

HammerIT: homopolymer-space Hamming clustering for IonTorrent read error correction

Anton Korobeynikov^{1,2}, Artem Tarasov¹

¹ Department of Statistical Modelling, St. Petersburg State University, St. Petersburg, Russia; ² Algorithmic Biology Laboratory, St. Petersburg Academic University, St. Petersburg, Russia

Introduction

Error correction of sequenced reads remains a difficult task, especially for data obtained using IonTorrent technology due to its higher error rate. The task is even more challenging in single-cell sequencing projects with extremely non-uniform coverage.

The existing error correction tools assume that the most sequencing errors in the data are mismatches and thus perform poorly on IonTorrent data with its prevailing errors due to homopolymer indels.

We introduce HAMMERIT — a novel error correction tool which is specifically tuned for IonTorrent sequencing errors.

IonTorrent error profile

Corrected flow signal intensities are available in BAM files produced by versions of Ion Torrent Suite prior to 3.4. Called homopolymer length is obtained as corrected flow signal intensity rounded to the nearest integer.

We have studied flow signal intensity distributions around insertion/deletion sites. File B7-295.bam, downloaded from the Ion Community website, contained 4.6M insertions, 5.0M deletions, and 1.5M mismatches.

The overwhelming majority of errors turned out to be insertions/deletions of length 1, occurring when flow signal intensity is approximately halfway between two adjacent integers.

Figure 1: Flow signal intensities at insertion sites

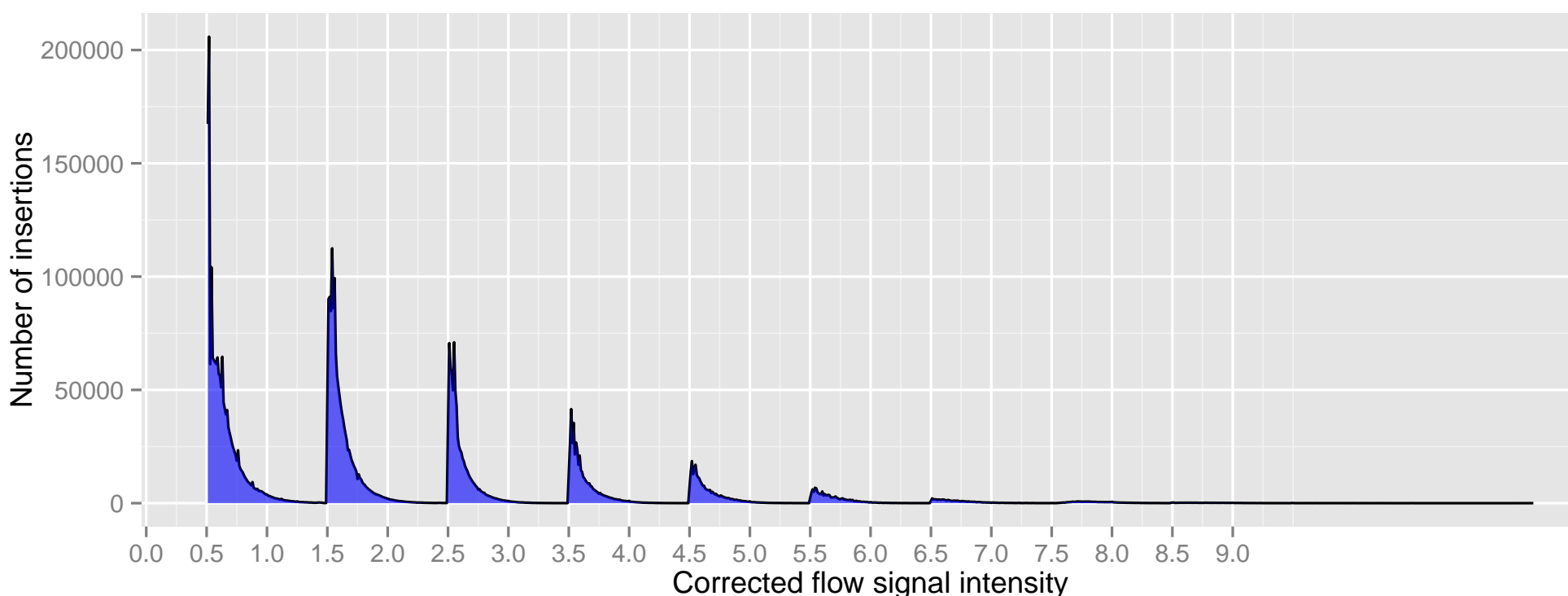
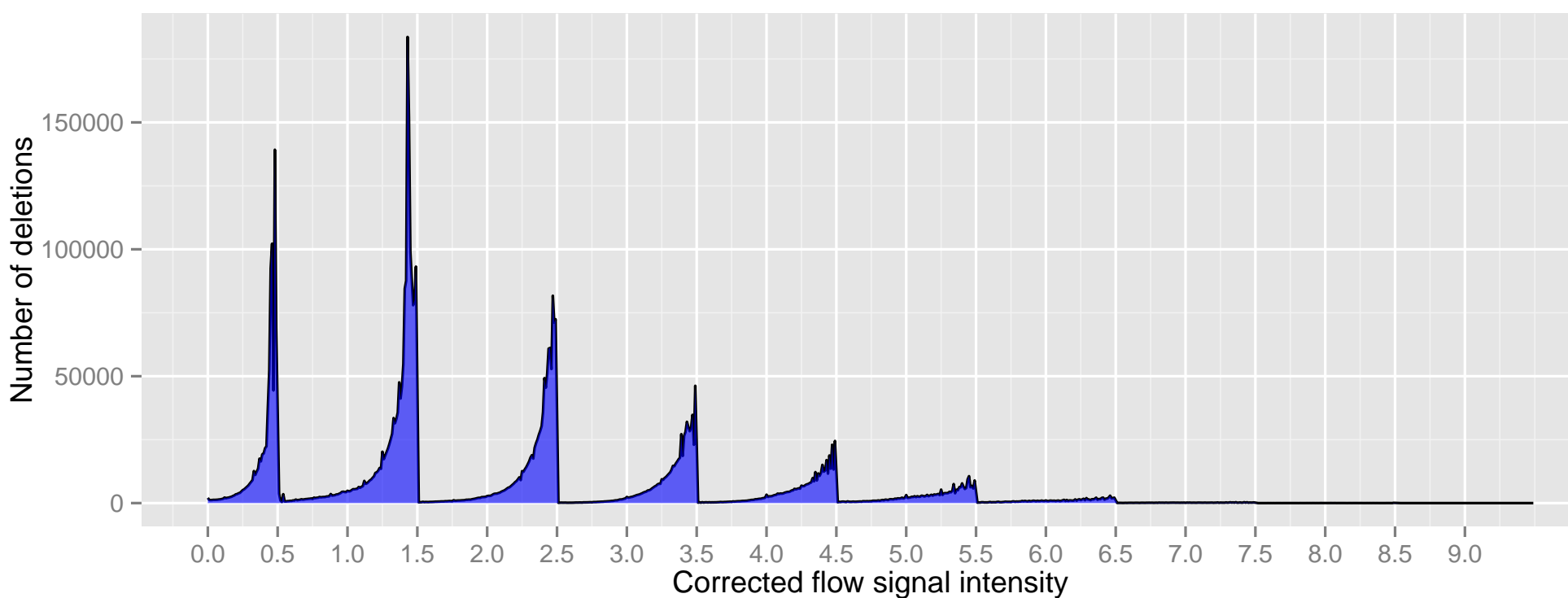


Figure 2: Flow signal intensities at deletion sites



Notation

Let \mathcal{N} denote the nucleotide alphabet $\{A, C, G, T\}$.

Define $\mathcal{H} = \mathcal{N} \times \mathbb{N}^+$.

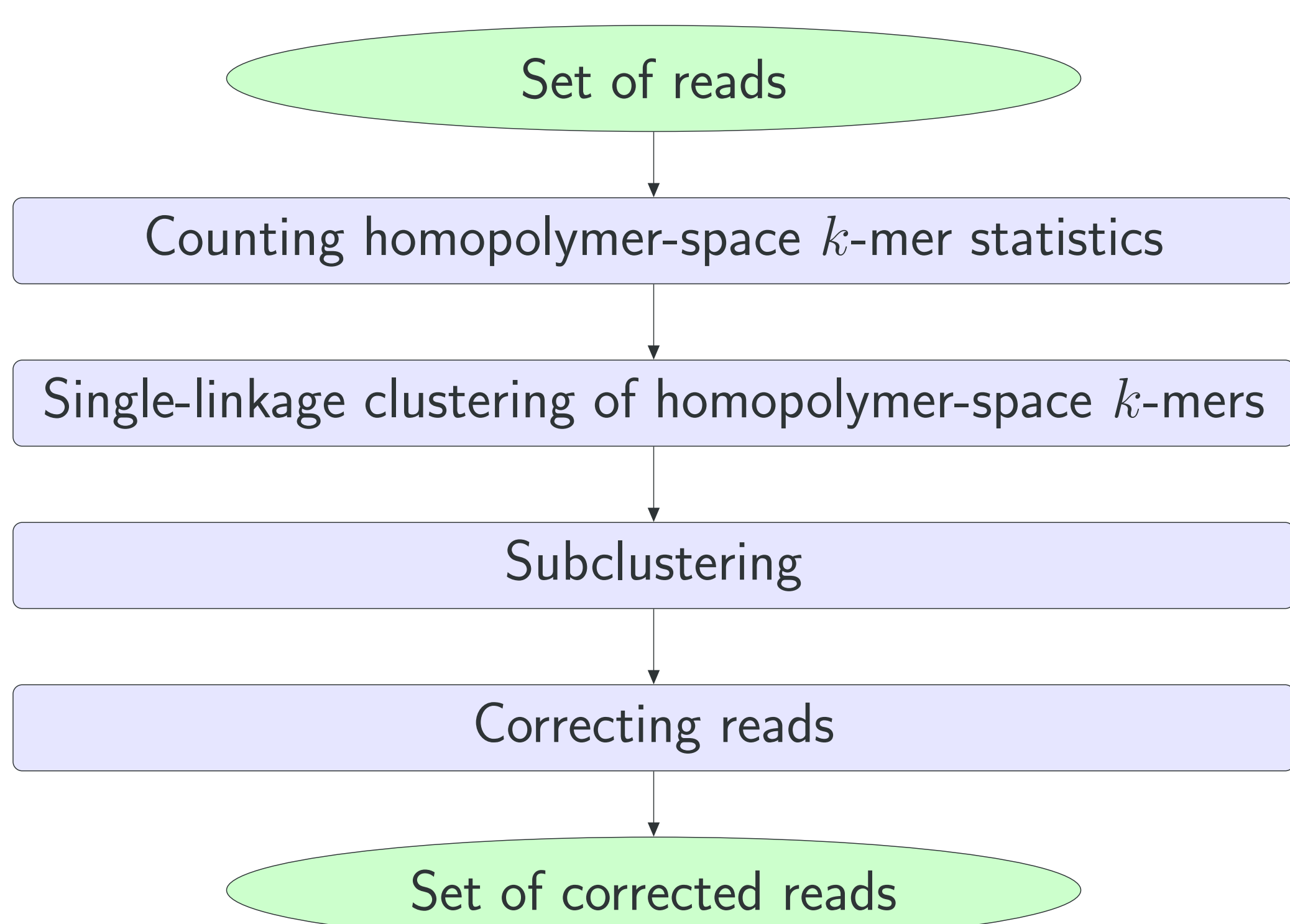
We call an element of the alphabet \mathcal{H} a *homopolymer run*, and an element of \mathcal{H}^k a *homopolymer-space k -mer*.

We use $x[k]$ to denote the k -th element of a sequence, using zero-based indexing; $x[k \dots l]$ to denote a subsequence $x[k]x[k+1] \dots x[l]$.

The length of x is denoted by $|x|$.

Distance between $i, j \in \mathcal{H}^k$ is defined as the minimum number of 1-base insertions/deletions/mismatches needed to align the common start of i and j . It is denoted by $\text{dist}(i, j)$.

HammerIT workflow



Pairwise distance calculation

We use 5-base lookahead to compute the distance between k -mers in homopolymer-space. A helper table stores precomputed values of

$$H_k : \mathcal{N}^k \times \mathcal{N}^k \rightarrow \{\text{Insertion, Deletion, Mismatch}\}, \quad k = 1, 2, 3, 4, 5.$$

The chosen value of 5 is a trade-off between accuracy and speed.

(Mapped reads from B7-295.bam dataset were used for training.)

G A G T A C A C T G T C G T C G
G T G T A C A T G T C G A T G C

$$H_5(\text{AGTAC}, \text{TGTAC}) = \text{Mismatch}$$

$$H_5(\text{CTGTG}, \text{TGTGCG}) = \text{Deletion}$$

$$H_3(\text{TCG}, \text{ATG}) = \text{Mismatch}$$

$$H_2(\text{CG}, \text{TG}) = \text{Mismatch}$$

Algorithm

Input: $x, y \in \mathcal{N}^+$ — homopolymer-space k -mers in nucleotide alphabet

Output: dist — approximate distance between x and y

$\text{pos.x} \leftarrow 0$; $\text{pos.y} \leftarrow 0$; $\text{dist} \leftarrow 0$;

while $\text{pos.x} < |x|$ and $\text{pos.y} < |y|$ **do**

if $x[\text{pos.x}] = y[\text{pos.y}]$ **then**

$\text{pos.x} \leftarrow \text{pos.x} + 1$; $\text{pos.y} \leftarrow \text{pos.y} + 1$;

else

$k \leftarrow \min(5, |x| - \text{pos.x}, |y| - \text{pos.y})$;

 adjust pos.x and pos.y according to

$$H_k(x[\text{pos.x} \dots \text{pos.x} + k - 1], y[\text{pos.y} \dots \text{pos.y} + k - 1]);$$

$\text{dist} \leftarrow \text{dist} + 1$;

Homopolymer-space k -mer clustering

The starting point is single-linkage clustering of homopolymer-space k -mers.

Two homopolymer-space k -mers belong to the same cluster if the distance between their nucleotide representations does not exceed one.

We reduce quadratic time requirements of the naive algorithm by noticing that if distance between two k -mers is less or equal to one, they share a common substring of length at least $\lfloor k/2 \rfloor$. This allows us to group k -mers into smaller blocks sharing a substring, and then use the quadratic algorithm for each block.

In order to detect all pairs of connected k -mers with such grouping, ranges $(0 \dots \lfloor k/2 \rfloor - 1)$, $(1 \dots \lfloor k/2 \rfloor)$, \dots , $(\lfloor (k+1)/2 \rfloor \dots k-1)$ should be examined. However, we use only the first and the last of the ranges, to speed up the clustering step. This increases the number of *singletons* — clusters consisting of only one k -mer, usually erroneous — but has little impact on correction performance because of adjacent read k -mers making much larger contribution into the consensus scores (the number of singleton occurrences is usually low).

$\text{partition}(\mathcal{K}, i, j) = \{B_s : \bigcup_{s \in \mathcal{H}(j-i+1)} B_s = \mathcal{K}, \forall x \in B_s \ x[i \dots j] = s\}$

Algorithm

$\text{kmers} \leftarrow$ homopolymer-space k -mers seen in the data

$\text{components} \leftarrow \{\{k\} : k \in \text{kmers}\}$

$\text{blocksL} \leftarrow \text{partition}(\text{kmers}, 0, \lfloor k/2 \rfloor - 1)$

$\text{blocksR} \leftarrow \text{partition}(\text{kmers}, \lfloor (k+1)/2 \rfloor, k-1)$

for each block in blocksL, blocksR do

for each $i \in \text{block}$ do

for each $j \in \text{block do}$

if $\text{dist}(i, j) \leq 1$ **then**

 join the components to which i and j belong

Error reduction results

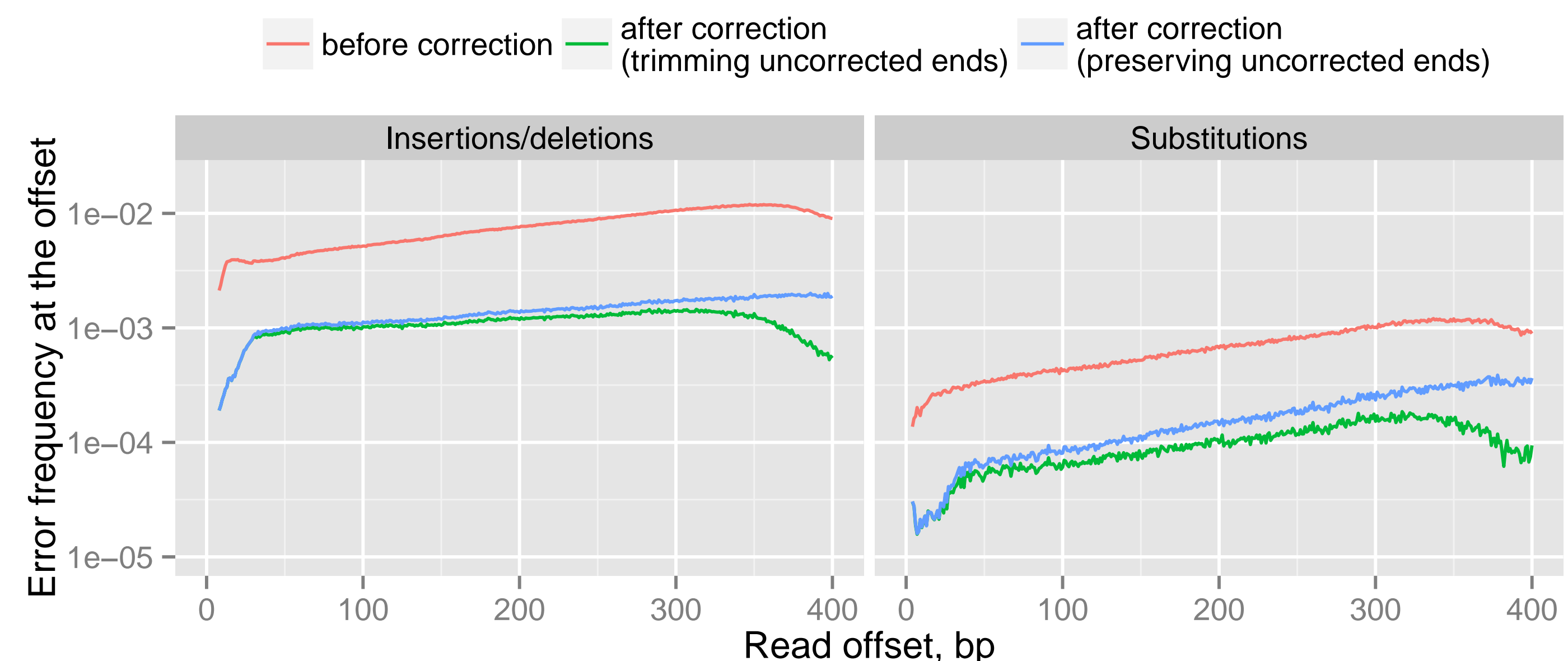
We evaluated HammerIT on 6 publicly available datasets, using the same pipeline as the authors of the recently published article (*Jünemann et al, Nat. Biotech., 2013; vol. 31, p.294–296*). In that article, error rate in four Ion Torrent datasets has been assessed. We used the same data plus two extra datasets from 314v2 chip, which recently became available on the Ion Community Portal.

Indel/mismatch error rates were calculated for uniquely mapped reads before and after correction. For each dataset, correction was done in two ways. In the first setup, trimming was done for read ends that couldn't be corrected due to lack of good k -mers, while in the second one such read ends were preserved in the output. The relative change in read coverage after correction stayed within 0.4% in all cases.

Figure 3: Error rates before and after correction



Figure 4: Error reduction by read position for Sakai 400bp reads



Subclustering

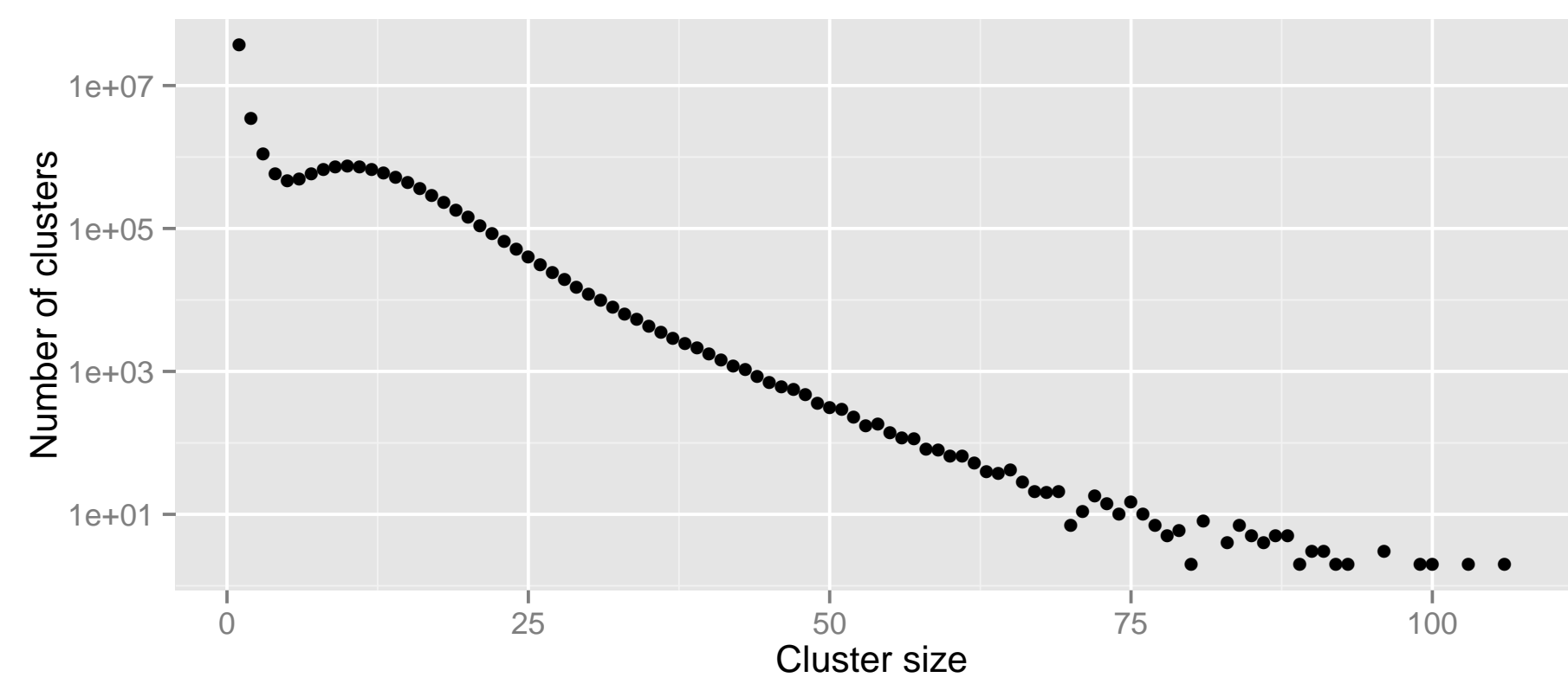
A cluster obtained from the initial process may contain $m \geq 2$ homopolymer-space k -mers from the genome. In this case, we select the m most occurring elements to be the subcluster centers. Each of the remaining elements is then assigned to the nearest center.

The subtle question is how to determine m . Currently we just set it to be the number of the cluster elements with quality within machine epsilon of 1, where *quality* of $x \in \mathcal{H}^k$ is defined as $\Pr(x \text{ is genomic})$. For the set of reads $\mathcal{R} \subset \mathcal{N}^+$ the quality of $x \in \mathcal{H}^k$ is computed as

$$1 - \prod_{r \in \mathcal{R}} \prod_{\substack{0 \leq m \leq |r| - |x'|, \\ r[m \dots m + |x'| - 1] = x'}} \left(1 - \prod_{m \leq n < m + |x'|} \Pr(r[n] \text{ is correct}) \right),$$

where $x' \in \mathcal{N}^+$ is the sequence of nucleotides in x .

Cluster size distribution



Typical cluster

homopolymer-space 16-mer	#	qual.
CCCGTTTGT-GCGCCC-GG--CATTT-GAT	562	1.00
CCCGTTTGT-GCGCCC-GG--CATT--GAT	13	1.00
CCCGTT-GT-GCGCCC-GG--CATTT-GAT	10	1.00
CCCGTTTGT-GCGCCC-GG--CATTT-GAT	6	1.00
CCCGTTTGT-GCGCCC-GG--CATTTGAT	5	0.99
CCCGTTTGT-GCGCCC-GG--CATTT-GAT	4	1.00
CCCGTTTGT-GCGCCC-GG--CATTT-GAT	4	0.98
CGTTTGT-GCGCCC-GG--CATTT-GAT	3	1.00
CCCGTTTGT-GCGCCC-GG--CATTT-GAT	3	0.98
CCCGTTTGT-GCGCCC-GG--CATTT-GAAT	3	0.98
CCCGTTTGT-GCGCCC-GGA-CATTT-GA	3	0.98
CCCGTTTGT-GCGCCC-GG--C-TTT-GATG	3	0.97
CCCGTT-GTTGCGCCC-GG--CATTT-GAT	2	0.94
CCCGTTTGTGCGCCC-GG--CATTT-GAT	2	0.94
CCCGTTTGTGCGCCC-GG--CATTT-GAT	2	0.93
CCCGTTTGT-GCGCCC-GG--CATTT-GAG	2	0.89
CCCGTTTGT-GCGCCC-GG--CATTTGGAT	2	0.74
CCCGTTTGTG-GCGTCG-GG--CATTT-G	1	0.97
CCCGTTTGTG-GCGCC-GG--CACTTTGA	1	0.87
CCCGTTTGTG-GCGCC-GG--CACTTT-GAT	1	0.85
CCCGTTTGTG-GCCC-G---CATTT-GAT	1	0.83
CCCGTTTGTG-GCCC-GG--CATTT-GTG	1	0.74
CCCGTTTGTGCGCCC-GG--CATTT-GAT	1	0.71
CCCGTTTGTG-GCGCCC-GG--CATTT-GAT	1	0.69
CCCGTTTGTG-GCGCCC-GGTACATTT-G	1	0.68
CCCGTTTGT-GCGCCC-GG--CATTTGGA	1	0.68
CCCGTTTGTGCGCCC-GG--CATTT-GAT	1	0.66
CCGTT-GT-GCGCC-GG--CATTT-GAT	1	0.66
CCCGTTTGTG-GCGCC-GG--CATTT-GAT	1	0.61
CCCGTTTGT-GCGCCC-GGACCATTT-GA	1	0.61
CCCGTTTGTG-GCCC-GG--CATTTGTA	1	0.58
CCCGTTTGTGCGCCC-GG--CATTT-GAT	1	0.57
CCCGTTTGTG-GCCC-GGA-CATTT-AG	1	0.56
CCCGTTTGT-GCGCCC-GG--CATTT-AGT	1	0.55
CCGTTTGTG-GCC-GG--CATTT-GAT	1	0.54
CCCGTT-GTGC-GCCC-GG--CATTT-GAT	1	0.44
CCCGTTTGTGCGCCCTGG--CATTT-GA	1	0.43
CCCGTTTGTGCGCCC-GG--CATTTGAT	1	0.36
CCCGTTTGTGCGCCCGG--CATTT-GAT	1	0.22

Assembly results

We have assembled *E. coli O157 H7 Sakai* 400bp reads before and after correction with SPAdes 2.4.0. Read coverage is ~240x, length of the reference genome is 5498450 bp.

Table 1: Assembly results. Contigs of length ≥ 500 are used.

	uncorrected	corrected
# contigs	266	242
Largest contig	316904	374932
Total length	5322223	5320566
NG50	106242	146551
NG75	41186	44318
# misassemblies	0	2
# local misassemblies	8	11
Misassembled contigs length	0	32563
Genome fraction (%)	93.880	93.966
# mismatches per 100 kbp	3.58	4.70
# indels per 100 kbp	6.21	5.88

Command-line parameters used for assembly:

--only-assembler -k 21,33,55,77,99

Acknowledgements

This work was supported by the Government of the Russian Federation (grant 11.G34.31.0018) and RFBR (grant 12-01-00747-a).

We would like to thank Sergey Nikolenko and Alla Lapidus for their insightful comments and fruitful discussions throughout the project.

Further information



SPAdes: <http://bioinf.spbau.ru/spades>