



ROBERT WOOD JOHNSON  
MEDICAL SCHOOL

University of Medicine & Dentistry of New Jersey

# **Next Generation Computational Chemistry Tools to Predict Toxicity of CWAs**

**William (Bill) Welsh**

**welshwj@umdnj.edu**

**Prospective Funding by DTRA/JSTO-CBD**



**A State-wide, Regional and National Resource**

**< [www.ebCTC.org](http://www.ebCTC.org) >**

Funded with support from the U.S. EPA

# Consortium Members



ROBERT WOOD JOHNSON  
MEDICAL SCHOOL

University of Medicine & Dentistry of New Jersey



U.S. Food and Drug  
Administration

Center for Toxicoinformatics, NCTR

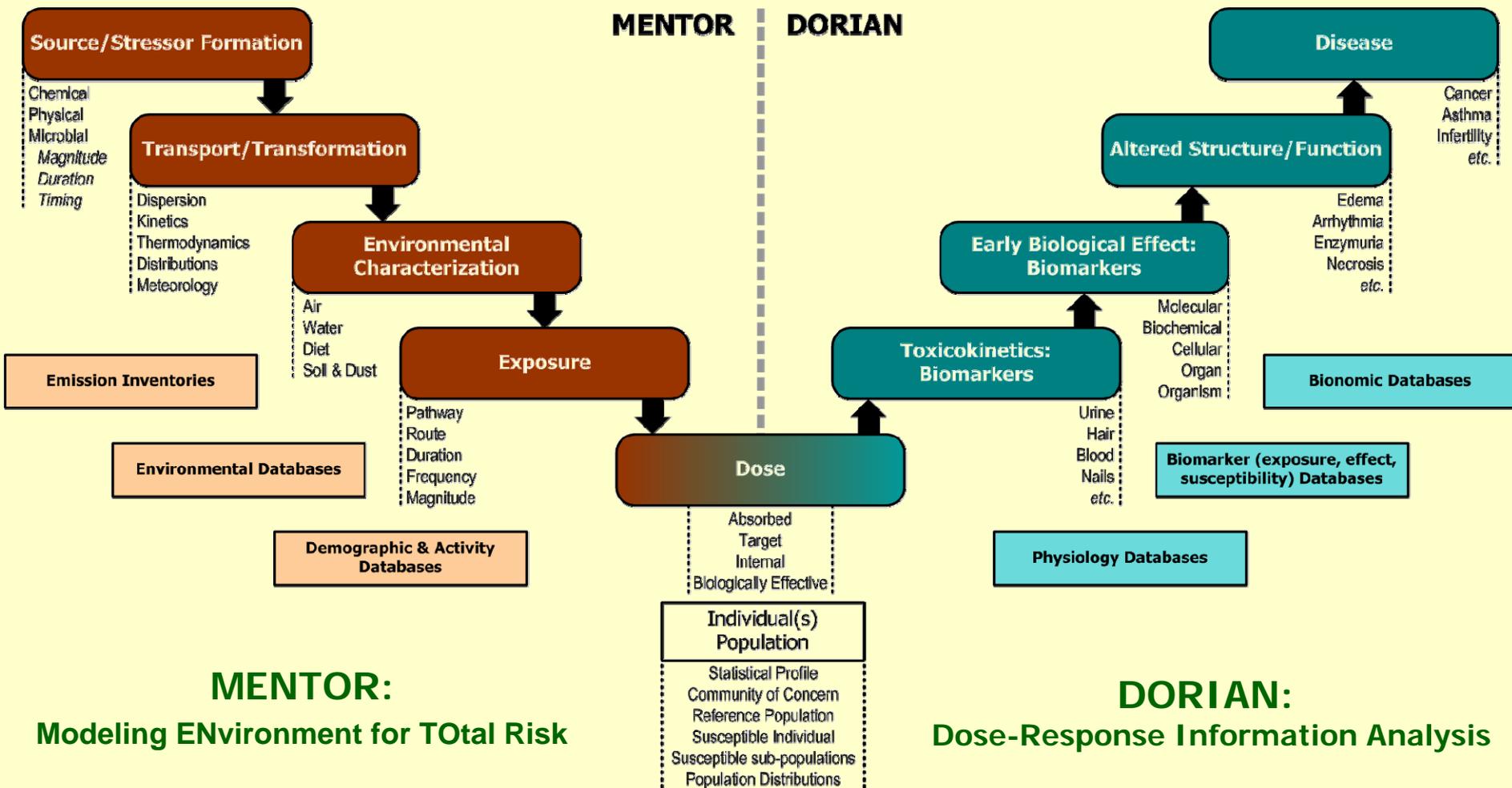


## Major Research Thrusts

- MENTOR-DORIAN Computational Toxicology System that spans the Source->Dose->Outcome continuum
- The Environmental Bioinformatics Knowledge Base (ebKB: [www.ebCTC.org](http://www.ebCTC.org))
- *ArrayTrack*: toxicological bioinformatics platform to process genomics, proteomics and metabonomics data
- Hepatocyte Metabolic Model for Xenobiotics
- **ChemTox, a suite of chem-informatics tools for toxicant identification & characterization**

# MENTOR & DORIAN

## Address the Source-to-Outcome Continuum



Adapted from chart by R. Calderon, USEPA/NHEERL, 2003

# ebKB environmental bioinformatics Knowledge Base

Computational Toxicology | Risk Assessment | Diagnostic Tools

## Genomics/Transcriptomics

### Databases

- Sequences/Maps
- Genetic Markers

The environmental bioinformatics Knowledge Base (ebKB) serves as a comprehensive compendium of tools, databases, and literature

### Search Environmental Bioinformatics Resources

## Proteomics

### Databases

- Shape/Structure
- Protein Markers
- NMR Spectra
- Classification
- Macromolecular Movements
- DNA-Protein Interactions

### Tools

- Data Management
- Search/Alignment
- Classification
- Structure Prediction
- Post Processing and Visualization

### Metadata

- Standards and Markup Languages
- Portals and Knowledge Libraries
- **Selected Literature**

## Metabonomics

### Databases

- Pathways

### Tools

- Data Management
- Simulators
- Pathway Profilers
- Post Processing and Visualization

### Metadata

- Standards and Markup Languages
- Portals and Knowledge Libraries
- **Selected Literature**

## Cytomics

### Databases

- Apoptosis
- Senescence
- Signal Transduction
- Oncology

### Tools

- Data Management
- Cell Simulators
- Post Processing and Visualization

### Metadata

- Standards and Markup Languages
- Portals and Knowledge Libraries
- **Selected Literature**

## Physiomics

### Databases

- Model Organisms

### Tools

- Data Management
- PBTK Models
- Post Processing and Visualization

### Metadata

- Standards and Markup Languages
- Portals and Knowledge Libraries
- **Selected Literature**

## CHEMINFORMATICS

### Databases

- QSAR
- Shape Signature

### Tools

- Data Management
- Virtual Synthesis
- Similarity Search
- Post Processing and Visualization

### Metadata

- Standards and Markup Languages
- Portals and Knowledge Libraries
- **Selected Literature**

## ENVIROINFORMATICS

### Databases

- Concentrations and Exposures
- Toxicity
- Demographics and Activities
- Biomonitoring

### Tools

- Data Management
- Environmental Fate and Transport Modeling
- Exposure Modeling
- Dose Modeling
- Integrated Modeling
- Post Processing and Visualization

### Metadata

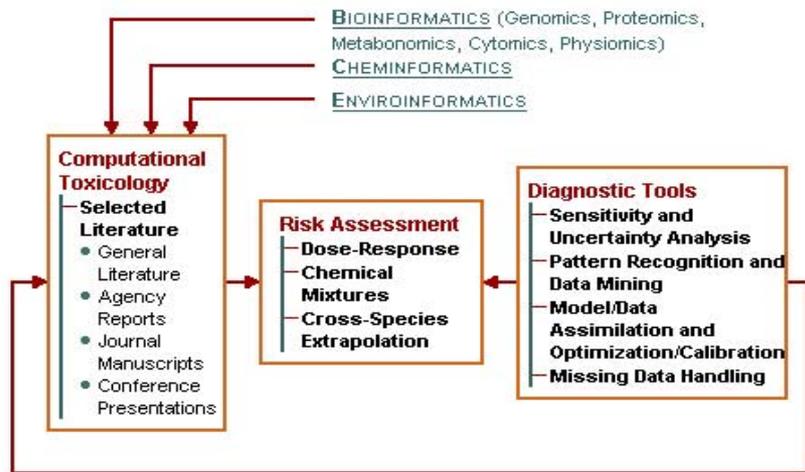
- Standards and Markup Languages
- Portals and Knowledge Libraries
- **Selected Literature**

## BIOINFORMATICS

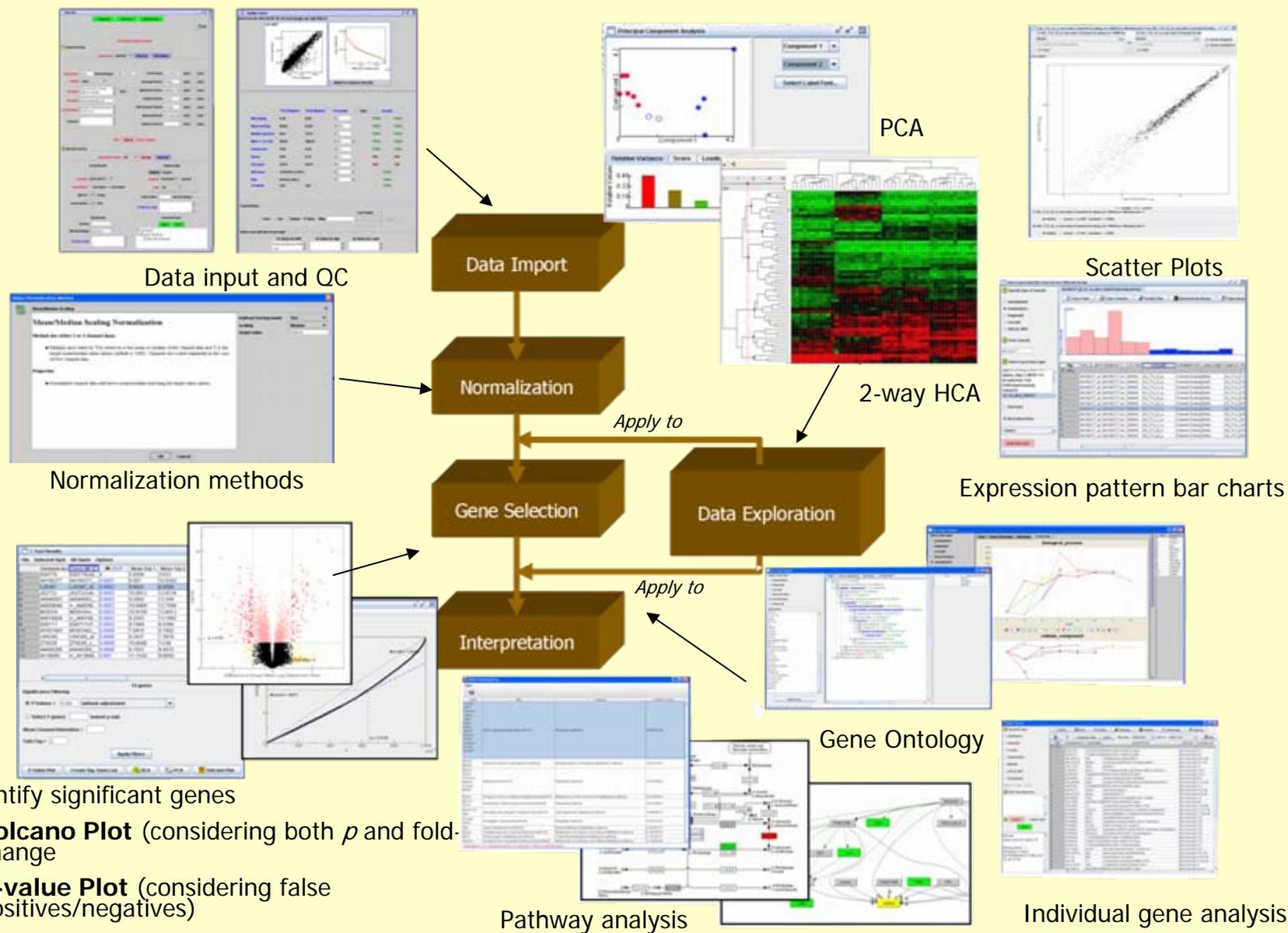
### Integrative Databases

### Integrative Tools

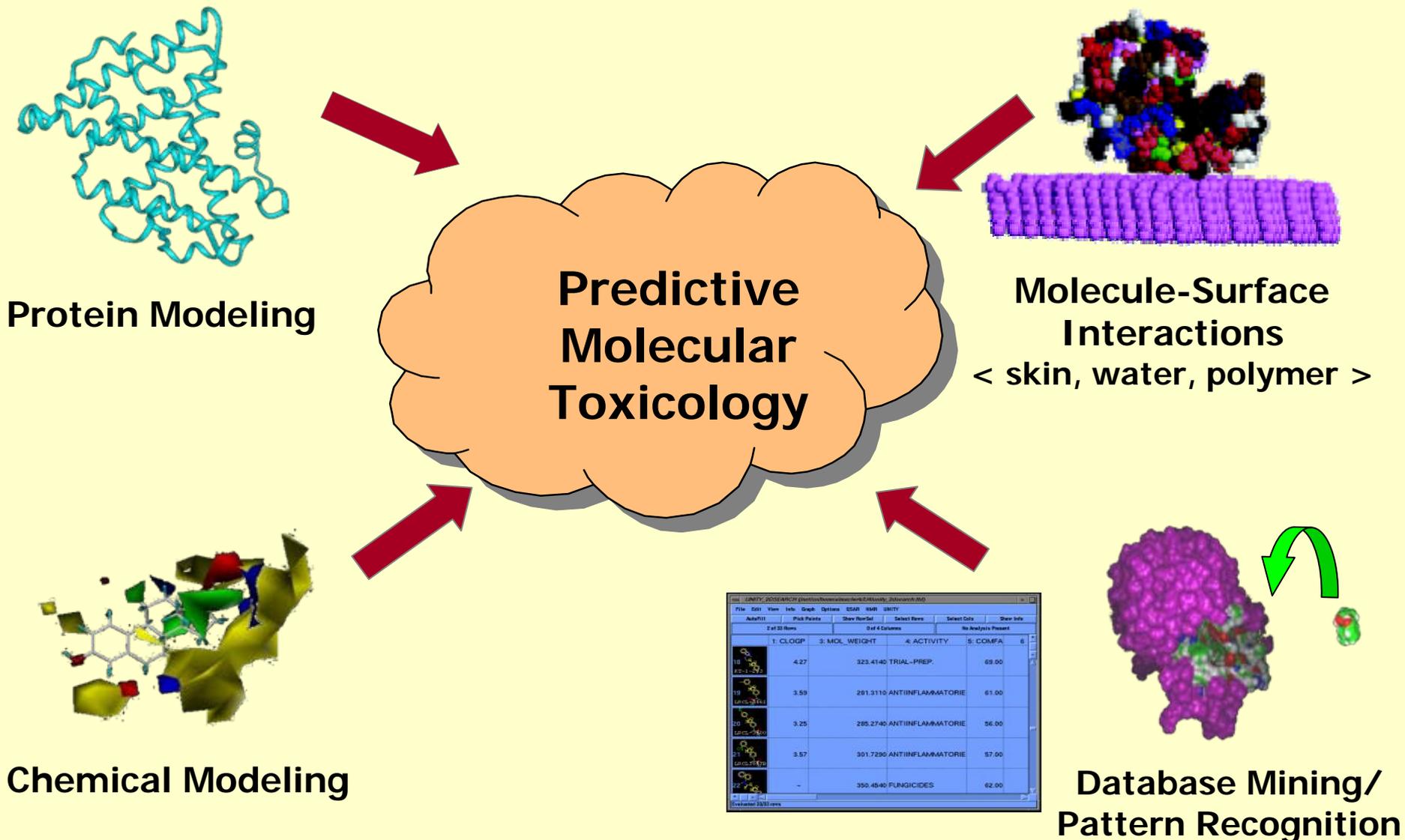
### Integrative Metadata



# ArrayTrack Suite of Bioinformatics Tools

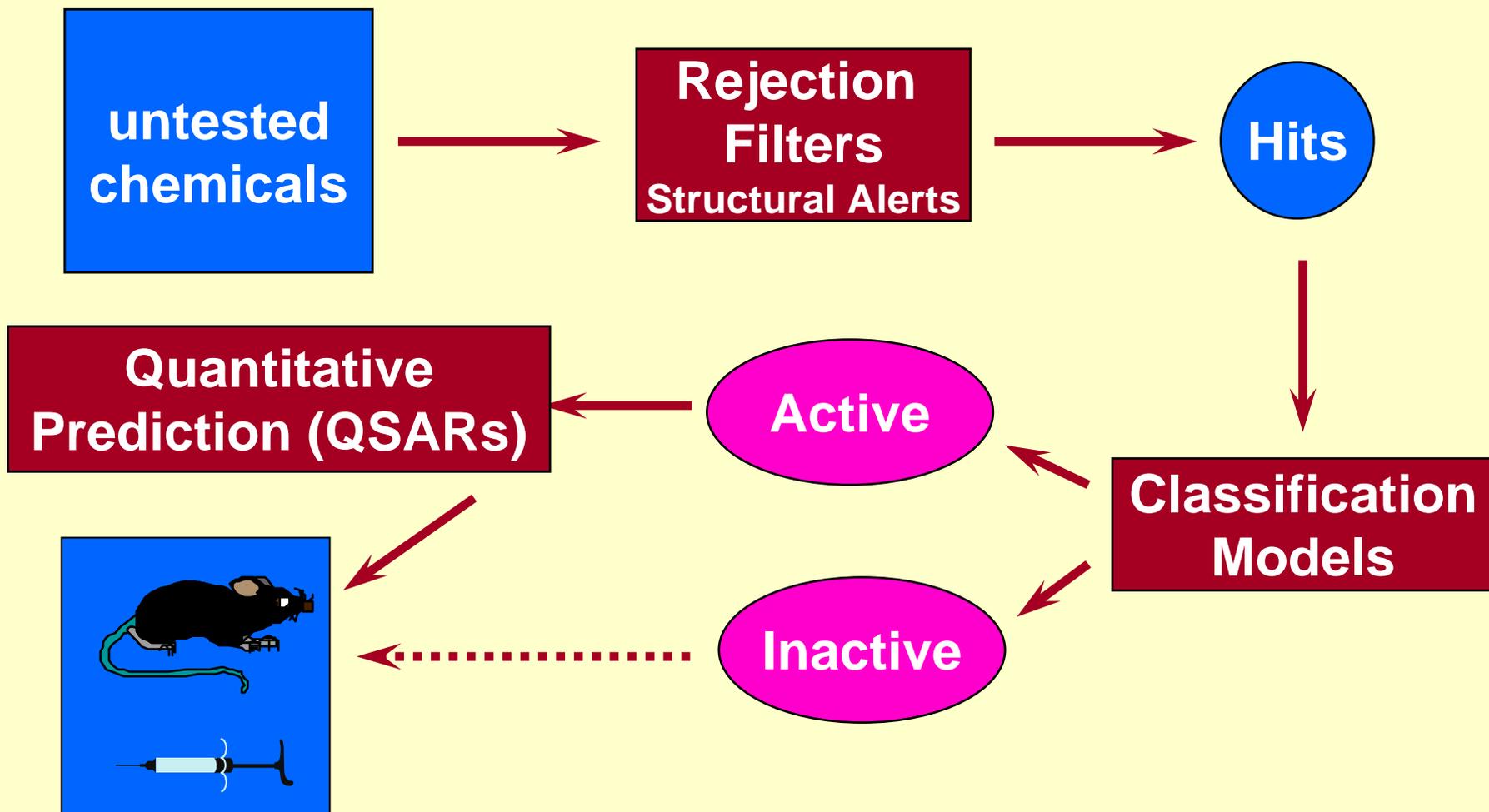


# ChemTox, an Integrated Suite of Cheminformatics Tools



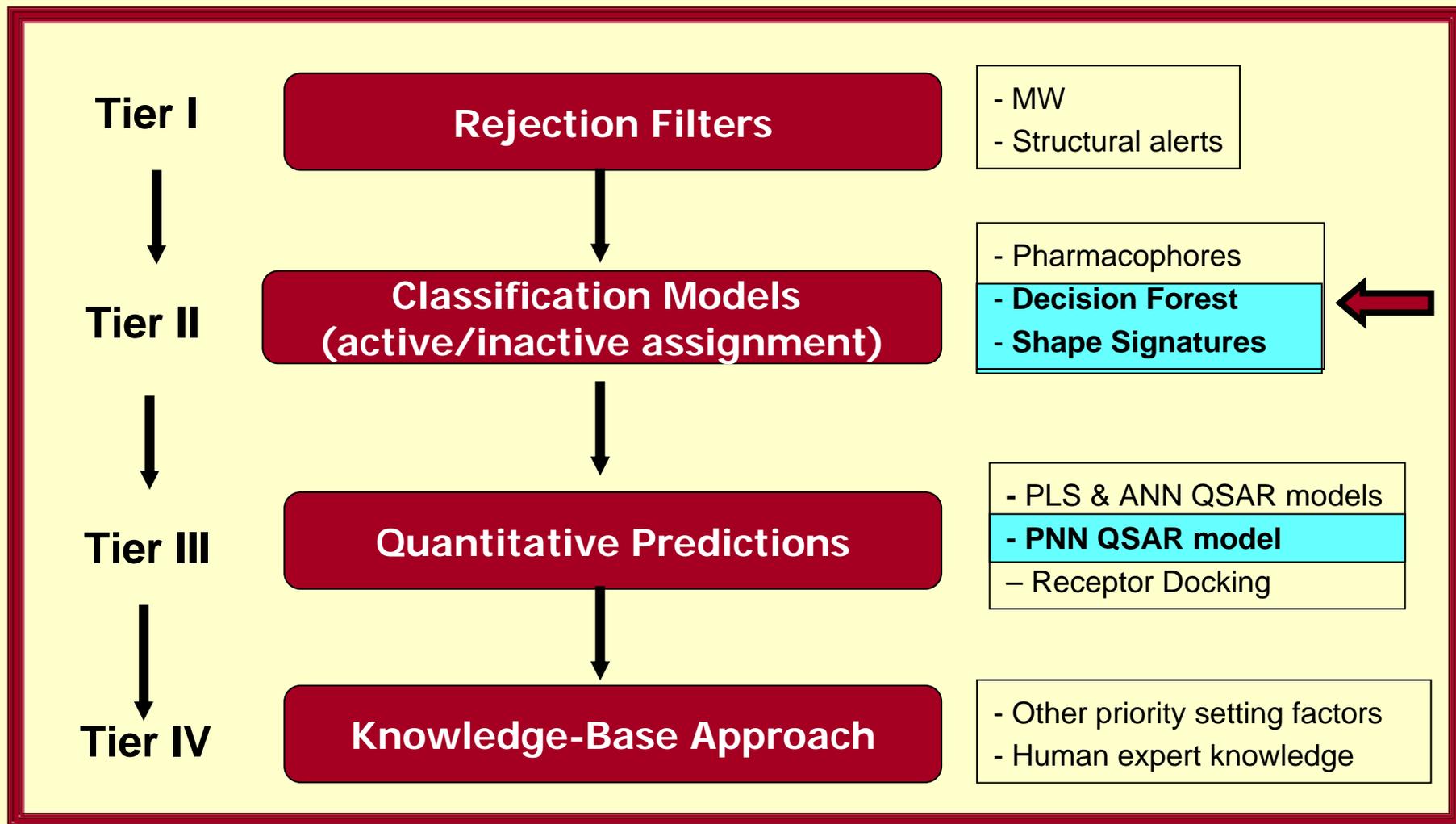
# Computational Screening Paradigm

## - Priority Setting -



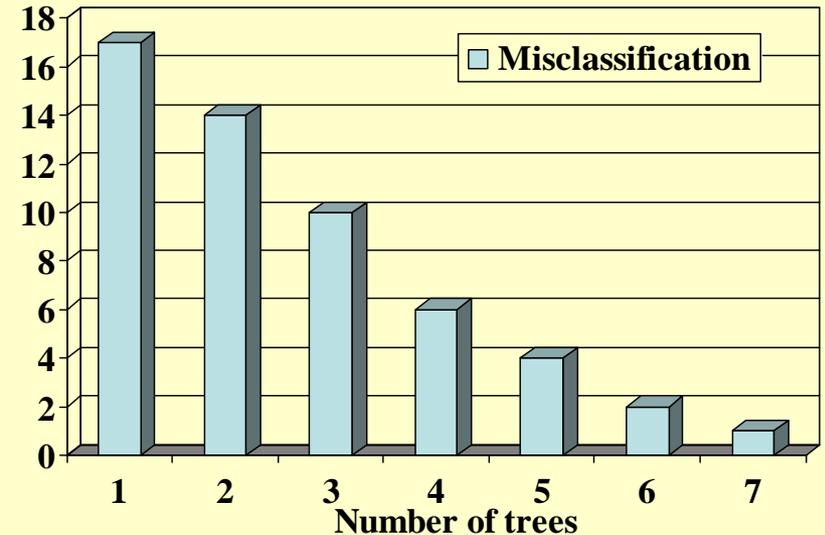
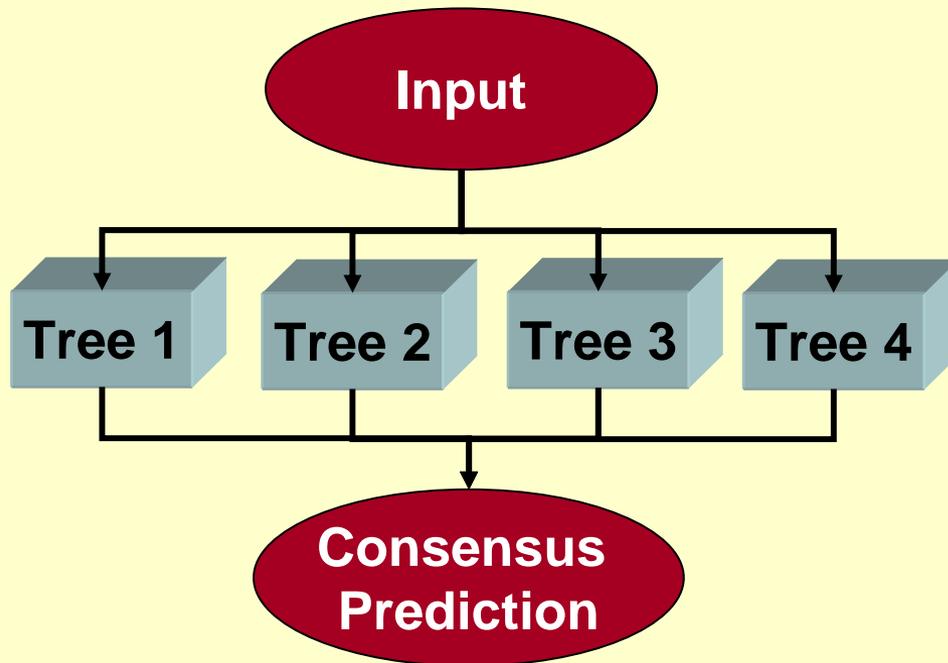
# Hierarchical Screening Framework

- addresses the need to minimize *false negatives* and *uncertainties*
- recognizes that no single computational model is adequate



# Decision Forest

- Improved classification by combining independent Decision Tree models -

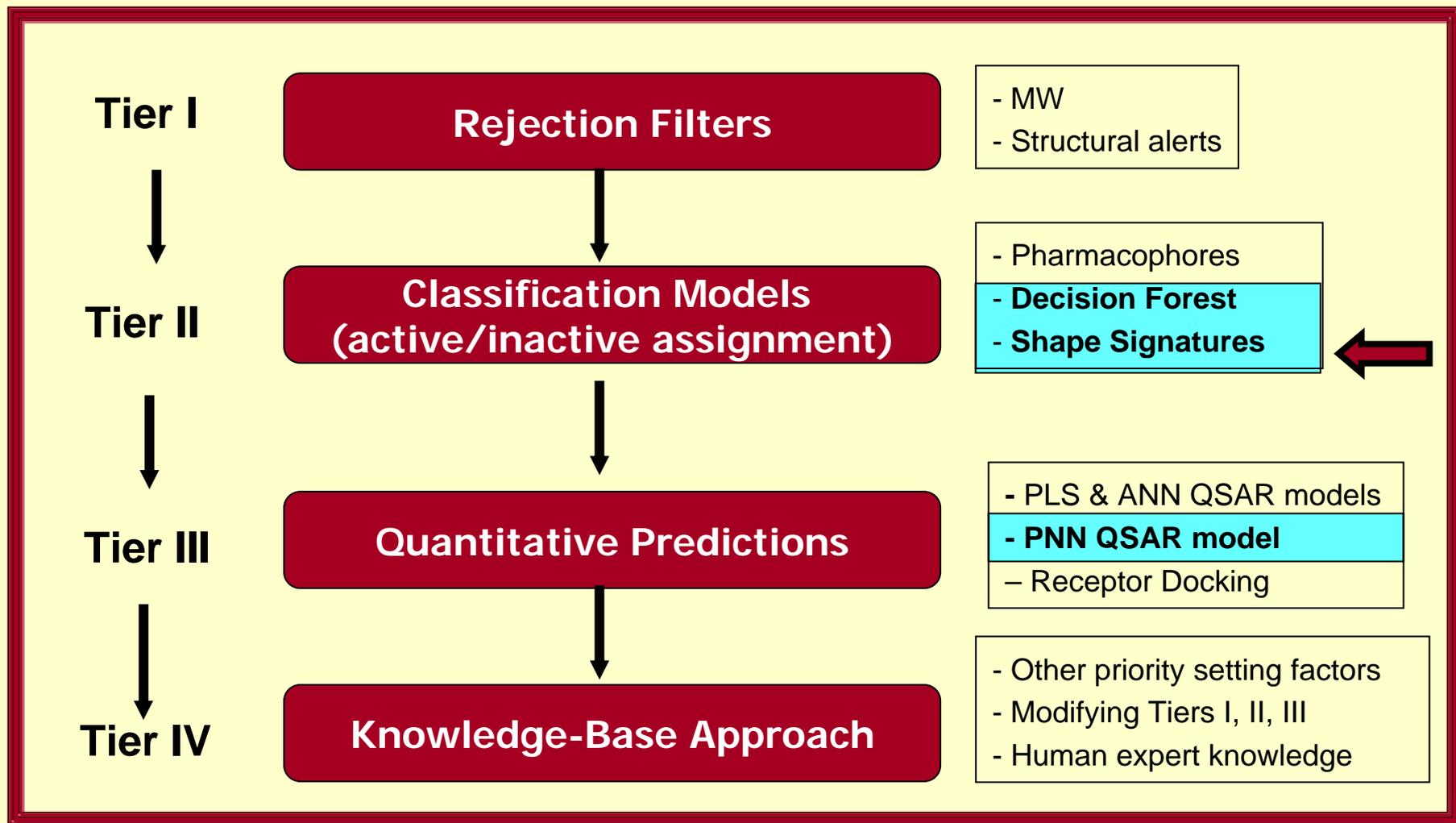


## Key Features

- Combining several independent yet predictive trees reduces misclassification
- DF structure permits assessment of prediction confidence
- Each tree consists of simple 'If-Then' branches, hence the DF is extremely fast

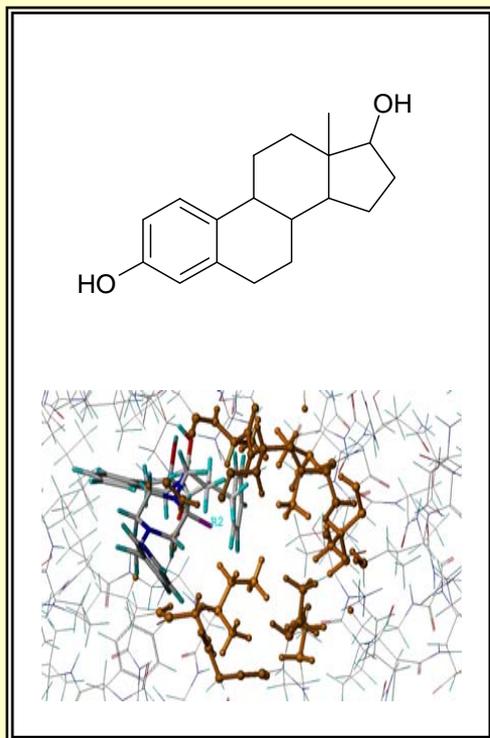
# Schematic of Hierarchical Framework

- addresses the need to minimize *false negatives* and *uncertainties* -



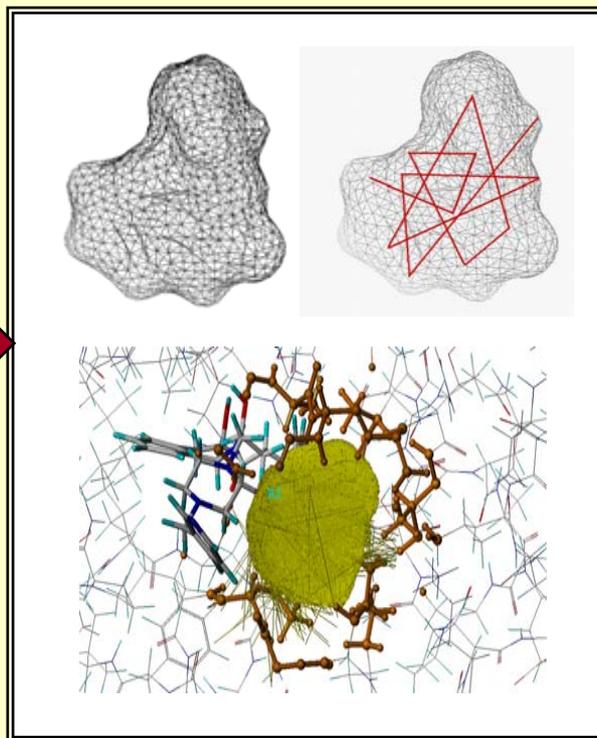
# Shape Signatures Tool

START



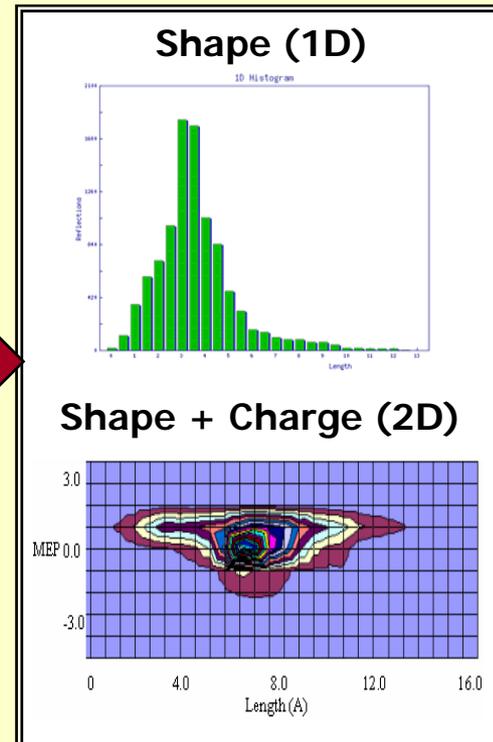
Small molecule or  
Protein binding pocket

PROCESSING



Ray tracing to  
generate the raw data

OUTPUT

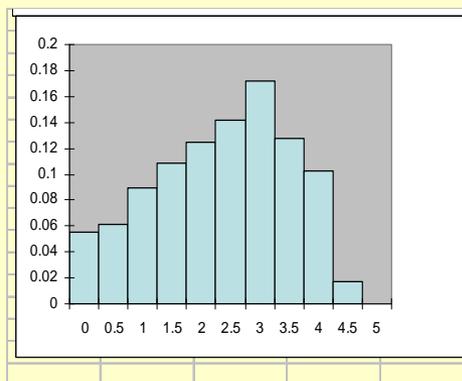
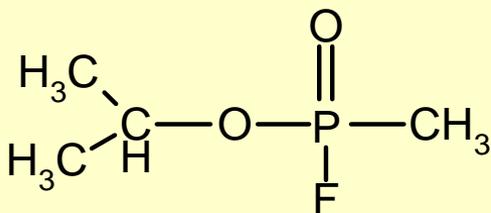


1D and 2D  
Shape Signatures

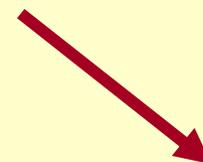
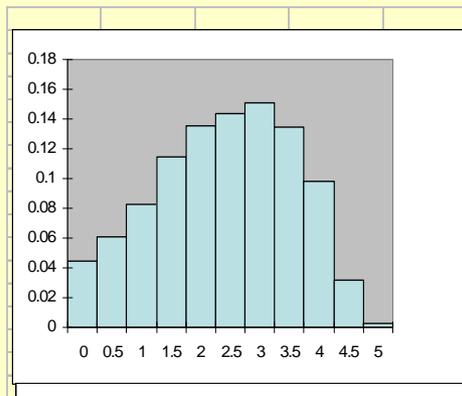
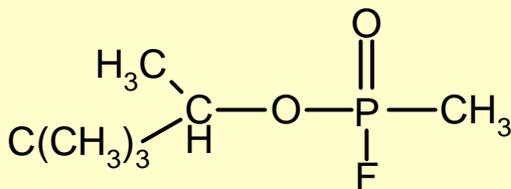
# Shape Signatures Tool

molecules are compared by subtracting their histograms

**Sarin**



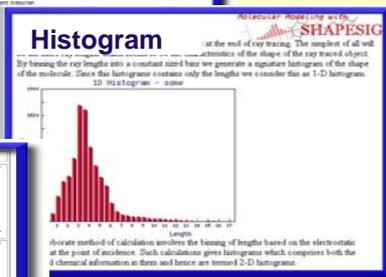
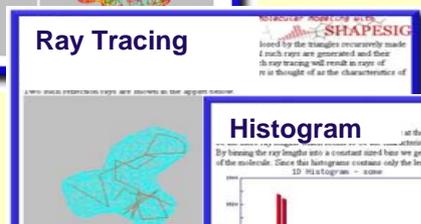
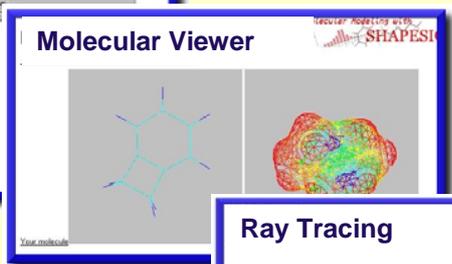
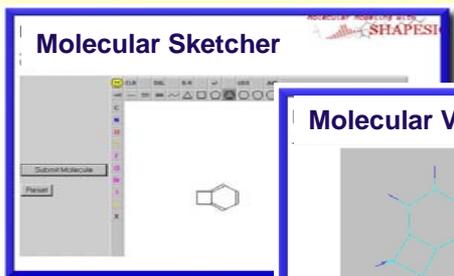
**Soman**

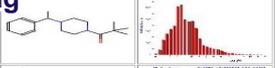
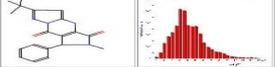
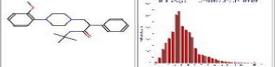
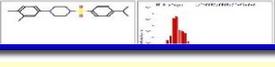


**Diff = 0.022**

Small *Diff* value means that two molecules have similar shape and polarity

# Shape Signatures Software Tool & Chemical Databases



Database Searching			HISTOGRAM
1	HTS_00651	MAYBRIDGE 0.0484	
2	HTS_08105	MAYBRIDGE 0.0551	
3	WAY-100135	WDI 0.0597	
4	ST4074848	GPCR 0.0615	

## Searchable *Shape Signatures* Databases

- 3+ million commercially available organic compounds
- 40,000 Natural Products
- Hazardous Chemicals (pesticides, nerve agents, mustards, psychotropic agents, other real or potential CWAs, TICs)
- PDB-extracted ligands

# Chemical → Target Protein → Mechanisms

Protein Data Bank (PDB): World Repository of ~35,000  
Protein-Ligand Crystal Structures (<http://www.rcsb.org/pdb/>)

In this page you can select organisms. The protein ID and the 2D images of the ligands will be displayed in a table form.

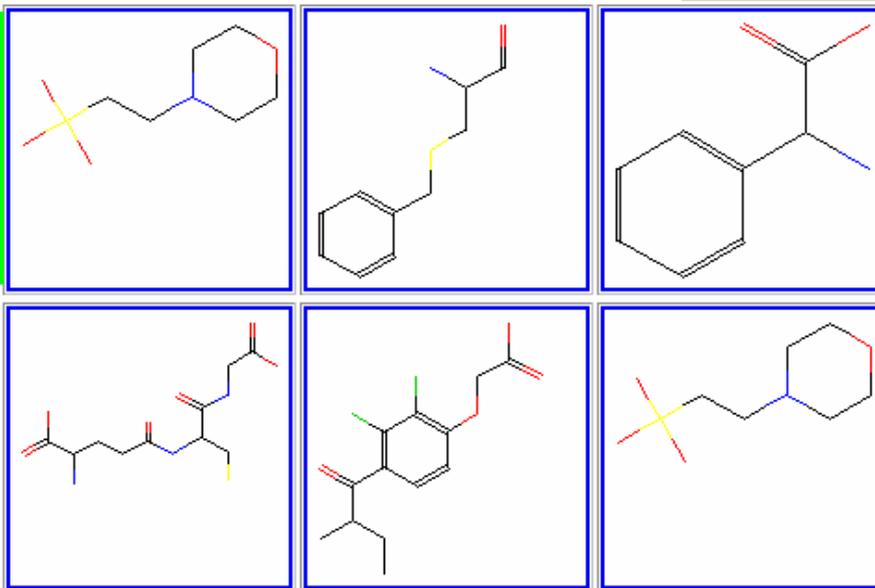
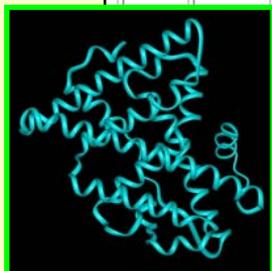
HUMAN (1751)

Submit Molecule

## Shape Signatures of PDB-extracted ligands

Here are the Results obtained by searching for **HUMAN** :

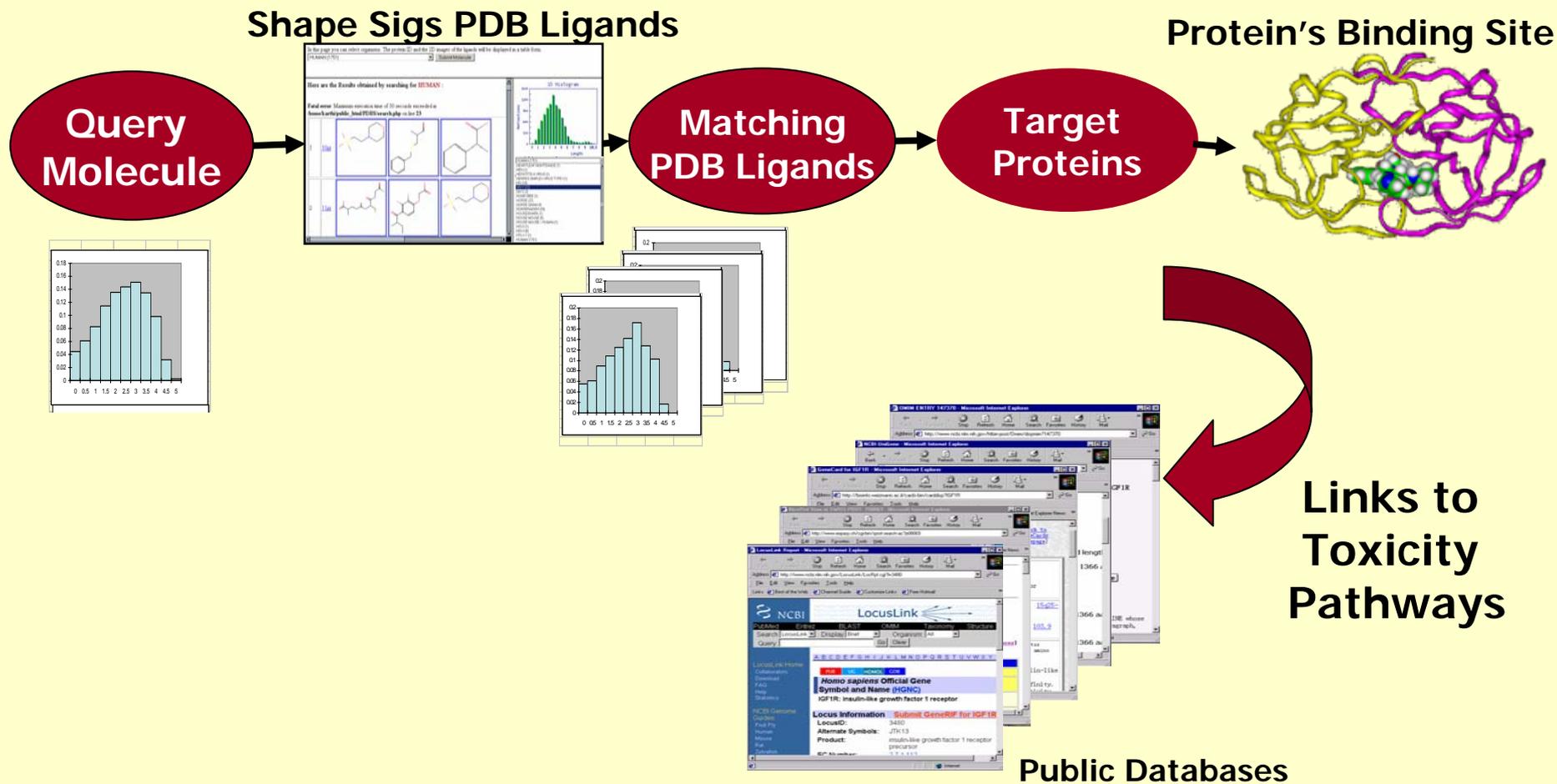
### Protein Structure



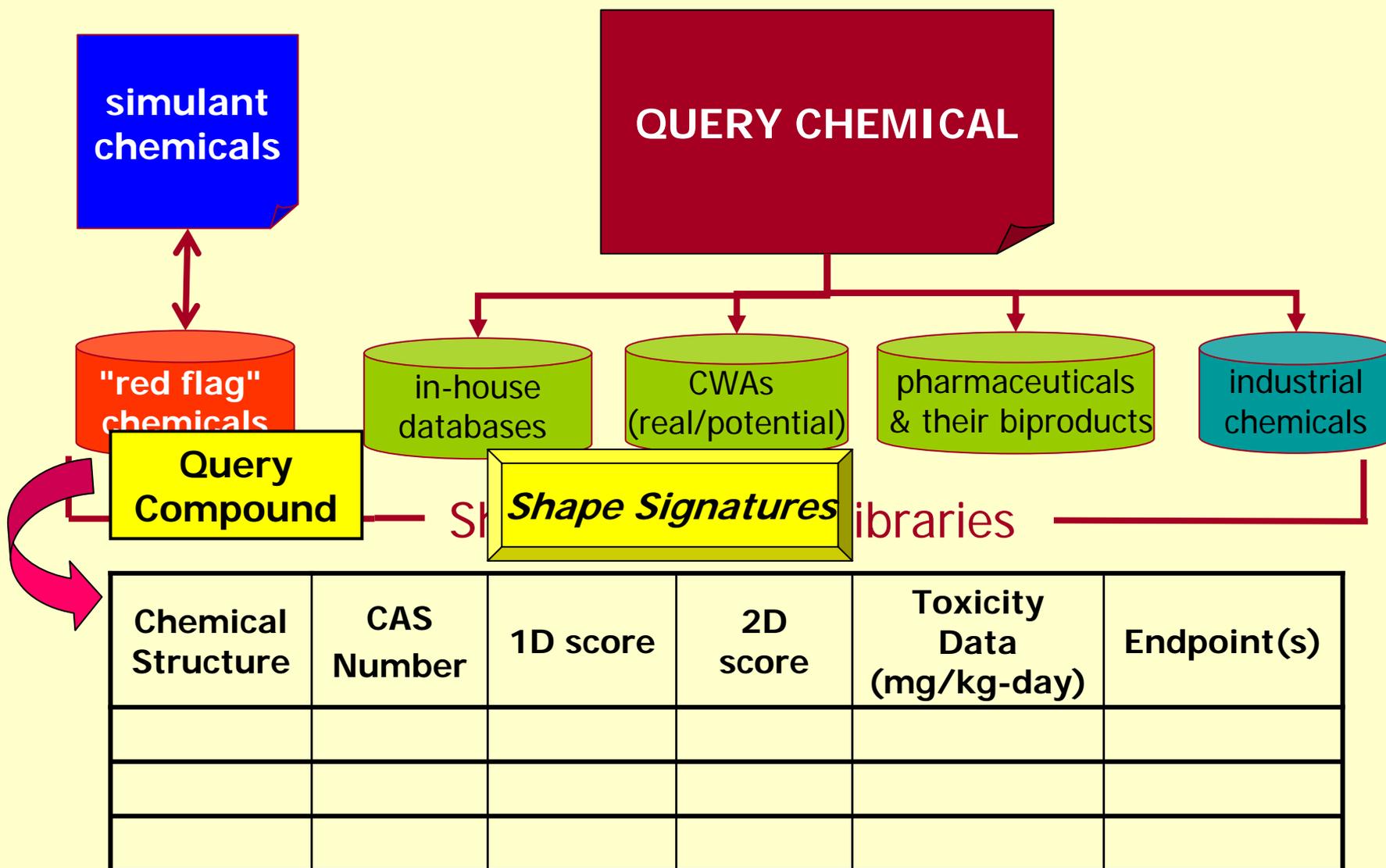
- HUMAN (1751)
- HEARTLEAF NIGHTSHADE (1)
- HEN (1)
- HEPATITIS A VIRUS (1)
- HERPES SIMPLEX VIRUS TYPE-1 (1)
- HIV (10)
- HIV-1 (72)**
- HIV-2 (3)
- HONEYBEE (1)
- HORSE (37)
- HORSE GRAM (4)
- HORSERADISH (28)
- HOUNDSHARK (1)
- HOUSE MOUSE (5)
- HOUSE MOUSE + HUMAN (1)
- HSV2 (1)
- HSV-1 (6)
- HTLV-1 (1)
- HUMAN (1751)

Species/Protein Family

# Molecules → Target Protein → Mechanism



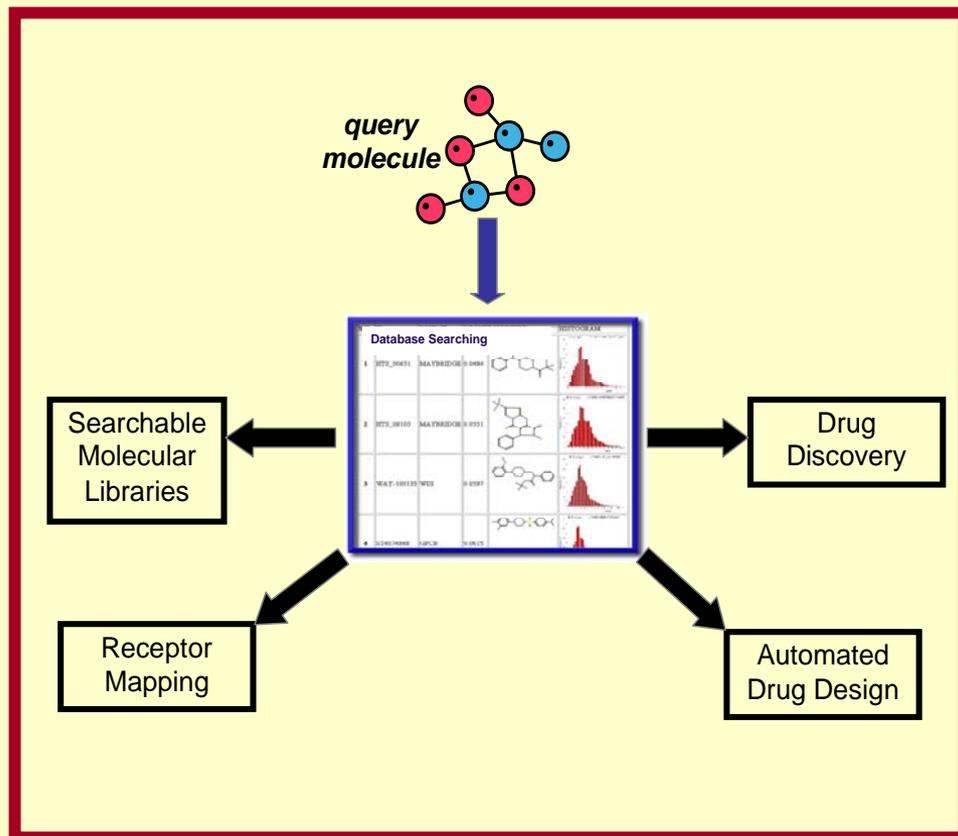
# Identifying Problem Chemicals



# Shape Signatures

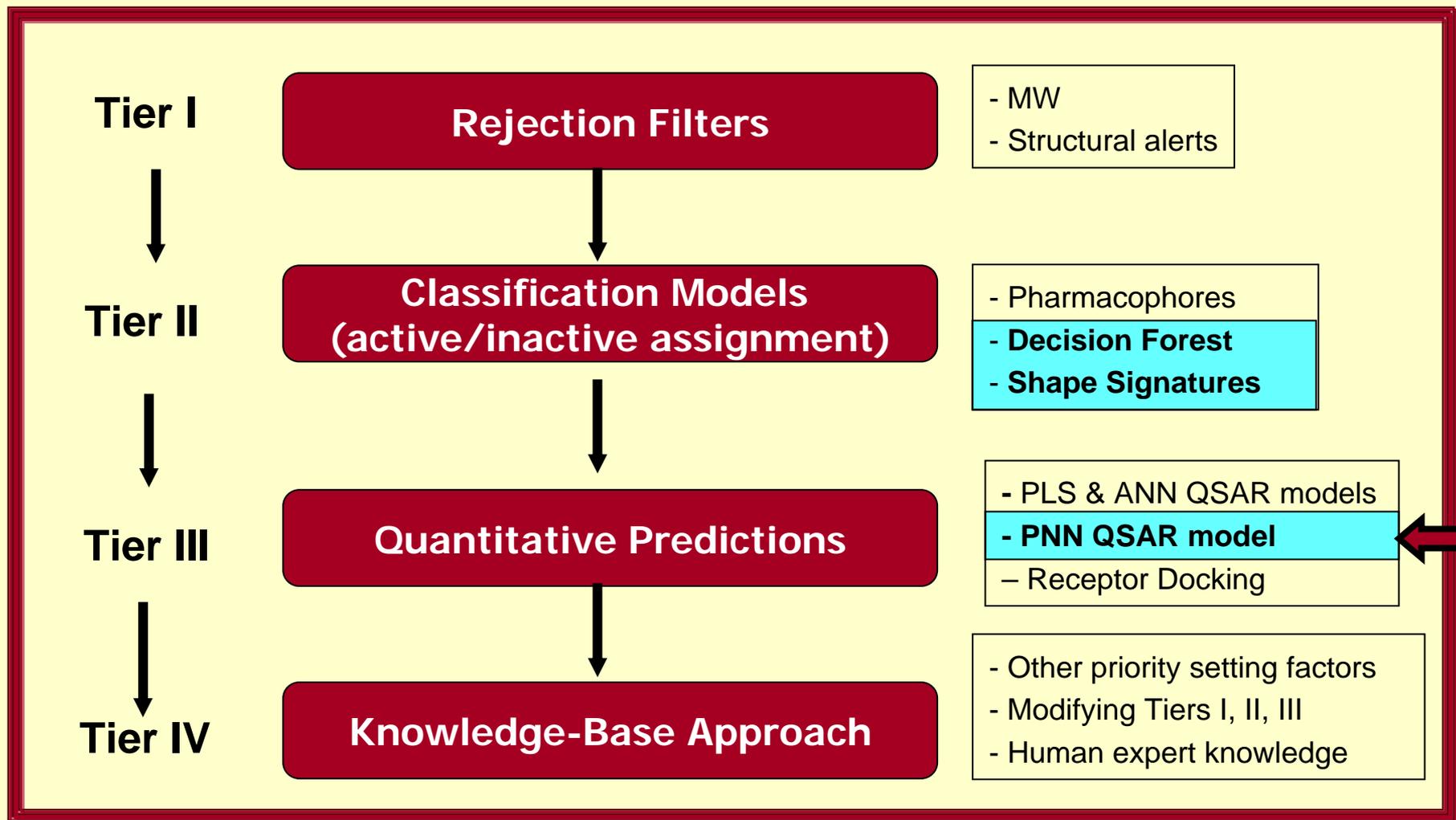
## - Key Features -

- Fast  
*screens large databases in secs*
- Extensible  
*works with any kind or number of molecular species*
- Portable  
*works on any platform*
- Versatile  
*broad utility, multiple databases*



# Schematic of Hierarchical Framework

- addresses the need to minimize *false negatives* and *uncertainties* -



# Building QSAR Models

target property  $\propto$  (molecular descriptors)

$$Y = f(X_i)$$

## Types of Molecular Descriptors

Type	Example
Constitutional	Molecular composition ( $M_w$ , # of atoms/bonds, # of H-bond donors/acceptors)
Topological	2-D structural formula (Kier-Hall indices, extent of branching)
Geometrical	3-D structure of molecule (molecular volume, solvent accessible surface area, polar and non-polar surface area)
Electrostatic	Charge distribution (atomic partial charges, electronegativities)
Quantum Mechanical	Electronic structure (HOMO-LUMO energies, band gap, dipole moment)

# Comparison of Regression Methods

## Desirable Features of Methods and Models

- predictions should be fast
- produces linear or non-linear models (i.e., relationship between obs toxicities and calc'd molecular features may be non-linear)
- models should be physically meaningful, interpretable, and assume parametric form

<b>Method</b>	<b>Speed</b>	<b>Linear Models?</b>	<b>Nonlinear Models?</b>	<b>Regression Equation?</b>	<b>Easy to Interpret?</b>
<b>PLS/MVR</b>	**	Yes	No	Yes	Yes
<b>ANN</b>	*	Yes	Yes	No	Yes
<b>PNN</b>	**	Yes	Yes	Yes	Yes

# Polynomial Neural Network (PNN)

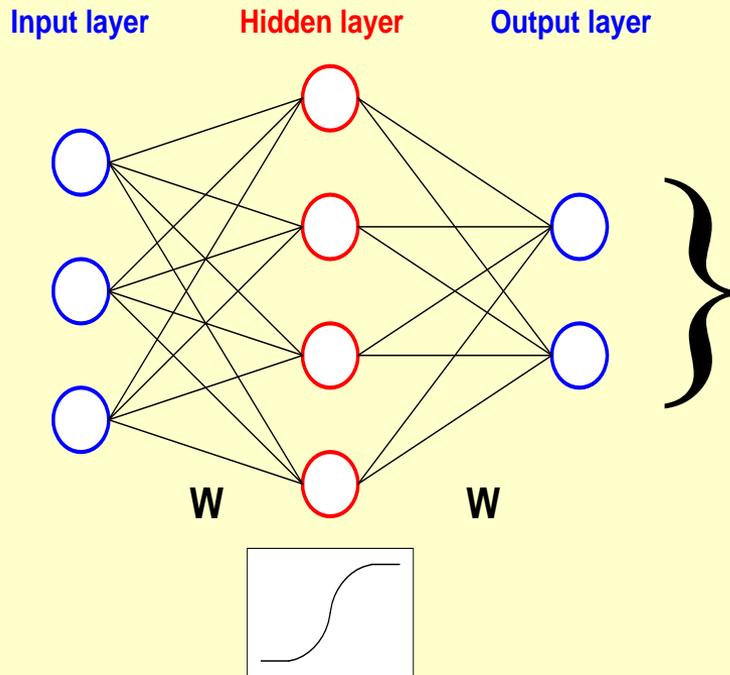
- combines best features of linear multivariate models (parametric form) and ANN models (nonlinearity) -

## Polynomial Neural Network

The screenshot displays the Pnn Discovery Client interface. At the top, there is a menu bar (File, View, Project, Tools, Help) and a toolbar. Below this is a data table with columns: diameter, pettjean, pettjeanSC, radius, VDistEq, VDistMa, weinerPath, weinerPol, a\_aro, a\_count, a\_IC, a\_ICM, a\_nH, and b\_trot. The table contains 15 rows of numerical data. A 'Project Settings' dialog box is open, showing options for 'Type of regression equation' (Linear), 'Applied criteria type' (RSS), 'Maximal terms number in equations' (40), 'Maximal iterations number' (50), 'Maximal number of saved models' (3), and 'Maximal equations degree' (2). Below the dialog box is a 'Settings' panel with a neural network diagram showing an input layer (blue circles), hidden layer (red circles), and output layer (blue circles), with weights 'W' indicated. The settings panel lists various parameters and their values, such as 'Type of regression equation: Linear equation', 'Applied criteria type: RSS', and 'Maximal iterations number: 50'.

- Produces linear or non-linear QSAR models in parametric form
- User control of model complexity
- Insensitive to irrelevant variables and outliers
- Yields predictive models, even for sparse or noisy data sets
- Trains rapidly, thus amenable to large data sets
- Automatically selects best models
- Customizable to fit user's needs

# Polynomial Neural Network (PNN)



1) PNN generates parametric solutions of any desired order 'n':

$$\text{Act.} = w_1(\text{SA}) + w_2(V) + w_3(\mu) + \dots$$

$$\text{Act.} = w_1(\text{SA}) + w_2(V)^2 + w_3(\mu)^3 + \dots$$

$$\text{Act.} = w_1(\text{SA})^2 + w_2(V) + w_3(\mu)^2 + \dots$$

$$\text{Act.} = w_1(\text{SA})^0 + w_2(V) + w_3(\mu)^2 + \dots$$

$$\text{Act.} = w_1(\text{SA}) + w_2(V)^2 + w_3(\mu)^2 + \dots$$

2) PNN selects best solutions:

$$\text{Act.} = w_1(\text{SA}) + w_2(V)^2 + w_3(\mu)^3 + \dots$$

$$\text{Act.} = w_1(\text{SA}) + w_2(V)^2 + w_3(\mu)^2 + \dots$$

***Thank You!***

***welshwj@umdnj.edu***