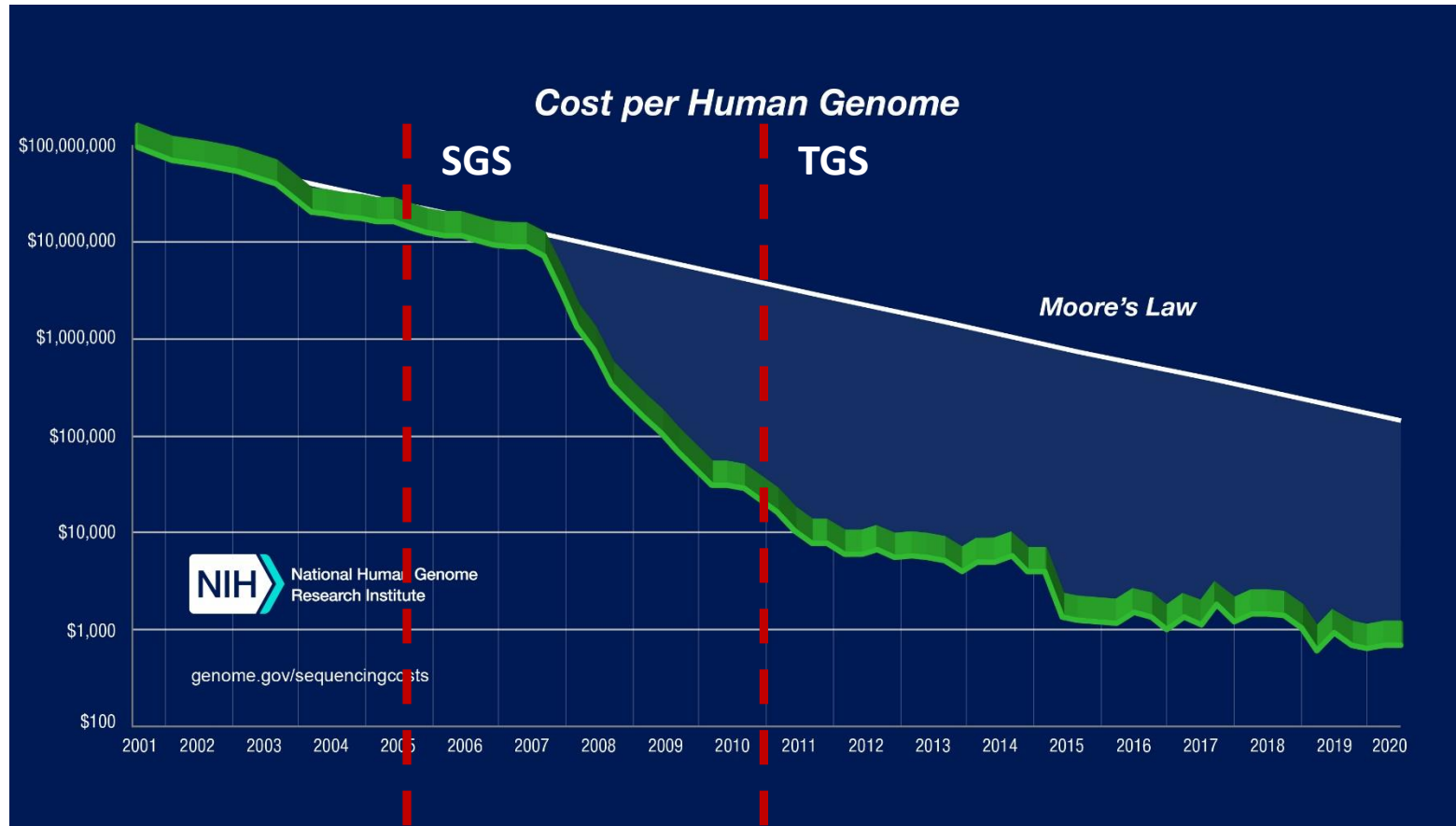# PGB2022

Dr. Agnieszka A. Golicz

agnieszka.golicz@agrar.uni-giessen.de
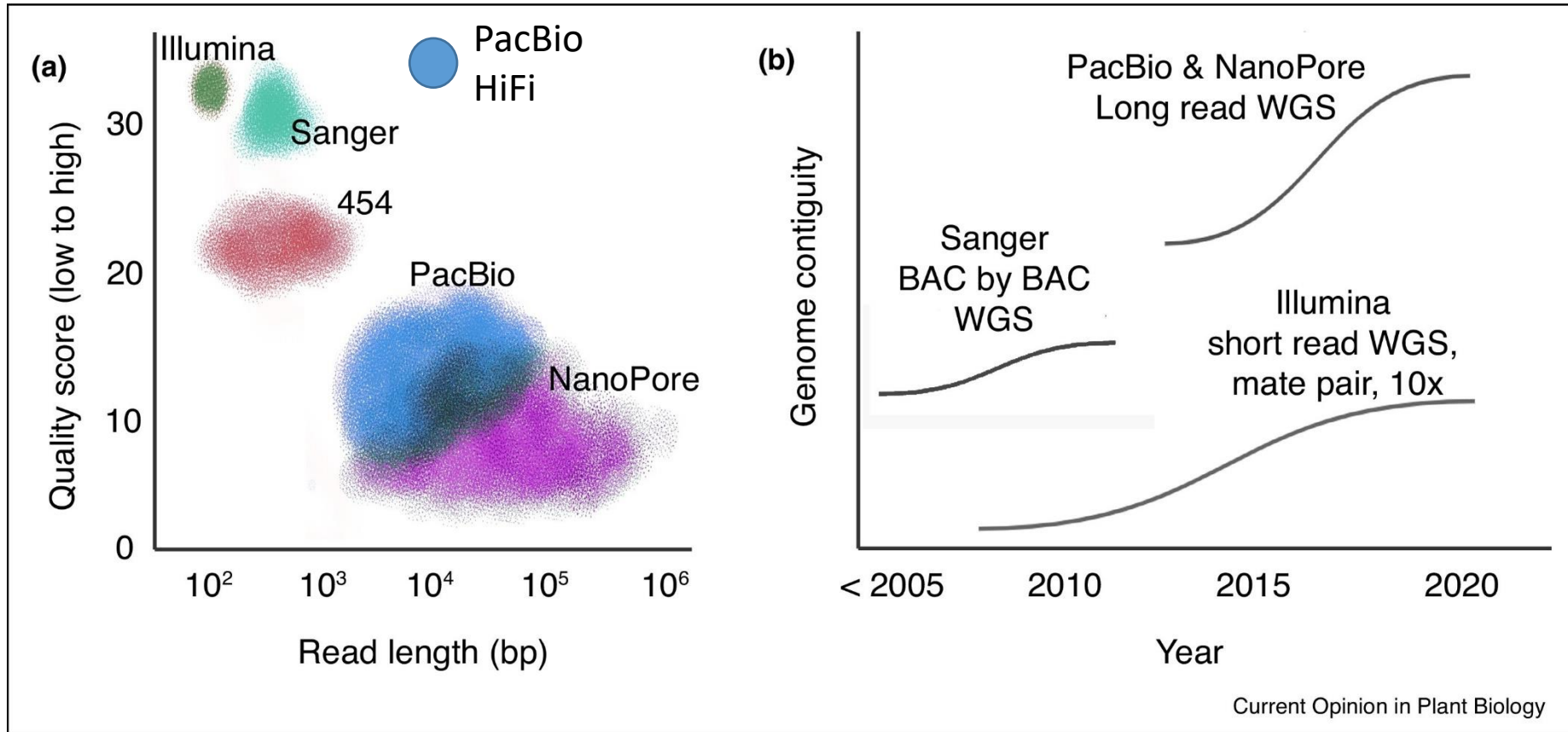
# Genome sequencing technology

Sanger sequencing
- Older technology
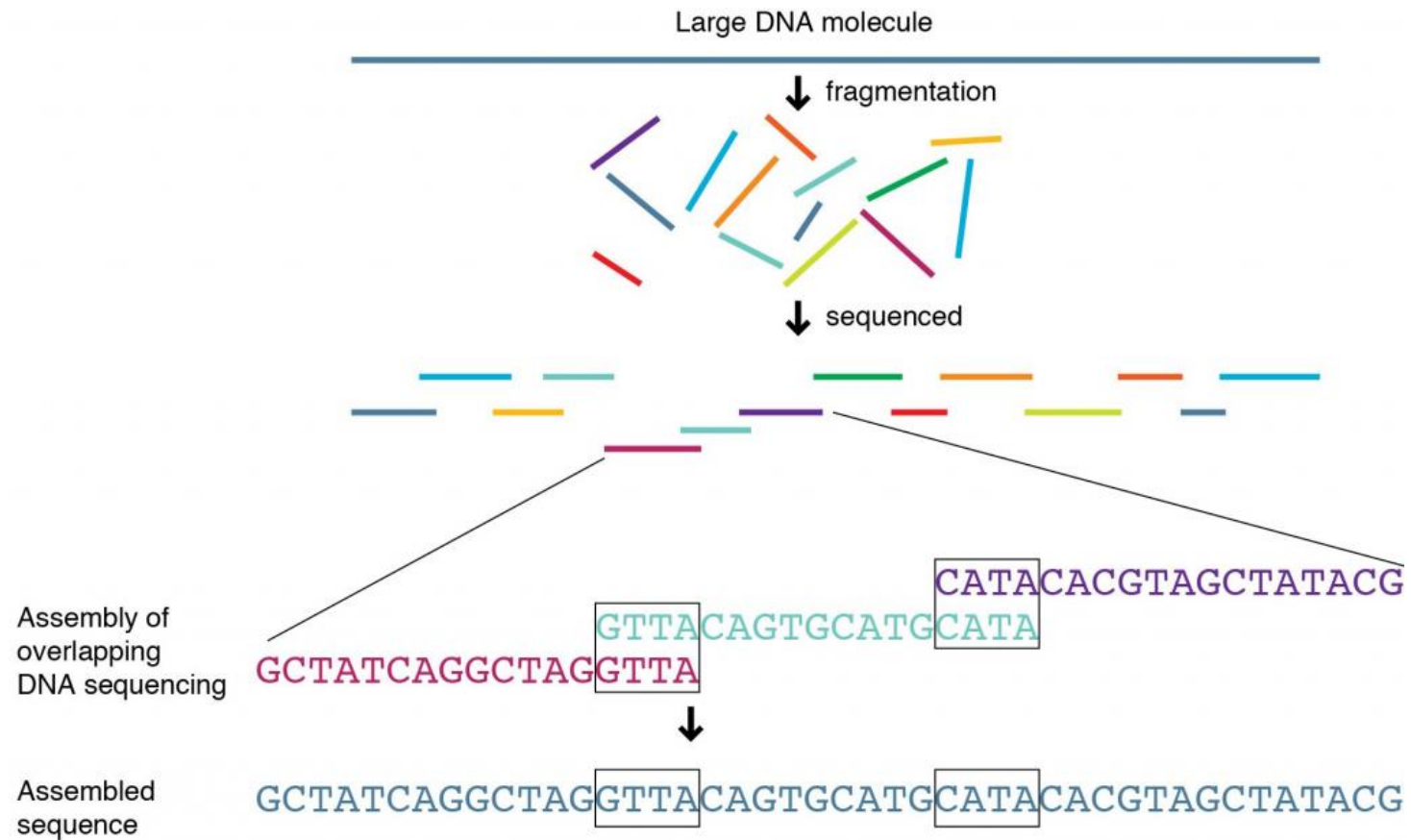- Reads ~1000 bp
- Very accurate

Second generation sequencing (Illumina)
- Reads ~200 bp
- More error prone
- But cheaper!

**Cost per Human Genome**

SGS

TGS

*Moore's Law*

$100,000,000

$10,000,000

$1,000,000

$100,000

$10,000

$1,000

$100

2001 2002 2003 2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020

NIH National Human Genome Research Institute

genome.gov/sequencingcosts

Technological limitation – we can not sequence (read) the whole chromosome at once
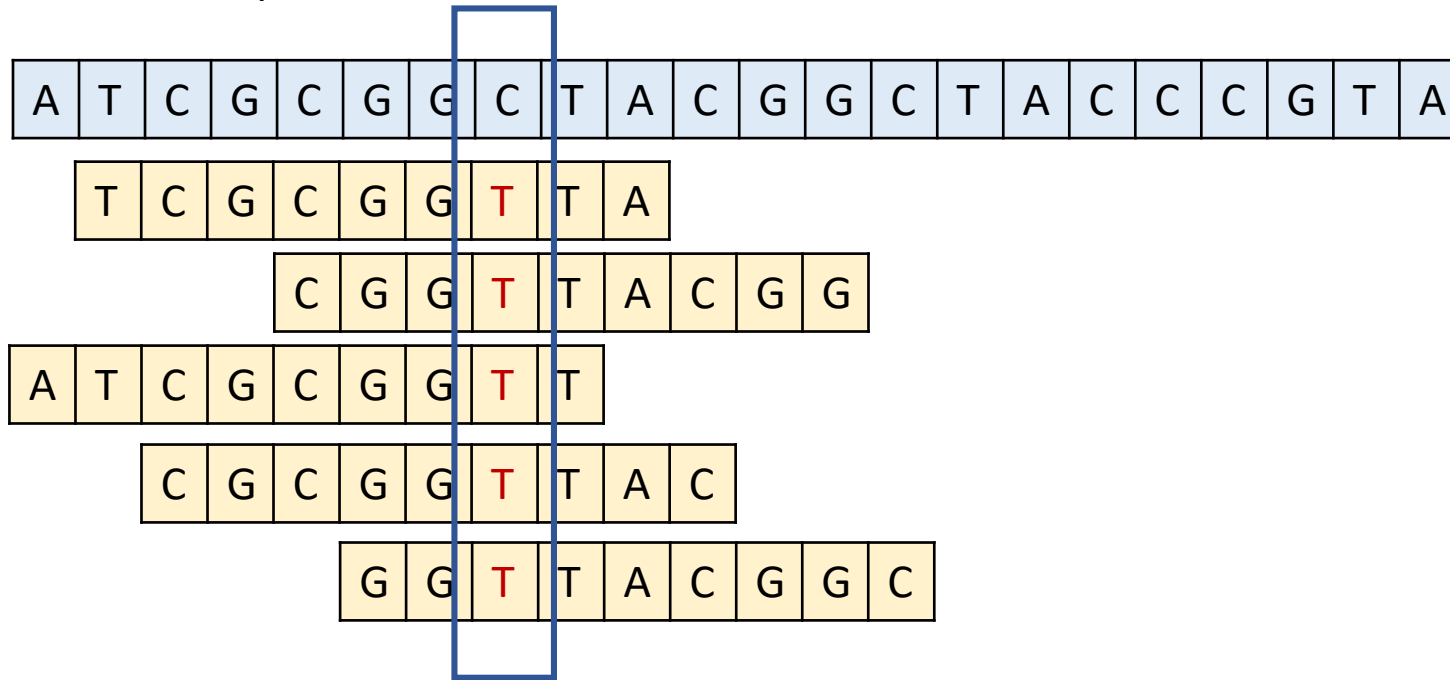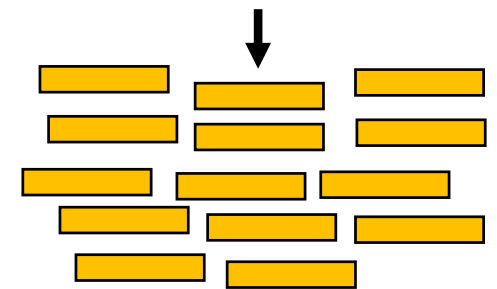Sequencing read – a short sequence representing a fragment of the DNA

https://www.genome.gov/about-genomics/fact-sheets/DNA-Sequencing-Costs-Data

# Genome sequencing

# Read mapping

Reference sequence

| A | T | C | G | C | G | G | C | T | A | C | G | G | C | T | A | C | C | C | G | T | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

| T | C | G | C | G | G | T | T | A |
|---|---|---|---|---|---|---|---|---|

| C | G | G | T | T | A | C | G | G |
|---|---|---|---|---|---|---|---|---|

| A | T | C | G | C | G | G | T | T |
|---|---|---|---|---|---|---|---|---|

| C | G | C | G | G | T | T | A | C |
|---|---|---|---|---|---|---|---|---|

| G | G | T | T | A | C | G | G | C |
|---|---|---|---|---|---|---|---|---|

Resequencing individuals

# Read mapping

# Bioinformatics file formats

➢ Dealing with large quantities of data
➢ Automated processes
➢ Require standardized file formats
    ➢ FASTA
    ➢ FASTQ
    ➢ SAM/BAM
    ➢ VCF
    ➢ GFF

# FASTA file format

```
>Sequence_1 assembly1
CCCTAAACCCTAAACCCTAAACCCTAAACCTCTGAATCCTTAATCCCTAAATCCCTAAAT
CTTTAAATCCTACATCCATGAATCCCTAAATACCTAATTCCCTAAACCCGAAACCGGTTT
CTCTGGTTGAAAATCATTGTGTATATAATGATAATTTTATCGTTTTTATGTAATTGCTTA
TTGTTGTGTGTAGATTTTTTAAAAATATCATTTGAGGTCAATACAAATCCTATTTCTTGT
GGTTTTCTTTCCTTCACTTAGCTATGGATGGTTTATCTTCATTTGTTATATTGGATACAA
GCTTTGCTACGATCTACATTTGGGAATGTGAGTCTCTTATTGTAACCTTAGGGTTGGTTT
ATCTCAAGAATCTTATTAATTGTTTGGACTGTTTATGTTTGGACATTTATTGTCATTCTT
>Sequence_2
CCCTAAACCCTAAACCCTAAACCCTAAACCTCTGAATCCTTAATCCCTAAATCCCTAAAT
CTTTAAATCCTACATCCATGAATCCCTAAATACCTAATTCCCTAAACCCGAAACCGGTTT
CTCTGGTTGAAAATCATTGTGTATATAATGATAATTTTATCGTTTTTATGTAATTGCTTA
TTGTTGTGTGTAGATTTTTTAAAAATATCATTTGAGGTCAATACAAATCCTATTTCTTGT
GGTTTTCTTTCCTTCACTTAGCTATGGATGGTTTATCTTCATTTGTTATATTGGATACAA
GCTTTGCTACGATCTACATTTGGGAATGTGAGTCTCTTATTGTAACCTTAGGGTTGGTTT
ATCTCAAGAATCTTATTAATTGTTTGGACTGTTTATGTTTGGACATTTATTGTCATTCTT
```
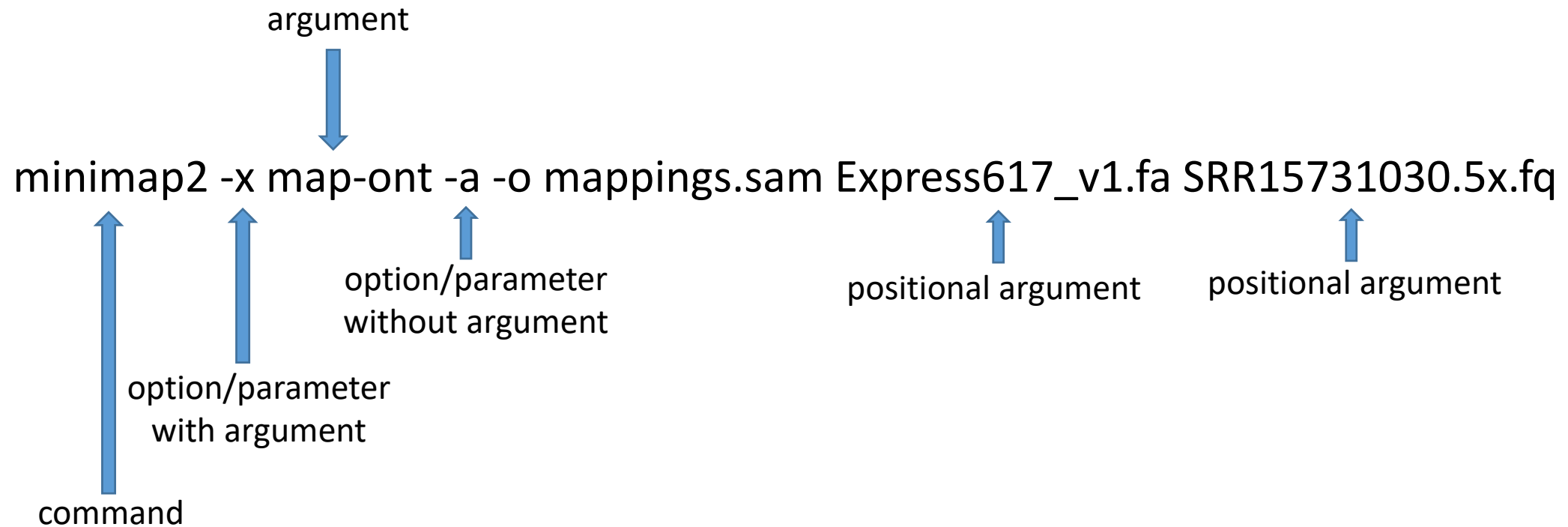
# FASTQ file format

Identifier ———————— `@HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1`

Sequence ———————— `TTAATTGGTAAATAAATCTCCTAATAGCTTAGATNTTACCTTNNNNNNNNNNTAGTTTCTTGAGA`

+ sign & identifier— `+HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1`

Quality scores ———— `efcfffffcfeefffcffffffddf`feed]`]_Ba_^__[YBBBBBBBBBRTT\]][]dddd``
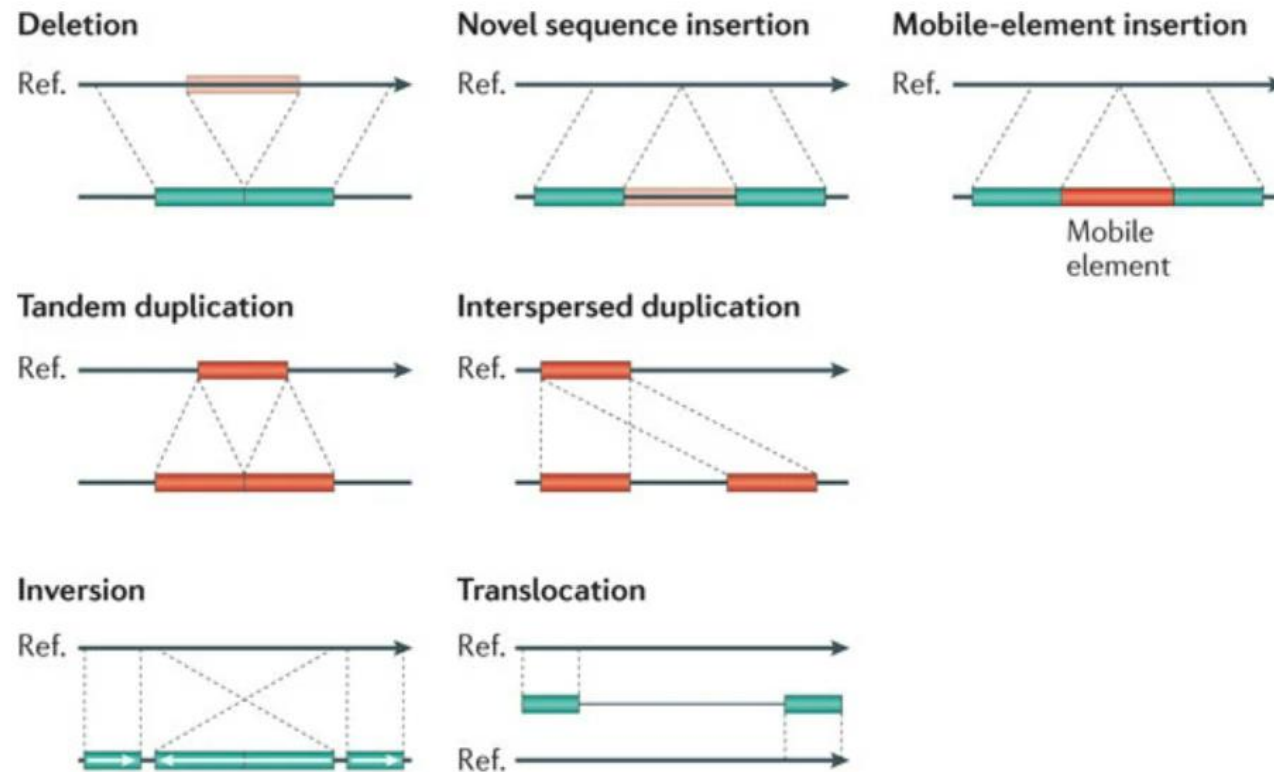
Base T
phred Quality  ] = 29

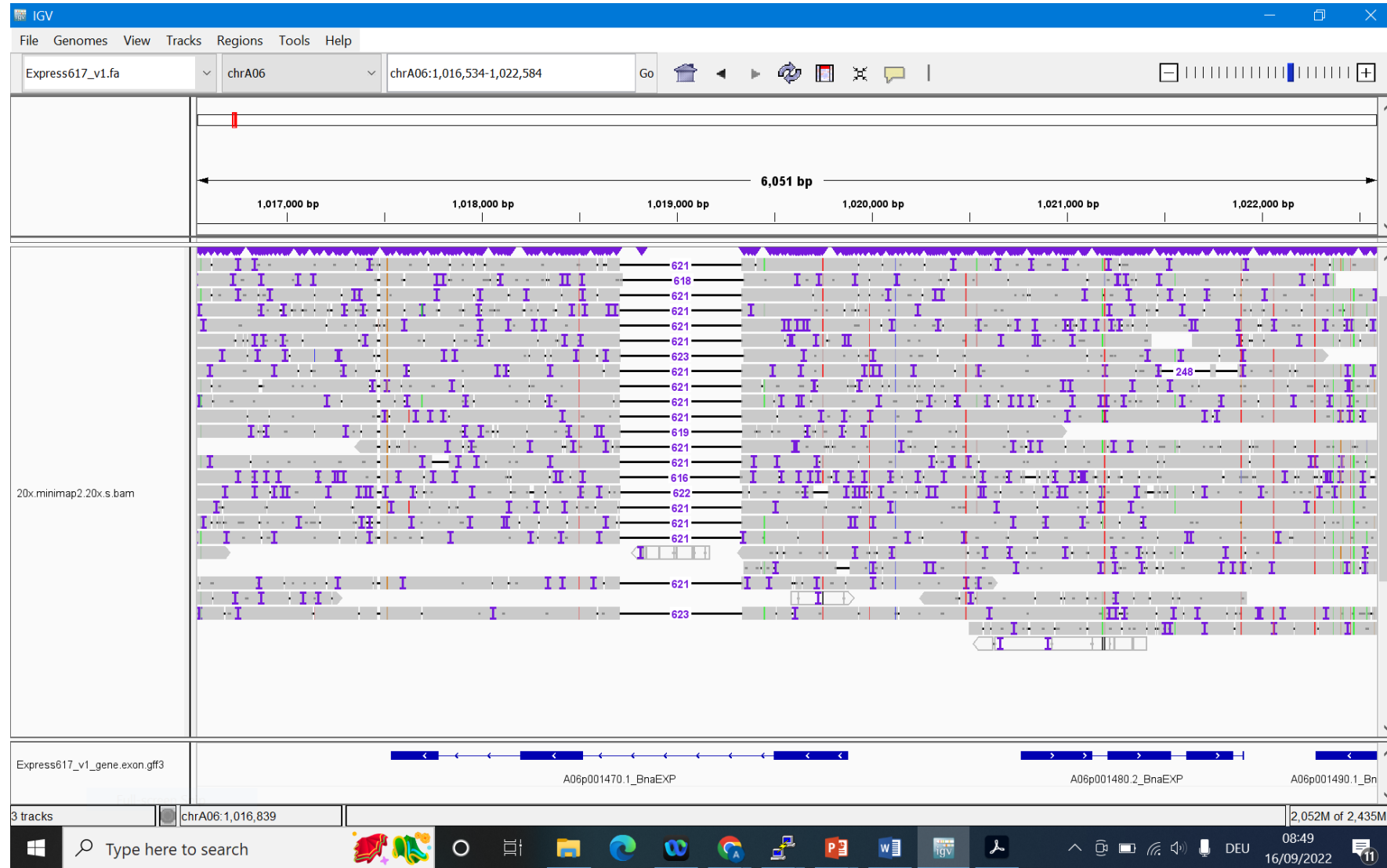# Using programs with command line options and arguments

argument

minimap2 -x map-ont -a -o mappings.sam Express617_v1.fa SRR15731030.5x.fq

option/parameter without argument

positional argument

positional argument

option/parameter with argument

command

# Comparing genomes – sequence variants

# Types of structural variants



Nature Reviews | Genetics
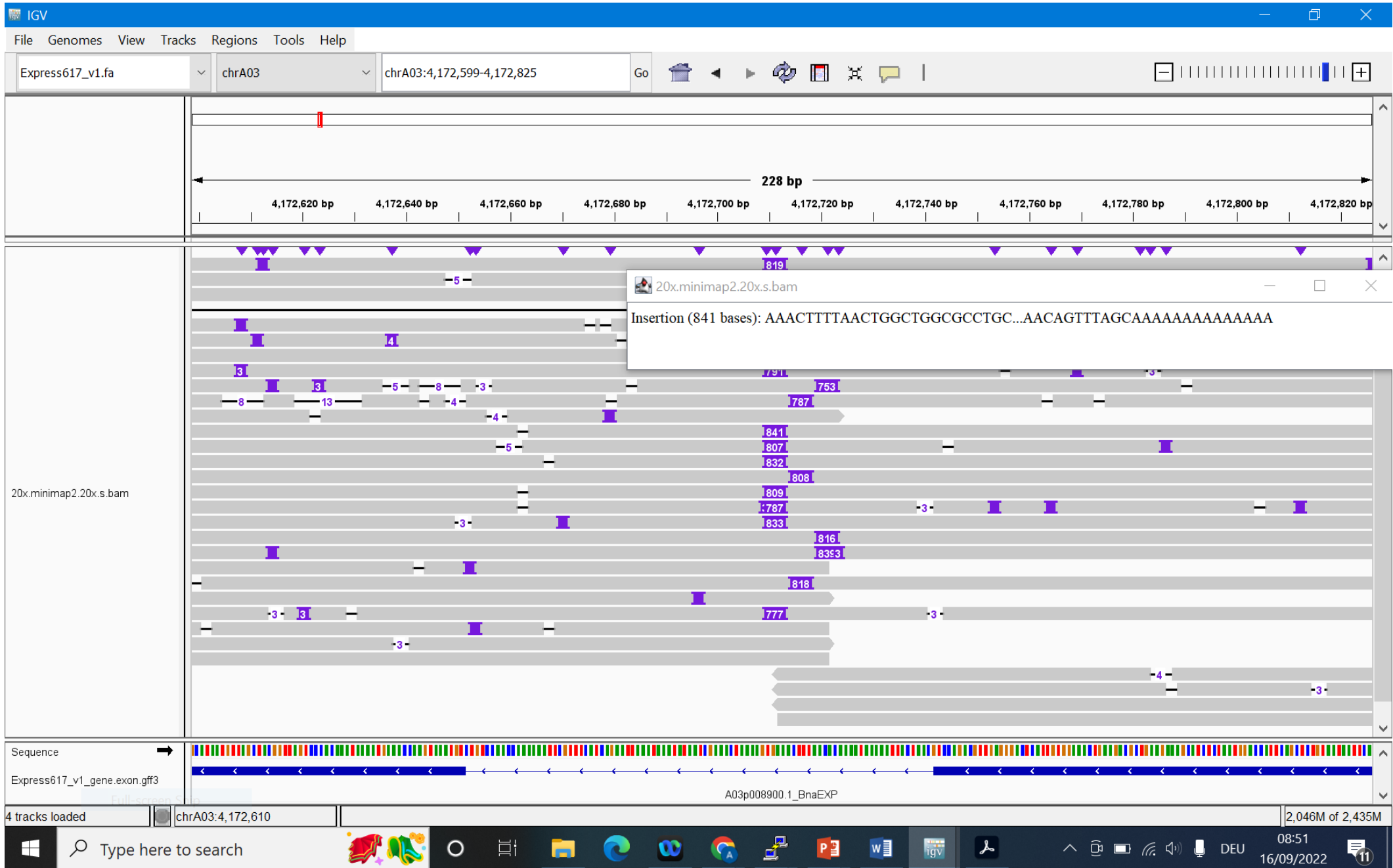
# Looking for structural variants with long reads

# Structural variant report

Pick one variant from SV.mRNA.overlap.tsv
Please note the ID of the variant and send it to me including your name in the email.

1. Short introduction - what are structural variants. Why are they important?
2. How was the structural variant identified?
3. Describe the structural variant. Include an image from IGV.
4. What are the possible consequences of this structural variant? Does it interrupt a gene, which part of the gene?

Please include at least 5-10 references.