



Assemblathon 1: A competitive assessment of de novo short read assembly methods

Dent A. Earl, Keith Bradnam, John St. John, et al.

Genome Res. published online September 16, 2011
Access the most recent version at doi:[10.1101/gr.126599.111](https://doi.org/10.1101/gr.126599.111)

P<P	Published online September 16, 2011 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
License	This manuscript is Open Access.
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Copyright © 2011, Cold Spring Harbor Laboratory Press

1 **TITLE PAGE**

2 **Assemblathon 1: A competitive assessment of de novo short read assembly methods**

4 **RUNNING TITLE: Assemblathon 1**

6 **KEYWORDS: Genome Assembly,**

8 **CONTACT:**

9 **Benedict Paten, benedict@soe.ucsc.edu**

12 Dent Earl^{1,2}, Keith Bradnam³, John St. John^{1,2}, Aaron Darling³, Dawei Lin^{3,4}, Joseph Fass^{3,4}, Hung
13 On Ken Yu³, Vince Buffalo^{3,4}, Daniel R. Zerbino², Mark Diekhans^{1,2}, Ngan Nguyen^{1,2}, Pramila
14 Nuwantha⁵, Ariyaratne Wing-Kin Sung^{5,6}, Zemin Ning⁷, Matthias Haimel⁸, Jared T. Simpson⁷,
15 Nuno A. Fonseca⁹, İnanç Birol¹⁰, T. Roderick Docking¹⁰, Isaac Y. Ho¹¹, Daniel S. Rokhsar^{11,12},
16 Rayan Chikhi^{13,14}, Dominique Lavenier^{13,14,15}, Guillaume Chapuis^{13,14}, Delphine Naquin^{14,15},
17 Nicolas Maillet^{14,15}, Michael C. Schatz¹⁶, David R. Kelley¹⁷, Adam M. Phillippy^{17,18}, Sergey
18 Koren^{17,18}, Shiao-Pyng Yang¹⁹, Wei Wu¹⁹, Wen-Chi Chou²⁰, Anuj Srivastava²⁰, Timothy I. Shaw²⁰,
19 J. Graham Ruby^{21,22,23}, Peter Skewes-Cox^{21,22,23}, Miguel Betegon^{21,22,23}, Michelle T. Dimon^{21,22,23},
20 Victor Solovyev²⁴, Igor Seledtsov²⁵, Petr Kosarev²⁵, Denis Vorobyev²⁵, Ricardo Ramirez-
21 Gonzalez²⁶, Richard Leggett²⁷, Dan MacLean²⁷, Fangfang Xia²⁸, Ruibang Luo²⁹, Zhenyu L²⁹,
22 Yinlong Xie²⁹, Binghang Liu²⁹, Sante Gnerre³⁰, Iain MacCallum³⁰, Dariusz Przybylski³⁰, Filipe J.
23 Ribeiro³⁰, Shuangye Yin³⁰, Ted Sharpe³⁰, Giles Hall³⁰, Paul J. Kersey⁸, Richard Durbin⁷, Shaun D.
24 Jackman¹⁰, Jarrod A. Chapman¹¹, Xiaoqiu Huang³¹, Joseph L. DeRisi^{21,22,23}, Mario Caccamo²⁶,
25 Yingrui Li²⁹, David B. Jaffe³⁰, Richard E. Green², David Haussler^{1,2,23}, Ian Korf^{3,32}, Benedict
26 Paten^{1,2}

27 (1) Center for Biomolecular Science and Engineering, University of California Santa Cruz, CA,
28 USA

29 (2) Biomolecular Engineering Department, University of California Santa Cruz, CA, USA

30 (3) Genome Center, University of California Davis, CA, USA

31 (4) Bioinformatics Core, Genome Center, University of California Davis, CA, USA

32 (5) Computational & Mathematical Biology Group, Genome Institute of Singapore, Singapore

33 (6) School of Computing, National University of Singapore, Singapore

34 (7) Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton, UK

35 (8) EMBL-EBI, Wellcome Trust Genome Campus, Hinxton, UK

36 (9) CRACS-INESC Porto LA, Universidade do Porto, Portugal

- 1 (10) Genome Sciences Centre, British Columbia Cancer Agency, Vancouver, British Columbia,
- 2 Canada
- 3 (11) DOE Joint Genome Institute, Walnut Creek, CA, USA
- 4 (12) UC Berkeley, Dept, of Molecular and Cell Biology, Berkeley, CA, USA
- 5 (13) Computer Science department, ENS Cachan/IRISA, Rennes, France
- 6 (14) CNRS/Symbiose, IRISA, Rennes, France
- 7 (15) INRIA, Rennes Bretagne Atlantique, Rennes, France
- 8 (16) Simons Center for Quantitative Biology, Cold Spring Harbor Laboratory, Cold Spring Harbor,
- 9 NY, USA
- 10 (17) Center for Bioinformatics and Computational Biology, University of Maryland, College Park,
- 11 MD, USA
- 12 (18) National Biodefense Analysis and Countermeasures Center, Fredrick, MD, USA
- 13 (19) Monsanto Company, 700 Chesterfield Parkway, Chesterfield, MO, USA
- 14 (20) Institute of Bioinformatics, University of Georgia, Athens, GA, USA
- 15 (21) Department of Biochemistry and Biophysics, University of California San Francisco, CA, USA
- 16 (22) Biological and Medical Informatics Program, University of California, San Francisco, CA,
- 17 USA
- 18 (23) Howard Hughes Medical Institute, Bethesda, MD, USA
- 19 (24) Department of Computer Science, Royal Holloway, University of London, UK
- 20 (25) Softberry Inc., 116 Radio Circle, Suite 400, Mount Kisco, NY, USA
- 21 (26) The Genome Analysis Centre, Norwich Research Park, Norwich, UK
- 22 (27) The Sainsbury Laboratory, Norwich Research Park, Norwich, UK
- 23 (28) Computation Institute, University of Chicago, IL, USA
- 24 (29) BGI-Shenzhen, Shenzhen 518083, China
- 25 (30) Broad Institute, Cambridge, MA, USA
- 26 (31) Department of Computer Science, Iowa State University, Ames, IA, USA
- 27 (32) Molecular and Cellular Biology, Genome Center, University of California Davis, CA, USA

28

29 **ABSTRACT**

30

31 Low cost short read sequencing technology has revolutionised genomics, though it is only just
32 becoming practical for the high quality de novo assembly of a novel large genome. We describe
33 the Assemblathon 1 competition, which aimed to comprehensively assess the state of the art in
34 de novo assembly methods when applied to current sequencing technologies. In a collaborative
35 effort teams were asked to assemble a simulated Illumina HiSeq dataset of an unknown,
36 simulated diploid genome. A total of 41 assemblies from 17 different groups were received. Novel

1 haplotype aware assessments of coverage, contiguity, structure, base calling and copy number
2 were made. We establish that within this benchmark (1) it is possible to assemble the genome to
3 a high level of coverage and accuracy, and that (2) large differences exist between the
4 assemblies, suggesting room for further improvements in current methods. The simulated
5 benchmark, including the correct answer, the assemblies and the code that was used to evaluate
6 the assemblies is now public and freely available from <http://www.assemblathon.org/>.

8 **INTRODUCTION**

9 Sequence assembly is the problem of merging and ordering shorter fragments, termed “reads,”
10 sampled from a set of larger sequences in order to reconstruct the larger sequences. The output
11 of an assembly is typically a set of “contigs,” which are contiguous sequence fragments, ordered
12 and oriented into “scaffold” sequences, with gaps between contigs within scaffolds representing
13 regions of uncertainty.

14
15 There are numerous subclasses of assembly problem that can be distinguished by, amongst
16 other things, the nature of (1) the reads, (2) the types of sequences being assembled and (3) the
17 availability of homologous (related) and previously assembled sequences, such as a reference
18 genome or the genome of a closely related species (Pop and Salzberg 2008) (Trapnell and
19 Salzberg 2009) (Chaisson et al. 2009). In this work we focus on the evaluation of methods for de
20 novo genome assembly using low cost “short read” technology, where the reads are
21 comparatively short in length but large in number, the sequences being assembled represent a
22 novel diploid genome and the nearest homologous genome to that being assembled is
23 significantly diverged.

24
25 In bioinformatics, the reads used in an assembly are derived from an underlying sequencing
26 technology. For a recent review of sequencing technologies see Metzker (2010). For the
27 assembly problem there are a number of key considerations, notably (1) the length of the reads,
28 (2) the error characteristics of the reads, (3) if and how the reads are “paired”, i.e. where reads
29 are produced in pairs separated by an approximately fixed length spacer sequence, and finally (4)
30 the number of reads produced for a given cost.

31
32 Sanger sequencing (Sanger et al. 1977) produces relatively long reads, typically between 300 to
33 1000 base pairs (bp) in length, with a low error rate, but which are comparatively expensive to
34 produce. After relying primarily on Sanger sequencing for decades the field of sequencing has
35 recently witnessed a diversification of competing technologies ((Margulies et al., 2005), (Bentley
36 2006), (Pandey et al. 2008), (Eid et al. 2009), (Pourmand et al. 2006)) and a rapid rate of overall

1 change. One direction of this development has been a move towards shorter reads, often less
2 than or equal to 150 bp, but at a much lower cost for a given volume of reads ((Bentley, 2006),
3 (Pandey et al. 2008), (Pourmand et al. 2006)).

4
5 As the field of sequencing has changed so has the field of sequence assembly, for a recent
6 review see Miller et al. (2010). In brief, using Sanger sequencing, contigs were initially built using
7 overlap or string graphs (Myers 2005) (or data structures closely related to them), in tools such
8 as Phrap (<http://www.phrap.org/>), GigAssembler (Kent and Haussler, 2001), Celera (Myers et al.
9 2000) (Venter et al. 2001), ARACHNE (Batzoglou et al. 2002), and Phusion (Mullikin and Ning
10 2003), which were used for numerous high quality assemblies such as human (Lander et al.
11 2001) and mouse (Mouse Genome Sequencing Consortium et al. 2002). However, these
12 programs were not generally efficient enough to handle the volume of sequences produced by the
13 next generation sequencing technologies, spurring the development of a new generation of
14 assembly software.

15
16 While some maintained the overlap graph approach, e.g. Edena (Hernandez et al. 2008) and
17 Newbler (<http://www.454.com/>), others used word look-up tables to greedily extend reads, e.g.
18 SSAKE (Warren et al. 2007), SHARCGS (Dohm et al. 2007), VCAKE (Jeck et al. 2007) and
19 OligoZip (<http://linux1.softberry.com/berry.phtml?topic=OligoZip>). These word look-up tables were
20 then extended into de Bruijn graphs to allow for global analyses (Pevzner et al. 2001), e.g. Euler
21 (Chaisson and Pevzner 2008), AllPaths (Butler et al. 2008) and Velvet (Zerbino and Birney 2008).
22 As projects grew in scale further engineering was required to fit large whole genome datasets into
23 memory ((ABYSS (Simpson et al. 2009), Meraculous (in submission)), (SOAPdenovo (Li et al.
24 2010), Cortex (in submission)). Now, as improvements in sequencer technology are extending the
25 length of “short reads”, the overlap graph approach is being revisited, albeit with optimized
26 programming techniques, e.g. SGA (Simpson and Durbin 2010), as are greedy contig extension
27 algorithms, e.g. PRICE (<http://derisilab.ucsf.edu/software/price/index.html>), Monument
28 (<http://www.irsia.fr/symbiose/people/rchkhi/monument.html>).

29
30 In general, most sequence assembly programs are multi stage pipelines, dealing with correcting
31 measurement errors within the reads, constructing contigs, resolving repeats (i.e. disambiguating
32 false positive alignments between reads) and scaffolding contigs in separate phases. Since a
33 number of solutions are available for each task, several projects have been initiated to explore the
34 parameter space of the assembly problem, in particular in the context of short read sequencing
35 ((Phillippy et al. 2008), (Hubis et al. 2011), (Alkan et al. 2011), (Narzisi and Mishra 2011), (Zhang

1 et al. 2011) and (Lin et al. 2011)). In this work we are concerned with evaluating assembly
2 programs as whole, with the aim of comprehensively evaluating different aspects of assemblies.

3
4 It is generally the case that the right answer to an assembly problem is unknown. Understandably
5 therefore, a common method for assessing assembly quality has been the calculation of length
6 summary statistics on the produced scaffold and contig sequences. Such metrics include various
7 weighted median statistics, such as the N50 defined below, as well as the total sequence lengths
8 and total numbers of sequences produced ((Lindblad-Toh et al. 2005), (Ming et al. 2008), (Liu et
9 al. 2009), (Church et al. 2009), (Li et al. 2010), (Locke et al. 2011) and (Colbourne et al. 2011)).

10
11 Other methods have been proposed for evaluating the internal consistency of an assembly, for
12 example, by analysing the consistency of paired reads, as in the clone-middle plot (Huson et al.
13 2001), by looking for variations in the depths of read coverage supporting a constructed assembly
14 (Phillippy et al. 2008) and looking at haplotype inconsistency (Lindblad-Toh et al. 2005).

15
16 To assess accuracy, assemblies may be compared to finished sequences derived from
17 independent sequencing experiments or to sequences held out of the assembly process. For the
18 dog genome, nine bacterial artificial chromosomes (BACs) were sequenced to finishing standards
19 and held out of the assembly (Lindblad-Toh et al. 2005), for the panda genome, which was
20 primarily an Illumina assembly, extra Sanger sequencing of BACs was performed (Li et al. 2010).
21 Additionally, if genetic mapping data is available such information can also be used to assess
22 scaffold quality, e.g. Church et al. (2009), which used a combination of linkage, radiation hybrid
23 and optical maps. Church et al. also demonstrate that transcriptome (the set of RNA molecules
24 for a given cell type) information, if available, can also be used to assess the validity of a genomic
25 assembly by checking the extent to which the assembly recapitulates the transcriptome.

26
27 When a reference genome or sequence is available a comparison between the assembly and
28 reference can be performed. This has previously been accomplished by studies using several
29 different genome alignment methods, including BLAST (e.g. Zang et al. 2011), LASTZ (e.g.
30 (Hubisz et al. 2011)) and Exonerate (Zerbino and Birney 2008, Hernandez et al. 2008). Given
31 such an alignment, most simply, the proportion of a reference's coverage can be reported (Li, et
32 al., 2010) (Zhang et al. 2011). Notably, Gnerre et al. (2011) compared novel short-read
33 assemblies to the human and mouse reference genomes,, and performed a comprehensive set of
34 analyses that encompassed coverage, contig accuracy and the long range contiguity of scaffolds.
35 Related to the work described here, Butler et al. (2008) described a graph based pattern analysis
36 using an assembly to reference alignment.

1

2 Comparison can also be made to a well-sequenced, related species. This can be done using the
3 complete genomic sequence of an out group, for example, Meader et al. (2010) presented an
4 assessment method based on patterns of insertions and deletions (indels) in closely related inter-
5 species genome alignments. Alternatively, specific genomic features can be studied, for example
6 Parra et al. examined the fraction of “core genes,” those present in all genomes, that could be
7 identified in draft genome assemblies (Parra et al. 2009).

8

9 Simulation has been a mainstay of genome assembly evaluation since assembly methodology
10 was first developed and with few exceptions (e.g. (Maccallum, et al., 2009), (Gnerre, et al., 2011))
11 is de rigueur when introducing new de novo assembly software (e.g. (Myers, et al., 2000),
12 (Lander et al. 2001), (Venter et al. 2001), (Batzoglou et al. 2002), (Warren et al. 2007), (Jeck et al.
13 2007), (Dohm et al. 2007), (Chaisson and Pevzner 2008), (Zerbino and Birney 2008), (Butler et al.
14 2008) etc). In this work we have also chosen to use simulations, utilising the new Evolver genome
15 evolution simulation tool (Edgar R, Aseminos G, Batzoglou S, Sidow Arend
16 <http://www.drive5.com/evolver/>) to produce a simulated diploid genome with parameters that
17 approximate that of a vertebrate genome, though at $\sim 1/10^{\text{th}}$ the scale.

18

19 From this novel genome we simulate reads, modelling an Illumina sequencing run, using a newly
20 developed read simulator. Assembly teams were asked to assemble this novel genome blind and
21 we present an analysis of the resulting assemblies. By using simulation we know a priori the
22 haplotype relationships; by a process of multiple sequence alignment (MSA) we assess the
23 relationships between the assemblies and the original haplotypes of our simulation. This novel
24 process allows us evaluate haplotype specific contributions to the assemblies. Additionally, as a
25 positive control for our results, we assess the generated assemblies using more traditional
26 BLAST (Zang et al. 2011), (Hubisz et al. 2011)) methods, and support our assessments by
27 making all the code and data from our assessments public and freely available
28 (<http://compbio.soe.ucsc.edu/assemblathon1/>, <http://www.assemblathon.org/>).

29

30 **RESULTS**

31

32 We start by giving an overview of the Assemblathon 1 dataset and its generation. We then
33 describe the assemblies before giving the results of different evaluations.

34

35 **Genome simulation**

36

1 Rather than use an existing reference genome for assessment, we opted to simulate a novel
2 genome. We did this primarily for three reasons. Firstly, it gave us a genome that had no
3 reasonable homology to anything other than out-group genomes that we generated and provided
4 to assemblers. This allowed for a fair, blind test in which none of the assembly contributors had
5 access to the underlying genomes during the competition. Secondly, we were able to precisely
6 tailor the proportions of the simulated genome to those desired for this experimental analysis, i.e.
7 to limit the size of the genome to less than that of a full mammalian genome and thus allow the
8 maximum number of participants, while still maintaining a size that posed a reasonable challenge.
9 Thirdly, we could simulate a diploid genome; we know of no existing diploid dataset (simulated or
10 real) in which the contributions of the two haplotypes are precisely and fully known. This allowed
11 us to assess a heretofore-unexplored dimension of assembly assessment.

12

13 To simulate the genome we used the Evolver suite of genome evolution tools. Evolver simulates
14 the forward evolution of multi-chromosome haploid genomes and includes models for
15 evolutionary constraint, protein codons, genes and mobile elements.

16

17 The input genome for the simulation, termed the *root genome*, was constructed by downloading
18 the DNA sequence and annotations (see methods) for human chromosome 13 (hg18/NCBI36,
19 95.6 non-N megabases (Mb)) from the UCSC Table Browser (Fujita et al. 2011) and dividing it
20 into four chromosomes of approximately equal length. Figure 1 shows the phylogeny used to
21 generate the simulated genomes, with branch lengths to scale. We first evolved the root genome
22 for ~200 million years (my) to generate the most recent common ancestor (MRCA) of the final leaf
23 genomes. We performed this long burn-in on the genome in order to reshuffle the sequence and
24 annotations present, thereby preventing simple discovery of the source of the root genome. The
25 simulation then proceeded along two independent lineages, generating both α and β , each ~50
26 my diverged from the MRCA. Finally, in both lineages we split the evolved genome into two sub-
27 lineages, termed *haplotypes*, and evolved these sub-lineages for a further ~1 my, to produce a
28 pair of diploid genomes $\alpha_{1,2}$ and $\beta_{1,2}$, each with a degree of polymorphism. The $\alpha_{1,2}$ genome's
29 haplotypes, α_1 and α_2 , each had three chromosomes and both haplotypes were 112.5 Mb in
30 total length with chromosome lengths of 76.3, 18.5 and 17.7 Mb.

31

32 The diploid $\alpha_{1,2}$ genome was used as the target genome for the assembly. The $\beta_{1,2}$ genome's
33 haplotypes, their common ancestor, β and their annotations were provided to the assemblers as
34 an out-group. Relatively few assemblers (see Table 1) reported using these sequences to assist
35 in the assembly process.

1

2 Table 2A provides a count of some of the events that took place along particular branches in the
3 phylogenetic tree during the course of the simulation. Table 2B provides a summary of the
4 pairwise differences between the α_1 and α_2 haplotypes and Supplementary Figure 1 shows a
5 dot-plot of their alignment. Supplementary Figures 2 and 3 show the length distribution of
6 annotations for the root, MRCA, internal node and leaf genomes, demonstrating these
7 annotations remained approximately static over the course of the simulation. We examined repeat
8 content of the simulated genomes (see Table 2 legend) and found a comparable portion of
9 annotated repeats to that in the original human chromosome 13, but a reduction of slightly more
10 than half in the proportion of repetitive 100-mers (Supplementary Figure 4). The simple
11 substitution model used by Evolver, which fails to capture some of the higher order dependencies
12 in substitution patterns that made the original human DNA sequence more repetitive, likely
13 explains this latter observation

14

15 **Read Simulation**

16

17 As mentioned, there are many competing technologies now available for sequencing, giving us
18 many possible options in designing the datasets for the first Assemblathon. However, we opted to
19 simulate just one combined short read dataset, with multiple read libraries, for the Illumina Hi-seq
20 platform (Illumina, Illumina HiSeq 2000), which is the current market leader for low cost de novo
21 sequencing on this scale. The advantage of this was chiefly (1) the avoidance of fragmentation in
22 the entries to the Assemblathon, thereby preventing categories with few or just one entry, and (2)
23 the ability to assess all the submitted assemblies with common sets of evaluations.

24

25 We needed a program that would generate short reads and model sources of error that the
26 Illumina protocols introduce. As we knew of no existing software that was capable of this (see
27 methods for a discussion of existing read simulators), we wrote our own short read simulator
28 called SimSeq (St. John J, <https://github.com/jstjohn/SimSeq>).

29

30 Abstractly, reads were sampled from the genome using one of two types of Illumina paired read
31 strategy, so called 'paired-end' and 'mate-pair' strategies, after which an error profile was applied
32 to each read in its proper orientation. In addition to generating reads from the target $\alpha_{1,2}$ genome,
33 three copies of an *Escherichia coli* sequence (gi 312944605) were added to the two haplotype
34 sequences to yield a ~5% bacterial contamination rate. Bacterial sequence was included as an
35 attempt model the sort of contamination occasionally present in data from sequencing centres,
36 though the specific choice of *E.coli* and the 5% level were arbitrary. Participants in the contest

1 were notified that some bacterial contamination was present in the data, though they were not
2 told about its precise nature nor explicitly told to remove it.

3
4 Multiple libraries were generated for both the paired-end and mate-pair strategies. Paired-end
5 libraries with 200 and 300 bp inserts contributed 80x, mate-pair libraries with separations of 3
6 kilobases (kb) and 10 kb contributed a further 40x, giving a total coverage of 120x for the sample.
7 Removing contamination reads gave an overall coverage of ~55x per haplotype.

8
9 A detailed description of the simulation method, the types of errors simulated, and the simulator's
10 limitations are given in the methods. Importantly, due to human error, the error model was
11 mistakenly reversed along the reads. This resulted in bases with a slightly higher error rate
12 tending to appear towards the beginning of the reads rather than towards the end of reads (see
13 Supplementary Figure 5). This issue only manifests itself if the reads are treated asymmetrically;
14 we surveyed participants on this matter and only one group, L'IRISA, indicated that their
15 methodology was possibly harmed more than other methods due to the mistake.

16 17 **Assemblies**

18
19 The competition started in January 2011 and teams were given just over a month to submit their
20 assemblies. Teams were allowed to submit up to five separate assemblies for consideration.
21 Additionally, assemblies were created by the organisers with popular assembly programs, using
22 default parameters, as a way of comparing naively generated assemblies with those that were
23 contributed by independent groups. Table 1 lists the evaluated assemblies, the main program
24 used to generate them and the groups that contributed them (see Supplementary section 8.2 for
25 detailed information on submissions). In total there were 59 assemblies, with 41 independently
26 contributed by 17 different groups using 15 different assembly programs and 18 generated by the
27 organisers using three popular programs.

28 29 **Evaluations**

30

1 We assessed all the contributed assemblies, full results for which can be found in the
2 supplementary material. However, to make the presentation succinct we choose to present only
3 the “top” assembly from each group in the following evaluations. To enable this we created a
4 ranking of the assemblies (see

5 Table 3, Supplementary Table 1), using the evaluations described below, and selected the
6 assembly from each group with the top overall ranking for inclusion. Full results for each
7 evaluation on every assembly in the main text can be found in the supplementary material.

8 9 **N50 and NG50**

10
11 A commonly used metric to assess assemblies is the N50 statistic. The N50 of an assembly is a
12 weighted median of the lengths of the sequences it contains, equal to the length of the longest
13 sequence s such that the sum of the lengths of sequences greater than or equal in length to s is
14 greater than or equal to half the length of the genome being assembled. As the length of the
15 genome being assembled is generally unknown, the normal approximation is to use the total
16 length of all of the sequences in an assembly as a proxy for the denominator. We follow this
17 convention for calculating N50, but additionally we define the *NG50* (G for genome). The NG50 is
18 identical to N50, except that we estimate the length of the genome being assembled as being
19 equal to the average of the length of the two haplotypes, α_1 and α_2 . Contig N50s and NG50s,
20 where the sequences are the set of assembly contigs, and scaffold N50s and NG50s, where the
21 sequences are the set of assembly scaffolds, are shown in Figure 2, Supplementary Figure 6 and
22 Supplementary Table 2.

23
24 The total span of most of the submitted assemblies was slightly larger than the haploid genome
25 size, primarily because of the degree of polymorphism of the two haplotypes. Thus the assembly-
26 specific N50s are in general smaller than the NG50s, with the median absolute difference
27 between contig NG50 and contig N50 being 599 bp (7.7%), and the median absolute difference
28 between scaffold NG50 and scaffold N50 being 1942 bp (3.6%). These differences are quite
29 small, though not negligible in every case, for example the CSHL assembly has a scaffold NG50
30 ~800 kb (31.6%) longer than scaffold N50.

31 32 **Multiple Sequence Alignment**

33
34 While N50 statistics give a sense of the scale and potential contiguity of an assembly they say
35 nothing necessarily about the underlying coverage or accuracy of an assembly. To compare each

1 assembly with the simulated genome and bacterial contamination we constructed a multiple
2 sequence alignment (MSA). The sequence inputs to the MSA were the two haplotypes, the
3 bacterial contamination and the scaffolds of the assembly. To generate the MSA we used an
4 adapted version (see methods) of the newly developed Cactus alignment program (Paten et al.
5 2011), a new MSA program able to handle rearrangements, copy number changes (duplications)
6 and missing data. The result of this alignment process was, for each assembly, a high specificity
7 map of the alignment of the assembly to the two haplotypes and the bacterial contamination.

8
9 As we aligned both the bacterial contamination and the two haplotypes together we used the
10 hypothetical existence of any alignments between the haplotypes and the bacterial contamination
11 as a negative control for non-specific alignment. We did not observe any such alignments. As an
12 additional control we replicated a similar, confirmatory analysis using a simple BLAST (Altschul et
13 al. 1990) strategy, details of which can be found in Supplementary section 7.1 and references to
14 which are made below.

15 16 **Coverage**

17
18 An MSA can be divided up into *columns*, each of which represents a set of individual base pair
19 positions in the input sequences that are considered homologous. We call columns that contain
20 positions within the haplotypes *haplotype columns*. We define the *overall coverage* of an MSA as
21 the proportion of haplotype columns that contain positions from the assembly. Similarly we can
22 define the coverage of *X*, where *X* is a specific haplotype or the bacterial contamination, as the
23 proportion of columns containing positions in *X* that also contain positions from the assembly.

24
25 Table 4 shows the overall, haplotype specific, and bacterial contamination specific coverage.
26 There is very little difference between the specific haplotype's coverage and the overall coverage,
27 and indeed little difference between many of the assemblies. The highest overall coverage was
28 the BGI assembly with 98.8%, but nearly all assemblies performed well in this metric with even
29 the assembly with 14th highest coverage, the CRACS assembly, providing 96.3% coverage.
30 However, there were huge differences in the coverage of the bacterial contamination
31 (Supplementary Figures 7 and 8), with many groups opting successfully to completely filter it out.
32 For example, the BGI assembly had no coverage of the bacterial sequence, while the ASTR
33 assembly had 100.0% coverage of the bacterial sequence.

34 35 **Blocks and Contig Paths**

1 Within an MSA a *block* is a maximal gapless alignment of a set of sequences, and is therefore
2 composed of a series of contiguous columns. The *length* of a block is equal to the number of
3 columns that it contains. We can use the block structure to define the *block NG50*, which is
4 exactly like the NG50, except that we use the distribution of block lengths rather than sequence
5 lengths. Supplementary Figures 9 and 10 show block coverage across the haplotypes. Alignment
6 of sequences that are very closely related are likely to contain fewer blocks with a greater base
7 pair length than sequences that are significantly diverged from one another. Unfortunately, the
8 two simulated haplotypes are sufficiently polymorphic with respect to one another that the block
9 NG50 of an alignment of just the two haplotypes is ~4 kb. As this length is much less than the
10 length of many sequences in the assemblies, assessing an assembly requires methods that do
11 not penalise the reconstruction of haplotype specific polymorphisms. This is evident by looking at
12 Figure 2, which shows that block NG50 is poorly discriminative. See Supplementary section 7.1.1
13 and Supplementary Figure 11 for supporting BLAST based analysis.

14

15 To extend our analysis we use a graph theoretic model of the alignments, which we now describe
16 in overview. An MSA can be described as a graph, and we call the simplest such graph an
17 *adjacency graph*. A formal description of the adjacency graph used here can be found in Paten et
18 al. (2011), it is closely related to the similarly name graph introduced in (Bergeron et al. 2006), but
19 also to a directed bigraph representation of a de Bruijn graph used in assembly (Medvedev and
20 Brudno 2009) and the multiple breakpoint graph used in the study of genome rearrangements
21 (Alekseyev and Pevzner 2009).

22

23 An adjacency graph G contains two kinds of edges, *block edges*, which represent the gapless
24 blocks of the alignment, and *adjacency edges*, which represent collections of connections
25 between the ends of segments of DNA. The nodes in the graph represent the ends of blocks of
26 aligned sequences. Figure 3 illustrates an example.

27

28 Each edge in G is labelled with the sub-sequences it represents, called *segments*, thus it is
29 possible to discern if the edge represents segments in the haplotypes, the assembly, the bacterial
30 contamination or some combination. As previously stated, no edges are contained in G that
31 represent segments in both the haplotypes and the bacterial contamination.

32

33 Within G a sequence is represented as a path of alternating adjacency and block edges, termed a
34 *thread*. We can assess the accuracy of assembly sequences by analysing their thread
35 representation in the adjacency graph. Let P be the thread representing an assembled sequence
36 in G . Any edge e in P is *consistent* if that edge is also labelled with segments from either or both

1 of the haplotypes. For any P a *contig path* is a maximal subpath of P in which all the edges are
2 consistent. Thus P can be divided up into a series of contig paths, possibly interspersed with
3 edges in P that are not contained in a contig path, see Figure 3 for an example. The *bp length* of
4 a contig path is equal to the sum of the bp lengths of the block edges it contains. Contig paths
5 represent maximal portions of the assembled sequence that are consistent with one or both of the
6 haplotypes and contain no assembly gaps, they can be thought of as portions of an assemblies'
7 contigs that perfectly follow a path through the graph of haplotype polymorphism.

8
9 Figure 2 shows *contig path NG50s*, defined analogously to block NG50; Supplementary Figures
10 12 and 13 show contig path coverage across the haplotypes, while Supplementary Figures 14
11 and 15 show in contrast the same plots, but instead use raw contig lengths. The contig path
12 NG50s are substantially larger than block NG50s, for example the BGI assembly has a contig
13 path NG50 1.5 orders of magnitude bigger than its block NG50. The difference between the
14 largest and smallest block NG50 is 2,556 bp (GACWT 1,351 bp to BGI 3,907 bp); the difference
15 between the largest and smallest contig path NG50 is 79,731 bp (GACWT 2,533 bp to BGI
16 82,264 bp). Thus the contig path NG50 results demonstrate that assemblies are able to
17 reconstruct substantial regions perfectly, and contig path NG50 appears to be a more
18 discriminative statistic than block NG50, as it indicates large differences between the assemblies.

19 20 **Scaffold paths**

21
22 To account for gaps within scaffolds, which we henceforth call *scaffold breaks*, we define *scaffold*
23 *paths*. Scaffold paths can be thought of as portions of the assemblies' scaffolds that perfectly
24 follow a path through the graph of haplotype polymorphism, but which are allowed to jump
25 unassembled sequences at *scaffold gaps*. Scaffold gaps are scaffold breaks (denoted as
26 contiguous runs of wildcard characters in an assembly) whose surrounding contig ends are
27 bridged by a path of haplotype representing edges within the adjacency graph, see Figure 3 for
28 an example and the methods for a formal definition.

29
30 Notably, our definition of a scaffold gap within the graph is permissive in that it allows (1) any
31 sequence of Ns to define a scaffold break and (2) the sequence of Ns that define the scaffold
32 break to be aligned within the ends of the block that sandwich the gap in the assembly. This
33 definition was sought because there is currently wide variation in the syntax used to define such
34 gaps within different assemblers and to be tolerant of alignment errors caused by the phenomena
35 of edge wander (Holmes and Durbin 1998), caused when the alignment of positions around a gap

1 has more than one equally probable scenario. As a scaffold path is a concatenation of contig
2 paths its bp length is just the sum of the bp lengths of the contig paths that it contains.

3
4 Figure 2 shows the scaffold path NG50, defined analogously to the block and contig path NG50s,
5 sorted with respect the scaffold NG50. In many cases the scaffold path NG50 is substantially
6 larger than the contig path NG50. Figure 4 shows a stack fill plot of the coverage of scaffold paths
7 along the three chromosomes of haplotype α_1 (see also Supplementary Figures 16-19). It
8 demonstrates the substantial differences between the assemblies and that large, megabase
9 regions of the haplotype can be reconstructed with assembly gaps, but without apparent error.

10 11 **Structural errors**

12
13 Despite the long lengths of many scaffold paths, for most assemblies the scaffold NG50 is
14 substantially larger than the scaffold path NG50, indicating that there were apparent errors that
15 broke scaffolds into smaller sets of scaffold paths. To analyse these errors we continued our
16 graph analysis, defining a number of subgraph types to represent them, which we formally define
17 in the methods. These subgraph definitions include erroneous intra and inter chromosomal joins,
18 insertions, deletions, simultaneous insertion and deletions and insertions at the ends of
19 assembled sequences. Table 5 (and Supplementary Table 3) shows the numbers of structural
20 errors for each assembly; Supplementary Figures 20 and 21 show structural errors across the
21 haplotypes. Many assemblies do not have categories of error to which they are particularly prone,
22 but a few do. In these cases there may be a systematic bias in the operation of the programs that
23 generated them or in the way that we interpreted them.

24
25 Insertion and deletion (indel) structural errors involve, respectively, the addition and removal of
26 contiguous run of bases. In Supplementary Figures 22 and 23 we investigate the size distribution
27 of such errors, using both the described MSA and supporting alignments from the
28 progressiveMauve program (Darling et al. 2010). We find that in almost all cases the size
29 distribution of the segments of inserted and deleted bases follows an approximately geometrically
30 decreasing distribution.

31
32 We also searched for subgraphs in which the assembly created a haplotype to contamination
33 chimera, but did not find any such subgraphs. To investigate this surprising result we searched
34 for such chimeras using BLAST with relaxed parameters (Supplementary section 7.1.3 and
35 Supplementary Figure 24). Using this approach we found 56 assemblies with no such chimeras,
36 7 assemblies with 1 chimera, 1 assembly with 2 chimeras and one outlier assembly with 26

1 chimeras. In each case we verified that these chimeras were missed in the graph approach due
2 to the stringent MSA parameters.

3 4 **Long-range contiguity**

5
6 The MSA graph theoretic analysis we have described is local in nature and quite strict, in that it
7 has no notion of large-scale contiguity, and refuses to stitch together paths that would be joined,
8 but for a small error. We thus sought a method to analyse the larger scale contiguity between
9 pairs of separated points in the genome. Formally, for two positions x_i and x_j in a haplotype
10 chromosome x , such that $i < j$, if there exists two positions y_k, y_l in an assembly scaffold y such
11 that (1) y_k is in the same column as x_i , (2) y_l is in the same column as x_j and (3) $k < l$, we say y_k
12 and y_l are *correctly contiguous*. Pairs may be correctly contiguous but not necessarily covered by
13 the same contig path or scaffold path, and indeed there may be arbitrary numbers of assembly
14 errors between two correctly contiguous positions.

15
16 Figure 5 shows the proportion of correctly contiguous pairs as a function of the pairs' separation
17 distance for each assembly. Taken at a high level, in all assemblies the proportion of correctly
18 linked pairs monotonically decreases with separation distance. Therefore we take the separation
19 distance at the 50th percentile, termed the *correct contiguity 50 (CC50)* as an essentially sufficient
20 statistic. See Supplementary section 7.1.2 and Supplementary Figures 25 and 26 for supporting
21 BLAST based analysis.

22 23 **Annotation analysis**

24
25 Evolver maintains annotations for a number of classes of simulated sequence, including genes,
26 which Evolver models as having exons, introns and untranslated regions (UTRs), and conserved
27 non-coding elements. Additionally, while Evolver does not track the history of individual repeat
28 elements following their insertion, it maintains a library of mobile elements, and thus, using
29 RepeatMasker (v1.25 <http://www.repeatmasker.org>) and tandem repeats finder (v4.0,
30 <http://tandem.bu.edu/trf/trf.html>) with this library, we identified a subset of repetitive sequence
31 within the haplotypes.

32
33 Continuing the previous MSA analysis, we define a *perfect path* as a maximal subpath of a
34 haplotype thread that is isomorphic to a subpath of an assembly thread. For a given assembly the
35 corresponding set of perfect paths reflects the regions of the haplotypes that are perfectly
36 reconstructed. Unlike contig and scaffold paths, perfect paths are intolerant of haplotype

1 polymorphism, but give a well defined set of intervals within $\alpha_{1,2}$ for comparison with a set of
2 annotations. Table 6 shows for each assembly the proportion of each annotation type contained
3 within perfect paths.

4
5 Both haplotypes of the $\alpha_{1,2}$ genome contain 176 protein-coding genes, Supplementary Figure 27
6 show the distribution of their lengths; we find that only a small proportion (max 11% of bp, min 2%
7 of bp) of these full length transcripts are perfectly reconstructed by the assemblies. Conversely,
8 we find that in the best assemblies almost all exons and a high proportion of UTRs are perfectly
9 reconstructed, for example, 99% of bp in exons and 84% of bp in UTRs of the BGI assembly. We
10 also find that most perfectly reconstructed genes are intron-less (data not shown), the assemblies
11 therefore fail mostly to reconstruct introns perfectly. To further characterise this failure we used
12 tBLASTN (Altschul et al. 1990) (see Supplementary section 7.1.4) to align the spliced transcripts
13 of α_1 (without introns and UTRs) to the scaffolds of the assemblies, counting a match if it
14 included 95% of the given transcript, see last column of Table 4 labelled *genic correctness* and
15 Supplementary Table 1. This more tolerant analysis reveals that in the best assemblies the
16 majority of exon chains are reconstructed contiguously (in the correct order and orientation) within
17 single scaffolds, e.g. the Broad assembly has 107 spliced transcripts (93.8% of bp) reconstructed
18 by this metric.

19
20 As expected, repeats were the least well reconstructed annotation types, with the best assembly,
21 BGI, reconstructing only 64% of repeats perfectly (see last column of Table 6). As these regions
22 are naturally difficult to align correctly within an MSA, we also performed BLAST based fragment
23 analysis, see Supplementary section 7.1.5, with similar results.

24
25 Finally, we also looked at conserved non-coding regulatory elements, which Evolver both models
26 and tracks. As these elements are short and relatively non-repetitive the majority (88% - 99% of
27 bp) were perfectly reconstructed by the assemblies.

28 **Substitution errors**

29
30
31 We have so far described assessments of structural correctness and contiguity, both overall and
32 for functional genic elements. We now turn to the assessment of somewhat orthogonal issues,
33 firstly by looking at base calling and then finally by analysing copy number errors.

34
35 Although we do not allow structural rearrangements within MSA blocks, blocks are tolerant of
36 substitutions. Let a (haplotype) column of aligned bases within a block that (1) contains a single

1 position from both haplotypes and (2) a single position from an assembly sequence be called
2 *valid*. We use these criteria because such columns unambiguously map a single assembled
3 sequence to a single position in the alignment of both haplotypes, while avoiding the issues of
4 paralogous alignment and multiple counting. We distinguish two types of valid columns: (1)
5 *homozygous columns*, those containing the same base pair from both haplotypes and (2)
6 *heterozygous columns*, those containing distinct base pairs from each haplotype. We also initially
7 considered columns that contain one but not both haplotypes, but found that the numbers of such
8 columns that we could consider reliably aligned was not sufficient for us to confidently compute
9 statistics.

10
11 Assemblers are free to use IUPAC ambiguity characters to call bases. To allow for this we use a
12 bit-score to score correct but ambiguous matches within valid columns (see methods). We say
13 there has been a *substitution error* if the position in the assembly sequence has an IUPAC
14 character that does not represent either of the haplotypes' base pair(s).

15
16 Some of the substitution errors that we observe are likely due to misalignments. These can occur
17 due to edge wander (Holmes and Durbin 1998), or the larger scale misalignment of an assembled
18 sequence to a paralog of its true ortholog. The sum of substitution errors over all valid columns is
19 therefore an upper bound on the substitution errors within valid columns. To obtain a higher
20 confidence set of substitution errors we select a subset of valid columns that meet the following
21 requirements: (1) Are part of blocks of at least 1 kb in length, avoiding errors within short indels,
22 (2) are not within 5 positions of the start and end of the block, avoiding edge wander, and (3) are
23 within blocks with 98% or higher sequence identity, ensuring the alignments are unlikely to be
24 paralogous. The sum of substitution errors within these high confidence valid columns represents
25 a reasonable lower bound of the number of substitution errors within valid columns.

26
27 Figure 6 (also Supplementary Figure 29 and Supplementary Tables 4 and 5) show, as might be
28 expected, that there are, in general, proportionally more errors made in heterozygous columns
29 than homozygous columns, though there are naturally much fewer overall heterozygous positions,
30 for example the WTSI-S and CRACS assemblies made no heterozygous errors. We find a strong
31 correlation between error rates in heterozygous and homozygous columns, with the exceptions of
32 the Broad and BCCGSC assemblies, which have proportionally higher rates of heterozygous
33 errors. The Broad result is explained by the large number of N ambiguity characters called at
34 heterozygous sites, which makes the number of errors per bit correspondingly higher, while the
35 BCCGSC result was due to a programmatic error in the assembler's pipeline that has since been
36 identified and resolved as a result of this analysis. Interestingly, we find considerable variation

1 between the programs in overall error rates. The strongest assembly, WTSI-S, makes one error
2 for every $15.3 \cdot 10^6$ to $2.94 \cdot 10^6$ correct bits, or approximately one every $7.7 \cdot 10^6$ to $1.49 \cdot 10^6$ bases,
3 while the weakest assembly, UCSF, makes an error for every $6.7 \cdot 10^3$ to $1.81 \cdot 10^4$ correct bits, or
4 approximately one in every $3.3 \cdot 10^3$ to $9.0 \cdot 10^3$ bases.

6 **Copy-number errors**

8 Within any haplotype column of the MSA the copy number of the simulated diploid genome can
9 be described by an interval $[min, max]$, where min is the minimum number of bases either of the
10 two haplotypes contributes and max is the maximum number of bases either of the two haplotype
11 contributes. To establish if assemblies were producing too many or too few copies of the
12 homologous positions within the two haplotypes we looked for haplotype columns where the copy
13 number of the assembly lay outside of this copy number interval. There are two possibilities,
14 either the number of copies in the assembly is less than min, in which case there is a deficiency
15 in the copy number, or the number of copies in the assembly is greater than max, in which case
16 there is an excess in the copy number.

18 Figure 7 (also Supplementary Figure 30) shows the proportions of haplotype columns with copy
19 number excesses and deficiencies. Again, to address contributions made by alignment errors we
20 choose to produce an upper and lower bound on these proportions. The lower bound is taken
21 over all haplotype columns in the alignments, while the upper bound is computed over only
22 haplotype columns that are part of blocks of at least 1 kb in length.

24 The deficiencies are dominated by cases where the assembly is not present; therefore copy
25 number deficiency is closely correlated with coverage. Unfortunately, there are not a sufficient
26 number of cases where the assembly is present but the copy number is deficient so that we may
27 make reliable inferences about this interesting category. This appears to be a consequence of the
28 genome simulation lacking sufficient numbers of recent duplications, and may be an indication
29 that the genome simulation is somewhat unrealistic, as other authors (Worley and Gibbs 2010)
30 (Alkan et al. 2011) have discussed that recent segmental duplications cause substantial problems
31 for assemblies generated with short reads.

33 We find that there are substantial numbers of copy number excesses, such that generally the
34 number of excesses was larger than the number of deficiencies. We find that excesses do not
35 correlate particularly well with deficiencies, particularly for programs with extremes of deficiency
36 or excess. We do find, however, that excesses correlate well with input assembly size (data not

1 shown). The best assembly, EBI, has excesses in between 0.0521% and 0.752% of haplotype
2 columns, while the least assembly, nABYSS has excesses in between 30.8% and 33.5% of
3 haplotype columns.

4 5 **DISCUSSION**

6
7 We have used simulation to create a novel benchmark dataset for de novo assembly. We have
8 evaluated a previously unprecedented 41 different assemblies from 17 different groups, making it
9 the largest short read de novo assembly evaluation to date. In summary, we have assessed
10 coverage, the lengths of consistent contig and scaffold paths, structural errors, long-range
11 contiguity, the assembly of specific annotated regions, including genes and repeats, base calling
12 errors and copy number errors; Table 7 conveniently summarises these evaluation metrics. This
13 benchmark dataset is freely available online at <http://www.assemblathon.org/> and is
14 supplemented by code that can take new assemblies and amalgamate the new result into the
15 analysis we present here. It is our hope that this standard will assist the assembly community
16 when introducing new methods by providing a large set of metrics and methods with which to
17 compare.

18
19 Given the degree of polymorphism within the $\alpha_{1,2}$ genome the haplotype aware evaluations
20 proved critical to the assessment. For example, the haplotype aware path analysis demonstrates
21 that methods are able to reconstruct multiple megabases, with scaffold breaks, essentially
22 perfectly. We chose to treat switches between the haplotypes of $\alpha_{1,2}$ permissively because the
23 assemblers were not asked to reconstruct the two haplotypes, but rather to produce a consensus
24 reference of the two. It is an open question if, with this dataset or one like it, an assembly could
25 produce phased variants of each scaffold. In section 7.3 of the supplement we test if there was
26 any evidence of teams phasing single nucleotide polymorphisms (SNPs) or structural variants by
27 preferentially choosing one haplotype, but we do not find convincing evidence for either, apart
28 from that inadvertently caused by a bias in the simulated reads (see Methods section: problems
29 with the error model used in Assemblathon 1).

30
31
32 Table 3 shows the rankings of each of the featured assemblies for each of the described
33 assessments; additionally in Supplementary Figures 31 and 32 we assess correlations between
34 the logs of different metrics. Intuitively, one might expect the path analysis metrics and the
35 contiguity assessments to be correlated to one another and inversely correlated with structural

1 errors. Indeed, this intuition proves partially correct. Contiguity (CC50) and scaffold path NG50
2 are strongly correlated ($R^2 = 0.77$, $P < 0.001$), while structural errors are inversely correlated with
3 scaffold path NG50s, with one explaining about half the variance of the other ($R^2 = 0.48$, $P <$
4 0.001). However, contig path NG50 is only weakly correlated with scaffold path NG50 (CPNG50 –
5 SPNG50 $R^2 = 0.38$, $P < 0.001$) and contiguity (CPNG50 - CC50 $R^2 = 0.31$, $P < 0.01$), suggesting
6 that the scaffolding process is more important in producing accurate long scaffolds than the prior
7 contigging process.

8
9 Given the popularity and simplicity of N50 statistics, it is perhaps reassuring how well these
10 metrics correlate with the path and contiguity metrics (SN50 – CC50 $R^2 = 0.98$, $P < 0.001$; SN50 –
11 SPNG50 $R^2 = 0.74$, $P < 0.001$; CN50 – CPNG50 $R^2 = 0.64$, $P < 0.001$), suggesting that one may
12 usefully compare N50 measurements between assemblies, and not just between assemblies by
13 the same program. Interestingly, the genic correctness measure also correlates with all the N50
14 measures, and most strongly with the correct contiguity 50 ($R^2 = 0.98$, $P < 0.001$) and scaffold
15 N50 measures ($R^2 = 0.98$, $P < 0.001$).

16
17 We do not find that substitution errors and copy number errors correlate substantially with
18 anything else, except for a correlation between substitution errors and structural errors ($R^2 = 0.45$,
19 $P < 0.001$). This is perhaps unsurprising, given the orthogonal basis of these metrics to each
20 other and the other evaluations. Perhaps surprisingly, coverage does not correlate strongly with
21 other measures, and in particular not with contig or scaffold N50 statistics, suggesting such naïve
22 measures are not good proxies for coverage.

23
24
25 Table 3 highlights that while the best assemblies are stronger in most categories than the
26 weakest assemblies, all the assemblies have areas in which they can improve relative to their
27 peers, if at a trade-off cost in other categories. For example, the BGI assembly, while having the
28 largest contig path NG50, has only the sixth largest scaffold path NG50, which is more than 4
29 times smaller than the strongest method in this category (WTSI-S), suggesting that its scaffolding
30 could be improved. Conversely, the WTSI-S and DOEJGI assemblies had large contiguity (CC50)
31 and scaffold path (SPNG50) measures and low numbers of structural errors, but relatively short
32 contig paths (CPNG50), suggesting their contigging could be made more aggressive, though
33 possibly with a corresponding increase in structural errors.

34

1 We have demonstrated in simulation that the best current sequence assemblers can reconstruct
2 at high coverage and with good accuracy large sequences of a substantial de novo genome. This
3 is concordant with other recent work that suggests that short read sequencing is becoming
4 competitive with capillary sequencing (Maccallum et al. 2009) (Gnerre et al. 2011). MacCallum et
5 al looked at five microbial genomes with sizes ranging from 2.8 Mb (*Staphylococcus aureus*) to
6 39.2 Mb (*Neurospora crassa*) and determined that with data from two paired libraries that the
7 ALLPATHS 2 program was able to produce assemblies with qualities that exceeded draft
8 assemblies using Sanger methods. Gnerre et al sequenced two genomes; a human cell line
9 (GM12878) and a mouse strain (C57BL/6J female) and assembled them using the ALLPATHS-
10 LG program. The authors found that with Illumina reads of 100x coverage in four library types that
11 their assemblies neared capillary sequencing quality in completeness, long range connectivity,
12 contiguity and accuracy.

13

14 There are a number of important limitations with the current work. Firstly, the use of simulation
15 makes it hard to know how applicable these results are to any other dataset; though this is
16 arguably true of any dataset, the simulation's limitations, in particular the noted issues with the
17 read simulation and with the low repeat content of the genome, likely influence the results.
18 Secondly, the limited size of the simulated genome means that some of the strategies employed
19 here may not work as effectively, or at all, on larger vertebrate scale genome datasets. Finally, as
20 our results derive from a single dataset in which no attempt was made to measure the variance of
21 our various metrics, it is questionable how reliable our measurements are. To address these
22 issues a second competitive evaluation, Assemblathon 2, is now underway.

23

24 Given the scale of challenges in making assessments, and to avoid fragmentation, we suggest
25 that Assemblathon 2 continue to focus on the assessment of complete pipelines, rather than
26 attempting to assess individual pipeline components, and that it continue to rely on the individual
27 assembly teams to compute their own assemblies, despite this making it difficult to compare the
28 computational requirements of the pipelines, given the self reported nature of such data and the
29 heterogeneous equipment upon which the assemblies are computed.

30

31 However, we conclude by making three distinguishing suggestions for Assemblathon 2 that would
32 sufficiently expand its scope from this initial competition. Firstly, it should feature at least one
33 mammalian genome scale data set, to test the scaling of the assembly pipelines. Secondly, it
34 should feature real data, to compare with the simulation results presented in this competition; this
35 may necessitate the use of a different set of evaluation metrics, where the "correct" answer is

1 unknown. Thirdly, it should be expanded to include other sequencing technologies, so that a
2 better comparative, unbiased understanding of available sequencing technologies can be made.

3 4 **METHODS**

5 6 **Genome simulation**

7
8 The Evolver simulation was managed by a set of scripts (Earl D, Paten B, Haussler D,
9 <https://github.com/dentearl/evolverSimControl/>), which control the execution of Evolver and allow
10 a general phylogeny to be simulated.

11
12 As well as a starting sequence, Evolver also requires a set of annotations that are used assign
13 sequences to element types that undergo differential evolution simulation. The following
14 annotations were used: UCSC Genes, UCSC Old Genes, CpG Islands, Ensembl Genes and
15 MGC Genes from the UCSC table browser (Fujita et al. 2011). The root genome was then
16 coupled with parameters and a mobile element library provided by Arend Sidow (pers. comm.) to
17 form the Evolver in-file set for the simulation.

18
19 Evolver proceeds iteratively by a series of discrete steps. We used an Evolver step length of 0.01
20 substitutions per site, meaning the initial branch length of 0.4 substitutions per site (~200 my
21 (Fujita et al. 2011) (Hedges et al. 2006)) from the root node to the most recent common ancestor
22 (MRCA) of the final leaf genomes node consisted of 40 separate Evolver cycles. The lineages
23 leading from the MRCA to α and β descend for a distance of 0.1 (~50 my) substitutions per site
24 in 10 Evolver cycles. The final splits into the lineages leading to the leaf genomes were each
25 performed in 1 Evolver cycle of 0.002 substitutions per site (~1 my), with parameters scaled
26 appropriately. An alignment between the α_1 and α_2 haplotypes is available on the project website
27 (<http://compbio.soe.ucsc.edu/assemblathon1/>).

28 29 **Read Simulation**

30
31 Prior to writing our own read simulator we considered several pre-existing tools. We first
32 considered wgsim (Li et al. 2010). Unfortunately, this program does not model mate-pair Illumina
33 reads, and it models error rates uniformly across the sequence. We note that this error rate
34 limitation is removed in dwgsim (<http://sourceforge.net/apps/mediawiki/dnaa/>). However, dwgsim
35 does not model chimeric mate-pair reads or paired-end contamination, which we wished to model.
36 We contacted Illumina and requested their in-house programs for simulating reads. The Illumina

1 software package was capable of modelling chimeric mate-pair reads, and it modelled error rates
2 by copying quality strings from a user supplied file of Illumina reads. Unfortunately, this method
3 did not allow us to model different error rates conditioned on different underlying bases, which we
4 felt was important. We also considered several other software packages for modelling Illumina
5 style reads, including metasim (Richter et al. 2008), PEMer (Korbel et al. 2009), ReSeqSim (Du et
6 al. 2009), SimNext (<http://evolution.sysu.edu.cn/english/software/simnext.htm>), Flux Simulator
7 (<http://flux.sammeth.net/index.html>), and Mason (part of the SeqAn package (Döring et al. 2008)),
8 all of which lack one or more of the criteria we desired.

9
10 Given these findings we wrote our own simulator, which combined the capabilities of the Illumina
11 supplied software to model chimeric mate-pair reads as well as standard paired-end reads, with
12 our own position and reference-base-specific empirical error model trained on Illumina data.

14 **Read Sampling Strategy**

15
16 For read sampling we employed two separate methods, one for mate-pair libraries and the other
17 for paired end libraries. Reads were first sampled uniformly across each sequence. Coverage
18 depth was kept approximately uniform by weighting the number of reads sampled from each
19 sequence by its length. Read fragments were sampled from either strand with equal probability.
20 Duplicates were produced with some probability before the error was applied to the reads. See
21 Supplementary Figure 33 for a density map of read depth across the haplotypes.

23 **Paired-End Sampling**

24
25 Illumina paired end sampling was the most straightforward strategy to simulate. It involved
26 randomly selecting fragments in the 150-500 base pair range uniformly across the genome until
27 the desired coverage was met (specific sizes below). Fragment size was sampled from a normal
28 distribution with a specified mean and variance. The reads were oriented facing each other and
29 were sampled from either strand with equal probability. The following paired-end libraries were
30 generated:

- 31 • 200 bp insert +/- 20 standard deviation
 - 32 ○ 2x 100 bp
 - 33 ○ 22,499,731 read pairs (~40x coverage of the diploid sequence)
 - 34 ○ 0.01 probability of being a duplicate
- 35 • 300 bp insert +/- 30 standard deviation.
 - 36 ○ 2x 100 bp

- 1 ○ 22,499,731 read pairs (~40x coverage)
- 2 ○ 0.01 probability of being a duplicate

3

4 **Mate-Pair Sampling**

5

6 Illumina mate-pair library construction differs from paired-end library construction in that it
7 introduces several unique types of error into the reads. In reality these libraries are constructed by
8 attaching a chemical tag onto the ends of a long sequence fragment, typically in the range of 2-10
9 kb, after which the fragment is circularized. The circularised product is then further fragmented
10 into sizes typically within the 200-500 bp range, which is the upper limit on fragment lengths for
11 Illumina sequencing. Finally the resulting mixture is purified for fragments that contain the
12 chemical tag, so that DNA from near the ends of the original 2-10 kb loop are what ideally get
13 sequenced.

14

15 There are three common types of error introduced in the mate-pair library preparation process,
16 and we modelled two of them. Firstly, when the fragments are circularized, there is a chance that
17 a loop will be formed between two non-related long fragments, resulting in chimeric reads
18 between two unrelated parts of the genome. We did not model this type of error. Assuming that
19 the fragment is properly circularised, the second type error is produced when a fragment that
20 does not contain the chemical tag is mistakenly sampled. When this happens, the loop join is not
21 part of the fragment, and a paired-end style read with a short insert is mixed in with the rest of the
22 library. We did model this type of error, and varied the probability of its occurrence with each
23 mate-pair library. The final major source of error is created during the random fragmentation
24 process and results in the loop join position occurring in the middle of a read rather than between
25 the two reads. We modelled this by assuming a uniform distribution of loop join sites across a
26 sampled loop fragment, which resulted in chimeric reads as a function of the size of the
27 fragmented loop piece, and the length of the reads. For example shorter reads and longer loop
28 fragmentation pieces were less likely to result in a chimeric read. The following mate-pair libraries
29 were generated:

- 30 • 3 kb loop length +/- 300 standard deviation
 - 31 ○ 2x 100 bp
 - 32 ○ 500 bp loop fragmentation size +/- 50 bp
 - 33 ○ 0.2 probability of sampling a PE fragment rather than an MP fragment
 - 34 ○ 11,249,866 read pairs (~20x coverage)
 - 35 ○ 0.05 probability of being a duplicate
- 36 • 10 kb loop length +/- 1 kb standard deviation

- 1 ○ 2x 100 bp
- 2 ○ 500 bp loop fragmentation size +/- 50 bp
- 3 ○ 0.3 probability of sampling a PE fragment rather than an MP fragment
- 4 ○ 11,249,866 read pairs (~20x coverage)
- 5 ○ 0.08 probability of being a duplicate

6

7 **Base-Level Error Model**

8

9 We utilized an error model that is dependent on the position within the read and the underlying
10 reference base. To generate this model we assembled a human mitochondrial genome using
11 reads from an Illumina HiSeq run (http://www.illumina.com/systems/hiseq_2000.ilmn) with the
12 reference guided assembler MIA (Green et al. 2010). We then took that assembly and mapped all
13 reads back to it using BWA with default settings, to do a paired end mapping to the sequence.
14 We kept all alignments with a mapq quality score over 10. We then iterated through the alignment
15 and built an empirical distribution of Phred (Ewing et al. 1998) scores and the probabilities of
16 observing one of A, C, G, T or N given the reference base, the position in the read, and the
17 reported Phred quality score. The error model was therefore conditioned on the Phred score,
18 position, and reference base, and did not assume that the Phred scores were an accurate
19 representation of the underlying error rates.

20

21 **Problems With the Error Model Used in Assemblathon 1**

22

23 The error model used was appealingly simple but has limitations that should be understood.
24 Firstly, in generating the error model we omitted many reads that had an error rate that was too
25 high to confidently map to the assembled mitochondria. In the future this could partially be
26 overcome by using the PhiX control lane
27 (http://www.illumina.com/products/multiplexing_sequencing_primers_and_phix_control_kit.ilmn),
28 where one can confidently force the vast majority of the reads to map back to the PhiX 174
29 genome (NCBI accession NC_001422.1), and do not have to be as sensitive to false positive
30 alignments.

31

32 Secondly, since we used a flat naive prior on the distribution of Phred scores, when training our
33 empirical model there was, due to noise, a mixture of good and poor quality bases at the ends of
34 the reads. Since each position was treated independently, the distribution of Phred scores was
35 therefore likely not typical, resulting in the likely relative failure of assembler heuristics used to
36 trim strings of bad Phred scores at the ends of reads.

1

2 Thirdly, since we wrote the simulator following the general algorithmic flow of the wgsim read
3 simulator (Li et al. 2010), reads were randomised within haplotype chromosomes, but not
4 between haplotype chromosomes, resulting in reads from each haplotype and chromosome being
5 clustered together separately in the data. Thankfully, an investigation of phasing bias in
6 Supplementary section 7.3 shows only a couple of assemblies showed evidence of any bias that
7 could likely be attributable to this.

8

9 **Cactus Alignment Assessment**

10

11 **Alignment Generation**

12

13 The Cactus program starts by using the Lastz pairwise alignment program
14 (<http://www.bx.psu.edu/~rsharris/lastz/>) to generate a set of pairwise alignments between all the
15 input sequences, including intra-sequence alignments that arise from recent duplications. In the
16 adapted version of Cactus used for the Assemblathon, which we hence forth call Cactus-A, we
17 used the following parameters to Lastz, after discussion with the program's author: --step=10 --
18 seed=match12 --notransition --mismatch=2,100 --match=1,5 --ambiguous=iupac --nogapped --
19 identity=98. This ensured that the resulting pairwise alignments were ungapped (without indels),
20 of minimum length of 100 bp and with an identity (sequence similarity) of 98% or greater, in
21 concordance with the evolutionary distance between the haplotypes. Cactus-A uses these
22 alignments to build a “sparse map” of the homologies between a set of input sequences. Once
23 this sparse map is constructed, in the form of a Cactus graph (Paten et al. 2011), a novel
24 algorithm is used to align together sequences that were initially unaligned in the sparse map. To
25 prevent sequences that are not homologous from being aligned in this process we set the
26 alignment rejection parameter, called γ , to 0.2, to filter positions from being aligned that are not
27 likely to have very recently been diverged. The results of Cactus-A are stored as MAF files
28 (Blanchette et al. 2004), one for each assembly, these are available in the supplementary
29 material.

30

31 **Scaffold Gaps, Error Subgraphs and Scaffold Paths**

32

33 Let P be a sequence of block edges $((x_1, x_2), (x_3, x_4) \dots (x_{n-1}, x_n))$ in a thread (thus ignoring the
34 alternating adjacency edges $(x_2, x_3), (x_4, x_5)$ etc.) representing an assembled sequence in the
35 adjacency graph. The *ambiguity* of a sequence is equal to the number of wildcard characters that
36 it contains (denoted as N_s). Similarly, the ambiguity of a subsequence of P is equal to the

1 ambiguity of the subsequence of the assembly sequence it represents. The prefix ambiguity of $(x_i,$
 2 $x_j)$ is equal to the number of wildcard characters in the first 5 bases of the assembly sequence
 3 that (x_i, x_j) represents, orienting the sequence from x_i to x_j . The *approximate ambiguity* of a
 4 subsequence $Q = ((x_i, x_{i+1}), (x_{i+2}, x_{i+3}) \dots (x_{i+j-1}, x_{i+j}))$ is equal to the ambiguity of Q plus the prefix
 5 ambiguity of (x_{i-1}, x_{i-2}) and (x_{i+j+1}, x_{i+j+2}) , if these edges exist. By using approximate ambiguity
 6 rather than just ambiguity we allow for wobble in the alignment caused by edge wander (Holmes
 7 and Durbin 1998) when denoting a scaffold gap.

8
 9 We say a thread is *empty* if it represents a sequence of zero length, else we say it is *non-empty*.

10
 11 Let a maximal thread of inconsistent adjacency edges and block edges that do not contain
 12 haplotypes or bacterial contamination segments be called a *joining thread*. A joining thread
 13 represents an unaligned portion of an assembly sequence. A scaffold gap or error subgraph is
 14 defined by a joining thread incident at one or both ends with blocks that contain haplotype
 15 segments. We classify such joining threads as follows:

16 **(A)** If the joining thread is not attached to anything at one end (i.e. it terminates)
 17 (Figure 8A):

- 18 • If it has approximate ambiguity then we classify it as a scaffold gap.
- 19 • Else we classify it as a hanging insert error.

20 **(B)** If the joining thread is attached at each ends to blocks, a and b , containing
 21 haplotype segments:

- 22 • If a and b are connected by a thread containing haplotype segments
 23 (Figure 8B):

24 i. If the joining thread has approximate ambiguity then it is a
 25 scaffold gap.

26 ii. Else it is an indel (insertion/deletion) error:

- 27 1. If the joining thread is empty then it is a deletion error, by
 28 definition all haplotype paths between a and b must be non-
 29 empty.

1 center) subaward on NHGRI grant no. U01HG004695 to the European Bioinformatics Institute;
2 ENCODE DCC (data coordination center) NHGRI grant no. U41HG004568; Browser (Center for
3 Genomic Science) NHGRI grant no. P41HG002371; Gencode subaward on NHGRI grant no.
4 U54HG004555 to the Sanger Center; NCI 1U24CA143858-01; NIH HG00064; PTDC/BIA-
5 BEC/100616/2008; PTDC/EIA-EIA/100897/2008; the Fundacao para a Ciencia e Tecnologia;
6 National Natural Science Foundation of China (30725008; 30890032;30811130531; 30221004); a
7 National Basic Research Program of China (973 program no. 2011CB809200); the Chinese 863
8 program (2006AA02Z177; 2006AA02Z334; 2006AA02A302;2009AA022707); NSF, Major
9 Research Instrumentation grant DBI 0821263 (University of Georgia Georgia Advanced
10 Computing Resource Center).

1 REFERENCES

- 2 Aird D, Ross M, Chen W-S, Danielsson M, Fennell T, Russ C, Jaffe DB, Nusbaum C, Gnirke A.
3 2011. Analyzing and minimizing PCR amplification bias in Illumina sequencing libraries. *Genome*
4 *Biology* 12 (2), R18.
- 5 Alekseyev M, Pevzner P. 2009. Breakpoint graphs and ancestral genome reconstructions.
6 *Genome Research* 19 (5), 943-57.
- 7 Alkan C, Sajjadian S, Eichler E. 2011. Limitations of next-generation genome sequence assembly.
8 *Nature methods* 8 (1), 61-5.
- 9 Altschul S, Gish W, Miller W, Myers E, Lipman, D. 1990. Basic local alignment search tool.
10 *Journal of molecular biology* 215 (3), 403-10.
- 11 Batzoglou S, Jaffe D, Stanley K, Butler J, Gnerre S, Mauceli E, Berger B, Mesirov JP, Lander ES.
12 2002. ARACHNE: a whole-genome shotgun assembler. *Genome research* 12 (1), 177-89.
- 13 Bentley D. 2006. Whole-genome re-sequencing. *Current opinion in genetics & development* 16
14 (6), 545-52.
- 15 Bergeron A, Mixtacki J, Stoye J. 2006. A unifying view of genome rearrangements. *Algorithms in*
16 *Bioinformatics LNBI 4175 (WABI 2006)*, 163-173.
- 17 Bergeron A, Mixtacki J, Stoye J. 2006. On sorting by translocations. *Journal of computational*
18 *biology : a journal of computational molecular cell biology* 13 (2), 567-78.
- 19 Blanchette M, Kent W, Riemer C, Elnitski L, Smit A, Roskin K, Baertsch R, Rosenbloom K,
20 Clawson H, Green ED, Haussler D, Miller W. 2004. Aligning multiple genomic sequences with the
21 threaded blockset aligner. *Genome research* 14 (4), 708-15.
- 22 Butler J, Maccallum I, Kleber M, Shlyakhter I, Belmonte M, Lander E, Nusbaum C, Jaffe DB. 2008.
23 ALLPATHS: de novo assembly of whole-genome shotgun microreads. *Genome research* 18 (5),
24 810-20.
- 25 Chaisson M, Pevzner P. 2008. Short read fragment assembly of bacterial genomes. *Genome*
26 *research* 18 (2), 324-30.
- 27 Chaisson M, Brinza D, Pevzner P. 2009. De novo fragment assembly with short mate-paired
28 reads: Does the read length matter? *Genome research* 19 (2), 336-46.
- 29 Church D, Goodstadt L, Hillier L, Zody M, Goldstein S, She X, Bult CJ, Agarwala R, Cherry JL,
30 DiCuccio M, et al. 2009. Lineage-specific biology revealed by a finished genome assembly of the
31 mouse. *PLoS Biology* 7 (5), e1000112.
- 32 Clarke J, Wu, H.-C., Jayasinghe L, Patel A, Reid S, Bayley H. 2009. Continuous base
33 identification for single-molecule nanopore DNA sequencing. *Nature nanotechnology* 4 (4), 265-
34 70.

1 Colbourne JK, Pfrender ME, Gilbert D, Thomas WK, Tucker A, Oakley TH, Tokishita S, Aerts A,
2 Arnold GJ, Basu MK, et al. 2011. The Ecoresponsive Genome of *Daphnia pulex*. *Science* 331
3 6017, 555-561.

4 Döring A, Weese D, Rausch T, Reinert K. 2008. SeqAn an efficient, generic C++ library for
5 sequence analysis. *BMC Bioinformatics* 9.

6 Darling A, Mau B, Perna N. 2010. progressiveMauve: multiple genome alignment with gene gain,
7 loss and rearrangement. *PloS one* 5 (6), e11147.

8 Dohm J, Lottaz C, Borodina T, Himmelbauer H. 2007. SHARCGS, a fast and highly accurate
9 short-read assembly algorithm for de novo genomic sequencing. *Genome research* 17 (11),
10 1697-706.

11 Du J, Bjornson R, Zhang Z, Kong Y, Snyder M, Gerstein M. 2009. Integrating sequencing
12 technologies in personal genomics: optimal low cost reconstruction of structural variants. *PLoS*
13 *Computational Biology* 5 (7).

14 Dunham A, Matthews LH, Burton J, Ashurst JL, Howe KL, Ashcroft KJ, Beare DM, Burford DC,
15 Hunt SE, Griffiths-Jones S, et al. 2004. The DNA sequence and analysis of human chromosome
16 13. *Nature* 428 6982, 522-8.

17 Eid J, Fehr A, Gray J, Luong K, Lyle J, Otto G, Peluso P, Rank D, Baybayan P, Bettman B, et al.
18 2009. Real-time DNA sequencing from single polymerase molecules. *Science* 323 5910, 133-8.

19 Ewing B, Hillier L, Wendl M, Green P. 1998. Base-calling of automated sequencer traces using
20 phred. I. Accuracy assessment. *Genome research* 8 (3), 175-85.

21 Fujita PA, Rhead B, Zweig AS, Hinrichs AS, Karolchik D, Cline MS, Goldman M, Barber GP,
22 Clawson H, Coelho A, et al. 2011. The UCSC Genome Browser database: update 2011. *Nucleic*
23 *acids research* 39 (Database issue), D876-82.

24 Gnerre S, Maccallum I, Przybylski D, Ribeiro FJ, Burton J, Walker BJ, Sharpe T, Hall G, Shea TP,
25 Sykes S, et al. 2011. High-quality draft assemblies of mammalian genomes from massively
26 parallel sequence data. *Proceedings of the National Academy of Sciences of the United States of*
27 *America* 108 (4), 1513-8.

28 Green R, Krause J, Briggs A, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W, Fritz
29 MH-Y, et al. 2010. A draft sequence of the Neandertal genome. *Science* 328 5979, 710-22.

30 Harris B, Miller, W. (n.d.). Retrieved from <http://www.bx.psu.edu/~rsharris/lastz/>

31 Hedges S, Dudley J, Kumar S. 2006. TimeTree: a public knowledge-base of divergence times
32 among organisms. *Bioinformatics* 22 (23), 2971-2.

33 Hernandez D, François P, Farinelli L, Osterås M, Schrenzel J. 2008. De novo bacterial genome
34 sequencing: millions of very short reads assembled on a desktop computer. *Genome research* 18
35 (5), 802-9.

- 1 Holmes I, Durbin R. 1998. Dynamic programming alignment accuracy. *Journal of computational*
2 *biology : a journal of computational molecular cell biology* 5 (3), 493-504.
- 3 Homer, N. (n.d.). *dwgsim*. Retrieved from <http://sourceforge.net/apps/mediawiki/dnaa/>
- 4 Hubisz M, Lin M, Kellis M, Siepel A. 2011. Error and error mitigation in low-coverage genome
5 assemblies. *PLoS ONE* 6 (2), e17034.
- 6 Huson D, Halpern A, Lai Z, Myers E, Reinert K, Sutton G. 2001. Comparing assemblies using
7 fragments and mate-pairs. In O. Gascuel, B. Moret (Ed.), *Proc. Workshop Algorithms in*
8 *Bioinformatics*. 2149, pp. 294-306. Aarhus, Denmark: Springer-Verlag.
- 9 Illumina. (n.d.). Illumina HiSeq 2000. Retrieved from
10 http://www.illumina.com/systems/hiseq_2000.ilmn
- 11 Illumina. (n.d.). Illumina PhiX. Retrieved from
12 http://www.illumina.com/products/multiplexing_sequencing_primers_and_phix_control_kit.ilmn
- 13 Jeck WR, Reinhardt JA, Baltrus DA, Hickenbotham MT, Magrini V, Mardis ER, Dangl JL, Jones
14 CD. 2007. Extending assembly of short DNA sequences to handle error. *Bioinformatics* 23 (21),
15 2942-4.
- 16 Kelley D, Schatz M, Salzberg S. 2010. Quake: quality-aware detection and correction of
17 sequencing errors. *Genome biology* 11 (11), R116.
- 18 Kent W, Haussler, D. (2001, Jan 1). *GigAssembler: An Algorithm for the Initial Assembly of the*
19 *Human Genome*. Technical Report UCSC-CRL-00-17 .
- 20 Korb J, Abyzov A, Mu X, Carriero N, Cayting P, Zhang Z, Snyder M, Gerstein MB. 2009.
21 PEMer: a computational framework with simulation-based error models for inferring genomic
22 structural variants from massive paired-end sequencing data. *Bioinformatics* 10 (2).
- 23 Kurtz S, Narechania A, Stein, J. C., Ware D. 2008. A new method to compute K-mer frequencies
24 and its application to annotate large repetitive plant genomes. *BMC Genomics* 9 (517),
25 doi:10.1168/1471-2164-9-517.
- 26 Lander E, Linton L, Birren B, Nusbaum C, Zody M, Baldwin J, Devon K, Dewar K, Doyle M,
27 FitzHugh W, et al. 2001. Initial sequencing and analysis of the human genome. *Nature* 409 6822,
28 860-921.
- 29 Li R, Fan W, Tian G, Zhu H, He L, Cai J, Huang Q, Cai Q, Li B, Bai Y, et al. 2010. The sequence
30 and de novo assembly of the giant panda genome. *Nature* 463 7279, 311-7.
- 31 Li R, Zhu H, Ruan J, Qian W, Fang X, Shi Z, Li Y, Li S, Shan G, Kristiansen K, et al. 2010. De
32 novo assembly of human genomes with massively parallel short read sequencing. *Genome*
33 *research* 20 (2), 265-72.
- 34 Lin Y, LI J, Shen H, Zhang L, Papasian, C. J., Deng, H.-W. 2011. Comparative studies of de novo
35 assembly tools for next-generation sequencing technologies. *Bioinformatics* 27 (15), 2031-2037.

- 1 Lindblad-Toh K, Wade C, Mikkelsen T, Karlsson E, Jaffe D, Kamal M, Clamp M, Chang JL,
2 Kulbokas EJ, Zody MC, et al. 2005. Genome sequence, comparative analysis and haplotype
3 structure of the domestic dog. *Nature* 438 7069, 803-19.
- 4 Liu Y, Qin X, Song, X.-Z., Jiang H, Shen Y, Durbin KJ, Lien S, Kent MP, Sodeland M, Ren Y, et al.
5 2009. *Bos taurus* genome assembly. *BMC Genomics* 10, 180.
- 6 Locke DP, Hillier LW, Warren WC, Worley KC, Nazareth LV, Muzny DM, Yang S-P, Wang Z,
7 Chinwalla AT, Minx P, et al. 2011. Comparative and demographic analysis of orang-utan
8 genomes. *Nature* 469 7331, 529-33.
- 9 Maccallum I, Przybylski D, Gnerre S, Burton J, Shlyakhter I, Gnirke A, Malek J, Mckernan K,
10 Ranade S, Shea TP, et al. 2009. ALLPATHS 2: small genomes assembled accurately and with
11 high continuity from short paired reads. *Genome Biology* 10 (10), R103.
- 12 Margulies M, Egholm M, Altman WE, Attiya S, Bader JS, Bemben LA, Berka J, Braverman MS,
13 Chen Y-J, Chen Z, et al. 2005. Genome sequencing in microfabricated high-density picolitre
14 reactors. *Nature* 437 7057, 376-80.
- 15 Meader S, Hillier L, Locke D, Ponting C, Lunter G. 2010. Genome assembly quality: assessment
16 and improvement using the neutral indel model. *Genome research* 20 (5), 675-84.
- 17 Medvedev P, Brudno M. 2009. Maximum likelihood genome assembly. *Journal of computational*
18 *biology : a journal of computational molecular cell biology* 16 (8), 1101-16.
- 19 Metzker, M. L. 2010. Sequencing technologies — the next generation. *Nature Reviews Genetics*
20 11 (1), 31-46.
- 21 Miller J, Koren S, Sutton G. 2010. Assembly algorithms for next-generation sequencing data.
22 *Genomics* 95 (6), 315-27.
- 23 Ming R, Hou S, Feng Y, Yu Q, Dionne-Laporte A, Saw JH, Senin P, Wang W, Ly BV, Lewis KLT,
24 et al. 2008. The draft genome of the transgenic tropical fruit tree papaya (*Carica papaya*
25 Linnaeus). *Nature* 452 7190, 991-6.
- 26 Monument. (n.d.). Retrieved from <http://www.irsia.fr/symbiose/people/rchkhi/monument.html>
- 27 Mullikin J, Ning Z. 2003. The phusion assembler. *Genome research* 13 (1), 81-90.
- 28 Myers EW. 2005. The fragment assembly string graph. *Bioinformatics* 21 Suppl 2, ii79-85.
- 29 Myers EW. (1995, Jan 1). Toward simplifying and accurately formulating fragment assembly.
30 *Journal of Computational Biology* .
- 31 Myers EW, Sutton GG, Delcher AL, Dew IM, Fasulo DP, Flanigan MJ, Kravitz SA, Mobarry CM,
32 Reinert KH, Remington KA,, et al. 2000. A whole-genome assembly of *Drosophila*. *Science* 287
33 5461, 2196-204.
- 34 Narzisi G, Mishra B. 2011. Comparing De Novo Genome Assembly: The Long and Short of It.
35 *PLoS ONE* 6 (4), e19175.
- 36 Newbler. (n.d.). Retrieved from <http://www.454.com/>

- 1 OligoZip. (n.d.). Retrieved from <http://linux1.softberry.com/berry.phtml?topic=OligoZip>
- 2 Pandey V, Nutter R, Prediger E. 2008. Applied Biosystems SOLiD System: Ligation-Based
3 Sequencing. *Next Generation Genome Sequencing: Towards Personalized Medicine* 29-41.
- 4 Parra G, Bradnam K, Ning Z, Keane T, Korf I. 2009. Assessing the gene space in draft genomes.
5 *Nucleic acids research* 37 (1), 289-97.
- 6 Paten B, Diekhans M, Earl D, St. John J, Ma J, Suh B, Haussler D. 2011. Cactus graphs for
7 genome comparisons. *J Comput Biol* (2011) vol. 18 (3) pp. 469-81
- 8 Paten B, Earl D, Nguyen N, Diekhans M, Zerbino D, Haussler D. 2011. Cactus: algorithms for
9 genome multiple sequence alignment. *Genome Res* (2011) vol. 21 (9) pp. 1512-28
- 10 Pevzner P, Tang H, Waterman M. 2001. An Eulerian path approach to DNA fragment assembly.
11 *Proceedings of the National Academy of Sciences of the United States of America* 98 (17), 9748-
12 53.
- 13 Phillippy A, Schatz M, Pop M. 2008. Genome assembly forensics: finding the elusive mis-
14 assembly. *Genome Biology* 9 (3), R55.
- 15 PHRAP. (n.d.). Retrieved from <http://www.phrap.org/>
- 16 Pop M, Salzberg, S. L. 2008. Bioinformatics challenges of new sequencing technology. *Trends in*
17 *Genetics* 24 (3), 142-149.
- 18 Pourmand N, Karhanek M, Persson, HH, Webb, CD, Lee, TH, Zahradnikova A, Davis RW. 2006.
19 Direct electrical detection of DNA synthesis. *PNAS* 103 (17), 6466-6470.
- 20 PRICE. (n.d.). Retrieved from <http://derisilab.ucsf.edu/software/price/index.html>
- 21 Richter D, Ott F, Auch A, Schmid R, Huson D. 2008. MetaSim: a sequencing simulator for
22 genomics and metagenomics. *PloS one* 3 (10), e3373.
- 23 Sanger F, Nicklen S, Coulson A. 1977. DNA sequencing with chain-terminating inhibitors.
24 *Proceedings of the National Academy of Sciences of the United States of America* 74 (12), 5463-
25 7.
- 26 SimNext. (n.d.). Retrieved from Laboratory of Evolution:
27 <http://evolution.sysu.edu.cn/english/software/simnext.htm>
- 28 Simpson J, Durbin R. 2010. Efficient construction of an assembly string graph using the FM-index.
29 *Bioinformatics* 26 (12), i367-73.
- 30 Simpson J, Wong K, Jackman S, Schein J, Jones S, Birol I. 2009. ABySS: a parallel assembler
31 for short read sequence data. *Genome research* 19 (6), 1117-23.
- 32 SOLiD. (n.d.). Retrieved from <http://www.appliedbiosystems.com>
- 33 St. John, J. (n.d.). Retrieved from <http://github.com/jstjohn/SimSeq>
- 34 Staden R, Beal K, Bonfield J. 2000. The Staden package, 1998. *Methods in molecular biology*
35 (Clifton, NJ) 132, 115-30.

1 Trapnell C, Salzberg SL, 2009. How to map billions of short reads onto genomes. Nature
2 Biotechnology 27 (5), 455-457.
3 Venter, JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, Smith HO, Yandell M, Evans
4 CA, Holt RA, et al. 2001. The sequence of the human genome. Science 291 5507, 1304-51.
5 Warren R, Sutton G, Jones S, Holt R. 2007. Assembling millions of short DNA sequences using
6 SSAKE. Bioinformatics 23 (4), 500-1.
7 Waterston RH, Lindblad-Toh K, Birney E, Rogers J, Abril JF, Agarwal P, Agarwala R, Ainscough
8 R, Alexandersson M, An P, et al. 2002. Initial sequencing and comparative analysis of the mouse
9 genome. Nature 420 6915, 520-62.
10 Worley K, Gibbs R. 2010. Genetics: Decoding a national treasure. Nature 463 7279, 303-4.
11 Zerbino D, Birney E. 2008. Velvet: algorithms for de novo short read assembly using de Bruijn
12 graphs. Genome research 18 (5), 821-9.
13 Zhang W, Chen J, Yang Y, Tang Y, Shang J, Shen B. 2011. A practical comparison of de novo
14 genome assembly software tools for next-generation sequencing technologies. PLoS ONE 6 (3),
15 e17915.

16

17 FIGURES

18

19 Figure 1 — The phylogeny of the simulated haploid genomes. The root genome derives from
20 human chromosome 13. The α_1 and α_2 haplotypes form the diploid genome from which we
21 generated reads. The β and β_2 haplotypes form a diploid out-group genome that was made
22 available to the assemblers.

23

24 Figure 2 — N50 statistics. Assemblies are sorted left to right in descending order by scaffold path
25 NG50. Data points for each assembly are slightly offset along the x-axis in order to show overlaps.

26

27

28 Figure 3 — An adjacency graph example demonstrating threads, contig paths and scaffold paths.
29 Each stack of boxes represents a block edge. The nodes of the graph are represented by the left
30 and right ends of the stacked boxes. The adjacency edges are groups of lines that connect the
31 ends of the stacked boxes. Threads are represented inset within the graph as alternating
32 connected boxes and coloured lines. There are three threads shown: (top to bottom) black, gray
33 and light gray. The black and gray threads represent two haplotypes, there are many alternative

1 haplotype threads that result from a mixture of these haplotype segments which are equally
 2 plausible given no additional information to de-convolve them. The light gray thread represents an
 3 assembly sequence. For the assembly thread consistent adjacencies are shown in solid red. The
 4 dashed red line between the right end of block g and the left end of block i represents a structural
 5 error (deletion). The dashed light gray line between the right end of block k and the left end of
 6 block m represents a scaffold gap, because the segment of the assembly in block n contains wild
 7 card characters. The example therefore contains three contig paths: (from left to right) blocks
 8 $a...g$ ACTGAAATCGGGACCCC; blocks i, j, k GGAAC; and block m CC. However the example
 9 contains only two scaffold paths because the latter two contig paths are concatenated to form one
 10 scaffold path.

11

12 Figure 4 — Assembly coverage along haplotype α_1 stratified by scaffold path length weighted
 13 overall coverage. The top 6 rows show density plots of annotations. CDS: coding sequence; UTR:
 14 untranslated region; NXE: non-exonic conserved regions within genes; NGE: non-genic
 15 conserved regions; island: CpG islands; repeats: repetitive elements. The remaining rows show
 16 the top ranked assembly from each group, sorted by scaffold path length weighted overall
 17 coverage. Each such row is a density plot of the coverage, with coloured stack fills used to show
 18 the length of scaffold paths mapped to a given location in the haplotype. For example, the left
 19 most light-orange block of the WTSI-S assembly row represents a region of haplotype α_1 that is
 20 almost completely covered by a scaffold path from the WTSI-S assembly greater than one
 21 megabase in length.

22

23 Figure 5 — The proportion of correctly contiguous pairs as a function of their separation distance.
 24 Each line represents the top assembly from each team. Correctly contiguous 50 (CC50) values
 25 are the lowest point of each line. The legend is ordered top to bottom in descending order of
 26 CC50. Proportions were calculated by taking 100 million random samples and binning them into
 27 2,000 bins, equally spaced along a \log_{10} scale, so that an approximately equal number of
 28 samples fell in each bin.

29

30 Figure 6 — Substitution (base) errors for the top assembly from each team. Top: substitution
 31 errors per correct bit within all valid columns, middle: substitution errors per correct bit within
 32 homozygous columns only, bottom: substitution errors per correct bit within heterozygous
 33 columns only. Assemblies are sorted from left to right in ascending order by the sum of

1 substitutions per correct bit. In each faceted plot each assembly is shown as an interval, giving
2 the upper and lower bounds on the numbers of substitution errors (see main text).

3

4 Figure 7 — Copy number errors for the top assembly from each team. Top: proportion of
5 haplotype containing columns with a copy number error, middle: proportion of haplotype
6 containing columns with an excess copy number error, bottom: proportion of haplotype containing
7 columns with an excess copy number error. Assemblies are sorted from left to right in ascending
8 order according to the proportion of haplotype containing columns with a copy number error. In
9 each faceted plot each assembly is shown as an interval, giving the upper and lower bounds on
10 the numbers of copy number errors (see main text).

11

12 Figure 8 — Scaffold gap and error subgraphs. Diagrams follow the format of Figure 3. The
13 rounded boxes represent extensions to the surrounding threads. Line ends not incident with the
14 edge of boxes represent the continuation of a thread unseen. In each diagram the right end of
15 block *a* and the left end of block *b* (if present) represent the ends of contig paths, the enclosed
16 gray thread represents the joining thread. The black thread represents a haplotype thread. The
17 gray thread represents either a haplotype or bacterial contamination thread. (A) Represents
18 (hanging) scaffold gaps and hanging insert errors. (B) Represents scaffold gaps and indel errors.
19 (C) Represents intra and inter chromosomal joining errors and haplotype to contamination joining
20 errors.

21

1 **TABLES**
2

1 Table 1 — Groups that submitted assemblies. The first 17 rows in the table correspond to entries
 2 submitted by participants in the competition. Assemblies with IDs beginning with “n,” (for naïve),
 3 were generated by organisers of the competition to demonstrate the performance of popular
 4 programs run with variations on their default parameters. * CSHL.1 used the β genome though
 5 that team’s top assembly, CSHL.2, which is referred to in the main paper as CSHL, did not.

6

7 Table 2 — Genome simulation statistics. (A) Event numbers are between the previous branch
 8 point and the named node. Mb: size of the genome in megabases; GC: percentage GC content;
 9 Reps: percent of the genome masked by the union of tandem repeats finder and RepeatMasker, *
 10 is the published value for chromosome 13 (Dunham et al. 2004); Reps 100mer: percent
 11 repetitiveness of the sequence and its reverse complement for 100-mers calculated with the
 12 tallymer tool (Kurtz et al. 2008); Chr: number of chromosomes; Subs: number of substitution
 13 events; Dels: number of deletion events; Inv: number of inversion events; Moves: number of
 14 translocations; Copy: number of DNA segmental duplications; Tandem: number of tandem repeat
 15 insertions; Chr Split: number of chromosome fission events; Chr Fuse: number of chromosome
 16 fusion events. (B) Differences between haplotypes α_1 and α_2 as determined by inspection of the
 17 Evolver pairwise alignment. SNPs: count of single nucleotide polymorphisms; Subs: count of
 18 substitutions, including SNPs; Σ Subs: sum of the lengths of all substitutions; Indels: count of
 19 insertion deletion events; Σ Indels: sum of the lengths of all insertion deletion events; Inv: the
 20 sum of number of inversions invoked in each of the α_1 and α_2 Evolver steps.

23

24 Table 3 — Rankings of the top assembly from each team in eight categories. For each category
 25 (listed below) all the received assemblies were ranked. The sum of the rankings from each
 26 category was then used to create an overall rank for the assemblies, the top (lowest number)
 27 ranked assembly from each group was then selected for inclusion in this manuscript. Numbers
 28 are ranks, with values shown in parentheses. Overall: sum of all rankings (possible range 8-160),
 29 CPNG50: Contig path NG50, SPNG50: Scaffold path NG50, Struct.: Sum of structural errors,
 30 CC50: length for which half of any two valid columns in the assembly are correct in order and
 31 orientation, Subs.: Total substitution errors per correct bit, Copy Num.: Proportion of columns with
 32 a copy number error, Cov. Tot.: Overall Coverage, Cov. Genic: Coverage within coding
 33 sequences.

34

1 Table 4 — Coverage statistics for the top assembly from each team. Hap Total: overall coverage,
2 Hap α_1 : percent coverage for Haplotype α_1 , Hap α_2 : percent coverage for Haplotype α_2 , Bac:
3 percent coverage of the bacterial contamination, Genic: percent coverage of the coding
4 sequences, Unmapped: number of unmapped bases, many corresponding to short contigs.

5

6

1 Table 5 — Structural error statistics for the top assembly from each team. Columns are defined in
2 the main text.

3

4

5 Table 6 — Inclusion of annotated features within perfect paths. Each annotation is represented as
6 a set of maximal non-overlapping intervals upon the haplotypes of $\alpha_{1,2}$. Each column represents
7 an annotation type, giving the number of bp contained within intervals of the type that are fully
8 contained within perfect paths, as a proportion of all bp in intervals of the type. Annotations from
9 left to right: Full length gene transcripts, exons, untranslated regions (UTRs), non-coding
10 conserved elements and repeats.

11

12

13

14 Table 7 — Summary of metrics used in the analysis.

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

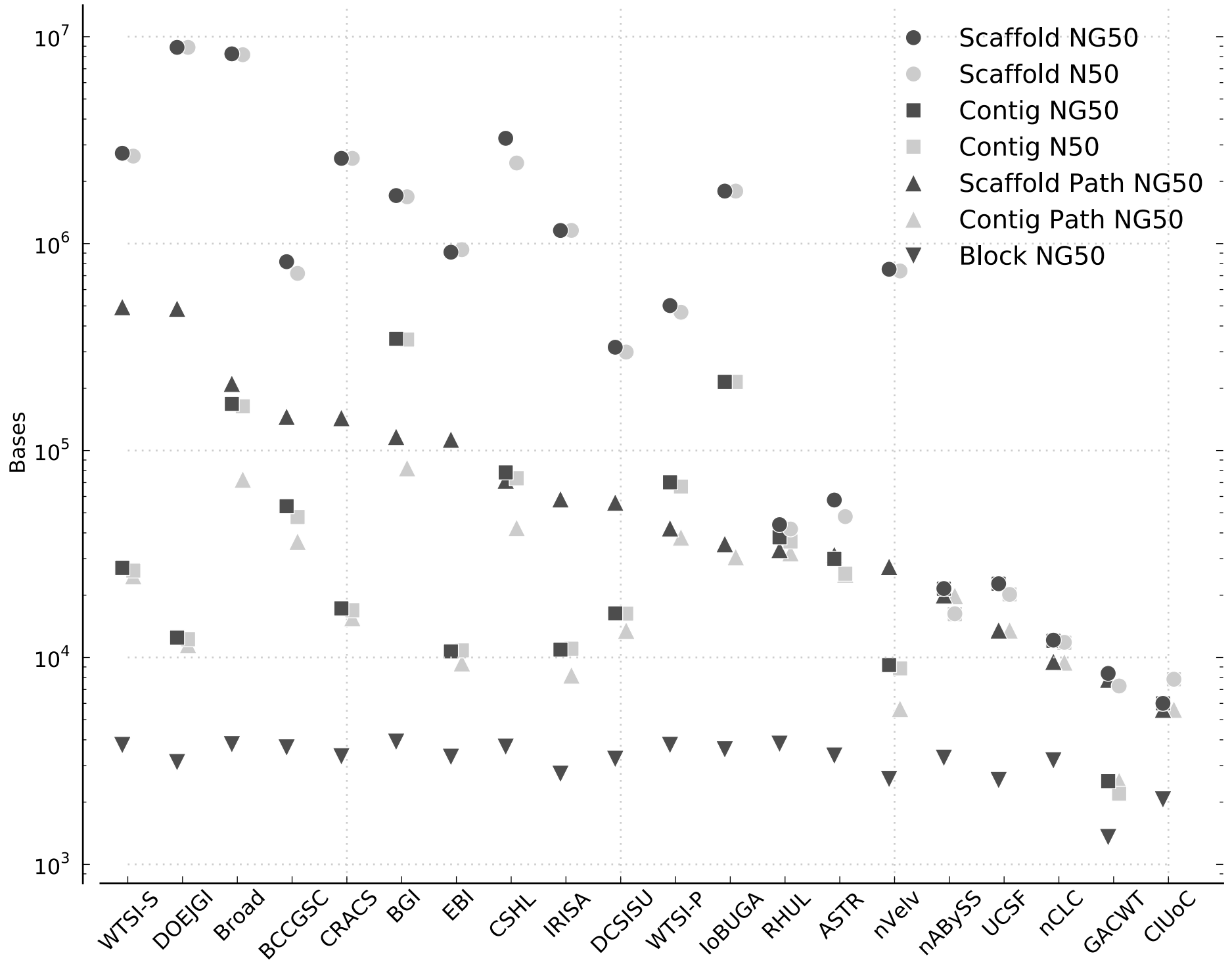
39

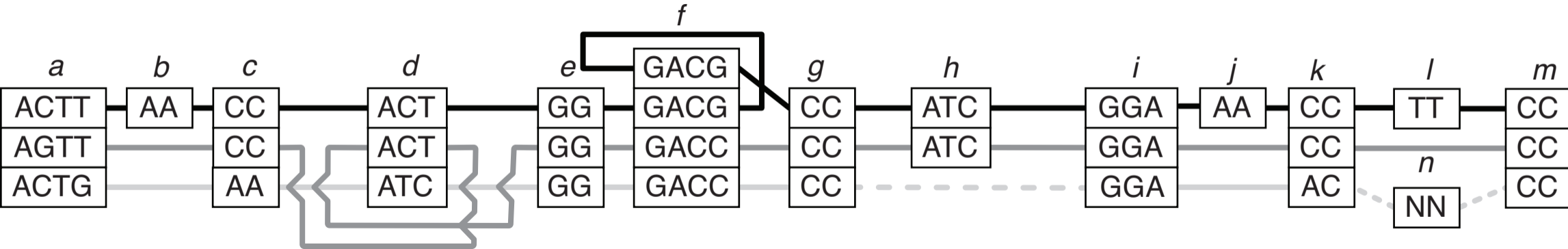
40

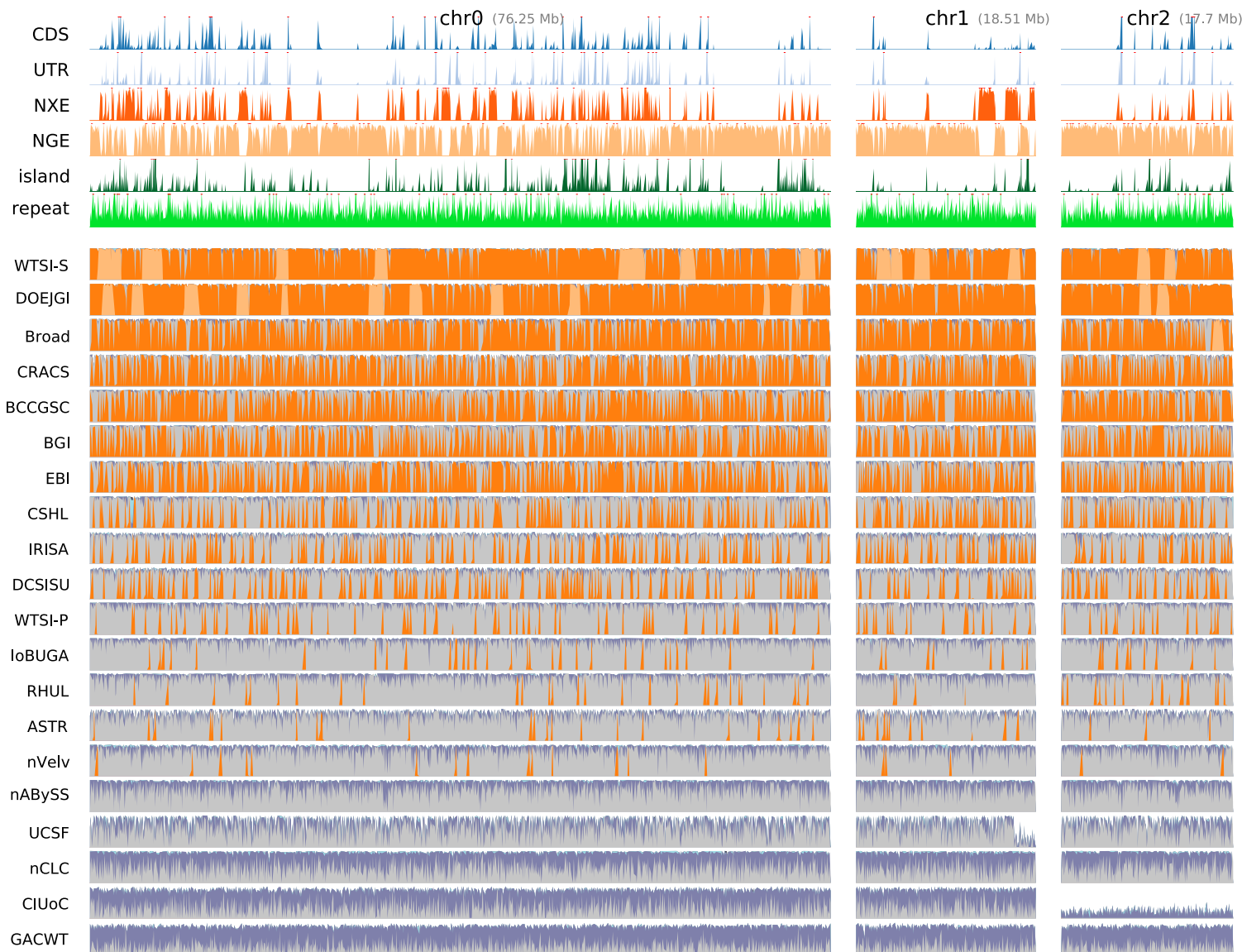
41



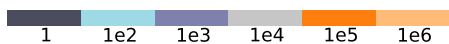
N50 Statistics



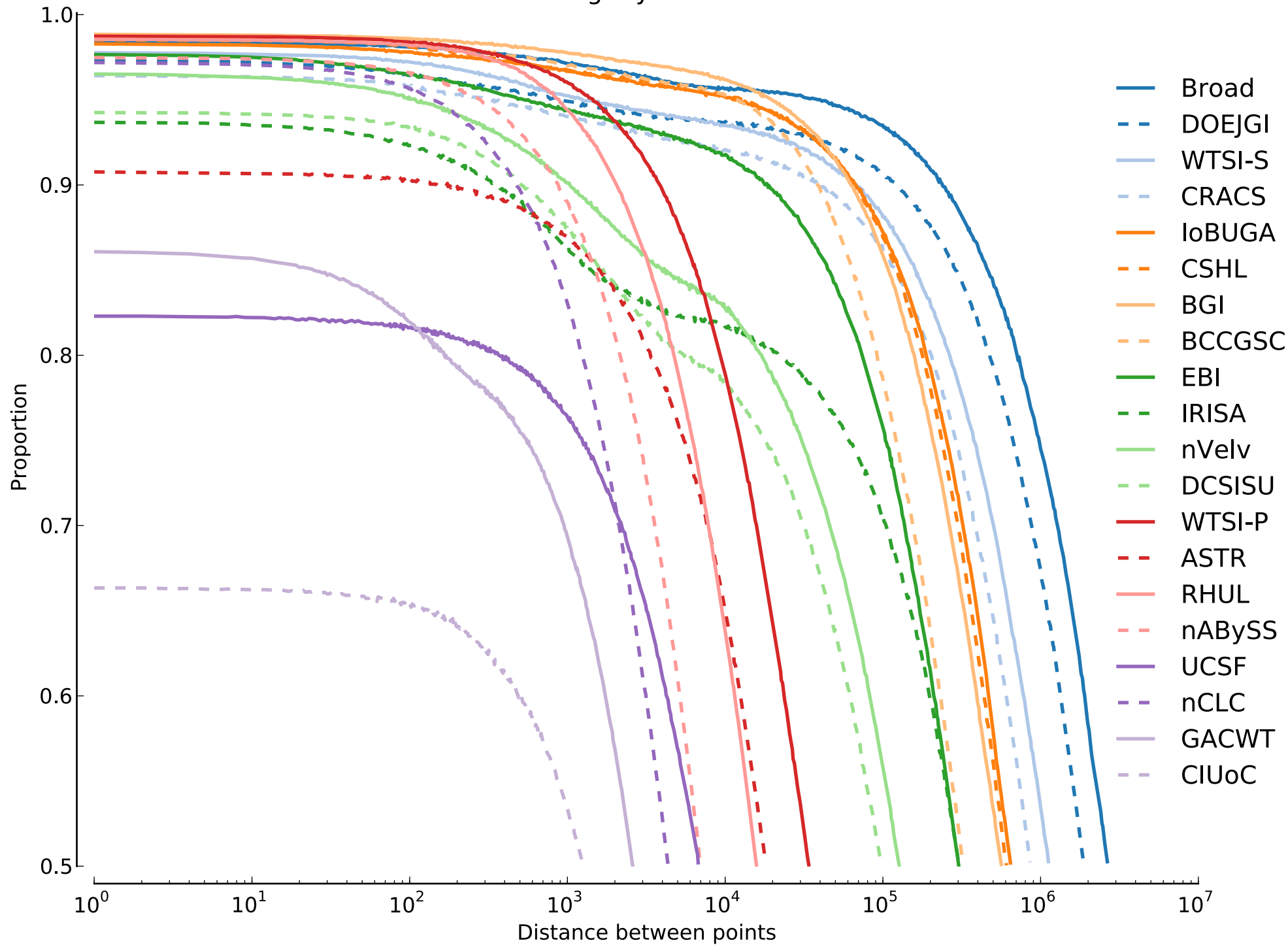




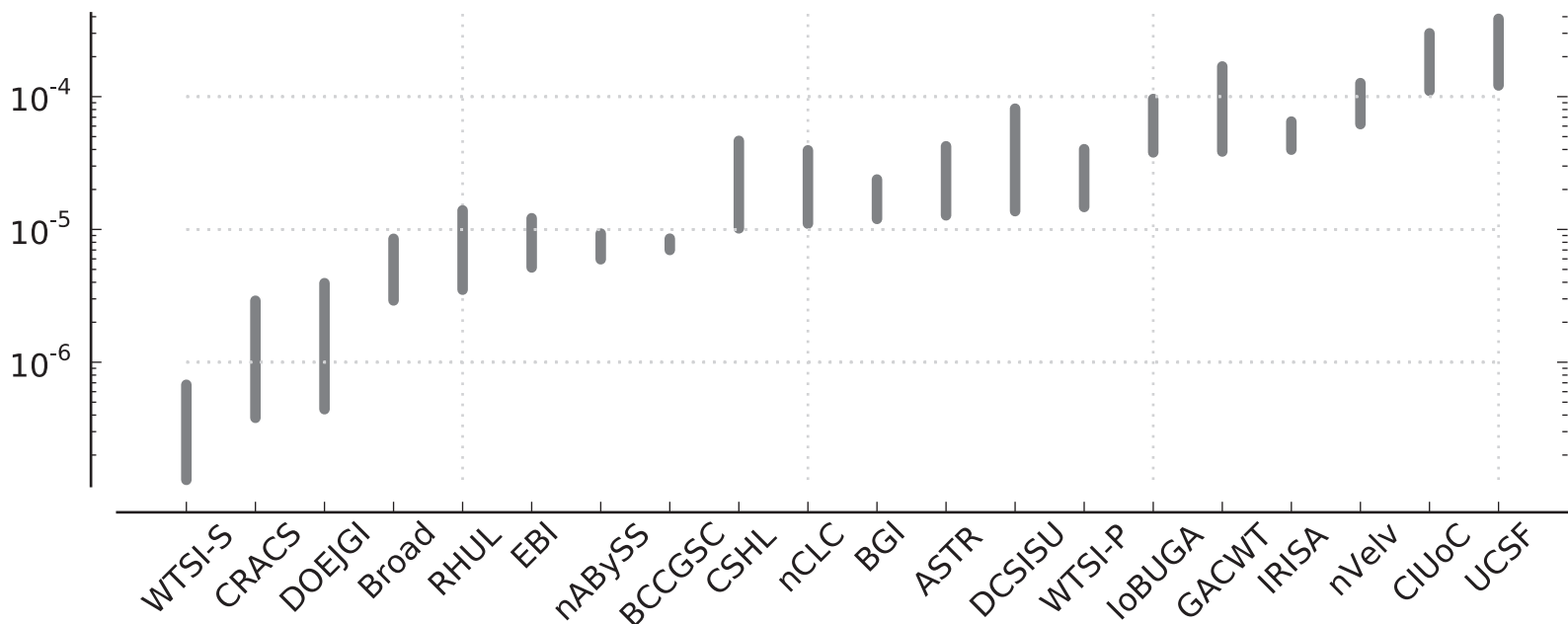
Fill Color Key
Item \geq



Contiguity Statistics



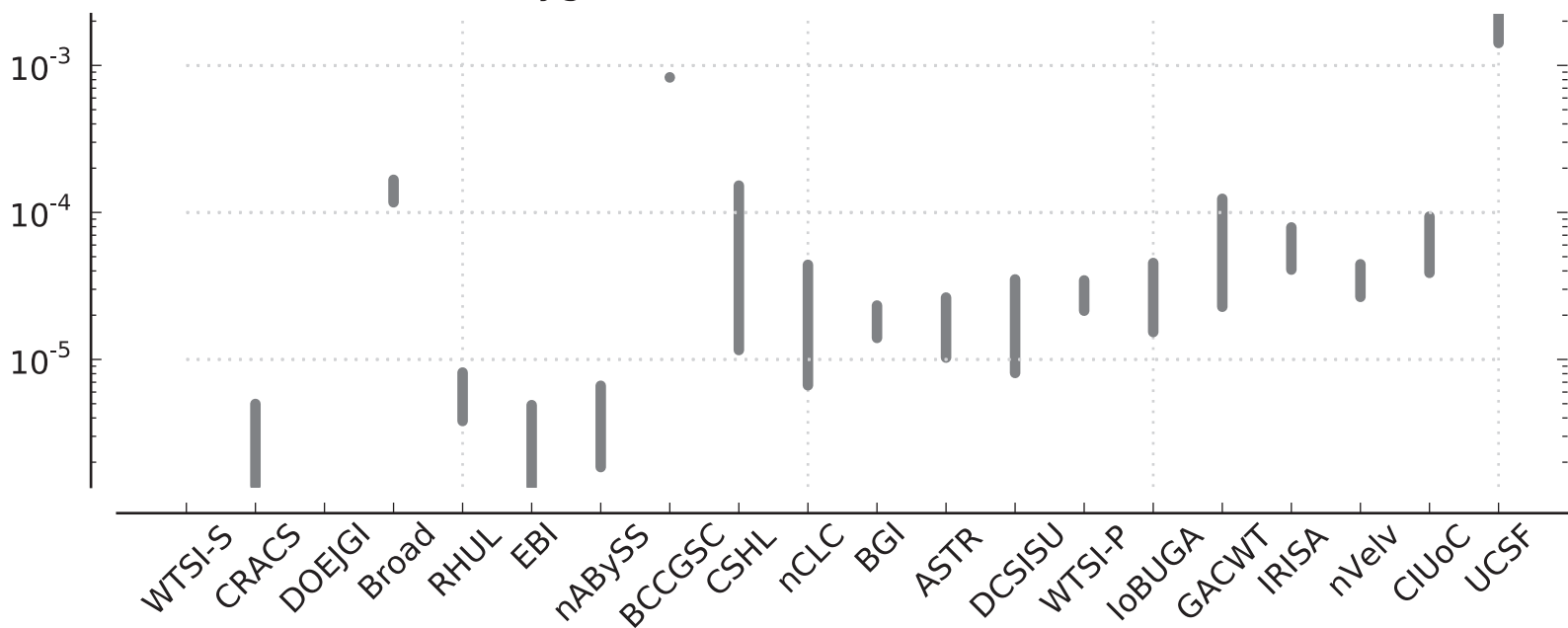
Sum of Substitution Errors / Correct (bits)



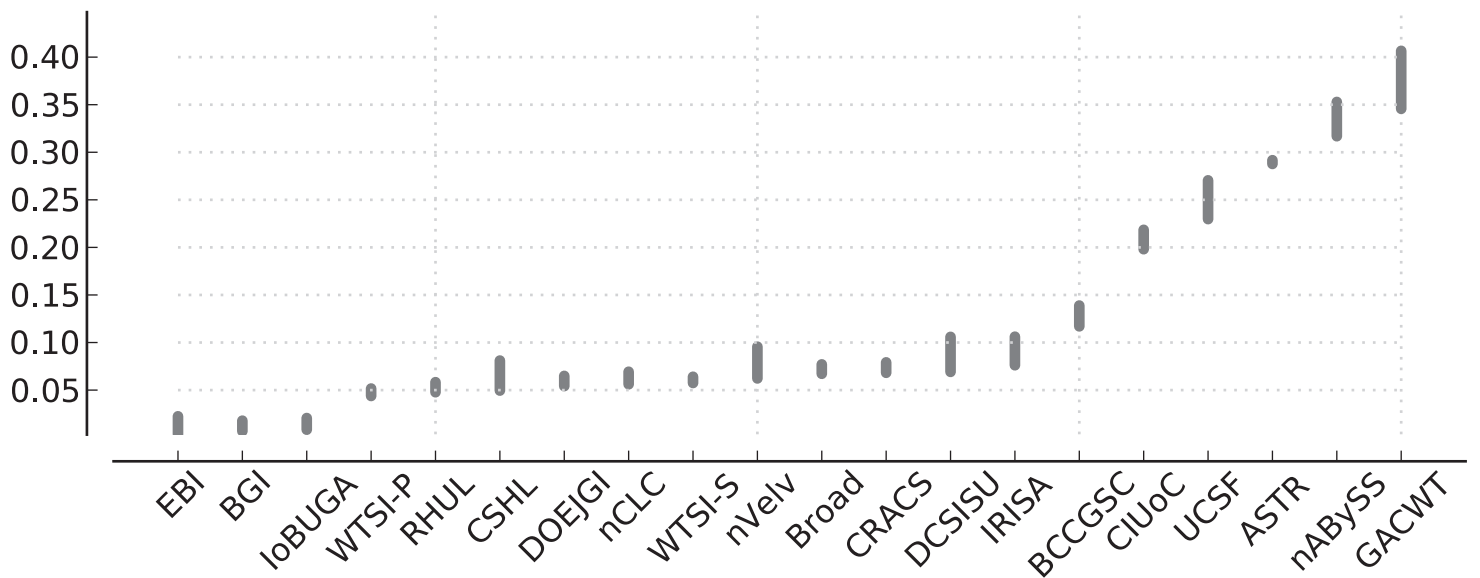
Homozygous Substitution Errors / Correct (bits)



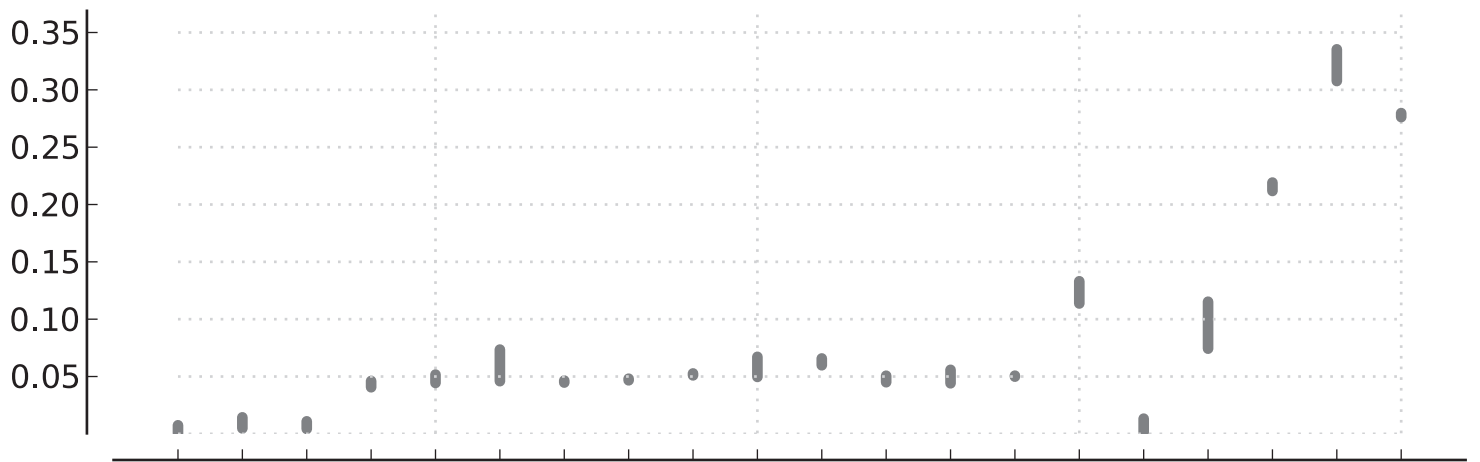
Heterozygous Substitution Errors / Correct (bits)



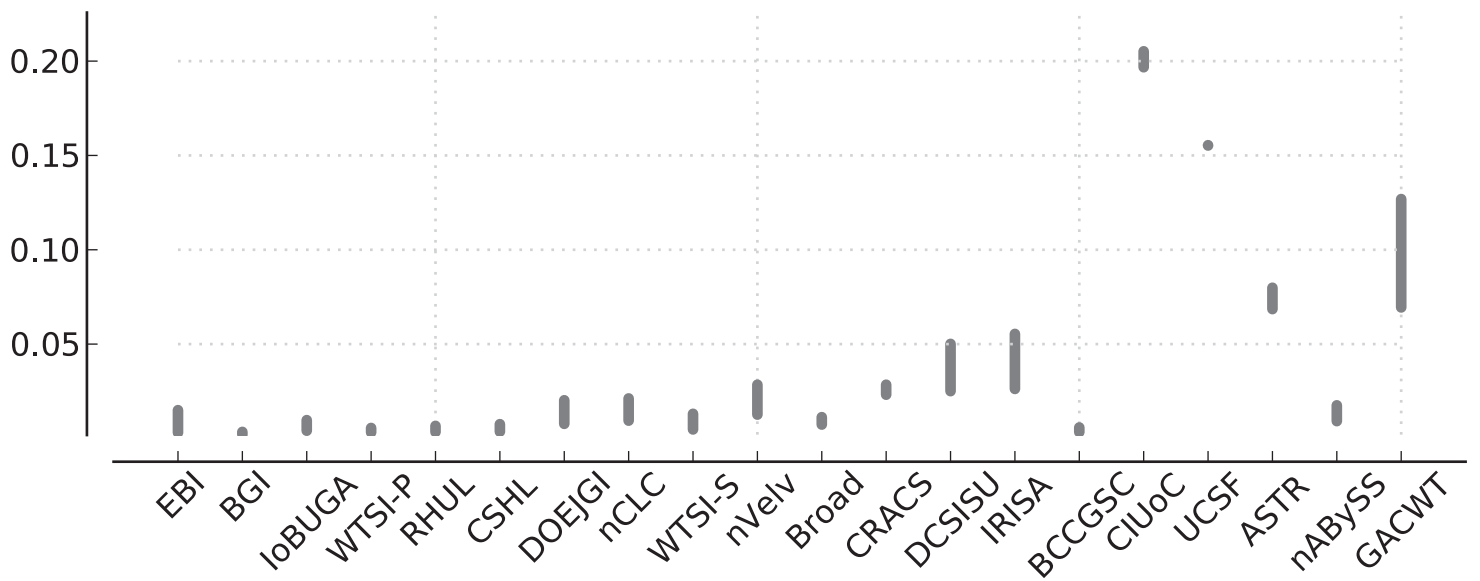
Sum of Proportional Copy Errors



Proportional Excess Copy Errors

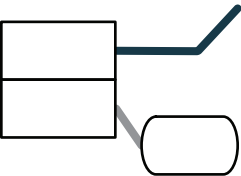


Proportional Deficient Copy Errors



(A)

a



(B)

a



b



(C)

a



b



ID	Affiliations	Entries	Software	Used β
ASTR	Agency for Science, Technology and Research, Singapore	1	PE-Assembler	No
WTSI-P	Wellcome Trust Sanger Institute, UK	2	Phusion2, phrap	No
EBI	European Bioinformatics Institute, UK	2	SGA, BWA, Curtain, Velvet	No
WTSI-S	Wellcome Trust Sanger Insitute, UK	4	SGA	No
CRACS	Center for Research in Advanced Computing Systems, Portugal	3	ABBySS	Yes
BCCGSC	BC Cancer Genome Sciences Centre, Canada	5	ABBySS, Anchor	No
DOEJGI	DOE Joint Genome Insitute, USA	1	Meraculous	No
IRISA	L'IRISA (Institut de recherche en informatique et systèmes aléatoires), France	5	Monument	No
CSHL	CSHL (Cold Spring Harbor Laboratory), USA	2	Quake, Celera, Bambus2	No*
DCISU	Department of Computer Science, Iowa State University	1	PCAP	No
IoBUGA	Computational Systems Biology Laboratory, University of Georgia, USA	3	Seqclean, SOAPdenovo	No
UCSF	UC San Francicso, USA	1	PRICE	Yes
RHUL	Royal Holloway, University of London, UK	5	OligoZip	No
GACWT	The Genome Analysis Centre, Sainsbury Laboratory, and Wellcome Trust Centre for Human Genetics, UK	3	Cortex_con_rp	No
CIUoC	Department of Computer Science, University of Chicago, USA	1	Kiki	No
BGI	BGI, Shenzhen China	1	SOAPdenovo	No
Broad	Broad Institute	1	ALLPATHS-LG	No
nVelv	—	6	Velvet	No
nCLC	—	9	CLC	No
nABBySS	—	6	ABBySS	No

Genome	Mb	GC (%)	Reps (%)	Reps 100mer (%)	Chr	Subs	Dels	Inv	Moves	Copy	Tandem	Chr Split	Chr Fuse
Input	95.6	38.8	7.1 / 42.3*	0.8	4	–	–	–	–	–	–	–	–
MRCAs	109.4	39.9	6.9	0.3	2	35.9e+06	2.47e+06	11,701	4,714	14,644	1.16e+06	2	4
α	112.4	40.0	7.5	0.3	3	9.70e+06	6.72e+05	3,325	1,369	4,151	3.13e+05	1	0
α_1	112.5	40.0	7.5	0.3	3	1.97e+05	13,528	54	34	83	6,436	0	0
α_2	112.5	40.0	7.5	0.3	3	1.97e+05	13,834	61	31	80	6,494	0	0
β	112.3	40.0	6.8	0.3	2	9.71e+06	6.74e+05	3,313	1,325	4,043	3.14e+05	0	0
β_1	112.4	40.0	6.8	0.3	2	1.97e+05	13,632	64	26	82	6,354	0	0
β_2	112.4	40.0	6.8	0.3	2	1.97e+05	13,621	71	35	79	6,445	0	0

Comparison	SNPs	Substitutions	Σ Subs	Indels	Σ Indels	Inversions
$\alpha_1 \alpha_2$	439,385	441,796	444,247	29,972	521,142	115

ID	Overall	CPNG50	SPNG50	Struct.	CC50	Subs.	Copy Num.	Cov. Tot.	Cov. Genic
Broad	31	2 (7.25e+04)	3 (2.11e+05)	3 (1244)	1 (2.66e+06)	4 (2.92e-06)	11 (6.71e-02)	6 (98.3)	1 (93.8)
BGI	37	1 (8.23e+04)	6 (1.17e+05)	6 (1878)	7 (5.66e+05)	11 (1.20e-05)	2 (6.75e-03)	1 (98.8)	3 (92.7)
WTSI-S	38	9 (2.48e+04)	1 (4.95e+05)	2 (475)	3 (1.14e+06)	1 (1.30e-07)	9 (5.74e-02)	8 (97.8)	5 (91.8)
DOEJGI	44	14 (1.15e+04)	2 (4.86e+05)	1 (456)	2 (1.89e+06)	3 (4.43e-07)	7 (5.42e-02)	11 (97.3)	4 (92.3)
CSHL	57	3 (4.23e+04)	8 (7.17e+04)	14 (5146)	6 (6.11e+05)	9 (1.02e-05)	6 (4.95e-02)	4 (98.5)	7 (89.1)
CRACS	58	11 (1.55e+04)	5 (1.44e+05)	4 (1666)	4 (8.61e+05)	2 (3.81e-07)	12 (6.82e-02)	14 (96.3)	6 (90.2)
BCCGSC	60	5 (3.63e+04)	4 (1.46e+05)	10 (2867)	8 (3.22e+05)	8 (7.00e-06)	15 (1.17e-01)	2 (98.7)	8 (88.9)
EBI	64	16 (9.39e+03)	7 (1.13e+05)	7 (2055)	9 (3.04e+05)	6 (5.17e-06)	1 (3.56e-03)	9 (97.7)	9 (88.5)
IoBUGA	65	7 (3.06e+04)	12 (3.54e+04)	15 (6310)	5 (6.47e+05)	15 (3.80e-05)	3 (8.38e-03)	6 (98.3)	2 (92.8)
RHUL	71	6 (3.20e+04)	13 (3.31e+04)	8 (2551)	15 (1.59e+04)	5 (3.52e-06)	5 (4.77e-02)	4 (98.5)	15 (67.4)
WTSI-P	74	4 (3.80e+04)	11 (4.21e+04)	13 (4895)	13 (3.41e+04)	14 (1.48e-05)	4 (4.38e-02)	2 (98.7)	13 (75.0)
DCSISU	99	12 (1.35e+04)	10 (5.61e+04)	12 (4319)	12 (9.75e+04)	13 (1.37e-05)	13 (6.91e-02)	15 (94.3)	12 (79.0)
nABySS	100	10 (1.99e+04)	16 (2.00e+04)	5 (1731)	16 (6.97e+03)	7 (5.96e-06)	19 (3.17e-01)	10 (97.5)	17 (57.2)
IRISA	103	17 (8.20e+03)	9 (5.82e+04)	11 (3725)	9 (3.04e+05)	17 (3.99e-05)	14 (7.61e-02)	16 (93.7)	10 (88.1)
ASTR	106	8 (2.52e+04)	14 (3.13e+04)	9 (2818)	14 (1.81e+04)	12 (1.28e-05)	18 (2.88e-01)	17 (90.9)	14 (68.5)
nVelv	114	18 (5.65e+03)	15 (2.75e+04)	18 (8626)	11 (1.27e+05)	18 (6.21e-05)	10 (6.22e-02)	13 (96.5)	11 (84.8)
nCLC	115	15 (9.47e+03)	18 (9.54e+03)	16 (7283)	18 (4.36e+03)	10 (1.11e-05)	8 (5.61e-02)	12 (97.2)	18 (55.4)
UCSF	138	12 (1.35e+04)	17 (1.35e+04)	20 (24987)	17 (6.84e+03)	20 (1.21e-04)	17 (2.30e-01)	19 (83.7)	16 (59.6)
GACWT	149	20 (2.53e+03)	19 (7.82e+03)	17 (8622)	19 (2.60e+03)	16 (3.86e-05)	20 (3.46e-01)	18 (86.4)	20 (48.0)
CIUoC	152	19 (5.60e+03)	20 (5.60e+03)	19 (11282)	20 (1.27e+03)	19 (1.11e-04)	16 (1.98e-01)	20 (78.5)	19 (48.9)

ID	Hap Total (%)	Hap α_1 (%)	Hap α_2 (%)	Bac (%)	Genic (%)	Unmapped
BGI	98.8	98.9	98.8	0.0	92.7	2.637e+05
BCCGSC	98.7	98.7	98.7	99.9	88.9	6.546e+06
WTSI-P	98.7	98.7	98.7	99.8	75.0	5.369e+06
RHUL	98.5	98.5	98.5	100.0	67.4	4.961e+06
CSHL	98.5	98.6	98.5	99.9	89.1	7.815e+06
Broad	98.3	98.4	98.3	68.9	93.8	3.538e+06
IoBUGA	98.3	98.3	98.3	4.8	92.8	7.822e+05
WTSI-S	97.8	97.8	97.8	99.1	91.8	4.948e+06
EBI	97.7	97.7	97.7	0.9	88.5	4.553e+05
nABySS	97.5	97.5	97.5	99.8	57.2	1.111e+07
DOEJGI	97.3	97.4	97.3	99.5	92.3	5.304e+06
nCLC	97.2	97.2	97.2	99.8	55.4	5.673e+06
nVelv	96.5	96.6	96.5	99.8	84.8	8.028e+06
CRACS	96.3	96.3	96.3	99.8	90.2	5.265e+06
DCSISU	94.3	94.3	94.2	99.5	79.0	6.259e+06
IRISA	93.7	93.7	93.7	99.7	88.1	5.426e+06
ASTR	90.9	90.9	90.9	100.0	68.5	5.175e+06
GACWT	86.4	86.4	86.4	0.0	48.0	2.053e+06
UCSF	83.7	83.7	83.7	0.0	59.6	1.822e+06
CIUoC	78.5	79.0	78.1	0.6	48.9	3.638e+05

ID	Intra chromosomal joins	Inter chromosomal joins	Insertions	Deletions	Insertion and deletion	Insertion at ends	Σ errors
DOEJGI	21	160	55	108	40	72	456
WTSI-S	6	191	56	76	19	127	475
Broad	75	161	524	379	9	96	1,244
CRACS	687	303	198	121	51	306	1,666
nABySS	17	48	208	188	63	1,207	1,731
BGI	368	288	355	639	98	130	1,878
EBI	458	563	127	547	53	307	2,055
RHUL	691	349	172	264	26	1,049	2,551
ASTR	2,065	200	109	227	73	144	2,818
BCCGSC	351	285	255	233	102	1,641	2,867
IRISA	147	203	925	1,593	116	741	3,725
DCSISU	1,410	956	330	954	109	560	4,319
WTSI-P	1,940	449	1,851	289	87	279	4,895
CSHL	396	337	417	3,287	223	486	5,146
IoBUGA	919	330	1,663	2,933	356	109	6,310
nCLC	23	64	2,359	2,237	68	2,532	7,283
GACWT	757	730	905	1,292	216	4,722	8,622
nVelv	2,885	455	1,473	2,838	306	669	8,626
CIUoC	1,205	684	1,189	2,026	65	6,113	11,282
UCSF	2,731	2,396	5,908	6,223	1,018	6,711	24,987

ID	COG-xcript (996,462)	CO-cds (562,627)	CO-utr (433,835)	CO-nxe+nge (21,292,660)	CO-repeat (14,475,489)
ASTR	0.11	0.92	0.82	0.92	0.56
WTSI-P	0.09	0.96	0.82	0.99	0.59
EBI	0.08	0.97	0.76	0.99	0.55
WTSI-S	0.07	0.89	0.75	0.99	0.56
CRACS	0.07	0.92	0.72	0.97	0.53
BCCGSC	0.08	0.94	0.79	0.99	0.59
DOEJGI	0.05	0.88	0.65	0.99	0.45
IRISA	0.06	0.89	0.66	0.97	0.37
CSHL	0.08	0.94	0.80	0.99	0.57
DCSISU	0.06	0.83	0.66	0.97	0.42
IoBUGA	0.08	0.97	0.81	0.99	0.58
UCSF	0.06	0.84	0.62	0.86	0.37
RHUL	0.09	0.96	0.81	0.99	0.59
GACWT	0.02	0.72	0.37	0.88	0.38
CIUoC	0.02	0.74	0.49	0.80	0.39
BGI	0.11	0.99	0.84	0.99	0.64
Broad	0.10	0.97	0.83	0.99	0.64
nVelv	0.04	0.88	0.69	0.99	0.34
nCLC	0.05	0.92	0.70	0.99	0.41
nABySS	0.06	0.91	0.73	0.99	0.46

Metric name	Units	Description
N50	—	A weighted median of the lengths of items, equal to the length of the longest item i such that the sum of the lengths of items greater than or equal in length to i is greater than or equal to half the length of all of the items. With regard to assemblies the items are typically contigs or scaffolds.
NG50	—	Whereas N50 sets the median in relation to the total length of all items in the set, we define NG50 to be normalised by the average length of the α_1 and α_2 haplotypes instead of the total length of all sequences as in N50, it is thus more reliable than N50 for comparison between assemblies.
CPNG50	bp	Contig path NG50. The weighted median of the lengths of contig paths. Contig paths represent maximal subsequences of contigs that are entirely consistent with $\alpha_{1,2}$.
SPNG50	bp	Scaffold path NG50. The weighted median of the lengths of scaffold paths. Scaffold paths represent maximal concatenations of contig paths and scaffold breaks that maintain correct order and orientation.
Structural errors	Counts	An error within a contig or scaffold. Errors include intra and inter chromosomal joins, insertions, deletions, simultaneous insertion and deletions and insertions at the ends of assembled sequences.
CC50	bp	Correct contiguity 50. The empirically sampled distance between two points in an assembly, where the two points are as likely to be correctly aligned as not.
Substitution errors	Counts per correct bits	Number of substitution errors per correct bit. Substitution errors are columns in the alignment where the α_1 and α_2 haplotypes contain either the same base (homozygous) or different bases (heterozygous) and the alignment contains a base (or IUPAC symbol) different from either α_1 or α_2 . The metric uses a bit score to allow for IUPAC symbols.
Copy number errors	Proportions	For a given haplotype column in the MSA the copy number of the simulated genome can be described as an interval $[min, max]$. Assemblies with copy number outside of this interval are classified either as an excess, for being above the interval, or a deficiency, for being below the interval.
Coverage	Percent	The coverage is the percent of columns in the MSA of the target (the whole genome, regions of a specific annotation type, etc) that contain positions of the assembly.
Genic correctness	Percent	The genic correctness is the percentage of bps in spliced transcripts from the haplotype sequences that align to the assembly with 95% coverage using WU-BLAST.