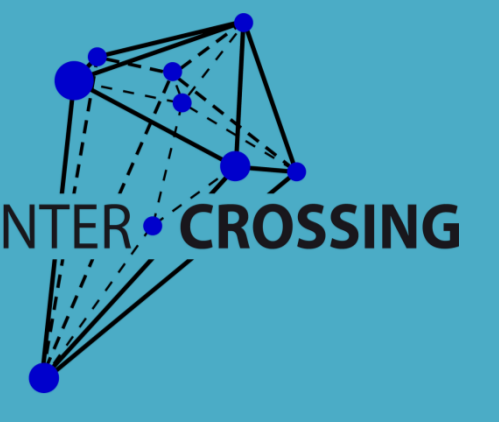


Assessing Allelic Configuration Models in Fixed Ploidy Variant Calling Using R - *Betula*



Jasmin Zohren¹, Igor Kardailsky², Kåre Lehmann Nielsen³, Anika Joecker², and Richard Buggs¹

1) School of Biological and Chemical Sciences, Queen Mary University of London, Mile End Road, London E1 4NS, UK; 2) CLC bio, a QIAGEN company, Silkeborgvej 2, Prismet, 8000 Aarhus C., DK; 3) Department of Chemistry and Bioscience, Sohngårdsholmsvej 49, 9000 Aalborg, DK



Genotyping polyploids

Genotyping of SNP loci in polyploids has always been challenging, but may become easier using high-coverage sequencing. The production of such data in turn requires the development of new methods and software that can be used to analyse it in user friendly software such as the CLC Genomics Workbench (GWB). The “Fixed Ploidy Variant Detection” tool of the CLC GWB already takes higher ploidy levels into account, but an explicit evaluation of the locus configuration is currently not provided. We have developed an algorithm that uses read counts and the average base quality to estimate the most likely allelic configuration at each variant position in polyploid samples.

The data



Solanum tuberosum (potato)
One tetraploid individual
Chromosome 11
50x coverage



Betula (birch)
213 individuals
Three different species
60x coverage

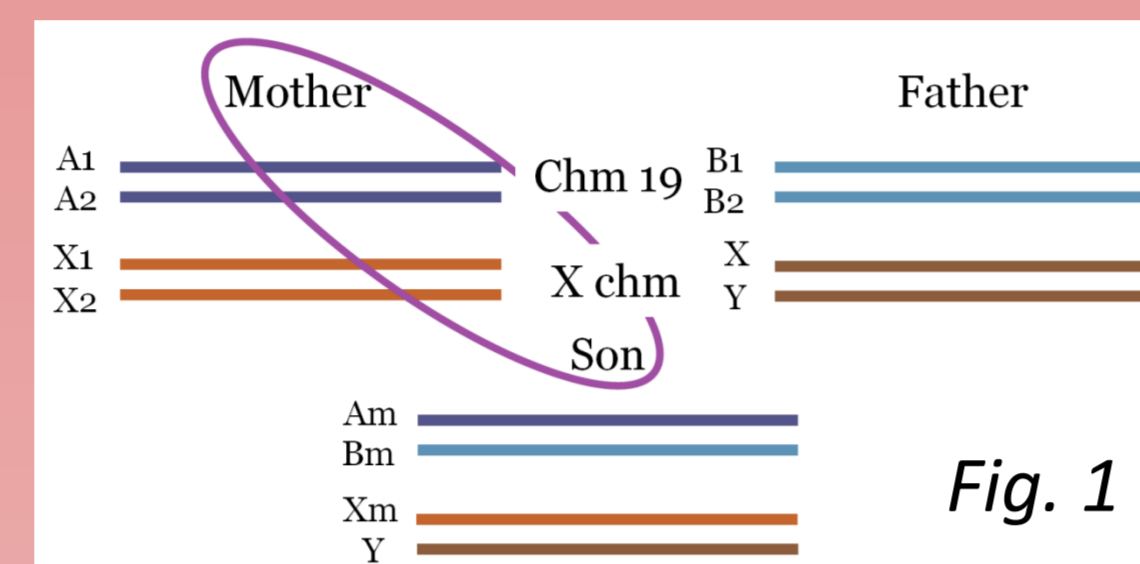


Homo sapiens
Illumina Platinum Genomes
Trio data set, chromosome 19
200x coverage

Tab. 1	Mother	Father	Father	Son	Mother	Son
Chm 19	A ₁ A ₂	B ₁ B ₂	B ₁ B ₂	B _m A _m	A ₁ A ₂	A _m B _m
X chm	X ₁ X ₂	XY	XY	X _m Y	X ₁ X ₂	X _m Y
Chm 19	A ₁ A ₂ B ₁ B ₂		B ₁ B ₂ B _m A _m		A ₁ A ₂ A _m B _m	
	4n (1:1:1:1)		3n (2:1:1)		3n (2:1:1)	
X chm	X ₁ X ₂ X		XX _m		X ₁ X ₂ X _m	
	3n (1:1:1)		2n (1:1)		2n (2:1)	

Allelic configurations of artificial polyploid humans. See figure 1 in “Tetraploid humans” section for details.

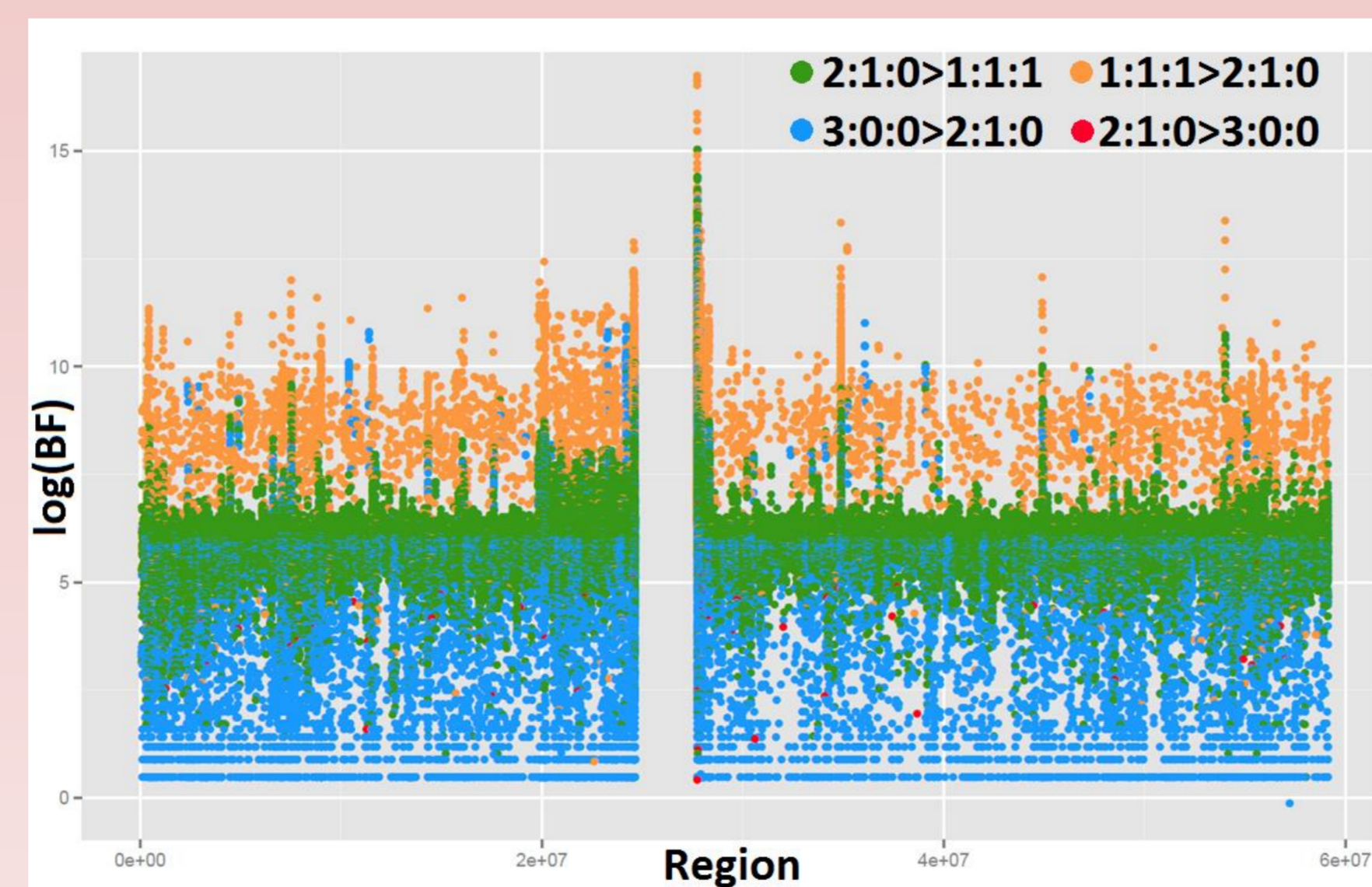
Tetraploid humans



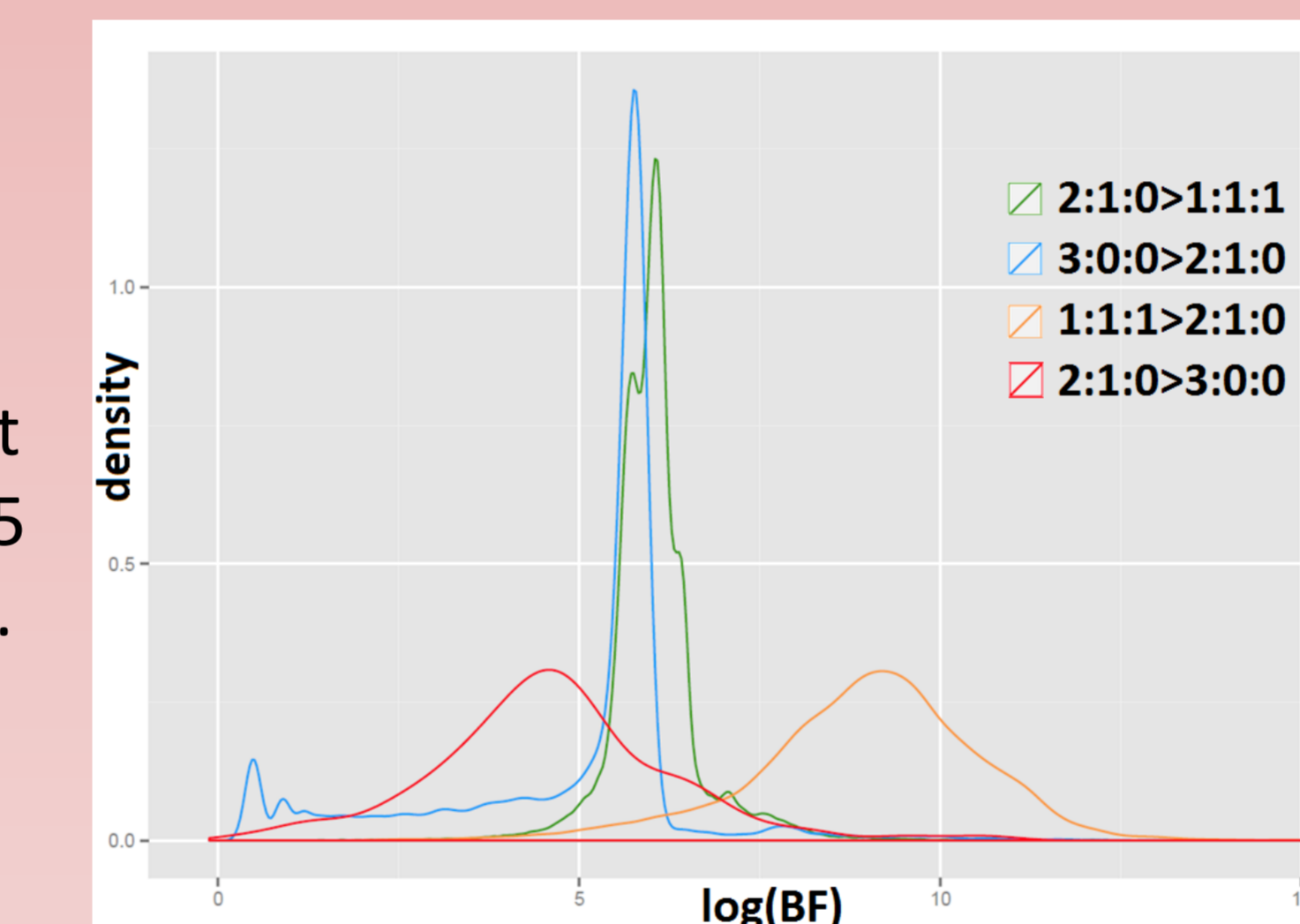
The reads of two individuals from the trio data set are merged, resulting in a polyploid human. The allele dosage differs depending on which individuals are merged and which chromosome is being analysed (see table 1 in “The data”). The data shown here is chromosome 19 of mother and son.

Best Model	#Hits	Second-best models	Mean (BF)	Median (BF)
3:0:0	31,560 (17.7%)	2:1:0	1,211	288
2:1:0	141,271 (79.2%)	3:0:0, 1:1:1*	926	406
1:1:1	5,615 (3.1%)	2:1:0	39,486	8,645
Total	178,446 (100%)			

Distribution of the Bayes factor (BF) for the different model comparisons. A Bayes factor of greater than 5 is considered to be significant.

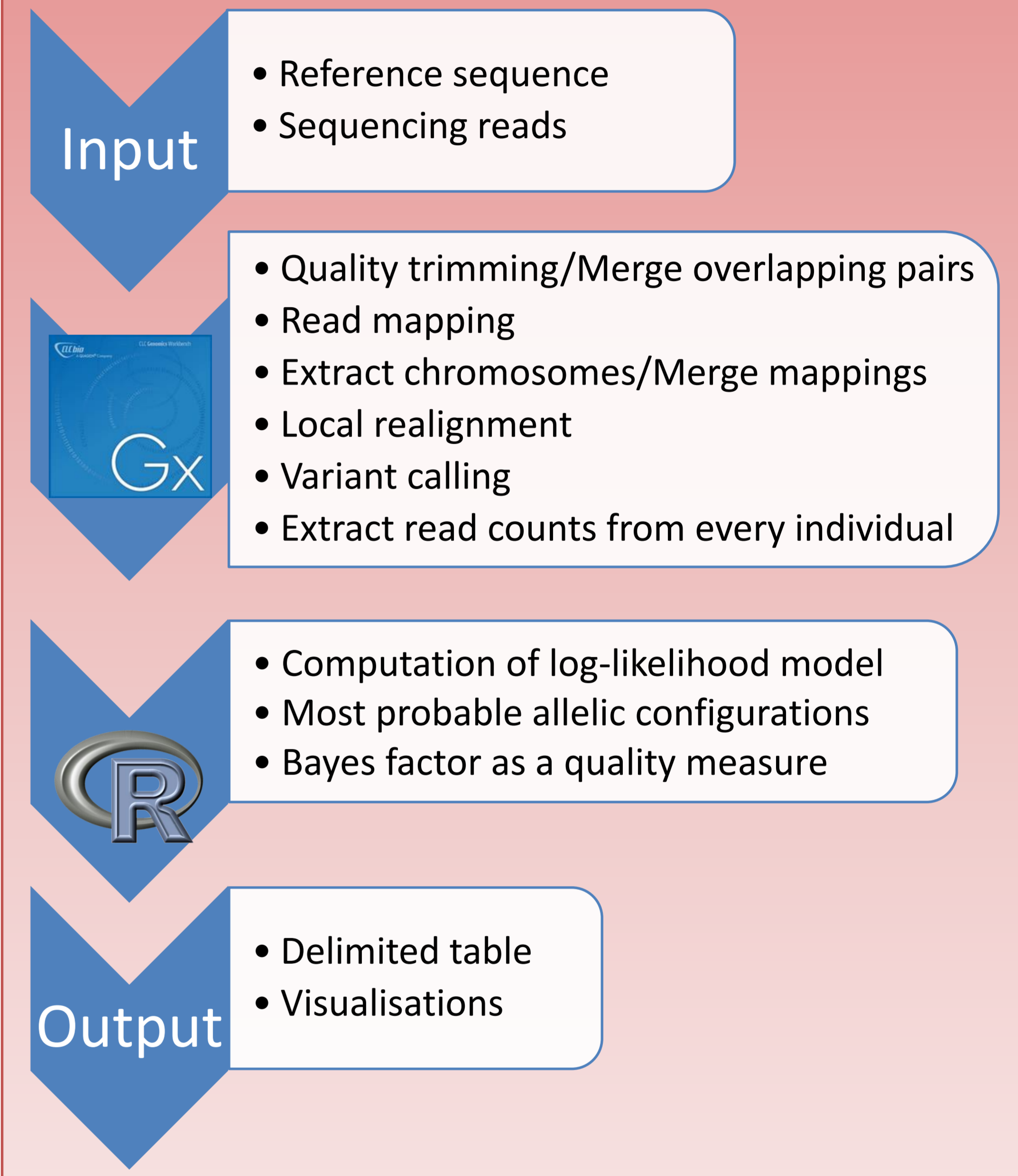


A total of 178,446 variants were called on the data set of chromosome 19 of mother and son. The “maximum expected ploidy” in the CLC variant caller was set to 3 to reflect the actual allele dosage (see table 1 in “The data” section).



Distribution of the different allele dosage models along chromosome 19. More “extreme” models are called with a higher certainty (i.e. Bayes factor) than others. A Bayes factor of greater than 5 is considered to be significant.

How it's done



Next steps

In order to validate the algorithm, more and different combinations of the trio data will be analysed (e.g. father and son and using reads mapping to the X chromosome). This might lead to a refinement of the mathematical model, for example to differentiate between allo- and auto-polyploids. We are also working on a sliding window approach to be able to detect regions with the same allele dosage and make an inference on sequence demography from sudden changes in these regions. The current version can already be used as a “ploidy detection” tool. Eventually, we intend to implement this as a tool within the CLC Genomics Workbench, probably as an addition to the already existing variant callers.

Funders

The Danish Council for Strategic Research

Jasmin Zohren

j.zohren@qmul.ac.uk

www.birchgenome.org