*PREPRINT*

# Ococo: an online variant and consensus caller

Karel Břinda [1,2]*, Valentina Boeva [3,4], and Gregory Kucherov [5,6]

[1]Center for Communicable Disease Dynamics, Department of Epidemiology, Harvard TH Chan School of Public Health, Boston MA; [2]Department of Biomedical Informatics, Harvard Medical School, Boston MA; [3]Institut Curie – Centre de Recherche, PSL Research University, Mines Paris Tech, INSERM U900, Paris, France; [4]Institut Cochin, INSERM U1016, CNRS UMR 8104, Université Paris Descartes UMR-S1016, Paris, France; [5]LIGM/CNRS, Université Paris-Est, Marne-la-Vallée, France; [6]Skolkovo Institute of Science and Technology (SkolTech), Moscow Region, Russia.

**ABSTRACT**

**Motivation:** Identifying genomic variants is an essential step for connecting genotype and phenotype. The usual approach consists of statistical inference of variants from alignments of sequencing reads. State-of-the-art variant callers can resolve a wide range of different variant types with high accuracy. However, they require that all read alignments be available from the beginning of variant calling and be sorted by coordinates. Sorting is computationally expensive, both memory- and speed-wise, and the resulting pipelines suffer from storing and retrieving large alignments files from external memory. Therefore, there is interest in developing methods for resource-efficient variant calling.

**Results:** We present Ococo, the first program capable of inferring variants in a real-time, as read alignments are fed in. Ococo inputs unsorted alignments from a stream and infers single-nucleotide variants, together with a genomic consensus, using statistics stored in compact several-bit counters. Ococo provides a fast and memory-efficient alternative to the usual variant calling. It is particularly advantageous when reads are sequenced or mapped progressively, or when available computational resources are at a premium.

**Availability:** `http://github.com/karel-brinda/ococo`

**Contact:** `kbrinda@hsph.harvard.edu`

## 1 INTRODUCTION

Identifying genomic variants is an essential step for connecting genotype and phenotype. The goal of *variant calling* is to identify genomic variants present in the sequenced individual or a population. Most commonly, variant calling proceeds by read mapping and then sliding a small window throughout the genome, collecting statistics for all reads aligned within the window and calculating the likelihood of variants observed in these alignments. We term this approach *offline variant calling* as it requires that all read alignments are available from the beginning. Offline calling is implemented in all major variant callers (see, e.g., (Bao et al., 2014)).

However, offline variant calling is highly time- and space-demanding. First, all alignments must be available and get sorted by coordinates prior to variant calling; this involves storing and retrieving large alignment files from external memory. Second, variant callers usually apply computationally expensive steps, such

as realignments, even for regions where this is not necessary. The resulting performance can be particularly limiting on portable devices, personal computers, or in a cloud environment with restricted resources.

Here, we introduce the concept of *online variant calling*, where variants are inferred in real time, as read alignments are fed in. We implement this approach in a program called Ococo, the first online variant caller. Ococo inputs unsorted alignments from an unsorted SAM/BAM stream (Li et al., 2009) and infers single-nucleotide variants, together with a genomic consensus, using statistics stored in compact several-bit counters. Ococo provides a fast and memory-efficient alternative to the usual variant calling, which is particularly advantageous when reads are sequenced or mapped progressively, or when available computational resources are at a premium.
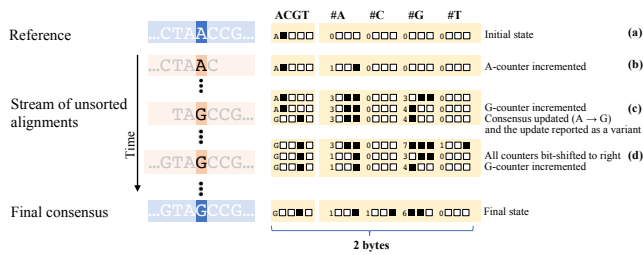
## 2 METHODS

**Overview.** Ococo calls variants and consensus directly from an unsorted SAM/BAM file, possibly provided in a stream. To do that, Ococo stores and maintains variant statistics for all genomic positions about previous alignments as well as a consensus sequence, which represents the current internal reference. The consensus can be initialized from a user-provided sequence, typically the same as used for read mapping. Whenever a new alignment is loaded, Ococo updates the statistics and assesses whether they are still concordant with the consensus. If not, the consensus is corrected and the corresponding substitution reported as a novel variant.

**Compact representation of variant statistics.** In the online approach reads can potentially map to any location. This is a fundamental difference from the offline calling, where reads are sorted and statistical inference uses a small sliding window, collecting information about locally overlapping alignments. Therefore, the main challenge of online calling is to design variant statistics for the whole genome that fit into main memory and at the same time be sufficiently informative for inferring variants. We propose using small, several-bit nucleotide counters and complementing them with fast bit operations.

The Ococo statistics consist of four integer counters per position, one per each nucleotide (**Figure 1**). Every counter represents the number of nucleotides aligned to that position; however only the most significant bits are stored. Whenever a new alignment is loaded, the corresponding counters are incremented (**Figure 1a**). If a counter is already saturated and yet is to be incremented, then all counters at the position are first bit-shifted, losing their rightmost (least significant) bit (**Figure 1c**). This mechanism makes it possible to compute nucleotide frequencies in a limited space and filter out randomly distributed sequencing errors. Ococo supports three counter

---

**Figure 1. Internal statistics of an online variant caller.** Example of update of OCOCO counters for a single position of the genome. The first 4 bits carrying the nucleotide consensus are followed by 4 nucleotide counters, each of them 3 bits long in this case. Vertical axis corresponds to time. The figure shows how the counters and consensus are updated based on the received alignments. **a)** At the beginning, the consensus base is initialized to the reference base (A). All counters are set to 0. **b)** The A counter is incremented; the statistics stay concordant with the consensus. **c)** The G counter is incremented, which triggers a consensus update and reporting a new variant (A→G). **d)** The G counter is to be incremented, but it is already saturated. Therefore, all counters of the position must be bit-shifted first.

configurations: 16, 32, and 64 bits per position corresponding to 4 bits for the consensus base and 3, 7, and 15 bits for each nucleotide counter, respectively.

**Variant and consensus calling strategy.** For the sake of speed, we propose considering individual genomic positions independently and keeping the consensus synchronized with the dominating base if such exists (**Figure 1b**). Formally defined, let $C_n$ denote the value of the counter for a nucleotide $n \in \{A, C, G, T\}$ at some fixed position, then the consensus is updated to $n$ if $2C_n > C_A + C_C + C_G + C_T$. This represents a simple instantiation of maximum likelihood estimation for haploid genomes and single-nucleotide variants without considering base qualities. Variants and the resulting consensus are reported in the VCF (Danecek et al., 2011) and FASTA format, respectively.

**Working modes.** OCOCO supports two modes of online calling: the *real-time* and *batch modes*. Whereas the *real-time mode* updates consensus and reports variants immediately after processing each read, the *batch mode* postpones reporting updates until all reads from the current batch have been processed.

**Implementation.** OCOCO is implemented in C++ and released under the MIT license. The software package is available from `http://github.com/karel-brinda/ococo`, BioConda (Grüning et al., 2018), and Zenodo (Břinda et al., 2017).

## 3 RESULTS

We performed two experiments demonstrating the accuracy and resource-efficiency of OCOCO (**Supplementary Figure 1**). Using RNFTOOLS (v0.3.1.3) (Břinda et al., 2016) with DWGSIM (v0.1.11), we simulated reads of length 100bp from the *Chlamydia trachomatis* genome (NC_021897.1, 1.046Mbp), with 20x coverage and 1% sequencing error rate. The genome was *in silico* mutated at the 2% level (with 15% of mutations being indels and with the 30% indel extension probability). The reads were then mapped back to the original reference using BWA-MEM (v0.7.17) (Li, 2013).

First, we evaluated how online variant calling progressed in time (**Supplementary Figure 1a**). We used the obtained alignments to call SNPs in the OCOCO (v0.1.2.6) streaming mode with

7-bit counters. After processing the first 50 thousand reads (corresponding to approximately 5x coverage), OCOCO neared a plateau of the edit distance between the simulated and inferred genomes. At this point, OCOCO had correctly identified 16,239 out of 21,914 single-edit variants (17,613 SNPs; indels of total length 4,301) inserted into the reference genome, i.e., 92,2% of SNPs were identified. After all the 209,742 reads have been processed, the proportion of correctly identified SNPs increased to 97,3%.

Second, we compared the speed of the OCOCO batch mode with a common variant-calling pipeline consisting of SAMTOOLS (v1.9) (Li et al., 2009) and VARSCAN (v2.4.3) (Koboldt et al., 2009) (**Supplementary Figure 1b**). For this pipeline, we measured the time required for sorting the alignments, computing an alignment pileup, and for the subsequent SNP calling. Times were measured in seconds, as the mean over 3 runs, using SNAKEMAKE (v5.3.0) (Köster and Rahmann, 2012) on an iMac 4.2 GHz Intel Core i7 with an SSD disk and 40 GB RAM. We observed that OCOCO provided 66x speedup for calling variants compared to SAMTOOLS and VARSCAN.

To assess scaling to larger genomes, we ran the same experiments on human chromosome 17 (HG18, 78,775Mbp) (**Supplementary Figure 2**). The obtained results were qualitatively similar: 91,3% and 96,1% of SNP were identified with 5x and 20x of reads, respectively, and OCOCO provided 55x speedup. The decrease in speed-up can be explained by more CPU cache misses due to the higher number of counters for longer genomes.

## 4 DISCUSSION

OCOCO brings several advantages over offline variant callers. First of all, it enables determining variants and consensus sequences directly from the output of a read mapper, avoiding heavy I/O operations. In many applications, OCOCO produces a sufficiently good-quality consensus supplemented by information about nucleotide frequencies at each position. This is especially relevant for bacterial genomics, where many statistical methods consider single-nucleotide variants only. Finally, online consensus calling is also an essential component of the *dynamic mapping approach* that we developed in a separate study (Břinda et al., 2016).

A key ingredient of OCOCO is the memory-efficient alignment statistics featuring small counters. We observe that they can be seen as an instance of approximate counting (Morris, 1978), but with a modified formulation: whereas traditional algorithms estimate *counts*, variant calling requires estimating *ratios* of the counts of nucleotides. Here, we provide an algorithmic description of the counter mechanics, although statistical properties of the resulting estimators are yet to be studied.

A major limitation of OCOCO is its restriction to single-nucleotide variants. Indels and more complicated variants present three challenges for future research. First, it is necessary to design an appropriate statistics, ideally counter-based, for each new variant type. Second, the statistics must be complemented with rules for triggering an update, which should be fast to evaluate. Finally, non-substitution updates entail changes in genomic and counter coordinates, which calls for a more sophisticated addressing and allocation system than the one implemented in OCOCO.
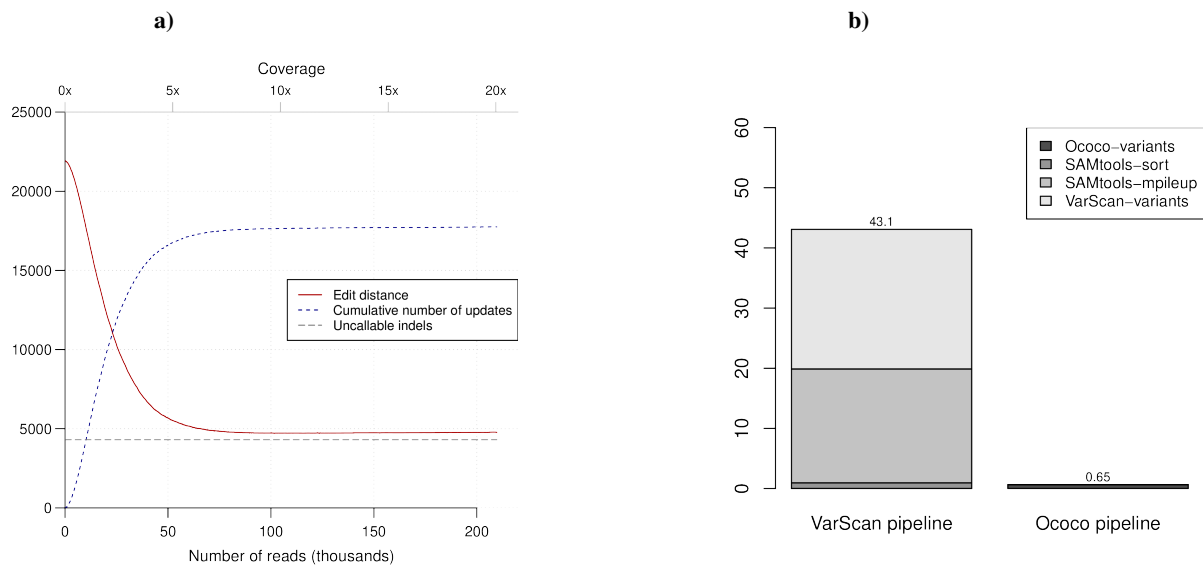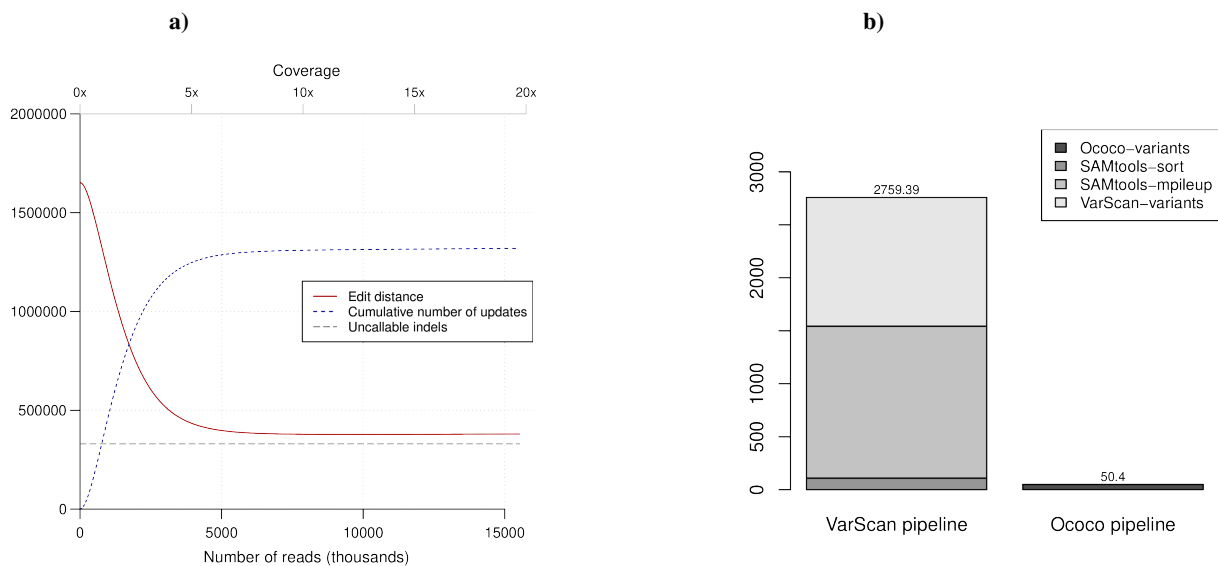
## 5 ACKNOWLEDGEMENTS

*Conflict of Interest*: None declared.

## REFERENCES

Bao, R., Huang, L., Andrade, J., Tan, W., Kibbe, W. A., Jiang, H., and Feng, G. (2014). Review of Current Methods, Applications, and Data Management for the Bioinformatics Analysis of Whole Exome Sequencing. *Cancer Informatics*, 13s2:CIN.S13779.

Břinda, K., Boeva, V., and Kucherov, G. (2016). Dynamic read mapping and online consensus calling for better variant detection. *arXiv:1605.09070*.

Břinda, K., Boeva, V., and Kucherov, G. (2016). RNF: a general framework to evaluate NGS read mappers. *Bioinformatics (Oxford, England)*, 32(1):136–9.

Břinda, K., Boeva, V., and Kucherov, G. (2017). karel-brinda/ococo: Ococo 0.1.2.6. *Zenodo*, 1095696. Available from https://doi.org/10.5281/zenodo.1095696.

Danecek, P., Auton, A., Abecasis, G., Albers, C. A., Banks, E., DePristo, M. A., Handsaker, R. E., Lunter, G., Marth, G. T., Sherry, S. T., McVean, G., and Durbin, R. (2011). The variant call format and VCFtools. *Bioinformatics*, 27(15):2156–2158.

Grüning, B., Dale, R., Sjödin, A., Chapman, B. A., Rowe, J., Tomkins-Tinch, C. H., Valieris, R., and Köster, J. (2018). Bioconda: sustainable and comprehensive software distribution for the life sciences. *Nature Methods*, 15(7):475–476.

Koboldt, D. C., Chen, K., Wylie, T., Larson, D. E., McLellan, M. D., Mardis, E. R., Weinstock, G. M., Wilson, R. K., and Ding, L. (2009). VarScan: variant detection in massively parallel sequencing of individual and pooled samples. *Bioinformatics*, 25(17):2283–5.

Köster, J. and Rahmann, S. (2012). Snakemake-a scalable bioinformatics workflow engine. *Bioinformatics*, 28(19):2520–2522.

Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv:1303.3997*, page 3.

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R., and 1000 Genome Project Data Processing Subgroup (2009). The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, 25(16):2078–9.

Morris, R. (1978). Counting large numbers of events in small registers. *Communications of the ACM*, 21(10):840–842.

**Supplementary Figure 1. Evaluation of OCOCO with Chlamydia Trachomatis (1.046Mbp)**. **a) Online variant calling as a function of time.** The blue curve shows the cumulative number of updates of the consensus as a function of the number of processed alignments (or the actual coverage). The red curve shows the edit distance from the simulated sequenced genome. **b) Speed comparison.** Comparison of time to completion of variant calling using OCOCO and a pipeline based on VARSCAN.



**Supplementary Figure 2. Evaluation of OCOCO with human chromosome 17 (78,775Mbp).** The figure is of the same format as **Supplementary Figure 1.**