# Hazard Assessment in the Era of Big Data

**Alexandra Maertens, Ph.D.**

**Toxicologist, Consortium for Environmental Risk Management**
**Green Toxicology @ Johns Hopkins Bloomberg School of Public Health**

# Changing The Paradigm

- **Regulatory approaches to chemical hazard have been focused on "black box" animal assays**

- **Protects against hazard but only very slowly - example asbestos, BPA**

- **Provides little mechanistic information to guide molecular design**

- **Assumes humans are 70 kg rats**

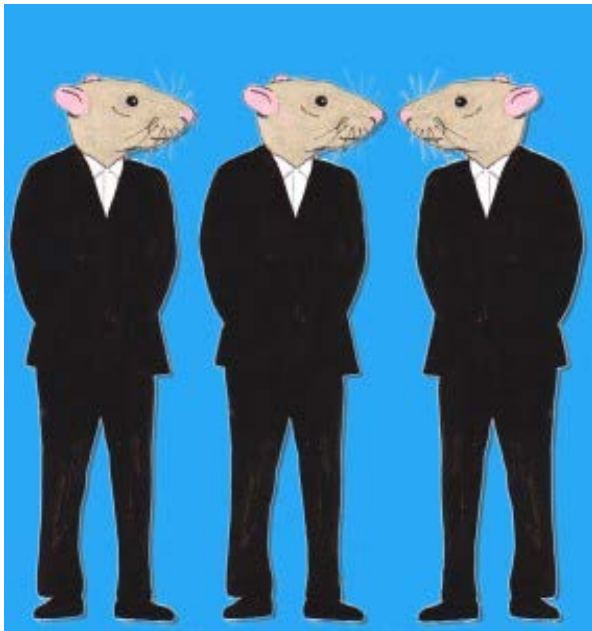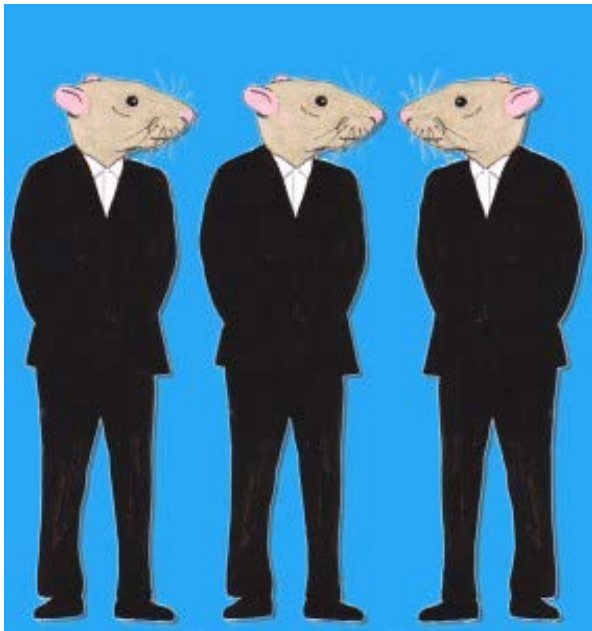- ■**Impossible to cope with the numbers of new chemicals being produced**

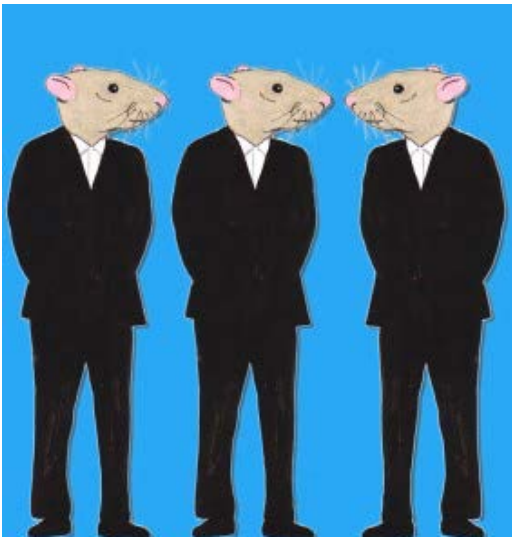Image from: http://antidote-europe.org

# Changing The Paradigm

- **Regulatory approaches to chemical hazard have been focused on "black box" animal assays**

- **Protects against hazard but only very slowly - example asbestos, BPA**

- **Provides little mechanistic information to guide molecular design**

- **Assumes humans are 70 kg rats**

- **Impossible to cope with the numbers of new chemicals being produced**

You mean it's 2019 and they are still using methods from the early 20th century?

Image from: http://antidote-europe.org

- In vitro testing
  - In vitro tests act as drop-in replacements for animal tests – BCOP for Draize tests
    - < 5 per cent of tests in ECHA
  - In vitro tests for mechanism - Ames test for mutagenicity
- In silico predictions
  - Several models exist that are widely used – Toxtree, SimCyp (FDA) to predict metabolism, but mostly for screening purposes
- Read-across/QSAR
  - 80 per cent of ECHA dossiers use read-across for at least one end point and approximately 20 per cent of endpoints
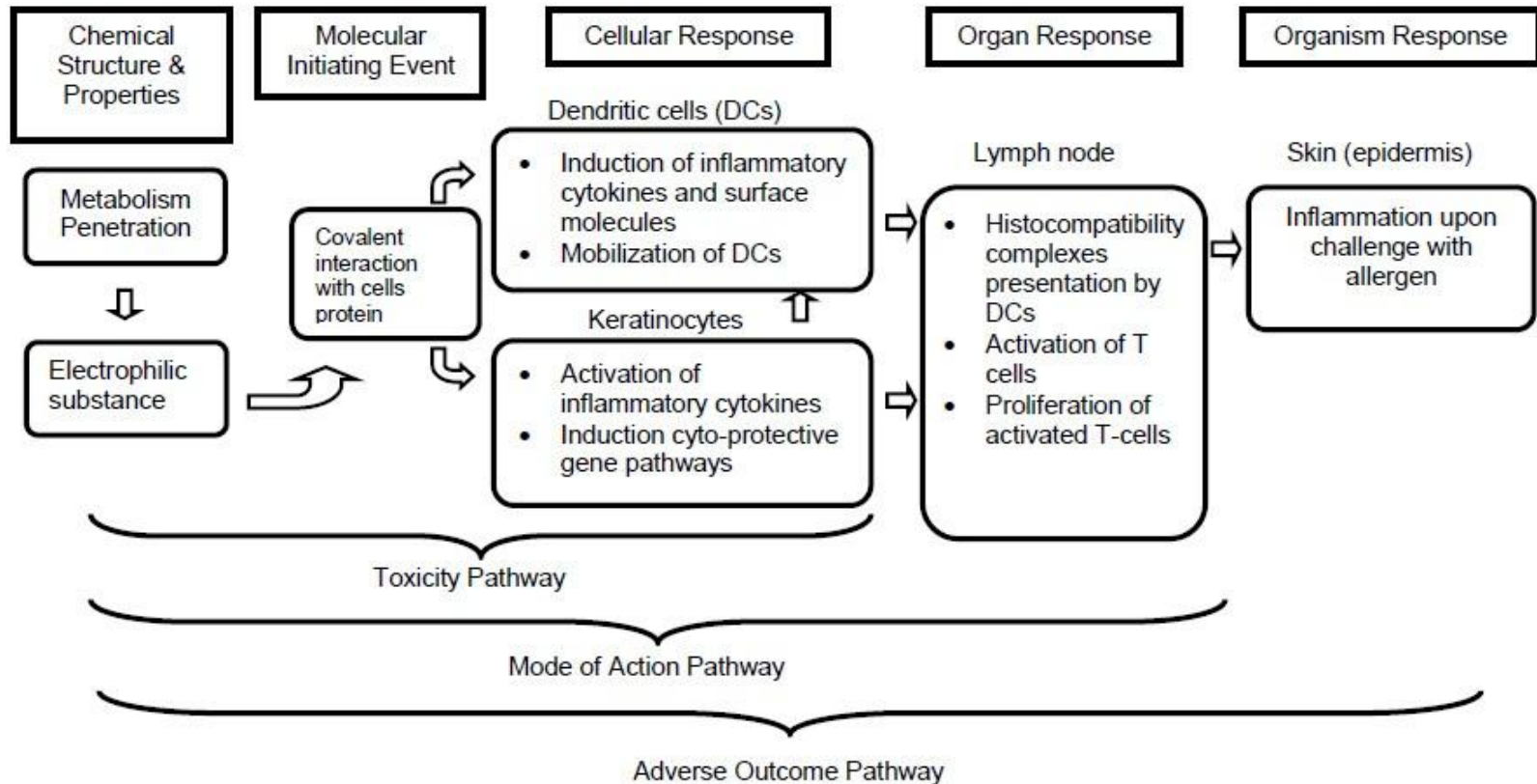  - 1 per cent use QSAR

# Case study: Skin Sensitization

Skin sensitization is an allergic response upon contact to a substance



Representation of skin irritation/sensitization

# What is sensitization and why is it important?

- 15-20% of the general population suffers from contact allergy (Thyssen et al., 2007) - prevalence may be increasing
- US Bureau of Labor Statistics (BLS) data shows that occupational skin diseases currently account for 10-15% of all occupational illness
- Occupational contact dermatitis is particularly prevalent in the personal services industry, with an estimated prevalence of 1.2% percent, e.g., in the beauty/ haircare industry (Warshaw et al., 2012),
- High prevalences in the petrochemical, rubber, plastic, metal and automotive industries (McDonald et al., 2006).
- Many chemicals used in occupational settings have NOT been tested for skin sensitization potential - prior to TSCA 21, only one in eight premarketing notifications to EPA is submitted with any toxicological data.
- However, EPA has of late been very proactive about sensitization

# Skin sensitization: Mechanism



OECD. (2012). The Adverse Outcome Pathway for Skin Sensitisation Initiated by Covalent Binding to Proteins Part 1: Scientific Evidence. [Series on Testing and Assessment No.168 ENV/JM/MONO(2012)10/PART1
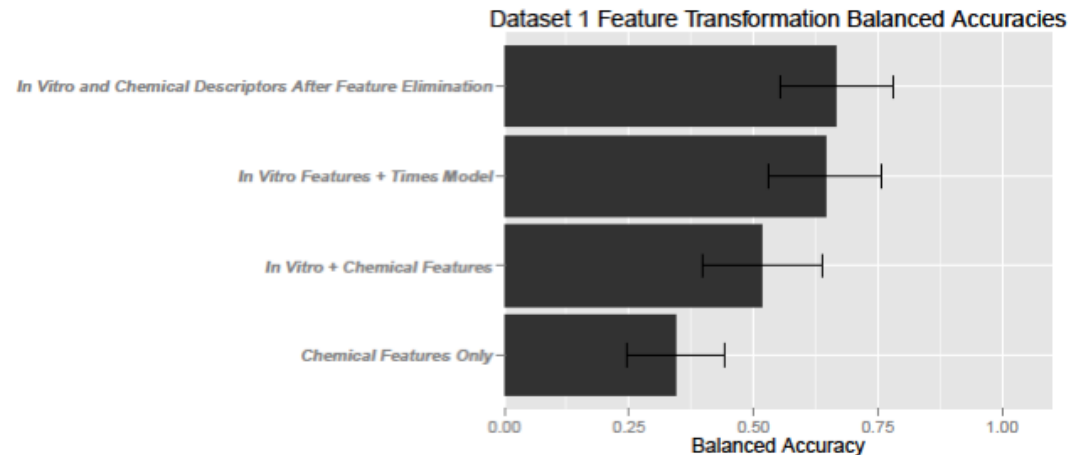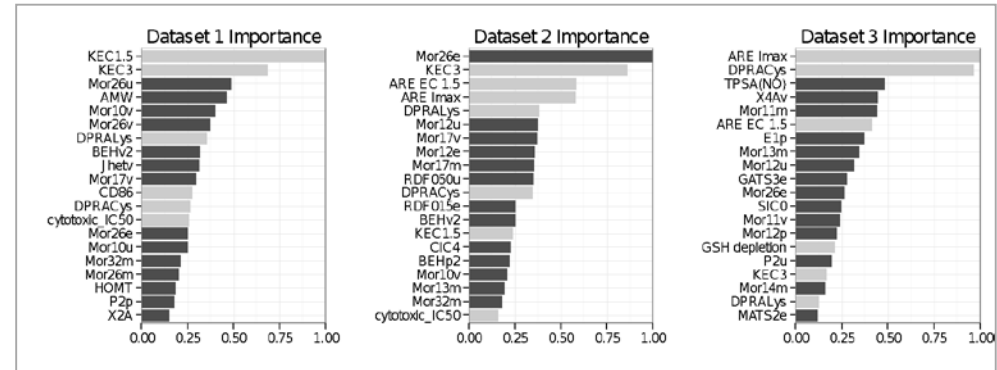
# Integrated Testing Strategy

- Use machine learning to combine multiple sources of information at each point of the AOP

- 145 chemicals with LLNA testing (reference classification) and additional *in vitro* and *in chemico* tests

  - ARE induction
  - Cytotoxicity
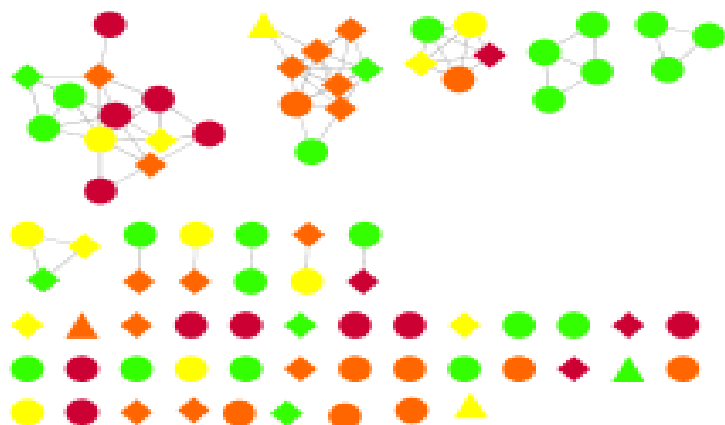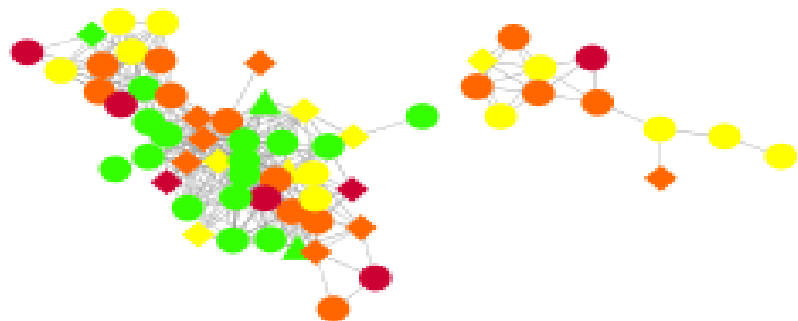  - Direct Peptide Reaction Assay

# Integrated Testing Strategy

- Variable importance consistently included *in vitro* assays in the top

- More data is not always better; accuracy improved by <u>reducing</u> features

- Using in vitro and chemicals descriptors achieved 80 percent accuracy
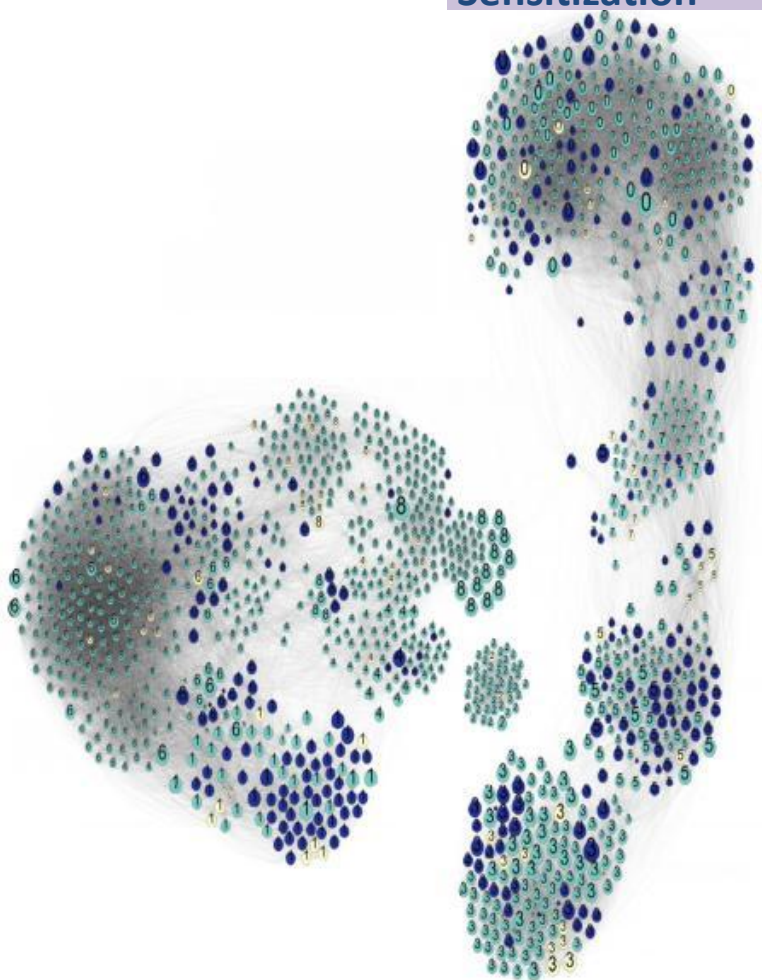
# Chemical Similarity and Skin Sensitization



- Tanimoto distance was used to calculate chemical similarity
- Each chemical with more than 70 percent similarity was linked in a map
- Many chemicals that are structurally similar have different potential as sensitizers
- Many chemicals in the dataset have no similar neighbors

Red = Extreme
Orange = Strong
Yellow = Moderate
Green = Weak/Non

# Chemical Similarity and Skin Sensitization



REACH Skin Sensitization

*Skin sensitization: Simple classification by nearest neighbor*

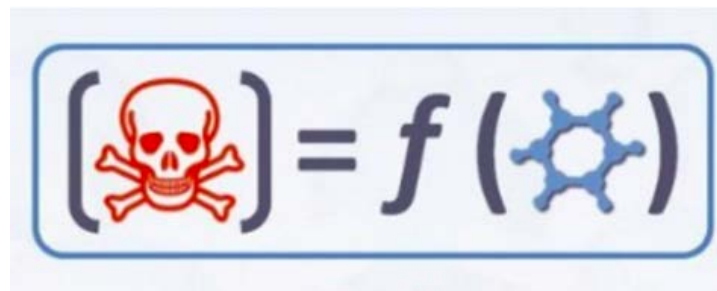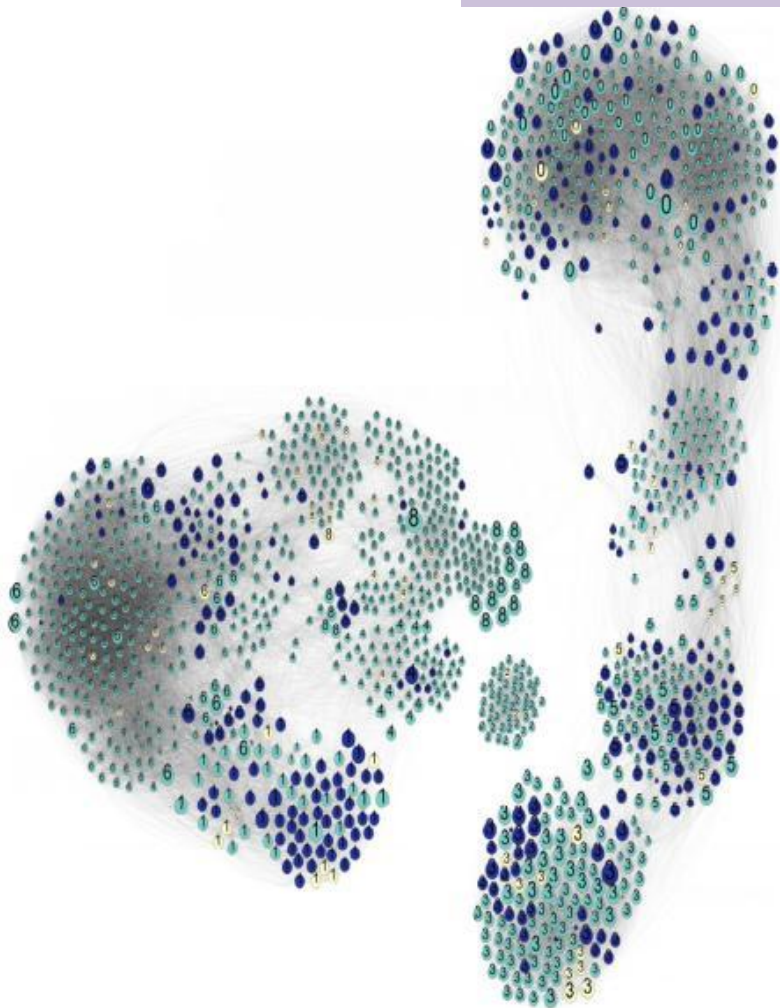*Skin sensitization:*
*Simple classification by nearest neighbor*

| Min. Similarity | Chemicals | Accuracy |
|---|---|---|
| 0.95 | 525 | 0.92 |
| 0.9 | 1189 | 0.85 |
| 0.85 | 1738 | 0.84 |
| 0.75 | 2288 | 0.80 |

**Accuracies better than different animal guideline tests against each other**

# Chemical Similarity and Skin Sensitization

# Chemical Similarity and Skin Sensitization

**Tab. 1: Classification agreement on chemicals with at least two sensitization studies in REACH dossiers from 2008-2014**
Studies found by searching for all studies with *studytype* = Buehler, GPMT, Patch-Test or LLNA.

| | Buehler | GPMT | Patch-test | LLNA |
|---|---|---|---|---|
| **Buehler** | 95.1% (344 chem.) | 91.8% (364 chem.) | 87.8% (58 chem.) | 76.8% (212 chem.) |
| **GPMT** | | 93% (624 chem.) | 90.5% (107 chem.) | 77.4% (403 chem.) |
| **Patch-test** | | | 92.1% (24 chem.) | 78.3% (40 chem.) |
| **LLNA** | | | | 88.5% (296 chem.) |

**Guinea pig**

**77%**

**Mouse**

**89%**

| SKIN SENSITIZATION MODELS | | | |
|---|---|---|---|
| **Model** | **Availability** | **Methodology** | **Results** |
| **PredSkin** | Free | QSAR | Predicts whether chemical will be a skin sensitizer<br>(Binary results for human prediction) |
| **Toxtree** | Free | Structural Alerts | Identifies structural alerts within target chemical |
| **OECD QSAR Toolbox** | Free | Read-Across/QSAR | Predicts whether chemical will be a skin sensitizer<br>(Binary results) |
| **Danish QSAR Database** | Free | Battery algorithm based on three individual QSAR models[#] | Predicts whether chemical will be a skin sensitizer<br>(Binary results for human prediction) |
| **CAESAR**<br>**(VEGA)** | Free | QSAR | Predicts whether chemical will be a skin sensitizer<br>(Binary results) |
| **UL's REACHAcross™** | Commercial | QSAR/<br>Read-Across | Predicts whether chemical will be a skin sensitizer<br>(Binary results) |
| **TIMES-SS** | Commercial | QSAR/<br>Skin Metabolism Simulation | Predicts whether chemical will be a skin sensitizer<br>(Results provided in potency scale: strong sensitizer, weak sensitizer, non-sensitizer) |

# Toxtree: Structural Alerts

- Advantages:
  - Any chemical can be analyzed for the key features
  - Provides a possible mechanism
- Disadvantages:
  - "Skin sensitization alert" or "Skin sensitization reactivity domain"?
- High amount of false positives
  - -approx. 30 percent of chemicals are Michaels acceptors

# PredSkin

# PredSkin

- Advantages:
  - Simple to use
  - Based on human data, but also provides prediction for LLNA as well as DPRA
  - Includes probability estimate
  - Provides batch mode
- Disadvantages:
  - Highly sensitive, not as specific

# OECD QSAR Toolbox

- QSAR Toolbox is very powerful, but there is a learning curve - Windows only
- Includes an automated workflow for data gap filling for skin sensitization
- OECD QSAR Toolbox has a large data set of analogs with LLNA and GPMT data



Workflow components

1. Input
2. Profiling
3. Category building
4. Data gathering
5. Sub-categorization
6. Data gap filling
7. Reporting

The algorithm

User defined target

Start of AW → Profiling

Category building
Data gathering

Sub-categorizaion

Data gap filling

End of AW

Reporting

# OECD QSAR Toolbox

- **Advantages:**
  - Highly accurate if you have similar analogs in the data set and an obvious chemical category
  - Includes potential metabolites
  - Batch mode
- **Disadvantages:**
  - Does not give a probability estimate
  - Minimal coverage

# Model Accuracy



| Model | Basketter et al. (2014) Balanced Accuracy (%) | HSDB Balanced Accuracy (%) |
|-------|-----------------------------------------------|----------------------------|
| PredSkin | 55.3409 | 51.8519 |
| Toxtree | 73.4951 | 60.0514 |
| QSAR Toolbox | 68.4169 | 66.2043 |
| UL ReachAcross | 83 | 78 |
| Danish QSAR Database | 66.4474 | 78.2682 |
| CESAR | 77.2527 | 56.3835 |
| TIMES-SS | 87.6543 | 72.4784 |
| DEREK | 86 | 70 |

# Which model should I use for a hazard assessment?

- You need to know about any *potential* alerts for skin sensitization from ToxTree but structural alerts are low information

- PredSkin can alert you to potential instances where human and LLNA data diverge

- If you there are adequate analogs, the OECD QSAR Toolbox will make a highly accurate prediction

- Accuracy is *highly dependent* on how similar your target chemical is to the data set

# Future Directions: Big Data

2015

2016

2019



From: Progress in Using Big Data in Chemical Toxicity Research at the National Center for Computational Toxicology
Anthony Williams/EPA

# Provider: T.E.S.T.

⬇ Download Summary ▾

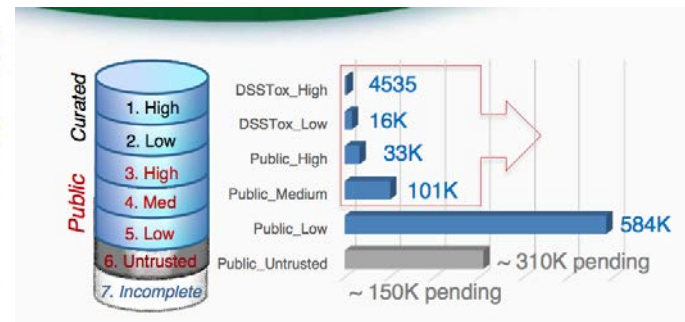| Property | Experimental Value | Consensus | Hierarchical clustering | Single model | Group contribution | Nearest neighbor |
|---|---|---|---|---|---|---|
| 96 hour fathead minnow LC50 | | 2.278 -Log10(mol/L) 590.987 mg/L | 2.455 -Log10(mol/L) 393.459 mg/L | 2.101 -Log10(mol/L) 887.682 mg/L | | |
| 48 hour D. magna LC50 | | | | | | |
| 48 hour T. pyriformis IGC50 | | | | | | 1.421 -Log10(mol/L) 4256.145 mg/L |
| Oral rat LD50 | | 2.410 -Log10(mol/kg) 435.705 mg/kg | 1.966 -Log10(mol/kg) 1213.489 mg/kg | | | 2.855 -Log10(mol/kg) 156.441 mg/kg |
| Bioaccumulation factor | | | | | | 0.639 Log10 4.351 |
| Developmental toxicity | | true | true | true | | true |
| Ames mutagenicity | | true | true | | | true |
| Estrogen Receptor RBA | | | | | | -3.440 Log10 3.631*10^-4 |
| Estrogen Receptor Binding | | false | false | false | false | false |
| Normal boiling point | | 185.1 °C | 232.6 °C | | 177.0 °C | 145.7 °C |
| Melting point | | 3.5 °C | -8.6 °C | | 49.6 °C | -30.7 °C |
| Flash point | | | | | | |
| Vapor pressure | | | | | | -0.579 Log10(mmHg) 0.264 mmHg |
| Density | | 1.425 g/cm³ | 1.506 g/cm³ | | 1.344 g/cm³ | 1.425 g/cm³ |

Type to enter a caption.

# Provider: T.E.S.T.

⬇ Download Summary ▾

| Property | Experimental Value | Consensus | Hierarchical clustering | Single model | Group contribution | Nearest neighbor |
|---|---|---|---|---|---|---|
| 96 hour fathead minnow LC50 | | 2.278 -Log10(mol/L) 590.987 mg/L | 2.455 -Log10(mol/L) 393.459 mg/L | 2.101 -Log10(mol/L) 887.682 mg/L | | |
| 48 hour D. magna LC50 | | | | | | |
| 48 hour T. pyriformis IGC50 | | | | | | 1.421 -Log10(mol/L) 4256.145 mg/L |
| Oral rat LD50 | | 2.410 -Log10(mol/kg) 435.705 mg/kg | 1.966 -Log10(mol/kg) 1213.489 mg/kg | | | 2.855 -Log10(mol/kg) 156.441 mg/kg |
| Bioaccumulation factor | | | | | | 0.639 Log10 4.351 |
| Developmental toxicity | | true | true | true | | true |
| Ames mutagenicity | | true | true | | | true |
| Estrogen Receptor RBA | | | | | | -3.440 Log10 3.631*10^-4 |
| Estrogen Receptor Binding | | false | fa | | | false |
| Normal boiling point | | 185.1 °C | 23 | | | 145.7 °C |
| Melting point | | 3.5 °C | -8.6 °C | | 49.6 °C | -30.7 °C |
| Flash point | | | | | | |
| Vapor pressure | | | | | | -0.579 Log10(mmHg) 0.264 mmHg |
| Density | | 1.425 g/cm³ | 1.506 g/cm³ | | 1.344 g/cm³ | 1.425 g/cm³ |

Based on specific data or warnings?

Type to enter a caption.

# Provider: T.E.S.T.

⬇ Download Summary ▾

| Property | Experimental Value | Consensus | Hierarchical clustering | Single model | Group contribution | Nearest neighbor |
|---|---|---|---|---|---|---|
| 96 hour fathead minnow LC50 | | 2.278 -Log10(mol/L)<br>590.987 mg/L | 2.455 -Log10(mol/L)<br>393.459 mg/L | 2.101 -Log10(mol/L)<br>887.682 mg/L | | |
| 48 hour D. magna LC50 | | | | | | |
| 48 hour T. pyriformis IGC50 | | | | | | 1.421 -Log10(mol/L)<br>4256.145 mg/L |
| Oral rat LD50 | | 2.410 -Log10(mol/kg)<br>435.705 mg/kg | 1.966 -Log10(mol/kg)<br>1213.489 mg/kg | | | 2.855 -Log10(mol/kg)<br>156.441 mg/kg |
| Bioaccumulation factor | | | | | | 0.639 Log10<br>4.351 |
| Developmental toxicity | | true | true | true | | true |
| Ames mutagenicity | | true | true | | | true |
| Estrogen Receptor RBA | | | | | | -3.440 Log10<br>$3.631 \times 10^{-4}$ |
| Estrogen Receptor Binding | | false | false | false | false | false |
| Normal boiling point | | 185.1 °C | 232.6 °C | | | 145.7 °C |
| Melting point | | 3.5 °C | -8.6 °C | | | -30.7 °C |
| Flash point | | | | | | |
| Vapor pressure | | | | | | -0.579 Log10(mmHg)<br>0.264 mmHg |
| Density | | 1.425 g/cm³ | 1.506 g/cm³ | | | 1.425 g/cm³ |



Type to enter a caption.

- ToxCast has data on over 1,800 chemicals- including industrial and consumer products, food additives, and potentially "green" chemicals that could be safer alternatives to existing chemicals.

- ToxCast screens chemicals in over 700 high-throughput assays that cover a range of high-level cell responses and approximately 300 signaling pathways.

- Data can be downloaded or accessed through the iCCS ToxCast Dashboard or Pubchem

- Can potentially evaluate mixtures (e.g. environmental contaminants of waste water)
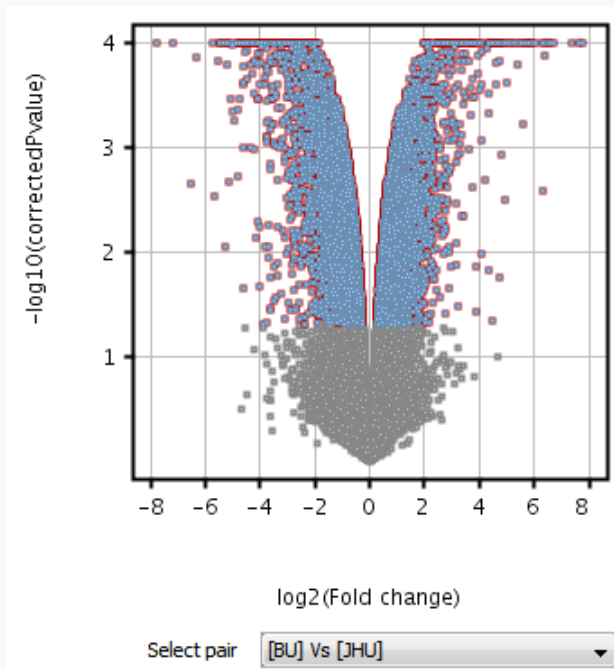
- Used very cautiously in hazard assessments

# HTS Strengths

- Human Assays
  - Eliminates need for animal to human extrapolation
  - Eliminates need for animal testing

- Bulk chemical testing and analysis

# HTS Weaknesses

- Prominent amount of false positives
  - Statistical Validations
  - Certainty of assays
  - Physical Properties - this is extremely important for nanomaterials

- Lacks ADME (Absorption, Distribution, <span style="color:red">Metabolism</span>, and Excretion)
  - Microsomal S9 Inserts

- Translating HTS to Adverse Outcome Pathway (AOP) problematic with regard to resources needed

- Cytotoxicity – false negatives/false positives
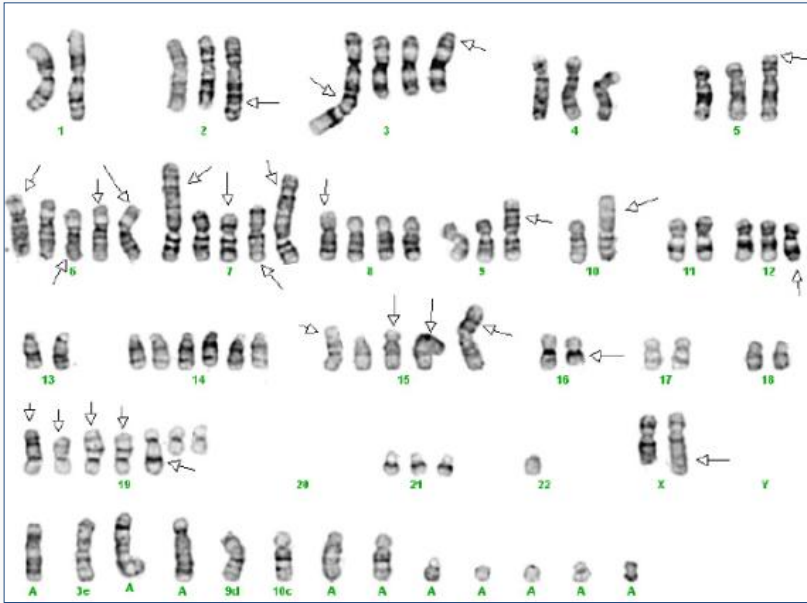
# Comparison of MCF-7 in two laboratories



**Same batch from ATCC**

**Method transfer**

**Transcriptomics**

**negative controls, 4h, gene level, *n* = 3 / group**

**Karyotyping**

## *Extent of deviations from normal genome*

| Classification | Kilobases | Percentage of genome |
|---|---|---|
| Losses | 4587603 | 51.2% |
| Deletions | 667374 | 7.5% |
| Amplifications | 26904 | 0.3% |
| Gains | 2587093 | 28.9% |
| Normal | 871166 | 9.7% |
| Centromeres | 217339 | 2.4% |
| Total Abberations | 7868974 | 87.8% |
| All Entries | 8957479 | |

**SurePrint G3 ISCA CGH+SNP Microarray Kit, 4x180K**
**115234 CGH features.2440 CGH replicate probes, 59647 SNP features**
**reference mapping: caucasian female human reference DNA**

**Kleensang et al., Nature Sci Rep, 2016**

# Future Directions

- Machine learning can already predict human skin sensitization with an accuracy similar to testing approaches just based on structural features; further improvement will require incorporating *in vitro* or much larger data sets
- It will be difficult to reproduce this success with other endpoints, although progress is being made
- The "Limiting Reagent" is human data - we can't test new chemicals on humans for ethical reasons
  - . . . except we do

# Future Directions

- We need better surveillance and reporting of human exposure before we can really model and predict human phenotypes

-  Existing models may be poorly suited for newer chemicals that are likely to be economically important
  - Biobased chemicals
  - Nanomaterials

- "Big data" based on models are already work well for screening; they will soon be key component of any hazard assessment
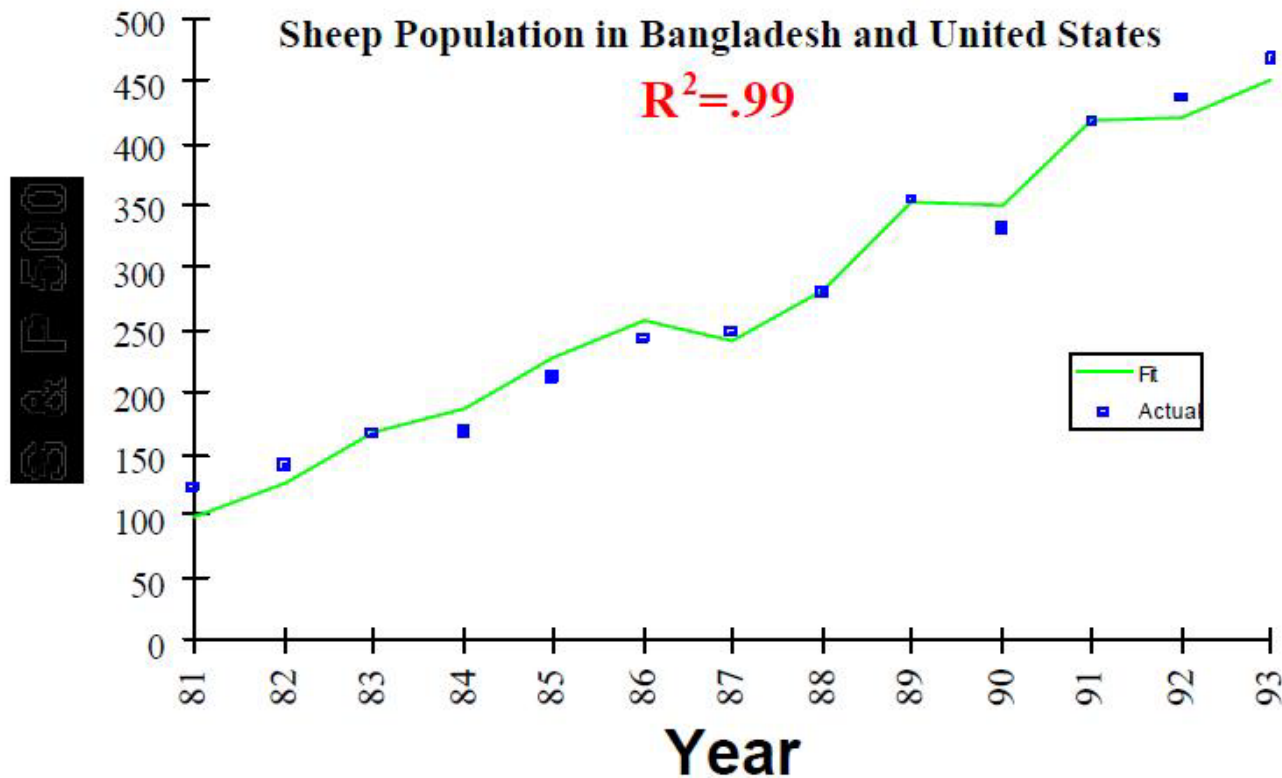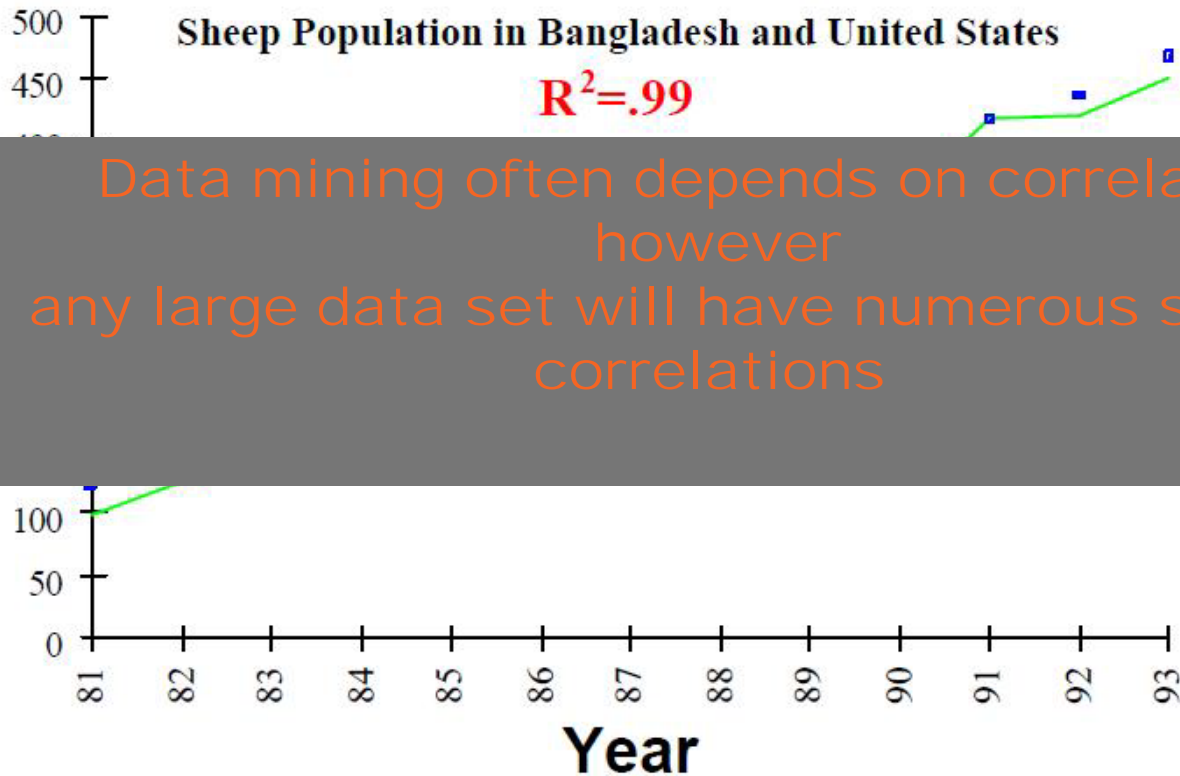
# More Data Is Not Always Better

# More Data Is Not Always Better



## Overfitting the S & P 500

**Butter Production in Bangladesh and United States**

**United States Cheese Production**

**Sheep Population in Bangladesh and United States**

$R^2 = .99$

Data mining often depends on correlations however
any large data set will have numerous spurious correlations

From: Leinweber, David J. Nerds on Wall Street: Math, machines and wired markets. John Wiley and Sons, 2009.

**Why You Should THROW OUT Your Microwave TODAY!**

1. THEY PROVIDE UNNECESSARY DAILY EXPOSURE TO RADIATION
2. THEY CREATE CARCINOGENIC COMPOUNDS IN CERTAIN FOODS
3. THEY DESTROY THE NUTRIENT VALUE OF YOUR FOOD
4. MICROWAVES ARE UNSAFE FOR BABY'S MILK
5. HAZARDOUS CHANGES HAVE BEEN FOUND IN THE BLOOD OF INDIVIDUALS CONSUMING MICROWAVED FOODS

LIVE LOVE FRUIT

# There are lies. . .

# There are damn lies. . .

# And then there are. . .

Lancet. 1998 Feb 28;351(9103):637-41.

### Ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children.

Wakefield AJ[1], Murch SH, Anthony A, Linnell J, Casson DM, Malik M, Berelowitz M, Dhillon AP, Thomson MA, Harvey P, Valentine A, Davies SE, Walker-Smith JA.

⊕ Author information

**Retraction in**
Retraction--Ileal-lymphoid-nodular hyperplasia, non-specific colitis, and pervasive developmental disorder in children. [Lancet. 2010]

**Partial retraction in**
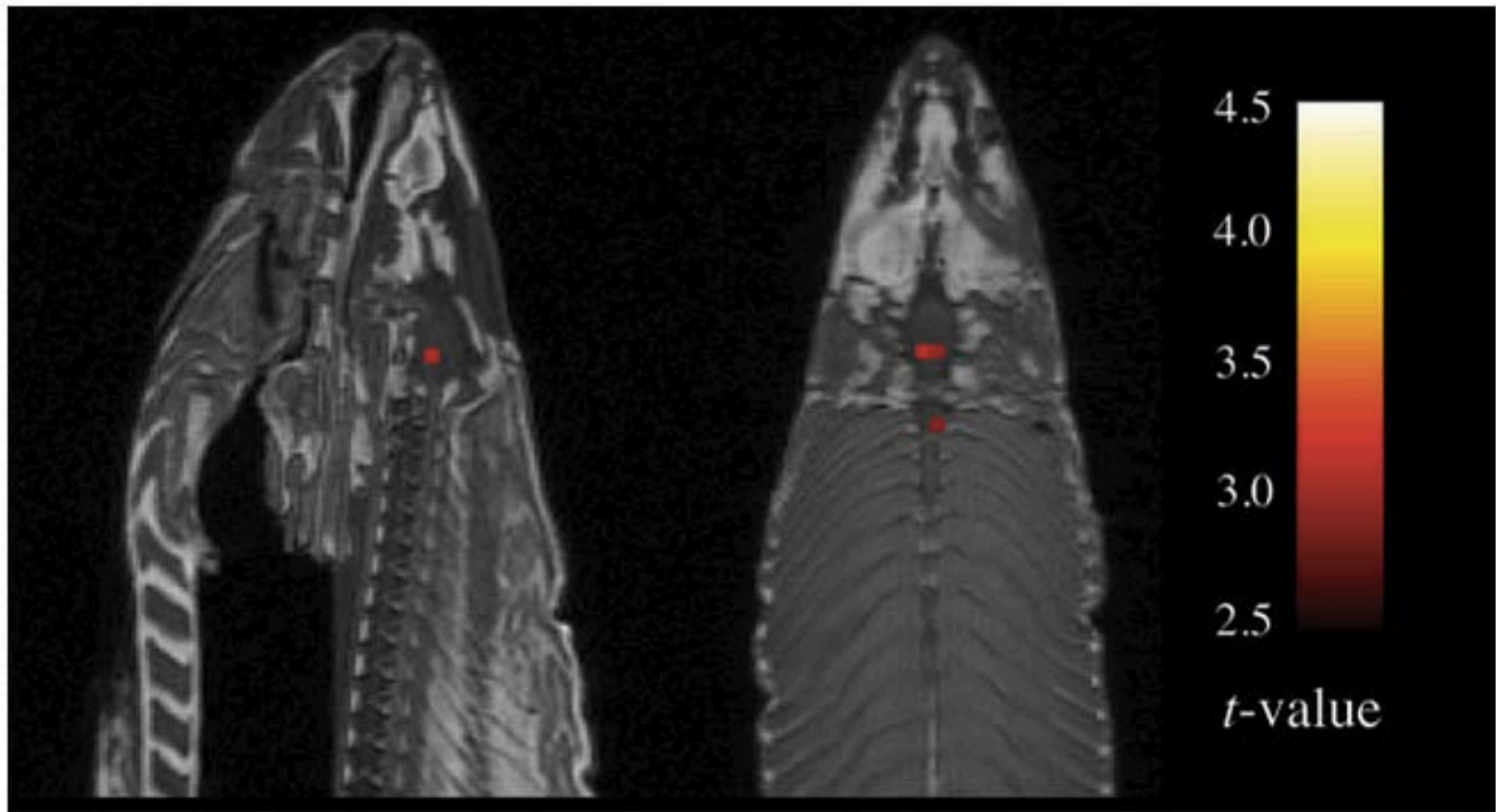Retraction of an interpretation. [Lancet. 2004]

**Abstract**
**BACKGROUND:** We investigated a consecutive series of children with chronic enterocolitis and regressive developmental disorder.

**METHODS:** 12 children (mean age 6 years [range 3-10], 11 boys) were referred to a paediatric gastroenterology unit with a history of normal development followed by loss of acquired skills, including language, together with diarrhoea and abdominal pain. Children underwent gastroenterological, neurological, and developmental assessment and review of developmental records. Ileocolonoscopy and biopsy sampling, magnetic-resonance imaging (MRI), electroencephalography (EEG), and lumbar puncture were done under sedation. Barium follow-through radiography was done where possible. Biochemical, haematological, and immunological profiles were examined.

**FINDINGS:** Onset of behavioural symptoms was associated, by the parents, with measles, mumps, and rubella vaccination in eight of the 12 children, with measles infection in one child, and otitis media in another. All 12 children had intestinal abnormalities, ranging from lymphoid nodular hyperplasia to aphthoid ulceration. Histology showed patchy chronic inflammation in the colon in 11 children and reactive ileal lymphoid hyperplasia in seven, but

# High dimensional data - lies with statistical significance and convinc

# Computers perform better than humans, but make *different* mistakes



[x] Human
[  ]  Not a human



[x] Human
[  ] Not a human

Thanks!  Questions?

amaertens@cermonline.com

Green Toxicology @ CAAT