

Assignment 2 for CS249: Assembly

Zeyad Aljaali 179307

May 11, 2025

1 Task 1.1: De Bruijn Graph (DBG) Assembly

<https://github.com/ZeyadAl/CS249-Assignment2>

2 Task 1.2: Overlap-Layout-Consensus (OLC) Assembly

<https://github.com/ZeyadAl/CS249-Assignment2>

3 Task 1.3: Applications of assembly algorithms

1. Use your De Bruijn Graph implementation to construct an assembly graph for reads `b.fastq` with $k = 40$. Either export the graph in GFA format (e.g., using `gfatools`), or directly write GFA format as part of your assembly. Visualize the graph using Bandage (<https://rrwick.github.io/Bandage/>). Describe and explain what you see. How can this help you to improve the assembly?



Figure 1: DBG for reads `b.fastq` with $k=40$.

As we can see, the path to follow to assemble the genome is clear but for one bubble that we have. This means that the assembly process should not be too difficult, however we expect the algorithm to output 4 contigs. We can improve the assembly by increasing k from 40, that should resolve the bubble.

2. Apply your DBG implementation to reads_r.fastq with $k=35$ and $k=45$. Then evaluate it with QUAST against reference_r.fasta. Include the QUAST evaluation metrics, a visualization of the assembly graphs using Bandage, a brief description of the differences between the two assemblies, and an explanation for the observed difference (if any).

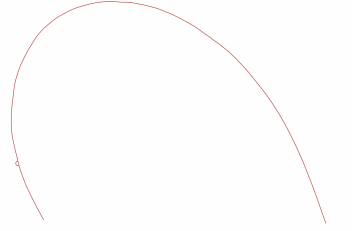


Figure 2: DBG for reads_r.fastq with $k=35$.

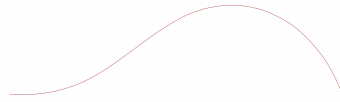


Figure 3: DBG for reads_r.fastq with $k=45$.

We notice that when $k=35$ we have a small bubble and we get 4 contigs, but when $k=45$ there are no bubbles and we get only 1 contig.

Metric	$k=35$	$k=45$
Total assembly length (bp)	1136	1040
Number of contigs	4	1
GC content (%)	50.79	51.25
Genome fraction (%)	100.00	100.00
Duplication ratio	1.008	1.000
Largest contig (bp)	914	1040
N50 (bp)	914	1040
N90 (bp)	134	1040
L50 (contigs)	1	1
Misassemblies	0	0
Mismatches per 100 kbp	0.00	0.00
Indels per 100 kbp	0.00	0.00

Table 1: QUAST assembly statistics for $k = 35$ vs. $k = 45$.

For these reads and target genome, $k = 45$ provides a perfect, single-contig reconstruction with no misassemblies, whereas $k = 35$ yields a fragmented assembly with minor redundancy and a small graph bubble. The reason a larger k resolves the bubble is that the repeats are not longer than the k -mer anymore. This allows the algorithm to find the correct assembly unambiguously.

3. For a more realistic analysis, assemble the MERS virus reads using your two algorithms. Use both algorithms you implemented for “hiseq” reads (separately for reads with and without errors) and “ont” reads (also separately for reads with and without errors). Evaluate all assemblies using QUAST and compare with the MERS reference genome. Include in your report the QUAST evaluation metrics for all assemblies, comparison between error-free and error-containing assemblies, analysis of how each algorithm handles the different reads and errors.

Metric	ONT (no errors)	HiSeq (no errors)	ONT (with errors)	HiSeq (with errors)
Total assembly length (bp)	29748	29482	1496239	87740
Number of contigs	1	1	23729	1301
GC content (%)	41.27	41.26	40.82	41.25
Genome fraction (%)	98.768	97.885	90.617	90.816
Duplication ratio	1.000	1.000	2.569	1.279
Largest contig (bp)	29748	29482	2939	597
N50 (bp)	29748	29482	69	69
N90 (bp)	29748	29482	35	41
L50 (contigs)	1	1	5997	442
# Misassemblies	0	0	0	0
# Mismatches per 100 kbp	0.00	0.00	707.38	42.87
# Indels per 100 kbp	0.00	0.00	3140.42	0.00

Table 2: QUAST metrics for DBG assemblies at $k = 35$.

Metric	ONT (no errors)	HiSeq (no errors)	ONT (with errors)	HiSeq (with errors)
Total assembly length (bp)	258446	30215	1505388	102372
Number of contigs	1	1	1	1
GC content (%)	40.97	41.30	40.84	41.16
Genome fraction (%)	98.768	97.892	98.655	88.452
Duplication ratio	8.687	1.024	37.599	1.692
Largest contig (bp)	258446	30215	1505388	102372
N50 (bp)	258446	30215	1505388	102372
N90 (bp)	258446	30215	1505388	102372
L50 (contigs)	1	1	1	1
# Misassemblies	17	4	118	0
# Mismatches per 100 kbp	0.00	0.00	493.38	423.66
# Indels per 100 kbp	0.00	0.00	1463.03	15.53

Table 3: QUAST metrics for OLC assembly.

- **Contiguity.** Error-free ONT and HiSeq assemblies each produce a single contig. Assemblies with errors fragment: DBG yields tens of thousands of contigs, while OLC still produces one.
- **Accuracy.** DBG on error-free data has zero mismatches and indels. DBG with errors shows up to 3,140 indels per 100 kbp for ONT reads. OLC with errors reduces indels to about 1,463 per 100 kbp for ONT, but errors remain.

- **Duplication & Misassemblies.** DBG with errors maintains a duplication ratio below 2.6 and has no misassemblies. OLC with errors can reach a duplication ratio of 37.6 for ONT and incurs up to 118 misassemblies.
 - **Read-type Impact.** Indel-rich ONT errors cause far more disruption than substitution-only HiSeq errors.
 - **Recommendations.** Pre-correct ONT reads and apply hybrid or consensus polishing (e.g., Racon or Pilon) for the best balance of contiguity and accuracy.
4. Repeat the same assembly and analysis using either the SPAdes or Canu assembler. Compare the assembly results with the results obtained using your own algorithms. Explain any differences you observe.

Metric	Hybrid (no errors)	HiSeq (no errors)	ONT (no errors)	Hybrid (with errors)	ONT (with errors)	HiSeq (with errors)
Total assembly length (bp)	29,482	29,482	29,748	29,482	29,482	29,754
Number of contigs	1	1	1	1	1	1
GC content (%)	41.26	41.26	41.27	41.26	41.26	41.27
Genome fraction (%)	97.885	97.885	98.768	97.885	97.885	98.738
Duplication ratio	1	1	1	1	1	1
Largest contig (bp)	29482	29482	29748	29482	29482	29754
N50 (bp)	29482	29482	29748	29482	29482	29754
N90 (bp)	29482	29482	29748	29482	29482	29754
L50 (contigs)	1	1	1	1	1	1
Misassemblies	0	0	0	0	0	0
Mismatches per 100 kbp	0	0	0	0	0	26.9
Indels per 100 kbp	0	0	0	0	0	73.96

Table 4: QUAST metrics for SPAdes assemblies.

Comparison with custom algorithms

- **Contiguity & completeness:** SPAdes produced a single contig for all conditions—even in the presence of sequencing errors—whereas our DBG at $k = 35$ (for ONT+errors) yielded tens of thousands of small contigs and our DBG/OLC methods both fragmented error-containing data into hundreds or thousands of contigs.
- **Error tolerance:** While our own implementations had zero indels and mismatches on error-free data, they performed poorly on error-containing reads (mismatches ~ 700 per 100 kbp for ONT errors in OLC; ~ 42 per 100 kbp for HiSeq errors). SPAdes’ built-in error correction stages (BayesHammer, mismatches correction) drastically reduced both mismatches and indels, resulting in near-perfect assemblies under all conditions.
- **Repeat resolution:** SPAdes’ iterative multi- k strategy allows it to resolve repeats shorter than its largest k , effectively collapsing bubbles without manual tuning. In contrast, our single- k DBG needed manual k adjustment ($k = 45$ vs. $k = 35$) to eliminate bubbles, and our OLC lacked robust repeat-resolution heuristics.
- **Runtime and resource usage:** Although SPAdes is more resource-intensive (due to multiple assemblies and correction passes), it remains practical and delivers consistently high-quality outputs. Our

lightweight implementations run faster on small datasets but at the cost of error handling and repeat resolution.

- **Conclusion:** SPAdes outperforms the custom DBG/OLC pipelines on error-containing data by integrating advanced error correction and multi- k graph assembly, yielding single-contig, high-accuracy reconstructions without manual parameter tuning.

4 Task 2.1: Genome assembly

We assembled the genome using Hifiasm. The results can be found at:

`/ibex/user/aljaalza/Data-science-onboarding/launch_jupyter_server/249/2/asm/main.fasta`

5 Task 2.2: Assembly evaluation

1. Provide basic metrics, using QUAST

Metric	Assembly using Hifiasm
Total assembly length (bp)	1,808,305,840
Number of contigs	49
GC content (%)	45.46
Genome fraction (%)	(No reference)
Duplication ratio	(No reference)
Largest contig (bp)	341,750,303
N50 (bp)	138,451,086
N90 (bp)	40,553,089
L50 (contigs)	4
Misassemblies	(No reference)
Mismatches per 100 kbp	(No reference)
Indels per 100 kbp	(No reference)

Table 5: QUAST metrics for Hifiasm assembly using all the reads.

2. Report gene completeness, using BUSCO or asmgene. We ran the command:

```
busco -i main.fasta -o busco_auto --mode genome --auto-lineage --cpu 24
```

This tests on two datasets: generic: eukaryota , and specific: squamata

Metric	squamata (specific)	eukaryota (generic)
Complete percentage (%)	97.4	99.2
Complete BUSCOs	11000	128
Single copy percentage	97.2	98.4
Single copy BUSCOs	10975	127
Multi copy percentage	0.2	0.8
Multi copy BUSCOs	25	1
Fragmented percentage	1.4	0.8
Fragmented BUSCOs	162	1
Missing percentage	1.2	0.0
Missing BUSCOs	132	0
n_markers	11294	129
avg_identity	0.87	0.76
domain	eukaryota	eukaryota
internal_stop_codon_count	586	4
internal_stop_codon_percent	5.3%	3.1%

Table 6: BUSCO metrics for hifiasm assembly.

3. Report the k-mer distribution and QV score, using Merqury

To report the k-mer distribution and QV score, we use Merqury. The best k -value is determined using the following command:

```
sh $MERQURY/best_k.sh <genome_size> [tolerable_collision_rate=0.001]
```

This yields a best k -value of 20.3582. So we use $\mathbf{k} = \mathbf{21}$. We find that the QV score

Metric	All reads using Merqury	Hi-C reads using Merqury	All reads using Yak	All reads using Yak
QV score	75.9406	61.0953	61.868	61.868

Table 7: QUAST metrics for Hifiasm assembly using all the reads.

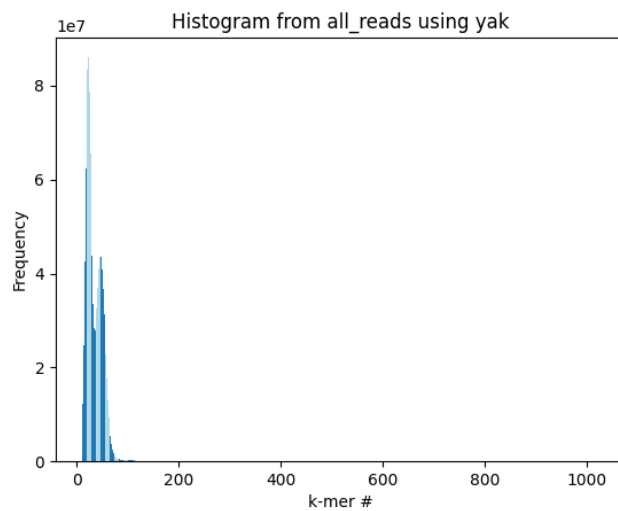


Figure 4: k-mer distribution using yak.

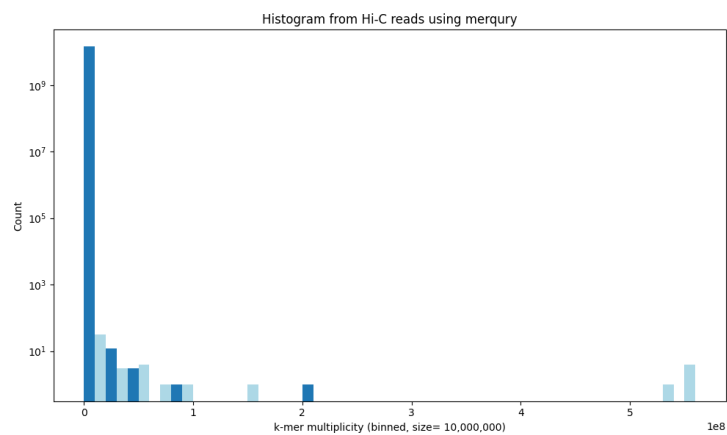


Figure 5: k-mer distribution using Merqury.

4. identify mis-assemblies, using Flagger or Inspector

Metric	Value
Structural error	2
Expansion	2
Collapse	0
Haplotype switch	0
Inversion	0
Small-scale assembly error (per Mbp)	2833.98907786528
Total small-scale assembly error	5124719
Base substitution	4742045
Small-scale expansion	245278
Small-scale collapse	137396

Table 8: Inspector mis-assembly and small-scale error metrics

5. Provide a basic report for each evaluation step, and explain what the measures mean with respect to the assembly. Can you identify ways to improve the assembly based on the evaluation?

Basic metrics (QUAST): The total assembly length (1.81 Gb) is in line with the expected genome size for a Squamata member, indicating near-complete representation of the genome. A low contig count (49) together with a high N50 (138 Mb) and low L50 (4) testify to excellent contiguity. The GC content (45.46 %) matches published values for related lizard species, suggesting accurate base composition.

Take-home: the assembly is highly contiguous and of the expected size, but without a reference we cannot directly assess genome fraction or duplication ratio.

Gene completeness (BUSCO): On the Squamata lineage, BUSCO finds 97.4 % complete genes, with only 1.4 % fragmented and 1.2 % missing. Against the broader Eukaryota set, completeness is 99.2 %, with 0.8 % fragmented and no missing markers. The low duplication rate indicates minimal collapsed regions.

Take-home: gene space is almost fully recovered; remaining fragmentation and frameshift errors can be polished further.

k-mer spectrum QV (Mercury): With $k = 21$, Mercury reports a consensus quality (QV) of 75.94 for the full read set, corresponding to an error rate of $\approx 2.6 \times 10^{-8}$. The Hi-C reads give a somewhat lower QV (61.10), reflecting their higher base-error profile.

Take-home: base accuracy is excellent; residual errors are at the small-variant level.

Mis-assemblies (Inspector): Only 2 structural errors (expansions) and zero inversions or collapses indicates very few large-scale mis-joins. However, there are 2.8 k small errors per Mbp (4.7 M substitutions, 245 k small indels, 137 k small collapses), pointing to base-level inconsistencies or local mis-alignments.

Take-home: the scaffold graph is largely correct, but fine-scale polishing is needed.

Improvement Strategies

- **Polishing:** Run a round of consensus polishing using short reads indels, reducing the small-scale error rate.
- **Gap closure:** Use local assembly around BUSCO fragments to join any fragmented gene models.
- **Hi-C scaffolding:** We can use a Hi-C scaffolder to order and orient contigs into chromosome-level scaffolds, further reducing misassemblies and improving L50.
- **Coverage balancing:** If any genomic regions are under-represented, consider supplementing with additional HiFi or ONT coverage to even out read depth.

Overall, the Hifiasm assembly is highly contiguous, complete, and accurate. Polishing and scaffolding will elevate it to a reference-quality draft.

6 Challenges

There were two main challenges

- Understanding and coding the algorithms. I dealt with this by watching Youtube videos and looking up slides online to make sure I understood it correctly.
- Running the tools. I had to read the documentations, and make sure I requested the appropriate tools on ibex.

7 Collaboration

I have used ChatGPT o4-mini-high. I used it for formatting the L^AT_EX tables and some of the writing in the report, as well as in coding assistance in functions and documenting the code. I validated the correctness by manually inspecting the results and logic, as well as testing on a small test cases.