

1 PAQman: reference-free ensemble evaluation of long-read 2 genome assemblies

3
4 Samuel O'Donnell^{1*}, Ningxiao Li^{2,4}, Jacob L. Steenwyk³, David M. Geiser⁴, Frank N. Martin²,
5 Emile Gluck-Thaler^{1,5*}

6
7 1: Department of Plant Pathology, University of Wisconsin-Madison, Madison, WI 53706, USA

8 2: Crop Improvement and Protection Research Unit, United States Department of Agriculture – Agricultural
9 Research Service, Salinas, CA 93905, USA

10 3: Howard Hughes Medical Institute and the Department of Molecular and Cell Biology, University of
11 California, Berkeley, CA, USA

12 4: Department of Plant Pathology and Environmental Microbiology, The Pennsylvania State University,
13 University Park, PA 16802, USA

14 5: Wisconsin Institute for Discovery, Madison, WI 53706, USA.

15
16 *Co-corresponding authors. Correspondence should be addressed to:
17 samuel.a.odonnell@gmail.com and gluckthaler@wisc.edu

18 **ORCID IDs:**

19 S O'Donnell: 0000-0003-2447-1993

20 N. Li: 0000-0002-7847-3567

21 J.L. Steenwyk: 0000-0002-8436-595X

22 D.M. Geiser: 0000-0002-1590-2045

23 F.N. Martin: 0000-0002-8050-1248

24 E. Gluck-Thaler: 0000-0003-0438-7495

25 **Running title:**

26
27 PAQman: quality assessment of long-read genomes

1
2 **Keywords:**
3 Assembly quality; Assembly evaluation; Oxford Nanopore Technologies, Pacific Biosciences
4

5 Abstract

6
7 Advances in long-read sequencing have made it easier and more cost effective to generate
8 high-quality genome assemblies. However, assessing assembly quality remains a challenge,
9 as existing tools often focus on a few metrics and/or require a reference assembly for
10 comparison. Furthermore, the number of available metrics and associated tools for genome
11 evaluation have expanded in recent years, making it more difficult for researchers to easily
12 use and develop comprehensive pipelines. To address this, we developed the Post-
13 Assembly Quality manager (PAQman), a tool that lowers the barrier to entry for assembly
14 quality assessment by measuring seven reference-free features of genome quality within a
15 single framework: Contiguity, Gene content, Completeness, Accuracy, Correctness,
16 Coverage, and Telomerality. PAQman integrates multiple commonly used tools alongside
17 custom scripts, requiring users to provide only a query genome assembly and its underlying
18 long-read data, while providing a streamlined and consistent framework for quality
19 assessment across datasets.

20 Introduction

21 The ease with which researchers can generate reference-quality assemblies has advanced
22 rapidly with the continued improvement of long-read sequencing technologies, namely
23 those offered by Oxford Nanopore Technologies (ONT) and Pacific Biosciences. In parallel,
24 assembly algorithms and pipelines have significantly improved in both accuracy and speed
25 (Li 2018; Kolmogorov et al. 2019; Giani et al. 2020; Espinosa et al. 2024). Despite these
26 advancements, assembly quality is often assessed by disparate software that individually
27 use a subset of available metrics, leading to inefficiencies in throughput and evaluation.
28 These different statistics often evaluate distinct genomic features and thus do not always
29 correlate with one another (Zhang et al. 2023). Determining the most appropriate metrics for
30 quality assessment is therefore challenging, especially when benchmarking software or
31 parameter settings to optimize assembly quality for a particular dataset or organism
32 (Bradnam et al. 2013; Zhang et al. 2022; Cosma et al. 2023). These difficulties are further

1 compounded in the absence of a reference genome, where defining the "best" assembly
2 becomes even more ambiguous.

3
4 Several tools have been developed to address the challenge of assembly quality
5 assessment, such as Genome Assembly Evaluation Process (GAEP) (Zhang et al. 2023),
6 GenomeQC (Manchanda et al. 2020) and the most commonly used QUAST pipeline
7 (Mikheenko et al. 2023). However, most are either sparsely documented, challenging to
8 install and/or use, or evaluate only a subset of relevant statistics. For example, GenomeQC,
9 if used without a reference assembly or gene annotations, only calculates a few assembly
10 stats such as contig N50 and gene content using BUSCO. Similarly, QUAST relies on
11 comparisons with a reference assembly to compute several statistics, such as structural
12 accuracy and completeness. Although GAEP evaluates a number of important statistics
13 such as contiguity, gene content, accuracy and structural errors, it does not assess
14 coverage, completeness or the presence of telomeric sequences, and is not packaged for
15 use as a software container or conda package for use in high-throughput computing
16 environments such as those commonly used for bioinformatics analyses.

17
18 To address these challenges, we developed the Post-Assembly Quality manager (PAQman),
19 a tool that calculates seven reference-free features of genome quality—Contiguity, Gene
20 content, Completeness, Accuracy, Correctness, Coverage and Telomericity—by integrating
21 multiple existing and popular tools and custom scripts. PAQman requires only the assembly
22 itself and the long-read data used to generate it, while offering users an accessible and
23 standardized approach for genome quality assessment in high-throughput computing
24 environments. Although PAQman was designed with eukaryotes in mind, there is no
25 limitation in organism type aside from the telomericity feature being designed for linear
26 chromosomes found primarily in eukaryotes.

27 Methods

28 The workflow summarized in Fig. 1 describes the two commands used in PAQman. The first,
29 *paqman*, performs the separate evaluation of seven features of assembly quality prior to
30 compiling the summary statistics table. By default, the minimum requirements for this
31 command are an assembly in FASTA format (-a) and long-reads in FASTQ format (-l).
32 Additional short paired-end reads may also be provided (-1 -2; FASTQ). Each feature and the
33 tool/s used for its evaluation are outlined below:

34

1 Feature 1: Contiguity. By far and above, the most common metrics evaluated post-assembly
2 are the number of contigs/scaffolds, the assembly size and the contig/scaffold N50.
3 PAQman runs QUAST (Mikheenko et al. 2023) to rapidly calculate these important metrics.
4 Generally speaking, a lower number of scaffolds and larger assembly size and contig N50
5 are associated with higher quality assemblies. Additional analyses from QUAST, such as GC
6 content and cumulative length plots, are kept in the output folder.

7
8 Feature 2: Gene content. Apart from contiguity, the most common tool used to evaluate
9 assemblies is BUSCO, which quantifies the recovery of various databases of universally
10 conserved single copy orthologs (Tegenfeldt et al. 2025). Using the presence/absence of
11 near-universal single-copy orthologs, BUSCO gives a proxy for assembly completeness in
12 terms of gene content. PAQman parameter `-b` allows users to provide the most appropriate
13 database for BUSCO to use given the taxonomic lineage of the focal organism (more
14 database information at https://busco.ezlab.org/busco_userguide.html#obtain-busco). The
15 percent recovery of single-copy orthologs (ranging from 0-100%) is positively associated with
16 assembly quality. Users can also provide a predownloaded or locally generated BUSCO
17 database and provide it to PAQman using `--localbuscodb`.

18
19 Feature 3: Completeness. Genomic regions containing a high density of protein-coding
20 genes are generally the easiest to assemble due to relatively low repetitiveness compared to
21 non-genic regions. Therefore, in addition to gene content, assembly completeness may also
22 be more globally evaluated in terms of whether k -mers present in the raw read inputs are
23 also present in the assembly. The k -mer-based completeness statistic is thus an accurate
24 representation of how much of the raw sequencing data is present in a final assembly,
25 including both repetitive and non-repetitive content. This k -mer-based measure of
26 completeness is therefore distinct from other uses of the term that are frequently used to
27 describe assembly quality, such as the gene content-based completeness measurement
28 implemented in BUSCO. PAQman first builds a k -mer distribution of the reads, selecting the
29 short-reads (`-1 -2`) if provided, using `meryl count` (Miller et al. 2008) then provides this `meryl`
30 database and the assembly to `Merqury` (Rhie et al. 2020) to calculate the k -mer based
31 completeness. Completeness is calculated as a percentage with a maximum of 100%.
32 Users may optionally provide a precomputed `meryl` database using `--meryldb`. A limitation
33 of this unbiased evaluation of raw reads is that the resulting k -mer based metrics can be
34 affected by contamination, mixed samples, or complex ploidy, complicating the
35 interpretation of completeness and QV estimates.

1
2 Feature 4: Accuracy. Using the k -mer distributions calculated for Feature 3, Merqury also
3 estimates the number of genome-wide sequencing errors and provides a Phred quality score
4 ($=-10 \times \log_{10}(Pe)$; Pe : estimated probability of error) for the assembly. This estimate indicates
5 how many errors may still be present in an assembly and may therefore be used to test
6 whether common assembly polishing tools/pipelines improve quality. In general,
7 assemblies with a Phred Quality greater than 30 are considered of high quality, although 45+
8 is common with more recent long-read chemistries and polishing methods. In rare cases
9 where no errors are detected, we used the rule of three to calculate a conservative Phred
10 score estimate, substituting the estimated probability of error by $3/n$, where n is the genome
11 size.

12
13 Feature 5: Correctness. As genome assemblies approach reference quality, assessment
14 must go beyond simple presence or absence of genomic content to include whether that
15 content is assembled in the correct genomic location and structure. CRAQ (Li et al. 2023) is
16 a tool that uses read mapping evidence to highlight regions within an assembly that have
17 potential assembly errors. This provides an estimate for genome-wide structural accuracy,
18 thereby providing evidence to support claims of assembly contiguity. Initially, PAQman
19 randomly down-samples the provided long-read sequencing data to a maximum of 30X
20 coverage by default (set by the `--maxcoverage` parameter) using Rasusa (Hall 2022), and then
21 provides the filtered read dataset to CRAQ in addition to short-reads, if provided. PAQman
22 then focuses on two versions of the Assembly Quality Index (AQI) summary statistic; the R-
23 AQI and S-AQI, that score the assembly quality based on the detection of small regional (R)
24 and large structural (S) errors respectively. The AQI scores are bound between 0-100, with
25 100 meaning the absence of detectable errors. Please refer to Li et al., 2023 for a complete
26 description on how both R- and S-AQI are calculated.

27
28 Feature 6: Coverage. Using read mapping, PAQman also evaluates whether coverage varies
29 in particular regions compared to the genome-wide median. One important use of this
30 evaluation metric is to determine whether multi-copy repeat regions have collapsed into a
31 single copy during assembly, as coverage will noticeably increase in collapsed regions.
32 Coverage may therefore help in understanding and locating more complex rearrangements
33 such as large duplications and aneuploidies. PAQman uses `bwa mem` (Li 2013) and/or
34 `minimap2` (Li 2018) for short- and long-read mapping respectively; `mosdepth` (Pedersen and
35 Quinlan 2018) and `bedtools` (Quinlan and Hall 2010) for pre- and post-processing; and

1 custom Rscripts with ggplot2 (Wickham 2016) for plotting the relative coverage normalised
2 by the genome-wide median (Fig. 1). PAQman uses the randomly downsampled long-reads
3 for mapping to improve speed. The parameters *-w* and *-s* control the size of the window and
4 slide, respectively, used to average coverage across the genome. For a final statistic,
5 PAQman calculates the percentage of the genome that is within two standard deviations of
6 the genome-wide median coverage. Ideally, larger percentage values are more desirable.

7
8 **Feature 7: Telomerality.** This term defines a set of statistics that indicate whether an
9 assembly contains assembled telomeric ends, which is a good indication that complex
10 subtelomeric regions have been assembled and that contigs may thus represent full
11 telomere-to-telomere (T2T) chromosomes. In most cases, this statistic will be used to
12 elevate assemblies ranging from very good quality to complete. This feature can be ignored
13 in organisms with circular DNA molecules. Organism-specific telomeric repeats are
14 provided through the *-r* parameter, quickly identified using seqkit *locate* (Shen et al. 2016)
15 and processed using bedtools *merge* for overlapping repeats with a max of one missing
16 repeat between. Repeats are considered telomeric if the aggregated repeat regions are at
17 least the size of two repeats. Contig/scaffold ends are considered to be capped by telomeres
18 if the distance of a repeat to an end is less than 75% of the size of the entire telomeric repeat
19 region. We recommend using TeloBase (Lyčka et al. 2024) in order to find the likely telomeric
20 sequence for your species, however this database does not allow for degenerate repeats. To
21 help with this, we are compiling a list
22 (https://github.com/SAMtoBAM/PAQman/blob/main/telomeric_sequences.md) of working
23 telomeric repeat patterns (both exact and regular expressions) found to work in manually
24 verified assemblies.

25
26 The second optional command in PAQman, *paqplot*, takes the summary statistics output by
27 *paqplot* from multiple assemblies and generates radar and lollipop plots to help benchmark
28 multi-assembly comparisons using both the raw and relative values (Fig. 2).

29 Results and Discussion

30 We tested PAQman on three different genome datasets with diverse genome characteristics:
31 *S. cerevisiae* (S288c), *Oryza sativa subsp. japonica* (Nipponbare) and Humans (CHM13);
32 encompassing a range of eukaryotic lineages (fungi, plants, and animals), genome sizes
33 (~12Mbp, ~380Mbp, and ~3Gbp), repeat content (~10%, ~40%, and ~55%) and ploidy
34 (haploid, diploid, and effectively haploid). Notably, each dataset contains multiple

1 independent and/or sequentially improved versions of the strain/cultivar/line with an
2 ultimately T2T version. For each dataset (S288c, Nipponbare and CHM13), a complete,
3 reproducible and detailed walkthrough can be found here
4 https://github.com/SAMtoBAM/PAQman/tree/main/pub_examples/.

5
6 For the *Saccharomyces cerevisiae* reference strain S288c, we tested PAQman on a dataset
7 of five publicly available assemblies (GCA_000146045.2, GCA_902192305.1;
8 GCA_022626425.2; GCA_002057635.1; GCA_016858165.1). Each was evaluated using the
9 same ONT dataset (SRR17374240) (Zhang et al. 2022) downsampled from 800x to 100x
10 coverage using Rasusa (*-b 1200000000*). *paqman.sh* was run with default options except '*-b*
11 *saccharomycetaceae*' and '*-r GGTGTG*' for the BUSCO db and telomeric repeat, respectively.
12 All summary stats were then evaluated by *paqplots.sh* using default settings (Fig. 2A-B).
13 PAQman shows that although assembly GCA022626425 (purple) contains good metrics for
14 the two most commonly evaluated features, Contiguity (fewest contigs and high contig N50)
15 and Gene Content (high single copy BUSCO content), it contains poorer metrics for other
16 less commonly used features such as Accuracy and Correctness. The lower AQI values in
17 GCA022626425, indicative of assembly errors, were confirmed by whole genome alignment
18 (data not shown). Similarly, GCA000146045 (teal), the canonical reference assembly,
19 contains good relative metrics for all features except Telomericity, with only 21 contig ends
20 identified as containing telomeric repeats, the same as previously detected (O'Donnell et al.
21 2023). These two examples highlight how each Feature is an independent dimension and
22 that only a comprehensive evaluation of all features allow us to accurately measure quality.

23
24 For the *Oryza sativa subsp. japonica* Nipponbare cultivar, we ran PAQman on five publicly
25 available assemblies (GCA_003865235.1, GCA_051403585.1, GCA_000005425.2,
26 GCA_001433935.1, GCA_034140825.1) which includes the previous reference assembly
27 (GCA_001433935.1) and the new gold standard T2T assembly (GCA_034140825.1). We
28 analysed all assemblies using PacBio HiFi (SRR25241090) reads (~85X) from the T2T project
29 (Shang et al. 2023). *paqman.sh* was run with default options except '*-b poales*' and '*-r*
30 *TTTAGGG*' for the BUSCO db and telomeric repeat respectively. All summary stats were then
31 evaluated by *paqplots.sh* using default settings (Fig. 2C-D). The results show that the
32 assemblies GCA051403585 (teal) and GCA034140825 (purple) have the best set of metrics
33 for all Features, including Accuracy ('qv-phred'). These two assemblies only differentiate
34 largely due to Telomericity, as the reportedly T2T assembly (GCA034140825; purple) indeed
35 contains a full set of 12 chromosomes all capped by telomeres ('T2T contigs') (Shang et al.
36 2023). This difference highlights the benefit of Telomericity as a feature of assembly quality.

1
2 For the human cell line CHM13, we ran PAQMan on five publicly available assemblies (T2T-
3 CHM13v0.7, T2T-CHM13v1.0, T2T-CHM13v1.1, T2T-CHM13v2.0 and GCA_002884485.1).
4 This includes 4 versions of the now telomere-2-telomere, complete assembly (T2T-
5 CHM13v2.0) and one much earlier assembly for the same cell line (GCA_002884485.1). We
6 also used PacBio HiFi (SRX5633451) reads from the T2T-CHM13 project
7 (<https://github.com/marbl/CHM13>) with approximately 24X coverage in total. *paqman.sh*
8 was run with default options except '*-b tetrapoda*' for the BUSCO db and *--coveragemax 0* to
9 skip downsampling due to low input coverage. All summary stats were then evaluated by
10 *paqplots.sh* using default settings (Fig. 2E-F). As expected we see that the Contiguity and
11 Telomerality features improve with subsequent versions of the T2T CHM13 versions and that
12 v0.7 already drastically improved upon the previous public accession (GCA002884485). This
13 once again shows the utility of Telomerality as an assembly quality Feature. Notably
14 CHM13v2 shows a decreased QV score compared to previous version, however this is only
15 due to the addition of the Y chromosome which comes from the HG002 cell line as was
16 combined with CHM13 (Rhie et al. 2023). Without the Y chromosome, the average QV value
17 of each chromosome was 75, comparable to the CHM13v1.1 calculated QV of 73.

18
19 For each run of *paqman.sh* we used varying numbers of threads (*--threads*) and calculated
20 the run time and maximum RAM usage using *time* and *sacct* (*--format=MaxRSS*) respectively
21 (Table 1; Table S1). For S288c using 16 threads, PAQman required on average 6 minutes per
22 assembly to complete and used a maximum of 5.5GB RAM. For Nipponbare, with a 31 times
23 larger genome (~12Mbp vs ~380Mbp) PAQman required on average 93 minutes and a
24 maximum of 48GB of RAM. For the human dataset CHM13 (~3Gbp), we increased the thread
25 count to 64 and saw that PAQman required on average 424 minutes and a maximum of 219
26 GB of RAM. In all datasets tested, speed improved with additional threads with minimal
27 changes to max RAM usage. Using the Nipponbare assembly GCA000005425 and 16
28 threads, we isolated each major step in PAQman and saw that 65% of time and the highest
29 RAM usage was during the steps containing read alignment with minimap2 and Correctness
30 estimation with CRAQ (Table S2).

31 Installation, Usage and Citation

32 PAQman was built for Linux operating systems. The easiest way to install PAQman is using
33 the conda package manager:

34

1 `conda config --append channels conda-forge`

2 `conda config --append channels bioconda`

3 `conda config --append channels pwwang`

4 `conda install samtobam::paqman`

5

6 It may also be run as an Apptainer/Singularity image:

7

8 `docker pull ghcr.io/samtobam/paqman:latest`

9

10 Both PAQman commands are run as below with default options (additional arguments
11 available for paqman.sh):

12

13 `paqman.sh -a path/to/assembly.fa -l path/to/longreads.fq`

14 `paqplots.sh -s path/to/combined_summary_stats.tsv`

15

16 The PAQman GitHub (<https://github.com/SAMtoBAM/PAQman>) contains descriptions of
17 installation/usage, commands, parameters, output files and both the features and their
18 metrics.

19

20 PAQman has many dependencies that are stand-alone programs in their own right. Please
21 cite both this manuscript in addition to its core dependencies. For example: “We used
22 PAQman v1.2.0 (O’Donnell et al. 2025) in conjunction with Quast v5.3.0 (Mikheenko et al.
23 2023), BUSCO v6.0.0 (Tegenfeldt et al. 2025), meryl v1.3 (Miller et al. 2008), Mercury v1.3
24 (Rhie et al. 2020), Rasusa 2.2.2 (Hall 2022), CRAQ v1.0.9 (Li et al. 2023), BWA v0.7.19 (Li
25 2013), minimap2 v2.30 (Li 2018), samtools v1.22.1 (Danecek et al. 2021), mosdepth v0.3.12
26 (Pedersen and Quinlan 2018), bedtools v2.31.1 (Quinlan and Hall 2010), seqkit v2.10.0
27 (Shen et al. 2016) and ggplot2 v3.5.2 (Wickham 2016) to assess and visualize assembly
28 quality.”

1 Conclusion

2 As the ease of producing reference-quality assemblies continues to advance, new tools are
3 needed to comprehensively evaluate assembly quality in a high-throughput manner.
4 PAQman provides a simple means of evaluating an assembly's overall quality by
5 streamlining the calculation of seven commonly used features in the absence of a reference
6 assembly. PAQman thus not only helps identify possible issues present within a specific
7 assembly but simplifies the process of benchmarking assembly software to find the set of
8 parameters that produce the “best” assembly for a given sequencing dataset and organism.

9 Data Availability

10 PAQman is distributed as a conda environment and an Apptainer image. Source code and documentation
11 are freely available through the GitHub repository at <https://github.com/samtobam/paqman> and a Zenodo
12 archive at <https://doi.org/10.5281/zenodo.16039705> for v1.2.0 used in this study.

13 Author Contributions

14 Samuel O’Donnell (Conceptualization, Investigation, Methodology, Software, Validation,
15 Visualization, Writing – original draft, Writing – review & editing)

16 Ningxiao Li (Data curation, Writing – review & editing)

17 Jacob L. Steenwyk (Conceptualization, Writing – review & editing)

18 David M. Geiser (Resources, Writing – review & editing)

19 Frank N. Martin (Data curation, Resources, Writing – review & editing)

20 Emile Gluck-Thaler (Conceptualization, Funding acquisition, Resources, Writing – review &
21 editing)

22 Study Funding

23 EGT and SO are supported by the Office of the Vice Chancellor for Research and Graduate
24 Education at the University of Wisconsin-Madison with funding from the Wisconsin Alumni
25 Research Foundation and the Department of Plant Pathology at the University of Wisconsin-
26 Madison. JLS is a Howard Hughes Medical Institute Awardee of the Life Sciences Research
27 Foundation.

1 Conflicts of Interest

2 JLS is an advisor to ForensisGroup Inc. JLS is a scientific consultant to Edison Scientific Inc.

3

4 References

5 Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, Boisvert S,
6 Chapman JA, Chapuis G, Chikhi R, et al. 2013. Assemblathon 2: evaluating de novo
7 methods of genome assembly in three vertebrate species. *GigaScience* 2:10.

8 Cosma B-M, Shirali Hossein Zade R, Jordan EN, van Lent P, Peng C, Pillay S, Abeel T.
9 2023. Evaluating long-read de novo assembly tools for eukaryotic genomes: insights
10 and considerations. *GigaScience* 12:giad100.

11 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane
12 T, McCarthy SA, Davies RM, et al. 2021. Twelve years of SAMtools and BCFtools.
13 *GigaScience* 10:giab008.

14 Espinosa E, Bautista R, Larrosa R, Plata O. 2024. Advancements in long-read genome
15 sequencing technologies and algorithms. *Genomics* 116:110842.

16 Giani AM, Gallo GR, Gianfranceschi L, Formenti G. 2020. Long walk to genomics:
17 History and current approaches to genome sequencing and assembly. *Comput.*
18 *Struct. Biotechnol. J.* 18:9–19.

19 Hall MB. 2022. Rasusa: Randomly subsample sequencing reads to a specified
20 coverage. *J. Open Source Softw.* 7:3941.

21 Kolmogorov M, Yuan J, Lin Y, Pevzner PA. 2019. Assembly of long, error-prone reads
22 using repeat graphs. *Nat. Biotechnol.* 37:540–546.

23 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with
24 BWA-MEM. *arXiv:1303.3997* [Internet]. Available from: <http://arxiv.org/abs/1303.3997>

25 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*
26 34:3094–3100.

27 Li K, Xu P, Wang J, Yi X, Jiao Y. 2023. Identification of errors in draft genome assemblies
28 at single-nucleotide resolution for quality assessment and improvement. *Nat.*
29 *Commun.* 14:6556.

- 1 Lyčka M, Bubeník M, Závodník M, Peska V, Fajkus P, Demko M, Fajkus J, Fojtová M.
2 2024. TeloBase: a community-curated database of telomere sequences across the
3 tree of life. *Nucleic Acids Res.* 52:D311–D321.
- 4 Manchanda N, Portwood JL, Woodhouse MR, Seetharam AS, Lawrence-Dill CJ,
5 Andorf CM, Hufford MB. 2020. GenomeQC: a quality assessment tool for genome
6 assemblies and gene structure annotations. *BMC Genomics* 21:193.
- 7 Mikheenko A, Saveliev V, Hirsch P, Gurevich A. 2023. WebQUAST: online evaluation of
8 genome assemblies. *Nucleic Acids Res.* 51:W601–W606.
- 9 Miller JR, Delcher AL, Koren S, Venter E, Walenz BP, Brownley A, Johnson J, Li K,
10 Mobarry C, Sutton G. 2008. Aggressive assembly of pyrosequencing reads with
11 mates. *Bioinformatics* 24:2818–2824.
- 12 O'Donnell S, Yue J-X, Saada OA, Agier N, Caradec C, Cokelaer T, De Chiara M, Delmas
13 S, Dutreux F, Fournier T, et al. 2023. Telomere-to-telomere assemblies of 142 strains
14 characterize the genome structural landscape in *Saccharomyces cerevisiae*. *Nat.*
15 *Genet.* 55:1390–1399.
- 16 Pedersen BS, Quinlan AR. 2018. Mosdepth: quick coverage calculation for genomes
17 and exomes. *Bioinformatics* 34:867–868.
- 18 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing
19 genomic features. *Bioinformatics* 26:841–842.
- 20 Rhie A, Nurk S, Cechova M, Hoyt SJ, Taylor DJ, Altemose N, Hook PW, Koren S,
21 Rautiainen M, Alexandrov IA, et al. 2023. The complete sequence of a human Y
22 chromosome. *Nature* 621:344–354.
- 23 Rhie A, Walenz BP, Koren S, Phillippy AM. 2020. Merqury: reference-free quality,
24 completeness, and phasing assessment for genome assemblies. *Genome Biol.*
25 21:245.
- 26 Shang L, He W, Wang T, Yang Y, Xu Q, Zhao X, Yang L, Zhang H, Li X, Lv Y, et al. 2023. A
27 complete assembly of the rice Nipponbare reference genome. *Mol. Plant* 16:1232–
28 1236.
- 29 Shen W, Le S, Li Y, Hu F. 2016. SeqKit: A Cross-Platform and Ultrafast Toolkit for
30 FASTA/Q File Manipulation. *PLOS ONE* 11:e0163962.
- 31 Tegenfeldt F, Kuznetsov D, Manni M, Berkeley M, Zdobnov EM, Kriventseva EV. 2025.
32 OrthoDB and BUSCO update: annotation of orthologs with wider sampling of
33 genomes. *Nucleic Acids Res.* 53:D516–D522.

1 Wickham H. 2016. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New
 2 York Available from: <https://ggplot2.tidyverse.org>

3 Zhang X, Liu C-G, Yang S-H, Wang X, Bai F-W, Wang Z. 2022. Benchmarking of long-
 4 read sequencing, assemblers and polishers for yeast genome. *Brief. Bioinform.*
 5 23:bbac146.

6 Zhang Y, Lu H-W, Ruan J. 2023. GAEP: a comprehensive genome assembly evaluating
 7 pipeline. *J. Genet. Genomics* 50:747–754.

8

9

10

11 Figure Legends

12 **Figure 1.** Schematic of the PAQman pipeline consisting of the *paqman* and *paqplots*
 13 commands. Grey boxes indicate files, blue boxes indicate computational steps calculating
 14 specific assembly metrics (with software names italicized below) and red boxes indicate
 15 outputs. Example figures output by PAQman are to the right. Brackets indicate the short-
 16 hand command line arguments for input files and parameters impacting that specific step.

17 **Figure 2.** Examples of output visualizations from *paqplot* outputs of various *paqman* Feature
 18 metrics. Three different datasets were used: (A-B) *S. cerevisiae* reference strain S288c, (C-
 19 D) *Oryza sativa* subsp. *japonica* Nipponbare cultivar, and the human cell line CHM13 (E-F);
 20 each with five different assemblies. (A, C, E) Radar and lollipop plots using the raw values of
 21 several metrics. (B, D, F) Radar plot using relative scales of metrics between assemblies. The
 22 individual metrics plotted represent all seven Features of assembly quality measured by
 23 PAQman; 1. Contiguity ('assembly size', 'assembly N50', 'assembly N90', 'largest contig',
 24 '#contigs' and '#contigs>10kb'); 2. Gene content ('BUSCO complete single (%)' and 'BUSCO
 25 complete (%)'); 3. Completeness ('kmer completeness (%)'); 4. Accuracy ('qv (phred)'); 5.
 26 Correctness ('R-AQI (%)' and 'S-AQI (%)'); 6. Coverage ('coverage normal (%)'); and 7.
 27 Telomericity ('Telomeric ends (%)', 'Telomeric ends' and 'T2T contigs').

28

1 Tables

2 **Table 1:** Average run time of *paqman.sh* across five assemblies for each dataset using an AMD EPYC 7513 32-Core
 3 Processor and x86_64 architecture. Run time (real) was calculated using the *time* command in Linux; Max RAM usage was
 4 calculated using *sacct* (mean \pm the standard deviation). Only 32 threads were tested on CHM13 with PacBio HiFi reads.

Input					Run results	
Dataset	Genome size	Reads	coverage	threads	real (minutes)	Max RAM usage (GB)
S288c	~12Mbp	ONT	~100X	8	6.92 \pm 0.1	6.19 \pm 0.78
				16	6.29 \pm 0.31	5.54 \pm 0.36
				32	5.3 \pm 0.14	4.73 \pm 0.42
Nipponbare	~380Mbp	HiFi	~85X	8	139.75 \pm 1.96	42.88 \pm 7.09
				16	92.86 \pm 1.71	47.864 \pm 12.55
				32	75.8 \pm 2.66	47.31 \pm 15.77
CHM13	~3Gbp	HiFi	~25X	64	423.59 \pm 23.32	218.67 \pm 32.16

5

6

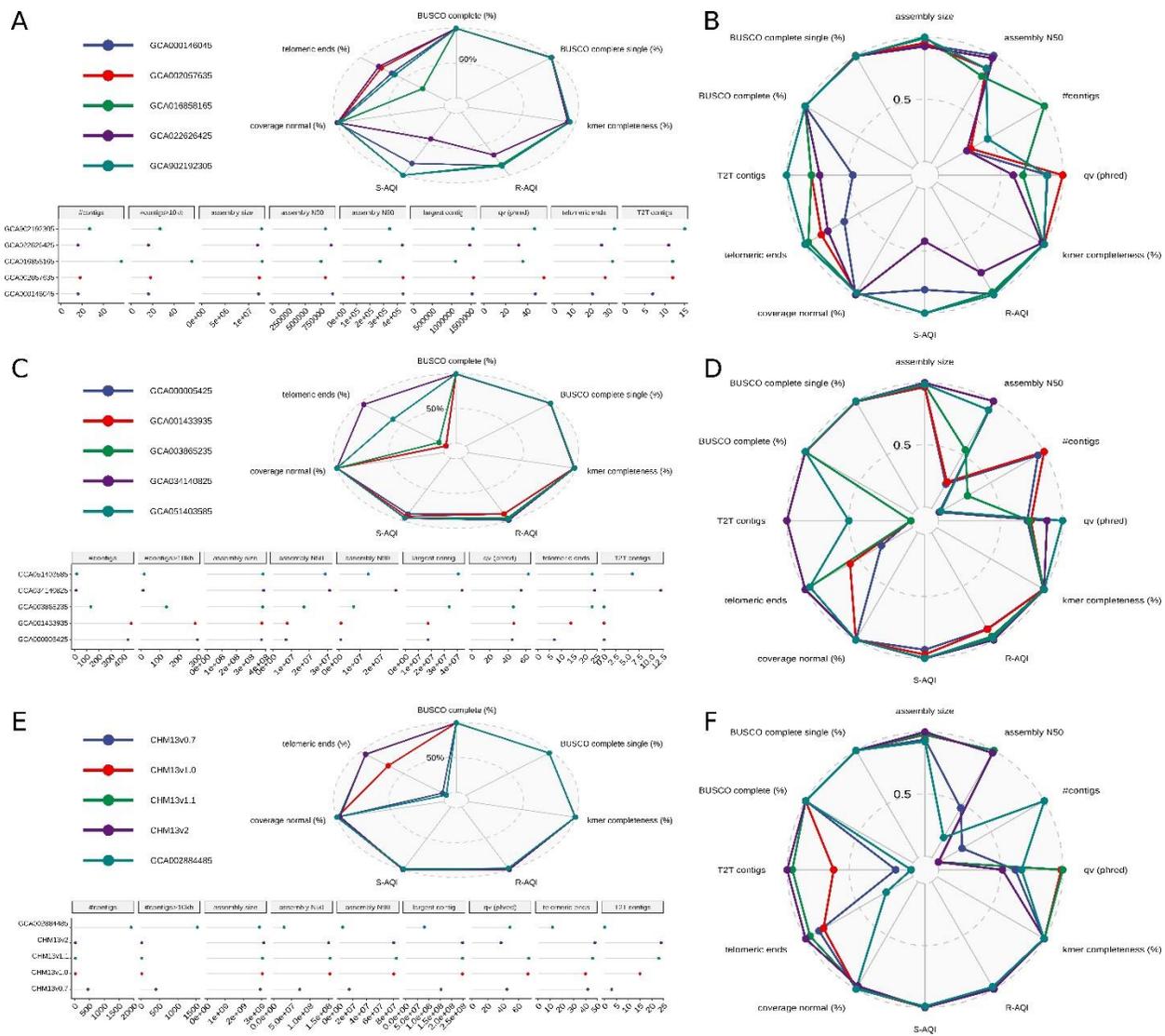


Figure 2
165x146 mm (DPI)

1
2
3