

Genomics

An introduction to genome sequencing, data processing, and analysis

#DatAlumni course series @ ZMT, July 6/7, 2021 – part II

Martin Helmkampf



```
martin — haex1482@hpc1004:~/data/shared/Serranus — ssh haex1482@carl.hpc.uni-oldenburg.de — 1...
[~/data/shared/Serranus]> zcat < G09452-L1_S5_L001_R1_001.fastq.gz | head -n 4
@J00124:43:HM5WNBBXX:1:1101:22932:1173 1:N:0:GNTTCG
ACCCANGGGAGAACCTGCAAACCTCCATACAGAAGTGTCCGAGCNGGATTGAACTCACGACCCAGGNTCNGAACACNCNTGCNGTAAGGCATGAGTGCTAACAC
TNCGCCACNTGGAGGCCCTCAAACNGAACAGTTAGCANAAN
+
AAFFF#JJFJJJJJJJJJJJJAJAJJFFFJJF<A7FJJFJ<JJ#FFJJ7AF<FJJJ<JJFJJJJFJJ#FJ#JFFJJF#J#FFF#FJF7FJJAJJJJ-FJJJ<JFFJ
F#JJJFJJJ#JJJFFJJJ<FFFFJJ#JFFAJFJJFFJFJA#FJ#
[~/data/shared/Serranus]>
```

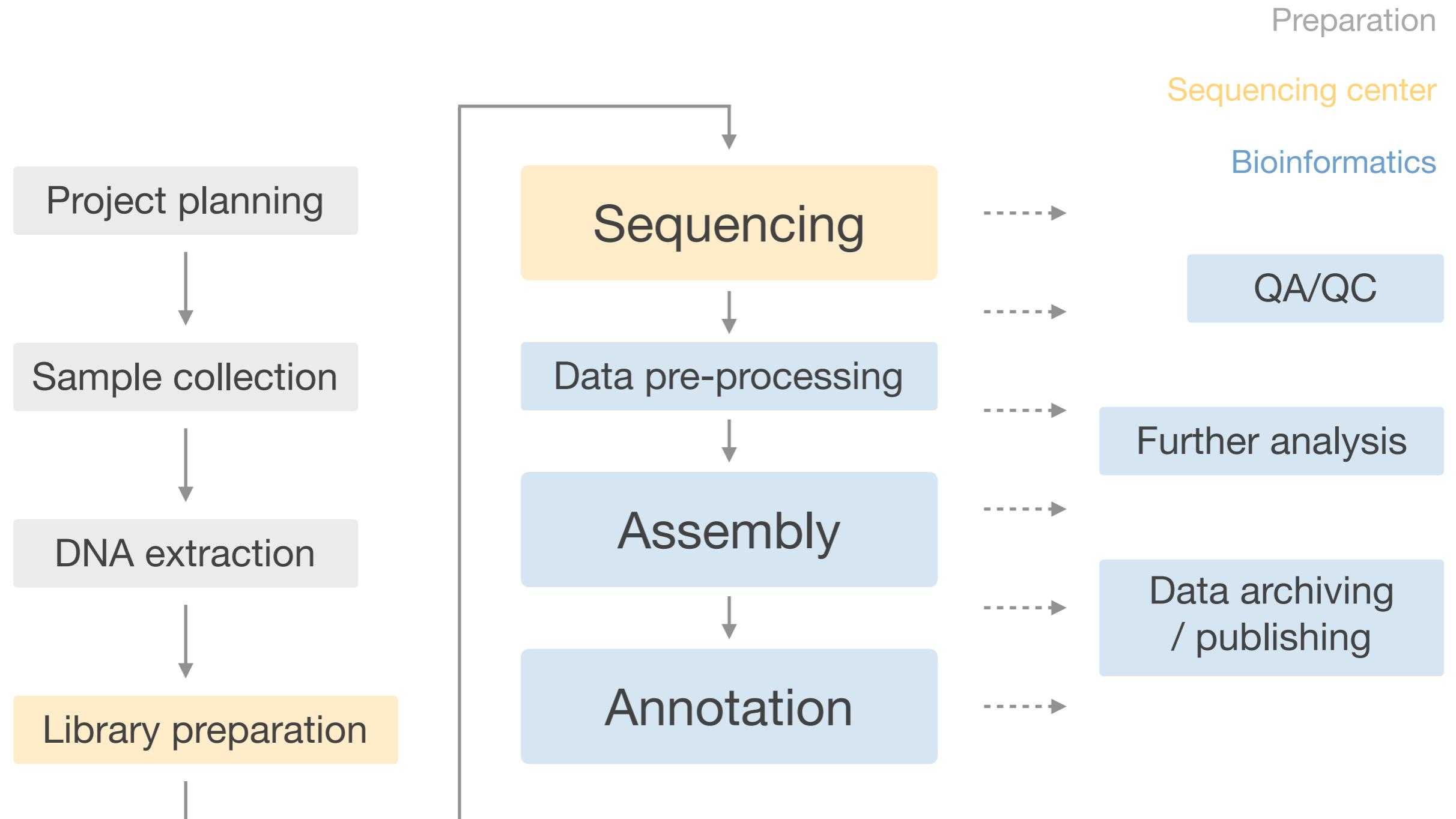
Workshop plan

Time	Block	Topics
Jul 06, 14:00–14:30	1. Background and introduction	<ul style="list-style-type: none">– Sequencing techniques and applications– Workflow overview– Planning a genome project
Jul 06, 14:30–16:00	2. Data pre-processing and QA/QC	<ul style="list-style-type: none">– Introduction to bioinformatics (exercise)– Raw sequencing data, trimming & filtering– Quality assurance / control
Jul 07, 14:00–15:00	3. Assembly and annotation	<ul style="list-style-type: none">– Genome assembly (exercise)– Genome re-sequencing– Genome annotation
Jul 06, 15:00–16:00	4. Data management	<ul style="list-style-type: none">– Documentation– Data organization– Publishing / sharing genome data

Objective: Convey basic knowledge of the steps required to sequence a genome, from project planning to data publication

De novo genome sequencing workflow

Legend



Recap: raw data

The FASTQ format

```
head -n 4 Hpu_e_raw300_F.fastq # display first 4 lines of file
```

1. @ followed by sequence id and optional info (e.g. instrument/run id, multiplex barcode)
 2. DNA sequence
 3. +, sometimes followed by sequence id
 4. per-base quality score (same length as sequence)

Recap: raw data

The FASTQ format

Phred **quality score**:

$$Q = -10 \log_{10} P$$

Common benchmark:

% bases with $Q \geq 30$

Quality score	P incorrect base call	Base call accuracy
10	1 in 10	90%
20	1 in 100	99%
30	1 in 1000	99.9%
40	1 in 10000	99.99%

FASTQ encoding (Illumina 1.8+):

ASCII Symbol:	!"#\$%&'()*)+,.-./0123456789:;=>?@ABCDEFGHIJ
Quality Score:	0.2.....10.....20.....30.....41

Recap: bash basics

Open shell (bash)

```
cd <dir>                                # change into workshop_zmt directory  
git pull                                  # update workshop example data via GitHub  
  
cd apps                                    # change into apps directory  
ls -l                                     # display directory content in long format
```

```
-rw-r--r--@ 1 martin staff 7902501 Jun 21 13:31 cutadapt-3.4.exe  
-rw-r--r--@ 1 martin staff 10249221 Jun 18 14:12 fastqc_v0.11.9.zip  
drwxr-xr-x 3 martin staff 96 Jul 7 10:44 test
```

d = directory, owner group file size modification file / dir name
permissions (bytes) time

Raw read quality check

Using the FastQC software

```
unzip fastqc_v0.11.9.zip          # decompress fastqc software  
FastQC/fastqc -h                  # confirm software works, show help  
  
cd ../../project/1_rawdata         # change into raw data folder  
../../apps/FastQC/fastqc Hpue_raw300_F.fastq    # run FastQC on file
```

```
Started analysis of Hpue_raw300_F.fastq  
Approx 5% complete for Hpue_raw300_F.fastq  
Approx 10% complete for Hpue_raw300_F.fastq  
...  
Analysis complete for Hpue_raw300_F.fastq
```

Results can also be found in `workshop_zmt/project/2_processing/qc`

Raw read quality check

Using the FastQC software



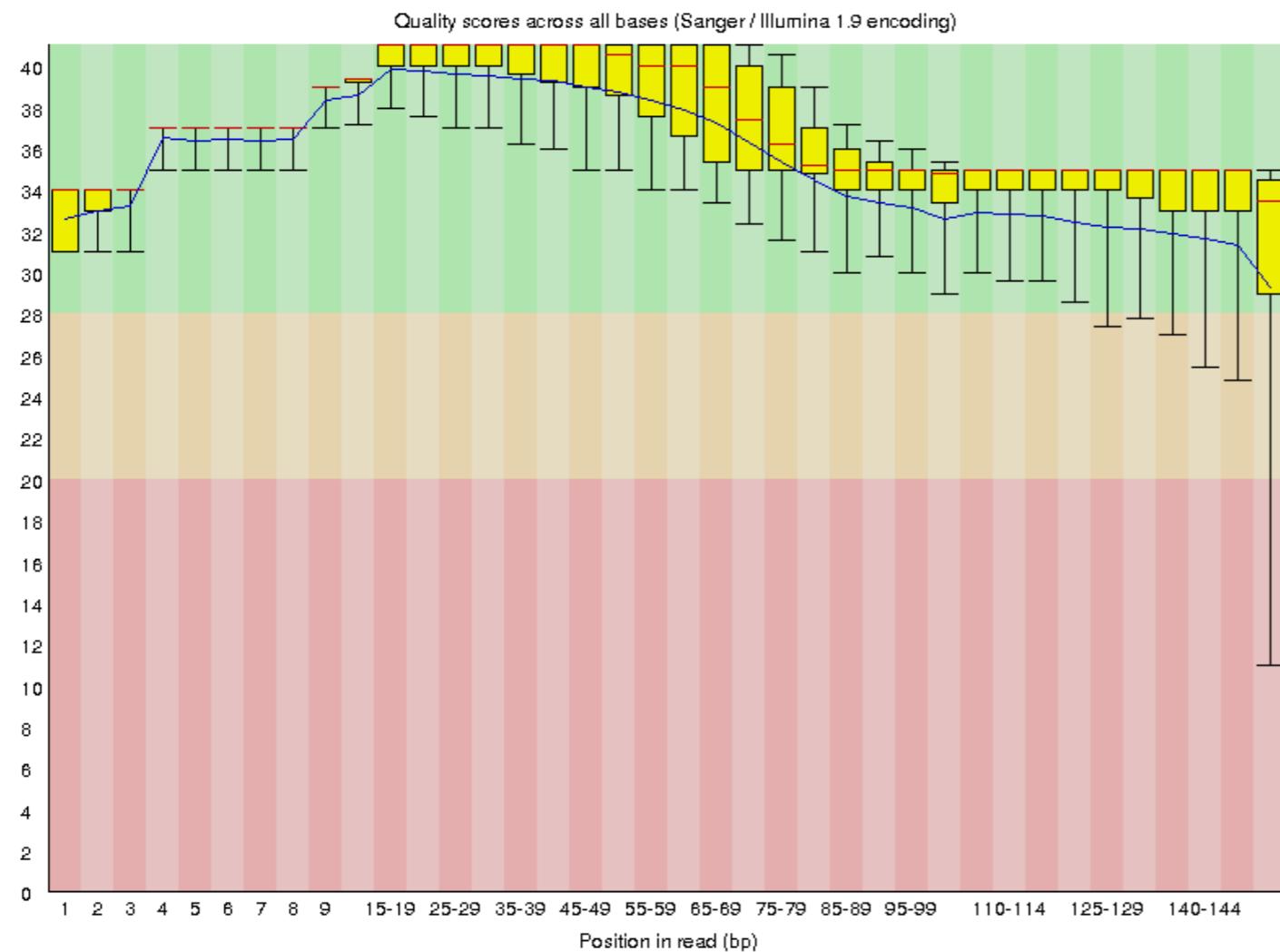
Basic Statistics

Measure	Value
Filename	Hpue_raw300_F.fastq
File type	Conventional base calls
Encoding	Sanger / Illumina 1.9
Total Sequences	75000
Sequences flagged as poor quality	0
Sequence length	151
%GC	40

Raw read quality check

Using the FastQC software

✓ Per base sequence quality



Raw read quality check

Using the FastQC software

⚠ Overrepresented sequences

Sequence	Count	Percentage	Possible Source
GATCGGAAGAGCACACGTCTGAACTCCAGTCACCGATGTATCTCGTATGC	640	0.8533333333333334	TruSeq Adapter, Index 2 (100% over 50bp)

Example for bad Illumina data:

https://www.bioinformatics.babraham.ac.uk/projects/fastqc/bad_sequence_fastqc.html

Read trimming & filtering

Using cutadapt for adapter removal and quality trimming

```
.../.../apps/cutadapt-3.4.exe -help          # confirm app works, show help (from 1_rawdata)

.../.../apps/cutadapt-3.4.exe -a AGATCGGAAGAGCACACGTCTGAACTCCAGTCA -q 10 -o
Hpue_raw300_F_trim.fastq Hpue_raw300_F.fastq
# run FastQC (-q 10: quality-filter last 10 bases, -a: trim Illumina adapter)
```

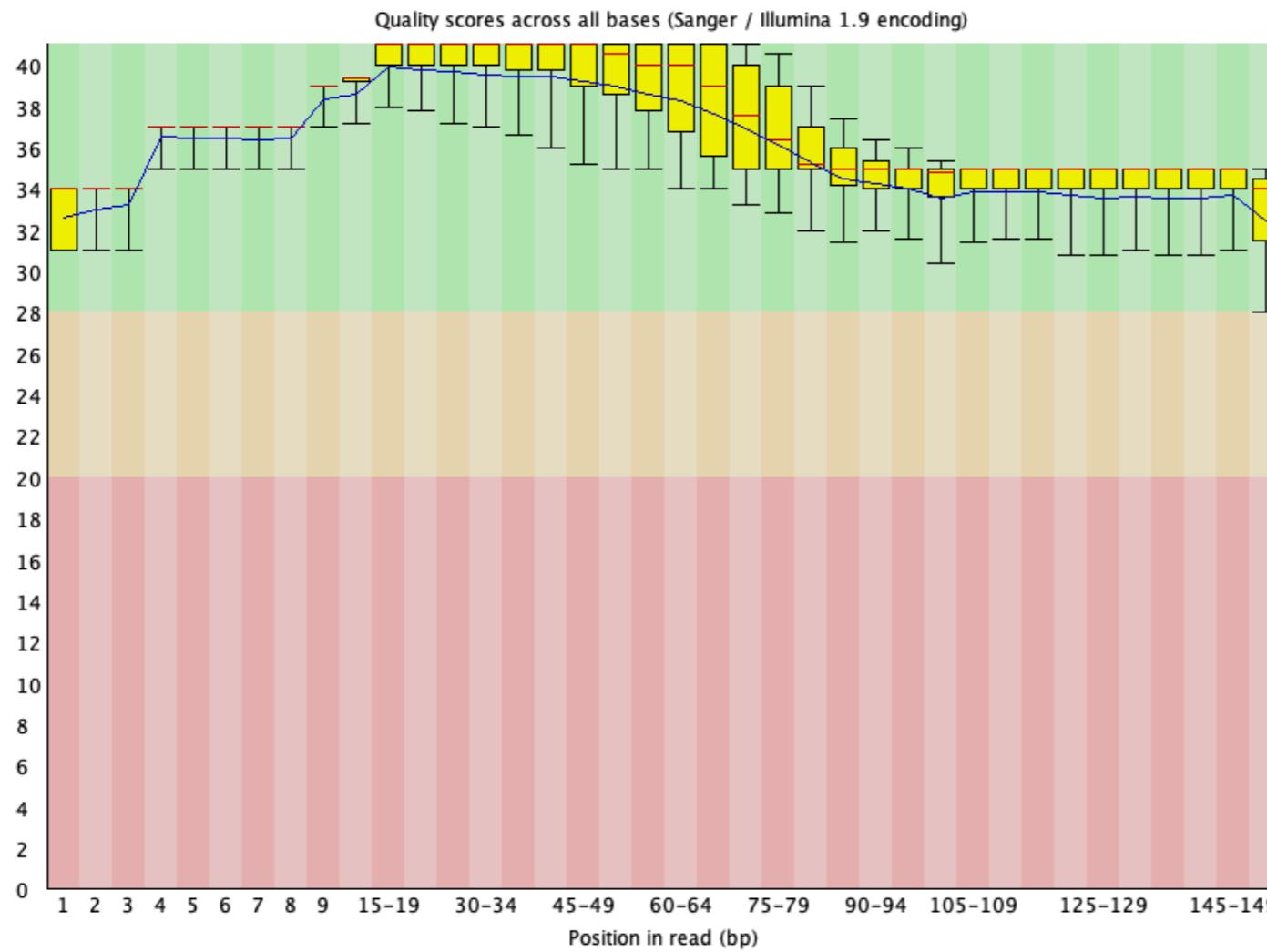
This is cutadapt 3.4 with Python 3.9.2
Command line parameters: ...
Processing reads ...

Rerun FastQC

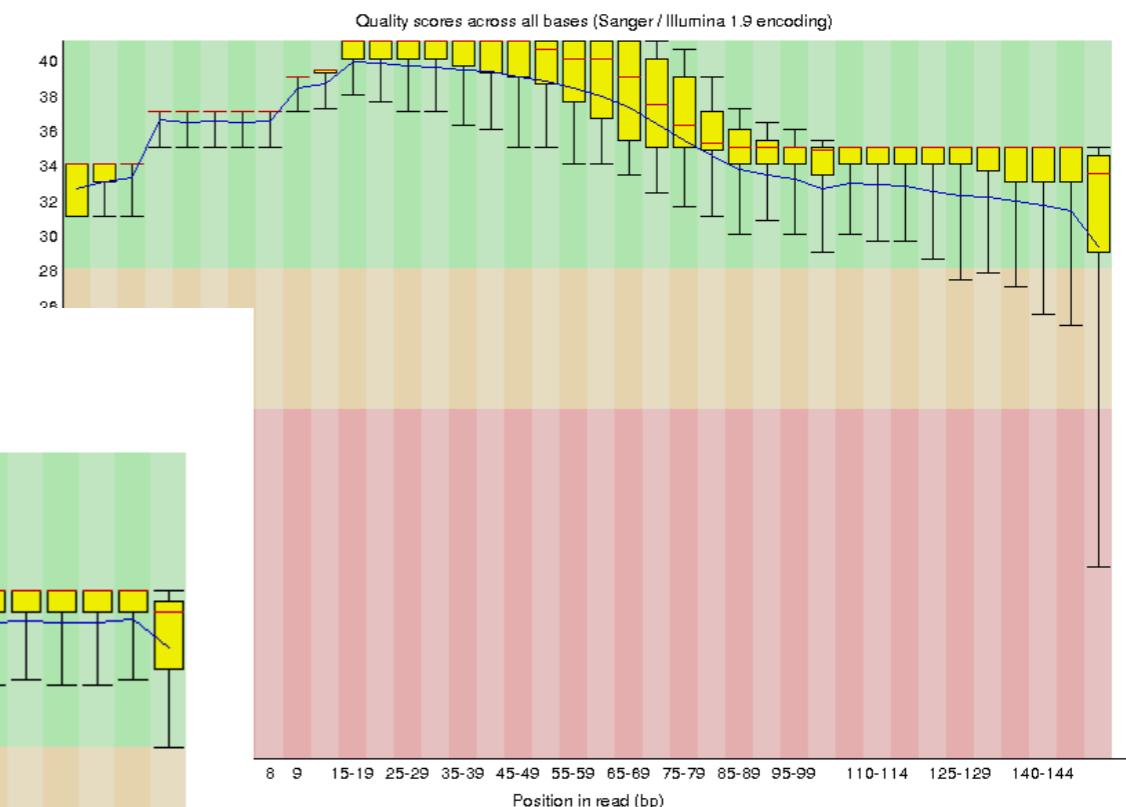
Filtered read quality check

Using the FastQC software

Per base sequence quality



Per base sequence quality



Raw reads

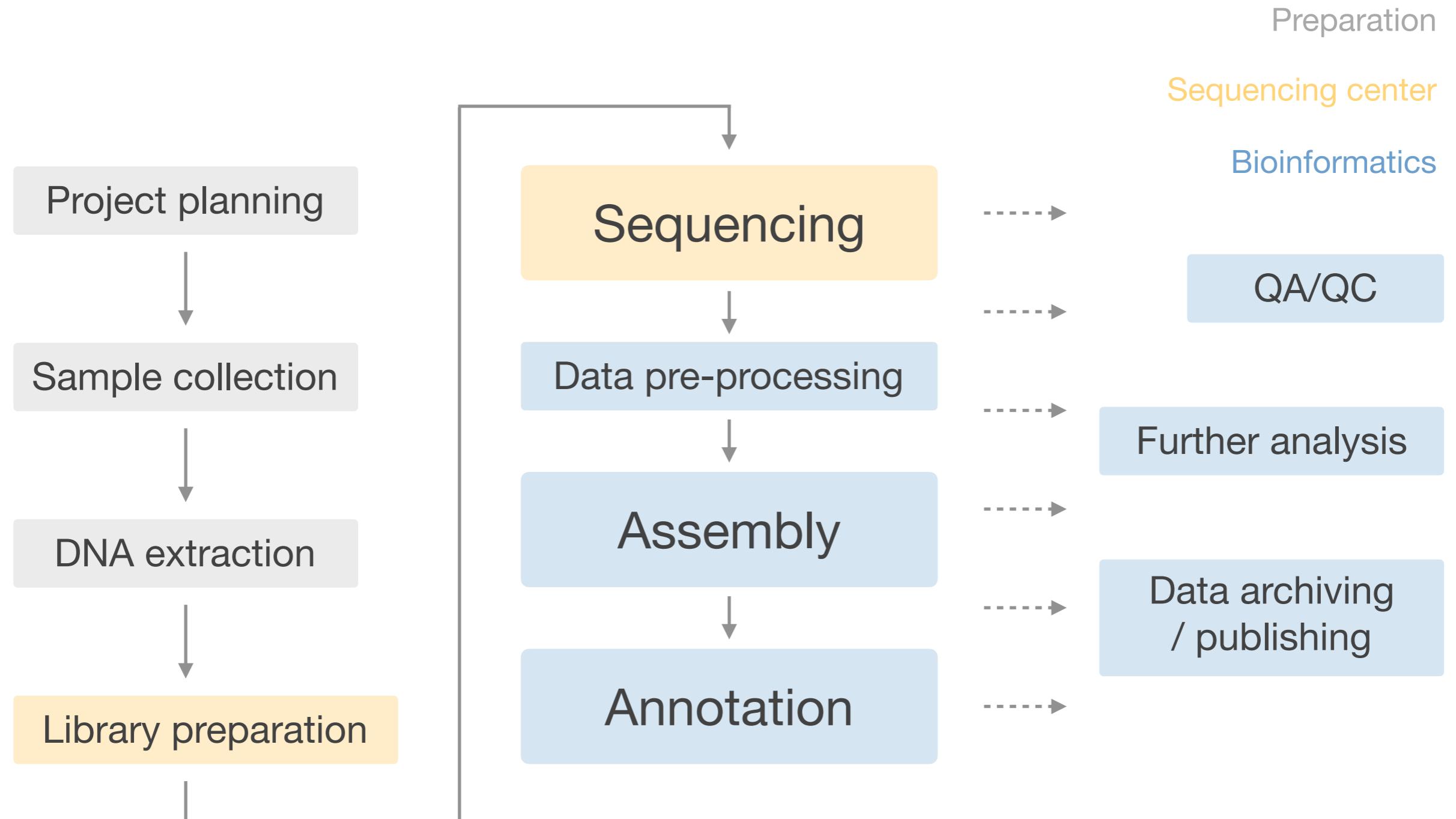
Trimmed reads

Data pre-processing summary

- Raw read quality assessment
- Adapter removal
- Quality filtering (trim low-quality bases / reads)
- Check for contamination (reads of viral, bacterial, human origin)
- Trimmed read quality assessment

De novo genome sequencing workflow

Legend



3. Assembly and annotation

- Genome assembly
- Metrics and assembly QA/QC
- Re-sequencing and variant calling
- Annotation

De novo genome assembly

- Reconstructing long, continuous sequence (up to chromosomes) from millions of fragments (reads)



- Computationally difficult due to repetitive (identical or highly similar) DNA
- Levels of assembly: Reads > contigs > scaffolds (> chromosomes)

= draft assembly = chromosome-level

Short read assembly

- Paired ends

(PE)

Forward read

Reverse read



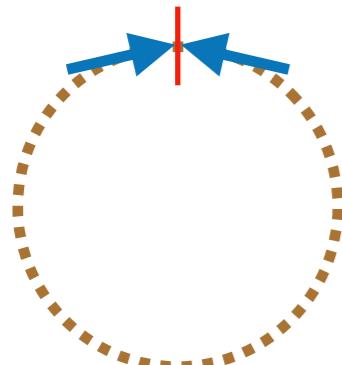
200–500 bp unknown sequence (insert)

Example file names

Hpue_raw300_F.fastq

Hpue_raw300_R.fastq

- Mate pairs



Forward read

Reverse read

2–20 kb unknown sequence (long insert)

Library type

Insert sizes

5 × pair

250–580 bp

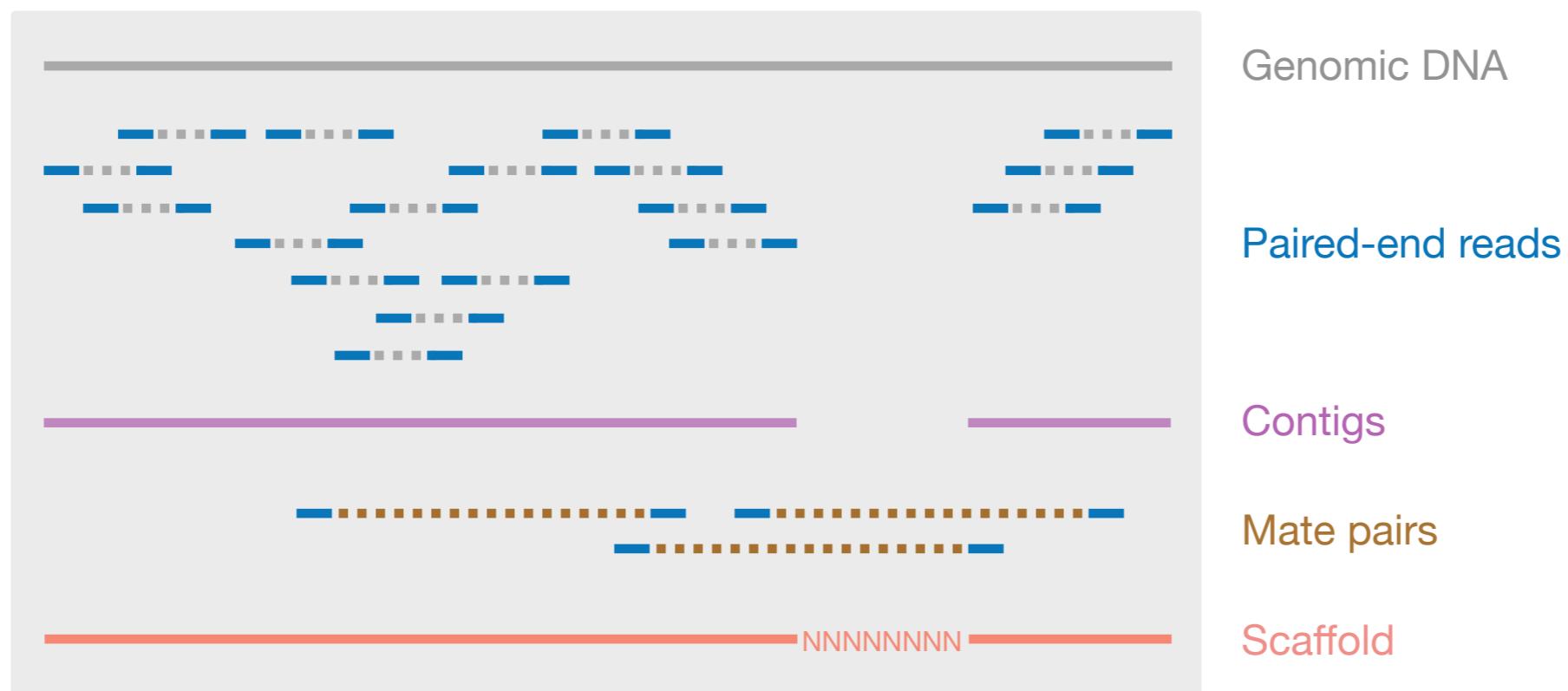
2 × mate

2.5, 4.3 kb

Example sequencing strategy

Short read assembly

- Overlapping reads > contigs
- Ordered, oriented contigs with gaps of known position, length > scaffolds
(gaps spanned by mate pairs)



Hybrid strategy

- Combine Illumina with long read sequencing
- Correct long reads (low coverage) with high coverage Illumina data
- Fill gaps in Illumina-based scaffolds with long reads (> up to chromosomes)
- Current state-of-the-art

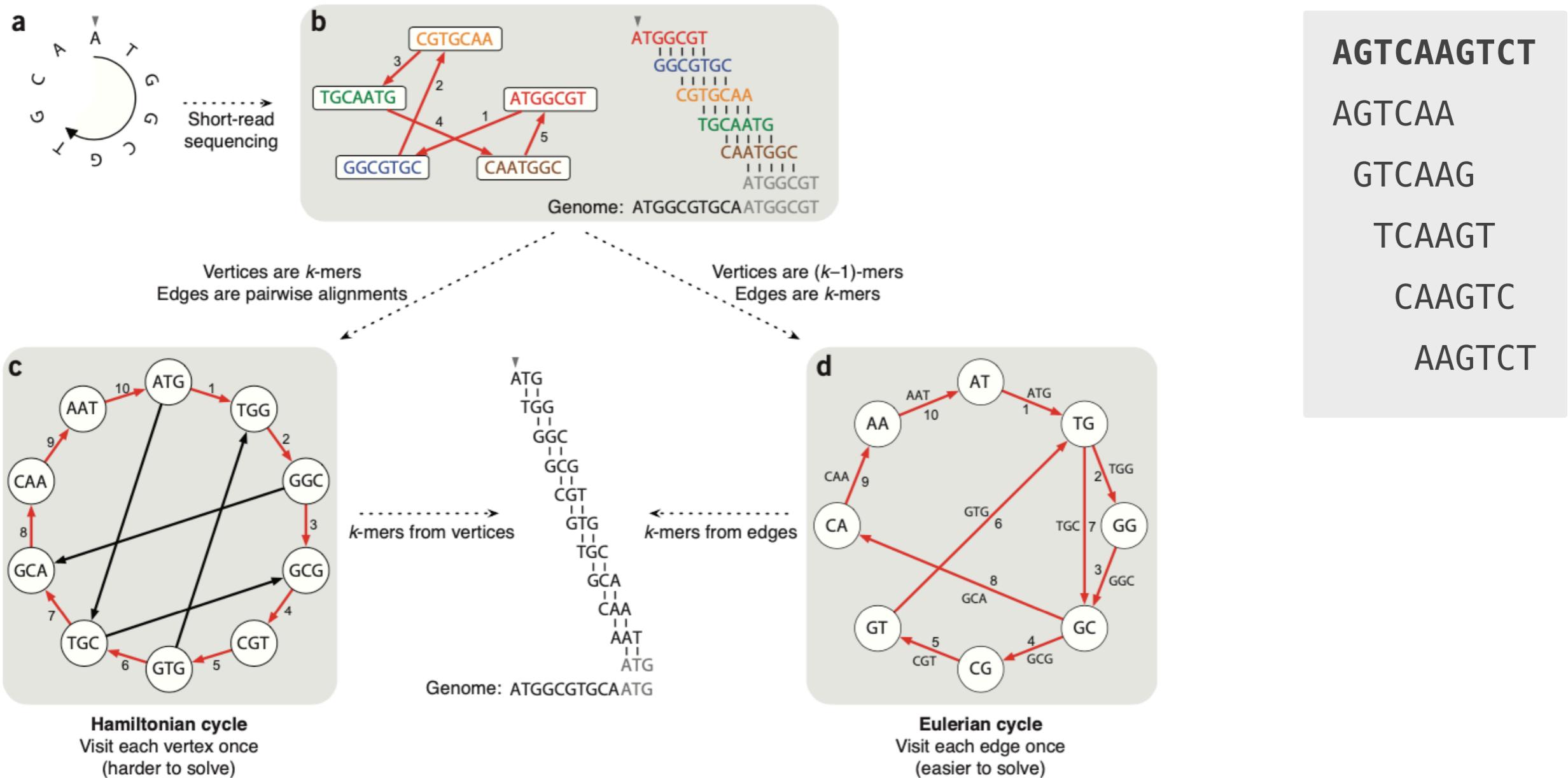
Assemblers

- Illumina: ALLPATHS, ABYSS, Discovar, Platanus, SOAPdenovo
- PacBio: Falcon, Peregrine, Canu (also Nanopore)
- 10x Genomics: Supernova
- Small genomes: SPAdes, Velvet, HGAP, MIRA
- Command-line based, though some graphical user interfaces are available for smaller genomes (< 100 Mb, e.g. Geneious)
- Expect run-times of multiple hours or days, depending on genome

De Bruijn graphs

Example:

6-mers of 10 bp-sequence



Compeau and Tessler 2011 (Nature Biotechnology)

De Bruijn graphs

Key assembly parameter: **k-mer size k**

- larger k : more specific, lower coverage, risk of under-assembly
- smaller k : less specific, higher coverage, risk of over-assembly

General advice

- Compare **alternate assemblies** (consider use case)
- Even high-quality assemblies should be **considered drafts**
- Expect lower quality in difficult (e.g. highly repetitive / heterozygous) regions

Assembly example

The FASTA format

```
cd workshop_zmt/project/3_assembly          # change directory  
zcat < Hpue_assembly_01_LG12.fas.gz | head -n 100    # look at compressed file (zcat),  
                                                       # display first 100 lines
```

```
>UXGA01000012.1 Hypoplectrus puella genome assembly, contig: LG12, whole genome shotgun  
sequence  
CCAAAAGTAGTTACAGTCCATGCGGCAAGTCCAGaactctctgcaccagactttTTTGAAGAAATGAGAGAGTTATGCG  
CCGGTCTAAACAAACTGAGGCTGCACTTTGAGCAGATTAAACAGCTCAGAAGGGACATGAGACCAAGATCAAGGCAAG  
ACTGGACCTCTGAAGAACCTTTGAACAACCTGTTCTCATCTCCTCAGCTGATGCCAGTGGACCTGTCTTCATCT  
...  
TGGGGATTAAGAGACGAAACACGTTAAACTCACATGAAACACTAAACTGATACTAGNNNNNNNNNNNNNNNNNNNN  
NNNTCGGTAGAAAAGTCTGTGCCTGTTGTCAGATGATTGAACTCATGAGCTGTCTACTGTCTGGTGTAT  
GTTAATCAAGTGTACTTCCGATAATGGAAACCTCTGATAACAGCTTTGTTGCCGCTGCTGATAAGCCCTGCAAAGC
```

Assessing assembly quality

- **Assembly metrics:** size distribution of contigs / scaffolds
- Gaps per Gb, unplaced contigs
- Base accuracy (Q score)
- Presence of false duplications
- **Gene completeness**
- Percentage of assembly assigned to chromosomes

Assembly metrics

- Total size (compare to expected genome size)
- Number of contigs / scaffolds
- Largest scaffold
- N50: contig / scaffold size where 50% of assembly is found on contigs / scaffolds of equal or larger size (measure for **sequence continuity**)

Scaffolds: 530, 760, 1050, 610, 450, 800, 220, and 1200 kb

Reorder: 1200, 1050, 800, 760, 610, 530, 450, 220 kb

Sum / 2: $5620 / 2 = 2810$

Add up until reached: $1200 + 1050 + 800 > 2810$

N50 = 800 kb

Example calculation

Assembly metrics

Using QUAST (command line tool or <http://cab.cc.spbu.ru/quast/>)

Report Example

E. coli single-cell assemblies

Aligned to "e.coli_reference" | 4 639 675 bp | 50.79% G+C | 884 operons

All statistics are based on contigs of size \geq 500 bp, unless otherwise noted (e.g., "# contigs (\geq 0 bp)" and "Total length (\geq 0 bp)" include all contigs).

	Worst	Median	Best	<input checked="" type="checkbox"/> Show heatmap
Statistics without reference				
# contigs	344	250	277	519
Largest contig	132 865	224 018	269 177	121 367
Total length	4 540 286	4 791 744	4 877 521	4 526 656
N50	33 616	96 947	106 927	20 445
Misassemblies				
# misassemblies	2	9	2	2
Misassembled contigs length	23 485	66 335	26 551	22 359
Mismatches				
# mismatches per 100 kbp	2.26	3.65	5.06	1.77
# indels per 100 kbp	0.7	0.2	0.7	0.92
# N's per 100 kbp	0	0	4.86	0
Genome statistics				
Genome fraction (%)	91.727	94.943	95.759	91.43
Duplication ratio	1.001	1.001	1.004	1.002
# genes	3767 + 160 part	4026 + 80 part	4046 + 102 part	3630 + 288 part
# operons	723 + 87 part	802 + 40 part	809 + 48 part	650 + 158 part
NGA50	32 051	96 947	110 539	19 791
Predicted genes				
# predicted genes (unique)	4258	4394	4417	4331
# predicted genes (\geq 0 bp)	4258	4394	4490	4331
# predicted genes (\geq 300 bp)	3643	3736	3784	3666
# predicted genes (\geq 1500 bp)	524	559	559	515
# predicted genes (\geq 3000 bp)	44	49	48	39

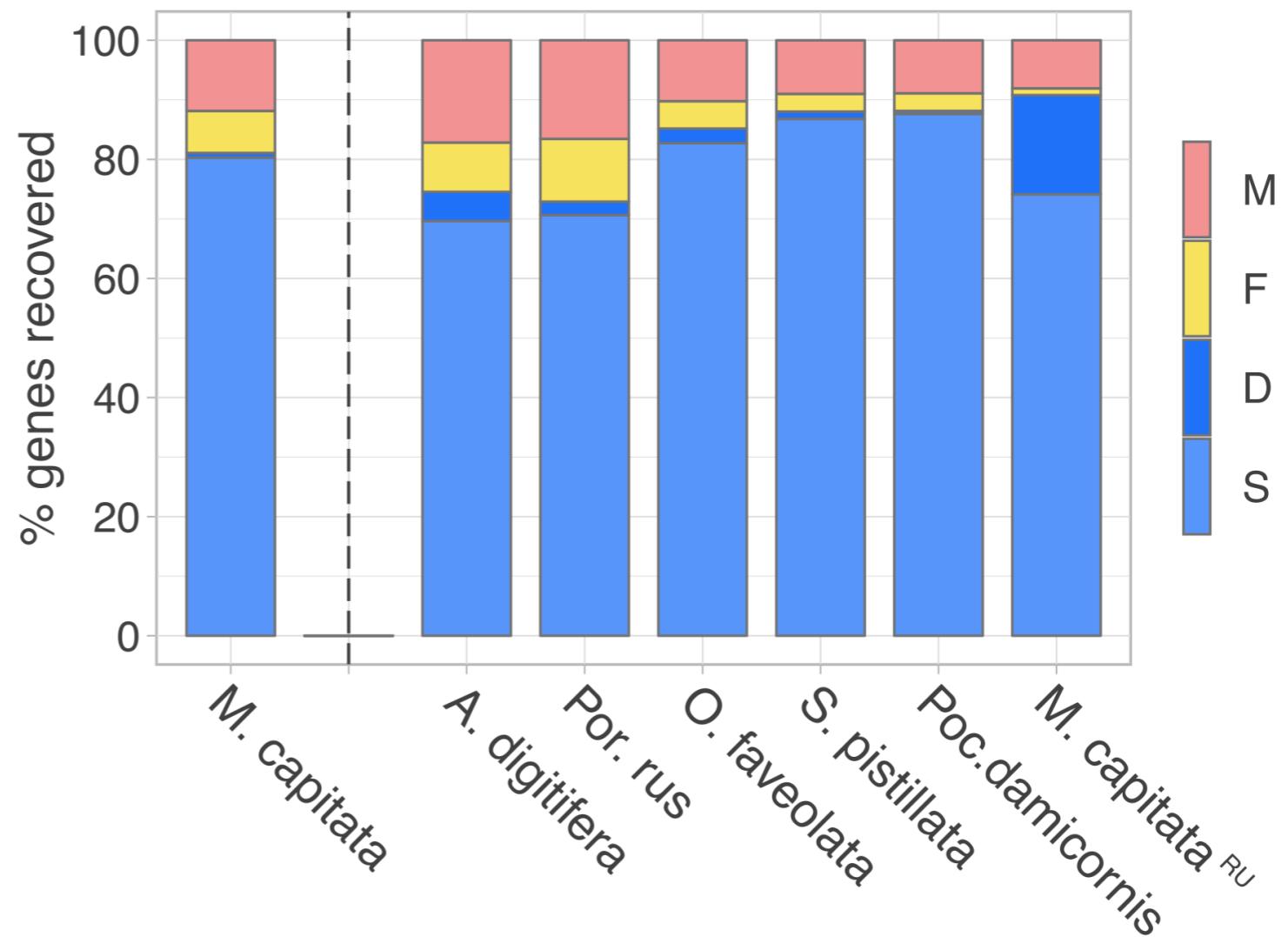
[Extended report](#)

Gene completeness

Using BUSCO

(<https://busco.ezlab.org>)

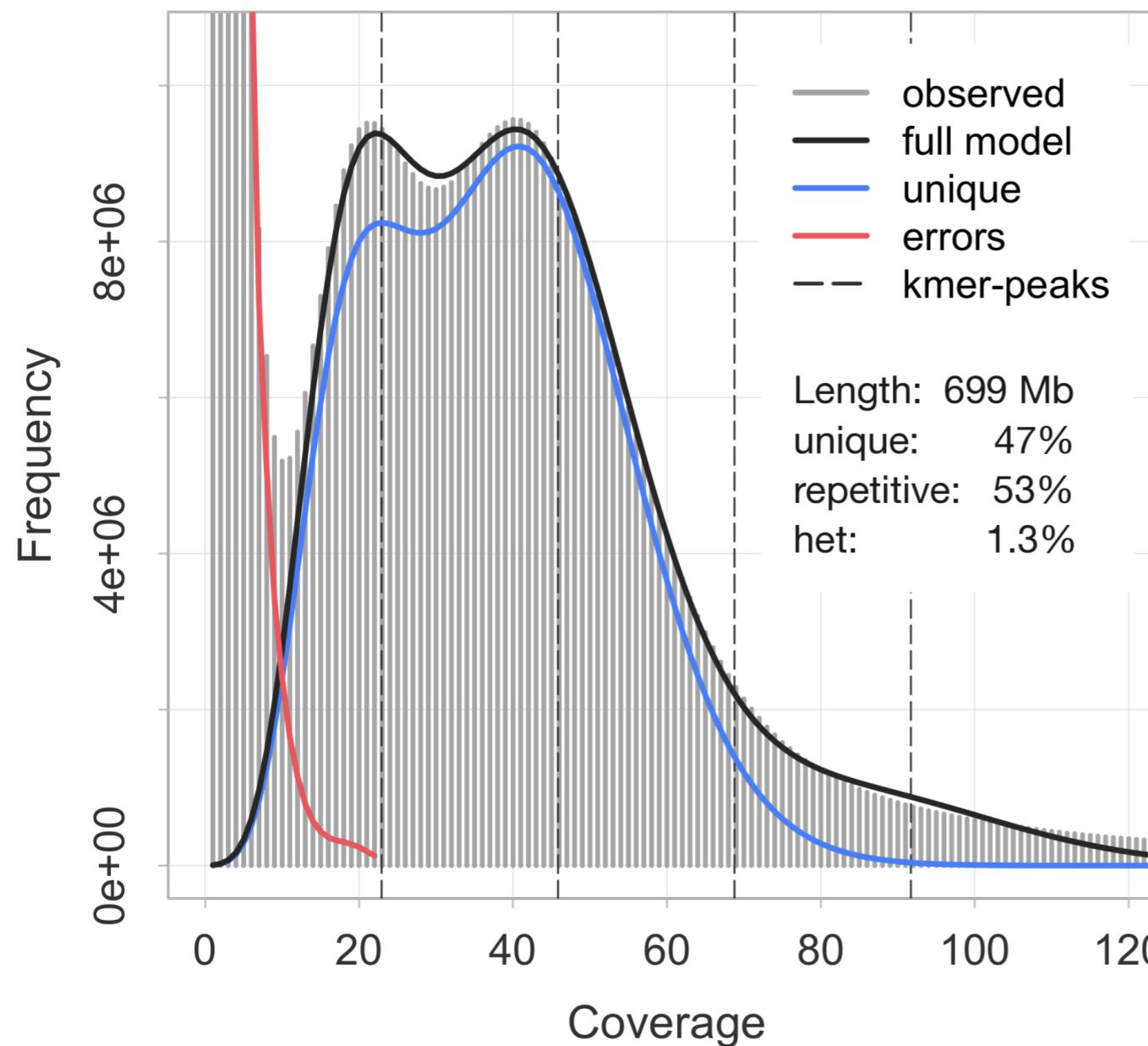
Quantifies assembly completeness based on presence of universal, highly conserved, single-copy genes (e.g. housekeeping genes)



Helmkampf et al. 2019 (Genome Biology and Evolution)

Genome assessment using k -mer profiles

Using GenomeScope (<http://qb.cshl.edu/genomescope/>)



Example:
6-mers of 10 bp-sequence

AGTCAAGTCT
AGTCAA
GTCAAG
TCAAGT
CAAGTC
AAGTCT

Helmkampf et al. 2019 (Genome Biology and Evolution)

Re-sequencing

- Prerequisite: high-quality reference genome
- Approach
 - low-coverage sequencing of multiple samples
 - aligning reads to reference (mapping)
 - identify differences > SNPs (variant calling / genotyping)



Re-sequencing

The VCF format

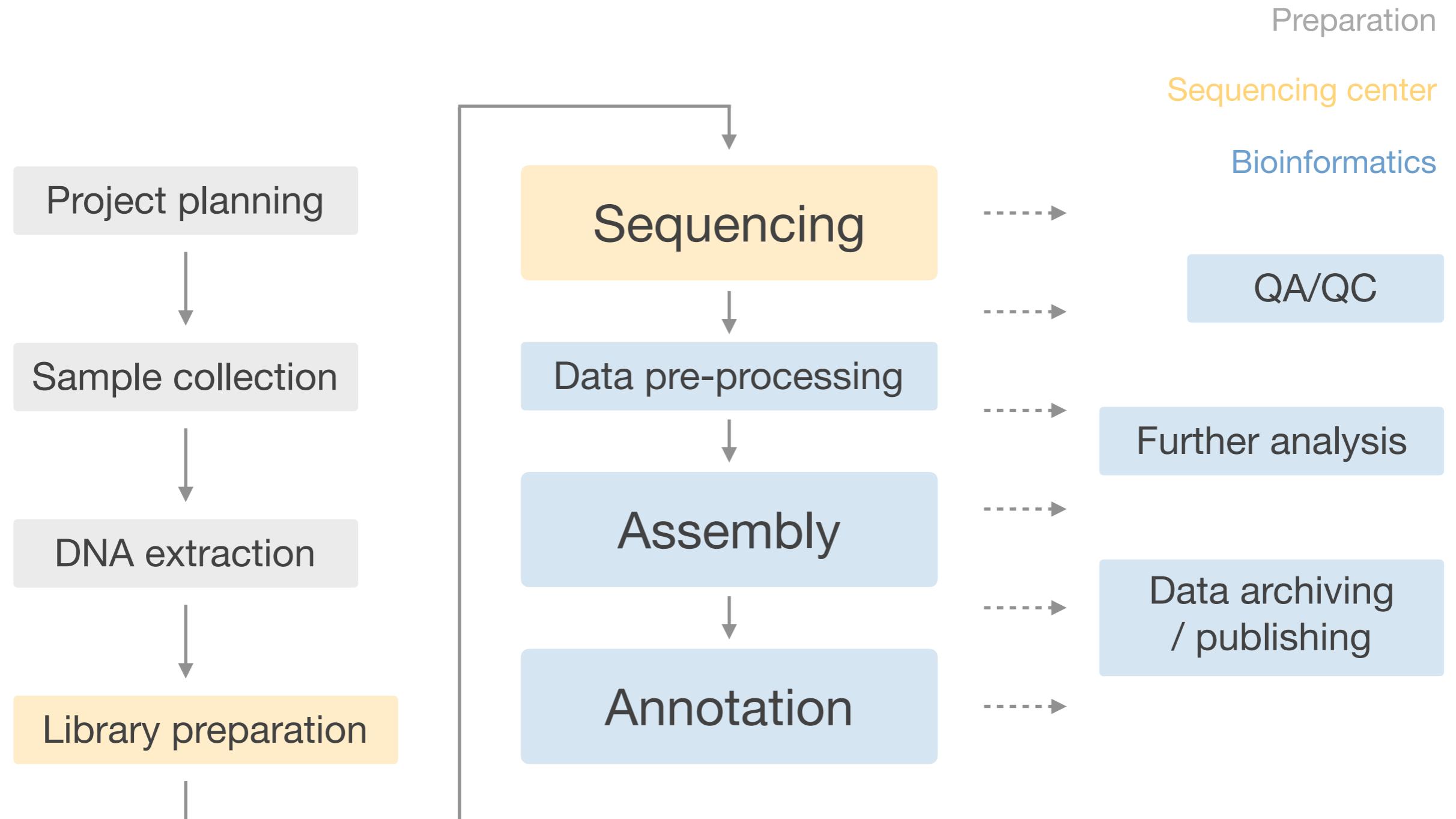
```
##fileformat=VCFv4.1
...
#CHROM POS ID REF ALT QUAL FILTER INFO FORMAT sampleA sampleB sampleC sampleD sampleE
LG01 1258 . C T . PASS . GT 0|0 0|0 0|0 0|0 0|0
LG01 1277 . G T . PASS . GT 0|0 0|0 0|0 0|0 0|0
LG01 1292 . G A . PASS . GT 0|0 0|0 0|0 0|0 0|0
LG01 1365 . C A . PASS . GT 0|0 0|0 0|0 0|0 0|0
LG01 1373 . C T . PASS . GT 0|1 0|1 1|1 0|0 1|1
LG01 1398 . G C . PASS . GT 0|1 0|1 1|1 0|0 1|1
LG01 1403 . G C . PASS . GT 0|0 0|0 0|0 0|0 0|0
LG01 1494 . C T . PASS . GT 1|0 0|0 0|0 1|1 0|1
LG01 1495 . G A . PASS . GT 0|1 1|0 1|1 0|0 1|0
...
```

A quick note on metagenomics

- Similar workflow: shotgun sequencing, read pre-processing, QA/AC
- Assembly stage: dedicated software like metaSPAdes, MetaQUAST
- Extra steps: **binning**, taxonomic profiling
- Annotation geared towards prokaryotes

De novo genome sequencing workflow

Legend



Genome annotation

Identifying the location, structure, and putative function of genome features

- repetitive and mobile DNA
- protein-coding genes
- regulatory elements

```
AGTCATTATTCTGCATCAACTTAACACACAACCTATTGGCTCGTTCAGCAGAGAAAGCAGATCAGTGAGCGTTC  
AAGGACCAATCATATTGAAAATCACTGAGCGTTACTGATCTGCTTCTCTGCTGAACGCAGCCATTCTGCGCTAA  
ATCCTTCTATATATTATTCTGTAATTCTGTTACTGGTATTACGGTCCGGTTAGGTCGGCCTCATTCTTCTTCAGCGT  
AGTTTCTAAGTATGTAGAGGATTAATAAATTATAAATGTATTACAGTCGAGGAAGCATAAAACCACAAAAATGGTA  
AAAATTCTATGCTTTCAATATATGCTAGATATGATCTTGTATATCATTATTATAATTGTTTAACTGTAATCGT  
GATTTGCTATTCTCTTTCTAAATTCTGCTCTATCTACTATGAGAAATATAGATTCTATTTCATTTCTCAAATCT  
AGTTAATCTAAATATATCAATGTATGTAATGAGTTTATTACCTGATCATATAGAAAAAGGCTCAGGAGAGCATTCAAC  
CTCCGCCAAGGCAAAGCTAAATTGGTCATCATTGCTAGCAATACGCCACCGCTAAGGTGAGACTGAGAAGTAGCACT  
TTAATGAATATAATCTTGATAAAATGTATACATATATACATATACATATTGTAAAAATTATATTCTATTGTAATT  
GTAATATTGATATAATCATCAATTGCAATTGGAAAACAAATAGGAGTACTATGGAGTAGTCTCCAGCTCTGCTTG  
GAAGTCGGAGATTGAATACTATGCAATGTTAGCGAAGACTGGTGTGCATCATTACACGGGAATAACATCGAACTGGG  
TATAAAATCTATGCCAACTGGTGTACAGCGTAATGTACAGTTTAATCCAATAAAATTCAAAACGTTTGATT  
TATTAATCTATATTAGTTATTTAGCAAGAGTAACATCAAAGATTCTCCTTAGATTTCACCGACTTGAAATA  
CTGAAGTTTCATTGTGTACATTAAACAAATCCGATTTCACGAATGCTAAGAATTGTACAATAAAATATAG  
AAAATCTTTAACATTATCCAAAAACTTTATACATATGTATATGTACATATACACCGCGCGCGCGTGT  
GTGCATCGTTCTTCAAATGATTGGTTGTTACAATAACGACAAATTGACCTTGCACGGACATGCAGTTACTGCA  
ATCTGCCACAGTTGTACATAGTCGCCAATGCCGACAACGTGATTGAATATGATTGTTCCACATTAGTACAA  
ACACCCCAGGCTGGACATCGTTCCGATCATTGCTGCCACATCGAACACATTCTATTATAGATTAGTCAGGAA  
ATAACAGCTAGATTGCTTTCAAAAACGGTATCCAATGTTCCCATGATTGTTATGGTAACTGCAGCTTCC  
TCGCCTCAACAGCAATGGACTGCAAGAAATATACAAATTATTAGAACATGCAATTGAGATTAGACATGATTAGCTA  
TTTTTATCTTCATACAAGTAATAGATATCAAACCTTAAAGCTACTTATGTTCTGCAATTCTCATTCA
```

Repeat annotation

- Most genomes contain large amounts of repetitive and mobile DNA
- Repeat classes:

Repetitive Elements Identified in the *Montipora capitata* Genome Assembly

	Number	Total Length (Mb)	Fraction of Assembly (%)
Tandem repeats	160,985	16.5	2.7
Interspersed repeats	1,028,006	257.4	41.9
DNA elements	31,667	11.6	1.9
LTR elements	14,753	10.0	1.6
Non-LTR elements	110,728	42.6	6.9
Unclassified	870,858	193.2	31.4

- Repeats may have regulatory or structural functions, but can also bias downstream analyses

Helmkampf et al. 2019 (Genome Biology and Evolution)

Repeat annotation

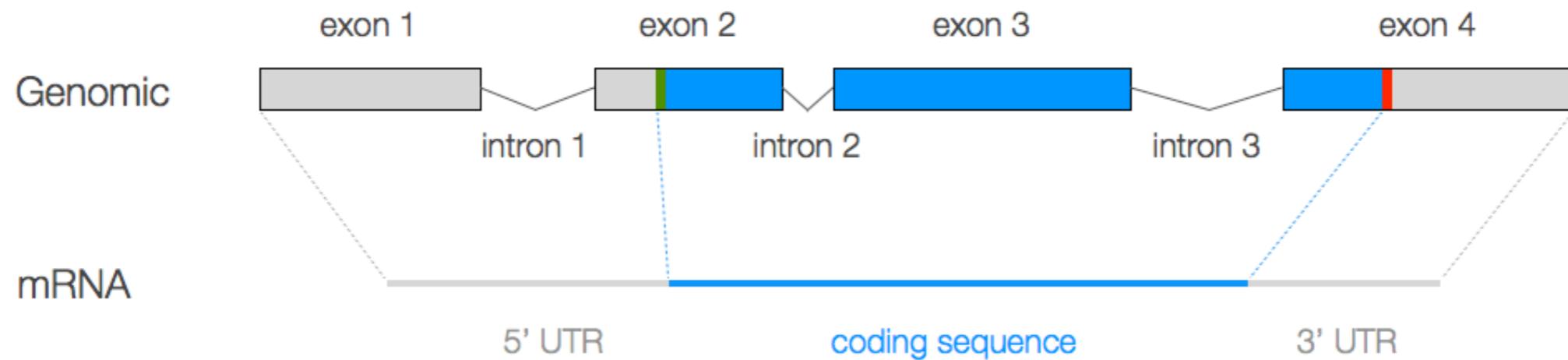
- **RepeatModeler** combines multiple repeat search engines and curated databases (e.g. Repbase) to identify repeats
- **RepeatMasker** marks repeats in assembly (replacement with Ns or soft-masking)

```
>Sc0000000
TAATGGCAGTCCAATCTGTACCCCTGATTATAATTCATTCAAATAATTCTTTTTATCCAATGATCAACCAAATGCATTGCTGCATTAAAAACTC
TTGTCAGCATTAAGTCCTCAGATATGGATCAACAAACACAGTAAAACCATTAAATCACTAAACTTAAAATAATAATTATTATTGAAGAATCTAAAtca
tcatttcttcacaggaacacacaagccaaacaaattgacctgttccaaacagattggcttcatalogtcagttggttctagcaattgcattcattggta
tcacagaggtaatgggttcaaattccattggagctgcctgaattttggtgtctTTTAAAGTGAUTCAATTCTATCTTCTCCCTGAATGCAGCG
TGTTAAGTATGCTTTTGAAACACAAGTTGAGAATCCTGTACTGCAAACATAATCCTTAGTTGCAGTGTATacacacacacacacacatacTGGAA
AATTACGGGTAAGAAATGTCAATCTCTCCTCATTCCAGTTCTCTCCAAAATAATGTTGAAGTTGAGGTAATGGTGGACCACGTGAATAGTGAATTCCT
GTTTCATCCTGTTATGTTACAATCACATGTCACATTGATTAATGGAAATATAATCCCCCTTAGAACCGATTGCAGACAGTCATTCTTAACTTG
...

```

Gene prediction

In eukaryotes, prediction of protein-coding genes is difficult due to coding sequence (**exons**) being interrupted by **introns**



Gene prediction

- Complex pipelines (e.g. AUGUSTUS) combine multiple approaches
 - *ab initio* prediction (pattern-based)
 - extrinsic evidence can include
 - RNA-seq (transcripts)
 - homologous protein sequences
 - Sensitivity and specificity are increased by multiple rounds of training
 - Repeats should be masked
- Computational gene prediction always requires careful quality control

Gene prediction

The GFF format

```
# start gene g29
Sc0000000 AUGUSTUS gene 656020 658053 0.9 - . g29
Sc0000000 AUGUSTUS transcript 656020 658053 0.9 - . g29.t1
Sc0000000 AUGUSTUS stop_codon 656020 656022 . - 0 transcript_id "g29.t1"; gene_id "g29";
Sc0000000 AUGUSTUS intron 656165 657273 1 - . transcript_id "g29.t1"; gene_id "g29";
Sc0000000 AUGUSTUS intron 657509 657956 1 - . transcript_id "g29.t1"; gene_id "g29";
Sc0000000 AUGUSTUS CDS 656020 656164 1 - 1 transcript_id "g29.t1"; gene_id "g29";
Sc0000000 AUGUSTUS CDS 657274 657508 1 - 2 transcript_id "g29.t1"; gene_id "g29";
Sc0000000 AUGUSTUS CDS 657957 658053 0.9 - 0 transcript_id "g29.t1"; gene_id "g29";
Sc0000000 AUGUSTUS start_codon 658051 658053 . - 0 transcript_id "g29.t1"; gene_id "g29";
# coding sequence = [atgaccgcgttcgaaagcgacgattctcaattggatgaagatagtgatcatggaaattgttcagcaaaacggttggcc
# caaaaagaaaaattcaggtcggtgcgcaggaaatcaaaagaatcatatgcgcctgggtcccaagggttgatgtaactccagcgagcgtacgatgaa
# ctggagcggaaacgtctttatgaactctggtaccttaaccagtcgagggtttccaaccctataccgtgagctccggagaggacgactttga
# tgaccttagtgtggagacatgaaaaagcgtcgagctaaaaacgaaagaattccgacaacgcgaagatacccacgatacacctgattgacagaca
# atgagggagagctgggttaacacctaagccgccttacgtcaacgaaggtgcccacacttcggatgatgccaatctgaactggaaaaggaa
# aaatga]
# protein sequence = [MTRFESDDSQLDEDSDHGNCSAKRLAPRKISGRLRRKSKESYRLGPKVDPASSDDELERETSFMNSGTFKPVEGF
# SNPIPVSSGEDDFDDLSVGDMKRRSSKTKEFRQREDTHDTPDLTDNEGELGFNTFKPASYVNEGAHTSDDAESELEKEK]
```

Gene prediction

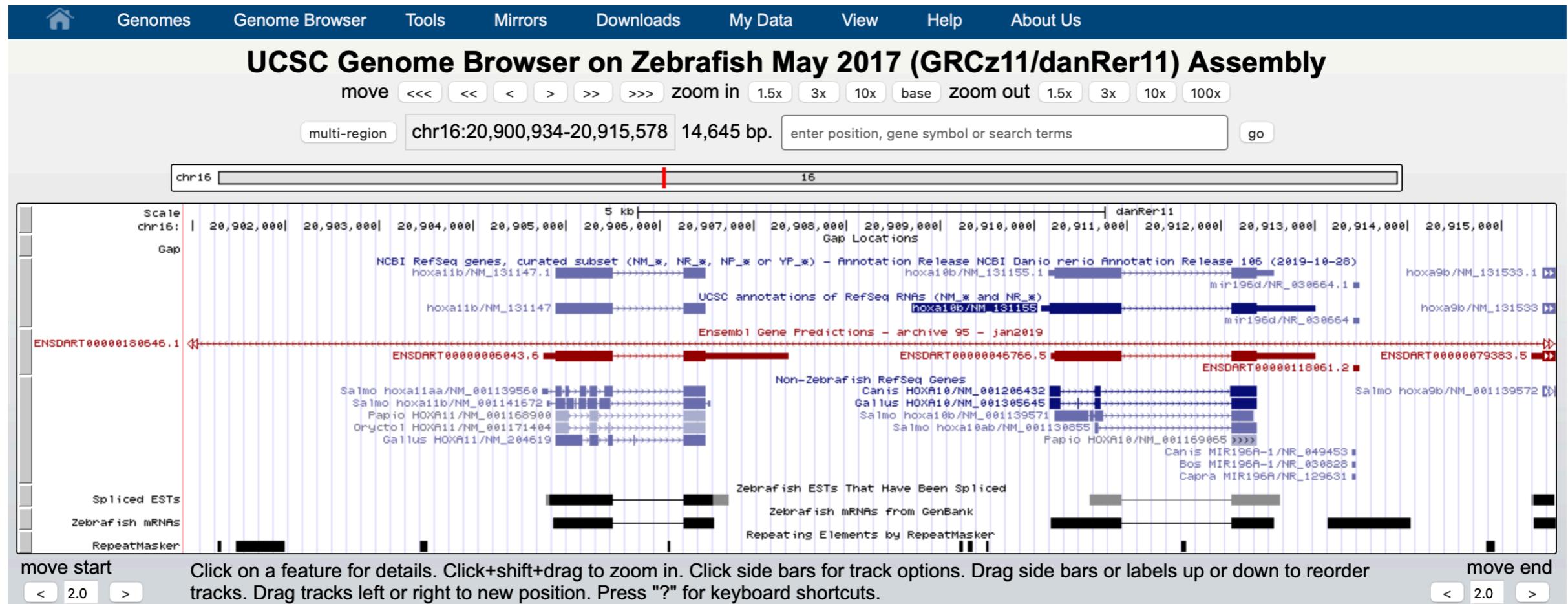
- Gene structure (UTRs, exons, introns)
- Coding sequence
- Transcription status and splice variants
- Putative identity and function

mostly homology-based, e.g. by BLAST versus UniProt (uniprot.org) database

- Computational + manually curated gene models = **Official Gene Set**

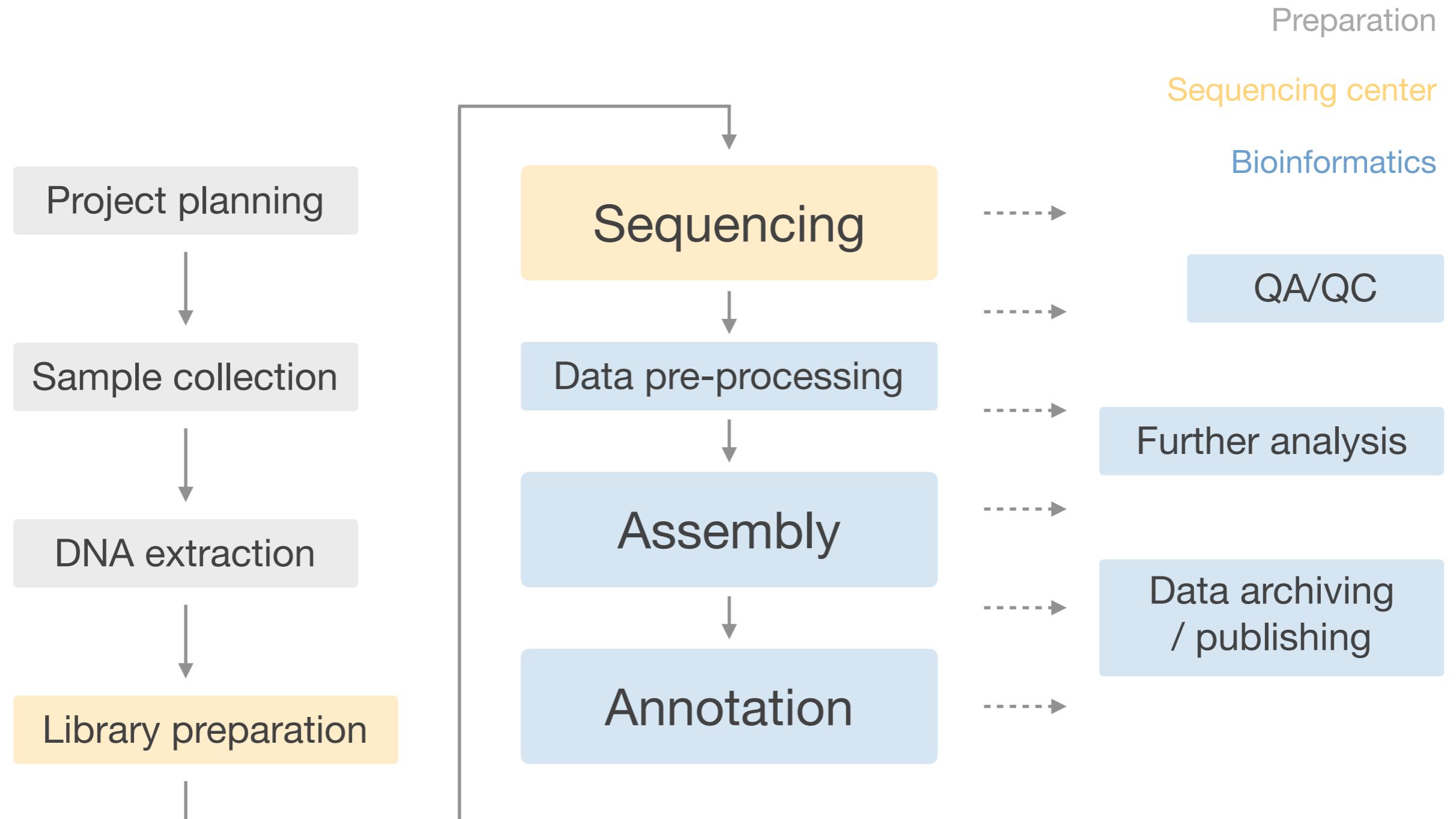
Visualization

GBrowse / JBrowse genome browser



De novo genome sequencing workflow

Legend

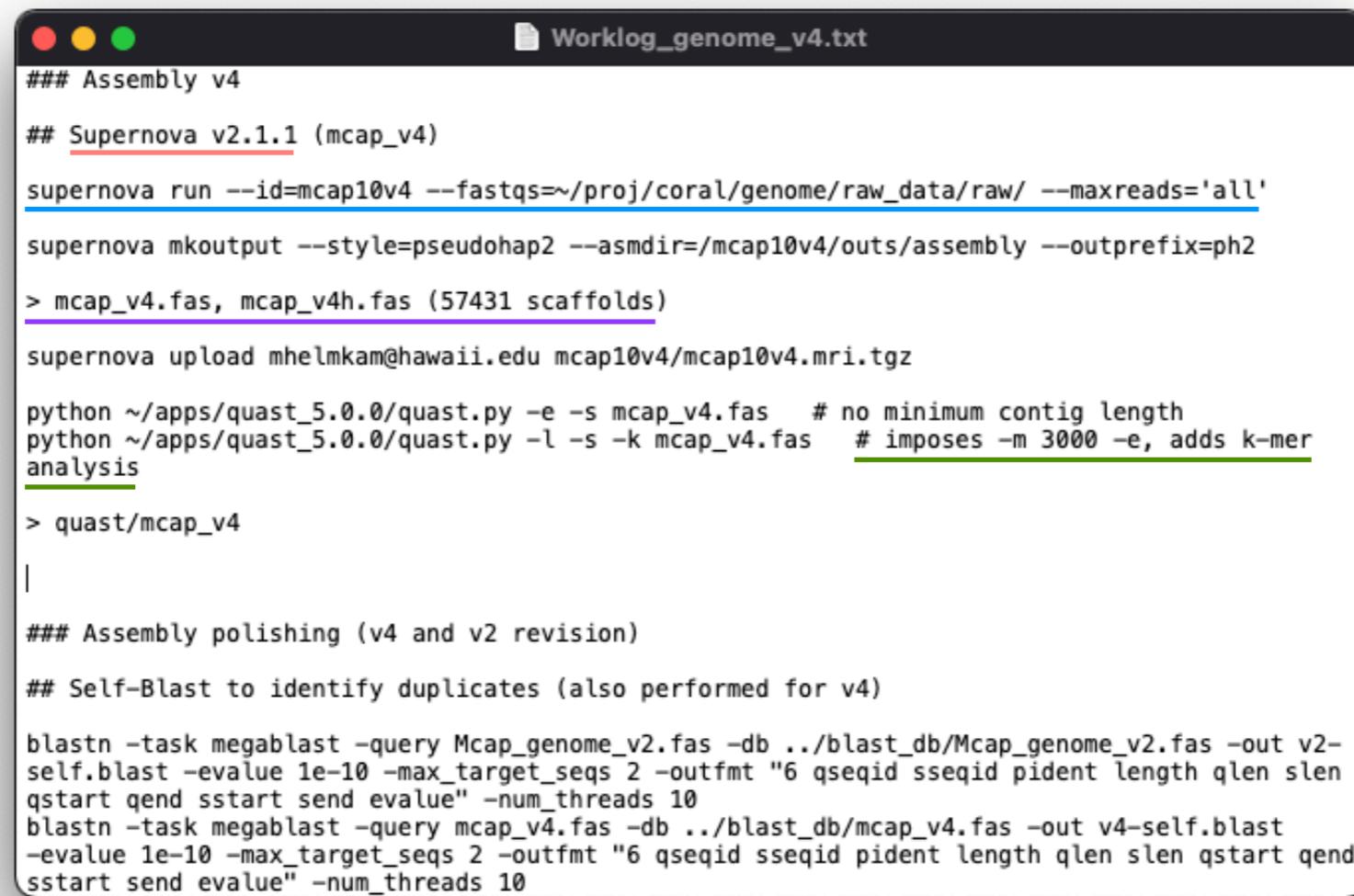


4. Data management

- Documentation
- Data organization
- Backing up / Version control
- Publishing / Sharing data

Documentation

- Document your workflow in plain text format (Atom /TextEdit / Notepad)
- Include **code**, **software versions**, **summary results**, and **comments**



```
## Assembly v4
## Supernova v2.1.1 (mcap_v4)
supernova run --id=mcap10v4 --fastqs=~/proj/coral/genome/raw_data/raw/ --maxreads='all'
supernova mkoutput --style=pseudohap2 --asmdir=/mcap10v4/outs/assembly --outprefix=ph2
> mcap_v4.fas, mcap_v4h.fas (57431 scaffolds)
supernova upload mhelmkam@hawaii.edu mcap10v4/mcap10v4.mri.tgz
python ~/apps/quast_5.0.0/quast.py -e -s mcap_v4.fas # no minimum contig length
python ~/apps/quast_5.0.0/quast.py -l -s -k mcap_v4.fas # imposes -m 3000 -e, adds k-mer
analysis
> quast/mcap_v4
|
### Assembly polishing (v4 and v2 revision)
## Self-Blast to identify duplicates (also performed for v4)

blastn -task megablast -query Mcap_genome_v2.fas -db ../blast_db/Mcap_genome_v2.fas -out v2-
self.blast -eval 1e-10 -max_target_seqs 2 -outfmt "6 qseqid sseqid pident length qlen slen
qstart qend sstart send evalue" -num_threads 10
blastn -task megablast -query mcap_v4.fas -db ../blast_db/mcap_v4.fas -out v4-self.blast
-eval 1e-10 -max_target_seqs 2 -outfmt "6 qseqid sseqid pident length qlen slen qstart qend
sstart send evalue" -num_threads 10
```

Data organization

Formats

- Standard file formats (e.g. plain text, pdf, png)
- Standard data formats (e.g. FASTA, FASTQ, VCF, GFF)

Data organization

Naming scheme

- Uniform, persistent
- Comprehensive, includes important information
- No spaces and special characters (use underscores instead)
- Consistent file extensions (e.g. .fastq, .vcf, .gff)
- May include date (YYYYMMDD) or version number

Data organization

Naming scheme examples

What is it?

Species acronym

Hpue_raw300_F.fastq.gz

File / data format

Mcap_genemodels_1.1_aa.fas

hyp155_a_0.33_mac4_5kb.raxml.log

Dataset name

Version number

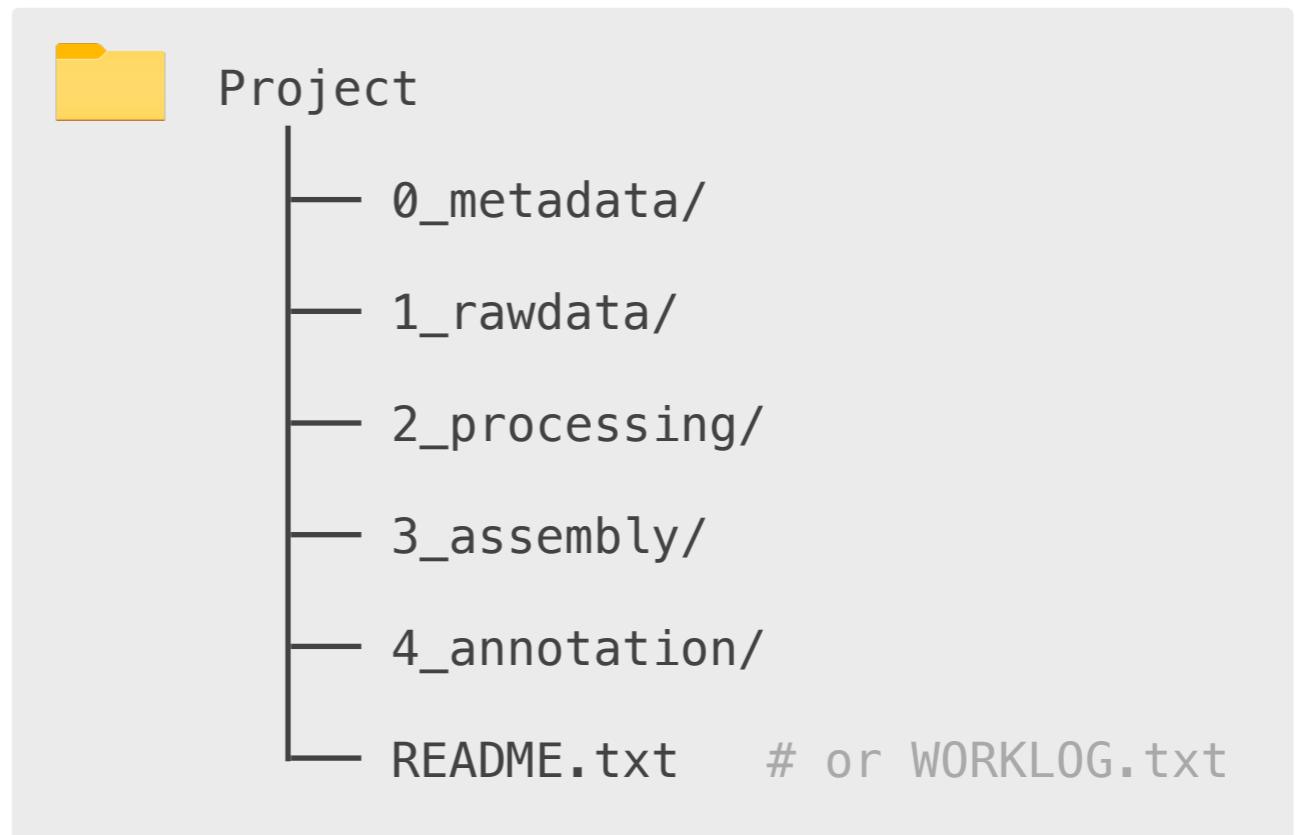
Important parameters

Data organization

Folder structure

- Be consistent
- Separate ongoing and completed work
- Consider links instead of duplicating files
- Include documentation (readme or worklog)

Example



Storing big data

- Keep raw data on dedicated hard drive
- Use file compression (e.g. [gzip](#))
- Transfer working copies to computing cluster or cloud service, if possible

Transferring big data

- Use improved copy tools (e.g. `rsync`) to avoid loss of progress or data corruption in case of connection issues
- Confirm data integrity using checksums after transfer

```
rsync -arvz <source_file> <destination_file>      # improved copy (-a: archive mode, -r:  
# include subdirectories, -z: compress  
# during transfer)  
  
md5 <file>                                         # calculate checksum (not working in Git  
# Bash)
```

Version control – git

- Tracks changes in files
- Can be used to back up and share data via GitHub
- Keep track of code, workflow / log files, result files, plots etc.
- Unsuitable to back up large data files

Version control – git

Typical git usage

```
git init                                # create git repository based on current folder  
git status                               # show repository status  
git add *                                 # track all files (exceptions in .gitignore)  
git commit -m "title"                      # record changes  
  
git push                                  # push to remote GitHub repository (if set up)
```

Example .gitignore file

```
1_rawdata/  
*.fastq.gz  
*.fastq
```

Publishing / sharing genome data

FAIR principles – **findable, accessible, interoperable, reusable**

Key attributes:

- unique and persistent identifier
- clear and detailed metadata
- standardized data format and metadata vocabulary
- accessible data usage license (e.g. open access)

Public sequence databases



Public sequence databases

NCBI Genbank (<https://www.ncbi.nlm.nih.gov/genbank/>)

The screenshot shows the NCBI Nucleotide database homepage. At the top, there's a search bar with 'Assembly' selected. A sidebar on the left lists 'Using Nucleotide' links like 'Quick Start Guide', 'FAQ', 'Help', 'GenBank FTP', and 'RefSeq FTP'. The main content area features a 'COVID-19 Information' box with links to CDC, NIH, SARS-CoV-2 data, HHS, and Spanish resources. Below it is a large image of DNA sequence data. To the right, there's a 'Nucleotide' section describing the database and its sources, along with 'Nucleotide Tools' (Submit to GenBank, LinkOut, E-Utilities, BLAST, Batch Entrez) and 'Other Resources' (GenBank Home, RefSeq Home, Gene Home, SRA Home, INSDC). A search bar at the bottom contains the term 'Scleractinia'.

Select database:
Assembly

Search term:
Scleractinia

Useful
search tags
[ACCN]
[ORGN]
[AUTH]
[TITL]

Public sequence databases

NCBI Genbank | Assembly report

The screenshot shows the NCBI Genbank assembly report for the organism *Montipora capitata* (stony corals). The assembly is named **Mcap_UHH_1.1**. Key details include:

- Organism name:** [Montipora capitata \(stony corals\)](#)
- Isolate:** Colony #1
- BioSample:** [SAMN10221787](#)
- BioProject:** [PRJNA495325](#)
- Submitter:** University of Hawaii
- Date:** 2019/07/02
- Assembly level:** Scaffold
- Genome representation:** full
- RefSeq category:** representative genome
- GenBank assembly accession:** GCA_006542545.1 (latest)
- RefSeq assembly accession:** n/a
- RefSeq assembly and GenBank assembly identical:** n/a
- WGS Project:** [RDEB01](#)
- Assembly method:** Supernova v. 2.1.1
- Expected final version:** yes
- Genome coverage:** 69.0x
- Sequencing technology:** Illumina HiSeq

IDs: 3590631 [UID] 11597468 [GenBank]

[History](#) ([Show revision history](#))

[Comment](#)

[Global statistics](#)

Total sequence length: 614,509,607

Global statistics

Total sequence length	614,509,607
Total ungapped length	571,886,739
Gaps between scaffolds	0
Number of scaffolds	27,865
Scaffold N50	185,537
Scaffold L50	747
Number of contigs	50,174
Contig N50	24,266
Contig L50	6,689
Total number of chromosomes and plasmids	0
Number of component sequences (WGS or clone)	27,865

Send to: [Download Assembly](#)

Access the data

- BLAST the assembly
- Run Primer-BLAST
- Full sequence report
- Statistics report
- FTP directory for GenBank assembly
- NCBI Datasets NEW

Assembly Information

- Assembly Help
- Assembly Basics
- NCBI Assembly Data Model

Related Information

- BioProject
- BioSample
- Genome
- PubMed
- Taxonomy

Public sequence databases

NCBI Genbank | BioSample

The screenshot shows a web browser window displaying the NCBI BioSample record for a sperm sample of *Montipora capitata*. The URL in the address bar is ncbi.nlm.nih.gov/biosample/SAMN10221787/.

Identifiers: BioSample: SAMN10221787; Sample name: Sample #1; SRA: SRS3943439

Organism: *Montipora capitata*
cellular organisms; Eukaryota; Opisthokonta; Metazoa; Eumetazoa; Cnidaria; Anthozoa; Hexacorallia; Scleractinia; Astrocoeniina; Acroporidae; Montipora

Package: MIGS: eukaryote; version 5.0

Attributes:

isolate	Colony #1
isolation source	sperm
collection date	2016-07-06
broad-scale environmental context	coral reef
local-scale environmental context	tidal pool
environmental medium	sea water
estimated size	700 Mb
geographic location	USA: Hawaii, Waiopae tide pools
isolation and growth condition	filtered sea water
latitude and longitude	19.4986 N 154.8183 W
number of replicons	2
ploidy	diploid
propagation	sexual

Description: Sperm sample from *Montipora capitata* gamete bundles released during spawning event (colony #1)
Keywords: GSC:MiS;MIGS:5.0

BioProject: PRJNA495325 *Montipora capitata* isolate:Colony #1
Retrieve [all samples](#) from this project

Related information: BioProject, SRA, Nucleotide, Assembly, Taxonomy

Recent activity: Turn Off, Clear

- Sperm sample of *Montipora capitata* biosample
- Montipora capitata* isolate:Colony #1 BioProject
- Mcap_UHH_1.1 - Genome - Assembly - NCBI Assembly
- montipora capitata AND (latest[filter] AND all[filter] NOT anomalous[filter]) (1) Assembly

[See more...](#)

Public sequence databases

NCBI Genbank | BioProject

The screenshot shows the NCBI BioProject page for the project PRJNA495325, which is a genome sequencing and assembly of *Montipora capitata*. The page includes a table of project details, a sequencing data table, and a sidebar with related information and recent activity.

Project Data:

Resource Name	Number of Links
SEQUENCE DATA	
Nucleotide (WGS master)	1
SRA Experiments	1
PUBLICATIONS	
PubMed	1
PMC	1
OTHER DATASETS	
BioSample	1
Assembly	1

Assembly details:

Assembly	Level	WGS	BioSample	Taxonomy
GCA_006542545.1	Scaffold	RDEB00000000	SAMN10221787	Montipora capitata

SRA Data Details:

Parameter	Value
Data volume, Gbases	55
Data volume, Mbytes	27719

Related information:

- Assembly
- BioSample
- Full text in PMC
- Genome
- Nucleotide
- PubMed
- SRA
- Taxonomy
- WGS master

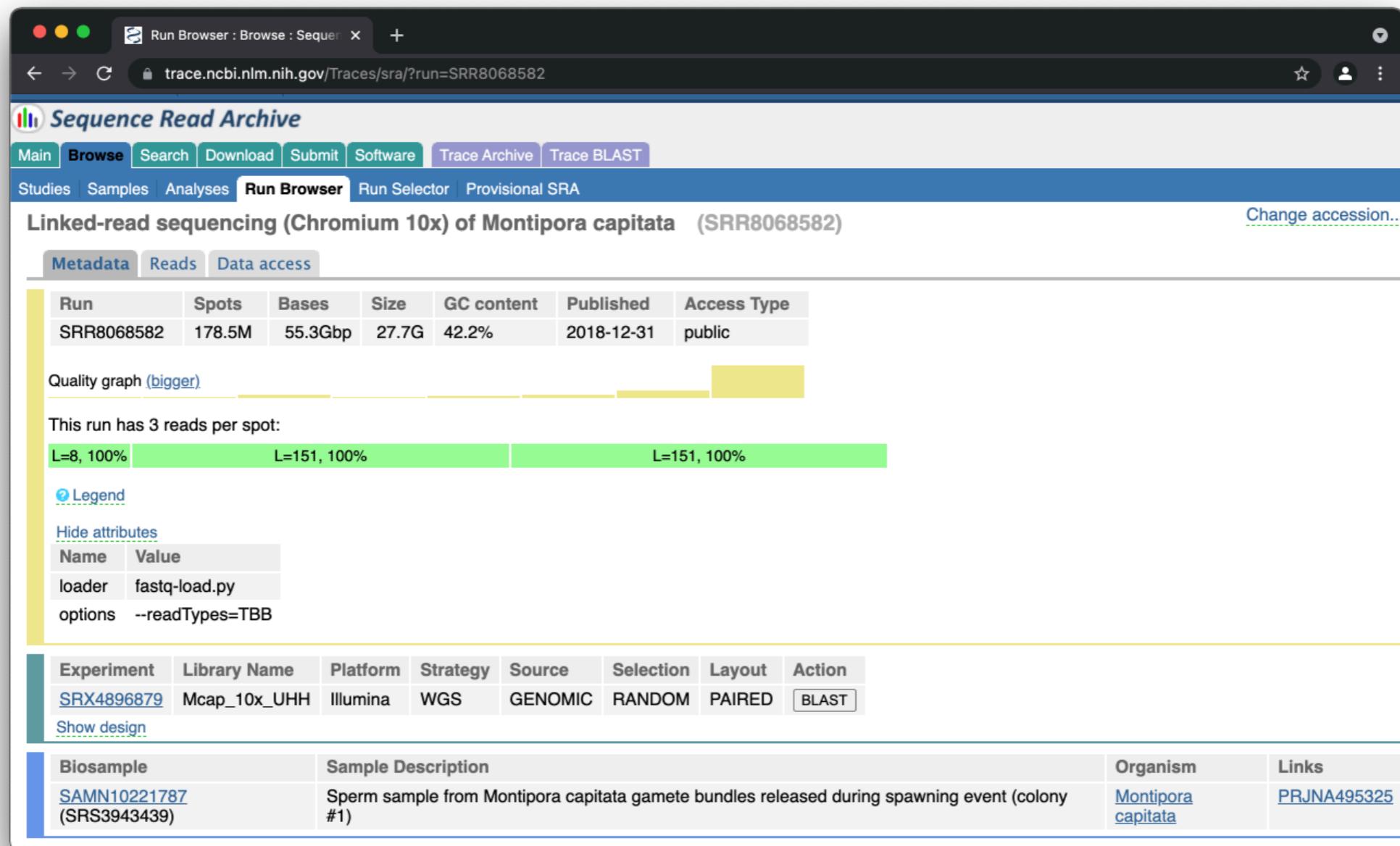
NAVIGATE ACROSS
5 additional projects are related by organism.

Recent activity:

- Montipora capitata isolate:Colony #1 (BioProject)
- Sperm sample of Montipora capitata (biosample)
- Mcap_UHH_1.1 - Genome - Assembly - NCBI (Assembly)
- montipora capitata AND (latest[filter] AND all[filter] NOT anomalous...) (1) (Assembly)

Public sequence databases

NCBI Genbank | Sequence Read Archive (SRA)



Genome submission

- NCBI genome submission guide:

<https://www.ncbi.nlm.nih.gov/genbank/genomesubmit/>

- Submission portal login (<https://submit.ncbi.nlm.nih.gov/subs/genome/>)
 - Log in with NCBI, Google, or ORCID etc. account
 - Start new or edit existing submission

Genome submission

The screenshot shows the 'Submission Portal' interface for the National Library of Medicine. The top navigation bar includes links for Home, My submissions (which is currently selected), Manage data, Groups, Templates, and My profile. The user 'mhelmkampf' is logged in. The main content area is titled 'Your submissions' and features a 'Start a new submission' section with options like GenBank, BioProject, Sequence Read Archive, BioSample, Genome, Supplementary Files, and API. Below this, a table lists 8 submissions:

Submission	Title	App	Group	Status	Updated
SUB4625258	Montipora capitata genome assembly	WGS		✓ Genomes: Processed (Details)	Jul 02 2019
SUB4632316	Montipora capitata Genome sequencing and assembly, Oct 15 '18	Sequence Read Archive (SRA)		✓ SRA: Processed SRR8068582 Download metadata file with SRA accessions View and manage my SRA submission data	Oct 17 2018
SUB4611388	Montipora capitata genome sequencing and assembly	BioProject		✓ BioProject: Processed PRJNA495325 : Montipora capitata genome sequencing and assembly (TaxID: 46704) Locus Tag Prefixes: • D8914	Oct 09 2018
SUB4611414	Sperm sample of Montipora capitata	BioSample		✓ BioSample: Processed (Details)	Oct 09 2018
SUB4139132	Corvus hawaiiensis diploid de novo genome assembly	WGS		✓ Genomes: Processed (Details)	Aug 16 2018

New submission
/ register:
[BioProject](#)
[BioSample](#)
[SRA](#)
[Genome \(WGS\)](#)

Genome submission

- Provide project description, sample metadata, technical details (method, coverage etc.), contact information, publication if applicable
- Upload assembly as single **FASTA file** (minimum contig length 200 bp)
 - Chromosome-level?
 - What do 'Ns' represent?
- Upload reads to SRA in FASTQ format
- Batch upload for multiple genomes

Genome submission

- Automated validation regarding:
 - Contamination
 - Adapter sequences
 - Duplicate contigs
- After successful review, **accession numbers** will be issued
- **Public release** at requested date or publication date, whichever comes first
- Data may run through database annotation pipelines

Dryad

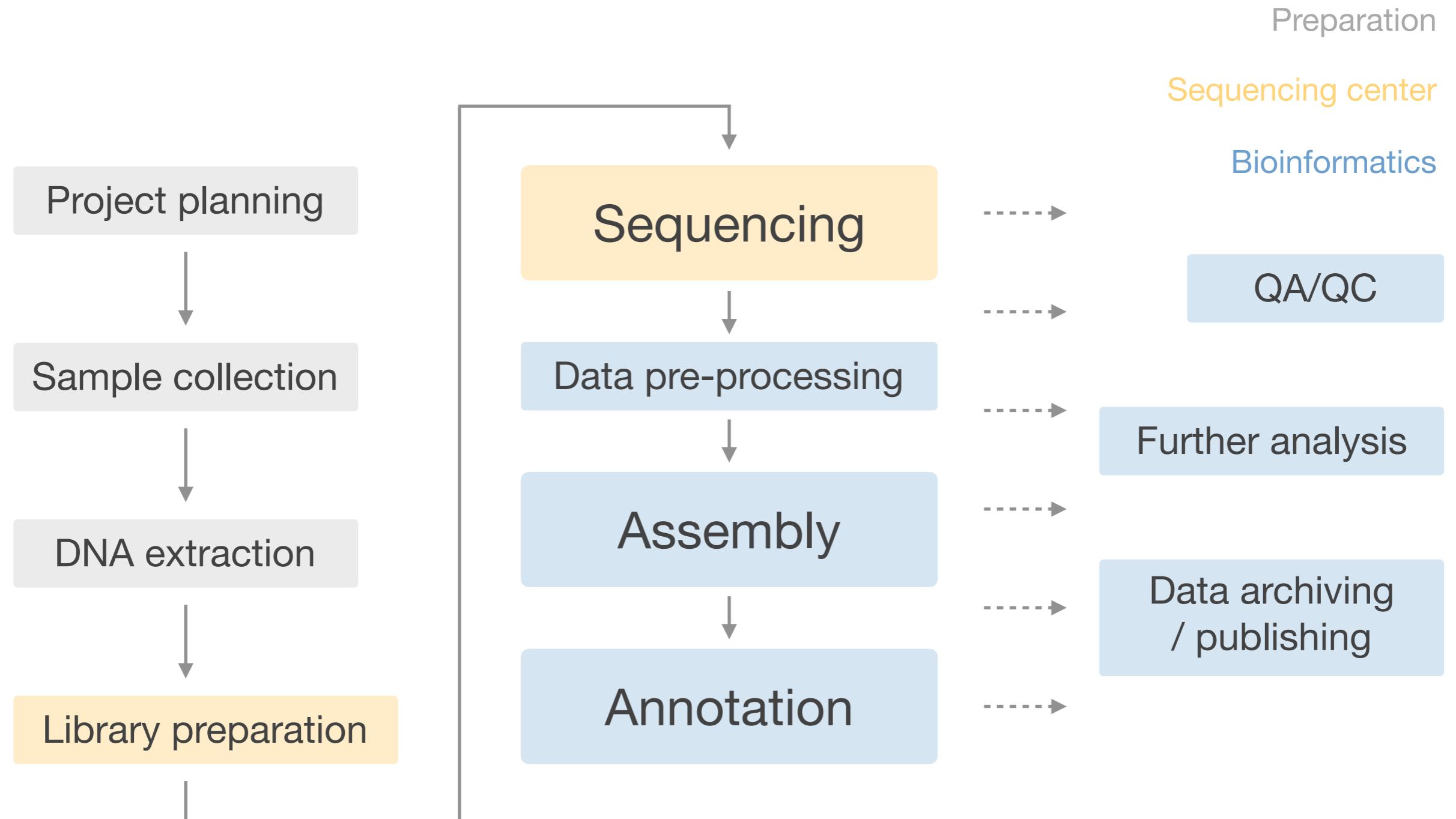
- <https://datadryad.org>
- Free open-access repository for data in any format
(e.g. GFF, official gene set FASTA, genotype data in VCF format)
- May or may not be linked to publication
- Data receive permanent DOI

Code and workflow documentation

- Key to reproducibility
- GitHub / GitHub Pages
- Example: https://k-hench.github.io/hamlet_radiation/

De novo genome sequencing workflow

Legend



Thank you!

Please leave feedback for this course here (Google form):

<https://tinyurl.com/9hucxfvh>

Questions?

Email me at martin.helmkampf@leibniz-zmt.de

