



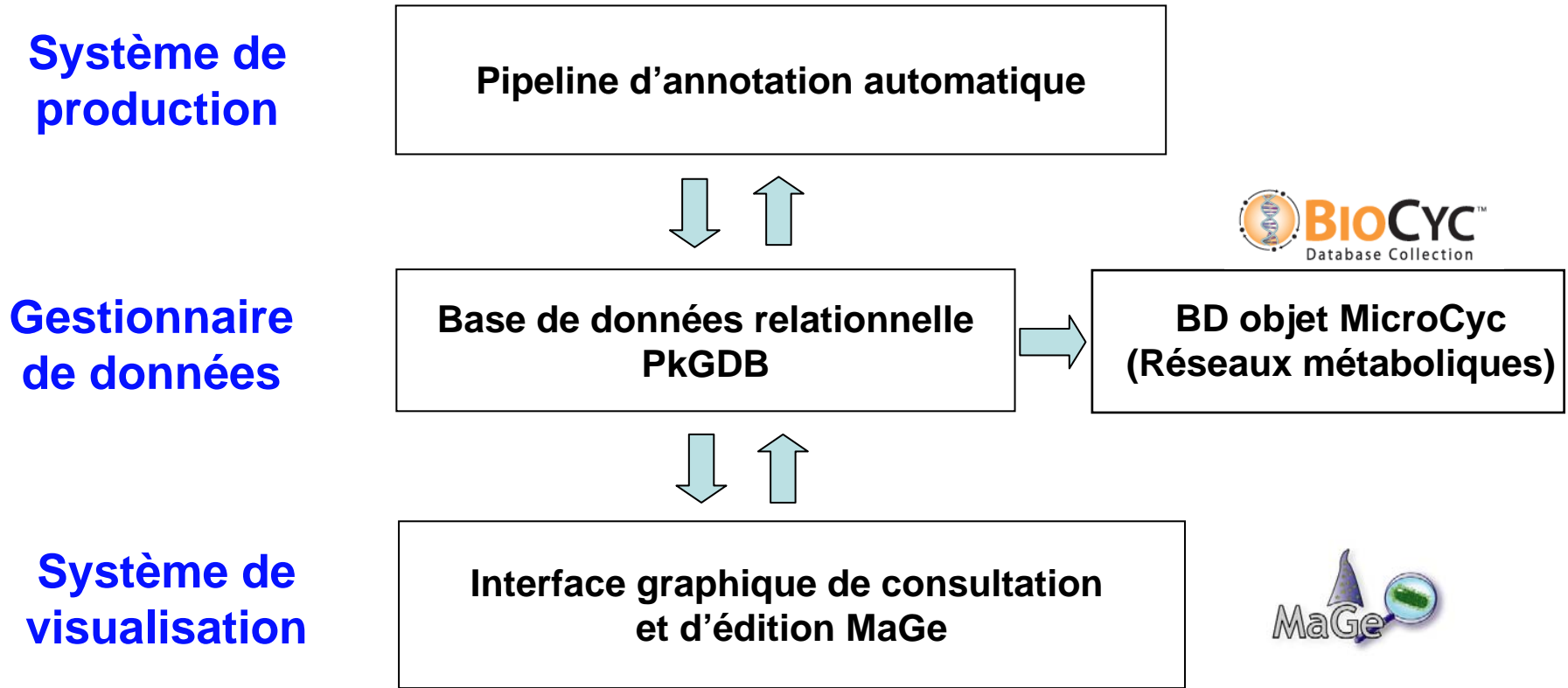
Automatisation du pipeline d'annotation de la plateforme MicroScope pour un passage à large échelle

Stéfan Engelen & David Vallenet

CEA/DSV/IG/Genoscope
CNRS-UMR 'Génomique Métabolique'
Laboratoire (Atelier)
de Génomique Comparative



Les composants de MicroScope



Explorer et éditer les connaissances contenues dans le gestionnaire de données :

- Outils de génomique comparative
- Annotation experte

Pipeline d'annotation : flux d'information

SYNTAXIQUE

Séquences
génomiques

- RNA and protein genes
- Transcription/translation start & stop
- Nucleotide composition and «Words»
 - Codon usage
 - Genomic islands

Gènes / Protéines

FONCTIONNEL

Calculs, recherche
de similarités

- Ortho/Para/Homologs
- Gene/protein families
- Subcellular localization
 - Motifs

Assignations
fonctionnelles

Protéines annotées

RELATIONNEL

Reconstruction
de processus

- Gene context, gene order
- Comparative genomics: PhyloProfile
- Gene fusion/fission
- Regulatory networks
- Protein interaction
- Metabolic networks

Processus
biologiques

EXPERTISE
HUMAINE

Visualisation
des données

Annotations
expertes

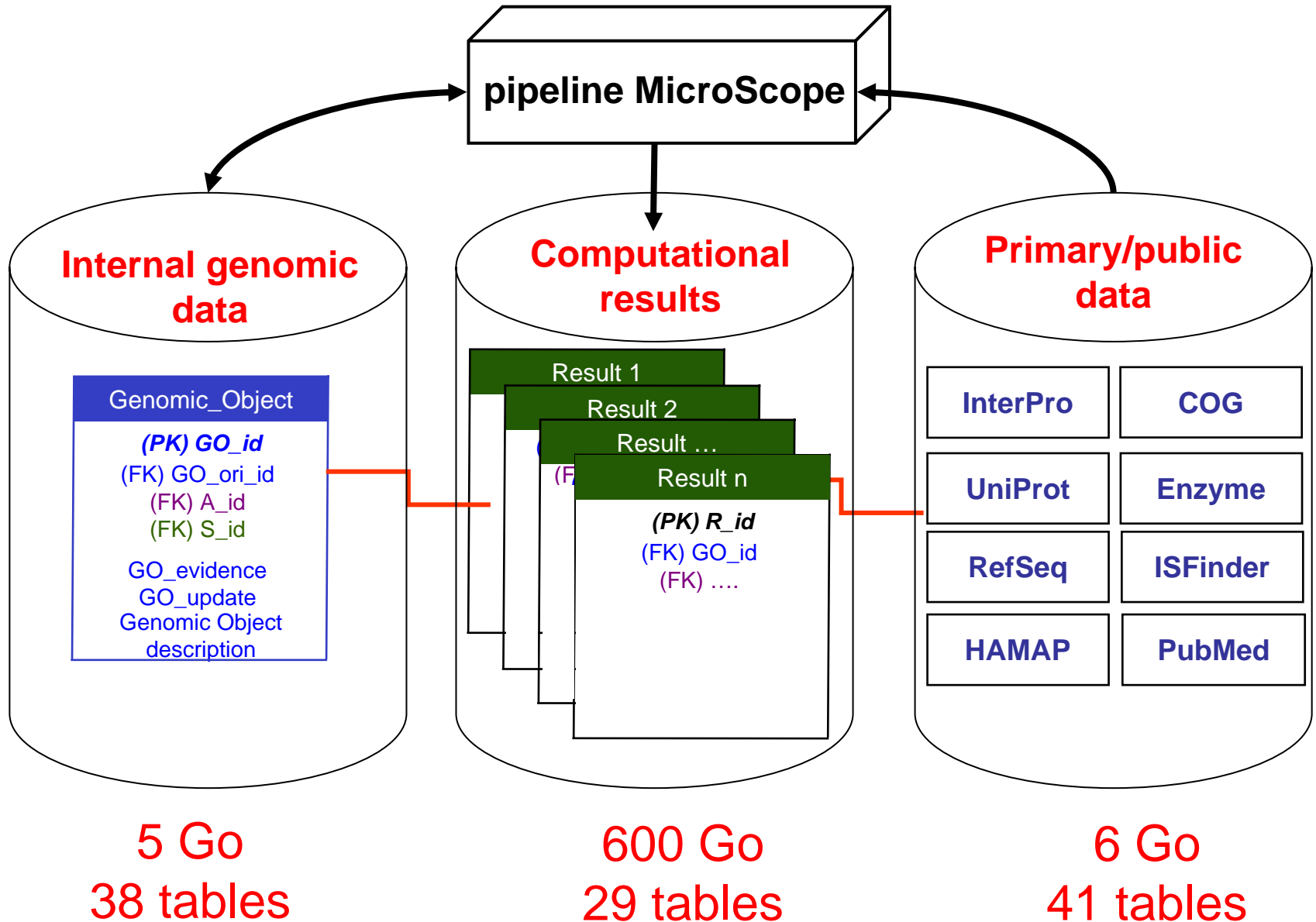
EXPERTISE
HUMAINE

Pipeline d'annotation : calculs

- ❑ Plus de **20 workflows différents** pour l'annotation syntaxique, fonctionnelle et relationnelle.
- ❑ **Temps CPU** (AMD optéron double-cœur) pour un génome (5000 gènes/protéines)
 - Blast sur **Uniprot** : 73 heures
 - Blast **PkGDB** + **Syntenie** : 16 heures
 - Blast **RefSeq** + **Syntenie** : 32 heures
 - **InterPro** : 305 heures
 - **Autres analyses** (COG, PRIAM, BioCyc, Rfam) : 10 heures

=> Total = 436 heures
- ❑ **Parallélisation sur 36 CPU** (72 cœurs)
 - pour **1 génome = 6 heures**
 - pour tous les génomes de PkGDB = 86 jours

La base PkGDB



Interface graphique : Magnifying Genomes



MicroScope project

Authentication

Help(s)

Options

Genome Overview

Export

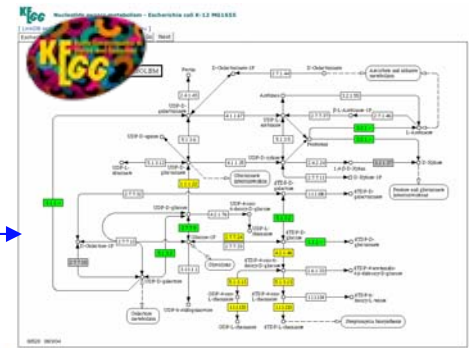
Gene Validation - ACIAD1137 (current annotation made by MaGe)

Type	Begin	End	Length	Frames	Mutation	Gene	Synonym	Date	Status
1137	54946	55943	998	1	0	ACIAD1137		2004-07-28 15:52:06	Valid

Product: ACIAD1137 protein (Acinetobacter sp.)

Comments: ACIAD1137 is a protein of 332 amino acids. It is a member of the ACIAD family of proteins. It is a member of the ACIAD family of proteins. It is a member of the ACIAD family of proteins.

Annotator editor

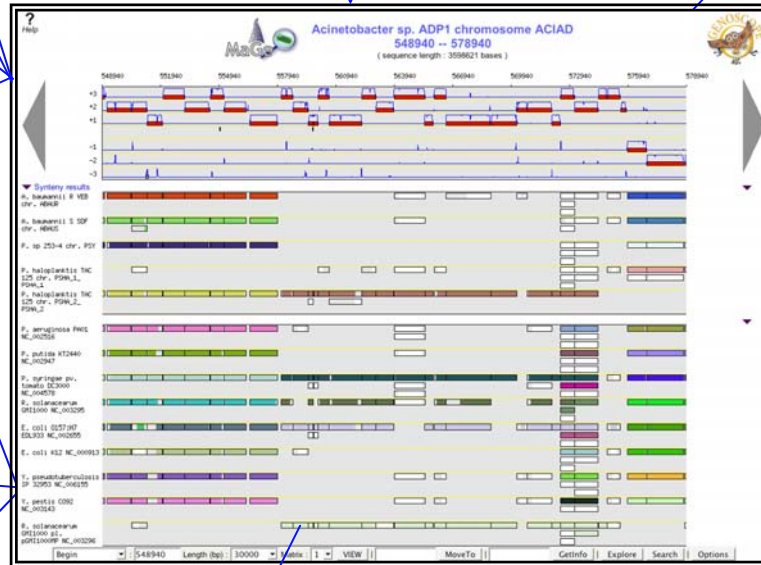


Pathway Hunter Tool

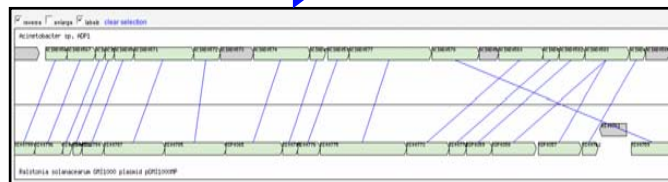
Metabolic pathways

- KeyWords
- Blast / Motif
- Phylogenetic profiles
- Fusions / Fissions
- Genomic islands
- Metabolic profiles

Exploration



Synteny map



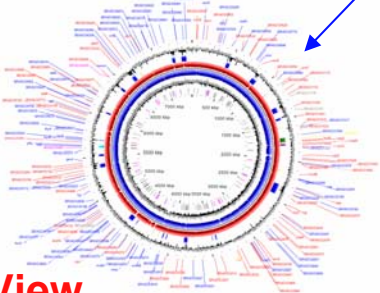
Synton visualization



Artemis



LinePlot

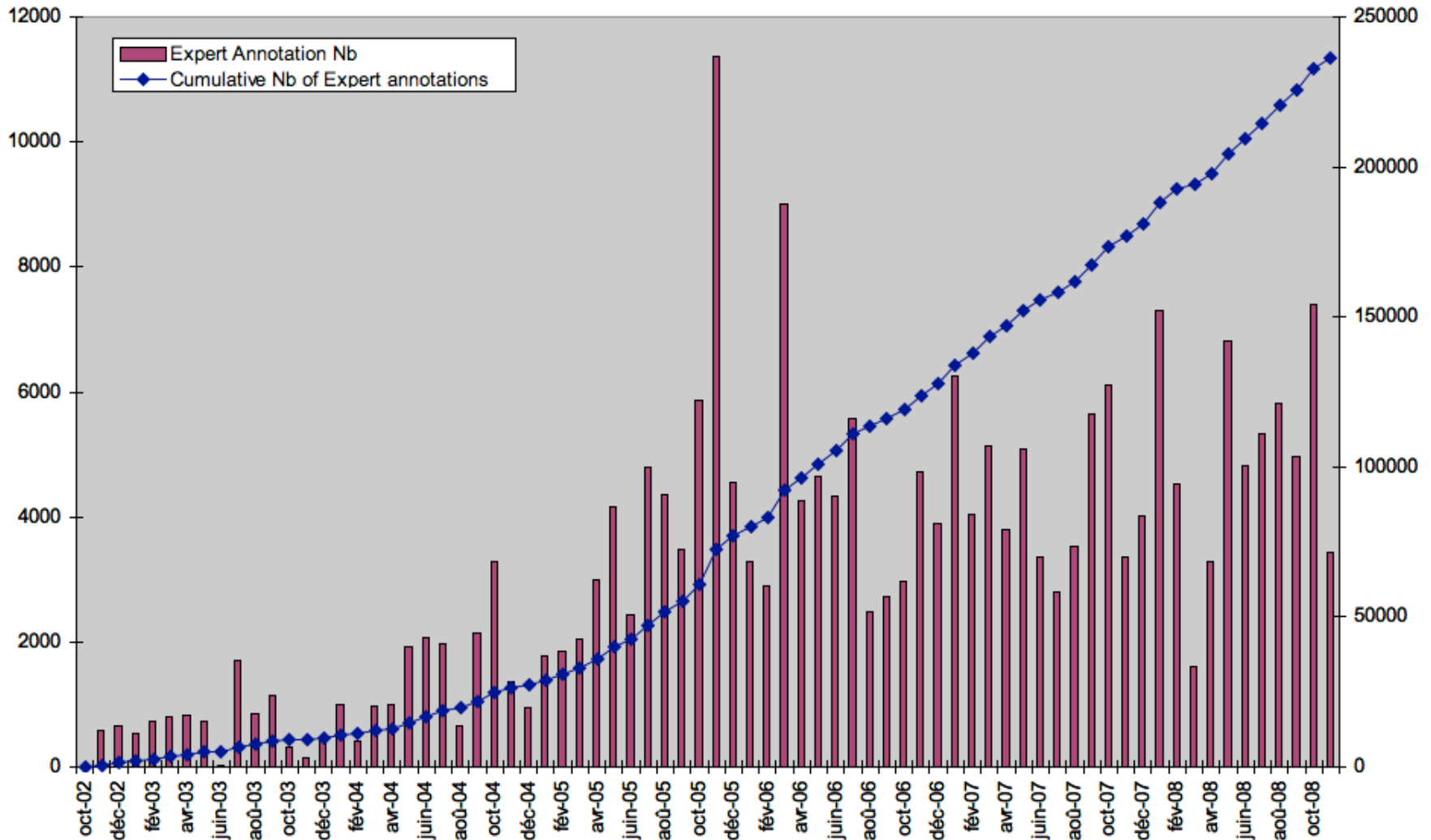


CGView

Taux d'utilisation de la plateforme

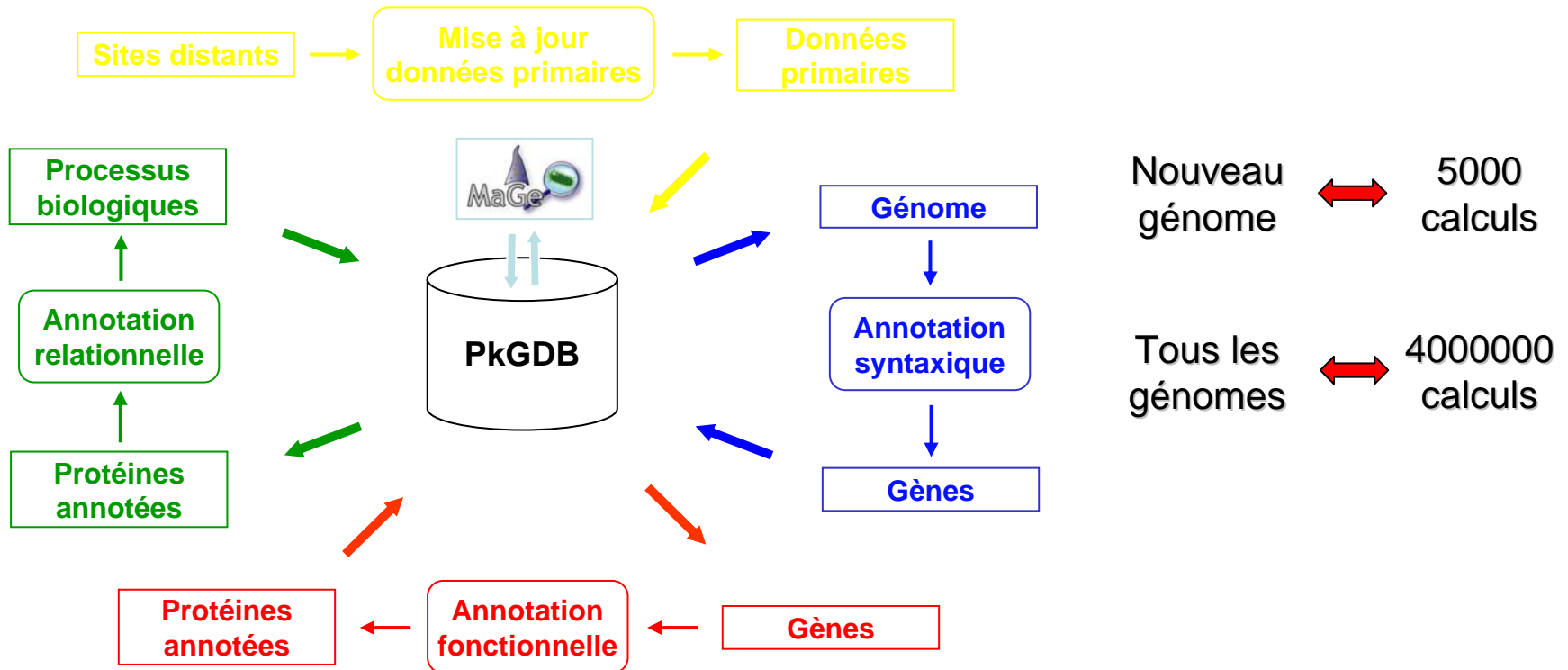
- ❑ **563 comptes** personnels {
 - 380 en France**
 - 76 en Europe**
 - 45 Etats-Unis + 62 autres pays**
- ❑ Projets collaboratifs avec une **trentaine de laboratoires en France** : Universités Paris, Strasbourg, Lyon, Marseille - CNRS - Institut Pasteur - CIRAD Montpellier – INRA - CEA.
- ❑ Projets collaboratifs avec **16 laboratoires hors France** : USA, Chine, Corée, Allemagne, Belgique, Portugal, ...
- ❑ **Partenaire(s) industriel(s)** : Sanofi-Aventis
- ❑ **70 authentications par jour** (env. 150 pages consultées par session)
 - ⇒ Efficacité du serveur Web MaGe (interface graphique)

Evolution du nombre d'annotations expertes



■ **236 233 annotations expertes** sont à ce jour enregistrées dans la base PkGDB

Pipeline d'annotation : objectifs ?



Objectifs

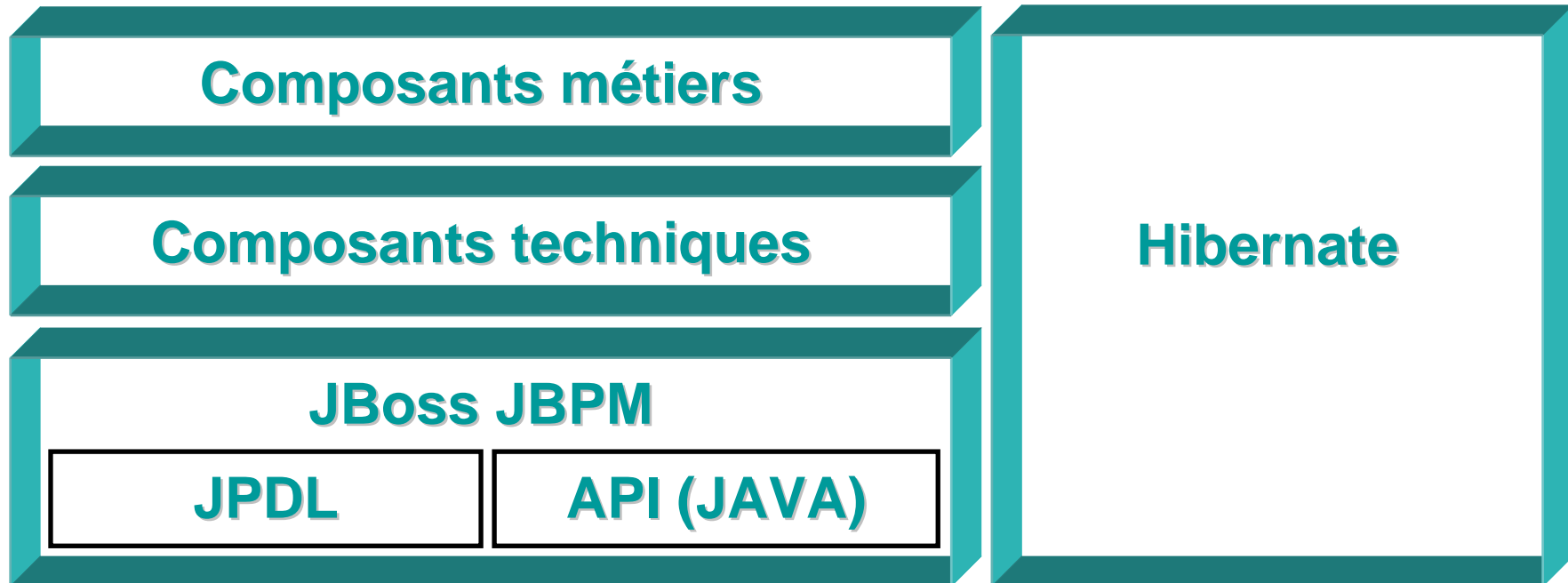
- ❑ Automatisation pour un passage à large échelle
- ❑ Garantir les performances d'accès (concurrency requêtes utilisateurs / chargement calculs)
- ❑ Maintenir les données primaires et calculs à jour

Identification des besoins

- ❑ **Définir / décrire** rigoureusement les processus métiers :
diagramme d'activités
- ❑ **Orchestrer / synchroniser** les activités humaines, calculs et processus systèmes l'aide d'un langage de description
- ❑ **Robustesse** : reprendre sur échec une activité
- ❑ **Contrôler et suivre** l'évolution des calculs en temps réel
- ❑ **Garantir la traçabilité** : sauvegarde de l'historique

JBPM - Java Business Process Management : solution technique à la mise en œuvre des processus métiers

- ❑ Solution open-source, souple, capable d'orchestrer un ensemble de programmes, quelque soit le langage d'implémentation
- ❑ JPD L : langage de description des activités (**Orchestrer / synchroniser**)
- ❑ API (JAVA) : développement de nouveaux composants techniques et métiers (**suivi / contrôle**)
- ❑ Hibernate : Persistance dans une BD (**robustesse / traçabilité**)



JBPM : cycle de développement

Identification des cas d'utilisation



Ecriture des scénarios



Elaboration du diagramme d'activités



Traduction avec JPDL

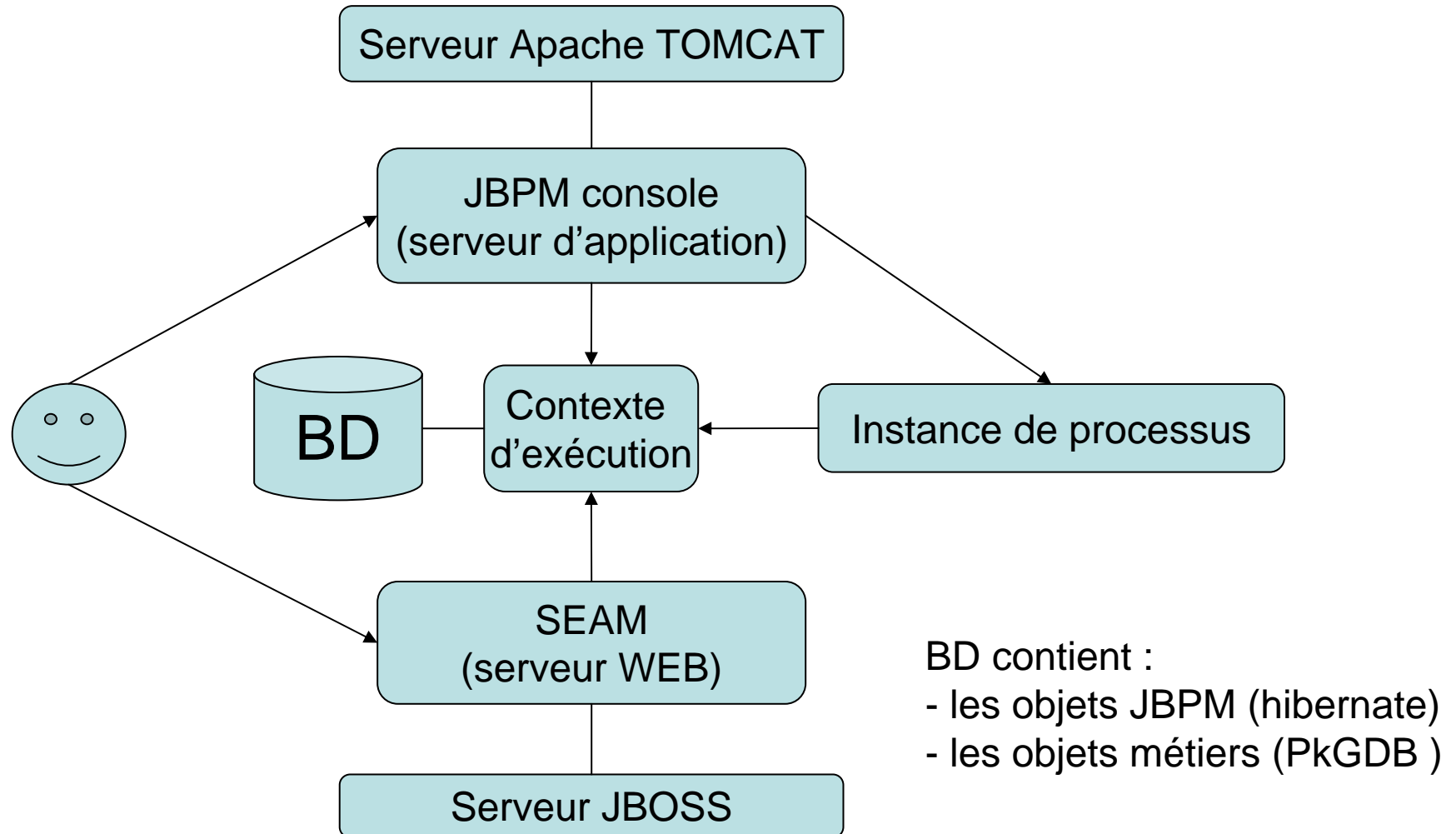


Adaptation des composants techniques



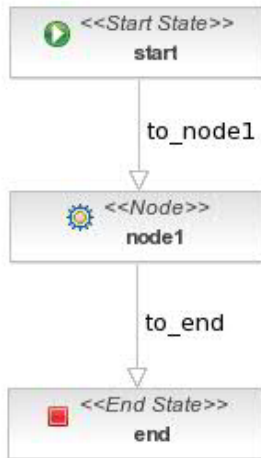
Développement des composants métiers

JBPM : architecture logicielle



Définition du processus

- ❑ On décrit le processus avec le langage **JPDL** (Java Process Definition Language) au format XML
- ❑ JPDL décrit un graphe et des propriétés dans chaque noeud et transition du graphe



```
<?xml version="1.0" encoding="UTF-8"?>
```

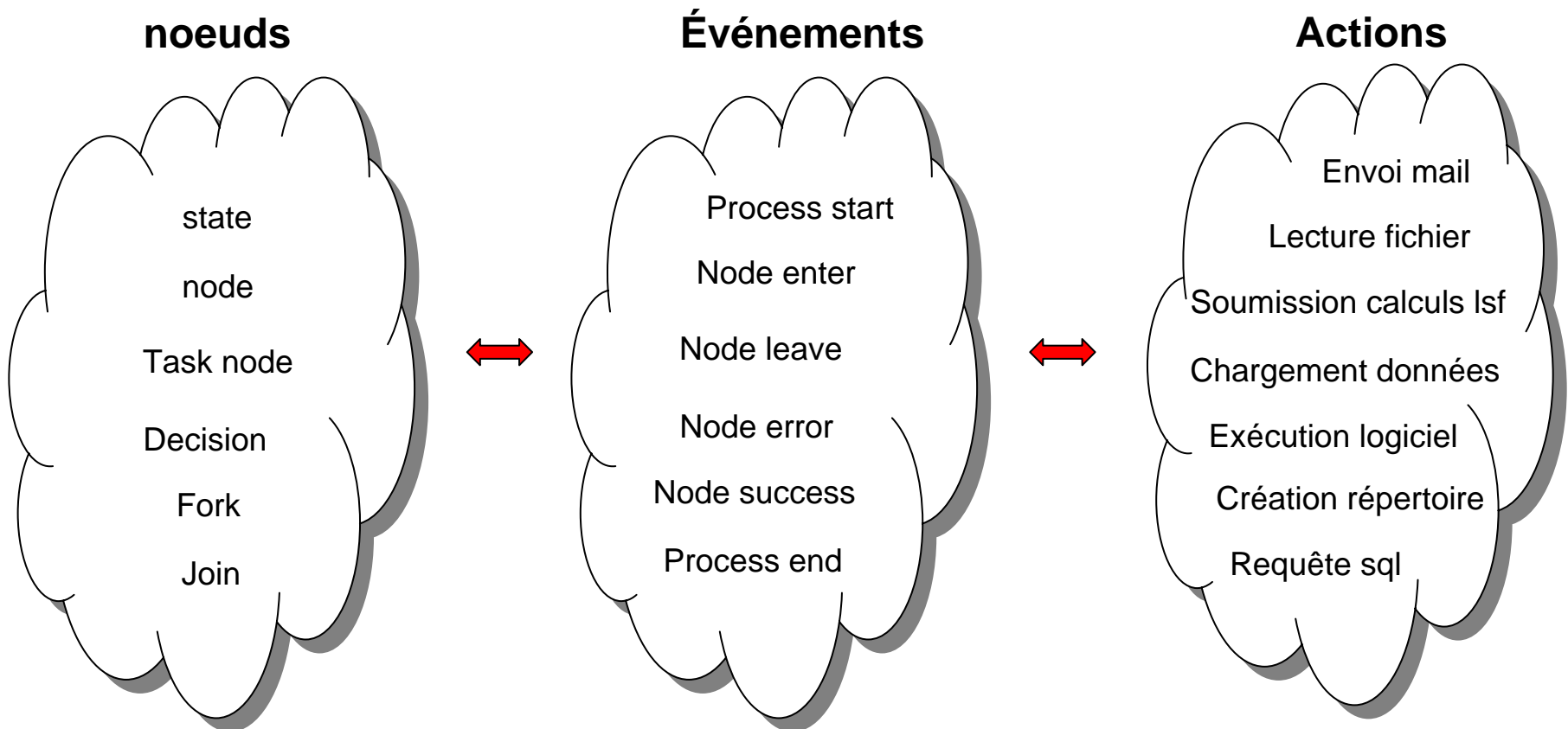
```
<process-definition xmlns="" name="Demo">  
  <start-state name="start">  
    <transition name="to_node1" to="node1"></transition>  
  </start-state>  
  <end-state name="end"></end-state>  
  <node name="node1">  
    <transition name="to_end" to="end"></transition>  
  </node>  
</process-definition>
```

Définition du processus

- ❑ Plusieurs type de nœuds dont l'activité dépend :
 - ❑ d'un comportement initial du nœud
 - ❑ **State** (état d'attente)
 - ❑ **Node** (état à comportement personnalisable)
 - ❑ **TaskNode** (assigne et crée des tâches)
 - ❑ **Decision** (dirige le flux en fonction d'un critère)
 - ❑ **Fork** (sépare le chemin d'exécution entrant)
 - ❑ **Join** (état bloquant/passant)
 - ❑ de l'exécution d'un certain nombre d'actions implémentées via l'API JBPM
 - ❑ de la synchronisation des actions par des évènements
- ❑ Un seul type de transition

Définition du processus

- ❑ L'activité d'un nœud correspond à l'exécution d'actions associées à des types d'évènements
- ❑ Action : composant métier reposant sur des composants techniques



Cas d'utilisation : workflow de calculs des synténies

Comparaison blast

Tache métier 1 : whichGenome

- génomes non à jour

Tache métier 2 : executeLSFscript

- comparer 2 à 2 n génomes avec blast
- 2ⁿ calculs parallélisés sur un cluster
- calculs terminés ou échoués

Tache métier 3 : loadDataAnalysis

- chargement dans PkGDB des calculs terminés à la tache 2

Calcul de synténies

Tache métier 1 : whichGenome

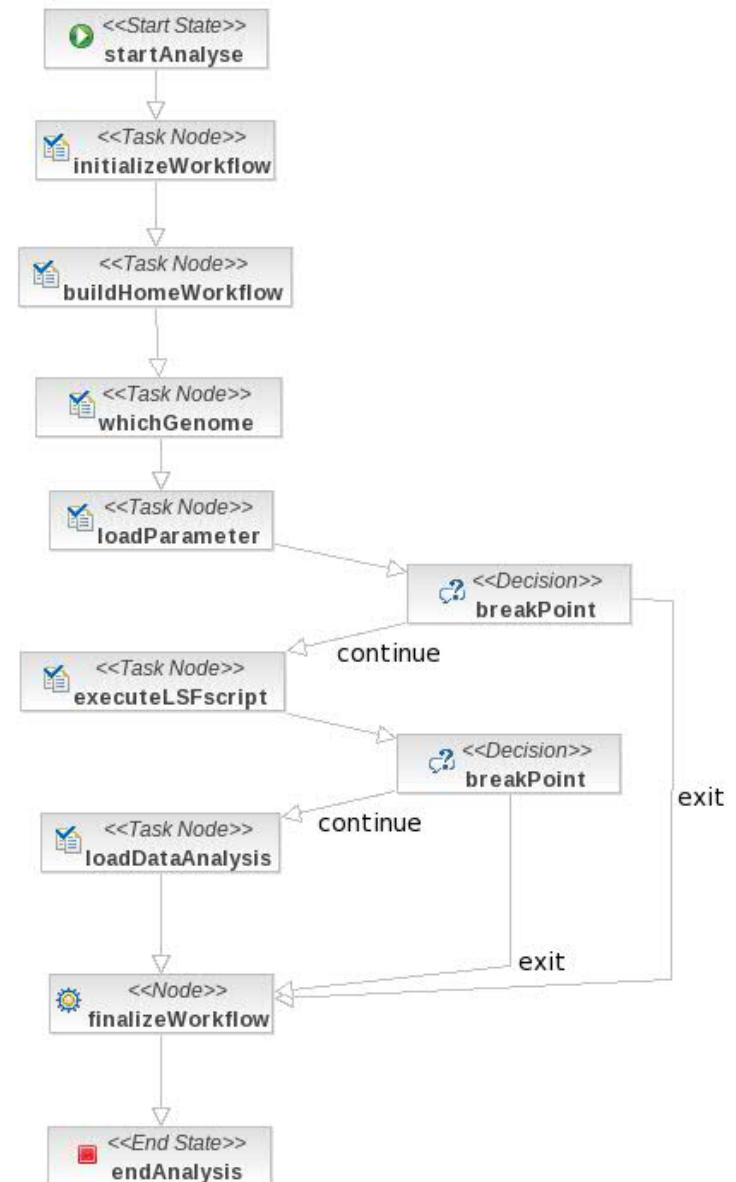
- génomes avec blast et sans synténies

Tache métier 2 : executeLSFscript

- calculs de synténies
- 2n calculs parallélisés sur un cluster
- calculs terminés ou échoués

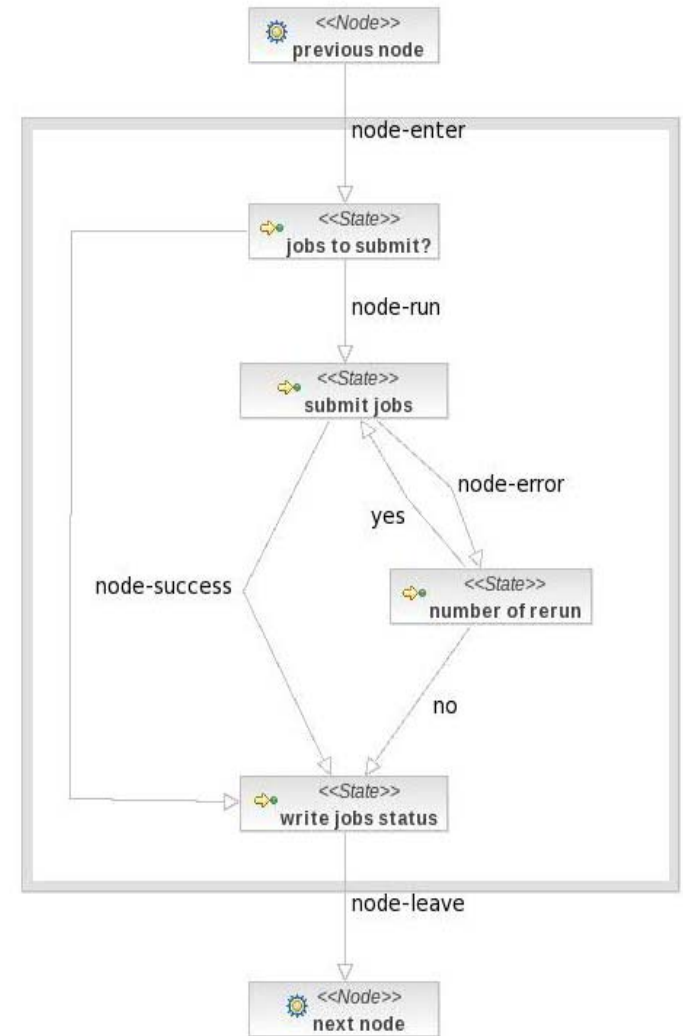
Tache métier 3 : loadDataAnalysis

- chargement des synténies terminés à la tache 2

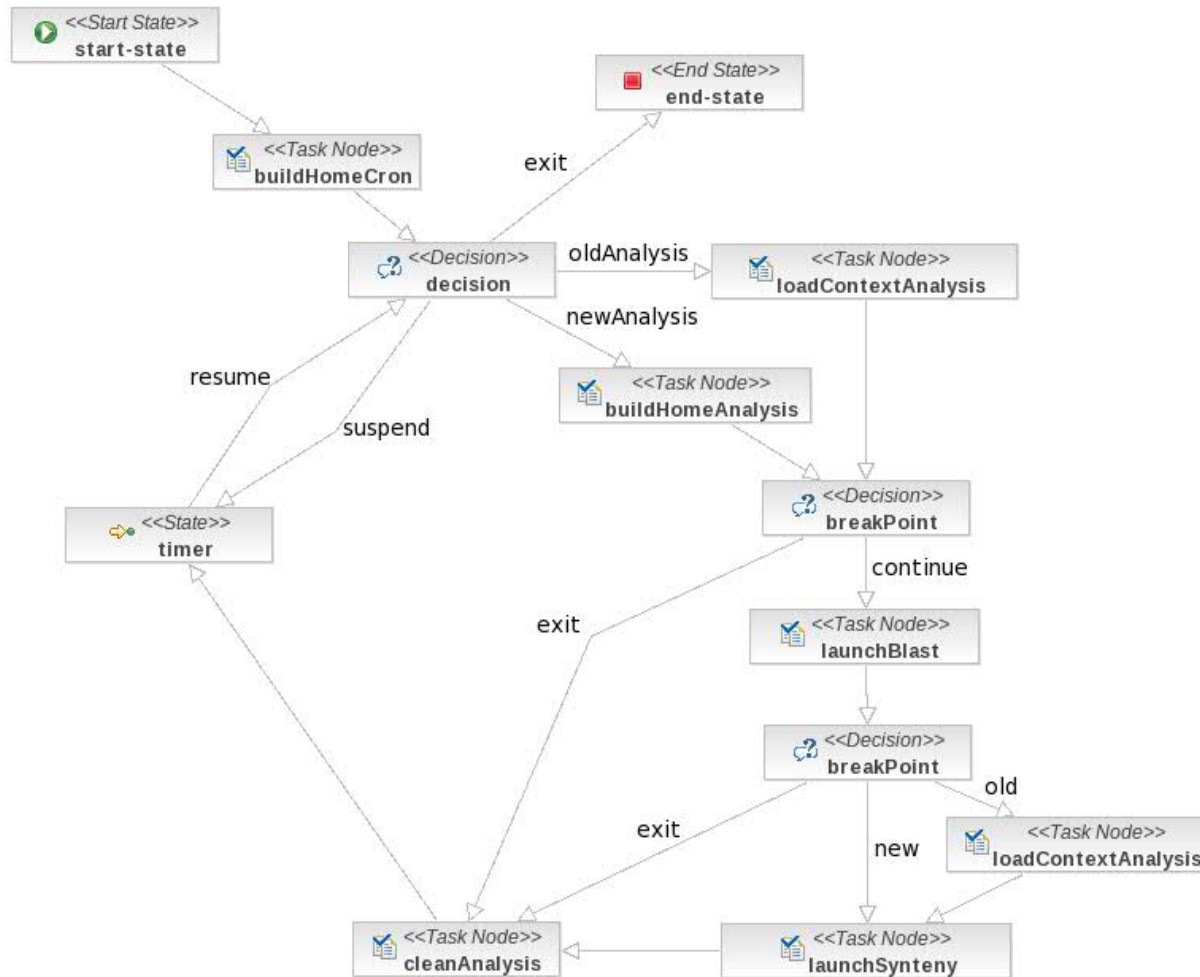


Tache métier : soumission de calculs blast

- ❑ Entrée : n génomes
- ❑ Formatage : 2^n lignes de commandes
 - o `blast -d genome_1 -i seq_1_1.fasta`
 - o `blast -d genome_n -i seq_m_n.fasta`
- ❑ Soumissions de 2^n calculs répartis en p paquets : LSF
 - o `bsub -o seq1.out -q big 'blast -d genome_1 -i seq_1_1.fasta'`
 - o `bsub -o seqN.out -q big 'blast -d genome_n -i seq_m_n.fasta'`
- ❑ Attente des résultats
- ❑ Parçage des p fichiers de résultats
- ❑ Resoumissions des calculs en erreurs
 - o nombre de resoumissions fixé
- ❑ Persistance des statuts des calculs
 - o N_e calculs en erreur
 - o N_s calculs en succès
- ❑ Sortie : $(2^n - N_e)$ alignements
- ❑ Reprise sur erreur
 - o base de donnée
 - o id workflows non terminés
 - o id taches échouées
 - o ligne de commande échouées
 - o données utilisées



CRON : Synchronisation de workflows métiers



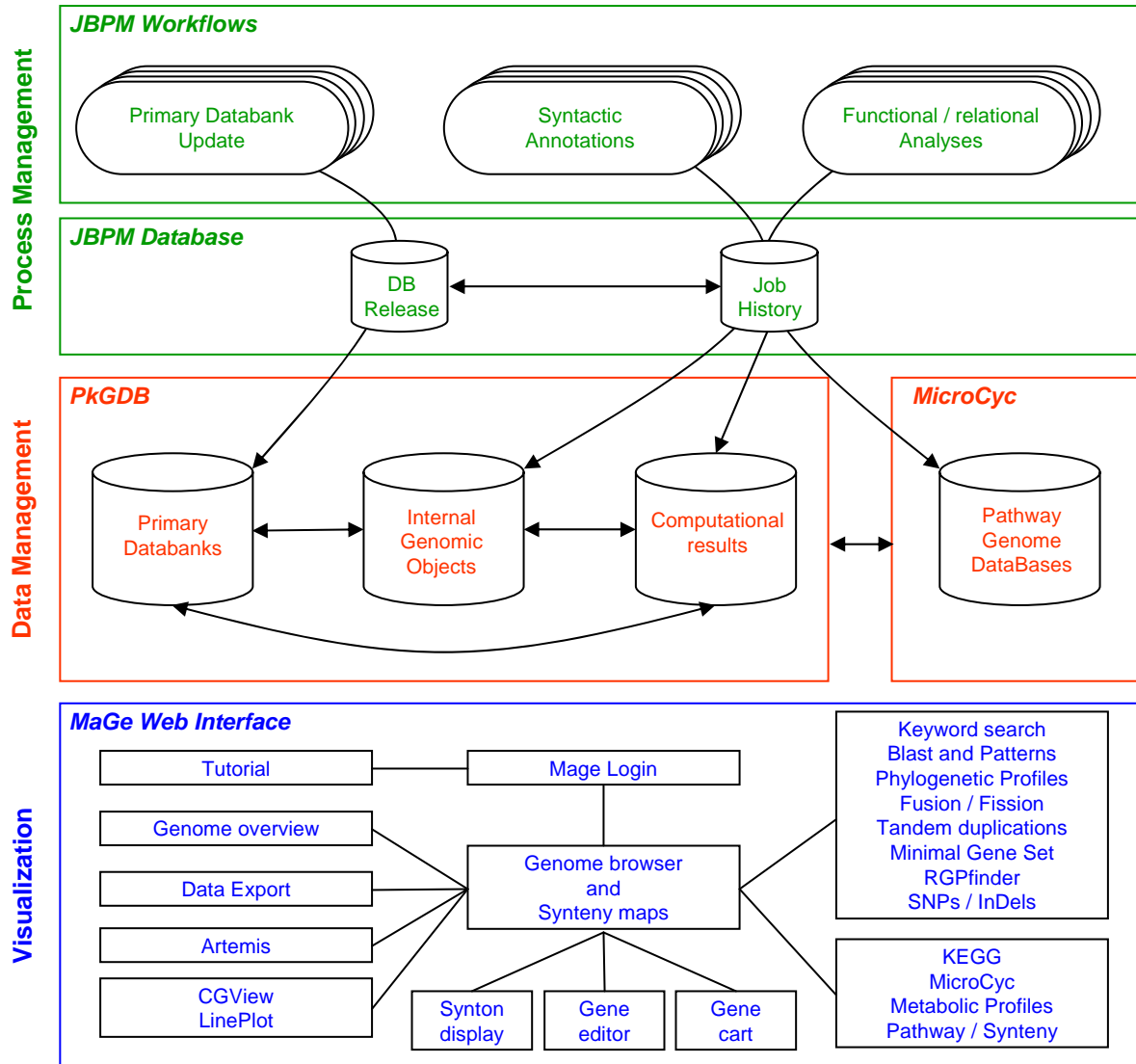
CRON

- Start
- Suspend
- Resume
- Stop
- new analysis
 - input sequences
 - update all sequences
- old analysis
 - id of process to rerun

Workflow métier

- Suspend jobs (bstop)
- Resume jobs (bresume)
- Stop jobs (bkill)
 - with load
 - without load
- Exit at breakpoint n

Architecture plateforme



Conclusion

❑ Automatisation :

- o 15 workflows de calculs de fréquence journalière
- o 9 workflows de mise à jour des données primaires

❑ Synchronisation :

- o Taches séquentielles, branchements conditionnels, états d'attente
- o Taches parallèles de chargement dans la base, backup des BD
- o Mise à jour données primaire / calculs

❑ Robustesse :

- o Reprise automatique en cas de panne d'un composant de la chaîne (base de données, cluster de calculs, panne électrique,...)
- o Mise en évidence d'erreurs et reprise des traitements sur erreurs

❑ Interaction en ligne de commande :

- o Suivi : avancement et historique des calcul, calculs en erreurs ...
- o Contrôle : priorité, reprise et arrêt de calculs, fréquence de mise à jour

❑ Volumétrie :

- o 2004-2009 : 2276387 calculs lancés et 200 GO chargés
- o depuis 6 mois : 4330795 calculs lancés et 400 GO chargés

Perspectives

❑ Court terme:

- o Intégration des méthodes d'analyses syntaxiques et relationnelles de la plateforme
- o Intégration de nouveaux workflows automatiques de mise à jour des données primaires différentielles ou complète

❑ Moyen terme:

- o Mise en place d'un service web « MaGe Light » pour la soumission anonyme de projets d'annotation de génomes
- o Mise en place et intégration de nouvelles méthodes d'analyses

❑ Long terme:

- o Utilisation de grilles pour les calculs lourds de la plateforme

Remerciements

- ❑ ANR PFTV MicroScope
- ❑ Equipe Informatique : Claude Scarpelli
 - ❑ Support informatique
 - ❑ Ludovic Fleury
- ❑ Equipe LGC : Claudine Medigue
 - ❑ Alexandra Calteau
 - ❑ Stéphane Cruveiller
 - ❑ David Vallenet
 - ❑ Zoé Rouy
 - ❑ Aurélie Lajus
 - ❑ Grégory Salvignol
 - ❑ David Roche
 - ❑ Alexander Smith
 - ❑ Damien Mornico