

# Quasi-species Identification by Model Based Clustering

Massively parallel sequencing

Bie Verbist  
Ghent University – Janssen

Potsdam - September 26th

Quasi-species identification by model-based clustering

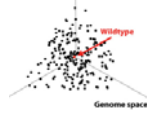
- Introduction
- Goal
- Massively parallel sequencing
  - Technology
  - Quality scores
- Model-based clustering
- Results

Potsdam - September 26th

Quasi-species identification by model-based clustering

## Introduction

- Identify and quantify quasi-species.

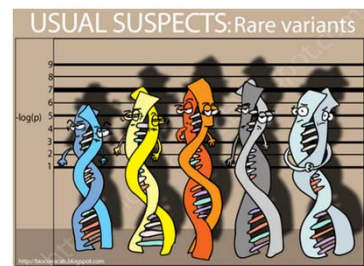


- Focus on low-frequency variants
- Deep sequencing: in-depth characterization of sequence variation.

Potsdam - September 26th

Quasi-species identification by model-based clustering

## Goal




Potsdam - September 26th

Quasi-species identification by model-based clustering

## Massively Parallel Sequencing

- illumina



Potsdam - September 26th

Quasi-species identification by model-based clustering

## Massively Parallel Sequencing

- Sequencing read:
 

```
@HWUSI-EAS1524:12:FC:1:1:18845:1091 1:N:0:
ATGACCCATCAAAGACTTAATAGCAGAAATACAGAAGCAGGGCAAGGCC
+
```
- ASCII-coded:
 

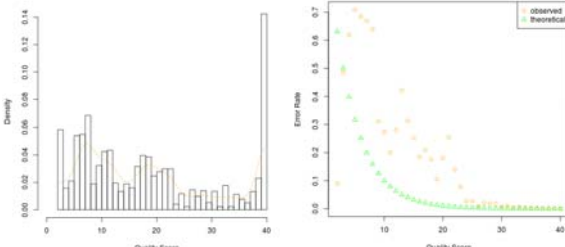
```
"#S%&'()*+,-./0123456789;=<=>?@ABCDEFGHI Qi ∈ [2,40]
```
- $Q_{ij} = -10 \log_{10}(p_{ij})$  with  $p_{ij}$  error probability

Potsdam - September 26th

Quasi-species identification by model-based clustering

## Massively Parallel Sequencing

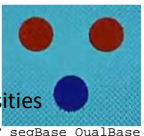
- Distribution of  $Q_{ij}$  at error positions
- $Q_{ij}$  doesn't reflect true error probabilities



Potsdam - September 26th

Quasi-species identification by model-based clustering

## Massively Parallel Sequencing

- Other metrics:
 
  - Corrected and censored raw intensities

lane	tile	cycle	x	y	A	C	G	T	seqBase	QualBase
1	1	1	4797	11926	2017	2	0	0	A	I
1	1	1	2727	14009	1114	573	450	443	A	I
1	1	1	2430	16718	42	73	0	10	C	:

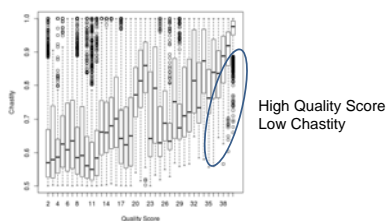
  - Second best base call

```
First Read (R) @HWUSI-EAS1524:12:FC:1:1:4797:11926
ATCTGCTCCTGTATCTAATAGAGCTTCCTTTAGTTGCCCCCTATCTTTAT
Second Read (S) @HWUSI-EAS1524:12:0:1:1:17484:13206
CAACAGCTTCTGTATCTTTATAGAGCCTTCACTCCTATTTTAGCCAACCCC
```

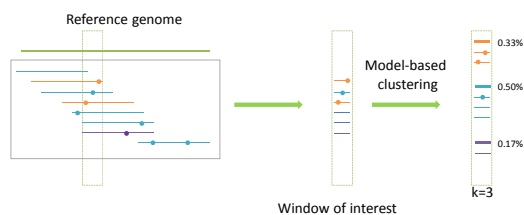
Potsdam - September 26th

## Massively Parallel Sequencing

- Chastity ( $C_{ii}$ ):  
max. intensity / sum of the 2 highest intensities



## Model-Based Clustering



- Clusters = haplotypes,  
- missing data problem: EM-algorithm

## Model-Based Clustering

- Notation:
  - $r_i$ : best base calls of read  $i$  ( $i=1 \dots n$ )
  - $s_i$ : second best base calls of read  $i$  ( $i=1 \dots n$ )
  - $z_{ij}$ :  $z_{ij}=1$  when read  $i$  belongs to haplotype  $j$  ( $j=1 \dots k$ )
  - $\tau_j$ : probability to belong to haplotype  $j$

- Complete Data Likelihood:

$$L = \prod_{i=1}^n \prod_{j=1}^k (f_j(r_i, s_i) \tau_j)^{z_{ij}} \quad \text{with } f_j(r, s) = \text{Prob}(r, s | \text{haplotype } j)$$

## Model-Based Clustering

$$L = \prod_{i=1}^n \prod_{j=1}^k (f_j(r_i, s_i) \tau_j)^{z_{ij}}$$

$$\text{with } f_j(r_i, s_i) = \prod_{l=1}^m \theta_{ril}^{I(r_{il}=h_{il})} \theta_{sil}^{I(s_{il}=h_{il})} \theta_{oil}^{(1-I(r_{il}=h_{il})) (1-I(s_{il}=h_{il}))}$$

$$\text{Prob}(\text{base } l \text{ of read } i = \text{base } l \text{ of haplotype } j) = \begin{cases} \theta_{ril} & \text{base} \\ \theta_{sil} & \text{second base} \end{cases}$$

$$\theta_{oil} = 1 - \theta_{ril} - \theta_{sil}$$

- Model is overidentified

## Model-Based Clustering

- Model  $\theta$  parameters:

$$\log \frac{\theta_{ril}}{\theta_{oil}} = \beta_{0r} + \beta_{1r} Q_{il} + \beta_{2r} Q_{i(i-1)} + \beta_{3r} C_{il} + \beta_{4r} H_{il}$$

$$\log \frac{\theta_{sil}}{\theta_{oil}} = \beta_{0s} + \beta_{1s} Q_{il} + \beta_{2s} Q_{i(i-1)} + \beta_{3s} C_{il} + \beta_{4s} H_{il}$$

- Complexity reduced to 10  $\beta$ -parameters
- Likelihood still depends on haplotype membership

$$L = \prod_{i=1}^n \prod_{j=1}^k (f_j(x_i, s_i) \tau_j)^{z_{ij}}$$

## Model-Based Clustering

- EM algorithm

- E step: update posterior probability

$$\hat{z}_{ij} = \text{Prob}(\text{read } i \text{ is of haplotype } j | \text{observed data})$$

- M step: update  $\beta$  parameters by maximizing expected likelihood

till convergence

## Model-Based Clustering

End result

- After conversion

- $\hat{z}_{ij} = \text{Prob}(\text{read } i \text{ is of haplotype } j | \text{observed data})$

- haplotype frequency  $\tau_j$

$$\hat{\tau}_j = \frac{\sum_{i=1}^n \hat{z}_{ij}}{n}$$

## Model-Based Clustering

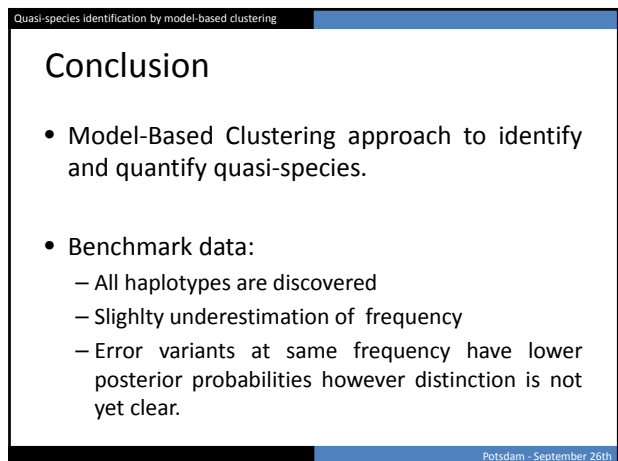
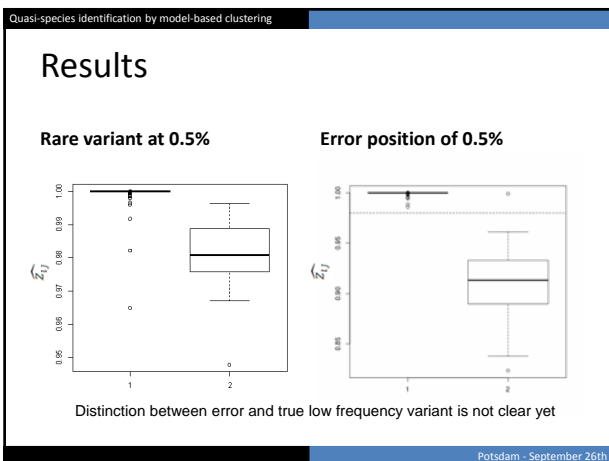
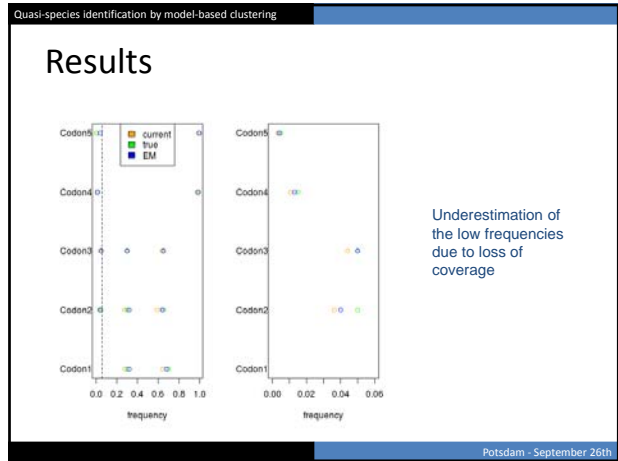
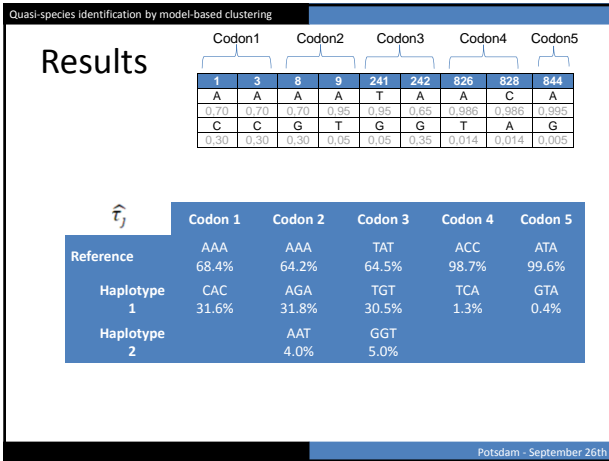
- Benchmark dataset

- 5 HIV-1 plasmids

SDM090122	SDM090126	SDM090141	SDM090147	SDM090151
0.63	0.30	0.050	0.014	0.0046

- 10 mutation positions

Pos	Codon1		Codon2		Codon3		Codon4		Codon5	
	1	3	8	9	199	241	242	826	828	844
Ref	A	A	A	A	A	T	A	A	C	A
	0.6988	0.6888	0.6988	0.9498	0.37	0.9498	0.6486	0.986	0.986	0.9954
Variant	C	C	G	T	G	G	G	T	A	G
	0.3012	0.3012	0.3012	0.0502	0.63	0.0502	0.3514	0.014	0.014	0.0046



## Conclusion

- Future research:
  - Check behaviour in error prone region of Illumina
  - Relax assumption number of haplotypes  $j$  is known
  - Allow for missing data in window
  - Model optimization at each codon position

## Acknowledgement



- Promoters: *Prof. Dr. O. Thas*<sup>1</sup> and *Dr. L. Bijns*<sup>2</sup>
- Translational genomics team
- Yves Wetzels, Tobias Verbeke, Joris Meys<sup>1</sup>



**BACK-UP**

## Corrected raw intensities

- Corrected for
  - signal amplitude: due to inequality of fluor emission intensities
  - spectral cross talk: overlap of emission frequency correction using the calculated matrix file
  - phasing and pre-phasing values: estimates of the proportion of molecules out of sync - behind or ahead

- Split in forward/reverse data frame
  - Perform analysis twice – same results??

Start position codon		
In algorithm (Codon level)	In excel list Koen (Ncl level)	In slides
598	2850	1
604	2856	7
838	3090	241
1423	3675	826
1441	3693	844
1258 (error pos)		

**FEEDBACK CONFERENCE**