**NanoSAR:** Structure-Activity Relationship Model for the Toxicity of nano particles
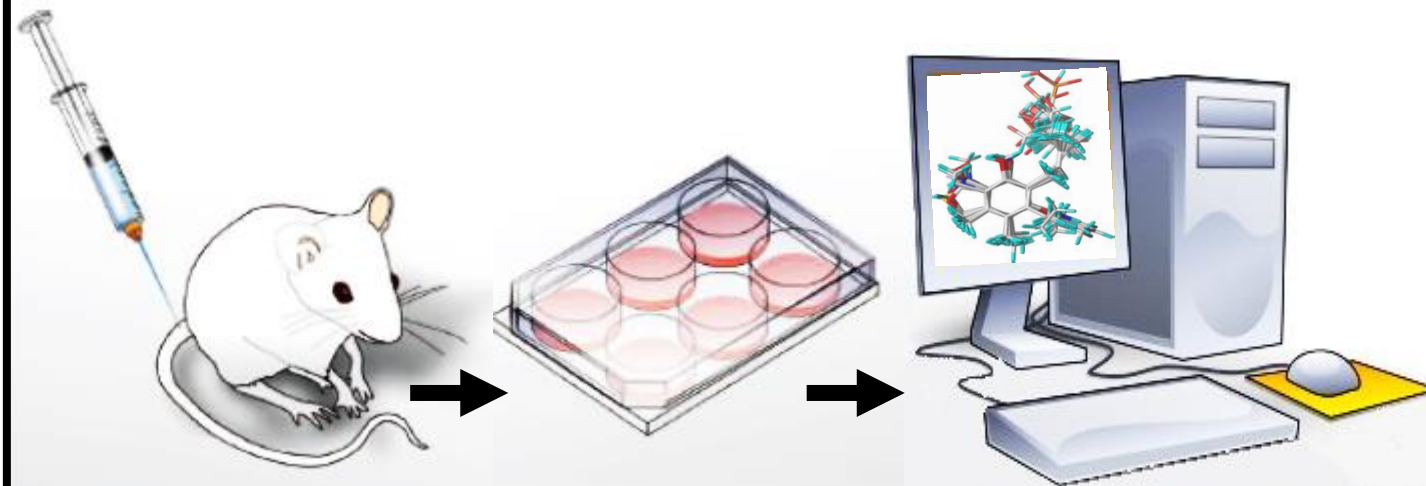
**<u>Ceyda OKSEL</u>**

**Xue Z Wang**

# Structure of the lecture

❖ BACKGROUND

- Why are things different at nanoscale ?

- Nanomaterial toxicity
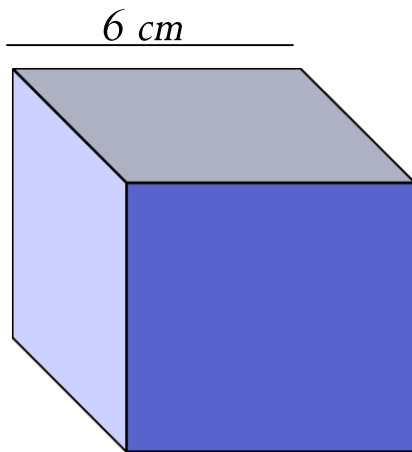
- Computational models for toxicity prediction

❖ COMPUTATIONAL MODELLING OF NANOMATERIAL TOXICITY

- What is (nano)QSAR ?
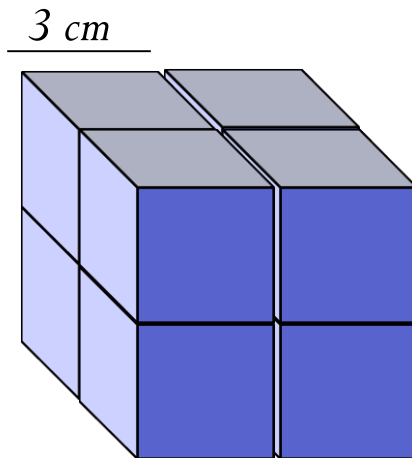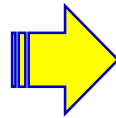
- 3 Case Studies

❖ CONCLUSIONS and FUTURE WORK

# Why are things different at nanoscale?
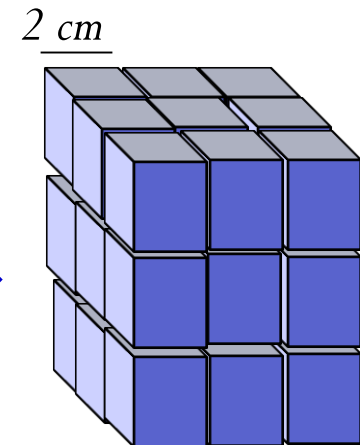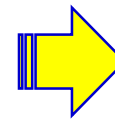
**Larger surface area**

6 cm

3 cm

2 cm

*Surface area*
=(6cm x 6cm x 6 faces x 1 cube)
=216cm²

*Surface area*
=(3cm x 3cm x 6 faces x 8 cubes)
=432cm²

*Surface area*
=(2cm x 2cm x 6 faces x 27 cubes)
=648cm²

**Quantum effects**

# Nanomaterial Toxicity

**New Technology**

**New Benefits**

opportunities

Uncertainties at early stages

**ACCEPTANCE**

**New Concerns**[1]

- Safety to human health and environment
- Suitability of risk management strategies

risks are identified

uncertainties are dealt

# Nano Particles, Mega Problems ?

*μm*

*nm*

**Reduction in particle size**

**Nano-specific Toxicity**

Change in toxicity of NPs

**Nano-products**

Increase in commercial nanoproducts

**Nanosafety Concerns**

Increase in nanosafety concerns

# Toxicity Testing

**IN SILICO**

(computational)

*predict*

*validate*

**IN VITRO**

(computational)

*predict*

*validate*

**IN VIVO**

(on living organism)

Experimental

Reduction in time, cost and animal testing

**QSAR**

**Q**uantitative **S**tructure-**A**ctivity **R**elationship models

# Why we need computational models?



**Innovations in Nanotechnology**

**Hazard Assessment of NMs**

Complexity

Knowledge gaps
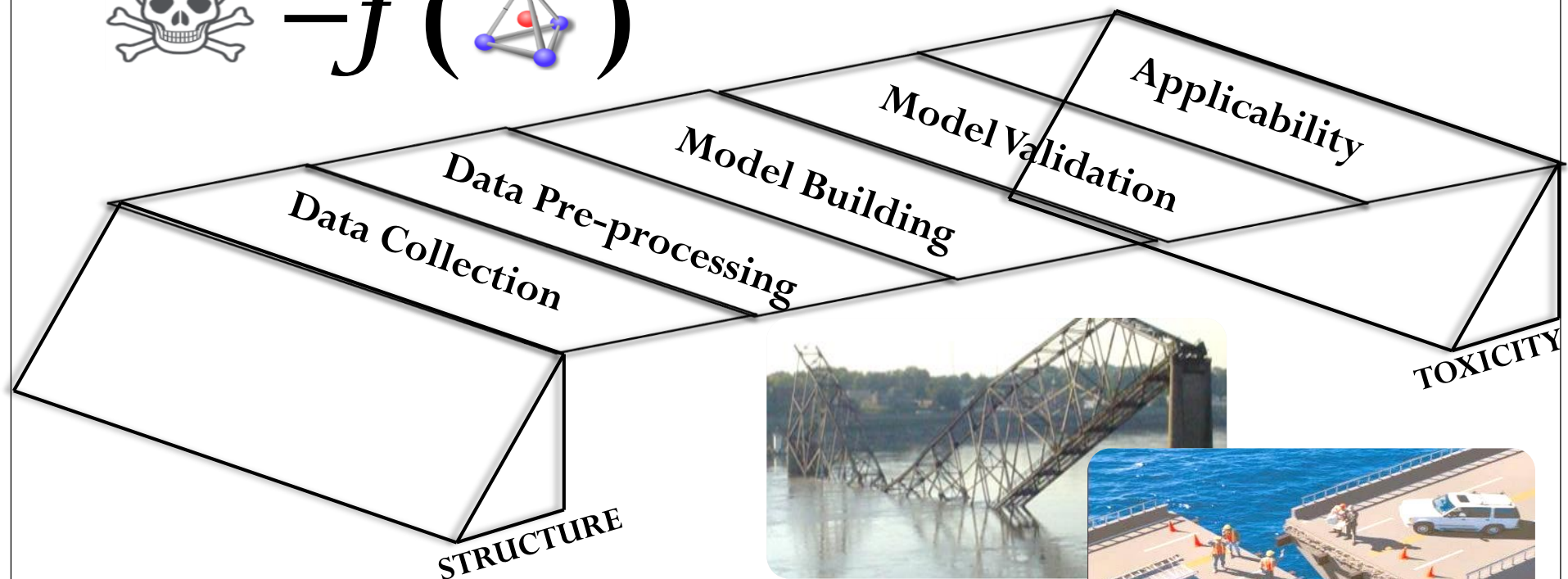
Effort, time and cost required

**NEED:** *The European REACH legislation promotes the use of non-animal testing methods*

**AIM: to satisfy this need!!!**

# What is nano-(Q)SAR ?

**A (Q)SAR is a statistical model that relates a set of physicochemical descriptors of a chemical compound to its biological activity.**



STRUCTURE

Data Collection

Data Pre-processing

Model Building

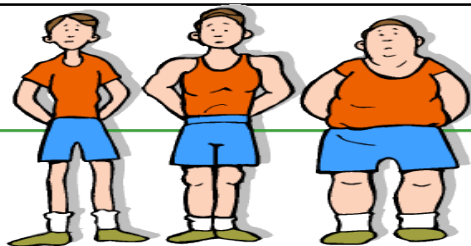Model Validation

Applicability

TOXICITY

Oksel, C., C.Y. Ma, and X. Z. Wang. "Current situation on the availability of nanostructure–biological activity data." *SAR and QSAR in Environmental Research* ahead-of-print (2015): 1-16.

Oksel, C., C.Y. Ma, J. J. Liu, T. Wilkins, X. Z. Wang, (2015) (Q)SAR modelling of nanomaterial toxicity: A critical review, Particuology, 10.1016/j.partic.2014.12.001
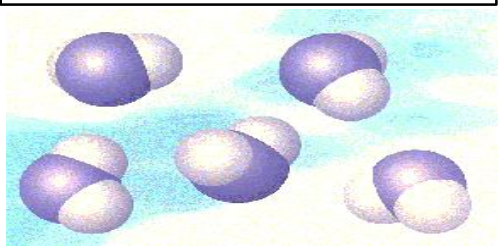
# Descriptors

| DESCRIBING A PERSON | DESCRIBING A MOLECULE | DESCRIBING A NANOPARTICLE |
|---|---|---|

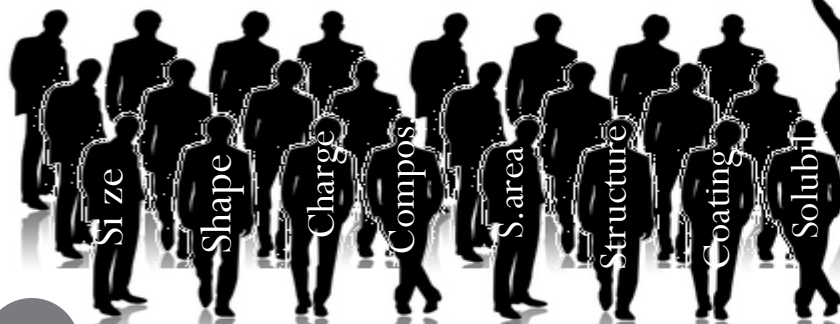

- ✓ Height
- ✓ Weight
- ✓ Attractiveness
- ✓ ...
- ✓ Eye
- ✓ Hair
- ✓ Build
- ✓ ...

- ✓ Molar mass
- ✓ Density
- ✓ Conductivity
- ✓ ...
- ✓ Atomic prop.
- ✓ Bonds
- ✓ Chirality
- ✓ ...

- ✓ Size
- ✓ Shape
- ✓ Composition
- ✓ ...
- ✓ Coating
- ✓ Charge
- ✓ Reactivity
- ✓ ...

**Experimental Descriptors**

**Descriptor Selection**
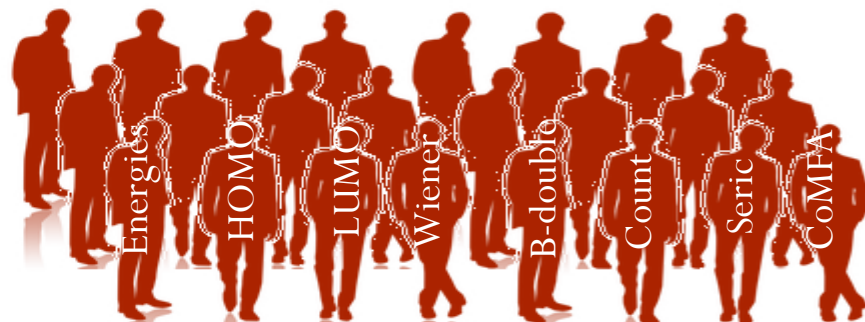Feature selection algorithms

**Theoretical Descriptors**

Size  Shape  Charge  Compos.  S.area  Structure  Coating  Solubil.

Energies  HOMO  LUMO  Wiener  B-doub.  Count  Seric  CoMFA

9

# Tree Induction From Genetic Programming

## GPTree: "in-house" software

Genetic Algorithms

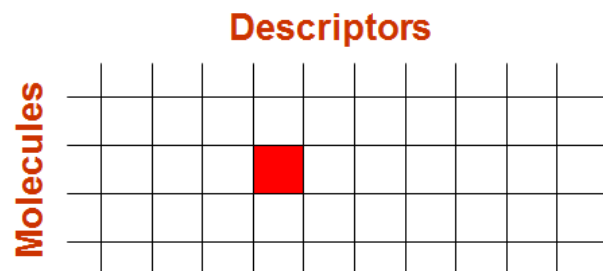| explore solution space | • Starts at random points<br>• Recombining (i.e., crossover)<br>• Optionally changing (i.e., mutation) |
|---|---|
| **Genetic Algorithm** | **(1) Randomly generate a pre-specified number of solutions, encoded as fixed size vectors.**<br><br>**(2) Either form a new generation or replace individuals in the population by**<br><br>2a. Selecting parents using the fitness function.<br><br>2b. Crossover the parents to form one or more offspring.<br><br>2c. Optionally mutate part of the solution.<br><br>**(3) Continue with Step 2 until a pre-specified number of generations or children have been grown, or until a good solution is found.** |

## GPTree: Methodology

- **DeLisle, R. K. and Dixon, S. L.** (2004) Induction of Decision Trees via Evolutionary Programming *Journal of Chemical Information and Computer Sciences,* **44,** 862-870.- **evolutionary programming of trees**
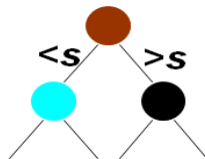
**1. Divide data into training and test sets**

**2. Generate the 1st population of trees**

- randomly choosing a row (i.e. a compound), and column (i.e. descriptor)

**Descriptors**

**Molecules**

- Using the value of the slot, *s,* to split, left child takes those data points with selected attribute values <= *s*, whilst the right child takes those > *s*.

<s   >s

# Tree Induction From Genetic Programming

## GPTree: Methodology

- If a child will not cover enough rows (e.g. 10% of the training rows), another combination is tried.

- A child node becomes a leaf node if pure/near pure, whilst the other nodes grow children.

-When all nodes either have two children or are leaf nodes, the tree is fully grown and added to the first generation.

-A leaf node is assigned to a class label corresponding to the majority class of points partitioned there.

### 3. Crossover and Mutation

# Tree Induction From Genetic Programming

**The key parameters**

| | |
|---|---|
| **y COL** | Column no containing the class of the data set. |
| **n Gen** | No of generations required |
| **n Trees** | No of treesrequired in each generation |
| **No. in tournament** | No of trees in the tournament to sort out the best for crossover operation |
| **Winn. Inc.** | Winners included (The N best trees are placed directly into the next generation, This was to allow ELITISM) |
| **L.I.I.A.T** | Low increase in accuracy tolerance (It forces a mutation for every tree if no improvement in the best accuracy has been seen for this many generations.) |
| **Mutation** | % age of mutation |
| **C in L.N** | Minimum no of cases in a leaf node |

# Case Study 1: Dataset

| **Compounds** | 75 Compounds |
| --- | --- |

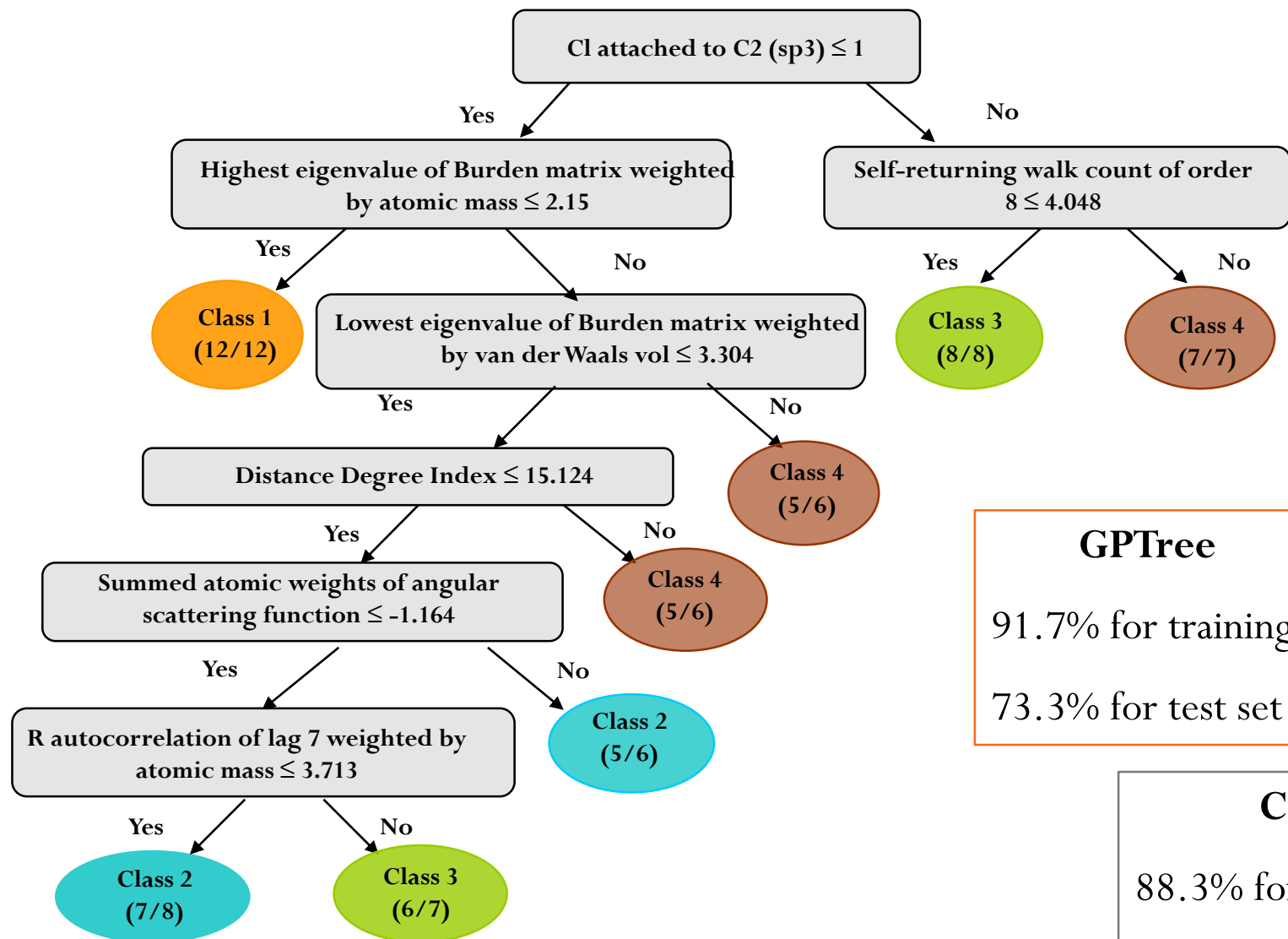| **Toxicity Data (4 classes)** | Concentration lethal to 50% of the population, LC50, 1/Log(LC50), of *vibrio fischeri, a biolumininescent bactorium* |
| --- | --- |

| **Descriptors** | 1069 molecular descriptors calculated by DRAGON |
| --- | --- |

## Parameters

| | |
| --- | --- |
| **y COL** | 1070 |
| **n Gen** | 60 |
| **n Trees** | 600 |
| **No. in tournament** | 16 |
| **Winn. Inc.** | 0 |
| **L.I.I.A.T** | 5 |
| **Mutation** | 66.7% |
| **C in L.N** | 2 |

# Case Study 1: Results

# Case Study 2: Dataset

| | |
|---|---|
| **Compounds** | 105 nanoparticles with different surface-modifying molecules |

| | |
|---|---|
| **Toxicity Data** | Cellular uptake in pancreatic cancer cell lines |

| | |
|---|---|
| **Threshold value** | Cellular uptake values:170-27 542 nanoparticles per cell<br>Threshold value: 10 000 nanoparticles per cell<br>18 nanoparticles with significant cellular uptake (CLASS 2)<br>87 nanoparticles with poor cellular uptake (CLASS 1) |

D. Fourches, D. Pu, C. Tassa, R. Weissleder, S.Y. Shaw, R.J. Mumper, and A. Tropsha, Quantitative nanostructure–activity relationship modeling, ACS Nano 4 (2010), pp. 5703–5712.

# Case Study 2: Dataset

## Descriptors

Nanoparticles $\longrightarrow$ Same core
Different surface-modifying molecules $\longrightarrow$ Conventional descriptors



INITIAL LIST OF SMILES

1. Removal of mixtures, inorganics (and eventually organometallics)

Structural conversion
Cleaning/removal of salts

2. Normalization of specific chemotypes

Treatment of tautomeric forms

3. Analysis/removal of duplicates

4. Manual inspection

CURATED DATASET

Fourches et al. (2010)

- **Data cleaning**
- **Structural Conversion**

SMILES strings $\longrightarrow$ 2D molecular graphs

(C=NC(=C(N=1)C1O)N=C(N=1)N)CNC(=CC=C(C1)C(O)=O)C

- **Manual inspection**

4 structure unmatched-excluded

- **Descriptor Calculation**

690 Dragon Descriptors

- **Descriptor Cleaning**

389 Dragon descriptors retained

Fourches, Denis, Eugene Muratov, and Alexander Tropsha. "Trust, but verify: on the importance of chemical structure curation in cheminformatics and QSAR modeling research." *Journal of chemical information and modeling* 50.7 (2010): 1189-1204.

# Case Study 2: Data Pre-processing

## Data splitting



## Pattern of splitting



**80% training**

**20% validation**

**The key parameters**

```
1    EPTREE Train.txt Test.txt 390 60 600 16 0 5 1 2 2
```

| Column no containing the class of the data set | 390 |
|---|---|
| No of generations required | 60 |
| No of trees in each generation required | 600 |
| No of trees in the tournament | 16 |
| Winners included | 0 |
| Low increase in accuracy tolerance | 5 |
| % age of mutation | 50% |
| Minimum no of cases in a leaf node | 2 |

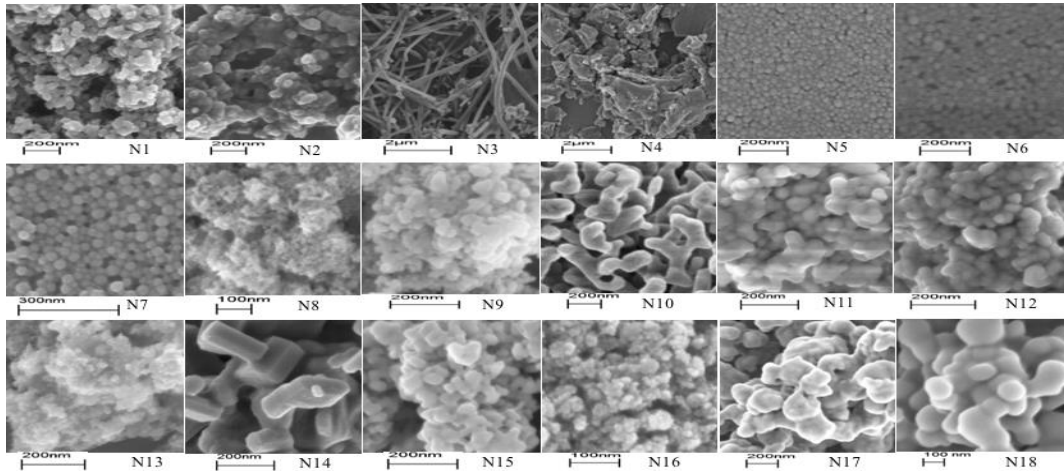# Case Study: Results

**GPTree Results**

# Case Study: Results



**Training accuracy: 96%**
**Test accuracy: 81%**

**9 descriptors out of 389**

# Case Study: Results

| DRAGON descriptor | Description | Block |
|---|---|---|
| JGI2 | mean topological charge index of order 2 | 2D autocorrelations |
| JGI5 | mean topological charge index of order 5 | 2D autocorrelations |
| ATSC8m | Centred Broto-Moreau autocorrelation of lag 8 weighted by mass | 2D autocorrelations |
| ATSC3v | Centred Broto-Moreau autocorrelation of lag 3 weighted by van der Waals volume | 2D autocorrelations |
| MATs6i | Moran autocorrelation of lag 6 weighted by ionization potential | 2D autocorrelations |
| GATS7s | Geary autocorrelation of lag 7 weighted by I-state | 2D autocorrelations |
| Eig05_EA(dm) | eigenvalue n. 5 from edge adjacency mat. weighted by dipole moment | Edge adjacency indices |
| SpMAD B(v) | spectral mean absolute deviation from Burden matrix weighted by van der Waals volume | 2D matrix-based descriptors |
| RBN | number of rotatable bonds | Constitutional indices |

# Case Study 3: Data Collection



Carbon Black **N1**

Diesel Exhaust **N2**

Japanese Nanotubes **N3**

Fullerene **N4**

Polystyrene Latex Beads **N5**

Polystyrene Latex Beads **N6**

Polystyrene Latex Beads **N7**

Aluminuim Oxide **N8**

Aluminuim Oxide **N9**

Aluminuim Oxide **N10**

Cerium Oxide **N11**

Nickel Oxide **N12**

Silicon Oxide **N13**

Zinc Oxide **N14**

Titanium Dioxide Rutile **N15**

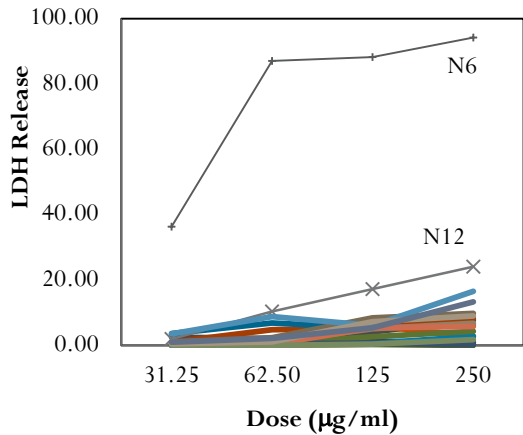Titanium Dioxide Anatase **N16**

Silver **N17**

Silver **N18**

## Characterization

- *Particle size and size distribution* were analysed using a Malvern MasterSizer 2000
- *Particle shape* was analysed using LEO 1530 Scanning Electron Microscope (SEM) or Philips CM20 Transmission Electron Microscope (TEM)
- *Surface area and porosity were measured using* TriStar 3000 BET
- *The free radical activities* were measured by EPR
- *Particle reactivity in solution*, the dithiothreitol (DTT) consumption
- *Metal Content* was measured
- *Charge:* z potential was measured using Malvern Instrument's Zetasizer Nano instrument
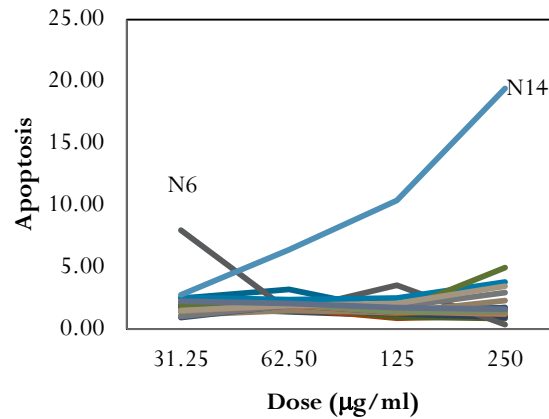
# Case Study 3: Data Collection
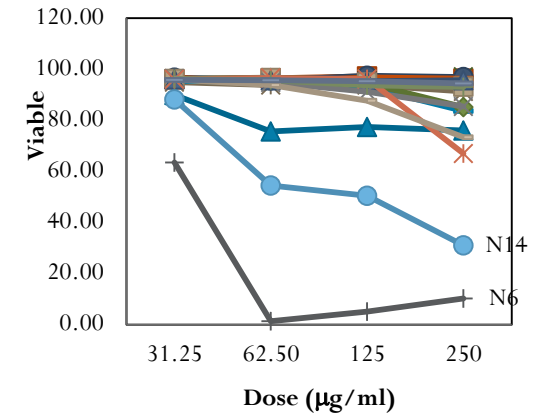
## Toxicological Evaluation

### LDH Release



### Apoptosis



### Viability



### Necrosis



### Pro-inflammation effects



### Haemolysis

# Case Study 3: Data Visualization

**Multidimensional data visualization:**
**Heat maps with hierarchical clustering**

# Case Study3: Model Development

**Clustering/Grouping based on Principal Component Analysis**



**Clustering based on toxicity data**

**Clustering based on Characterization data**

Wang, Xue Z., et al. "Principal component and causal analysis of structural and acute in vitro toxicity data for nanoparticles." Nanotoxicology 8.5 (2014): 465-476.

# Conclusions

- In LEEDS, we have developed a decision tree software which can be successfully employed for nano-(Q)SAR investigations
- (Q)SAR tools are useful for identifying the properties that influence the toxicity

- Many potential profits:
  - An alternative, fast and cheap way of hazard assessment
  - Risk Reduction
  - Safety-by-design

# Future Work

| No | Dataset | Nanomaterials | Toxicity Endpoint | Characterization |
|----|---------|---------------|-------------------|------------------|
| 1 | Wang et al. (2014) | 18 NMs (carbon-based and metal oxides) | LDH release, apoptosis, pro-inflammatory effects, haemolysis, MTT, DiOC6, cell morphology assay | size, surface area, morphology, metal content, reactivity, free radical generation and zeta potential |
| 2 | Shaw et al. (2008) | 50 NMs with diverse core structures | ATP content, reducing equivalents, apoptosis, mitochondrial membrane potential | core composition, coating type, surface modification, size, relaxivities and zeta potential |
| 3 | NANOMMUNE project | 18 NMs | In vitro assays | core, coating, 2 sizes and zeta potential |
| 4 | Puzyn et al. (2011) | 17 metal oxide NMs | Cytotoxicity (EC50) | 12 different quantum-mechanical descriptors |
| 5 | MARINA project | 9 NMs | In vitro assays | experimental descriptors |
| 6 | Weissleder et al. (2005) | 109 NMs with the same core but different surface modifiers | Cellular uptake | theoretical descriptors |
| 7 | B. Yan (private communication) | 80 surface-modified MWCNTs | Protein binding activities, cell viability, nitrogen oxide generation | theoretical descriptors |
| 8 | Liu et al. (2011) | 9 metal oxide NMs | Cytotoxicity (PI uptake) | a set of 10 descriptors |
| 9 | Sayes and Ivanov (2010) | 42 NMs with two cores (differing in concentrations) | Cellular membrane damage (LDH release) | primary particle size, size in water and buffered solutions, concentration and zeta potential |
| 10 | ENPRA project | 10 NMs | In vitro/in vivo assays | size, dustiness, surface area and impurities |
| 11 | Gajewicz et al. (2014) | 18NMs | Cellular viability (LC50) | 18 quantum mechanical descriptors, 11 image descriptors, 3 experimental descriptors |

NANO-FEAR

SUSTAINABILITY of NANOTECHNOLOGY

Thank you !