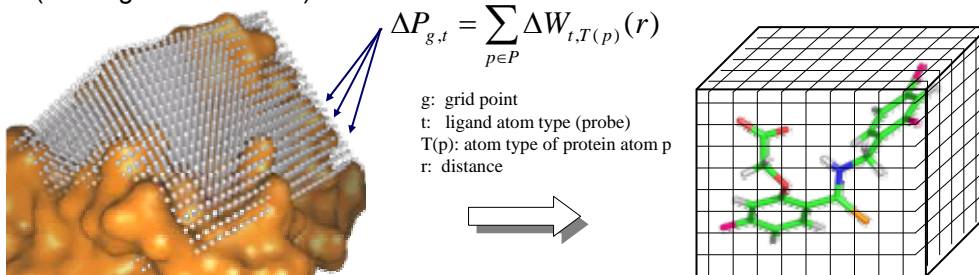
   Interaction points and vectors

   ($\rightarrow$ Programs LUDI, FlexX)



3. Underline{Grid representations}

   Interaction potentials of probe atoms are mapped to grid points

   ($\rightarrow$ Program AutoDock)

$$\Delta P_{g,t} = \sum_{p \in P} \Delta W_{t,T(p)}(r)$$

g: grid point
t: ligand atom type (probe)
T(p): atom type of protein atom p
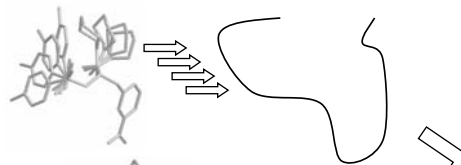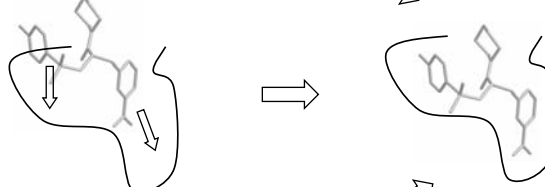r: distance



---

# B.) Ligand        major problem: conformational flexibility

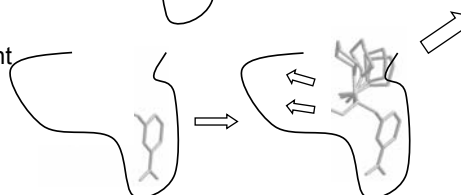$\Rightarrow$ | strategies for flexible ligand docking |

1.) rigid docking of conformers

2.) simultaneous optimization
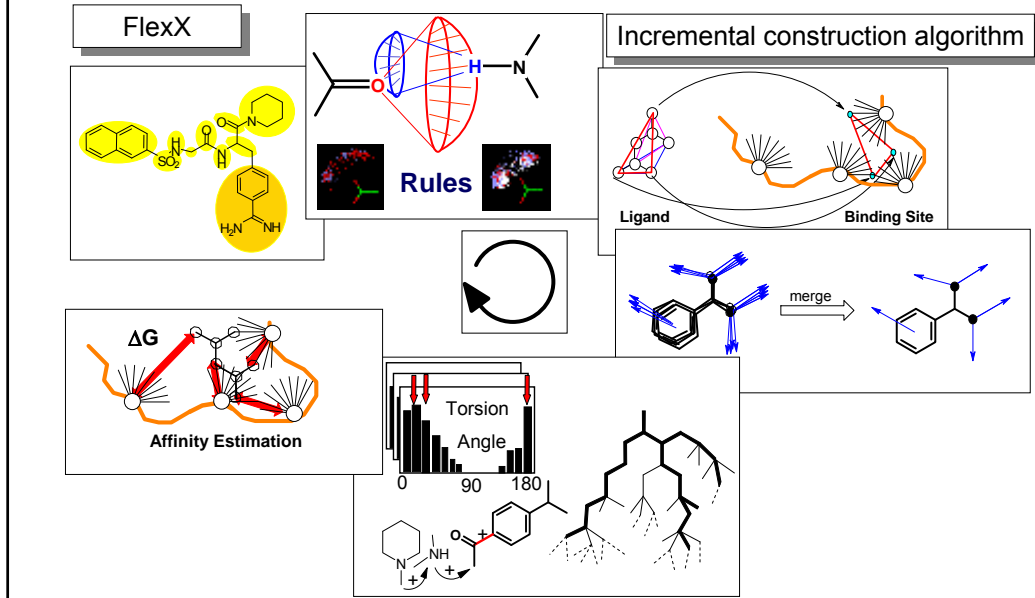    of orientation and conformation

3.) placement of a base fragment
    followed by incremental
    construction

# „Searching": Search algorithms for docking procedures

## 1.) Rule-based: geometric-combinatorial methods

FlexX

Rules

Incremental construction algorithm

Ligand   Binding Site

merge

ΔG

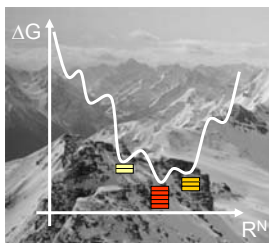**Affinity Estimation**

Torsion Angle

0   90   180

---

## 2.) Energy-driven: stochastic optimization methods

general assumption:

experimentally determined complex structure corresponds to global minimum of $\Delta G_{bind}$

Docking = optimization problem

$\Delta G$

$R^N$

- Search for $\min(\Delta G_{bind})$-binding mode
- $\Delta G_{bind}$ approximated by scoring function
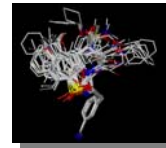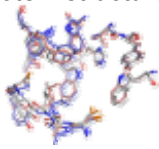- „rugged", multi-dimensional energy landscape

Monte-Carlo methods, genetic algorithms

examples: AutoDock, ICM, GOLD

## Before docking ...

... take care of the setup!

- Protein structures:
  - Protonation states and H-bonding networks
  - Quality and completeness of structural data
  - Location of binding site
  - Experimental data about water molecules and flexible regions

- Ligand structures:
  - Protonation states (influenced by protein!)
  - Tautomers
  - Conformers

- Docking program:
  - Choose suitable parameters
  - Validate, validate, validate (in particular for your system)!

- Know you program!

- Check structures and setup visually!

- Critically assess the quality of automated setup routines!

---

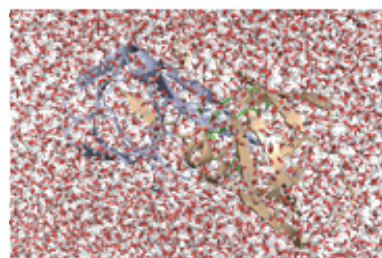## „Evaluating" / Scoring: Why is affinity prediction a challenge?

1.) Protein-ligand complexes are dynamic systems in aqueous solution

- huge number of particles
- simultaneous, unperiodic, continuously changing interactions

⟹ Simulation methods required!

Statistical thermodynamics: Calculation of ΔG° needs integration over entire phase space!

⟹ Computationally very expensive!

2.) The prediction methods need to be fast

Database screens: ~ $10^3 - 10^6$ molecules need to be compared

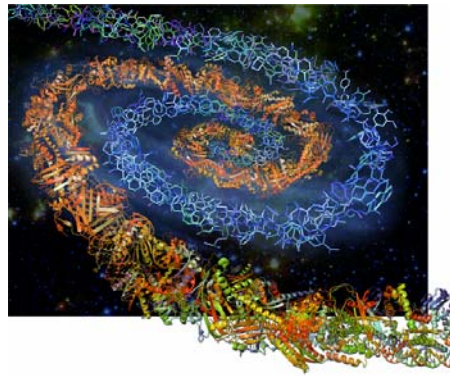Docking runs: ~ $10^7 - 10^9$ configurations need to be evaluated

⟹ „Scoring functions" required:

Fast, simplified, heuristic methods for prediction of binding strength

## Scoring functions: Goals

The ultimate goals of an ideal function:

- accurate within less than 1 $pK_D$ unit (<1.4 kcal/mol)

- generally valid (not system specific; large affinity range)

- robust (tolerant with respect to structural uncertainties)

- widely applicable (docking, virtual screening)

- physically meaningful (interpretable)

- fast and easy to compute



Aus: Klebe, Wirkstoffdesign, 2. Aufl. © Spektrum Akademischer Verlag GmbH, 2009

## Scoring functions: Tasks and types
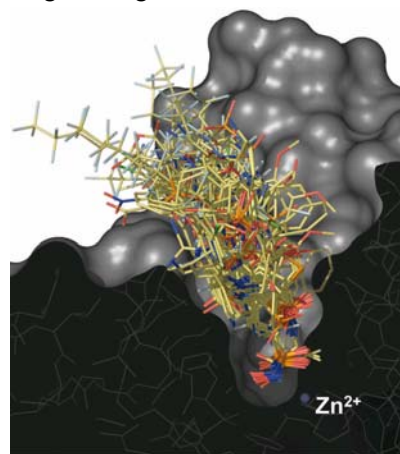
Application tasks:

A) Identification of the correct binding mode for a given ligand
   *Pose prediction in docking*

B) Identification of new active ligands
   *Virtual screening*

C) Affinity ranking for compound series
   *Ligand design, lead optimization*

Available approaches:

- Force field-based methods

- Knowledge-based scoring functions
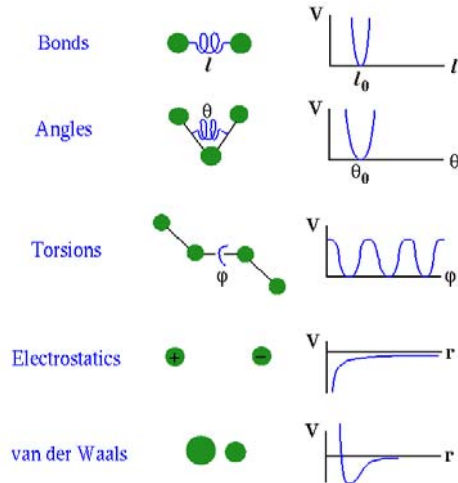
- Empirical scoring functions

# Force field-based methods

Molecular Mechanics (MM):

- atoms $\rightarrow$ charged spheres

- bonds $\rightarrow$ springs

- classical potentials

- no electrons $\rightarrow$ no bond formation / cleavage

- typically parameterized to reproduce
  molecular potential energy surface
  ($\rightarrow$ conformational $\Delta H$ in the gas phase!)

⟹ Scoring protein-ligand complexes:

  **+** for pose prediction in docking

  **−** for ligand ranking by affinity

⟹ Terms accounting for (de)solvation & entropic factors required (cf. MM-PBSA)

Bonds

Angles

Torsions

Electrostatics

van der Waals

---

# Knowledge-based scoring functions
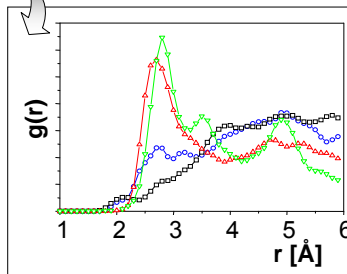
Derivation from
crystal-structure data

$$P_{ij}(r) = -\ln \frac{g_{ij}(r)}{g_{ref}}$$

$P_{ij}$: distance-dependent pair potential

$g_{ij}$: frequency distribution of atom-atom contacts
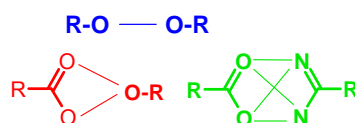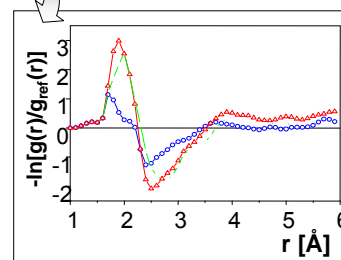
$g_{ref}$: reference distribution

**Relibase**

Frequency of occurrence

*No experimental affinities used!*

Statistical potential

g(r)

r [Å]

$-\ln[g(r)/g_{ref}(r)]$

r [Å]

R-O —— O-R

# Empirical scoring functions

Regression-based:

$$pKi = \Sigma \, pKi_n \, f_n(\text{structure})$$

affinity    weighting factors    structure descriptors

determined via regression analysis (MLR, PLS)

Data:

Experimental
binding affinities

Experimental
structures

# Where do we stand with docking & scoring?

A not too unusual result ...



Correlation with affinity
for a test set of 800
known complexes:

*in general,*
r < 0.55  (r$^2$ < 0.3)

Wang et al., *J. Chem. Inf.
Comp. Sci.* 44 (2004), 2114

So, what is possible and what is not?

# I. Docking

Preface: „Comparing protein-ligand docking programs is difficult"

Cole et al., *Proteins* 60 *(2005)*, 325

- Test sets needs to be carefully selected to
  - ensure sufficient diversity
  - provide good experimental reliability
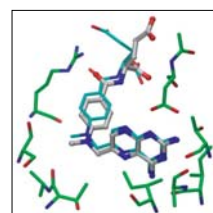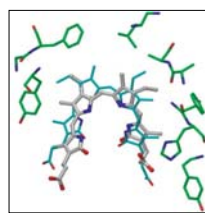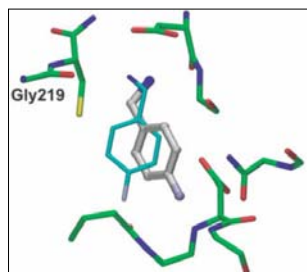  - avoide crystal packing effects
- Consider search complexity and timings
- RMSD values can be misleading
- Tests may cover different aspects, e.g.
  - redocking
  - crossdocking
  - blind docking
  - blind predictions



---

## Comparative evaluations of docking programs

Warren et al., *J. Med. Chem.* 49 (2006), 5912

Success rate for reproducing exp. binding mode on top rank with RMSD < 2Å

Compiled by Moitessier et al., *Br. J. Pharmacol.* 153 (2008), S7

- best approaches typically around 60%
- individual success rates up to 90%
- no approach consistently best
- highly target-dependent

Similar general conclusions by recent studies:

- Cross et al., *J. Chem. Inf. Model.* 49 (2009), 1455
  (68 complexes; DOCK, FlexX, Glide, ICM, PhDOCK, Surflex)
- Li et al., *J. Comput. Chem.* 31 (2010), 2109
  (195 complexes; Glide, GOLD, LigandFit, Surflex)

Figure 4 Co
co-workers, w

## A critical issue: Conformational flexibility!

• Complex reconstruction from **rigid binding partners**: Essentially a solved problem!

  e.g.: RosettaLigand, 85 complexes (Astex diverse): 99% success rate; av. RMSD <1Å

  Davis et al., *J.Mol.Biol.* 385 *(*2009), 381

• **Flexible ligand** – rigid protein docking: Standard, but not without problems

  - docking success rate drops for more flexible ligands (>7-8 rotatable bonds)
  - danger of insufficient sampling (correct conformation and pose is not generated)

• Flexible ligand – **flexible protein** docking: Active field of development

  Modeling of protein flexibility:  • *before* ligand placement (e.g., ensemble docking)

  • *after* ligand placement (e.g., complex refinement)

  • *during* ligand placement (e.g., MC/MD techniques)

  ⇨  But: even „simple" conformational changes can be out of reach!

---

## Example 1: TGT - Successful docking and ...



X-ray confirmed a perfect binding-mode-prediction for a new virtual-screening hit!

## ... surprises out of reach for any docking program



form H-bond with Leu231

**Backbone flip at Leu231**

**and water molecule**

**mediate formation of new**

**H-bond interaction!**

Brenk et al.,
*J.Med.Chem.* 46 *(*2003), 1133

---

## ... surprises out of reach for any docking program



Virtually impossible to predict
with current protein-flexibility
docking approaches
(unless alternative conformation
is experimentally known in advance)

**Backbone flip at Leu231**

**and water molecule**

**mediate formation of new**

**H-bond interaction!**

Brenk et al.,
*J.Med.Chem.* 46 *(*2003), 1133

I. Docking: What is possible and what is not?

## Example 2: Aldose Reductase - docking to multiple pocket conformers



**Sorbinil pocket**

**Tolrestat pocket**

**IDD594 pocket**

*„In-situ"* Cross-Docking
of new pyridazinone inhibitor
to multiple pocket conformations

Zentgraf et al. *ChemMedChem.* 1 (2006), 1355

---

I. Docking: What is possible and what is not?

## Example 2: Aldose Reductase - docking to multiple pocket conformers



**binds to
IDD594 pocket!**

Docking *vs.* X-ray
RMSD = 0.49 Å

Steuber et al., *J. Mol. Biol.* 356 (2006), 45

## Example 2: Aldose Reductase - docking to multiple pocket conformers



**binds to
IDD594 pocket!**

Docking *vs.* X-ray
RMSD = 0.49 Å



Steuber et al., *J. Mol. Biol.* 356 (2006), 45

Reason for successful prediction:

- ligand binds to protein conformer known from previous X-ray structures

- scoring function correctly scores the true binding mode much better than binding modes in alternative protein conformers

⟹   despite protein flexibility:

„easy task" for common docking tools

---

## Design of new inhibitors: Tolrestat analogues



**tolrestat**

**1**

**2**

Da Settimo et al., *J. Med. Chem.* 48 (2005), 6897.
*Naphtho[1,2-d]isothiazole acetic acid derivatives as
a novel class of selective aldose reductase inhibitors.*

**Do the new compounds adopt**

**the same binding mode as tolrestat?**

**I. Docking – Example 2: Aldose Reductase**

Docking of **1** to three different binding pocket conformers, using AutoDock

-10.0 kcal/mol — -9.1 kcal/mol — -9.9 kcal/mol

sorbinil — tolrestat — idd 594



**I. Docking – Example 2: Aldose Reductase**

Docking of **1** to three different binding pocket conformers, using AutoDock

-10.0 kcal/mol

sorbinil

Preferred binding mode of **1**:

- sorbinil-like, not tolrestat-like

- closed specificity pocket

- 4-COO⁻ binds to catalytic site (!)

tolrestat **1** **2**

Docking result of **1** in comparison with crystal structure



• specificity pocket closed

• 4-COO⁻ in catalytic site

***But:***

Unexpected
conformational changes!

• Trp 20 rotated by 35°

• Lys 21 salt bridge broken

• Trp 219 disordered

***Unpredictable***
***with docking methods!***
***(incl. FlexX, GOLD, Glide)***

Docking of **2** to three different binding pocket conformers, using AutoDock



-12.0 kcal/mol          -10.4 kcal/mol          -11.4 kcal/mol

sorbinil          tolrestat          idd 594

Docking of **2** to three different binding pocket conformers, using AutoDock



-12.0 kcal/mol

sorbinil

Preferred binding mode of **2**:

- sorbinil-like, not tolrestat-like

- closed specificity pocket

- 2-COO⁻ binds to catalytic site

tolrestat    1    2

Zentgraf et al., *Angew.Chem. Int. Ed.* 46 (2007), 3575

---

Docking result of **2** in comparison with crystal structure



- specificity pocket closed
- 2-COO⁻ in catalytic site
- no conformational changes!

**But:**

Water molecules

immobilized in binding pocket!

3 very „similar" ligands
lead to
3 very different binding modes!

AutoDock results obtained when using the „correct" binding-site conformer



***Bindung mode exactly reproduced in both cases!***

---

I. Docking: What is possible and what is not?

Protein flexibility and docking

<u>What´s already possible:</u>

• simultaneous docking to multiple protein conformers of arbitrary difference

• correct predictions if multiple protein conformers are known

• support by MD: generation of relevant conformers

• docking with explicit side-chain flexibility

<u>What remains a problem:</u>

• predicting:  - the details

          - backbone mobility

          - large conformational changes

• fast estimate of energetic contributions from protein

• explicit consideration of full protein dynamics upon ligand binding

## II. Scoring

Application tasks:

A) Identification of the correct binding mode for a given ligand
   *Pose prediction in docking*

B) Identification of new active ligands
   *Virtual screening*

C) Affinity ranking for compound series
   *Ligand design, lead optimization*

Available approaches:

- Force field-based methods

- Knowledge-based scoring functions

- Empirical scoring functions



Zn²⁺

---

## A) Pose prediction in docking

Identification of near-native binding pose
among a set of geometric decoys

- Test set of 195 complexes of 65 different targets
- 100 low-energy poses per complex (0-10 Å rmsd)
- 29 scoring functions tested

• **native poses can be detected fairly well**

• **success rates of up to ~80%**

• **knowledge-based approaches work best**



Success rate for identifying
best-scored ligand binding pose
with

- ▮ (yellow) rmsd < 1.0 Å
- ▮ (orange) rmsd < 2.0 Å
- ▮ (blue) rmsd < 3.0 Å

Cheng et al., *J. Chem. Inf. Model.* 49 (2009), 1079

DrugScoreX^CSD    93%

19

# B) Virtual screening

Detection of active compounds in screening databases

Problem: Testing scoring function performance in virtual screening is not trivial!



- significant enrichment can be obtained

- not always for the right reasons

- no function performs consistently well

Compiled by Moitessier et al., *Br. J. Pharmacol.* 153 (2008), S7

---

# C) Affinity prediction

Correlation of scores with experimental binding affinities

and ranking of compounds



Compiled by Moitessier et al.,
*Br. J. Pharmacol.* 153 (2008), S7

- poor correlation for generic data sets

- hardly possible to obtain correct ranking

- of limited use for ligand optimization

## II. Scoring: What is possible and what is not?

*Since all methods are of empirical nature:*

Do more and „better" experimental data

lead to better functions?



---

SFCscore  empirical scoring functions

SFC: Scoring Function Consortium

➡ Data collection from public & industry sources

• affinity data from literature for PDB complexes

• „diversity" from PDB, SAR series from industry

• unique data format and encoding for industry data

• up to 58 complexes per target, 28 series, mostly $IC_{50}$ (!)

➡ Raw data in total (public + industrial):

| | | |
|---|---|---|
| complexes from PDB: | 440 | filtered: 290 |
| complexes from industry: | 618 | filtered: 565 |
| total: | 1058 | 855 |

## SFCscore Training sets: Regression statistics

| Function | Method | $N$ | $k$ | $r$ | $r^2$ | $s$ | $F$ |
|----------|--------|-----|-----|-------|-------|-------|-------|
| sfc_290m | MLR | 290 | 7 | 0.843 | 0.711 | 1.085 | 99.2 |
| sfc_229m | MLR | 229 | 7 | 0.842 | 0.709 | 1.098 | 76.9 |
| sfc_frag | MLR | 130 | 4 | 0.810 | 0.656 | 0.973 | 59.8 |
| sfc_855 | PLS | 855 | 6 | 0.770 | 0.593 | 0.994 | 205.9 |
| sfc_ser | PLS | 466 | 4 | 0.843 | 0.711 | 0.952 | 284.0 |
| sfc_met | PLS | 341 | 4 | 0.844 | 0.713 | 1.046 | 208.9 |
| sfc_290p | PLS | 290 | 5 | 0.867 | 0.751 | 1.005 | 171.3 |
| sfc_229p | PLS | 229 | 6 | 0.875 | 0.766 | 0.982 | 121.2 |

$N$, number of complexes in the training set; $k$, number of components for PLS functions, number of variables for MLR functions; $r$ and $r^2$, correlation coefficient and its square; $s$, standard error; $F$, $F$-value.

Sotriffer et al., *Proteins* 73 (2008), 395

## SFCscore Training sets: Internal cross validation

| Function | $Q^2$ | $s_{PRESS}$ |
|----------|-------|-------------|
| sfc_290m | 0.692 | 1.121 |
| sfc_229m | 0.683 | 1.147 |
| sfc_frag | 0.627 | 1.015 |
| sfc_855 | 0.572 | 1.033 |
| sfc_ser | 0.692 | 1.028 |
| sfc_met | 0.688 | 1.135 |
| sfc_290p | 0.722 | 1.080 |
| sfc_229p | 0.723 | 1.086 |

For the functions derived by MLR, leave-one-out (LOO) cross-validation was used (lines highlighted in italics); for PLS functions, 10-fold cross-validation (20 runs) was applied and the average $Q^2$ and $s_{PRESS}$ of the 20 runs are reported.

Sotriffer et al., *Proteins* 73 (2008), 395

## Comparison with other scoring functions

| | R | R$^2$ | s | F | Q$^2$ | s$_{PRESS}$ |
|---|---|---|---|---|---|---|
| SFCscore: sfc_290m ( k = 7, n = 290 ) | 0.843 | 0.711 | 1.09 | 99.2 | 0.692 | 1.12 |
| X-CSCORE eq3 (Wang 2002): ( k = 4, n = 200 ) | 0.756 | 0.571 | 1.41 | 70.4 | 0.551 | 1.47 |
| Chemscore (Eldridge 1997): ( k = 4, n = 82 ) | 0.843 | 0.710 | 1.40 | 47.1 | 0.658 | 1.52 |
| Score2 (Böhm 1998): ( k = 7, n = 82 ) | 0.890 | 0.792 | 1.27 | 40.3 | | |
| Score1 (Böhm 1994): ( k = 4, n = 45 ) | 0.873 | 0.762 | 1.38 | 32.0 | | |

---

| Function | R$_P$ | SD | ME |
|---|---|---|---|
| SFCscore:: sfc_met | 0.585 | 1.80 | 1.37 |
| SFCscore:: sfc_ser | 0.572 | 1.82 | 1.40 |
| SFCscore:: sfc_855 | 0.570 | 1.82 | 1.40 |
| X-Score::HMScore | 0.566 | 1.82 | 1.42 |
| SFCscore:: sfc_290p | 0.564 | 1.83 | 1.39 |
| SFCscore:: sfc_229p | 0.553 | 1.85 | 1.41 |
| SFCscore:: sfc_229m | 0.534 | 1.87 | 1.44 |
| SFCscore:: sfc_290m | 0.525 | 1.89 | 1.45 |
| SFCscore:: sfc_frag | 0.523 | 1.89 | 1.46 |
| X-Score::HPScore | 0.514 | 1.89 | 1.47 |
| X-Score::HSScore | 0.506 | 1.90 | 1.48 |
| Sybyl::ChemScore | 0.499 | 1.91 | 1.50 |
| DrugScore:Pair/Surf | 0.476 | 1.94 | 1.50 |
| DrugScore: Pair | 0.473 | 1.94 | 1.51 |
| DrugScore: Surf | 0.463 | 1.95 | 1.53 |
| Cerius2:: PLP1 | 0.458 | 1.96 | 1.52 |
| Sybyl:: G-Score | 0.443 | 1.98 | 1.56 |
| Cerius2:: LigScore | 0.406 | 2.00 | 1.57 |
| Cerius2:: LUDI2 | 0.379 | 2.04 | 1.62 |
| GOLD:: GoldScore_opt | 0.365 | 2.06 | 1.63 |
| HINT | 0.330 | 2.08 | 1.65 |
| Cerius2:: PMF | 0.253 | 2.13 | 1.71 |
| Sybyl:: F-Score | 0.141 | 2.19 | 1.77 |

Testing on external data set and comparison with other functions

800 PDB complexes with exp. pK$_i$

Wang et al., *J. Chem. Inf. Comp. Sci.* 44 (2004), 2114

⟹ improvement, but still only moderate correlation

---

New, carefully compiled test set of 195 PDB complexes with exp. pK$_i$:

Cheng et al., *J. Chem. Inf. Model.* 49 (2009), 1079

Best functions:

| | R$_P$ | SD |
|---|---|---|
| SFCscore:: sfc_met | 0.646 | 1.82 |
| X-Score::HMScore | 0.644 | 1.83 |

Why did many functions in the past appear more successful than they are?

⇒ Very small external test sets of limited diversity

cf. how many of the now available complexes are well predicted by SFCscore!

| | | Residual $< 1.5$ | | |
| --- | --- | --- | --- | --- |
| Function | N | $R_P$ | $r^2_{pred}$ | $SE_{pred}$ |
| sfc_290m | 551 | 0.874 | 0.763 | 0.809 |
| sfc_229m | 546 | 0.879 | 0.769 | 0.803 |
| sfc_frag | 417 | 0.915 | 0.818 | 0.835 |
| sfc_855 | 555 | 0.850 | 0.720 | 0.820 |
| sfc_ser | 558 | 0.876 | 0.765 | 0.806 |
| sfc_met | 553 | 0.872 | 0.759 | 0.790 |
| sfc_290p | 559 | 0.875 | 0.765 | 0.796 |
| sfc_229p | 531 | 0.887 | 0.784 | 0.790 |

⇒ CAVE with any conclusions derived from too small test sets!

---

Have the limits of empirical approaches been reached?

⇒ Consider quality and comparability of experimental data!

• Structural data (mainly X-ray) of protein-ligand complexes

- multiple conformations (highly dynamic systems)

- hydrogen atom positions (protonation states) not observable

- side-chain orientation may be ambiguous (Asn, Gln, His)

- water molecules are only partially observable

- binding modes may depend on crystallization conditions and crystal packing

• Affinity data of protein-ligand complexes

- may highly depend on pH, buffer, salt concentration, temperature

- enyzme kinetics: inhibition mechanism must be known

- $IC_{50} \leftrightarrow K_i \leftrightarrow K_d$

Knowledge-based and empirical scoring methods
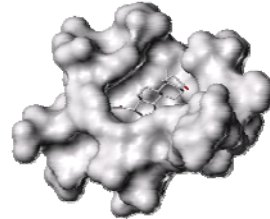cannot be better than the exp. data they are based on!

II. Scoring: What is possible and what is not?

## Have the limits of empirical approaches been reached?

⟹ For the development of generic scoring functions:

Problems difficult to overcome, even by concerted efforts!

⟹ Focus on target- or target-class-specific functions!

• Target-specific adaptation of existing functions

• Better comparability of experimental data

• Definition of standards for acquisition of new affinity data possible

## Recommendations ...

... for approaching the scoring problem:

1) Validate the scoring function for your system of interest

2) Train the scoring function for your system („Tailored scoring function")

3) Try applying multiple scoring functions („Consensus Scoring")

4) Tackle the problem with additional pre- and postfiltering steps

## Further developments required ...

... to overcome the most serious simplifications in scoring functions:

> **„Flexibility – Water – Entropy"**

- single configuration of the binding partners in the complex

- no consideration of the unbound state

- no or simplified consideration of the solvent

- focused on enthalpic contributions and interaction descriptors

- additivity of interaction terms

*A single model may not be sufficient to capture*
*the complex interplay of*
*residual mobility, desolvation, and interaction quality in protein-ligand complexes!*
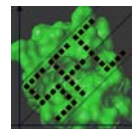
## The wrong conclusion ...

## Acknowledgement



UNIVERSITÄT WÜRZBURG

**David Zilian**
Daniel Cappel
Monika Nocker
Ulrich Peinz
Benjamin Schaefer
Martin Sippel
Christine Topf
Constanze Waltenberger
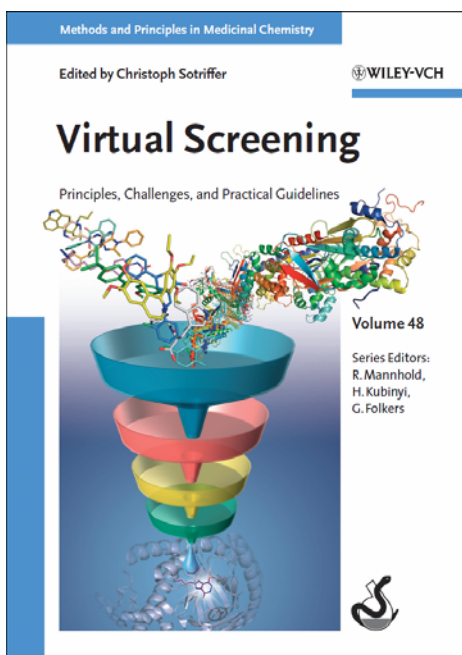Armin Welker

© Rolf Nachbar

**S**coring
**F**unction
**C**onsortium

| Astra | Aventis | |
|-------|---------|------|
| BASF | Boehringer | |
| Glaxo | Novo Nordisk | |
| Pfizer | Agouron | |
| Roche | Schering | CCDC |

Hans Matter  (Sanofi-Aventis)

Gerhard Klebe (Univ. of Marburg)
Paul Sanschagrin

---

Methods and Principles in Medicinal Chemistry

Edited by Christoph Sotriffer

WILEY-VCH

# Virtual Screening

Principles, Challenges, and Practical Guidelines

Volume 48

Series Editors:
R. Mannhold,
H. Kubinyi,
G. Folkers

available in autumn 2010