



Central European Institute of Technology
BRNO | CZECH REPUBLIC

Moderní metody analýzy genomu - analýza

Mgr. Nikola Tom

Brno, 2.4.2014



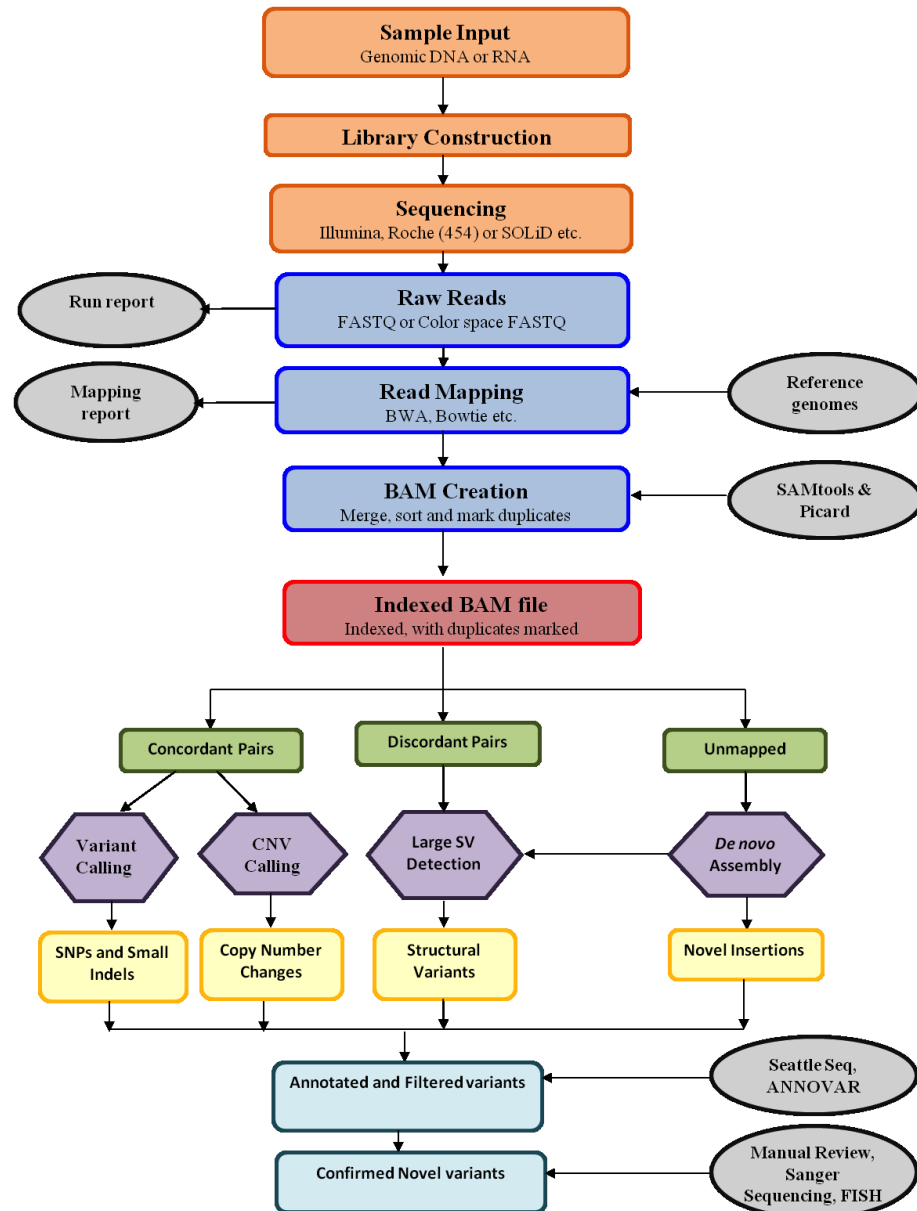
EUROPEAN UNION
EUROPEAN REGIONAL DEVELOPMENT FUND
INVESTING IN YOUR FUTURE



OP Research and
Development for Innovation



Workflow

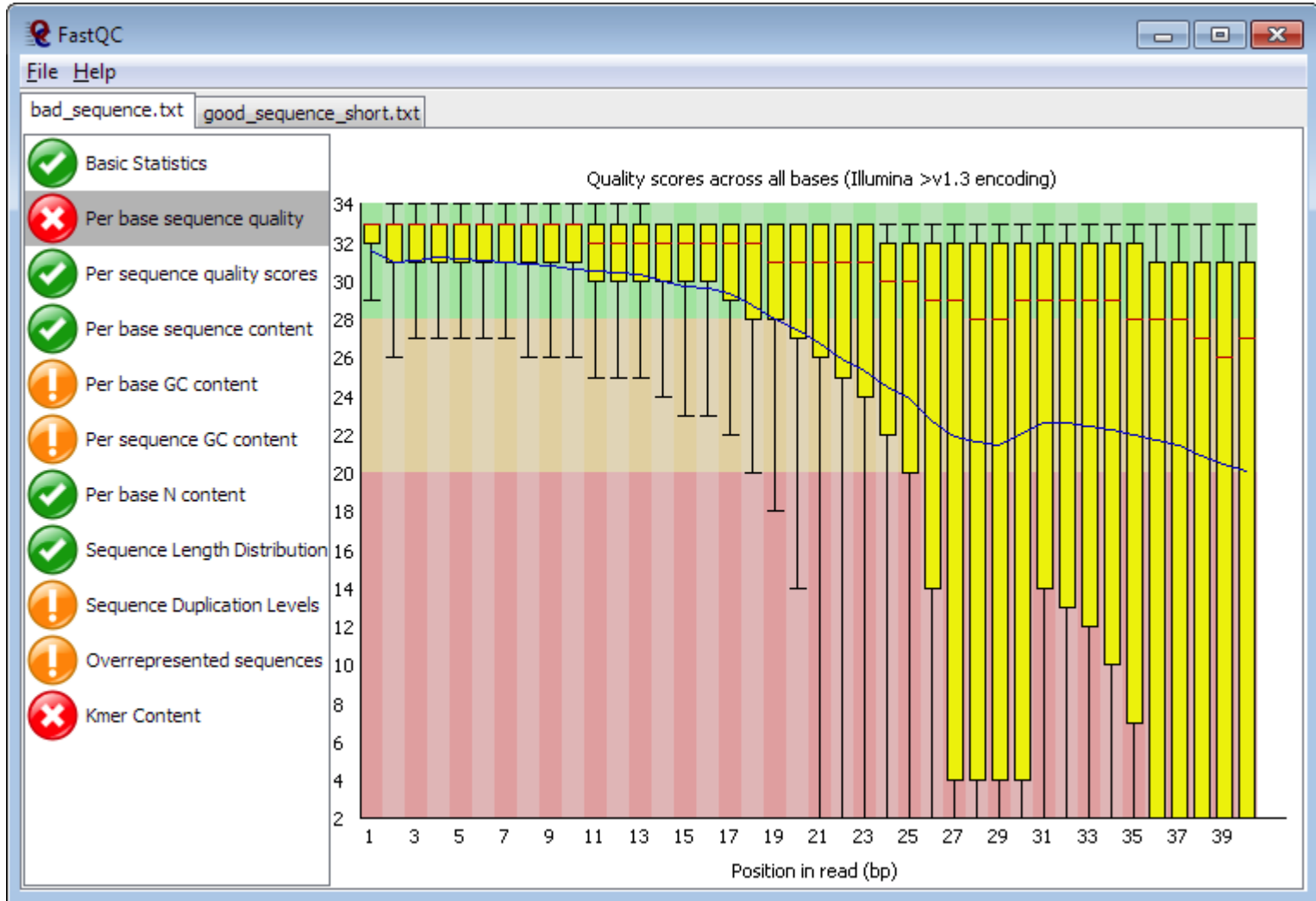


Raw sequence = fastq

- Biological sequence
- Corresponding quality scores
- ASCII character
- (fasta+ qual, csfasta + csqual, sff)

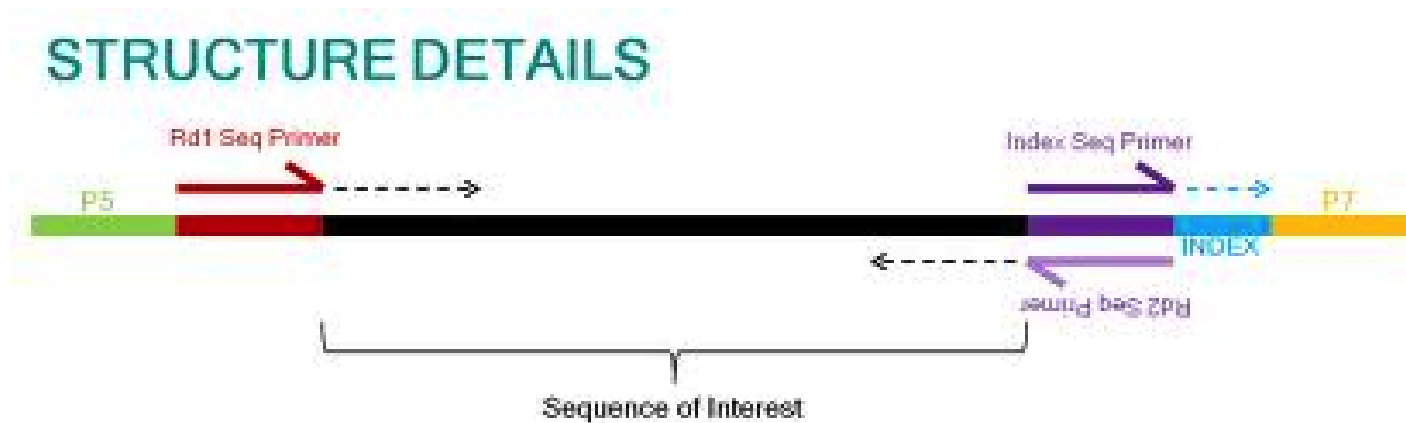
```
@  
SEQ_ID GATTTGGGGTTCAAAGCAGTATCGATCAAATAGTAAATCCATTTGTTCAACTCACAGTTT  
+  
!"*((( (**+))%%%+)(%%%) .1***-+*")**55CCF>>>>>CCCCCCC65
```

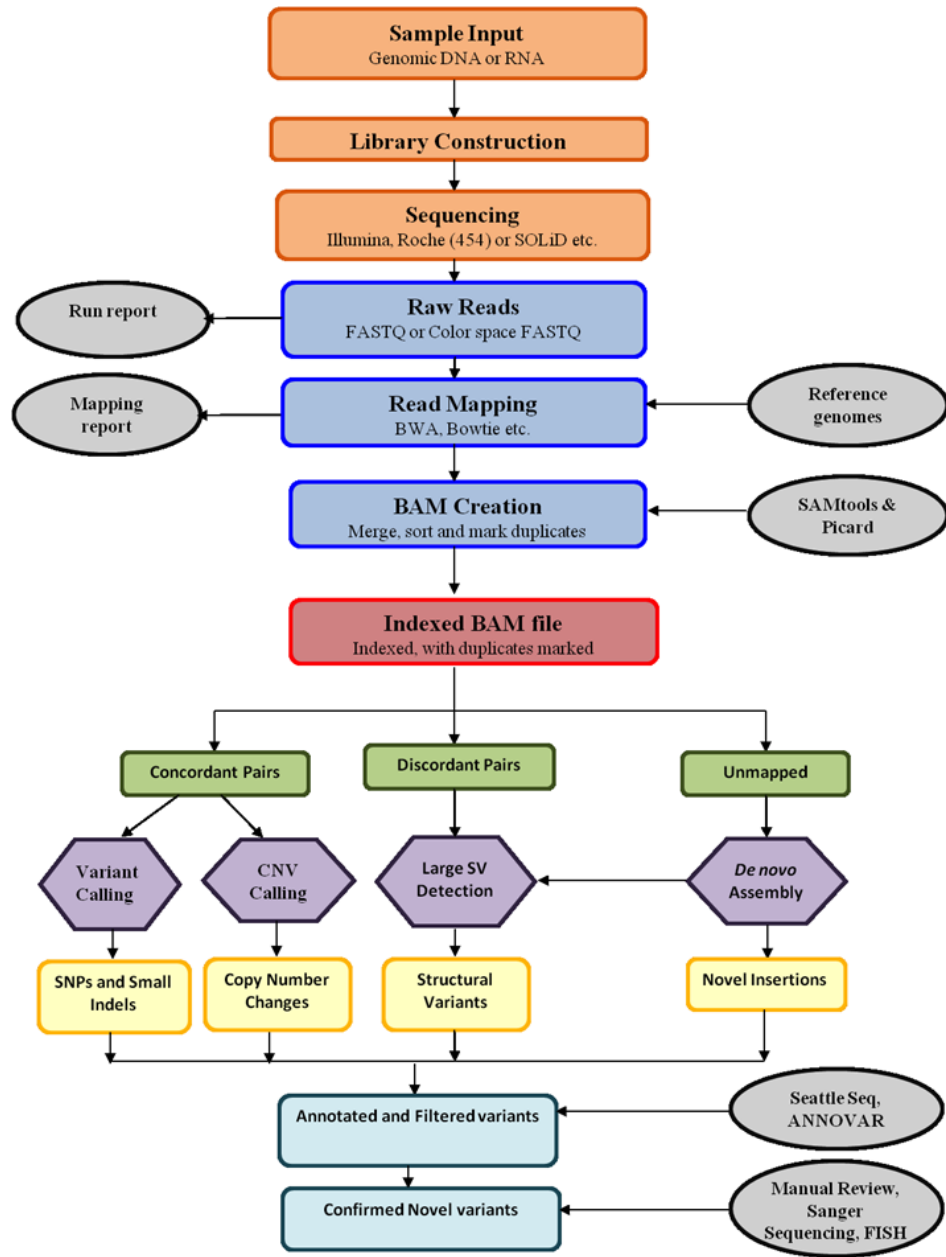
FastQC



Cutadapt

- Adaptor trimming (miRNA)
- Quality filtering
- Length filtering





Read mapping => SAM, BAM

- Usually mapping reads on reference
- miRNA - special case
 - Grouping and annotate against mirBase
- **DNA**
 - BWA, Bowtie, Bfast, SHRiMPclc
- **RNA**
 - TopHat (*de novo* splice aligner)
- **Commercial**
 - CLC Genomics Workbench

SAM

Headers

```
1 @HD VN:1.0 SO:unsorted
2 @SQ SN:gil110640213|refINC_008253.1 LN:4938920
3 @PG ID:bowtie2 PN:bowtie2 VN:2.1.0
```

Alignments

```
4 gil110640213|refINC_008253.1_418_952_1:0:0_1:0:0_0/1 0 gil110640213|refINC_
008253.1 418 42 70M * 0 0 CCAGGCAGTGGCAGGTGGCCACCGTCCTCTCTGCCCCCGCCAAA
ATCACCAACCATCTGGTAGCGATGAT 2222222222222222222222222222222222222222222
22222222222222222222 AS:i:-3 XN:i:0 XM:i:1 XO:i:0 XG:i:0 NM:i:1 MD:Z:8G61
YT:Z:UU
5 gil110640213|refINC_008253.1_31_476_0:0:0_0:0:0_1/1 16 gil110640213|refINC_
008253.1 407 42 70M * 0 0 GGAAAGCAATGCCAGGCAGGGCAGGTGGCCACCGTCCTCTCTG
CCCCCGCCAAAATCACCAACCATCTG 2222222222222222222222222222222222222222222
22222222222222222222 AS:i:0 XN:i:0 XM:i:0 XO:i:0 XG:i:0 NM:i:0 MD:Z:70 YT:Z
-III
6 gil110640213|refINC_008253.1_210_743_2:0:0_1:1:0_2/1 0 gil110640213|refINC_
008253.1 210 42 70M * 0 0 CATTACCACCACCATCACCATTACCACAGGAAACGGTGCGGGCT
GACGCGTACAGGAAACACCGAAAAAA 2222222222222222222222222222222222222222222
22222222222222222222 AS:i:-6 XN:i:0 XM:i:2 XO:i:0 XG:i:0 NM:i:2 MD:Z:30T31A7
YT:Z:UU
```

Each row describes a single alignment of a raw read against the reference genome. Each alignment has 11 mandatory fields, followed by any number of optional fields.

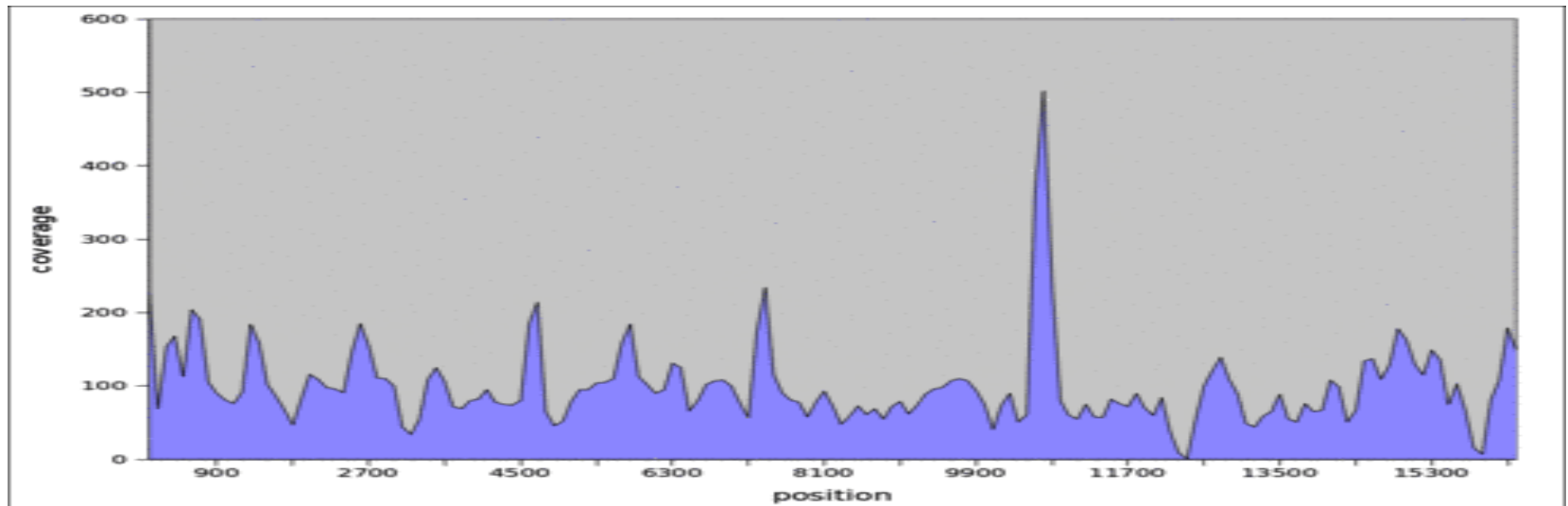
Alignment

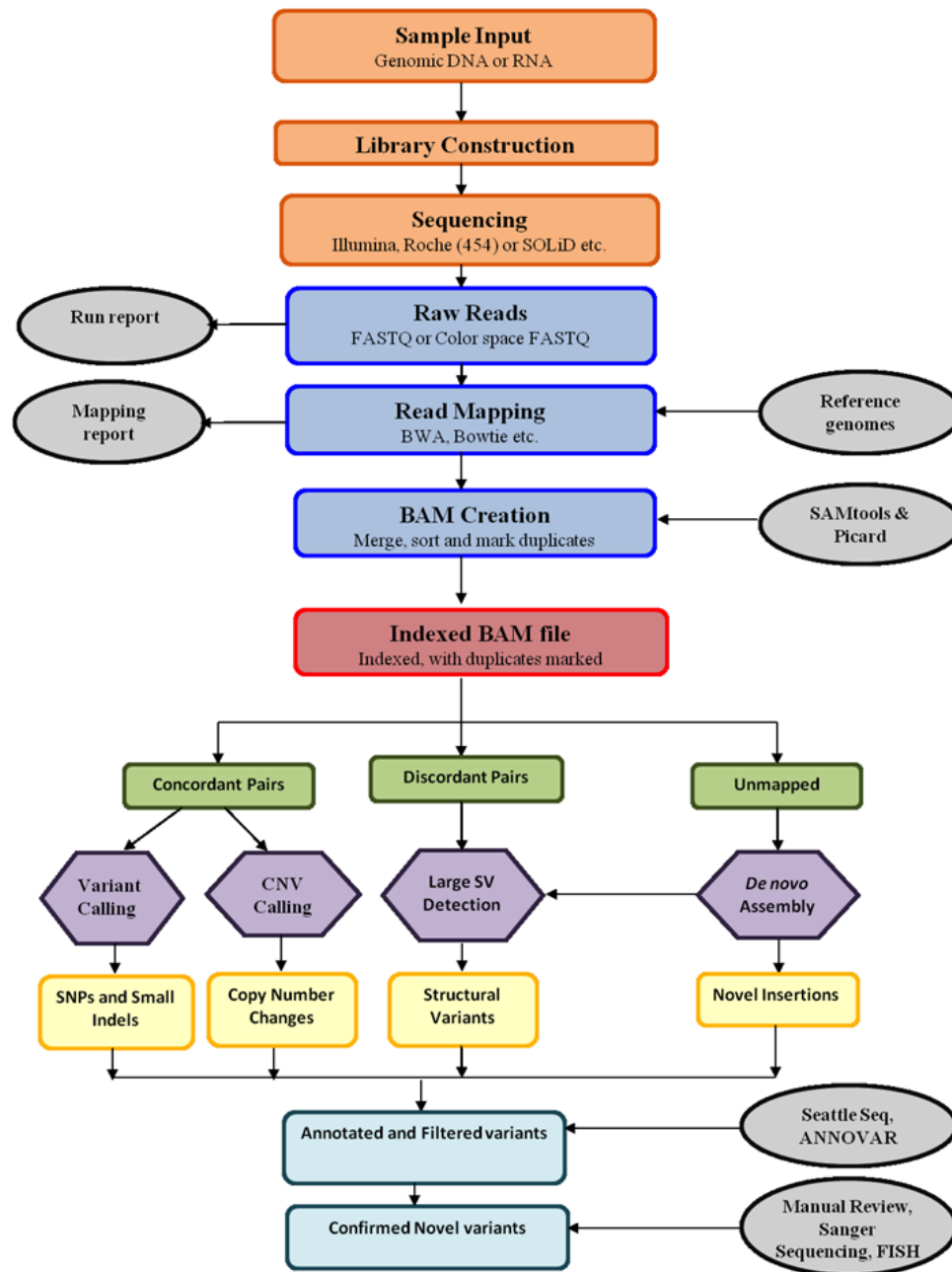


Aligned reads	<pre> ACGCGATTCAGGTTACCAAG GCGATTCAGGTTACCAAGG GATTCAGGTTACCAAGCGTA TTCAGGTTACCAAGCGTAGC CAGGTTACCAAGCGTAGCGC GGTTACCAAGCGTAGCGCAT TTACCAAGCGTAGCGCATT ACCACGCGTAGCGCATTACA CACGCGTAGCGCATTACACA CGCGTAGCGCATTACACAGA CGTAGCGCATTACACAGATT TAGCGCATTACACAGATTAG </pre>
Consensus contig	<pre> ACGCGATTCAGGTTACCAAGCGTAGCGCATTACACAGATTAG </pre>

Mapping, Coverage reports

- Important checkout for lab protocol
- Specificity of PCR
- Normalization
- Settings of variant calling threshold, CNV





SNV and small InDel Calling

- Coverage
- Frequency
- Base quality

- !!!
- Genomic context (homopolymers)
- Nucleotide type
- Position in read (errors at the read end)
- Alignment errors

CNV variations

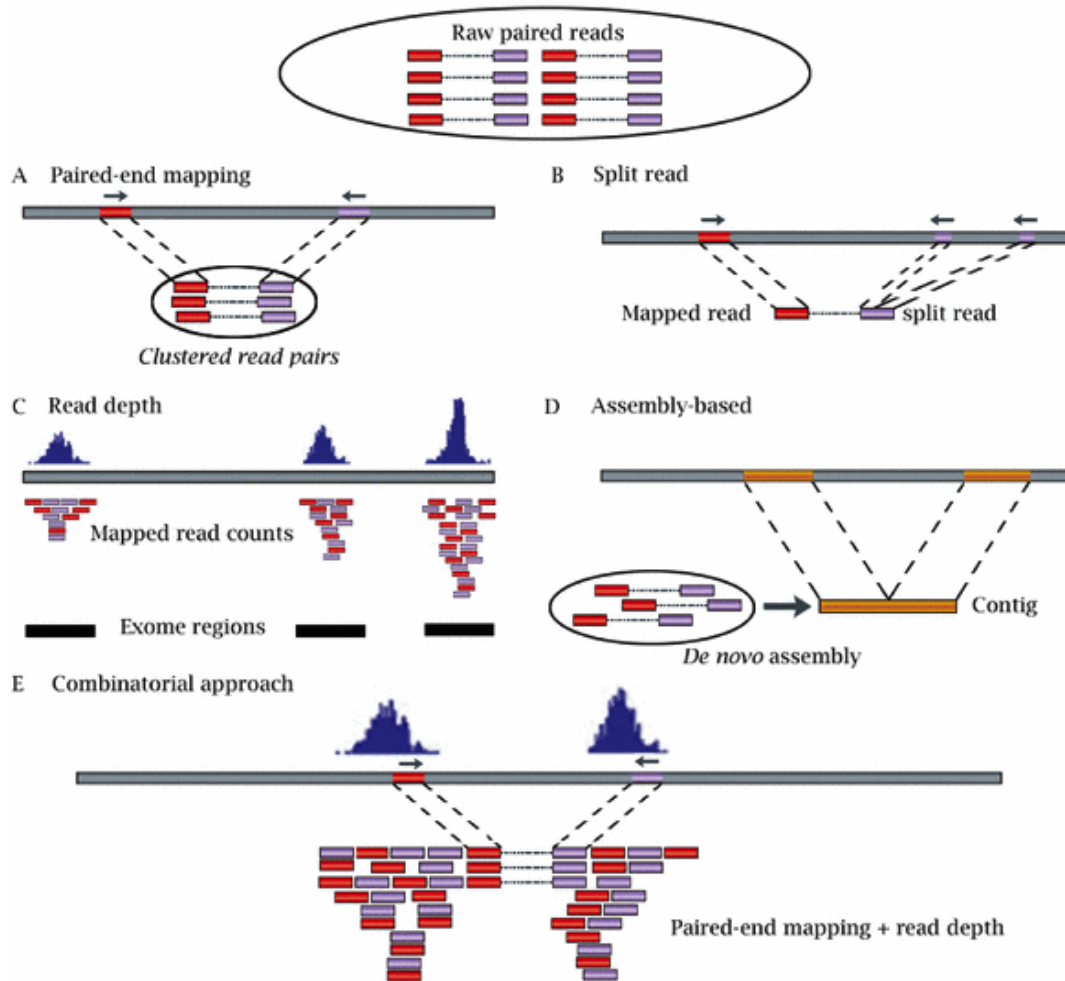
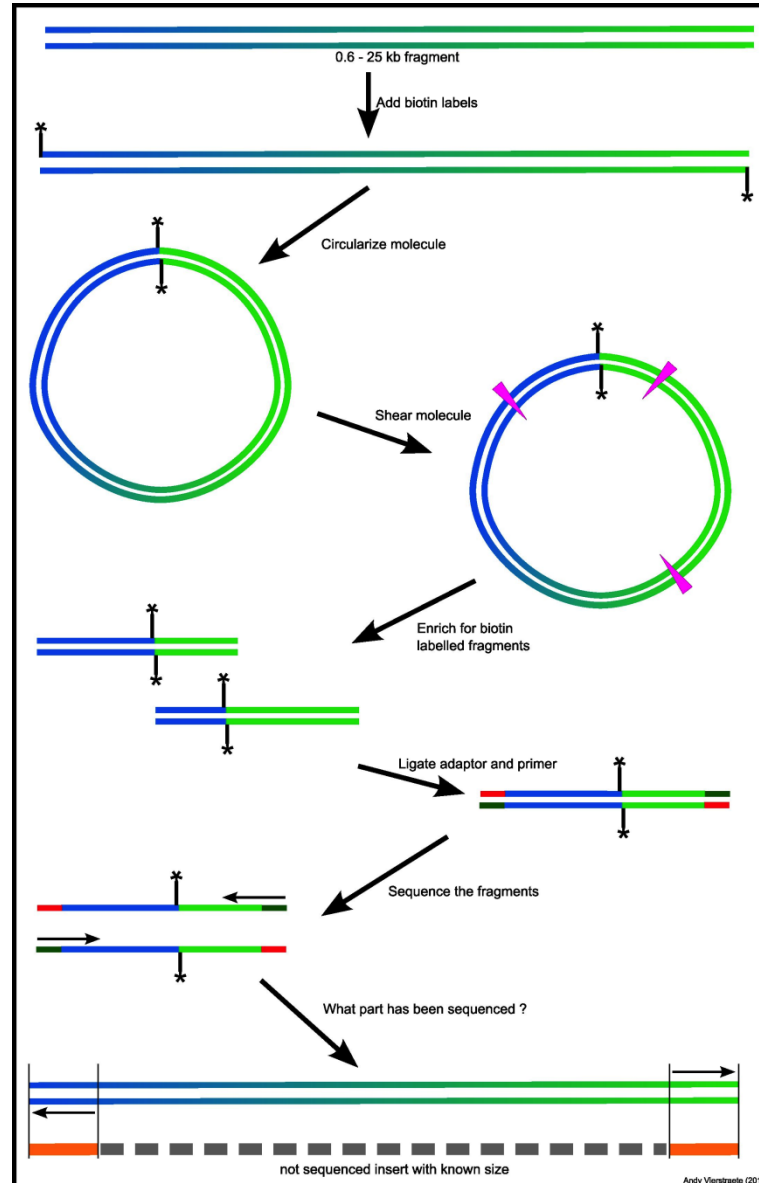
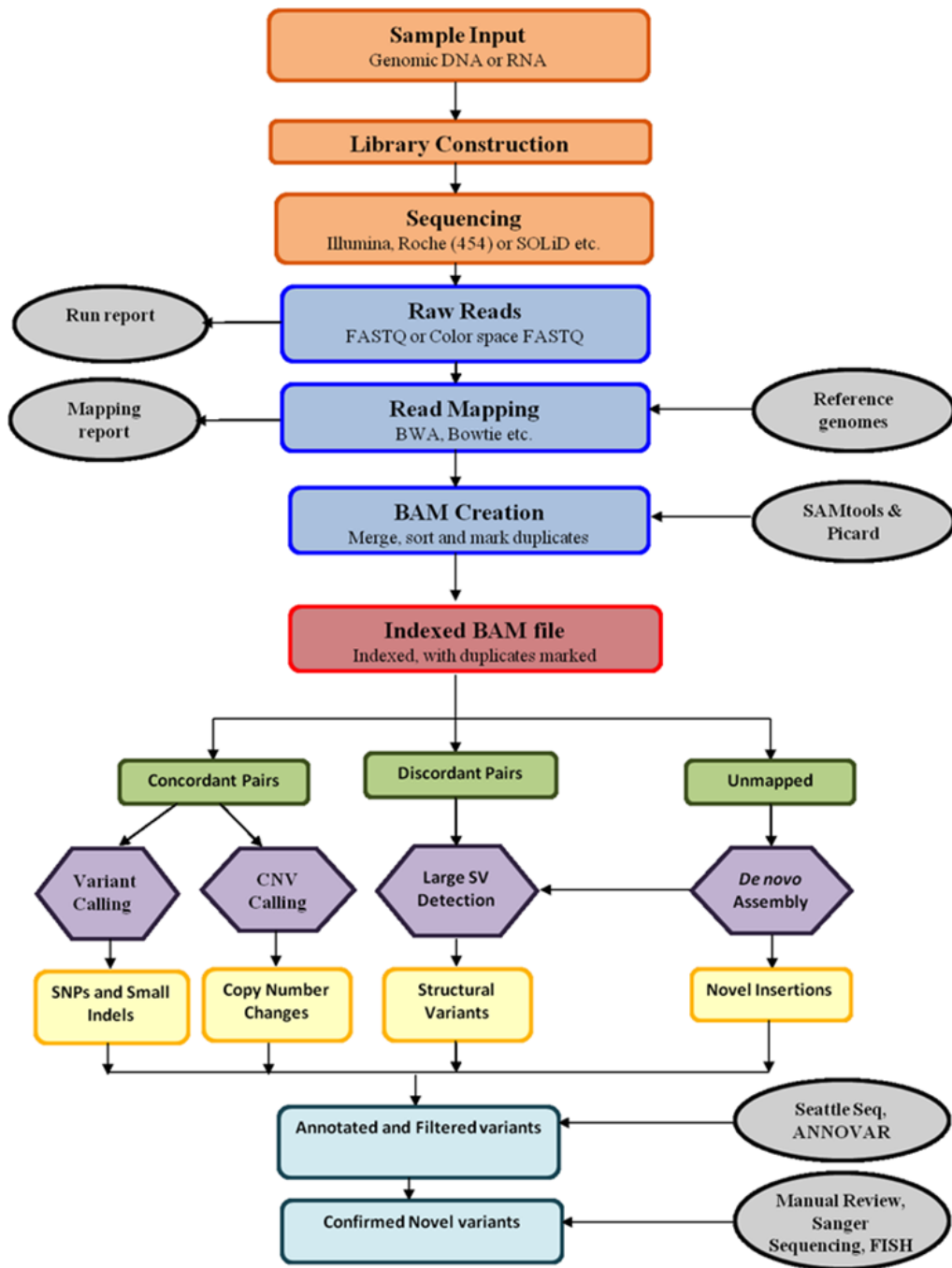


Figure 1 Five approaches to detect CNVs from NGS short reads. A. Paired-end mapping (PEM) strategy detects CNVs through discordantly mapped reads. A discordant mapping is produced if the distance between two ends of a read pair is significantly different from the average insert size. B. Split read (SR)-based methods use incompletely mapped read from each read pair to identify small CNVs. C. Read depth (RD)-based approach detects CNV by counting the number of reads mapped to each genomic region. In the figure, reads are mapped to three exome regions. D. Assembly (AS)-based approach detects CNVs by mapping contigs to the reference genome. E. Combinatorial approach combines RD and PEM information to detect CNVs.

Structural variations

- Mate-pair library
- Long InDel
- Translocation





Annotating and filtering

- Gene
- Transcript
- dbSNP
- Regulation
- Comparative genomics
- Repeats
- Functional
- Gene ontology
- miRNA targets
- Etc.