



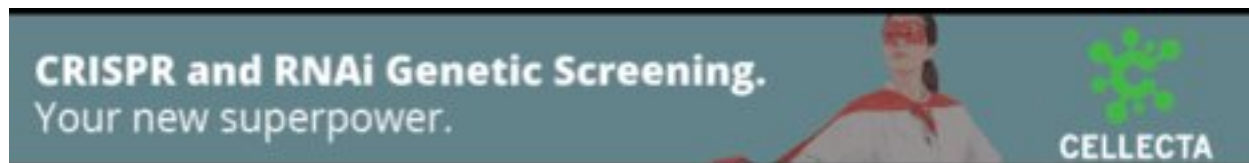
Alignathon: A competitive assessment of whole genome alignment methods

Dent Earl, Ngan K Nguyen, Glenn Hickey, et al.

Genome Res. published online October 1, 2014

Access the most recent version at doi:[10.1101/gr.174920.114](https://doi.org/10.1101/gr.174920.114)

| | |
|---------------------------------|--|
| P<P | Published online October 1, 2014 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. |
| Creative Commons License | This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at http://creativecommons.org/licenses/by/4.0/ . |
| 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>

Published by Cold Spring Harbor Laboratory Press

Alignathon: A competitive assessment of whole genome alignment methods

Dent Earl^{1,2}, Ngan Nguyen^{1,2}, Glenn Hickey^{1,2}, Robert S. Harris³, Stephen Fitzgerald⁴, Kathryn Beal⁴, Igor Seledtsov⁵, Vladimir Molodtsov⁵, Brian J. Raney¹, Hiram Clawson¹, Jaebum Kim⁶, Carsten Kemena^{7,8,9}, Jia-Ming Chang^{7,8,10}, Ionas Erb^{7,8}, Alexander Poliakov¹¹, Minmei Hou¹², Javier Herrero^{4,13}, William James Kent^{1,2}, Victor Solovyev⁵, Aaron E. Darling¹⁴, Jian Ma¹⁵, Cedric Notredame^{7,8}, Michael Brudno^{16,17}, Inna Dubchak^{11,18}, David Haussler^{1,2,19}, Benedict Paten^{1,2}

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

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

(3) Department of Biology, The Pennsylvania State University, PA, USA

(4) EMBL-European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, CB10 1SD, UK

(5) Softberry Inc., USA

(6) Department of Animal Biotechnology, Konkuk University, Seoul 143-701, Korea

(7) Centre For Genomic Regulation (CRG), Carrer del Doctor Aiguader 88, 08003 Barcelona, Spain

(8) Universitat Pompeu Fabra (UPF), 08003 Barcelona, Spain

(9) Westfalian Wilhelms University, Institute of Evolution and Biodiversity, Huefferstr. 1, 48149 Muenster, Germany

(10) Institute of Human Genetics, 141 rue de la Cardonille, 34396 Montpellier, France

(11) Department of Energy Joint Genome Institute, 2800 Mitchell Drive, Walnut Creek, CA, 94598, USA

(12) Department of Computer Science, Northern Illinois University, USA

(13) The Genome Analysis Centre, Norwich Research Park, Norwich, NR4 7UH, UK

(14) ithree Institute, University of Technology Sydney, Sydney, Australia

(15) Department of Bioengineering and Institute for Genomic Biology, University of Illinois at Urbana-Champaign, USA

(16) Department of Computer Science and the Donnelly Centre, University of Toronto 10 King's College Rd, Toronto, ON M5S 3G4, Canada

(17) Centre for Computational Medicine and the Genetics and Genome Biology Program, Hospital for Sick Children, 555 University Avenue, Toronto, ON M5G 1X8 Canada

(18) Lawrence Berkeley National Laboratory, 1 Cyclotron Road, Berkeley, CA 94710, USA

(19) Howard Hughes Medical Institute, Bethesda, MD, USA

Abstract

Multiple sequence alignments (MSAs) are a prerequisite for a wide variety of evolutionary analyses. Published assessments and benchmark datasets for protein and, to a lesser extent, global nucleotide MSAs are available, but less effort has been made to establish benchmarks in the more general problem of whole genome alignment (WGA).

Using the same model as the successful Assemblathon competitions we organized a competitive evaluation in which teams submitted their alignments and then assessments were performed collectively after all the submissions were received. Three datasets were used; two were simulated and based on primate and mammalian phylogenies and one was comprised of 20 real fly genomes. In total 35 submissions were assessed, submitted by ten teams using 12 different alignment pipelines.

We found agreement between independent simulation-based and statistical assessments indicating that there are substantial accuracy differences between contemporary alignment tools. We saw considerable difference in the alignment quality of differently annotated regions and found few tools aligned the duplications analysed. We found many tools worked well at shorter evolutionary distances, but fewer performed competitively at longer distances. We provide all datasets, submissions and assessment programs for further study and provide, as a resource for future benchmarking, a convenient repository of code and data for reproducing the simulation assessments.

Introduction

Given a set of sequences, a multiple sequence alignment (MSA) is a partitioning of the residues in the sequences, be they amino-acids or nucleotides, into related sets. Here we are interested in the relationship of evolutionary homology. In other contexts residues may be aligned with a different aim, as in structural alignments, where residues are aligned if located at the same point in a shared crystal structure (Kolodny et al. 2005). MSA is a fundamental problem in biological sequence analysis because it is a prerequisite for most phylogenetic and evolutionary analyses (Felsenstein 2003, Wallace et al. 2005, Edgar & Batzoglou 2006, Notredame 2007). Most MSAs are termed “global”, made of sequences assumed to be related through the mutational processes of residue substitution, subsequence insertion and subsequence deletion (collectively insertions and deletions are termed indels); for a review see Notredame 2007. The availability of whole genome sequences has led to an interest in MSAs for complete genomes, including all sequences: genes, promoters, repetitive regions, etc. Termed whole genome alignment (WGA), this requires the aligner to additionally consider genome rearrangements, such as inversions, translocations, chromosome fusions, chromosome fissions and reciprocal translocations. Some tools for WGA are also capable of modelling unbalanced rearrangements that lead to copy number change, such as tandem and segmental duplications (Blanchette et al. 2004, Miller et al. 2007, Paten et al. 2008, Angiuoli & Salzberg 2011, Paten et al. 2011). WGA methods have been critical to understanding the selective forces acting across genomes, allowing evolutionary analysis of many potential functional elements (The ENCODE Project Consortium 2012), and in particular, the identification of conserved non-coding functional elements (*Drosophila* 12 Genomes Consortium 2007, Lindblad-Toh et al. 2011), including cis-regulatory elements (Kellis et al. 2003), enhancers and noncoding RNAs.

The lack of accepted gold standard reference alignments has made it hard to objectively assess the relative merits of WGA methods. Previous evaluations of MSAs can be split into roughly four types: those using simulation, those using expert information, those using direct statistical assessments and finally those that assess how well an alignment functions for a downstream analysis. We briefly describe and review these approaches; for a more comprehensive review see Iantorno et al. 2014.

In simulation evaluations a set of sequences and an alignment is generated using a model of evolution. Alignments are created from the simulated sequences and the resulting predictions are compared to the “true” simulated alignment. There are two basic types of simulators for DNA sequence evolution, coalescent simulators and non-coalescent forward-time simulators (Carvajal-Rodríguez 2010). While useful for modelling populations, coalescent simulators can not yet efficiently model general sequence evolution, and as a result MSA simulators currently use forward-time approaches. There are numerous forward-time simulators useful for assessing global MSA tools (Stoye et al. 1997, Blanchette et al. 2004, Cartwright 2005, Varadarajan et al. 2008). However, the simulation options for assessing WGA have until recently been absent, essentially because to do so requires modelling both low level sequence evolution and higher level genome rearrangements — a formidable challenge given the large and complex parameter space that potentially encompasses all aspects of genome evolution. The sgEvolver simulator (Darling et al. 2004, 2010) is used to generate simulated genome alignments, though it lacks an explicit model for sequence translocation or mobile element evolution. EvolSimulator is a genome simulator but it has a somewhat simple model of evolution and a focus on ecological parameters (Beiko & Charlebois 2007). Another option, the ALF simulator (Dalquen et al. 2012), models gene

and neutral DNA evolution. For this study we used the EVOLVER software, which can simulate full-sized, multi-chromosome genome evolution in forward time (Edgar et al. 2009). EVOLVER models an explicitly haploid genome and lacks a population model; its framework and expert-curated extensive parameter set are intended to produce “reference-like” genomes, i.e., haploid genomes. EVOLVER models DNA sequence evolution with sequence annotations; a gene model; a base-level evolutionary constraint model; chromosome evolution, including inter- and intra-chromosomal rearrangements; tandem and segmental duplications; and mobile element insertions, movements, and evolution.

An alternative approach to assessing MSA is to use expert biological information not available to the aligner. While interpreting the results of simulations is made difficult by the uncertainty to which they approximate reality, the clear advantage of using expert information is that it can be used to assess alignments of actual biological sequences. For protein and RNA alignment there are several popular benchmarks that provide either reference structural alignments or expertly curated alignments (Wilm et al. 2006, Blackshields et al. 2006, Kemena et al. 2013). Nontranscribed DNA alignments are, however, much harder to assess since one lacks an external criterion to assemble objective gold standard references (Kemena & Notredame 2009). This explains why untranslated DNA alignments are usually evaluated using more ad hoc expert information (Margulies et al. 2007, Paten et al. 2008). The main strength of these procedures is that they provide an objective evolutionary context when evaluating the alignment. The difficulty with relying upon such expert information is that it may address only a small fraction of the alignment (e.g. in the referenced papers, coding exons and ancient repeats), and may itself rely on other forms of inference (e.g. ancient repeat analyses has an explicit dependence on the sequence alignment procedures used to determine ancestral repeat relationships) and have unknown variance, generality and discriminative power.

The third approach addresses alignments by statistical measures. For global MSA there are several options, e.g. the T-Coffee CORE index (Notredame & Abergel 2003), Heads or Tails (HoT) (Landan & Graur 2008), GUIDANCE (Penn et al. 2010a, 2010b) and StatSigMA-w (Chen & Tompa 2010). For this work we expand on the probabilistic sampling-based alignment reliability (PSAR, Kim & Ma 2011) method, which samples pairwise suboptimal alignments to assess the reliability of MSAs. Statistical measures are attractive because they can be used with the complete alignments of real sequences. However without a gold standard to compare against they are only a proxy to a true assessment of accuracy.

The final category of common assessment methods addresses how well a program generates alignments for a given computational task. This is typically the assessment made by a biologist in choosing an alignment program — i.e. how well does it perform in practice, according to intuition or analysis? Unfortunately, these assessments, often being one-offs, rarely make it into the literature, and are difficult if not impossible to generalize from because these assessments are made for the purposes of a given analysis. Notably for WGAs, Bradley et al. assessed how much alignment methods influenced de novo ncRNA predictions and Margulies et al. analysed the effect of different WGAs on the prediction of conserved elements (Bradley et al. 2009, Margulies et al. 2007).

There have been relatively few independent or community organised assessments of WGA pipelines. Notably, as part of the ENCODE Pilot Project (Margulies et al. 2007), four pipelines were assessed across a substantial number of regions, and Tompa et al. later compared those alignments using the StatSigMA-w tool (Chen & Tompa 2010). The Alignathon is an attempt to

perform a larger and more comprehensive evaluation. It is a natural intellectual successor to the Assemblathon collaborative competitions (Earl et al. 2011, Bradnam et al. 2013). The starting point of the Alignathon is to assume that the problem of genome assembly is largely a solved problem. Though we admit this is currently a dubious assumption, it appears that the problem of genome assembly will shrink in size in the coming years as new sequencing technologies become available and existing assembly software is perfected to take advantage of more numerous, longer and less error-prone reads (Branton et al. 2008, Schreiber et al. 2013, Laszlo et al. 2014). With this future as a starting point, the question a biologist faces changes from a proximate one of “how do I best assemble the genome of my favorite species?” to a higher level question of “how is my favorite species related to the pantheon of other sequenced species?” Such a question is answered through a WGA. If organized community efforts to sequence large numbers of genomes, such as the Genome 10K for vertebrates and insect 5K for insects, are to maximally fulfill their promise by revealing and refining the evolutionary history of all of their species then it is vital that we have the best possible methods for WGA (Genome 10K Community of Scientists 2009, i5K Consortium 2013).

Results

Of the four discussed strategies to assess alignments we pursue two: simulations and statistical assessment. We now describe the Alignathon datasets, the submissions we received, how the submissions were processed and the evaluations that were performed.

Datasets

The Alignathon used three test sets. Two of the test sets were created by way of forward-time simulation, using the EVOLVER tool, starting from a $\sim 1/20$ th scale mammalian genome, a genome size of 120 megabases (Mb), based upon a subset of hg19/GRCh37 (chromosomes 20, 21 and 22; see Methods). The first simulated dataset models a great ape phylogeny consisting of genomes with the same evolutionary relationships as humans, chimpanzees, gorillas and orangutans (Fig. 1). The second simulated dataset is based upon a mammalian phylogeny consisting of genomes with the same evolutionary relationships as humans, mice, rats, cows and dogs (Fig. 1). On a gross level, the summary statistics of the two simulated datasets are shown in Tables 1 and S1. After an initial burn-in phase to shuffle the original input sequences and ensure the simulation had reached stationarity (see Methods), the primate phylogeny contained, amongst other changes, one chromosomal fusion and over three million substitutions in the lineage from the most recent common ancestor (MRCA) to the simulated human. The mammal phylogeny contained, amongst other changes, two chromosomal splits, one fusion and over 27 million substitutions in the lineage from the MRCA to the simulated human.

Recognizing the limitations of simulations, our third test set consisted of 20 real fly genomes (Fig. 1). The fly genomes were available in various states of completion from near-finished in the case of *Drosophila melanogaster* (dm3 assembly, chromosome sequences) to fragmentary in the case of *D. rhopaloa* (droRho assembly, 34,000 contigs) (Table 2).

Competition Organization and Submissions

The initial datasets were released in December 2011, and teams were given until February 2012 to submit their entries. The initial simulated datasets included the truths and information on where to obtain (and optionally contribute to) the analysis software. As in the Assemblathons, none of the

teams had access to the datasets until their initial release. The Alignment received 35 submissions, 13 for the primate simulation, 13 for the mammal simulation and nine for the fly dataset (Tables 3, S2, S3 and S4). The pipelines that were used to generate the alignments represent those used by genome browsers to generate their WGAs: VISTA-LAGAN for the VISTA Browser (Frazer et al. 2004, Dubchak et al. 2009); MULTIZ for the UCSC Genome Browser (Miller et al. 2007, Meyer et al. 2013) and Pecan and EPO for the Ensembl Browser (Paten et al. 2008, Paten et al. 2009, Flicek et al. 2013). In addition, we tested a fairly broad set of standalone WGA tools, including: progressiveMauve (Darling et al. 2010); TBA (Blanchette et al. 2004); Cactus (Paten et al. 2011); Mugsy (Angiuoli & Salzberg 2011), which was designed for closely related genomes; a meta-WGA tool, Robusta (Notredame 2012), which combines results from multiple standalone tools; and a realignment tool, PSAR-Align, which was used to realign MULTIZ based alignments in this competition, but can in principle refine alignments from any multiple alignment tool. We also tested pairwise WGAs from the GenomeMatch team. As might be expected, not all algorithms / pipelines were run for all test sets. Participants cited limitations of the methods applied (e.g. inability to handle the scale of the fly dataset), and of resources (time, person-hours, funding, etc.) as reasons for not participating in all datasets. Descriptions generated by the teams of the computation of each submission are given in the supplement, as are details on runtimes and computational resources used.

Genomewide Comparison to Simulated Genome Alignments

All submissions were received in multiple alignment format (MAF) (<http://genome.ucsc.edu/FAQ/FAQformat.html#format5>). A suite of MAF comparison tools was developed for the project (mafTools, see Methods) including a comparator tool, so-called because it compares two alignment files. We call the set of aligned pairs of residues within an alignment its *alignment relation*. The comparator tool works by taking two input MAF files, A and B, and comparing their alignment relations. For the simulated datasets, if A is the predicted alignment created by a tool and B is the simulated truth, then the ratio of the number of pairs in the intersection of A and B to the number of pairs in A is the *precision* of the prediction. Conversely, the ratio of the number of pairs in the intersection of A and B to the number of pairs in B is the *recall* of the prediction. One standard method for combining precision and recall into a single value is the balanced F-score, which is simply the harmonic mean of precision and recall:

$$F\ score = 2 * \frac{precision * recall}{precision + recall}$$

(Beitzel 2006).

The cardinality of the alignment relation of the considered WGAs is exceedingly large, e.g. approximately 1.7 billion pairs for the simulated mammalian alignment. This made complete comparison impractical. Instead, for each pair of MAFs compared we sampled (see Methods) a subset of the alignment relation of one and checked if any or all elements of the subset were present in the alignment relation of the other. 10 million pairs were sampled for each direction of a MAF pair comparison, and variance between sampling runs was negligible (data not shown).

For the simulated datasets we performed analyses both with respect to the entire genome and to areas of the genome subsetted by annotation type (genic, neutral and repetitive; see Methods). Results are shown in Fig. 2 and Tables S5-8. We find that many of the submissions were able to align the primate dataset with both relatively high recall and precision, and, with the exception of the GenomeMatch submissions, which had lower values in the repetitive regions, performance

was consistently high across annotation types, e.g. the top eight submissions differed by only 0.007 in F-score, and all had recall and precision above 0.98.

For the mammal simulations we found a much wider spread of results, both between aligners and within different annotation classes. The strongest submission, Cactus, had an F-score 0.081 points higher than its nearest competitor. Looking at the mammal results by annotation type, generally (and predictably) submissions performed the best in genic regions where simulated selection, which lead to strong conservation, was presumably highest. Performance was intermediate in neutral regions and submissions generally performed most poorly in repetitive regions. Generally submissions retained their ranking across annotation regions, that is to say that the submissions ranked 1 and 2 overall were also ranked 1 and 2 in genic regions. However, this trend did not hold for repetitive regions, and surprisingly several submissions performed slightly better in repetitive regions than in the neutral regions (Mugsy, Pecan, EPO, Robusta).

As phylogenetic distance between species grows the number of unobserved mutation events increases and the alignment problem naturally becomes more difficult (Holmes & Durbin 1998, Wong et al. 2008, Landan & Graur 2008). To see this, we stratified the results by phylogenetic distance (path length between leaves in the simulated phylogenies) between all pairs of species; see Fig. 3. Longer distances are indeed observed to lead to lower precision and recall values, and therefore lower F-score values. For reference based aligners, which use one species as a reference (here simHuman), there is a clear dip in performance for non-reference pairs (pairs not including the reference sequence). This is especially prevalent in Fig. 3.B for the PSAR-Align submission, which used the MULTIZ program, and for the MULTIZ and AutoMZ submissions, which also rely upon the MULTIZ program.

Evaluating Genome Alignments in the Absence of a True Alignment

We used the recently developed PSAR statistical alignment tool (Kim & Ma 2011) to compare to the simulation results and to assess the fly dataset. PSAR assesses an alignment by removing a sequence, sampling suboptimal alignments between the removed sequence and the remaining alignment using the forward algorithm with a pairwise hidden Markov model (pair-HMM), and then checking to see how well the newly sampled alignments match the original alignment. By repeatedly performing this sampling with every possible sequence PSAR is able to calculate an alignment reliability score, termed a *PSAR pair score*, for every pair of matched residues in the alignment. Each *PSAR pair score* is similar to the posterior probability that a given pair of residues in the input alignment are aligned (Durbin et al. 1998), i.e. it can be thought of as a proxy to a local measure of accuracy that factors in the edit matrix surrounding the pair of aligned residues.

To deal with its limited alignment model — which is appropriate for global MSA, allowing only substitutions, insertions and deletions — and to make it computationally feasible to assess the alignments we ran PSAR on subsections of the datasets. For each of the datasets we ran PSAR on five sampled half-megabase subregions (see Methods). Subregion alignments were converted to make them appropriate for PSAR (e.g. removing duplications, ordering rows, etc; see Methods). For a pair of genomes we define the *PSAR-precision* as the average of the PSAR pair scores of their aligned residues. The overall *PSAR-precision* for the complete alignment is the average of PSAR-precision for genome pairs including the reference. The PSAR-precision scores are analogous to the precision measures calculated from the simulations, because they

estimate the expected number of pairs in the alignment that are correctly aligned.

To complement our proxy to precision we used a simple proxy to recall: coverage. For a pair of genomes A and B, the proportion of residues in A aligned to a residue from B is the *coverage of B on A*. The *overall coverage* (where we drop the “overall” when it is clear from the context) is the average of coverages for all pairs of distinct species. Hypothesizing that PSAR-precision can be used to approximate precision and that coverage can be used as an estimate of recall, the natural statistical analog to F-score is the harmonic mean of PSAR-precision and coverage, which we call the *pseudo F-score*.

To see how consistent our statistical measures were to the measures derived from the simulations we calculated them for the user-generated simulated primate and mammal alignments (Fig. 4 and Tables S5 and S6). To check for any bias created by the use of a set of adjusted regional alignments in calculating the PSAR-precision and coverages values we calculated regional precision and recall values using the regional alignments (see Methods) and found good correlations between the regional and overall versions of these numbers (see Methods, Fig. 4 and Tables S5 and S6).

In the simulated primates all the values were uniformly high, and hence saturated. However, looking at the simulated mammals we find a very good, linear correlation between recall and coverage ($r^2 = 0.984$), but no linear correlation between PSAR-precision and precision. In particular PSAR reports relatively consistent, high scores for all the different alignment programs, suggesting that at a local, residue level on aggregate the alignments look equivalently reasonable between alignment programs. Despite the lack of linear correlation between precision and PSAR-precision, we find that because of the excellent recall and coverage correlation the F-score and pseudo F-score results linearly correlate strongly ($r^2 = 0.975$ in simulated mammals). This appears to be because the more limiting factor in many of the alignments’ performance was not precision, but rather a lack of relative recall/coverage, something particularly affecting the GenomeMatch, Mugsy, and, to a lesser extent, Pecan, EPO and Robusta submissions.

Fig. 5 and Table S9 shows the overall PSAR-precision, coverage and pseudo F-score results for the fly dataset, and Fig. S3 shows pseudo F-score stratified by phylogenetic distance for the fly dataset. For the teams that submitted alignments for both datasets we see good concordance between the fly and simulated results. Again, the difference between the aligners is dominated by coverage, with uniformly high (all greater 0.97) average PSAR-precision values that mostly lie within the regional standard deviations of one another, with the exception of the GenomeMatch alignments, which have very high PSAR-precision values but relatively low coverage. Surprisingly given their reference assisted nature, we find that, along with Cactus and TBA, MULTIZ and AutoMZ had high relative coverage and pseudo F-scores, even when factoring that coverage was calculated over all pairs, not just reference containing pairs. Plotting the pairwise coverages between all pairs of species (Fig. 6) we see that all the programs had higher relative coverage for pairs involving the reference, partially this is an artifact of the structure of the phylogeny (Fig. 1). The reference based aligners (here MULTIZ and AutoMZ) indeed did have the highest coverage for reference pairs, and the strongest non-reference based aligners by these metrics, TBA and Cactus, showed a smaller separation between reference and non-reference species pairs.

Visualizing and Analysing Regional Accuracy

We have demonstrated that the performance of alignments varied regionally according to the

simulation of annotation types. Having developed scoring metrics that can be applied across simulated and non-simulated genomes, we corroborated this analysis by visualizing how the scores vary across the sampled subregions. To view a complete subregion at approximately this level of resolution we binned the reference sequence into 1 kb non-overlapping intervals and calculated the F-score (for the simulated datasets) and pseudo F-score for each bin, calculating the score for a bin as if it represented the complete alignment, and for the simulated comparison restricting the true alignment to just those pairs involving residues in the reference interval that defined the bin. Fig. 7, S1 and S2 visualizes how the scores vary across example subregions of, respectively, the simulated mammalian, real fly and simulated primate alignments. It is clear that the “best” alignments by these measures differs substantially from the poorest, and that for many submissions there is considerable regional variation. Looking across all the simulated regions, the F-score and pseudo F-score measures correlate reasonably bin-by-bin (Fig. S4-6, $r^2 = 0.671$), indicating that pseudo F-score can be used as a reasonable proxy to F-score at this regional level of resolution (e.g. see Fig. S7, the equivalent to Fig. S2, but using pseudo F-score instead of F-score). It should be noted that the correlation is imperfect; in particular it appears that the pseudo F-scores saturate at high values, while the corresponding F-scores still discriminate alignment quality. I.e. pseudo F-scores do not always discriminate between good and very good alignments.

Comparing the Submissions Directly

Several of the pipelines used some of the same underlying programs. To see how these commonalities affected the alignments, for each dataset we calculated the Jaccard distance between the alignment relations of each of the submissions (Fig. 8). As predicted by the earlier analyses, the primate submissions are relatively similar to one another, while the mammalian and fly submissions prove much more divergent. The inter-dataset commonality between some submissions is striking, with the same patterns being repeated across the three datasets, and fits well with the programmatic commonalities that the pipelines share. The results indicate that some of the programmatic commonalities between the alignment pipelines are perhaps more important than others. For example, sharing the same synteny block generator (Mercator or MULTIZ), appears to have had a greater effect on the results than sharing the same synteny block aligner. In particular the EPO and Pecan submissions both use the Pecan program (Paten et al. 2008, Paten & Birney 2009) to align sets of syntenic sequences, and the Cactus program (Paten et al. 2011) uses the same pairwise-HMM to generate much of its multiple alignment as Pecan, but these submissions were relatively different from one another.

Assessing Similarity vs. Homology

The MAF specification allows sequence residues to appear in one or more columns of the file. This allows a program to report sequence similarity, which is not a transitive property, but allows inconsistency if the intention is to report homology. E.g. if a residue x is aligned to a residue y and y is aligned to a residue z , then x and z should be aligned, because it is not possible for x and y to share a common ancestor and for y and z to share a common ancestor while x and z do not share a common ancestor. Comparing two alignments, one of which is transitively closed, and one of which is not, based only on the aligned pairs contained in a MAF potentially gives an advantage to the non-transitively closed submission. This is because the transitively closed submission must align all residues transitively connected by alignments, which may lower the overall precision of the set of aligned pairs. To see how different the results would be if we were to have enforced transitive closure — and therefore a strict homology assumption — we created a tool

(mafTransitiveClosure) which computes the transitive closure of a MAF (a linear time operation). Fig. S8 shows the results for the simulated datasets. We found that the progressiveMauve, Cactus, Pecan, Robusta, EPO and Mugsy programs produced WGAs that were transitively closed, and therefore unaffected by transformation. As predicted, those submissions that were not initially transitively closed, such as the pairwise and reference based alignments, all saw their precision performance decline, in some cases very substantially, and no submission saw a large boost in recall.

Missing Duplications

To find duplications within the alignments we used a simple metric, duplicative coverage. For a pair of genomes A and B, the duplicative coverage of B on A is the proportion of residues in A aligned to two or more distinct residues in B. This assessment is made complicated by the lack of transitive closure in some submissions, because a single residue may align to two or more residues in a genome, but in separate columns of the file. To avoid this complexity we assessed the submissions after computing the transitive closure (which also made the computational task significantly easier). To avoid misrepresenting submissions we dropped submissions from the assessment for which the transitive closure adversely (> 0.05 change) affected the F-score or pseudo F-score. Figure S9 shows the results; in short, we find that only Cactus had significantly non-zero duplicative coverages, e.g. just over 3% of all fly genome bases were found to be duplicated, on average, when looking at any other genome.

A Code and Data Repository to Reproduce the Simulation Results

To aid future assessment we have created an easy to evaluate benchmarking pipeline (available at <http://compbio.soe.ucsc.edu/alignathon/>). Unfortunately the PSAR analysis involved using a compute cluster, making it expensive for outside groups to repeat this assessment. However, given a MAF file of one of the simulated datasets, the benchmarking pipeline can be used to make a performance assessment. The user can download the analysis repository, compile the necessary software, download the requisite data, place their alignment in a specified subdirectory and type “make” in the terminal window to launch the analysis. This approach will hopefully spur future development and assessment upon this resource.

Discussion

With the explosion in sequencing delivering ever larger numbers of near complete genome assemblies, WGA is an essential and increasingly important task. We have tested a total of 35 submissions from 12 different pipelines across three different datasets to produce the largest and most comprehensive assessment of WGA to date. The assessment purposefully chose test genomes in the 100-200 Mb range. The decision to use datasets at this size scale, rather than at the scale of larger vertebrate genomes, was balanced, just as in the first Assemblathon, by the desire to attract the largest possible number of entrants while still creating a reasonable challenge.

The primate simulations indicated that for closely related genome sequences aligners can find the vast majority of homologies accurately. In contrast, the simulated mammalian alignments showed a broader distribution of results. In concordance with this, we find that accuracies were substantially higher between more closely related genomes, and higher in more conserved regions, even in alignments also involving more distantly related genomes — this was apparent

both in looking at F-scores in the simulated mammals (Fig. 3), and pseudo F-scores in the flies (Fig. S3). We also find via simulation that more highly conserved sequence is easier to align, and that duplications are poorly represented in current WGA.

Testing using both simulations and real data we find a clear concordance between the rankings. In addition, using the simulated datasets we were able to demonstrate reasonable linear correlations, both overall and regionally, between F-scores and pseudo F-scores. This indicates that the high-level aggregate differences we highlight between the submissions can be found by two entirely independent means. We did not find a linear correlation between precision and the statistical measure of precision we used (PSAR-precision, see below), but we did find a very strong correlation between recall and coverage. Importantly, for the submissions we received on both flies and simulated mammals, differences in recall were overall greater than differences in precision, and therefore more critical in determining the observed performance differences. We speculate, given the various overlaps in the tools used between the pipelines and the resulting similarities between the submissions, that the larger differences in recall were largely due to differences in synteny mapping, and that this is one area where there is clear room for improvement. Visualising the data regionally, we were able to observe local differences in performance that fit well with these overall results.

In the simulations, for simplicity of interpretation, homologies that predated the MRCA of the extant genomes were not included in the “true” simulated alignments, therefore some ancient homologies captured by the aligners are considered false by the benchmarking pipeline. Additionally, EVOLVER does not track the alignments it generates when creating simulated mobile element (e.g. transposon) insertions, thus two highly similar transposon copies from separate insertion events are not considered homologous in the true alignment. For these reasons, the reported simulated precision values may be considered a lower bound that may, for some purposes (detecting ancient and mobile element mediated alignments), underestimate the accuracy of the alignments. This may partly explain the lack of correlation between precision and PSAR-precision, because such repeat regions might appear to be reasonably alignable to PSAR, but false homologies according to the simulation.

There were some dependencies between the assessments and the assessed programs. EVOLVER simulations were used to benchmark the Cactus aligner in its initial publication, though at approximately 1/250th the scale used here (Paten et al. 2011), and as part of two separate simulation assessments in that paper. It is therefore difficult to know if its substantial increase in relative performance is partially an artifact of training Cactus to the EVOLVER evolutionary model, though we note Cactus also performed well in the independent fly assessments. Similarly, PSAR uses the same pair-HMM alignment model as used by the PSAR-Align team in generating their alignments. We might expect therefore that the PSAR-Align alignments would be judged most accurate by PSAR, though we actually found a number of other programs earned equivalently high results.

The use of the MAF format for submissions made apples-to-apples comparison somewhat difficult because the format does not force transitivity of homology. However, this permissiveness in format allowed us to assess a variety of WGAs, some of which are naturally not transitively closed, such as the reference based MULTIZ aligners and the GenomeMatch pairwise submissions. To make comparisons under a strict homology assumption, we tried taking the transitive closure of such alignments, but this does not generally result in a reasonable WGA. In general, when performing consistent evolutionary analyses the nature of the alignment relation

— similarity or homology — bears consideration.

In theory several of the tools should have been able to align together duplicated regions. Unfortunately we received only submission for the simulation and not fly datasets for some duplication aware tools, such as EPO and Mugsy, but across the pipelines the lack of aligned duplications indicates there is likely significant room for WGA tools to improve in detecting duplicative homologies.

Some individual results on the simulations were surprising. For example, the EPO results had particularly low coverage on the simulations — substantially lower than that pipeline achieved in the genome alignments available from Ensembl (Paten et al. 2008). One possibility put forward by the authors of EPO is that the tools' reliance on using highly conserved sequences (“anchors”) for constructing a synteny map was not well suited to the simulations, which though modeled constraint, were likely different to the vertebrate genomes to which the EPO pipeline is normally applied.

Conclusion

Robust WGA tools are critical for the future of comparative genomics, but to make objective progress we must agree on assessments or risk not knowing when a genuine advance has been made. As is typical in bioinformatics, much prior assessment of WGAs has been made as part of the publication of a novel tool. Naturally these assessments tend to present results that favour the new tool. The few independent assessments that have been made of WGA, while useful, are several years old and assessed only a fraction of the available methods (Margulies et al. 2007, Chen & Tompa 2010). In comparison the Alignathon has been a success; it has leveraged a collaborative-competition model which has had broad community involvement and led to a broader set of WGAs being compared than in any prior attempt and certainly more than any single research group would likely have had the patience or expertise to handle.

Just as in any area where there is no accepted ground truth, comparing WGAs is hard, and each of the assessments types categorized in the introduction has flaws. Here we put an emphasis on two independent methods for assessment, and showed some consistency between them across datasets. For this reason, and that many of the results met with our prior expectation, we have some confidence in the results. Indeed, it is possible to compare the various F-score and pseudo F-scores (Figures 3-5, 7, S1-3, S7) and see some pipelines performed particularly strongly. However, given the uncertainty about the realism of the simulations and the apparent limited resolution of our statistical metrics we caution against over-interpretation. Therefore, more independent lines of assessment need to be developed: more simulations, more statistical assessments, and more assessments at different scales (e.g. full mammalian genomes), etc.

Assessments like the Alignathon are useful to spur community activity. However, these kind of benchmarking exercises risk becoming one-offs whose results are not comparable with the next generation of tools. To avoid this we have tried to make the simulation assessments developed here easily reusable so that they might be included in future publications. There is then a risk that tools may become overfit to these benchmarks, therefore updating the benchmarks periodically is essential.

In the supplement each of the teams describes how they computed each submission, which should be useful for reproducing their results. As the submissions were computed independently

by each team and each team had a different hardware environment we can not fairly compare the computational cost of the different pipelines. It would be useful for future efforts to assess this aspect, perhaps by getting groups to run their aligners on a common platform, such as Amazon EC2 or Microsoft Azure, where a controlled comparison could be made. This may prove to be an optimistic goal though, as many WGA pipelines are designed and implemented at individual institutions by researchers whose goal is the sharing of the results of the pipeline but not the pipeline itself. The computation-environmental peculiarities of individual institutions can thereby be reflected in their pipelines through unintentional design.

We have left a number of questions unresolved. For example, we have not attempted to determine how tools for WGA compare to methods for other types of MSA, such as protein aligners, or how the quality of the input genome assemblies affects WGAs. In summary, we very much hope that the Alignment will help pave the way for subsequent efforts with more datasets, comparisons, accurate statistical assessments and benchmarking exercises with even broader scope.

Methods

Simulations

As in the Assemblathon 1 project (Earl et al. 2011) simulated genomes were generated using the EVOLVER suite of tools forward-time whole genome evolution simulation tools (Edgar et al. <http://www.drive5.com/evolver/>). Specific parameter files used to create the simulations are available on the project website. EVOLVER has a model for proteins, genes and base-level evolutionary constraints. EVOLVER uses a two step process for simulating a single forward step in a simulation: the first step is an intra-chromosomal evolution step, and the second an inter-chromosomal step. The intra-chromosomal step allows events such as substitutions, insertions and deletions, duplications, translocations, and inversions, according to rates distributed according to the length of the event. The inter-chromosomal step allows chromosome fusions, fissions, segment copying, segment movement, reciprocal translocations and non-reciprocal translocations. Additionally EVOLVER keeps a separate mobile element library that can insert mobile element DNA into the simulated genome; this library is itself also undergoing simulated evolution. EVOLVER logs all evolutionary events that take place during a cycle and keeps track of the relationships between residues in the parent and child genomes.

EVOLVER as distributed is only capable of performing a single cycle of evolution. In order to run the arbitrary phylogenies necessary for this project we used the `evolverSimControl` and `evolverInfileGeneration` tools available at <https://github.com/dentearl/evolverSimControl> and <https://github.com/dentearl/evolverInfileGeneration/>, respectively. These extra tools, along with `mafJoin` <https://github.com/dentearl/mafJoin/> were used to construct MAF files containing the entire simulated true evolutionary relationships of all of the genomes: leaves, internal nodes, and the root.

As in the Assemblathon project, we initiated the simulation using a subset of the well annotated human genome, hg19/GRCh37. Complete chromosome sequences for chromosomes 20, 21 and 22 along with annotations for those chromosomes from the UCSC Genome Browser tracks `mgcGenes`, `knownGene`, `knownGeneOld5`, `cpgIslandExt`, and `ensGene` were obtained from the UCSC golden path download site. The tool suite `evolverInfileGeneration` was used to take the raw data and make it into an EVOLVER infile dataset. This starting dataset was then put through the EVOLVER simulator for a distance of 1.0 neutral substitutions per site, an evolutionary time of

approximately 500 million years of vertebrate evolution (Hedges et al. 2006, Earl et al. 2011, Fujita et al. 2011). This process, which we term a burn-in, shuffles the sequences, genes and chromosomes of the genome. The resulting genome was termed the most recent common ancestor (MRCA) because it was used as the starting point for both the primate and mammalian simulations. It has been previously ascertained that distributions on the numbers and lengths of tracked annotation types in EVOLVER simulations stay stationary over time (Earl et al. 2011), so this burn-in process, from a simulation point of view, does not adversely affect the nature of the simulated genomes.

EVOLVER rediscovers the tandem repeat sequence annotation at every step of a simulation by calling tandem repeats finder (Benson 1999, v4.0). RepeatMasker (Smit et al. 2010, Smit & Hubley 2010, v1.25) and tandem repeats finder were used to identify and soft-mask repetitive sequence in the final leaf genomes.

The primate simulation was described by the phylogenetic tree (in newick format; Fig. 1):

```
((simGorilla:0.008825,(simHuman:0.0067,simChimp:0.006667)sHuman-sChimp:0.00225)sG-sH-sC:0.00968,simOran
g:0.018318);
```

The mammal simulation was described by the phylogenetic tree (in newick format; Fig. 1):

```
((simCow:0.18908,simDog:0.16303)sCow-sDog:0.032898,(simHuman:0.144018,(simMouse:0.084509,simRat:0.0915
89)sMouse-sRat:0.271974)sH-sM-sR:0.020593);
```

We used the EVOLVER produced repetitive element library from the simHuman genome as an input library for RepeatMasker. Following each simulation, the EVOLVER mobile element library from the simHuman leaf node genome was used as an input into the repetitive sequence finder RepeatMasker. RepeatMasker was then used to mask simple repeats and repeats from the provided library in the other non-human simulated genomes.

Complete sequence and annotations of the leaf genomes and the MRCA genome were provided to participants.

Flies Dataset

The phylogeny was created by merging the phylogeny provided in the modENCODE (The modENCODE Consortium et al. 2010) comparative genomics white paper (http://www.genome.gov/Pages/Research/Sequencing/SeqProposals/modENCODE_ComparativeGenomics_WhitePaper.pdf accessed 15 October 2013) courtesy of Artyom Kopp (UC Davis) and the phylogeny used by UCSC for the 15-way insect alignment. The Kopp tree lacked droSim1 and droSec1 which were added by normalizing the branch lengths between the dm3 branches on the two trees. Extraneous species were trimmed using [tree_doctor](#) from PHAST. This tree was provided for progressive aligners that need a guide tree. This phylogeny corresponds to the newick tree (Fig. 1):

```
((droGri2:0.183954,droVir3:0.093575):0.000000,(droMoj3:0.110563,(((droBip:0.034265,droAna3:0.042476):0.121927,
(droKik:0.097564,((droFic:0.109823,(((dm3:0.023047,(droSim1:0.015485,droSec1:0.015184):0.013850):0.016088,(dro
Yak2:0.026909,droEre2:0.029818):0.008929):0.047596,(droEug:0.102473,(droBia:0.069103,droTak:0.060723):0.0158
55):0.005098):0.010453):0.008044,(droEle:0.062413,droRho:0.051516):0.015405):0.046129):0.018695):0.078585,(dro
Per1:0.007065,dp4:0.005900):0.185269):0.068212,droWil1:0.259408):0.097093):0.035250);
```

To create the fly sequence dataset we took 12 flies available from the UCSC golden path server

on 14 December 2011 (droAna3, dreEre2, droGri2, droMoj3, dp4, droVir3, droWil1, dm3, droSim1, droPer1, droSec1, droYak2) and eight flies from NCBI on 25 January 2012 (droBia, droBip, droEle, droFic, droKik, droTak, droRho, droEug).

MafTools

Participants submitted their predictions of alignments in MAF files. To process the submissions we wrote a suite of open-source tools called mafTools, available at <https://github.com/dentearl/mafTools/> to perform the majority of transformations, manipulations and analyses. Scripts to perform the analyses described can be found in the analysis repository.

MAF Comparisons

Exhaustively checking all pairs of aligned residues between alignments is computationally impractical, so instead we developed a method, termed *mafComparator*, based upon sampling pairs of aligned residues. Sampling is performed by reading each input MAF file twice, once to count the total number of pairs present in the file, such that given a user specified number of pairs to sample we can calculate the probability of picking a given pair at random. The MAF file is then read a second time.

During the second pass we iterate over every block in the MAF and then every column in the block. We calculate the number of pairs present in the block, call this k , and then make a draw from a binomial distribution with probability s / m (where s is the number of samples taken, here 10,000,000, and m is the total number of pairs present in the MAF) to see how many (if any) pairs to sample from that column. If x many pairs are to be sampled we then sample x times from a discrete uniform $[0, k - 1]$ decrementing the range of the distribution with each sample, without replacement, and then map those integers to pairs using a bijective function. This allows us to efficiently sample pairs without iterating through each and every pair.

Regional Alignments & PSAR

To accommodate PSAR, which processes global MSAs in which the alignment is represented as a 2D matrix where the aligned sequences, interspersed with gaps, are the rows and the columns represent the equivalence classes of aligned bases, we constructed sub-alignments of sampled regions.

Regions were randomly sampled by picking intervals of a chosen reference genome (for the flies *D. melanogaster*, dm3, and for the simulations simHuman). For each of the three test sets, regional intervals were selected by sampling five different starting values from a discrete uniform distribution $[0, g - 1 - 500,000]$, where g is the total length of the reference genome and 500,000 is the length of the interval. Sampled values were then mapped back to individual chromosomes. All alignments containing any positions of the reference within these intervals were extracted from the submitted alignments. Though this model of sampling does not prevent overlapping regions, no overlapping regions were sampled. Likewise this model of sampling does not prevent regions that cross between chromosomes, but no such bridged regions were sampled.

Details of how we adjusted each submission for regional analysis with PSAR are in the supplement.

Data Access

The project website is available at <http://compbio.soe.ucsc.edu/alignathon/>. The website links to all of the datasets, submissions and benchmarking code.

Acknowledgements

We would like to thank the Howard Hughes Medical Institute, Dr. and Mrs. Gordon Ringold, NIH grant 2U41 HG002371-13 and NHGRI/NIH grant 5U01HG004695 for providing funding. We would like to thank the Genome 10K organizers for providing a venue to discuss an early version of these findings. We thank three anonymous referees who provided a great deal of helpful feedback and suggestions.

Disclosure Declaration

We have no conflicts of interest.

Figure Captions

FIGURE 1 — The phylogenies of the three test sets: primate simulation, mammal simulation and fly real data set. Branch lengths are in units of neutral substitutions per site.

FIGURE 2 — Simulated primate and mammal F-score results. Recall as a function of precision is shown for primates in A and mammals in C. GenomeMatch-3 is omitted from plot A as both of its values are low (see its overall F-score in B). B shows the primate F-score results isolated to different annotation types: overall, genes, neutral and repetitive regions. D shows the mammal version of C. Legends for B and D are ordered as in the overall category and this order is maintained in the genes and neutral annotations.

FIGURE 3 — Primate and Mammal simulation F-score results stratified by phylogenetic distance. For each subplot the vertical axis shows the F-score and the horizontal axis shows 13 individual submissions ordered from left to right (descending) by average overall F-score. Horizontal gray lines show the overall F-score of the submission taking into account all sequence pairs. Horizontal black lines show the overall F-score of the submission taking into account only sequences pairs including the reference. Submissions are comprised of points connected by a line where the points are in ascending order of phylogenetic distance (all possible pairs are shown).

FIGURE 4 — Simulated mammal results comparing simulation values to statistical values. Shown is precision and PSAR precision; recall and coverage; F-score and pseudo F-score. Each column represents the results of one submission, columns are in descending order of overall (full genome) F-score value. The horizontal line is, respectively, the overall precision, recall or F-score value; the up-pointing triangle with a vertical line is the regional precision, recall or F-score mean value, +/- the

regional standard deviation; the down-pointing triangle with a vertical line is the PSAR-precision, coverage or pseudo F-score mean value, +/- standard deviation for values that were computed using regional sub-alignments.

FIGURE 5 — Fly results: values of PSAR-precision; average overall coverage between all pairs; pseudo F-score. Columns are in descending order of mean pseudo F-score value. For each metric, each submission is made up of an orange circle with a vertical line representing the regional mean +/- stdev.

FIGURE 6 — Overall pairwise coverage values in the flies dataset. Submissions are ordered left to right (descending) by overall coverage. Gray points are non-reference pairs, black points contain the reference. The horizontal gray line shows the average coverage of the submission for all points, the horizontal black line shows the average coverage of the submission for just pairs containing the reference. Beneath the pairwise coverage plot is a barcode plot showing the phylogenetic distances of all pairs, shorter gray lines are non-reference pairs and longer black lines are reference containing pairs.

FIGURE 7 — Region 2 of *D. melanogaster* (dm3) with respect to *D. grimshawi* (droGri2) of the regional analysis of the mammal simulation data set. Region 2 is defined as bases 12,450,223 - 12,950,222 of dm3 chromosome 3R (horizontal axis). Rows are: the relative abundance of genes within the region; the relative abundance of repetitive sequence in the region; submissions in descending order of average pseudo F-score. Each submission row shows the pseudo F-score of the submission in black. The vertical axis of each row is the same scale, as labeled in the bottom row. In grey, in the background, is shown the pseudo F-score value of the top submission for this region (TBA).

FIGURE 8 — The Jaccard distance (1 - Jaccard similarity coefficient) matrix and accompanying hierarchical clustering (UPGMA) of submissions for each of the three test sets. Higher values indicate that the sets of aligned pairs of two submissions are more dissimilar, lower values indicate similarity.

Figures

FIGURE 1

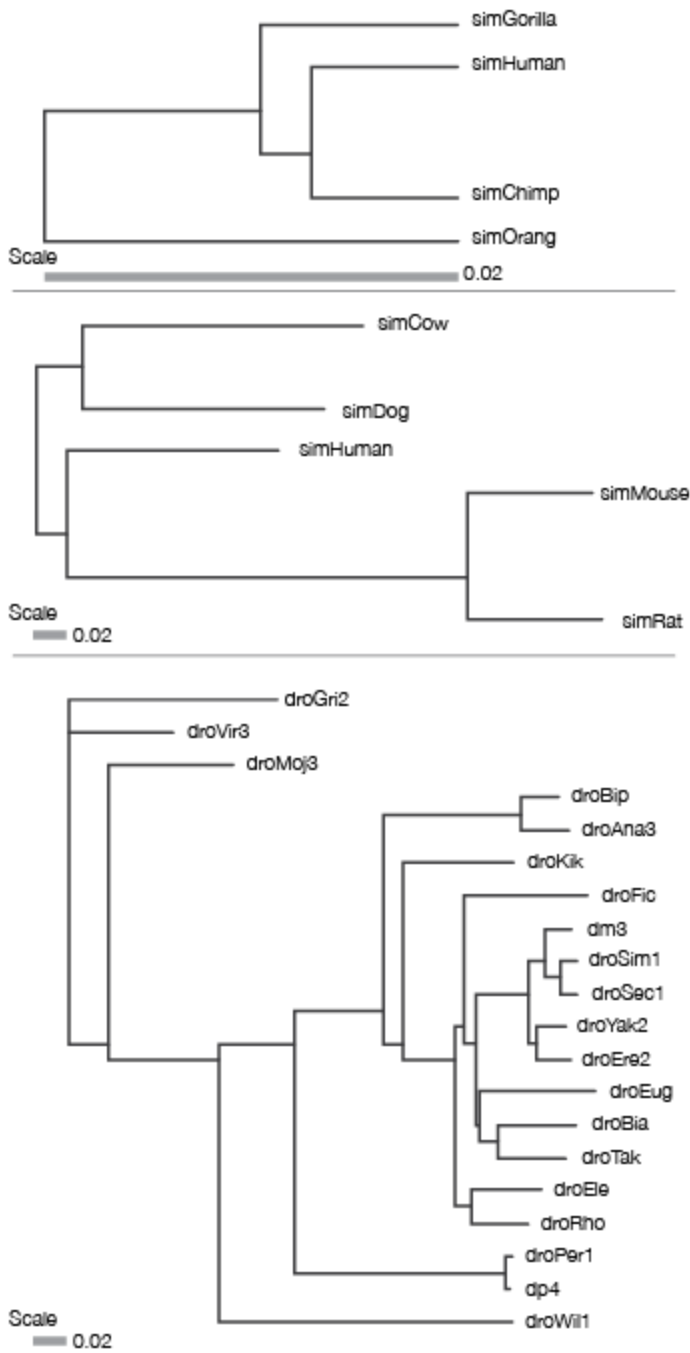


FIGURE 2

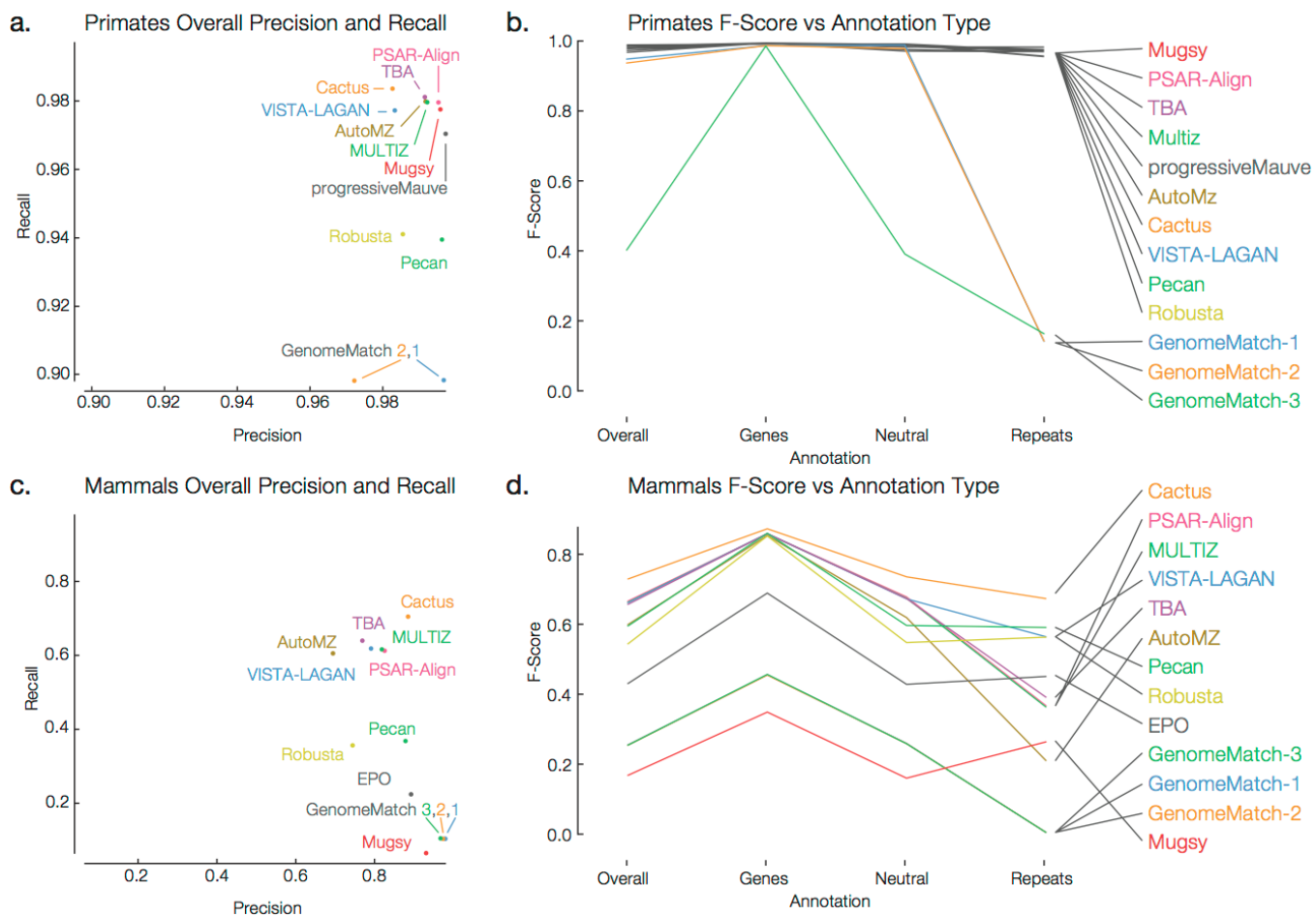


FIGURE 3

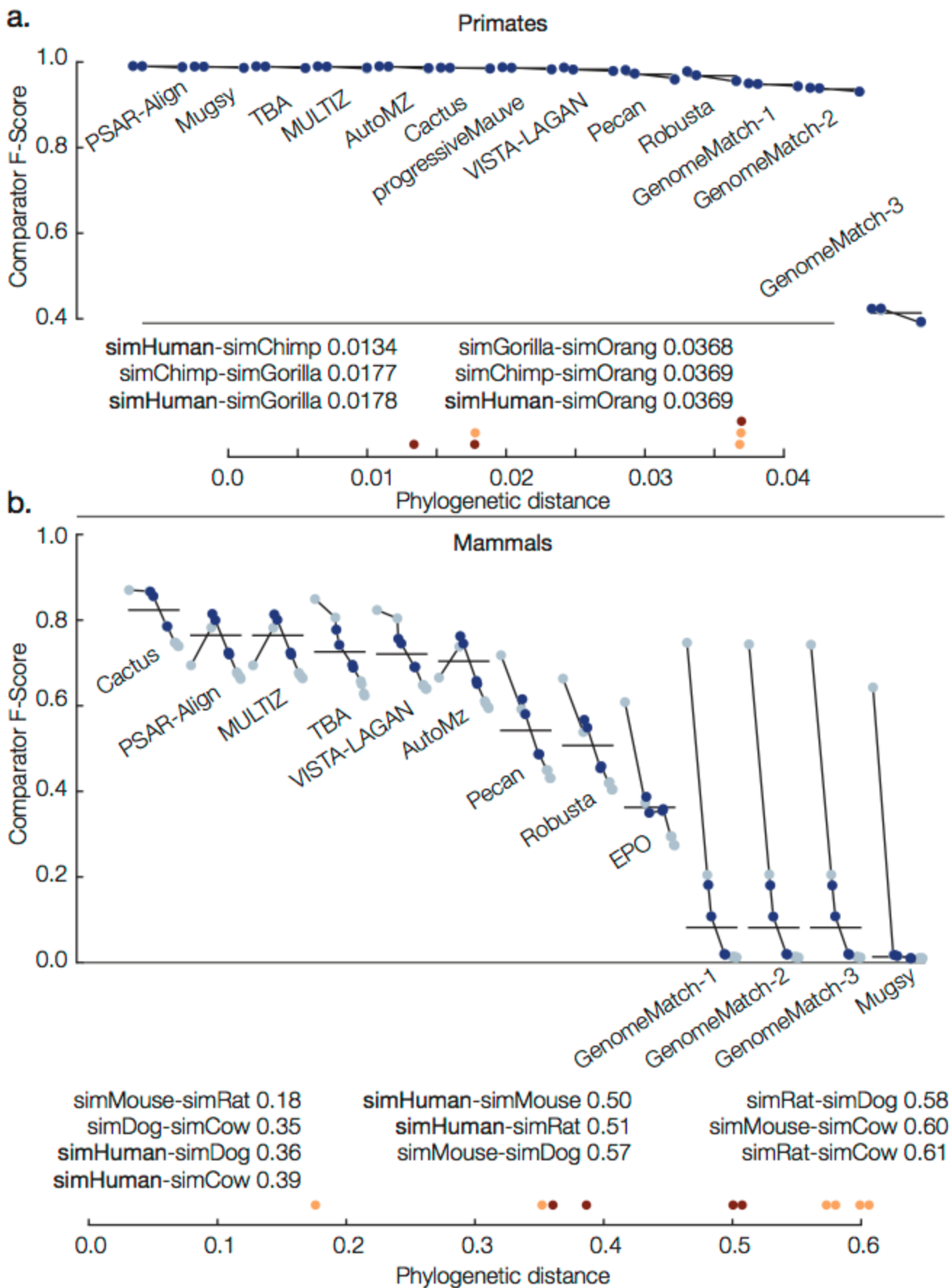


FIGURE 4

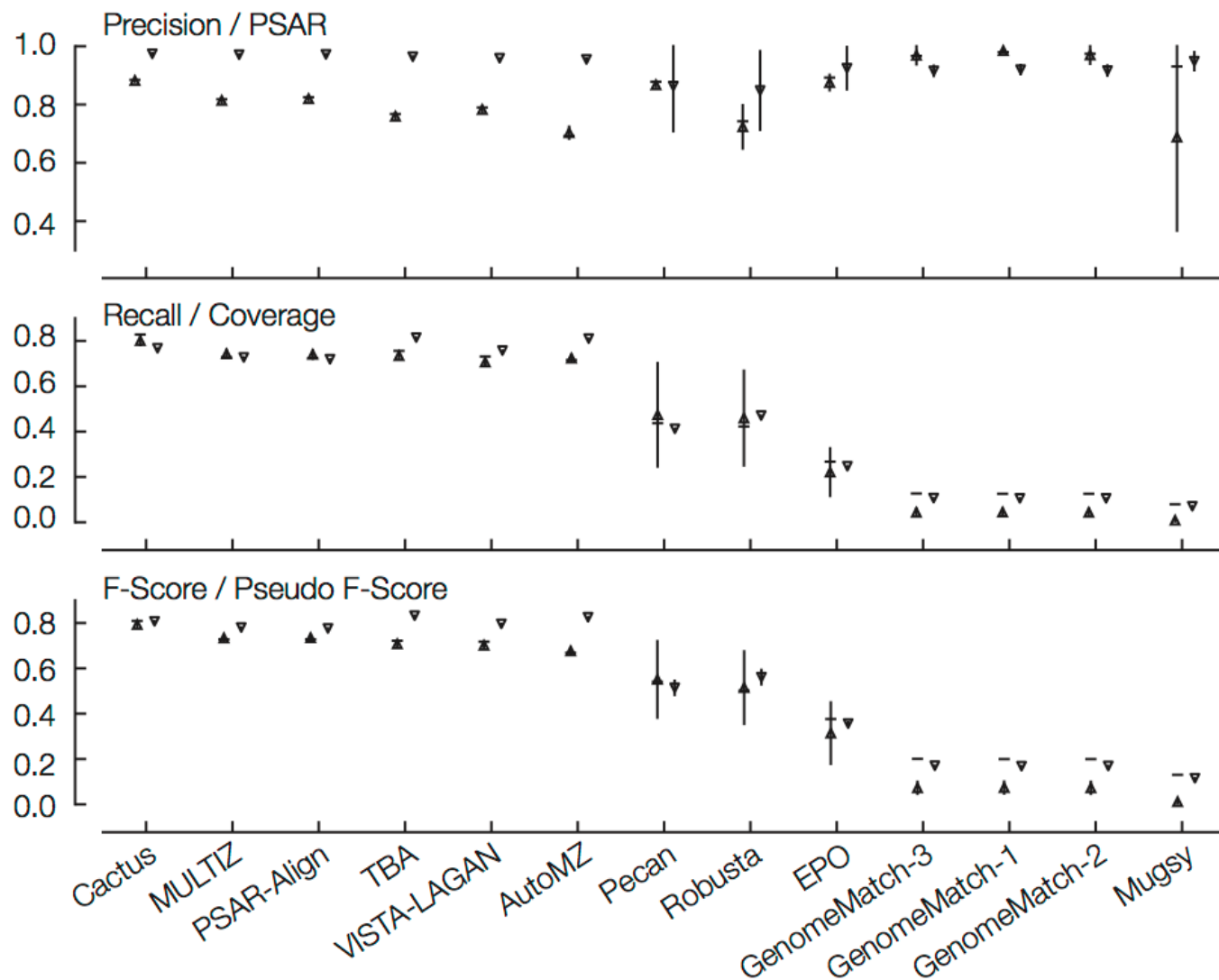


FIGURE 5

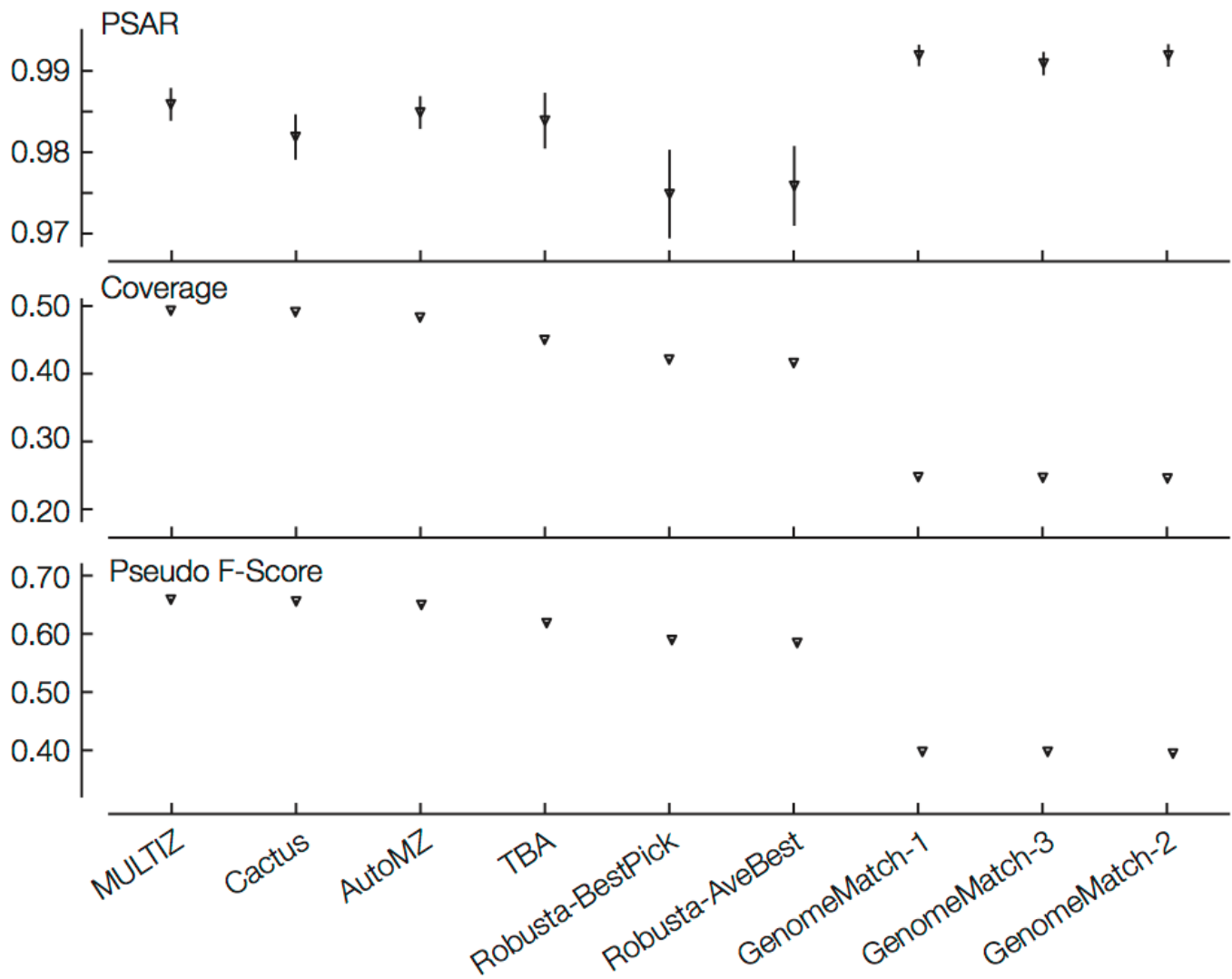


FIGURE 6

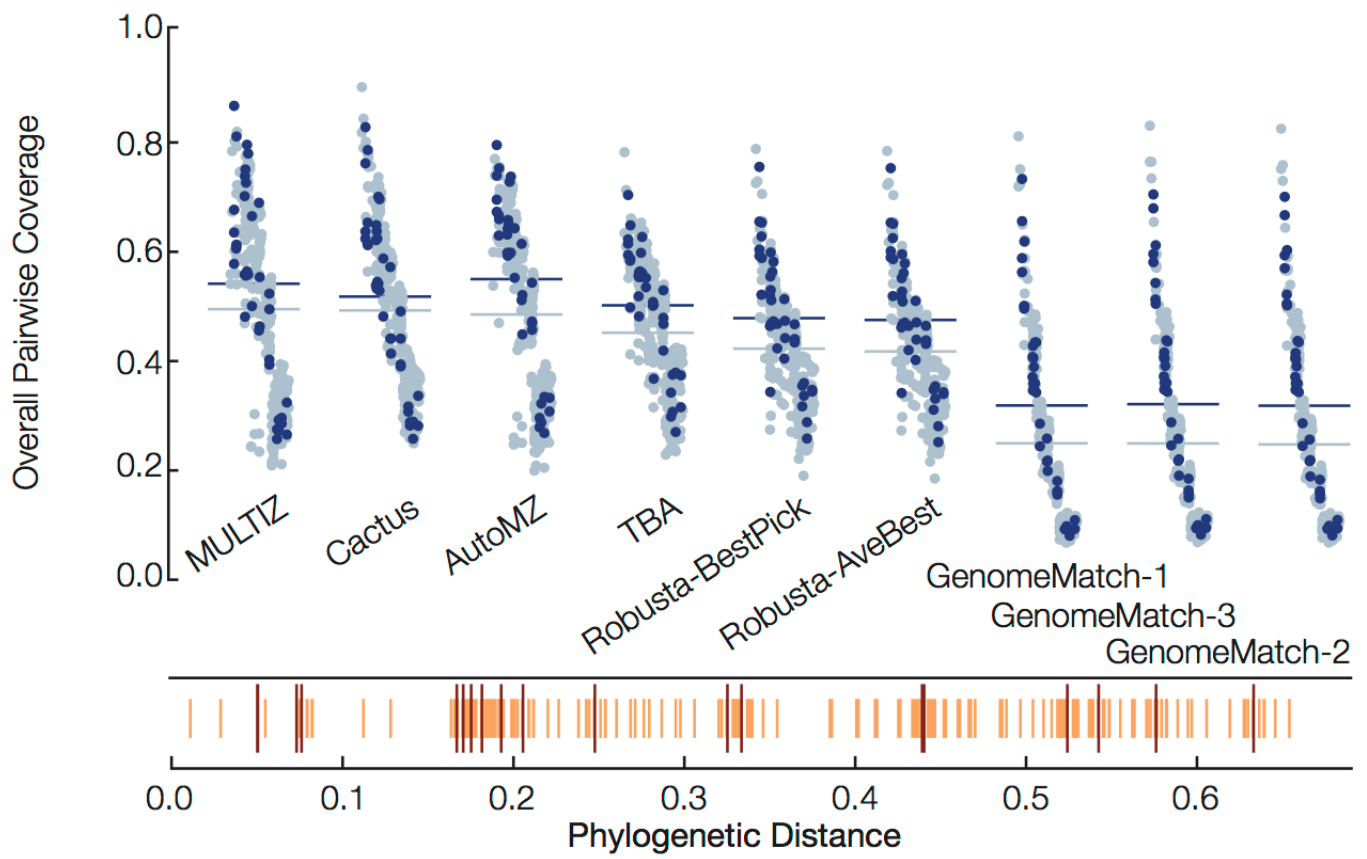


FIGURE 7

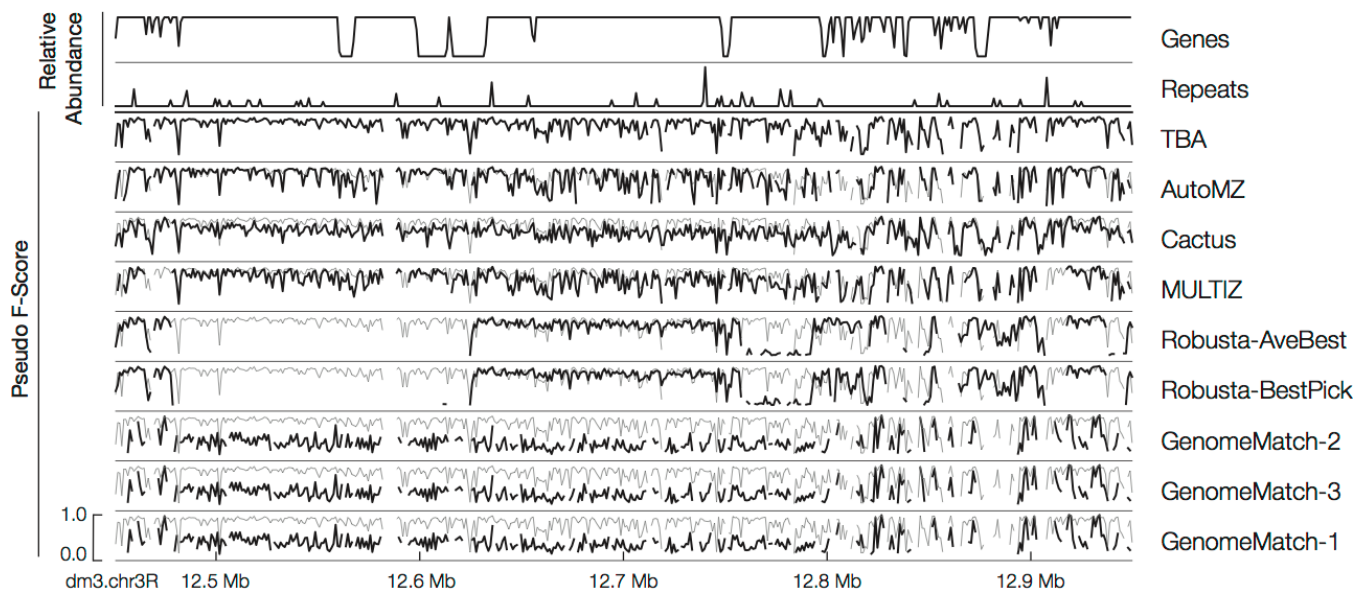
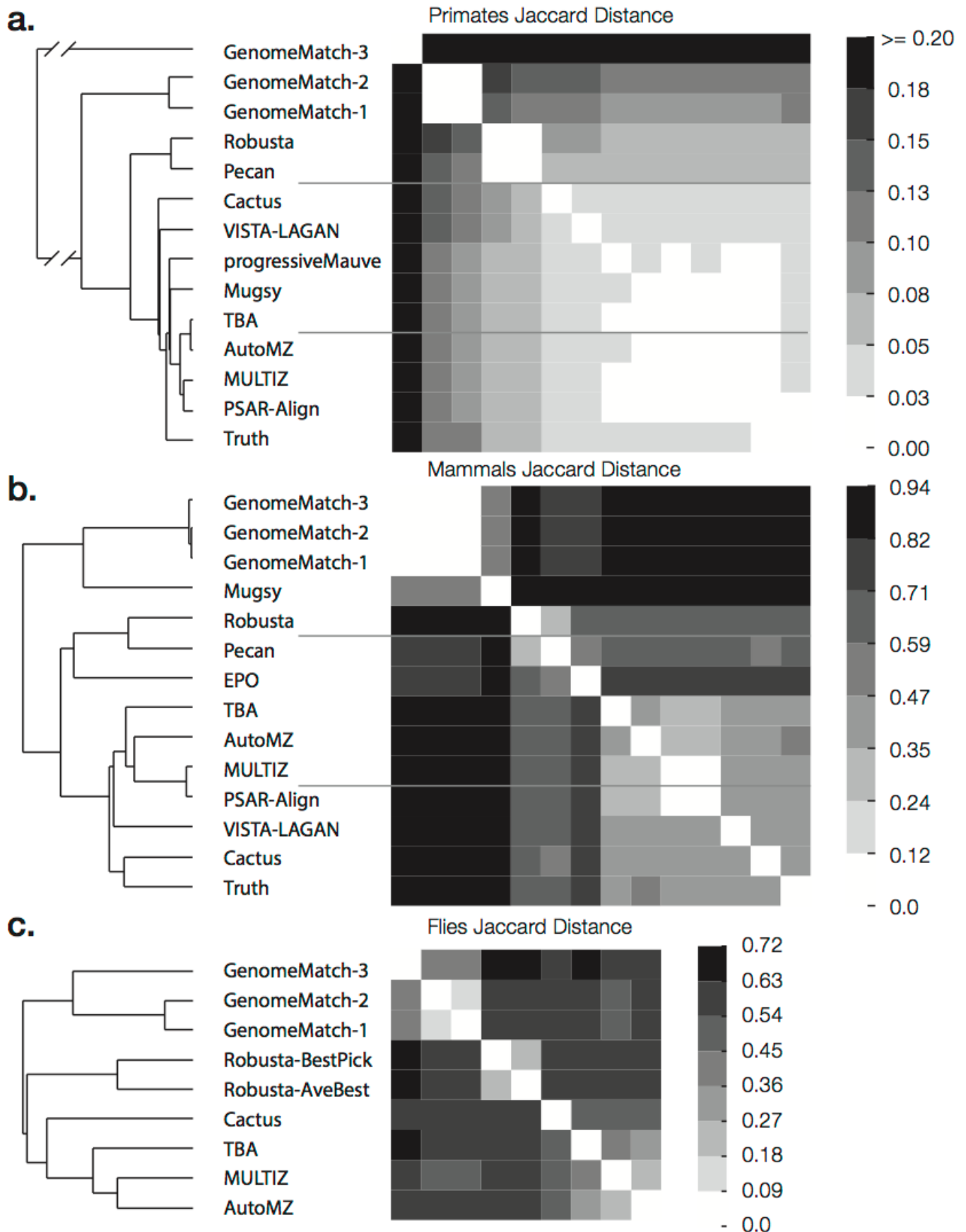


FIGURE 8



Tables

| Simulation | Genome | #Chrs | Max | Min | Mean | Median | GC% | Length (Mb) | Dist to ref |
|------------|------------|-------|------------|------------|------------|--------------|-------|-------------|-------------|
| Primate | simChimp | 4 | 85,778,862 | 10,574,168 | 46,284,070 | 44,391,624.5 | 41.3% | 185.1 | 0.016637 |
| | simGorilla | 4 | 85,848,133 | 10,570,608 | 46,298,606 | 44,387,841.0 | 41.3% | 185.2 | 0.017775 |
| | simHuman | 4 | 85,835,872 | 10,572,275 | 46,286,362 | 44368649.50 | 41.3% | 185.1 | 0 |
| | simOrang | 5 | 85,903,762 | 10,564,720 | 370,67,732 | 35,683,973.0 | 41.3% | 185.3 | 0.036948 |
| Mammal | simCow | 5 | 86,443,571 | 6,172,747 | 38,605,187 | 33,408,597.0 | 41.1% | 193.0 | 0.386589 |
| | simDog | 6 | 64,906,724 | 5,551,284 | 32,033,832 | 30,919,174.0 | 41.1% | 192.2 | 0.360539 |
| | simHuman | 6 | 88,398,963 | 2,880,482 | 31,799,332 | 22,095,903.5 | 41.1% | 190.8 | 0 |
| | simMouse | 7 | 71,158,916 | 3,949899 | 28,417,345 | 16,897,397.0 | 41.1% | 198.9 | 0.500501 |
| | simRat | 5 | 88,137,694 | 4,060,565 | 39,740,941 | 45,269,609.0 | 41.1% | 198.7 | 0.507581 |

TABLE 1 — Summary statistics for the simulated genomes. Rows are leaf genomes generated by simulation. Columns are different metrics: #Chrs is the number of chromosomes; Max is the longest chromosome; Min is the shortest chromosome; Mean is the average length of a chromosome; median is the median chromosome length; GC% is the percent GC composition of the genome; Length is the total length of the genome in megabases (Mb); Dist to ref is the phylogenetic distance from the leaf to the reference species (named simHuman in both simulations).

| Genome | #Seqs | Max | Min | Mean | Median | GC% | Est. Length (Mb) | NG50 | Dist to dm3 |
|------------------------------------|--------|------------|--------|------------|-----------|-------|------------------|------------|-------------|
| <i>Drosophila melanogaster</i> dm3 | 15 | 29,004,656 | 19,517 | 12,027,821 | 6,668,899 | 41.7% | 176** | 23,011,544 | 0 |
| <i>D. pseudoobscura</i> dp4 | 4896 | 30,794,189 | 101 | 31,202 | 1,734 | 45.2% | 161** | 11,692,001 | 0.439806 |
| <i>D. ananassae</i> droAna3 | 13,749 | 23,697,760 | 55 | 16,801 | 1,517 | 42.0% | 190** | 5,790,199 | 0.334455 |
| <i>D. biarmipes</i> droBia | 7,864 | 2,773,931 | 109 | 21,437 | 1,756 | 42.0% | 195** | 354,307 | 0.176787 |
| <i>D. bipectinata</i> droBip | 8,681 | 1,366,469 | 69 | 19,169 | 4,348 | 42.0% | 195** | 79,801 | 0.326244 |
| <i>D. elegans</i> droEle | 8,393 | 1,324,493 | 74 | 20,318 | 1,713 | 40.3% | 190** | 166,643 | 0.183046 |
| <i>D. erecta</i> droEre2 | 5,124 | 26,641,161 | 154 | 29,808 | 1,729 | 42.2% | 156** | 18,748,788 | 0.077882 |
| <i>D. eugracilis</i> droEug | 7,566 | 1,857,654 | 71 | 20,660 | 1,927 | 40.9% | 225** | 61,086 | 0.194302 |
| <i>D. ficusphila</i> droFic | 9,164 | 1,642,143 | 112 | 16,484 | 1,649 | 41.9% | 180** | 163,015 | 0.207007 |
| <i>D. grimshawi</i> droGri2 | 17,440 | 24,565,398 | 80 | 11,494 | 1,702 | 38.0% | 231* | 6,267,026 | 0.633146 |
| <i>D. kikkawai</i> droKik | 8,344 | 1,394,654 | 70 | 19,593 | 1,871 | 41.4% | 205** | 133,834 | 0.248921 |
| <i>D. mojavensis</i> droMoj3 | 6,841 | 34,148,556 | 101 | 28,336 | 1,654 | 39.5% | 166** | 26,866,924 | 0.524505 |
| <i>D. persimilis</i> droPer1 | 12,838 | 11,822,988 | 206 | 14,674 | 1,671 | 44.9% | 180** | 1,930,428 | 0.440971 |
| <i>D. rhopaloa</i> droRho | 34,038 | 561,403 | 65 | 5,696 | 1,465 | 40.0% | — | 19,476† | 0.172149 |
| <i>D. sechellia</i> droSec1 | 14,730 | 21,120,651 | 207 | 11,309 | 1,710 | 42.1% | 171** | 2,104,621 | 0.052081 |
| <i>D. simulans</i> droSim1 | 18 | 27,517,382 | 14,972 | 8,371,773 | 2,996,586 | 42.5% | 152** | 19,596,830 | 0.052382 |
| <i>D. takahashii</i> droTak | 9,700 | 1,026,890 | 96 | 18,661 | 3,477 | 40.0% | 195** | 109,442 | 0.168407 |
| <i>D. virilis</i> droVir3 | 13,530 | 25,233,164 | 43 | 15,216 | 1,215 | 40.0% | 332** | 510,240 | 0.542767 |
| <i>D. willistoni</i> droWill1 | 14,927 | 16,660,200 | 868 | 15,857 | 1,508 | 37.2% | 222* | 4,707,319 | 0.576257 |
| <i>D. yakuba</i> droYak2 | 21 | 28,832,112 | 16,019 | 8,427,486 | 2,539,874 | 42.3% | 166** | 22,324,452 | 0.074973 |

TABLE 2 — Summary statistics for the fly genomes. For each species provided to participants the following information is shown: number of sequences that comprise the genome; the maximum length of a sequence; the minimum length of a sequence; the mean length of all sequences; the median length of all sequences; the percent GC content of the genome; the estimated size in megabases; the NG50 value of the genome (Earl et al. 2011); the phylogenetic distance to *D. melanogaster* (dm3). *Genome size estimates from Bosco et al. 2007. **Genome size estimates from Gregory & Johnston 2008. †N50 instead of NG50 value due to lack of genome size estimate.

| Submission | Tools | Submitter(s) | Primate Simulation | Mammal Simulation | Fly data set (real) |
|------------------|------------------------|---|--------------------|-------------------|---------------------|
| AutoMZ | AutoMZ | Minmei Hau | 1 | 1 | 1 |
| Cactus | Cactus | Benedict Paten, Glenn Hickey | 1 | 1 | 1 |
| EPO | Enredo, Pecan, Ortheus | Stephen Fitzgerald, Kathryn Beal, Javier Herrero | — | 1 | — |
| Pecan | Mercator, Pecan | Stephen Fitzgerald, Kathryn Beal, Javier Herrero | 1 | 1 | — |
| GenomeMatch | Genome Match | Igor Seledtsov, Vladimir Molodtsov, PI: Victor Solovyev | 3 | 3 | 3 |
| Mugsy | Mugsy | Aaron E. Darling | 1 | 1 | — |
| MULTIZ | MULTIZ | Brian J. Raney | 1 | 1 | 1 |
| PSAR-Align | PSAR-Align | Jaebum Kim, Jian Ma | 1 | 1 | — |
| progressiveMauve | progressiveMauve | Aaron E. Darling | 1 | — | — |
| Robusta | Robusta | Carsten Kemena, Jia-Ming Chang, Ionas Erb, Cedric Notredame | 1 | 1 | 2 |
| TBA | TBA | Minmei Hau | 1 | 1 | 1 |
| VISTA-LAGAN | VISTA-LAGAN | Alexander Poliakov, Michael Brudno, Inna Dubchak | 1 | 1 | — |
| TOTAL | | | 13 | 13 | 9 |

TABLE 3 — Submissions to the Alignathon. Each row shows a tool with: the name of the submission as used in this paper; the names of the submitters; the number of submissions from the tool for the primate data set; the number of submissions from the tool for the mammal data set; the number of submissions from the tool for the fly data set.

References

Angiuoli SV., Salzberg SL. 2011. Mugsy: fast multiple alignment of closely related whole genomes. *Bioinformatics* 27(3):334-42.

Beiko RG, Charlebois RL. 2007. A simulation test bed for hypotheses of genome evolution. *Bioinformatics*. 23:825-831.

Beitzel S M. 2006. On Understanding and Classifying Web Queries (Ph.D. thesis). IIT. CiteSeerX: 10.1.1.127.634

Benson G. 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Research* 27(2):573-580.

Blackshields G, Wallace IM, Larkin M, Higgins DG. 2006. Analysis and comparison of benchmarks for multiple sequence alignments. *In Silico Biology* 6: 321-339.

Blanchette M., Kent W.J., Riemer C., Elnitski L., Smit A.F., Roskin K.M., Baertsch R., Rosenbloom K., Clawson H., Green E.D., et al. 2004. Aligning multiple genomic sequences with the threaded blockset aligner. *Genome Res.* 14: 708-715.

Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, Boisvert S, Chapman JA, Chapuis G, Chikhi R, et al. 2013. Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species. *Gigascience.* 2(1):10.

- Branton D, Deamer DW, Marziali A, Bayley H, Benner SA, Butler T, Di Ventra M, Garaj S, Hibbs A, Huang X, et al. 2008. The potential and challenges of nanopore sequencing. *Nature Biotechnology* 26(10): 1146-53.
- Cartwright RA. 2005. DNA Assembly With Gaps (Dawg): Simulating Sequence Evolution. *Bioinformatics* 21 (Suppl. 3): iii31-iii38.
- Carvajal-Rodríguez A. 2010. Simulation of genes and genomes forward in time. *Curr Genomics*. 11(1):58-61. doi: 10.2174/138920210790218007.
- Chen X., Tompa M. 2010. Comparative assessment of methods for aligning multiple genome sequences. *Nature Biotechnology*, 28(6): 567-572. doi:10.1038/nbt.1637
- Dalquen D.A., Anisimova M., Gonnet G.H., Dessimoz C. 2012. ALF—A Simulation Framework for Genome Evolution. *Mol Biol Evol* 29 (4): 1115-1123 first published online December 8, 2011 doi:10.1093/molbev/msr268
- Darling ACE, Mau B, Blattner FR, Perna NT. 2004. Mauve: multiple alignment of conserved genomic sequence with rearrangements. *Genome research* 14, no. 7: 1394-403.
- Darling ACE, Mau B, Perna NT. 2010. progressiveMauve: multiple genome alignment with gene gain, loss and rearrangement. *PLoS One*. 5(6):e11147. doi: 10.1371/journal.pone.0011147.
- Dewey CN. 2007. Aligning multiple whole genomes with Mercator and MAVID. *Methods Mol Biol*. 2007;395:221-36.
- Drosophila 12 Genomes Consortium. 2007. Evolution of genes and genomes on the Drosophila phylogeny. *Nature* 450(7167)203-18. doi:10.1038/nature06341
- Dubchak I, Poliakov A, Kislyuk A, Brudno M. 2009. Multiple whole-genome alignments without a reference organism. *Genome Research* 19(4):682-9.
- Durbin R, Eddy S, Krogh A, Mitchison G. 1998. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge University.
- The ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the human genome. *Nature*. 2012 Sep 6;489(7414):57-74.
- Earl D, Bradnam K, St John J, Darling A, Lin D, Fass J, Yu HO, Buffalo V, Zerbino DR, Diekhans M, et al. 2011. Assemblathon 1: a competitive assessment of de novo short read assembly methods. *Genome Res*. 21(12):2224-41. doi: 10.1101/gr.126599.111.
- Edgar RC, Batzoglou S. 2006. Multiple sequence alignment. *Current Opinion in Structural Biology*. 16(3):368-373.
- Edgar R, Asimenos G, Batzoglou S, Sidow A. 2009. EVOLVER, <http://www.drive5.com/evolver/> accessed 2 January 2014.

- Felsenstein, J. 2003. *Inferring Phylogenies*. Sinauer Associates, Sunderland, Massachusetts.
- Flicek P, Ahmed I, Amode MR, Barrell D, Beal K, Brent S, Carvalho-Silva D, Clapham P, Coates G, Fairley S, et al. 2013. Ensembl 2013. *Nucleic Acids Res.* 41(Database issue):D48-55. doi: 10.1093/nar/gks1236. Epub 2012 Nov 30.
- Frazer KA, Pachter L, Poliakov A, Rubin EM, Dubchak I. 2004. VISTA: computational tools for comparative genomics. *Nucleic Acids Research* 32(Web Server issue):W273-9.
- Fujita PA, Rhead B, Zweig AS, Hinrichs AS, Karolchik D, Cline MS, Goldman M, Barber GP, Clawson H, Coelho A, et al. 2011. *Nucleic Acids Research* 39(Database issue):D876-82.
- Genome 10K Community of Scientists. 2009. Genome 10K: a proposal to obtain whole-genome sequence for 10,000 vertebrate species. *J Hered.* 100(6):659-74. doi: 10.1093/jhered/esp086. Epub 2009 Nov 5.
- Hedges SB, Dudley J & Kumar. 2006. TimeTree: a public knowledge-base of divergence times among organisms. *Bioinformatics* 22:2971-2972.
- Holmes I, Durbin R. 1998. Dynamic programming alignment accuracy. *J Comput Biol.* 5(3):493-504.
- i5K Consortium. 2013. The i5K Initiative: advancing arthropod genomics for knowledge, human health, agriculture, and the environment. *J Hered.* 104(5):595-600. doi: 10.1093/jhered/est050.
- Iantorno S, Gori K, Goldman N, Gil M, Dessimoz C. 2014. Who watches the watchmen? An appraisal of benchmarks for multiple sequence alignment. *Methods Mol Biol.* 1079:59-73. doi: 10.1007/978-1-62703-646-7_4.
- Kellis M, Patterson N, Endrizzi M, Birren B, Lander ES. 2003. Sequencing and comparison of yeast species to identify genes and regulatory motifs. *Nature* 423(6937):241-54
- Kemena C, Notredame C. 2009. Upcoming challenges for multiple sequence alignment methods in the high-throughput era. *Bioinformatics* 27(19): 2455-2465. doi:10.1093/bioinformatics/btp452.
- Kemena C, Bussotti G, Capriotti E, Marti-Renom MA, Notredame C. 2013. Using tertiary structure for the computation of highly accurate multiple RNA alignments with the SARA-Coffee package. *Bioinformatics.* 29(9):1112-9. doi: 10.1093/bioinformatics/btt096. Epub 2013 Feb 28.
- Kim J, Ma J. 2011. PSAR: measuring multiple sequence alignment reliability by probabilistic sampling. *Nucleic Acids Res.* 2011 Aug;39(15):6359-68. doi: 10.1093/nar/gkr334.
- Kim J, Ma J. 2013. PSAR-Align: improving multiple sequence alignment using probabilistic sampling. *Bioinformatics.* 2013 Nov 29. [Epub ahead of print]
- Kolodny R, Koehl P, Levitt M. 2005. Comprehensive evaluation of protein structure alignment methods: scoring by geometric measures. *J Mol Biol.* 346(4):1173-1188.

- Landan G, Graur D. 2008. Local reliability measures from sets of co-optimal multiple sequence alignments. *Pac Symp Biocomput.* 2008:15-24.
- Laszlo AH, Derrington IM, Ross BC, Brinkerhoff H, Adey A, Nova IC, Craig JM, Langford KW, Samson JM, Daza R, et al. 2014. *Nature Biotechnology* 32(8): 829-33.
- Lindblad-Toh K, Garber M, Zuk O, Lin M F, Parker B J, Washietl S, Kheradpour P, Ernst J, Jordan G, Mauceli E, Ward L D, et al. 2011. A high-resolution map of human evolutionary constraint using 29 mammals. *Nature.* 478(7370) 476-482. doi:10.1038/nature10530
- Margulies EH, Cooper GM, Asimenos G, Thomas DJ, Dewey CN, Siepel A, Birney E, Keefe D, Schwartz AS, Hou M, et al. 2007. Analyses of deep mammalian sequence alignments and constraint predictions for 1% of the human genome. *Genome Res.* 17(6):760-74.
- Meyer LR, Zweig AS, Hinrichs AS, Karolchik D, Kuhn RM, Wong M, Sloan CA, Rosenbloom KR, Roe G, Rhead B, et al. 2013. The UCSC Genome Browser database: extensions and updates 2013. *Nucleic Acids Res.* 41(Database issue):D64-9. doi: 10.1093/nar/gks1048. Epub 2012 Nov 15.
- Miller W, Rosenbloom K, Hardison RC, Hou M, Taylor J, Raney B, Burhans R, King DC, Baertsch R, Blankenberg D, et al. 2007. 28-way vertebrate alignment and conservation track in the UCSC Genome Browser. *Genome Res.* 2007 Dec;17(12):1797-808.
- The modENCODE Consortium et al. 2010. Identification of Functional Elements and Regulatory Circuits by *Drosophila* modENCODE. *Science* 330:1787-1797.
- Notredame C, and Abergel C. 2003. Using Multiple Alignment Methods to Assess the Quality of Genomic Data Analysis, in *Bioinformatics and Genomes: Current Perspectives*, M. Andrade, Editor. 2003, Horizon Scientific Press. p. 30-50
- Notredame C. 2007. Recent evolutions of multiple sequence alignment algorithms. *PLoS Comput Biol.* 2007 Aug;3(8):e123.
- Notredame C. 2012. Robusta Home Page, <http://www.tcoffee.org/Projects/robusta/> accessed 15 January 2014.
- Paten B, Herrero J, Beal K, Fitzgerald S, Birney E. 2008. Enredo and Pecan: genome-wide mammalian consistency-based multiple alignment with paralogs. *Genome Res.* 2008 Nov;18(11):1814-28. doi: 10.1101/gr.076554.108.
- Paten B, Herrero J, Beal K, Birney E. 2009. Sequence progressive alignment, a framework for practical large-scale probabilistic consistency alignment. *Bioinformatics.* 25(3):295-301. doi: 10.1093/bioinformatics/btn630.

- Paten B, Earl D, Nguyen N, Diekhans M, Zerbino D, Haussler D. 2011. Cactus: Algorithms for genome multiple sequence alignment. *Genome Res.*; 21(9):1512-28. doi: 10.1101/gr.123356.111. Epub 2011 Jun 10.
- Penn O, Privman E, Landan G, Graur D, Pupko T. 2010a. An alignment confidence score capturing robustness to guide tree uncertainty. *Mol Biol Evol.* 27(8):1759-67. doi: 10.1093/molbev/msq066.
- Penn O, Privman E, Ashkenazy H, Landan G, Graur D, Pupko T. 2010b. GUIDANCE: a web server for assessing alignment confidence scores. *Nucleic Acids Res.* 38(Web Server issue):W23-8. doi: 10.1093/nar/gkq443.
- Schreiber J, Wescoe ZL, Abu-Shumays R, Vivian JT, Baatar B, Karplus K, Akeson M. 2013. Error rates for nanopore discrimination among cytosine, methylcytosine, and hydroxymethylcytosine along individual DNA strands. *PNAS* 110(47):18910-5.
- Smit AFA, Hubley R, Green P. 2010. RepeatMasker Open-3.0. <<http://www.repeatmasker.org>>
- Smit AFA, Hubley R. 2010. RepeatModeler Open-1.0. <<http://www.repeatmasker.org>>
- Stoye J, Evers D, Meyer F. 1997. Generating benchmarks for multiple sequence alignments and phylogenetic reconstructions. *Proc Int Conf Intell Syst Mol Biol*, 5:303–6.
- Varadarajan A, Bradley RK, Holmes I. 2008. Tools for simulating evolution of aligned genomic regions with integrated parameter estimation. *Genome Biology.* 9(10):R147 doi:10.1186/gb-2008-9-10-r147.
- Wallace IM, Blackshields G, Higgins DG. 2005. Multiple sequence alignments. *Current Opinion in Structural Biology.* 15(3):261-6.
- Wilm A, Mainz I, Steger G. 2006. An enhanced RNA alignment benchmark for sequence alignment programs. *Algorithms Mol Biol.* 1:19.
- Wong KM, Suchard MA, Huelsenbeck JP. 2008. Alignment uncertainty and genomic analysis. *Science.* 319(5862):473-6. doi: 10.1126/science.1151532.