# Integration of in-house and external data in practice

**D. Domine, C. Merlot, M. Ibberson, and M. de Francesco**
**Drug Discovery Informatics**
**Merck Serono Geneva Research Center**

# Outline

◆ **Introduction: Drivers and Needs in KM**

◆ **Merck Serono Knowledge Space**
  - ◆ **In-house data integration**
  - ◆ **External data integration**
  - ◆ **Workflows**
  - ◆ **Project management**
  - ◆ **Entity aggregation**

◆ **Examples**
◆ **Concluding remarks**

# Drivers

- **Informed decision-making at different levels: Fast, easy, well-informed.**

- **Increased complexity, variety, quantity of information/technologies**

- **Broader scope (e.g Discovery, Clinics, Manufacturing, etc)**

- **More data outside than inside. Need better integration with in-house information**

- **All project team members and management need to understand the global status and goal**

- **Navigation between different levels of abstraction.**

- **Know-how capitalization**

- **Optimize resources**

- **Enable mergers**

# Needs

- **In-house experimental data integration. Never ending**
  - Constantly evolving
  - Broader scope (Chem., bio., pharm., clin., gen., manuf., IP, Compet., ….)
  - Capture interpretations vs experimental data, results

- **External data/Prediction tools integration. Starting**
  - Need to be structured
  - Requires manual curation, understanding
  - Lack of standards

- **Know-how, unstructured data acquisition, mining, interpretation (in-house and external). Starting**

- **Data and Information Mapping/Scoring. Starting**

- **Different levels of understanding/reporting (lab., project, Management). Starting**

- **Automated annotation/interpretation to turn data into information and knowledge Possible at all?**

# Some existing technical solutions

- **Integration: Function of needs and company history**
    - Requisite for automated interpretation
    - Links – Cross-references
    - Indices (sequence dbs, DiscoveryGate).
    - Metalayers (description of underlying data sources and business rules)
    - Data marts (duplication of information)
    - Standards: Unique definition of entities and relations (e.g. Semantic web)

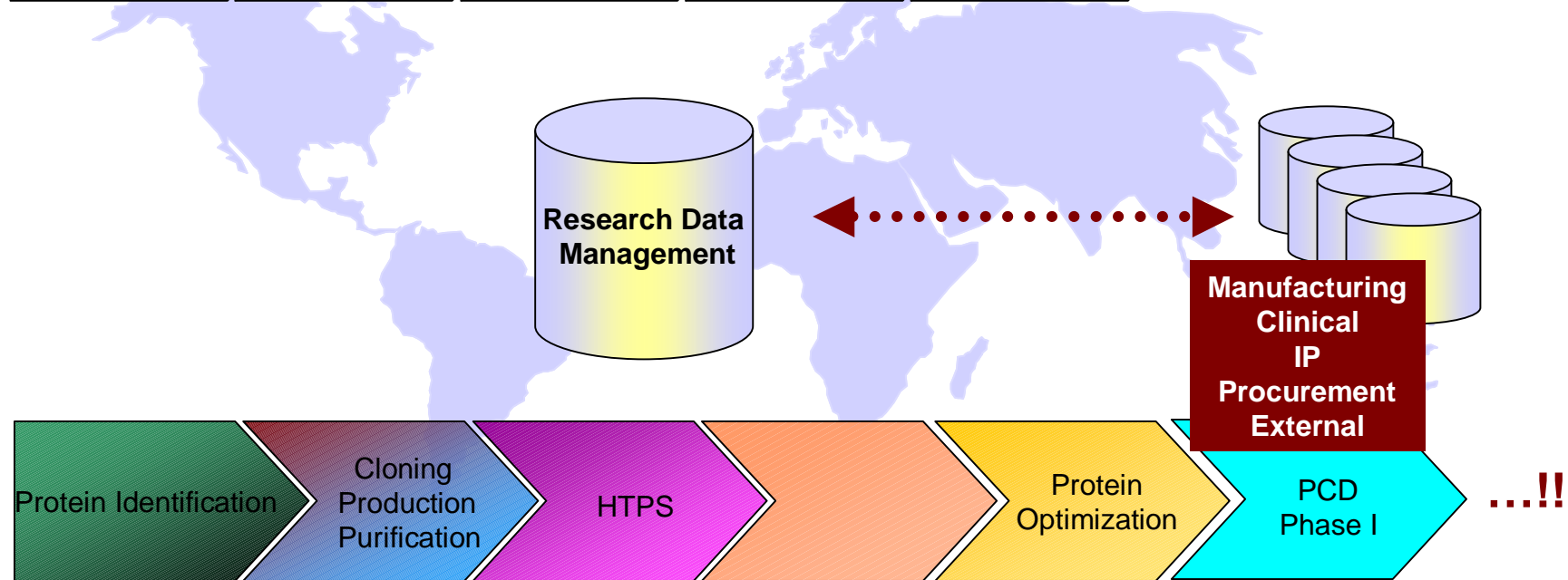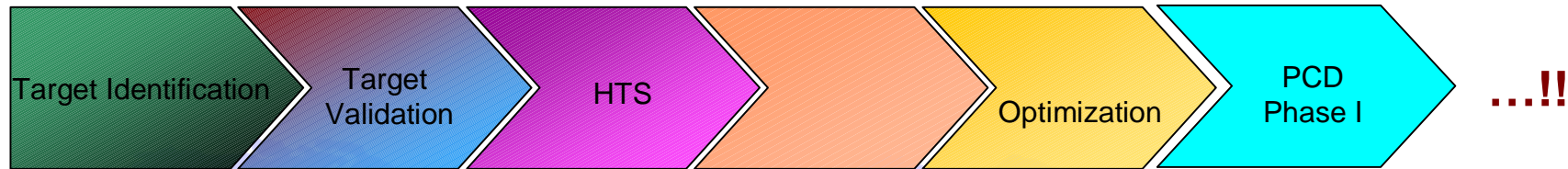- **Know-how acquisition**
    - ELNs
    - Documentation management

- **Unstructured data mining**
    - Text mining

- **Information mapping/scoring**
    - Mostly scoring
    - Viewers (e.g. sequence viewers)

# Research Context

**MERCK SERONO**

## Small molecules

| Target Identification | Target Validation | HTS | | Optimization | PCD Phase I | …!! |

**Research Data Management**

**Manufacturing Clinical IP Procurement External**

| Protein Identification | Cloning Production Purification | HTPS | | Protein Optimization | PCD Phase I | …!! |

## Proteins and Antibodies

**Right (Right) information to the right people at the right time in the right context**

# A Culture of Sharing

- **One place** for all biological, chemical & pharmacological information
- No boundaries **between sites, teams, projects**

**MERCK SERONO**

**Research Corporate Db**

Access to well integrated information:
- to make people "knowledgeable"
- to facilitate information flows
- to facilitate information tracking

# An Enabling Architecture for Integration

# Research Knowledge Space - Framework

**MERCK SERONO**

**User-friendly web interfaces**

**in silico knowledge generation tools  e.g. Text mining,  Data Analysis**

Corporate Database Core

**Integrated Oracle-Based Systems**

- **- Workflows**
- **- Aggregation**
- **- Projects**

**Drill-down**

E-notebook

**Publish/Control Posting APIs**

LIMS QC

LIMS QC

LIMS QC

Other Serono Sources

External Sources

**Specialized lab-oriented user interfaces & External sources**

# Research Knowledge Space

**External Results
And annotations**
Public and Commercial
Dbs

**Experimental Results
and annotations**
HTS, MTS, HCS,
Physchem, Preformulation,
Analytical, Selectivity,
In vivo pharmacology
ADMET, progression
**Extended to**
Expression, Phenotypes,
Disease Association,

**Relations**

**Entities**
Therapeutics:
SM, Proteins, mAbs
**Extended to**
Tissues
Targets
Primers
Plasmids
Cell lines
Mouse collection

**Corporate
Database**

Inventory
And
Laboratory
Data

# Knowledge Space – Framework – Molecule centric



**System comprising of:**

- a central gateway for handling entities (e.g.: genes, proteins, chemicals): the Merck Serono Wide Index (MS-Wii)

- specialized databases to handle large amount of raw data, each specific of a technology, referring to entities in the SWI

- a system to manage project-relevant information. Also, an easy-to-understand entry point for more detailed results and raw data (EDS)

- a tracking of production workflows

- same system for external data sources -> One single entry point

# Knowledge Space – External

- **Several categories:**
  - low-throughput: Must be curated, but high value. eg.: NTP reports database, public FDA databases
  - high-throughput: Large databases integrating heterogeneous data. Public large scale experiments e.g.: ENSEMBL, MDL MDDR, McKay
- **Integration**
  - Index in MS-Wii
  - Filter information relevant to the company's objectives.
- **Some issues:**
  - Some servers might not accept to full download
  - Need to build an appropriate analysis pipeline to interpret raw data
  - Requires prior analysis of structure and contents to design parsing scripts
    - "Thought-to-be structured"
- **Examples:**
  - 1- Sequences: new releases parsed and integrated in MS-Wii.
    - 1 s to get all in-house and external information for a given sequence whether it is a gene, transcript or protein
    - Information on past releases kept
    - Initiated same work on small molecules
  - 2- Toxicity data and predictions.

# Workflow Integration: Bioprocessing Pipeline

**Request follow-up**

**Production workflow**

Cloning — Expression — Purification

**Integration with the other data**

**Plasmids**

**Protein registration**

**Corporate Database**

# Screening data pipeline



**Raw data Storage**

**Processing QC**

**Corporate Db**

**Posting**

Data is transferred automatically, reducing human mistakes. Screeners can focus on QC.

# Electronic Laboratory Notebook



**MERCK SERONO**

**Queries, Reports (e.g. patent writing that automatically re-builds all reaction scheme)**

Inventory Reagent db

Corporate Database

Inventory

Analytical data

**+ Parallel Synthesis**
**+ Analytical modules**

Automatic structure import and calculations, registration in corporate db and inventory ease the chemists' work

# Workflow Integration

- Project-relevant information must remain in the Corporate DB

- Rule of thumb:
  - Corporate DB is a repository for all entities
  - LIMS processes raw data into information
  - LIMS associates information to entities in the Corporate DB

# Aggregation: Sequences

# Application – Protein Database

**Single access to the biological information:**
- **Description**
- **Gene expression**
- **Orthologs**
- **Disease link**
- **…**

**Grouping under a single protein all the annotations about related genes, transcripts or orthologs**

**Access to public and corporate data (e.g. Screening, production)**

# Reporting sequence related data



**Results are associated to a plasmid**

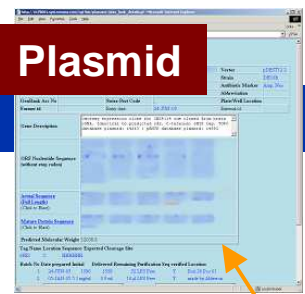**Results are displayed with said plasmid**

**They are also displayed with the protein's other results**

- Most results are associated to a protein (cell assays, in vivo models) but some are associated to other entities (plasmids,…).
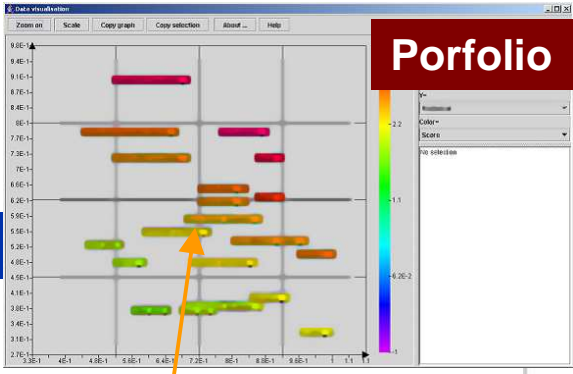
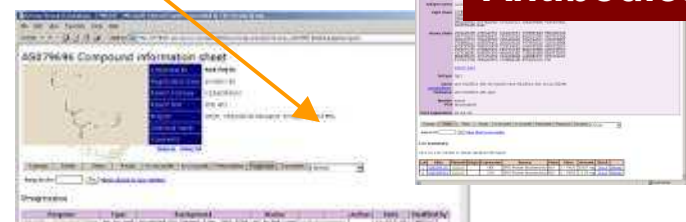- All must appear in the different contexts.

# From Target to ….

**Porfolio**

**Plasmid**

**Projects Progression**

**Genetic results**

**Target**

**Public database**

**Proteins Antibodies**

**Small molecules**

**Progression**

**MERCK SERONO**

# Molecule Centric Integration: One Compound, All Information Workflow + Aggreg



**Batch info**

**Screening results**

**Link to raw data**

**In vivo profile**

**Link to ELN or bioprocessing pipeline**

**Link to assay description**

**Catalogs Predictions Pubchem …**

**Inventory**

**In vitro profile**

# Structure Integration

# Result Integration

# FDA Toxicity Database



Venn diagram of the FDA toxicity databases

# "Structured" information sources

◆ **Integration requires prior in-depth analysis of contents**

◆ **Database (relational) <> structured.**
  - ◆ Would need to expose model and contents structure
  - ◆ No standards nor unique definitions

◆ **E.g. Toxicity.**
  - ◆ Endpoint definition, experimental details.
  - ◆ Models: Might predict the same activity but for different reasons.
    - ◆ Two different fragments responsible for same activity
    - ◆ Activity extrapolated from different experiments

# Needs

- **Integrate**
- **Monitor**
- **Informed decision-making**



**Research Corporate Db**

structure → phys-chem → in vitro → in vivo → in man

# Combining Toxicity Prediction Methods - Theoretical Study

- **Combining predictions increases either sensitivity <u>OR</u> specificity.**

- **Hypothesis:**

  - A substance is toxic if it reaches then binds/reacts to its target. These two requirements are modeled separately. We assume the substance reaches its target.

  - The binding event is due to the arrangement of some atoms in the molecule. The rest of them is not required but might interfere (eg: steric hindrance)

- **Method 1:** results for a single method
  - false positive rate:      30.97%
  - false negative rate:      12.04%

- **Method 2:** predicted positive if A and B are positive with a common descriptor
  - false positive rate:       0.25%      **=> highest specificity**
  - false negative rate:      22.41%

- **Method 3:** predicted positive when both A and B are positive
  - false positive rate:      12.10%
  - false negative rate:      17.87%

- **Method 4:** predicted positive when either A or B is positive
  - false positive rate:      48.93%
  - false negative rate:       6.52%      **=> highest sensitivity**

Simulation based on evaluating an outcome for 100000 molecules with different theoretical methods. Principle: different methods know some arrangements that cause toxicity, they miss some of them (=> false negatives) and include arrangements that are not a root cause (=> false positives).

# chemXML – sharing alerts



- Extends CML to substructure queries
  http://cml.sourceforge.net

- Handles almost all substructure query features described in MOL files, Derek and SMARTS

- Text-based format for sharing easily query definition

- Most of the validation is done by the XML parser (eg: Xerces)

# Concluding remarks

- **Unique Id for entities**
  - Advantages: Fast and reliable. "Easy" integration for in-house information
  - Drawbacks:
    - Indefinitely scalable?:
      - chemical space estimated to be 10^100 individual structures. Not incl mixtures.
      - Nessea: Developed algorithm to avoid the enumeration of Combichem Libraries
- **Production workflow tracking**
- **Project relevant information in corporate database**
- **Horizontal and vertical navigation**
- **Main trends:**
  - ELNs for know-how capitalization.
  - Information scoring/mapping
  - Extend scope. Cover whole life of entities.
  - Exploit data from outside
- **Main issue to allow integration and automated of parsing/interpretation:**
  - Lack of unique definitions/standards to depict entities, relations, contents

# Backup slide: Project Management

- Integrating project-related data
- Tracking product's life across pipeline