



LEIDEN UNIVERSITY MEDICAL CENTER

The Dutch Variant Database.

Jeroen F. J. Laros

Leiden Genome Technology Center

Department of Human Genetics

Center for Human and Clinical Genetics



The Dutch Variant Database.

Goal: Share variants for annotation purposes.

The Dutch Variant Database.

Goal: Share variants for annotation purposes.

Prototype by:

- Hubrecht / UMCU Utrecht.
- UMCN Nijmegen.
- LUMC Leiden.

The Dutch Variant Database.

Goal: Share variants for annotation purposes.

Prototype by:

- Hubrecht / UMCU Utrecht.
- UMCN Nijmegen.
- LUMC Leiden.

Other participants:

- UMCG Groningen.
- AMC Amsterdam.
- EMC Rotterdam.
- VUMC Amsterdam.
- NBIC.

In *exome sequencing*, we select genomic regions of interest using a *target-enrichment strategy*.

- PCR.
- On array capture.
- **In-solution capture.**

In *exome sequencing*, we select genomic regions of interest using a *target-enrichment strategy*.

- PCR.
- On array capture.
- **In-solution capture.**

Overview of an in-solution capture.

- Fragmentation.
- Size selection.
- Linker ligation.
- Capture.

In *exome sequencing*, we select genomic regions of interest using a *target-enrichment strategy*.

- PCR.
- On array capture.
- **In-solution capture.**

Overview of an in-solution capture.

- Fragmentation.
- Size selection.
- Linker ligation.
- Capture.

These regions are then *sequenced*.

- ABI SOLiD 5500.
- Illumina HiSeq 2000.

Introduction

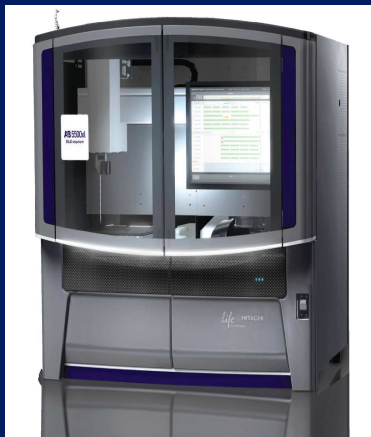


Figure 1: SOLiD 5500XL.



Figure 2: HiSeq 2000.

Paired end, high throughput, cheap.

Raw data acquisition.



Figure 3: Raw image.

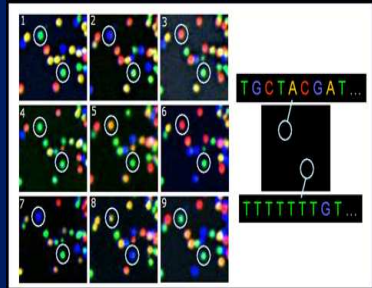


Figure 4: Base calling.

Raw data acquisition.

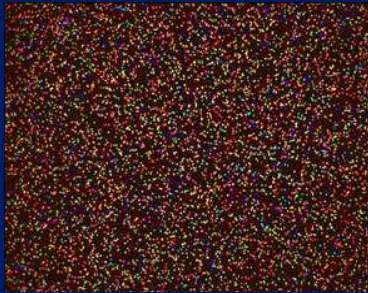


Figure 3: Raw image.

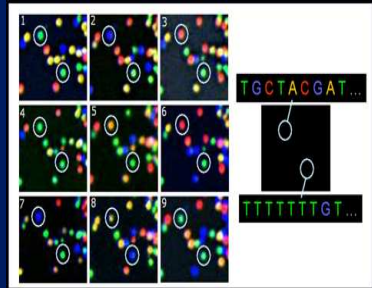


Figure 4: Base calling.

Output:

- Long list of short DNA sequences.
- Quality scores per base.

Exome sequencing pipelines can roughly be divided in five steps.

Exome sequencing pipelines can roughly be divided in five steps.

1. Pre-alignment.
 - Data cleaning.

Exome sequencing pipelines can roughly be divided in five steps.

1. Pre-alignment.
 - Data cleaning.
2. Alignment.

Exome sequencing pipelines can roughly be divided in five steps.

1. Pre-alignment.
 - Data cleaning.
2. Alignment.
3. Variant calling.
 - Description of difference with the reference.

Exome sequencing pipelines can roughly be divided in five steps.

1. Pre-alignment.
 - Data cleaning.
2. Alignment.
3. Variant calling.
 - Description of difference with the reference.
4. Filtering.
 - Variant calling has a high false positive rate.

Exome sequencing pipelines can roughly be divided in five steps.

1. Pre-alignment.
 - Data cleaning.
2. Alignment.
3. Variant calling.
 - Description of difference with the reference.
4. Filtering.
 - Variant calling has a high false positive rate.
5. Annotation.

We use the Trimmomatic / FASTX toolkit for data cleaning.

- Remove linker sequences.
- Clip low quality reads at the end of the read.
- Judge the read that is left over.

We use the Trimmomatic / FASTX toolkit for data cleaning.

- Remove linker sequences.
- Clip low quality reads at the end of the read.
- Judge the read that is left over.

The FASTQC toolkit is used for quality control (both before and after the data cleaning step).

- GC content.
- GC distribution.
- Quality scores distribution.
- ...

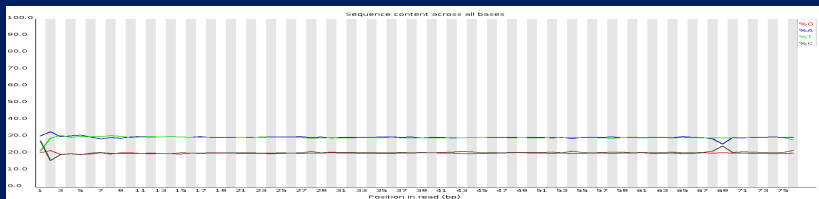


Figure 5: Per base sequence content.

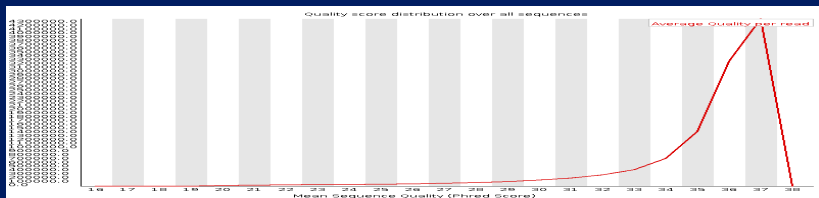


Figure 6: Per sequence quality.

Stampy: A statistical algorithm for sensitive and fast mapping of Illumina sequence reads.

Stampy: A statistical algorithm for sensitive and fast mapping of Illumina sequence reads.

Some features:

- Base quality recalibration.
 - First map 1% of the input.
 - Recalibrate the Fastq quality scores.
 - Redo the alignment with the recalibrated scores.

Stampy: A statistical algorithm for sensitive and fast mapping of Illumina sequence reads.

Some features:

- Base quality recalibration.
 - First map 1% of the input.
 - Recalibrate the Fastq quality scores.
 - Redo the alignment with the recalibrated scores.
- Uses BWA for the hard work.
 - Switches to its accurate built in aligner when BWA fails.

Burrows-Wheeler Aligner (BWA) is a short read aligner that allows small insertions and deletions.

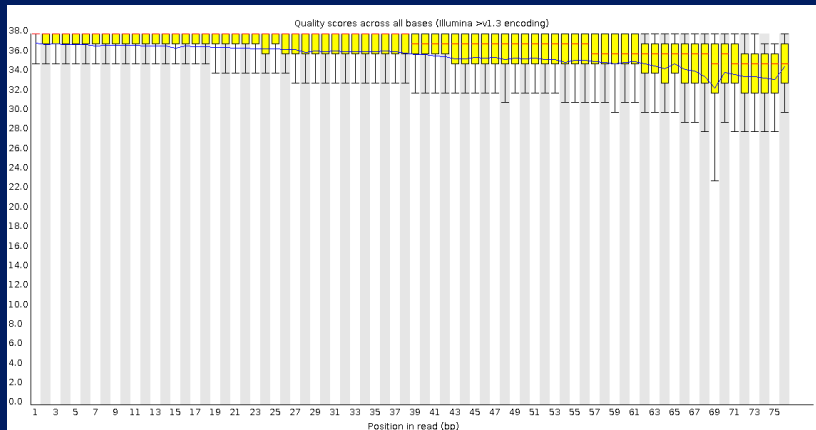


Figure 7: Base quality recalibration.

Variant calling

```

ACCGTTAAGACCAAGTCTTTCGGACTCTCGA X 4
ACCGTTAAGACCAAGTCTTTCGGACTCTCGAC X 2
ACCGTTAAGACCAAGTCTTTCGGACTCTCGAGC X 2
CGGTTAAGACCAAGTCTTTCGGACTCTCGACT X 1
GTTTANGACCAAGTCTTTCGGACTCTCGACTC X 2
GTTTANGACCAAGTCTTTCGGACTCTCGACTCG X 1
GTTTANGACCAAGTCTTTCGGACTCTCGACTCG X 1
TTAAGACCAAGTCTTTCGGACTCTCGACTCGA X 2
TTAAGACCAAGTCTTTCGGACTCTCGACTCGA X 1
TAAGACCAAGTCTTTCGGACTCTCGACTCGAA X 2
TAAGACCAAGTCTTTCGGACTCTCGACTCGAA X 2
TAAGACCAAGTCTTTCGGACTCTCGACTCGAA X 1
TAAGACCAAGTCTTTCGGACTCTCGACTCGAA X 1
GACCAAGTCTTTCGGACTCTCGACTCGAACCT X 1
GACCAAGTCTTTCGGACTCTCGACTCGAACCT X 1
ACCAAGTCTTTCGGACTCTCGACTCGAACCT X 1
CCAAAGTCTTTCGGACTCTCGACTCGAACCT X 1
TAAAGTCTTTCGGACTCTCGACTCGAACCTTTA X 1
CAAGTCTTTCGGACTCTCGACTCGAACCTTTA X 1
AAGTCTTTCGGACTCTCGACTCGAACCTTTAG X 1
AAGTCTTTCGGACTCTCGACTCGAACCTTTAG X 1
AGTCTTTCGGACTCTCGACTCGAACCTTTAG X 1
GTCTTTCGGACTCTCGACTCGAACCTTTAGT X 1
GTCTTTCGGACTCTCGACTCGAACCTTTAGT X 1
TCCTTTCGGACTCTCGACTCGAACCTTTAGGT X 2
TCCTTTCGGACTCTCGACTCGAACCTTTAGGT X 1
CTTTCGGACTCTCGACTCGAACCTTTAGGT X 1
CTTTCGGACTCTCGACTCGAACCTTTAGGT X 1
TTTTCGGACTCTCGACTCGAACCTTTAGGTGA X 2
TTTTCGGACTCTCGACTCGAACCTTTAGGTGA X 1
TTTTCGGACTCTCGACTCGAACCTTTAGGTGA X 2
TCGGACTCTCGACTCGAACCTTTAGGTGAAA X 5
CGGACTCTCGACTCGAACCTTTAGGTGAAA X 1
CGGACTCTCGACTCGAACCTTTAGGTGAAA X 1
GGACTCTCGACTCGAACCTTTAGGTGAAA X 1
GACTCTCGACTCGAACCTTTAGGTGAAA X 1
GACTCTCGACTCGAACCTTTAGGTGAAA X 1
CTTTCGGACTCTCGAACCTTTAGGTGAAA X 1
GTTTCGGACTCTCGAACCTTTAGGTGAAA X 1
CTCGACTCGAACCTTTAGGTGAAA X 1
TCGACTCGAACCTTTAGGTGAAA X 2
TCGACTCGAACCTTTAGGTGAAA X 1
CGGACTCGAACCTTTAGGTGAAA X 1
TTGGGAATCTGGTTGAGAAAGCTGAGAACCGACTGGAAATCCAGATTTCTCTGGCTGC

```

Figure 8: Variant calling.

Variant calling

Variant calling is done by Samtools, BCFtools / VCFutils.

The output of most modern aligners is in *Sequence Alignment / Map* (SAM) format.

Variant calling

Variant calling is done by Samtools, BCFtools / VCFutils.

The output of most modern aligners is in *Sequence Alignment / Map* (SAM) format.

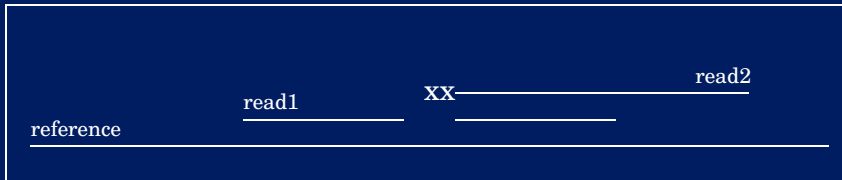
Mainly file format conversions.

- **SAM** → BAM.
- BAM → BAM.sorted.
- BAM.sorted → BAM.sorted.index.
- BAM.sorted → mpileup (**BAQ realignment**).
- BAM.sorted → BCF.
- BCF → **VCF**.

We end up with a list in *Variant Call Format* (VCF).

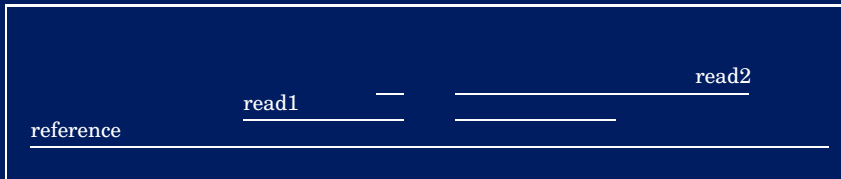
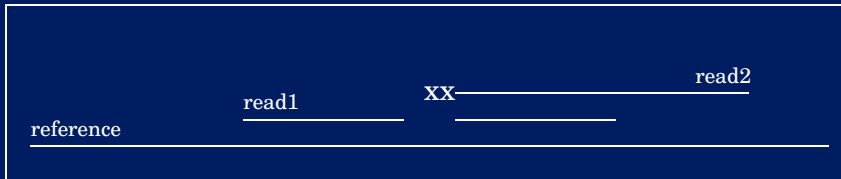
Variant calling

Base Alignment Quality (BAQ) realignment:
 Remove SNPs around indels.



Variant calling

Base Alignment Quality (BAQ) realignment:
Remove SNPs around indels.



Samtools varfilter.

- Minimum coverage threshold.
- Strand bias.
- Quality scores.
- **Maximum coverage threshold.**
 - Copy number variation.
 - Alignment artefacts.

Samtools varfilter.

- Minimum coverage threshold.
- Strand bias.
- Quality scores.
- **Maximum coverage threshold.**
 - Copy number variation.
 - Alignment artefacts.

Transition transversion rates:

- Transitions: **A** \Leftrightarrow **G** or **C** \Leftrightarrow **T**.
- Transversions: **A** \Leftrightarrow **C** or **G** \Leftrightarrow **T** or **A** \Leftrightarrow **T** or **C** \Leftrightarrow **G**.
- Transitions are more frequent.

Samtools varfilter.

- Minimum coverage threshold.
- Strand bias.
- Quality scores.
- **Maximum coverage threshold.**
 - Copy number variation.
 - Alignment artefacts.

Transition transversion rates:

- Transitions: **A \Leftrightarrow G or C \Leftrightarrow T.**
- Transversions: **A \Leftrightarrow C or G \Leftrightarrow T or A \Leftrightarrow T or C \Leftrightarrow G.**
- Transitions are more frequent.

- Around 2.1 human full genome.
- Around 2.8 to 3.0 human exome.

We use different annotation sources and an in-house database.

We use different annotation sources and an in-house database.

A selection of SeattleSeq annotation:

- Is the variant known?
- Does it hit a gene?

We use different annotation sources and an in-house database.

A selection of SeattleSeq annotation:

- Is the variant known?
- Does it hit a gene?
 - Is it in an intron?
 - Does it hit a splice site?

We use different annotation sources and an in-house database.

A selection of SeattleSeq annotation:

- Is the variant known?
- Does it hit a gene?
 - Is it in an intron?
 - Does it hit a splice site?
 - Is it in the coding region?
 - Is there a gain/loss of a stop codon?
 - Does the variant result in a frameshift?
 - ...

We use different annotation sources and an in-house database.

A selection of SeattleSeq annotation:

- Is the variant known?
- Does it hit a gene?
 - Is it in an intron?
 - Does it hit a splice site?
 - Is it in the coding region?
 - Is there a gain/loss of a stop codon?
 - Does the variant result in a frameshift?
 - ...
 - Is it in the 5'/3' UTR of a gene?
 - ...

We use different annotation sources and an in-house database.

A selection of SeattleSeq annotation:

- Is the variant known?
- Does it hit a gene?
 - Is it in an intron?
 - Does it hit a splice site?
 - Is it in the coding region?
 - Is there a gain/loss of a stop codon?
 - Does the variant result in a frameshift?
 - ...
 - Is it in the 5'/3' UTR of a gene?
 - ...
- Is it in a regulatory region?
- ...

Combining all these tools in a pipeline:

```
1 bwa aln -t 8 $reference $i > $i.sai
2 bwa samse $reference $i.sai $i > $i.sam
3 samtools view -bt $reference -o $i.bam $i.sam
```

Listing 1: Shell script

Combining all these tools in a pipeline:

```

1  bwa aln -t 8 $reference $i > $i.sai
2  bwa samse $reference $i.sai $i > $i.sam
3  samtools view -bt $reference -o $i.bam $i.sam

```

Listing 1: Shell script

```

1  %.sai: %.fq
2      $(BWA) aln -t $(THREADS) $(call MKREF, $@) $< > $@
3
4  %.sam: %.sai %.fq
5      $(BWA) samse $(call MKREF, $@) $^ > $@
6
7  %.bam: %.sam
8      $(SAMTOOLS) view -bt $(call MKREF, $@) -o $@ $<

```

Listing 2: Makefile

Description of the data.

Exome sequencing output:

- List of variants (VCF format).
- Regions with sufficient coverage (BED format).

Description of the data.

Exome sequencing output:

- List of variants (VCF format).
- Regions with sufficient coverage (BED format).

We need the coverage information to calculate the frequencies correctly.

- We need to see who *could* have seen a variant.

Advanced annotation and filtering.

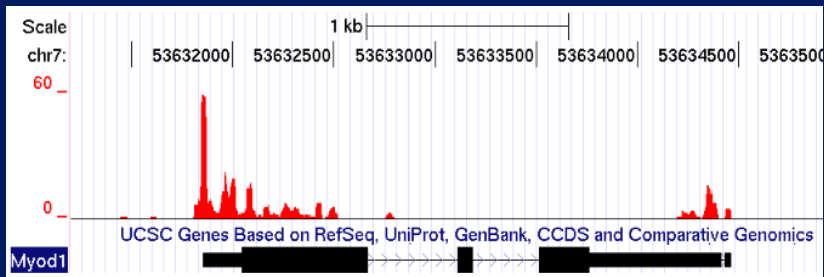


Figure 9: Coverage / depth histogram.

Advanced annotation and filtering.

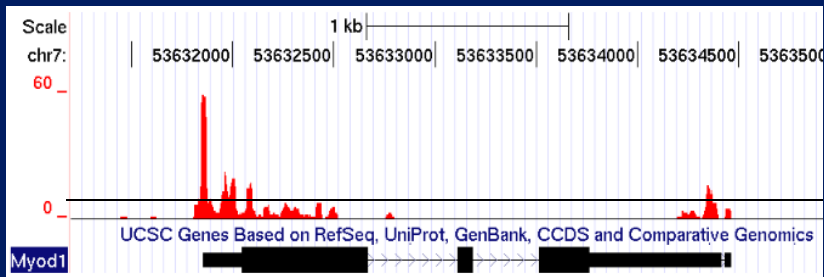


Figure 9: Coverage / depth histogram.

Advanced annotation and filtering.

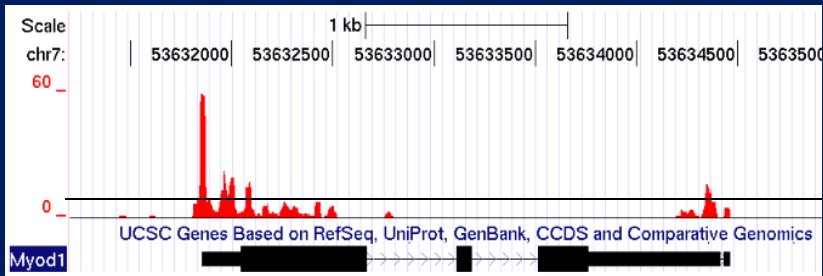


Figure 9: Coverage / depth histogram.



Figure 10: Coverage summary.

Proposed setup:

For privacy reasons:

- No browsing.
- Pre-defined interface.
- Encrypted connection.
- Authentication.

Proposed setup:

For privacy reasons:

- No browsing.
- Pre-defined interface.
- Encrypted connection.
- Authentication.

General design:

- A SOAP interface for interaction with pipelines.
- SOAP over HTTPS (SSL) for encryption.
- Trivial authentication.
 - Send username / password in each call.

Prototype running in Utrecht:

- MySQL database.
- SOAP interface.

Prototype running in Utrecht:

- MySQL database.
- SOAP interface.

The exome data can be enriched with other datasets, not available anywhere else:

- Latest 1000 Genomes Project results.
- Genome of the Netherlands (GoNL).
- RefSeq false positives.

National variant database

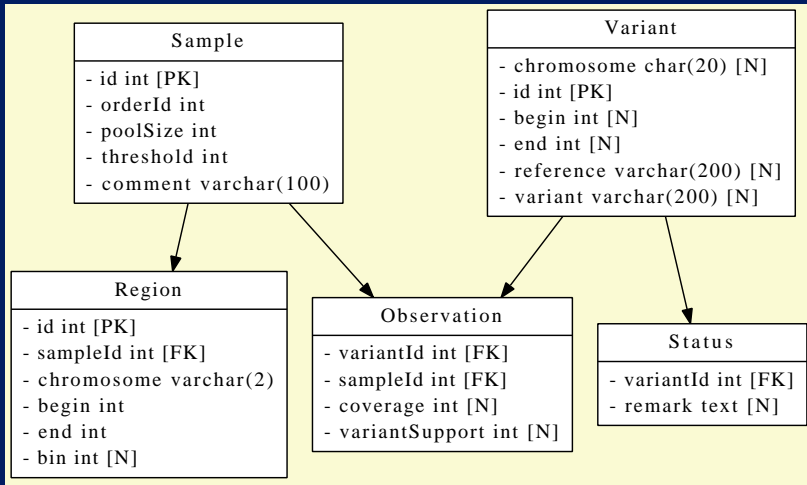


Figure 11: Database schema.

Client perspective:

In principle, one program that handles all:

- Upload a VCF and a BED.
- Server responds with an ID.
- Use the ID to poll for the results.

Client perspective:

In principle, one program that handles all:

- Upload a VCF and a BED.
- Server responds with an ID.
- Use the ID to poll for the results.

Client characteristics:

- Around 30 lines of code.
- Programming language independent.
- Operating system independent.

Server perspective:

Uploading:

- Receive a VCF and a BED.
- Respond with an ID.
- A process is forked that starts annotation and importing.

Server perspective:

Uploading:

- Receive a VCF and a BED.
- Respond with an ID.
- A process is forked that starts annotation and importing.

Annotation:

- Receive an ID.
- If the annotation for this ID is finished, return the result.

Server perspective:

Uploading:

- Receive a VCF and a BED.
- Respond with an ID.
- A process is forked that starts annotation and importing.

Annotation:

- Receive an ID.
- If the annotation for this ID is finished, return the result.

Note that if the annotation part is finished, the result can already be returned.

Column	Name	Description
1	CHROM	Name of the chromosome.
2	POS	Position on the chromosome.
3	REF	Reference allele.
4	ALT	Genotype.
5	Frequency	Frequency of this genotype.
6	Occurrences	Number of hits.
7	IDs (> 100)	If the frequency is low: a list of sample IDs.
8	IDs (< 100)	A list with hits in reserved samples; 1000 genomes, GoNL, HGMD, etc.

Table 1: Output format.

Stakeholder meeting.

Stakeholder meeting.

- The proposed interface.

Stakeholder meeting.

- The proposed interface.
- More sensitive issues.
 - Who gets access?

Stakeholder meeting.

- The proposed interface.
- More sensitive issues.
 - Who gets access?
 - What if someone submits the same data twice?

Stakeholder meeting.

- The proposed interface.
- More sensitive issues.
 - Who gets access?
 - What if someone submits the same data twice?
 - Can we get all variants in a gene?
 - Maybe only for public samples.

Acknowledgements:

Christian Gilissen
Victor Guryev
Frans-Paul Ruzius
Nienke Wieskamp
Hailiang Mei
Johan den Dunnen

<https://www.mutalyzer.nl/svn/NGSdata/branches/dvd/>