



LEIDEN UNIVERSITY MEDICAL CENTER

# Research selection

**Jeroen F.J. Laros**

**Leiden Genome Technology Center**

**Department of Human Genetics**

**Center for Human and Clinical Genetics**



*About me*

The last decade:

- Master Computer Science.
- Ph.D. Mathematics and Natural Sciences.
- Post-doctoral researcher LUMC.
- Senior researcher.
  - Coördinator Bioinformatics LGTC.

*About me*

The last decade:

- Master Computer Science.
- Ph.D. Mathematics and Natural Sciences.
- Post-doctoral researcher LUMC.
- Senior researcher.
  - Coördinator Bioinformatics LGTC.

Currently active in:

- Variant databases / formal descriptions.
- Next Generation Sequencing.
- Metagenomics.
- Forensics.

## *Activities*

### Collaborations:

- Within the LUMC: Human Genetics, Clinical Genetics, Forensic lab, Medical Microbiology, Dermatology, Hematology, ICT.
- Leiden University: Leiden Institute of Advanced Computer Science, Faculty of Social Sciences.
- Hogeschool Leiden.
- SURFSara.
- Commercial: GenomeScan, BaseClear, PhenoSystems.

## *Activities*

Other activities:

- Communication with external partners.
- Innovation / research and development.
- Implementation in diagnostics.
- Design and policy HPC infrastructure.
- Design and policy good research practice.
- Algorithm design.
- Audits.
- Education.

## *Sequencers*



Figure 1: HiSeq 2500.



Figure 2: Ion proton.

## *Next generation sequencing data*

```

1  @SGGPP : 4 : 101
2  TTCGGGGGCTGGCAAATCCACTTCCGTGACACGCTACCATTCGCTGGTGGT
3  +
4  -' +4589, 53330-0&07+03 : 54/2362-+ . 488587>@/25440++0 (+
5  @SGGPP : 4 : 102
6  CGGTAAACCACCCTGCTGACGGAACCCTAATGCGCCTGAAAGACAGCGTTC
7  +
8  34/--0' + . 000 ( . 55 : ; : 99 ( 0 (+2 (22 (0316; 185; ; 0; : <<>=AA59
9  @SGGPP : 4 : 106
10 TCGTTAACGACTTTGTTTCGCCACCGCAACCGCCTGTTTCGGGTCACAGGCA
11 +
12 09875; 5?<; ?@A4?B : BBB<AA>CCC>C>BB0 . ->=0488+3444 : @5@<
13 @SGGPP : 4 : 112
14 TTGATGAATATATTATTTTCAGGGAATAATTATGACACCTTTAGAACGCATT
15 +
16 70<<@ : : 5 : <; ==7; >>/79< : . : 494 . 8 ( , , 8 : 753/5@5??C>B??B7

```

Listing 1: A FastQ file.

*Data analysis*

Very diverse.

Align to a reference genome (resequencing):

- Variant detection.
- Phylogenetic reconstruction.



*Data analysis*

Very diverse.

Align to a reference genome (resequencing):

- Variant detection.
- Phylogenetic reconstruction.

Or to multiple references:

- Antibiotic resistance testing.

## *Data analysis*

Very diverse.

Align to a reference genome (resequencing):

- Variant detection.
- Phylogenetic reconstruction.

Or to multiple references:

- Antibiotic resistance testing.

*De novo* assembly:

- First step to make a reference genome.
- Finding large rearrangements.

## *Metrics for NGS data files*

Comparing  $k$ -mer profiles.

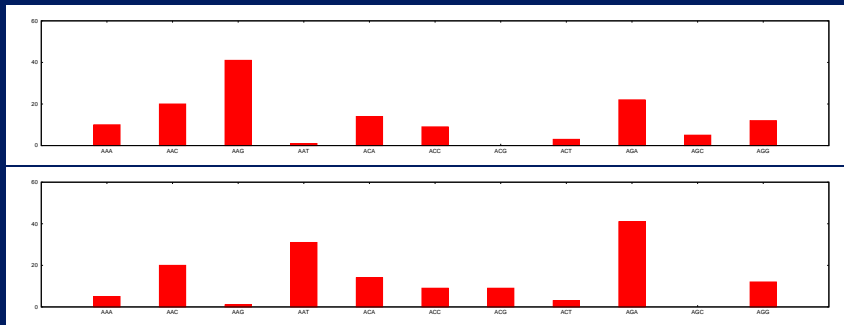


Figure 3: Two  $k$ -mer profiles.

*Forensics*

## STR profiling:

- Look deeper into STRs by using sequencing.
- Semi-global alignment of flanking sequences.
- Regular expressions for known alleles.
- Classification of new alleles.

*Forensics*

## STR profiling:

- Look deeper into STRs by using sequencing.
- Semi-global alignment of flanking sequences.
- Regular expressions for known alleles.
- Classification of new alleles.

## SNP profiling:

- Highly variable regions in a certain population.
- Easier to work with than with STRs.

## *Effect prediction*

LUMC Mutalyzer
☰

---

**Variant description**

AB026906.1:c.[274G>T;300del;341\_342del]

Example: AB026906.1:c.274G>T

---

**Reference protein**

```

1  MAVLWRLSAV CGALGGRALL LRTPVVRPAH ISAFLLQDRPI PEWCGVQHIH LSPSHHSGSK
61 AASLHWTSER VVSVLLLGLL PAAYLNPCSA MDYSLAAAL T LHGHMGLGQV VTDVYVHGDL
121 QKAAKAGLLA LSALTFAGLC YFNYHDVIGIC KAVAMLWKL*
```

---

**Protein predicted from variant coding sequence**

```

1  MAVLWRLSAV CGALGGRALL LRTPVVRPAH ISAFLLQDRPI PEWCGVQHIH LSPSHHSGSK
61 AASLHWTSER VVSVLLLGLL PAAYLNPCSA MYISLAAAL T FMVTGALDKL LLTVHGDLQ
121 KAAKAGLLAL SALTFAGLCY FNYHDVIGICK AVAMLWKL*
```

Figure 4: Mutalyzer.

*A “human” way of finding a description*

Observation:

- There is always a default way of describing a variant (**delins**).
- A **delins** may be split in smaller parts.

## A “human” way of finding a description

### Observation:

- There is always a default way of describing a variant (**delins**).
- A **delins** may be split in smaller parts.

### Outline:

- Find the *area of change*.
- Describe this as a **delins**.
- Find the largest overlap in this area of change, splitting the area in two.
- Describe the two sub areas, and see whether this description is smaller than the one we have.



## *Outline of the algorithm*

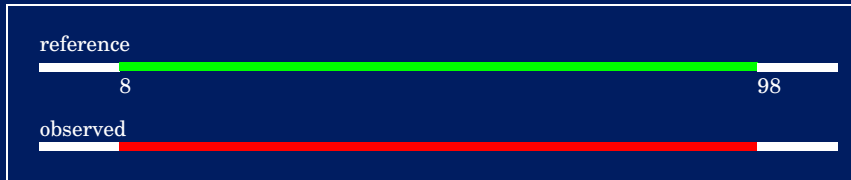


Figure 5: How would a human do it?

8\_98delinsAGATGCGATAGATTAGCTATATAGGATCG . . .

## *Outline of the algorithm*

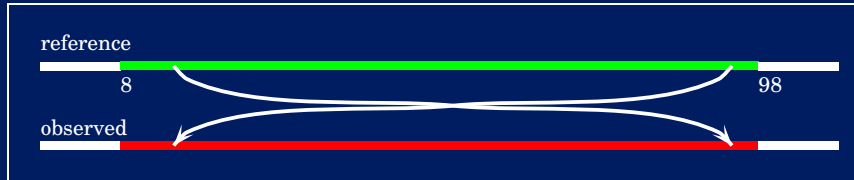


Figure 5: How would a human do it?

8\_98delinsAGATGCGATAGATTAGCTATATAGGATCG...

## *Outline of the algorithm*

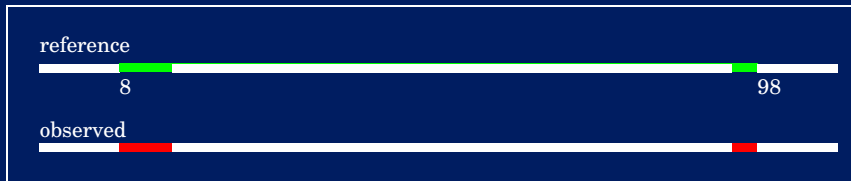


Figure 5: How would a human do it?

```
8_98delinsAGATGCGATAGATTAGCTATATAGGATCG . . .
[8_12delinsAGATG;13_96inv;97_98delinsTG]
```

## *Outline of the algorithm*

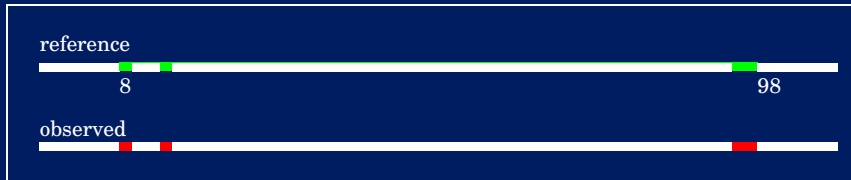


Figure 5: How would a human do it?

```
8_98delinsAGATGCGATAGATTAGCTATATAGGATCG . . .
[8_12delinsAGATG;13_96inv;97_98delinsTG]
[8G>A;12C>G;13_96inv;97_98delinsTG]
```

## *Finding common sub strings*

How would a computer do it?

	A	T	G	A	G	C	G
A	1	0	0	1	0	0	0
T	0	2	0	0	0	0	0
C	0	0	0	0	0	1	0
A	1	0	0	1	0	0	0
G	0	0	1	0	2	0	1
C	0	0	0	0	0	3	0
A	1	0	0	1	0	0	0

Table 1: LCS dynamic programming.

*Accuracy vs. speed*AGAGGACG  $\Rightarrow$  AG AG GA CGGAGGACA  $\Rightarrow$  GA AG GG GA AC CA

## *Accuracy vs. speed*

AGAGGACG  $\Rightarrow$  AG AG GA CG

GAGGACA  $\Rightarrow$  GA AG GG GA AC CA

	A	A	G	C
	G	G	A	G
GA	0	0	1	0
AG	1	1	0	0
GG	0	0	0	0
GA	0	0	2	0
AC	0	0	0	0
CA	0	0	0	0

Table 2: Rough method to find large strings.

## *Accuracy vs. speed*

AGAGGACG  $\Rightarrow$  AG AG GA CG

GAGGACA  $\Rightarrow$  GA AG GG GA AC CA

	A	A	G	C
	G	G	A	G
GA	0	0	1	0
AG	1	1	0	0
GG	0	0	0	0
GA	0	0	2	0
AC	0	0	0	0
CA	0	0	0	0

Table 2: Rough method to find large strings.



## *Accuracy vs. speed*

	A	A	G	C
	G	G	A	G
GA	0	0	1	0
AG	1	1	0	0
GG	0	0	0	0
GA	0	0	2	0
AC	0	0	0	0
CA	0	0	0	0

Table 3: “Zoom out”  $k = 2$ .

	A	G
	G	G
	A	A
GAG	0	0
AGG	0	0
GGA	0	1
GAC	0	0
ACA	0	0

Table 4: “Zoom out”  $k = 3$ .

## *Accuracy vs. speed*

	A	A	G	C
	G	G	A	G
GA	0	0	1	0
AG	1	1	0	0
GG	0	0	0	0
GA	0	0	2	0
AC	0	0	0	0
CA	0	0	0	0

Table 3: “Zoom out”  $k = 2$ .

	A	G
	G	G
	A	A
GAG	0	0
AGG	0	0
GGA	0	1
GAC	0	0
ACA	0	0

Table 4: “Zoom out”  $k = 3$ .

We find all common sub strings larger than  $k$ .

## *Accuracy vs. speed*

	A	A	G	C
	G	G	A	G
GA	0	0	1	0
AG	1	1	0	0
GG	0	0	0	0
GA	0	0	2	0
AC	0	0	0	0
CA	0	0	0	0

Table 3: “Zoom out”  $k = 2$ .

	A	G
	G	G
	A	A
GAG	0	0
AGG	0	0
GGA	0	1
GAC	0	0
ACA	0	0

Table 4: “Zoom out”  $k = 3$ .

We find all common sub strings larger than  $k$ .

The length of these strings are at least  $\ell k$  and at most  $\ell k + (k - 1)$  long.



## Acknowledgements:

Martijn Vermaat

Jonathan Vis

Lusine Khachatryan

Yahya Anvar

Kristiaan van der Gaag

Peter de Knijff

Johan den Dunnen