

High Performance Computing for Haplotyping: Models and Platforms

Andrea Tangherloni¹(✉), Leonardo Rundo^{1,5}, Simone Spolaor¹,
Marco S. Nobile^{1,6}, Ivan Merelli², Daniela Besozzi¹, Giancarlo Mauri^{1,6},
Paolo Cazzaniga^{3,6}, and Pietro Liò⁴

¹ Department of Informatics, Systems and Communication,
University of Milano-Bicocca, Milan, Italy
`andrea.tangherloni@disco.unimib.it`,

² Institute of Biomedical Technologies,
Italian National Research Council, Segrate, Italy

³ Department of Human and Social Sciences,
University of Bergamo, Bergamo, Italy

⁴ Computer Laboratory,
University of Cambridge, Cambridge, United Kingdom

⁵ Institute of Molecular Bioimaging and Physiology,
Italian National Research Council, Cefalù (PA), Italy

⁶ SYSBIO.IT Centre of Systems Biology, Milano, Italy

Abstract. The reconstruction of the haplotype pair for each chromosome is a hot topic in Bioinformatics and Genome Analysis. In Haplotype Assembly (HA), all heterozygous Single Nucleotide Polymorphisms (SNPs) have to be assigned to exactly one of the two chromosomes. In this work, we outline the state-of-the-art on HA approaches and present an in-depth analysis of the computational performance of GenHap, a recent method based on Genetic Algorithms. GenHap was designed to tackle the computational complexity of the HA problem by means of a *divide-et-impera* strategy that effectively leverages multi-core architectures. In order to evaluate GenHap's performance, we generated different instances of synthetic (yet realistic) data exploiting empirical error models of four different sequencing platforms (namely, Illumina NovaSeq, Roche/454, PacBio RS II and Oxford Nanopore Technologies MinION). Our results show that the processing time generally decreases along with the read length, involving a lower number of sub-problems to be distributed on multiple cores.

Keywords: Future-generation sequencing · Genome analysis · Haplotype assembly · High Performance Computing · Master-Slave paradigm

1 Introduction

The advent of second-generation sequencing technologies revolutionized the field of genomics, enabling a more complete view and understanding of the genome

of different species. However, despite their great contribution to the field, the data produced by these technologies are still unsuitable for several applications, including Haplotype Assembly (HA). This problem consists in assigning all heterozygous Single Nucleotide Polymorphisms (SNPs) to exactly one of the two homologous chromosomes, leveraging data from sequencing experiments. The short length of the reads produced by second-generation sequencing technologies might be not long enough to span over a relevant number of SNP positions, leading to the reconstruction of short haplotype blocks [1, 2] and ultimately hindering the possibility of reconstructing the full haplotypes.

In recent years, however, a third-generation of sequencing technologies was developed and paved the way to the production of sequencing data characterized by reads covering hundreds of kilobases, thus able to span different SNP loci at once [3–5]. Unfortunately, the increase in length comes at the cost of a decrease in the accuracy of the reads, compared to the short and precise ones produced by second-generation sequencing technologies, such as NovaSeq (Illumina Inc., San Diego, CA, USA) [6]. In order to compensate for this inadequacy, there is a need to increase read coverage. Formally, the coverage of a sequencing experiment is the average number of times that each nucleotide is expected to be covered by a read. This value is given by the following relationship:

$$\text{cov} = (L \cdot N)/G, \quad (1)$$

where cov stands for the coverage, L for the read length, N for the number of reads and G for the length of the haploid region of genome on which the reads are mapped [7]. Eq. (1) shows that longer reads or a higher amount of reads are needed to increase the coverage. In practice, an average coverage higher than $30\times$ is the *de facto* standard for accurate SNP detection [8]. Along with the HA issues, novel challenges—e.g., poliploidity, metagenomics, analysis of cancer cell heterogeneity and chromosomal capture experiments—require sequencing data with a high coverage.

In this paper, we briefly describe the state-of-the-art on haplotype computational tools, providing a taxonomy based on the employed computational techniques. Then, we focus on GenHap [9], an evolutionary method enabled for High Performance Computing (HPC). We show how GenHap performs on data produced by four different sequencing platforms, namely:

- Illumina NovaSeq (Illumina Inc., San Diego, CA, USA) [6]: the most used and widespread platform belonging to the class of second-generation sequencing technologies, able to produce a huge number of short and precise reads (up to 150bp);
- Roche/454 (Roche AG, Basel, Switzerland) [10]: a second-generation sequencing technology able to produce accurate and slightly longer reads than Illumina sequencers (up to 700bp);
- PacBio RS II (Pacific Biosciences of California Inc., Menlo Park, CA, USA) [3, 4]: a third-generation sequencing technology able to produce long reads (up to 30000bp);

- Oxford Nanopore Technologies (ONT) MinION (ONT Ltd., Oxford, United Kingdom) [5, 11, 12]: the latest developed third-generation sequencing technology, able to produce reads that are tens of kilobases long.

The manuscript is structured as follows. Sect. 2 describes and classifies the most used HA approaches, focusing on HPC potential provided by GenHap. The achieved results, in terms of scalability and efficiency on multi-core architectures, are shown and analyzed in Sect. 3. Finally, future directions and possible fruitful connections with other research fields, such as machine learning and security in distributed computing, are mentioned in Sect. 4.

2 HPC in Haplotype Assembly

Current human Whole Genome Sequencing (WGS) approaches do not generally provide phasing information, limiting the identification of clinically-relevant samples, estimation of compound heterozygosity as well as population-level phenomena, including haplotype diversity and Linkage Disequilibrium patterns that could help to resolve migratory patterns and mutation origins [13].

Several computational HA approaches for human genome phasing have been proposed in literature [13]. Most of these methods solve the NP-hard Minimum Error Correction (MEC) problem, which aims at inferring the haplotype pair that yields two disjoint sets of the sequencing reads characterized by the minimum number of SNP values to be corrected [14]. An additional variant of MEC exists, called weighted MEC (wMEC) [15], which takes into account also the information concerning the quality of the reads.

In what follows, we concisely describe the most diffused HA methods and graphically represent them by means of a “phylogenetic tree”-like diagram (Fig. 1). Then, we focus on the functioning of the distributed GenHap implementation on multi-core architectures [9].

2.1 Related Work

Beagle [16] is one of the earliest heuristic approaches based on Hidden Markov Models (HMMs). Considering the genotype information of an individual, *Beagle* finds the most likely haplotype pair among different possible haplotype solutions. It has a quadratic computational complexity with respect to the input data.

SHAPEIT [17] starts from genotyping data from a population and, given the genotype data of an individual, exploits an HMM-based approach to estimate the haplotype pair. The population data are used to apply constraints on the graph, which denotes all possible haplotypes compatible with the input data, in order to determine the haplotype of that individual. At each iteration, *SHAPEIT* has a linear complexity with respect to the number of haplotypes.

Eagle2 [18] is a phasing algorithm that exploits the Burrows-Wheeler transform to encode the information from large external reference panels. It relies on an HMM to explore only the most relevant phase paths among all possible paths. The authors showed that *Eagle* is 20 times faster than *SHAPEIT* [17].

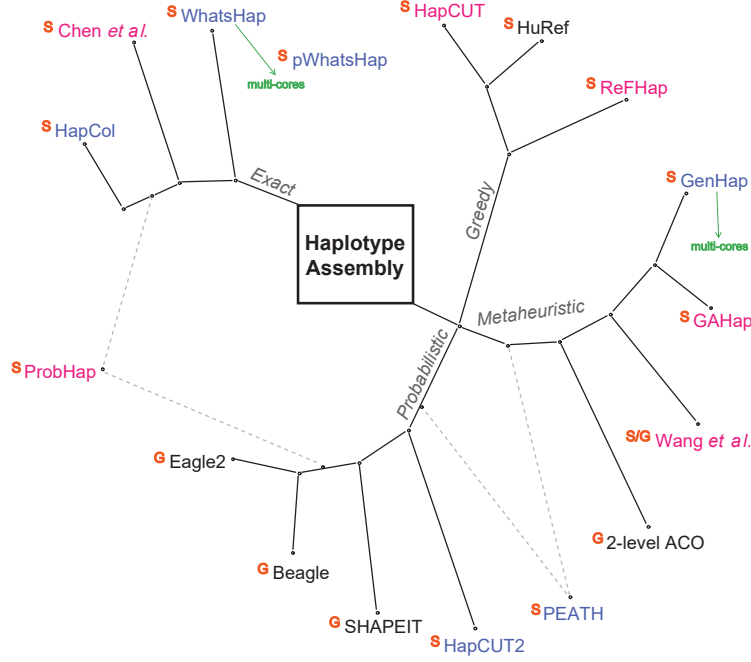


Fig. 1. The “phylogeny” of haplotyping methods. Over the past few years, the repertoire of tools for haplotyping has rapidly expanded. A “phylogenetic tree”-like diagram is used here to depict the division of the algorithms in 4 different classes, namely: exact, greedy, probabilistic, metaheuristic. Hybrid methods are connected with dashed lines to the implemented multiple computational techniques. The orange superscript denotes the analyzed data: sequencing (S) and genotyping (G). Methods that solve either the MEC or the wMEC problem are denoted with blue or magenta, respectively. Finally, the HA methods that exploit HPC are highlighted with a green arrow directed to the used computational resources. See text in Sect. 2.1 for descriptions of the most common software representatives of branches, and for the definitions of abbreviations.

HapCUT [19] leverages sequencing data (i.e., the entire set of reads is considered) instead of population genotypes. It infers the haplotype pair of an individual by partitioning the set of reads solving the MEC problem. The MEC problem is reduced to the max-cut problem, which is greedily solved over the graph representation of the input instance.

HapCUT2 [20] is a recent heuristic approach that exploits a haplotype likelihood model for the sequencing reads. A partial likelihood function is used to evaluate the likelihood of a subset of the fragments. Differently from its previous version (*HapCUT* [19]), which is based on a max-cut algorithm, *HapCUT2* optimizes the likelihood to find a max-cut in graph representation of the input instance.

ProbHap [21] relies on an exact likelihood optimization technique to solve a generalized version of the MEC problem. It exploits a dynamic programming

algorithm capable of exactly optimizing a likelihood function, which is specified by a probabilistic graphical model that generalizes the MEC problem.

ReFHap [22] is based on a heuristic algorithm to find the max-cut. ReFHap solves the Maximum Fragments Cut (MFC) problem instead of the classic MEC problem. The max-cut problem is reduced to the MFC problem, which is addressed using a greedy approach.

HuRef [23] is a heuristic approach that aims at inferring the heterozygous variants of an individual. It is based on a greedy algorithm that iteratively refines the initial partial haplotype solutions. The authors leveraged this HA approach to study non-SNP genetic alterations considering the diploid nature of the human genome.

Chen et al. (2013) [24] proposed an exact approach for the MEC problem using an integer linear programming solver. First, the fragment matrix is decomposed into small independent sub-matrices. Each of these sub-matrices is used to define an integer linear programming problem that is then exactly solved.

WhatsHap [25] is an exact method relying on a dynamic programming algorithm used to solve wMEC. It implements a fixed parameter tractable algorithm, where the fixed parameter is the maximum coverage of the input instance, to deal with the NP-hardness of the wMEC problem. This method does not assume the all-heterozygosity of the phased positions.

pWhatsHap [26] is an efficient version of WhatsHap [25], which was designed to leverage multi-core architectures in order to obtain a relevant reduction of the execution time required by WhatsHap. The proposed implementation exploits the physical shared memory of the underlying architecture to avoid data communication among threads.

HapCol [27] implements a dynamic programming algorithm to solve an alternative version of the wMEC problem, called k -MEC, which is used to take into account the distribution of sequencing errors of future-generation technologies. In this strategy, the number of corrections per column is bounded by the parameter k . No all-heterozygous assumption is required.

Two-Level ACO [28] is based on the Ant Colony Optimization (ACO) technique, which is a metaheuristic designed to deal with combinatorial problems on graphs generated starting from the genotyping data given as input. This approach is based on the pure parsimony criterion to find the smallest set of distinct haplotypes that solves the HA problem.

Probabilistic Evolutionary Algorithm with Toggling for Haplotyping (PEATH) [29] is based on the Estimation of Distribution Algorithm (EDA), which is a metaheuristic suitable for continuous problems. During each iteration of EDA, the promising individuals are used to build probabilistic models that are sampled to explore the search space. This metaheuristic is exploited to deal with noisy sequencing reads, aiming at inferring one haplotype, under the all-heterozygous assumption.

Wang et al. (2005) [14] relies on Genetic Algorithms (GAs), which are a family of metaheuristics designed to tackle combinatorial and discrete problems.

This method was proposed to address an extended version of the MEC problem in which genotyping data are considered during the SNP correction process.

GAHap [30] uses GAs to infer the haplotype pair of an individual working on nucleotide strings. During the optimization, GAHap solves the MEC problem by means of a majority rule that takes into account allele frequencies. No all-heterozygous assumption is required.

GenHap [9] is a novel computational method based on GAs to solve the wMEC problem. This method exploits a *divide-et-impera* approach to partition the entire problem into smaller and manageable overlapped sub-problems. In order to solve in parallel the sub-problems, GenHap was developed using a Master-Slave approach to leverage multi-core architectures.

2.2 GenHap: A Distributed Computing Implementation for HA

Hereafter, we briefly recall the peculiarities of GenHap [9], by focusing on the HPC implementation. GenHap tackles the HA problem by solving the wMEC problem, exploiting an approach based on GAs. Since the execution time and the problem difficulty increase with the number of reads and SNPs of the input data, GenHap follows a *divide-et-impera* approach [31] in which the wMEC problem is efficiently solved by splitting the fragment matrix \mathbf{M} into $\Pi = \lfloor m/\gamma \rfloor$ sub-matrices consisting of γ reads (where γ depends on the coverage value and on the nature of the sequencing technology). So doing, the problem difficulty is reduced by solving the sub-problems by means of independent GA executions that eventually converge to solutions having two sub-haplotypes with the least number of corrections to the SNP values. Finally, these sub-haplotypes are combined to achieve the complete haplotype pair. It is worth noting that GenHap considers all phased positions [32] as heterozygous during the optimization phase with GAs. As soon as the sub-haplotypes are obtained, all possible uncorrected heterozygous sites are removed and the correct value is assigned by checking the columns of the sub-partitions.

GenHap makes use of a Master-Slave distributed programming paradigm [33] to speed up the overall execution (Fig. 2). It was developed using the C++ programming language and the Message Passing Interface (MPI) specifications to leverage multi-core Central Processing Units (CPUs). The Master-Slave strategy of GenHap consists of the following phases: **1**) the Master process (MPI rank 0) proceeds by (i) allocating the necessary resources, (ii) partitioning the matrix into Π sub-matrices, and (iii) offloading the data onto the available Σ Slave processes. Each Slave σ (with MPI rank $1 \leq \sigma \leq \Sigma$) proceeds by randomly generating the initial population of the GA; **2**) each Slave executes the assigned wMEC sub-task by means of an independent GA instance. If multiple cores are available, the Slave processes are executed in a parallel fashion; **3**) as soon as the wMEC sub-tasks are terminated, the Master process recombines the sub-solutions received from the Slaves, and yields the complete wMEC solution.

According to the GA settings analysis provided in [9], we used here 100 individuals, tournament selection with size equal to 10 individuals, crossover and mutation rates equal to 0.9 and 0.05, respectively. Finally, the elitism strategy is

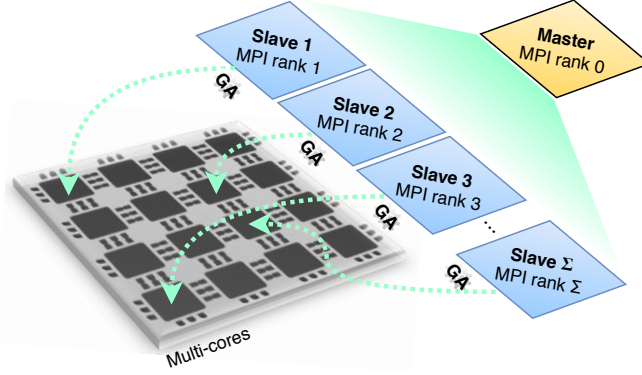


Fig. 2. Graphical representation of the Master-Slave approach implemented by GenHap: the Master process handles all the Σ Slaves sending one or more sub-partitions to each Slave, which then solves the assigned wMEC sub-task leveraging a core.

exploited to copy the best individual from the current population into the next one without undergoing the genetic operators.

3 Test Battery and Results

In what follows, we present some computational results obtained by considering different sequencing technologies, namely: Illumina NovaSeq, Roche/454, PacBio RS II, and ONT MinION. In [9], GenHap was shown to be faster than HapCol achieving up to $20\times$ speed-up on PacBio RS II instances, reconstructing haplotypes characterized by a very low haplotype error rate. Moreover, GenHap was capable of solving in about 10 minutes a real PacBio RS II instance characterized by $\#\text{SNPs} \simeq 28000$ and $\#\text{reads} \simeq 140000$, with average and maximum coverages equal to 29 and 565, respectively. Notice that a direct comparison with the only other parallel method, pWhatsHap [26], was not possible since the source code of that tool is no longer publicly available.

In order to assess the computational performance of GenHap, we used the General Error-Model based SIMulator (GemSIM) toolbox [34] to generate different instances of synthetic (yet realistic) data, compliant with these sequencing technologies. GemSIM generates the instances relying on empirical error models and distributions learned from real NGS data. A detailed description of the whole pipeline is described in [9]. For each sequencing technology, we generated a single instance varying the following parameters: (i) $\#\text{SNPs} \in \{500, 1000, 5000, 10000, 20000\}$; (ii) $\text{cov} \in \{\sim 30\times, \sim 40\times, \sim 50\times, \sim 60\times\}$; (iii) average $f_{\text{SNPs}} = 200$ (i.e., one SNP every 200bp exists [35, 36]).

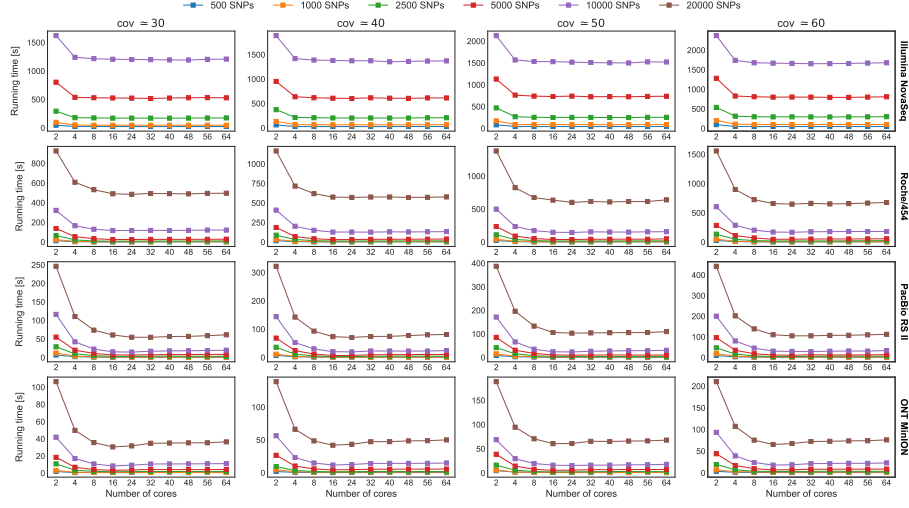


Fig. 3. Comparison of the running time required by GenHap on sequencing data generated by four sequencing technologies (Illumina NovaSeq, Roche/454, PacBio RS II, ONT MinION) by varying the coverage values. Note that the instances generated using the Illumina NovaSeq technology and characterized by $\#SNPs = 20000$ required more RAM than the amount of memory available on the computing nodes used for the tests. The tests were executed by increasing the number of cores exploited to run GenHap, to evaluate the scalability of the implementation based on distributed computing.

These instances were used to evaluate the scalability of GenHap by varying the number of cores, that is, $\#cores \in \{2, 4, 8, 16, 24, 32, 40, 48, 56, 64\}$. All tests were performed on the MARCONI supercomputer, which is based on Lenovo NeXtScale System[®] platform (Morrisville, NC, USA), provided by the Italian inter-university consortium CINECA (Bologna, Italy). Three different partitions running on CentOS 7.2 are available on this supercomputer:

- A1 Broadwell (BDW) partition consists of 720 compute nodes, each one equipped with 2 Intel[®] Xeon[®] E5-2697 v4 (18 cores at 2.30GHz) and 128 GB RAM;
- A2 Knights Landing (KNL) partition consists of 3600 compute nodes, each one equipped with an Intel[®] Knights Landing (68 cores at 1.40GHz and 16 GB MCDRAM), which is the next-generation of the Intel[®] Xeon Phi[™] product family many-core architecture, and 93 GB RAM;
- A3 Skylake (SKL) partition consists of 92 compute nodes, each one equipped with 2 Intel[®] Xeon[®] 8160 CPU (18 cores at 2.10GHz) and 192 GB RAM.

Our analysis was carried out by using the computing nodes of partition A2, which was chosen due to the availability of a higher number of computing cores.

Fig. 3 depicts the running times required by GenHap to infer the pairs of haplotypes. As expected, the processing time generally decreases along with the read length: according to Eq. (1), the same coverage can be obtained by means

of long reads coupled with a lower number of reads. This circumstance leads to a lower number of sub-problems to be solved, reducing the necessary computational effort. Moreover, the lowest running time is achieved on the instances generated relying on the ONT MinION, which is capable of producing long reads (up to 6000bp) with accuracy greater than 92%. As a matter of fact, the amount of SNPs to be corrected decreases when reads characterized by high accuracy are taken into account, allowing the GA instances to have a fast convergence to the optimal solutions. The results obtained for each sequencing platform are summarized as follows: (i) Illumina NovaSeq: independently from the coverage, the lowest running time is achieved by exploiting 24 cores to parallelize the GA instances when $\#SNPs = 10000$. When $\#SNPs < 10000$, 16 or 24 cores require the minimum running time to infer the haplotype pairs; (ii) Roche/454: when $\#SNPs \geq 5000$, the best GenHap's performance is achieved by exploiting 16 or 24 cores, otherwise the best choice is 24 cores; (iii) PacBio RS II: in every test, the fastest executions are generally obtained by exploiting 24 cores to parallelize the GA instances, except when $\#SNPs = 500$ is taken into account. In this case, the running time decreases when 16 cores are exploited to effectively distribute the GA instances on multiple cores. Since the reads generated by relying on this technology have a low accuracy (approximately 87%), which makes the problem more difficult to be solved (i.e., the amount of SNPs to be corrected increases), the scalability of GenHap is emphasized; (iv) ONT MinION: in all tests, the best choice is 16 cores that allow for efficiently distributing the computational load. In all tests, a number of cores greater than 24 does not reduce the running time since the overhead introduced by MPI is not entirely mitigated by the required computational load. Furthermore, when the number of sub-problems is lower than the number of available cores, our Master-Slave approach exploits a number of cores equal to the number of sub-problems. On the one hand, when technologies producing short reads are considered, the number of haplotype blocks increases along with $\#SNPs$. Since these blocks are solved sequentially and are generally characterized by a number of sub-problems lower than the available cores, 16 or 24 cores allow for balancing the computational load. On the other hand, technologies producing long reads generate a small number of reads that lead to a low number of sub-problems to be solved. Notice that exploiting the accuracy of the reads produced using Illumina NovaSeq, Roche/454 and ONT MinION, the GA instances have a fast convergence to the optimal solutions requiring only a dozen of generations.

4 Conclusion and Future Trends

In this paper, we presented a complete overview on the currently available HA computational tools, focusing on the potential of HPC. In particular, we investigated the computational performance of GenHap [9], a recent evolutionary method leveraging multi-core architectures.

As a future development, we plan to extend GenHap to deal with HA in organisms characterized by different ploidity. Differently from diploid organisms

having two copies of each chromosome set, polyploid organisms have multiple copies of their chromosome sets. Polyploidy has gained scientific interest in the study of the ongoing species diversification phenomena [37]. This characteristic is mainly present in plant genomes, but also in animals (such as salmonid fishes and African clawed frogs) [38]. In these comparative genomic studies, haplotype-aware assemblies play a crucial role in elucidating genetic and epigenetic regulatory evolutionary aspects. Unfortunately, the computational burden of the HA problem is emphasized in the case of polyploid haplotypes with respect to diploids [39]. Therefore, HPC represents a key element for efficient, accurate, and scalable methods for HA of both diploid and polyploid organisms.

An interesting future trend in Genome Analysis is related to its connection with machine learning. As a matter of fact, deep learning has been successfully applied in population genetic inference and learning informative features of data [40]. Combining population genetics inference and HA can provide insights on patterns regarding the genetic diversity in DNA polymorphism data, especially for rapid adaptation and selection [41].

An additional issue worth of notice is that, although the integration of various types of information (e.g., electronic health records and genome sequences) conveys a wealth of information, it is giving rise to unique challenges in bioinformatics analysis even in terms of secure genomic information sharing [42]. With reference to secure Genome-Wide Association Study (GWAS) in distributed computing environments, multi-party computation schemes based on conventional cryptographic techniques achieves limited performance in practice [13]. Therefore, HPC could become an enabling factor also in this context.

Acknowledgment

We acknowledge the CINECA for the availability of High Performance Computing resources and support.

References

1. Zhang, K., Calabrese, P., Nordborg, M., Sun, F.: Haplotype block structure and its applications to association studies: power and study designs. *Am. J. Hum. Genet.* **71**(6) (2002) 1386–1394
2. Daly, M.J., Rioux, J.D., Schaffner, S.F., Hudson, T.J., Lander, E.S.: High-resolution haplotype structure in the human genome. *Nat. Genet.* **29**(2) (2001) 229
3. Rhoads, A., Au, K.F.: PacBio sequencing and its applications. *Genomics Proteomics Bioinformatics* **13**(5) (2015) 278–289
4. Roberts, R.J., Carneiro, M.O., Schatz, M.C.: The advantages of SMRT sequencing. *Genome Biol.* **14**(6) (2013) 405
5. Jain, M., Fiddes, I.T., Miga, K.H., Olsen, H.E., Paten, B., Akeson, M.: Improved data analysis for the MinION Nanopore sequencer. *Nat. Methods* **12**(4) (2015) 351

6. Quail, M.A., Kozarewa, I., Smith, F., Scally, A., Stephens, P.J., Durbin, R., Swerdlow, H., Turner, D.J.: A large genome center's improvements to the Illumina sequencing system. *Nat. Methods* **5**(12) (2008) 1005
7. Lander, E.S., Waterman, M.S.: Genomic mapping by fingerprinting random clones: a mathematical analysis. *Genomics* **2**(3) (1988) 231–239
8. Sims, D., Sudbery, I., Illott, N.E., Heger, A., Ponting, C.P.: Sequencing depth and coverage: key considerations in genomic analyses. *Nat. Rev. Genet.* **15**(2) (2014) 121
9. Tangherloni, A., Spolaor, S., Rundo, L., Nobile, M.S., Cazzaniga, P., Mauri, G., Liò, P., Merelli, I., Besozzi, D.: GenHap: A novel computational method based on genetic algorithms for haplotype assembly. *BMC Bioinform.* (2018) Accepted.
10. Luo, C., Tsementzi, D., Kyrpides, N., Read, T., Konstantinidis, K.T.: Direct comparisons of Illumina vs. Roche 454 sequencing technologies on the same microbial community DNA sample. *PloS one* **7**(2) (2012) e30087
11. Jain, M., Koren, S., Miga, K.H., Quick, J., Rand, A.C., Sasani, T.A., Tyson, J.R., Beggs, A.D., Diltney, A.T., Fiddes, I.T., et al.: Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat. Biotechnol.* **36**(4) (2018) 338
12. Senol Cali, D., Kim, J.S., Ghose, S., Alkan, C., Mutlu, O.: Nanopore sequencing technology and tools for genome assembly: computational analysis of the current state, bottlenecks and future directions. *Brief. Bioinform.* (2018) bby017.
13. Choi, Y., Chan, A.P., Kirkness, E., Telenti, A., Schork, N.J.: Comparison of phasing strategies for whole human genomes. *PLoS Genet.* **14**(4) (2018) e1007308
14. Wang, R., Wu, L., Li, Z., Zhang, X.: Haplotype reconstruction from SNP fragments by minimum error correction. *Bioinformatics* **21**(10) (2005) 2456–2462
15. Greenberg, H.J., Hart, W.E., Lancia, G.: Opportunities for combinatorial optimization in computational biology. *INFORMS J Comput* **16**(3) (2004) 211–231
16. Browning, S.R., Browning, B.L.: Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am. J. Hum. Genet.* **81**(5) (2007) 1084–1097
17. Delaneau, O., Marchini, J., Zagury, J.F.: A linear complexity phasing method for thousands of genomes. *Nat. Methods* **9**(2) (2012) 179
18. Loh, P.R., Danecek, P., Palamara, P.F., Fuchsberger, C., Reshef, Y.A., Finucane, H.K., Schoenherr, S., Forer, L., McCarthy, S., Abecasis, G.R., et al.: Reference-based phasing using the haplotype reference consortium panel. *Nature Genet.* **48**(11) (2016) 1443
19. Bansal, V., Bafna, V.: HapCUT: an efficient and accurate algorithm for the haplotype assembly problem. *Bioinformatics* **24**(16) (2008) i153–i159
20. Edge, P., Bafna, V., Bansal, V.: HapCUT2: robust and accurate haplotype assembly for diverse sequencing technologies. *Genome Res.* **27**(5) (2017) 801–812
21. Kuleshov, V.: Probabilistic single-individual haplotyping. *Bioinformatics* **30**(17) (2014) i379–i385
22. Duitama, J., Huebsch, T., McEwen, G., Suk, E., Hoehe, M.: ReFHap: a reliable and fast algorithm for single individual haplotyping. In: *Proceedings of the First ACM International Conference on Bioinformatics and Computational Biology*, ACM (2010) 160–169
23. Levy, S., Sutton, G., Ng, P.C., Feuk, L., Halpern, A.L., Walenz, B.P., Axelrod, N., Huang, J., Kirkness, E.F., Denisov, G., et al.: The diploid genome sequence of an individual human. *PLoS Biol.* **5**(10) (2007) e254
24. Chen, Z.Z., Deng, F., Wang, L.: Exact algorithms for haplotype assembly from whole-genome sequence data. *Bioinformatics* **29**(16) (2013) 1938–1945

25. Patterson, M., Marschall, T., Pisanti, N., Van Iersel, L., Stougie, L., Klau, G.W., Schönhuth, A.: WhatsHap: weighted haplotype assembly for future-generation sequencing reads. *J. Comput. Biol.* **22**(6) (2015) 498–509
26. Bracciali, A., Aldinucci, M., Patterson, M., Marschall, T., Pisanti, N., Merelli, I., Torquati, M.: pWhatsHap: efficient haplotyping for future generation sequencing. *BMC Bioinform.* **17**(Suppl 11) (2016) 342
27. Pirola, Y., Zaccaria, S., Dondi, R., Klau, G., Pisanti, N., Bonizzoni, P.: HapCol: accurate and memory-efficient haplotype assembly from long reads. *Bioinformatics* **32**(11) (2015) 1610–1617
28. Benedettini, S., Roli, A., Di Gaspero, L.: Two-level ACO for haplotype inference under pure parsimony. In: *Proceedings of the International Conference on Ant Colony Optimization and Swarm Intelligence*. Volume 5217 of LNCS., Springer (2008) 179–190
29. Na, J.C., Lee, J.C., Rhee, J.K., Shin, S.Y.: PEATH: Single individual haplotyping by a probabilistic evolutionary algorithm with toggling. *Bioinformatics* (2018) bty012
30. Wang, T.C., Taheri, J., Zomaya, A.Y.: Using genetic algorithm in reconstructing single individual haplotype with minimum error correction. *J. Biomed. Inform.* **45**(5) (2012) 922–930
31. Maisto, D., Donnarumma, F., Pezzulo, G.: Divide et impera: subgoalng reduces the complexity of probabilistic inference and problem solving. *J. R. Soc. Interface* **12**(104) (2015) 20141335
32. Kuleshov, V., Xie, D., Chen, R., Pushkarev, D., Ma, Z., Blauwkamp, T., Kertesz, M., Snyder, M.: Whole-genome haplotyping using long reads and statistical methods. *Nat. Biotechnol.* **32**(3) (2014) 261–266
33. Tangherloni, A., Rundo, L., Spolaor, S., Cazzaniga, P., Nobile, M.S.: GPU-powered multi-swarm parameter estimation of biological systems: a master-slave approach. In: *Proceedings of the 26th Euromicro Conference on Parallel, Distributed, and Network-Based Processing (PDP)*, Cambridge, UK, March 21–23, 2018, Euromicro (2018) In press.
34. McElroy, K.E., Luciani, F., Thomas, T.: GemSIM: general, error-model based simulator of next-generation sequencing data. *BMC Genomics* **13**(1) (2012) 74
35. Nachman, M.W.: Single nucleotide polymorphisms and recombination rate in humans. *Trends Genet.* **17**(9) (2001) 481–485
36. Gabriel, S.B., Schaffner, S.F., Nguyen, H., Moore, J.M., Roy, J., Blumenstiel, B., Higgins, J., DeFelice, M., Lochner, A., Faggart, M., et al.: The structure of haplotype blocks in the human genome. *Science* **296**(5576) (2002) 2225–2229
37. Otto, S.P., Whitton, J.: Polyploid incidence and evolution. *Annu. Rev. Genet.* **34**(1) (2000) 401–437
38. Rodriguez, F., Arkhipova, I.R.: Transposable elements and polyploid evolution in animals. *Curr. Opin. Genet. Dev.* **49** (2018) 115–123
39. Das, S., Vikalo, H.: SDhaP: haplotype assembly for diploids and polyploids via semi-definite programming. *BMC Genomics* **16**(1) (2015) 260
40. Sheehan, S., Song, Y.S.: Deep learning for population genetic inference. *PLoS Comput. Biol.* **12**(3) (2016) e1004845
41. Hermisson, J., Pennings, P.S.: Soft sweeps and beyond: understanding the patterns and probabilities of selection footprints under rapid adaptation. *Methods Ecol. Evol.* **8**(6) (2017) 700–716
42. Bianchi, L., Liò, P.: Opportunities for community awareness platforms in personal genomics and bioinformatics education. *Brief. Bioinform.* **18**(6) (2016) 1082–1090