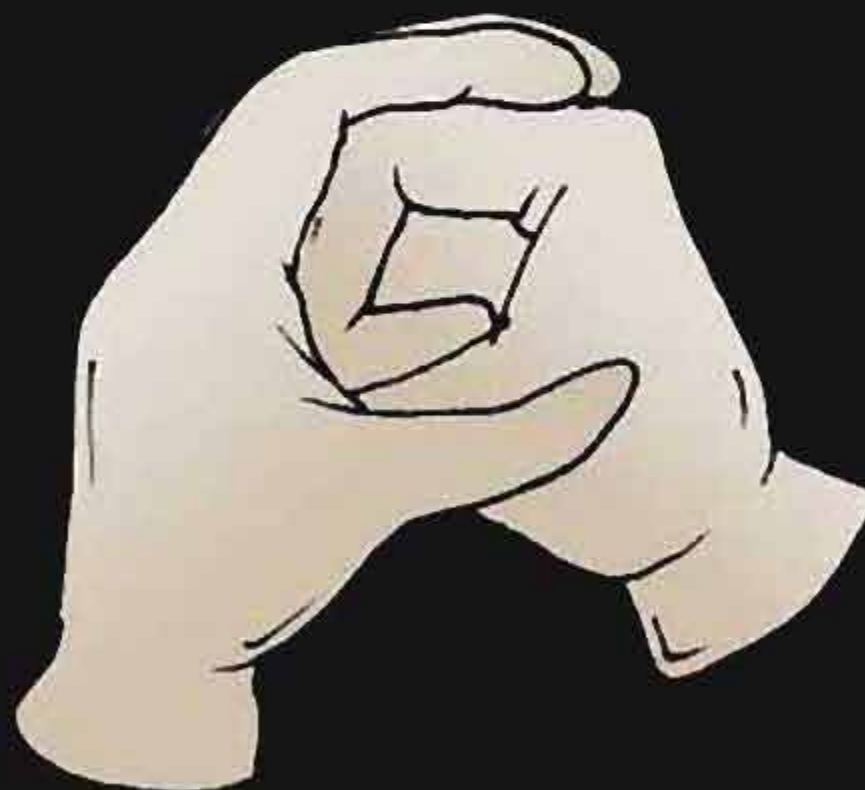## J3.1.1 Similarity and Complementarity-Based Drug Design

Ligand-based and structure-based drug design are the two major approaches in rational drug design. Ligand-based drug design exploits the likeness between molecules that is expressed in terms of "molecular similarity" whereas structure-based drug design exploits the detailed 3D recognition features between a ligand and its receptor, and the concept used is "molecular complementarity". The molecular similarity concept and its applications are discussed in detail in this chapter.

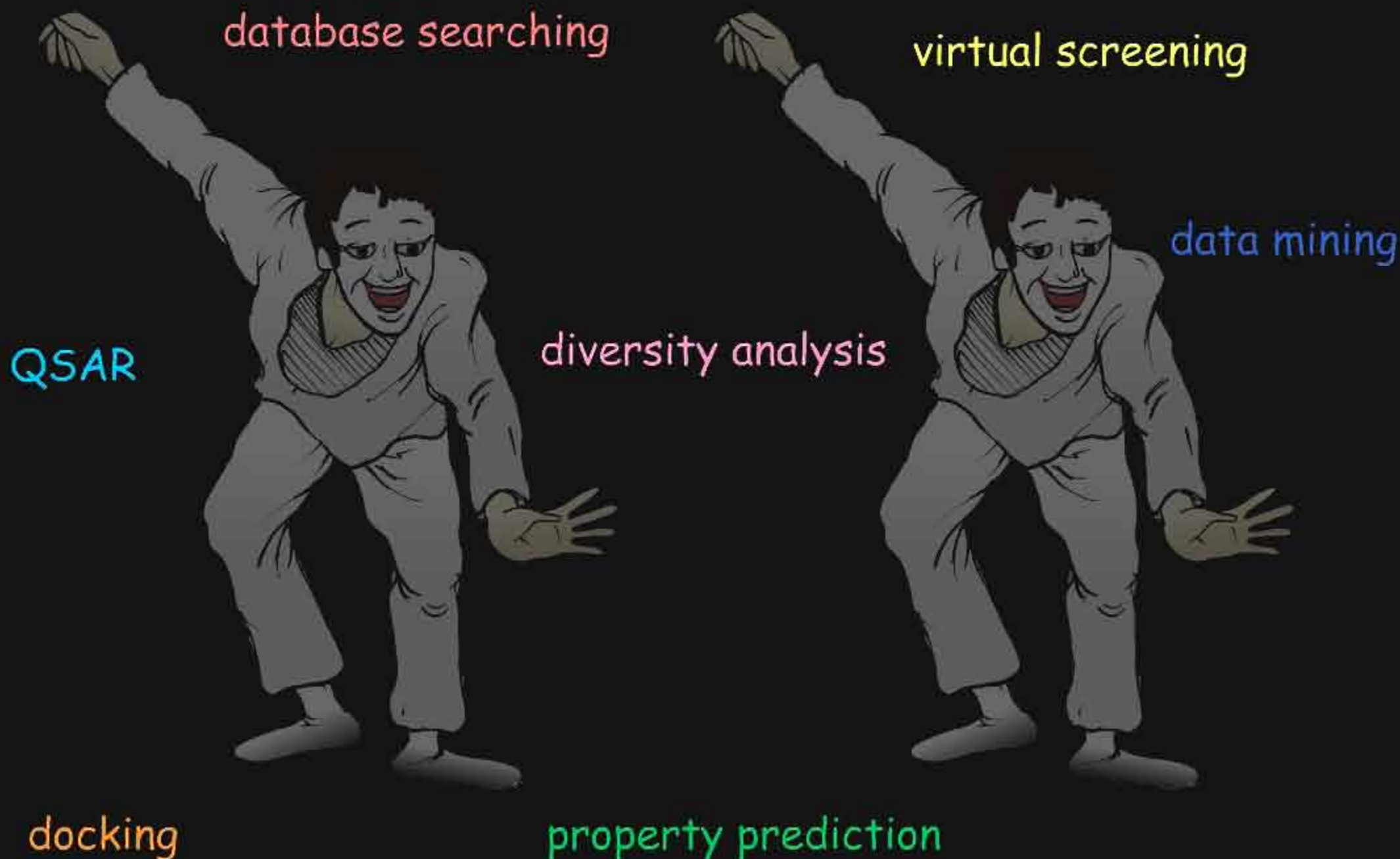ligand-based drug design          structure-based drug design

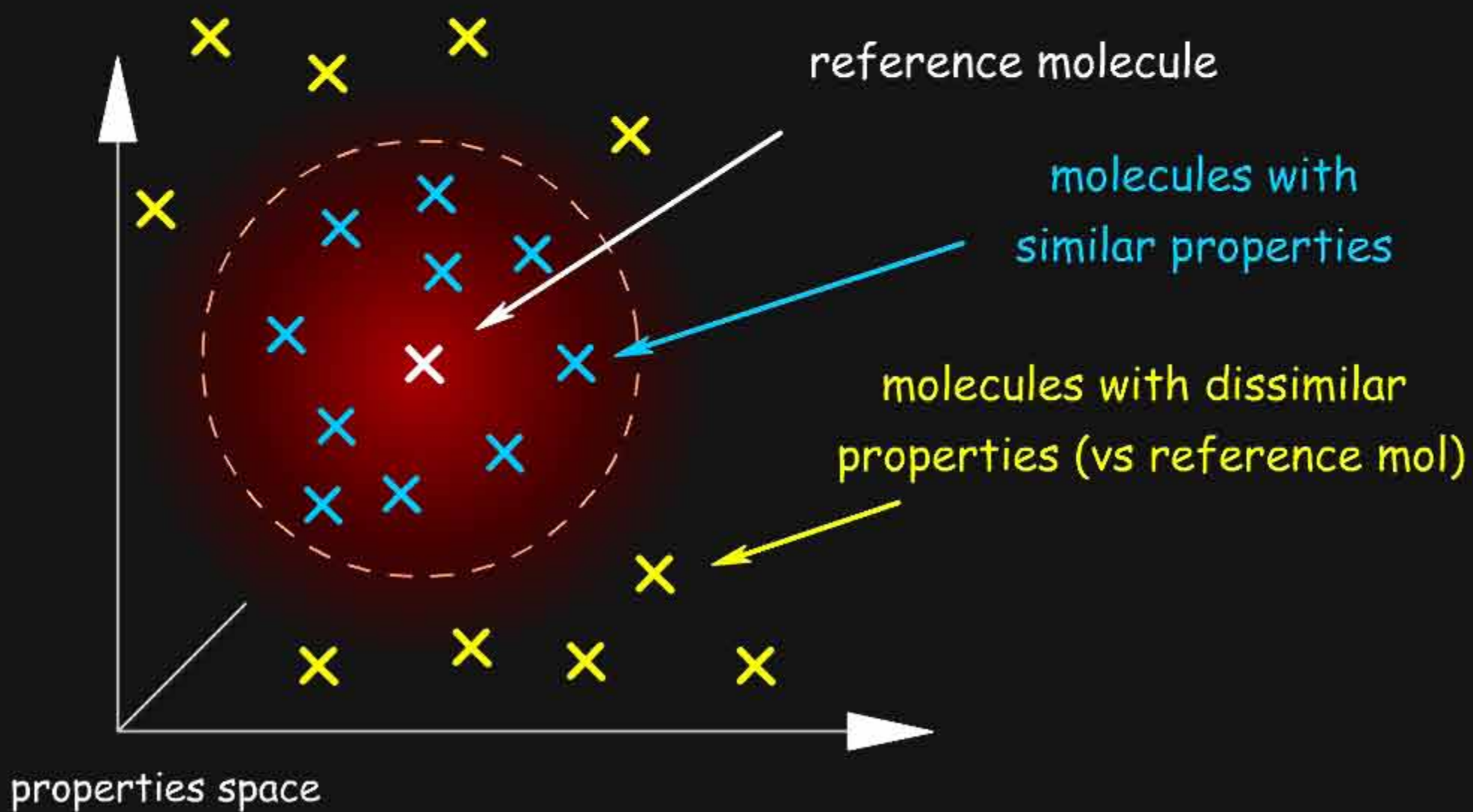molecular similarity ?          molecular complementarity ?

## J3.1.2 Comparing Molecules: a Central Issue in Drug Discovery

Assessing the similarity between molecules is a central issue in drug discovery. The molecular similarity concept has created a broad range of cheminformatics tools that have proven useful in drug design for finding new lead compounds.
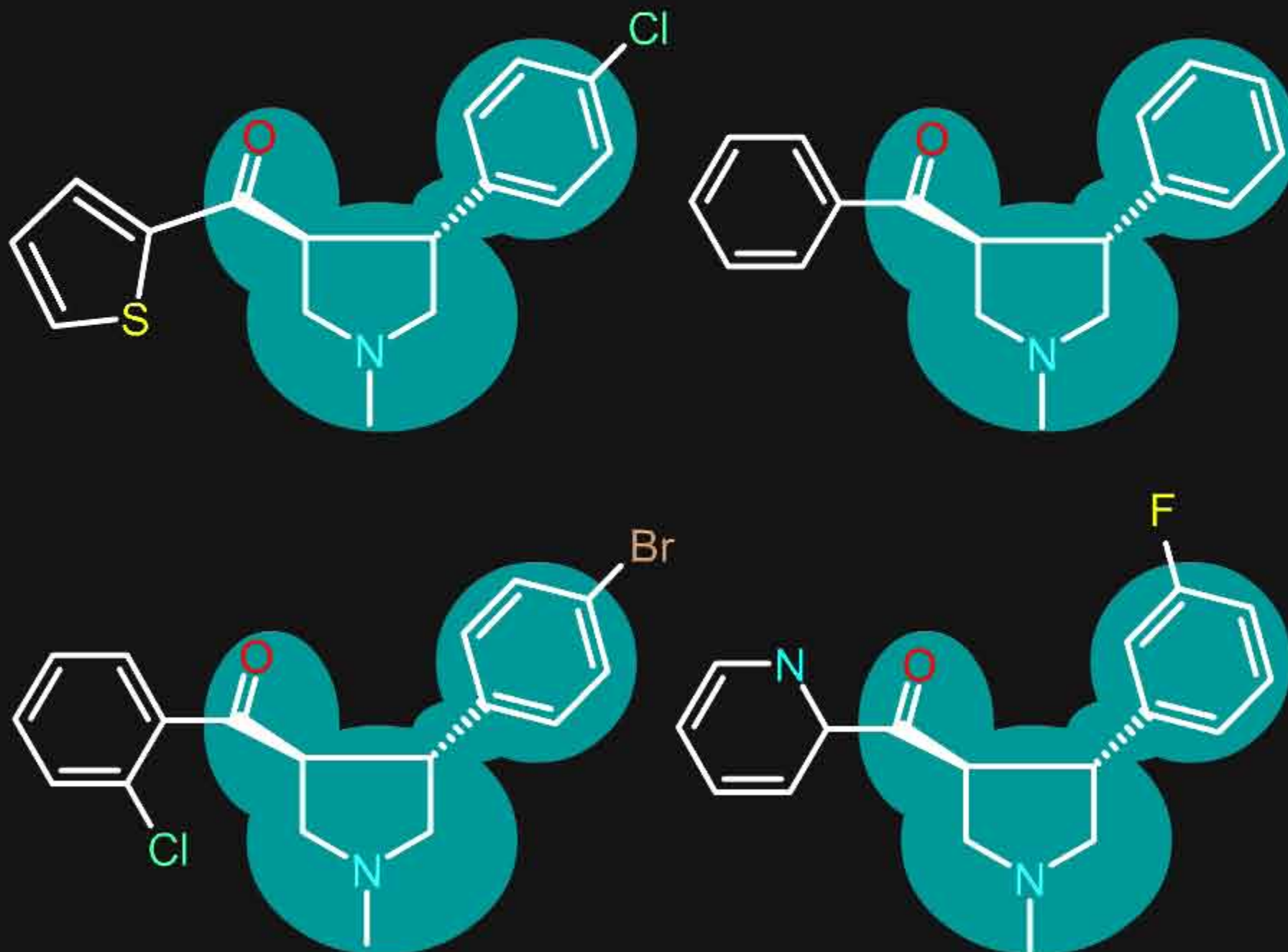
# J3.1.3 The Molecular Similarity Principle

The "molecular similarity principle", which is also known as the "similar property principle", is the underlying concept of virtually all ligand-based drug design methods. Its underlying assumption states that similar molecules tend to behave similarly, while more dissimilar molecules exhibit more distinct properties.

similar molecules ➡ similar properties

dissimilar molecules ➡ dissimilar properties

reference molecule

molecules with similar properties

molecules with dissimilar properties (vs reference mol)

properties space

# J3.1.6 2D-Structure Similarity

Since chemists are very familiar with structural chemical formulas, a straightforward way to compare structures is based on their 2D connectivity i.e. the kinds of atoms which are bonded to each other, in which topology. In the structure below there are four structures which present identical bioactivity, and they possess very similar 2D structures.
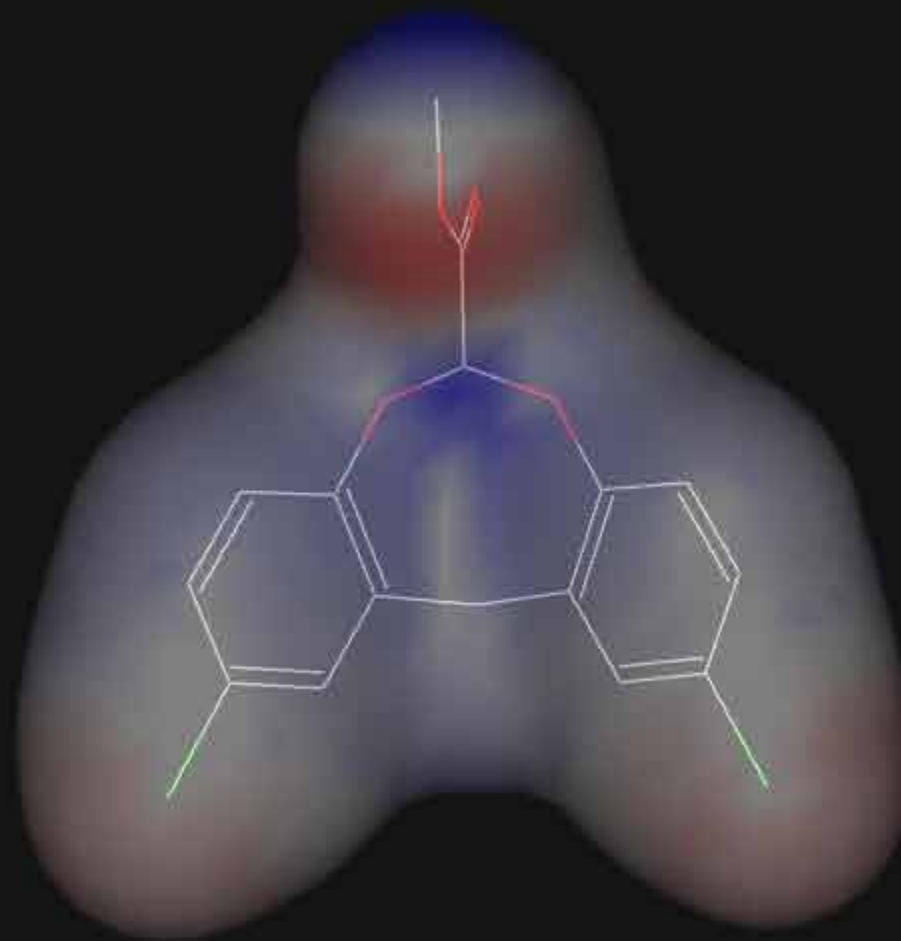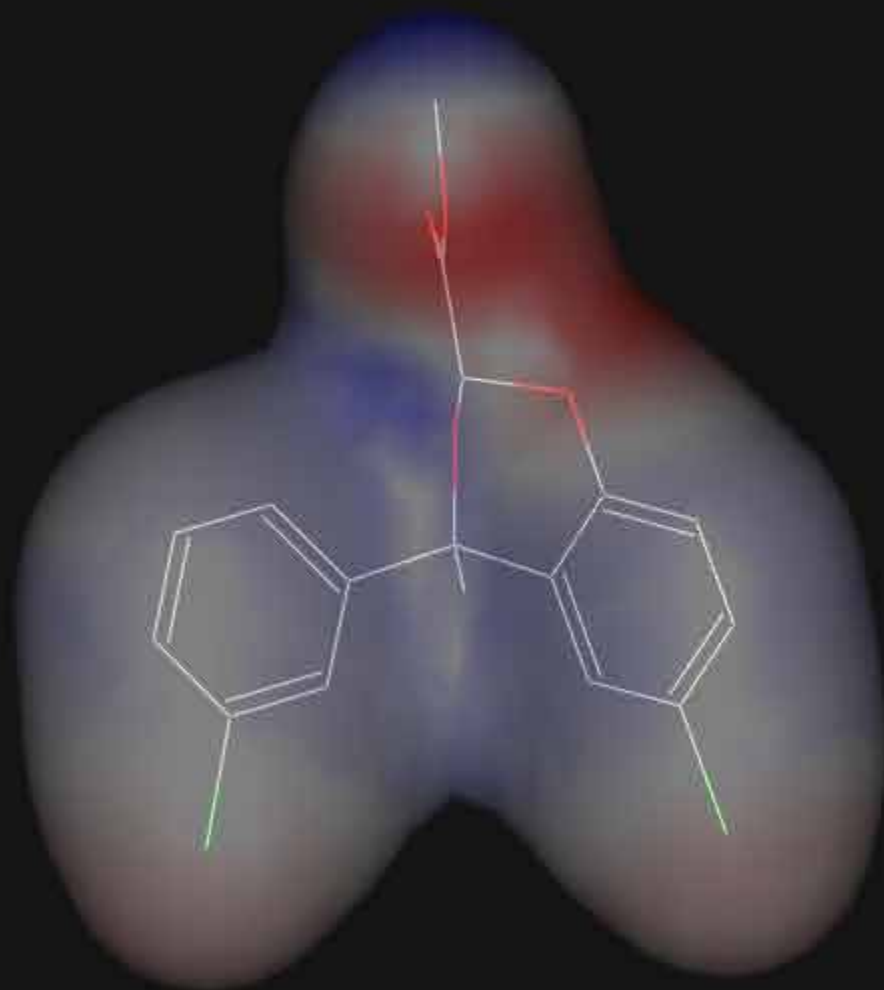
## J3.1.7  Shape Similarity

Another way to define similarity between molecules is shape similarity. Molecular shape is an important determinant in the biological activity of a molecule. Below two structures are shown which look very different with respect to their 2D similarity. Still, if they are rendered in three dimensions you can see that their shapes are similar, which accounts for their similar biological profile.

⦿ Shape                    ⦿ 2D
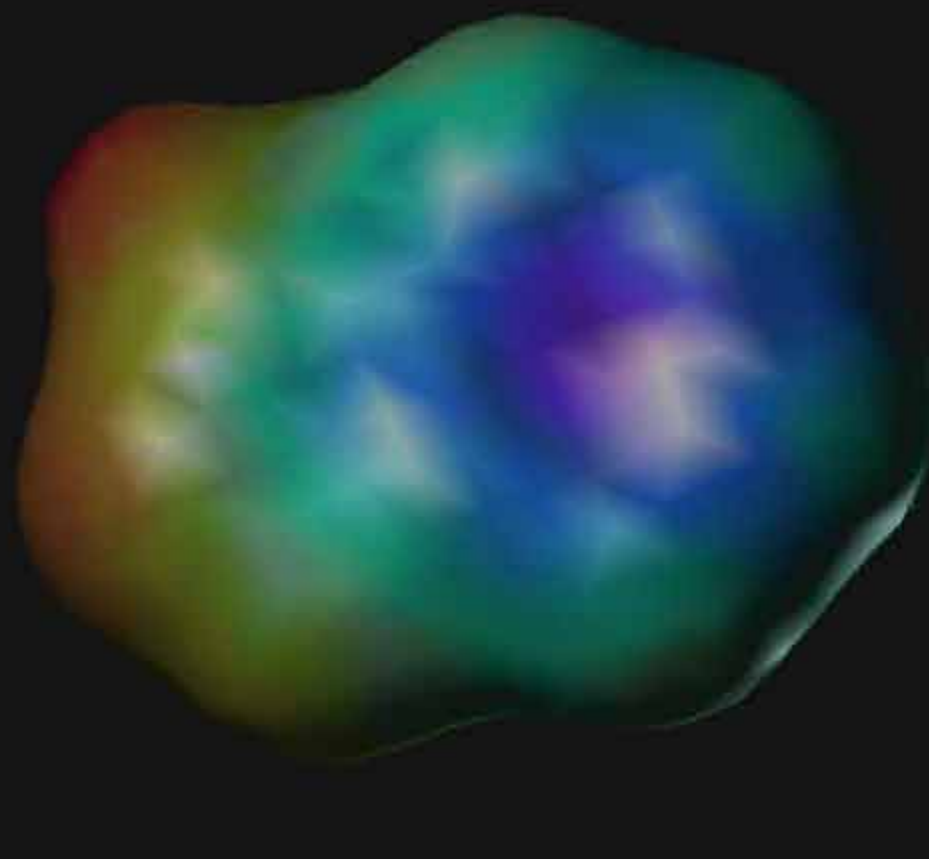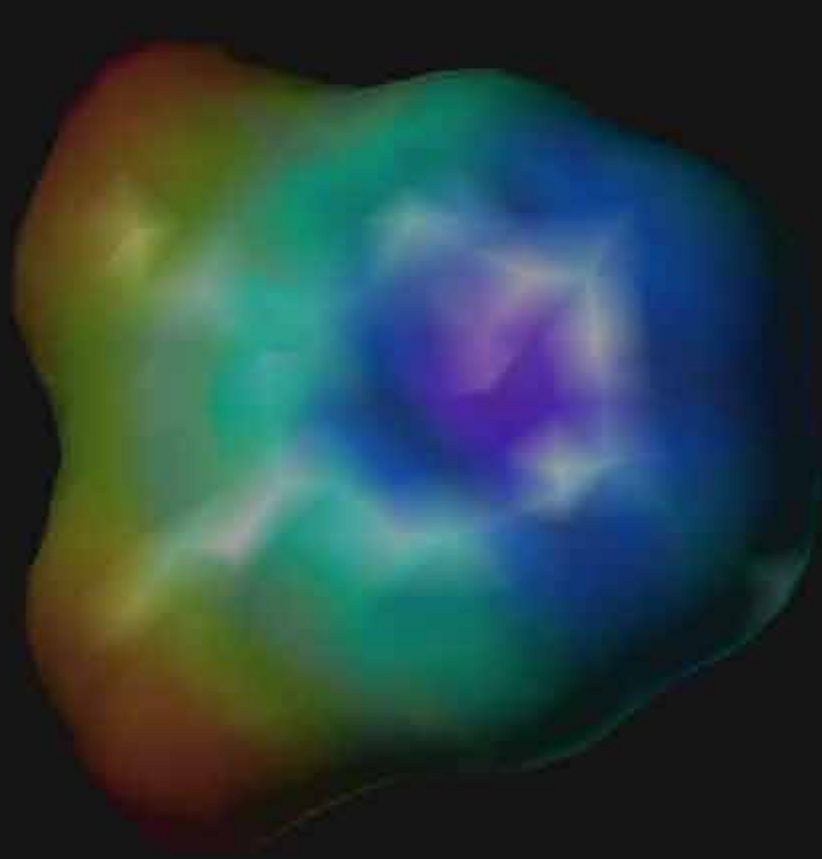


RU25961                    Treloxinate

# J3.1.8  Surface Physicochemical Similarity

Apart from shape, surface property similarity is very important. Properties such as atomic charges, electrostatic potentials, hydrophobicity, polarizability can be represented and compared on the molecular surfaces. Below you can see the surface electrostatic potential of two molecules: a catechol and its structurally quite dissimilar bioisosteric replacement containing a second nitrogen heterocycle instead of the two original hydroxyl groups. Despite the dissimilarity with respect to their 2D representations, both structures show similar electrostatic potential on the surface, which in turn might result in similar bioactivity.

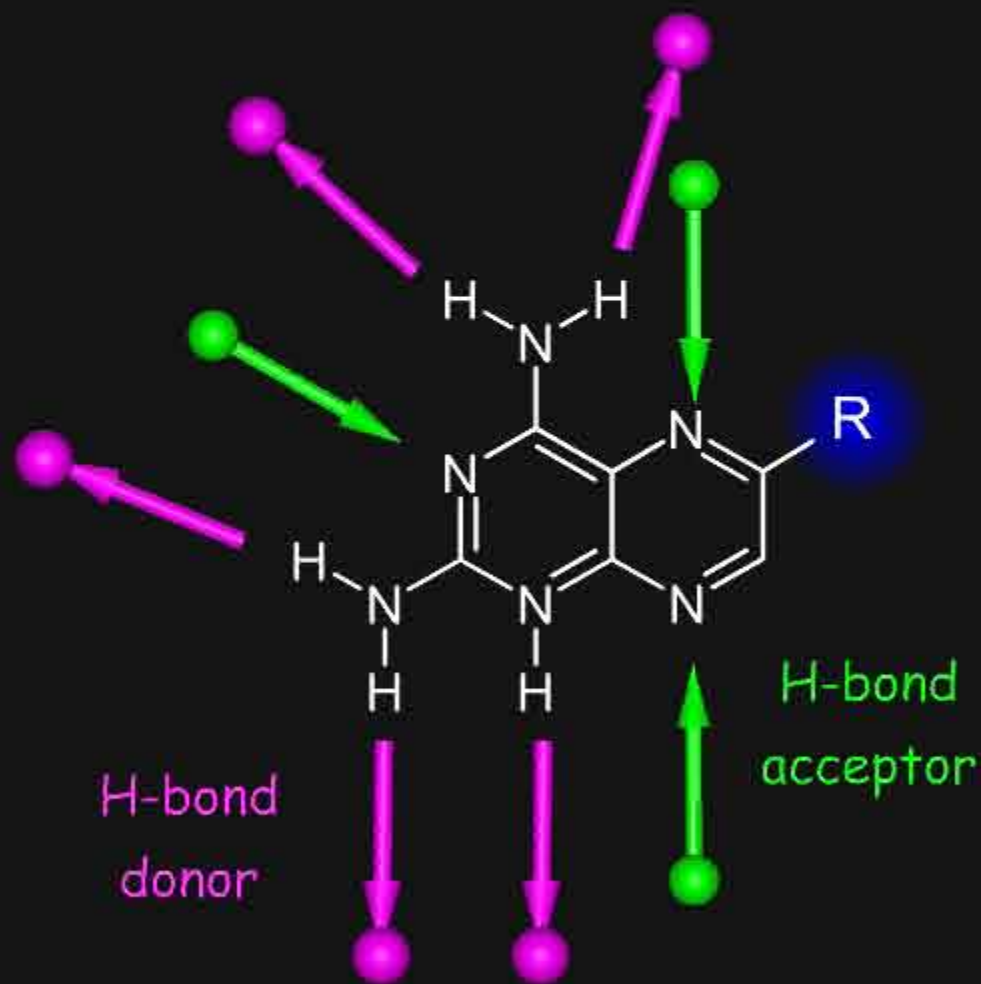● 2D                                                    ◉ Surface
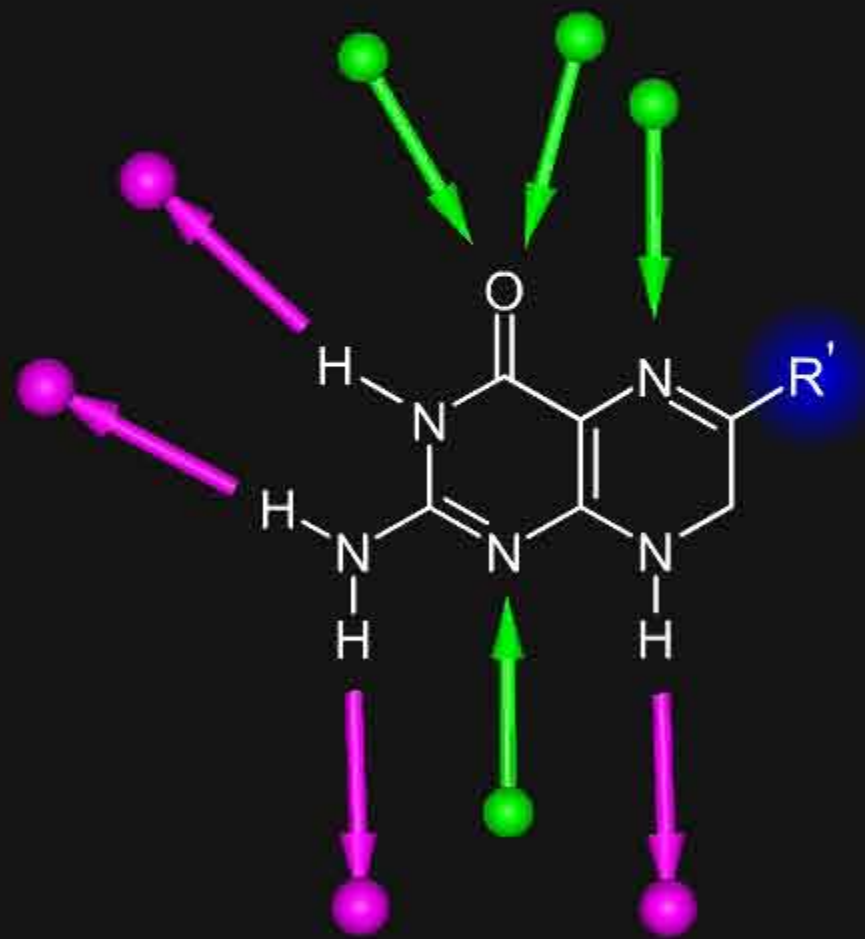
# J3.1.9  H-Bond Similarity

Since the hydrogen bonds are generally essential for the selectivity of a ligand, molecules can be compared in terms of their hydrogen bond pattern similarity. Both methotrexate and dihydrofolate bind to dihydrofolate reductase; by looking at their chemical structures it is reasonable to assume that they bind in an orientation as indicated in '2D alignment'. However, X-ray studies reveal that each molecule binds upside down, relative to the other, which makes sense when looking at their H-bond similarities (see 'H-bond alignment').
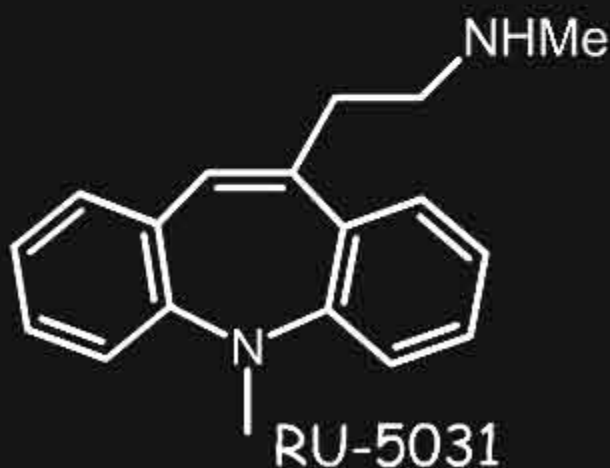
○ 2D alignment
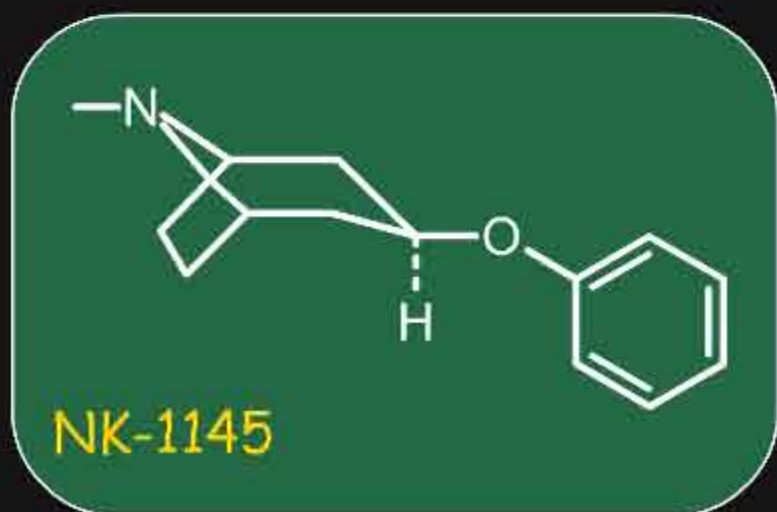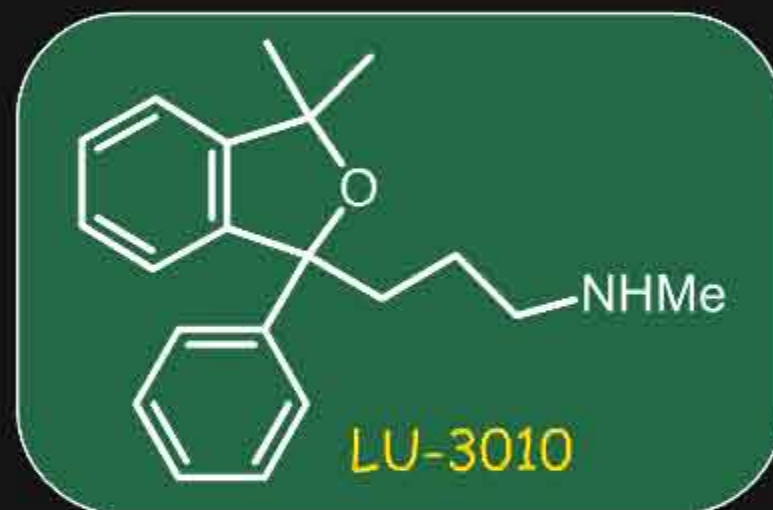
◉ H-bond alignment



H-bond acceptor

H-bond donor

methotrexate

dihydrofolate

## J3.1.10 Absence of Particular Features

Are a turtle and an elephant similar because neither of them do not have wings? While this is partly a philosophical question it needs to be mentioned since there are molecular similarity methods that, in addition to judging similarity based on features that indeed do match, also take features into account which do not match. The molecules shown below are all antidepressants. NK-1145 and LU-3010 do not have the tricyclic moiety normally found in the other structures, so in this context they are similar.



Imipramine

NMe₂

Amitriptyline

NMe₂

LU-3010

NHMe

NK-1145

RU-5031

NHMe

Protriptyline

NHMe

# J3.1.11 Pharmacophore Similarity

The pharmacophore is the three-dimensional (spatial) arrangement of ligand features responsible for ligand-target interactions. Comparing molecules in terms of their pharmacophore pattern focuses on considering only the essential parts of these molecules. In the example below the pharmacophore consists of a fluoro-phenyl, an amino group and a carboxyl moiety.

○ 2D                                                        ○ Superimposition

# J3.1.14  Relevant Characteristics: What is Important?

The ideal characteristic used for a comparison should be like an ideal witness in court, who states 'the truth, the whole truth, and nothing but the truth". In other words, the characteristic captures all the relevant aspects of the property we are attempting to predict ("the whole truth"), but at the same time does not add noise ("nothing but the truth"). Specific descriptors can be good or bad, depending on the situation.

- Relevant property

- Capture the truth

## Molecule-1

- 2D structure
- Shape
- H-Bond pattern
- Pharmacophore
- Physico-chemical
- Electronic property
- etc...

## Relevant Properties

## Molecule-2

- 2D structure
- Shape
- H-Bond pattern
- Pharmacophore
- Physico-chemical
- Electronic
- etc...
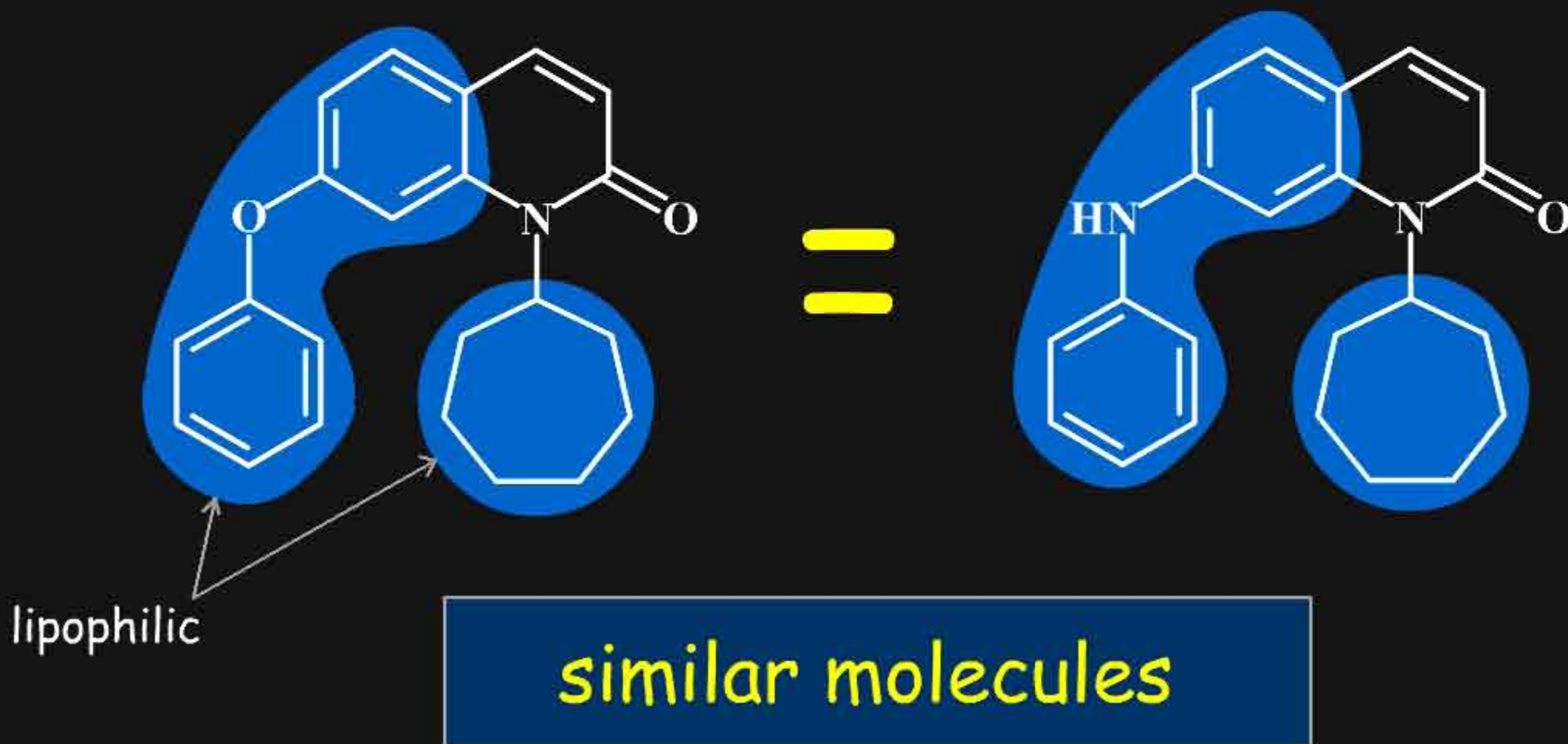
# J3.1.15 Relativity of Relevant Properties

In molecular similarity calculations the relevant descriptors differ from case to case. In the example below: in terms of lipophilicity the replacement of an oxygen linker (–O–) by a secondary amine (–NH–) does not introduce major changes; however this modification may have radical repercussions if the group is involved in specific hydrogen bond interactions with the receptor.

○ Lipophilicity context                    ◉ H-bond context
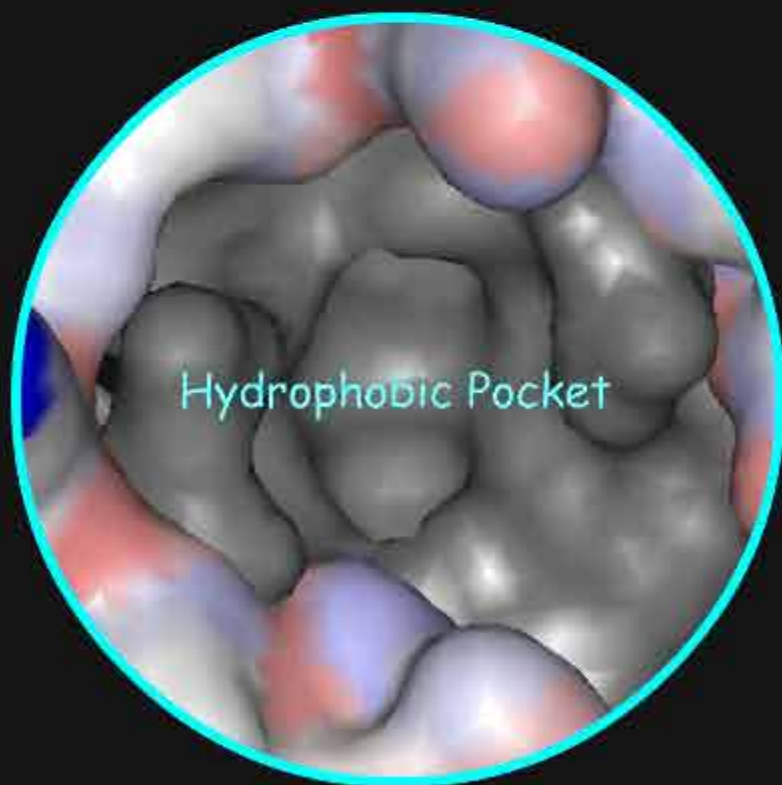


lipophilic

**similar molecules**

## J3.1.16 Interpretable Characteristics

It is a good rule of thumb to aim for properties that are interpretable because they provide insights into the exact content of the molecules at hand. However they only capture what chemists are trained to capture. Hard-to-interpret descriptors try to capture alternative aspects of the molecular structure that are elusive to the human mind, but can be related to observable properties.



Lipophilicity

Hydrophobic Pocket

Magnetic Susceptibility

$\chi$ ?

## J3.1.17  Global and Local Characteristics

Some molecular characteristics are 'global' and provide a very short description of the molecule, as for example the LogP. Other characteristics can be 'local', describing the properties of some regions, fragments, atoms or even a point in space. Local regions can be compared individually, enabling local similarities.
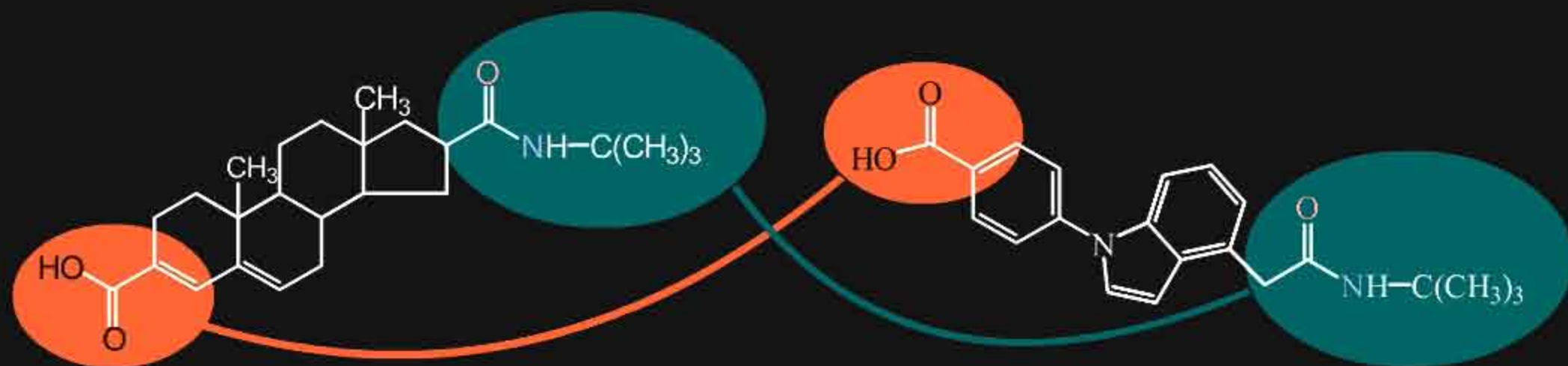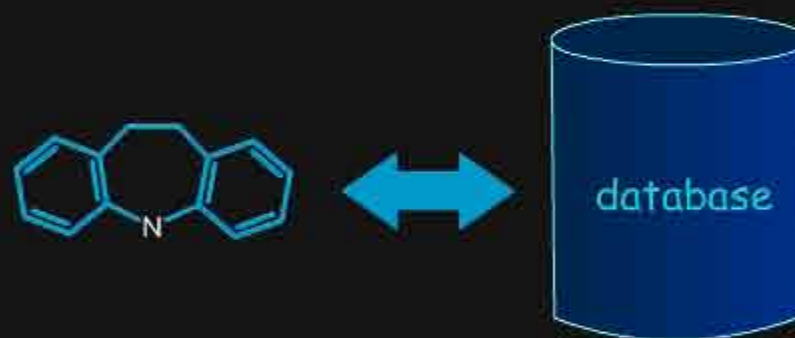
# J3.1.21  Cheminformatics

With the advent of computers, the definition of "molecular similarity" became more explicit, and two major approaches have emerged. The first is involved in similarity searching as the database implementation of the similarity concept. The second deals with the development of models for predicting molecular properties. Each of these approaches will be discussed in detail in the present chapter.

molecular similarity in cheminformatics

1. similarity searching

2. development of models    $P = f\left( \ \right)$

In the previous section we saw how to handle molecular structures and properties in the computer by using molecular descriptors. Now we need to establish a numerical similarity measure between molecules which assigns a single number (the "similarity index") to structure pairs.

# J3.6.2 Similarity Coefficients of Relevant Properties

As already mentioned, similarity is a subjective concept that requires relevant descriptors that can represent the properties to be compared (e.g. biological activities). A poor choice of the relevant physicochemical properties or descriptors leads to similarity values that are meaningless.

# J3.6.3  Binary and Distance-Based Formulas

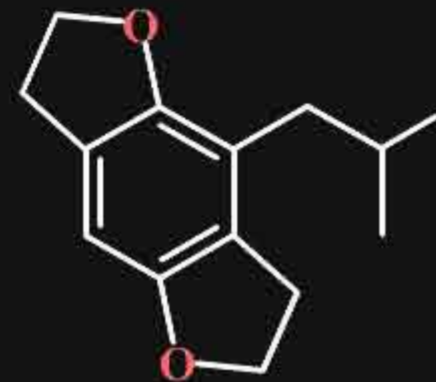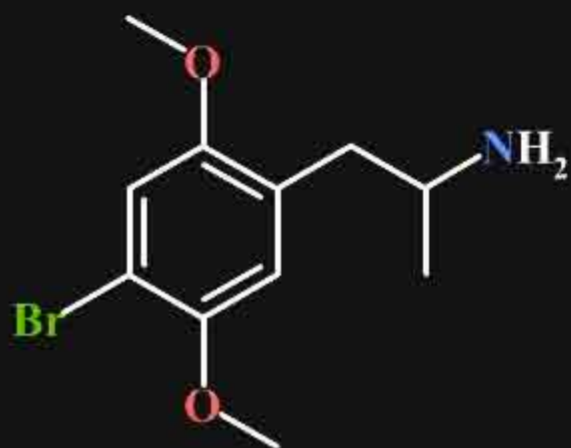Depending on how molecules are encoded (using either binary strings or numerical descriptors), the similarity indexes are calculated from either binary or numerical values. They are assigned the same name (e.g. Tanimoto) but have different mathematical forms. When the calculation are based on binary values the numerical similarity measure is called the "Similarity Coefficient" and when the calculations are based on numerical values it is called the "Distance Coefficient".

Binary Descriptor

Numerical Descriptor

| 1 | 0 | 0 |

| 0.3 | 5.2 | 51 |

Formula for binary variables

Formula for continuous variables

$$\text{Tanimoto (binary)} = \frac{C}{A + B - C}$$

$$\text{Tanimoto (numerical)} = \frac{(x_1 x_2) + (y_1 y_2)}{(x_1^2 + x_2^2) + (y_1^2 + y_2^2) - (x_1 x_2) + (y_1 y_2)}$$

Similarity Coefficient

Distance Coefficient

# J3.6.9 Examples of Similarity and Distance Coefficients

In the following slides we will only discuss the most popular similarity and distance-based similarity coefficients.

- The Tanimoto coefficient

- Dice and Cosine

- Tversky coefficient

- Euclidiean

- Hamming

Due to their binary nature, the distribution of similarity coefficients usually does not tend to be random. Simple ratios such as 1/2, 1/3 and 2/3, etc. are much more frequent than other ratios with larger numerators and denominators. For small fingerprints it was shown that 1/3 is the most frequent similarity value for an infinitely long bit string which is half occupied, a value that should be kept in mind to have a feeling for similarity values.

## J3.10.9 Size-Bias: Favoring Large Molecules

Similarity coefficients can usually not be compared directly; the size of the query has a major influence on the distribution of similarities. In the following figure is illustrated the distributions of Tanimoto coefficients between a query and all molecules of a database. The larger the query, the higher the average similarity, and the wider the distribution of values. This is also known as a 'size-bias' of the Tanimoto coefficient: in similarity calculations larger compounds are favored.
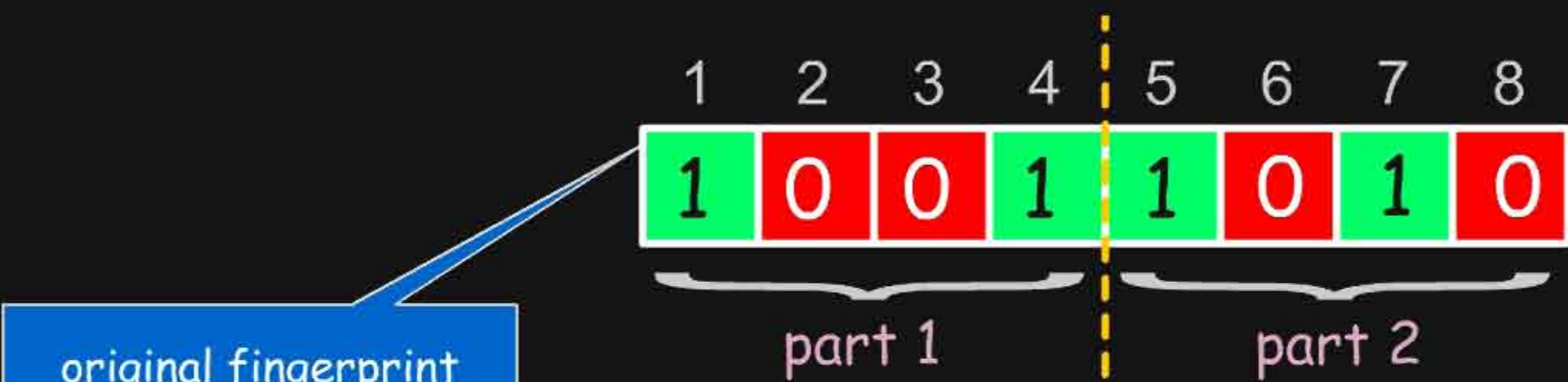
# J3.10.5 Folding of Fingerprints

It is sometimes useful to compress the information encoded in long fingerprints. In this case, the whole fingerprint is said to be "folded". It is cut into two halves, which are combined by a logical OR (or other) operation to give a new fingerprint of only half the size of the original fingerprint. In the resulting fingerprint, bits are set at position n if they were present in position n of the original fingerprint OR position (length/2)+n (which corresponds to bit n of the "second half" of the fingerprint). This process can also be repeated.
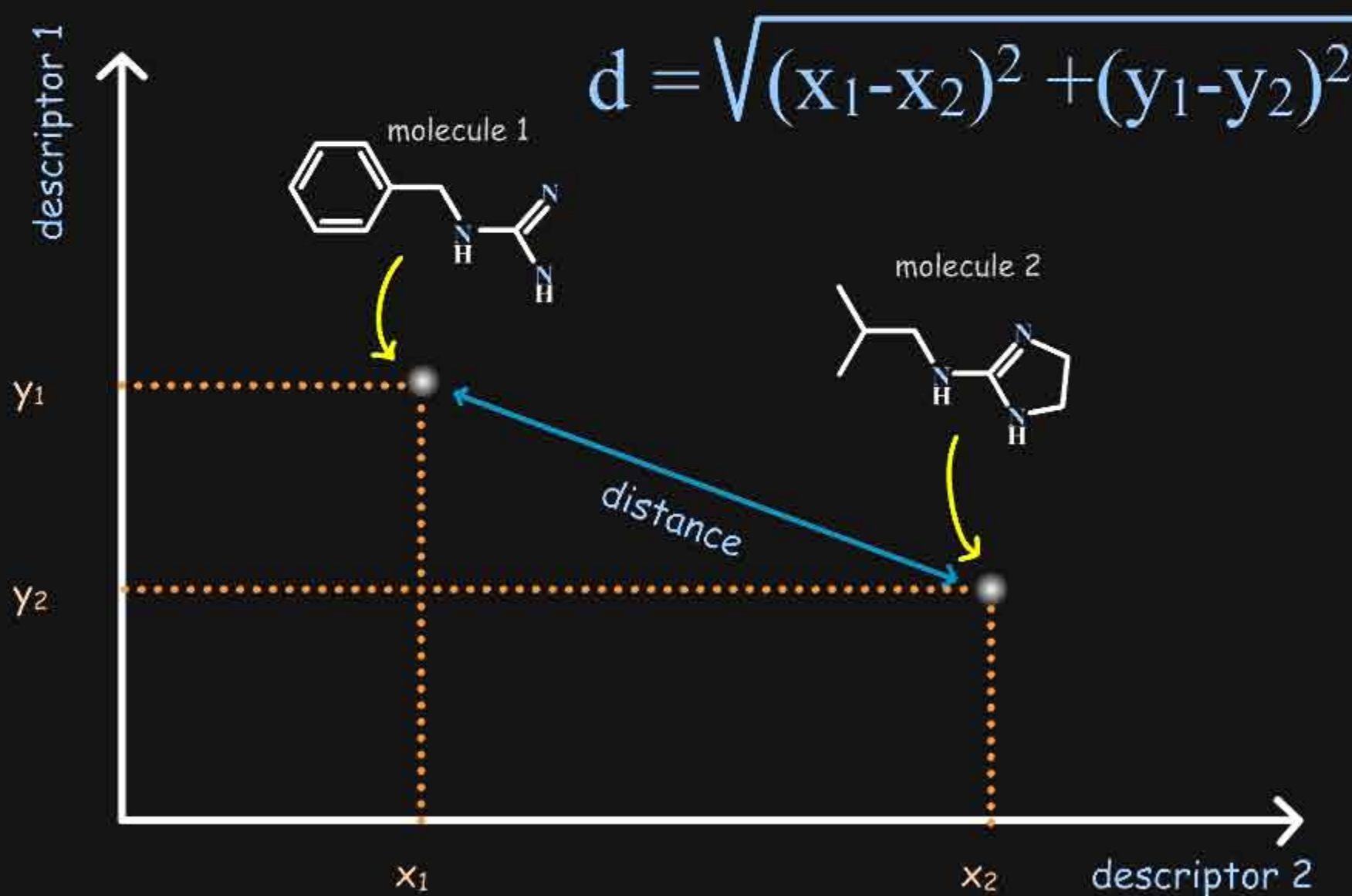
# J3.6.4 Distance Coefficients

Analogous to establishing the distance of objects in space, distance coefficients establish the distance between two descriptor representations of molecules. Distance coefficients range between 0 and infinite. One well-known distance coefficient is the Euclidean Distance, which treats the descriptor entries as dimensions in Euclidean space and is calculated as given below.



identical

0 < distance coefficient < infinite

dissimilar

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2}$$

descriptor 1

molecule 1

molecule 2

$y_1$

distance

$y_2$

$x_1$

$x_2$

descriptor 2

While similarity coefficients can be defined in many ways, one question is whether the coefficient behaves symmetrically, in other words is "A is similar to B" as "B is similar to A".
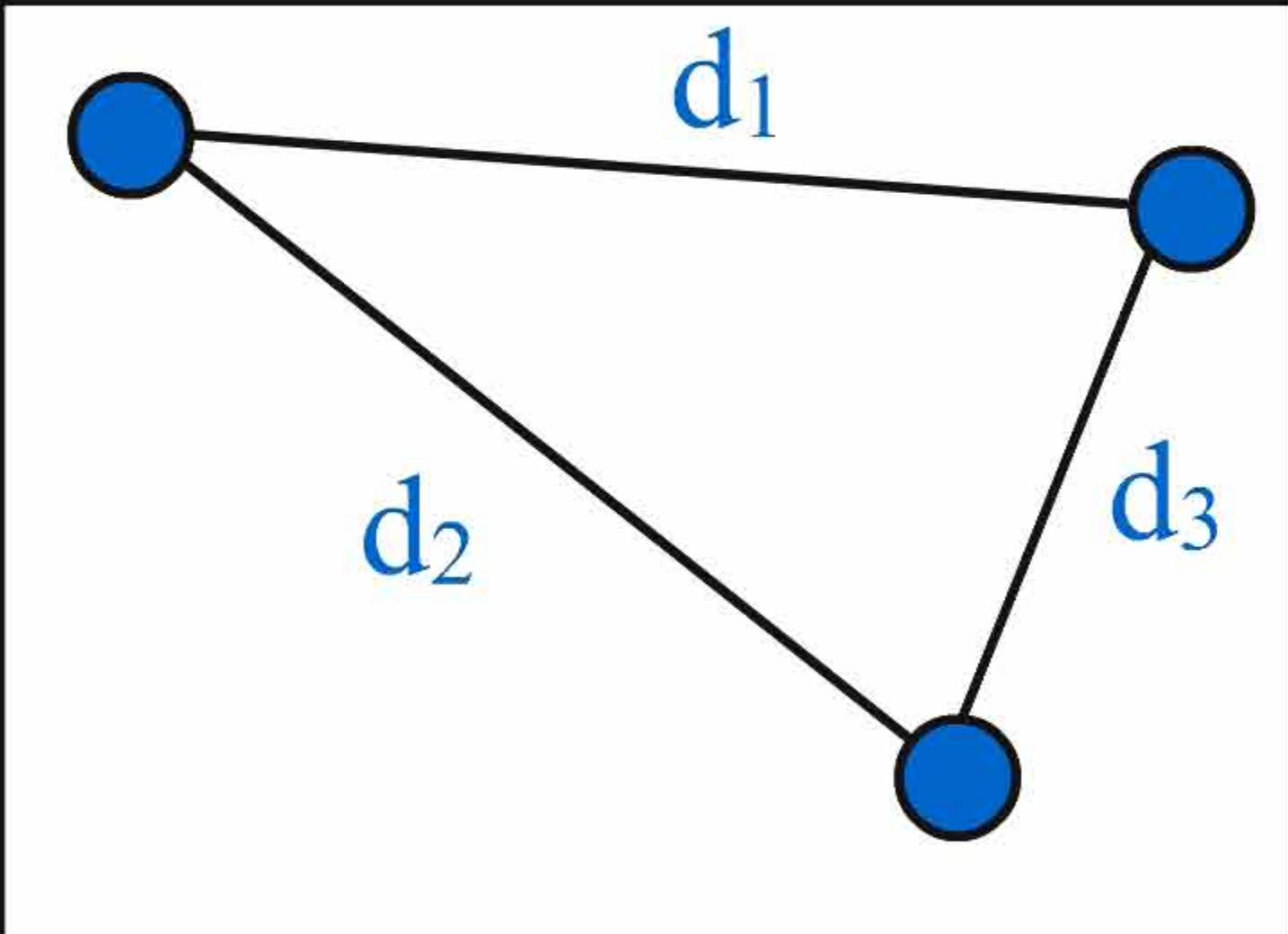
A

B

same distance ?

# J3.10.6 The Concept of Binning

If descriptors of molecules involve the distance of intermolecular features in space and one wants to determine whether identical (or equivalent) features are present in two molecules, exact matching of a particular distance virtually never occurs. In order to still be able to compare spatial arrangements of features a concept called "binning" is often employed, where ranges of intra-feature distances are assigned to the same "distance range", also called a "bin".
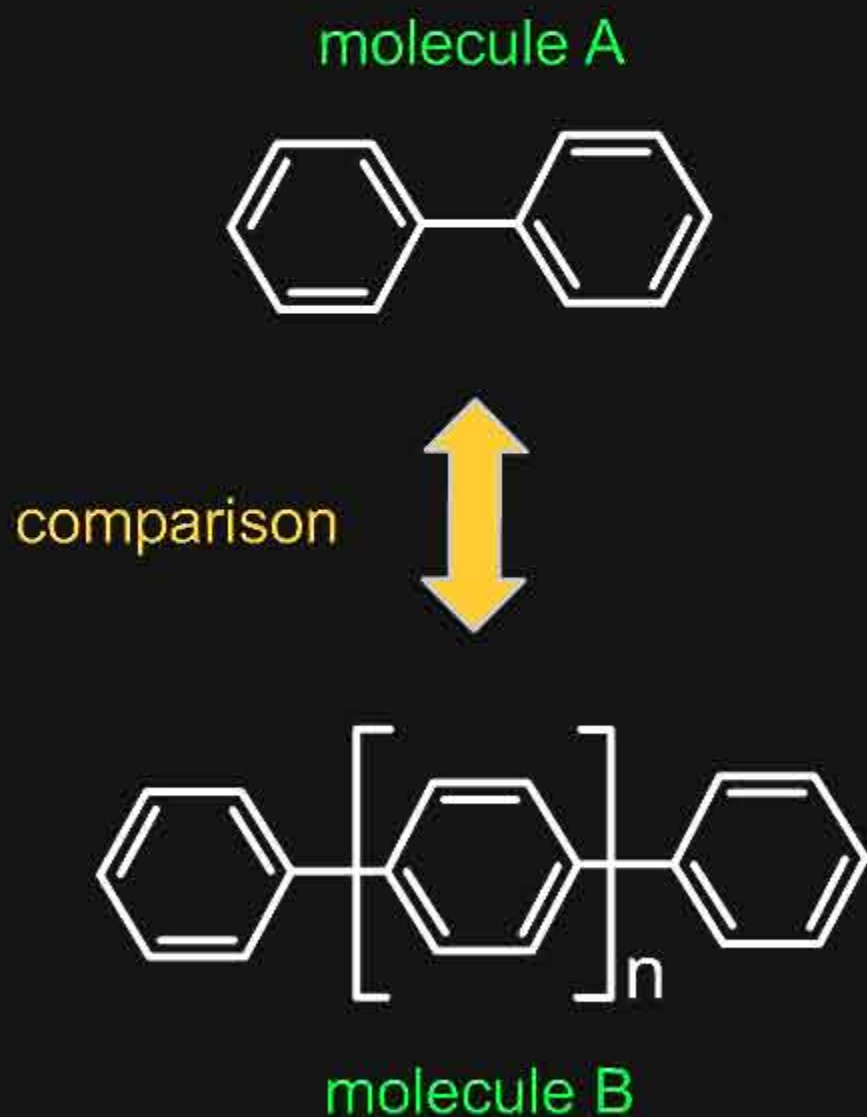
## J3.10.8 Size-Bias of the Tanimoto Similarity Coefficient

Both the molecular representation in descriptor space and the comparison via similarity coefficients emphasize certain molecular features or properties while they neglect others. One of the most important properties of the Tanimoto coefficient is its size-bias as illustrated in the biphenyl query shown below. The Tanimoto coefficient reaches a limit and does not decrease when the size of the second molecule is increased.

molecule A

comparison

molecule B

| n | Tanimoto |
|---|---|
| 0 | 1 (identity) |
| 1 | 0.932 |
| 2 | 0.674 |
| 3 | 0.633 |
| 4 | 0.608 |
| 5 | 0.608 |
| .... | 0.608 |
| 2000 | 0.608 |

# J3.10.7 The Concept of "Fuzzy" Descriptors

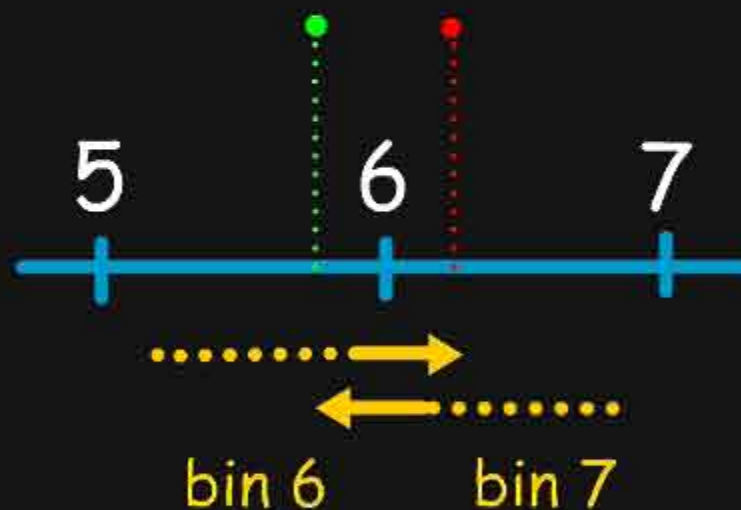If distances between features are binned, what means that they are assigned to a discrete distance range (for example according to the rule "a distance between 6 and 7 Angstroms corresponds to bin 7"), very similar distances might still be assigned to different bins. This happens around the bin borders; in the example above a distance of 5.9 Angstroms would be assigned to bin 6, but a (very similar) distance of 6.1 Angstroms would be assigned to bin 7. "Fuzzy" descriptors also increase the counts of neighboring bins, thereby alleviating the problem of discretizing real-valued distances. An example is given in the applications section.

## J3.10.1  Neighborhood Behavior

If descriptors are employed for similarity searching, the goal is to determine property similarity, with descriptor similarity being the practical means towards this goal. Thus, correlation between the distance in property space (what one is really interested in) and in descriptor space (what one is able to determine computationally) is highly relevant. The concept of "Neighborhood Behavior" describes high correlations between descriptor similarity and property similarity which is crucial for practically every relevant molecular descriptor.

## J3.10.2 Back-Projectability

One aim of property models is to make sensible predictions about novel compounds. Of additional interest is the following question: which structural features actually make this particular compound so active, soluble or toxic? Back-projectable descriptors try to answer this question by projecting features found to be important back onto the molecular structure: they are back-projectable.



identification of structural
features responsible for binding

Training Set:
Known active and
inactive molecules

Back-Projectability

Model of the
property

Predict the property
of novel compounds
using the model as
a black box

## J3.10.10 2D vs. 3D Descriptors

Two-dimensional and three-dimensional descriptors also have specific advantages and disadvantages. While each descriptor also shows individual differences, the following general rules can be given: 2D descriptors perform well where invariant topological features constitute a given activity, they perform less well in the case of topologically more diverse molecules which all show the same property. 3D descriptors on the other hand are able to identify structures with different topologies that still show the same activity, a capability often referred to as "scaffold hopping". On the other hand, they are slower and often found to be less efficient for example in virtual screening runs.

Analogs                                    Unrelated

2D descriptors perform well (computationally efficient)

3D descriptors perform well (computationally slow)

# J3.11.1 Unique Content of Each Similarity Coefficient?

While many similarity coefficients can be defined, and indeed have been, the question arises as to whether they contain different information from each other or whether they behave similarly; in other words, whether they cluster into coefficients with very similar properties. While the first step is to investigate whether or not similarity coefficients show similar behavior, we might also ask a second question: if similarity coefficients indeed show different properties, might it be advantageous to combine them? Both questions will be explored in the following pages.
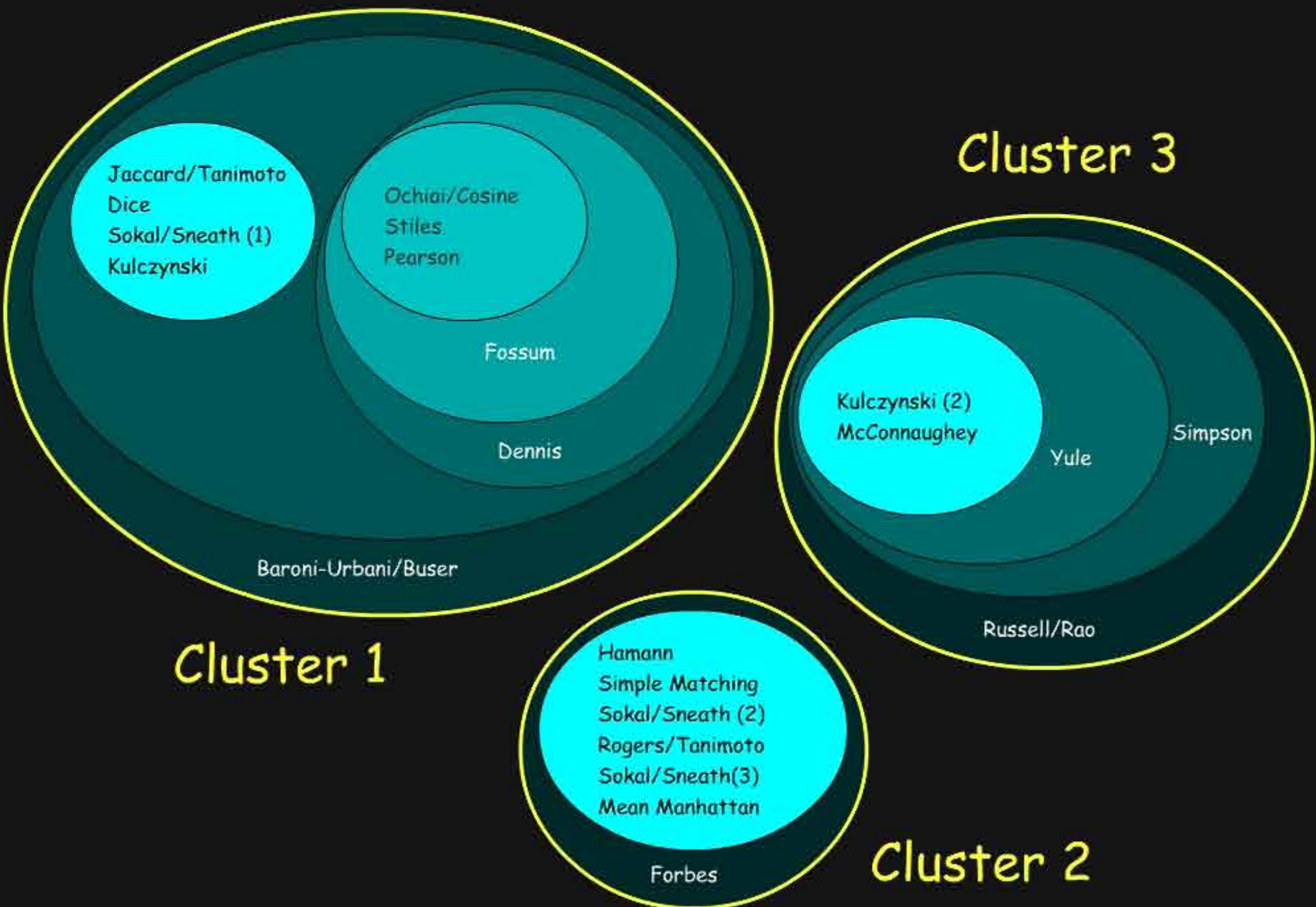
Jaccard/Tanimoto

Dice

Yule

Kulczynski (2)

Simpson

Russell/Rao

Rogers/Tanimoto

Ochiai/Cosine

Stiles

Sokal/Sneath (1)

Are they correlated ?

Hamann

Sokal/Sneath(3)

Forbes

Simple Matching

Pearson

Fossum

McConnaughey

Dennis

Kulczynski

Sokal/Sneath (2)

Baroni-Urbani/Buser

Mean Manhattan

# J3.11.2  Clustering Similarity Coefficients

In a recent work 22 similarity coefficients were analyzed with respect to their behavior, and 3 distinct clusters could be identified. This may suggest that many of them were redundant. A graphic representation of the correlation between descriptors is visualized below (the smaller the distance, the higher the correlation between the descriptors).
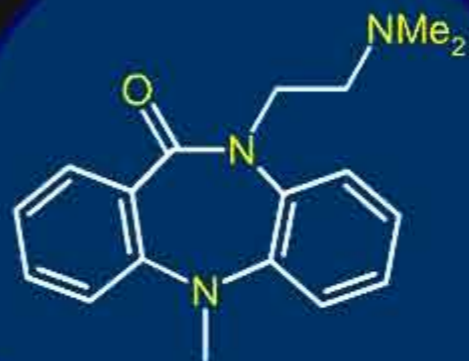
○ 22 similarity coefficients          ○ 3 clusters



Cluster 1

Jaccard/Tanimoto
Dice
Sokal/Sneath (1)
Kulczynski

Ochiai/Cosine
Stiles
Pearson

Fossum

Dennis

Baroni-Urbani/Buser

Cluster 2

Hamann
Simple Matching
Sokal/Sneath (2)
Rogers/Tanimoto
Sokal/Sneath(3)
Mean Manhattan

Forbes

Cluster 3

Kulczynski (2)
McConnaughey

Yule

Simpson

Russell/Rao

Imagine you are on a TV quiz show and have to answer a set of questions ranging from any area of sports to any area of culture. You have two people who are allowed to help you; one is better at sports (and okay at culture), one is better at culture (and okay at sports). But you don't know who's better in which area since the persons are anonymous. Whom do you trust? One option is to ask both people and take both opinions into account. This is precisely the concept of "consensus scoring". This concept can be applied in Similarity Analysis.

consensus score = 6/10

4/10

How would you score the similarity between these two molecules?

8/10

## J3.11.4  Why Does Consensus Scoring Improve the Results?

Mathematically speaking, consensus scoring improves the results in cases where individual errors are independent of each other: they cancel out. This is not true in every case, but on average better results are obtained. On the other hand, if the predictions errors are not (at least partly) independent of each other, consensus scoring is unable to improve results.
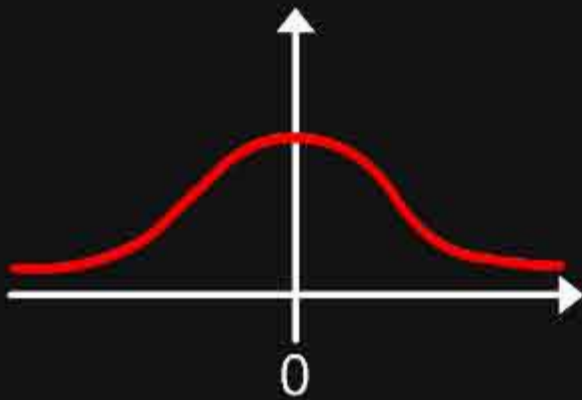
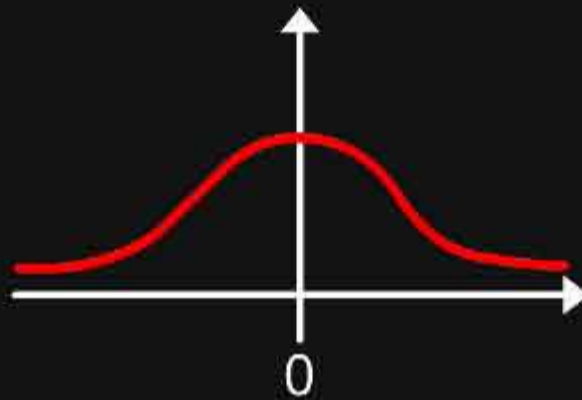◉ Dependent errors                    ◉ Independent errors

large standard deviation ⟷ mutually dependent errors

distribution of errors

method 1                    method 2                    consensus
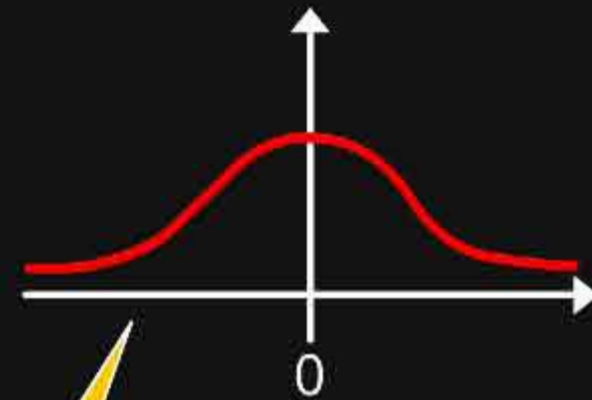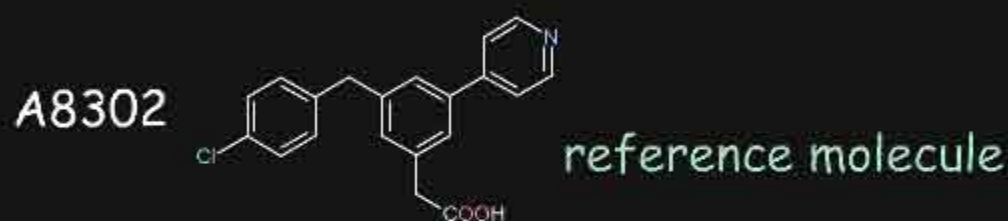
0                          0                          0

no improvement over single methods

# J3.11.5 What Algorithms Exist for Consensus Scoring?

One can merge the predictions from multiple classifiers in many ways: One can use the sum of the individual predictions, the relative ranks or the highest prediction. In practice, the sum of individual predictions has been found to perform best in some cases. In others the highest score assigned to each individual prediction performed best. It should be noted that no generally applicable rules exist here unfortunately.

A8302

reference molecule

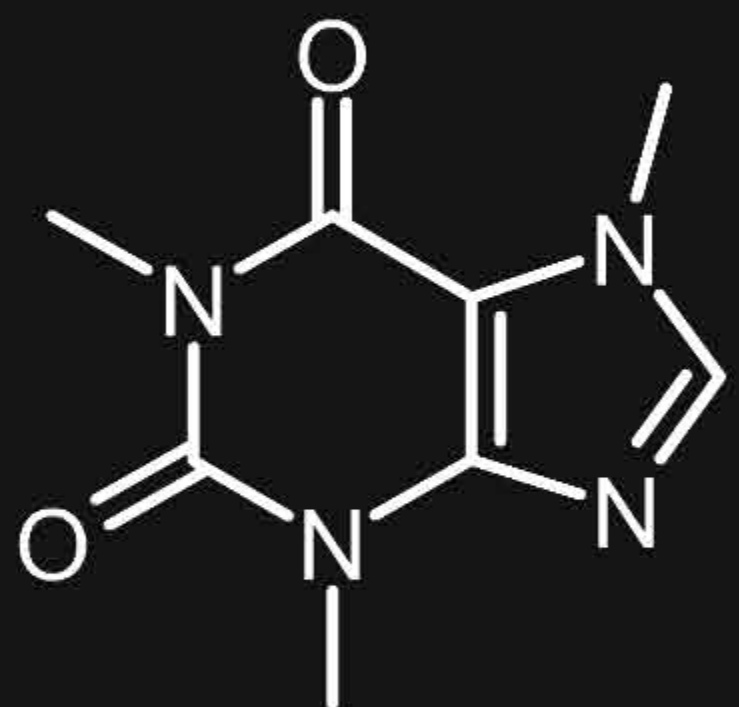| Compound | Tanimoto | Cosine | Dice | Consensus |
|----------|----------|--------|------|-----------|
| A1001 | 0.812 | 0.763 | 0.791 | 0.788 |
| A246 | 0.731 | 0.620 | 0.706 | 0.685 |
| A7052 | 0.520 | 0.480 | 0.532 | 0.510 |
| A5031 | 0.792 | 0.742 | 0.813 | 0.782 |
| A9458 | 0.412 | 0.420 | 0.450 | 0.427 |
| A5674 | 0.910 | 0.924 | 0.873 | 0.902 |

## J3.12.1 Limitation of Ligand-Based Approaches

The "molecular similarity principle" is generally used when the biological receptor (or the mechanism of action) is not known, and this lack of knowledge represents the most fundamental limitation (but also its strength!) of all ligand-based approaches: when focusing only on the ligands, everything cannot be explained. Ignoring the receptor introduces a bias with the risk of incorrect interpretations. This is discussed in the following pages.
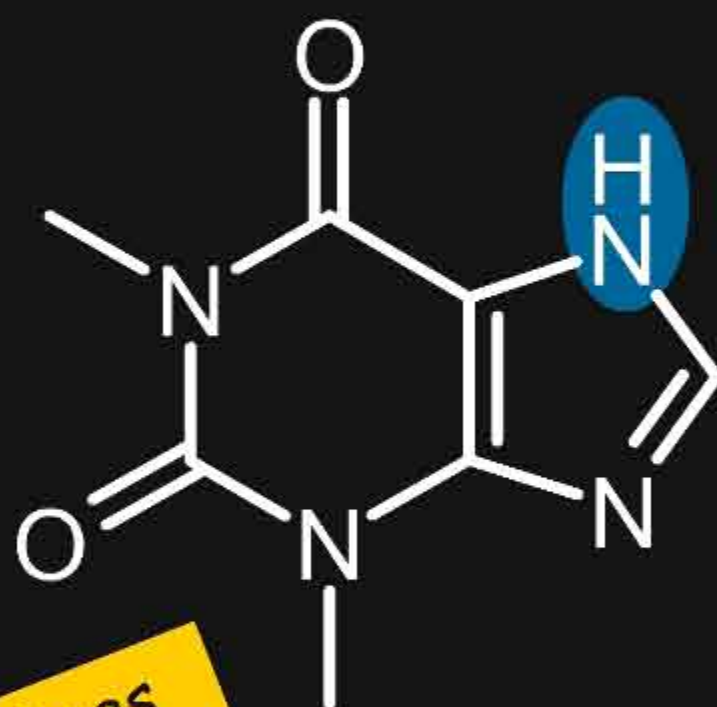
## J3.12.2 Example of Ligand-Based Approach Limitation

Below, two very similar structures are shown Caffeine differs from Theophyllin purely by a methyl group, a very small change overall. Still, RNA binding is decreased by a factor of 10,000, which would not be expected based on the molecular similarity principle. If we ignore the interactions with the receptor, we cannot explain this difference.



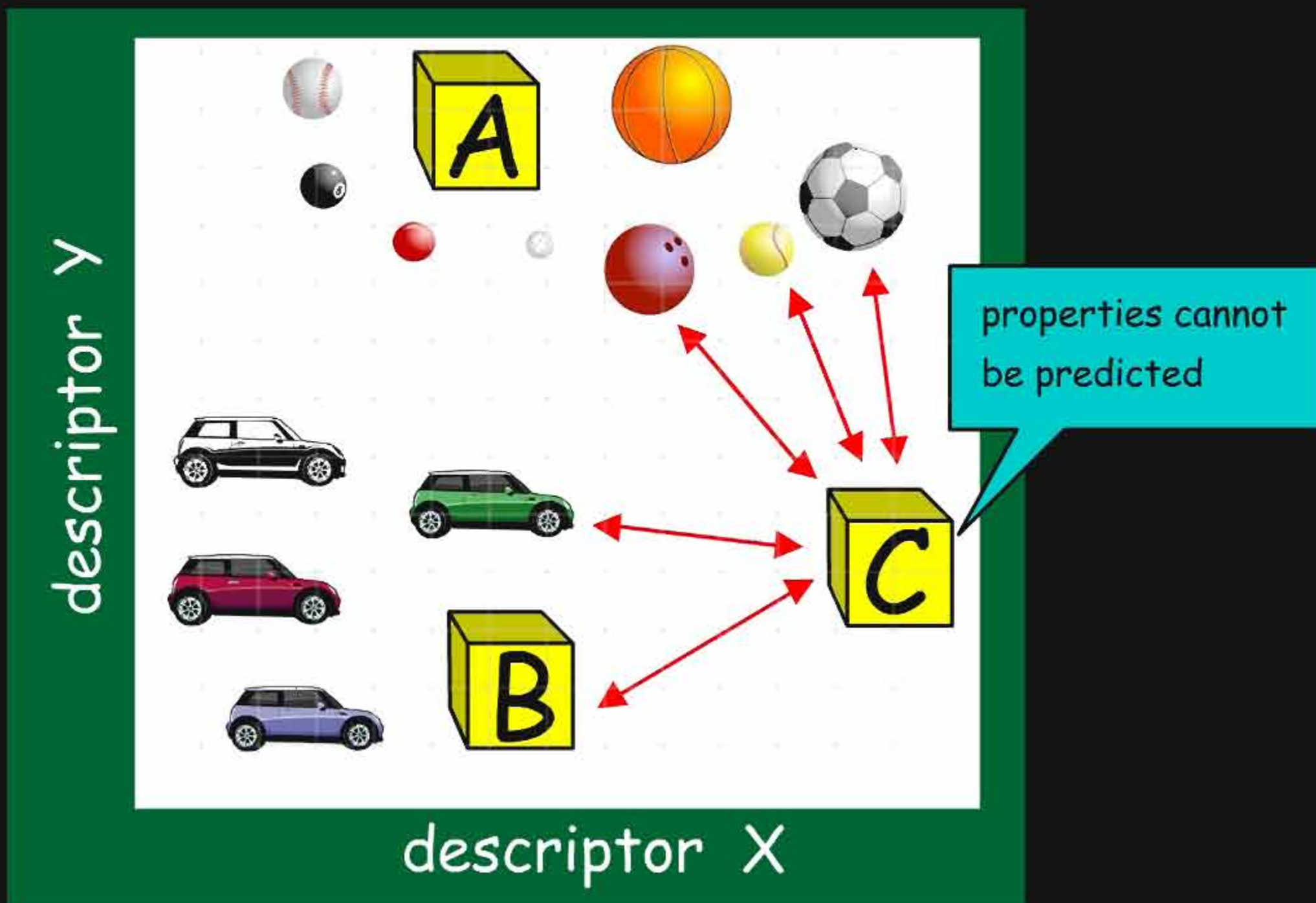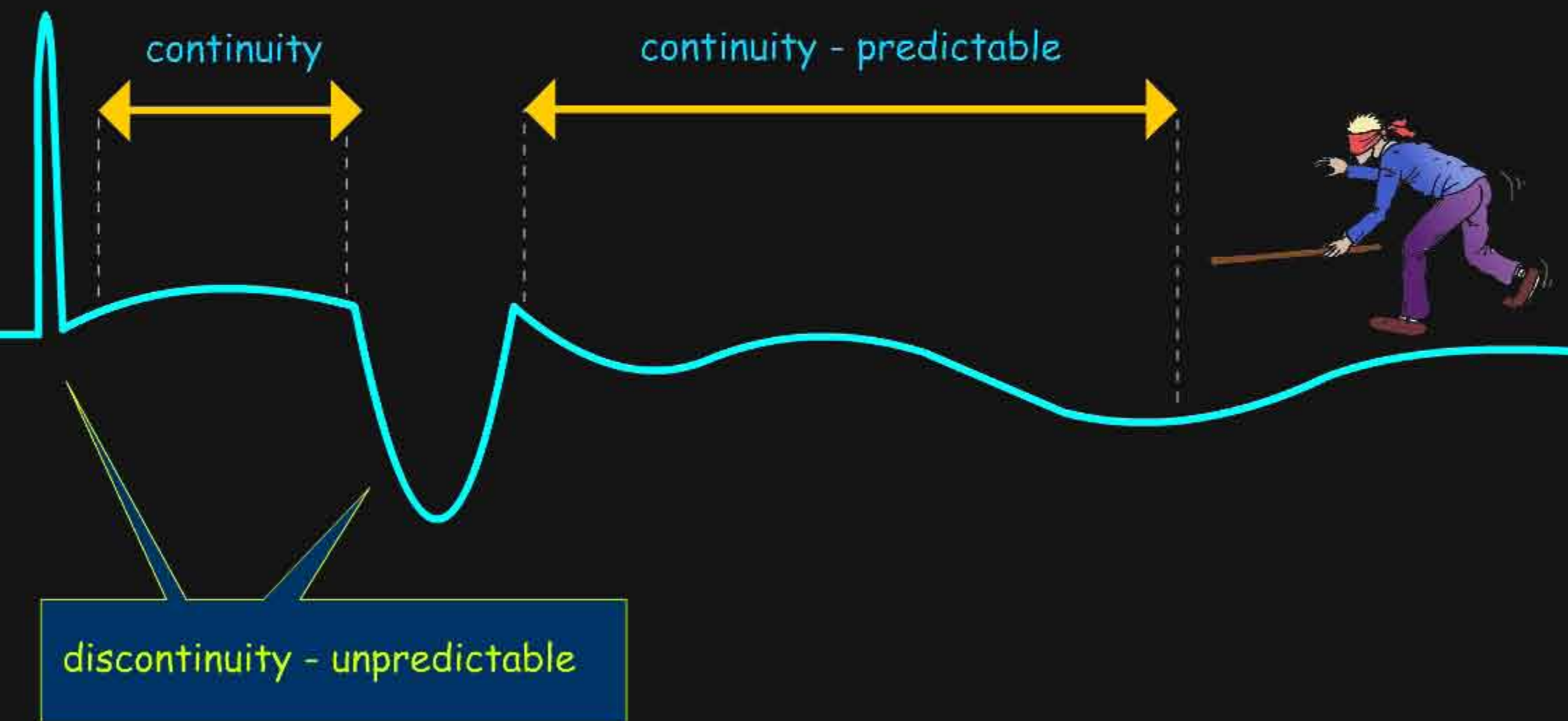10000 times more potent

Caffeine

Theophilline

## J3.12.3  Limitation Due to Extrapolations

The first limitation of the principle of similarity concerns the distribution of the molecules of the initial dataset. While it may straightforward to extrapolate the properties of closely related objects (e.g. A or B), the extrapolation for C becomes highly debatable. The only way to overcome this limitation is to cover a large range of values for those descriptors believed to be relevant to the property at hand.

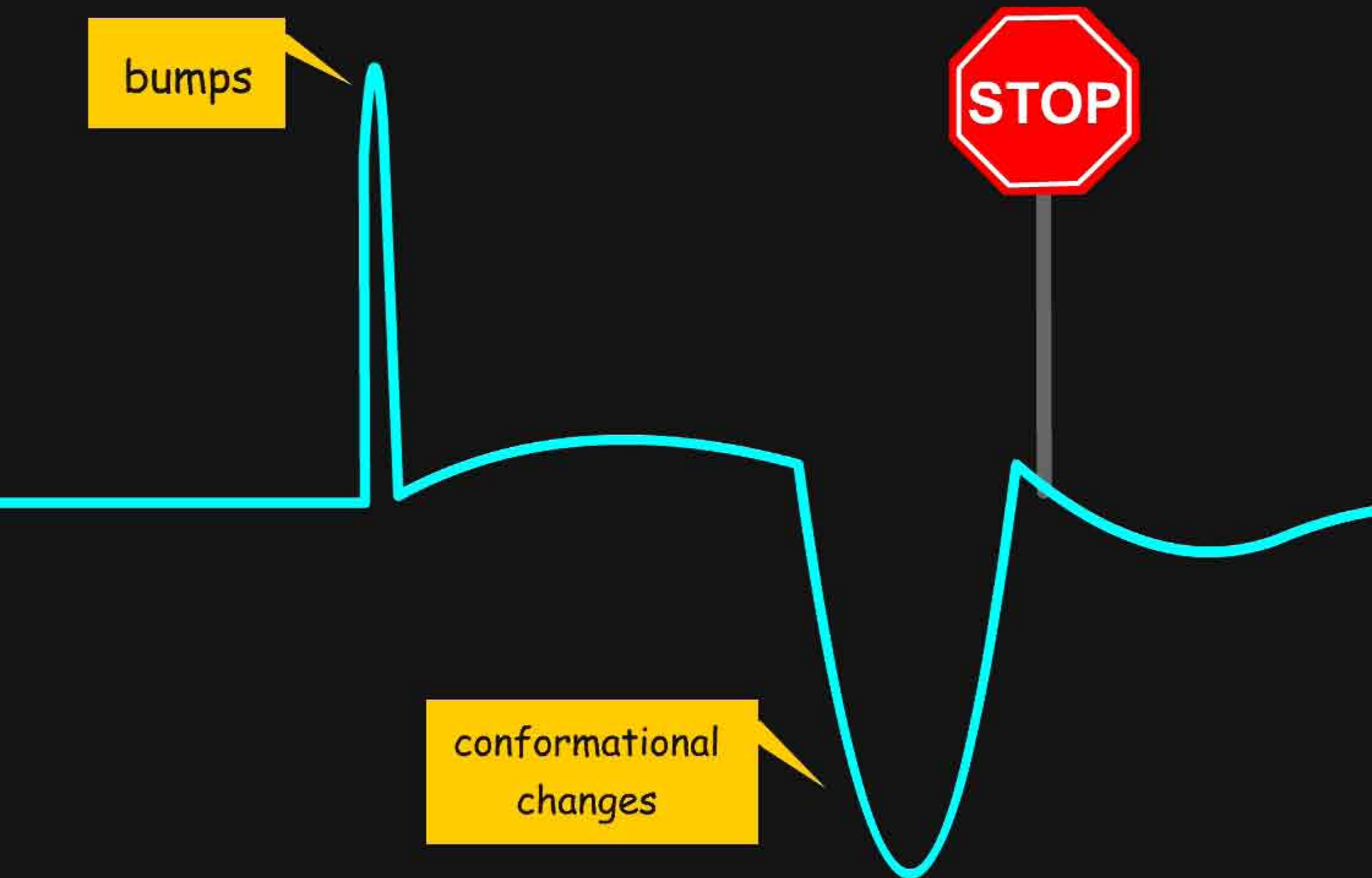## J3.12.5  Principle of Continuity

Because of our ignorance of the interactions with the receptor, the molecular similarity principle (based on ligand-based methods) relies on the principle of "continuity", assuming that changes are gradual in the recognition of molecules. Therefore, the similarity principle is not applicable when abrupt changes occur in the system.

# J3.12.6  Discontinuity in Molecular Recognition

Examples of cases where abrupt changes occur and the similarity principle is not applicable include the following: bumps, ligand or receptor conformational changes, flip in binding modes, discontinuity in ligand properties or in receptor function.

A small change in the structure of a ligand can result in dramatic changes of the biological properties. For example in the structure below, the simple move of one nitrogen atom in the pyrimidine ring transforms a potent inhibitor into an inactive molecule. This may be due to bumps with the receptor, a discontinuity in the activity that in the first place, is not foreseeable (see button "Explanation"). Here, the similarity principle is not applicable.

◉ SAR                    ◉ Explanation



⚠ similarity principle not applicable
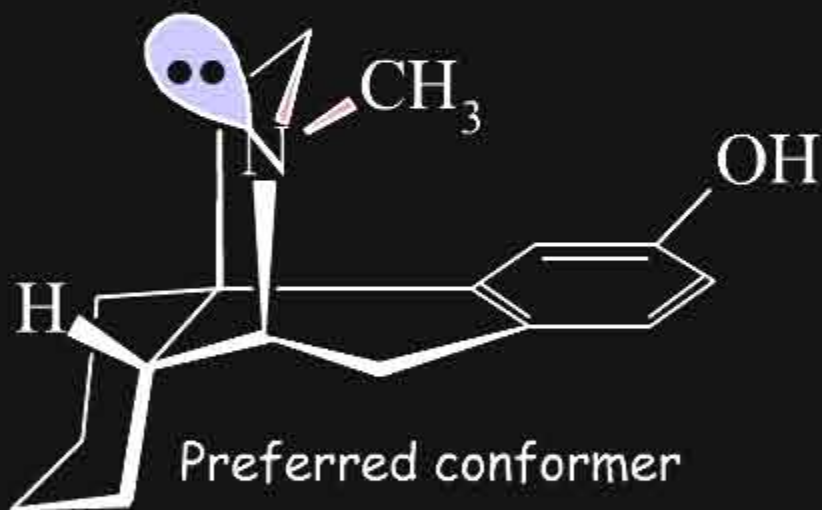
Ki = 30 nM                    Ki > 10,000 nM

# J3.12.8 Ligand Conformational Change

A change in the structure of a ligand can result in an important change in the conformation of the molecule, and the principle of continuity cannot be invoked in such cases. Examples of conformational changes are shown below, where only one carbon atom is added to the initial structure. The similarity principle is not applicable.
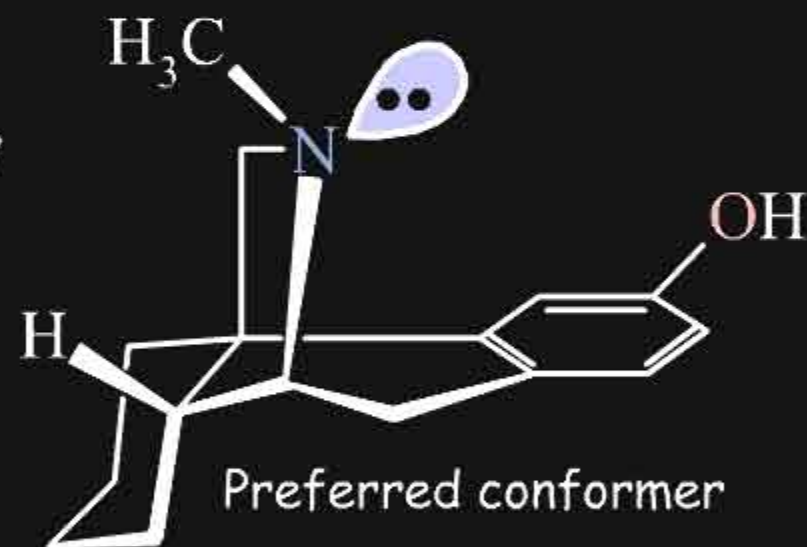
◉ Morphinan          ◉ Nifedipine          ◉ Amide



Morphinan

Preferred conformer

**Morphinan**

similarity principle not applicable
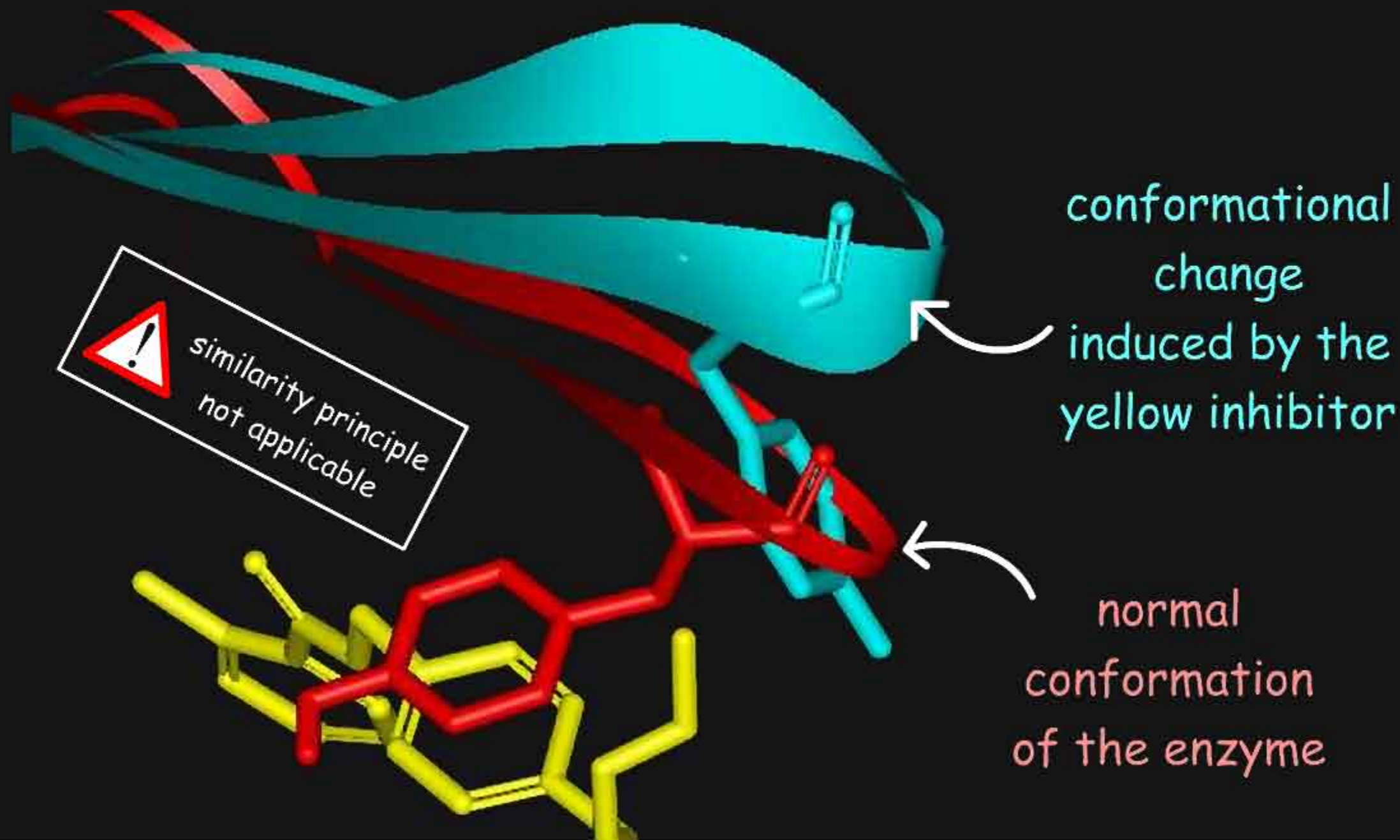
Preferred conformer

**D-nor morphinan**

## J3.12.9  Receptor Conformational Changes

For reasons that are not predictable, the receptor may undergo substantial conformational changes. In the example below the enzyme renin adopts different conformations (blue and red) depending on the inhibitors. Note the change of a key tyrosine residue, which exhibits favorable interactions with the inhibitor (yellow); while with the red conformation it clashes. The similarity principle is not applicable when this is overlooked.
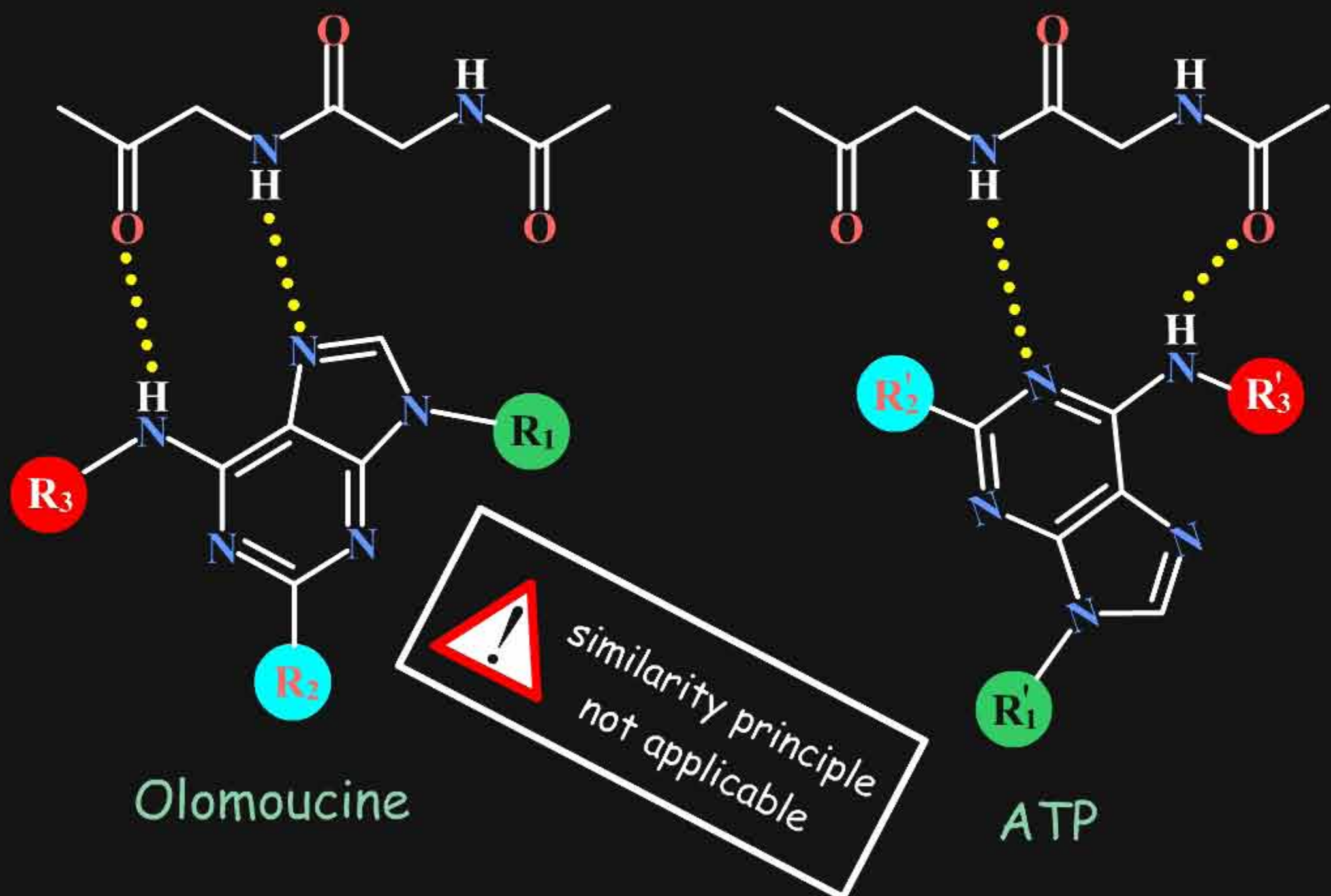
◉ Sketch                                                        ◉ 3D



conformational change induced by the yellow inhibitor

⚠ similarity principle not applicable

normal conformation of the enzyme

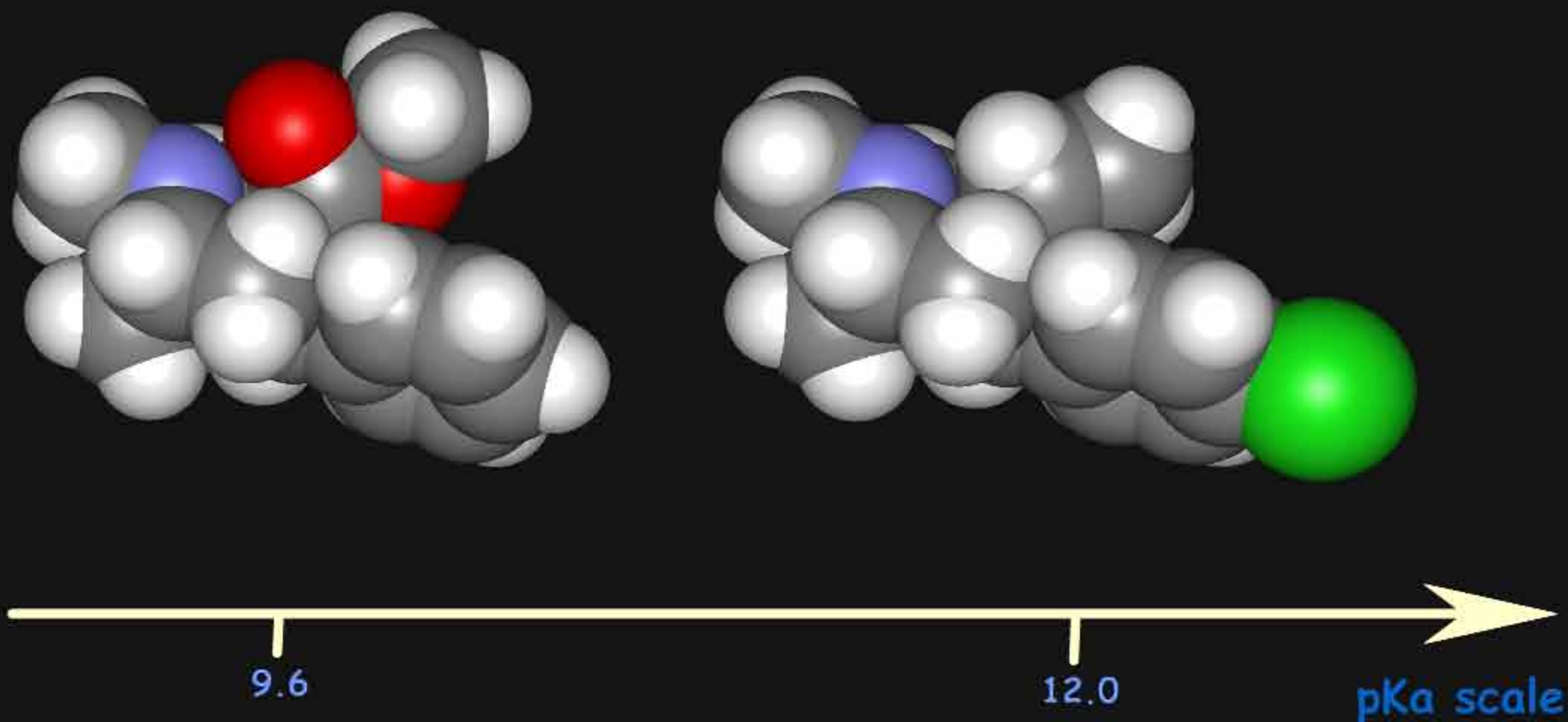## J3.12.10 Flip in Binding Mode

When an X-ray structure of a complex is available for a given ligand, the design normally includes the assumption that the binding mode of the analog will be similar to that of the reference molecule. This has proven to be valid in many cases. However, there are examples where important flips were observed (in some cases a 180 degree rotation of the structure), and this is another limitation of the above mentioned continuity principle. The similarity principle is not applicable.



Olomoucine

ATP

similarity principle not applicable

# J3.12.11 Discontinuity in Ligand Property

A small alteration of a structure may introduce substantial changes in specific physico-chemical properties that are not apparent in the first place. For example, the replacement of a carbon atom by an oxygen, or the change in the number of carbon atoms in a molecule can affect many properties of the substance. Examples of changes in $pK_a$ and logP are presented on the following pages.
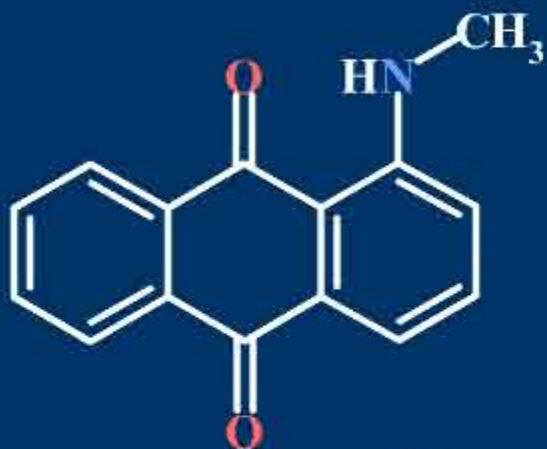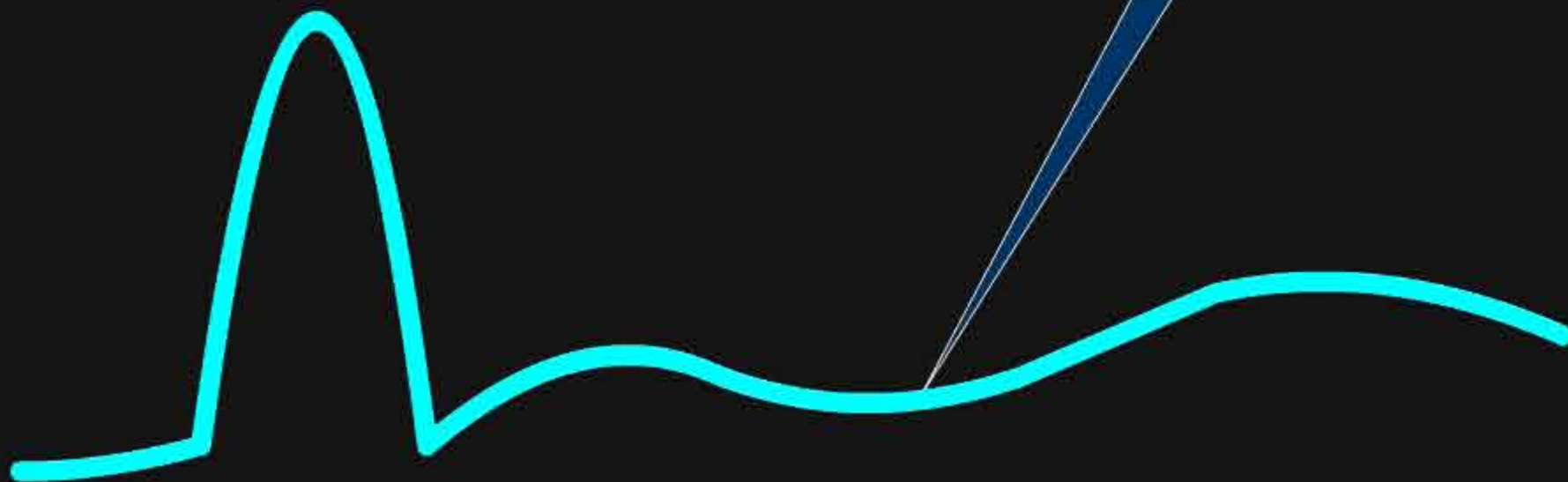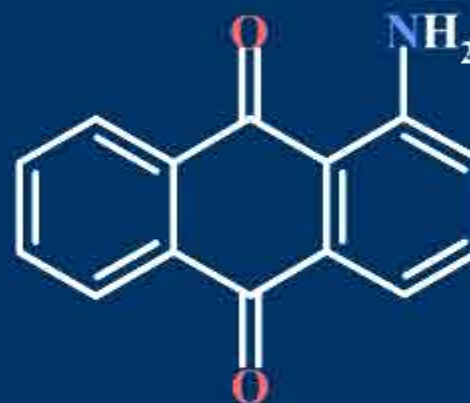


9.6                    12.0

pKa scale

# J3.12.13 LogP

When changes are introduced in a structure, their impact on the global characteristics of the modified compounds must be taken into account. The example below illustrates a change in logP that is not negligible.

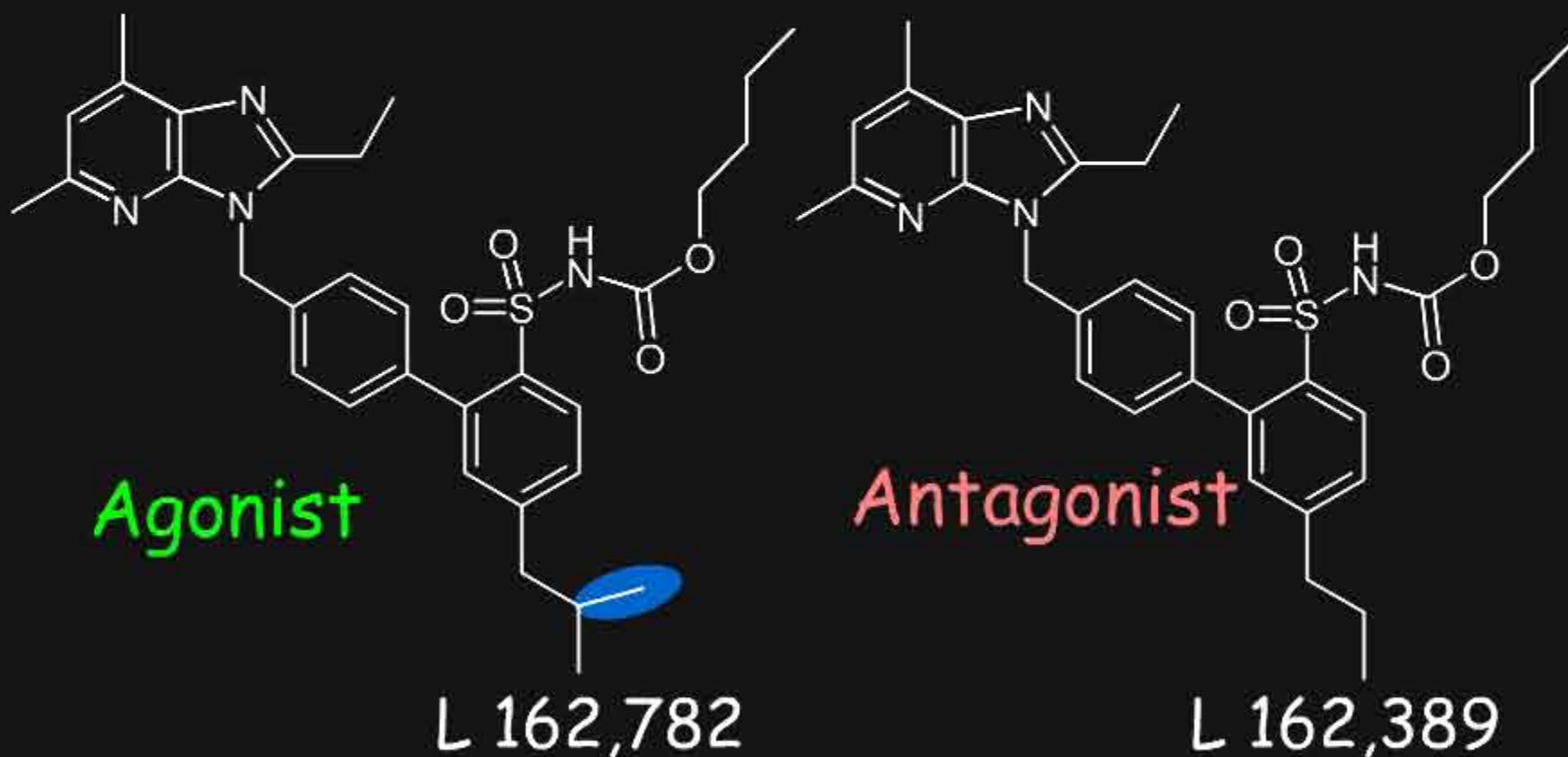# J3.12.14  Discontinuity in the Function of the Receptor

A small change of a ligand can result in dramatic alteration of effect with a receptor. For example L 162,389 and L 162,782 differ only by one methyl group. The former is an antagonist and the latter an agonist of the angiotensin AT1 receptor. Appreciate the dramatic change in receptor binding that occurs for all the examples shown below. A conformational change of the receptor is initiated by the agonist but not the antagonist: a discontinuity that is very hard to predict.
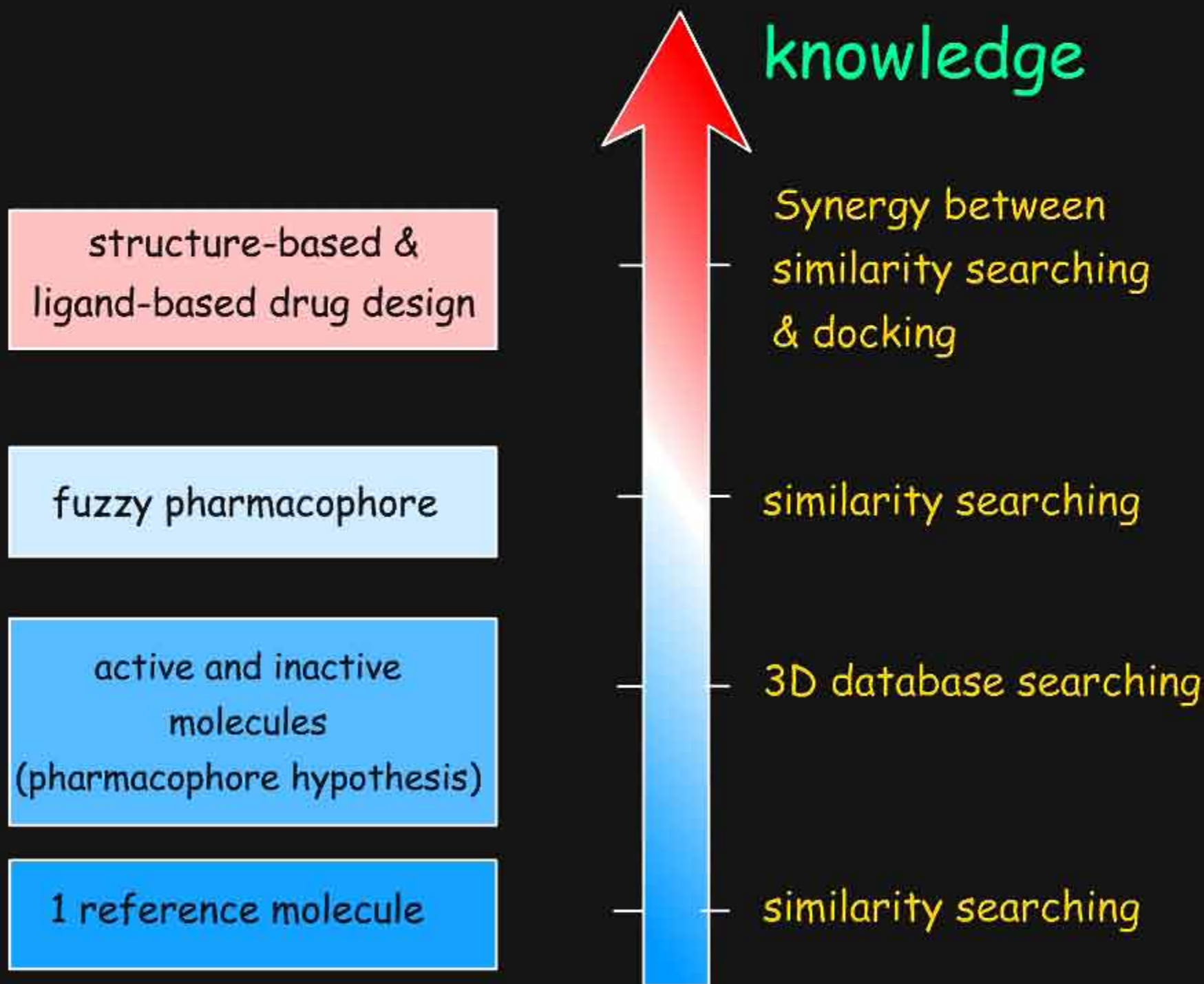
● Angiotensin AT1          ● GABA          ● Benzodiazepine



Binding to the angiotensin II AT1 receptor

Agonist

L 162,782

Antagonist

L 162,389

## J3.13.1 How Does "Molecular Similarity" Fare Today?

Despite of the limitations of ligand-based design, "Molecular Similarity" is a concept routinely applied in pharmaceutical and related industries today. Molecular similarity methods are very much of importance in the early stages of drug discovery. They are most often applied when little knowledge is available (e.g. only one active compound). With increasing knowledge, their importance generally decreases and vanishes when the structure-based level is reached.
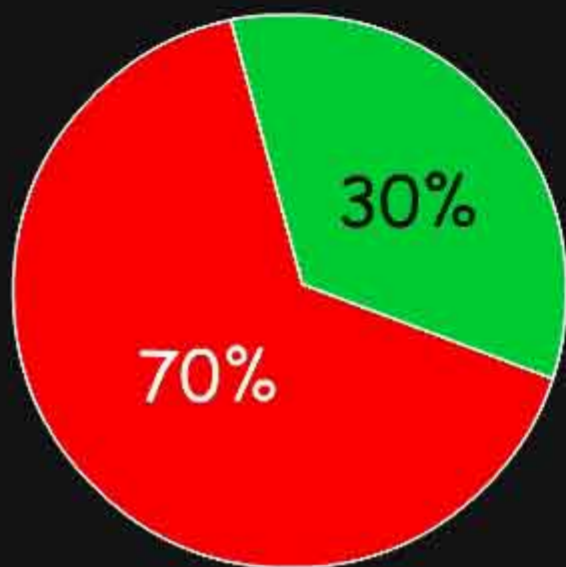
In an attempt to answer the question "do structurally similar molecules have similar biological activity?" a follow-up study to 115 high-throughput screening assays shows that at high Tanimoto similarity values (>0.85), only 30% of compounds similar to an active molecule are themselves active. This reveals the complexity of molecular similarity. A list of possible factors explaining such observations is given below.

for only 30% of molecules the molecular similarity principle is true

Possible reasons

30%

70%

- relevant descriptors with noise
- discontinuity in the properties
- invalid extrapolation
- bad choice of similarity index
- experimental errors
- model too simplistic