# Next Generation Sequencing

# and

# Bioinformatics Analysis Pipelines

**Adam Ameur**

**National Genomics Infrastructure**

**SciLifeLab Uppsala**
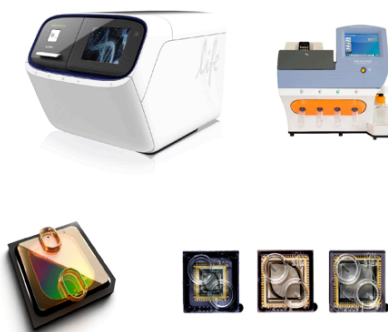
**adam.ameur@igp.uu.se**

# Today's lecture

- Management of NGS data at NGI/SciLifeLab

- Examples of analysis pipelines:

    - Human exome/genome sequencing
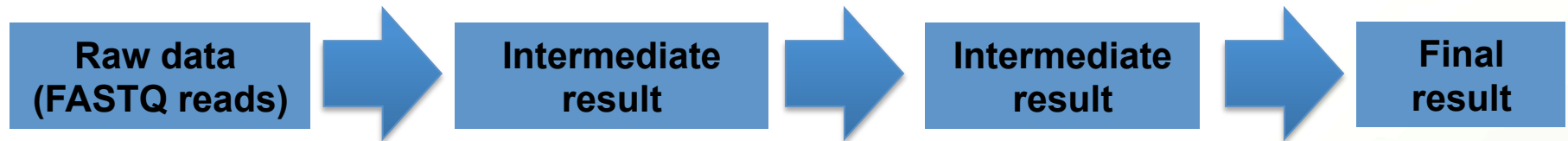
    - Assembly using long reads
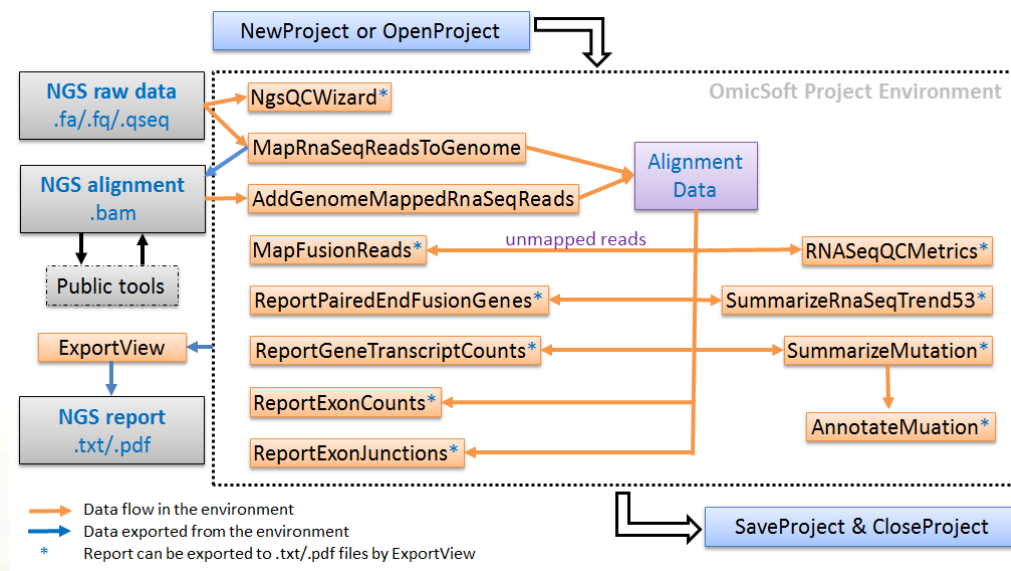
    - Clinical routine sequencing

# What is an analysis pipeline?

- Basically just a number of steps to analyze data
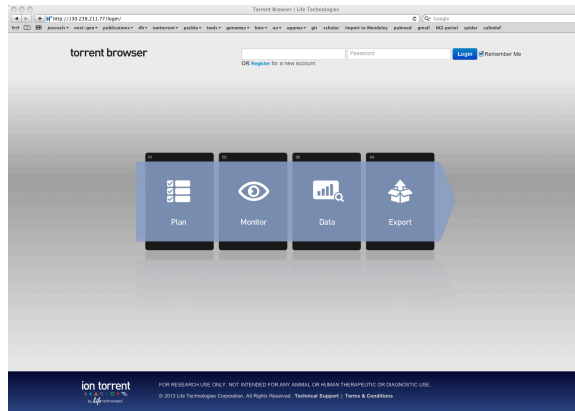


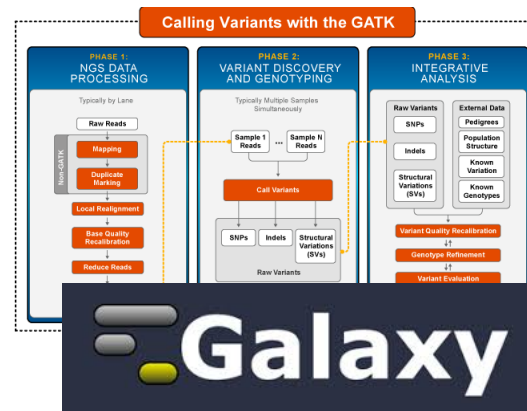- Pipelines can be simple or very complex…
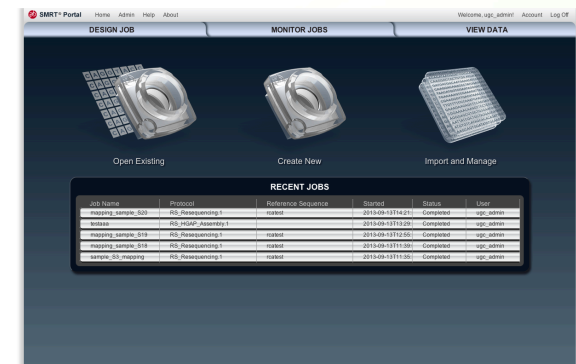
# Some analysis pipelines for NGS data

**Ion Torrent**
Torrent Suite Software



**Illumina**
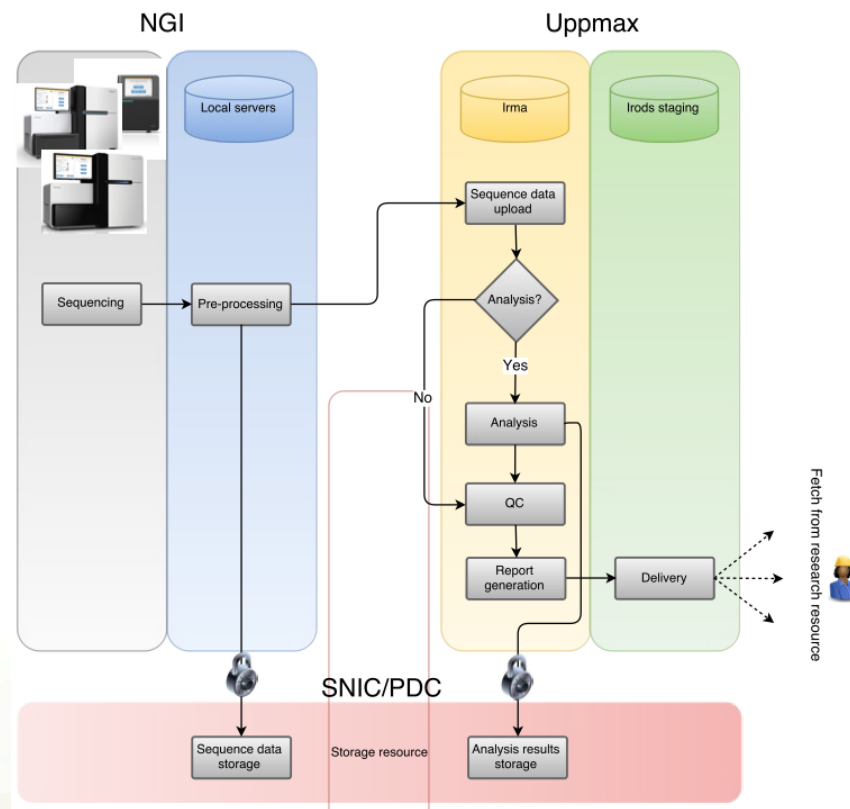GATK, Galaxy,…



**PacBio**
SMRT analysis portal



- Enables variant calling, de novo assembly, RNA expression analyses, …

- Many other tools exists, also from commercial vendors

# Data processing at NGI

- Raw data from is processed in automated pipelines
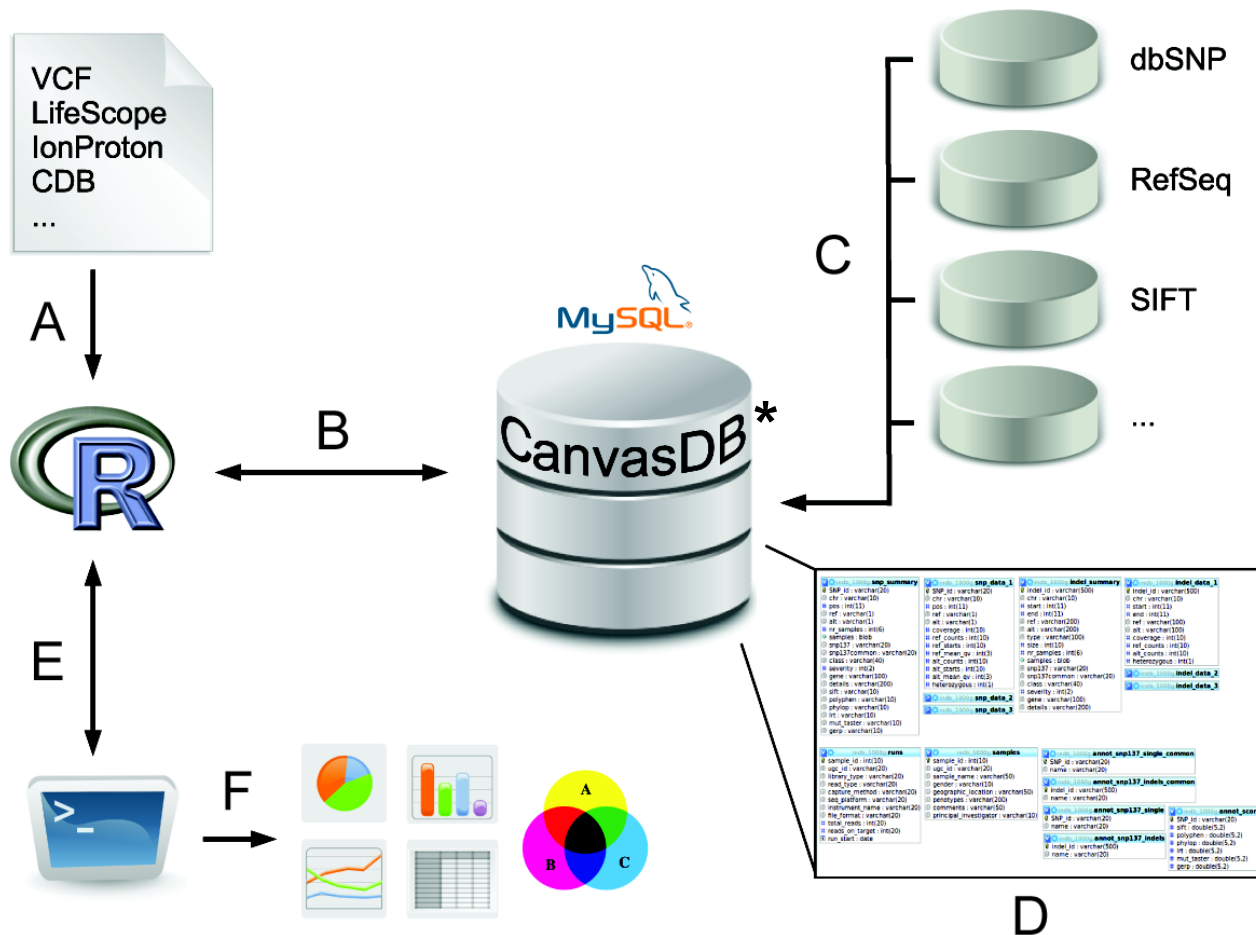- Delivered to user accounts at UPPNEX

# In-house development of pipelines

- In some cases NGI develops own pipelines

- But only when we see a need for a specific analysis

Some examples follows:

**I.  Building a local variant database (WES/WGS)**

**II.  Assembly of genomes using long reads**

**III. Clinical sequencing – Leukemia Diagnostics**

# Example I:
# Computational infrastructure for exome-seq data

# Background: exome-seq

- ## Main application of exome-seq
  - Find disease causing mutations in humans

- ## Advantages
  - Allows investigate all protein coding sequences
  - Possible to detect both SNPs and small indels
  - Low cost (compared to WGS)
  - Possible to multiplex several exomes in one run
  - Standardized work flow for data analysis

- ## Disadvantage
  - All genetic variants outside of exons are missed (~98%)

# Exome-seq throughput

- ## We are producing a lot of exome-seq data
  - 4-6 exomes/day on Ion Proton
  - In each exome we detect
    - Over 50,000 SNPs
    - About 2000 small indels
  => Over 1 million variants/run!
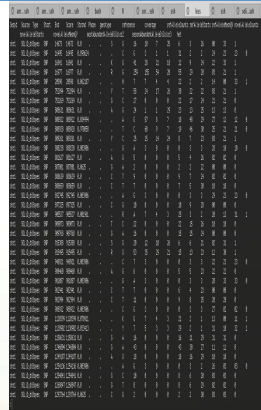    - In plain text files

# How to analyze this?

- Traditional analysis - A lot of filtering!
  - Typical filters
    - Focus on rare SNPs (not present in dbSNP)
    - Remove FPs (by filtering against other exomes)
    - Effect on protein: non-synonymous, stop-gain etc
    - Heterozygous/homozygous
  - This analysis can be automated (more or less)

Start:
All identified SNPs

Result:
A few candidate causative SNP(s)!

# Why is this not optimal?

- Drawbacks
  - Work on one sample at time
    - Difficult to compare between samples
  - Takes time to re-run analysis
    - When using different parameters
  - No standardized storage of detected SNPs/indels
    - Difficult to handle 100s of samples


- Better solution
  - A database oriented system
    - Both for data storage and filtering analyses

# Analysis: In-house variant database



*CANdidate Variant Analysis System and Data Base*

*Ameur et al., Database Journal, 2014*

# CanvasDB - Filtering

# CanvasDB - Filtering speed

- Rapid variant filtering, also for large databases

# A recent exome-seq project

- Hearing loss: 2 affected brothers
  - Likely a rare, recessive disease
  => Shared homozygous SNPs/indels

- Sequencing strategy
  - TargetSeq exome capture
  - One sample per PI chip

| nr reads | (% mapped) | 76M-89M (97%) |
|---|---|---|
| mapped reads | (% on target) | 73M-88M (83%) |
| SNPs | (% in dbSNP) | 85k-93k (93%) |
| Indels | (% in dbSNP) | 5k-6k (48%) |

# Filtering analysis

- *CanvasDB* filtering for a variant that is…
  - rare
    - at most in 1% of ~700 exomes
  - shared
    - found in both brothers
  - homozygous
    - in brothers, but in no other samples
  - deleterious
    - non-synonymous, frameshift, stop-gain, splicing, etc..



D recessive variant

```
> cand <- filterRecessive(c("up_001_1","up_001_2"), outfile="cand.txt")
Total time for filtering: 27.012s
```

# Filtering results

- ## Homozygous candidates
  - 2 SNPs
    - stop-gain in *STRC*
    - non-synonymous in *PCNT*
  - 0 indels



D recessive variant

- ## Compound heterozygous candidates (lower priority)
  - in 15 genes

| sample_name | class | chr | pos | ref | alt | snp137 | gene | ref_counts | alt_counts |
|---|---|---|---|---|---|---|---|---|---|
| up_001_1 | stopgain | chr15 | 43896948 | G | A | rs144948296 | STRC | 3 | 58 |
| up_001_2 | stopgain | chr15 | 43896948 | G | A | rs144948296 | STRC | 5 | 55 |
| up_001_1 | nonsynonymous | chr21 | 47808772 | G | A | rs35044802 | PCNT | 0 | 21 |
| up_001_2 | nonsynonymous | chr21 | 47808772 | G | A | rs35044802 | PCNT | 1 | 14 |

=> Filtering is fast and gives a short candidate list!

# STRC - a candidate gene

## STRC

From Wikipedia, the free encyclopedia

**Stereocilin** is a protein that in humans is encoded by the *STRC* gene.[1][2][3]

This gene encodes a protein that is associated with the hair bundle of the sensory hair cells in the inner ear. The hair bundle is composed of stiff microvilli called stereocilia and is involved with mechanoreception of sound waves. This gene is part of a tandem duplication on chromosome 15; the second copy is a pseudogene. Mutations in this gene cause autosomal recessive non-syndromic deafness.[3]

=> Stop-gain in STRC is likely to cause hearing loss!

# IGV visualization: Stop gain in STRC

# STRC, validation by Sanger

- Sanger validation



- Does not seem to be homozygous..
  - Explanation: difficult to sequence STRC by Sanger
    - Pseudo-gene with very high similarity

- New validation showed mutation is homozygous!!

# CanvasDB – some success stories

**Solved cases, exome-seq - Niklas Dahl/Joakim Klar**

| | |
|---|---|
| *Neuromuscular disorder* | *NMD11* |
| *Artrogryfosis* | *SKD36* |
| *Lipodystrophy* | *ACR1* |
| *Achondroplasia* | *ACD2* |
| *Ectodermal dysplasia* | *ED21* |
| *Achondroplasia* | *ACD9* |
| *Ectodermal dysplasia* | *ED1* |
| *Arythroderma* | *AV1* |
| *Ichthyosis* | *SD12* |
| *Muscular dystrophy* | *DMD7* |
| *Neuromuscular disorder* | *NMD8* |
| *Welanders myopathy (D)* | *W* |
| *Skeletal dysplasia* | *SKD21* |
| *Visceral myopathy (D)* | *D:5156* |
| *Ataxia telangiectasia* | *MR67* |
| *Exostosis* | *SKD13* |
| *Alopecia* | *AP43* |
| *Epidermolysis bullosa* | *SD14* |
| *Hearing loss* | *D:9652* |

*Success rate >80% for recent Proton projects!*

# CanvasDB - Availability

- CanvasDB system now freely available on GitHub!

## Installation of the CanvasDB system

This section describes how to download and install CanvasDB on your local computer. Make sure that MySQL, R and ANNOVAR are running on your computer before starting the installation.

Step 1. Download code from github

```
$ git clone https://github.com/UppsalaGenomeCenter/CanvasDB.git
$ cd CanvasDB
```

Step 2. Set the current path to 'rootDir' in canvasDB.R

# Next Step: Whole Genome Sequencing

- New instruments at SciLifeLab for human WGS…



**Capacity of HiSeq X Ten: 320 whole human genomes/week!!!**

- More work on pipelines and databases needed!!!

# Analysis of WGS data @ SciLifeLab

We have a working group for WGS at SciLifeLab!

## wgs-toolbox@scilifelab.se

Contacts with Genomics England initiated for analyses

# The SciLifeLab Human WGS Initiative

- WGS of patient cohorts (n=10,000 ind/year)

- Genetic Variant Database for the Swedish Population (n=1000)

# The Swedish Genetic Variant Project

A. Identify a cohort that reflects the genetic structure of the Swedish population

B. Generate WGS data using short- and long-read MPS technologies

C. Establish a user-friendly database to make information available to the research community (association analyses) and clinical genetics laboratories.

# The Swedish Twin registry

- Inclusion based on twinning

- Distribution like population density

- General population-prevalence of any disease

- 10,000 individuals have been analysed with SNP arrays

- Identify 1,000 individuals based on genetic structure  and diversity across Sweden.

# Principal components of European samples from 1,000 genomes project and 10,000 Swedish samples

# European samples from 1,000 genomes project and 1,000 selected Swedish samples

# WGS of Swedish control cohort

**Step 1:** 30X Illumina data of the 1,000 individuals (Q2 2016)

**Step 2:** PacBio *de novo* sequencing of 3 individuals (Q2 2016)

Ref genome individuals



Principal components generated in individuals selected for WGS

**Step 3:** Sequencing of HLA and other clinically relevant loci

# Example II:

# Assembly of genomes using Pacific Biosciences

# Genome assembly using NGS

- Short-read *de novo* assembly by NGS
  - Requires mate-pair sequences
    - Ideally with different insert sizes
  - Complicated analysis
    - Assembly, scaffolding, finishing
    - Maybe even some manual steps
  => Rather expensive and time consuming


- Long reads really makes a difference!!
  - We can assemble genomes using PacBio data only!

# HGAP *de novo* assembly

- HGAP uses both long and shorter reads

**Short reads** →

**Long reads (seeds)** →

# PacBio – Throughput & read lengths

- >10kb average read lengths! (run from April 2014)



- ~ 1 Gb of sequence from one PacBio SMRT cell

# PacBio assembly analysis

- Simple -- just click a button!!

# PacBio assembly, example result

- Example: Complete assembly of a bacterial genome



>70% A/T rich bacteria
1 contig: 2,024,078 bp

# PacBio assembly – recent developments

- Also larger genomes can be assembled by PacBio..



**2013**

**2014**

Spinach[5]
**1 Gb**
Contig N50
**531 kb**

Drosophila[4]
**170 Mb**
Contig N50
**4.5 Mb**

Arabidopsis[3]
**120 Mb**
Contig N50
**7.1 Mb**

Yeast[2]
**12 Mb**
Resolve most
**chromosomes**

Bacteria[1]
**1-10 Mb**
**Finished**
**Genomes**

Human[6]
**3.2 Gb**
Contig N50
**4.4 Mb**
**Max=44 Mb**
(Assembly
powered by
**Google Cloud**)

# Assembly of large genomes

- A computational challenge!!

## Data Release: ~54x Long-Read Coverage for PacBio-only De Novo Human Genome Assembly

We are pleased to make publicly available a new shotgun sequence dataset of long PacBio® reads from a human DNA sample. We previously released sequence data using Single Molecule, Real-Time (SMRT®) Sequencing of ~10x coverage of this sample, sufficient for reference-based detection of structural variation. Today we expand on that release with additional data that increases the total sequencing coverage to ~54x. This long-read data has enabled the generation of the first *de novo* human genome assembly from PacBio-only sequence reads.
Download the 54x long-read coverage dataset.

**405,000 CPUh used on Google Cloud!**

# De novo WGS of Swedish cohort

Establish Swedish reference genome sequences by *de novo* assembly of long-read PacBio data (+10X Genomics?)

Ref genome individuals



Principal components generated in individuals selected for WGS

# First Swedish PacBio WGS

- 20 kb library

- 157 SMRT cells

- 140 Gb data (~45X)

- FALCON assembly

|  | First PacBio Assembly |
|---|---|
| # of contigs (>=0 bp) | 7708 |
| # of contigs (>=1000 bp) | 7653 |
| Total length (>=0 bp) | 2844 Mb |
| Total length (>=1000 bp) | 2844 Mb |
| No of contigs | 7692 |
| Largest contig | 19.5 Mb |
| Total contig length | 2844 Mb |
| N50 | 4.35 Mb |
| N75 | 1.97 Mb |

# Why clinical WGS using long reads?

**Precision medicine requires high-quality genome sequences!**

- Resolving repetitive and complex regions

- Annotation of unknown genomic regions

- Haplotype phasing

- …

*Jim Lupski: "The Goal Is De Novo Assembly in the Clinic"*



Jim Lupski, Baylor

# Example III:

# Clinical sequencing for Leukemia Treatment

# Chronic Myeloid Leukemia

- BCR-ABL1 fusion protein – a CML drug target



The BCR-ABL1 fusion protein can acquire resistance mutations following drug treatment

www.cambridgemedicine.org/article/doi/10.7244/cmj-1355057881

# BCR-ABL1 workflow – PacBio Sequencing

From sample to results: < 1 week

1 sample/SMRT cell

*Cavelier et al., BMC Cancer, 2015*

# BCR-ABL1 mutations at diagnosis

PacBio sequencing generates ~10 000X coverage!



Sample from time of diagnosis:

# BCR-ABL1 mutations in follow-up sample



BCR

ABL1

Sample 6 months later

Mutations acquired in fusion transcript.
Might require treatment with alternative drug.

# BCR-ABL1 dilution series results

- Mutations down to 1% detected!

# Mutations mapped to protein structure

# BCR-ABL1 - Multiple isoforms in one individual!

# Clinical Diagnosis of BCR-ABL1 mutations

**Clinical Genetics**



- Collection of samples
- Seq library preparation

**Sequencing Facility**



- SMRT sequencing
- Mutational analysis

**IT developers**



- Web server for results

- Ongoing routine service, 0-4 samples/week
- Over 120 patient samples run so far
- 100% of Sanger-positive mutations recovered
- Developments: Detect low frequency mutations down to 0.1%

# Web system for result sharing

| Details | Sample ID | Run ID | Unresolved (count) | Unknown (count) | M244V | Q252H | Y253H | E255K | E255V | K262N | D276G | T277A | L298V | T315I | T315A | M351T | F359V | L387M | E450G | E453G | E459G | M472I | E499E | Date |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 91 | R12021 | cba_011_2 | | | | | | | | | | | | | | | | | | | | | | 2015-09-07 |
| 92 | R12023 | cba_011_3 | | | | | | | | | | | | | | | | | | | | | | 2015-09-07 |
| 93 | R12026 | cba_011_4 | | | | | | | | | | | | | | | | | | | | | | 2015-09-07 |
| 94 | R12091 | cba_012_1 | | | | | | | | | | | | | | | | | | | | | | 2015-09-17 |
| 95 | R12092 | cba_012_2 | | | | | | | | | | | | | | | | | | | | | | 2015-09-17 |
| 96 | R12093 | cba_012_3 | | | | | | | | | | | | | | | | | | | | | | 2015-09-17 |
| 97 | R12095 | cba_012_4 | | | | 45.2 | | | | | | | | | | | | | | | | | | 2015-09-17 |
| 98 | R12124 | cba_013_1 | | | | | | | | | | | | | | | | | | | | | | 2015-09-23 |
| 99 | R12125 | cba_013_2 | | | | | | | | | | | | | | | | | | | | | | 2015-09-23 |
| 100 | R12123 | cba_013_3 | | | | | | | | | | | | | | | | | | | | | | 2015-09-23 |
| 101 | R12126 | cba_014_1 | | | | | | | | | | | | | | | | | | | | | | 2015-09-29 |
| 102 | R12149 | cba_014_2 | | | | | | | | | | | | | | | | | | | | | | 2015-09-29 |
| 103 | R12165 | cba_015_1 | | | | | | | | | | | | | | | | | | | | | | 2015-10-07 |
| 104 | R12143 | cba_016_1 | | | | | | | | | | | | | | | | | | | | | | 2015-11-04 |
| 105 | R12281 | cba_017_1 | | | | | | | | | | | | | | | | | | | | | | 2015-11-12 |
| 106 | R12282 | cba_017_2 | | | | | | | | | | | | | | | | | | | | | | 2015-11-12 |
| 107 | R12222 | cba_018_1 | | | | | | | | | | | | | | | | | | | | | | 2015-11-18 |
| 108 | R12291 | cba_019_1 | | | | | | | | | | | | | | | | | | | | | | 2015-12-02 |
| 109 | R12355 | cba_019_2 | | | | | | | | | | | | | | | | | | | | | | 2015-12-02 |
| 110 | R12200 | cba_020_1 | | | | | | | | | | | | | | | | | | | | | | 2015-12-16 |

96    **Sample 97**    98      New Search

| Sample ID | Run ID | Date |
|---|---|---|
| R12095 | cba_012_4 | 2015-09-17 |

**Downloads:**

Results          Sequence          Details          Clonal distribution

| mutation | sequence | wt_reads | mut_reads | other_reads | freq | detection |
|---|---|---|---|---|---|---|
| M351T | CACTCAGATCTCGTCAGCCA[T/C]GGAGTACCTGGAGAAGAAAA | 16176 | 19154 | 3 | 0.542 | positive |
| Q252H | CACAAGCTGGGCGGGGGCCA[G/C]TACGGGGAGGTGTACGAGGG | 12918 | 10686 | 16 | 0.452 | positive |
| K262N | GTGTACGAGGCGCGTGTGGAA[G/T]AAATACAGCCTGACGGTGGC | 25673 | 7035 | 16 | 0.215 | positive |
| M244V | TGGAACGCACGGACATCACC[A/G]TGAAGCACAAGCTGGGCGGG | 32901 | 33 | 2 | 0.001 | negative |
| K247K | GGACATCACCATGAAGCACA[A/G]GCTGGGCGGGGGCCAGTACG | 27186 | 32 | 9 | 0.001 | negative |
| L248V | ACATCACCATGAAGCACAAG[C/G]TGGGCGGGGGCCAGTACGGG | 27214 | 3 | 17 | 0 | negative |
| G250E | CATGAAGCACAAGCTGGGCG[G/A]GGGCCAGTACGGGGAGGTGT | 23601 | 8 | 3 | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0.001 | negative |
| | | | | | 0.002 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| | | | | | 0 | negative |
| V299L | AGATCAAACACCCTAACCTG[G/T]TGCAGCTCCTTGGGGTCTGC | 30283 | 2 | 9 | 0 | negative |
| F311V | TCTGCACCCGGGAGCCCCCG[T/G]TCTATATCATCACTGAGTTC | 27076 | 1 | 35 | 0 | negative |

| | Frequency | Reads |
|---|---|---|
| M351T | 49.9 % | 9268 |
| Q252H | 23.8 % | 4418 |
| Q252H    K262N | 17.4 % | 3245 |
| | 8.69 % | 1613 |

Karolinska Institutet   KTH VETENSKAP OCH KONST   Stockholms universitet   UPPSALA UNIVERSITET

SciLifeLab

# Ion Torrent – Ongoing developments

## Ion S5 XL system

## Ion Proton PII chip (EA)

# PacBio - Ongoing developments

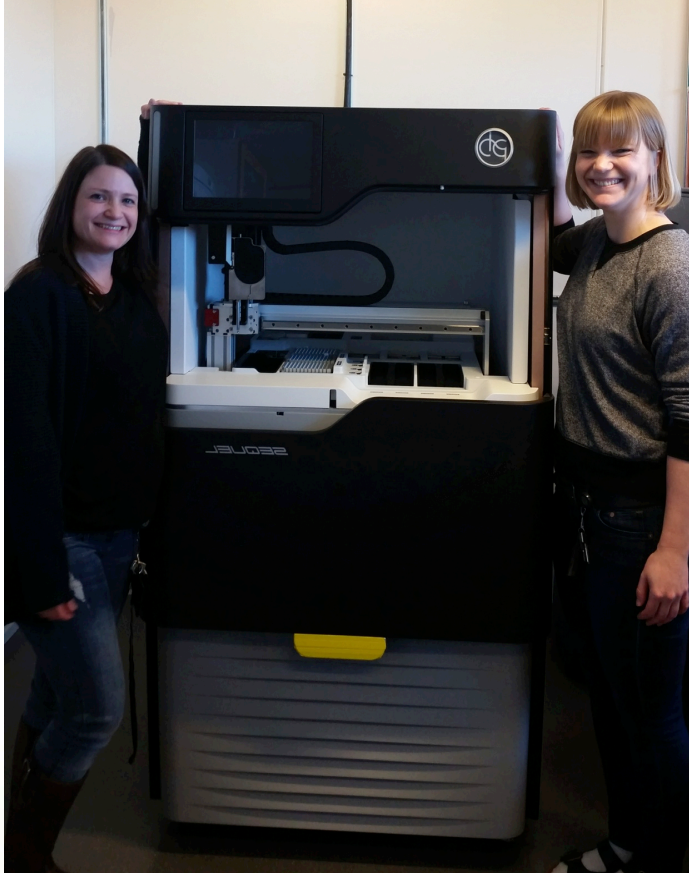**Sequel** - New instrument with higher throughput!



**7x more data per SMRT cell!**

Installation at NGI during 2016

# Who does the sequencing?



**Ulf Gyllensten**
Platform director

**Inger Jonasson**
Facility manager

**Olga Vinnere Pettersson**
Project coordinator

**Adam Ameur**
Bioinformatician, NGS

**Ignas Bunikis**
Bioinformatician, NGS

**Christian Tellgren-Roth**
Bioinformatician, NGS

**Susana Häggqvist**
Research engineer
NGS

**Ida Höijer**
Research engineer
NGS

**Cecilia Lindau**
Research engineer
NGS

**Maria Schenström**
Research engineer
NGS

**Magdalena Andersson**
Research engineer
NGS

**Ulrika Broström**
Research engineer
NGS

**Nina Williams**
Research engineer
NGS

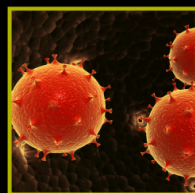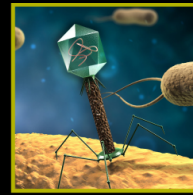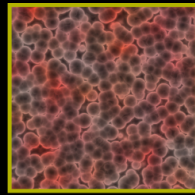**Carolina Ilbäck**
Research engineer
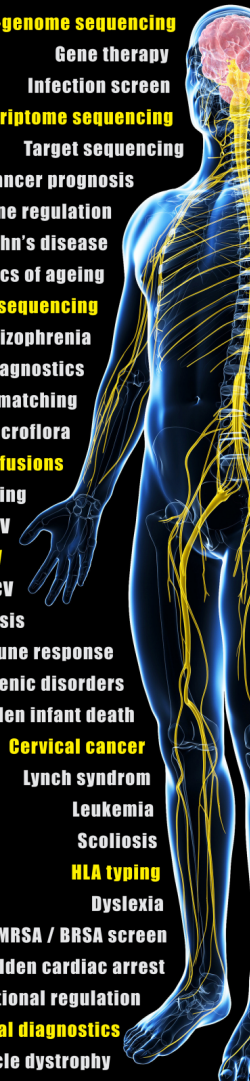NGS

**Anna Petri**
Research engineer
Sequencing Service

**Anne-Christine Lindström**
Research engineer
Sequencing Service

# What we sequence at NGI /

SciLifeLab

THANK YOU

Diabetes
Alzheimer's disease
Whole-genome sequencing
Gene therapy
Infection screen
Whole-transcriptome sequencing
Target sequencing
Cancer prognosis
Gene regulation
Crohn's disease
Genomics of ageing
Exome sequencing
Schizophrenia
Cancer diagnostics
Organ donor matching
Gut microflora
Gene fusions
RNA editing
HIV
HPV
HCV
Scoliosis
Immune response
Monogenic disorders
Sudden infant death
Cervical cancer
Lynch syndrom
Leukemia
Scoliosis
HLA typing
Dyslexia
MRSA / BRSA screen
Sudden cardiac arrest
Transcriptional regulation
Prenatal diagnostics
Muscle dystrophy
Individualised cancer therapy
and much more...