

http://get.genotoul.fr



http://bioinfo.genotoul.fr

Gestion des données de l'HiSeq 2000

Plateformes GenoToul Bioinformatique et GeT-PlaGe Ch. Klopp / G. Salin







Contexte NGS

- Collaboration entre PF GeT-PlaGe et Bioinfo depuis 2007
 - PF GeT-PlaGe : production de données (454, HiSeq2000)
 - PF Bioinfo : moyens d'analyse et de mise à disposition
- CPER/IBiSA communs pour les investissements dans le NGS
- Acteurs du projet France Génomique (Investissements d'avenir)



L'HiSeq2000, côté PlaGe

Quelques dates

Validation de l'application Installation « séquençage d'ADNg »

Sept. 2010

Déc. 2010

Avr. 2011

Oct. 2011

Validation technique par Illumina

Triplement du débit

Projets pilotes ADNg

Projets pilotes RNA-Seq



L'HiSeq2000, côté PlaGe

Quelques chiffres

36 flowcells en 2x101b depuis le 23/11/2010 + 8 HS

- 288 lanes autant d'analyses nG6
 - 227 ADNg
 - 61 ARN
- 28 projets différents
- 366 librairies
- ~54To de données brutes transférées
- 19 flowcells planifiées



L'HiSeq2000, côté PlaGe

Notre expérience

- -- : Problèmes réactifs + machine (8 flowcells repassées)
- : Dépendance vis-à-vis de l'infrastructure informatique
- + : Beaucoup de problèmes logiciels corrigés par Illumina
- ++: Un débit plus élevé qu'annoncé (+40%)
 - 130 M de reads / Lane attendus selon spécifications Illumina en Paired-ends
 - 185 M de reads / Lane en moyenne
- ++: Retour positif des utilisateurs



Une infrastructure en conséquence

- Réseau à 1Gb/s avec connexion directe sur un NAS
- 14 To pour le transfert des données brutes
- 8 To pour la sauvegarde des données analysées
- Cluster avec 49 nœuds de calcul (> 400 coeurs)
- nG6, un système d'information dédié au NGS



Historique de nG6

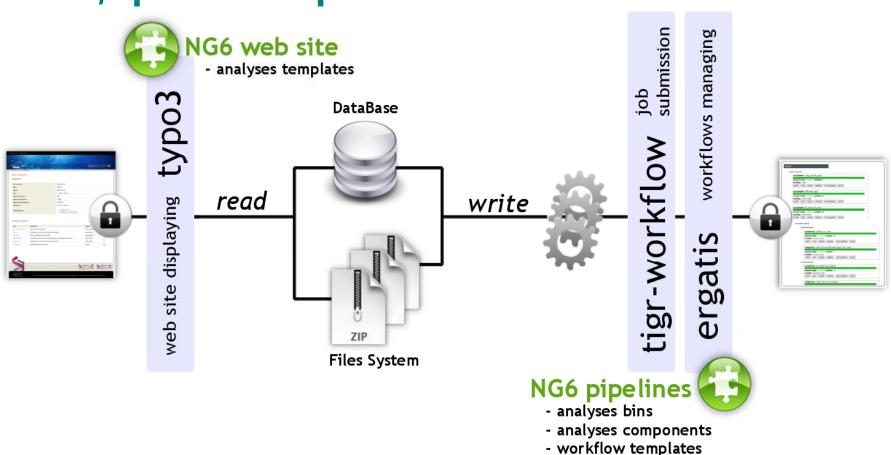
- Projet commun Bioinfo/PlaGe depuis 2008
- Pour la gestion des données Roche 454 :
 - Traçabilité des données et analyses associés
 - Portail sécurisé d'accès aux données pour les utilisateurs
 - Sauvegarde des données analysées
- Intégration des premières données HiSeq fin 2010

Sites web

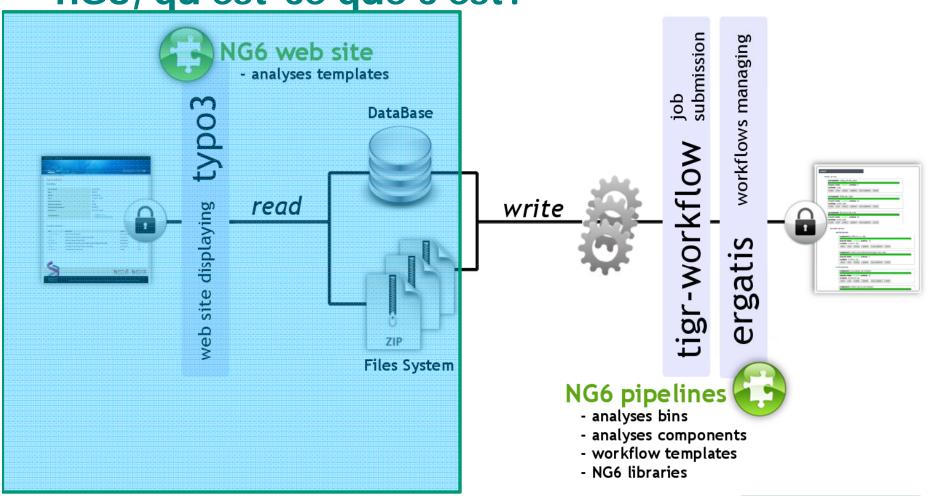
- Démo : http://ng6.toulouse.inra.fr
- Forge: http://mulcyber.toulouse.inra.fr/projects/ng6/



- NG6 libraries

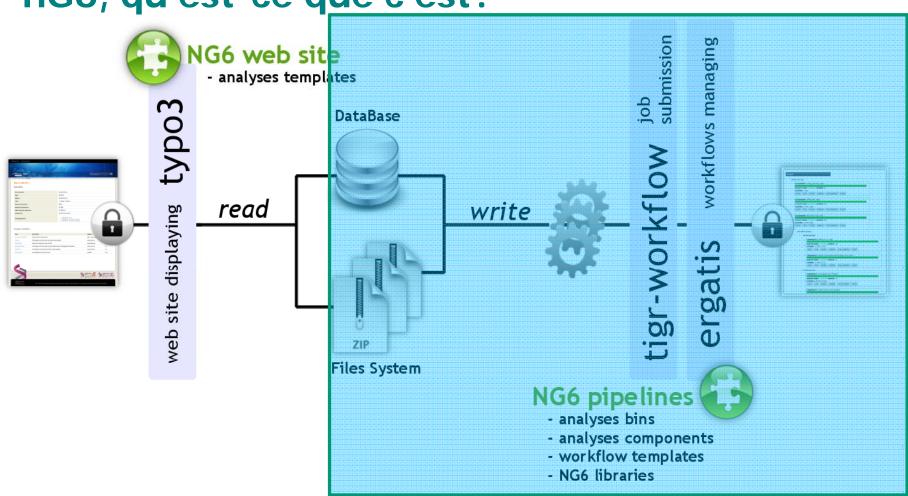








L'HiSeq2000 et nG6





Frontend, site web utilisateur

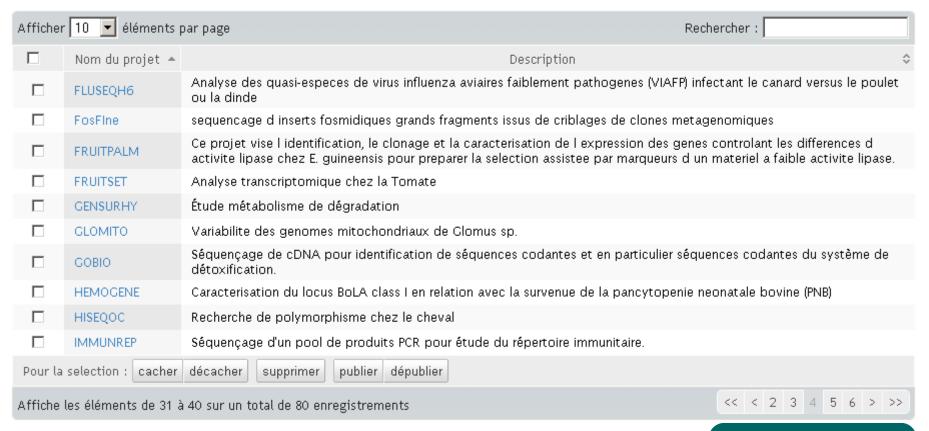
 Accès aux résultats d'analyse par projet ou par run + module de recherche



Liste des projets auxquels vous avez accès :

Vous avez accès à 80 projets.

L'ensemble des données brutes et des résultats d'analyses occupent 4.65 Tb d'espace disque pour l'ensemble des projets.



FR4486041658-

BESNEBUCK

CARTOSEQ 23-08-11

Bovin

ADNa

1/8

Flowcell

A - lane

163 355

506

ech

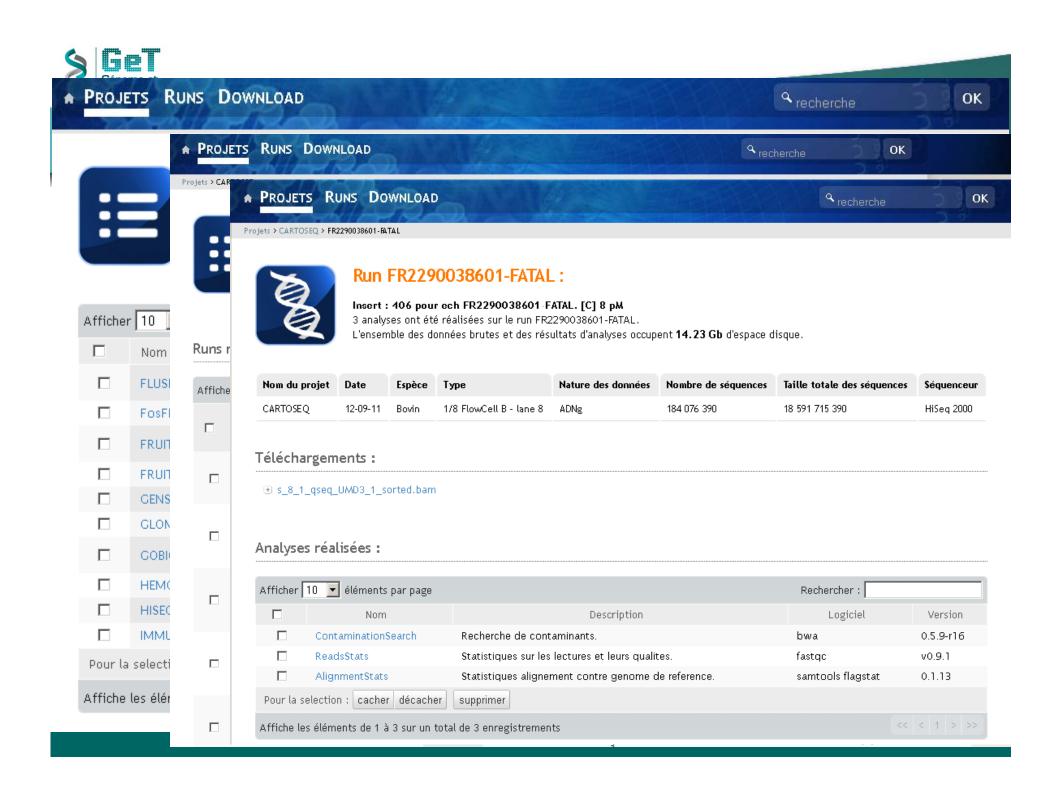
FR4486041658-

BESNBUCK. [C] 7

HiSeq 2000

16 498

906 106





Frontend, site web utilisateur

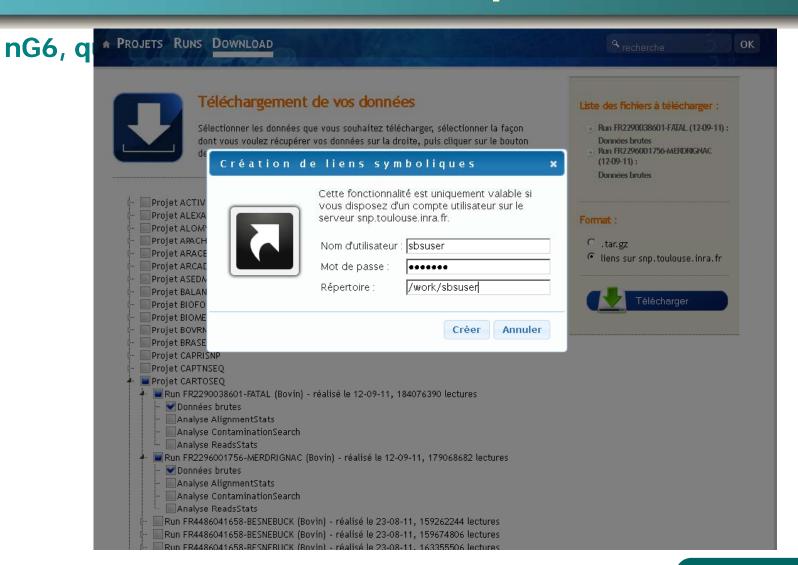
- Analyses proposées en routine pour des données 454 :
 - Statistiques de base (longueur, qualité, N...)
 - Recherche de contaminants (blastall 2.2.25)
 - Nettoyage sur les biais connus du 454 (Pyrocleaner 1.2)
 - Assemblage (newbler 2.6)
- Analyses proposées en routine pour des données HiSeq2000:
 - Statistiques de base (FastQC 0.9.1)
 - Recherche de contaminants (bwa 0.5.9)
 - Alignement contre génome de ref (bwa 0.5.9)



Frontend, site web utilisateur

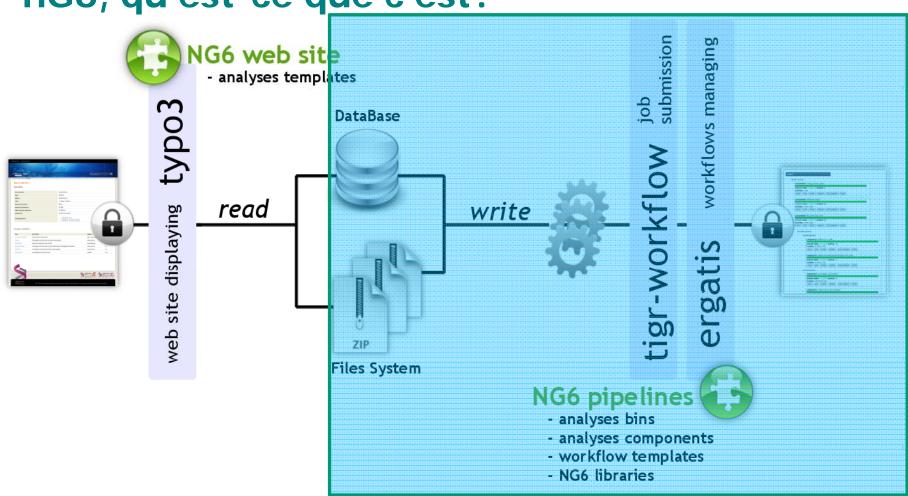
 Possibilité de télécharger les données « brutes» (sff, bam, fastq.gz) après archivage ou de créer un lien symbolique sur serveur Bioinfo pour analyse sur le cluster







L'HiSeq2000 et nG6





Backend, Ergatis

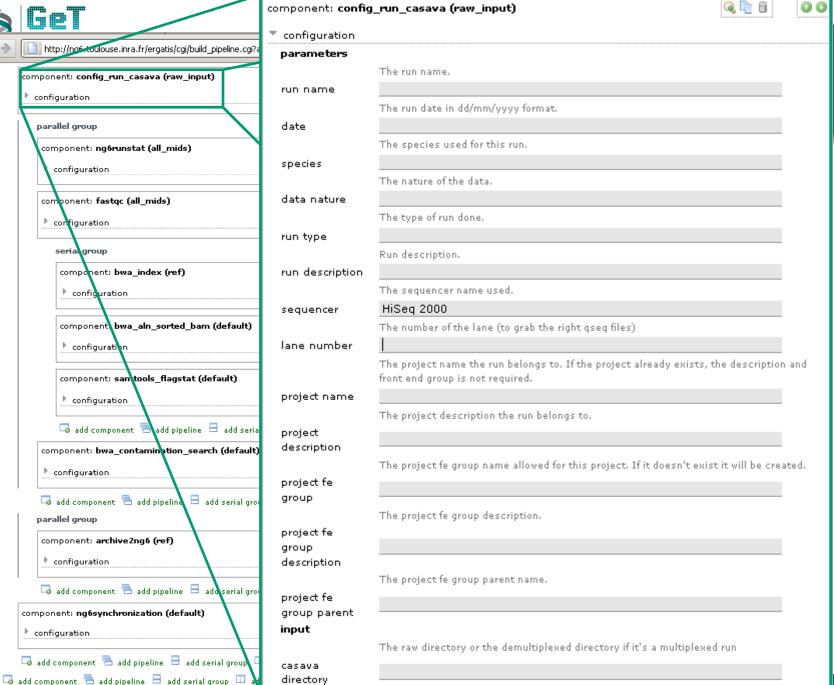
Pipelines préconfigurés suivant

- le séquenceur
 - 454
 - HiSeq2000, GA
- le type d'échantillon
 - ADNg / ARN
 - Multiplexé ou non
 - Métagénomique
- Le format d'entrée des données
 - sff
 - fasta
 - fastq.gz
 - qseq





2000, côté Bioinfo







L'HiSeq2000, les Hommes

nG6, qui sont les petites mains?

PF Bioinfo:

- C. Klopp: Responsable technique
- J. Mariette : l'expert, développeur principal
- D. Laborie : administrateur système

PF GeT-*:

- J. Lluch, E. Robe, S. Ruzafa, E. Lhuillier, N. Marsaud, D. Esquerre : équipes 454 et HiSeq, exécutent les analyses qualité sur ergatis et rédigent les comptes-rendus
- O. Bouchez : responsable NGS GeT-PlaGe, valide les analyses qualité
- F. Escudié et moi : gestion des données, suivi des analyses qualité et développement, interface PF Bioinfo



L'HiSeq2000, les Hommes

nG6, qui sont les petites mains?

Groupes de travail sur différentes thématiques visant à créer des pipelines à intégrer à nG6 (RNASeq, variations, miRNA, ChIP-SEQ, méthylation...)

https://mulcyber.toulouse.inra.fr/projects/ngspipelines > Onglet MediaWiki



L'HiSeq2000, futur PlaGe

Et dans l'année à venir...

PF GeT-PlaGe:

- Triplement du débit de l'HiSeq2000 dès maintenant
- Doublement du débit annoncé courant d'année prochaine
- 1 voire 2 HiSeq supplémentaires (130To de données à sauvegarder)



L'HiSeq2000, futur Bioinfo

Et dans l'année à venir...

PF Bioinfo:

- Passage en production des pipelines CASAVA1.8
- Continuer le travail d'optimisation de l'espace dans la gestion des données
- Nouveaux pipelines d'analyse + nouvelles fonctionnalités nG6
- Formations NGS (polymorphisme SNP / RNA-Seq)
- Mise en place nouvelle infrastructure en janvier 2012 (CPER)
 - Calcul + + +
 - Stockage+



Remerciements / Questions



• Toute l'équipe PF Bioinfo



L'EIC Toulouse



• Les équipes de recherche pour leurs retours

Toute l'équipe PlaGe...et plus particulièrement les biologistes moléculaires qui mettent la main au clavier







































