# GEPETTO : OPEN-SOURCE FRAMEWORK FOR GENE PRIORITIZATION

*Laboratoire de Bioinformatique et de Génomique Intégratives (LBGI)*
*Illkirch, France*

*Team leader : Olivier POCH*
*Project Manager : Hoan NGUYEN*
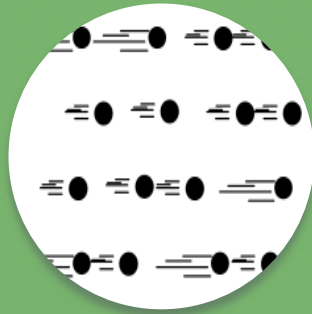*Presented by : Vincent WALTER*
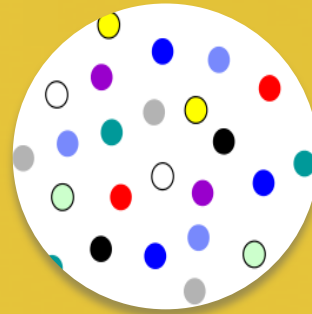
# BIG DATA CONTEXT



**VOLUME**

**Data at rest**

Terabytes to exabytes of existing data to process
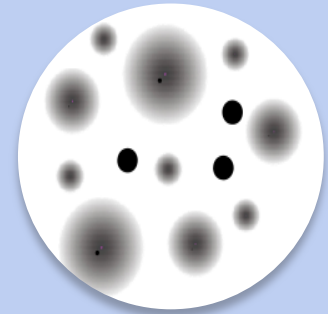
**VELOCITY**

**Data in motion**

Streaming data, queries take milliseconds to seconds to respond

**VARIETY**

**Data in many forms**

Structured, unstructured, text, multimedia

**VERACITY**

**Data in doubt**

Uncertainly due to data inconsistency and incompleteness, ambiguities, latency, model approximations

# SM2PH – BIG DATA MANAGEMENT

SM2PH (**S**tructural **M**utation to **P**athology **P**henotypes in **H**uman) http://decrypthon.igbmc.fr/sm2ph/

```xml
SRP9.xml
1   <?xml version="1.0" encoding="UTF-8" standalone="yes"?>
2   <GeneCard xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://decrypthon.igbmc.fr/sm2ph/profiles/ sm2ph.xsd">
3       <GeneInformation>
4           <GeneName>
5               <GeneSymbol>SRP9</GeneSymbol>
6               <GeneId>6726</GeneId>
7           </GeneName>
8           <Localisation>
9               <Chromosom>1</Chromosom>
10              <Cytoband>1q42.12</Cytoband>
11          </Localisation>
12      </GeneInformation>
13      <ProteinList>
14          <Number>1</Number>
15          <Protein>
16              <ProteinInformation>
17                  <UniprotAccession>
18                      <Accession>P49458</Accession>
19                      <Synonym>Q6NVX0</Synonym>
20                      <Synonym>Q8WTW0</Synonym>
21                  </UniprotAccession>
22                  <ProteinName>
23                      <EntryName>SRP09_HUMAN</EntryName>
24                      <Name>Signal recognition particle 9 kDa protein</Name>
25                      <Synonym>SRP9</Synonym>
26                  </ProteinName>
27                  <Comments>Signal-recognition-particle assembly has a crucial role in targeting secretory proteins to the rough endoplasmic r
28              </ProteinInformation>
29              <OrthologList>
30                  <Ortholog specie="RAT">
31                      <Symbol>Srp9</Symbol>
32                      <Name>Srp9</Name>
33                      <Chromosom>13q26</Chromosom>
34                      <SourceId>ENTREZ:690345</SourceId>
35                  </Ortholog>
36                  <Ortholog specie="MOUSE">
37                      <Symbol>Srp9</Symbol>
```
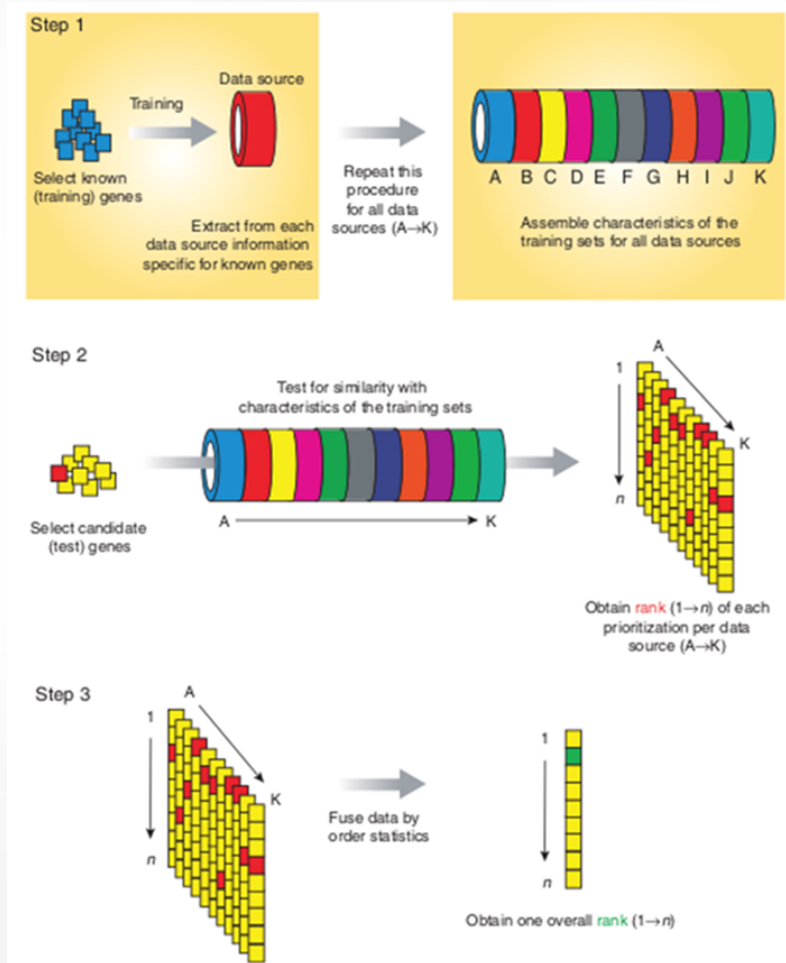
# PRIORITIZATION PROCESS

- Identification of the **most promising** feature associated to a question (ex. biological process, pathology, network) in dynamic systems

- **3 STEPS PROCESS**

  - Building model for training features

  - **Local** prioritizations
    - Candidate feature profile
    - **Evaluation** according to the similarity with the model (= scoring)
    - **Ranking**

  - **Global** prioritization

**GEPETTO** = dedicated to human gene prioritization



*Aerts, S., et al. **Gene prioritization through genomic data fusion.** Nature Biotechnology 24, 537–544.*

# WHY AN OPEN-SOURCE FRAMEWORK?

**SEVERAL GENE PRIORITIZATION TOOLS**
- New ones published regularly
- Most of them are web GUI

**BUT MANY DRAWBACKS**
- **Maintenance** – short period
- **Availability** - couple of months

- **Queryability** – not for high-throughput process
- **Integrability**– no libraries for gene prioritization

- **Extensibility** – not open-source
- **Data** – key aspects are not integrated (genomic context, 3D-structure,…)

| Tools | Functional annotations | Expression | Text (co-citation) | Text (functional) | Interactions | Pathways | Sequence | Phenotype | Conservation/homology | Disease probabilities |
|---|---|---|---|---|---|---|---|---|---|---|
| Candid | | X | | X | X | | X | | X | X |
| DGP | | | | | | | X | | X | |
| Endeavour | X | X | | X | X | X | X | | | X |
| GeneRank | X | X | | | | | | | | |
| GeneRanker | X | | X | X | X | | | X | | |
| GeneSeeker | | X | | X | | | | X | X | |
| PolySearch | X | | | X | X | X | | | | |
| PosMed | X | | X | X | X | | | X | | |
| SNPs3D | X | | | X | X | X | X | X | | |
| ToppGene | X | X | X | | X | X | X | X | | |

# GEPETTO FRAMEWORK : SPECIFICATIONS

**OPEN-SOURCE**
- Available on SourceForge platform **http://sourceforge.net/projects/gepetto/**
- **Supported** by the community
- **Improved** by the community
- **Completed** by the community
  - Adding **your own** datasources

**ACCESSIBILITY**
- **Biologists and clinicians**
  - **Web inteface**
- **Bioinformaticians**
  - **Standalone application (command line), Java API**

**STANDARDIZATION / FLEXIBILTY**
- **Facilitated** maintenance
- **Modularity** - 1 JAR project by local prioritization
- **Scalability / Customizability**

**VERSIONING/ CONTRIBUTION**
- **SVN** (via SourceForge)
- GitHub and Mercurial (in progress)

# GEPETTO FRAMEWORK : TECHNOLOGIES

**PROGRAMMING LANGUAGES**

- **Java : modules libraries / core of the application**
- **Python : command line launcher**
- **R : statistical tests for local prioritization**

**PLATFORMS**

- **Only UNIX systems are supported today**
  - **Debian-based distributions (Ubuntu,...)**
    - Débian package to install easily
  - **Redhat-based (Fedora,...)**
    - Source tarball (tar.gz archive)

**INTEROPERABILITY**

- **Galaxy integration**
  - **I/O : Use file as input and can generate file as output**
  - **Queryable by command line**

# LOCAL PRIORITIZATION MODULES

## GENERIC MODULES

- **Protein sequence** (Alignment)
  - E-value of the best hit with BLASTp

- **Evolutionary barcodes** (Evolucode)
  - Evolutionary histories in vertebrates (16) for all human genes
  - Based on K-means clustering

- **Genomic context** (Geco)
  - Nested repeats, open chromatin, PolII,...

- **Transcriptomic** (GxDb)
  - Tissular expression in 79 human tissues
  - Pearson Correlation and Fisher's omnibus analysis

- **Protein-protein interactions** (String)
  - Physical / functional interactions
  - Size of the overlap with the training model

## SPECIFIC MODULES

- **Hereditary disease gene probability (IDGP)**
  - Combination of the recessive / dominant probability of involvment

- **Genome-wide association study (GWAS-AMD)**
  - Double GC value

Linard, B., et al. (2011). EvoluCode: **Evolutionary Barcodes as a Unifying Framework for Multilevel Evolutionary Data.** Evolutionary Bioinformatics 61.

8

# GLOBAL PRIORITIZATION

- Global scoring and ranking
  - Data fusion
  - Based on the different local ranking

Aerts S, et al. **Gene prioritization through genomic data fusion. Nat Biotechnol. 2006 May;24(5): 537-44.**

Britto, R., et al. (2012). **GPSy: a cross-species gene prioritization system for conserved biological processes--application in male gamete development.** Nucleic Acids Research.

### ORDER STATISTICS

- Able to handle features with missing values.
- Minimizes the bias for known or well-characterized features

- $Q(r_1, r_2, ..., r_N) = N! \int_0^{r_1} \int_0^{r_2} ... \int_0^{r_n} ds_N \, ds_{N-1} ... ds_1$

### ROBUST RANK AGGREGATION

- Similar to order statistiss
- More significant result for global prioritization

- $P(r_k' \leq r_k) = \sum_{i=k}^{n} \binom{n}{i} \cdot (r_k)^i (1 - r_k)^{n-i}$

### MALLOWS MODEL

- Looks for a consensus ranking

- $(\pi_1, ..., \pi_n) = 1/n \cdot \sum_{k=1}^{n} 1_{\{i < \pi_k^j\}}$

### GPSy

- Assigns weights to the different criteria to find the optimal weight combination

- $r_i = \sum_{j=1}^{n} w_j \cdot r_{ij} / \sum_{j=1}^{n} w_j$

# GEPETTO WORKFLOW USING JBPM

- A division of an american multinational software company : RedHat
- Specialized in writting and supporting open-source middleware software
- Provides tools for Java applications

- Open-source workflow engine
- Since 2003 (Don't reinvent the wheel)
- Complete workflow engine
  - Manages information flows
  - Makes the bridge between biologists (business analysts and end users) and bioinformaticians (developers).
  - Provides features not provided by other workflow engines
- Based on BPMN2 standard notation (Business Process Model and Notation)
- Describes the process step by step

# ADVANTAGES OF JBPM

| | |
|---|---|
| Field usage | • Universal / Not oriented<br>• More flexible |
| Installation | • Import .JAR libraries<br>• Very easy |
| Skills required | • No computer programming competences required |
| Workflow – possibilities | • Complex<br>• Parallel gateway or  exclusive gateway (decisionnal workflow) |
| Workflow – definition | • BPMN2 (Business Process Managment Notation) – standard notation<br>• XML |
| Workflow – graphical editor | • Web Service for JbossAS (Application Server)<br>• Web Service for Apache Tomcat |
| Workflow – passing data mode between nodes | • Java / Object (directly in memory)<br>• Not necessary to handle many files |
| Workflow – languages or programs | • All languages / programs wrapped into WorkItemHandler |

# JBPM : EXAMPLE OF SPECIFIC FEATURES

**PARALLEL GATEWAY**
- Used to split / synchronize the respectively incoming or outgoing sequence flow
- Executes task in parallel

*GEPETTO : Used for local prioritization*
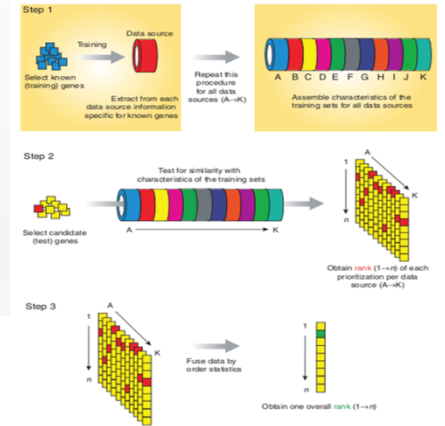
```
<process id="parallelGateway" name="BPMN2 example parallel gatewar">

  <startEvent id="Start" />

  <sequenceFlow id="flow1" name="fromStartToSplit"
    sourceRef="Start"
    targetRef="parallelGatewaySplit"  />

  <parallelGateway id="parallelGatewaySplit" name="Split"
    gatewayDirection="diverging"/>

  <sequenceFlow id="flow2a" name="Leg 1"
    sourceRef="parallelGatewaySplit"
    targetRef="prepareShipment" />

  <userTask id="prepareShipment" name="Prepare shipment"
    implementation="other" />

  <sequenceFlow id="flow2b" name="fromPrepareShipmentToJoin"
```

**EXCLUSIVE GATEWAY**
- Guides the workflow progression using any p

*GEPETTO : Used to select the global prioritizatio*

```
<process id="exclusiveGateway" name="BPMN2 Example exclusive gateway">

  <startEvent id="start" />

  <sequenceFlow id="flow1" name="fromStartToExclusiveGateway"
    sourceRef="start" targetRef="decideBasedOnAmountGateway" />

  <exclusiveGateway id="decideBasedOnAmountGateway" name="decideBasedOnAmount" />

  <sequenceFlow id="flow2" name="fromGatewayToEndNotEnough"
    sourceRef="decideBasedOnAmountGateway" targetRef="endNotEnough">
    <conditionExpression xsi:type="tFormalExpression">
      ${amount < 100}
    </conditionExpression>
  </sequenceFlow>

  <sequenceFlow id="flow3" name="fromGatewayToEnEnough"
    sourceRef="decideBasedOnAmountGateway" targetRef="endEnough">
    <conditionExpression xsi:type="tFormalExpression">
      ${amount <= 500 && amount >= 100}
    </conditionExpression>
  </sequenceFlow>

  <endEvent id="endNotEnough" name="not enough" />

  <endEvent id="endEnough" name="enough" />

  <endEvent id="endMoreThanEnough" name="more than enough" />

</process>
```

12

# GEPETTO WORKFLOW USING JBPM



**STEP 2**

**LOADER**

**GLOBAL PRIORITIZATION**

**STEP 1**

**LOCAL PRIORITIZATION**

**STEP 3**

# ADD NEW LOCAL PRIORITIZATION MODULE

1. Download GEPETTOModel.jar

2. Create new JAR project

3. Implements class that inherits GEPETTOModel interfaces using Features provided by GEPETTOModel library.

4. Create/Edit MANIFEST.mf file

5. Update

# ADD NEW LOCAL SCORING METHOD

- Local scoring method are managed using Polymorphism
- To add a new scoring method
  - Implement a new class which inherit from GenePrioritization or ProteinPrioritization
  - Add this class to the ModulePrioritizationFactory
  - Change the method used in Spring Application Context File (XML file)



- Spring Framework

  - Open-source Java application framework
  - Inversion Of Control (IOC) container

  - Facilitates development and test of Java (Web) Applications
  - Standardization of source code : more robust and easier to maintain

  - Reduce code of SM2PH-DAO(Data Access Object)
  - GEPETTO uses *BeanFactory* /Spring IOC (Dependency Injection)

# APPLICATIONS – AMD USE CASE

- **AMD** (Age-related Macular Degeneration)
  - Macula degeneration
  - Affects **old people** (over 50 years)
  - Causes a significant **weakening** of the **visual capabilities**

- Confidential GWAS-data
  - Provided by the AMD Gene Consortium
    - T. Léveillard, et al.
  - 8.000 patients
  - 8.000 controls

- Data sets
  - Genes involved in AMD (12 SNPs, 14 genes in KEGG, Pubmed)
  - Recently validated new genes (7 SNPs, 9 genes)
  - Candidate genes out of the limit of GWAS detection (21 SNPs, 29 genes)

AMD Gene Consortium, Poch O, Ripp R, Léveillard T among 156 writers, **Seven new loci associated with age-related macular degeneration**. Nat Genet. 2013 Apr;45(4):433-9. doi: 10.1038/ng.2578. Epub 2013 Mar 3.

# Applications – AMD use case

**APPLICATION 1**

- Evaluating the ability to **successfully detect** AMD validated genes

- Comparison against Endeavour / ToppGene
  - Method : **ROC-AUC**
  - Test set : AMD known and new genes (23) + Genes closed to GWAS limits (29)

| TRAINING SET | AUC (prioritization of the 23 AMD genes) | | |
|:---:|:---:|:---:|:---:|
| | **GEPETTO** | **ToppGene** | **Endeavour** |
| **Known genes (14)** | *0,649* | 0,720 | 0,701 |
| **New genes (9)** | *0,825* | 0,479 | 0,643 |
| **Mixed genes (11)** | *0,887* | 0,854 | 0,903 |

- Conclusion
  - **0,64 < AUC < 0,89** : GEPETTO uses a good model.
  - Able to detect genes that are experimentally validated (over the limit of GWAS detection)

# APPLICATIONS – AMD USE CASE

**APPLICATION 2**

- Evaluate the ability to **discriminate** AMD target genes and Retinitis Pigmentosa (RP) known genes

- Compared to Endeavour
  - Method : **ROC-AUC**
  - Training set : AMD known genes
  - Test set : All AMD candidates + RP known genes

| DISCRIMINATION CRITERIA | AUC | |
|---|---|---|
| | **GEPETTO** | **Endeavour** |
| TP = AMD known genes ; FP = Other genes (AMD/RP) | 0,967 | 1,000 |
| TP = AMD validated genes ; FP = Other genes (AMD/RP) | 0,859 | 0,892 |
| TP = AMD candidates genes ; FP = RP genes | 0,752 | 0,564 |

- Conclusion
  - **0,75 < AUC < 0,97** : GEPETTO uses a good model.
  - Very discriminant in some cases

# PERSPECTIVES

| | |
|---|---|
| **Optimisation** | • Implementation of the 3 new methods of global ranking<br>• Implementation of new methods of local ranking |
| **New parameters** | • Using SNP integration from MSV3d (**MisSense Variants mapped to 3D-structures**) / Bayesian networks |
| **Extends to other organisms** | • Human-centric gene prioritization<br>• Add the possibility to prioritize genes for other species |
| **Pattern extraction** | • Extraction of patterns for hereditary disease causing genes |
| **Application** | • Development of modules devoted to ciliopathies or other rare diseases |

# ACKNOWLEDGEMENT

PARSEC
Poster No44

Olivier POCH                  Alexis ALLOT
Julie THOMPSON                Kirsley CHENNEN
Hoan NGUYEN                   Etienne GOFFINET
Odile LECOMPTE               Alan LAHURE
Luc MOULINIER                Can MOREL
Jean MULLER                  Anaïs NICOL
Laetitia POIDEVIN            Hélène POLVECHE
Raymond RIPP                 Raphaël SCHNEIDER
Wolfgang RAFFELSBERGER  Carlos BERMJO-DAS-NEVES