# Computational Platform for Compound Identification

*Aurelien Monge\*, Stephane Cano, Stephane Albrecht, Pavel Pospisil, Elyette Martin\**

*Philip Morris International, R&D*

*19th June 2013*

*ACD/Labs Symposium on Laboratory Intelligence (EUM), Neuchatel Switzerland*
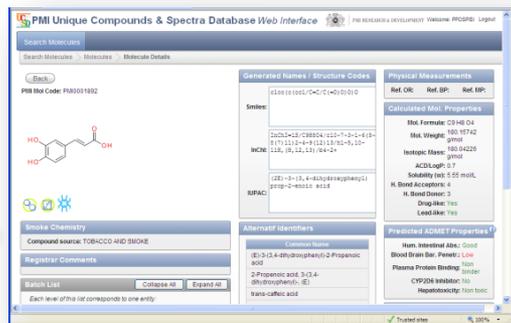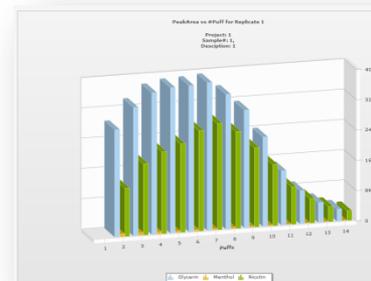
# Introduction

Chemoinformatics challenges at PMI:

- Compound identification from complex matrices

- Efficient and/or automatized compound registration (including registration of mixtures and stereochemical isomers)

- Managing spectral data

- Associating toxicity data to compounds and mixtures

- Inserting chemical data in corporate Scientific Data Warehouse

- Reporting chemical data in a relevant way

- Building R&D chemoinformatics platform using Rapid Development Tools

PMI RESEARCH & DEVELOPMENT

*A. Monge & E. Martin, PMI R&D*

# PMI Chemoinformatics platform



**1. Unique Compounds & Spectra Database**

**2. Computer Assisted Structure Identification**

**3. Chemoinformatics Knowledge Base**

**Scientific Data Warehouse**

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT
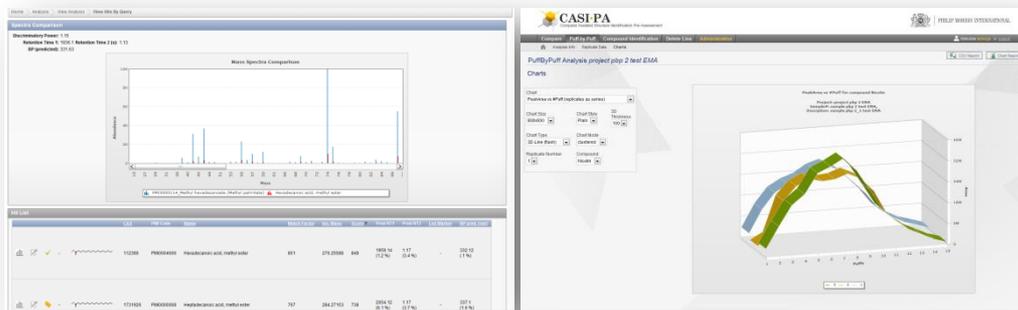
# 1. Unique Compounds & Spectra Database



**1. Unique Compounds & Spectra Database**

**3. Chemoinformatics Knowledge Base**

**Scientific Data Warehouse**

**2. Computer Assisted Structure Identification**

*A. Monge & E. Martin, PMI R&D*
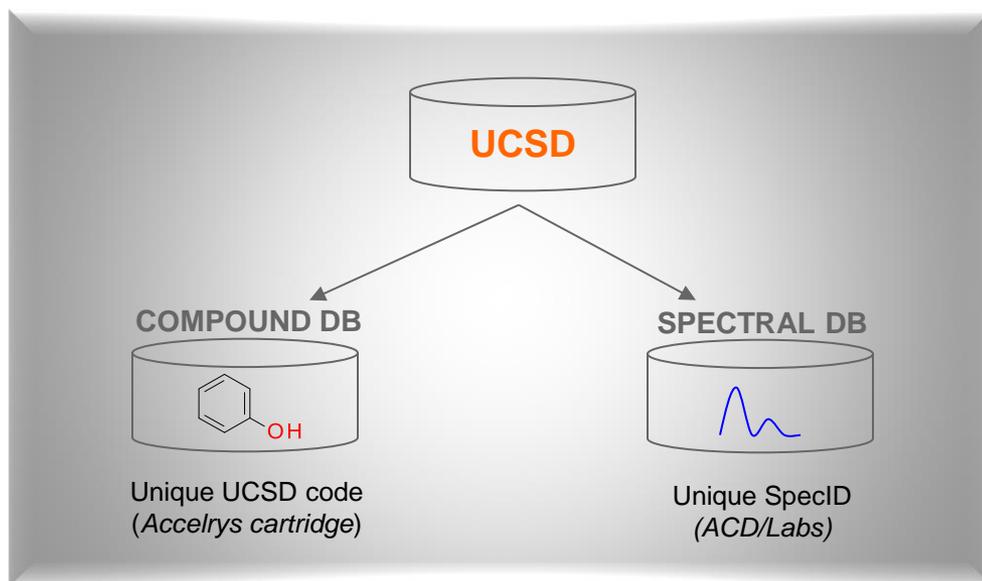
PMI RESEARCH & DEVELOPMENT

# Concept of UCSD



**Unique Compound and Spectra Database:** To assemble R&D chemical information into a central repository of chemical substances and analytical spectra.

*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# Concept of UCSD

- UCSD is an internal database with no external access

- Accelrys chemical cartridge with enhanced stereochemistry

- ACD/Labs for the analytical spectra part

- Web interface developed in Oracle Application Express and automatic processes in Accelrys Pipeline Pilot



*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# UCSD Roles



**1. Submitter**

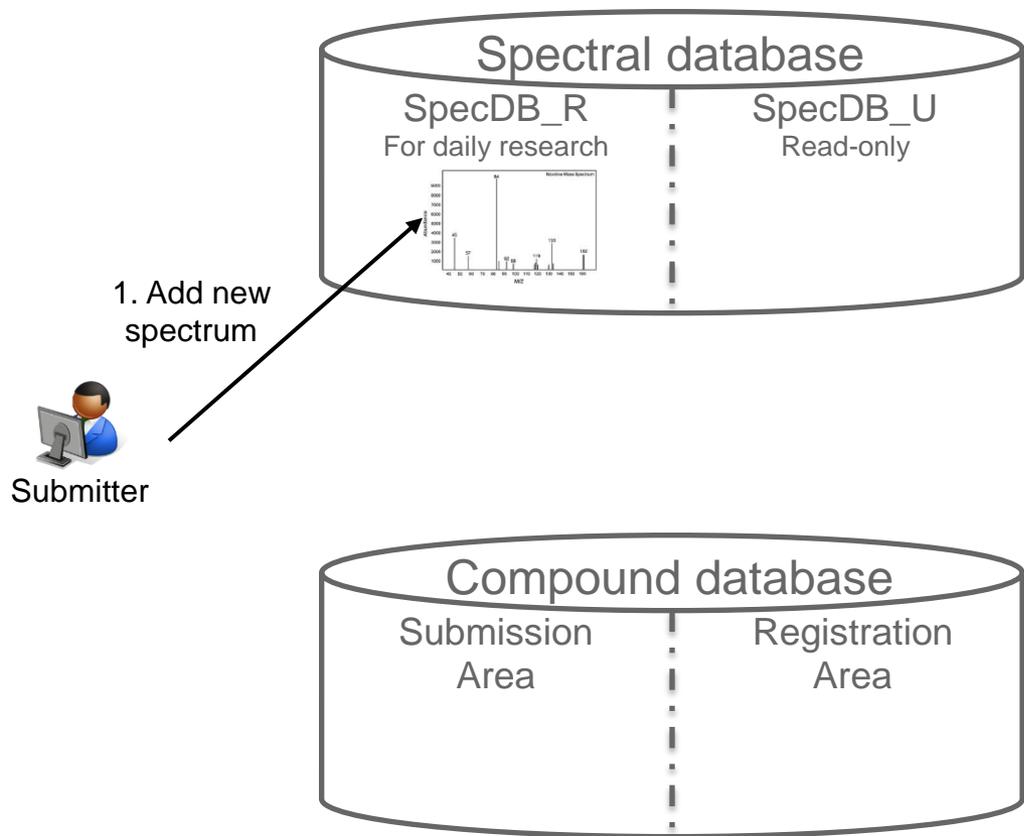A chemist who can **search**, **view** data and also **submit** molecules and spectra.



**2. Registrar**

A chemoinformatician that can **search, view, submit** and also check and **register** data.



**Viewer**

User can search and **view** registered data (print, export, list creation)

*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# UCSD Workflow



Spectral database

SpecDB_R
For daily research

SpecDB_U
Read-only

1. Add new spectrum

Submitter

Compound database

Submission Area

Registration Area

*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# SpecDB_R

## ACD/ChemFolder Enterprise module



*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# UCSD Workflow



**Spectral database**

SpecDB_R
For daily research

SpecDB_U
Read-only

**1. Add new spectrum**

Submitter

**2. Add new batches with SpecID**

**Compound database**

Submission Area

Registration Area

**3. Validate new record**

Note: submitted batches can also correspond to 'known unknowns' that once identified, can be later registered as known molecules.

PMI RESEARCH & DEVELOPMENT

*A. Monge &  E. Martin, PMI R&D*

# Concept of Molecule, Substance and Batch Assignment
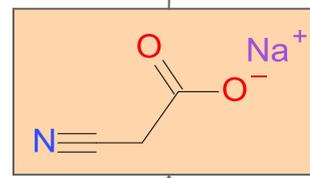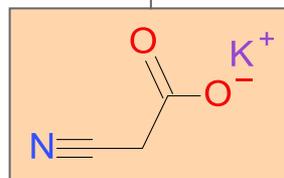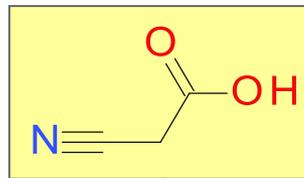
Unique means there is no redundancy, ensured by uniqueness of chemical structure and systematic and accurate registration process.



**Molecule =**
Neutral, unique chemical entity

Every molecule is assigned a unique company code:
e.g., **PMI01234567**

**Substance =**
Molecule + Salt

**Substance ID = e.g.,** PMI01234567-A · PMI01234567-B · PMI01234567-C

**Batch =**
Substance physically present or project-relevant as cited in literature

Batch ID =  e.g.,

| | | |
|---|---|---|
| BC000002152 | BC000000122 | BC000000121 |
| BC000008641 | BC000000154 | BC000000158 |
| BC000000560 | BC000000176 | BC000003174 |
| BC000000320 | BC000002504 | BC000002598 |
| | BC000000894 | |

PMI RESEARCH & DEVELOPMENT

*A. Monge &  E. Martin, PMI R&D*

# UCSD Compound Part



• To each newly registered compound we developed automatic processes that:

- attach SMILES strings, InChI codes and IUPAC names

- calculate structure-derived physico-chemical and ADMET properties

⇒ Naming and calculation of compounds after registration (uniqueness check) reduces further ambiguity

See more in Martin et al., *Journal of Cheminformatics* 2012, 4:11

*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# UCSD Compound Part



See more in Martin et al., *Journal of Cheminformatics* 2012, 4:11

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# In waiting for Spectrus

Search by spectra similarity is not possible in ChemFolder module ➜ implementation of a script that copies each entry in Spectrum Database module



**Add** ➜ **ChemFolder Enterprise** **Copy** ➜ **Spectrum Database**

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# Process for bulk import

## INPUT for SpecDB_R



**CSV file**

| | A | B | C | D | E | F | G | |
|---|---|---|---|---|---|---|---|---|
| Title | | Submitter | Scientist | Project | Confidence | Ionization type | Analyzer type | Pro |
| 2-Nonanone-1,1,1,3,3-d5 | | emartin | Dossin Eric | Chemistry to Exposure | Reference | EI | QTOF | C9D |
| Nitrobenzene-d5 | | emartin | Dossin Eric | Chemistry to Exposure | Reference | EI | QTOF | C6D |
| O-cresol-d8 | | emartin | Dossin Eric | Chemistry to Exposure | Reference | EI | QTOF | C7D |

**JDX file**

```
##TITLE= 2-Nonanone-1,1,1,3,3-d5
##DATA TYPE= MASS SPECTRUM
##MW= 147.1671489294
##MOLFORM=C9D5H13O
##CAS REGISTRY NO=1398065-76-3
##CAS NAME=2-Nonanone-1,1,1,3,3-d5
##NAMES=PMI0008955
##NPOINTS= 87
##XYDATA= (XY..XY)
30.0629, 0.996579903650516 31.0520, 1.03540117877503
##END=
##TITLE= Nitrobenzene-d5
##DATA TYPE= MASS SPECTRUM
##MW= 128.0634121394
##MOLFORM=C6D5NO2
##CAS REGISTRY NO=4165-60-0
```

## OUTPUT

| | A | B | C | D | E | F | G | H | |
|---|---|---|---|---|---|---|---|---|---|
| Title | | SpecID | Submitter | Scientist | Project | Confidence | Ionization type | Analyzer type | Pro |
| 2-Nonanone-1,1,1,3,3-d5 | | SPR000000411 | emartin | Dossin Eric | Chemistry to Exposure | Reference | EI | QTOF | C9D |
| Nitrobenzene-d5 | | SPR000000412 | emartin | Dossin Eric | Chemistry to Exposure | Reference | EI | QTOF | C6D |
| O-cresol-d8 | | SPR000000413 | emartin | Dossin Eric | Chemistry to Exposure | Reference | EI | QTOF | C7D |

**SD file**

**Merge and import in the compound part via Pipeline Pilot protocol**



SD Reader | Get Dicts Values | Get User & Scientist | Get User & Submitter | Get SUBID | Insert Molecule and Salts | Insert Batch | Insert Substance Properties | GPROPERTIES | ODBC Publish All

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# 2. Computer Assisted Structure Identification



**1. Unique Compounds & Spectra Database**

**2. Computer Assisted Structure Identification**

**3. Chemoinformatics Knowledge Base**

**Scientific Data Warehouse**

*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# Computer Assisted Structure Identification

Our in-house developed Computer Assisted Structure Identification platform is constituted of two software:





Software 1 "CASI-PA": assists the processing of MS data for given analyses:

- Compound Identification

- Comparison of 2 samples

- Puff-by-Puff

Automated non-target processing of MS data.

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# Puff-by-Puff Analysis of Cigarette Aerosols

Analysis of the concentration of given chemical compounds per puff.



*A. Monge & E. Martin, PMI R&D*

# Computer Assisted Structure Identification

Our Computer Assisted Structure Identification platform is constituted of two software:





Help for the processing of MS data for given analyses:

- Compound Identification

- Puff-by-Puff

- Comparison of two products

Software 2: Automated non-target processing of MS data.

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# Computer Assisted Structure Identification

Automated platform to accelerate and standardize identification of chemical structures of aerosol constituents with high confidence power.

*Smoke of a conventional cigarette, measured by GCxGC-EI-TOF-MS*



Compound?

A. Knorr et al., Anal. Chem. 2013, submitted for publication

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# CASI Concept



*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# Predictive QSPR Models for KI and 2DrelRT

### Kovats Index



### 2DrelRT



*Validation n=60*          r$^2$ = 0.981

r$^2$ = 0.855

GA – linear regression, 7 Molecular Descriptors

GA – Support Vector Regression, 12 molecular descriptors

*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# Ranking of hits using CASI Score vs. NIST MS Match Factor



CASI Score improves NIST MS Search ranking.

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# CASI 2 Software

# 3. Chemoinformatics Knowledge Base



**1. Unique Compounds & Spectra Database**

**3. Chemoinformatics Knowledge Base**

**Scientific Data Warehouse**

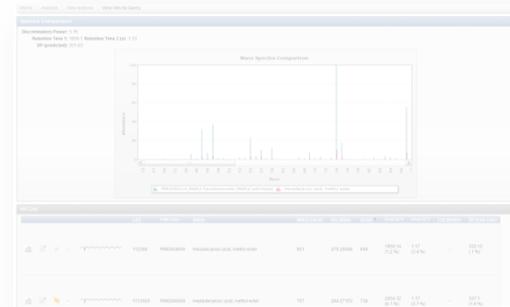**2. Computer Assisted Structure Identification**

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# CIKB Scope

Central platform to **report** on the presence and concentration of constituents in smoke aerosols and their known or measured toxicities to support the toxicological assessment.

Information of interest to be combined with chemical – analytical information of UCSD and CASI platforms:

**Product Lifecycle:** Complete specification of developed and tested products.

**LIMS:** Quantitative smoke characterization.

***In-vitro* Assays:** Standardized *in vitro* assays to assess toxicity.

**Toxicities:** Known toxicological information from literature.

**Metabolites:** Known metabolites from literature.

A. Monge & E. Martin, PMI R&D

PMI RESEARCH & DEVELOPMENT

# Chemistry Based Reporting



**1. Unique Compounds & Spectra Database**

**3. Chemoinformatics Knowledge Base**

**2. Computer Assisted Structure Identification**

**Scientific Data Warehouse**

**Product Lifecycle**   **LIMS**   **Assays**   **Toxicities**   **Metabolites**

*A. Monge & E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT

# Used tools

- Rapid Application Development tools were used to develop software to support the processing of analytical chemistry data (mass spectrometry) and thus to answer quickly to business needs.

| Web Application | Methodology |
|---|---|
| **Web Application**<br>User friendly and no desktop installation required | **Methodology** |

**Web Application**
User friendly and no desktop installation required

**Oracle® Apex**
Used for database centric developments (e.g. chemical registration system).

**Grails**
Used to develop custom data processing application for analytical chemists.

**Accelrys® Pipeline Pilot**
Used for complex manipulation of chemical structures and data.

**Accelrys ® Direct cardridge**
Used to manage chemical structures in compound database.

**ACD/Labs modules**
Used to manage spectra in spectral database.

**Oracle® 11gR2**
Database Management System

**Methodology**

✓ Iterative development

✓ Users involved during all the process

✓ 21 CFR Part 11 requirements

✓ R&D Computerized Systems Development best practices

PMI RESEARCH & DEVELOPMENT

*A. Monge & E. Martin, PMI R&D*

# Conclusion



**UCSD** developed to facilitate and automatize processes of the unambiguous registration of compounds and their analytical spectra
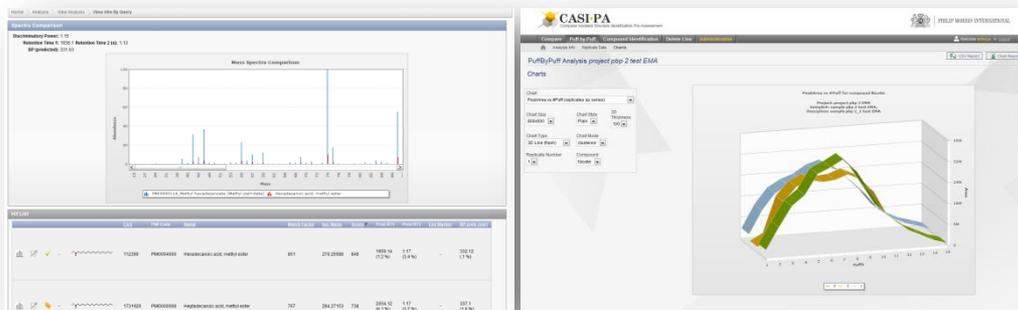
**CIKB** in development to report on product chemical constituents and the associated toxicity and metabolic data

**Scientific Data Warehouse**

**CASI** platforms developed to increase the confidence in the identification of compounds using GC-MS methods and facilitates to users the comparison of analyses

PMI RESEARCH & DEVELOPMENT

*A. Monge & E. Martin, PMI R&D*

# Acknowledgment

- *Stephane Cano, Stephane Albrecht, Pavel Pospisil*

- *PMI R&D chemistry analytics teams*

- *PMI R&D HPC team*

- *PMI R&D SIS team*

- *Accelrys Symyx consultants*

- *ACD/Labs consultants*

*Thank you for your attention.*

*A. Monge &  E. Martin, PMI R&D*

PMI RESEARCH & DEVELOPMENT