# Assessing efficiency of third generation sequencing read correction

Medhat Mahmoud
Department of Protein Biosynthesis (IBCh)
&
Department of Computational Biology (AMU)
Poznan 24-05-2017

# The goal

Compare sensitivity, specificity and performance of Single molecule real-time  (SMRT) PacBio reads correction tools

• • •

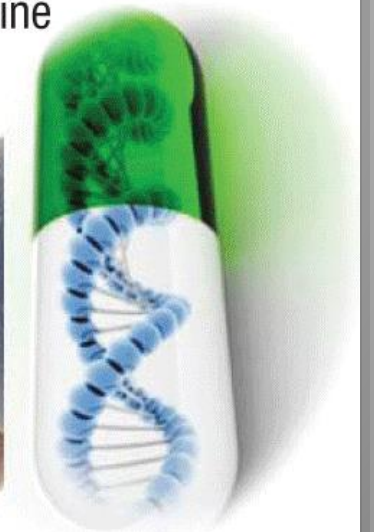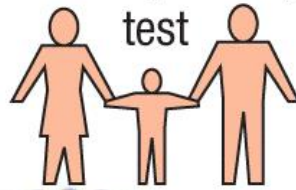Genetic ancestry test

Genetic paternity test

Genetic compatibility test

DNA

Genetic fingerprinting

%

Genetic disease risk

Personalized medicine

# Large DNA molecule

fragmentation

sequenced

Assembly of
overlapping
DNA sequencing

CATACACGTAGCTATACG

GTTACAGTGCATGCATA

GCTATCAGGCTAGGTTA

Assembled
sequence

GCTATCAGGCTAGGTTACAGTGCATGCATACACGTAGCTATACG
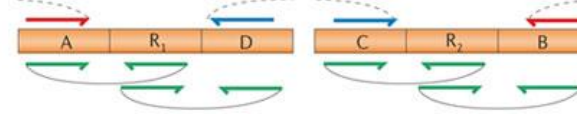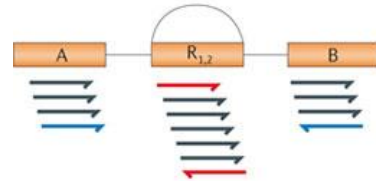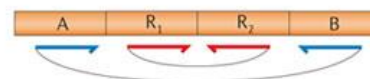
Aa Assembly graph

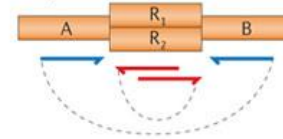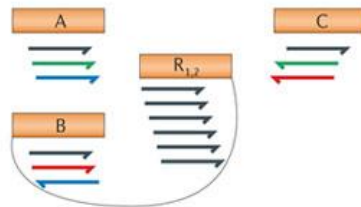Ab Correct assembly

Ac Misassembly

Ba Assembly graph

Bb Correct assembly
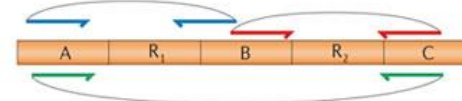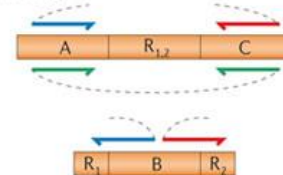
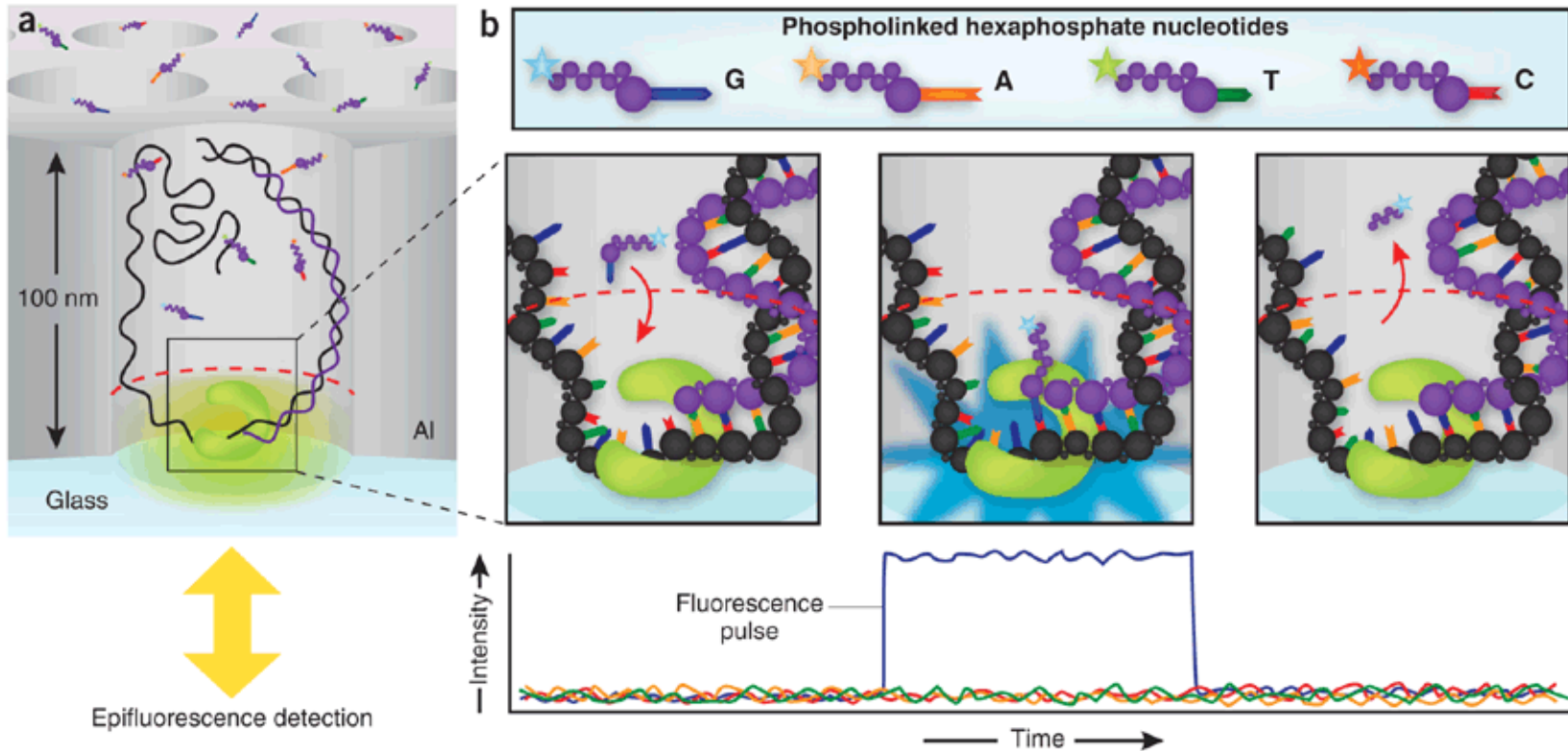Bc Misassembly

Ca Assembly graph

Cb Correct assembly

Cc Misassembly

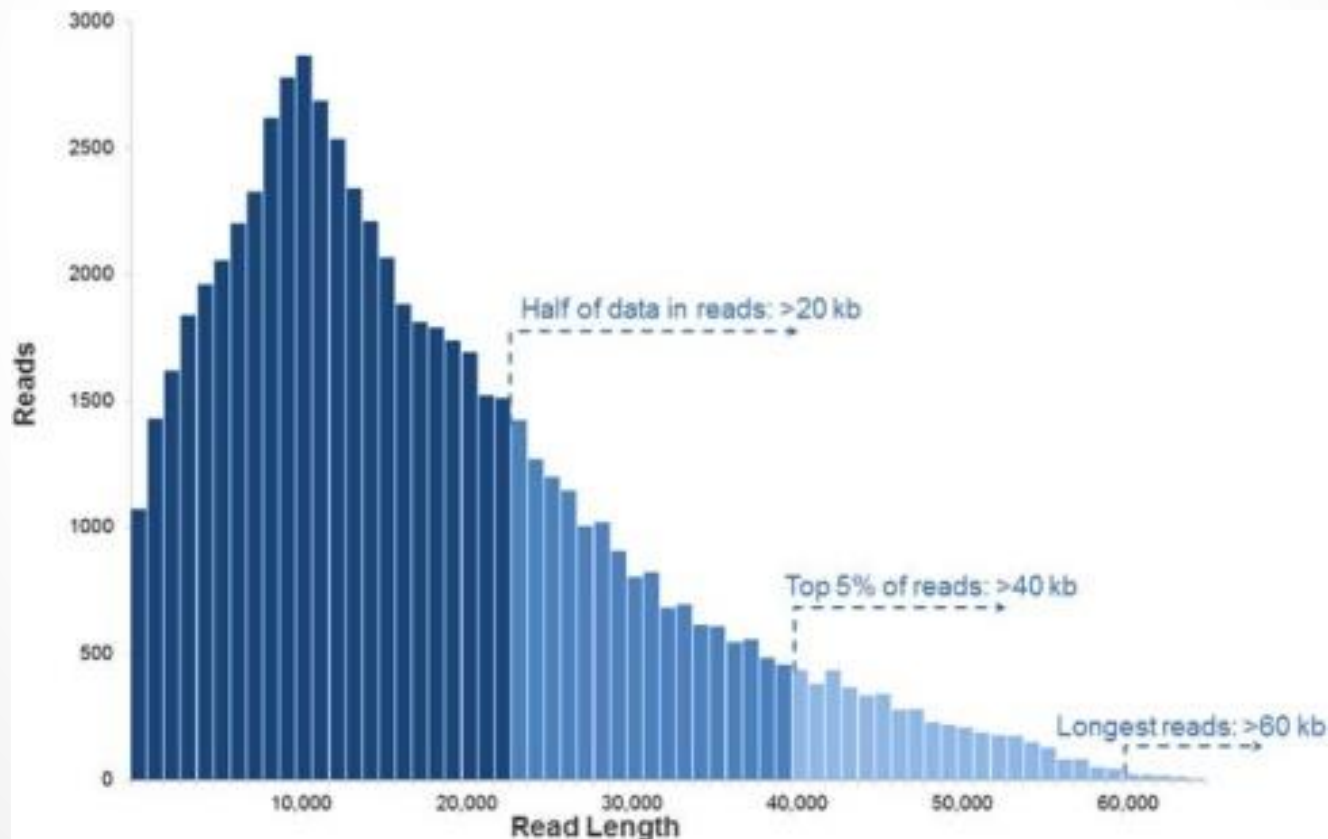Nature Reviews | Genetics

# SMRT PacBio sequencing

● ● ●

# The idea

# PacBio read advantages

- offers much longer read lengths and faster runs than SGS (reaches 60,000 bp)
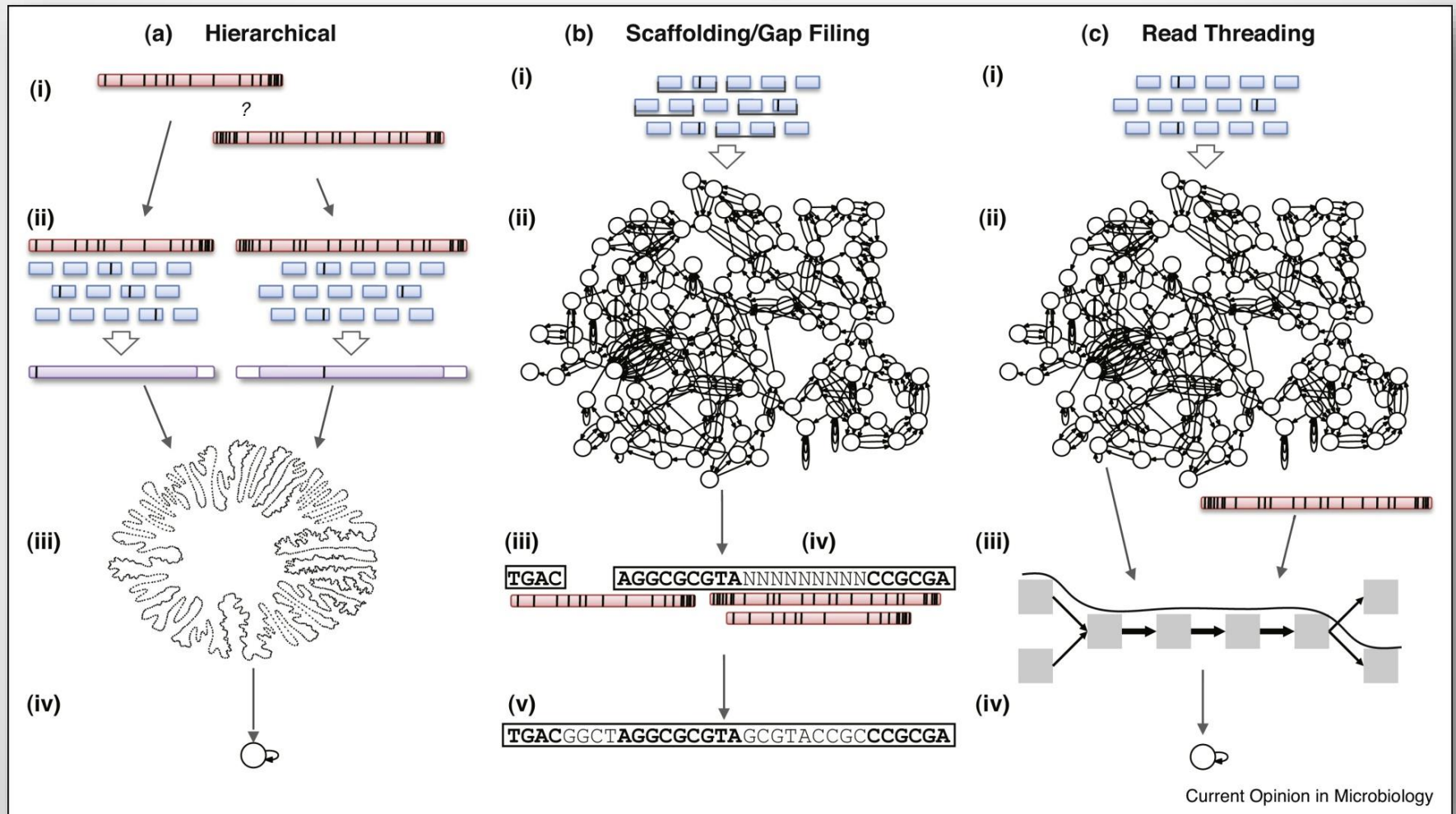
repetitive structural mutations variation hybrid detect sequencing gaps novel assemblies gene novo isoforms annotated close regions methylation resolve

# Where is the catch ?!

• • •

higher error rate, and higher cost per base

# PacBio assembly usage example



Current Opinion in Microbiology

# Which tool to choose?

• • •

# Assessment steps

**Organism selection**
- Human
- Rice
- Yeast
- Trypanosoma
- E.coli

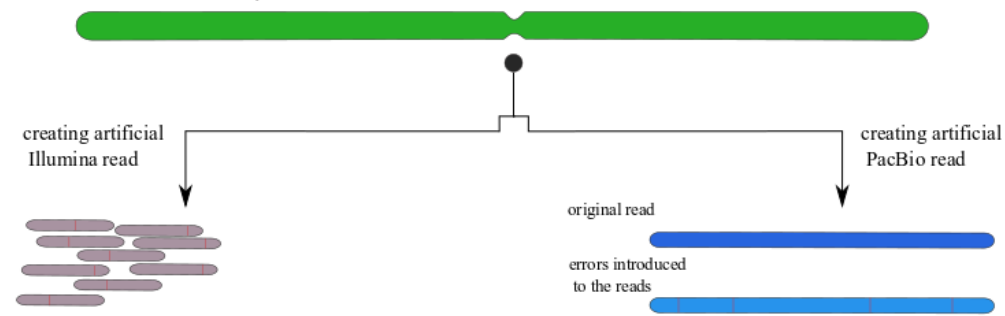**Creation of artificial short and long reads**
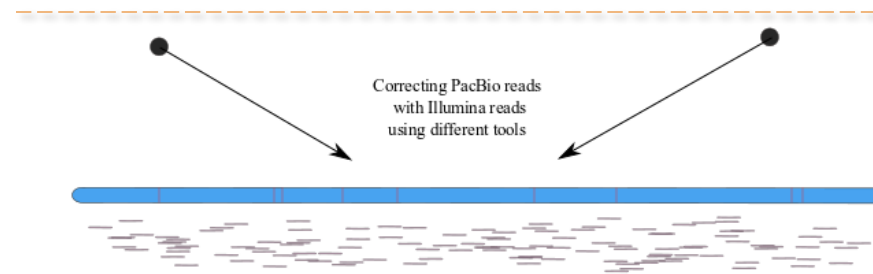
**Tool selection**
- LoRDEC
- Proovread
- LSC
- PBcR

**Assessment of:**
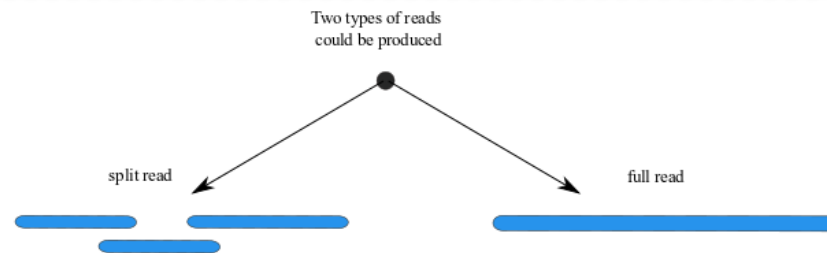- Efficiency
- Data loss
- Read length

chromosome under study

creating artificial
Illumina read

creating artificial
PacBio read

original read

errors introduced
to the reads

Correcting PacBio reads
with Illumina reads
using different tools

Two types of reads
could be produced

split read

full read

both will be aligned
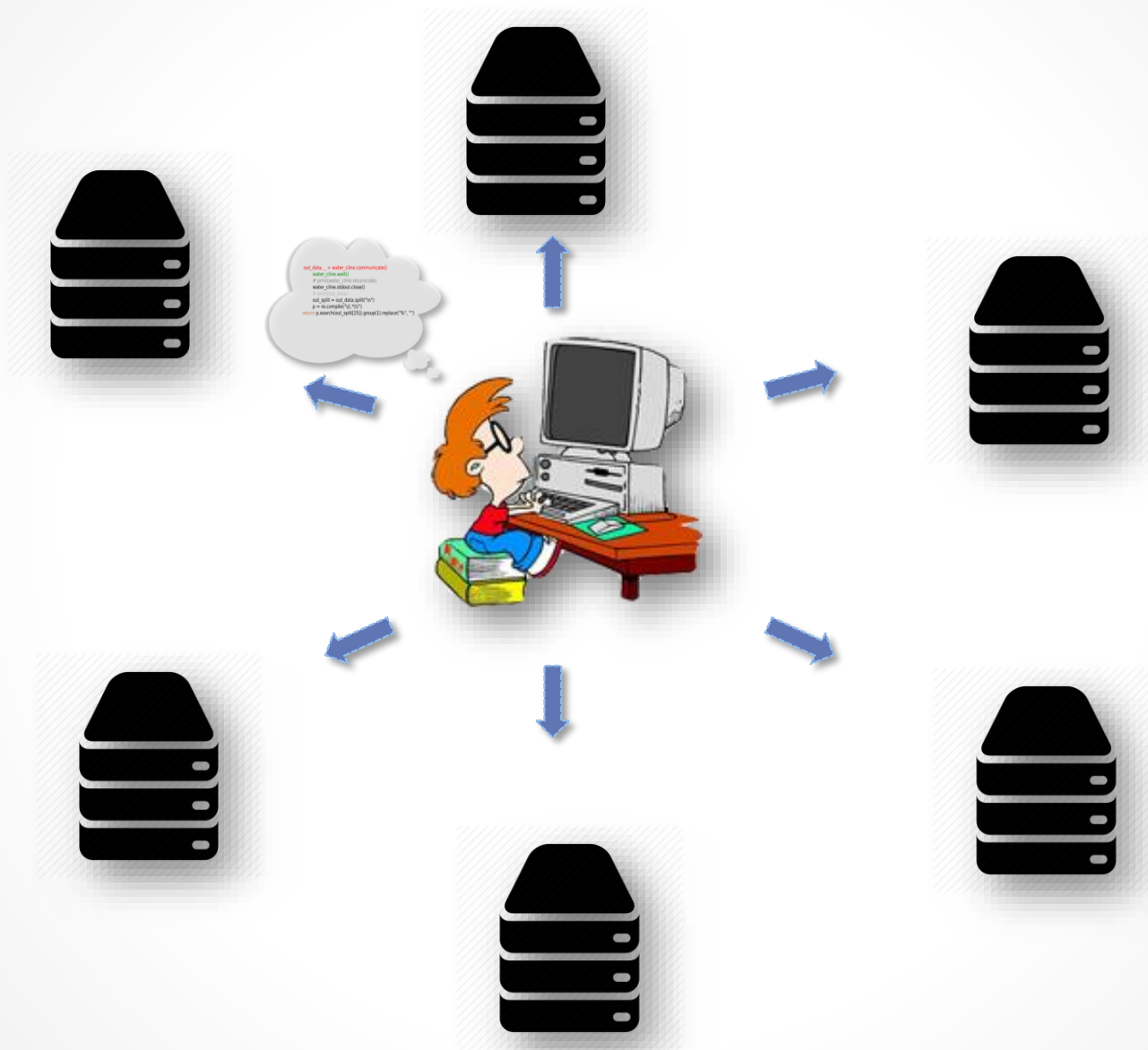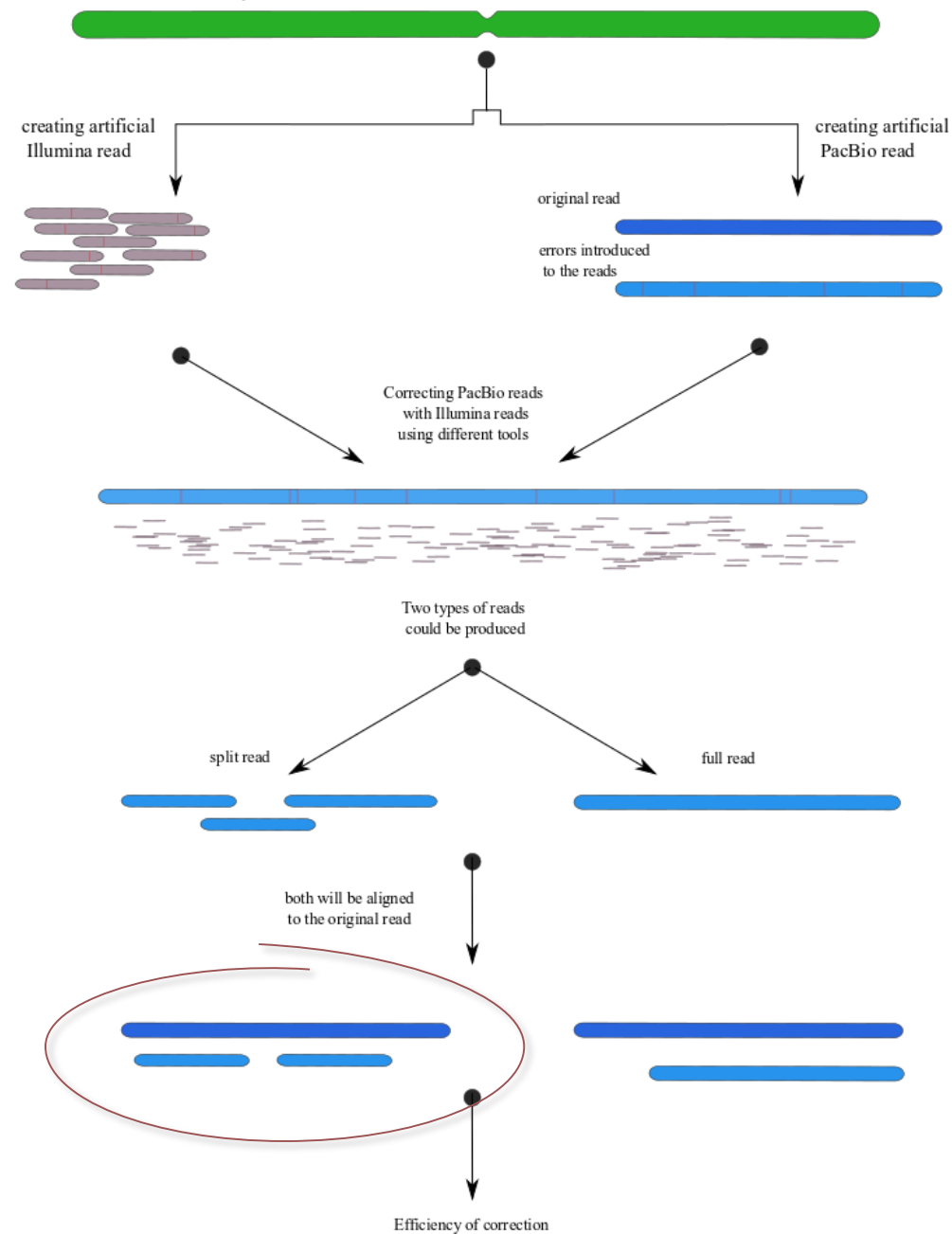to the original read

Efficiency of correction

1

2

3

4

15

# Eagle
• • •

# The result

• • •

For the split reads

# chromosome under study

creating artificial
Illumina read

creating artificial
PacBio read

original read

errors introduced
to the reads

Correcting PacBio reads
with Illumina reads
using different tools

Two types of reads
could be produced

split read

full read

both will be aligned
to the original read

Efficiency of correction

# Data loss
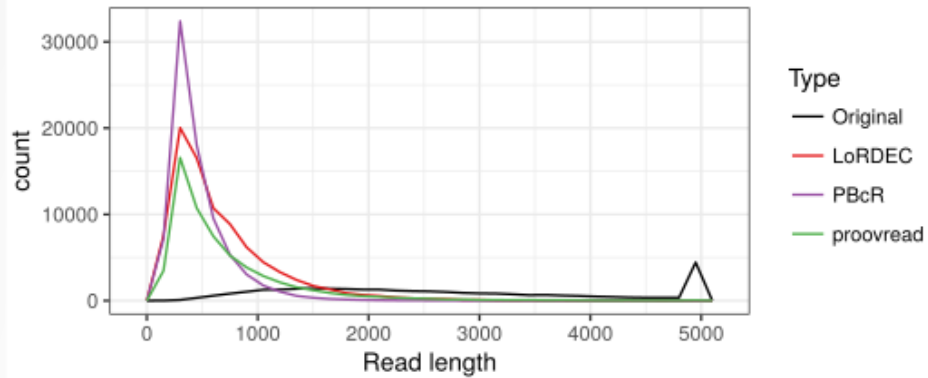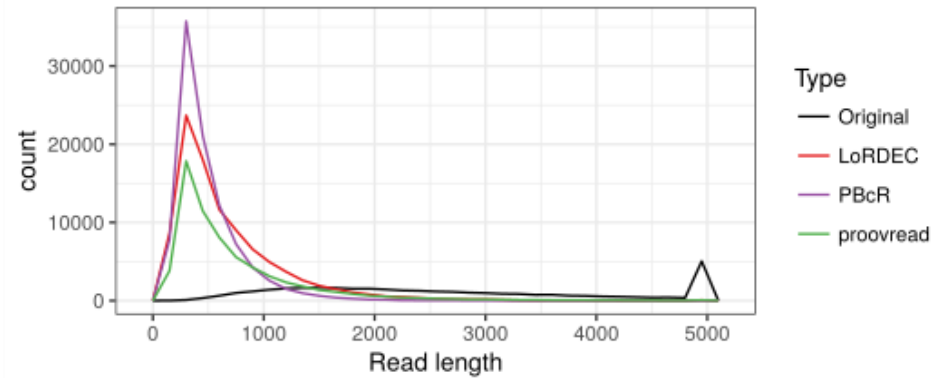
Clipped sequence distribution

# Read length count
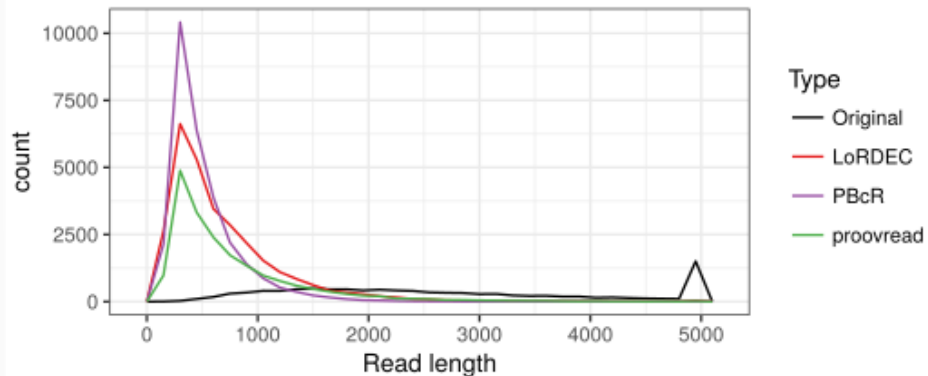
Correction length with original

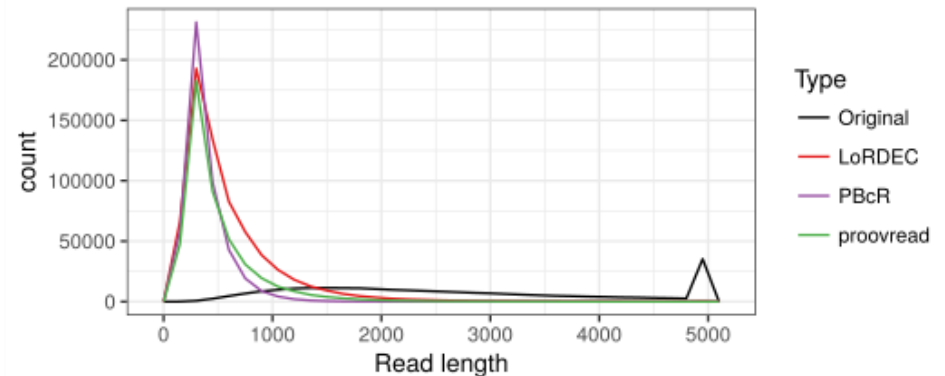E.coli Distribution of length after correction with the original length

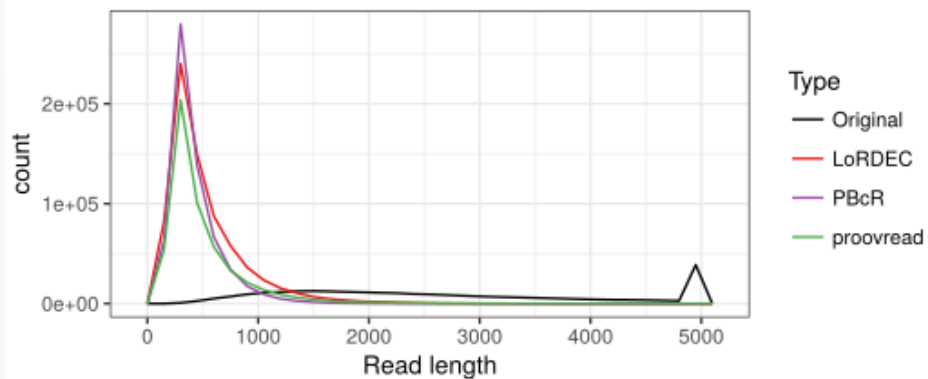Trypanosoma Distribution of length after correction with the original length

Yeast Distribution of length after correction with the original length

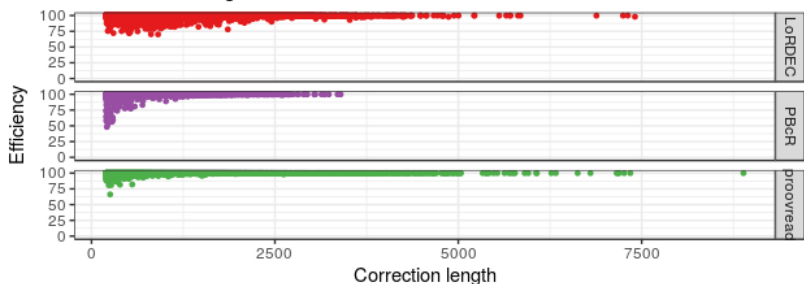Rice Distribution of length after correction with the original length

Human Distribution of length after correction with the original length
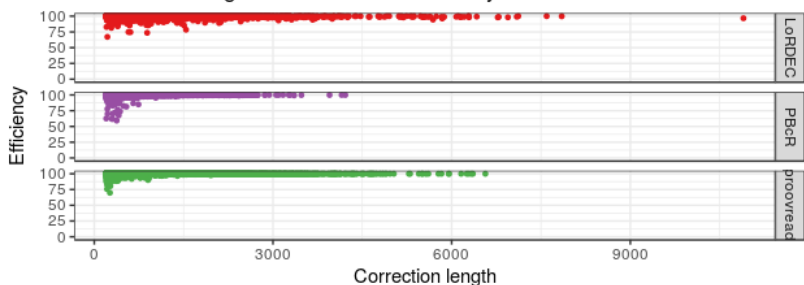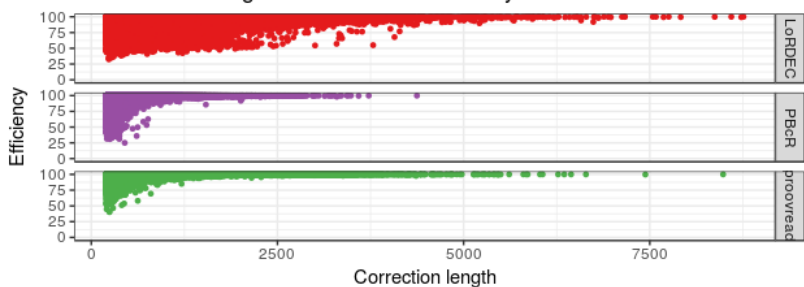
# Correction efficiency

Correction vs Length local

# Split reads result conclusion

# The result

· · ·

For the full reads

creating artificial
Illumina read

creating artificial
PacBio read

original read

errors introduced
to the reads

Correcting PacBio reads
with Illumina reads
using different tools

Two types of reads
could be produced

split read

full read

both will be aligned
to the original read

Efficiency of correction

E.coli corrected vs not corrected reads similarity

Trypanosoma corrected vs not corrected reads similarity

Yeast corrected vs not corrected reads similarity
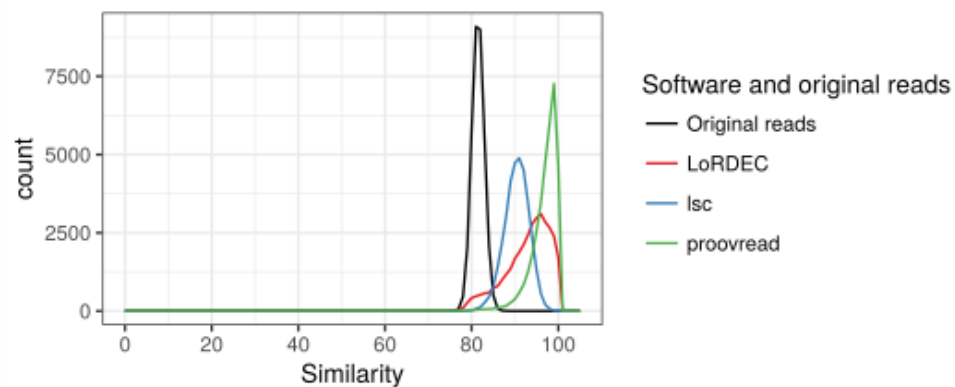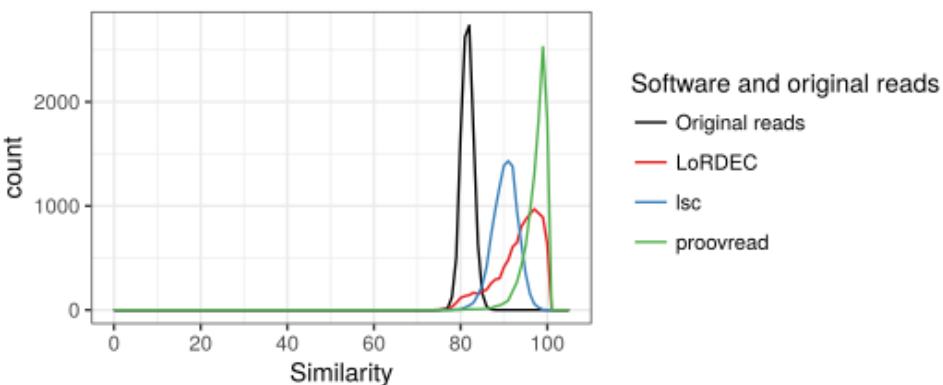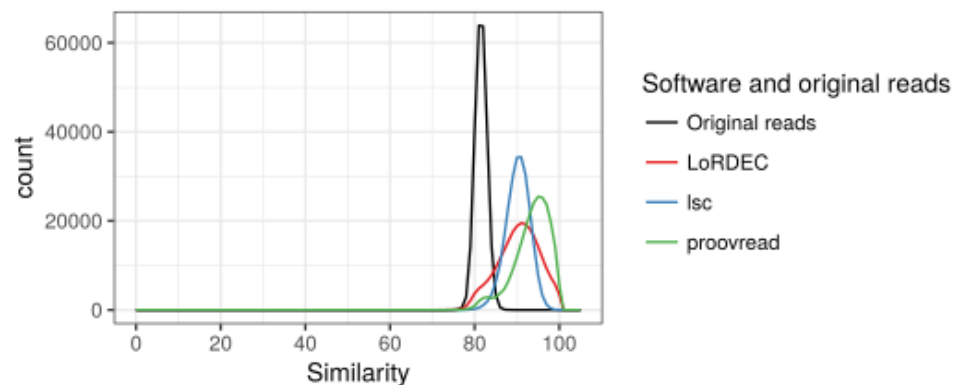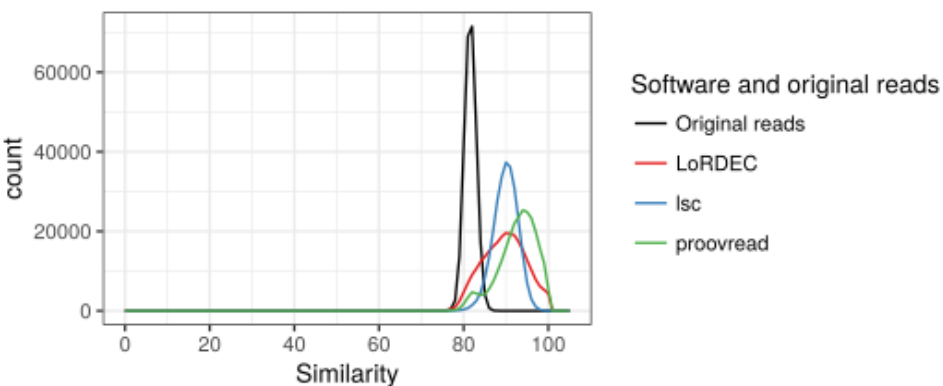
Rice corrected vs not corrected reads similarity

Human corrected vs not corrected reads similarity

# Full reads result conclusion

# Performance

| | User | System | CPU | Total |
|---|---|---|---|---|
| Proovread | 49302.87s (13.69h) | 6811.14s (1.89h) | 676% | 2:18:11.62 |
| LoRDEC | 2943.50s (0.81h) | 532.17s (0.14h) | 790% | 7:19.84 |
| LSC | 58944.97s (16.37h) | 195.09s (0.05h) | 934% | 1:45:27.54 |
| PBcR | 20538.56s (5.70h) | 1140.17s (0.31h) | 383% | 1:34:11.75 |

# Chimera

• • •

- Read correction

- Genome assembly

**Department of Computational Biology (AMU):**

Prof. dr hab. Wojciech Karłowski

Dr Marek Żywicki

**Department of Protein Biosynthesis (IBCh):**

Prof. dr hab. Tomasz Twardowski

(The head of the project)

Dr Agata Tyczewska

Dr Joanna Gracz

www.combio.pl