

Session 2– Lecture 1

Alignment to reference genomes

Qi Sun

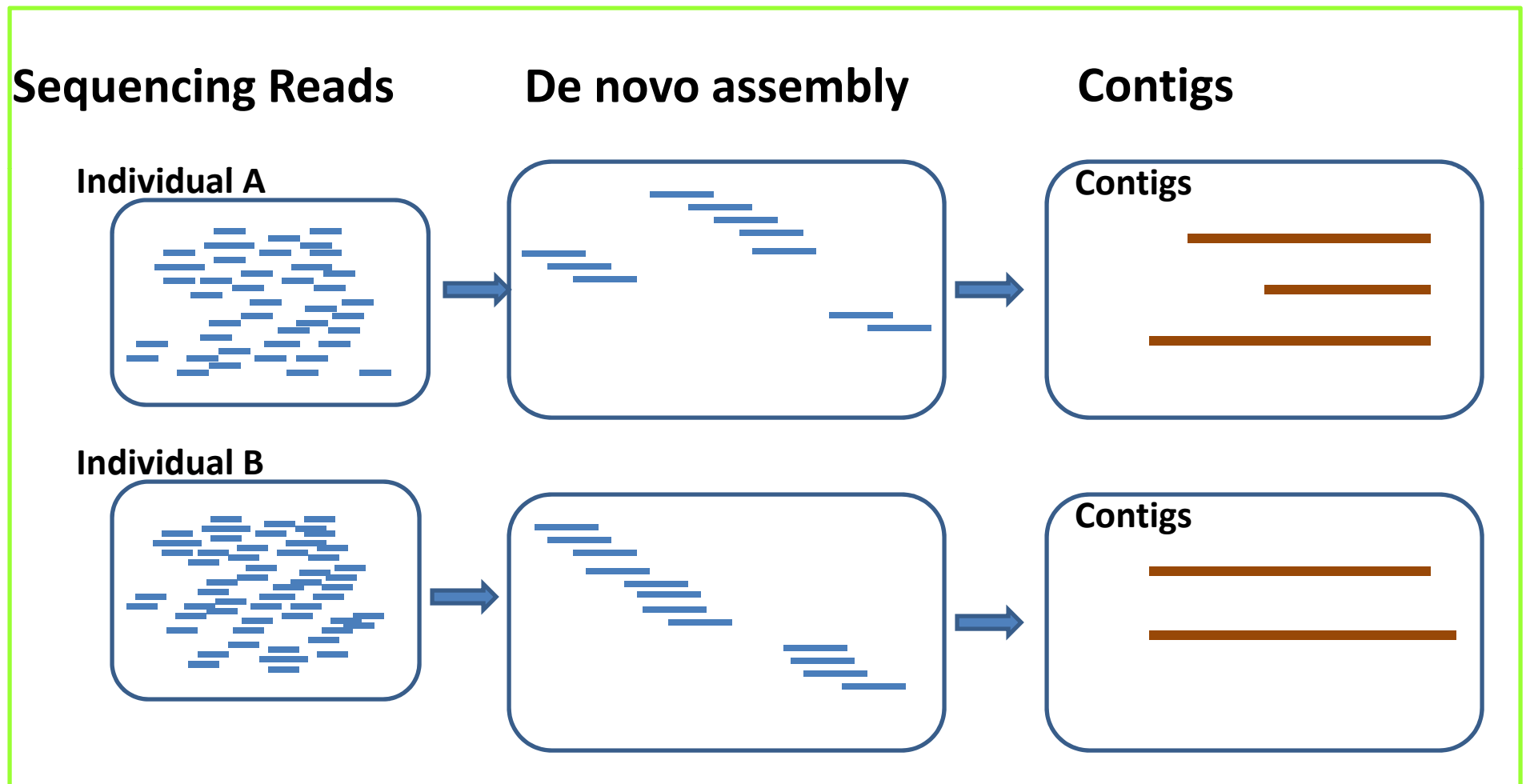
**Computational Biology Service Unit
Cornell University**

Outline

- 1. Alignment and assembly**
- 2. How alignment software work**
- 3. Commonly used alignment software**
- 4. Standard output file format**

Difference between alignment and assembly

Assembly process

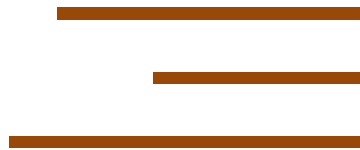


Difference between alignment and assembly

Assembly process

Individual A

Contigs



Individual B

Contigs



Align the contigs

Contig 3 from Individual A



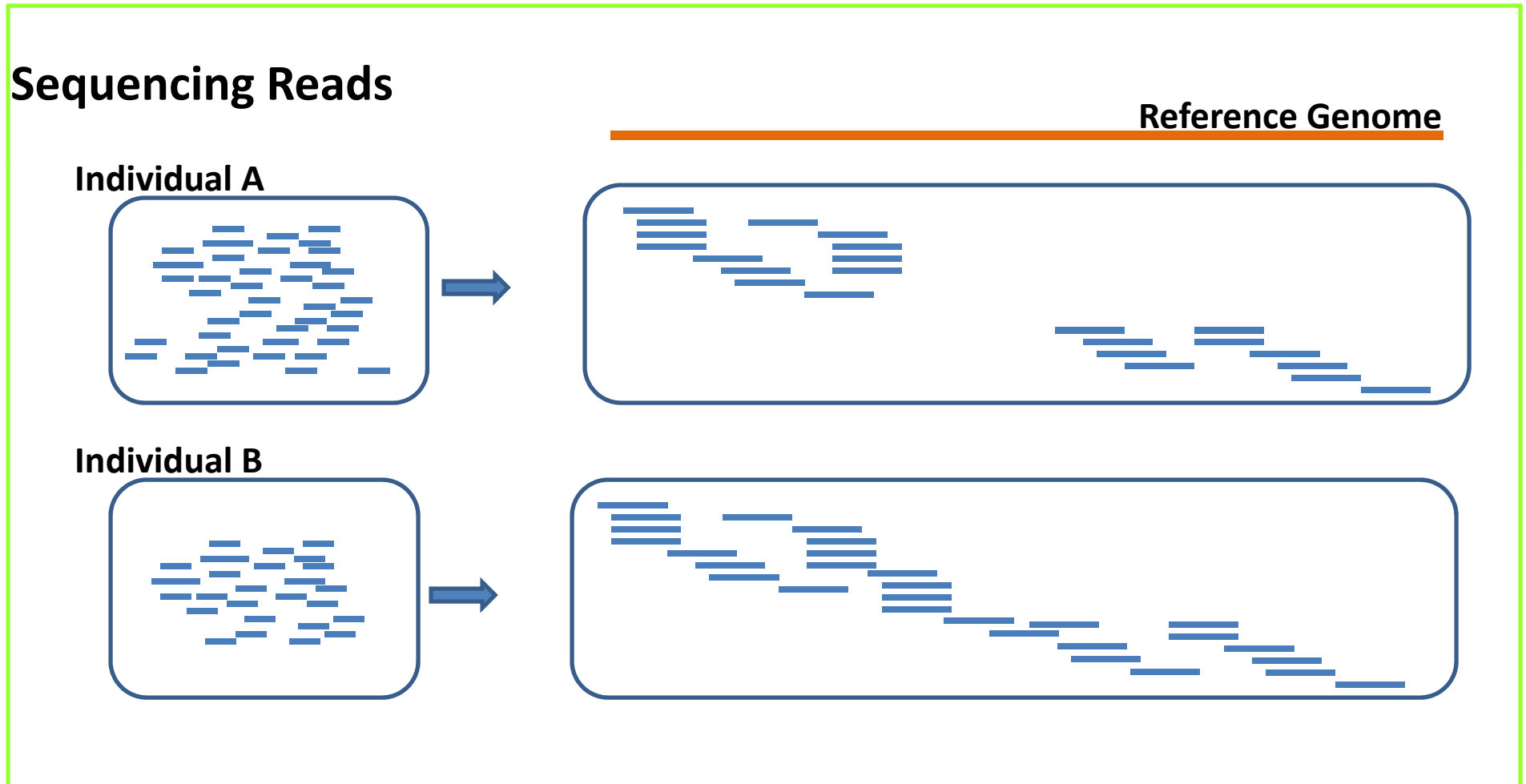
Contig 1 from Individual B



*** Most often, assemblies are followed by gene annotation, and comparison in the annotated genes.**

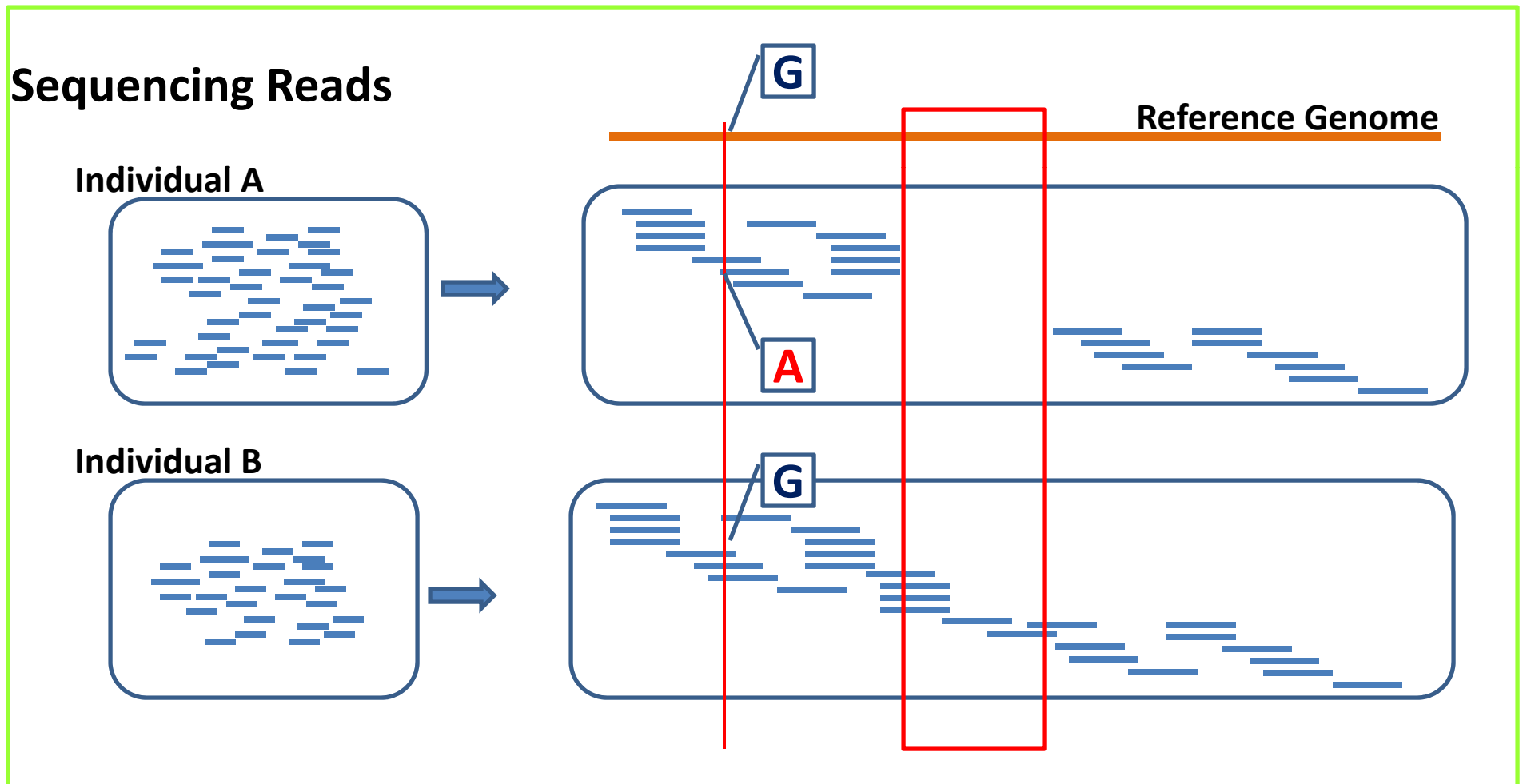
Difference between alignment and assembly

Alignment process



Difference between alignment and assembly

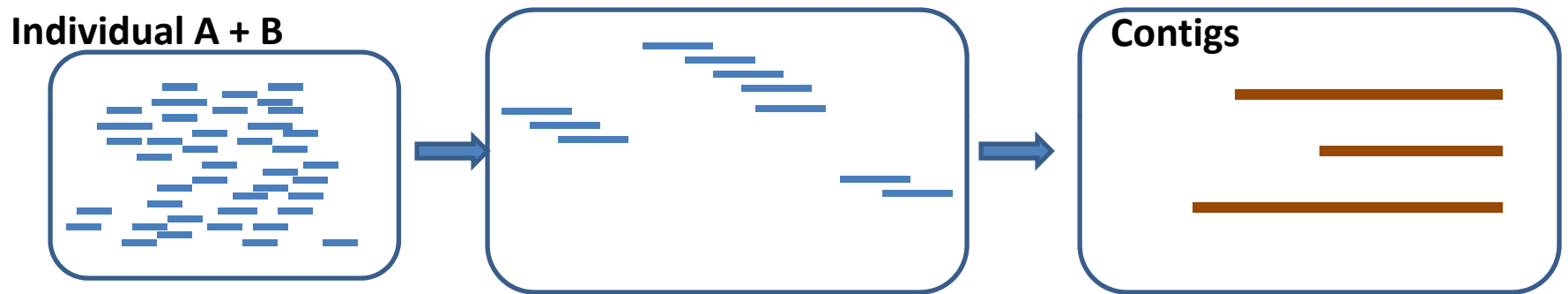
Alignment process



Combining the two approaches

Step1 : Assembly

Mix the reads from sample A and B, creating an assembly



Step2 : Alignment

Use the assembled contigs as reference, and align the reads from A and B to the reference.

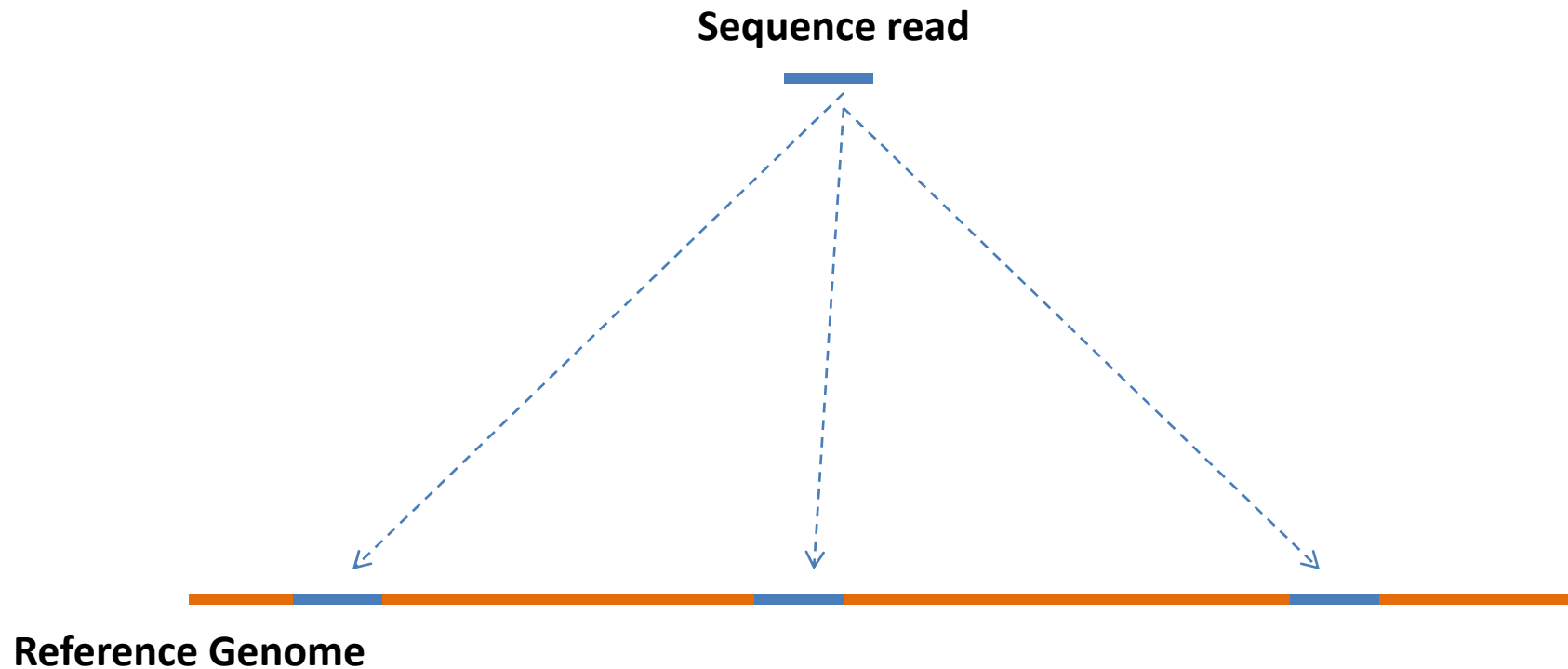
Limitation of alignment approach

Restricted by reference genome. Highly polymorphic regions or large insertions cannot be detected.

Alignment based approach is rarely used in comparative bacterial genome projects.

2-Step Alignment Strategy

Step 1: Mapping. Goal: Quickly identify candidates of hits



2-Step Alignment Strategy

Step 1: Mapping. Goal: Quickly identify candidates of hits

Sequence read

1. Hash table based;

2. Burrows Wheeler transform (BWT)-based;

Reference Genome

Issues in the mapping step

1. The mapping step is heuristic. Not all reads that should have been aligned would be aligned, especially for reads in highly polymorphic regions, repetitive regions, or reads with indels.
2. Balance between alignment speed and alignment accuracy.

2-Step Alignment Strategy

Step 2: Alignment and reporting. Goal: Score the alignment

Read: AGG**t**CCGGG**A**TACCGGGGAC

Candidate 1 (chr1): AGG**G**CCGGGA**A**ACCGGGGAC Score: -2

Candidate 2 (chr2): AGG**G**CCGGG**A**TACCGGGGAC Score: -1

Issues in the alignment/reporting step

- 1. Some software would use the base quality score to evaluate the alignment. Others do not.**
- 2. Software parameters relevant to this step: 1) Maximum mismatches that would be reported; 2) Reporting unique hits or multiple hits.**

Features in alignment software

1. Gapped vs ungapped alignment:

| | |
|-----------|-------------------------------------|
| Reference | GATGGACCCTTA--GTACGC..... . |
| Read | GATGGACCCTTACGGTACGC |

Issues with gapped alignments

Issues with gapped alignments

- Indels located at the edge of a read

Reference GATGGACCCTTA--GTACGC..... .

Read 1 GATGGACCCTTACGGTACGC

Read 2 gccaatGATGGACCCTTAC

Reported as
mismatch

- Reported positions of the INDELS

..... . GATGGACCCTT**AAA**GTACGC..... .

GATGGACCCTT--**AA**CGGTACGC

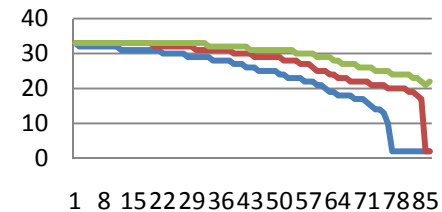
Features in alignment software

1. Global vs Local alignment

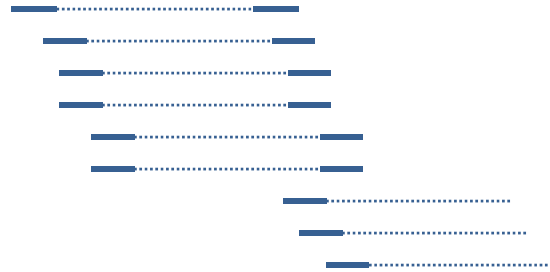
| | |
|-----------|-----------------------------------|
| Reference | ...CTACTGATGGACCCTTACGGTTGAG..... |
| Read | GGGGGATGGACCCTTACGGTACGC |

Situations where only part of the reads can be aligned

- **Low quality part of the read**
- **Reads span breaking points in a chromosome re-arrangement event**
- **Reads span splicing junctions**



Paired – end alignment



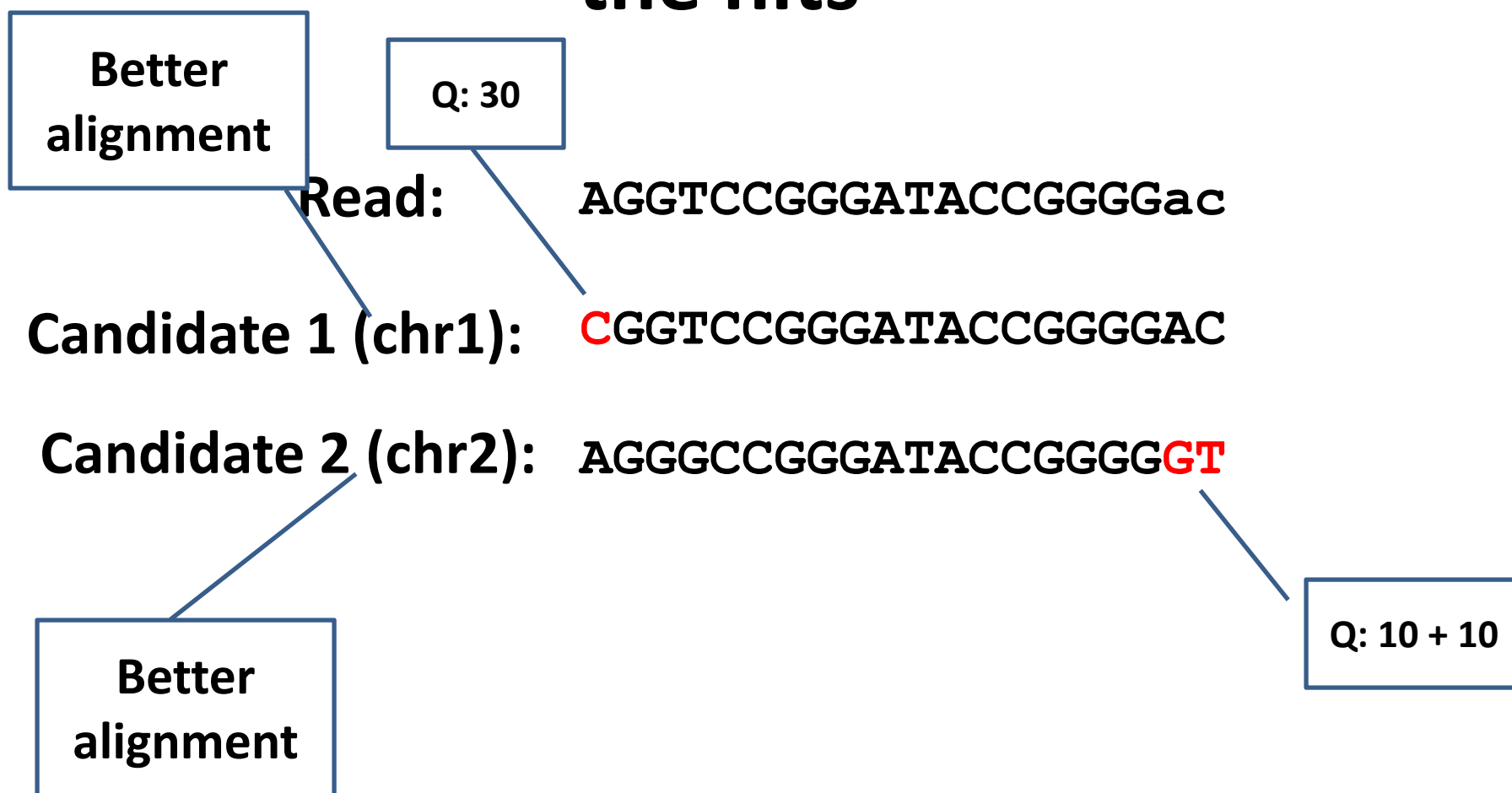
Software supporting paired-end alignment requires user input:

- 1. Minimum and maximum size of the inserts**
- 2. Options of reporting unpaired alignment**

Options of reporting the alignments

- **Unique alignment only**
- **Reporting ambiguous hits**
 - a) **Randomly report one**
 - b) **Report all alignment above cutoff parameters**
 - c) **Report all alignments with best alignment scores**

Using base qualities to evaluate the hits



Hash size selection/ Guaranteed number of mismatches

- **Hash size too big: miss the reads with many mismatches.**
- **Hash size too small: slow; too many hits in the mapping step and fail to be aligned**

BWA

- BWA aln:

- 1) Short reads up to 200 bp with errors <5%
- 2) global alignment;
- 3) Gapped alignment;
- 4) Base quality is not used in evaluating hits;
- 5) Can do paired-end;
- 6) Report ambiguous hits;

- BWASW

- 1) Longer reads with more mismatches;
- 2) Local alignment;
- 3) Slow and less accurate, does not work with paired -end

Bowtie

- **One of the fastest alignment software for short reads**
- **Not gapped-alignment;**
- **Base quality can be used evaluating alignment;**
- **Paired-end;**
- **Flexible reporting mode**

Other software:

- **ELAND** **part of the Illumina CASAVA software**
- **SOAP** **developed at BGI**
- **MAQ** **developed by Heng Li**

Specialized alignment software

- **Tophat: splicing junction alignment**
- **BSMAP: C->T tolerant alignment**

Commercial Alignment Software

Common Features:

Gapped alignment tools;
Reporting ambiguous hits;
Supporting paired-end alignment;
flexible reporting modes;
Not open source.

• SlimSearch by RealTimeGenomics

• Available in April, 2010. Cost: Unknown. CBSU is in its early access program.

• Novoalign by Novocraft

• Available now. Free for academic users.

SAM/BAM

at

Alignment section

| | | | |
|----|-------|--|----------------|
| 1 | QNAME | Query (pa | |
| 2 | FLAG | bitwise FL | |
| 3 | RNAME | Reference | |
| 4 | POS | 1-based l | ipped sequence |
| 5 | MAPQ | MAPping | |
| 6 | CIAGR | extended CIGAR string | |
| 7 | MRNM | Mate Ref | e as RNAME) |
| 8 | MPOS | 1-based M | |
| 9 | ISIZE | Inferred i | |
| 10 | SEQ | query SEQUENCE on the same strand as the reference | |
| 11 | QUAL | query QU | se quality) |
| 12 | OPT | variable C | VTYP:VALUE |

Strand;
Paired-end;
et al.

Map position

Indels; Junctions;
et al

Read sequence &
base qualities

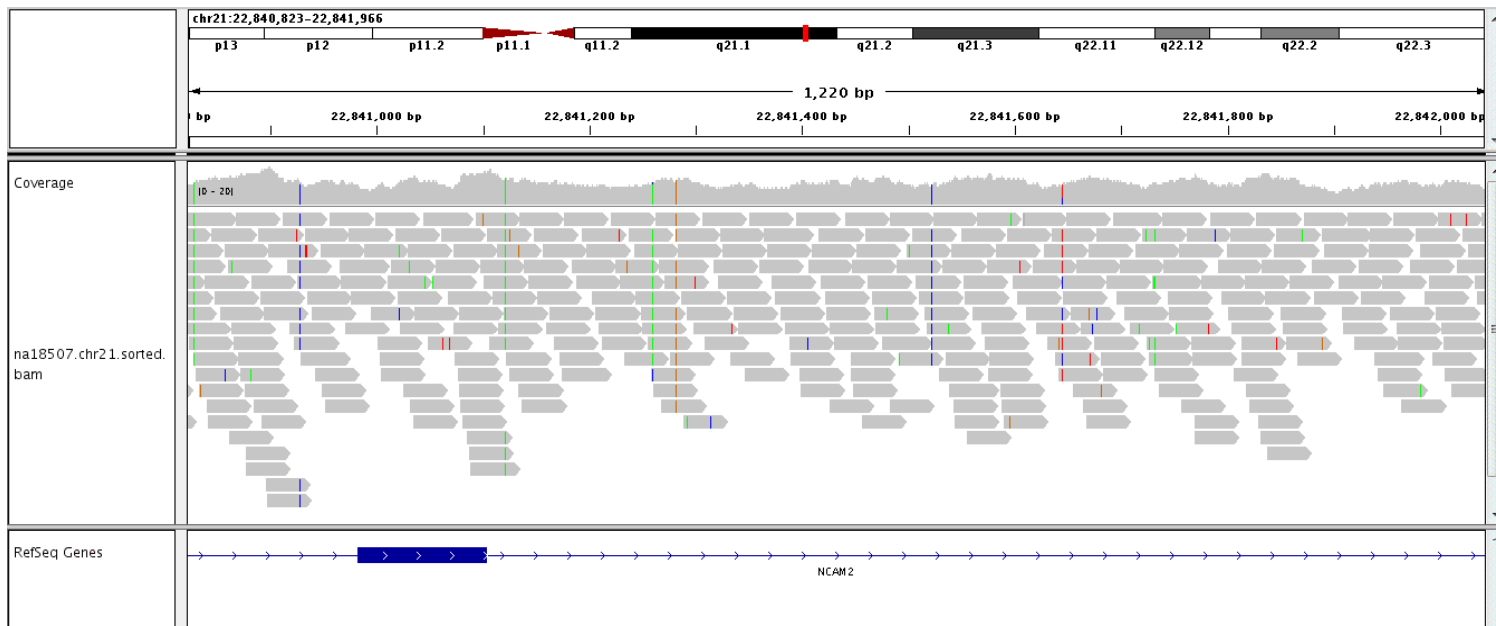
Departures from the standard

How to represent reads with multiple hits in SAM file?

BOWTIE/TOPHAT represents multiple hits with multiple lines, one hit per line.

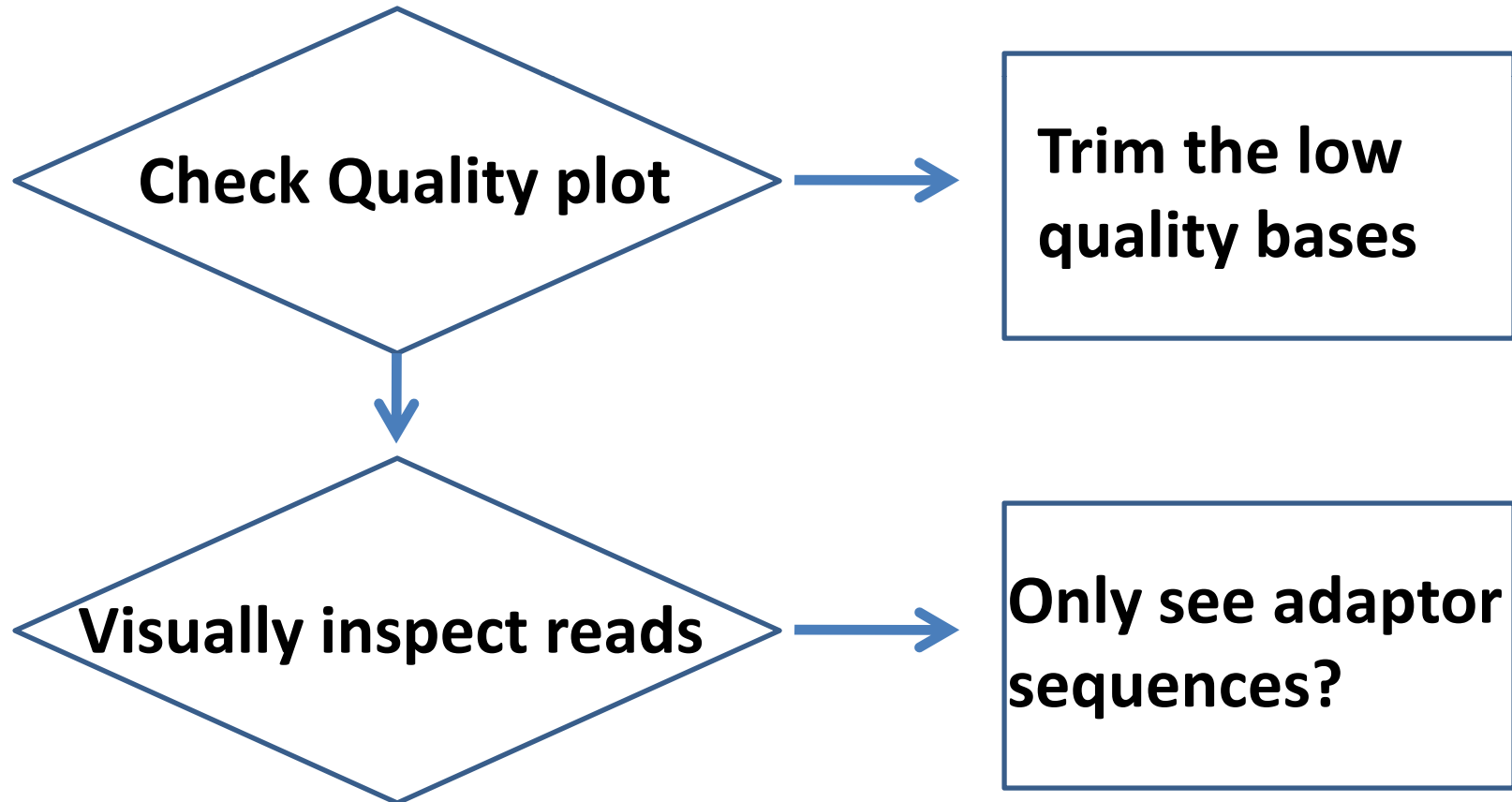
BWA represents multiple hits with single. The alternative alignments are reported with optional XA tag.

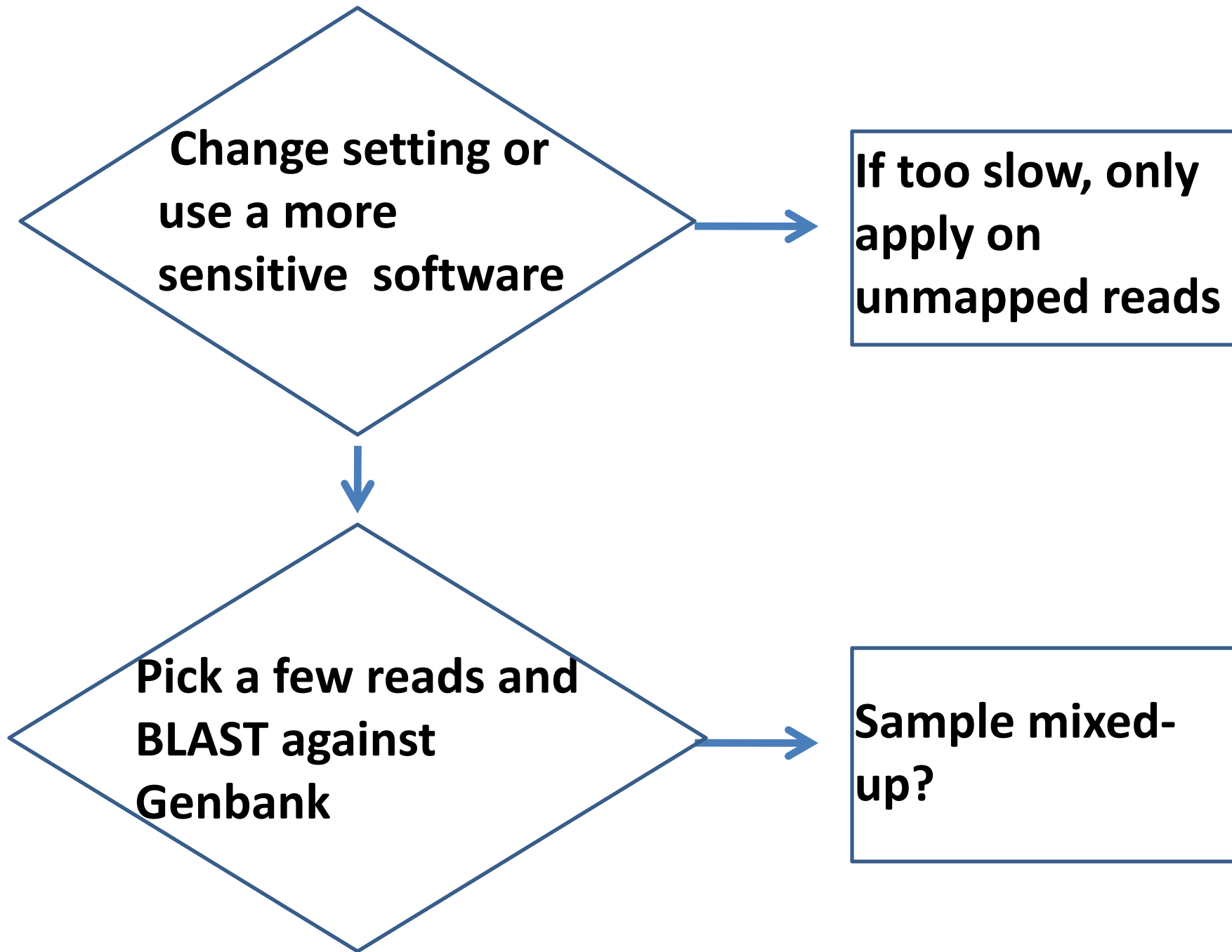
Visualization of BAM file through IGV



Summary

- Use a fast alignment tool in the first try. (BWA, Bowtie, Tophat, et al.). Normally $\frac{3}{4}$ of the reads should align to the genome.





Commercial solutions:

- 1) Alignment software like SlimSearch and Novoalign;**
- 2) Specialized cloud computing service with combined hardware and software solutions, GenomeQuest et al.**