

# Next-Generation Sequencing and Methylation Profiling.

Alignment to a reference of BS-treated sequences

#### Alberto Policriti

(joint work with C. Del Fabbro, E. De Paoli, N. Prezza, and F. Vezzi)

Napoli — September 26, 2012

Alignment of BS-treated sequences









Alignment of BS-treated sequences

# **Dot-plot**



**Dot-plot** 

Alignement

Alignment of BS-treated sequences

#### A familiar alignment "program"



- extends to strings a valuation given on the alphabet by a cost matrix
- gaps are allowed but it is not clear how much to charge for them

Alignment of BS-treated sequences

## Dot-plot

#### A familiar alignment "program"



- extends to strings a valuation given on the alphabet by a cost matrix
- gaps are allowed but it is not clear how much to charge for them

#### Aligning

Computing a *sliding* dot-plot

Introduction

Alignement

Alignment of BS-treated sequences

## **Distances and Scores**

# A "negative" or "positive" point of view Distance ↓ Score ↑

Alignment of BS-treated sequences

## **Distances and Scores**

| A "negative" or "positive" point of view |         |  |
|--|---------|--|
| Distance $\Downarrow$                    | Score ↑ |  |

 $\alpha$  and  $\beta$  strings

Definition (Levensthein (edit) distance)

 $d_L(\alpha,\beta)$ : min number of Insertion Deletion Substitution

to convert a  $\alpha$  into  $\beta$ 

Alignment of BS-treated sequences

## **Distances and Scores**

| A "negative" or "positive" point of view |         |  |
|--|---------|--|
| Distance $\Downarrow$                    | Score ↑ |  |

 $\alpha$  and  $\beta$  strings

Definition (Levensthein (edit) distance)

 $d_L(\alpha,\beta)$ : min number of Insertion Deletion Substitution

to convert a  $\alpha$  into  $\beta$ 

What are we searching for?

 $Alignment \equiv Program$ 

Alignment of BS-treated sequences

## Let us simplify our lives

#### **Definition (Hamming distance)**

 $d_H(\alpha, \beta)$ : min number of Substitution

to convert a  $\alpha$  into  $\beta$ 

Alignment of BS-treated sequences

## Let us simplify our lives

#### **Definition (Hamming distance)**

 $d_H(\alpha, \beta)$ : min number of Substitution

to convert a  $\alpha$  into  $\beta$ 

#### **Complexity and gain**

Hamming  $\Rightarrow |\alpha|$  (the length of  $\alpha$ ) must be equal to  $\beta$ Levensthein  $\Rightarrow |\alpha|$  and  $|\beta|$  may be different

# Matrix-based alignment and *fast* alignments

#### Matrix-based alignments (Levensthein)

A weighted dot-plot computation of the *Pattern* against the *Reference* along the entire R



# Matrix-based alignment and *fast* alignments

#### Matrix-based alignments (Levensthein)

A weighted dot-plot computation of the *Pattern* against the *Reference* along the entire R



#### Cost

 $3G \times 100$  cells to scan (for Human)!

Introduction

Alignement ○●○○○ Alignment of BS-treated sequences

# Matrix-based alignment and fast alignments

#### Fast alignments (Hamming)

#### A scan of R with P stored in memory



Introduction

Alignement ○●○○○ Alignment of BS-treated sequences

# Matrix-based alignment and *fast* alignments

Fast alignments (Hamming)

#### A scan of R with P stored in memory



time necessary to load the reference

Alignment of BS-treated sequences

## **Data structures**

#### Observation

We can structure up R while loading it

#### Advantages: example

Ordering a set of 1G ( $\approx 2^{30})$  numbers allows to search *any* element in no more than 30 steps

Alignment of BS-treated sequences

## **Data structures**

#### Observation

We can structure up R while loading it

#### Advantages: example

Ordering a set of 1G ( $\approx 2^{30})$  numbers allows to search any element in no more than 30 steps

#### Question

Can we order a text (like R)?

Introduction

Alignement

Alignment of BS-treated sequences

## Two important techniques

#### Suffixes

R-suffixes can be ordered lexicographically

#### Hashing

*R*-blocks (*k*-mers) can be used as numbers to *index* an array

# Two important techniques

#### Suffixes

| 1  | AGGTTGCCAGTGT | 1  | AGGTTGCCAGTGT |
|----|---------------|----|---------------|
| 2  | GGTTGCCAGTGT  | 9  | AGTGT         |
| 3  | GTTGCCAGTGT   | 8  | CAGTGT        |
| 4  | TTGCCAGTGT    | 7  | CCAGTGT       |
| 5  | TGCCAGTGT     | 6  | GCCAGTGT      |
| 6  | GCCAGTGT      | 2  | GGTTGCCAGTGT  |
| 7  | CCAGTGT       | 12 | GT            |
| 8  | CAGTGT        | 10 | GTGT          |
| 9  | AGTGT         | 3  | GTTGCCAGTGT   |
| 10 | GTGT          | 13 | Т             |
| 11 | TGT           | 5  | TGCCAGTGT     |
| 12 | GT            | 11 | TGT           |
| 13 | Т             | 4  | TTGCCAGTGT    |

Alignment of BS-treated sequences

# Two important techniques



Alignment of BS-treated sequences

# Compression

Ordered *R*-suffixes can be *compressed*: Burrows-Wheeler

AGGTTGCCAGTGT AGTGT CAGTGT CCAGTGT GCCAGTGT GGTTGCCAGTGT GT GTGT **GTTGCCAGTGT** т TGCCAGTGT TGT TTGCCAGTGT

Alignment of BS-treated sequences

## Compression

Ordered *R*-suffixes can be *compressed*: Burrows-Wheeler

\$ AGGTTGCCAGTGT\$ AGTGT\$ CAGTGT\$ CCAGTGT\$ GCCAGTGT\$ **GGTTGCCAGTGT\$** GT\$ GTGT\$ GTTGCCAGTGT\$ **T**\$ TGCCAGTGT\$ TGT\$ TTGCCAGTGT\$

Alignment of BS-treated sequences

## Compression

Ordered *R*-suffixes can be *compressed*: Burrows-Wheeler

\$AGGTTGCCAGTGT AGGTTGCCAGTGT\$ AGTGT\$AGGTTGCC CAGTGT\$AGGTTGC CCAGTGT\$AGGTTG GCCAGTGT\$AGGTT GGTTGCCAGTGT\$A GT\$AGGTTGCCAGT GTGT\$AGGTTGCCA GTTGCCAGTGT\$AG T\$AGGTTGCCAGTG TGCCAGTGT\$AGGT TGT\$AGGTTGCCAG TTGCCAGTGT\$AGG

Alignment of BS-treated sequences

## Compression

Ordered *R*-suffixes can be *compressed*: Burrows-Wheeler

\$AGGTTGCCAGTGT AGGTTGCCAGTGT\$ AGTGT\$AGGTTGCC CAGTGT\$AGGTTGC CCAGTGT\$AGGTTG GCCAGTGT\$AGGTT GGTTGCCAGTGT\$A GT\$AGGTTGCCAGT GTGT\$AGGTTGCCA GTTGCCAGTGT\$AG T\$AGGTTGCCAGTG TGCCAGTGT\$AGGT TGT\$AGGTTGCCAG TTGCCAGTGT\$AGG

# Aligning BS-treated sequences: the problem

## Determine Methylated C's (i.e. $C^{m}$ 's) along the reference R

- Sodium Bisulphite converts un-methylated cytosines to uraciles
- Methylated cytosines remain un-converted
- ⇒ after PCR un-methylated cytosines appear as thymines and methylated ones remain unaltered

Alignment of BS-treated sequences

## Aligning BS-treated sequences: the problem







# Aligning BS-treated sequences: the problem

### Determine Methylated C's (i.e. $C^{m}$ 's) along the reference R

- Sodium Bisulphite converts un-methylated cytosines to uraciles
- Methylated cytosines remain un-converted
- ⇒ after PCR un-methylated cytosines appear as thymines and methylated ones remain unaltered

#### Two (main) issues

- some mismatch conveys information
- we are working with a 5-characters alphabet

Alignment of BS-treated sequences

## **Characters and mismatches**

#### "Natural" approach

Ignore T - C mismatches ( $\Rightarrow$  identify *T*'s and *C*'s  $\Rightarrow$  use a 3-characters alphabet!)

Alignment of BS-treated sequences

## **Characters and mismatches**

#### "Natural" approach

Ignore T - C mismatches

 $(\Rightarrow$  identify T's and C's  $\Rightarrow$  use a 3-characters alphabet!)

#### Hamming distance to "render" methylated C's

 $d_{GH}$ : the Generalized Hamming distance for BS-treated sequences, assigns

0 distance to T-read/C-reference

mismatches (as well as G-read/A-reference)

## ERNE-BS5

## $d_{GH}$ indirectly reduces the alphabet to 3 characters

- Many reads do not reach a threshold to be reliably aligned
- Many misalignment

#### Idea: use the *fifth* character to improve alignment

Methylated C's can be used to disambiguate multiple alignments

# **ERNE-BS5**

 $N_i^C$  number of *C*'s read at position *i*;  $N_i^T$  number of *T*'s read at position *i*;

**Definition (Methylation level)** 

$$\mu(i) = \frac{N_i^C}{N_i^C + N_i^T}$$

#### Definition (Context-aware Hamming distance $d_{\alpha}$ )

Watson strand:

$$\alpha(i, \mathbf{x}) = \begin{cases} 1 - \mu(i) \\ \mu(i) \\ 0 \end{cases}$$

if 
$$R[i] = C \land x = C$$
;  
if  $R[i] = C \land x = T$ ;  
otherwise.

Crick strand analogous.

## **ERNE-BS5**

#### The full strategy

step 1Align reads using  $d_{GH}$ step 2On-line compute  $d_{\alpha}$ step 3Align reads multiply aligning using  $d_{\alpha}$ step 4If new alignments are found, go to step 2

# Results on Arabidopsis thaliana, real data



# **Further work**

- Extensively test ERNE-BS5
- Encode more in  $d_{\alpha}$
- Integrate more knowledge
- Study the two methylation's patterns
- Compress