

# La donnée d'imagerie microscopique en biologie structurale à l'IGBMC

Jonathan Michalon, 31 janvier 2017

Journées scientifiques « Equip@Meso »

# Contexte

- IGBMC (Institut de Génétique et de Biologie Moléculaire et Cellulaire), 4 départements, près de 750 personnes
- département de biologie structurale intégrative
- problématiques informatiques de stockage et de calcul
  - principal microscope : 2 To de données brutes par jour
  - autres systèmes
- Instruct (européen), FRISBI (French Infrastructure For Integrated Structural Biology)

# Infrastructure

## Locale

- 1,2 Po de capacité brute disque
- 5 serveurs GPGPU
- 4 serveurs « haute mémoire » avec 1To de RAM et 64 cœurs de calcul chacun

## Méso-centre

- 1,3 Po de capacité disque en ligne
- 6500 cœurs sur plus de 550 nœuds
- mutualisation IGBMC (648 cœurs, 8 Go de RAM par cœur) :
  - 8 machines achetées par l'IGBMC, 256 cœurs
  - 14 machines dimensionnées IGBMC (via CPER), 392 cœurs

# Organisation et flux de données

# Acquisition

- Cristallographie (rayons X), microscopie photonique et électronique
  - Enjeux : microscopie électronique (jeux de données pouvant atteindre 2To)
- Réalisée par les systèmes d'acquisition sur une passerelle dédiée :
  - tampon pendant l'acquisition (jusqu'à 10 To)
  - isole le système d'acquisition (parfois encore Windows XP !)
  - formulaire web simplifié pour transfert du dossier d'acquisition vers une zone de travail sur le stockage central
  - interconnexion 10 gbps

## Accès pendant le traitement

Le stockage central est découpé en « zones », qui sont des dossiers soit personnels (aux chercheurs), soit pour la plateforme, soit spécifiques à certains systèmes.

Il est directement accessible sur toutes les machines de traitement locales, que ce soit les pré-traitements, le GPGPU, etc.

Ces serveurs ont aussi du stockage local, pour les fichiers temporaires principalement.

Les mêmes zones sont aussi réexportées en Samba/CIFS, ce qui permet d'accéder à l'ensemble des données depuis n'importe quel poste de travail de l'institut (Linux / Windows).

# Méso-centre

Directement depuis le stockage central, chacun peut à tout moment transférer des données via le lien 10 gbps depuis et vers le méso-centre. Une petite interface web permet d'initier la copie, qui est ensuite effectuée en parallèle à l'aide de globus-url-copy.

Toutes les tâches de gros calcul parallèle y sont effectuées, directement par les chercheurs, lesquels disposent d'un compte personnel sur le méso-centre.

# Archivage

- Données brutes (pas systématiquement)
- Une fois le projet terminé (le dataset a été exploité, il faut garder pour la publication ou autre)

Sur cassette (LTO6) :

- dossiers complets
- une description est ajoutée au début du tar
- Via un petit utilitaire en shell (interface zenity) :
  - monte la zone
  - remplit la description
  - ajoute les données à la bande

La personne est alors responsable de sa bande, et peut à tout moment reverser le contenu sur les zones de la même manière.



# Pipeline « Titan Krios »

# Topologie

Le microscope est contrôlé par un ordinateur (Windows 7) gérant le fonctionnement général et une des caméra (Falcon 2).

Une nouvelle caméra à comptage d'électron, la Gatan K2 Summit est contrôlée via son propre ordinateur (Windows 7 : 1,2 To SSD).

La passerelle (Linux : 10 To) est accessible depuis ces deux PC d'acquisition, par montage Samba/CIFS.

Une machine de pré-traitement dotée de 4 GPU GTX 1080 sert à toutes les tâches automatiques, avec un lien direct sur la passerelle et le stockage central monté à demeure.

# Arrivée des données

Les données sont générées par l'un ou l'autre des PC gérant les caméras. L'un peut directement acquérir sur le montage, l'autre utilise les SSD comme premier tampon.

Une fois les données sur la passerelle Linux, nous pouvons commencer à automatiser.

- script shell, plus facile pour gérer des processus
- inotifywait, pour réagir aux événements du système de fichier

[http ://cbi-dev.igbmc.fr/cbi/titan-autoprocess](http://cbi-dev.igbmc.fr/cbi/titan-autoprocess)

# Tâches exécutées

Est exécuté à la volée :

- Empilement des micrographes dans le cas d'une acquisition en mode "movie"
- Alignement

Nécessite de déterminer un moment de déclenchement :

- Tentative de sélection des particules automatique
- Premier modèle initial

Cette seconde partie qui est largement dépendante de l'échantillon observé, n'est pas encore fonctionnelle mais utilisera la machine de pré-traitement.

# Détermination de structure

Pour obtenir cette structure le workflow sera . . .

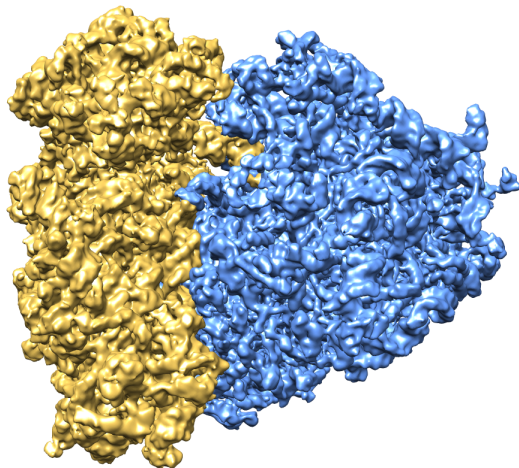


FIGURE 1: Structure 3D d'un ribosome

## Étapes pour obtenir l'image de ribosome précédente :

- Acquisition : microscope Titan Krios
  - plus de 1To de micrographes sur la passerelle-tampon
  - transfert parallélisé via formulaire web vers le stockage centralisé
- Picking : sélection des particules
  - eman2 avec visualisation à distance, directement sur le stockage
- Classification 2D : groupement des particules en classes
  - eman2 sur le méso-centre, transferts aller-retour en 10gbps
- Reconstruction 3D : affinage itératif de la structure
  - relion sur R920

# Conclusion

- Infrastructure informatique forte ;
- répondant aux divers problématiques inhérentes aux travaux des chercheurs ;
- se recentrer sur les traitements, pipeline, développements en étant assuré d'une base de travail solide et évolutive.



Merci pour votre attention.<sup>1</sup>

---

1. Support de présentation réalisé en Markdown, utilisant la classe Beamer de  $\text{\LaTeX}$  via Pandoc.