

# Aspetti algoritmici di base nell'analisi di dati NGS

## Basic Algorithmic Aspects of NGS data analysis

NGS Milan Meeting, 17 luglio 2012

Gianluca Della Vedova

Università degli Studi di Milano-Bicocca

[gianluca.dellavedova@unimib.it](mailto:gianluca.dellavedova@unimib.it)

# Efficient Algorithms and NGS

Perfect match

Lots of data

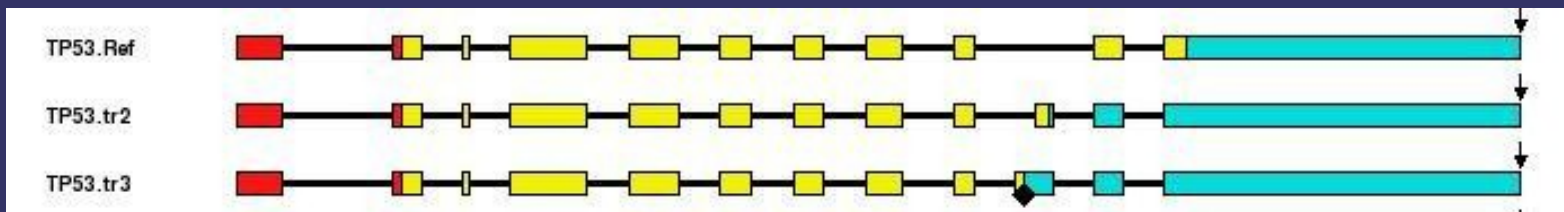
- short fragments
- Different error distribution
  - Too many errors

You cannot be too fast

You cannot be too space efficient

# NGS in Transcriptomics

- **Genome Annotation (with RNA-seq data)**
  - splice junctions (exon-intron boundaries)
  - reconstruction of full-length isoforms
  - expression level of known isoforms
  - validation of predicted isoforms



<http://t.caspur.it/ASPicDB/index.php>

# The Most Basic Problem: Alignment

- Sequence comparison
- Pattern matching
- Strategy:
  - Find *seeds* (short perfect matches)
  - Filter seeds
  - Extend seeds

# NGS Alignment

**INPUT:** a set  $R$  of NGS data (genomic reads or RNA-seq reads) and a reference sequence

**OUTPUT:** for each read  $r$  in  $R$ , find the reference regions very similar to  $r$

**APPROACH:** align each read to the reference sequence

# Basic problem: Alignment

- Gapped/ungapped alignment of a **huge amount** of reads to a **long** reference sequence (genome)
- **Goal:** speed and accuracy
- **Strategy:** find seeds → time-space efficient data structures for indexing reads and/or reference

# Index

- Index = data structure for pattern matching
  - Seed Hash Table
  - Suffix Tree
  - Suffix Array
  - Burrows-Wheeler Transform (BWT) and FM-index

# Seed Hash Table

S = acgtgacgtg

acgt	1, 6
cgtg	2, 7
gtga	3
tgac	4
gacg	5

Spaced seed (1101):

ac\*t, cg\*g,

gt\*a, tg\*c,

ga\*g

**Human Genome >12GB**

SOAP, MAQ



# Suffix Tree

S = dabdac

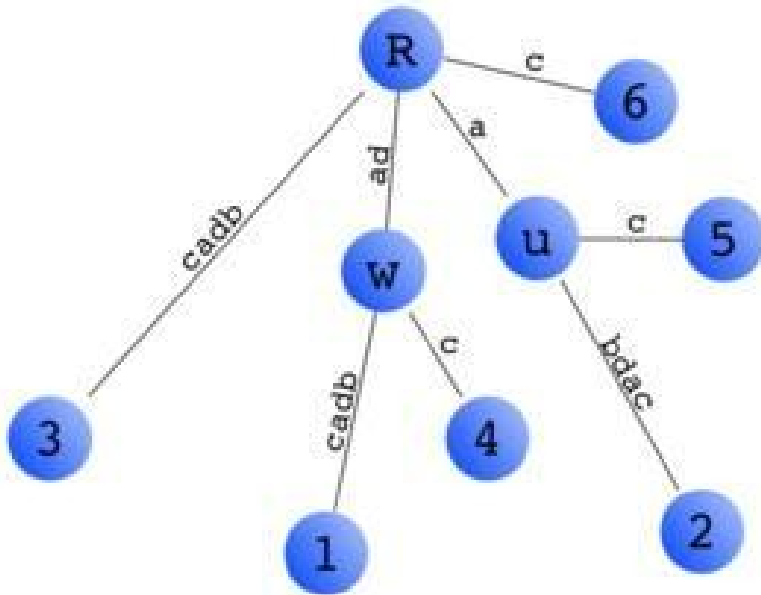


Image from Wikipedia

**Human Genome >35GB**

MUMmer (not for NGS)

Gianluca Della Vedova – Aspetti algoritmici di base nell'analisi di dati NGS

# Suffix array

S = acgtgacgtg\$

acgtgacgtg\$	1
cgtgacgtg\$	2
gtgacgtg\$	3
tgacgtg\$	4
gacgtg\$	5
acgtg\$	6
cgtg\$	7
gtg\$	8
tg\$	9
g\$	10
\$	11

SA

\$	11
acgtg\$	6
acgtgacgtg\$	1
cgtg\$	7
cgtgacgtg\$	2
g\$	10
gacgtg\$	5
gtg\$	8
gtgacgtg\$	3
tg\$	9
tgacgtg\$	4

Human Genome ~12GB

Vmatch, Segemehl

Gianluca Della Vedova – Aspetti algoritmici di base nell'analisi di dati NGS

# Burrows-Wheeler Transform (BWT)

T = acgtgacgtg\$

# Burrows-Wheeler Transform (BWT)

T = acgtgacgtg\$

acgtgacgtg\$
cgtgacgtg\$a
gtgacgtg\$ac
tgacgtg\$acg
gacgtg\$acgt
acgtg\$acgtg
cgtg\$acgtga
gtg\$acgtgac
tg\$acgtgacg
g\$acgtgacgt
\$acgtgacgtg

# Burrows-Wheeler Transform (BWT)

T = acgtgacgtg\$

acgtgacgtg\$
cgtgacgtg\$a
gtgacgtg\$ac
tgacgtg\$acg
gacgtg\$acgt
acgtg\$acgtg
cgtg\$acgtga
gtg\$acgtgac
tg\$acgtgacg
g\$acgtgacgt
\$acgtgacgtg

\$acgtgacgtg
acgtg\$acgtg
acgtgacgtg\$
cgtg\$acgtga
cgtgacgtg\$a
g\$acgtgacgt
gacgtg\$acgt
gtg\$acgtgac
gtgacgtg\$ac
tg\$acgtgacg
tgacgtg\$acg

# Burrows-Wheeler Transform (BWT)

T = acgtgacgtg\$

BWT

acgtgacgtg\$
cgtgacgtg\$a
gtgacgtg\$ac
tgacgtg\$acg
gacgtg\$acgt
acgtg\$acgtg
cgtg\$acgtga
gtg\$acgtgac
tg\$acgtgacg
g\$acgtgacgt
\$acgtgacgtg

\$acgtgacgtg
acgtg\$acgtg
acgtgacgtg\$
cgtg\$acgtga
cgtgacgtg\$a
g\$acgtgacgt
gacgtg\$acgt
gtg\$acgtgac
gtgacgtg\$ac
tg\$acgtgacg
tgacgtg\$acg

# Burrows-Wheeler Transform (BWT)

T = acgtgacgtg\$

acgtgacgtg\$
cgtgacgtg\$a
gtgacgtg\$ac
tgacgtg\$acg
gacgtg\$acgt
acgtg\$acgtg
cgtg\$acgtga
gtg\$acgtgac
tg\$acgtgacg
g\$acgtgacgt
\$acgtgacgtg

F

BWT

\$acgtgacgtg
acgtg\$acgtg
acgtgacgtg\$
cgtg\$acgtga
cgtgacgtg\$a
g\$acgtgacgt
gacgtg\$acgt
gtg\$acgtgac
gtgacgtg\$ac
tg\$acgtgacg
tgacgtg\$acg

# Inverting BWT

BWT
g
g
\$
a
a
t
t
c
c
g
g



# Inverting BWT

<b>F</b>	<b>BWT</b>
\$	g
a	g
a	\$
c	a
c	a
g	t
g	t
g	c
g	c
t	g
t	g

# Inverting BWT

F	BWT	BWT	F'
\$	g	g	\$
a	g	g	a
a	\$	\$	a
c	a	a	c
c	a	a	c
g	t	t	g
g	t	t	g
g	c	c	g
g	c	c	g
t	g	g	t
t	g	g	t

# Inverting BWT

F	BWT	BWT	F	F2
\$	g	g	\$	\$a
a	g	g	a	ac
a	\$	\$	a	ac
c	a	a	c	cg
c	a	a	c	cg
g	t	t	g	g\$
g	t	t	g	ga
g	c	c	g	gt
g	c	c	g	gt
t	g	g	t	tg
t	g	g	t	tg

# Inverting BWT

F
\$
a
a
c
c
g
g
g
g
t
t

BWT
g
g
\$
a
a
t
t
c
c
g
g

BWT	F
g	\$
g	a
\$	a
a	c
a	c
t	g
t	g
c	g
c	g
g	t
g	t

BWT	F2
g	\$a
g	ac
\$	ac
a	cg
a	cg
t	g\$
t	ga
c	gt
c	gt
g	tg
g	tg

# BWT: pattern searching

BWT

S = acgtgacgtg

P = acg

\$acgtgacgtg
acgtg\$acgtg
acgtgacgtg\$
cgtg\$acgtga
cgtgacgtg\$a
g\$acgtgacgt
gacgtg\$acgt
gtg\$acgtgac
gtgacgtg\$ac
tg\$acgtgacg
tgacgtg\$acg

# Burrows-Wheeler Transform (BWT) vs Suffix Array

	SA
\$	11
acgtg\$	6
acgtgacgtg\$	1
cgtg\$	7
cgtgacgtg\$	2
g\$	10
gacgtg\$	5
gtg\$	8
gtgacgtg\$	3
tg\$	9
tgacgtg\$	4

	BWT
\$acgtgacgtg	g
acgtg\$acgtg	g
acgtgacgtg\$	\$
cgtg\$acgtga	a
cgtgacgtg\$a	a
g\$acgtgacgt	t
gacgtg\$acgt	t
gtg\$acgtgac	c
gtgacgtg\$ac	c
tg\$acgtgacg	g
tgacgtg\$acg	g

# FM-index

- Compressed genome
- Compressed information for pattern matching
  - $C(x)$  = number of characters of  $S$  smaller than  $x$
  - $\text{Occ}(x,i)$  = occurrence of  $x$  in  $\text{BWT}[1:i]$
  - $\text{SA}(T)$  = suffix array

## Proprietà:

$$\text{begin}(aP) = C(a) + \text{Occ}(a, \text{begin}(P)) - 1$$

$$\text{end}(aP) = C(a) + \text{Occ}(a, \text{end}(P)) - 1$$

$$\text{BWT}[i] = T[\text{SA}(i) - 1]$$

# BWT e FM-index

- Bowtie
  - ungapped alignment (no indels)
  - human genome in 1.3 GB
- BWA
  - gapped alignment (with indels)
- SOAP2
- Bowtie2 (gapped alignment)



# Thanks



# Research group [www.algolab.eu](http://www.algolab.eu)

- ⇒ Paola Bonizzoni
- ⇒ Gianluca Della Vedova
- ⇒ Raffaella Rizzi
- ⇒ Yuri Pirola (Post-doc)
- ⇒ Stefano Beretta (Ph.D. student)
- ⇒ And more...

# References

- Indexing compressed text. Paolo Ferragina, Giovanni Manzini.  
<http://doi.acm.org/10.1145/1082036.1082039>
- A block sorting lossless data compression algorithm. M. Burrows and D. Wheeler.  
<http://www.hpl.hp.com/techreports/Compaq-DEC/SRC-RR-124.pdf>
- Fast and accurate short read alignment with Burrows–Wheeler transform. Heng Li and Richard Durbin, doi: 10.1093/bioinformatics/btp324
- Exploring single-sample SNP and INDEL calling with whole-genome de novo assembly. Heng Li. doi: 10.1093/bioinformatics/bts280
- Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Ben Langmead, Cole Trapnell, Mihai Pop and Steven L Salzberg.  
doi:10.1186/gb-2009-10-3-r25

Quest' opera è distribuita con licenza Creative Commons Attribuzione - Condividi allo stesso modo 3.0 Italia (<http://creativecommons.org/licenses/by-sa/3.0/it/>).