# Synteny Portal: a web-based application portal for synteny block analysis

**Jongin Lee[1], Woon-young Hong[1], Minah Cho[1], Mikang Sim[1], Daehwan Lee[1], Younhee Ko[2] and Jaebum Kim[1],***

[1]Department of Animal Biotechnology, Konkuk University, Seoul 05029, South Korea and [2]Department of Clinical Genetics, Department of Pediatrics, Yonsei University College of Medicine, Seoul 03722, South Korea

## ABSTRACT

**Recent advances in next-generation sequencing technologies and genome assembly algorithms have enabled the accumulation of a huge volume of genome sequences from various species. This has provided new opportunities for large-scale comparative genomics studies. Identifying and utilizing synteny blocks, which are genomic regions conserved among multiple species, is key to understanding genomic architecture and the evolutionary history of genomes. However, the construction and visualization of such synteny blocks from multiple species are very challenging, especially for biologists with a lack of computational skills. Here, we present Synteny Portal, a versatile web-based application portal for constructing, visualizing and browsing synteny blocks. With Synteny Portal, users can easily (i) construct synteny blocks among multiple species by using prebuilt alignments in the UCSC genome browser database, (ii) visualize and download syntenic relationships as high-quality images, (iii) browse synteny blocks with genetic information and (iv) download the details of synteny blocks to be used as input for downstream synteny-based analyses, all in an intuitive and easy-to-use web-based interface. We believe that Synteny Portal will serve as a highly valuable tool that will enable biologists to easily perform comparative genomics studies by compensating limitations of existing tools. Synteny Portal is freely available at http://bioinfo.konkuk.ac.kr/synteny_portal.**

## INTRODUCTION

Recent advances in next-generation sequencing technologies and genome assembly algorithms, together with the results of various large-scale genome projects, such as the Genome 10K Project (1), the Bird 10K Project (2) and i5K

(3), have enabled the accumulation of a huge volume of genomic sequences from various species. These genome sequences are usually stored in public databases, such as the NCBI Reference Sequence database (4), the Genome On-Line Database (5), the UCSC genome browser database (6) and Ensembl Genomes (7); they have been used to uncover genomic similarities and differences between different species and the functional consequences of such similarities and differences. Comparative genomics studies have contributed to a better understanding of the molecular-level mechanisms of species diversity and genome evolution (8–12).

Synteny blocks, which are genomic regions that are conserved among multiple species, play a pivotal role in comparative genomics. Various algorithms and tools have been introduced for the construction and utilization of synteny blocks, such as CYNTENATOR (13), DRIMM-Synteny (14), i-ADHoRe 3.0 (15), inferCars (16), OSfinder (17) and Sibelia (18). They are stand-alone and command-line tools, and users need to download and install them on a users' machine. Therefore, the construction of synteny blocks using such tools requires computer resources, and advanced bioinformatics skills to properly prepare input data and running command-line programs.

To alleviate these difficulties, more user-friendly tools for utilizing and visualizing synteny blocks have been developed. For example, GenomeMatcher (19), Mauve (20), MizBee (21) and SyMap (22) are stand-alone synteny browsers that support graphical user interface. Gbrowse_syn (23) and Sybil (24) are web-based software packages for comparative genomics that need to be installed and configured in users' own machines. GSV (25) and mGSV (26) are web-based synteny viewers that help to visualize user-provided synteny information. Cinteny (27) is a web server for identifying synteny blocks and analyzing genome rearrangements. CoGe (28), VISTA (29) and Ensembl SyntenyView (30) are web-based platforms that perform various analyses for comparative genomic studies. Genomicus (31) is a web browser that displays the homologous genomic contexts of different species. C-Sibelia

*To whom correspondence should be addressed. Tel: +82 2 450 0456; Fax: +82 2 455 1044; Email: jbkim@konkuk.ac.kr

(32) is a web server for comparing assemblies to a reference genome together with the annotation of variants in genes. GRIMM-Synteny (33) is a web server for analyzing genome rearrangements by using the list of genomic markers represented by numbers with the annotation of variants in genes. However, different tools still have different weaknesses, such as only supporting pairwise comparison, requiring the generation of synteny information by users, limited features for visualization and not supporting the saving or downloading of high-quality images (Supplementary Data for tool comparison).

Here, we present Synteny Portal, a versatile web-based application portal that allows for construction, visualization and browsing of synteny blocks by compensating limitations of existing tools. With Synteny Portal, users can easily (i) construct synteny blocks among multiple species by using prebuilt alignments in the UCSC genome browser database (6), (ii) visualize and download syntenic relationships as high-quality images, (iii) browse synteny blocks together with genetic information and (iv) download the details of synteny blocks that can be used as input for downstream synteny-based analyses, all in an intuitive and easy-to-use web interface.

## MATERIALS AND METHODS

### Data preparation

*Data collection.* We collected genomic information on well-studied reference species, including human, mouse, cow, dog, horse, pig, rat, rhesus, chicken and zebra fish, from the UCSC genome browser database (6), which includes genome assembly sequences, gene annotation, whole-genome pairwise alignments of the reference species in the chain and net format (34), cytobands and assembly sizes of the reference species. In total, 458 whole-genome pairwise alignments for 67 different species, including 124 different assembly versions and 33 assembly sequences for the ten reference species with different assembly versions are available on our web server. Gene/protein identifiers annotated with seven different criteria, such as Wiki Genes (35), RefSeq genes and proteins (4), Entrez genes (36), Ensembl genes, transcripts and proteins (7), allow for searching and browsing user-provided gene/protein identifiers. The mapping information using different identifiers was obtained from the Uniprot website (http://www.uniprot.org/) and the biomaRt package from Bioconductor (37).

*Synteny block construction and visualization.* Synteny blocks were built using the inferCars program package (16) with four different resolutions (150, 300, 400 and 500 Kbp). Specifically, pairwise collinear genome segments (PCGSs) of two species (i.e. a reference species and another species) are first collected from pairwise whole-genome alignments. Then, to create multiple collinear genome segments (MCGSs), two PCGSs are merged by splitting the PCGSs of one species at positions corresponding to the boundaries of the PCGSs of the other species based on the coordinates of the reference genome. This merging process repeats until the PCGSs of all species are merged, and the final MCGSs, which are larger than the given resolution,

are reported as synteny blocks. If outgroup species are provided, which is optional, genomic regions of the outgroup species matched to each of the calculated synteny blocks are added. Syntenic relationships are visualized through the Circos program (38) and converted to an interactive version on the website with JavaScript. The BLAST program (39) is used to search for a top-matched region in the reference genome based on the user-provided DNA or protein sequence with default parameter values (Supplementary Data).

*Implementations.* The Synteny Portal website is implemented with the HTML5 and JavaScript libraries, including jQuery (http://jquery.com), d3.js (http://d3js.org/) and GenomeD3Plot (40) for the client-side user interface, and by PHP scripts (5.3.3) for interactive data processing. The server-side processes are generated by Perl scripts (v5.22.0).

## RESULTS

### Synteny Portal overview

Synteny Portal provides four main web applications: SynCircos, SynBrowser, SynSearcher and SynBuilder. SynCircos visualizes synteny blocks against a user-selected reference and target species in an interactive Circos plot. SynBrowser displays more details of synteny blocks with annotated reference genes. SynSearcher finds syntenic regions in other species based on query sequences. Finally, SynBuilder constructs synteny blocks among multiple chosen species. These four applications generate the high-quality Circos plots or records of synteny blocks that can be downloaded in various formats, including SVG, PNG, JPEG and PDF.

### SynCircos

SynCircos depicts synteny blocks in an interactive Circos plot. Users can select a reference species, which is the reference of pairwise whole-genome alignment, multiple target species, target chromosomes and a resolution for the minimum size of a homologous reference block. The syntenic relationships among chosen species are visualized as an interactive Circos plot (Figure 1A), which can be downloaded in various formats. Ribbons connect homologous blocks in different species that belong to the same synteny block and users can highlight a specific syntenic relationship by placing a mouse pointer on a specific ribbon.

*Inputs to SynCircos.* SynCircos requires user-selected reference and target species, genome assembly version of the selected species, chromosome number, a resolution for the synteny blocks and an image format for visualization in the Circos plot.

*Outputs of SynCircos.* The outputs of SynCircos is an interactive Circos plot. Tracks and ribbons in the Circos plot represent the chromosomes and cytobands (if any) of the chosen species and the syntenic relationships among different species, respectively. Figure 1A shows a highlighted synteny block in human, mouse and cow.
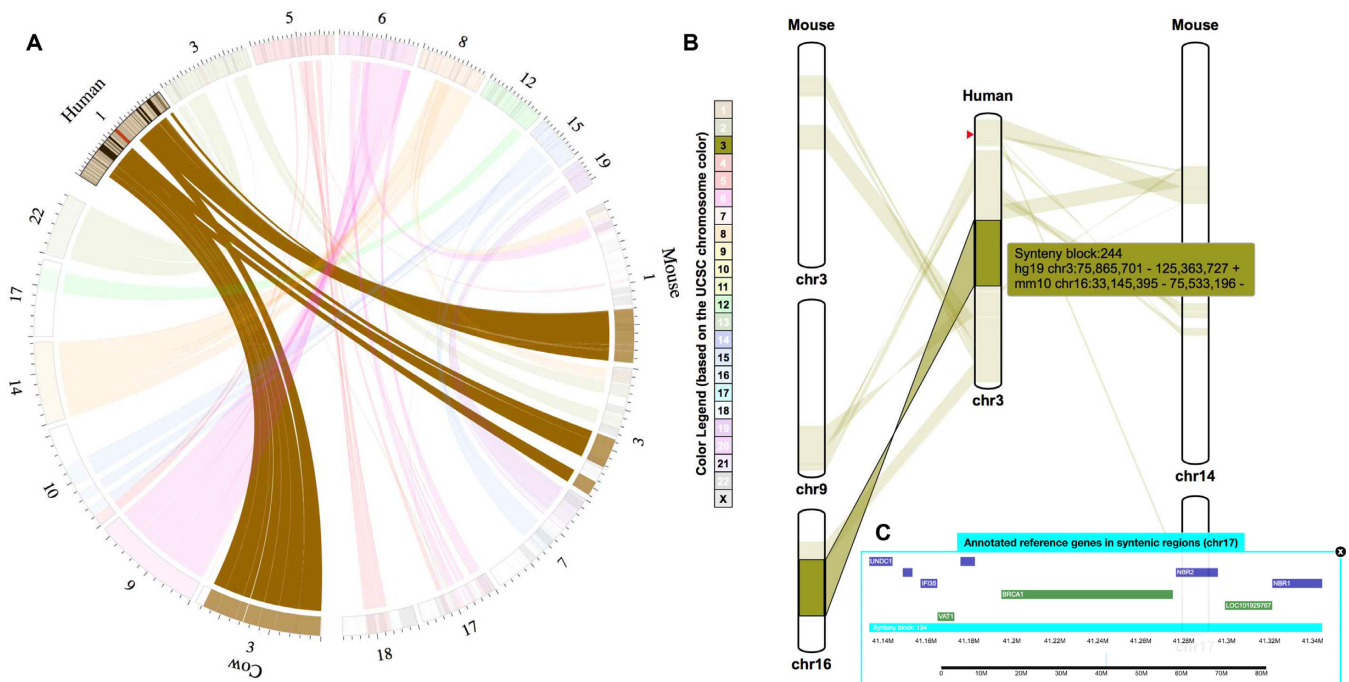
**Figure 1.** Examples of SynCircos and SynBrowser outputs. (**A**) A representative Circos plot generated by SynCircos, which compared three species and their specific chromosomes (Human: chr1, chr3, chr5, chr6, chr8, chr12, chr15, chr19; Mouse: chr1, chr3, chr7, chr17, chr18; and Cow: chr3, chr9, chr10, chr14, chr17, chr22). The highlighted ribbons represent syntenic relationships between human chromosome 1 and the chromosomes of the other species. (**B**) A representative image generated by SynBrowser which shows the relationships between human chromosome 3 and the related mouse chromosomes. The coordinates of homologous blocks in each synteny block can be seen when a specific genomic region is selected (the khaki-colored box in the center of the image). The red triangle indicates the position of the user-queried gene. (**C**) A representative gene browser produced by SynBrowser which shows the genes of the reference species (human in this image) in synteny blocks. Details of the annotated genes are provided via links to the UCSC genome browser.

## SynBrowser

SynBrowser displays more details of syntenic relationships between two species (i.e. a reference and a target species), whereas SynCircos shows a global view of syntenic relationships among multiple species. SynBrowser presents homologous blocks between the two species in an interactive graphical form (Figure 1B and C). Users can easily highlight specific homologous blocks via exact genomic coordinates in a floating window. SynBrowser also provides two ways of navigating synteny blocks: (i) chromosome-level navigation, by clicking on a specific reference chromosome in a chromosome legend on the left of a webpage, and (ii) region-level navigation, by specifying a coordinate of a reference chromosome. Genes within a synteny block can also be navigated using a gene browser through a direct search for a specific gene with a variety of gene identifiers, such as RefSeq, Wiki Gene and Entrez Gene. The images in SynBrowser can also be downloaded in multiple formats.

*Inputs to SynBrowser.* To visualize the pairwise synteny blocks, SynBrowser needs only four user-selected parameters: a reference species, a reference chromosome, a target species and a synteny block resolution. Synteny blocks containing a specific gene or protein can be found by searching for a gene or protein identifier. Users can directly move to a chosen position in the reference genome by providing specific coordinates.

*Outputs of SynBrowser.* The outputs of SynBrowser are a plot showing connections between homologous blocks between two species and a gene browser displaying reference genes within synteny blocks. Figure 1B and C show highlighted homologous blocks between human and mouse in the same synteny block, along with their coordinates and reference genes within the synteny block. The red triangle in Figure 1B indicates the location of a gene searched for by users.

## SynSearcher

SynSearcher searches for a reference genome sequence and finds the best-matched region based on a user-provided DNA or protein sequence by BLAST query; synteny blocks encompassing the top-matched reference region are retrieved (Figure 2A). The results are provided to users in two different formats: (i) the details of BLAST search results and the coordinates of synteny blocks as a text file, and (ii) an interactive Circos plot highlighting the retrieved syntenic relationships (Figure 2B). The Circos plot can be downloaded in various formats, and the text file can also be downloaded for downstream analyses with minor modifications.

*Inputs to SynSearcher.* A user-provided DNA or protein sequence can be directly inserted in a text box or uploaded as a file in a FASTA format. Users can also select a reference and target species with a chosen resolution.
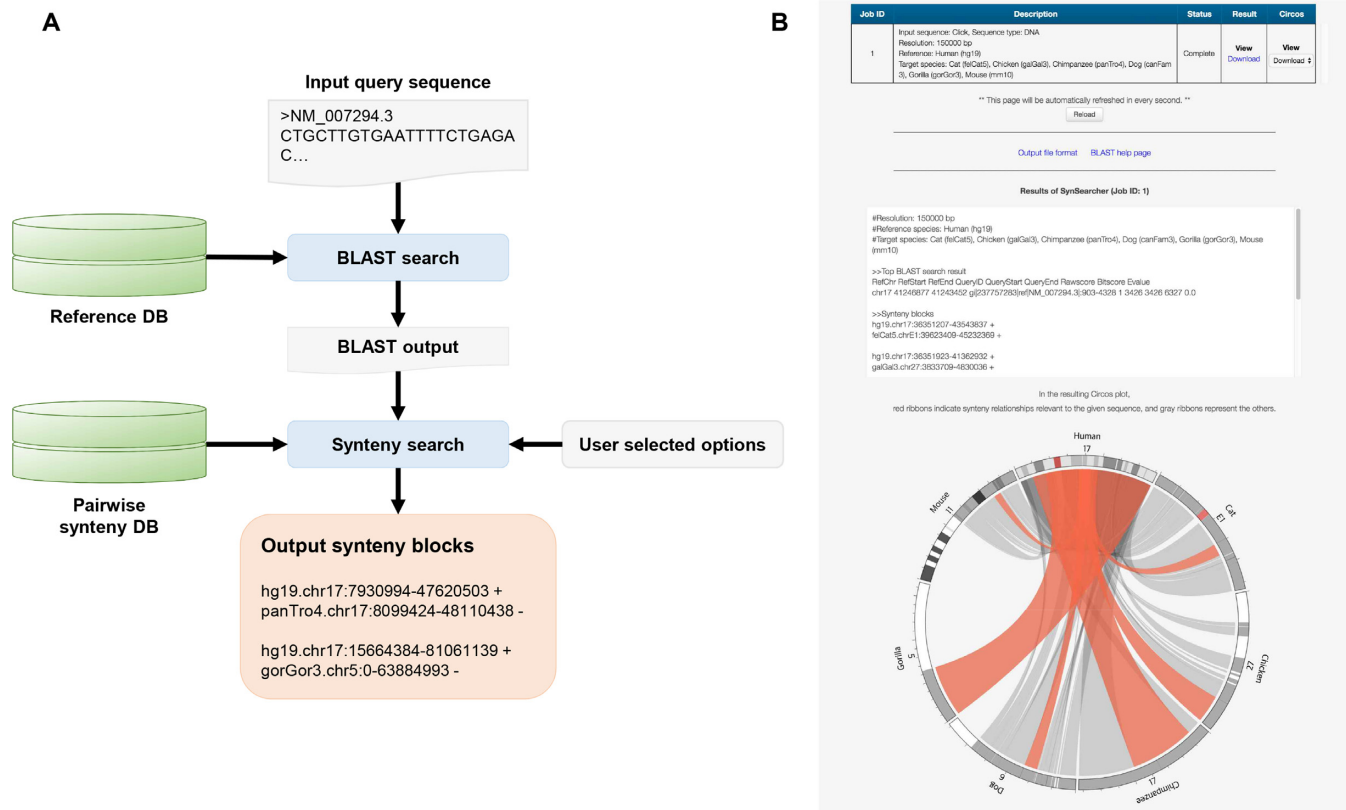
**Figure 2.** Flowchart and outputs of SynSearcher. (**A**) Given a DNA or protein sequence in the FASTA format, SynSearcher finds matching regions in a reference species using a BLAST search and returns pairwise synteny blocks containing the matched regions of the reference species. (**B**) The outputs of SynSearcher (white box, BLAST search results and the coordinates of pairwise synteny blocks; below, the Circos plot) for seven species (human, chimpanzee, gorilla, mouse, cat, dog and chicken) given a user-provided sequence.

*Outputs of SynSearcher.* The results consist of (i) the details of the BLAST output for the top-matched alignment result, (ii) the coordinates of identified pairwise synteny blocks and (iii) a Circos plot visualizing the synteny blocks. Figure 2B shows synteny blocks obtained from seven species (human, chimpanzee, mouse, gorilla, chicken, cat and dog) given a user-provided sequence.

**SynBuilder**

SynBuilder creates synteny blocks for user-selected species (a reference, target and optional outgroup species) with a chosen resolution using the inferCars program (16). Specifically, pairwise whole-genome alignments between the reference and the target species are first collected and then homologous genomic regions are identified, and, finally, co-linear genomic regions among the chosen species are grouped into synteny blocks (Figure 3A). If outgroup species are selected, matched genomic regions between the outgroup species and a reference species are added to the synteny blocks. The results of SynBuilder consist of (i) the coordinates of the synteny blocks as a text file, and (ii) an interactive Circos plot depicting syntenic relationships (Figure 3B). The Circos plot can be downloaded in various formats, and the text file can also be downloaded and used for downstream analyses with minor modifications.

*Inputs to SynBuilder.* SynBuilder requires the selection of a reference, target and optional outgroup species, and a resolution for the synteny blocks.

*Outputs of SynBuilder.* The results of SynBuilder are (i) the coordinates of synteny blocks and (ii) an interactive Circos plot depicting syntenic relationships among the chosen species.

**CONCLUSION**

Synteny Portal facilitates studies on comparative genomics through easy construction of synteny blocks, intuitive graphical representation with a high-quality image format and easy-to-use querying and browsing functionalities. Synteny Portal enables a user with a lack of computational skills to perform comparative genomic analyses. The current version only works with prebuilt whole-genome alignments because of long computational time and high computational resource requirements for online creation of whole-genome sequence alignments. However, whole-genome alignment data will be periodically updated via mirroring of data in the UCSC genome browser database. In addition, our website will be updated to support the whole-genome alignment of user-provided genome sequences with additional reference species in future. We believe that Synteny Portal will serve
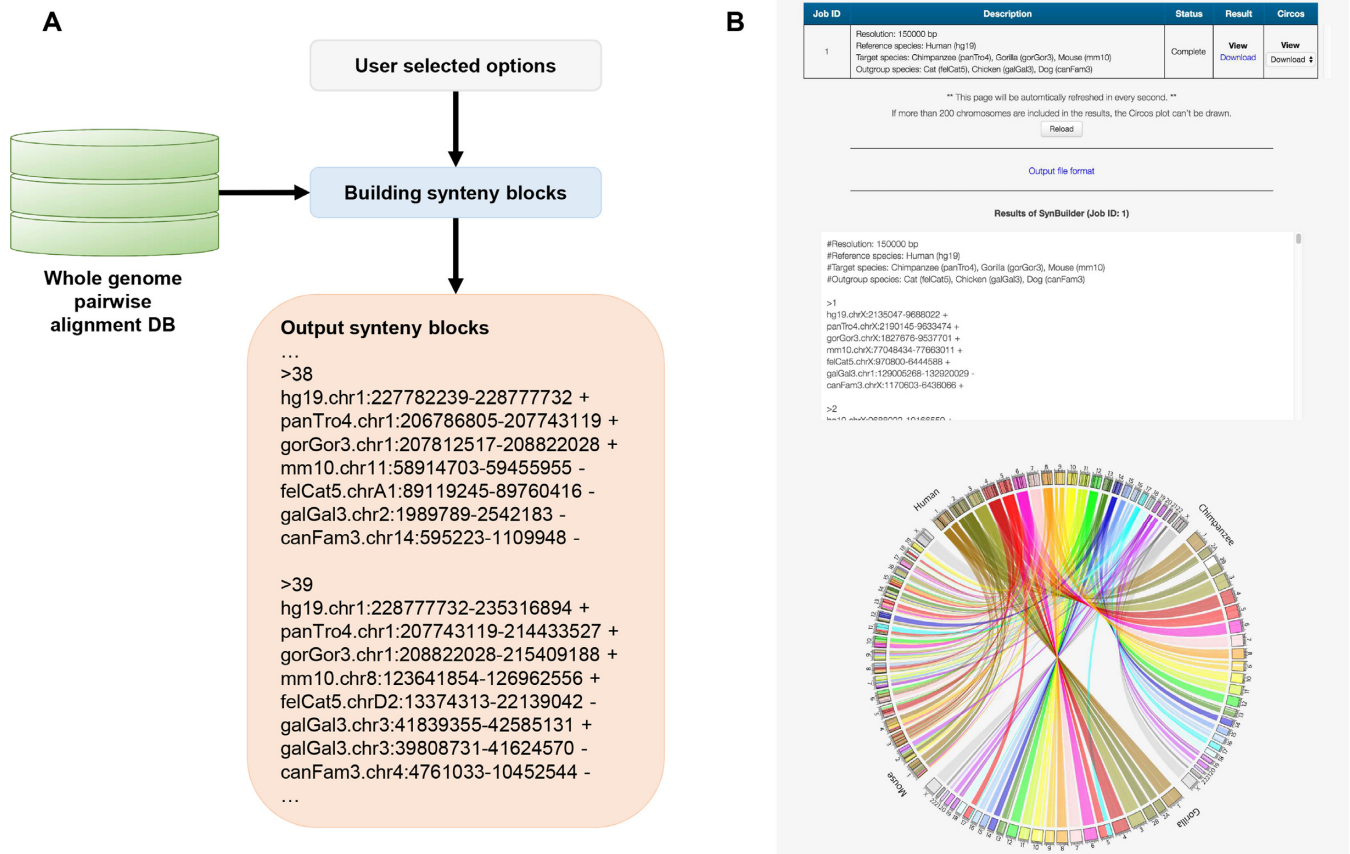
**Figure 3.** Flowchart and outputs of SynBuilder. (**A**) Given options (a reference species, target species, and optional outgroup species, and a resolution) selected by users, SynBuilder creates synteny blocks using pairwise whole-genome alignments between the reference and the other species. (**B**) The output of SynBuilder (white box, coordinates of synteny blocks; below, the Circos plot) for four species (human, chimpanzee, gorilla and mouse) with three outgroup species (cat, dog and chicken).

as a highly valuable tool that will enable biologists to easily perform comparative genomics studies.

## AVAILABILITY

The Synteny Portal web server is available at http://bioinfo. konkuk.ac.kr/synteny_portal/.

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## FUNDING

## REFERENCES

1. Koepfli,K.P., Paten,B. and O'Brien,S.J. (2015) The Genome 10K Project: a way forward. *Annu. Rev. Anim. Biosci.*, **3**, 57–111.
2. Zhang,G., Rahbek,C., Graves,G.R., Lei,F., Jarvis,E.D. and Gilbert,M.T. (2015) Genomics: bird sequencing project takes off. *Nature*, **522**, 34.
3. i5K Consortium. (2013) The i5K Initiative: advancing arthropod genomics for knowledge, human health, agriculture, and the environment. *J. Hered.*, **104**, 595–600.
4. Pruitt,K.D., Tatusova,T. and Maglott,D.R. (2007) NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.*, **35**, D61–D65.
5. Reddy,T.B., Thomas,A.D., Stamatis,D., Bertsch,J., Isbandi,M., Jansson,J., Mallajosyula,J., Pagani,I., Lobos,E.A. and Kyrpides,N.C. (2015) The Genomes OnLine Database (GOLD) v.5: a metadata management system based on a four level (meta)genome project classification. *Nucleic Acids Res.*, **43**, D1099–D1106.
6. Speir,M.L., Zweig,A.S., Rosenbloom,K.R., Raney,B.J., Paten,B., Nejad,P., Lee,B.T., Learned,K., Karolchik,D., Hinrichs,A.S. *et al.* (2016) The UCSC Genome Browser database: 2016 update. *Nucleic Acids Res.*, **44**, D717–D725.
7. Kersey,P.J., Allen,J.E., Armean,I., Boddu,S., Bolt,B.J., Carvalho-Silva,D., Christensen,M., Davis,P., Falin,L.J., Grabmueller,C. *et al.* (2016) Ensembl Genomes 2016: more genomes, more complexity. *Nucleic Acids Res.*, **44**, D574–D580.
8. Guiliano,D.B., Hall,N., Jones,S.J., Clark,L.N., Corton,C.H., Barrell,B.G. and Blaxter,M.L. (2002) Conservation of long-range synteny and microsynteny between the genomes of two distantly related nematodes. *Genome Biol.*, **3**, RESEARCH0057.
9. Bowers,J.E., Chapman,B.A., Rong,J. and Paterson,A.H. (2003) Unravelling angiosperm genome evolution by phylogenetic analysis of chromosomal duplication events. *Nature*, **422**, 433–438.

10. International Chicken Genome Sequencing Consortium. (2004) Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. *Nature*, **432**, 695–716.

11. Chimpanzee Sequencing and Analysis Consortium. (2005) Initial sequence of the chimpanzee genome and comparison with the human genome. *Nature*, **437**, 69–87.

12. Alfoldi,J. and Lindblad-Toh,K. (2013) Comparative genomics as a tool to understand evolution and disease. *Genome Res.*, **23**, 1063–1068.

13. Rodelsperger,C. and Dieterich,C. (2010) CYNTENATOR: progressive gene order alignment of 17 vertebrate genomes. *PLoS One*, **5**, e8861.

14. Pham,S.K. and Pevzner,P.A. (2010) DRIMM-Synteny: decomposing genomes into evolutionary conserved segments. *Bioinformatics*, **26**, 2509–2516.

15. Proost,S., Fostier,J., De Witte,D., Dhoedt,B., Demeester,P., Van de Peer,Y. and Vandepoele,K. (2012) i-ADHoRe 3.0–fast and sensitive detection of genomic homology in extremely large data sets. *Nucleic Acids Res.*, **40**, e11.

16. Ma,J., Zhang,L., Suh,B.B., Raney,B.J., Burhans,R.C., Kent,W.J., Blanchette,M., Haussler,D. and Miller,W. (2006) Reconstructing contiguous regions of an ancestral genome. *Genome Res.*, **16**, 1557–1565.

17. Hachiya,T., Osana,Y., Popendorf,K. and Sakakibara,Y. (2009) Accurate identification of orthologous segments among multiple genomes. *Bioinformatics*, **25**, 853–860.

18. Minkin,I., Patel,A., Kolmogorov,M., Vyahhi,N. and Pham,S. (2013) Sibelia: a scalable and comprehensive synteny block generation tool for closely related microbial genomes. In: Darling,A and Stoye,J (eds). *Algorithms in Bioinformatics*. Springer, Berlin Heidelberg, pp. 215–229.

19. Ohtsubo,Y., Ikeda-Ohtsubo,W., Nagata,Y. and Tsuda,M. (2008) GenomeMatcher: a graphical user interface for DNA sequence comparison. *BMC Bioinformatics*, **9**, 376.

20. Darling,A.C., Mau,B., Blattner,F.R. and Perna,N.T. (2004) Mauve: multiple alignment of conserved genomic sequence with rearrangements. *Genome Res.*, **14**, 1394–1403.

21. Meyer,M., Munzner,T. and Pfister,H. (2009) MizBee: a multiscale synteny browser. *IEEE Trans. Vis. Comput. Graph.*, **15**, 897–904.

22. Soderlund,C., Bomhoff,M. and Nelson,W.M. (2011) SyMAP v3.4: a turnkey synteny system with application to plant genomes. *Nucleic Acids Res.*, **39**, e68.

23. McKay,S.J., Vergara,I.A. and Stajich,J.E. (2010) Using the Generic Synteny Browser (GBrowse_syn). *Curr. Protoc. Bioinformatics*, **31**, 9.12.1–9.12.25.

24. Crabtree,J., Angiuoli,S.V., Wortman,J.R. and White,O.R. (2007) Sybil: methods and software for multiple genome comparison and visualization. *Methods Mol. Biol.*, **408**, 93–108.

25. Revanna,K.V., Chiu,C.C., Bierschank,E. and Dong,Q. (2011) GSV: a web-based genome synteny viewer for customized data. *BMC Bioinformatics*, **12**, 316.

26. Revanna,K.V., Munro,D., Gao,A., Chiu,C.C., Pathak,A. and Dong,Q. (2012) A web-based multi-genome synteny viewer for customized data. *BMC Bioinformatics*, **13**, 190.

27. Sinha,A.U. and Meller,J. (2007) Cinteny: flexible analysis and visualization of synteny and genome rearrangements in multiple organisms. *BMC Bioinformatics*, **8**, 82.

28. Lyons,E. and Freeling,M. (2008) How to usefully compare homologous plant genes and chromosomes as DNA sequences. *Plant J.*, **53**, 661–673.

29. Frazer,K.A., Pachter,L., Poliakov,A., Rubin,E.M. and Dubchak,I. (2004) VISTA: computational tools for comparative genomics. *Nucleic Acids Res.*, **32**, W273–W279.

30. Clamp,M., Andrews,D., Barker,D., Bevan,P., Cameron,G., Chen,Y., Clark,L., Cox,T., Cuff,J., Curwen,V. *et al.* (2003) Ensembl 2002: accommodating comparative genomics. *Nucleic Acids Res.*, **31**, 38–42.

31. Louis,A., Muffato,M. and Roest Crollius,H. (2013) Genomicus: five genome browsers for comparative genomics in eukaryota. *Nucleic Acids Res.*, **41**, D700–D705.

32. Minkin,I., Pham,H., Starostina,E., Vyahhi,N. and Pham,S. (2013) C-Sibelia: an easy-to-use and highly accurate tool for bacterial genome comparison. *F1000Res.*, **2**, 258.

33. Pevzner,P. and Tesler,G. (2003) Human and mouse genomic sequences reveal extensive breakpoint reuse in mammalian evolution. *Proc. Natl. Acad. Sci. U.S.A.*, **100**, 7672–7677.

34. Kent,W.J., Baertsch,R., Hinrichs,A., Miller,W. and Haussler,D. (2003) Evolution's cauldron: duplication, deletion, and rearrangement in the mouse and human genomes. *Proc. Natl. Acad. Sci. U.S.A.*, **100**, 11484–11489.

35. Hoffmann,R. (2008) A wiki for the life sciences where authorship matters. *Nat. Genet.*, **40**, 1047–1051.

36. Brown,G.R., Hem,V., Katz,K.S., Ovetsky,M., Wallin,C., Ermolaeva,O., Tolstoy,I., Tatusova,T., Pruitt,K.D., Maglott,D.R. *et al.* (2015) Gene: a gene-centered information resource at NCBI. *Nucleic Acids Res.*, **43**, D36–D42.

37. Durinck,S., Spellman,P.T., Birney,E. and Huber,W. (2009) Mapping identifiers for the integration of genomic datasets with the R/Bioconductor package biomaRt. *Nat. Protoc.*, **4**, 1184–1191.

38. Krzywinski,M., Schein,J., Birol,I., Connors,J., Gascoyne,R., Horsman,D., Jones,S.J. and Marra,M.A. (2009) Circos: an information aesthetic for comparative genomics. *Genome Res.*, **19**, 1639–1645.

39. Camacho,C., Coulouris,G., Avagyan,V., Ma,N., Papadopoulos,J., Bealer,K. and Madden,T.L. (2009) BLAST+: architecture and applications. *BMC Bioinformatics*, **10**, 421.

40. Laird,M.R., Langille,M.G. and Brinkman,F.S. (2015) GenomeD3Plot: a library for rich, interactive visualizations of genomic data in web applications. *Bioinformatics*, **31**, 3348–3349.