

Olorin: combining gene flow with exome sequencing in large family studies of complex disease

James A. Morris* and Jeffrey C. Barrett

Wellcome Trust Sanger Institute, Hinxton, Cambridge, CB10 1HH, UK

Associate Editor: Inanc Birol

ABSTRACT

Motivation: The existence of families with many individuals affected by the same complex disease has long suggested the possibility of rare alleles of high penetrance. In contrast to Mendelian diseases, however, linkage studies have identified very few reproducibly linked loci in diseases such as diabetes and autism. Genome-wide association studies have had greater success with such diseases, but these results explain neither the extreme disease load nor the within-family linkage peaks, of some large pedigrees. Combining linkage information with exome or genome sequencing from large complex disease pedigrees might finally identify family-specific, high-penetrance mutations.

Results: Olorin is a tool, which integrates gene flow within families with next generation sequencing data to enable the analysis of complex disease pedigrees. Users can interactively filter and prioritize variants based on haplotype sharing across selected individuals and other measures of importance, including predicted functional consequence and population frequency.

Availability: <http://www.sanger.ac.uk/resources/software/olorin>

Contact: olorin@sanger.ac.uk

Received on June 20, 2012; revised on October 4, 2012; accepted on October 8, 2012

1 INTRODUCTION

Next generation sequencing has rapidly become the standard approach for identifying mutations responsible for Mendelian diseases (Bamshad *et al.*, 2011). Although software and file formats for the processing of raw sequence data are relatively robust (Danecek *et al.*, 2011; Li *et al.*, 2009), there is currently a lack of easy-to-use software for downstream analysis of these data. For some study designs, such as focused analysis of fully penetrant *de novo* mutations or autosomal recessive inheritance, exome sequence data can be analysed and filtered relatively simply. Increasingly, however, sequence-based approaches are being applied to complex diseases, which are unlikely to follow a simple genetic model, such as autism (Neale *et al.*, 2012), and to more complicated scenarios, such as large pedigrees with incomplete penetrance. These studies require new tools to enable the diverse community of researchers working on such families to interactively and comprehensively analyze next generation sequence data. Figure 1 shows how our new program, Olorin, integrates within-family linkage analysis with exome sequencing in a user-friendly package.

2 FEATURES

2.1 File formats

Olorin uses four types of data file: two that provide information about the gene flow calculated by MERLIN (Abecasis *et al.*, 2002), one defining the pedigree structure, and a list of variants identified by sequencing. MERLIN's haplotyping functionality is used to compute haplotype inheritance within the pedigree. Details of the genomic markers used in the estimation of haplotypes, and pedigree information about the relationships between individuals and their disease status are read from standard.map and.ped MERLIN format files. All variants identified from sequencing across samples need to be provided as a single variant call format (VCF) file (version 4.0 or greater) (Danecek *et al.*, 2011).

2.2 Workflow

2.2.1 Selecting individuals On loading data, Olorin automatically generates an interactive pedigree using standard conventions for information such as sex and disease status. Users can

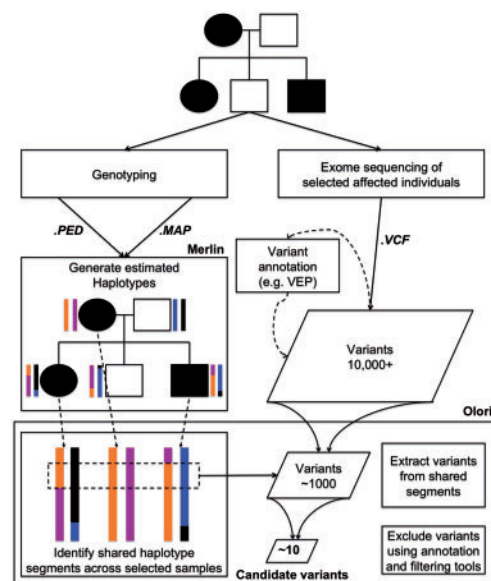


Fig. 1. Olorin uses patterns of gene flow estimated by MERLIN to identify genomic regions shared by affected individuals in large pedigrees. This information is combined with next generation sequence data, and only those variants that lie within shared regions are analysed. Users can further refine the list of variants using Olorin's realtime filtering tools

*To whom correspondence should be addressed.

