

Next Generation Sequencing and Bioinformatics Analysis Pipelines

Adam Ameer
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 & whole genome sequencing
 - Assembly using long reads
 - Clinical routine sequencing

illumina®



ThermoFisher
SCIENTIFIC
life
ion torrent
5 4 3 2 1



PACIFIC
BIOSCIENCES®

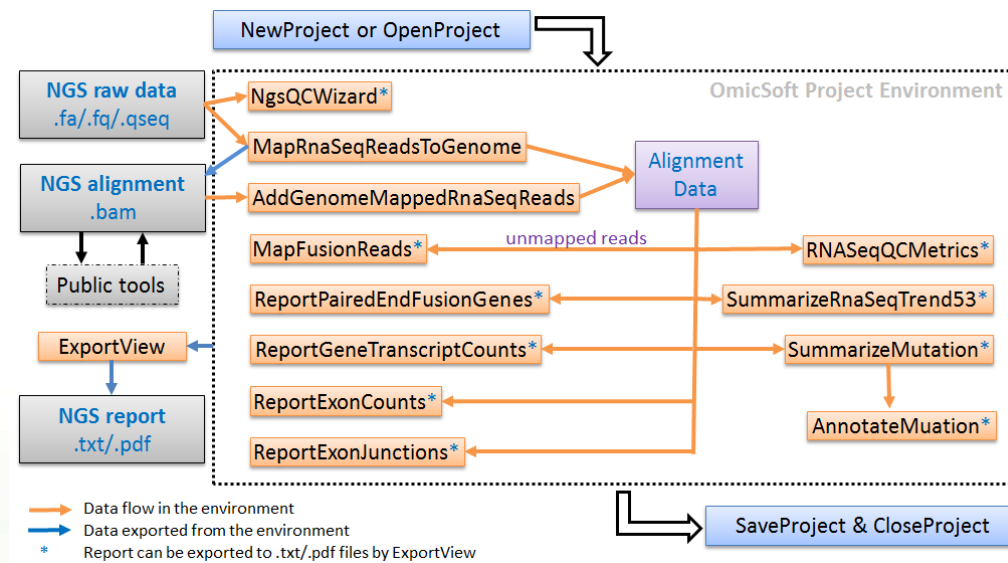


What is an analysis pipeline?

- Basically just a number of steps to analyze data

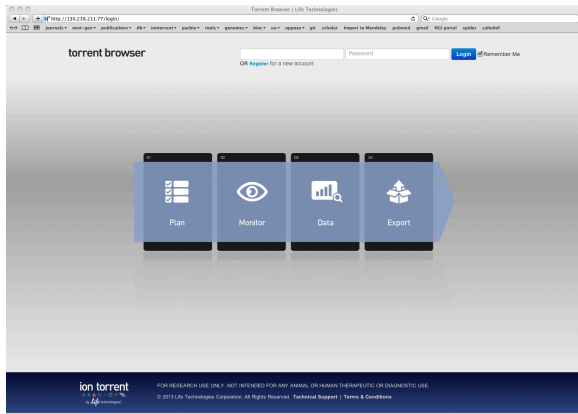


- Pipelines can be simple or very complex...

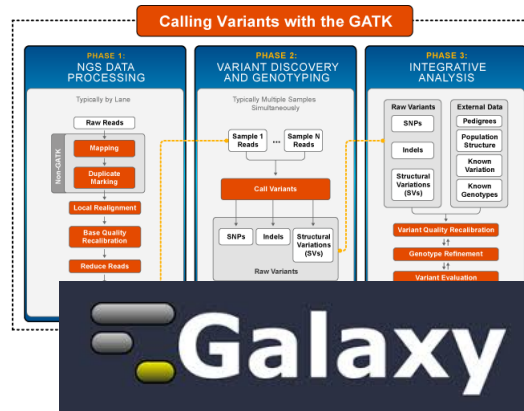


Some analysis pipelines for NGS data

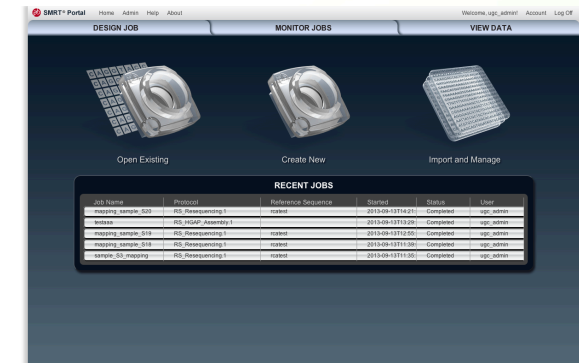
Ion Torrent Torrent Suite Software



Illumina GATK, Galaxy,...



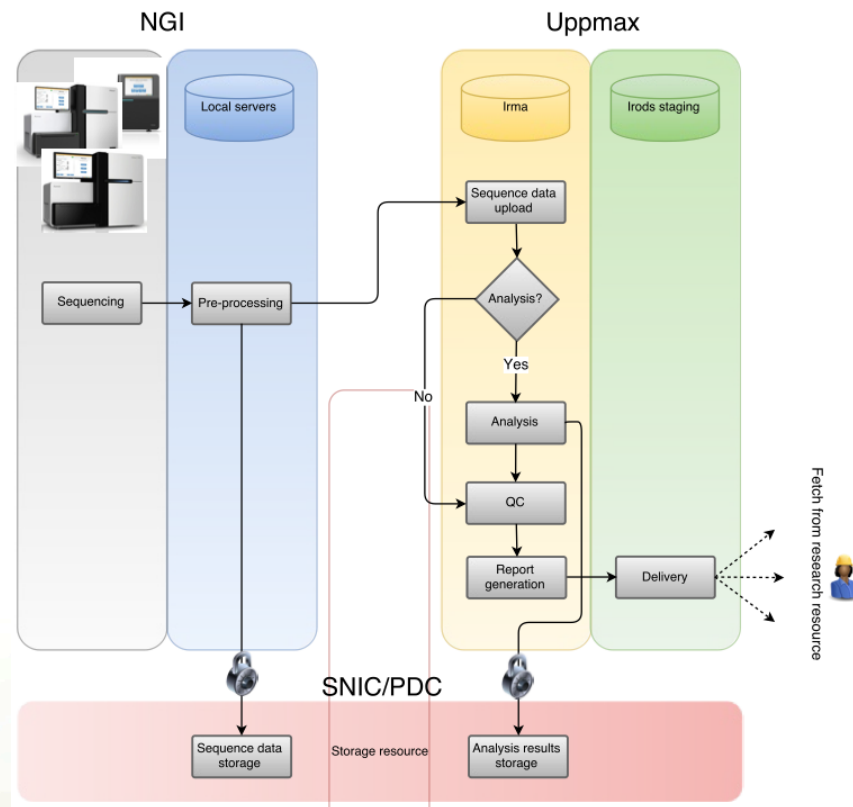
PacBio SMRT analysis portal



- 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



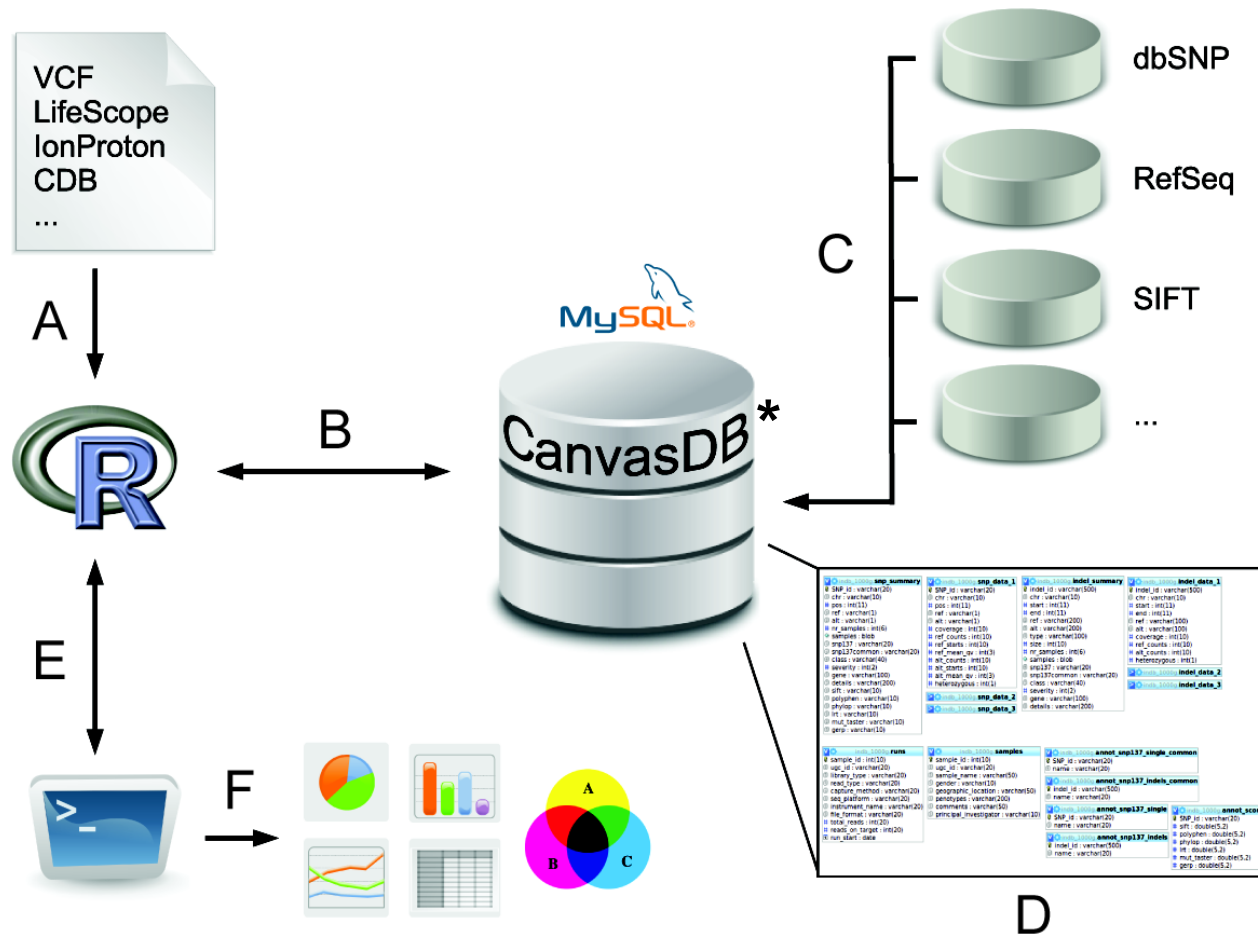
In-house development of pipelines

- 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 databases (WES/WGS)**
- II. Assembly of genomes using long reads**
- III. Clinical sequencing – Leukemia Diagnostics**

Example I: Computational infrastructure for exome-seq data



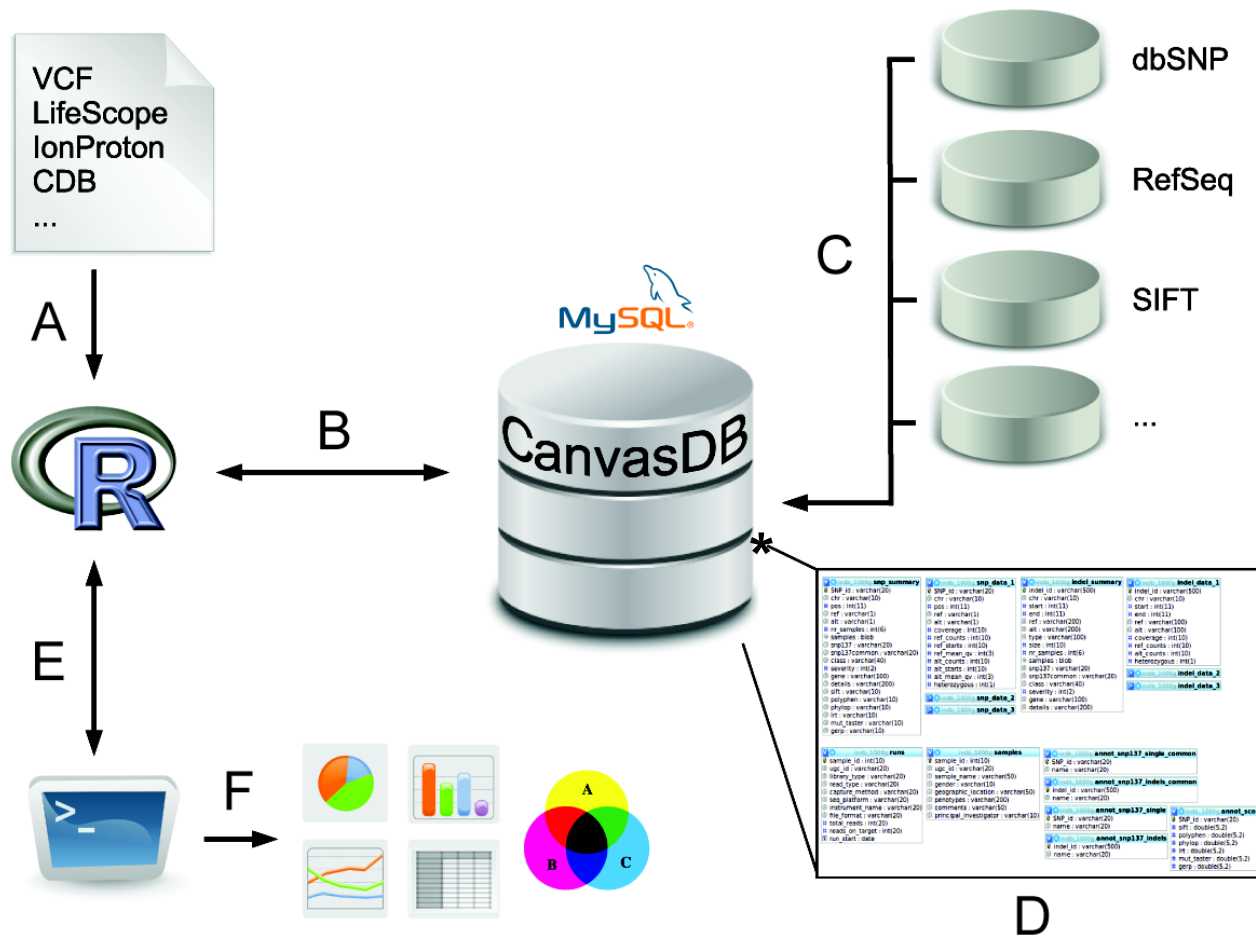
Background: exome-seq

- 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%)

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

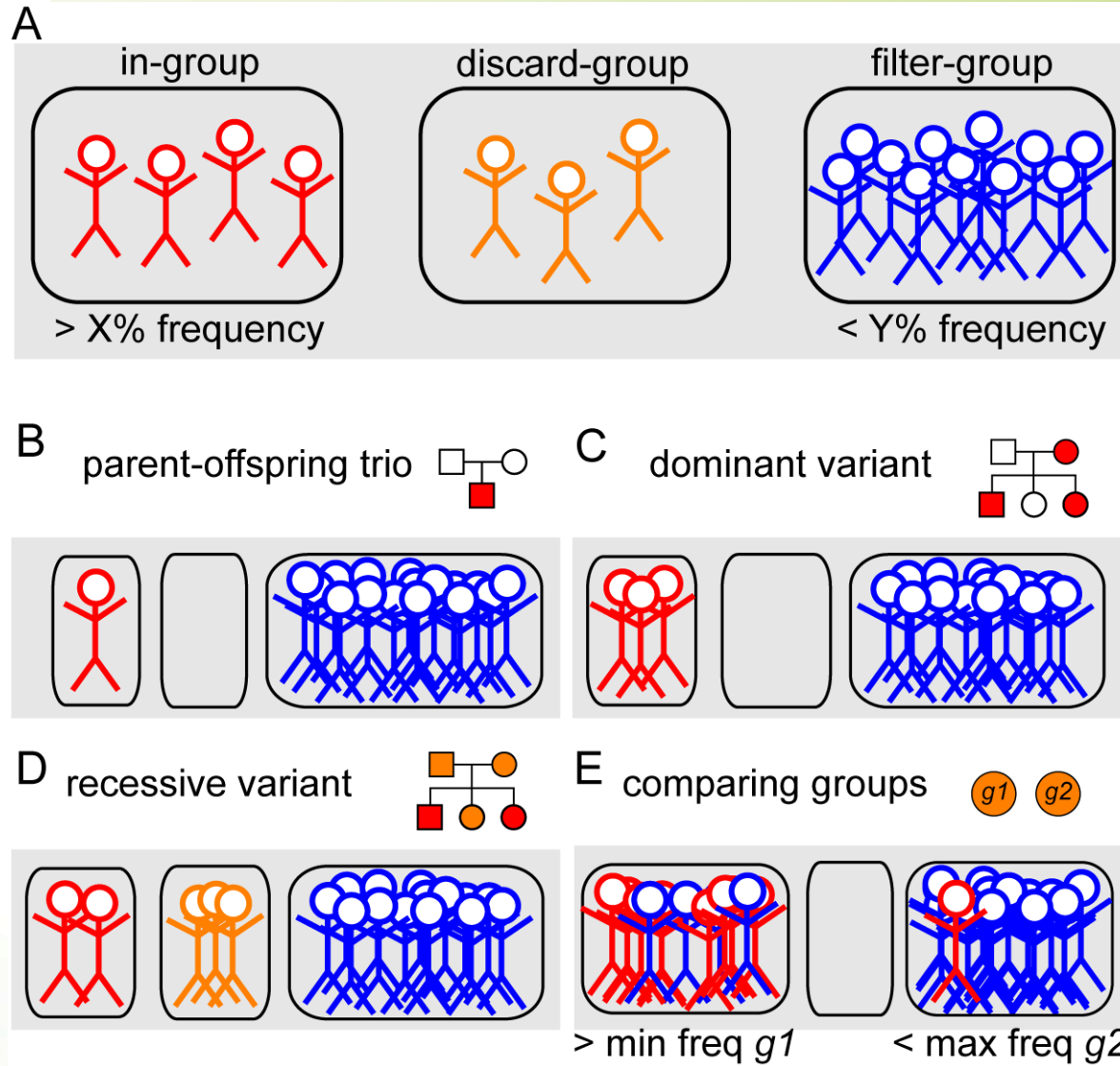
Analysis: In-house variant database



***CANdicate Variant Analysis System and Data Base**

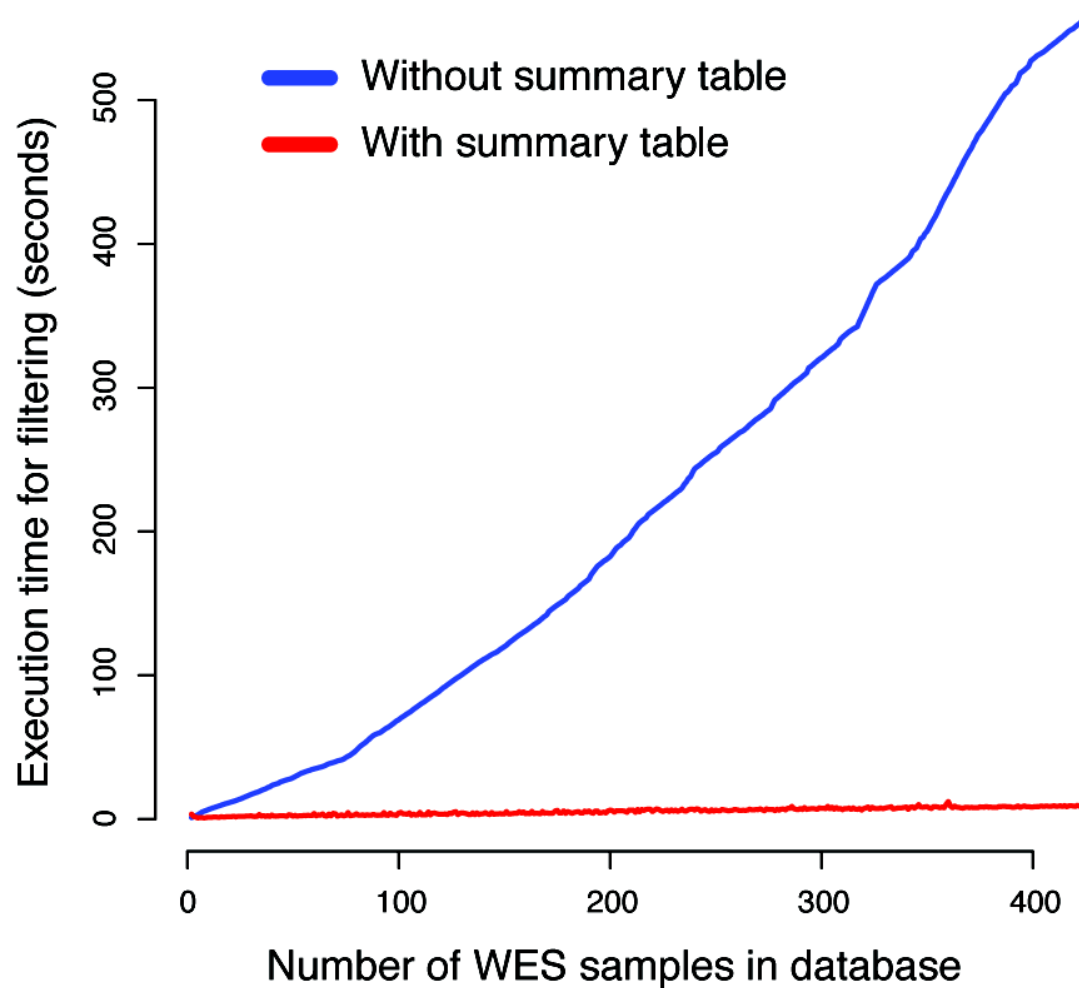
Ameur et al., Database Journal, 2014

CanvasDB - Filtering



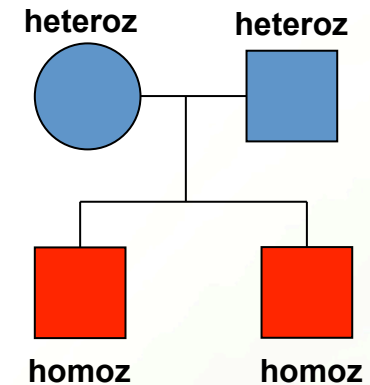
CanvasDB - Filtering speed

- 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

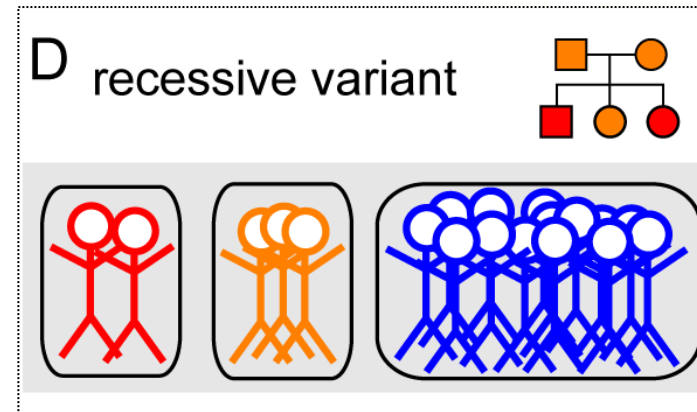
- found in both brothers

- homozygous

- in brothers, but in no other samples

- deleterious

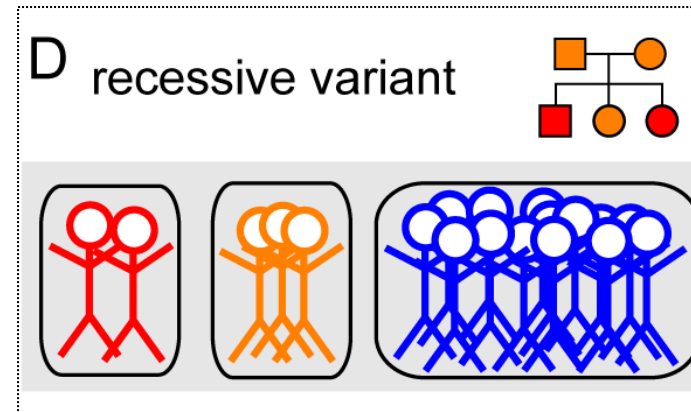
- non-synonymous, frameshift, stop-gain, splicing, etc..



```
> 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
- 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

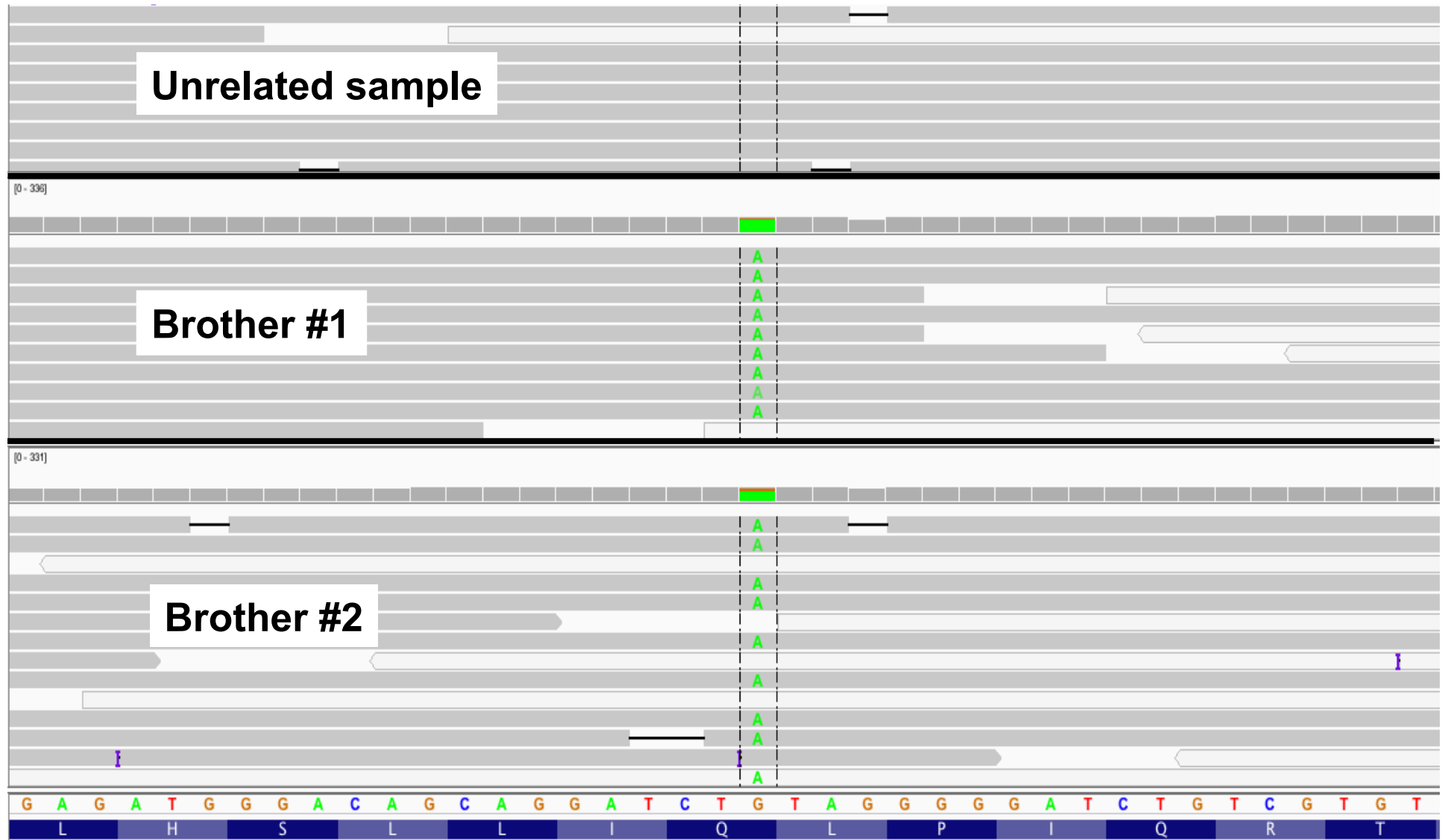
From Wikipedia, the free encyclopedia

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]

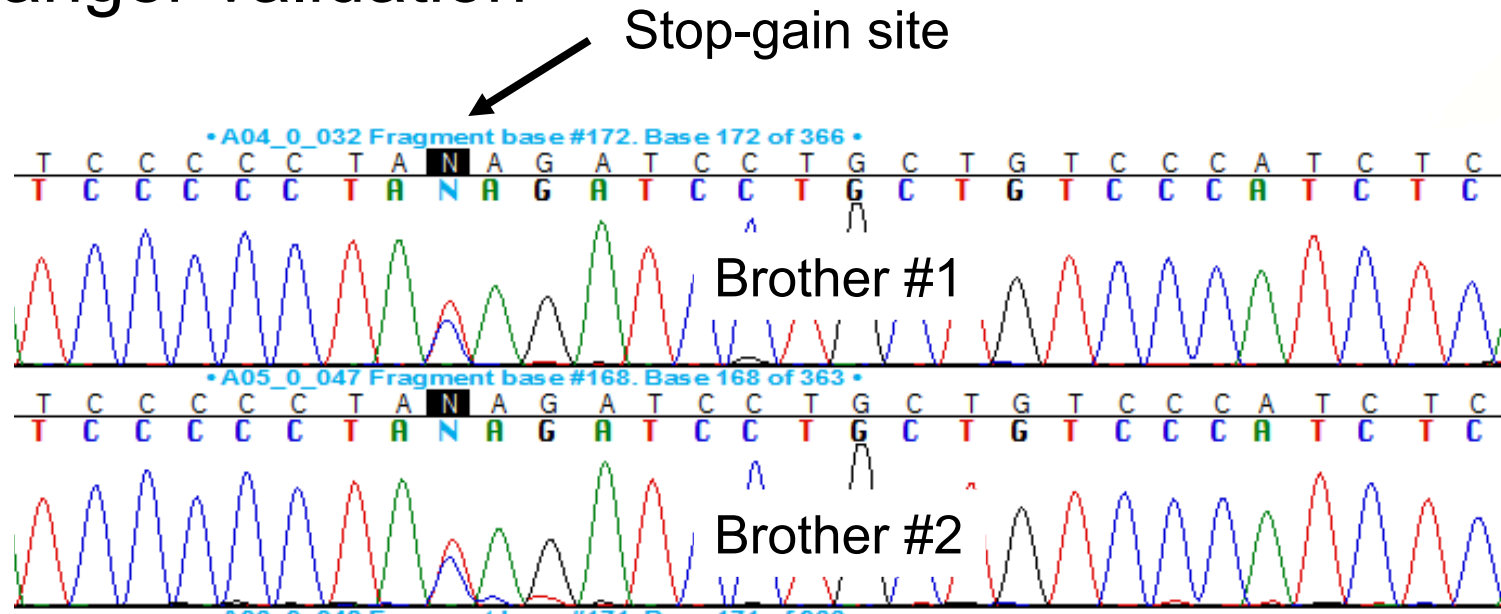
=> Stop-gain in STRC is likely to cause hearing loss!

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

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

Success rate >80% for recent Proton projects!

CanvasDB - Availability

- CanvasDB system 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



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

Genomics
england



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.

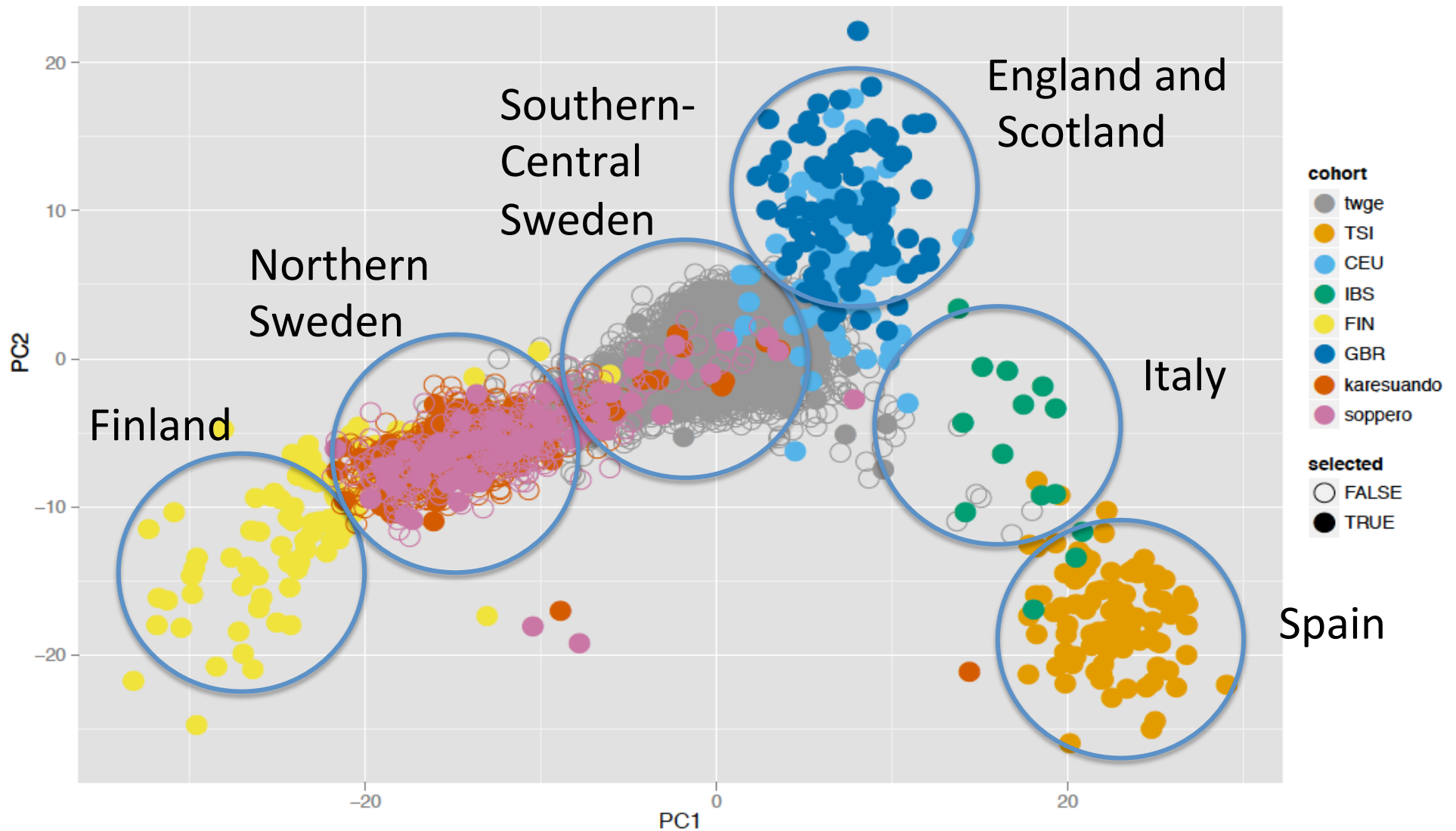
Twin Registry samples used as control cohort

- Inclusion based on twinning
- Distribution like population density
- General population-prevalence of disease
- 10,000 individuals have been analysed with SNP arrays



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

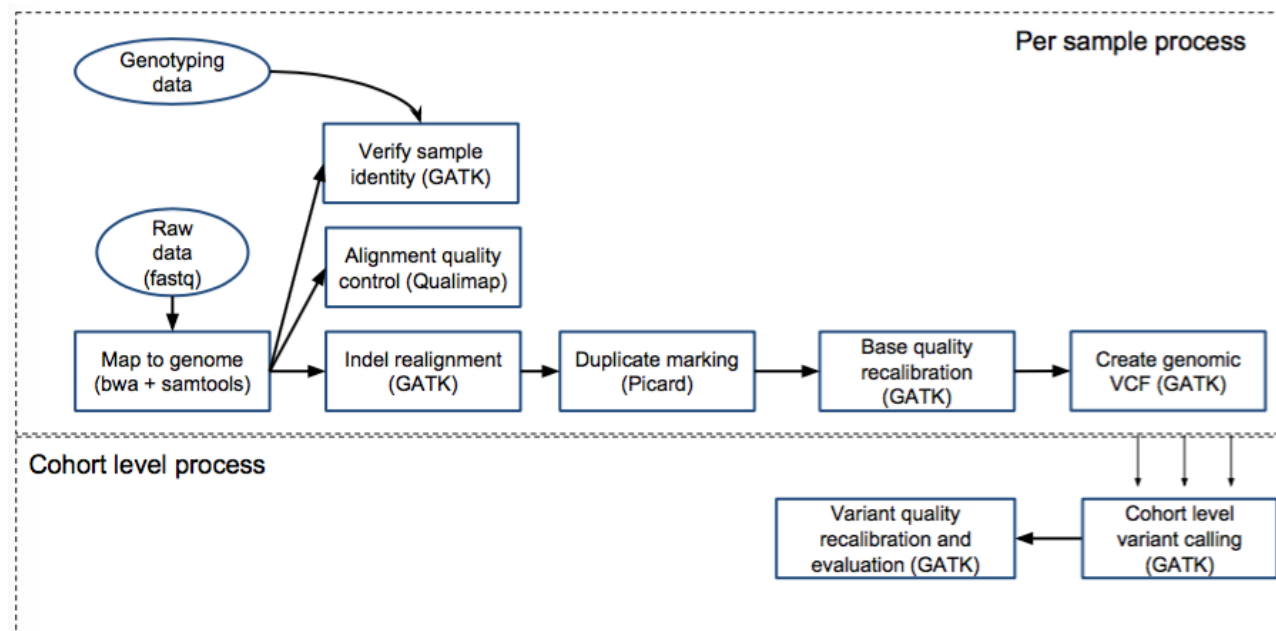
Principal components of European samples from 1,000 genomes project and 10,000 Swedish samples



Illumina WGS of Swedish control cohort

Step 1: 30X Illumina data of the 1,000 individuals

Step 2: Mapping and variant calling



Step 3: Making genotype frequencies available for download

A web server for 'SweGen' data

SweFreq

About

Terms of use

Data Beacon

ExAC Browser

Download Data

Admin

Adam Ameer [Logout](#)
adam.ugc@gmail.com

SweGen Variant Frequency Database

This server hosts whole-genome variant frequencies for 1000 Swedish individuals generated within the SweGen project. The frequency data is intended to be used as a resource for the research community and clinical genetics laboratories. Individual positions in the genome can be viewed using the Data Beacon or ExAC Browser by clicking the links above. To access the variant frequency file you need to register.

Please note that the 1000 individuals included in the SweGen project represent a cross-section of the Swedish population and that no disease information has been used for the selection. The frequency data may therefore include genetic variants that are associated with, or causative of, disease.

We request that any use of data from the SweGen project cite [this preprint on bioRxiv](#).



Released on Oct 19th! Data available from: swefreq.nbis.se



SciLifeLab

Example II: Assembly of genomes using Pacific Biosciences



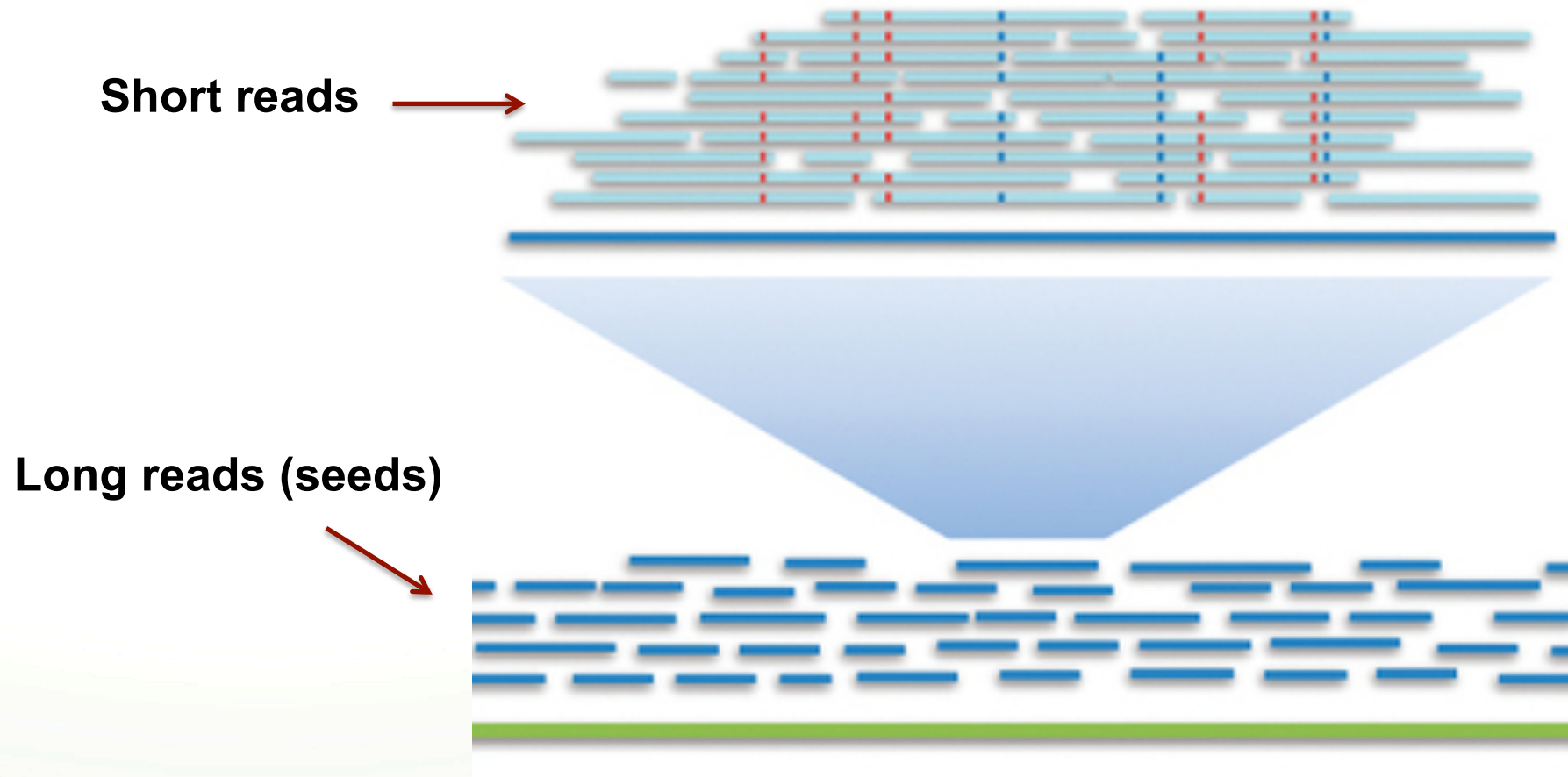
Genome assembly using NGS

- 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



PacBio assembly analysis

- Simple -- just click a button!!

The screenshot displays the PacBio SMRT Portal interface for job assembly analysis. The browser address bar shows the URL `127.0.0.1:8080/smrtportal/#/Design-Job/Details-of-Job/16497`. The page title is "Details of Job assembly". The interface includes a navigation bar with "DESIGN JOB", "MONITOR JOBS", and "VIEW DATA" tabs. The "MONITOR JOBS" tab is active. The job name is "assembly". The protocols are set to "RS_HGAP_Assembly.3" and the reference is "[None selected]".

There are two tables displayed:

SMRT Cells Available (Viewing 1 - 31 of 31)

Sample	Version	User	Groups	Started	Uri
Pb9_frax 21	2.0.2		all	2014-02-20T19:28:20+0000	/home/pacbio/...
Pb9_frax 44	2.0.2		all	2014-02-20T19:28:20+0000	/home/pacbio/...
Pb9_frax 63	2.0.2		all	2014-02-20T19:28:20+0000	/home/pacbio/...
Pb33_1	2.0.2		all	2014-02-20T19:28:20+0000	/home/pacbio/...
Pb33_2	2.0.2		all	2014-02-20T19:28:20+0000	/home/pacbio/...
Pb 33-5	2.0.2		all	2014-02-24T13:48:09+0000	/home/pacbio/...
Pb 33-7	2.0.2		all	2014-02-24T13:48:09+0000	/home/pacbio/...
Pb 33-6	2.0.2		all	2014-02-24T13:48:09+0000	/home/pacbio/...
Pb 33-3	2.0.2		all	2014-02-24T13:48:09+0000	/home/pacbio/...
Pb 33-9	2.0.2		all	2014-02-24T13:48:09+0000	/home/pacbio/...
Pb 33-8	2.0.2		all	2014-02-24T13:48:09+0000	/home/pacbio/...
Pb 33-4	2.0.2		all	2014-02-24T13:48:09+0000	/home/pacbio/...
Pb 33-10	2.0.2		all	2014-02-24T13:48:09+0000	/home/pacbio/...
Pb55_f2rpt	2.1.0		all	2014-05-09T10:48:14+0000	/home/pacbio/...
Pb_46_3_repeat	2.1.0		all	2014-05-09T10:48:14+0000	/home/pacbio/...
Pb55_f2rpt	2.1.0		all	2014-05-09T10:48:14+0000	/home/pacbio/...
Pb_46_9	2.1.0		all	2014-05-09T10:48:14+0000	/home/pacbio/...
Pb_46_10	2.1.0		all	2014-05-09T10:48:14+0000	/home/pacbio/...
Pb46_3	2.1.0		all	2014-05-08T11:08:49+0000	/home/pacbio/...
Pb46_5	2.1.0		all	2014-05-08T11:08:49+0000	/home/pacbio/...

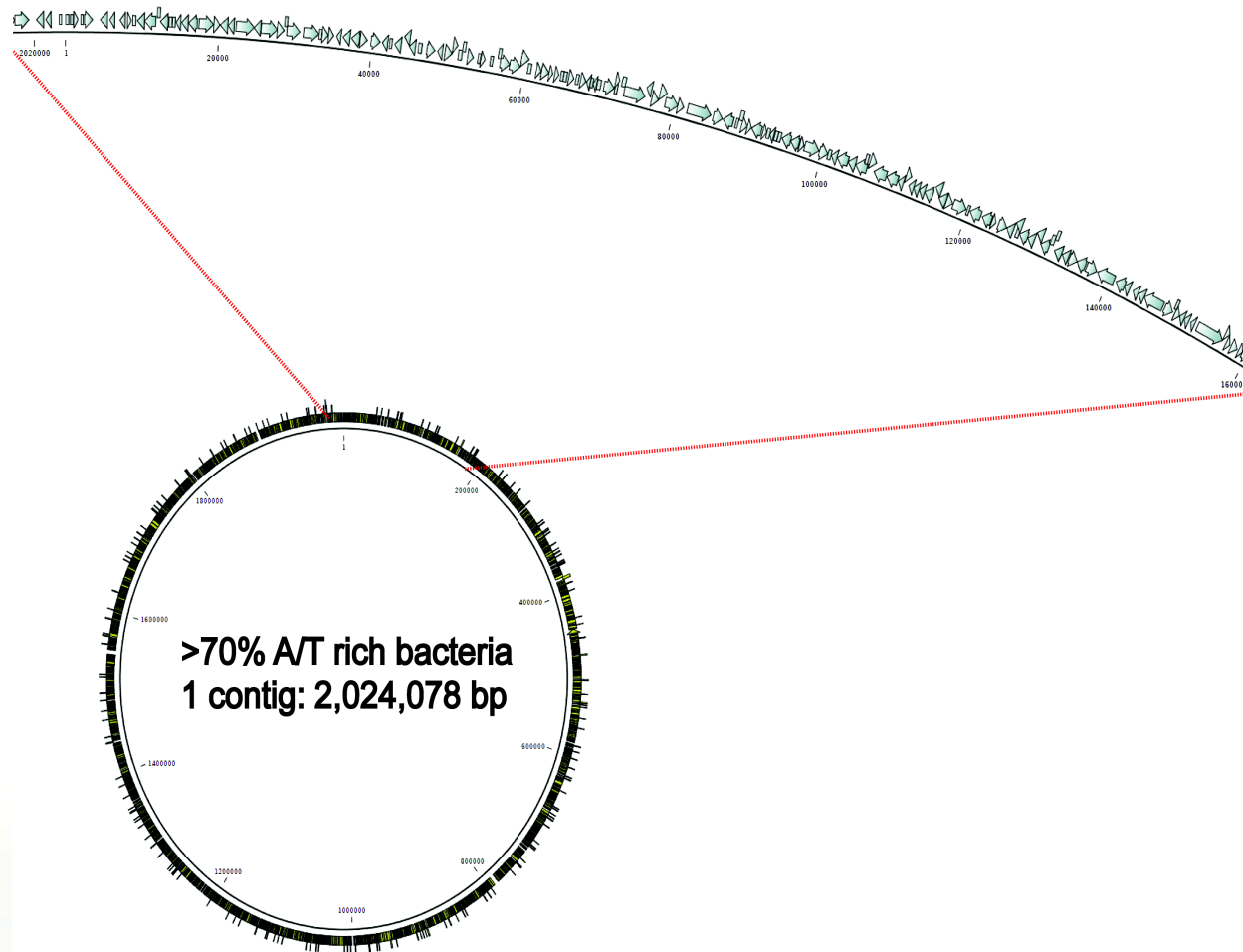
SMRT Cells in Job (Viewing 1 - 1 of 1)

Sample	Version	User	Groups	Uri
Pb33_1	2.0.2		all	/home/pacbio/DATA/adam/Pb_33_F...

Buttons at the bottom: Start, Save, Copy, Cancel.

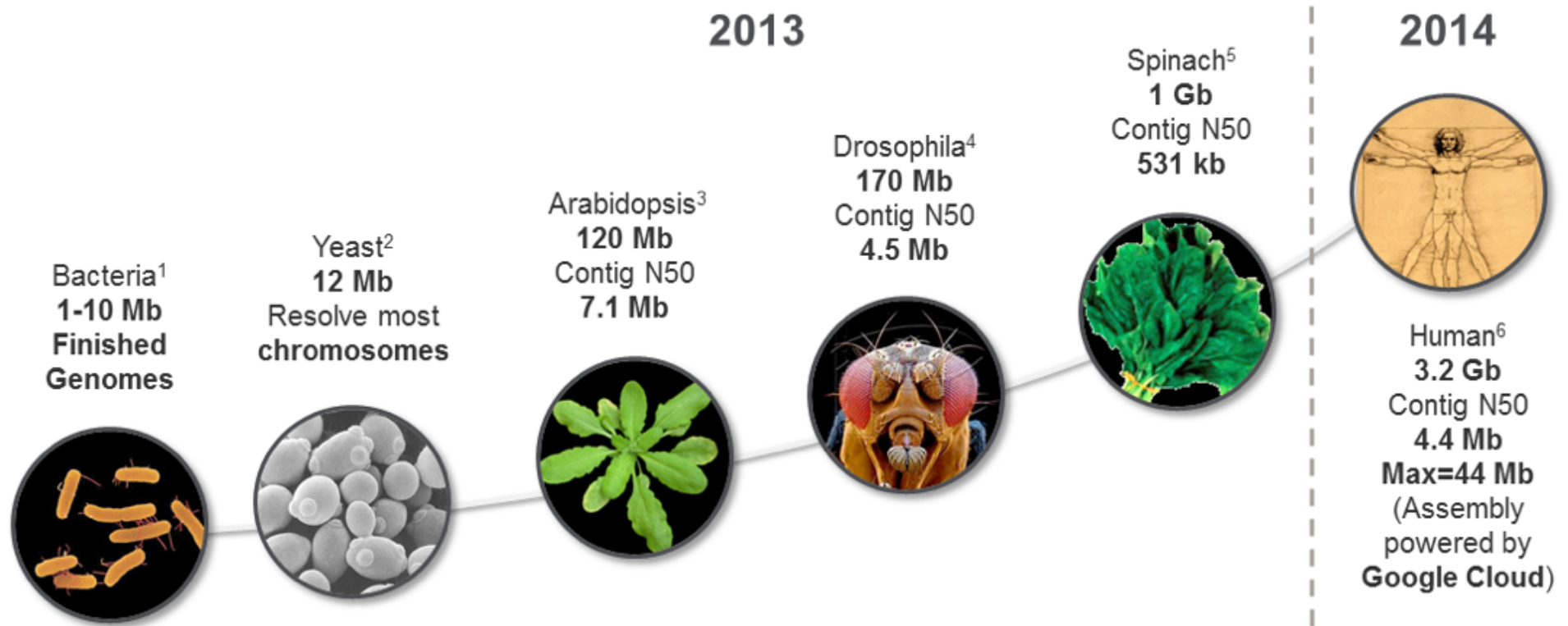
PacBio assembly, example result

- Example: Complete assembly of a bacterial genome



PacBio assembly – recent developments

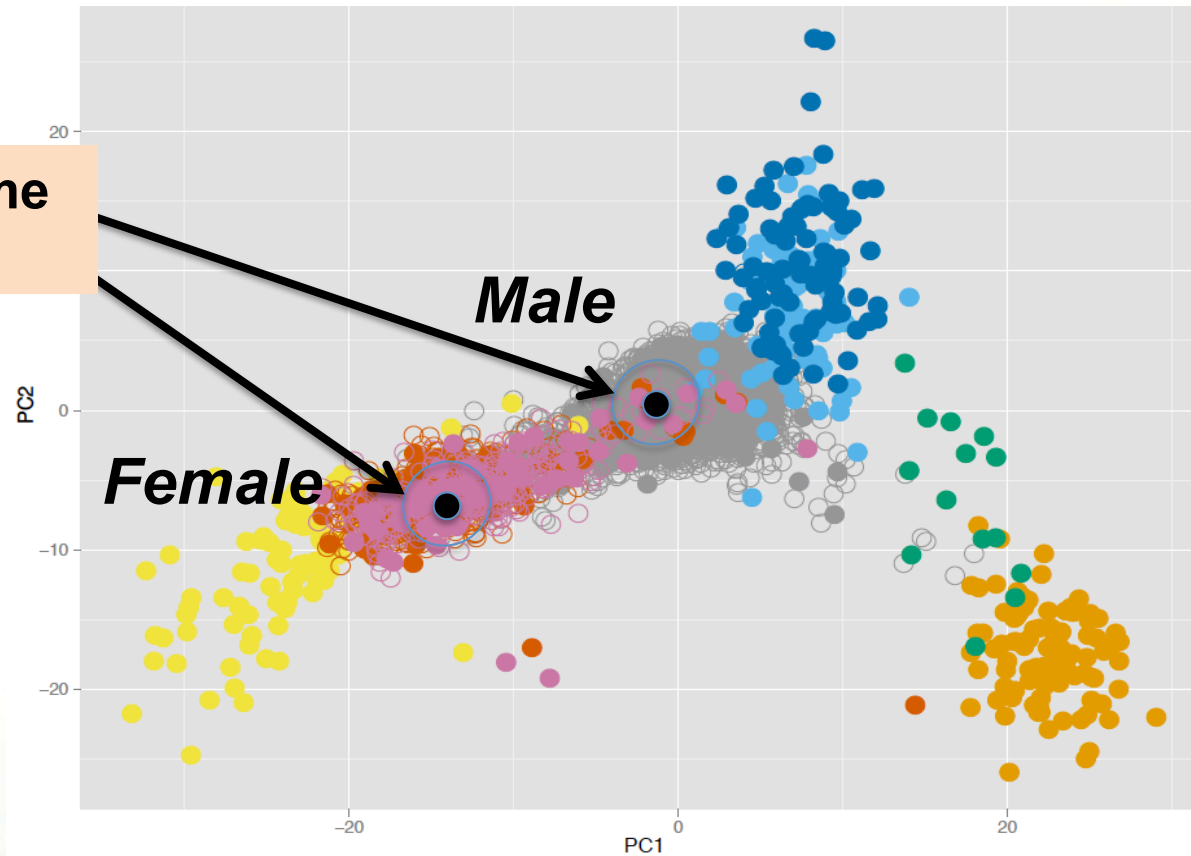
- Also larger genomes can be assembled by PacBio..



De novo WGS of Swedish cohort

Establish Swedish reference genome sequences by *de novo* assembly of long-reads: ***PacBio+BioNano+10X Genomics***

Reference genome individuals

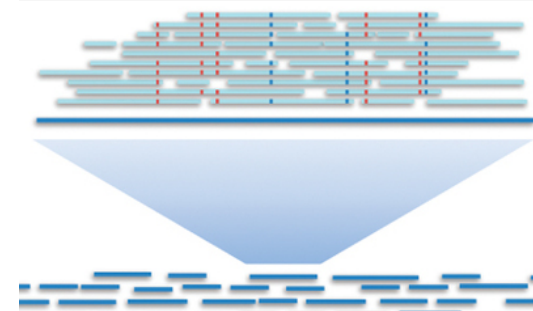
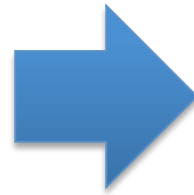


De novo assembly of 75X PacBio data

Assembly (FALCON)



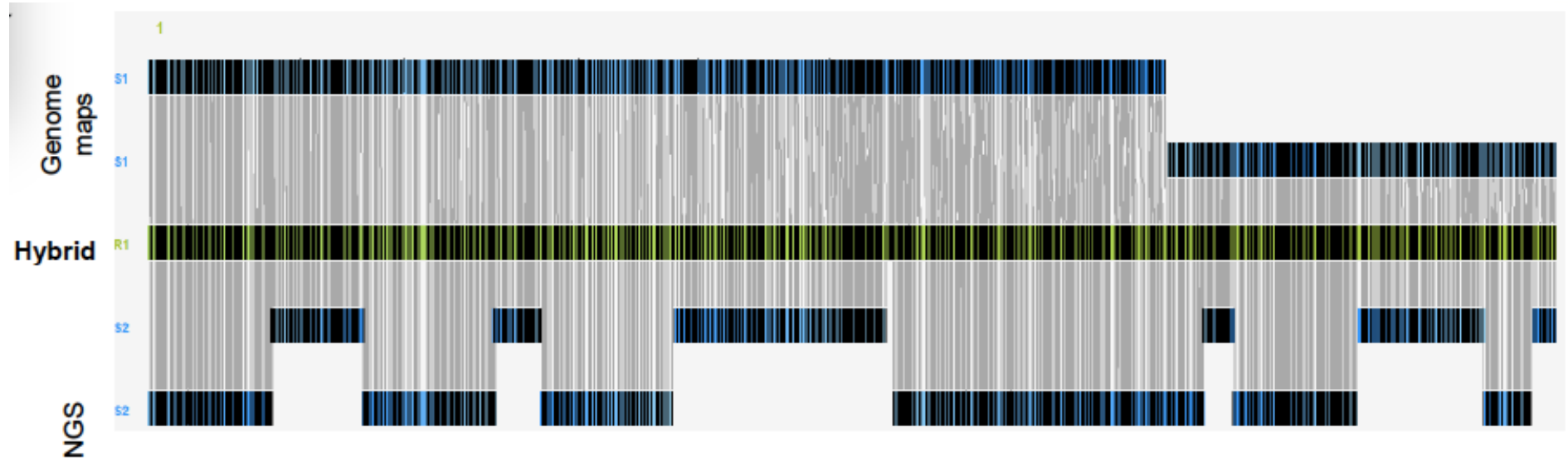
Error correction (2 x Quiver)



Analysis time: 1 month/genome

	Individual 1	Individual 2
Assembly size	3,039,619,582	3,024,752,299
Nr contigs	11,249	11,601
Longest contig	36,8 Mb	54,1 Mb
N50	8,9 Mb	8,3 Mb

Hybrid scaffolding, PacBio + BioNano

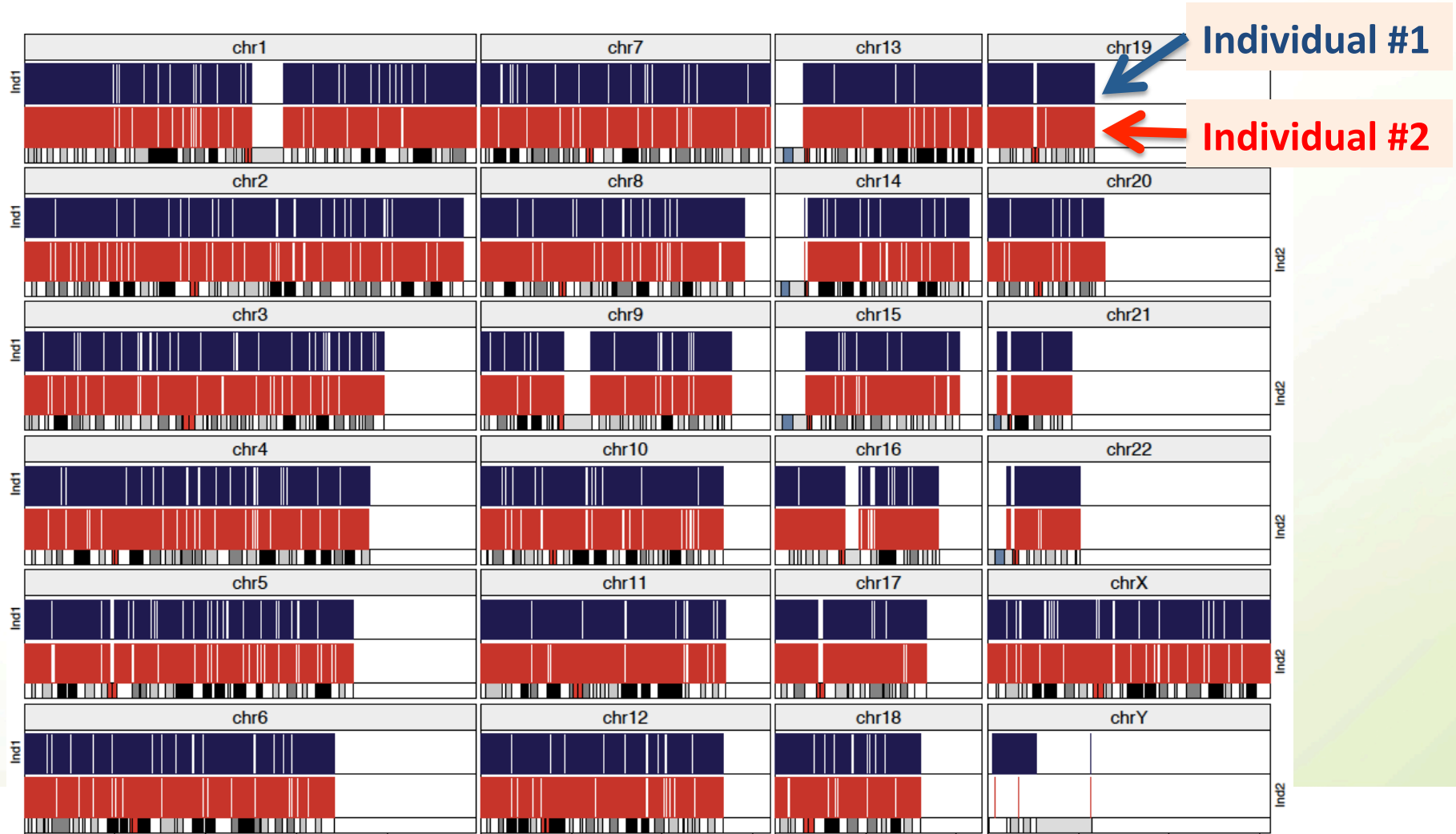


Hybrid scaffolding with two labellings resulted in

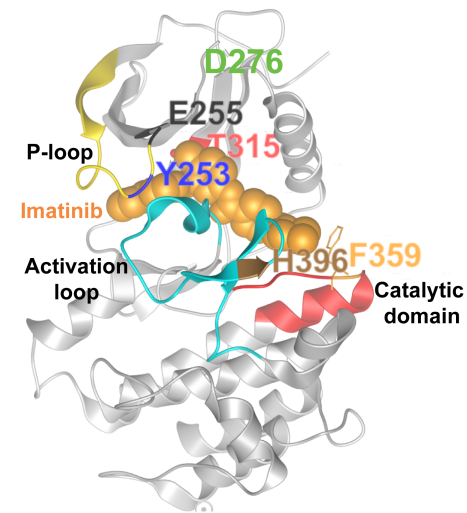
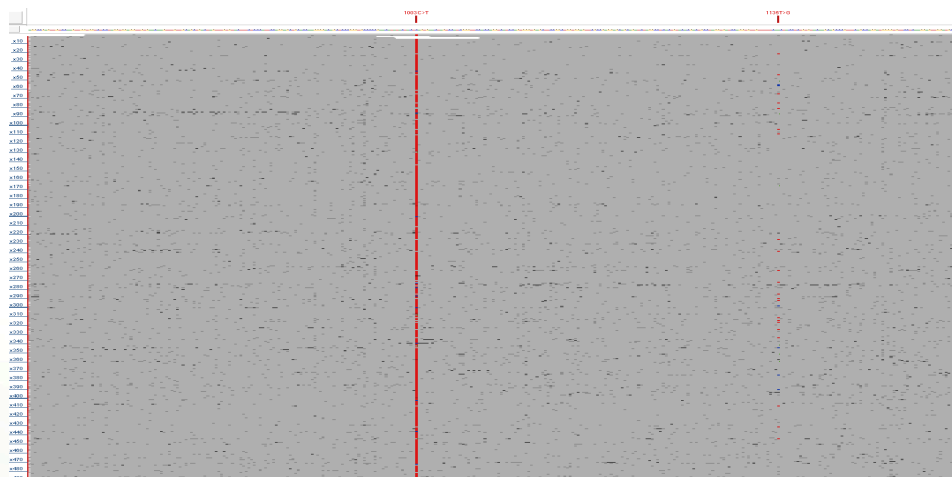
- **3,1 Gb** assembly, **51 Mb** N50 (for individual #1)
- **3,1 Gb** assembly, **46 Mb** N50 (for individual #2)

Aligning contigs to human reference

> 99% of bases can be aligned to human reference (hg38)

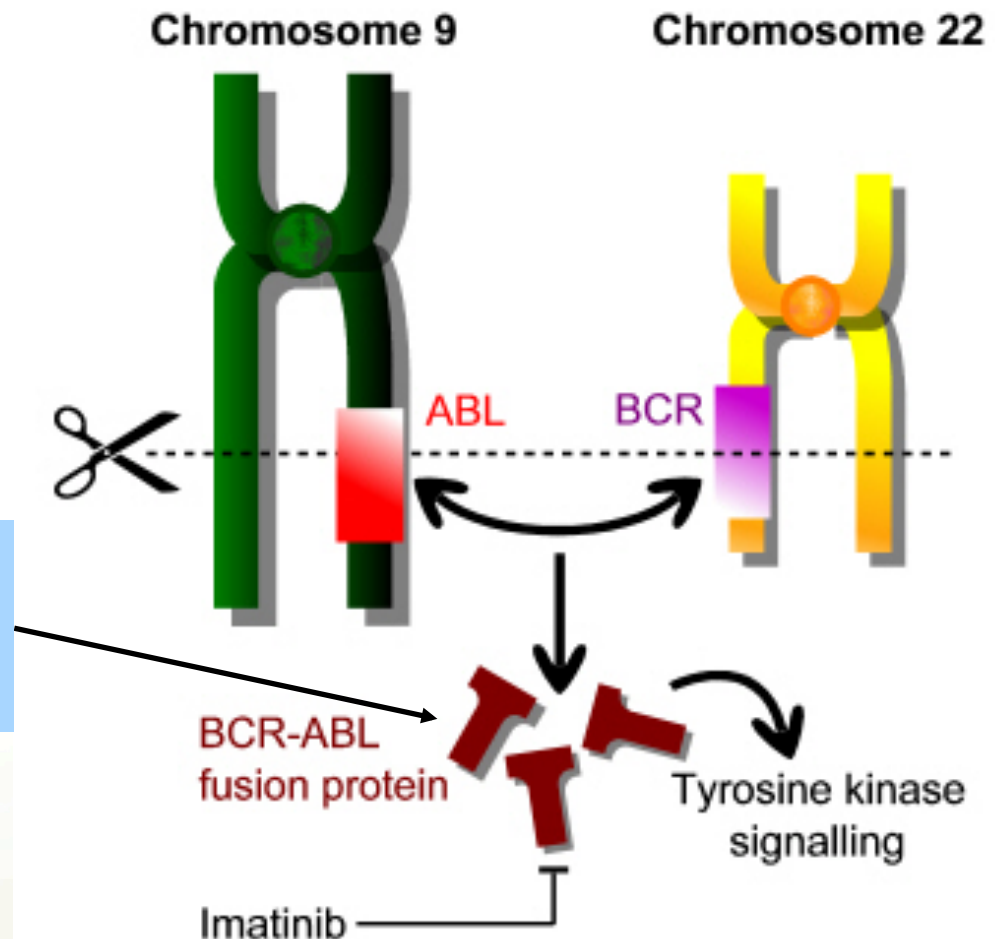


Example III: Clinical sequencing for Leukemia Treatment



Chronic Myeloid Leukemia

- BCR-ABL1 fusion protein – a CML drug target

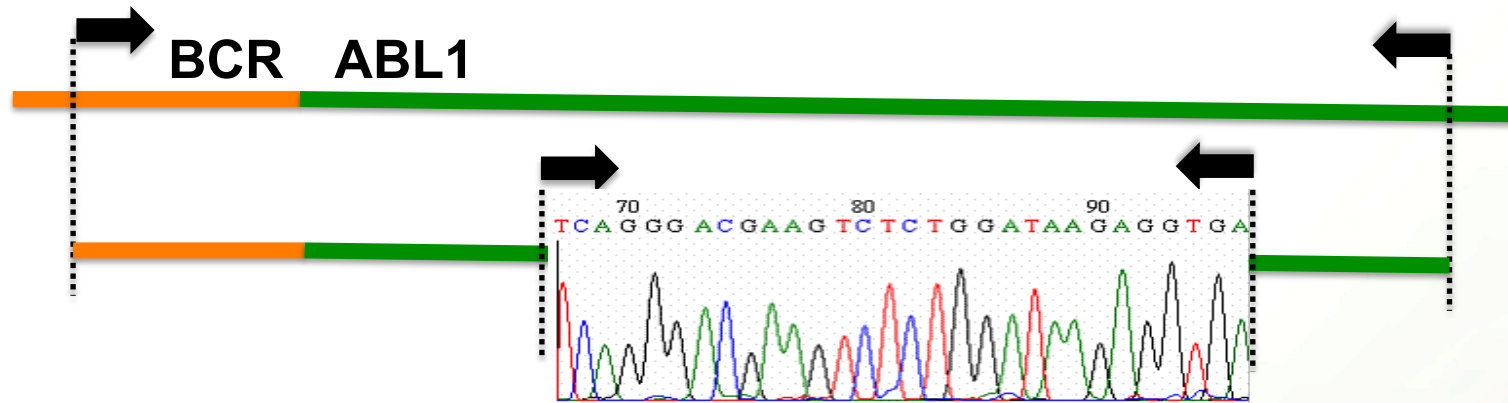


The BCR-ABL1 fusion protein can acquire resistance mutations following drug treatment

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

Traditional mutation screening in BCR-ABL1

Nested PCR and Sanger sequencing:



Limitations:

- Mutations at frequencies below 10-20% not seen
- Biases may be introduced by nested PCR
- Whole BCR-ABL1 fusion transcript not sequenced
- Clonal composition of mutations not determined

Our clinical diagnostics pipeline for BCR-ABL1

Cavelier et al. *BMC Cancer* (2015) 15:45
DOI 10.1186/s12885-015-1046-y



RESEARCH ARTICLE

Open Access

Clonal distribution of *BCR-ABL1* mutations and splice isoforms by single-molecule long-read RNA sequencing

Lucia Cavelier^{1*}, Adam Ameur^{1†}, Susana Häggqvist¹, Ida Höjjer¹, Nicola Cahill¹, Ulla Olsson-Strömberg² and Monica Hermanson¹

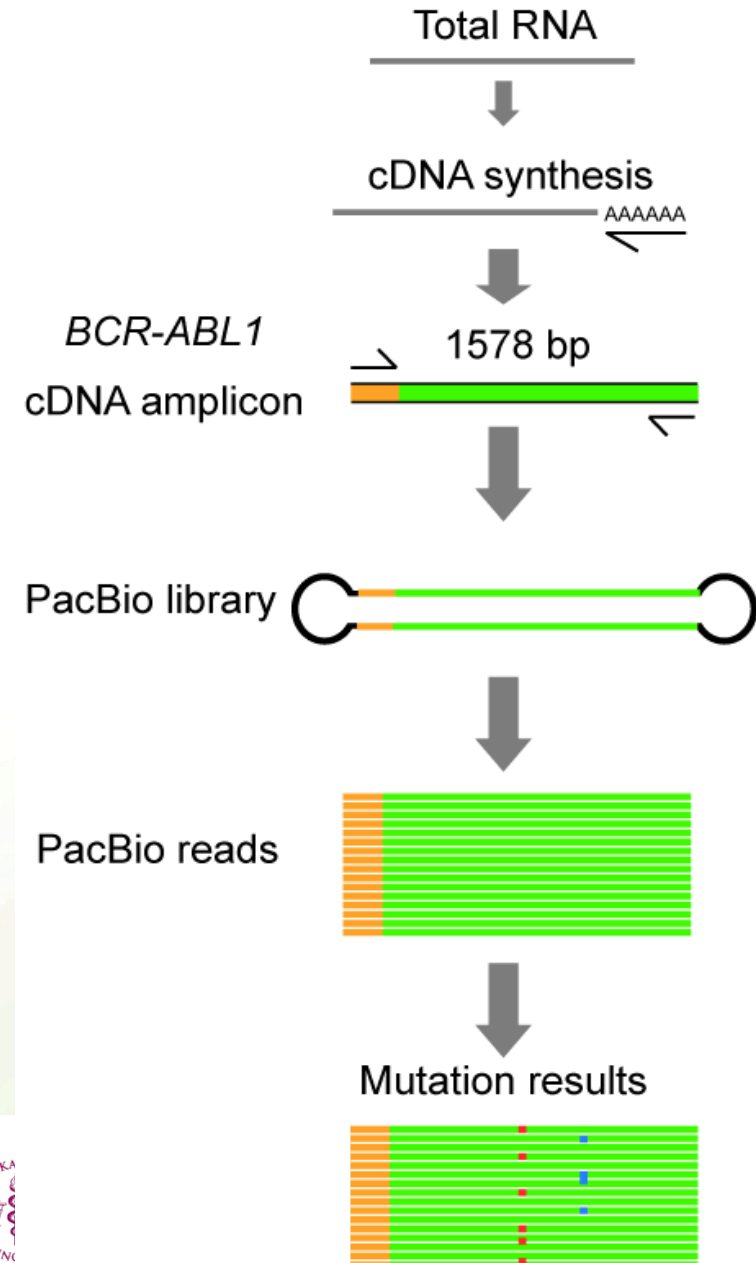
Abstract

Background: The evolution of mutations in the *BCR-ABL1* fusion gene transcript renders CML patients resistant to tyrosine kinase inhibitor (TKI) based therapy. Thus screening for *BCR-ABL1* mutations is recommended particularly in patients experiencing poor response to treatment. Herein we describe a novel approach for the detection and surveillance of *BCR-ABL1* mutations in CML patients.

Methods: To detect mutations in the *BCR-ABL1* transcript we developed an assay based on the Pacific Biosciences (PacBio) sequencing technology, which allows for single-molecule long-read sequencing of *BCR-ABL1* fusion transcript molecules. Samples from six patients with poor response to therapy were analyzed both at diagnosis and follow-up. cDNA was generated from total RNA and a 1,6 kb fragment encompassing the *BCR-ABL1* transcript was amplified using long range PCR. To estimate the sensitivity of the assay, a serial dilution experiment was performed.

Results: Over 10,000 full-length *BCR-ABL1* sequences were obtained for all samples studied. Through the serial dilution analysis, mutations in CML patient samples could be detected down to a level of at least 1%. Notably, the assay was determined to be sufficiently sensitive even in patients harboring a low abundance of *BCR-ABL1* levels. The PacBio sequencing successfully identified all mutations seen by standard methods. Importantly, we identified several mutations that escaped detection by the clinical routine analysis. Resistance mutations were found in all but one of the patients. Due to the long reads afforded by PacBio sequencing, compound mutations present in the same molecule were readily distinguished from independent alterations arising in different molecules. Moreover, several transcript isoforms of the *BCR-ABL1* transcript were identified in two of the CML patients. Finally, our assay allowed for a quick turn around time allowing samples to be reported upon within 2 days.

Conclusions: In summary the PacBio sequencing assay can be applied to detect *BCR-ABL1* resistance mutations in both diagnostic and follow-up CML patient samples using a simple protocol applicable to routine diagnosis. The method besides its sensitivity, gives a complete view of the clonal distribution of mutations, which is of importance when making therapy decisions.

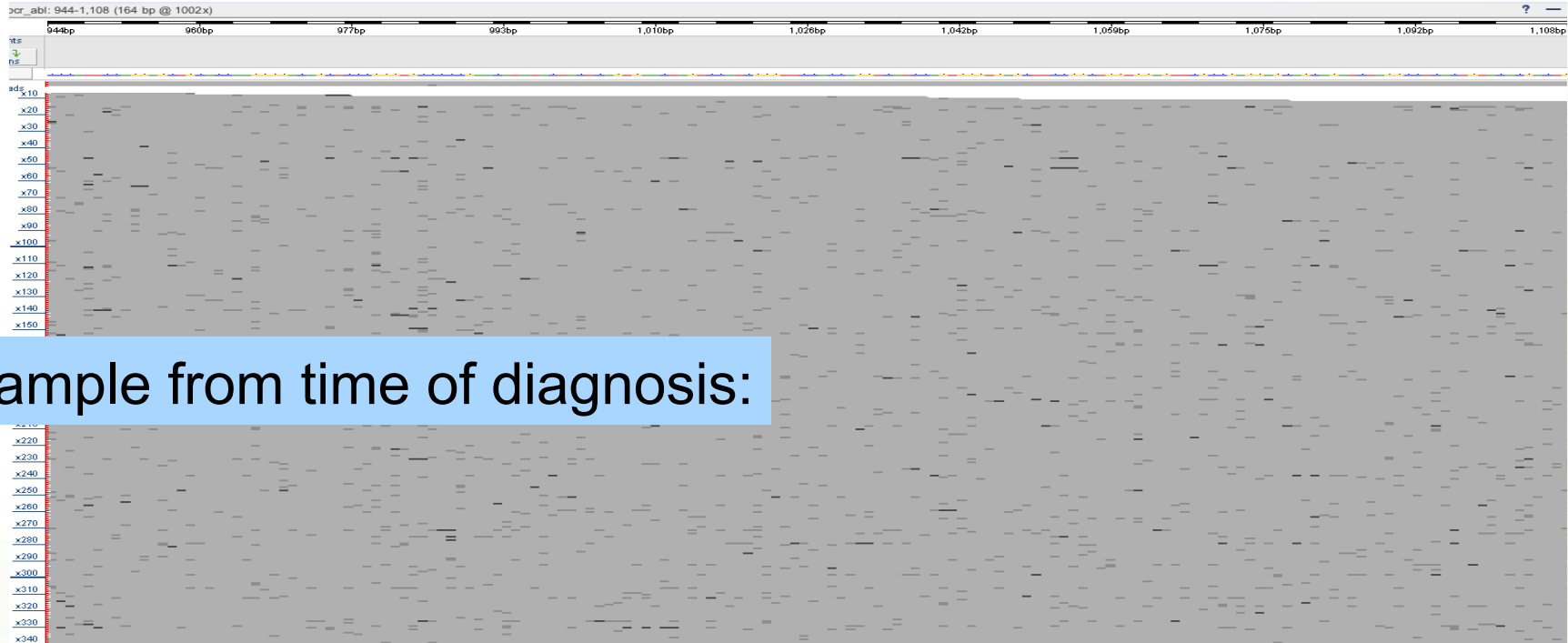


BCR-ABL1 mutations at diagnosis

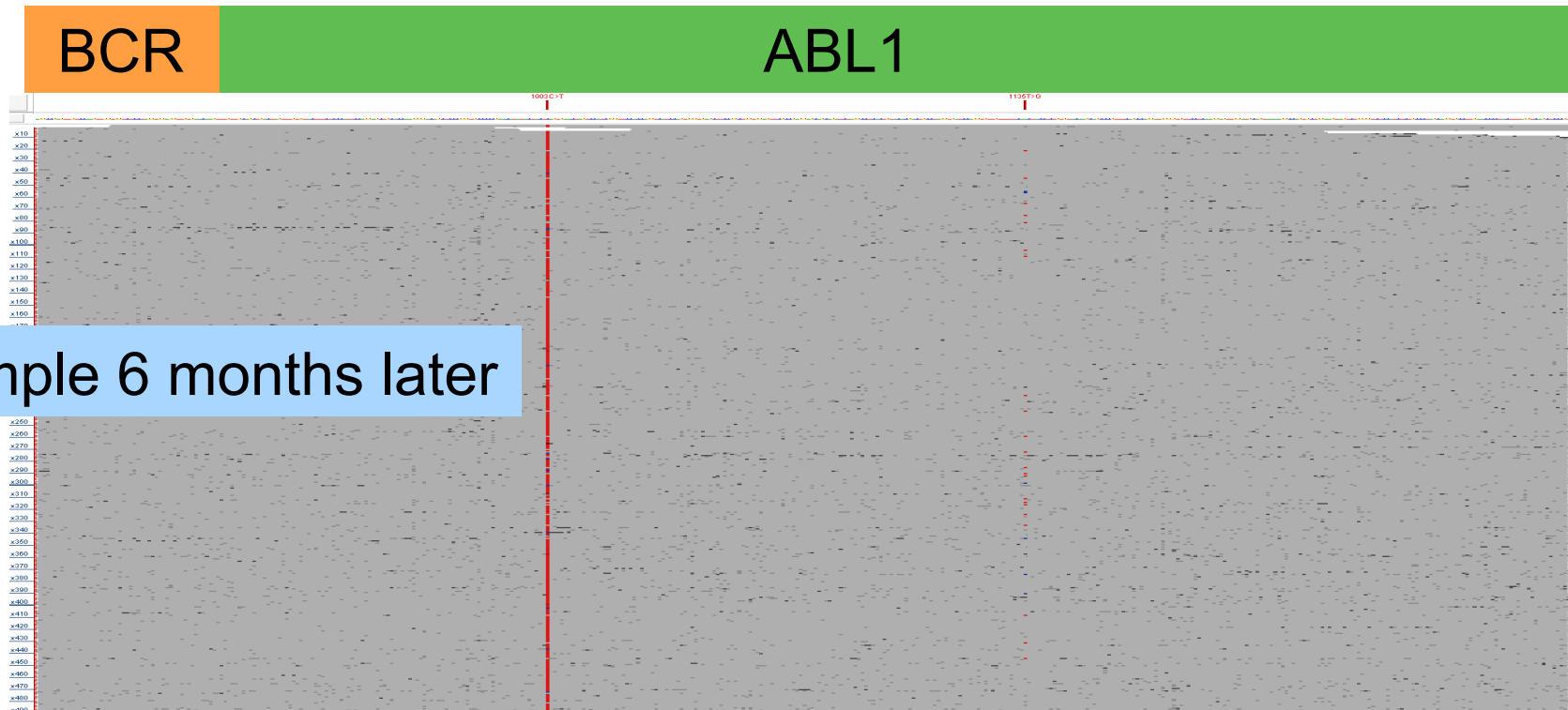
PacBio sequencing generates ~10 000X coverage!

BCR

ABL1



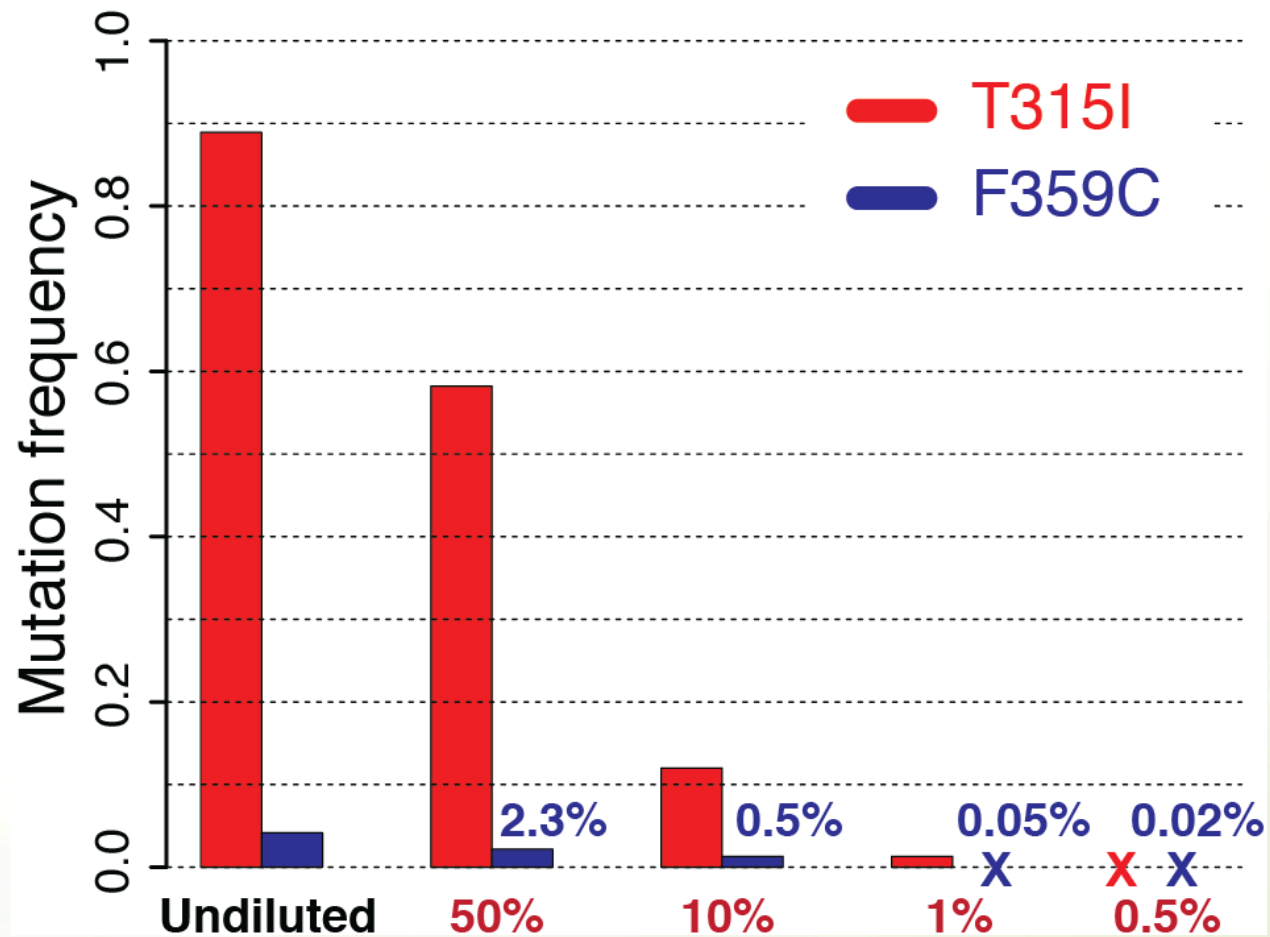
BCR-ABL1 mutations in follow-up sample



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

BCR-ABL1 dilution series results

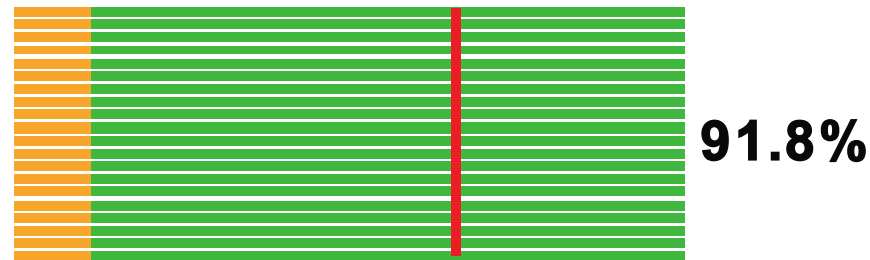
- Mutations down to 1% detected!



BCR-ABL1 - Compound mutations

49 months

T315I

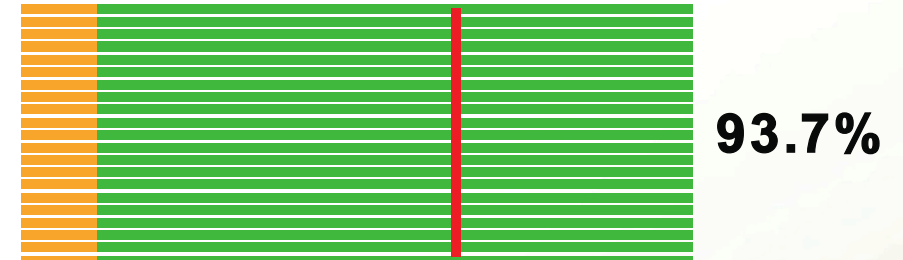


F359C



55 months

T315I



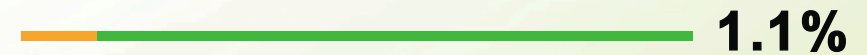
D276G



F359C



H396R



Analysis method for BCR-ABL1 mutations

- Create CCS reads and screen all known resistance mutations
- CAVA analysis - count number of WT and MUT sequences

WT sequence: **TATATCATCACTGAGTTCATG**

MUT sequence: **TATATCATCA**T**TGAGTTCATG**



BCR-ABL1 resistance mutation

- Classify each mutation
 - Less than 500X coverage => **Unresolved**
 - At least 0.5% mutation frequency => **Positive**
 - Otherwise => **Negative**

Clinical Diagnosis of BCR-ABL1 mutations

Clinical Genetics



- Collection of samples
- Seq library preparation

Sequencing Facility



- SMRT sequencing
- CAVA analysis

IT developers

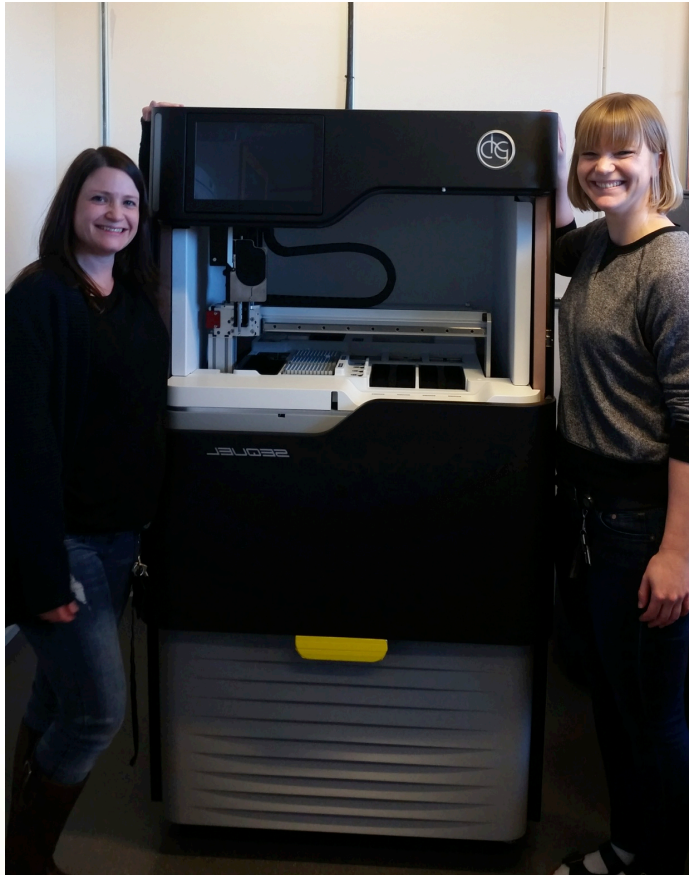


- Web server for results

- Ongoing routine service, 0-4 samples/week.
- Over 150 patient samples run so far
- 100% consistency with Sanger results

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 Höjjer
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



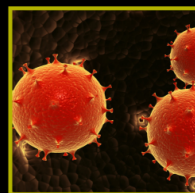
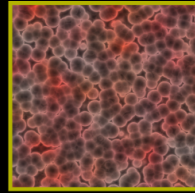
Anna Petri
Research engineer
Sequencing Service



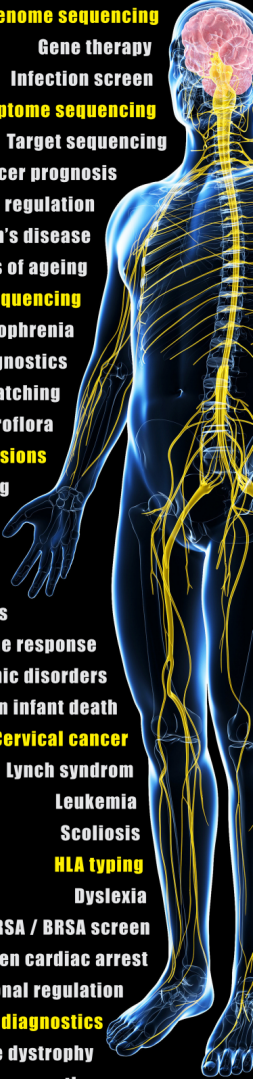
Anne-Christine Lindström
Research engineer
Sequencing Service

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...