

Computational Pan-Genomics: Status, Promises and Challenges

Tobias Marschall^{1,2}, Manja Marz^{3,4,5,6}, Thomas Abeel⁷, Louis Dijkstra^{8,9}, Bas E. Dutilh^{10,11,12}, Ali Ghaffaari^{1,2}, Paul Kersey¹³, Wigard P. Kloosterman¹⁴, Veli Mäkinen¹⁵, Adam M. Novak¹⁶, Benedict Paten¹⁶, David Porubsky¹⁷, Eric Rivals^{18,19}, Can Alkan²⁰, Jasmijn Baaijens²¹, Paul I. W. De Bakker¹⁴, Valentina Boeva^{22,23,24,25}, Francesca Chiaromonte²⁶, Rayan Chikhi²⁷, Francesca D. Ciccarelli²⁸, Robin Cijvat²⁹, Erwin Datema³⁰, Cornelia M. Van Duijn³¹, Evan E. Eichler^{32,33}, Corinna Ernst³⁴, Eleazar Eskin^{35,36}, Erik Garrison³⁷, Mohammed El-Kebir^{21,38,39}, Gunnar W. Klau²¹, Jan O. Korbel^{13,40}, Eric-Wubbo Lameijer⁴¹, Benjamin Langmead⁴², Marcel Martin⁴³, Paul Medvedev^{44,45,46}, John C. Mu⁴⁷, Pieter Neerincx⁴¹, Klaasjan Ouwers^{48,49}, Pierre Peterlongo⁵⁰, Nadia Pisanti^{51,52}, Sven Rahmann³⁴, Ben Raphael³⁹, Knut Reinert⁵³, Dick de Ridder⁵⁴, Jeroen de Ridder⁷, Matthias Schlesner⁵⁵, Ole Schulz-Trieglaff⁵⁶, Ashley Sanders⁵⁷, Siavash Sheikhezadeh⁵⁴, Carl Schneider⁵⁸, Sandra Smit⁵⁴, Daniel Valenzuela¹⁵, Jiayin Wang^{59,60,61}, Lodewyk Wessels⁶², Ying Zhang^{29,21}, Victor Guryev^{17,14}, Fabio Vandin^{63,39}, Kai Ye^{64,65,61} and Alexander Schönhuth²¹

¹Center for Bioinformatics, Saarland University, Saarbrücken, Germany; ²Max Planck Institute for Informatics, Saarbrücken, Germany; ³Bioinformatics and High Throughput Analysis, Faculty of Mathematics and Computer Science, Friedrich Schiller University Jena, Leutragraben 1, 07743 Jena, Germany; ⁴FLI Leibniz Institute for Age Research, Beutenbergstraße 11, 07745 Jena, Germany; ⁵Michael Stifel Center Jena, Ernst-Abbe-Platz 2, 07743 Jena, Germany; ⁶German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig; ⁷The Delft Bioinformatics Lab, Department of Intelligent Systems, Delft University of Technology, Mekelweg 4, 2628 CD, Delft, The Netherlands; ⁸Computational Science Lab, University of Amsterdam, Amsterdam, 1098XG, The Netherlands; ⁹Department of High Performance Computing, ITMO University, Saint Petersburg, 197101, Russia; ¹⁰Radboud Institute for Molecular Life Sciences, Center for Molecular and Biomolecular Informatics, Radboud University Medical Center, Nijmegen, Netherlands; ¹¹Theoretical Biology and Bioinformatics, Utrecht University, Utrecht, Netherlands; ¹²Department of Marine Biology, Institute of Biology, Federal University of Rio de Janeiro, Rio de Janeiro, Brazil; ¹³EMBL-European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, CB10 1SD, UK; ¹⁴Department of Genetics, Center for Molecular Medicine, University Medical Center Utrecht, 3584 CG, Utrecht, The Netherlands; ¹⁵HIIT and Department of Computer Science, University of Helsinki, Finland; ¹⁶Genomics Institute, University of California Santa Cruz, Santa Cruz, CA 95064, USA; ¹⁷European Research Institute for the Biology of Ageing, University Medical Center Groningen, University of Groningen, Antonius Deusinglaan 1, AV Groningen 9713, The Netherlands; ¹⁸LIRMM, CNRS and Université de Montpellier, Montpellier, France; ¹⁹Institut de Biologie Computationnelle, CNRS and Université de Montpellier, Montpellier, France; ²⁰Department of Computer Engineering, Bilkent University, Bilkent, Ankara, 06800, Turkey; ²¹Life Sciences Group, Centrum Wiskunde & Informatica, Amsterdam, 1098XG, The Netherlands; ²²Institut Curie, Centre de Recherche, Inserm U900, F-75005 Paris, France; ²³Mines ParisTech, F-77305 cedex Fontainebleau, France; ²⁴PSL Research University, F-75005 Paris, France; ²⁵Institut Cochin, Inserm U1016, CNRS UMR 8104, Université Paris Descartes UMR-S1016, F-75014 Paris, France; ²⁶Departments of Statistics, The Pennsylvania State University, University Park, PA, 16802; ²⁷CNRS, Univ. Lille, UMR 9189 CRISTAL, F-59000 Lille, France; ²⁸Division of Cancer Studies, King's College London, London SE11UL, UK; ²⁹MonetDB Solutions, Amsterdam, The Netherlands; ³⁰KeyGene N.V., Agro Business Park 90, 6708 PW Wageningen, The Netherlands; ³¹Department of Epidemiology, Erasmus Medical Center, Rotterdam, The Netherlands; ³²Department of Genome Sciences, University of Washington, Seattle, WA, 98195, USA; ³³Howard Hughes Medical Institute, University of Washington, Seattle, WA, 98195, USA; ³⁴Genome Informatics, Institute of Human Genetics, University Hospital Essen, University of Duisburg-Essen, Essen, Germany; ³⁵Department of Computer Science, University of California, Los Angeles, USA; ³⁶Department of Human Genetics, University of California, Los Angeles, USA; ³⁷Wellcome Trust Sanger Institute, Cambridge, UK; ³⁸Centre for Integrative Bioinformatics VU (IBIVU), VU University Amsterdam, De Boelelaan 1081A, 1081 HV Amsterdam, The Netherlands; ³⁹Center for Computational Molecular Biology and Department of Computer Science, Brown University, Providence, RI 02912, USA; ⁴⁰European Molecular Biology Laboratory (EMBL), Genome Biology Unit, Meyerhofstrasse 1, 69117 Heidelberg, Germany; ⁴¹Genomics Coordination Center, University of Groningen, University Medical Center Groningen, Groningen, 9700RB, The Netherlands; ⁴²Department of Computer Science and Center for Computational Biology, Johns Hopkins University, Baltimore, Maryland; ⁴³Science for Life Laboratory, Dept. of Biochemistry and Biophysics, Stockholm University, Box 1031, SE-17121 Solna, Sweden; ⁴⁴Department of Computer Science and Engineering, The Pennsylvania State University, USA; ⁴⁵Department of Biochemistry and Molecular Biology, The Pennsylvania State University, USA; ⁴⁶Genomic Sciences Institute of the Huck, The Pennsylvania State University, USA; ⁴⁷Bina Technologies, Roche Sequencing, Redwood City, CA 94065, USA; ⁴⁸Biological Psychology, Vrije Universiteit Amsterdam, The Netherlands; ⁴⁹Genalix BV, Harderwijk, The Netherlands; ⁵⁰INRIA Campus de Beaulieu- Rennes, Rennes Cedex 35042, France; ⁵¹Dipartimento di Informatica, University degli Studi di Pisa, Pisa, Italy; ⁵²Erable Team, INRIA; ⁵³Department of Mathematics and Computer Science, Freie Universität Berlin, Berlin, Germany; ⁵⁴Bioinformatics Group, Wageningen University, Droevendaalsesteeg 1, 6708 PB, Wageningen, The Netherlands; ⁵⁵Division of Theoretical Bioinformatics, German Cancer Research Center (DKFZ), Im Neuenheimer Feld 280, 69120 Heidelberg, Germany; ⁵⁶Illumina Cambridge Ltd, Chesterford Research Park, Little Chesterford, Essex CB10 1XL, UK; ⁵⁷Terry Fox Laboratory, BC Cancer Agency, Vancouver, British Columbia, Canada; ⁵⁸Leiden Observatory, Leiden University, P.O. Box 9513, 2300 RA Leiden, The Netherlands; ⁵⁹School of Management, Xi'an Jiaotong University, Xi'an, Shaanxi 710049, China; ⁶⁰Institute of Data Science and Information Quality, Xi'an Jiaotong University; ⁶¹Shaanxi Engineering Research Center of Medical and Health Big Data, Xi'an Jiaotong University; ⁶²Netherlands Cancer Institute (NKI), Amsterdam, the Netherlands; ⁶³Department of Mathematics and Computer Science, University of Southern Denmark, Campusvej 55, DK-5230, Odense M, Denmark; ⁶⁴School of Electronic and Information Engineering, Xi'an Jiaotong University, Xi'an, Shaanxi 710049, China; ⁶⁵Kai Ye Young Scientist Studio, Xi'an Jiaotong University

March 29, 2016

Abstract

Many disciplines, from human genetics and oncology to plant and animal breeding, microbiology and virology, commonly face the challenge of analyzing rapidly increasing numbers of genomes. In case of *Homo sapiens*, the number of sequenced genomes will approach hundreds of thousands in the next few years. Simply scaling up established bioinformatics pipelines will not be sufficient for leveraging the full potential of such rich genomic datasets. Instead, novel, qualitatively different computational methods and paradigms are needed. We will witness the rapid extension of *computational pan-genomics*, a new sub-area of research in computational biology. In this paper, we examine already available approaches to construct and use pan-genomes, discuss the potential benefits of future technologies and methodologies, and review open challenges from the vantage point of the above-mentioned biological disciplines. As a prominent example for a computational paradigm shift, we particularly highlight the transition from the representation of reference genomes as strings to representations as graphs. We outline how this and other challenges from different application domains translate into common computational problems, point out relevant bioinformatics techniques and identify open problems in computer science. In this way, we aim to form a computational pan-genomics community that bridges several biological and computational disciplines.

1 Introduction

In 1995, the complete genome sequence for the bacterium *Haemophilus influenzae* was published, followed by the sequence for the eukaryote *Saccharomyces cerevisiae* in 1996 and the landmark publication of the human genome in 2001. These sequences, and many more that followed, have served as *reference genomes*, which formed the basis for both major advances in functional genomics and for studying genetic variation by re-sequencing other individuals from the same species. The advent of rapid and cheap “next-generation” sequencing technologies since 2006 has turned re-sequencing into one of the most popular modern genome analysis workflows. As of today, an incredible wealth of genomic variation within populations has already been detected, permitting functional annotation of many such variants, and it is reasonable to expect that this is only the beginning.

With the number of sequenced genomes steadily increasing, it makes sense to re-think the idea of a *reference genome*. Such a reference sequence can take a number of forms, including:

- the genome of a single selected individual,
- a consensus drawn from an entire population,
- a “functional” genome (without disabling mutations in any genes), or
- a maximal genome that captures all sequence ever detected.

Depending on the context, each of these alternatives may make sense. However, many early reference sequences did not represent any of the above. Instead, they consisted of collections of sequence patches, assayed from whatever experimental material had been available, often from a relatively unstructured mix of individual biological sources. Only lately has the rapid spread of advanced sequencing technologies allowed the reasonably complete determination of many individual genome sequences from particular populations, taxonomic units, or environments. To take full advantage of these data, a good “reference genome” should have capabilities beyond the alternatives listed above. This entails a paradigm shift, from focusing on a single reference genome to using a *pan-genome*, that is, a representation of all genomic content in a certain species.

1.1 Definition of Computational Pan-Genomics

The term *pan-genome* was first used by Sigaux [121] to describe a public database containing an assessment of genome and transcriptome alterations in major types of tumors, tissues, and experimental models. Later, Tettelin et al. [131] defined a microbial pan-genome as the combination of a *core genome*, containing genes present in all strains, and a *dispensable genome* (also known as flexible or accessory genome) composed of genes absent from one or more of the strains. A generalization of such a representation could contain not only the genes, but also other variations present in the collection of genomes. The idea of transitioning to a human pan-genome is also gaining more and more attention¹.

Here, we generalize the above definitions and use the word *pan-genome* to refer to any collection of genomic sequences to be analyzed jointly or to be used as a reference. These sequences can be linked in a graph-like structure, or simply constitute sets of (aligned or unaligned) sequences. Questions about efficient data structures, algorithms and statistical methods to perform bioinformatic analyses of pan-genomes give rise to the discipline of *computational pan-genomics*.

While being aware that the above definition of a pan-genome is general, we argue that it is instrumental for identifying common computational problems that occur in different disciplines. Our notion of computational pan-genomics therefore intentionally intersects with many other bioinformatics disciplines. In particular it is related to *metagenomics*, which studies the entirety of genetic material sampled from an environment; to *comparative genomics*, which is concerned with retracing evolution by analyzing genome sequences; and to *population genetics*, whose main subject is the change of a population’s genetic composition in response to various evolutionary forces and migration. All these fields have developed their own algorithms and data structures to represent sets of genomes and can therefore contribute to the pan-genomics toolbox. By advocating *computational pan-genomics*, we hope to increase awareness of

¹See <http://www.technologyreview.com/news/537916/rebooting-the-human-genome>, for an example of recent media coverage

common challenges and to generate synergy among the involved fields.

At the core of pan-genomics is the idea of replacing traditional, linear reference genomes by richer data structures. The paradigm of a single reference genome has endured in part because of its simplicity. It has provided an easy framework within which to organize and think about genomic data; for example, it can be readily visualized through a genome browser. With the currently rapidly growing number of sequences we have at our disposal, this approach increasingly fails to fully capture the information on variation, similarity, frequency, and functional content implicit in the data. Although pan-genomes promise to be able to represent this information, there is not yet a conceptual framework or a toolset for working with pan-genomes that has achieved widespread acceptance. For many biological questions, it is not yet established how to best extract the relevant information from any particular pan-genome representation, and even when the right approach can be identified, novel bioinformatics tools often need to be developed in order to apply it.

In this paper, we explore the challenges of working with pan-genomes, and identify conceptual and technical approaches that may allow us to organize such data to facilitate its application in (green, blue, red, and white [65]) biotechnology and fundamental research.

1.2 Goals of Computational Pan-Genomics

On a high level, desirable features of a pan-genome include *completeness*, or containing all functional elements and enough of the sequence space to serve as a reference for the analysis of additional individuals; *stability*, or having uniquely identifiable features that can be studied by different researchers and at different points in time; *comprehensibility*, or facilitating understanding of the complexities of genome structures across many individuals or species; and *efficiency*, or organizing data in such a way as to accelerate downstream analysis.

These desiderata highlight the breadth of challenges facing pan-genomics as a field, some of which go beyond scientific questions. Reaching *completeness*, for instance, requires the necessary (financial and technical) resources to collect and sequence a

sufficient number of genomes for a particular tissue, organism, species, other taxonomic unit, ecological community, or geospatial niche of interest to be accurately represented. The availability data sharing mechanisms will greatly influence how quickly *completeness* can be achieved. Issues of data sharing include technical ones (mostly due to the data being big), political ones, and ethical/privacy concerns [53], as well as issues related to the interplay of these three areas. Achieving *stability* requires a central, recognized authority equipped with the long-term resources for curating reference pan-genomes. Besides this organizational component, achieving stability also requires reaching consensus about ways to define coordinate systems on pan-genomes. The goal of *comprehensibility* is mostly a biological problem. What it means exactly can differ substantially between application domains, as we outline below. The goal of *efficiency*, on the other hand, is in the domain of computer science. Aligning the needs of researchers in the application domains with efforts to develop algorithms and statistical methods is key to designing efficient solutions. With this paper, we hope to contribute significantly to this communication process.

2 Applications

2.1 Microbes

Bacteria and fungi are widely studied—and applied—in fields including biology, medicine, and biotechnology. A full understanding of the functional and evolutionary repertoire of microbial genomes thus not only is interesting from a scientific point of view, but also opens up possibilities for developing therapies and engineering applications.

For a number of microorganisms, pan-genome sequence data is already available; refer to [79, 41] for examples. Microbes provide a unique opportunity for pan-genome construction: the size of their genomes is relatively small, and for many species there are multiple fully closed genome sequences available. Furthermore, for some clinically interesting bacterial species, up to thousands of sequenced strains are available at sufficient depth to create draft genome assemblies. This has enabled pan-genome studies at the gene level [135], for which es-

tablished workflows and mature software are available, as reviewed in [144]. With the current data, however, we are in a position to create a pan-genome at the sequence level, as in e.g. [99]. In this context, a pan-genome is a representation that encodes the complete sequence information of many individual strains.

From an evolutionary point of view, microbial pan-genomes support comparative genomics studies. These are particularly interesting due to most microorganisms' potential for horizontal gene exchange. This means that not all genes in a genome adhere to the same phylogenetic sub-tree [36]. Thus, the evolution of microorganisms including bacteria, but also higher organisms [26], is more naturally represented as a phylogenetic network, rather than a phylogenetic tree [62]. We envision that these phylogenetic networks can be encoded in the structure of the pan-genome.

Applying genome-wide association studies (GWAS) to microbes is an emerging field [39, 112], promising to pinpoint genetic variables that correlate with relevant traits such as drug resistance or secondary metabolism. Such studies can operate at the level of individual variants—such as single-nucleotide polymorphisms (SNPs), insertion/deletion variants (indels), and structural variants (SVs)—or at the level of absence or presence of whole genes, annotated functions, or mobile genetic elements such as integrons or prophages. Computational pan-genomic approaches could be applied at each of these levels. Important challenges amenable to a pan-genomic approach include establishing reliable data processing pipelines to deliver variant calls, extracting gene absence/presence signals from NGS data, annotation for hypothetical genes and proteins, and specifically computational challenges such as the definition of a coordinate system to identify sequence loci on pan-genomes or the handling nested variation, such as SNP positions in large insertions. By addressing these challenges, computational pan-genomics has the potential to substantially contribute to the success of microbial GWAS.

2.2 Metagenomics

Metagenomics studies the genomic composition of microorganisms sampled from an environment. Abundant metagenomic data is currently being

generated from various environments such as human hosts [133, 76], the world's oceans [142, 18], and soil [59]. One main advantage of this approach lies in allowing the sampling of *all* microorganisms in an environment, not only those that can be cultured. This however comes at the cost of having to untangle the sequencing data generated from such a mixture computationally. A first question often asked is about the taxonomic composition of the sample. Other relevant questions that can be approached with metagenomic data include ascertaining the presence of certain gene products or whole pathways, and determining which genomes these functional genes are associated with.

Metagenomics can be applied to gain insights on human health and disease. Metagenome-wide association studies that aim to associate the microbial composition in the human gut with diseases such as type 2 diabetes are an example [110]. Metagenomics has also been shown to be capable of revealing the genomes of entire species, and tracing them through environments, as in the example of the shiga-toxigenic *Escherichia coli* being responsible for a recent major outbreak in Germany [81].

In the metagenomic setting, the set of genomic sequences underlying a pan-genome is not defined by ancestral relationships, but by co-occurrence in an environment. This presents both a challenge and an opportunity. On the one hand, constructing such a pan-genome and drawing robust conclusions from it is difficult, especially when sequencing reads are short. On the other hand, it presents the chance to reveal common adaptations to the environment as well as co-evolution of interactions.

2.3 Viruses

Viruses are notorious mutation machines. A viral quasi-species is a cloud of viral haplotypes that surround a given master virus [35]. Although viral genomes are comparatively short (RNA viruses range from 3–30kb, DNA viruses are usually not larger than 3Mb), their high sequence variability makes it challenging to assemble full viral genomes *de novo*. There are two major sequencing approaches for viruses: sequencing isolated viral clones, and metagenomic sequencing. The latter usually identifies a metapopulation consensus genome sequence rather than a single haplotype [40], and includes confounding genetic se-

quences such as the genome of other community members and of the cellular virus host. Thus far, the obvious approach of viral particle sorting by Fluorescence-Activated Cell Sorting (FACS), followed by single virus sequencing, has remained elusive due to their small genome size [5, 88]. New long-read technologies (e.g. PacBio, Oxford Nanopore) are now providing the first promising results in the sequencing of complete viral genomes [80, 136]. Currently, error rates in these third-generation long read sequencing technologies still far exceed the frequencies of rare strains or haplotypes. However, as sequencing chemistry and technologies progress, such techniques are likely to become key tools for the construction of viral pan-genomes.

Low frequency strains are hardly detectable, especially for fast evolving RNA viruses with a replication mutation rate of about the sequencing error rates. Reliable viral haplotype reconstruction is not fully solved, although to date many promising approaches have been presented [13]. Haplotype resolution techniques such as Strand-Seq [47] are not applicable for small virus particles.

One of the goals of pan-genomics, both in virology and in medical microbiology, will be to fight infectious disease. We expect that computational pan-genomics will assist GWAS approaches, which may allow the prediction of crucial parameters such as the exact diagnosis, staging, and suitable therapy selection from a given patient's viral pan-genome. For example, several studies have shown relationships between genetic diversity and disease progression, pathogenesis, immune escape, effective vaccine design, and drug resistance in HIV [74, 11, 20]. Thus, computational pan-genomics promises to be useful when studying the response of the quasi-species to the host immune system, in the context of personalized medicine.

The molecular interactions between pathogens and their hosts lead to a genetic arms race that allows virus-host interactions to be predicted [30]. In this context, metagenomics techniques can also be applied [43, 96]. The large metagenomic datasets mentioned in Section 2.2 can serve as input for such studies. We expect that computational pan-genomics will allow increased power and accuracy, for example by allowing the pan-genome structure of a viral population to be directly compared with that of a susceptible host population.

2.4 Plants

Genomic hybridization of accessions of crops or flowers has been exploited for over a century to create offspring with desirable traits. Genes found in wild varieties that improve important properties of crops, such as appearance, nutrient content, resistance to certain pests or diseases, or tolerance for stresses such as drought or heat, are now routinely bred into commercial crop varieties.

Large-scale genomics projects to characterize the genetic diversity in plants are already ongoing, not only for the model plant *Arabidopsis thaliana* [138], but also for crops [10]. Examples include the resequencing of hundreds to thousands of varieties of rice [61], maize [64], sorghum [83], and tomato [132]. Future projects aim to sequence many more varieties, e.g. 100,000 varieties of rice². Mining and leveraging the sequence data in such large-scale projects requires a pan-genomic approach. Particularly challenging is the fact that many plant genomes are large, complex (containing many repeats) and often polyploid.

A pan-genome structure has multiple advantages over a single, linear reference genome sequence in plant breeding applications. Having a pan-genome available for a given crop that includes its wild relatives provides a single coordinate system to anchor all known variation and phenotype information, and will allow for identification of novel genes from the available germplasm that are not present in the reference genome(s). Moreover, the pan-genome will reveal chromosomal rearrangements between genotypes which can hinder the introgression of desired genes. It also provides a compact representation of polyploid genomes and, in case of autopolyploids, allows for the quantitation of allele dosage between individuals.

2.5 Rare Genetic Diseases

Mutations that are causal for rare Mendelian diseases have successfully been discovered during the past few years using whole exome and genome sequencing [50, 9]. Key for these studies is the availability of databases of common and rare genetic variants present in control populations that do not carry the disease. Current resources such as

²<http://irri.org/our-work/research>

the Exome Variant Server (EVS)³ and the ExAC Browser [46] have amassed large amounts of rare variants found in the human exome. Yet caution is needed when relating a genetic variation to a rare human disease, because any genome sequence contains many potentially functional rare variants that could result in false-positive associations to disease [82]. Particularly for non-coding genetic variants, assessing pathogenicity is a challenging task given the lack of knowledge on predicting their functional consequences.

An important step towards strengthening the identification of disease-causing genetic variants will be efforts to aggregate and categorize large amounts of common and rare (population-specific) genetic variation into a fully annotated pan-genome data structure. Such a centralized data structure would serve as a general baseline and circumvent the need for comparing variant calls from patients to several different variant resources generated by a variety of consortia, for instance [120, 1, 49]. Furthermore, encapsulating all possible genetic variations into a reference sequence would greatly improve variant detection following read alignment, or even as part of read alignment against a pan-genome reference structure. This is especially relevant for structural genomic changes, which play an important role in rare disease genetics [126].

Despite tremendous efforts to capture structural variation based on discordant mapping of short reads, a major fraction remains undetected in large part because of their complexity and due to the incompleteness of the current reference genome [4, 21]. Incorporating fully resolved high-quality structural variation data into the reference, preferably from long-read sequencing data, would greatly improve the genotyping of known structural variations and limit false-positives among novel structural variation calls. This will be highly relevant in the clinical setting, where genome sequencing is expected to replace array-based copy number variation profiling within a few years.

Addition of variation in genome structure to a human reference genome sequence provides a more complete gene complement for the human population including resolved paralogous genes, genome assembly collapses, redundant regions and population-specific genes. Inclusion of the pheno-

typic effects of these variants will also pave the way for methods to *diagnose* a patient *in silico*. This is essential for disease-gene identification for patients with a hitherto unexplained genetic disease.

2.6 Cancer

Cancer is caused mostly by somatic DNA alterations that accumulate during an individual's lifetime [128]. Somatic mutations in different individuals arise independently, and recent large cancer studies have uncovered extensive *inter-patient heterogeneity* among somatic mutations, with any two tumors presenting a different complement of hundreds to tens of thousands of somatic mutations [66, 72]. Heterogeneity also manifests *intra-patient*, with different populations of cells presenting different complements of mutations in the same tumor [90, 91].

Inter-patient and intra-patient heterogeneity pose several challenges to the detection and the interpretation of somatic mutations in cancer. The availability of a pan-genome reference would greatly improve the detection of somatic mutations in general, through improved quality of read mapping to polymorphic regions, and in particular in cases when matched normal tissue is not available or when only a reduced sequence coverage can be obtained.

In addition to a pan-genome reference, a somatic cancer pan-genome, representing the variability in the observed as well as inferred background alteration rate across the genome and for different cohorts of cancer patients, would enhance the identification of genomic alterations related to the disease (*driver events*) based on their recurrence across individuals. Even more important would be the availability of a somatic pan-genome describing the general somatic variability in the human population, which would provide an accurate baseline for assessing the impact of somatic alterations.

For the medium and long term future, we envision a comprehensive cancer pan-genome to be built for each tumor patient, comprising single-cell data, haplotype information as well as sequencing data from circulating tumor cells and DNA. Such a pan-genome will most likely constitute a much better basis for therapy decisions compared to current cancer genomes which mainly represent the most abundant cell type.

³<http://evs.gs.washington.edu/EVS>

2.7 Phylogenomics

Phylogenomics reconstructs the evolutionary history of a group of species by using their complete genome sequences, and can exploit various signals such as sequence or gene content [38, 124]. Computational pan-genomics will allow genomic features with an evolutionary signal to be rapidly extracted, such as gene content tables, sequence alignments of shared marker genes, genome-wide SNPs, or internal transcribed spacer (ITS) sequences, depending on the level of relatedness of the included organisms. This will facilitate evolutionary analyses ranging from the reconstruction of species phylogenies, where heterogeneity between genomes is high [23], to tracing epidemic outbreaks and cancer lineages within a patient, where heterogeneity between genomes is low. For example, the yeast dataset described in [93] allowed a phylogenetic classification based on the presence and location of mobile elements in several strains of *S. cerevisiae*. Computational pan-genomics would also enhance such mobilomics analysis when the pan-genome is built from a set of distinct strains of the same species.

Unambiguous phylogenomic trees of organismal or cellular lineages form invaluable input data for applications in various biomedical fields, for example to map the evolutionary dynamics of mutation patterns in genomes [16] or to understand the transfer of antibiotic resistance plasmids [31]. At the same time, the size of the pan-genome often hampers the inference of such a “tree of life” computationally as well as conceptually. One clear bonus offered by the pan-genome, is that for traditional phylogenomics only the best aligned, and most well behaved residues of a multiple sequence alignment can be retained. In contrast, the pan-genomic representation of multiple genomes allows for a clear encoding of the various genomic mutations in a model of the evolutionary events. This leads to the possibility for radical new evolutionary discoveries in fields including the origin of complex life [141], the origin of animals [97] and plants [147], or the spread of pathogens [44, 57], but also inferring the relationships between cancer lineages within a single patient [52, 24].

2.8 Gene Regulation

All genomes contain functional elements such as genes, but also numerous signals to control their expression and to ensure the genome’s maintenance. Transcription factors that need to bind to specific motifs in a regulatory region of a gene in order to initiate transcription are one example of such a mechanism [25]. Likewise, the replication and sharing of the DNA between two daughter cells is a crucial and hence highly regulated mechanism. Again, specific sequence motifs control order and origins of genome replication [32].

Computational tasks related to understanding gene regulation therefore include detecting a given motif or to find all its occurrences in certain genomic context, or in a discovery mode, to extract or infer such over-represented motifs from positive and negative examples; refer to [116, 29] for reviews. These tasks are relevant for all species, because the underlying biological phenomena are universal.

Clearly, for such sequence analysis tasks, the use of a pan-reference and of a dedicated data structure for interrogating multiple genomes efficiently promises substantial improvements both in terms of sensitivity and efficiency. Indeed, a pan-genome structure gives access to all variants for a given genomic locus: searching the motif against it informs us about the presence or absence of the motif, but also on its frequency in the sampled population, strains, or species—depending on the phylogenetic level at which the pan-genome was computed. This information is useful for estimating the statistical significance of occurrences, but also to gain evolutionary insights on the mechanisms under investigation.

Overall, a pan-genome search increases the chance to detect occurrences across strains or species. Improvement in efficiency follows from the fact that a pan-genome search avoids searching each genome individually, and offers a common coordinate system. Furthermore, it might facilitate an easy and integrated use of information on both sequence and evolutionary conservation.

3 Impact of Sequencing Technology on Pan-Genomics

Next-generation short-read sequencing has contributed tremendously to the increase in the known number of genetic variations in genomes of many species. The inherent limitations of commonly used short-read sequencing are three-fold. First, the short read lengths prohibit the interrogation of genomic regions consisting of repetitive stretches, the direct phasing of genetic variants [51], and the detection of large structural variations [21]. Second, non-random errors hamper the detection of genetic variations [6]. Third, there is a non-uniform distribution of sequencing coverage [115] due to various factors including biases in PCR amplification, polymerase processivity, and bridge amplification.

Establishing pan-genome sequences ideally requires a complete set of *phased* – that is, haplotype resolved – genetic variations. Experimental techniques to capture such linkage information have witnessed significant progress recently, as reviewed by [125]. Ultimately, specialized protocols for haplotype-resolved sequencing will be rendered obsolete once sufficiently long sequencing reads are routinely available.

The most promising developments in sequencing technology involve single-molecule real-time sequencing of native DNA strands. Currently, SMRT sequencing (Pacific Biosciences) is widely used for variation discovery and genome assembly [21]. The MinION device (Oxford Nanopore Technologies) [118] provides even longer reads of single DNA molecules, but has been reported to exhibit GC biases [71]. Data generated on the MinION platform has been successfully used for assembly of small genomes and for unraveling the structure of complex genomic regions [7, 84].

Despite this progress, sequencing reads are not yet sufficiently long to traverse and assemble all repeat structures and other complementary technologies are necessary to investigate large, more complex variation. Presently, array comparative genomic hybridization (arrayCGH), synthetic long reads (Moleculo [70], 10X Genomics [146]), chromatin interaction measurements [19] and high-throughput optical mapping [130, 54, 85] all aid the detection of structural variation.

Beyond interrogating genomes, sequencing tech-

nologies also serve to measure various other signals that can be seen as additional layers of information to be stored and analyzed in a pan-genome framework. Most notably, specialized protocols exist to measure transcriptomes, DNA-protein interaction, 3D genome structure, epigenetic information, or translatoemes. In all these cases a current challenge consists in transitioning from bulk to single-cell sequencing.

We expect that novel technologies will continue to greatly improve all mentioned applications in genomics and beyond. Nonetheless, further decreasing costs and conducting appropriate benchmark studies that illustrate specificity and sensitivity are problems yet to be tackled.

4 Data Structures

4.1 Design Goals

Different applications give rise to different requirements for data structures that represent pan-genomes. Refer to Figure 1 for a schematic overview. Depending on the specific application, a pan-genome data structure may need to offer any of the following capabilities:

Construction and Maintenance. Pan-genomes should be constructable from different independent sources, such as (1) existing linear reference genomes and their variants, (2) haplotype reference panels, and (3) raw reads, either from bulk sequencing of complex mixtures or from multiple samples sequenced separately. The data structure should allow dynamic updates of stored information without rebuilding the entire data structure, including local modifications such as adding a new genetic variant, insertions of new genomes, deletion of contained genomes.

Coordinate System. A pan-genome defines the space in which (pan-)genomic analyses take place. It should provide a “coordinate system” to unambiguously identify genetic loci and (potentially nested) genetic variants. Desirable properties of such a “coordinate system” include that nearby positions should have similar coordinates, paths representing genomes should correspond to monotonic

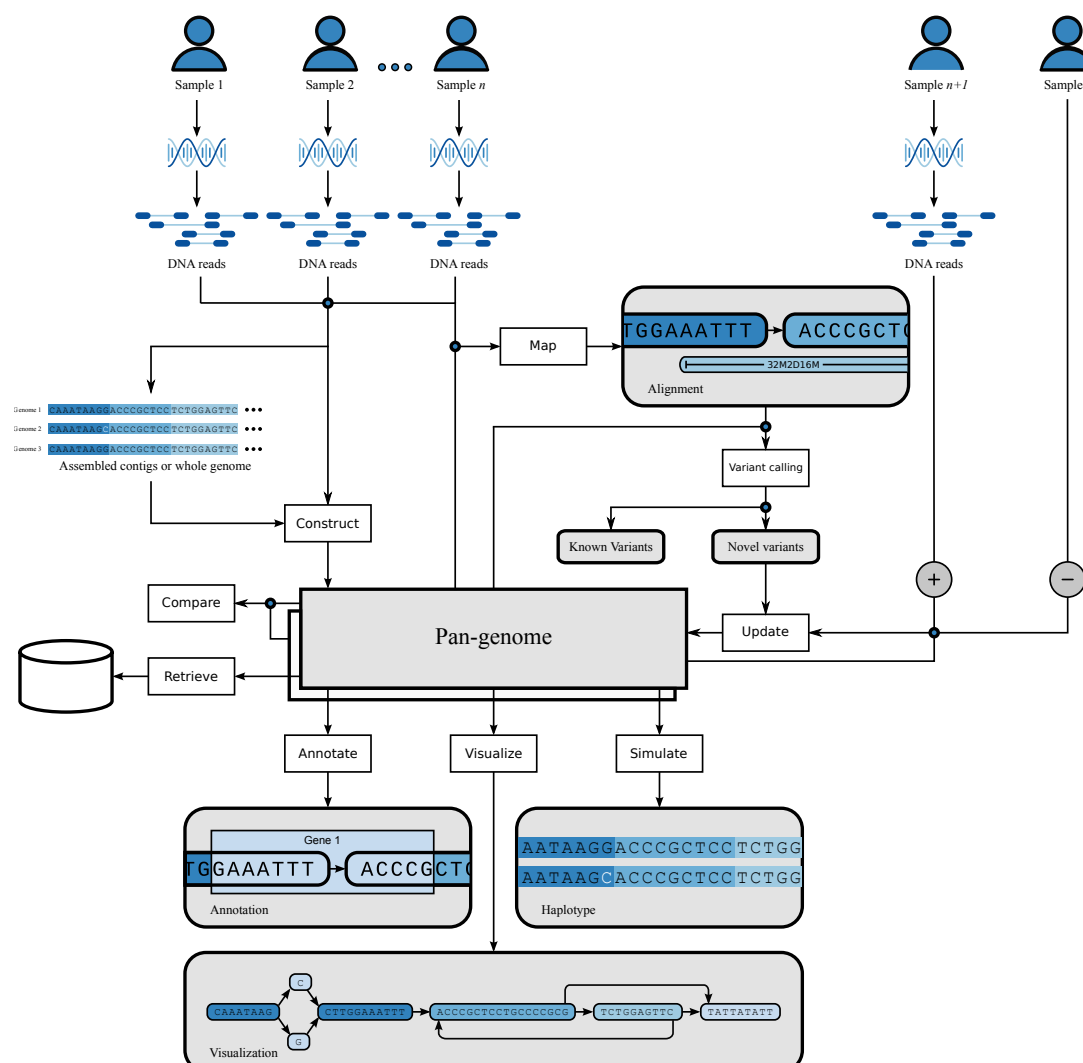


Figure 1: Illustration of operations to be supported by a pan-genome data structure.

sequences of coordinates where possible, and coordinates should be concise and interpretable.

Biological Features and Computational Layers. Annotation of biological features should be coherently provided across all individual genomes. Computationally these features represent additional layers on top of pan-genomes. This includes information about (1) genes, introns, transcription factor binding sites; (2) epigenetic properties; (3) linkages, including haplotypes; (4) gene regulation; (5) transcriptional units; (6) genomic 3D structure and (7) taxonomy among individuals.

Data Retrieval. A pan-genome data structure should provide positional access to individual genome sequences, access to all variants and to the corresponding allele frequencies. Haplotypes should be reconstructable including information about all maximal blocks and linkage disequilibrium between two variants.

Searching within Pan-Genomes. Comparisons of short and long sequences (e.g. reads) with the pan-genome ideally results in the corresponding location and the best matching individual genome(s). This scenario may occur for transcrip-

tomic data as well as for DNA re-sequencing data, facilitating the identification of known variants in new samples (variant calling).

Comparison among Pan-Genomes. Given any pair of genomes within a pan-genome, we expect a data structure to highlight differences, variable and conserved regions, as well as common syntenic regions. Beyond that, a global comparison of two (or more) pan-genomes, e.g. with respect to gene content or population differentiation, should be supported.

Simulation. A pan-genome data structure should support the generation (sampling) of individual genomes similar to the genomes it contains.

Visualization. All information within a data structure should be easily accessible for human eyes by visualization support on different scales. This includes visualization of global genome structure, structural variants on genome level and local variants on nucleotide level, but also biological features and other computational layers (see above) should be represented.

Efficiency. We expect a data structure to use as little space on disk and memory as possible, while being compatible to computational tools with a low running time. Supporting specialized hardware, such as general purpose graphics processing units (GPGPUs) or field-programmable gate arrays (FPGAs), is partly an implementation detail. Yet, in some cases, the target platform can influence data structure design significantly.

4.2 Approaches

There are natural trade-offs between some of the desiderata discussed above. For instance, the capability to allow dynamic updates might be difficult to achieve while using only small space and allowing for efficient indexing. It is one of the core challenges of computational pan-genomics to design data structures that support (some of) the above query types efficiently. While desirable in principle, we consider it difficult, if not impossible, to develop a solution that meets *all* the listed requirements at

once. Therefore, future research should aim to delineate the compromises that may have to be made and thereby provide guidance on which solution is suitable for which application scenario. As the field matures, additional queries will appear, and data structures will need to adapt to support them.

In the following, we discuss traditional approaches to meet fundamental requirements for genome analysis, first extensions for pan-genomes, as well as future challenges.

Unaligned Sets of Sequences. The conceptually simplest representation of a pan-genome consists of a set of individual sequences (Figure 2a), which might be either whole genomes or parts of it. The traditional view of a species' pan-genome as the set of all genes [135], which is prevalent in microbiology, can be considered an example for this. Unaligned whole genome sequences on the other hand are, in general, of limited utility for most applications, especially when the genomes are long. So we consider collections of individual genomes mostly as input to build the more advanced representations discussed in the following.

Multiple-Sequence-Alignment Based Representations. Pan-genomes can be represented by alignments of multiple genomes. In a multiple sequence alignment (MSA), the input sequences are aligned by inserting gap characters into each sequence (Figure 2b). The result is a matrix, where each column represents putatively homologous characters. Refer to [42, 102] for reviews on current methods and remaining challenges. Such classical colinear alignments are not able to capture larger rearrangements like inversions and translocations well and hence only apply to short genomic regions such as single genes or to very similar genomes.

One advantage of using an MSA as a representation of a pan-genome is that it immediately defines a coordinate system across genomes: a column in the alignment represents a location in the pan-genome. MSAs furthermore support many comparison tasks.

All approaches designed for linear reference genomes can, in principle, be extended to multiple alignments at the expense of adding bookkeeping data structures to record where the gaps are.

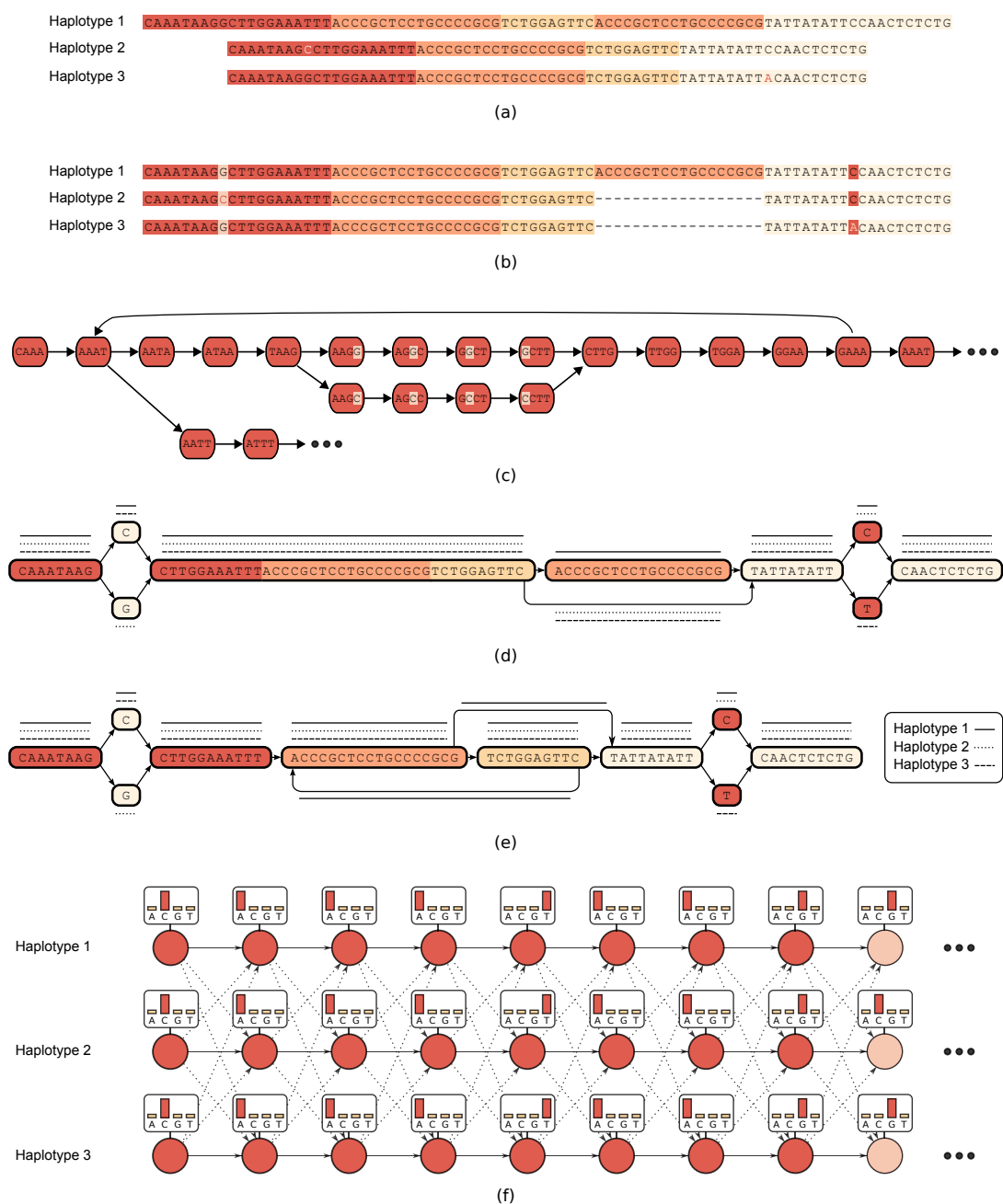


Figure 2: Selected examples of pan-genome representations: (a) three unaligned sequences, colors highlight similarities; (b) a multiple sequence alignment of the same three sequences; (c) the De Bruijn graph of the first (red) sequence block; (d) acyclic sequence graph, paths representing the three haplotypes shown as solid/dashed/dotted lines; (e) cyclic sequence graph; (f) Li-Stephens model of the first nine characters with states indicated by circles, emission distributions given in boxes and transitions given by arrows; dashed arrows indicate the (less likely) "recombination" transitions.

Efficient data structures for prefix sum, rank, and select queries exist [98], which can be used for the purpose of doing projections to and from a sequence and its gapped version as a row of an MSA. Multiple sequence alignments can be compactly represented by journaled string trees [111]. This data structure also allows for efficiently executing sequential algorithms on all genomes in the MSA simultaneously. One example for such a sequential algorithm is online pattern matching, that is, searching all genomes for the exact or approximate occurrence of a pattern without building an index structure first.

When aligning two or more whole genomes, structural differences such as inversions and translocations need to be taken into account. Standard methods for a colinear MSA are therefore not applicable. Instead, one aims to partition the input genomes into blocks such that sequences within blocks can be aligned colinearly. Creating such a partitioning is a non-trivial task in itself and mostly approached through graph data structures that represent local sequence similarities and adjacencies. On the one hand, such graphs therefore facilitate whole genome alignment. On the other hand, they can be understood as representations of the pan-genome. Concrete realizations of this idea include A-Bruijn graphs [108], Enredo graphs [106] and Cactus graphs [104, 105]. For detailed definitions and a comparison of these concepts we refer the reader to the review [68].

Block-based multiple sequence alignments can also serve as the basis for a coordinate system on a pan-genome: by numbering blocks as well as numbering columns inside each colinearly aligned block, a notion of a position in a pan-genome can be defined. This idea is explored by Herbig et al. [56], who furthermore show how it can serve as a foundation for visualization.

k-mer-Based Approaches. Starting from either assembled genomes, contigs, or just collections of (error-corrected) reads, a pan-genome can also be represented as a collection of k -mers, i.e. strings of length k . The task of efficiently counting all k -mers occurring in an input sequence has been studied extensively in recent years and many solutions are available, including Jellyfish [92], DSK [114] and KMC2 [33]. Such a k -mer collection is a rep-

resentation of the corresponding de Bruijn Graph (DBG), illustrated in Figure 2c. DBGs were introduced in the context of sequence assembly [94], but can be used as pan-genome representations supporting many applications beyond assembly. When k -mer neighborhood queries are sufficient, and no k -mer membership queries are required, then even more space-efficient data structures for DBGs exist [22].

When building DBGs for multiple input samples, one can augment each k -mer by the set of samples containing it. This idea is realized in colored DBGs where we color each k -mer according to the input samples it occurs in. Colored DBGs have been used successfully for reference-free variant calling and genotyping [63]. Recently, Holley et al. [58] introduced bloom filter tries, a data structure able to efficiently encode such colored DBGs.

For k -mer based representations of pan-genomes, the length k is obviously an important parameter and picking the right value depends on the intended application. Data structures able to represent a pan-genome at different granularities (i.e. at different values of k) are hence an interesting research topic. For instance, Minkin et al. [95] show that iteratively increasing k helps to capture nested synteny structure.

Pan-genomes encompassing many species can be encoded as a mapping between k -mers and clades: given a phylogenetic tree, each k -mer is mapped to the lowest common ancestor of all genomes containing it. This technique was introduced by Wood and Salzberg [143], who show that it efficiently supports the task of analyzing the composition of metagenomic samples.

Advantages of k -mer-based representations include simplicity, speed, and robustness: it is not necessary to produce an assembly or an alignment, which both can be error-prone, and very efficient data structures to store them exist. However, they do not explicitly represent structural information at distances greater than the k -mer length. For applications where such information is needed, DBGs can sometimes serve as a basis to design richer data structures. Colored DBGs [63, 58] are an example of this since they store information about occurrence in individual genomes on top of each k -mer.

Further Sequence Graphs. Building on the above ideas, more general approaches conceptualize a pan-genome as an (edge- or node-labeled) graph of generic pieces of sequence. Such graphs are not necessarily constructed using an MSA and the constituting sequences are not necessarily fixed-length k -mers. Figures 2d and 2e show examples of a directed and an undirected sequence graph, respectively. Individual genomes can be represented as paths in such graphs and node identifiers can serve as a “coordinate system”.

Compressed DBGs (also called compacted DBGs), which collapse chains of non-branching nodes in a DBG into a single node, are an example of this. Marcus et al. [89] show how such compressed DBGs can be constructed for a pan-genome by first identifying maximal exact matches using a suffix tree, by-passing uncompressed DBGs. Beller and Ohlebusch [14] and Baier et al. [8] show how the same can be achieved more efficiently, using an FM index resp. compressed suffix trees and the Burrows-Wheeler transform.

Useful data structures for pan-genomes may combine some of the basic approaches discussed so far. For example, PanCake [45] uses a graph-based structure to represent common genomic segments and uses a compressed multiple-alignment based representation in each node of the graph. Dilthey et al. [34] propose a generative model by representing sequence variation in a k -mer-emitting HMM.

Further examples of implementations of sequence graphs include the Global Alliance for Genomics and Health (GA4GH) “side graph” data model and the FASTG format⁴. Side graphs represent a pan-genome as a set of sequences and an additional set of joins, each of which defines an extra adjacency between the sides of two bases within the sequences. The GA4GH graph tools⁵ allow side graphs and embeddings of individual sampled genomes in that graph to be made available over the Internet, for data distribution and remote analysis.

Haplotype-Centric Models. When a fixed set of (non-nested) sequence variants is considered, every haplotype in a population can be represented as a string of fixed length. The character at position

k reflects the status of the k -th variant. When all variants are bi-allelic, then these haplotype strings are formed over a binary alphabet. Such collections of haplotypes are often referred to as *haplotype panels*. This representation is favorable for many population genetic analyses since it makes shared blocks of haplotypes more easily accessible, for instance compared to sets of paths in a graph.

A recent data structure to represent haplotype panels, termed Positional Burrows-Wheeler Transform (PBWT) [37], facilitates compression and supports the enumeration of maximal haplotype matches.

One of the most widely used haplotype-based models is the Li-Stephens model [77]. In a nutshell, it can be viewed as a hidden Markov model (HMM) with a grid of states with one row per haplotype and one column per variant, as sketched in Figure 2f. Transitions are designed in a way such that staying on the same haplotype is likely but jumping to another one is also possible with less probability. It hence is a generative probabilistic model for haplotypes that allows for sampling new individuals and provides conditional probabilities for new haplotypes given the haplotypes contained in the model.

5 Computational Challenges

Pan-genomic data have all of the standard properties of big data — in particular, volume, variety, velocity and veracity. Especially due to the sheer size of generated sequencing data, extreme heterogeneity of data and complex interaction on different levels, pan-genomics comes with big challenges for algorithm and software development [12]. The International Cancer Genome Consortium (ICGC) has amassed a dataset in excess of two petabytes in just five years with the conclusion to store data generally in clouds, providing an elastic, dynamic and parallel way of processing data in a cheap, flexible, reliable and secure manner [127].

Currently large high computing infrastructure providers and large public repositories (e.g. NCBI/EBI/DDBJ) are completely separated. We need hybrids that offer both large public repositories as well as the computing power to analyze these in the context of individual samples/data. We consider it desirable to bring the computation as close

⁴<http://fastg.sourceforge.net>

⁵<https://github.com/ga4gh/server> and <https://github.com/ga4gh/schemas>

as possible to the data by uploading queries or in-database computing.

These general Big-Data-related challenges apply to all individual computational problems we discuss below.

5.1 Read Mapping

Given a set of reads sequenced from a donor, *read mapping* consists in identifying parts of the reference genome matching each read. Read mapping to a pan-genome has a potential to improve alignment accuracy and subsequent variant calling, especially in genomic regions with a high density of (complex) variants.

For a single reference sequence, the read mapping problem has mostly been solved by indexing the reference into a data structure that supports efficient pattern search queries. Most successful approaches use k -mer based or Burrows-Wheeler transform based indexes, as reviewed in [75]. Indexing a pan-genome is more complicated.

Efficient indexing of a set of reference genomes for read mapping was first studied in [86, 87]. The approach uses compressed data structures, exploiting the redundancy of long runs of the same letter in the Burrows-Wheeler transform of a collection of similar genomes. This approach yields a reasonably compressed representation of the pan-genome, but read alignment efficiency is hampered by the fact that most reads map to all of the references, and that extraction of these occurrence locations from a compressed index is potentially slow. More recently, approaches based on Lempel-Ziv compression have been proposed to speed-up the reporting of occurrences, as reviewed in [48].

The earliest approach to index a *sequence graph* (see Section 4.2) was proposed in [117], where k -mer indexing on the paths of such a graph was used; instead of a full sequence graph, a *core* sequence graph was used where columns were merged in regions of high similarity (core genome) to avoid extensive branching in the graph. After finding seed occurrences for a read in this graph, the alignment was refined locally using dynamic programming. Similar k -mer indexing on sequence graphs has since been used and extended in several read mapping tools such as MuGI [27], BGREAT [78]

and VG⁶.

Instead of k -mer indexing, one can also use Burrows-Wheeler-based approaches, based on appending extracted contexts around variations to the reference genome [60]. Context extraction approaches work only on limited pattern length, as with long patterns they suffer from a combinatorial explosion in regions with many variants; the same can happen with a full sequence graph when all nearby k -mer hit combinations are checked using dynamic programming. There is also a special Burrows-Wheeler transform and an index based on that for a sequence graph [122, 123]. This approach works on any pattern length, but the index itself can be of exponential size in the worst case; best case and average case bounds are similar to the run-length compressed indexes for set of references like [87]. The approach is also likely to work without exponential growth on a core sequence graph of [117], but as far as we know, this combination has not been explored in practice. A recent implementation⁷ avoids the worst case exponential behavior by stopping the construction early; if this happens, the approach also limits the maximum read length. This implementation has been integrated into VG as an alternative indexing approach. HISAT2⁸ implements an index structure that is also based on [122], but builds many small index structures that combinedly cover the whole genome.

In summary, a number of approaches to perform read mapping against a pan-genome reference under various representation models exist, and efficient implementations for daily usage are under active development. However, we consider this field as being far from saturated and still expect considerable progress in both algorithmic and software engineering aspects. To reach the full potential of these developments, the interactions between read mapping and variant calling methods need to be considered.

5.2 Variant Calling and Genotyping

The task of determining the differences between a sequenced donor genome and a given (linear) reference genome is commonly referred to as *variant*

⁶<https://github.com/ekg/vg>

⁷<https://github.com/jltsiren/gcsa2>

⁸<https://ccb.jhu.edu/software/hisat2/index.shtml>

calling. In case of diploid or polyploid organisms, we additionally want to determine the corresponding *genotype*. In the face of pan-genome data structures, variant calling becomes decomposed into two steps: identifying *known* variants already represented in the data structure and calling *novel* variants. Refer to Schneeberger et al. [117] for an early work on pan-genome variant calling. They do not only show the feasibility of short read alignment against a graph representing a pan-genome reference (see Section 5.1) but also demonstrate its positive impact on variation calling in the frame of the Arabidopsis 1001 Genomes Project.

Known Variants. By using a pan-genome reference, one merges read mapping and calling of known variants into a single step. Read alignments to sequence variants encapsulated in our pan-genome data structure indicate the presence of these variants in the donor genome. In particular, this applies not only to small variants which can be covered by a single read (such as SNPs and indels), but also to larger structural variants such as inversions or large deletions. Integrating those steps potentially decreases overall processing time and, more importantly, removes read-mapping biases towards the reference allele and hence improves accuracy of calling known variants. One important challenge is to statistically control read mapping ambiguity on a pan-genome data structure. Leveraging the associated statistical models for estimating genotype likelihoods is expected to lead to significant improvements in genotyping.

As a first major step in that direction, Dilthey et al. [34] cast the (diploid) variant calling problem into finding a pair of paths through a pan-genome reference represented as a k -mer-emitting Hidden Markov Model. They demonstrate that this leads to substantially improved performance in the variation-rich MHC region.

Novel Variants. Detecting variants not present in a pan-genome data structure is similar to traditional variant calling with respect to a linear reference genome. Still, differences exist that require special attention. The most straightforward way to use established variant calling methods is to use read alignments to a pan-genome and project them onto a linear sequence. For small variants such as

SNPs and indels, that are contained within a read, this approach is likely to be successful. Methods to characterize larger structural variation (SV) need to be significantly updated. SV calling methods are usually classified into four categories based on the used signal: read pair, read depth, split read, and assembly, as reviewed by Alkan et al. [4]. Each of these paradigms has its merits and shortcomings and state-of-the-art approaches usually combine multiple techniques [129]. Each of these ideas can and should be translated into the realm of pan-genomes. For split-read and assembly based approaches, the problem of aligning reads and contigs, respectively, to a pan-genome data structure (while allowing alignments to cross SV breakpoints) needs to be addressed. In case of read pair methods, a different notion of “distance” is implied by the pan-genome model and has to be taken into account. For read depth methods, statistical models of read mapping uncertainty on pan-genomes have to be combined with models for coverage (biases). Developing standards for reporting and exchanging sets of potentially nested variant calls is of great importance.

Somatic Mutations. Calling somatic mutations from paired tumor/normal samples is an important step in molecular oncology studies. Refer to Section 2.6 for details and to [3] for a comparison of current work flows. Calling somatic variants is significantly more difficult compared to calling germline variants, mostly due to tumor heterogeneity, the prevalence of structural variants, and the fact that most somatic variants will be novel. Pan-genome data structures promise to be extremely useful in cancer studies for the stable detection of somatic variants. A conceivable approach for leveraging pan-genome data structures in this context would be to map reads from the matched normal sample to the pan-reference, call germline mutations, create a restricted pan-genome with detected variants and map tumor reads to that pan-reference for calling somatic mutations. There are many more potential applications including building a pan-genome representation of a heterogeneous tumor to be used as a starting point for retracing tumor evolution.

Storing Variants. Storing and exchanging variant calls genotyped in a large cohort of samples increasingly becomes a bottleneck with growing cohort sizes. Some improvement is achieved by adopting binary instead of text-based data formats for variant calls, i.e. using BCF instead of VCF⁹, but more efficient approaches are urgently needed. Organizing data by individual rather than by variant while sorting variants by allele frequency has proven beneficial for compression and some retrieval tasks [73]. We expect the question of storing, querying and exchanging variant data to remain an active and relevant field of research in the coming years.

5.3 Haplotype Phasing

Humans are diploid, that is, each chromosome comes in two copies, one inherited from the mother and one inherited from the father. The individual sequences of these two chromosomal copies are referred to as *haplotypes*, where one often restricts the attention to polymorphic sites. The process of assigning each allele at heterozygous loci to one of the two haplotypes is referred to as *phasing*. Plants are often polyploid. For example, wheat can be tetra- (= 4 copies) or hexaploid (= 6 copies), while certain strawberries are even decaploid (= 10 copies). As an extreme, the “ploidy” of viral quasispecies, that is the number of different viral strains that populate an infected person (see Section 2.3) is usually unknown and large. The same applies to heterogeneous tumors, as discussed above.

Pan-genome data structures have the potential to, on the one hand, store haplotype information and, on the other hand, be instrumental for phasing. Currently, several approaches for obtaining haplotype information exist. *Statistical phasing* [17] uses genotype information of large cohorts to reconstruct haplotypes of all individuals based on the assumption that haplotype blocks are conserved in a population. Once sets of haplotypes, called reference panels, are known, additional individuals can be phased by expressing the new haplotypes as a mosaic of the already known ones. The question of how to best organize and store reference panels is open. To this end, Durbin [37] has proposed the aforementioned PWB index struc-

ture. We consider marrying reference panels to pan-genome data structures an important topic for future research.

To determine haplotypes of single individuals, including rare and de novo variants, statistical approaches are not suitable and experimental techniques to measure linkage are needed. Such techniques include specialized protocols and emerging long-read sequencing platforms, as discussed in Section 3. Currently, first approaches for haplotype-resolved local assembly are being developed [113]. More literature exists on the problem of phasing from aligned long reads, e.g. [107, 109, 69]. In practice, this technique is hampered by insufficient alignment quality of long error-prone reads. Since phasing is based on heterozygous loci, avoiding allelic biases during read mapping by means of pan-genome data structures can contribute to solving this problem. Combining the virtues of read-based phasing with statistical information from reference panels is an active area of research [70]. Leveraging pan-genome data structures that encode reference haplotypes towards this goal constitutes a promising research direction.

These problems are amplified when phasing organisms or mixtures of higher or unknown ploidy such as plants, viral quasispecies or tumors. Algorithms with manageable runtime on polyploid organisms [2, 15] and for the reconstruction of quasispecies [145, 134] require the use of specialized techniques (especially when allele frequencies drop below sequencing error rates). Extending these approaches to pan-genome data structures, as outlined above for the diploid case, is another challenging topic for future research.

5.4 Visualization

Pan-genomics introduces new challenges for data visualization. Fundamentally, the problems relate to how to usefully view a large set of genomes and their homology relationships, and involve questions of scale and useful presentation in the face of huge volumes of information.

At a high-level of abstraction, pan-genome bag-of-genes approaches can be visualized using methods for comparing sets, such as Venn diagrams, flower plots, and related representations. For example, the recent tool Pan-Tetris visualizes a gene-based pan-genome in a grid [55], color-coding ad-

⁹<http://samtools.github.io/hts-specs/>

ditional annotation. For divergent genomes, as in bacterial- and meta- pan-genomics, and where complete assembly is not possible, such approaches provide useful summary information.

For the viewing of individual, assembled genomes or sequences, genome browsers and applications frequently display an individual sequence along a linear or circular axis upon which other genomics information is visualized, as reviewed in [101]. This trope, which is popular and widely understood, forces interpretation through the lens of one chosen genome. When this genome is a distantly related reference genome there is a visual reference bias which may lead to misinterpretation.

Pan-genome displays can potentially help to alleviate this visual bias. One option is to aim to improve linear visualizations: either the chosen individual reference sequence can be replaced by a more visually useful imputed pan-genome reference, or the pan-genome data structures which relate different genomes in the population can be used to translate information to the most closely related genome possible for the display. In the former case, a pan-genome display can be made more inclusive than any single genome [100]. At the base level such inclusive displays are somewhat analogous to popular multiple sequence alignment displays such as Mauve [28] or Jalview [137] that focus on displaying all the differences between a set of sequence as clearly as possible. The latter case, translation, where a pan-genome alignment is used to show information on the most closely related genome possible, is likely to become more popular as the number of available personal genomes grows, see [99] for an early example of such an approach.

More adventurously than linear layouts, pan-genome displays can attempt to visualize graphs of variation. This has the flexibility of allowing arbitrary genome variation within a clean semantic model, but can prove visually complex for even small, non-trivial examples. For example, a graph of a few dozen bacterial strains contains tens to hundreds of thousands of nodes and edges. So far graph visualizations have proved popular for assemblies, and the visualization of heterozygosity, for example DISCOVAR [139] contains a module that allows you to visualize subsets of an assembly graph in a figure. One popular tool is Cytoscape [119], which is a generic biological graph/network visualization tool, but lacks scalability and semantic nav-

igation. Another tool, Bandage [140], visualizes *de novo* assembly graphs specifically.

A number of challenges exist moving forwards. In a useful visualization it will be possible to navigate and to zoom in and out on pan-genome structures. Zooming should be done semantically, i.e. different zoom levels can use different representations of the data to convey biologically relevant information. The upper scales should give information about global genome structure. Zooming in the visuals should focus on structural variants in a genomic region and the most zoomed in views should enable exploration of local variants on nucleotide level. Furthermore these visuals need to be put in the context of the phylogeny, e.g. the relation of the various samples that went into the pan-genome. This will enable rapid identification and interpretation of observed variants. Finally, any pan-genome graph visualization should offer the same basic features that current reference based genome browsers have. There should be visual ways to indicate biologically interesting features such as gene annotations and position based continuous valued signals such as wiggle tracks in the UCSC genome browser. Basic analytical capabilities would be beneficial to visually highlight interesting biologically relevant mutations. For example, it would be useful to have different visual representations for different types of mutations: indels, (non)-synonymous SNPs, structural variants, repeats etc.

5.5 Data Uncertainty Propagation

One of the computational (and modeling) challenges facing the field of pan-genomics is how to deal with data uncertainty propagation through the individual steps of analysis pipelines. In order to do so, the individual processing steps need to be able to take uncertain data as input and to provide a ‘level of confidence’ for the output made. This can, for instance, be done in the form of posterior probabilities. Examples where this is already common practice include read mapping qualities and genotype likelihoods.

Computing a reasonable confidence level commonly relies on weighing alternative explanations for the observed data. In the case of read mapping for example, having an extensive list of alternative mapping locations aids in estimating the probability of the alignment being correct. A pan-genome

expands the space of possible explanations and can, therefore, facilitate the construction of fairer and more informative confidence levels.

As an illustration, consider a pipeline including read mapping, variant calling and genotyping, phasing and association testing. Substantial uncertainty and sequence composition biases are already inherent to the input data generated by next-generation sequencing [67]. The following read alignment step adds uncertainty in read placement, leading to uncertain coverage and uncertain fragment lengths. These uncertainties translate into uncertainties in variant calling and genotyping, and further into uncertainties in phasing. This, finally, results in uncertainties in association testing in genome-wide association studies. The precise quantification of the propagation of these effects is largely unclear. The advent of ever larger and refined panels, supported by appropriate pan-genome data structures, bears the promise of making quantification and alleviation of such effects possible.

6 Conclusions

Already today, the DNA having been sequenced for many biologically coherent ensembles—such as certain taxonomic units or virus populations—likely captures the majority of their frequently occurring genetic variation. Still, the pace at which genomes are currently sequenced is on a steep rise, thanks to accumulation of sequencers in laboratories and frequent, significant advances in sequencing technology. Therefore, capturing *all of genomes*, in terms of genetic variation content and abundance, is no longer wishful thinking, but will materialize soon for many species, populations and cancer genomes. In other words, life sciences have entered the era of *pan-genomics*, which is characterized by knowing *all* major genetic variation of a collection of genomes of interest. In this white paper, we have been addressing how to arrange and analyze this incredible wealth of knowledge and also how to deal with some of the consequences in downstream analyses.

The computational aspects that need to be considered fan out across a large variety of particular challenges, usually governed by the realm of application they stem from. We have listed the many facets of pan-genomes in terms of func-

tionality, annotational detail, computational efficiency issues and visualization. We have discussed how the availability of well-arranged pan-genomes will affect population genetics, cancer genomics, pathogen research, plant breeding, phylogenomics, functional genomics as well as genetic disease research and genome-wide association studies. We have surveyed the impact of sequencing technology advances on the field of pan-genomics, and we have considered also the complications that come along with these advances. We have put particular emphasis on data structures and supporting algorithms that make it possible to consistently work with pan-genomes. One of the currently most evident processes in computational pan-genome research is the move away from linear reference genomes towards reference systems that are rooted in graph theory in some form. The effort of the Data Working Group of the Global Alliance for Genomics and Health (GA4GH) is a prominent example for this. We have also discussed how the transition in terms of data structures will affect operations such as read mapping, variant discovery, genotyping and phasing, all of which are at the core of modern genomics research. Last but not least, we have analyzed the issues that arise in visualizing pan-genomes, and we have also briefly discussed future issues in uncertain data handling, recently an ever recurring theme in genome data analysis, often arising from the repetitive structure of many genomes.

We have concentrated on computational challenges of pan-genomics in this survey. We are aware that there are also political challenges that have to be addressed that concern data sharing and privacy. Clearly, the usefulness of any pan-genomic representation will increase with the number of genomes it represents, strengthening its expressive and statistical power. Unfortunately, however, only a fraction of the sequenced data is currently publicly available. This is partly due to the confidential nature of human genetic data, but also, to a large extent, by missing policies and incentives to make genomic data open access or to prevent intentional withholding of data. Funding agencies like the National Institutes of Health (NIH) in the US have started to address these issues [103]¹⁰.

Overall, we have provided a broad overview of

¹⁰see also <http://www.nih.gov/news-events/news-releases/nih-issues-finalized-policy-genomic-data-sharing>

computational pan-genomics issues, which we hope will serve as a reference for future research proposals and projects. However, so far, we have mostly been addressing how to deal with genomes as sequences, that is from a “one-dimensional” point of view, and so we have been focusing on storing and analyzing sequences and the mutual relations of particular subsequence patches, like variant alleles and their interlinkage, genes and/or transcripts. We have done this because we believe that at this point in genomics history, only the consistent exploration and annotation of exhaustive amounts of sequence information can lay the solid foundation for additional “pan-genomics oriented” steps.

Yet, even after having resolved the corresponding issues—and we are hopeful that, at this point, our summary has helped to consistently structure these—there is more to follow. New approaches have already appeared on the horizon that will benefit from the cornerstone provided by primarily sequence-driven pan-genomics. For example, it can be expected that one can lift pan-genomes into three dimensions in the mid-term future, thanks to rapidly developing technologies that allow to infer their three-dimensional conformation. This will mean that future, three-dimensional pan-genomes will not only represent all sequence variation applying for species or populations, but also encode their spatial organization as well as their mutual relationships in that respect.

Epigenomics topics have not been exhaustively addressed here either, but will need to be addressed as soon as the first “primary” pan-genomes stand. Technologies, by which to not only monitor sequential and three-dimensional arrangement, but also additional biochemical modifications have likewise been on a steep rise recently. Most importantly, we will be in position to link sequential pan-genomes to maps that indicate hypo- and hypermethylated regions relatively soon. Likely, the integration of such basic biochemical modification will serve as template for further, often more complex elements of biochemical genomic maps.

In summary, the emergence of computational pan-genomics as a field is an expression of a major advance in contemporary genomics research. For the first time, we have entered an era that holds the promise to close large gaps in global maps of genomes and to draw the full picture of their vari-

ability. We therefore believe that we can expect to witness amazing, encompassing insights about extent, pace, and nature of evolution in the mid-term future.

Acknowledgments

We are deeply grateful to the Lorentz Center for hosting the workshop “Future Perspectives in Computational Pan-Genomics” (June 8–12, 2015), which gave rise to this paper. In particular, we like to thank the Lorentz Center staff, who turned organizing and attending the workshop into a great pleasure. The workshop received additional financial support by KNAW, Bina Technologies, ERIBA, PacBio, and Genalix. Bas E. Dutilh was supported by the Netherlands Organization for Scientific Research (NWO) Vidi grant 864.14.004 and CAPES/BRASIL. Veli Mäkinen and Daniel Valenzuela were funded by the Academy of Finland, grant 284598 (CoECGR). Louis Dijkstra wishes to acknowledge partial funding by the Russian Scientific Foundation, under grant #14-11-00826. Eric Rivals thanks Défi MASTODONS from CNRS and the French ANR-12-BS02-0008 Colib’read project. Alexander Schönhuth was supported by the Netherlands Organization for Scientific Research (NWO) Vidi grant 639.072.309.

Competing Interests

Ole Schulz-Trieglaff is an employee of Illumina Inc. and receives stocks as part of his compensation. Illumina is a public company that develops and markets systems for genetic analysis.

References

- [1] 1000 Genomes Project Consortium, Goncalo R. Abecasis, Adam Auton, Lisa D. Brooks, Mark A. DePristo, Richard M. Durbin, Robert E. Handsaker, Hyun Min Kang, Gabor T. Marth, and Gil A. McVean. An integrated map of genetic variation from 1,092 human genomes. *Nature*, 491(7422):56–65, November 2012.
- [2] Derek Aguiar and Sorin Istrail. Haplotype assembly in polyploid genomes and identical by descent shared tracts. *Bioinformatics*, 29(13):i352–i360, July 2013.
- [3] Tyler S. Alioto, Ivo Buchhalter, Sophia Derdak, Barbara Hutter, Matthew D. Eldridge, Eivind Hovig, Lawrence E. Heisler, Timothy A. Beck, Jared T. Simpson, Laurie Tonon, Anne-Sophie Sertier, Ann-Marie Patch, Natalie Jäger, Philip Ginsbach, Ruben Drews, Nagarajan Paramasivam, Rolf Kabbe, Sasithorn Chotewutmontri, Nicolle Diessl, Christopher Previti, Sabine Schmidt, Benedikt Brors, Lars Feuerbach, Michael Heinold, Susanne Gröbner, Andrey Korshunov, Patrick S. Tarpey, Adam P. Butler, Jonathan Hinton, David Jones, Andrew Menzies, Keiran Raine, Rebecca Shepherd, Lucy Stebbings, Jon W. Teague, Paolo Ribeca, Francesc Castro Giner, Sergi Beltran, Emanuele Raineri, Marc Dabad, Simon C. Heath, Marta Gut, Robert E. Denroche, Nicholas J. Harding, Takafumi N. Yamaguchi, Akihiro Fujimoto, Hidewaki Nakagawa, Víctor Quesada, Rafael Valdés-Mas, Sigve Nakken, Daniel Vodák, Lawrence Bower, Andrew G. Lynch, Charlotte L. Anderson, Nicola Waddell, John V. Pearson, Sean M. Grimmond, Myron Peto, Paul Spellman, Minghui He, Cyriac Kandoth, Semin Lee, John Zhang, Louis Létourneau, Singer Ma, Sahil Seth, David Torrents, Liu Xi, David A. Wheeler, Carlos López-Otín, Elías Campo, Peter J. Campbell, Paul C. Boutros, Xose S. Puente, Daniela S. Gerhard, Stefan M. Pfister, John D. McPherson, Thomas J. Hudson, Matthias Schlesner, Peter Lichter, Roland Eils, David T. W. Jones, and Ivo G. Gut. A comprehensive assessment of somatic mutation detection in cancer using whole-genome sequencing. *Nature Communications*, 6:10001, 2015.
- [4] Can Alkan, Bradley P. Coe, and Evan E. Eichler. Genome structural variation discovery and genotyping. *Nature Reviews Genetics*, 12(5):363–376, May 2011.
- [5] Lisa Zeigler Allen, Thomas Ishoey, Mark A. Novotny, Jeffrey S. McLean, Roger S. Lasken, and Shannon J. Williamson. Single Virus Genomics: A New Tool for Virus Discovery. *PLoS ONE*, 6(3):e17722, March 2011.
- [6] Manuel Allhoff, Alexander Schönhuth, Marcel Martin, Ivan G. Costa, Sven Rahmann, and Tobias Marschall. Discovering motifs that induce sequencing errors. *BMC Bioinformatics (Proceedings of RECOMB-seq)*, 14(Suppl 5):S1, apr 2013.
- [7] Philip M. Ashton, Satheesh Nair, Tim Dallman, Salvatore Rubino, Wolfgang Rabsch, Solomon Mwaigwisya, John Wain, and Justin O’Grady. MinION nanopore sequencing identifies the position and structure of a bacterial antibiotic resistance island. *Nature Biotechnology*, 33(3):296–300, March 2015.
- [8] Uwe Baier, Timo Beller, and Enno Ohlebusch. Graphical pan-genome analysis with compressed suffix trees and the burrows-wheeler transform. *Bioinformatics*, AOP, 2015.
- [9] Michael J. Bamshad, Sarah B. Ng, Abigail W. Bigham, Holly K. Tabor, Mary J. Emond, Deborah A. Nickerson, and Jay Shendure. Exome sequencing as a tool for Mendelian disease gene discovery. *Nature Reviews Genetics*, 12(11):745–755, November 2011.
- [10] Delfina Barabaschi, Davide Guerra, Katia Lacrima, Paolo Laino, Vania Michelotti, Simona Urso, Giampiero Valè, and Luigi Cativelli. Emerging knowledge from genome sequencing of crop species. *Molecular Biotechnology*, 50(3):250–266, March 2012.
- [11] István Barthá, Jonathan M. Carlson, Chanson J. Brumme, Paul J. McLaren, Zabrina L.

- Brumme, Mina John, David W. Haas, Javier Martinez-Picado, Judith Dalmau, Cecilio López-Galíndez, Concepción Casado, Andri Rauch, Huldrych F. Günthard, Enos Bernasconi, Pietro Vernazza, Thomas Klimkait, Sabine Yerly, Stephen J. O'Brien, Jennifer Listgarten, Nico Pfeifer, Christoph Lippert, Nicolo Fusi, Zoltán Kutalik, Todd M. Allen, Viktor Müller, P. Richard Harrigan, David Heckerman, Amalio Telenti, and Jacques Fellay. A genome-to-genome analysis of associations between human genetic variation, HIV-1 sequence diversity, and viral control. *eLife*, 2:e01123, October 2013.
- [12] C. Beckstein, S. Böcker, M. Bogdan, H. M. Bücker H. Bruelheide, J. Denzler, P. Dittrich, I. Grosse, A. Hinneburg, B. König-Ries, F. Löffler, M. Marz, M. Müller-Hannemann, M. Winter, and W. Zimmermann. Explorative analysis of heterogeneous, unstructured, and uncertain data: A computer science perspective on biodiversity research. In M. Helfert, A. Holzinger, O. Belo, and C. Francalanci, editors, *Proceedings of the 3rd International Conference on Data Management Technologies and Applications, DATA 2014, Vienna, Austria*, pages 251–257. SCITEPRESS, August 29–31 2014.
- [13] Niko Beerenwinkel, Huldrych F. Günthard, Volker Roth, and Karin J. Metzner. Challenges and opportunities in estimating viral genetic diversity from next-generation sequencing data. *Virology*, 3:329, 2012.
- [14] Timo Beller and Enno Ohlebusch. Efficient Construction of a Compressed de Bruijn Graph for Pan-Genome Analysis. In Ferdinando Cicalese, Ely Porat, and Ugo Vaccaro, editors, *Combinatorial Pattern Matching*, number 9133 in Lecture Notes in Computer Science, pages 40–51. Springer International Publishing, June 2015.
- [15] Emily Berger, Deniz Yorukoglu, Jiang Peng, and Bonnie Berger. Haptree: a novel Bayesian framework for single individual polyplotyping using NGS data. *PLoS Computational Biology*, 10(3):e1003502, 2014.
- [16] Brigitte Boeckmann, Marina Marcet-Houben, Jonathan A. Rees, Kristoffer Forslund, Jaime Huerta-Cepas, Matthieu Muffato, Pelin Yilmaz, Ioannis Xenarios, Peer Bork, Suzanna E. Lewis, Toni Gabaldón, and the Quest for Orthologs Species Tree Working Group. Quest for Orthologs Entails Quest for Tree of Life: In Search of the Gene Stream. *Genome Biology and Evolution*, 7(7):1988–1999, July 2015.
- [17] Sharon R. Browning and Brian L. Browning. Haplotype phasing: existing methods and new developments. *Nature Reviews Genetics*, 12(10):703–714, October 2011.
- [18] Jennifer R. Brum, J. Cesar Ignacio-Espinoza, Simon Roux, Guilhem Doulier, Silvia G. Acinas, Adriana Alberti, Samuel Chaffron, Corinne Cruaud, Colomban de Vargas, Josep M. Gasol, Gabriel Gorsky, Ann C. Gregory, Lionel Guidi, Pascal Hingamp, Daniele Iudicone, Fabrice Not, Hiroyuki Ogata, Stéphane Pesant, Bonnie T. Poulos, Sarah M. Schwenck, Sabrina Speich, Celine Dimier, Stefanie Kandels-Lewis, Marc Picheral, Sarah Searson, Tara Oceans Coordinators, Peer Bork, Chris Bowler, Shinichi Sunagawa, Patrick Wincker, Eric Karsenti, and Matthew B. Sullivan. Patterns and ecological drivers of ocean viral communities. *Science*, 348(6237):1261498, May 2015.
- [19] Joshua N. Burton, Andrew Adey, Rupali P. Patwardhan, Ruolan Qiu, Jacob O. Kitzman, and Jay Shendure. Chromosome-scale scaffolding of *de novo* genome assemblies based on chromatin interactions. *Nature Biotechnology*, 31(12):1119–1125, 2013.
- [20] Jonathan M. Carlson, Chanson J. Brumme, Eric Martin, Jennifer Listgarten, Mark A. Brockman, Anh Q. Le, Celia K. S. Chui, Laura A. Cotton, David J. H. F. Knapp, Sharon A. Riddler, Richard Haubrich, George Nelson, Nico Pfeifer, Charles E. DeZiel, David Heckerman, Richard Apps, Mary Carrington, Simon Mallal, P. Richard Harrigan, Mina John, Zabrina L. Brumme, and the International HIV Adaptation Collaborative. Correlates of Protective Cellular Immunity

- Revealed by Analysis of Population-Level Immune Escape Pathways in HIV-1. *Journal of Virology*, 86(24):13202–13216, December 2012.
- [21] Mark J. P. Chaisson, John Huddleston, Megan Y. Dennis, Peter H. Sudmant, Maika Malig, Fereydoon Hormozdiari, Francesca Antonacci, Urvashi Surti, Richard Sandstrom, Matthew Boitano, Jane M. Landolin, John A. Stamatoyannopoulos, Michael W. Hunkapiller, Jonas Korlach, and Evan E. Eichler. Resolving the complexity of the human genome using single-molecule sequencing. *Nature*, 517(7536):608–611, January 2015.
- [22] Rayan Chikhi, Antoine Limasset, Shaun Jackman, Jared T. Simpson, and Paul Medvedev. On the Representation of de Bruijn Graphs. In Roded Sharan, editor, *Research in Computational Molecular Biology*, volume 8394 of *Lecture Notes in Computer Science*, pages 35–55. Springer International Publishing, April 2014.
- [23] Francesca D. Ciccarelli, Tobias Doerks, Christian von Mering, Christopher J. Creevey, Berend Snel, and Peer Bork. Toward automatic reconstruction of a highly resolved tree of life. *Science (New York, N.Y.)*, 311(5765):1283–1287, March 2006.
- [24] Colin S. Cooper, Rosalind Eeles, David C. Wedge, Peter Van Loo, Gunes Gundem, Ludmil B. Alexandrov, Barbara Kremeyer, Adam Butler, Andrew G. Lynch, Niedzica Camacho, Charlie E. Massie, Jonathan Kay, Hayley J. Luxton, Sandra Edwards, Zsofia Kote-Jarai, Nening Dennis, Sue Merson, Daniel Leongamornlert, Jorge Zamora, Cathy Corbishley, Sarah Thomas, Serena Nik-Zainal, Manasa Ramakrishna, Sarah O’Meara, Lucy Matthews, Jeremy Clark, Rachel Hurst, Richard Mithen, Robert G. Bristow, Paul C. Boutros, Michael Fraser, Susanna Cooke, Keiran Raine, David Jones, Andrew Menzies, Lucy Stebbings, Jon Hinton, Jon Teague, Stuart McLaren, Laura Mudie, Claire Hardy, Elizabeth Anderson, Olivia Joseph, Victoria Goody, Ben Robinson, Mark Maddison, Stephen Gamble, Christopher Greenman, Dan Berney, Steven Hazell, Naomi Livni, the ICGC Prostate Group, Cyril Fisher, Christopher Ogden, Pardeep Kumar, Alan Thompson, Christopher Woodhouse, David Nicol, Erik Mayer, Tim Dudderidge, Nimish C. Shah, Vincent Gnanapragasam, Thierry Voet, Peter Campbell, Andrew Futreal, Douglas Easton, Anne Y. Warren, Christopher S. Foster, Michael R. Stratton, Hayley C. Whitaker, Ultan McDermott, Daniel S. Brewer, and David E. Neal. Analysis of the genetic phylogeny of multifocal prostate cancer identifies multiple independent clonal expansions in neoplastic and morphologically normal prostate tissue. *Nature Genetics*, 47(4):367–372, April 2015.
- [25] Antoine Coulon, Carson C. Chow, Robert H. Singer, and Daniel R. Larson. Eukaryotic transcriptional dynamics: from single molecules to cell populations. *Nature Reviews Genetics*, 14(8):572–584, August 2013.
- [26] Alastair Crisp, Chiara Boschetti, Malcolm Perry, Alan Tunnacliffe, and Gos Micklem. Expression of multiple horizontally acquired genes is a hallmark of both vertebrate and invertebrate genomes. *Genome Biology*, 16(1):50, March 2015.
- [27] Agnieszka Danek, Sebastian Deorowicz, and Szymon Grabowski. Indexes of Large Genome Collections on a PC. *PLoS ONE*, 9(10):e109384, October 2014.
- [28] Aaron C. E. Darling, Bob Mau, Frederick R. Blattner, and Nicole T. Perna. Mauve: Multiple Alignment of Conserved Genomic Sequence With Rearrangements. *Genome Research*, 14(7):1394–1403, July 2004.
- [29] Modan Das and Ho K. Dai. A survey of DNA motif finding algorithms. *BMC Bioinformatics*, 8(Supplement 7):S21, 2007.
- [30] Matthew D. Daugherty and Harmit S. Malik. Rules of Engagement: Molecular Insights from Host-Virus Arms Races. *Annual Review of Genetics*, 46(1):677–700, December 2012.
- [31] Mark de Been, Val F. Lanza, María de Toro, Jelle Scharringa, Wietske Dohmen, Yu Du,

- Juan Hu, Ying Lei, Ning Li, Ave Tooming-Klunderud, Dick J. J. Heederik, Ad C. Fluit, Marc J. M. Bonten, Rob J. L. Willems, Fernando de la Cruz, and Willem van Schaik. Dissemination of Cephalosporin Resistance Genes between *Escherichia coli* Strains from Farm Animals and Humans by Specific Plasmid Lineages. *PLoS Genet*, 10(12):e1004776, December 2014.
- [32] Sarah De Val, Neil C. Chi, Stryder M. Meadows, Simon Minovitsky, Joshua P. Anderson, Ian S. Harris, Melissa L. Ehlers, Pooja Agarwal, Axel Visel, Shan-Mei Xu, Len A. Pennacchio, Inna Dubchak, Paul A. Krieg, Didier Y. R. Stainier, and Brian L. Black. Combinatorial Regulation of Endothelial Gene Expression by Ets and Forkhead Transcription Factors. *Cell*, 135(6):1053–1064, December 2008.
- [33] Sebastian Deorowicz, Marek Kokot, Szymon Grabowski, and Agnieszka Debudaj-Grabysz. KMC 2: fast and resource-frugal k-mer counting. *Bioinformatics*, 31(10):1569–1576, May 2015.
- [34] Alexander Dilthey, Charles Cox, Zamin Iqbal, Matthew R. Nelson, and Gil McVean. Improved genome inference in the MHC using a population reference graph. *Nature Genetics*, 47(6):682–688, June 2015.
- [35] Esteban Domingo. Quasispecies Theory in Virology. *Journal of Virology*, 76(1):463–465, January 2002.
- [36] W. Ford Doolittle. Phylogenetic Classification and the Universal Tree. *Science*, 284(5423):2124–2128, June 1999.
- [37] Richard Durbin. Efficient haplotype matching and storage using the positional Burrows–Wheeler transform (PBWT). *Bioinformatics*, 30(9):1266–1272, May 2014.
- [38] B. E. Dutilh, V. van Noort, R. T. J. M. van der Heijden, T. Boekhout, B. Snel, and M. A. Huynen. Assessment of phylogenomic and orthology approaches for phylogenetic inference. *Bioinformatics (Oxford, England)*, 23(7):815–824, April 2007.
- [39] Bas E. Dutilh, Lennart Backus, Robert A. Edwards, Michiel Wels, Jumamurat R. Bayjanov, and Sacha A. F. T. van Hijum. Explaining microbial phenotypes on a genomic scale: GWAS for microbes. *Briefings in Functional Genomics*, 12(4):366–380, July 2013.
- [40] Bas E. Dutilh, Martijn A. Huynen, and Marc Strous. Increasing the coverage of a metapopulation consensus genome by iterative read mapping and assembly. *Bioinformatics*, 25(21):2878–2881, November 2009.
- [41] Bas E. Dutilh, Cristiane C. Thompson, Ana C. P. Vicente, Michel A. Marin, Clarence Lee, Genivaldo G. Z. Silva, Robert Schmieder, Bruno G. N. Andrade, Luciane Chimetto, Daniel Cuevas, Daniel R. Garza, Iruka N. Okeke, Aaron Oladipo Aboderin, Jessica Spangler, Tristen Ross, Elizabeth A. Dinsdale, Fabiano L. Thompson, Timothy T. Harkins, and Robert A. Edwards. Comparative genomics of 274 *Vibrio cholerae* genomes reveals mobile functions structuring three niche dimensions. *BMC genomics*, 15(1):654, August 2014.
- [42] Robert C Edgar and Serafim Batzoglou. Multiple sequence alignment. *Current Opinion in Structural Biology*, 16(3):368–373, June 2006.
- [43] Robert A. Edwards and Forest Rohwer. Viral metagenomics. *Nature Reviews Microbiology*, 3(6):504–510, 2005.
- [44] Mark Eppinger, Talima Pearson, Sara S. K. Koenig, Ofori Pearson, Nathan Hicks, Sonia Agrawal, Fatemeh Sanjar, Kevin Galens, Sean Daugherty, Jonathan Crabtree, Rene S. Hendriksen, Lance B. Price, Bishnu P. Upadhyay, Geeta Shakya, Claire M. Fraser, Jacques Ravel, and Paul S. Keim. Genomic Epidemiology of the Haitian Cholera Outbreak: a Single Introduction Followed by Rapid, Extensive, and Continued Spread Characterized the Onset of the Epidemic. *mBio*, 5(6):e01721–14, December 2014.
- [45] Corinna Ernst and Sven Rahmann. Pancake: A Data Structure for Pangenomes. In Tim Beißbarth, Martin Kollmar, Andreas Lehmann, Burkhard Morgenstern, Anne-Kathrin

- Schultz, Stephan Waack, and Edgar Wingen-der, editors, *German Conference on Bioinformatics 2013*, volume 34 of *OpenAccess Series in Informatics (OASIS)*, pages 35–45, Dagstuhl, Germany, 2013. Schloss Dagstuhl–Leibniz-Zentrum fuer Informatik.
- [46] Exome Aggregation Consortium, Monkol Lek, Konrad Karczewski, Eric Minikel, Kaitlin Samocha, Eric Banks, Timothy Fennell, Anne O'Donnell-Luria, James Ware, Andrew Hill, Beryl Cummings, Taru Tukiainen, Daniel Birnbaum, Jack Kosmicki, Laramie Duncan, Karol Estrada, Fengmei Zhao, James Zou, Emma Pierce-Hoffman, David Cooper, Mark DePristo, Ron Do, Jason Flannick, Menachem Fromer, Laura Gauthier, Jackie Goldstein, Namrata Gupta, Daniel Howrigan, Adam Kiezun, Mitja Kurki, Ami Levy Moonshine, Pradeep Natarajan, Lorena Orozco, Gina Peloso, Ryan Poplin, Manuel Rivas, Valentin Ruano-Rubio, Douglas Ruderfer, Khalid Shakir, Peter Stenson, Christine Stevens, Brett Thomas, Grace Tiao, Maria Tusie-Luna, Ben Weisburd, Hong-Hee Won, Dongmei Yu, David Altshuler, Diego Ardisino, Michael Boehnke, John Danesh, Elosua Roberto, Jose Florez, Stacey Gabriel, Gad Getz, Christina Hultman, Sekar Kathiresan, Markku Laakso, Steven McCarroll, Mark McCarthy, Dermot McGovern, Ruth McPherson, Benjamin Neale, Aarno Palotie, Shaun Purcell, Danish Saleheen, Jeremiah Scharf, Pamela Sklar, Sullivan Patrick, Jaakko Tuomilehto, Hugh Watkins, James Wilson, Mark Daly, and Daniel MacArthur. Analysis of protein-coding genetic variation in 60,706 humans. *bioRxiv*, 2015.
- [47] Ester Falconer, Mark Hills, Ulrike Naumann, Steven S. S. Poon, Elizabeth A. Chavez, Ashley D. Sanders, Yongjun Zhao, Martin Hirst, and Peter M. Lansdorp. DNA template strand sequencing of single-cells maps genomic rearrangements at high resolution. *Nature Methods*, 9(11):1107–1112, November 2012.
- [48] Travis Gagie and Simon J. Puglisi. Searching and indexing genomic databases via kernelization. *Bioinformatics and Computational Biology*, 3:12, 2015.
- [49] Genome of the Netherlands Consortium. Whole-genome sequence variation, population structure and demographic history of the Dutch population. *Nature Genetics*, 46(8):818–825, August 2014.
- [50] Christian Gilissen, Jayne Y. Hehir-Kwa, Djie Tjwan Thung, Maartje van de Vorst, Bregje W. M. van Bon, Marjolein H. Willemssen, Michael Kwint, Irene M. Janssen, Alexander Hoischen, Annette Schenck, Richard Leach, Robert Klein, Rick Tearle, Tan Bo, Rolph Pfundt, Helger G. Yntema, Bert B. A. de Vries, Tjitske Kleefstra, Han G. Brunner, Lisenka E. L. M. Vissers, and Joris A. Veltman. Genome sequencing identifies major causes of severe intellectual disability. *Nature*, 511(7509):344–347, July 2014.
- [51] Gustavo Glusman, Hannah C Cox, and Jared C Roach. Whole-genome haplotyping approaches and genomic medicine. *Genome Medicine*, 6(9):73, September 2014.
- [52] Chris D. Greenman, Erin D. Pleasance, Scott Newman, Fengtang Yang, Beiyuan Fu, Serena Nik-Zainal, David Jones, King Wai Lau, Nigel Carter, Paul A. W. Edwards, P. Andrew Futreal, Michael R. Stratton, and Peter J. Campbell. Estimation of rearrangement phylogeny for cancer genomes. *Genome Research*, 22(2):346–361, February 2012.
- [53] Richard J. Hall, Jenny L. Draper, Fiona G. G. Nielsen, and Bas E. Dutilh. Beyond research: a primer for considerations on using viral metagenomics in the field and clinic. *Frontiers in Microbiology*, 6(224), March 2015.
- [54] Alex R. Hastie, Lingli Dong, Alexis Smith, Jeff Finklestein, Ernest T. Lam, Naxin Huo, Han Cao, Pui-Yan Kwok, Karin R. Deal, Jan Dvorak, Ming-Cheng Luo, Yong Gu, and Ming Xiao. Rapid genome mapping in nanochannel arrays for highly complete and accurate *de novo* sequence assembly of the

- complex *aegilops tauschii* genome. *PLoS One*, 8(2):e55864, 2013.
- [55] André Hennig, Jörg Bernhardt, and Kay Nieselt. Pan-Tetris: an interactive visualisation for Pan-genomes. *BMC Bioinformatics*, 16(Suppl 11):S3, August 2015.
- [56] A. Herbig, G. Jäger, F. Battke, and K. Nieselt. GenomeRing: alignment visualization based on SuperGenome coordinates. *Bioinformatics*, 28(12):i7–i15, June 2012.
- [57] Matthew T. G. Holden, Li-Yang Hsu, Kevin Kurt, Lucy A. Weinert, Alison E. Mather, Simon R. Harris, Birgit Strommenger, Franziska Layer, Wolfgang Witte, Herminia de Lencastre, Robert Skov, Henrik Westh, Helena Žemličková, Geoffrey Coombs, Angela M. Kearns, Robert L. R. Hill, Jonathan Edgeworth, Ian Gould, Vanya Gant, Jonathan Cooke, Giles F. Edwards, Paul R. McAdam, Kate E. Templeton, Angela McCann, Zhemin Zhou, Santiago Castillo-Ramírez, Edward J. Feil, Lyndsey O. Hudson, Mark C. Enright, Francois Balloux, David M. Aanensen, Brian G. Spratt, J. Ross Fitzgerald, Julian Parkhill, Mark Achtman, Stephen D. Bentley, and Ulrich Nübel. A genomic portrait of the emergence, evolution, and global spread of a methicillin-resistant *Staphylococcus aureus* pandemic. *Genome Research*, 23(4):653–664, April 2013.
- [58] Guillaume Holley, Roland Wittler, and Jens Stoye. Bloom filter trie - a data structure for pan-genome storage. In *Proceedings of WABI*, volume 9289 of *LNBI*, pages 217–230, 2015.
- [59] Adina Chuang Howe, Janet K. Jansson, Stephanie A. Malfatti, Susannah G. Tringe, James M. Tiedje, and C. Titus Brown. Tackling soil diversity with the assembly of large, complex metagenomes. *Proceedings of the National Academy of Sciences*, 111(13):4904–4909, 2014.
- [60] Lin Huang, Victoria Popic, and Serafim Batzoglou. Short read alignment with populations of genomes. *Bioinformatics*, 29(13):i361–i370, July 2013.
- [61] Xuehui Huang, Nori Kurata, Xinghua Wei, Zi-Xuan Wang, Ahong Wang, Qiang Zhao, Yan Zhao, Kunyan Liu, Hengyun Lu, Wenjun Li, Yunli Guo, Yiqi Lu, Congcong Zhou, Danlin Fan, Qijun Weng, Chuanrang Zhu, Tao Huang, Lei Zhang, Yongchun Wang, Lei Feng, Hiroyasu Furuumi, Takahiko Kubo, Toshie Miyabayashi, Xiaoping Yuan, Qun Xu, Guojun Dong, Qilin Zhan, Canyang Li, Asao Fujiyama, Atsushi Toyoda, Tingting Lu, Qi Feng, Qian Qian, Jiayang Li, and Bin Han. A map of rice genome variation reveals the origin of cultivated rice. *Nature*, 490(7421):497–501, October 2012.
- [62] Daniel H. Huson and Celine Scornavacca. A Survey of Combinatorial Methods for Phylogenetic Networks. *Genome Biology and Evolution*, 3:23–35, January 2011.
- [63] Zamin Iqbal, Mario Caccamo, Isaac Turner, Paul Flicek, and Gil McVean. De novo assembly and genotyping of variants using colored de Bruijn graphs. *Nature Genetics*, 44(2):226–232, February 2012.
- [64] Yinping Jiao, Hainan Zhao, Longhui Ren, Weibin Song, Biao Zeng, Jinjie Guo, Baobao Wang, Zhipeng Liu, Jing Chen, Wei Li, Mei Zhang, Shaojun Xie, and Jinsheng Lai. Genome-wide genetic changes during modern breeding of maize. *Nature Genetics*, 44(7):812–815, July 2012.
- [65] Paweł Kafarski. Rainbow code of biotechnology. *Chemik*, 66(8):811–816, 2012.
- [66] Cyriac Kandoth, Michael D. McLellan, Fabio Vandin, Kai Ye, Beifang Niu, Charles Lu, Mingchao Xie, Qunyu Zhang, Joshua F. McMichael, Matthew A. Wyczalkowski, Mark D. M. Leiserson, Christopher A. Miller, John S. Welch, Matthew J. Walter, Michael C. Wendl, Timothy J. Ley, Richard K. Wilson, Benjamin J. Raphael, and Li Ding. Mutational landscape and significance across 12 major cancer types. *Nature*, 502(7471):333–339, October 2013.
- [67] Pınar Kavak, Bayram Yüksel, Soner Aksu, M. Oguzhan Kulekci, Tunga Güngör, Faraz Hach, S. Cenk Şahinalp, Turkish Human

- Genome Project, Can Alkan, and Mahmut Şamil Sağiroğlu. Robustness of massively parallel sequencing platforms. *PLoS ONE*, 10(9):e0138259, September 2015.
- [68] Birte Kehr, Kathrin Trappe, Manuel Holtgrewe, and Knut Reinert. Genome alignment with graph data structures: a comparison. *BMC Bioinformatics*, 15(1):99, April 2014.
- [69] Volodymyr Kuleshov. Probabilistic single-individual haplotyping. *Bioinformatics*, 30(17):i379–i385, September 2014.
- [70] Volodymyr Kuleshov, Dan Xie, Rui Chen, Dmitry Pushkarev, Zhihai Ma, Tim Blauwkamp, Michael Kertesz, and Michael Snyder. Whole-genome haplotyping using long reads and statistical methods. *Nature Biotechnology*, 32(3):261–266, March 2014.
- [71] T. Laver, J. Harrison, P. A. O’Neill, K. Moore, A. Farbos, K. Paszkiewicz, and D. J. Studholme. Assessing the performance of the Oxford Nanopore Technologies MinION. *Biomolecular Detection and Quantification*, 3:1–8, March 2015.
- [72] Michael S. Lawrence, Petar Stojanov, Paz Polak, Gregory V. Kryukov, Kristian Cibulskis, Andrey Sivachenko, Scott L. Carter, Chip Stewart, Craig H. Mermel, Steven A. Roberts, Adam Kiezun, Peter S. Hammerman, Aaron McKenna, Yotam Drier, Lihua Zou, Alex H. Ramos, Trevor J. Pugh, Nicolas Stransky, Elena Helman, Jaegil Kim, Carrie Sougnez, Lauren Ambrogio, Elizabeth Nickerson, Erica Shefler, Maria L. Cortés, Daniel Auclair, Gordon Saksena, Douglas Voet, Michael Noble, Daniel DiCara, Pei Lin, Lee Lichtenstein, David I. Heiman, Timothy Fennell, Marcin Imielinski, Bryan Hernandez, Eran Hodis, Sylvan Baca, Austin M. Dulak, Jens Lohr, Dan-Avi Landau, Catherine J. Wu, Jorge Melendez-Zajgla, Alfredo Hidalgo-Miranda, Amnon Koren, Steven A. McCarroll, Jaume Mora, Ryan S. Lee, Brian Crompton, Robert Onofrio, Melissa Parkin, Wendy Winckler, Kristin Ardlie, Stacey B. Gabriel, Charles W. M. Roberts, Jaclyn A. Biegel, Kimberly Stegmaier, Adam J. Bass, Levi A. Garraway, Matthew Meyerson, Todd R. Golub, Dmitry A. Gordenin, Shamil Sunyaev, Eric S. Lander, and Gad Getz. Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature*, 499(7457):214–218, July 2013.
- [73] Ryan M. Layer, Neil Kindlon, Konrad J. Karczewski, Exome Aggregation Consortium, and Aaron R. Quinlan. Efficient genotype compression and analysis of large genetic-variation data sets. *Nature Methods*, advance online publication, November 2015.
- [74] Thomas Lengauer and Tobias Sing. Bioinformatics-assisted anti-HIV therapy. *Nature Reviews Microbiology*, 4(10):790–797, October 2006.
- [75] Heng Li and Nils Homer. A survey of sequence alignment algorithms for next-generation sequencing. *Briefings in Bioinformatics*, 11(5):473–483, September 2010.
- [76] Junhua Li, Huijue Jia, Xianghang Cai, Huanzi Zhong, Qiang Feng, Shinichi Sunagawa, Manimozhiyan Arumugam, Jens Roat Kultima, Edi Prifti, Trine Nielsen, Agnieszka Sierakowska Juncker, Chaysavanh Manichanh, Bing Chen, Wenwei Zhang, Florence Levenez, Juan Wang, Xun Xu, Liang Xiao, Suisha Liang, Dongya Zhang, Zhaoxi Zhang, Weineng Chen, Hailong Zhao, Jumanana Yousuf Al-Aama, Sherif Edris, Huanming Yang, Jian Wang, Torben Hansen, Henrik Bjørn Nielsen, Søren Brunak, Karsten Kristiansen, Francisco Guarner, Oluf Pedersen, Joel Doré, S. Dusko Ehrlich, MetaHIT Consortium, Peer Bork, and Jun Wang. An integrated catalog of reference genes in the human gut microbiome. *Nature Biotechnology*, 32(8):834–841, August 2014.
- [77] Na Li and Matthew Stephens. Modeling Linkage Disequilibrium and Identifying Recombination Hotspots Using Single-Nucleotide Polymorphism Data. *Genetics*, 165(4):2213–2233, December 2003.
- [78] Antoine Limasset, Bastien Cazaux, Eric Rivals, and Pierre Peterlongo. Read Mapping on de Bruijn graph. *arXiv*, 1505.04911, 2015.

- [79] Gianni Liti, David M. Carter, Alan M. Moses, Jonas Warringer, Leopold Parts, Stephen A. James, Robert P. Davey, Ian N. Roberts, Austin Burt, Vassiliki Koufopanou, Isheng J. Tsai, Casey M. Bergman, Douda Bensasson, Michael J. T. O’Kelly, Alexander van Oudenaarden, David B. H. Barton, Elizabeth Bailes, Alex N. Nguyen, Matthew Jones, Michael A. Quail, Ian Goodhead, Sarah Sims, Frances Smith, Anders Blomberg, Richard Durbin, and Edward J. Louis. Population genomics of domestic and wild yeasts. *Nature*, 458(7236):337–341, March 2009.
- [80] Nicholas J. Loman, Joshua Quick, and Jared T. Simpson. A complete bacterial genome assembled de novo using only nanopore sequencing data. *Nature Methods*, 12(8):733–735, August 2015.
- [81] Loman NJ, Constantinidou C, Christner M, and et al. A culture-independent sequence-based metagenomics approach to the investigation of an outbreak of shiga-toxigenic *escherichia coli* o104:h4. *JAMA*, 309(14):1502–1510, 2013.
- [82] D. G. MacArthur, T. A. Manolio, D. P. Dimmock, H. L. Rehm, J. Shendure, G. R. Abecasis, D. R. Adams, R. B. Altman, S. E. Antonarakis, E. A. Ashley, J. C. Barrett, L. G. Biesecker, D. F. Conrad, G. M. Cooper, N. J. Cox, M. J. Daly, M. B. Gerstein, D. B. Goldstein, J. N. Hirschhorn, S. M. Leal, L. A. Pennacchio, J. A. Stamatoyannopoulos, S. R. Sunyaev, D. Valle, B. F. Voight, W. Winckler, and C. Gunter. Guidelines for investigating causality of sequence variants in human disease. *Nature*, 508(7497):469–476, April 2014.
- [83] Emma S. Mace, Shuaishuai Tai, Edward K. Gilding, Yanhong Li, Peter J. Prentis, Lianle Bian, Bradley C. Campbell, Wushu Hu, David J. Innes, Xuelian Han, Alan Cruickshank, Changming Dai, Céline Frère, Haikuan Zhang, Colleen H. Hunt, Xianyuan Wang, Tracey Shatte, Miao Wang, Zhe Su, Jun Li, Xiaozhen Lin, Ian D. Godwin, David R. Jordan, and Jun Wang. Whole-genome sequencing reveals untapped genetic potential in Africa’s indigenous cereal crop sorghum. *Nature Communications*, 4:2320, July 2013.
- [84] Mohammed-Amin Madoui, Stefan Engelen, Corinne Cruaud, Caroline Belser, Laurie Bertrand, Adriana Alberti, Arnaud Lemainque, Patrick Wincker, and Jean-Marc Aury. Genome assembly using Nanopore-guided long and error-free DNA reads. *BMC Genomics*, 16(1):327, April 2015.
- [85] Angel C. Y. Mak, Yvonne Y. Y. Lai, Ernest T. Lam, Tsz-Piu Kwok, Alden K. Y. Leung, Annie Poon, Yulia Mostovoy, Alex R. Hastie, William Stedman, Thomas Anantharaman, Warren Andrews, Xiang Zhou, Andy W. C. Pang, Heng Dai, Catherine Chu, Chin Lin, Jacob J. K. Wu, Catherine M. L. Li, Jing-Woei Li, Aldrin K. Y. Yim, Saki Chan, Justin Sibert, Željko Džakula, Han Cao, Siu-Ming Yiu, Ting-Fung Chan, Kevin Y. Yip, Ming Xiao, and Pui-Yan Kwok. Genome-wide structural variation detection by genome mapping on nanochannel arrays. *Genetics*, 202(1):351–362, January 2016.
- [86] Veli Mäkinen, Gonzalo Navarro, Jouni Sirén, and Niko Välimäki. Storage and Retrieval of Individual Genomes. In Serafim Batzoglou, editor, *Research in Computational Molecular Biology*, number 5541 in Lecture Notes in Computer Science, pages 121–137. Springer Berlin Heidelberg, 2009.
- [87] Veli Mäkinen, Gonzalo Navarro, Jouni Sirén, and Niko Välimäki. Storage and retrieval of highly repetitive sequence collections. *Journal of Computational Biology: A Journal of Computational Molecular Cell Biology*, 17(3):281–308, March 2010.
- [88] Christine M. Malboeuf, Xiao Yang, Patrick Charlebois, James Qu, Aaron M. Berlin, Monica Casali, Kendra N. Pesko, Christian L. Boutwell, John P. DeVincenzo, Gregory D. Ebel, Todd M. Allen, Michael C. Zody, Matthew R. Henn, and Joshua Z. Levin. Complete viral RNA genome sequencing of ultra-low copy samples by sequence-

- p>independent amplification.
- Nucleic Acids Research*
- , 41(1):e13–e13, January 2013.
- [89] Shoshana Marcus, Hayan Lee, and Michael C. Schatz. SplitMEM: a graphical algorithm for pan-genome analysis with suffix skips. *Bioinformatics*, 30(24):3476–3483, December 2014.
 - [90] Andriy Marusyk, Vanessa Almendro, and Kornelia Polyak. Intra-tumour heterogeneity: a looking glass for cancer? *Nature Reviews. Cancer*, 12(5):323–334, May 2012.
 - [91] Nicholas McGranahan and Charles Swanton. Biological and Therapeutic Impact of Intra-tumor Heterogeneity in Cancer Evolution. *Cancer Cell*, 27(1):15–26, January 2015.
 - [92] Páll Melsted and Jonathan K. Pritchard. Efficient counting of k-mers in DNA sequences using a bloom filter. *BMC Bioinformatics*, 12(1):333, August 2011.
 - [93] Giulia Menconi, Giovanni Battaglia, Roberto Grossi, Nadia Pisanti, and Roberto Marangoni. Mobilomics in *saccharomyces cerevisiae* strains. *BMC Bioinformatics*, 14:102, 2013.
 - [94] Jason R. Miller, Sergey Koren, and Granger Sutton. Assembly Algorithms for Next-Generation Sequencing Data. *Genomics*, 95(6):315–327, June 2010.
 - [95] Ilya Minkin, Anand Patel, Mikhail Kolmogorov, Nikolay Vyahhi, and Son Pham. Sibelia: A Scalable and Comprehensive Synteny Block Generation Tool for Closely Related Microbial Genomes. In Aaron Darling and Jens Stoye, editors, *Algorithms in Bioinformatics*, number 8126 in Lecture Notes in Computer Science, pages 215–229. Springer Berlin Heidelberg, 2013.
 - [96] John L Mokili, Forest Rohwer, and Bas E Dutilh. Metagenomics and future perspectives in virus discovery. *Current Opinion in Virology*, 2(1):63–77, 2012.
 - [97] Leonid L. Moroz, Kevin M. Kocot, Mathew R. Citarella, Sohn Dosung, Tigran P. Norekian, Inna S. Povolotskaya, Anastasia P. Grigorenko, Christopher Dailey, Eugene Berezikov, Katherine M. Buckley, Andrey Ptitsyn, Denis Reshetov, Krishanu Mukherjee, Tatiana P. Moroz, Yelena Bobkova, Fahong Yu, Vladimir V. Kapitonov, Jerzy Jurka, Yuri V. Bobkov, Joshua J. Swore, David O. Girardo, Alexander Fodor, Fedor Gusev, Rachel Sanford, Rebecca Bruders, Ellen Kittler, Claudia E. Mills, Jonathan P. Rast, Romain Derelle, Victor V. Solovyev, Fyodor A. Kondrashov, Billie J. Swalla, Jonathan V. Sweedler, Evgeny I. Rogaev, Kenneth M. Halanych, and Andrea B. Kohn. The ctenophore genome and the evolutionary origins of neural systems. *Nature*, 510(7503):109–114, June 2014.
 - [98] Gonzalo Navarro and Veli Mäkinen. Compressed Full-text Indexes. *ACM Comput. Surv.*, 39(1):61, April 2007.
 - [99] Ngan Nguyen, Glenn Hickey, Brian J. Raney, Joel Armstrong, Hiram Clawson, Ann Zweig, Donna Karolchik, William James Kent, David Haussler, and Benedict Paten. Comparative assembly hubs: Web-accessible browsers for comparative genomics. *Bioinformatics*, 30(23):3293–3301, December 2014.
 - [100] Ngan Nguyen, Glenn Hickey, Daniel R. Zerbino, Brian Raney, Dent Earl, Joel Armstrong, W. James Kent, David Haussler, and Benedict Paten. Building a Pan-Genome Reference for a Population. *Journal of Computational Biology*, 22(5):387–401, January 2015.
 - [101] Cydney B. Nielsen, Michael Cantor, Inna Dubchak, David Gordon, and Ting Wang. Visualizing genomes: techniques and challenges. *Nature Methods*, 7:S5–S15, 2010.
 - [102] Cédric Notredame. Recent Evolutions of Multiple Sequence Alignment Algorithms. *PLoS Comput Biol*, 3(8):e123, August 2007.
 - [103] Dina N Palttoo, Laura Lyman Rodriguez, Michael Feolo, Elizabeth Gillanders, Erin M Ramos, Joni L Rutter, Stephen Sherry, Vivian Ota Wang, Alice Bailey, Rebecca Baker, Mark Caulder, Emily L Harris, Kristofor Langlais, Hilary Leeds, Erin Luetkemeier, Taunton Paine, Tamar Roomian, Kimberly

- Tryka, Amy Patterson, and Eric D Green. Data use under the NIH GWAS Data Sharing Policy and future directions. *Nature Genetics*, 46(9):934–938, August 2014.
- [104] Benedict Paten, Mark Diekhans, Dent Earl, John St. John, Jian Ma, Bernard Suh, and David Haussler. Cactus Graphs for Genome Comparisons. *Journal of Computational Biology*, 18(3):469–481, March 2011.
- [105] Benedict Paten, Dent Earl, Ngan Nguyen, Mark Diekhans, Daniel Zerbino, and David Haussler. Cactus: Algorithms for genome multiple sequence alignment. *Genome Research*, 21(9):1512–1528, September 2011.
- [106] Benedict Paten, Javier Herrero, Kathryn Beal, Stephen Fitzgerald, and Ewan Birney. Enredo and Pecan: Genome-wide mammalian consistency-based multiple alignment with paralogs. *Genome Research*, 18(11):1814–1828, November 2008.
- [107] Murray Patterson, Tobias Marschall, Nadia Pisanti, Leo van Iersel, Leen Stougie, Gunnar W. Klau, and Alexander Schönhuth. WhatsHap: Weighted Haplotype Assembly for Future-Generation Sequencing Reads. *Journal of Computational Biology*, 22(6):498–509, February 2015.
- [108] Pavel A. Pevzner, Haixu Tang, and Glenn Tesler. De Novo Repeat Classification and Fragment Assembly. *Genome Research*, 14(9):1786–1796, September 2004.
- [109] Yuri Pirola, Simone Zaccaria, Riccardo Dondi, Gunnar W. Klau, Nadia Pisanti, and Paola Bonizzoni. HapCol: Accurate and memory-efficient haplotype assembly from long reads. *Bioinformatics*, 2015. Advance access.
- [110] Junjie Qin, Yingrui Li, Zhiming Cai, Shenghui Li, Jianfeng Zhu, Fan Zhang, Suisha Liang, Wenwei Zhang, Yuanlin Guan, Dongqian Shen, Yangqing Peng, Dongya Zhang, Zhuye Jie, Wenxian Wu, Youwen Qin, Wenbin Xue, Junhua Li, Lingchuan Han, Donghui Lu, Peixian Wu, Yali Dai, Xiaojuan Sun, Zesong Li, Aifa Tang, Shilong Zhong, Xiaoping Li, Weineng Chen, Ran Xu, Mingbang Wang, Qiang Feng, Meihua Gong, Jing Yu, Yanyan Zhang, Ming Zhang, Torben Hansen, Gaston Sanchez, Jeroen Raes, Gwen Falony, Shujiro Okuda, Mathieu Almeida, Emmanuelle LeChatelier, Pierre Renault, Nicolas Pons, Jean-Michel Batto, Zhaoxi Zhang, Hua Chen, Ruifu Yang, Weimou Zheng, Songgang Li, Huanming Yang, Jian Wang, S. Dusko Ehrlich, Rasmus Nielsen, Oluf Pedersen, Karsten Kristiansen, and Jun Wang. A metagenome-wide association study of gut microbiota in type 2 diabetes. *Nature*, 490(7418):55–60, 2012.
- [111] René Rahn, David Weese, and Knut Reinert. Journaled String Tree - A scalable data structure for analyzing thousands of similar genomes on your laptop. *Bioinformatics*, 30(24):3499–3505, July 2014.
- [112] Timothy D. Read and Ruth C. Massey. Characterizing the genetic basis of bacterial phenotypes using genome-wide association studies: a new direction for bacteriology. *Genome Medicine*, 6(11):109, November 2014.
- [113] Andy Rimmer, Hang Phan, Iain Mathieson, Zamin Iqbal, Stephen R. F. Twigg, Wgs500 Consortium, Andrew O. M. Wilkie, Gil McVean, and Gerton Lunter. Integrating mapping-, assembly- and haplotype-based approaches for calling variants in clinical sequencing applications. *Nature Genetics*, 46(8):912–918, August 2014.
- [114] Guillaume Rizk, Dominique Lavenier, and Rayan Chikhi. DSK: k-mer counting with very low memory usage. *Bioinformatics*, 29(5):652–653, March 2013.
- [115] Michael G. Ross, Carsten Russ, Maura Costello, Andrew Hollinger, Niall J. Lennon, Ryan Hegarty, Chad Nusbaum, and David B. Jaffe. Characterizing and measuring bias in sequence data. *Genome Biology*, 14(5):R51, May 2013.
- [116] Geir K. Sandve and Finn Drabløs. A survey of motif discovery methods in an integrated framework. *Biology Direct*, 1(1):11, April 2006.

- [117] Korbinian Schneeberger, Jörg Hagmann, Stephan Ossowski, Norman Warthmann, Sandra Gesing, Oliver Kohlbacher, and Detlef Weigel. Simultaneous alignment of short reads against multiple genomes. *Genome Biology*, 10(9):R98, September 2009.
- [118] Grégory F. Schneider and Cees Dekker. DNA sequencing with nanopores. *Nature Biotechnology*, 30(4):326–328, April 2012.
- [119] Paul Shannon, Andrew Markiel, Owen Ozier, Nitin S. Baliga, Jonathan T. Wang, Daniel Ramage, Nada Amin, Benno Schwikowski, and Trey Ideker. Cytoscape: A Software Environment for Integrated Models of Biomolecular Interaction Networks. *Genome Research*, 13(11):2498–2504, November 2003.
- [120] S. T. Sherry, M. H. Ward, M. Kholodov, J. Baker, L. Phan, E. M. Smigielski, and K. Sirotkin. dbSNP: the NCBI database of genetic variation. *Nucleic Acids Research*, 29(1):308–311, January 2001.
- [121] François Sigaux. [Cancer genome or the development of molecular portraits of tumors]. *Bulletin de l'Académie nationale de médecine*, 184(7):1441–1449; discussion 1448–1449, October 2000.
- [122] J. Sirén, N. Välimäki, and V. Mäkinen. Indexing Graphs for Path Queries with Applications in Genome Research. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 11(2):375–388, March 2014.
- [123] Jouni Sirén, Niko Välimäki, and Veli Mäkinen. Indexing Finite Language Representation of Population Genotypes. In Teresa M. Przytycka and Marie-France Sagot, editors, *Algorithms in Bioinformatics*, number 6833 in Lecture Notes in Computer Science, pages 270–281. Springer Berlin Heidelberg, 2011.
- [124] Berend Snel, Martijn A. Huynen, and Bas E. Dutilh. Genome trees and the nature of genome evolution. *Annual Review of Microbiology*, 59:191–209, 2005.
- [125] Matthew W. Snyder, Andrew Adey, Jacob O. Kitzman, and Jay Shendure. Haplotype-resolved genome sequencing: experimental methods and applications. *Nature Reviews Genetics*, 16(6):344–358, June 2015.
- [126] Paweł Stankiewicz and James R. Lupski. Structural variation in the human genome and its role in disease. *Annual Review of Medicine*, 61:437–455, 2010.
- [127] Lincoln D. Stein, Bartha M. Knoppers, Peter Campbell, Gad Getz, and Jan O. Korbel. Data analysis: Create a cloud commons. *Nature*, 523(7559):149–151, July 2015.
- [128] Michael R. Stratton, Peter J. Campbell, and P. Andrew Futreal. The cancer genome. *Nature*, 458(7239):719–724, April 2009.
- [129] Lorenzo Tattini, Romina D’Aurizio, and Alberto Magi. Detection of genomic structural variants from next-generation sequencing data. *Frontiers in Bioengineering and Biotechnology*, 3:92, 2015.
- [130] Brian Teague, Michael S. Waterman, Steven Goldstein, Konstantinos Potamouisis, Shiguo Zhou, Susan Reslewic, Deepayan Sarkar, Anton Valouev, Christopher Churas, Jeffrey M. Kidd, Scott Kohn, Rodney Runnheim, Casey Lamers, Dan Forrest, Michael A. Newton, Evan E. Eichler, Marijo Kent-First, Urvasi Surti, Miron Livny, and David C. Schwartz. High-resolution human genome structure by single-molecule analysis. *Proceedings of the National Academy of Sciences*, 107(24):10848–10853, June 2010.
- [131] Hervé Tettelin, Vega Massignani, Michael J. Cieslewicz, Claudio Donati, Duccio Medini, Naomi L. Ward, Samuel V. Angiuoli, Jonathan Crabtree, Amanda L. Jones, A. Scott Durkin, Robert T. DeBoy, Tanja M. Davidsen, Marirosa Mora, Maria Scarselli, Immaculada Margarit y Ros, Jeremy D. Peterson, Christopher R. Hauser, Jaideep P. Sundaram, William C. Nelson, Ramana Madupu, Lauren M. Brinkac, Robert J. Dodson, Mary J. Rosovitz, Steven A. Sullivan, Sean C. Daugherty, Daniel H. Haft, Jeremy Selengut, Michelle L. Gwinn, Liwei

- Zhou, Nikhat Zafar, Hoda Khouri, Diana Radune, George Dimitrov, Kisha Watkins, Kevin J. B. O'Connor, Shannon Smith, Teresa R. Utterback, Owen White, Craig E. Rubens, Guido Grandi, Lawrence C. Madoff, Dennis L. Kasper, John L. Telford, Michael R. Wessels, Rino Rappuoli, and Claire M. Fraser. Genome analysis of multiple pathogenic isolates of *Streptococcus agalactiae*: Implications for the microbial "pan-genome". *Proceedings of the National Academy of Sciences of the United States of America*, 102(39):13950–13955, September 2005.
- [132] The 100 Tomato Genome Sequencing Consortium, Saulo Aflitos, Elio Schijlen, Hans de Jong, Dick de Ridder, Sandra Smit, Richard Finkers, Jun Wang, Gengyun Zhang, Ning Li, Likai Mao, Freek Bakker, Rob Dirks, Timo Breit, Barbara Gravendeel, Henk Huits, Darush Struss, Ruth Swanson-Wagner, Hans van Leeuwen, Roeland C.H.J. van Ham, Laia Fito, Laëtitia Guignier, Myrna Sevilla, Philippe Ellul, Eric Ganko, Arvind Kapur, Emmanuel Reclus, Bernard de Geus, Henri van de Geest, Bas te Lintel Hekkert, Jan van Haarst, Lars Smits, Andries Koops, Gabino Sanchez-Perez, Adriaan W. van Heusden, Richard Visser, Zhiwu Quan, Jiumeng Min, Li Liao, Xiaoli Wang, Guangbiao Wang, Zhen Yue, Xinhua Yang, Na Xu, Eric Schranz, Erik Smets, Rutger Vos, Johan Rauwerda, Remco Ursem, Cees Schuit, Mike Kerns, Jan van den Berg, Wim Vriezen, Antoine Janssen, Erwin Datema, Torben Jahrman, Frederic Moquet, Julien Bonnet, and Sander Peters. Exploring genetic variation in the tomato (*Solanum section Lycopersicon*) clade by whole-genome sequencing. *The Plant Journal*, 80(1):136–148, October 2014.
- [133] The Human Microbiome Project Consortium. Structure, function and diversity of the healthy human microbiome. *Nature*, 486(7402):207–214, June 2012.
- [134] Armin Töpfer, Tobias Marschall, Rowena A. Bull, Fabio Luciani, Alexander Schönhuth, and Niko Beerenwinkel. Viral Quasispecies Assembly via Maximal Clique Enumeration. *PLoS Comput Biol*, 10(3):e1003515, March 2014.
- [135] George Vernikos, Duccio Medini, David R Riley, and Hervé Tettelin. Ten years of pan-genome analyses. *Current Opinion in Microbiology*, 23:148–154, February 2015.
- [136] Jing Wang, Nicole Elizabeth Moore, Yi-Mo Deng, David A. Eccles, and Richard J. Hall. MinION nanopore sequencing of an influenza genome. *Virology*, 6:766, August 2015.
- [137] Andrew M. Waterhouse, James B. Procter, David M. A. Martin, Michèle Clamp, and Geoffrey J. Barton. Jalview Version 2—a multiple sequence alignment editor and analysis workbench. *Bioinformatics*, 25(9):1189–1191, May 2009.
- [138] Detlef Weigel and Richard Mott. The 1001 genomes project for *Arabidopsis thaliana*. *Genome Biology*, 10(5):107, May 2009.
- [139] Neil I. Weisenfeld, Shuangye Yin, Ted Sharpe, Bayo Lau, Ryan Hegarty, Laurie Holmes, Brian Sogoloff, Diana Tabbaa, Louise Williams, Carsten Russ, Chad Nusbbaum, Eric S. Lander, Iain MacCallum, and David B. Jaffe. Comprehensive variation discovery in single human genomes. *Nature Genetics*, 46(12):1350–1355, December 2014.
- [140] Ryan R. Wick, Mark B. Schultz, Justin Zobel, and Kathryn E. Holt. Bandage: interactive visualization of de novo genome assemblies. *Bioinformatics*, 31(20):3350–3352, 2015.
- [141] Tom A. Williams, Peter G. Foster, Cymon J. Cox, and T. Martin Embley. An archaeal origin of eukaryotes supports only two primary domains of life. *Nature*, 504(7479):231–236, December 2013.
- [142] Shannon J. Williamson, Douglas B. Rusch, Shibu Yooseph, Aaron L. Halpern, Karla B. Heidelberg, John I. Glass, Cynthia Andrews-Pfannkoch, Douglas Fadrosh, Christopher S. Miller, Granger Sutton, Marvin Frazier, and J. Craig Venter. The Sorcerer II Global Ocean Sampling Expedition: Metagenomic

- Characterization of Viruses within Aquatic Microbial Samples. *PLoS ONE*, 3(1):e1456, January 2008.
- [143] Derrick E. Wood and Steven L. Salzberg. Kraken: ultrafast metagenomic sequence classification using exact alignments. *Genome Biology*, 15(3):R46, March 2014.
- [144] Jingfa Xiao, Zhewen Zhang, Jiayan Wu, and Jun Yu. A Brief Review of Software Tools for Pangenomics. *Genomics, Proteomics & Bioinformatics*, 13(1):73–76, February 2015.
- [145] Osvaldo Zagordi, Arnab Bhattacharya, Nicholas Eriksson, and Niko Beerenwinkel. ShoRAH: estimating the genetic diversity of a mixed sample from next-generation sequencing data. *BMC Bioinformatics*, 12(1):119, April 2011.
- [146] Grace X.Y. Zheng, Billy T. Lau, Michael Schnall-Levin, Mirna Jarosz, John M. Bell, Christopher M. Hindson, Sofia Kyriazopoulou-Panagiotopoulou, Donald A. Masquelier, Landon Merrill, Jessica M. Terry, Patrice A. Mudivarti, Paul W. Wyatt, Rajiv Bharadwaj, Anthony J. Makarewicz, Yuan Li, Phillip Belgrader, Andred D. Price, Adam J. Lowe, Patrick Marks, Gerard M. Vurens, Paul Hardenbol, Luz Montesclaros, Melissa Luo, Lawrence Greenfield, Alexander Wong, David E. Birch, Steven W. Short, Keith P. Bjornson, Pranav Patel, Erik S. Hopmans, Christina Wood, Sukhvinder Kaur, Glenn K. Lockwood, David Stafford, Joshua P. Delaney, Indira Wu, Heather S. Ordonez, Susan M. Grimes, Stephanie Greer, Josephine Y. Lee, Kamila Belhocine, Kristina M. Giorda, William H. Heaton, Geoffrey P. McDermott, Zachary W. Bent, Francesca Meschi, Nikola O. Kondov, Ryan Wilson, Jorge A. Bernate, Shawn Gauby, Alex Kindwall, Clara Bermejo, Adrian N. Fehr, Adrian Chan, Serge Saxonov, Kevin D. Ness, Benjamin J. Hindson, and Han-lee P. Ji. Haplotyping germline and cancer genomes with high-throughput linked-read sequencing. *Nature Biotechnology*, AOP, 2016.
- [147] Bojian Zhong, Linhua Sun, and David Penny. The origin of land plants: A phylogenomic perspective. *Evolutionary Bioinformatics Online*, 11:137–141, 2015.