# The needs for chemistry standards, database tools and data curation at the chemical-biology interface

*Antony Williams*, Kamel Mansouri, Ann Richard and Chris Grulke*

This work was reviewed by EPA and approved for presentation but does not necessarily reflect official Agency policy.

*January 25th, 2016*
*SLAS 2016, San Diego, CA*

EPA
United States
Environmental Protection
Agency

- Ensuring Chemical Structure, Biological Data and Computational Model Quality - Sean Ekins

- Bioassay Variability and Reliability in the Published and Patent Literature - John Overington

- Machine Learning to Optimize Experiments Why Have One Model When You Could Have Thousands? - Alex Clark

- And the rest of you experience it in so many ways…

# Why does Quality Matter?

# What is the Correct Structure of Vitamin K?

# Experiences over the years

## Dedicating Christmas Time to the Cause of Curating Wikipedia

My overall conclusions so far…my estimate is that about 2–3% of the structure records online have errors. What's an error?

1) The structure does not match what it "should be" based on review of many other sources.

2) Systematic nomenclature can be poor…if the name displayed on Wikipedia is converted to a structure then sometimes it is inconsistent with the actual structure displayed

3) Sometimes the formula or mass displayed in the ChemBox are inconsistent with the **actual** mass or formula of the structure displayed

4) The SMILES or InChI String associated with the structure can produce a different structure when converted.

5) The registry number matches either a different structure or a different "form" of the structure. For example, the structure shown is a neutral form of the compound but the registry number is for the salt.

# Experiences over the years

## Aspirin

### Systematic (IUPAC) name
2-(acetoxy)benzoic acid

### Clinical data

**Pronunciation** acetylsalicylic acid
/əˈsiːtəl sælɪˈsɪlɪk/

**AHFS/Drugs.com** monograph

**MedlinePlus** a682878

**Pregnancy category** AU: C
US: C (Risk not ruled out) D
in the 3rd trimester

### Pharmacokinetic data

| | |
|---|---|
| **Bioavailability** | 80–100%[1] |
| **Protein binding** | 80–90%[2] |
| **Metabolism** | Hepatic, (CYP2C19 and possibly CYP3A), some is also hydrolysed to salicylate in the gut wall.[2] |
| **Biological half-life** | Dose-dependent; 2–3 hours for low doses, 15–30 hours for large doses.[2] |
| **Excretion** | Urine (80–100%), sweat, saliva, feces[1] |

### Identifiers

| | |
|---|---|
| **CAS Number** | 50-78-2 ✓ |
| **ATC code** | A01AD05 B01AC06, N02BA01 |
| **PubChem** | CID: 2244 |
| **IUPHAR/BPS** | 4139 |
| **DrugBank** | DB00945 ✓ |
| **ChemSpider** | 2157 ✓ |
| **UNII** | R16CO5Y76E ✓ |
| **KEGG** | D00109 ✓ |
| **ChEBI** | CHEBI:15365 ✓ |
| **ChEMBL** | CHEMBL25 ✓ |
| **Synonyms** | 2-acetoxybenzoic acid acetylsalicylate acetylsalicylic acid O-acetylsalicylic acid |
| **PDB ligand ID** | AIN (PDBe, RCSB PDB) |

### Chemical data

| | |
|---|---|
| **Formula** | $C_9H_8O_4$ |
| **Molecular mass** | 180.157 g/mol |
| **SMILES** [hide] | O=C(Oc1ccccc1C(=O)O)C |
| **InChI** [hide] | InChI=1S/C9H8O4/c1-6(10)13-8-5-3-2-4-7(8)9(11)12/h2-5H,1H3,(H,11,12) ✓
Key:BSYNRYMUTXBXSQ-UHFFFAOYSA-N ✓ |

### Physical data

| | |
|---|---|
| **Density** | 1.40 g/cm$^3$ |
| **Melting point** | 135 °C (275 °F) |
| **Boiling point** | 140 °C (284 °F) (decomposes) |
| **Solubility in water** | 3 mg/mL (20 °C) |

**Sustainable chemistry relates to the design and use of chemicals that minimize impacts to human health, ecosystems and the environment.**

**Sustainable Chemistry**

**Human Health Hazard**

**Environmental Persistence**

**Ecosystems Hazard**

➢ How can predictive tools be used to screen for potential impacts of chemicals early in the product development process?

➢ How can chemical and bioassay data inputs be combined to screen and prioritize testing of thousands of existing chemicals lacking data?

➢ How do chemical transformations impact the hazard potential and environmental persistence?

# NCCT Publicly Available Data
http://www.epa.gov/ncct



**ACToR**

**ToxRefDB**

**DSSTox**

*Chemical Structure Annotation*
*Chemical Management*
*Chemical QC*

**ToxCast**

**Tox21**

# Chemistry as a Data Foundation



*MOA & Knowledge-informed features & chemotypes*

*Chemical search by CAS-names SMILES Structures*

*Structure file downloads*

*PK/ADME model inputs & outputs*

CAS look-up
New list error-checking

**Structure "cleaning"**
**SAR-ready structure files**
**Structure-based predictions**
*Similarity searching*
*Fingerprints, feature sets*

*Phys-chem calc & measured properties*

Chemical Structures

CASRN & Chemical Names

Chemical Inventories & Data Sources (ACTOR)

**1:1 mapping**
**CAS:Name:structure**

*Quality scores*

>700K substances

*EPA SRS (TSCA, EcoTox)*

*ToxCast/Tox21*

*ToxRef*

*EDSP21*

*CPCAT*

**QSAR models; phys-chem properties**

CSS Dashboards

ToxCast/Tox21 HTS activities

# Diverse quality in databases

- ## Our challenges in assembling a database
    - – Sourcing high quality data – sorting wheat from chaff
    - – How to mesh data together – based on structure? On name? On CAS number? Other identifier?
    - – Checking self-consistency of data?
    - – Structure validation versus property value validation – very different challenges

# iCSS Chemistry Dashboard

- Deliver a web-based platform hosting prediction models – both in-house and "community models"

- Initial set of models to be based on EPISuite's PHYSPROP datasets – logP, BP, MP, Wsol etc.

- A question of "data quality" - a project to review data initiated.

## 6.2. KOWWIN Estimation Accuracy & Domain

### 6.2.1. Training Accuracy

The following graph illustrates the correlation between the experimental and KOWWIN estimated log Kow values:



KOWWIN Training Set (2447 compounds)

Total Training Set Statistics:
```
    number in dataset      = 2447
    correlation coef (r²)  = 0.982
    standard deviation     = 0.217
    absolute deviation     = 0.159
    avg Molecular Weight   = 199.98
```

# Overall Predicted vs. Experimental

(>15k chemicals)

- Sourced the SDF files online as downloads

- Basic Manual review searching for errors:
  - hypervalency, charge imbalance, undefined stereo
  - deduplication
  - mismatches between identifiers

    - CAS Numbers not matching structure
    - Names not matching structure
    - Collisions between identifiers

# Incorrect valences

# Differences in Values (MP)

| Structure | Formula | FW | CAS | NAME | MP | EstMP | ErrorMP |
|---|---|---|---|---|---|---|---|
|  | $C_{11}H_{14}O_2$ | 178.2277 | 000093-15-2 | METHYLEUGENOL | -4.000000000000000e+000 | 3.285000000000000e+001 | 3.685000000000000e+001 |
|  | $C_{11}H_{14}O_2$ | 178.2277 | 006380-24-1 | 4-Allyl-1,2-dimethoxy-benzene | 7.000000000000000e+001 | 3.285000000000000e+001 | -3.715000000000000e+001 |

# Salts



PHYSPROP database

Appropriate Depiction

- Same structure depictions (Molfiles) - different CAS, different names, different SMILES

| | | | | | |
|---|---|---|---|---|---|
|  | $C_{24}H_{40}O_4$ | 392.5720 | 000083-49-8 | HYODEOXYCHOLIC ACID | 3.080000000000000e+000 |
|  | $C_{24}H_{40}O_4$ | 392.5720 | 000128-13-2 | URSODEOXYCHOLIC ACID | 3.000000000000000e+000 |

# Identifiers

- Many chemical names are truncated
- Many chemicals don't have CAS Numbers

| | |
|---|---|
| SRC000-02-7 | Ethanaminium, N,N,N-trimethyl-2-[(1-oxo-2-propen |
| SRC000-04-3 | Guanidine, N-hydroxy-N''-[4-(methylthio)benzeneme |
| SRC000-04-4 | Hydrazinecarboximidamide, N'-[4-(methylthio)benz |
| SRC000-04-5 | NNN5-TeMe-N-(3FuranMe),ammon Br |
| SRC000-04-6 | Benzenamine, 4-bromo-N,N-bis(2,2,2-trifluoroethy |
| SRC000-04-7 | 2-Propenoic acid, 3-(2-chlorophenoxy)-, methyl e |
| SRC000-05-1 | 9H-Purine-9-acetaldehyde, a-(1-formyl-2-hydroxye |
| SRC000-05-2 | N1-Pr-N2-CN-N3-Me guanidine |
| SRC000-05-3 | 1-(2-OHEt)-2-Me imidazoline HCL |

# KNIME workflow to evaluate dataset

## 15,809 chemicals in the KOWWIN dataset

- CAS Checksum: 12163 valid, 3646 invalid

- Invalid names: 555

- Invalid SMILES 133

- Valence errors: 322 Molfile, 3782 SMILES

- Duplicates check:
  - 31 DUPLICATES MOLFILE
  - 626 DUPLICATES SMILES
  - 531 DUPLICATES NAMES

- SMILES vs. Molfiles (structure check)
  - 1279 differ in stereochemistry
  - 362 "Covalent Halogens"
  - 191 differ as tautomers
  - 436 are different compounds

# Quality flags inserted into data

- 4 Stars ENHANCED: 4 levels of consistency with stereo information
- 4 Stars: 4 levels of consistency, stereo ignored.
- 3 Stars Plus: 3 out of 4 levels. The 4th is a tautomer.
- 3 Stars ENHANCED: 4 levels of consistency with stereo information
- 3 Stars: 3 levels of consistency, stereo ignored.
- 2 Stars PLUS : 2 out of 4 levels. The 3rd is a tautomer.
- 1 Star - What's left.

# Standardization of structures
## Make QSAR ready file for modeling



- Explicit hydrogen removal
- Removal of chirality info, isotopes and pseudo-atoms
- Aromatization + add explicit hydrogen atoms
- Standardize Nitro groups
- Other tautomerization/mesomerization checking
- Neutralize (when possible)

# Mesomerization/tautomerization

- Azide mesomers
- Exo-enol tautomers
- Enamine-Imine tautomers
- Ynol-ketene tautomers
- ….

# Neutralize Structures

# Building the Models

- Interested in "impact of quality vs. quantity of data on prediction models"
- Is it *worth* the effort to clean and enhance data underpinning the models?
- QSAR ready forms for modeling – standardize tautomers, remove stereochemistry, no salts
- Compare 3 STAR and BETTER models with entire dataset

# What about the "cluster"?



QSAR ready dataset

'''SpMAD_DzZ''' Spectral mean absolute deviation from Barysz matrix / weighted by atomic number
'''AATS1m''' Average Broto-Moreau autocorrelation - lag 1 / weighted by mass
'''ATS4v''' Average Broto-Moreau autocorrelation - lag 4 / weighted by van der Waals volumes
'''ATSC4m''' Centered Broto-Moreau autocorrelation - lag 4 / weighted by mass
'''nHBint5''' Count of E-State descriptors of strength for potential Hydrogen Bonds of path length 5
'''minsOH''' Minimum atom-type E-State: -OH
'''ETA_Beta''' A measure of electronic features of the molecule
'''MPC6''' Molecular path count of order 6
'''MPC9''' Molecular path count of order 9

# Weighted kNN Model, 10 descriptors

**weighted KNN model 10 descriptors**



Weighted 5-nearest neighbors
Training: 11251 chemicals
Test set: 2798 chemicals

5 fold cross validation:
R2: 0.87
RMSE: 0.67

'''CrippenLogP''' Crippen's LogP
'''GATS2c''' Geary autocorrelation - lag 2 / weighted by charges
'''LipoaffinityIndex''' Lipoaffinity index
'''AATS1p''' Average Broto-Moreau autocorrelation - lag 1 / weighted by polarizabilities
'''ATSC1i''' Centered Broto-Moreau autocorrelation - lag 1 / weighted by 1st ionization potential
'''ETA_EtaP''' Composite index Eta relative to molecular size
'''MLFER_S''' Combined dipolarity/polarizability
'''nN''' Number of nitrogen atoms
'''ETA_Beta''' A measure of electronic features of the molecule
salt _ index  salt info
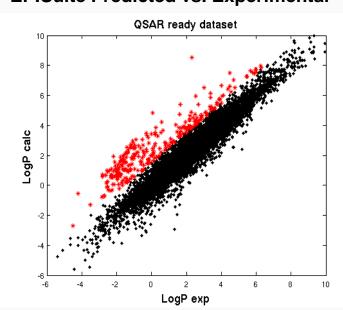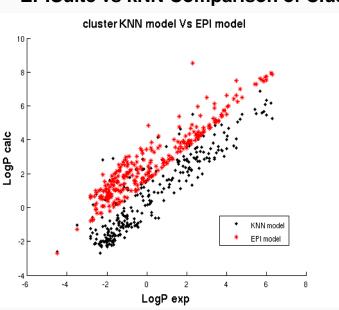
**Applicability Domain of the
Original EPISuite was an issue…**

**EPISuite Predicted vs. Experimental**

**EPISuite vs kNN Comparison of Cluster**

# Bringing together PhysChem Data

- Additional Experimental Data being assembled
    - *PubChem*
    - *eChemPortal*
    - *Open Data Sources*
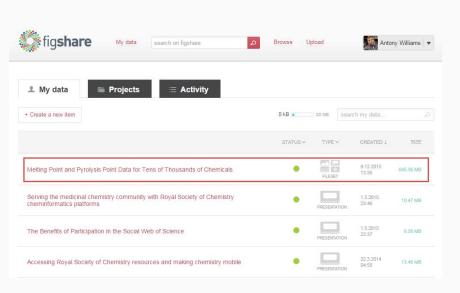
# Open Data Set - example

- 200,000 melting points and 13,000 pyrolysis data points

**Research article**

**Open Access**

**The development of models to predict melting and pyrolysis point data associated with several hundred thousand compounds mined from PATENTS**

**Igor V. Tetko**[1,2]*, **Daniel M. Lowe**[3] and **Antony J. Williams**[4]

- Modeling this much data is a challenge!
- Accessibility to data?
  – Figshare.com
  – DataDryad.com
  – Institutional Repositories
  – Publishers?

figshare   My data   search on figshare   Browse   Upload   Antony Williams ▾

My data   Projects   Activity

+ Create a new item   0 kB   20 GB   search my data...

| | STATUS | TYPE | CREATED ↓ | SIZE |
|---|---|---|---|---|
| Melting Point and Pyrolysis Point Data for Tens of Thousands of Chemicals | ● | FILESET | 9.12.2015 13:35 | 645.36 MB |
| Serving the medicinal chemistry community with Royal Society of Chemistry cheminformatics platforms | ● | PRESENTATION | 1.5.2015 23:46 | 10.47 MB |
| The Benefits of Participation in the Social Web of Science | ● | PRESENTATION | 1.5.2015 23:37 | 6.58 MB |
| Accessing Royal Society of Chemistry resources and making chemistry mobile | ● | PRESENTATION | 22.3.2014 04:55 | 13.46 MB |

31

# iCSS Chemistry Dashboard Releasing in March 2016

# iCSS Chemistry Dashboard
# Releasing in March 2016

# iCSS Chemistry Dashboard Releasing in March 2016

- Links external resources: EPA, NIH, property predictors

# …Integrate Toxicity Estimation Software Tool

http://www.epa.gov/chemical-research/toxicity-estimation-software-tool-test



| Option |
| --- |
| Fathead minnow LC50 (96 hr) |
| Daphnia magna LC50 (48 hr) |
| T. pyriformis IGC50 (48 hr) |
| Oral rat LD50 |
| Bioaccumulation factor |
| Developmental Toxicity |
| Mutagenicity |
| Normal boiling point |
| Vapor pressure at 25°C |
| Melting point |
| Flash point |
| Density |
| Surface tension at 25°C |
| Thermal conductivity at 25°C |
| Viscosity at 25°C |
| Water solubility at 25°C |
| Molecular Descriptors |

# Toxicity Estimation Software Tool



**Toxicity prediction results for 333-41-5 for Hierarchical clustering method**

Prediction results

| Endpoint | Experimental value CAS: 333-41-5 Source: ECOTOX | Predicted value[a] | Prediction interval |
|---|---|---|---|
| Fathead minnow $LC_{50}$ (96 hr) -Log(mol/L) | 4.81 | 5.39 | $4.54 \leq Tox \leq 6.24$ |
| Fathead minnow $LC_{50}$ (96 hr) mg/L | 4.70 | 1.23 | $0.17 \leq Tox \leq 8.71$ |

[a]Note: the test chemical was present in the external test set.

Cluster model predictions and statistics

| Cluster model | Test chemical descriptor values | Prediction interval -Log(mol/L) | $r^2$ | $q^2$ | #chemicals |
|---|---|---|---|---|---|
| 1296 | Descriptors | $6.010 \pm 1.136$ | 0.793 | 0.733 | 101 |
| 1300 | Descriptors | $5.458 \pm 1.312$ | 0.729 | 0.645 | 111 |
| 1301 | Descriptors | $5.136 \pm 1.169$ | 0.747 | 0.718 | 294 |
| 1302 | Descriptors | $4.922 \pm 1.182$ | 0.774 | 0.751 | 641 |

Cluster models with violated constraints

| Cluster Model | $r^2$ | $q^2$ | # chemicals | Message |
|---|---|---|---|---|
| 1121 | 0.810 | 0.576 | 10 | Rmax constraint not met |
| 1209 | 0.799 | 0.574 | 11 | Fragment constraint not met |
| 1247 | 0.919 | 0.647 | 20 | Fragment constraint not met |
| 1264 | 0.869 | 0.781 | 22 | Fragment constraint not met |
| 1268 | 0.675 | 0.553 | 24 | Fragment constraint not met |

Descriptor values for test chemical

## Model # 1296

| Parameter | Value |
|---|---|
| Endpoint | Fathead minnow $LC_{50}$ (96 hr) |
| $r^2$ | 0.793 |
| $q^2$ | 0.733 |
| #chemicals | 101 |
| Model | Model # 1296 |



**Model fit results**

36

# Conclusion

- Improved chemistry support and data will impact every NCCT project

- Creating high quality data is difficult & iterative

- Public domain and open data is valuable *but* requires validation

- Modern cheminformatics approaches can help validate and prepare data for modeling

- **There are increasing efforts to source and release high quality, open data to the world**

# Acknowledgements

## Contributors to the iCSS Chemistry Dashboard

- Jeff Edwards
- Jeremy Fitzpatrick
- Jordan Foster
- Chris Grulke
- Jason Harris

- Dave Lyons
- Kamel Mansouri
- Ann Richard
- Aria Smith
- Jennifer Smith
- Indira Thillainadarajah