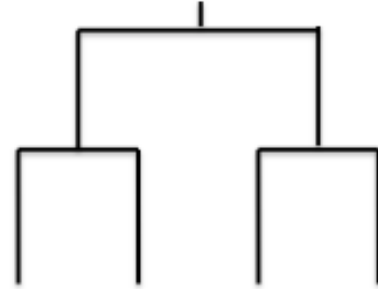# Bacterial Genomics Workshop

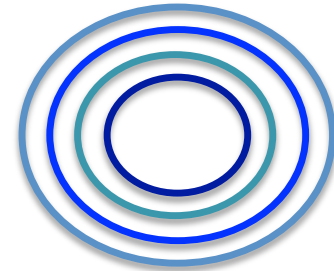## March 2021

# Goals of workshop

- Get an overview of steps in microbial genomics pipeline

- Get exposure to common file formats and terminology in genomics

- Get hands on experience with a set of tools that could compose a genomics pipeline

- Get experience working in a high-performance computing environment
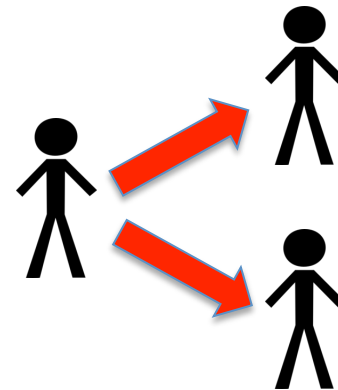
# So you want to sequence some bacteria?

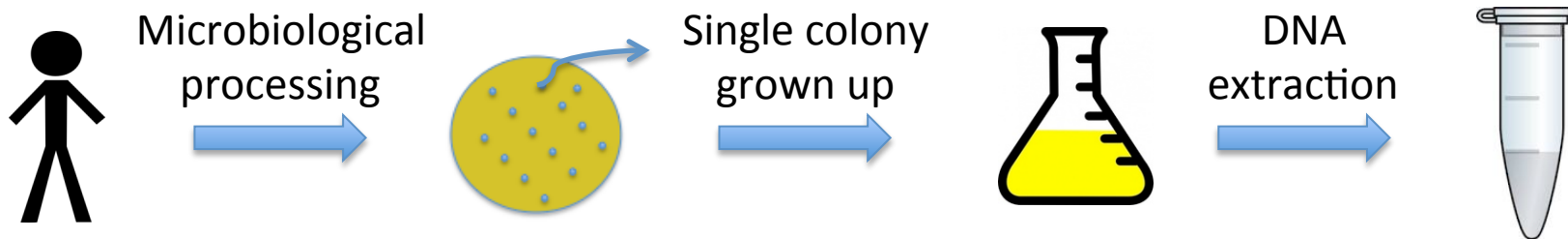- Microbial phylogenetics

- Comparative genomics

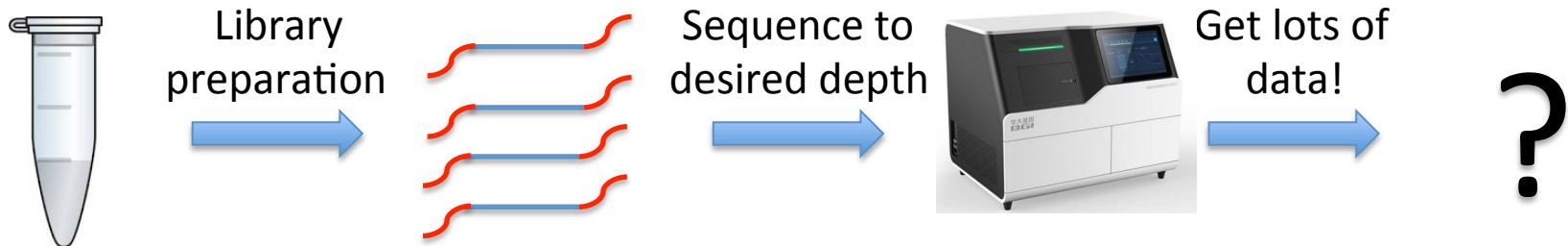- Population genomics

# DNA and library preparation

**1. Sample Preparation**

Microbiological processing → Single colony grown up → DNA extraction →

**2. Sequencing**

Library preparation → Sequence to desired depth → Get lots of data! → **?**

# Mile-high view of a genomics pipeline

Variant
Annotation

I -> L    gyrA
A -> *    mutS

Variant
Calling

Quality
Control

Phylogenetic
Analysis

Genome
Assembly

Comparative
Genomics

Genome
Annotation

# Sequencing quality control

**Forward reads**

```
@seq1_1
ACTGCACT
+
8-8,,+@+
@seq2_1
TGCATCTA
+
@+@E++BF
.
.
.
```

**Reverse reads**

```
@seq1_2
TCTATCGA
+
A<-9BFCFF
@seq2_2
CTAGTTAA
+
**>D7?7=.
.
.
.
```

**Raw fastq files**

**FastQC**

1. Contaminants
2. Aberrant quality

## FastQC Report

**Summary**

✅ Basic Statistics
⚠️ Per base sequence quality
✅ Per tile sequence quality
✅ Per sequence quality scores
❌ Per base sequence content
⚠️ Per sequence GC content
✅ Per base N content
⚠️ Sequence Length Distribution
✅ Sequence Duplication Levels
❌ Overrepresented sequences
✅ Adapter Content
❌ Kmer Content

**Trimmomatic**

1. Filter reads
2. Trim reads

**Forward reads**

```
@seq1_1
ACTGCACT
+
8-8,,+@+
.
.
.
.
.
.
```

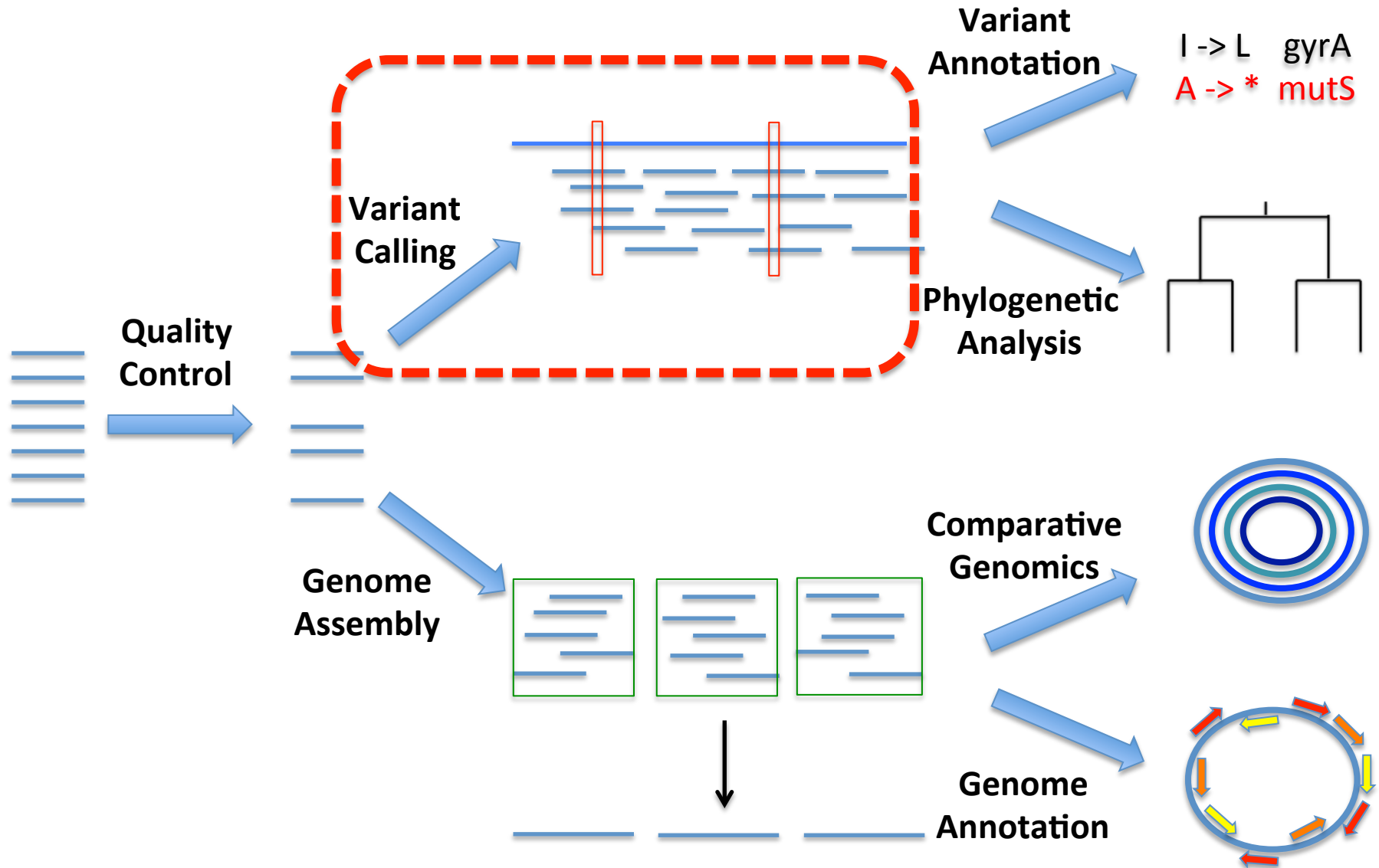**Reverse reads**

```
@seq1_2
TCTATCGA
+
A<-9BFCFF
.
.
.
.
.
.
```

**Clean fastq files**

# Mile-high view of a genomics pipeline

**Variant Annotation**

I -> L    gyrA
A -> *    mutS

**Variant Calling**

**Phylogenetic Analysis**

**Quality Control**

**Genome Assembly**

**Comparative Genomics**
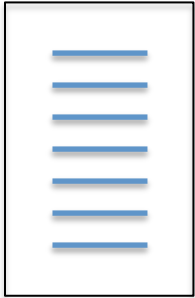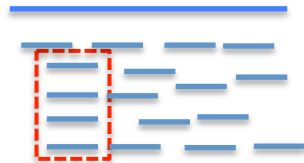
**Genome Annotation**

# Variant identification



**Forward reads**

**Reverse reads**

**bwa**

Read mapping

**Picard**

Remove duplicates

**samtools + bcftools**

Call variants

| Ref | Var |
|-----|-----|
| A | T |
| C | A |
| G | A |
| C | - |

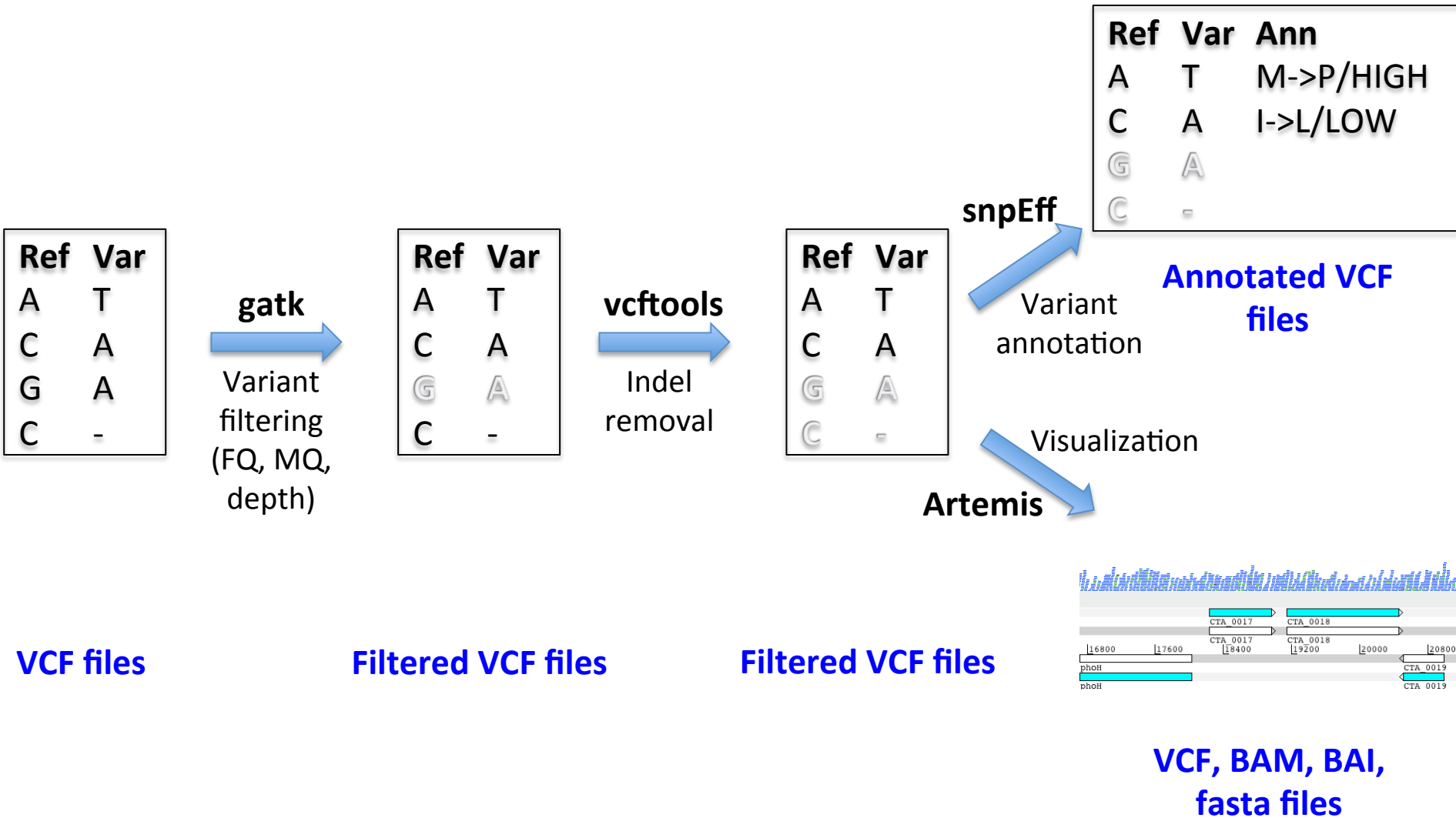**Clean fastq files**        **SAM/BAM files**        **SAM/BAM files**        **Raw VCF files**

# Variant filtering and annotation

| Ref | Var | Ann |
|-----|-----|-----|
| A | T | M->P/HIGH |
| C | A | I->L/LOW |
| G | A | |
| C | - | |

**snpEff**

**Annotated VCF files**

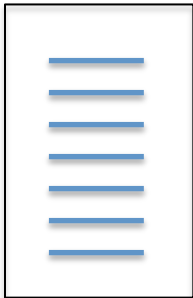| Ref | Var |
|-----|-----|
| A | T |
| C | A |
| G | A |
| C | - |

**gatk**

Variant filtering (FQ, MQ, depth)

| Ref | Var |
|-----|-----|
| A | T |
| C | A |
| G | A |
| C | - |

**vcftools**

Indel removal

| Ref | Var |
|-----|-----|
| A | T |
| C | A |
| G | A |
| C | - |

Variant annotation

Visualization

**Artemis**

**VCF files**

**Filtered VCF files**

**Filtered VCF files**



**VCF, BAM, BAI, fasta files**

**KEY**

- Input files
- Individual genome analysis
- Comparative genome analysis
- Future implementation

*Optional step

Raw DNA sequencing read/s (**.fastq.gz**)

Reference sequence (**.fasta**)

*De novo* assembly (**Velvet**)

Read alignment against reference (**BWA**)

Export of unaligned reads (**.bam**)

Conversion to .bam format (**SAMTools**)

*De novo* assembly (**Velvet**)

Removal of optical duplicates (**Picard**)

Local realignment of high mismatch regions (**GATK**)

Locus presence/absence analysis (**BEDTools**)

Variant (SNP and indel) calling (**GATK**)

Merged presence/absence matrix (**SPANDx, MS Excel**)

Annotated reference* (**.gbk, .gff**)

Variant filtering (**GATK, SPANDx**)

Variant annotation* (**SNPEff**)

Filtered variant outputs (**.vcf**)

Merged SNPs (**GATK**)

Merged indels (**GATK**)

Downstream phylo-genetic analysis (**PAUP\*, PHYLIP, RAxML**)

Reinterrogate ambiguous variants (**GATK**)

Filtered SNP matrix (**.nex**)

Remove non-orthologous and low-quality SNPs (**SPANDx**)

Merge "clean" SNPs (**VCFtools**)

Merge "clean" indels (**VCFtools**)

Annotated reference* (**.gbk, .gff**)

Merge annotated variants* (**SNPEff, SPANDx**)

# Mile-high view of a genomics pipeline

**Variant Annotation**

I -> L    gyrA
A -> *    mutS

**Variant Calling**

**Quality Control**

**Phylogenetic Analysis**

**Genome Assembly**

**Comparative Genomics**

**Genome Annotation**
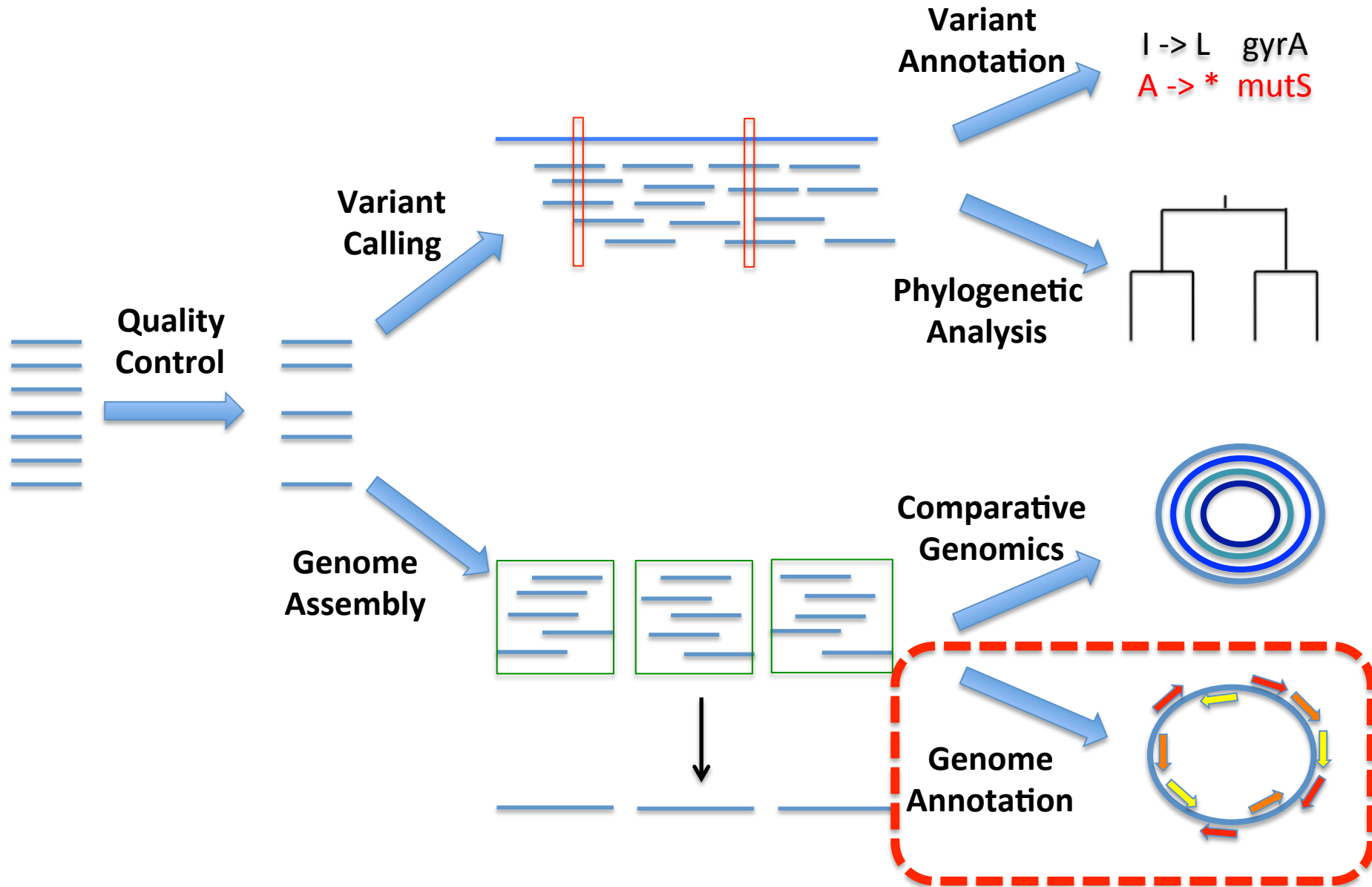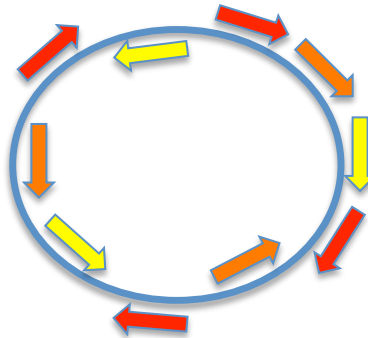
# Genome assembly

# Mile-high view of a genomics pipeline

# Genome annotation



```
>pseudo-molecule
ATCGTCGTGCTGC
TGCTGTCGTGCTG
CAGTGCATGTGCTA
GACTGTCGATGCTA
AGCTGTACCGATG
ACTGCTGACTGAC

.
```
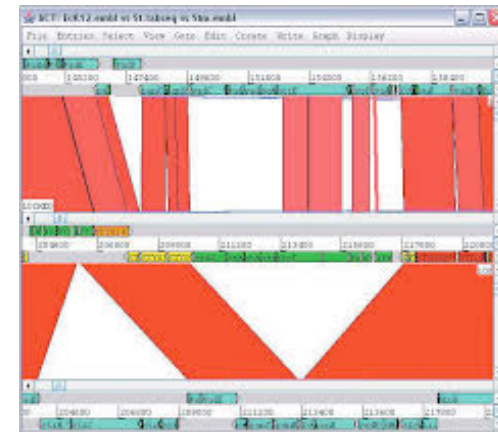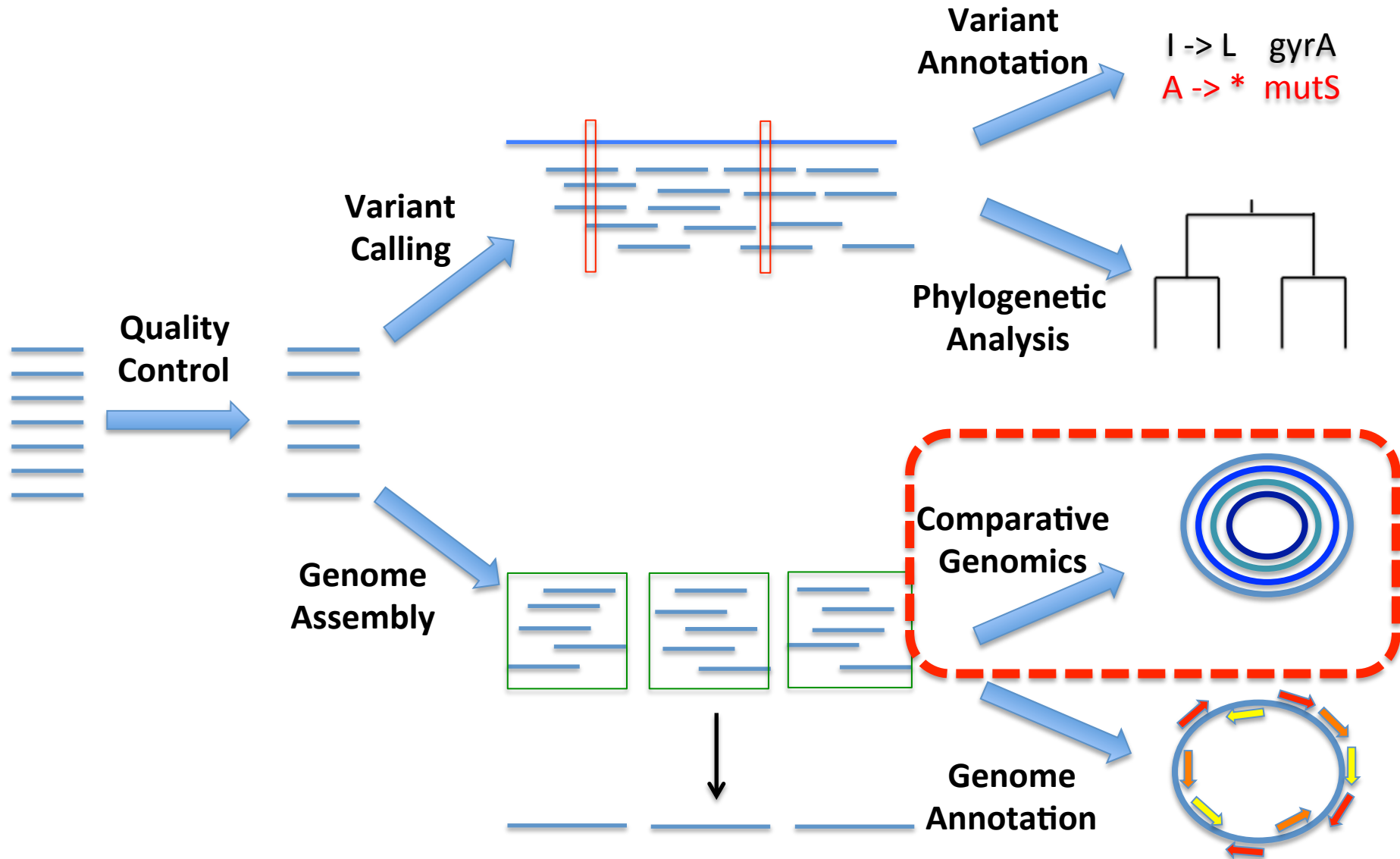
**Prokka**

1) Gene finding
2) Basic annotation

**ACT**

Visualization

**Fasta file**

**Genbank file**

**Genbank files, alignment files**

# Mile-high view of a genomics pipeline

Variant Annotation

I -> L    gyrA
A -> *    mutS

Variant Calling

Phylogenetic Analysis

Quality Control

Genome Assembly

Comparative Genomics

Genome Annotation

# Comparative genomics



**BLAST**

Genome mining

**LS-BSR**

Pan-genome analysis

**ACT**

Structural variants

**Fasta, genbank and/or pep**

|  | Gene 1 | Gene 2 | Gene 3 | Gene 4 | Gene 5 |
|---|---|---|---|---|---|
| Genome 1 |  | ■ | ■ |  | ■ |
| Genome 2 | ■ | ■ |  |  | ■ |
| Genome 3 |  | ■ | ■ | ■ |  |

Genome 1
Genome 2
Genome 3

# Mile-high view of a genomics pipeline

# Phylogenetics



**Orthofiner iqtree**

Tree construction

**iTOL**

Beautification

>Genome 1
ATCGTCGTGCTGC
TGCTGTCGTGCTG

>Genome 2
CAGTGCATGTGCTA
GACTGTCGATGCTA

>Genome 3
AGCTGTACCGATG
ACTGCTGACTGAC
.

Genome 1    Genome 2    Genome 3    Genome 4

**Multi-fasta file**

**Nexus file**

**Gubbins**
Recombination filtering

**Recipient**

**HGT**

**Donor**

**Seaview /ape**
Tree construction