

Examples of algorithms used in genomics

Michel Koskas

January 10-th 2018

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Brujin graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Bruijn graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion

Mapping. From: A reference Genome

0 1 2 3 4 5 6 7
012345678901234567890123456789012345678901234567890123456789012345
G=TGGCTTGTTGGTCATGGGGTCGACTAGGAGACCATTATTTAGTGCC TTCACAATATCAACAGGCCCTCCATTTCT

and a set of Reads:

```
GGGGTCGACT
GAGACGTACTGT
CCGAGTCCCGGT
GCAGAAACCAG
CAAAGTCC
GGAAAGCTTGAC
GATTTGCCTCAG
TGCTTCCAAAA
R = GCAAGTTCTCTA
AATCACCCATA
ATCTTCTAAGGA
GGCTCACAAGTA
ACGAATTTAA
ATGGTTTTCACT
GTCACCCAATCA
GGTAACGGTGA
```

Exact Match:

0 1 2 3 4 5 6 7
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5
G=TGGCTTGTGGTCATG**GGGGTCGACT**AGGAGACCATTATTTAGTGCCTTCACAATATCAACAGGCCCTCCATTTCT

GGGGTCGACT
GAGACGTA CTGT
CCGAGTCCCGGT
GCAGAAACCAG
CAAAGTCC
GGAAAGCTTGAC
GATTTGCCTCAG
TGCTTCCAAAA
R = GCAAGTTCTCTA
AATCACCCATA
ATCTTCTAAGGA
GGCTCACAAGTA
ACGAATTTAA
ATGGTTTTCACT
GTCACCCAATCA
GGTAACGGTGA

Mismatch:

0 1 2 3 4 5 6 7
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5
G=TGGCTTGTGGTCATG**GGGG****CGACT**AGGAGACCATTATTTAGTGCCTTCACAATATCAACAGGCCCTCCATTTCT

*GGGG***T***CGACT*

GAGACGTACTGT

CCGAGTCCCGGT

GCAGAAACCAG

CAAAGTCC

GGAAAGCTTGAC

GATTTGCCTCAG

TGCTTCCAAAA

R = GCAAGTTCTCTA

AATCACCCATA

ATCTTCTAAGGA

GGCTCACAAGTA

ACGAATTTAA

ATGGTTTTCACT

GTCACCCAATCA

GGTAACGGTGA

InDel:

0 1 2 3 4 5 6 7
 01234567890123456789012345678901234567890123456789012345678901234567890123456789012345
 G=TGGCTTGTGGTCATGGGGG-CGACTAGGAGACCATTATTTAGTGCCTTCACAATATCAACAGGCCCTCCATTTCT

GGGGTCGACT

GAGACGTACTGT

CCGAGTCCCGGT

GCAGAAACCAG

CAAAGTCC

GGAAAGCTTGAC

GATTTGCCTCAG

TGCTTCCAAAA

R = GCAAGTTCTCTA

AATCACCCATA

ATCTTCTAAGGA

GGCTCACAAGTA

ACGAATTTAA

ATGGTTTTCACT

GTCACCCAATCA

GGTAACGGTGA

InDel:

0 1 2 3 4 5 6 7
012345678901234567890123456789012345678901234567890123456789012345
G=TGGCTTGTGGTCATGGGGTCGACTAGGAGACCATTATTTAGTGCCTTCACAATATCAACAGGCCCTCCATTTCT

GGGG-CGACT
GAGACGTACTGT
CCGAGTCCCGGT
GCAGAAACCAG
CAAAGTCC
GGAAAGCTTGAC
GATTTGCCTCAG
TGCTTCCAAAA
R = GCAAGTTCTCTA
AATCACCCATA
ATCTTCTAAGGA
GGCTCACAAAGTA
ACGAATTTAA
ATGGTTTTCACT
GTCACCCAATCA
GGTAACGGTGA

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Bruijn graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion

A puzzle

A NGS machine produces billions of pieces of DNA.
These pieces have to be put together to find DNA back.

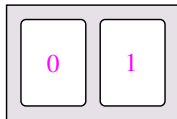
For instance how should we assemble the following *reads*?
BRAC, CADA, ACAD, ABRA, ABRA, ADAB ?

A puzzle

A NGS machine produces billions of pieces of DNA.
These pieces have to be put together to find DNA back.

For instance how should we assemble the following *reads*?
BRAC, CADA, ACAD, ABRA, ABRA, ADAB ?

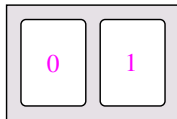
The de Bruijn Graphs: An erudite burglar.



A digital lock has a three digits code.
When trying 0100, one tries 010 and 100.

Question: is it possible to try all the possibilities in a minimal time
(minimal number of keys pressed)
(Such a sequence should be of length $3 + 2^3 - 1 = 10$)

The de Bruijn Graphs: An erudite burglar.

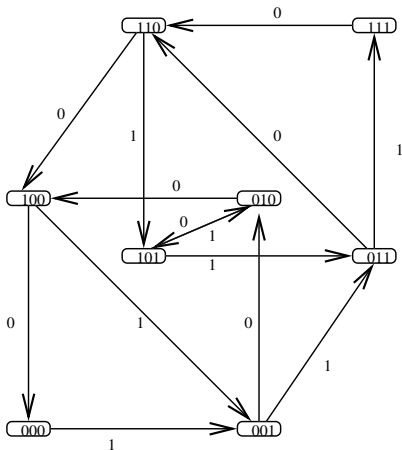


A digital lock has a three digits code.
When trying 0100, one tries 010 and 100.

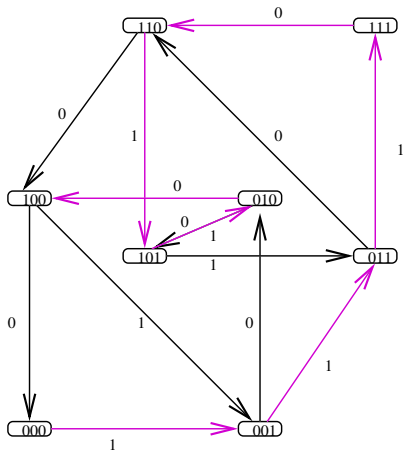
Question: is it possible to try all the possibilities in a minimal time
(minimal number of keys pressed)

(Such a sequence should be of length $3 + 2^3 - 1 = 10$)

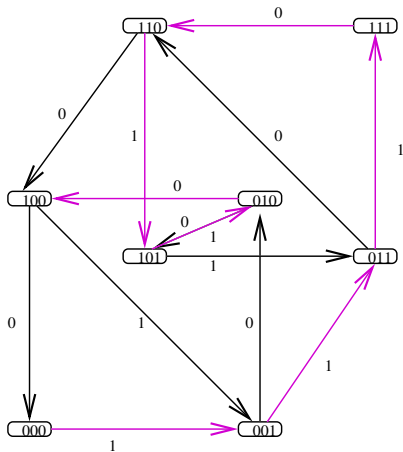
The question is to find a hamiltonian path among this graph:



The question is to find a hamiltonian path among this graph:



The question is to find a hamiltonian path among this graph:



0001110100

Now the first goal is to build the de bruijn graph of the reads
 This is the problem we address here.
 ABRA, ACAD, ABRA, ADAB, BRAC, CADA

	ABRA	ACAD	ABRA	ADAB	BRAC	CADA
ABRA	ABRA	ABRACAD	ABRA	ABRADAB	ABRAC	ABRA_CADA
ACAD	ACAD_ABRA	ACAD	ACAD_ABRA	ACADAB	ACAD_BRAC	ACADA
ABRA	ABRA	ABRACAD	ABRA	ABRADAB	ABRAC	ABRA_CADA
ADAB	AD ABRA	ADAB_ ACAD	AD ABRA	ADAB	ADAB BRAC	ADAB_CADA
BRAC	BRAC_ABRA	BR ACAD	BRAC_ABRA	BRAC_ADAB	BRAC	BRAC CADA
CADA	CAD ABRA	CAD ACAD	CAD ABRA	CADAB	CADA_ BRAC	CADA

Now the first goal is to build the de bruijn graph of the reads
 This is the problem we address here.

ABRA, ACAD, ABRA, ADAB, BRAC, CADA

(Cleaning up ...)

	ABRA	ACAD	ABRA	ADAB	BRAC	CADA
ABRA		ABRACAD	ABRA	ABRADAB	ABRAC	ABRA_CADA
ACAD	ACAD_ABRA		ACAD_ABRA	ACADAB	ACAD_BRAC	ACADA
ABRA	ABRA	ABRACAD		ABRADAB	ABRAC	ABRA_CADA
ADAB	ADABRA	ADAB_ACAD	ADABRA		ADABRAC	ADAB_CADA
BRAC	BRAC_ABRA	BRACAD	BRAC_ABRA	BRAC_ADAB		BRACADA
CADA	CADABRA	CADACAD	CADABRA	CADAB	CADA_BRAC	
BRAC	BRAC_ABRA	BRACAD	BRAC_ABRA	BRAC_ADAB	BRAC	BRACADA

Now the first goal is to build the de bruijn graph of the reads

This is the problem we address here.

ABRA, ACAD, ABRA, ADAB, BRAC, CADA

(Cleaning up ...)

	ABRA	ACAD	ABRA	ADAB	BRAC	CADA
ABRA		ABRACAD	ABRA	ABRADAB	ABRAC	
ACAD				ACADAB		ACADA
ABRA	ABRA	ABRACAD		ABRADAB	ABRAC	
ADAB	ADABRA		ADABRA		ADABRAC	
BRAC		BRACAD				BRACADA
CADA	CADABRA	CADACAD	CADABRA	CADAB		

Now the first goal is to build the de bruijn graph of the reads

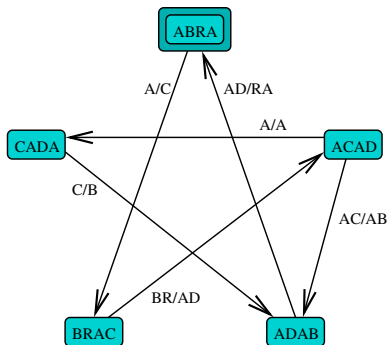
This is the problem we address here.

ABRA, ACAD, ABRA, ADAB, BRAC, CADA

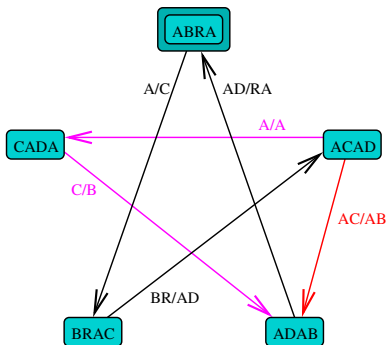
(Cleaning up ...)

	ABRA	ACAD	ABRA	ADAB	BRAC	CADA
ABRA			ABRA		ABRAC	
ACAD				ACADAB		ACADA
ABRA	ABRA				ABRAC	
ADAB	ADABRA		ADABRA			
BRAC		BRACAD				
CADA				CADAB		

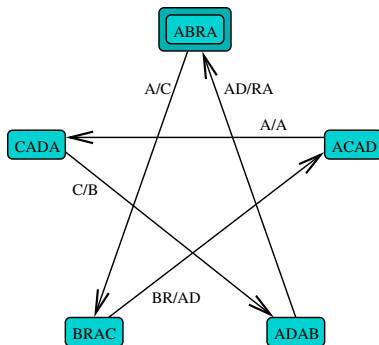
Then one reconstitutes the de Bruijn graph:



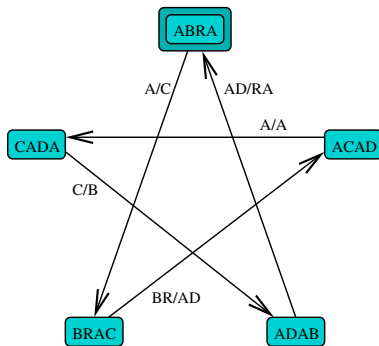
Then one reconstitutes the de Bruijn graph:



Then one reconstitutes the de Bruijn graph:

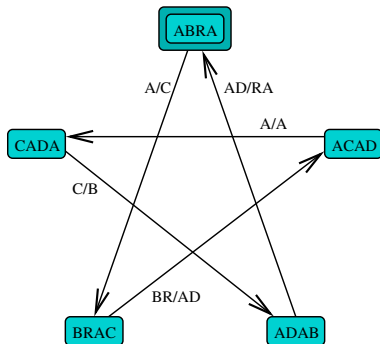


Then one reconstitutes the de Bruijn graph:



The frequencies allow us to reconstitute the sequence
ABRACADABRA

Then one reconstitutes the de Bruijn graph:



The frequencies allow us to reconstitute the sequence
ABRACADABRA

Repetitions in the full sequence is a very serious
issue... **ABRACADABRA** (Or was it
ABRABRABRABRABRABRACADABRA ?)

de Novo Sequencing (exact matching)

The problem:

Two sets of reads, Left and Right

	<i>CT...AGT</i>		<i>GT... A GT</i>
	<i>AA...CGA</i>		<i>AT... C GT</i>
<i>Left =</i>	<i>:::·.:::</i>	<i>Right =</i>	<i>:::·.:::</i>
	<i>GT...CGT</i>		<i>GA T ...GA</i>
	<i>CC...GAT</i>		<i>CC... G AT</i>

(Actually, when the reads are stored for instance in a FATSFA or FASTQ file we have more information that is disregarded for the sake of clarity)

de Novo Sequencing (exact matching)

The problem:

and one looks for exact matching overlaps between the Left and Right sets (prefix of a Left side read matches a suffix of a right side read) ...

<i>CT...AGT</i>	<i>GT... A GT</i>
<i>AA...CGA</i>	<i>AT... C GT</i>
<i>Left = :::.:::</i>	<i>Right = :::.:::</i>
<i>GT...CGT</i>	<i>GA T ...GA</i>
<i>CC...GAT</i>	<i>CC... G AT</i>

de Novo Sequencing (exact matching)

The problem:

... for any length

$CT...AGT$	$GT... A GT$
$AA...CGA$	$AT... C GT$
$Left = \vdots \vdots \cdot \vdots \vdots \vdots$	$Right = \vdots \vdots \cdot \vdots \vdots \vdots$
$GT...CGT$	$GAT...GA$
$CC...GAT$	$CC... GAT$

Remark: *Left* and *Right* may be the same set.

Remark: Of course reads might have different lengths, but for the sake of clarity we suppose they have the same length. Otherwise it is easy to adapt the description to the general case.

<i>Left</i> =	<i>C</i>	<i>T</i>	...	<i>A</i>	<i>G</i>	<i>T</i>		<i>G</i>	<i>T</i>	...	<i>A</i>	<i>G</i>	<i>T</i>
	<i>A</i>	<i>A</i>	...	<i>C</i>	<i>G</i>	<i>A</i>		<i>A</i>	<i>T</i>	...	<i>C</i>	<i>G</i>	<i>T</i>
	⋮	⋮	⋮	⋮	⋮	⋮	<i>Right</i> =	⋮	⋮	⋮	⋮	⋮	⋮
	<i>G</i>	<i>T</i>	...	<i>C</i>	<i>G</i>	<i>T</i>		<i>G</i>	–	<i>T</i>	...	<i>G</i>	<i>A</i>
	<i>C</i>	<i>C</i>	...	<i>G</i>	<i>A</i>	<i>T</i>		<i>C</i>	<i>C</i>	...	<i>G</i>	<i>A</i>	<i>T</i>

Of course there might be mismatches and/or indels in the overlaps...

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Bruijn graph
- 3 Usual technologies**
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion

Usual technologies use “Seed & Grow” techniques:

Find an exact match for a k mer of a read (the seed)

Then try to extend the seed, accepting a maximum number of mismatches and/or InDel.

Usual technologies use “Seed & Grow” techniques:
Find an exact match for a k mer of a read (the seed)
Then try to extend the seed, accepting a maximum number of mismatches and/or InDel.

Hash Tables

A Hash table is an efficient in-memory storage of key-value couples enabling one

- To add an element in a constant average time
- find the value for a given key in a constant time.

Hash the k -mers of the reference genome (key: hashed k -mer, value: indexes of occurrences)

$0 \qquad 1 \qquad 2 \qquad 3 \qquad 4 \qquad 5 \qquad 6 \qquad 7$
 $012345678901234567890123456789012345678901234567890123456789012345$
 $G=$ *TGG***C**TTG**TGG**TC**ATGG**GGGCCGACTAGGAGACCATTATTTAGTGCCTTCACAATATCAACAGGCCCTCCATTTCT

Key	Value
TGG	0, 8, 14

<i>SeedSize</i>	<i>TableSize</i>	<i>Storagespace(GB)</i>
32	18,446,744,073,709,551,616	137,438,953,472
28	72,057,594,037,927,936	536,870,912
24	281,474,976,710,656	2,097,152
20	1,099,511,627,776	8,192
18	68,719,476,736	512
16	4,294,967,296	32
12	16,777,216	128MB

(The storage space is the *maximal* storage space, each number being stored on 8 bytes.)

Long seeds will demand more space,
while short seeds will require more computation time...

<i>SeedSize</i>	<i>TableSize</i>	<i>Storagespace(GB)</i>
32	18,446,744,073,709,551,616	137,438,953,472
28	72,057,594,037,927,936	536,870,912
24	281,474,976,710,656	2,097,152
20	1,099,511,627,776	8,192
18	68,719,476,736	512
16	4,294,967,296	32
12	16,777,216	128MB

(The storage space is the *maximal* storage space, each number being stored on 8 bytes.)

Long seeds will demand more space,

while short seeds will require more computation time...

<i>SeedSize</i>	<i>TableSize</i>	<i>Storagespace(GB)</i>
32	18,446,744,073,709,551,616	137,438,953,472
28	72,057,594,037,927,936	536,870,912
24	281,474,976,710,656	2,097,152
20	1,099,511,627,776	8,192
18	68,719,476,736	512
16	4,294,967,296	32
12	16,777,216	128MB

(The storage space is the *maximal* storage space, each number being stored on 8 bytes.)

Long seeds will demand more space,
while short seeds will require more computation time. . .

Suffix Arrays

0 1
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9

$G = TGGCTTGTTGGTCATGGGGG$

Build a suffix array consists in

- 1 Add a letter \$ at the end of the word (lexicographically smaller

than any letter) 0 1 2
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
 $G = TGGCTTGTTGGTCATGGGGG\$$

- 2 compute its conjugates
- 3 Sort them

Suffix Arrays

0 1
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9

G=TGGCTTGTTGGTCATGGGGG

Build a suffix array consists in

- 1 Add a letter \$ at the end of the word (lexicographically smaller than any letter)

0 1 2
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
G=TGGCTTGTTGGTCATGGGGG\$

- 2 compute its conjugates

0 1 2
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0

G=TGGCTTGTTGGTCATGGGGG\$ (0)
GGCTTGTTGGTCATGGGGG\$T (1)
GCTTGTTGGTCATGGGGG\$TGT (2)
CTTGTTGGTCATGGGGG\$TGG (3)
TTGTTGGTCATGGGGG\$TGGC (4)
TGTGTTGGTCATGGGGG\$TGGCT (5)
GTTGGTCATGGGGG\$TGGCTT (6)
TTGGTCATGGGGG\$TGGCTTG (7)
TGGTCATGGGGG\$TGGCTTGT (8)
GGTCATGGGGG\$TGGCTTGTT (9)
GTCATGGGGG\$TGGCTTGTTG (10)
TCATGGGGG\$TGGCTTGTTGG (11)
CATGGGGG\$TGGCTTGTTGGT (12)
ATGGGGG\$TGGCTTGTTGGTC (13)
TGGGGG\$TGGCTTGTTGGTCA (14)
GGGGG\$TGGCTTGTTGGTCAT (15)
GGGG\$TGGCTTGTTGGTCATG (16)
GGG\$TGGCTTGTTGGTCATGG (17)
GG\$TGGCTTGTTGGTCATGGG (18)
G\$TGGCTTGTTGGTCATGGGG (19)

Suffix Arrays

0 1
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9

G=TGGCTTGTGGTCATGGGG

Build a suffix array consists in

- 1 Add a letter \$ at the end of the word (lexicographically smaller than any letter)

0 1 2
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
G=TGGCTTGTGGTCATGGGGG\$

- 2 compute its conjugates
- 3 Sort them

0 1 2
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
G=\$TGGCTTGTGGTCATGGGGG (20)
ATGGGGG\$TGGCTTGTGGTC (13)
CATGGGGG\$TGGCTTGTGGT (12)
CTTGTGGTCATGGGGG\$TGG (3)
GCTTGTGGTCATGGGGG\$TG (2)
G\$TGGCTTGTGGTCATGGGG (19)
GG\$TGGCTTGTGGTCATGGG (18)
GGCTTGTGGTCATGGGGG\$T (1)
GGG\$TGGCTTGTGGTCATGG (17)
GGGG\$TGGCTTGTGGTCATG (16)
GGGGG\$TGGCTTGTGGTCAT (15)
GGTCATGGGGG\$TGGCTTGT (9)
GTCATGGGGG\$TGGCTTGTG (10)
GTTGGTCATGGGGG\$TGGCTT (6)
TCATGGGGG\$TGGCTTGTGG (11)
TGGCTTGTGGTCATGGGGG\$ (0)
TGGGGG\$TGGCTTGTGGTCA (14)
TGGTCATGGGGG\$TGGCTTGT (8)

looking for the seed GGG

Suffix Arrays

0 1
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9

G=TGGCTTGTGGTCATGGGGG

Build a suffix array consists in

- 1 Add a letter \$ at the end of the word (lexicographically smaller than any letter)

0 1 2
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
G=TGGCTTGTGGTCATGGGGG\$

- 2 compute its conjugates

- 3 Sort them

0 1 2
0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0

G=\$TGGCTTGTGGTCATGGGGG (20)
 ATGGGGG\$TGGCTTGTGGTC (13)
 CATGGGGG\$TGGCTTGTGGT (12)
 CTTGTTGGTCATGGGGG\$TGG (3)
 GCTTGTGGTCATGGGGG\$TG (2)
 G\$TGGCTTGTGGTCATGGGG (19)
 GG\$TGGCTTGTGGTCATGGG (18)
 GGCTTGTGGTCATGGGGG\$T (1)
 GGG\$TGGCTTGTGGTCATGG (17)
 GGGG\$TGGCTTGTGGTCATG (16)
 GGGGG\$TGGCTTGTGGTCAT (15)
 GGTCATGGGGG\$TGGCTTGT (9)
 GTCATGGGGG\$TGGCTTGTG (10)
 GTTGGTCATGGGGG\$TGGCTT (6)
 TCATGGGGG\$TGGCTTGTGG (11)
 TGGCTTGTGGTCATGGGGG\$ (0)
 TGGGGG\$TGGCTTGTGGTCA (14)
 TGGTCATGGGGG\$TGGCTTGT (8)

looking for the seed GGG

The Suffix Tree may be seen as an attempt to save space. . .

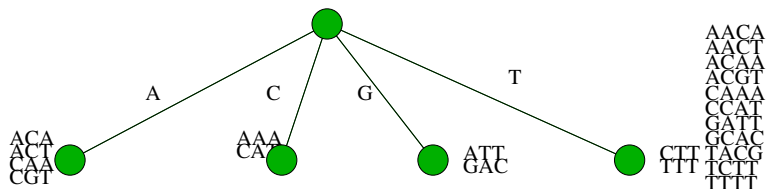
AACA
AACT
ACAA
ACGT
CAAA
CCAT
GATT
GCAC
TACG
TCTT
TTTT

The Suffix Tree may be seen as an attempt to save space. . .

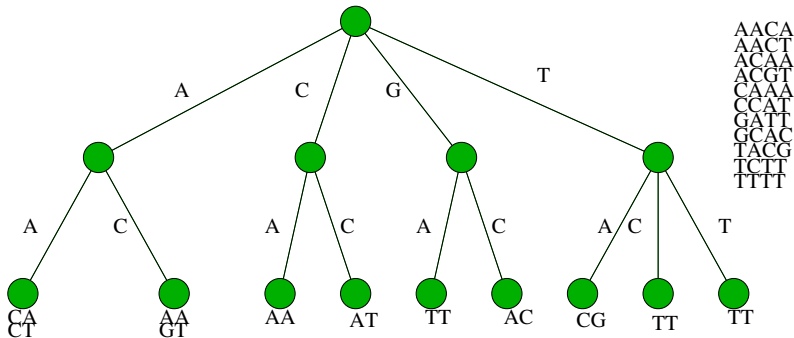


AACA
AACT
ACAA
ACGT
CAAA
CCAT
GATT
GCAC
TACG
TCCT
TTTT

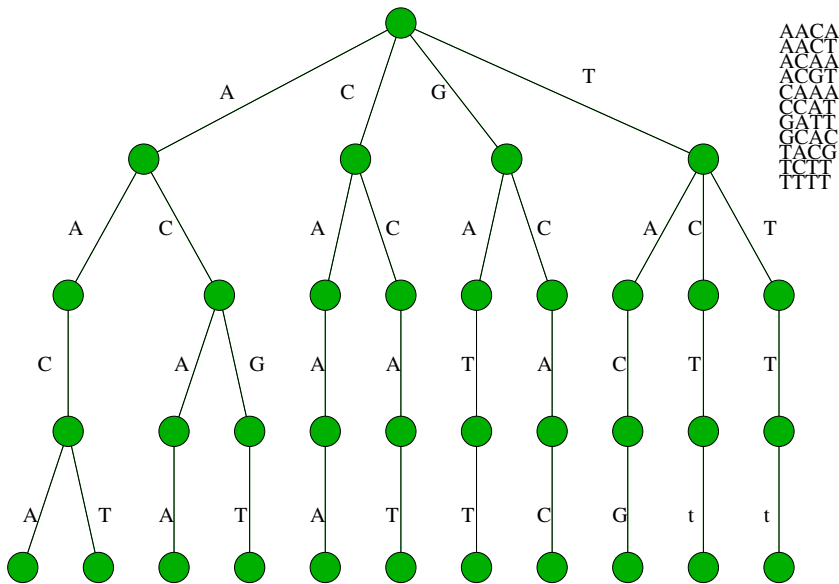
The Suffix Tree may be seen as an attempt to save space...



The Suffix Tree may be seen as an attempt to save space...



The Suffix Tree may be seen as an attempt to save space...



The Burrows-Wheeler Transform consists in keeping only the last column of the preceding array:

```

      0           1           2
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
G=TGGCTTGTGGTCATGGGGG$
      0           1           2
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
G=$TGGCTTGTGGTCATGGGGG (20)
  ATGGGGG$TGGCTTGTGGTC (13)
  CATGGGGG$TGGCTTGTGGT (12)
  CTTGTGGTCATGGGGG$TGG (3)
  GCTTGTGGTCATGGGGG$TG (2)
  G$TGGCTTGTGGTCATGGGG (19)
  GG$TGGCTTGTGGTCATGGG (18)
  GGCTTGTGGTCATGGGGG$T (1)
  GGG$TGGCTTGTGGTCATGG (17)
  GGGG$TGGCTTGTGGTCATG (16)
  GGGGG$TGGCTTGTGGTCAT (15)
  GGTCATGGGGG$TGGCTTGT (9)
  GTCATGGGGG$TGGCTTGTG (10)
  GTTGGTCATGGGGG$TGGCTT (6)
  TCATGGGGG$TGGCTTGTGG (11)
  TGGCTTGTGGTCATGGGGG$ (0)
  TGGGGG$TGGCTTGTGGTCA (14)
  TGGTCATGGGGG$TGGCTTGT (8)
  TGTTGGTCATGGGGG$TGGCT (5)
  TTGGTCATGGGGG$TGGCTTG (7)
  TTGTTGGTCATGGGGG$TGGC (4)
      0           1           2
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
GCTGGGGTGGTTGTG$ATTGC

```

The Burrows-Wheeler Transform consists in keeping only the last column of the preceding array:

```
      0           1           2
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
G=TGGCTTGTTGGTCATGGGGG$
      0           1           2
      0 1 2 3 4 5 6 7 8 9 0 1 2 3 4 5 6 7 8 9 0
GCTGGGGTGGTTGTG$ATTGC
```

It is easy to reconstitute the whole array with this single line.

Seeds with FM-Indexes:

One builds two arrays

Relative Indexes & move from relative to Absolute indexes:

	0	1	2	3	4	5	6	7	8	9	0	1	2	3	4	5	6	7	8	9	2
	G	C	T	G	G	G	G	T	G	G	T	T	G	T	G	\$	A	T	T	G	C
\$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
C	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2
G	1	1	1	2	3	4	5	5	6	7	7	7	8	8	9	9	9	9	9	10	10
T	0	0	1	1	1	1	1	2	2	2	3	4	4	5	5	5	5	6	7	7	7
	(20)	(13)	(12)	(3)	(2)	(19)	(18)	(1)	(17)	(16)	(15)	(9)	(10)	(6)	(11)	(0)	(14)	(8)	(5)	(7)	(4)

c	\$	A	C	G	T
First c Ind	0	1	2	4	14

Seeds with FM-Indexes:

One builds two arrays

Relative Indexes & move from relative to Absolute indexes:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
	G	C	T	G	G	G	G	T	G	G	T	T	G	T	G	\$	A	T	T	G	C
\$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
C	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2
G	1	1	1	2	3	4	5	5	6	7	7	7	8	8	9	9	9	9	9	10	10
T	0	0	1	1	1	1	1	2	2	2	3	4	4	5	5	5	5	6	7	7	7
	(20)	(13)	(12)	(3)	(2)	(19)	(18)	(1)	(17)	(16)	(15)	(9)	(10)	(6)	(11)	(0)	(14)	(8)	(5)	(7)	(4)

c	\$	A	C	G	T
First c Ind	0	1	2	4	14

And notices that $[L(aP), R(aP)] = Abs(a) + [Rel_a(L(P) - 1), Rel_a(R(P)) - 1]$.

$\varepsilon \rightarrow [0, 20]$.

Hence G $\rightarrow 4 + [0, 10 - 1] = [4, 13]$

Hence GG $\rightarrow 4 + [2, 8 - 1] = [6, 11]$

Hence GGG $\rightarrow 4 + [4, 7 - 1] = [8, 10]$

Hence GGG appears at indexes $\{17, 16, 15\}$.

Seeds with FM-Indexes:

One builds two arrays

Relative Indexes & move from relative to Absolute indexes:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
	G	C	T	G	G	G	G	T	G	G	T	T	G	T	G	\$	A	T	T	G	C
\$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
C	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2
G	1	1	1	2	3	4	5	5	6	7	7	7	8	8	9	9	9	9	9	10	10
T	0	0	1	1	1	1	1	2	2	2	3	4	4	5	5	5	5	6	7	7	7
	(20)	(13)	(12)	(3)	(2)	(19)	(18)	(1)	(17)	(16)	(15)	(9)	(10)	(6)	(11)	(0)	(14)	(8)	(5)	(7)	(4)

<i>c</i>	\$	A	C	G	T
<i>First c Ind</i>	0	1	2	4	14

And notices that $[L(aP), R(aP)] = Abs(a) + [Rel_a(L(P) - 1), Rel_a(R(P)) - 1]$.

$\varepsilon \rightarrow [0, 20]$.

Hence G $\rightarrow 4 + [0, 10 - 1] = [4, 13]$

Hence GG $\rightarrow 4 + [2, 8 - 1] = [6, 11]$

Hence GGG $\rightarrow 4 + [4, 7 - 1] = [8, 10]$

Hence GGG appears at indexes $\{17, 16, 15\}$.

Seeds with FM-Indexes:

One builds two arrays

Relative Indexes & move from relative to Absolute indexes:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
	G	C	T	G	G	G	G	T	G	G	T	T	G	T	G	\$	A	T	T	G	C
\$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
C	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2
G	1	1	1	2	3	4	5	5	6	7	7	7	8	8	9	9	9	9	9	10	10
T	0	0	1	1	1	1	1	2	2	2	3	4	4	5	5	5	5	6	7	7	7
	(20)	(13)	(12)	(3)	(2)	(19)	(18)	(1)	(17)	(16)	(15)	(9)	(10)	(6)	(11)	(0)	(14)	(8)	(5)	(7)	(4)

<i>c</i>	\$	A	C	G	T
<i>First c Ind</i>	0	1	2	4	14

And notices that $[L(aP), R(aP)] = Abs(a) + [Rel_a(L(P) - 1), Rel_a(R(P)) - 1]$.

$\varepsilon \rightarrow [0, 20]$.

Hence $G \rightarrow 4 + [0, 10 - 1] = [4, 13]$

Hence $GG \rightarrow 4 + [2, 8 - 1] = [6, 11]$

Hence $GGG \rightarrow 4 + [4, 7 - 1] = [8, 10]$

Hence GGG appears at indexes $\{17, 16, 15\}$.

Seeds with FM-Indexes:

One builds two arrays

Relative Indexes & move from relative to Absolute indexes:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
	G	C	T	G	G	G	G	T	G	G	T	T	G	T	G	\$	A	T	T	G	C
\$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
C	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2
G	1	1	1	2	3	4	5	5	6	7	7	7	8	8	9	9	9	9	9	9	10
T	0	0	1	1	1	1	1	2	2	2	3	4	4	5	5	5	5	6	7	7	7
	(20)	(13)	(12)	(3)	(2)	(19)	(18)	(1)	(17)	(16)	(15)	(9)	(10)	(6)	(11)	(0)	(14)	(8)	(5)	(7)	(4)

<i>c</i>	\$	A	C	G	T
<i>First c Ind</i>	0	1	2	4	14

And notices that $[L(aP), R(aP)] = Abs(a) + [Rel_a(L(P) - 1), Rel_a(R(P)) - 1]$.

$\varepsilon \rightarrow [0, 20]$.

Hence $G \rightarrow 4 + [0, 10 - 1] = [4, 13]$

Hence $GG \rightarrow 4 + [2, 8 - 1] = [6, 11]$

Hence $GGG \rightarrow 4 + [4, 7 - 1] = [8, 10]$

Hence GGG appears at indexes $\{17, 16, 15\}$.

Seeds with FM-Indexes:

One builds two arrays

Relative Indexes & move from relative to Absolute indexes:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
	G	C	T	G	G	G	G	T	G	G	T	T	G	T	G	\$	A	T	T	G	C
\$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
C	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2
G	1	1	1	2	3	4	5	5	6	7	7	7	8	8	9	9	9	9	9	9	10
T	0	0	1	1	1	1	1	2	2	2	3	4	4	5	5	5	5	6	7	7	7
	(20)	(13)	(12)	(3)	(2)	(19)	(18)	(1)	(17)	(16)	(15)	(9)	(10)	(6)	(11)	(0)	(14)	(8)	(5)	(7)	(4)

<i>c</i>	\$	A	C	G	T
<i>First c Ind</i>	0	1	2	4	14

And notices that $[L(aP), R(aP)] = Abs(a) + [Rel_a(L(P) - 1), Rel_a(R(P)) - 1]$.

$\varepsilon \rightarrow [0, 20]$.

Hence $G \rightarrow 4 + [0, 10 - 1] = [4, 13]$

Hence $GG \rightarrow 4 + [2, 8 - 1] = [6, 11]$

Hence $GGG \rightarrow 4 + [4, 7 - 1] = [8, 10]$

Hence GGG appears at indexes $\{17, 16, 15\}$.

Seeds with FM-Indexes:

One builds two arrays

Relative Indexes & move from relative to Absolute indexes:

	0	1	2	3	4	5	6	7	8	9	10	1	2	3	4	5	6	7	8	9	20	
	G	C	T	G	G	G	G	T	G	G	T	T	G	T	G	\$	A	T	T	G	C	
\$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1	
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	
C	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2	
G	1	1	1	2	3	4	5	5	6	7	7	7	8	8	9	9	9	9	9	9	10	10
T	0	0	1	1	1	1	1	2	2	2	3	4	4	5	5	5	5	6	7	7	7	
	(20)	(13)	(12)	(3)	(2)	(19)	(18)	(1)	(17)	(16)	(15)	(9)	(10)	(6)	(11)	(0)	(14)	(8)	(5)	(7)	(4)	

<i>c</i>	\$	A	C	G	T
<i>First c Ind</i>	0	1	2	4	14

And notices that $[L(aP), R(aP)] = Abs(a) + [Rel_a(L(P) - 1), Rel_a(R(P)) - 1]$.

$\varepsilon \rightarrow [0, 20]$.

Hence $G \rightarrow 4 + [0, 10 - 1] = [4, 13]$

Hence $GG \rightarrow 4 + [2, 8 - 1] = [6, 11]$

Hence $GGG \rightarrow 4 + [4, 7 - 1] = [8, 10]$

Hence GGG appears at indexes $\{17, 16, 15\}$.

Seeds with FM-Indexes:

One builds two arrays

Relative Indexes & move from relative to Absolute indexes:

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20
	G	C	T	G	G	G	G	T	G	G	T	T	G	T	G	\$	A	T	T	G	C
\$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1	1
A	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	1	1	1	1
C	0	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	2
G	1	1	1	2	3	4	5	5	6	7	7	7	8	8	9	9	9	9	9	10	10
T	0	0	1	1	1	1	1	2	2	2	3	4	4	5	5	5	5	6	7	7	7
	(20)	(13)	(12)	(3)	(2)	(19)	(18)	(1)	(17)	(16)	(15)	(9)	(10)	(6)	(11)	(0)	(14)	(8)	(5)	(7)	(4)

c	\$	A	C	G	T
First c Ind	0	1	2	4	14

(Each integer should be stored on at least 4 bytes so a naive storage of this array would take at least $17 \times |Ref|$)

The Smith-Waterman Algorithm

Now that seeds are found, let's grow them! (dynamic programming)

At Index 1 we find:

GGCTT

How does it match *GGACT*?

For instance we suppose that

An exact match is $+3$, A Mismatch is -3 , Insertion or Deletion is -2

Let's find the best matching. . .

		G	G	C	T	T
	0	0	0	0	0	0
G	0					
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	-3				
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	-3				
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	3				
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	3				
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	3	0			
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	3	-3			
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	3	3			
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	3	3			
G	0					
A	0					
C	0					
T	0					

		G	G	C	T	T
	0	0	0	0	0	0
G	0	3	3	0	-3	-3
G	0	3	6	3	0	-3
A	0	0	3	3	0	-3
C	0	-3	0	6	3	0
T	0	-3	-6	3	9	6

		G	G	C	T	T
	0	0	0	0	0	0
G	0	3	3	0	-3	-3
G	0	3	6	3	0	-3
A	0	0	3	3	0	-3
C	0	-3	0	6	3	0
T	0	-3	-6	3	9	6

G G - C T T
G G A C T -

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Brujin graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines**
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion



- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Brujin graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem**
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion

We have maize lines

<i>Ind</i>	<i>M/A</i>	<i>l</i> ₁	<i>l</i> ₂	<i>l</i> ₃	...	<i>l</i> _{<i>n</i>}
<i>L</i> ₁	<i>A/C</i>	<i>AA</i>	<i>CC</i>	<i>CC</i>	...	<i>AA</i>
<i>L</i> ₂	<i>T/A</i>	<i>TT</i>	<i>TT</i>	<i>AA</i>	...	<i>TT</i>
⋮	⋮	⋮	⋮	⋮	⋮	⋮
<i>L</i> _{<i>L</i>}	<i>G/C</i>	<i>GG</i>	<i>CC</i>	<i>GG</i>	...	<i>GG</i>

We have maize lines

		l_1	l_2	l_3	...	l_n
L_1	AC	A	C	C	...	A
L_2	TA	T	T	A	...	T
\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
L_L	GC	G	C	G	...	G

We have maize lines

		l_1	l_2	l_3	\dots	L_n
L_1	<i>AC</i>	1	0	0	\dots	1
L_2	<i>TA</i>	1	1	0	\dots	1
\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
L_L	<i>GC</i>	1	0	1	\dots	0

BUT some of the contents are unknown

		l_1	l_2	l_3	...	l_n
L_1	<i>AC</i>	1	?	0	...	1
L_2	<i>TA</i>	?	?	0	...	1
\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
L_L	<i>GC</i>	1	0	?	...	?

and our mission (if we accept it) is to fill missing data by using statistical methods.

(this file format is a very simplified VCF file format)

BUT some of the contents are unknown

		l_1	l_2	l_3	...	l_n
L_1	<i>AC</i>	1	?	0	...	1
L_2	<i>TA</i>	?	?	0	...	1
\vdots	\vdots	\vdots	\vdots	\vdots	\ddots	\vdots
L_L	<i>GC</i>	1	0	?	...	?

and our mission (if we accept it) is to fill missing data by using statistical methods.

(this file format is a very simplified VCF file format)

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Brujin graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)**
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion

A classical resolution consists in:

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	0	?	1	1	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

A classical resolution consists in:

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	0	?	1	1	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

Identifying haplotypes blocks

A classical resolution consists in:

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	0	0	1	1	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	0	?	1	1	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

Then apply majority among blocks *having the same known data*

This is a “natural” way in the sense that it respects biology and its accuracy is very high. But

- 1 the computation is very long
- 2 It predicts a few heterozygous loci
- 3 when the missing data rate is very high, the accuracy drops down

Can we approach this accuracy with much less computation time?

This is a “natural” way in the sense that it respects biology and its accuracy is very high. But

- 1 the computation is very long
- 2 It predicts a few heterozygous loci
- 3 when the missing data rate is very high, the accuracy drops down

Can we approach this accuracy with much less computation time?

This is a “natural” way in the sense that it respects biology and its accuracy is very high. But

- 1 the computation is very long
- 2 It predicts a few heterozygous loci
- 3 when the missing data rate is very high, the accuracy drops down

Can we approach this accuracy with much less computation time?

This is a “natural” way in the sense that it respects biology and its accuracy is very high. But

- 1 the computation is very long
- 2 It predicts a few heterozygous loci
- 3 when the missing data rate is very high, the accuracy drops down

Can we approach this accuracy with much less computation time?

This is a “natural” way in the sense that it respects biology and its accuracy is very high. But

- 1 the computation is very long
- 2 It predicts a few heterozygous loci
- 3 when the missing data rate is very high, the accuracy drops down

Can we approach this accuracy with much less computation time?

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Brujin graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains**
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion

This is a purely computer-oriented approach, blind to biology

We first choose a number m (the “Markov Chain Order”).

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	0	?	1	1	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

m may be interpreted as a “recall length”
(In this example, $m = 3$).

Then for each and every locus:

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	0	?	1	1	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

then at this given locus one numbers the occurrences of consecutive given letters
(number the 000 followed by 0, 000 by 1, 001 by 0, 001 by 1 ...)

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	0	?	1	1	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

... this for each Ind

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	0	?	1	1	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

For instance here 011 \rightarrow 0 has been seen +1 time here (same locus, Ind varies)

But what if there is missing data?

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	0	?	1	1	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

But what if all data is not seen? ?0? → 1 is turned into

000 → 1 : +0.25

001 → 1 : +0.25

100 → 1 : +0.25

101 → 1 : +0.25

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	?	?	?	?	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

But what if all data is not seen? ??? \rightarrow ? could cause increments by $\frac{1}{2^{m+1}}$ which seems useless (won't change relative orders).

	L32	L33	L34	L35	L36	L37	L38	L39	L40	L41	L42	L43	L44	L45
I27	?	0	1	?	0	0	1	1	0	?	0	1	1	0
I28	1	?	?	1	?	?	0	0	?	1	1	?	?	0
I29	?	0	1	1	0	0	0	?	1	?	?	0	0	?
I30	1	?	0	?	1	?	1	0	0	0	?	1	?	0
I31	?	?	?	?	?	0	?	1	1	0	0	?	1	1
I32	?	1	0	0	1	0	0	?	1	?	1	0	?	1
I33	?	?	1	0	1	?	0	1	?	0	?	1	1	?
I34	1	1	1	0	0	1	0	1	1	0	0	1	?	1
I35	?	0	0	1	1	0	?	0	?	1	?	1	0	1
I36	1	1	?	0	?	1	1	0	0	0	1	0	0	0

But what if all data is not seen? ??? → ? could cause increments by $\frac{1}{2^{m+1}}$ which seems useless (won't change relative orders).

We chose to do nothing (no increment)

These matrices of transitions are a markov chain of order m .
How should we choose m ? First it has a lot to do with the necessary memory for the process.

If the number of loci is L , the necessary amount of memory is $(L - m)2^{m+3}$ bytes.

That is to say that if $L = 300,000$ and $m = 10$ you will need 2.5 GB of memory.

The amount of necessary memory will *double* if m is *incremented by 1!*
usually $m \leq 4$.

But in our case we will make it run with bigger values of m .

These matrices of transitions are a markov chain of order m .
How should we choose m ? First it has a lot to do with the necessary memory for the process.

If the number of loci is L , the necessary amount of memory is $(L - m)2^{m+3}$ bytes.

That is to say that if $L = 300,000$ and $m = 10$ you will need 2.5 GB of memory.

The amount of necessary memory will *double* if m is *incremented by 1!*
usually $m \leq 4$.

But in our case we will make it run with bigger values of m .

These matrices of transitions are a markov chain of order m .
How should we choose m ? First it has a lot to do with the necessary memory for the process.

If the number of loci is L , the necessary amount of memory is $(L - m)2^{m+3}$ bytes.

That is to say that if $L = 300,000$ and $m = 10$ you will need 2.5 GB of memory.

The amount of necessary memory will *double* if m is *incremented by 1!*
usually $m \leq 4$.

But in our case we will make it run with bigger values of m .

These matrices of transitions are a markov chain of order m .
How should we choose m ? First it has a lot to do with the necessary memory for the process.

If the number of loci is L , the necessary amount of memory is $(L - m)2^{m+3}$ bytes.

That is to say that if $L = 300,000$ and $m = 10$ you will need 2.5 GB of memory.

The amount of necessary memory will *double* if m is *incremented by 1!*
usually $m \leq 4$.

But in our case we will make it run with bigger values of m .

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Brujin graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing**
- 9 Performances
- 10 An R package
- 11 Conclusion

At this point we may consider that for a given Ind, for every locus we have the probability that m consecutive letters are followed by a 0 or a 1.

This probability depends on the locus. It is a local property.

These probabilities spread through the whole sequence.

We compute $P(I, L) = 1$ for each locus of a given Ind, taking into account the spreading of the probabilities

Given a threshold, one can now impute the data itself.

For instance if threshold = 0.9, if $P(I, L) > 0.9$ we impute with a 1. If $P(I, L) < 1 - 0.9 = 0.1$, we impute with 0.

In between we do not impute and the unknown char remains unknown.

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Brujin graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances**
- 10 An R package
- 11 Conclusion

Known & Full VCF files consisted in 381 lines and 81,678 loci (chromosome 1).

From these files were created VCF files with missing data at rates (10%, 20%, . . . , 90%).

Then we imputed these files and compared results with Beagle.

We had several ways to remove known data:

- 1 choose Ind and locus at random
- 2 (coming soon): choose a number of untouched Ind, then choose some loci at random and erase all data for these loci for all Ind (except the first ones)

Time computation: (computations have been performed on average computers)

depending on the missing data rate, Beagle averagely needed between 2 and 3 hours to impute a chromosome

The same work was made by LinelImputer in a time $< 20''$ (Forward backward and Viterbi imputations both included).

Beware: these times are a rough guide only: softwares did not run on the same machines.

BUT:

Time computation: (computations have been performed on average computers)

depending on the missing data rate, Beagle averagely needed between 2 and 3 hours to impute a chromosome

The same work was made by LinelImputer in a time $< 20''$ (Forward backward and Viterbi imputations both included).

Beware: these times are a rough guide only: softwares did not run on the same machines.

BUT:

Time computation: (computations have been performed on average computers)

depending on the missing data rate, Beagle averagely needed between 2 and 3 hours to impute a chromosome

The same work was made by LinImputer in a time $< 20''$ (Forward backward and Viterbi imputations both included).

Beware: these times are a rough guide only: softwares did not run on the same machines.

BUT:

Time computation: (computations have been performed on average computers)

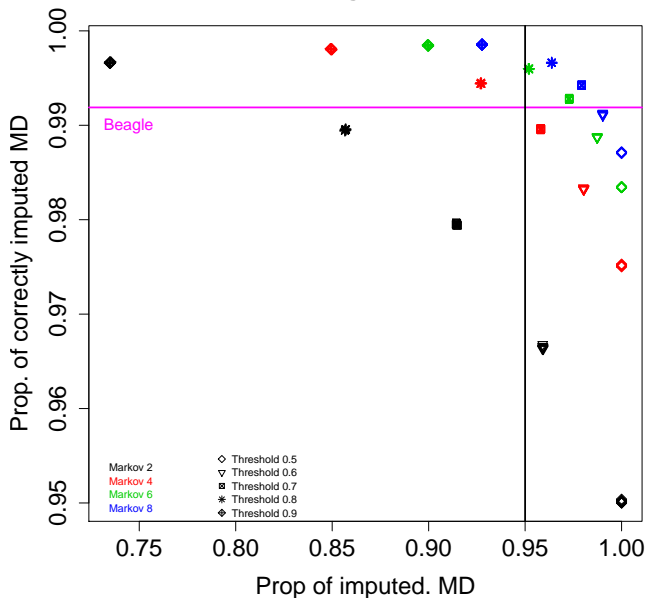
depending on the missing data rate, Beagle averagely needed between 2 and 3 hours to impute a chromosome

The same work was made by LinImputer in a time $< 20''$ (Forward backward and Viterbi imputations both included).

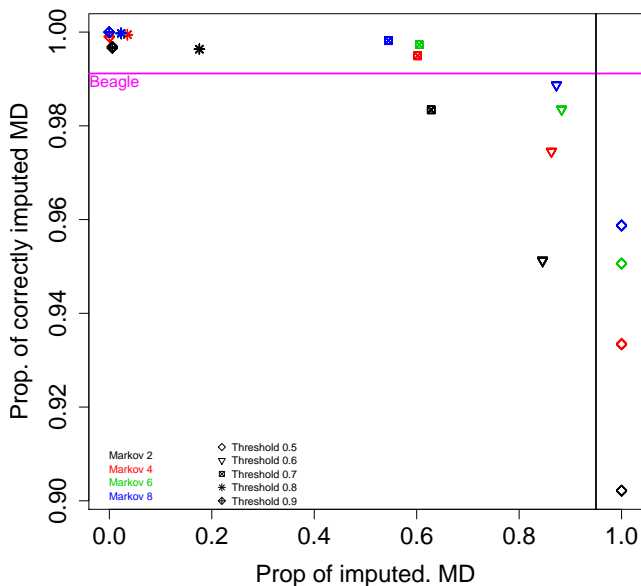
Beware: these times are a rough guide only: softwares did not run on the same machines.

BUT:

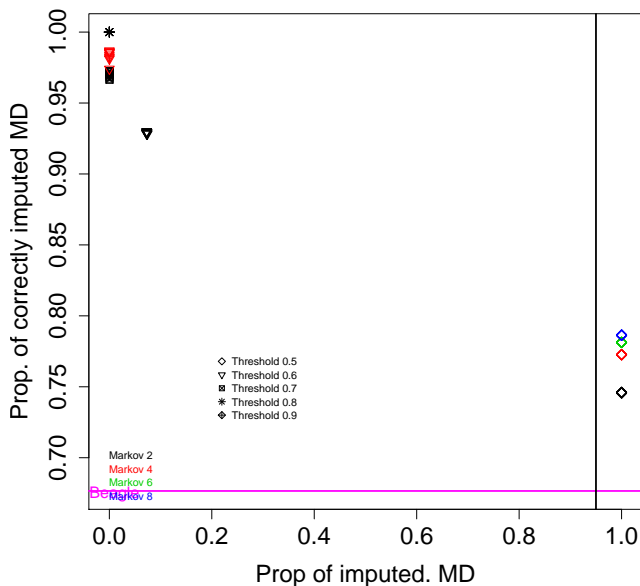
Missing rate: 0.1



Missing rate: 0.5



Missing rate: 0.9



- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Brujin graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package**
- 11 Conclusion

(Coming soon): an *R* package allowing one to

- 1 merge VCF files (linear time and space for 2 files)
- 2 Impute a VCF file. One chooses $2 \leq m \leq 12$ (the Markov chain order) and the Threshold. It produces a file containing probabilities and imputed files using forward/backward and/or Viterbi imputed files (linear time and space complexity)

(Coming soon): an *R* package allowing one to

- 1 merge VCF files (linear time and space for 2 files)
- 2 Impute a VCF file. One chooses $2 \leq m \leq 12$ (the Markov chain order) and the Threshold. It produces a file containing probabilities and imputed files using forward/backward and/or Viterbi imputed files (linear time and space complexity)

- 1 Mapping problem
- 2 De Novo Sequencing (exact match)
 - The de Bruijn graph
- 3 Usual technologies
 - Seed...
 - ... & Grow!
- 4 Imputing maze lines
- 5 The Problem
- 6 Classical Resolution (Beagle)
- 7 Markov Chains
- 8 Imputing
- 9 Performances
- 10 An R package
- 11 Conclusion**

LineImputer is (will soon be actually) an R package based on a Markov chain approach.

It allows one to merge VCF files.

It imputes provided the user specifies m and the threshold

It is very fast

It's accuracy may be good (depending on the choices of the Markov Chain Order and the Threshold) without outbesting Beagle, except for very high missing data rates

The quality of imputation decreases (possibly fast) when the threshold decreases towards 0.5, depending on the choice of m .

LineImputer is (will soon be actually) an R package based on a Markov chain approach.

It allows one to merge VCF files.

It imputes provided the user specifies m and the threshold

It is very fast

It's accuracy may be good (depending on the choices of the Markov Chain Order and the Threshold) without outbesting Beagle, except for very high missing data rates

The quality of imputation decreases (possibly fast) when the threshold decreases towards 0.5, depending on the choice of m .

LineImputer is (will soon be actually) an R package based on a Markov chain approach.

It allows one to merge VCF files.

It imputes provided the user specifies m and the threshold

It is very fast

It's accuracy may be good (depending on the choices of the Markov Chain Order and the Threshold) without outbesting Beagle, except for very high missing data rates

The quality of imputation decreases (possibly fast) when the threshold decreases towards 0.5, depending on the choice of m .

LineImputer is (will soon be actually) an R package based on a Markov chain approach.

It allows one to merge VCF files.

It imputes provided the user specifies m and the threshold

It is very fast

It's accuracy may be good (depending on the choices of the Markov Chain Order and the Threshold) without outbesting Beagle, except for very high missing data rates

The quality of imputation decreases (possibly fast) when the threshold decreases towards 0.5, depending on the choice of m .

LineImputer is (will soon be actually) an R package based on a Markov chain approach.

It allows one to merge VCF files.

It imputes provided the user specifies m and the threshold

It is very fast

It's accuracy may be good (depending on the choices of the Markov Chain Order and the Threshold) without outbesting Beagle, except for very high missing data rates

The quality of imputation decreases (possibly fast) when the threshold decreases towards 0.5, depending on the choice of m .

LineImputer is (will soon be actually) an R package based on a Markov chain approach.

It allows one to merge VCF files.

It imputes provided the user specifies m and the threshold

It is very fast

It's accuracy may be good (depending on the choices of the Markov Chain Order and the Threshold) without outbesting Beagle, except for very high missing data rates

The quality of imputation decreases (possibly fast) when the threshold decreases towards 0.5, depending on the choice of m .

Future work: try to improve accuracy, for instance by:

- 1 choose a value of m (quite big) and a big threshold
- 2 impute (this will necessitate a very smart management of hardware for big values of m , in particular the available memory)
- 3 use the output file as an input file
- 4 decrease m (do not touch the threshold)
- 5 if a lot remains to impute go back to step 2
- 6 else end the process

The goal would be to get comparable performances to the ones of Beagle but in much less time

Future work: try to improve accuracy, for instance by:

- 1 choose a value of m (quite big) and a big threshold
- 2 impute (this will necessitate a very smart management of hardware for big values of m , in particular the available memory)
- 3 use the output file as an input file
- 4 decrease m (do not touch the threshold)
- 5 if a lot remains to impute go back to step 2
- 6 else end the process

The goal would be to get comparable performances to the ones of Beagle but in much less time

Future work: try to improve accuracy, for instance by:

- 1 choose a value of m (quite big) and a big threshold
- 2 impute (this will necessitate a very smart management of hardware for big values of m , in particular the available memory)
- 3 use the output file as an input file
- 4 decrease m (do not touch the threshold)
- 5 if a lot remains to impute go back to step 2
- 6 else end the process

The goal would be to get comparable performances to the ones of Beagle but in much less time

Future work: try to improve accuracy, for instance by:

- 1 choose a value of m (quite big) and a big threshold
- 2 impute (this will necessitate a very smart management of hardware for big values of m , in particular the available memory)
- 3 use the output file as an input file
- 4 decrease m (do not touch the threshold)
- 5 if a lot remains to impute go back to step 2
- 6 else end the process

The goal would be to get comparable performances to the ones of Beagle but in much less time

Future work: try to improve accuracy, for instance by:

- 1 choose a value of m (quite big) and a big threshold
- 2 impute (this will necessitate a very smart management of hardware for big values of m , in particular the available memory)
- 3 use the output file as an input file
- 4 decrease m (do not touch the threshold)
- 5 if a lot remains to impute go back to step 2
- 6 else end the process

The goal would be to get comparable performances to the ones of Beagle but in much less time

Future work: try to improve accuracy, for instance by:
use Line Imputer to impute a vast majority of Data
and then use Beagle to impute the unimputed data

Future work: try to improve accuracy, for instance by:

simply use bigger values for m (take care of the use of memory in this case)

Thank you for your attention.
Questions?