

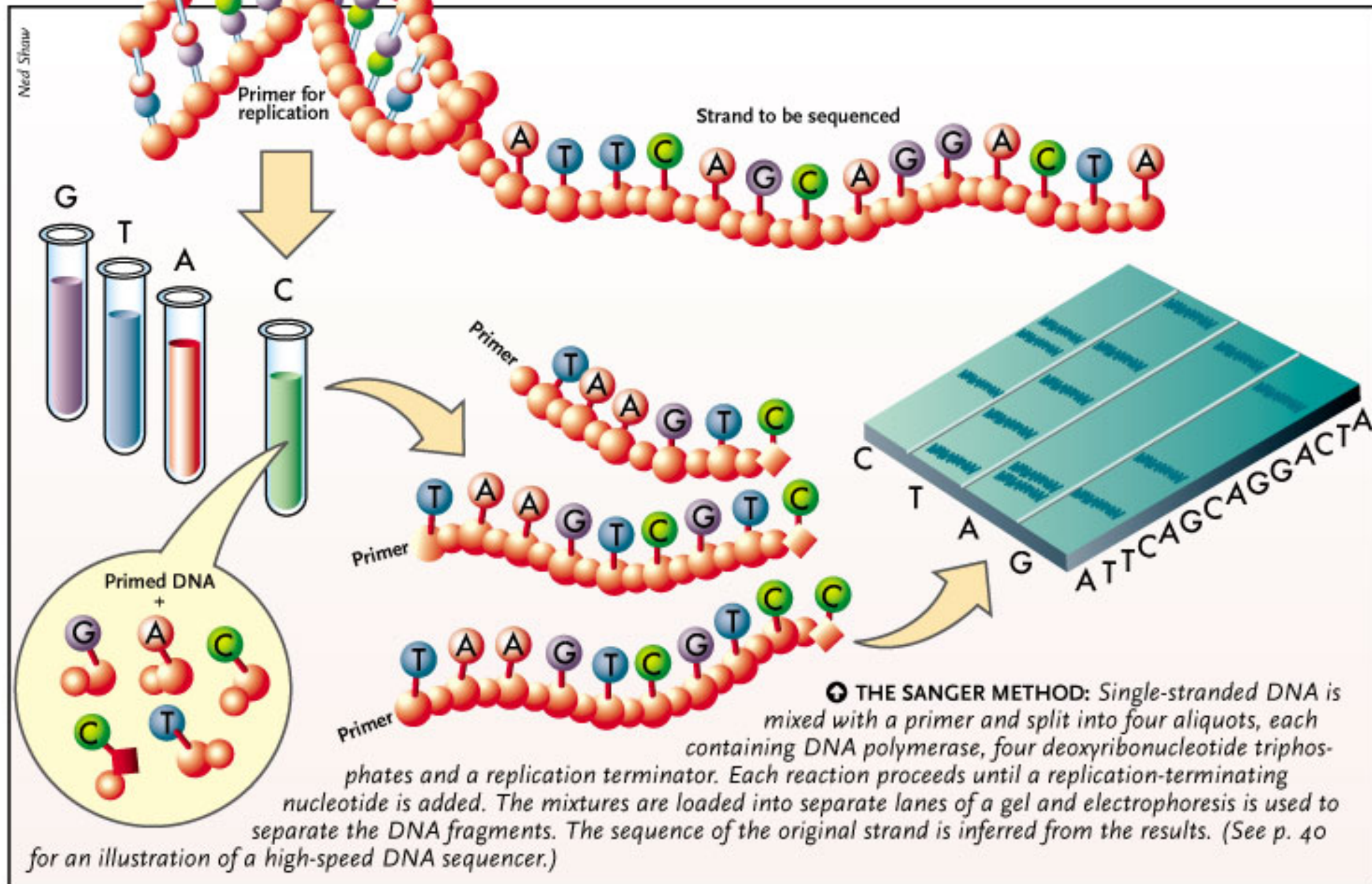
Detecting Gene Mutations in Cancer

androwis abumoussa | janice spence phd | john spence phd | richard burack md/phd

Background

sanger sequencing | next generation sequencing

DNA Sequencing



DNA Sequencing

biology

g c g t a a g a c c g c g t a a g a c c g c g t
g c a t t c t g g c g c a t t c t g g c g c a

alignment

[illegible]

Reference Position :

1

2

3

4

5

6

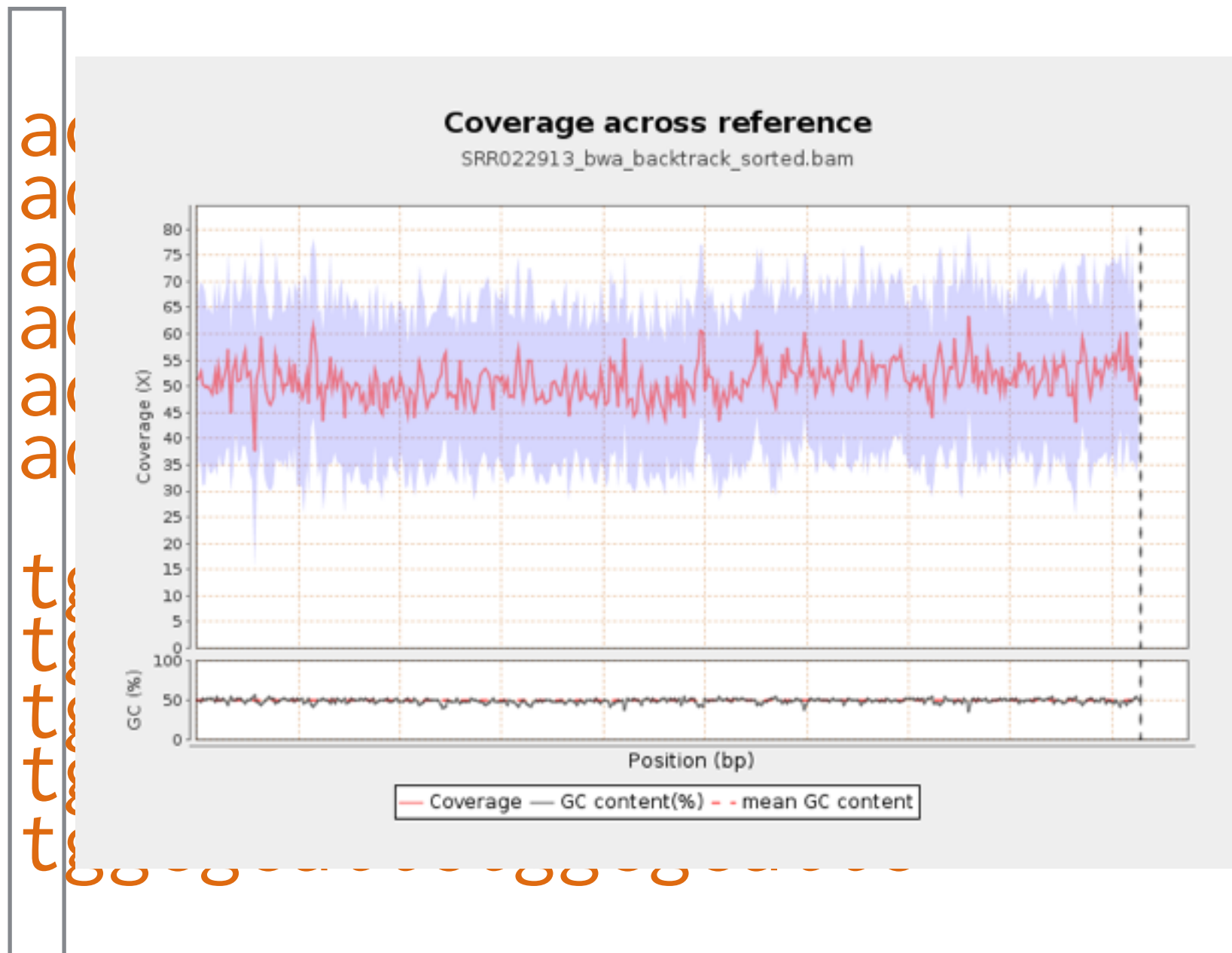
7

.....

[illegible]

DNA Sequencing

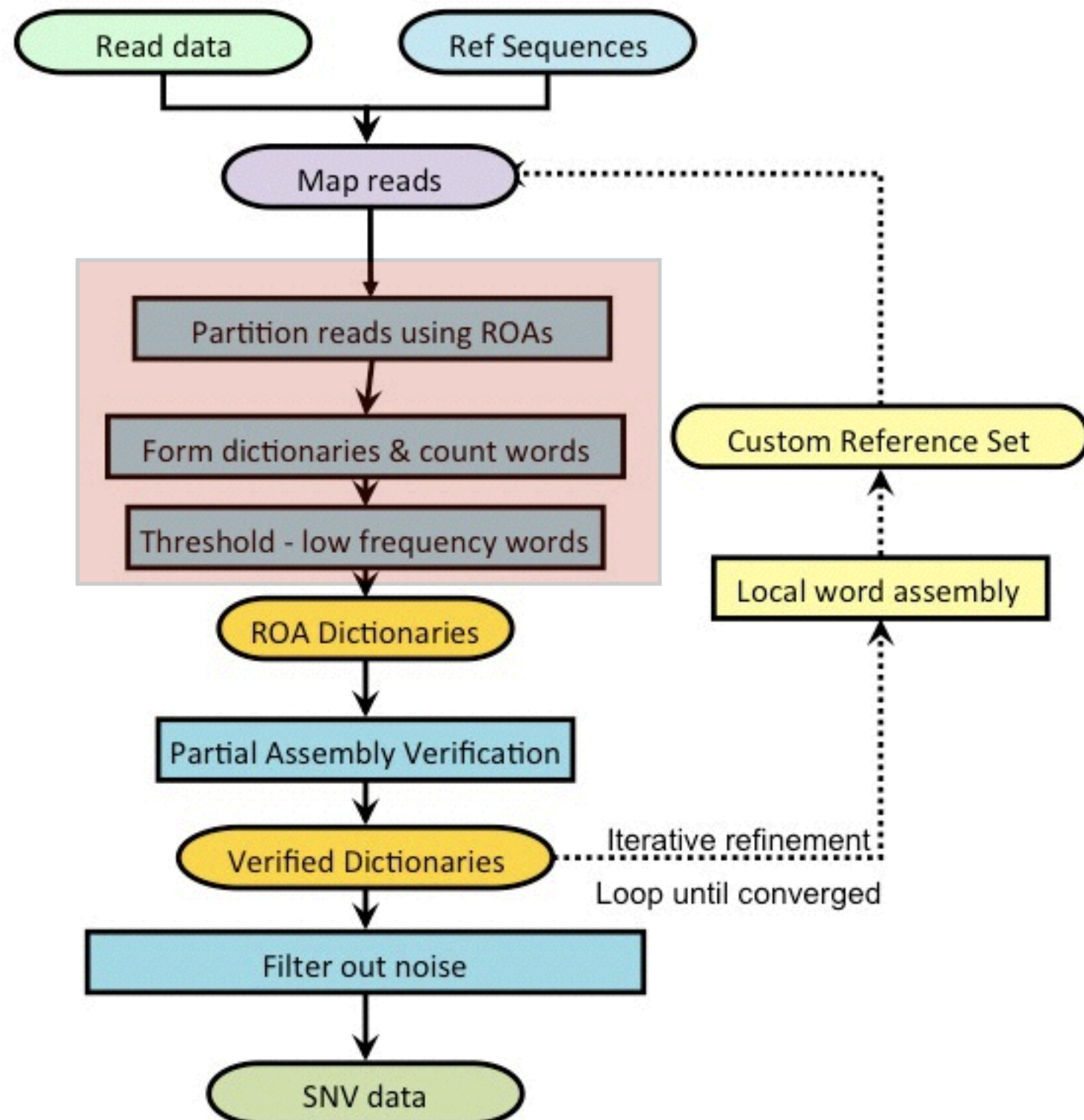
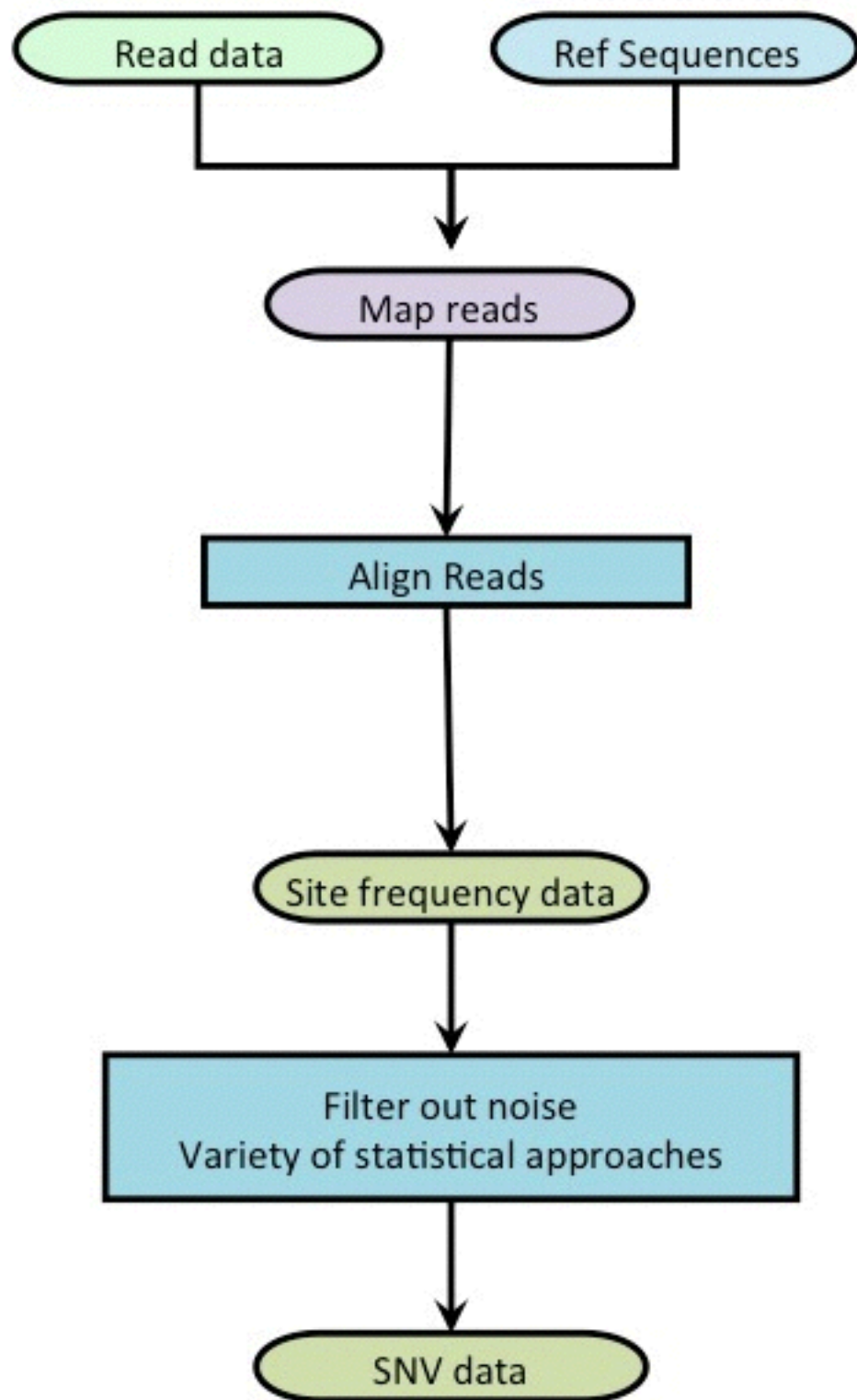
map reduce



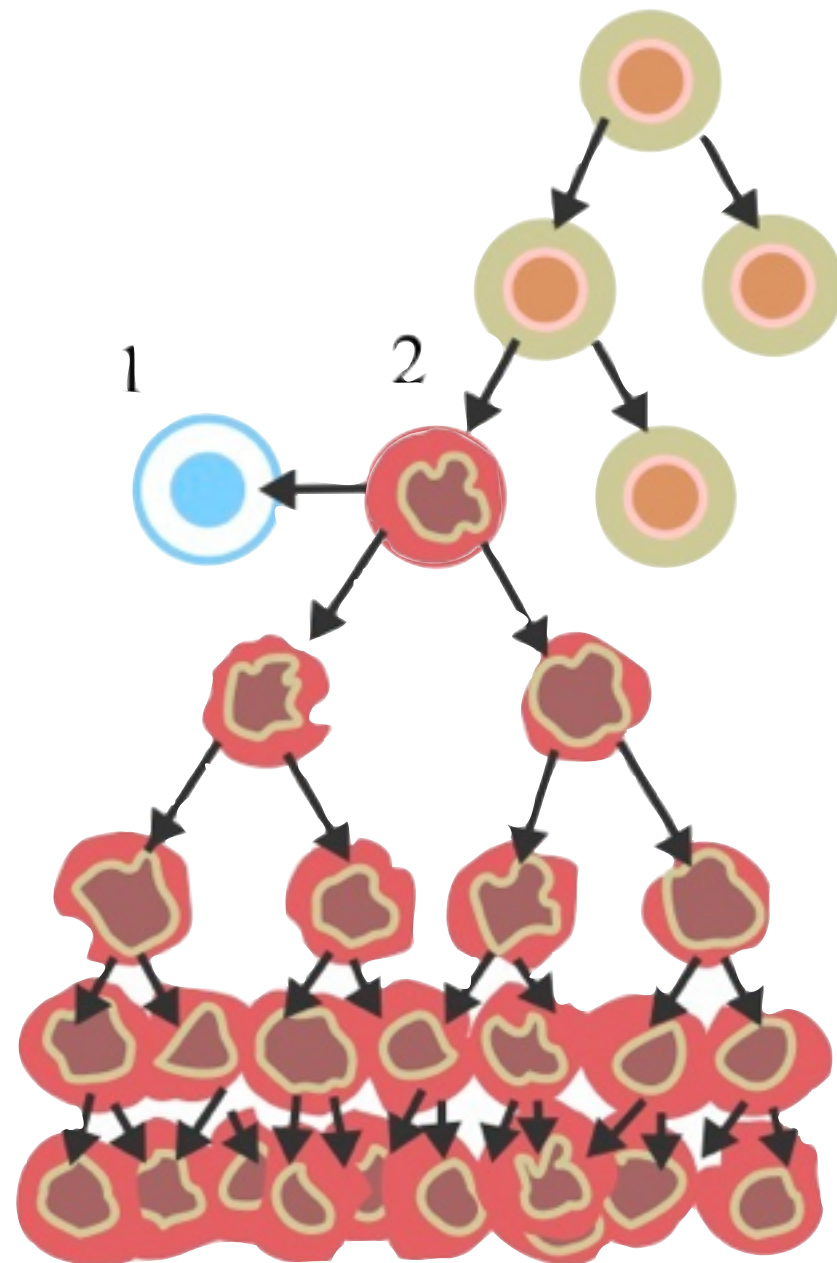
DNA Sequencing | Alignment

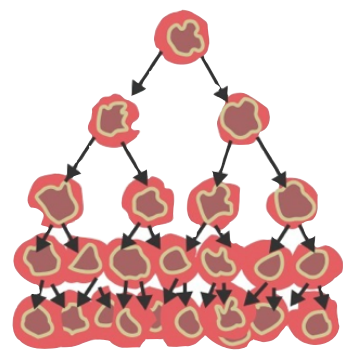
- Maximize Similarity
- Bound Search Space
- Make assumptions to do this

Targeted Resequencing



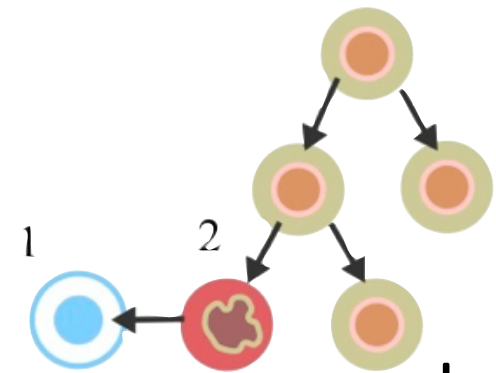
Cell Pathways



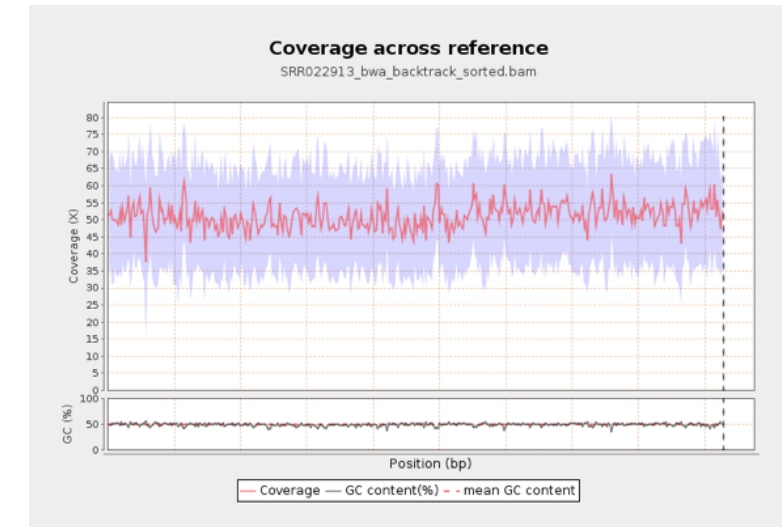
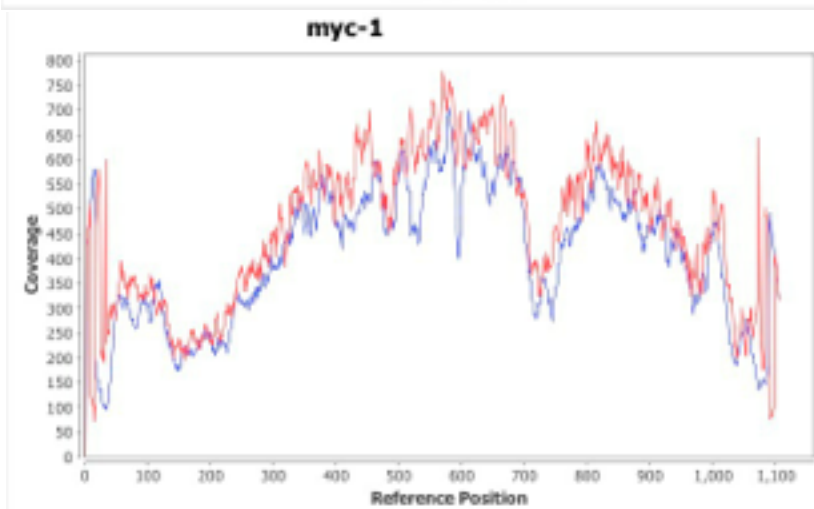
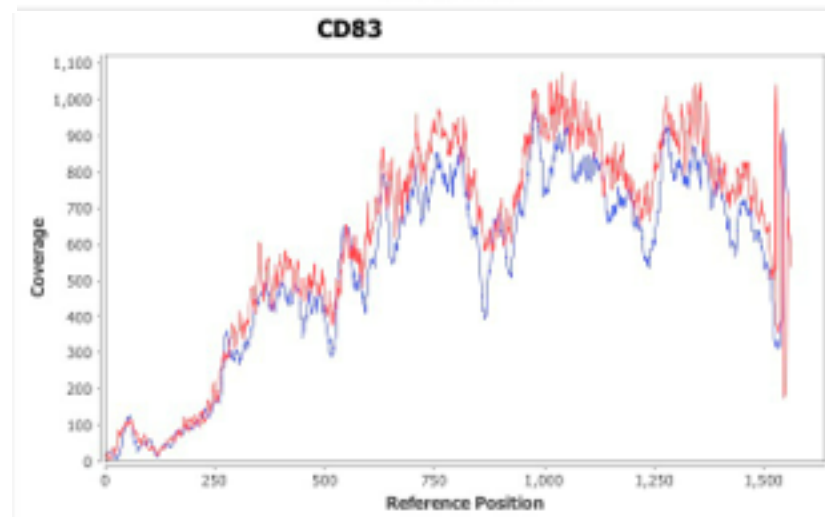
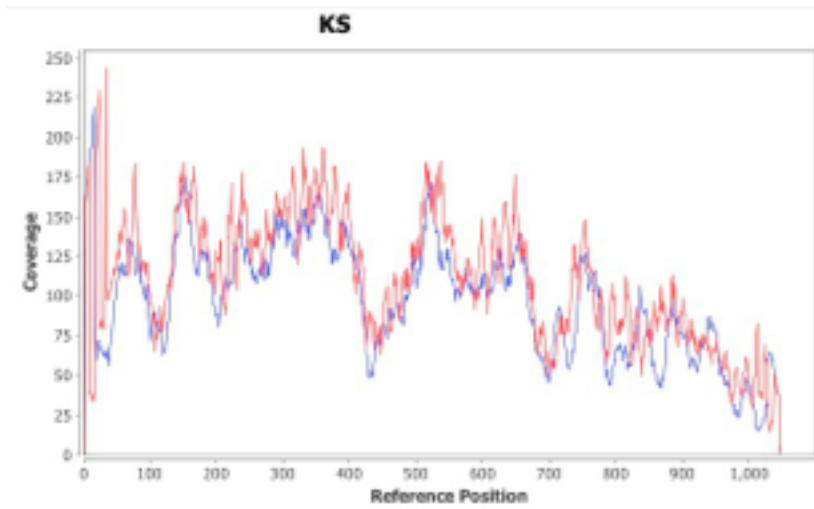


cancer

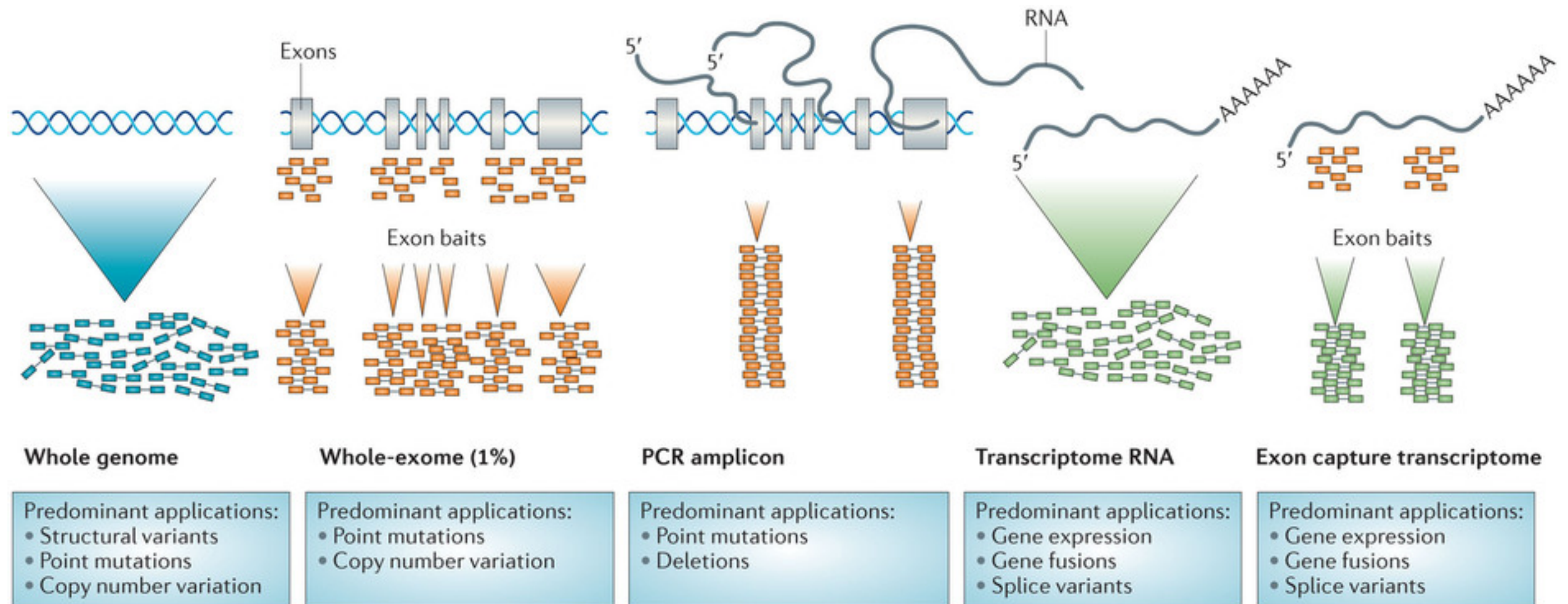
Cell Pathways



normal



Next Gen Sequencing



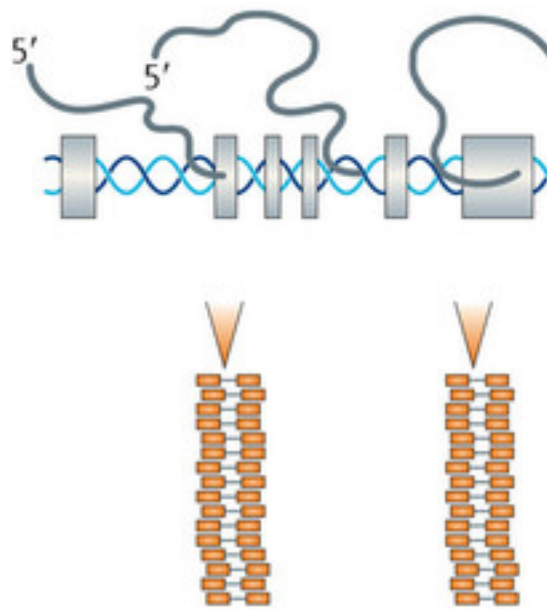
Big Data in Genetics

Parallel Data Structures | Map Reduce | Distributed Computing

Parallel Data Structures

trie | burrow-wheeler's algorithm

Trie Formation



PCR amplicon

Predominant applications:

- Point mutations
- Deletions

```
read : accgcgtaag.....
read : actgcgttag.....
read : actgcgtatg.....
read : accgcgtaag.....
read : actgcgttag.....
read : actgcgtatg.....
read : accgcgtaag.....
read : accgcgtaag.....
read : actgcgttag.....
read : actgcgtatg.....
read : accgcgtaag.....
read : actgcgttag.....
read : actgcgtatg.....
read : accgcgtaag.....
```


Trie Formation

read : accgcgtaag

Step 1 : create suffixes

read : accgcgtaag\$
ccgcgtaag\$
cgcgtaag\$
gcgtaag\$
cgtaag\$
gtaag\$
taag\$
aag\$
ag\$
g\$
\$

Step 2 : alphabetize suffixes

read : aag\$
accgcgtaag\$
ag\$
ccgcgtaag\$
cgcgtaag\$
cgtaag\$
gcgtaag\$
gtaag\$
g\$
taag\$
\$

Trie Formation

```
read  : accgcgtaag
hash  : 1346579028
```

Step 2 : alphabetize suffixes

```
read : 0 aag$
```

1 accgcgtaag\$

2 ag\$

3 ccgcgtaag\$

4 cgcgtaag\$

5 cgtaag\$

6 gcgtaag\$

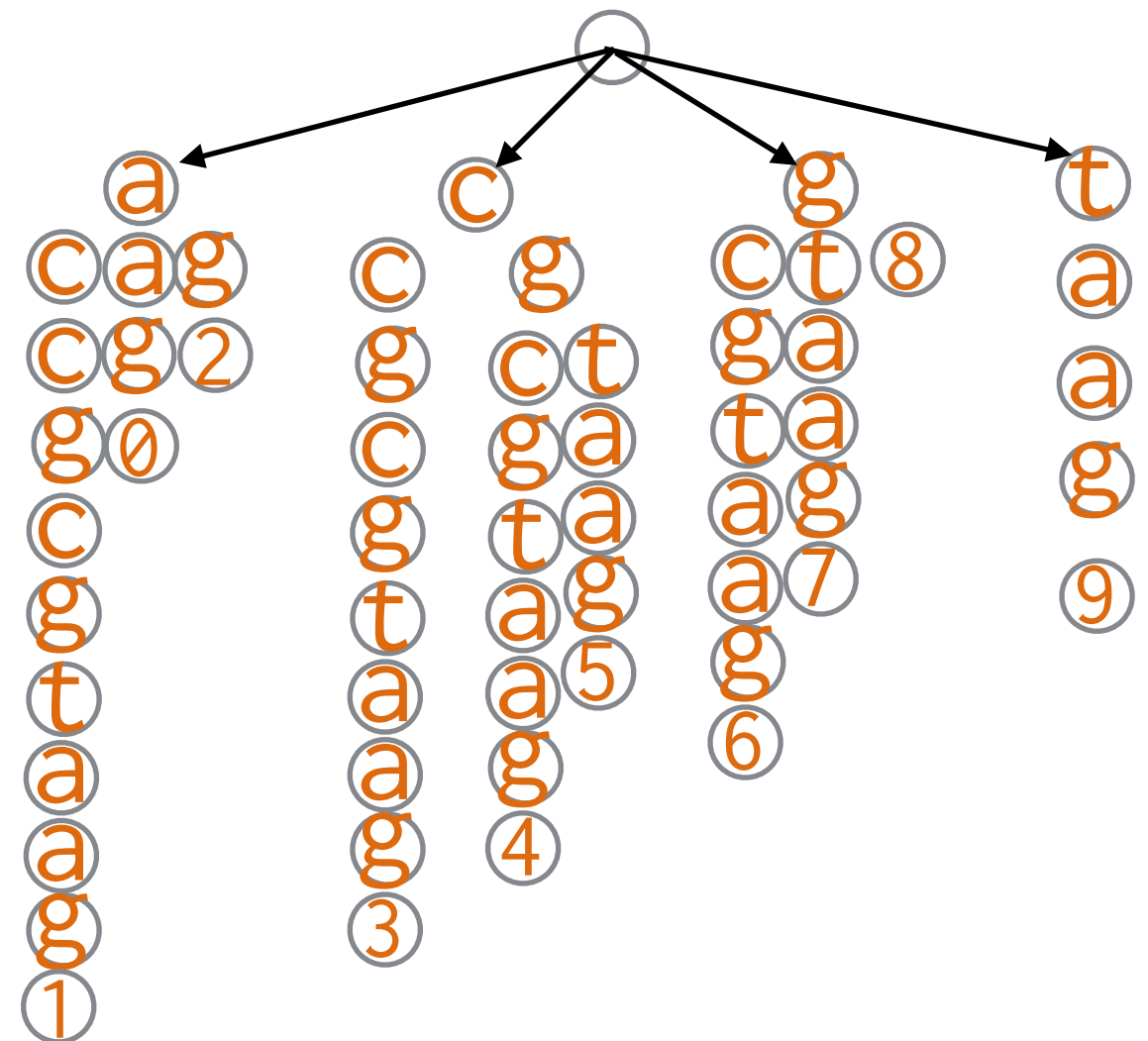
7 gtaag\$

8 g\$

9 taag\$

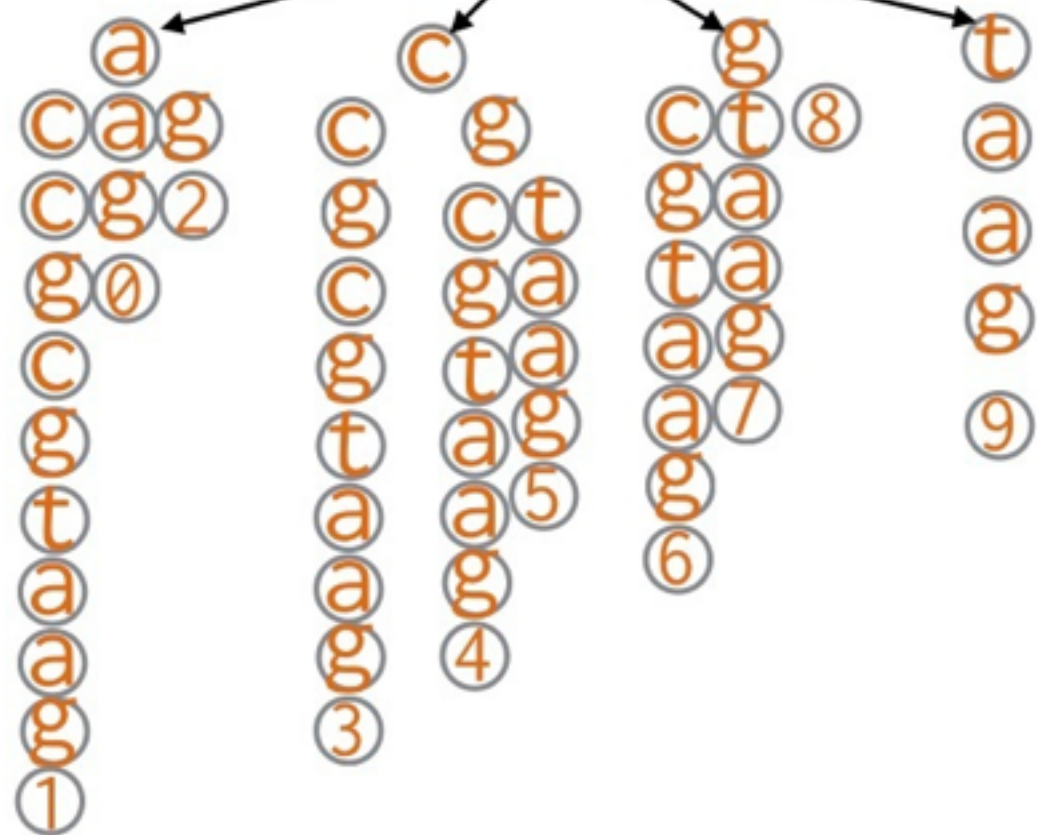
\$

Step 3 : build trie



Burrows-Wheeler

Reference : ...ccgcgtaagaccgcgtaagaccgcgtaaga...



Burrows-Wheeler (actual)

Reference : ...cgtaagaccgcgtaaga...



Map-Reduce in Genetics

Regions of Analysis

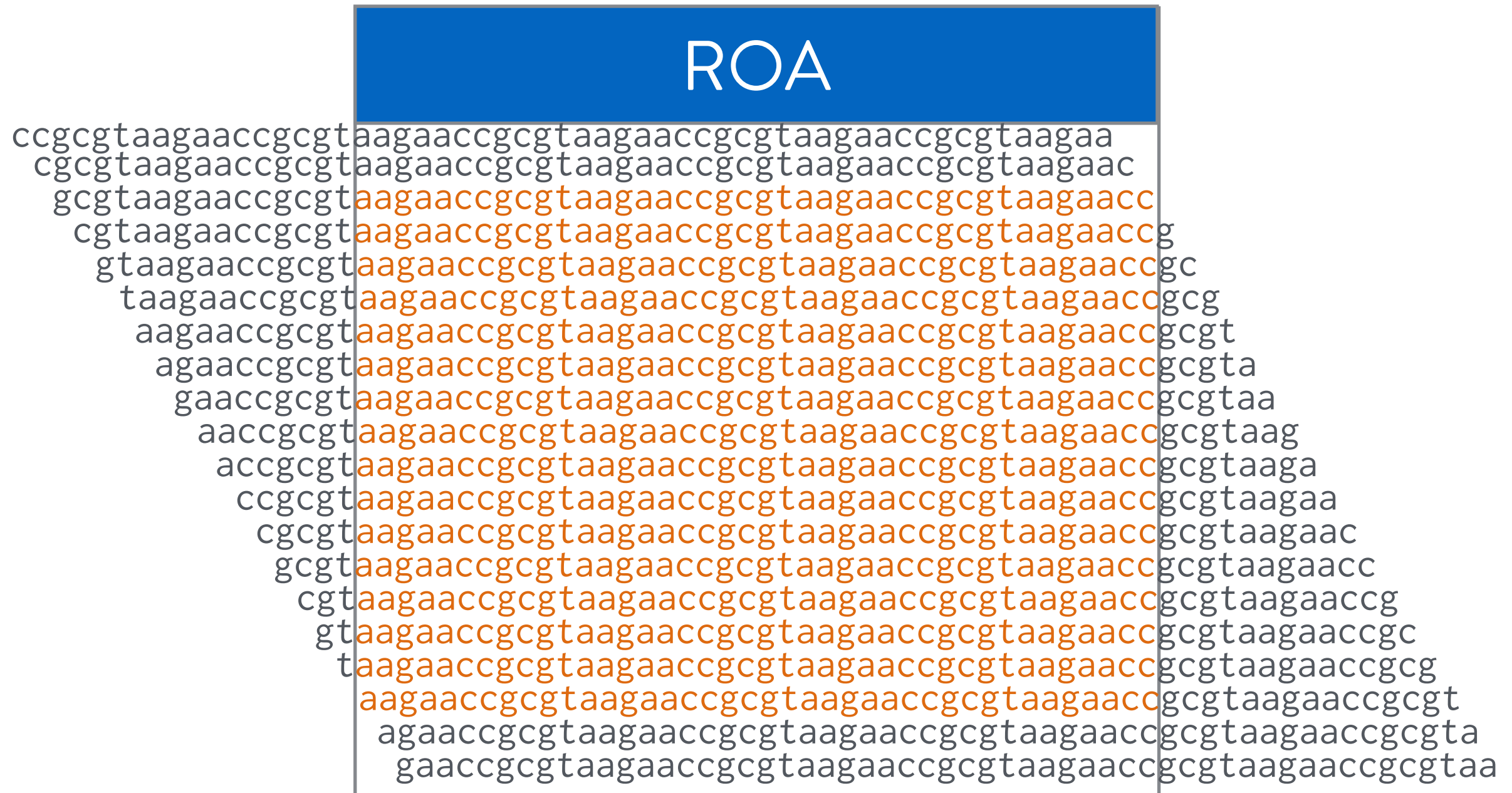
Map Reduce

accgcgtaagaccgcgtaag
accgcgtaagaccgcgtaag
accgcgtaagaccgcgtaag
accgcgtaagaccgcgtaag
accgcgtaagaccgcgtaag
accgcgtaagaccgcgtaag

tggcgcatcttctggcgcatct
tggcgcatcttctggcgcatct
tggcgcatcttctggcgcatct
tggcgcatcttctggcgcatct
tggcgcatcttctggcgcatct

ROA Partitioning

Ref : ...taagaaccgcgtaagaaccgcgtaagaaccgcgtaagaaccgcgtaagaaccgc...



Length : 34 BP

Dictionary Formation

Ref : ...taagaaccgcgtaagaaccgcgtaagaaccgcgtaagaaccgcgtaagaaccgc...

ROA
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc
aagaaccgcgtaagaaccgcgtaagaaccgcgtaagaacc

Length : 34 BP

Dictionary Formation

ROA
acgtggttacctgtacgtttgggaccaatgca
acgtggttac t gtacgtttgggaccaatgca
acgtggttacctgtat t gtttgggaccaatgca
ac a tggttacctgtacgtttgggaccaatgca
ac a tggttacctgtat t gtttgggacacgtggt
acgtggttacctgtacgtttgggac t aatgca
acgtggttacctgtacgtttgggaccaatgca
acgtggttac t gtacgtttgggaccaatgca
acgtggttacctgtat t gtttgggac t aatgca

REF:	8657
Word 1	1033
Word 2	427
Word 3	98
Word 4	2

Dictionary Thresholding

ROA
acgtggttacctgtacgtttgggaccaatgca
acgtggttacttggtacgtttgggaccaatgca
acgtggttacctgtatgtttgggaccaatgca
acatggttacctgtacgtttgggaccaatgca
acatggttacctgtatgtttgggacacgtggt
acgtggttacctgtacgtttgggaccaatgca
acgtggttacttggtacgtttgggaccaatgca
acgtggttacctgtatgtttgggactaatgca

REF:	8657
Word 1	1033
Word 2	427
Word 3	98

Baseline Reconstruction

ROA 18-51		ROA 52-85		ROA 85-118	
		ROA 35-68		ROA 69-92	
acgt	acgtggttacctg	tacgttttgggac	caatgca		
cgt	acgtggttacctg	tacgttttgggac	caatgcat		
gt	acgtggttacctg	tacgttttgggac	caatgcatt		
t	acgtggttacctg	tacgttttgggac	caatgcattg		
	gtggttacctg	tacgttttgggac	caatgcattgcaa		
	tggttacctg	tacgttttgggac	caatgcattgcaact		
	gttacctg	tacgttttgggac	caatgcattgcaactga		
	tacctg	tacgttttgggac	caatgcattgcaactgac		
		acgttttgggac	caatgcattgcaactgatccctgat	cg	
		ttgggac	caatgcattgcaactgatccctgat	cg	t
		gggac	caatgcattgcaactgatccctgat	cg	ta
		gac	caatgcattgcaactgatccctgat	cg	tag

Dictionary Reconstruction

Initial dictionary

Reference Sequence

$t = 1$

Word (68)	Reference (50)	Word (68)	Reference (50)	Word (68)
-----------	----------------	-----------	----------------	-----------

...

	Word (68)	Reference (50)	Word (68)	Reference (50)	
--	-----------	----------------	-----------	----------------	--

Reference (50)	Word (68)	Reference (50)	Word (68)	Reference (50)
----------------	-----------	----------------	-----------	----------------

...

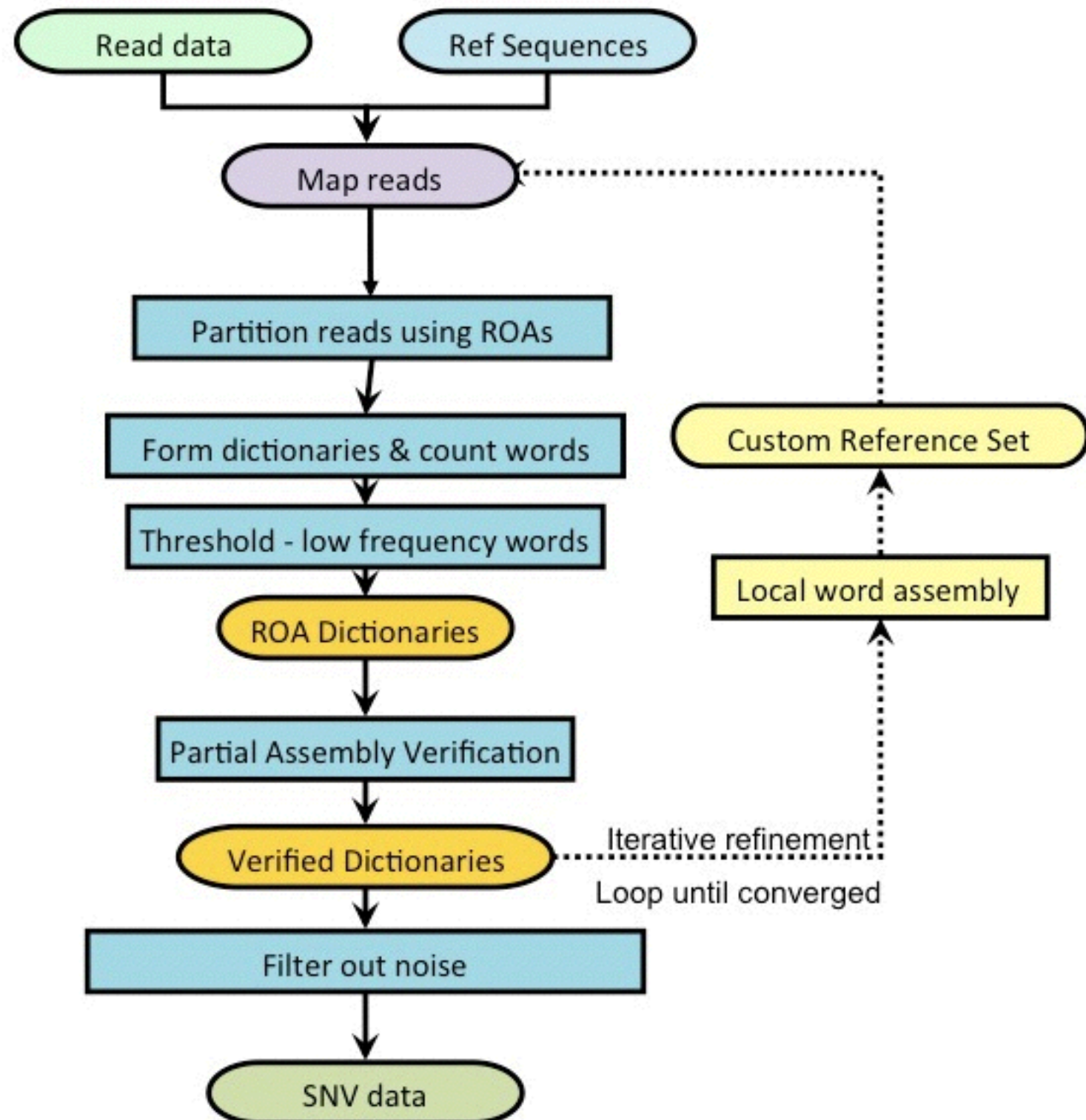
$t = n$

	Word (68)		Word (68)		Word (68)
--	-----------	--	-----------	--	-----------

	Word (68)		Word (68)	
--	-----------	--	-----------	--

		Word (68)		Word (68)	
--	--	-----------	--	-----------	--

Repeat



Distributed Systems in Genetics

Distributed Systems

