1    **HapFlow: Visualising haplotypes in sequencing data**

2    Mitchell J. Sullivan, Nathan L. Bachmann, Peter Timms, Adam Polkinghorne*


3    Faculty of Science, Health, Education and Engineering, University of the Sunshine Coast,
4    Sippy Downs 4556, Australia

5

6    Running Head: HapFlow: Visualising haplotypes

7

8    *To whom correspondence should be addressed. Tel: +61 7 5456 5578; Email:
9    apolking@usc.edu.au

10

11

12

13    *Abstract*

14    **Summary:** HapFlow is a python application for visualising haplotypes present in high-

15    throughput sequencing data. HapFlow identifies nucleotide variant profiles in raw read

16    sequences and creates an abstract visual representation of these profiles to make haplotypes

17    easier to identify.

18    **Availablity:** HapFlow is freely available (under a GPL license) for download (for Mac OS X,

19    Unix and Microsoft Windows) from github (http://mjsull.github.io/HapFlow).

20    **Contact:** apolking@usc.edu.au

21

## Introduction

The emergence of high-throughput sequencing has enabled new experimental approaches such as the sequencing of bacterial populations. Infections frequently contain multiple strains of the same species (Darch, et al., 2015; Taylor, et al., 1995). This has important implications for detecting transmission events (Bachmann, et al., 2015) and determining treatment outcomes (Cohen, et al., 2012). Several methods have been developed to analyse mixed-strain populations. ShoRAH (Zagordi, et al., 2011) reconstructs a minimal set of global haplotypes and estimates the frequency of inferred haplotypes. It requires variants be dense enough to be linked by overlapping reads. A two-step maximum likelihood approach has also been described to identify the portion of infection rising from dominant and minor strains (Eyre, et al., 2013). This approach does not rely on variant density but is unable to infer local or global haplotypes. A tool that visualises haplotypes in sequencing data is needed to identify the best strategy for genomic analysis of multiple strains of the same bacteria within a sample.

Many excellent read alignment visualisation tools exist including Savant (Fiume, et al., 2010), Tablet (Milne, et al., 2010) and Consed (Gordon and Green, 2013). These tools arrange reads in a linear fashion with each read represented as a line, or row of bases. This layout is satisfactory for identifying variants or misaligned reads, however, it is not ideal for identifying haplotypes present in reads. Reads are packed tightly together making it difficult to determine whether distant variants are located on the same read pair. Additionally, reads are not grouped by haplotype making it difficult to identify how frequently a haplotype is represented in the sequencing data.

HapFlow addresses these problem by abstracting read alignment data to make the haplotypes

47  present easier to identify. HapFlow can be used to help identify potential sites of

48  recombination, identify the minimum number of strains present in a sample and determine

49  whether defining local or global haplotypes is possible using sequence data alone.

50  ***Implementation***

51  HapFlow is a python tool that uses the Tkinter windows system. It is available as a Python

52  script or using the package manager PIP. It contains two parts: HapFlow-generator, a process

53  for creating a flow file, which contains the count of reads with each haplotype profile and

54  HapFlow-viewer, a tool for visualising the flow file.

55

56  **HapFlow-generator** can be executed from the GUI or the command-line. It takes a VCF file

57  of called variants and an indexed BAM file of aligned reads as input. Pysam is used to create

58  a profile of variants present in each read of the alignment. This profile consists of which

59  variant or variants are present in the read, on which pair each variant is present and the

60  direction of the read. If the variant profile is unique, a flow (profile of variants in a read) is

61  created. If the flow already exists in another read, the count of the flow is incremented by

62  one.

63

64  **HapFlow-viewer** displays the created flow file on the canvas of the GUI. An orange

65  rectangle within the blue rectangle represents the portion of the genome currently being

66  displayed. Underneath, an orange rectangle with vertical lines represents where the variants

67  are located within the displayed section of the genome, these lines are extended below and

68  spaced an equal distance apart in the area where the flows are viewed. Each flow consisting

69  of one or more reads is represented as one or more arrows overlapping each variant line that

70  the reads of the flow align to. Width of the arrow represents the number of reads within that

71  flow. Variants on the same read of a pair are joined by a solid line, variants on different reads

72 of a pair are joined by a dotted line. Arrows grouped at the top of the canvas represent the

73 most common variant, the second group of flows represents the second most common variant

74 and so on. The last group represents potential sequencing, alignment or variant calling errors

75 as they have sequence that has not been called as a variant. Information about the sequence of

76 each variant is represented underneath the flows. The canvas is scrollable and zoomable

77 allowing the user to easily navigate through whole genomes.

78

79 **Results and discussion**

80 To demonstrate the application of HapFlow, reads from the recent sequencing of a *Chlamydia*

81 *pecorum* PCR-positive swab sample collected from the urogenital tract of a koala with mixed

82 *C. pecorum* infections were analysed. *C. pecorum* DNA was extracted directly from the host

83 cell contaminants using Sure-Select RNA probes and sequenced using an Illumina Hi-Seq to

84 produce 101bp paired-end reads, as previously described (Bachmann, et al., 2015). These

85 reads were then mapped back to E58 using Bowtie-2 and then variant calling was performed

86 using FreeBayes. Exploration of the HapFlow diagram identified several regions in low

87 complexity areas where non-chlamydial DNA had been captured. Importantly, several regions

88 where read coverage in the dominant strain dropped below that of the minor strain were

89 identified (Figure 1). This was not unexpected as sequence capture is less efficient at

90 capturing DNA in areas where the sequence of the strain varies significantly from the probe.

91 This meant that any method of consensus calling that relied on coverage would result in a

92 chimeric genome not representative of either strain. Due to the proximity of variants, a

93 linkage approach was used to determine the sequence for large regions of both strains.
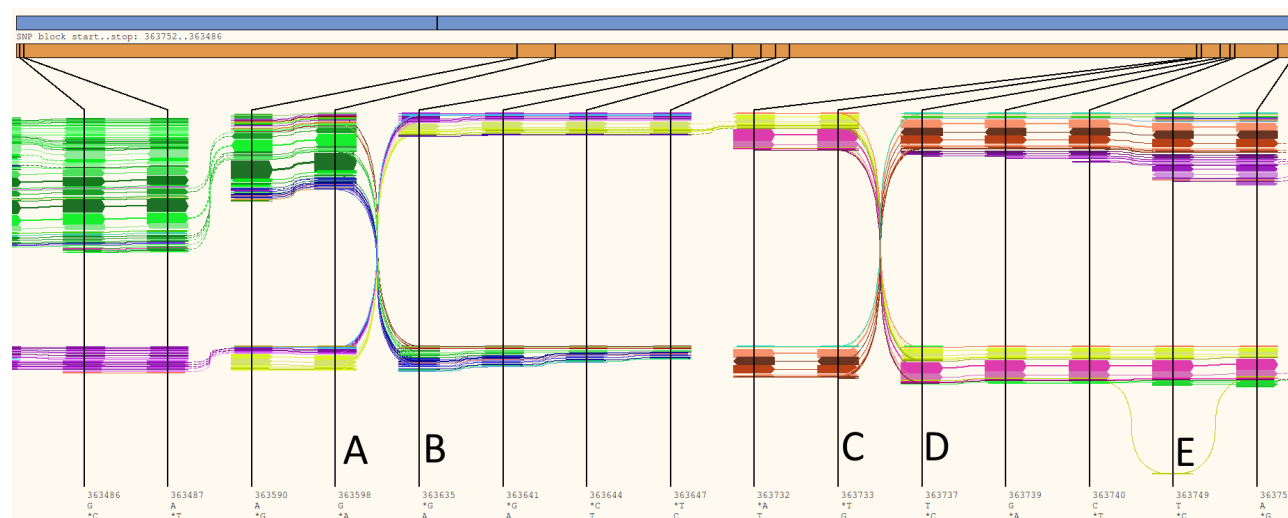
98



**Figure 1: HapFlow diagram of sequencing data from a urogenital tract infection in a koala mapped to *C. pecorum*.** Flows containing only the dominant variant group at the top, flows containing only the minor variant group in the second row while mixed flows switch between top and middle. A flow containing an alignment or sequencing error can be seen at site E. All reads with the most common variant at site A have the least common variant at site B. Similarly all reads with the least common variant at site A, have the most common variant at site B. This pattern is repeated at sites C and D.

*References*

Bachmann, N.L.*, et al.* (2015) Culture-independent genome sequencing of clinical samples reveals an unexpected heterogeneity of infections by *Chlamydia pecorum*. *J Clin Microbiol* **5,** 1-9.

Cohen, T.*, et al.* (2012) Mixed-strain mycobacterium tuberculosis infections and the implications for tuberculosis treatment and control. *Clin Microbiol Rev* **25,** 708-719.

Darch, S.E.*, et al.* (2015) Recombination is a key driver of genomic and phenotypic diversity in a *Pseudomonas aeruginosa* population during cystic fibrosis infection. *Sci. Rep* **5,** 7649.

Eyre, D.W.*, et al.* (2013) Detection of mixed infection from bacterial whole genome sequence data allows assessment of its role in *Clostridium difficile* transmission. *PLoS Comput Biol* **9,** e1003059.

Fiume, M.*, et al.* (2010) Savant: genome browser for high-throughput sequencing data. *Bioinformatics* **26,** 1938-1944.

Gordon, D. and Green, P. (2013) Consed: a graphical editor for next-generation sequencing.

121    *Bioinformatics* **29,** 2936-2937.
122    Milne, I.*, et al.* (2010) Tablet—next generation sequence assembly visualization.
123    *Bioinformatics* **26,** 401-402.
124    Taylor, N.S.*, et al.* (1995) Long-term colonization with single and multiple strains of
125    Helicobacter pylori assessed by DNA fingerprinting. *J Clin Microbiol* **33,** 918-923.
126    Zagordi, O.*, et al.* (2011) ShoRAH: estimating the genetic diversity of a mixed sample from
127    next-generation sequencing data. *BMC Bioinformatics* **12,** 119.
128
129