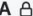# Analysis of long read data at the Flemish Super Computer

Erika Souche

Laboratory for Cytogenetics and Genome Research

# Laboratory for Cytogenetics and Genome Research

- Aims
  - Implement genomic technologies to improve genetic diagnostic testing
  - Map the causes and mechanisms underlying chromosomal rearrangements
  - Improve preimplantation, prenatal and postnatal genetic testing of rare diseases
  - Use liquid biopsies as a biomarker to map genic and non-genic diseases

# One nextflow pipeline



Center for Human Genetics

# Use of external software

# Variant calling from mapping

Long reads

Reference genome

Small variants
- Single Nucleotide Variants (SNVs)
- Indels

Repeat expansions

Structural Variants (SVs)
- Insertions
- Deletions
- Duplications
- Inversions
- Translocations

# Variant calling from *de novo* assembly



Long reads

Contig 1

Contig 2

Contig 3

Contig 4

Structural Variants (SVs)
- Insertions
- Deletions
- Duplications
- Inversions
- Translocations

# Project 1: Developmental disorders

- Whole Exome Sequencing / Whole Genome Sequencing (WGS)

    ~ 10-15 % Copy Number Variants (CNVs)

    ~ 25-30 % SNVs & indels

# Project 1: Developmental disorders

- Whole Exome Sequencing / Whole Genome Sequencing (WGS)
    - ~ 10-15 % Copy Number Variants (CNVs)
    - ~ 25-30 % SNVs & indels

- ONT sequencing of unsolved trios
    - 27 trios
    - Structural variants (SVs)

# SVs

~ 23,000 SVs per individual

# SVs

~ 23,000 SVs per individual

# SVs

~ 23,000 SVs per individual

~ 100 *de novo* SVs

# SVs

~ 23,000 SVs per individual

~ 100 *de novo* SVs

~ 7 rare *de novo* SVs

# SVs

~ 23,000 SVs per individual

~ 100 *de novo* SVs

~ 7 rare *de novo* SVs



Patient

Father

Mother

# Complex SVs

- Inversion on chrX in a male patient

# Complex SVs

- *de novo* assembly



> 91.3% of hg38

401st contig:
3 099 922 541, Reference
**2 832 815 483, DevDis_P23_flye_assembly_polished_racon**

# Complex SVs

- *de novo* assembly

- Mapping of contigs to hg38

# Complex SVs

- *de novo* assembly
- Mapping of contigs to hg38

# Complex SVs

- SV locus not covered by *de novo* assembly

# Project 1: Developmental disorders

- ONT sequencing of unsolved trios
  - Structural variants (SVs)
  - Repeat expansions (STRs)
  - Epigenic modifications
  - Mosaïcism
  - …

# Project 1: Developmental disorders

- ONT sequencing of unsolved trios
  - Structural variants (SVs)
  - Repeat expansions (STRs)
  - Epigenic modifications
  - Mosaïcism
  - …

- Need for longer-term storage
for active data

# Project 1: Developmental disorders

- ONT sequencing of unsolved trios
  - Structural variants (SVs)
  - Repeat expansions (STRs)
  - Epigenic modifications
  - Mosaïcism
  - …

- Need for longer-term storage for active data

Tier1 data

# Tier1 data

- Active data

- Data workflow

- Data structure

=> More sustainable for research laboratories

# Tier1 data

- Metadata

- 4-year plan
  - Number of samples
  - Resource estimation

- Exit strategy

# Resource monitoring

# Project 2: Episignature

- Neurodevelopmental disorders with methylation disturbances



Illumina Infinium/Epic array
disease A
control
disease B
Aref-Esghi *et al.* 2020

Talk "Episignature in patients with developmental disorders" by Benjamin Huremagic

# Project 2: Episignature

- Neurodevelopmental disorders with methylation disturbances

- Proof of concept
  - 20 patients
  - 40 control samples



**Illumina Infinium/Epic array**

disease A

control

disease B

Aref-Esghi *et al.* 2020

**Oxford Nanopore sequencing**

Current study

Talk "Episignature in patients with developmental disorders" by Benjamin Huremagic

# Project 2: Episignature

- Neurodevelopmental disorders with methylation disturbances

- Proof of concept
  - 20 patients
  - 40 control samples



Illumina Infinium/Epic array

disease A

control

disease B

Aref-Esghi *et al.* 2020

Oxford Nanopore sequencing

Current study

- Extend model

  => Processing of 210 samples

Talk "Episignature in patients with developmental disorders" by Benjamin Huremagic

KU LEUVEN

# Project 2: Episignature

- Neurodevelopmental disorders with methylation disturbances

- Proof of concept
  - 20 patients
  - 40 control samples

- Extend model
  => Processing of 210 samples



Illumina Infinium/Epic array

disease A

control

disease B

Aref-Esghi *et al.* 2020

Oxford Nanopore sequencing

Current study

Need GPU processing to get methylation calls

Talk "Episignature in patients with developmental disorders" by Benjamin Huremagic
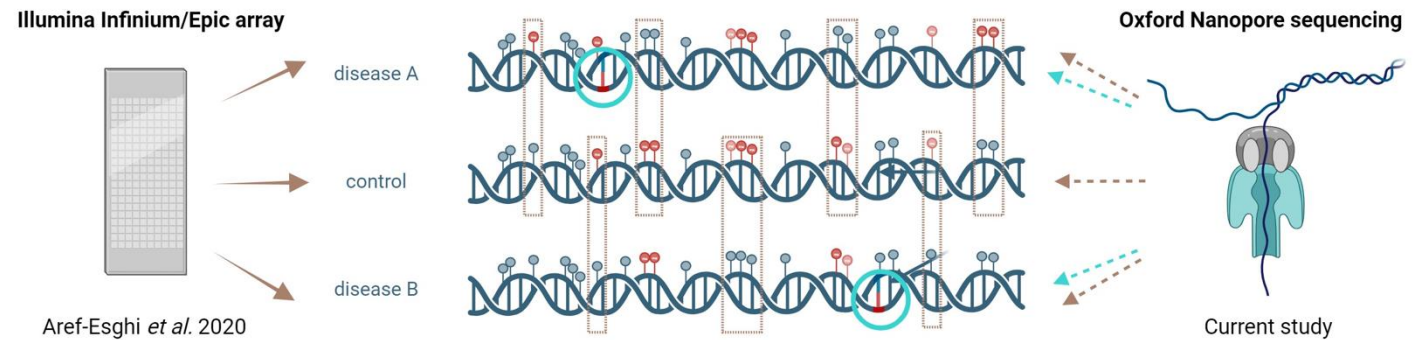
# Project 2: Episignature

- Neurodevelopmental disorders with methylation disturbances

- Proof of concept
  - 20 patients
  - 40 control samples



Illumina Infinium/Epic array

disease A

control

disease B

Aref-Esghi *et al.* 2020

Oxford Nanopore sequencing

Current study

Need GPU processing to get methylation calls

- Extend model
  => Processing of 210 samples

Tier1 compute

Talk "Episignature in patients with developmental disorders" by Benjamin Huremagic

# Tier1 compute

- Pipeline description

```
   FAST5  ──pod5──▶   POD5   ──dorado──▶   BAM
   files               files                file
```

- Resource optimisation
  => Starting grant

# Resource optimisation

- FAST5 to POD5

| Number of nodes | Total number of cores | Wall clock time (s) | Speed-up (w.r.t. baseline) | Efficiency |
|---|---|---|---|---|
| 1 | 1 | 16128 | Not applicable | Not applicable |
| 1 | 2 | 12566 | 1 | 1 |
| 1 | 4 | 6778 | 1.85 | 0.93 |
| 1 | 8 | 2708 | 4.64 | 1.16 |
| 1 | 16 | 2407 | 5.22 | 0.65 |
| 1 | 32 | 1445 | 8.70 | 0.54 |
| 1 | 64 | 1449 | 8.67 | 0.27 |
| 1 | 128 | 1548 | 8.12 | 0.13 |

# Resource estimation

- Based on
  - Results of efficiency tests
  - Number of samples to process

### Tier1 compute

9,240 CPU hours
8,316 GPU hours
21 TB scratch volume

# Resource monitoring

# Acknowledgements

Mathilde Geysens

Benjamin Huremagic

Chiara Campanelli

Greet Peeters

Senne Meynants

Natalia Olszewska

Kris Van Den Bogaert

Joris Vermeesch

Qiang Fu

Tatjana Jatsenko

Marta Sousa Santos

Kate Elizabeth Stanley

Olga Tsuiko

Stefania Tuveri

Yan Zhao

Jonas Demeulemeester

KU LEUVEN

# Tools used

**Mapping**

Minimap2

**Quality control**

mosdepth

NanoPlot

QUAST

*de novo* assembly

shasta

racon

flye

hifiasm

Verkko

HapDup

**Structural variant calling**

PBSV

SVIM

Sniffles2

cuteSV

QDNAseq

DipDiff

SURVIVOR

**Repeat expansion calling**

straglr
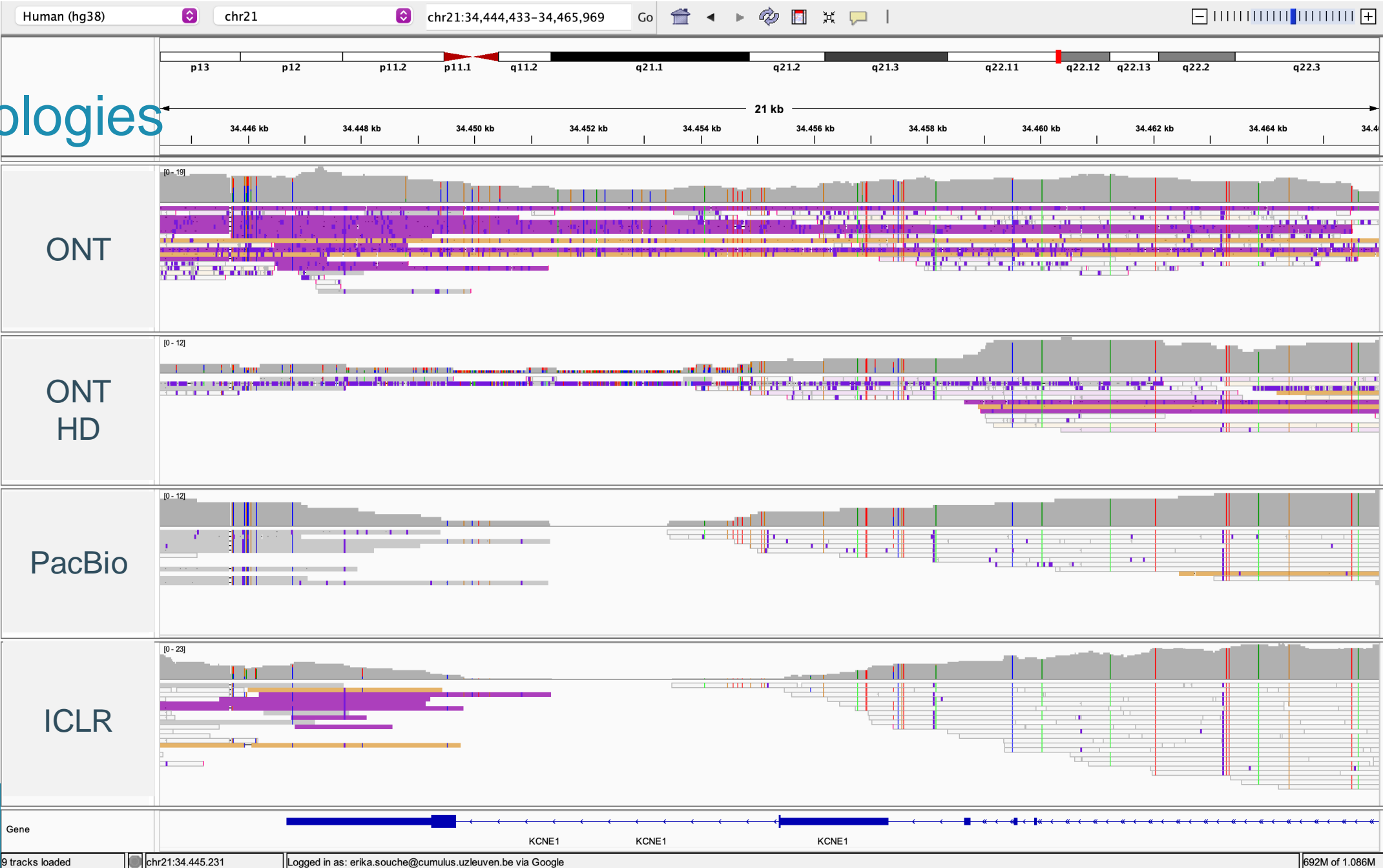
TRGT

**Small variant calling**

Clair3

DeepVariant

WhatsHap

# Data footprint

| | Format | ONT | PacBio | ICLR |
|---|---|---|---|---|
| Raw data | FAST5/POD5<br>FASTQ.gz | 1 TB<br>100 GB | -<br>50 GB | -<br>40 GB * |
| Aligned reads | BAM<br>CRAM | 100 GB<br>50 GB | 50 GB<br>15 GB | 25 GB<br>5 GB |
| Small variants | VCF.gz | 100 MB | 100 MB | 100 MB |
| Structural variants | VCF | 50 MB | 50 MB | 50 MB |
| Repeat expansions | VCF | 10 MB | 10 MB | 10 MB |
| *De novo* assembly | FASTA | 3 GB | 3 GB | 3 GB |

\* After processing by Illumina
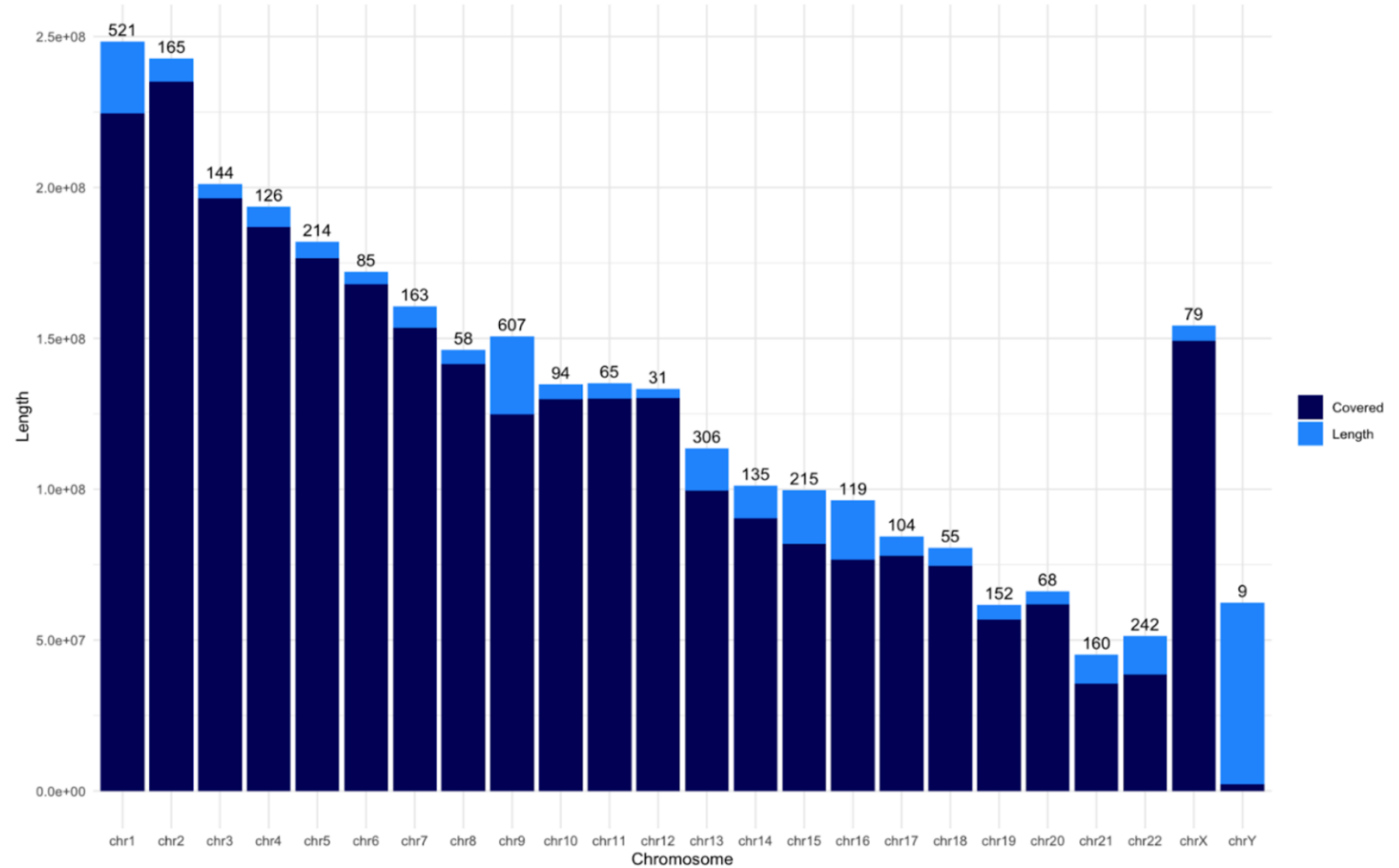
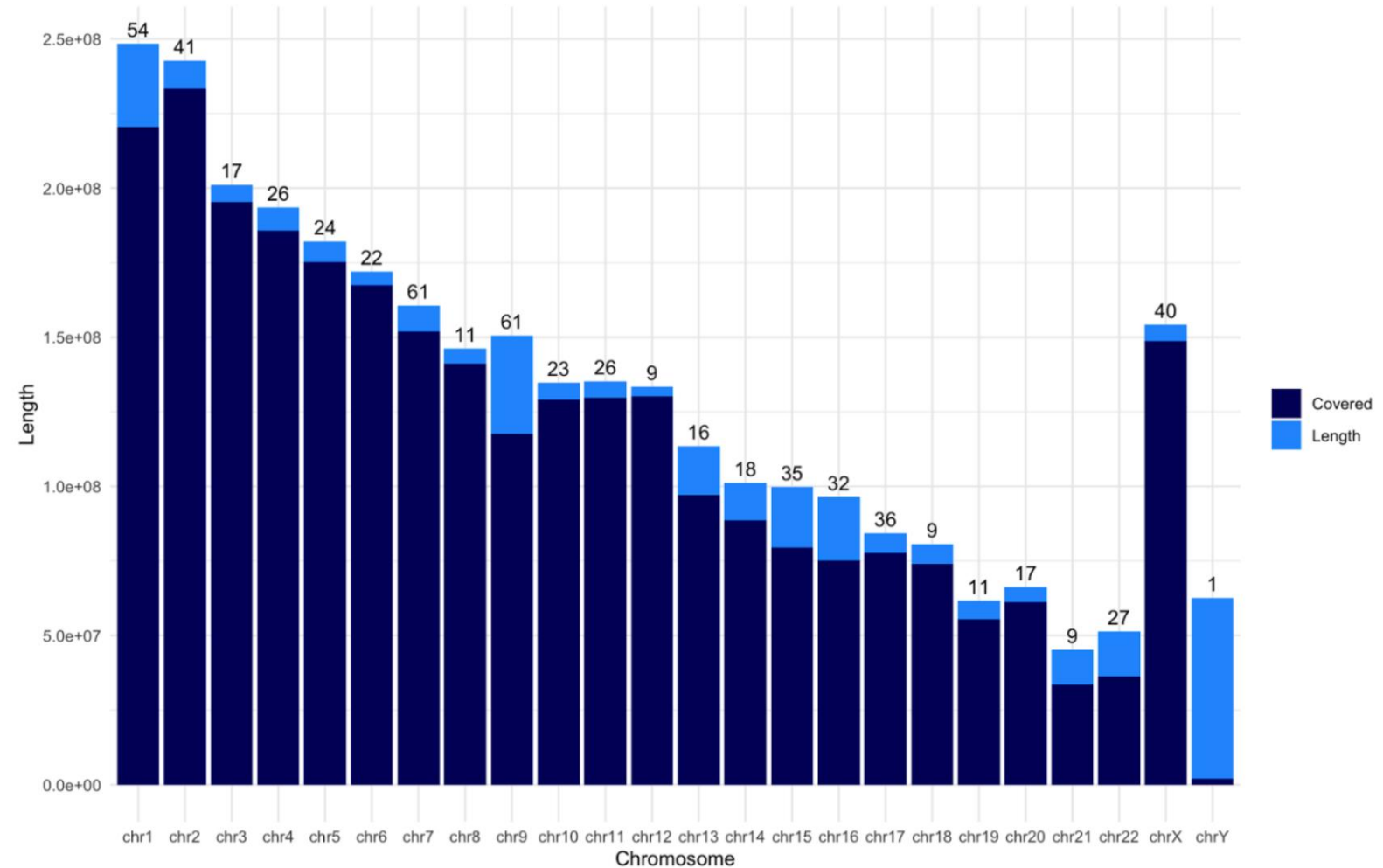KU LEUVEN

# All technologies

All technologies

# *de novo* assembly
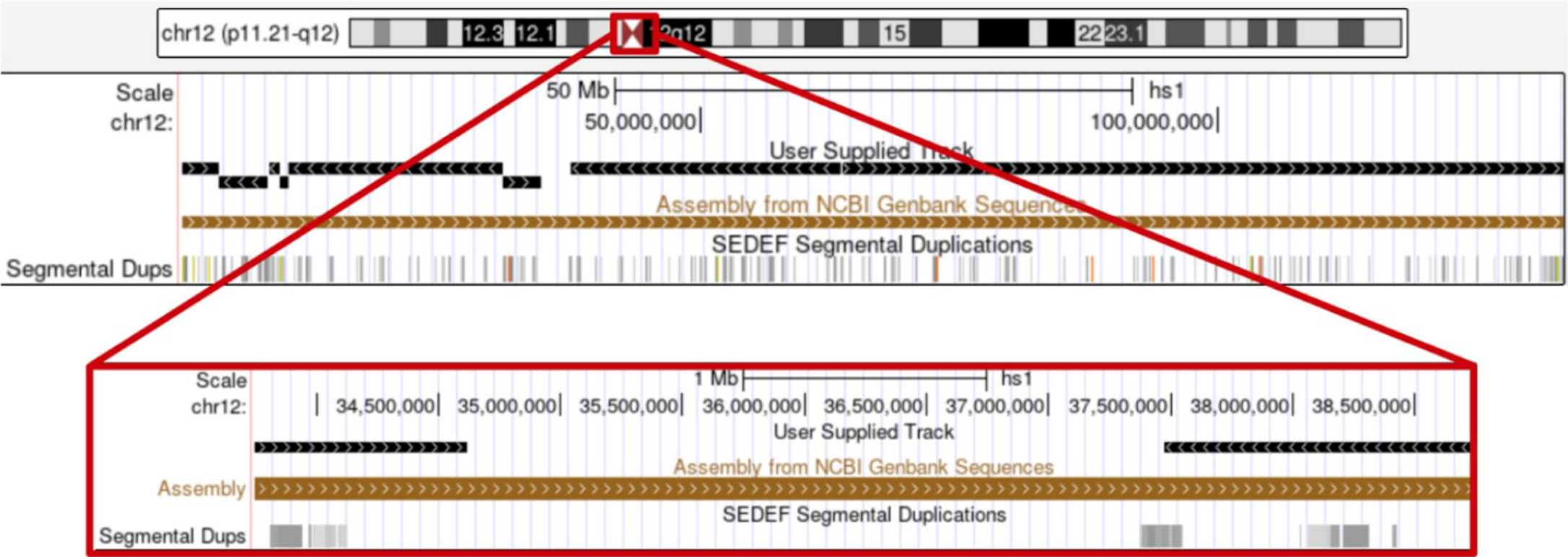
- 2,727 contigs
- ~86.7% of T2T covered

# *de novo* assembly

- 626 contigs > 100 kb
- ~86.67% of T2T covered

# *de novo* assembly

# *de novo* assembly