

Collapsed haplotype pattern method for linkage analysis of next-generation sequencing data

Gao T. Wang, Di Zhang, Biao Li, Hang Dai and Suzanne M. Leal[†]

[†]Center for Statistical Genetics, Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, Texas, USA



Motivation

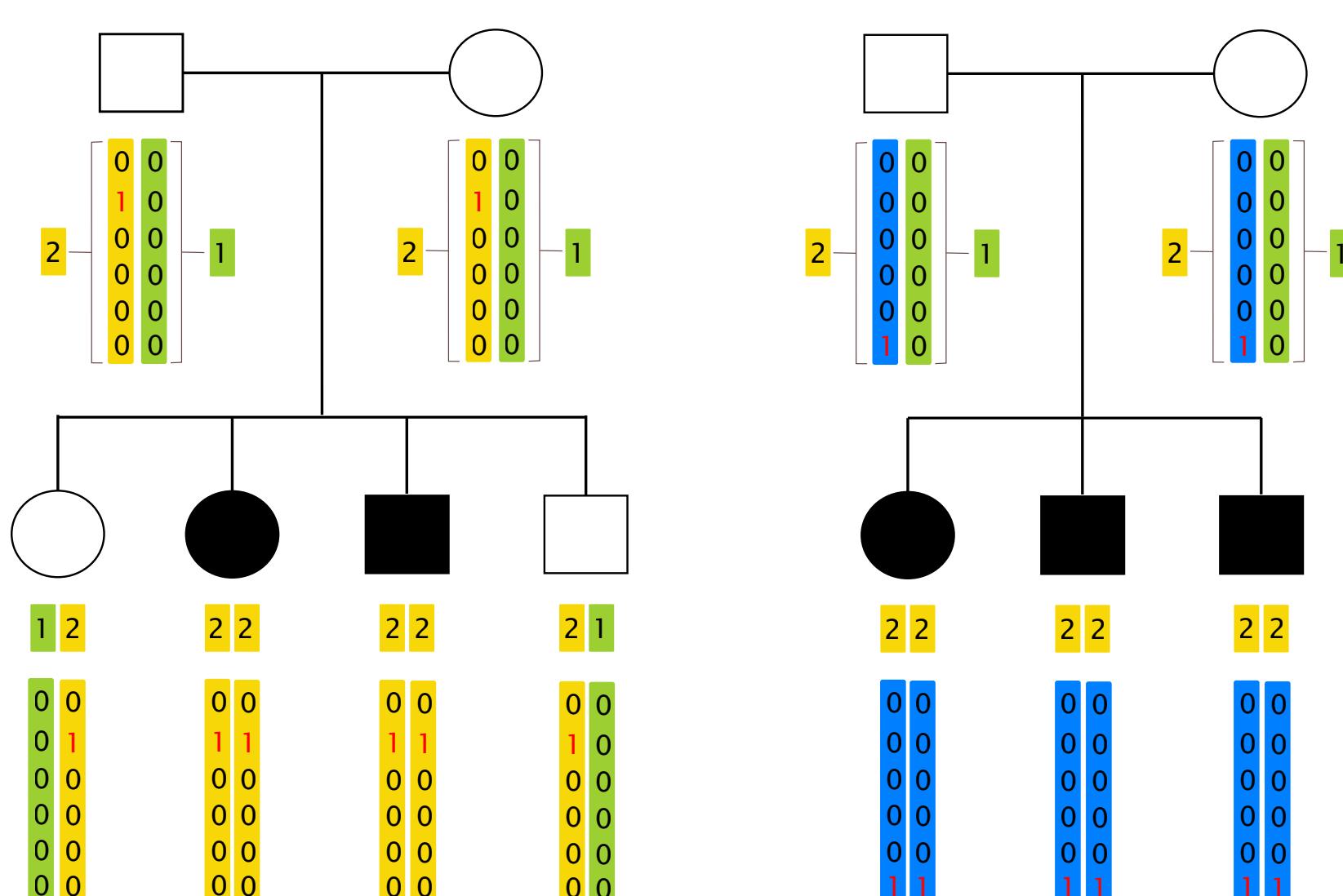
- Variant filtering using next-generation sequencing (NGS) data of families has successfully identified many causal mutations for Mendelian diseases, yet such approach
 - May result in many variants to follow-up
 - Is sensitive to mis-classification
 - Sample swaps, phenocopies, reduced penetrance
 - Do not provide statistical evidences
- Linkage analysis has been a powerful approach to map Mendelian disease loci using genetic marker data, but is under-powered when applied to sequence data, mainly due to lack of heterogeneity in single-nucleotide variants
- We developed a **Collapsed Haplotype Pattern** (CHP) method and a **SEQLinkage** software to utilize NGS data of multiple families for powerful linkage analysis

The CHP Method

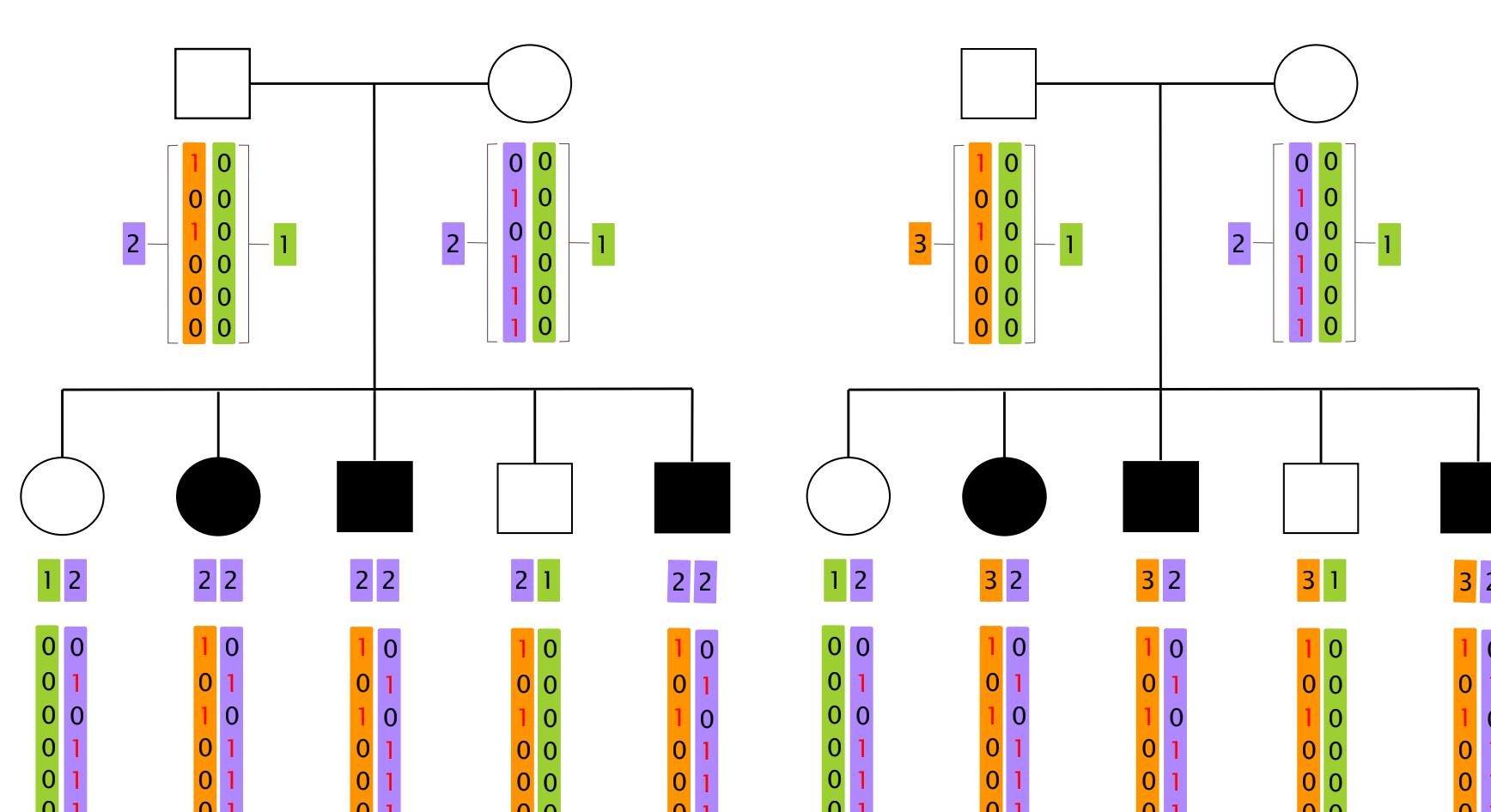
Inspired by rare variant "burden test" for association, the CHP method tests for linkage with a genetic region rather than with individual variants

- Regional markers** are generated for given genomic units
 - e.g. genes for exome sequence data
- CHP regional markers are **more heterozygous**
 - More informative than single variants in tracking the transmission of disease mutations within families, and when analyzing multiple pedigrees in the presence of allelic heterogeneity
- Tolerates missing data on causal variant sites
- No LD pruning required, all variants in data are used

Complete Collapsing



Complete vs. No Collapsing



Implementation

1. Family NGS data preprocessing

- Mendelian inconsistency check
- Haplotype reconstruction via genetic phasing

2. Regional marker construction

- Complete collapsing and LD based collapsing themes
- Sub-regional markers by recombination events

3. Annotation and statistics for linkage analysis

- Genetic distance interpolation via Rutgers Map
- Calculation of marker allele frequencies for linkage analysis with samples having no genotype data

4. Parametric test for linkage using regional markers

- LOD/HLOD scores computed over multiple pedigrees

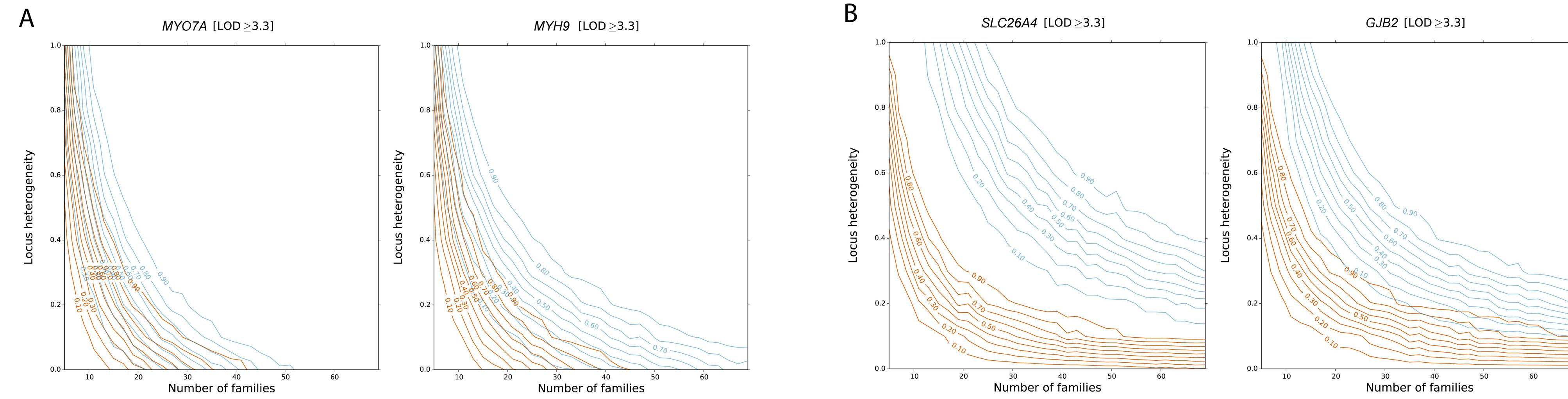
Power analyses and sample size estimations

Simulation study design for power analyses

4 nonsyndromic hearing impairment gene sequences simulated, using sequences from European American samples in NHLBI Exome Sequencing Project; causal variants determined by NCBI-Clinvar database; **two-generational pedigrees** generated with 3 to 8 offspring based on USA population demographic data; **full-penetrance** for causal variants, allowing for **allelic heterogeneity** and varying degrees of **locus heterogeneity**; two-point linkage analysis performed comparing **CHP vs. single variant linkage** methods; empirical power evaluated via 500 replicates

Power comparisons

CHP (orange curves) outperforms single variant linkage method (blue curves) under both dominant (A) and compound recessive (B) models. With 50% locus heterogeneity it requires **12 families** for CHP to achieve a power of 90% for *SLC26A4* at a genome-wide α level of 0.05, while single variant linkage method requires over **50 families**



Sample size estimations

Number of families required to achieve desired power are evaluated; **50% locus heterogeneity** is assumed for all scenarios; power calculations are based on HLOD score instead; impact of **missing causal variant** in families is evaluated

Required Power	Gene	MOI	CHP ^a	SNV ^b	CHP-M75% ^c	SNV-M75%
0.8	<i>SLC26A4</i>	recessive	11	40	39	160
0.9	<i>SLC26A4</i>	recessive	13	45	46	180
0.8	<i>SLC26A4</i>	compound recessive	11	50	39	200
0.9	<i>SLC26A4</i>	compound recessive	13	55	46	220
0.8	<i>GJB2</i>	recessive	12	23	44	92
0.9	<i>GJB2</i>	recessive	14	28	52	112
0.8	<i>GJB2</i>	compound recessive	12	25	44	100
0.9	<i>GJB2</i>	compound recessive	14	34	52	136
0.8	<i>MYO7A</i>	dominant	12	16	31	64
0.9	<i>MYO7A</i>	dominant	14	20	36	80
0.8	<i>MYH9</i>	dominant	11	13	32	52
0.9	<i>MYH9</i>	dominant	14	18	41	72

a. minimum number of families required to achieve desired power for CHP method

b. minimum number of families required to achieve desired power for single variant linkage method

c. "M75%": minimum number of families requirement when the causal variant in 75% families are missing

The SEQLinkage software

- Written in C++ with a Python parallel computing interface to rapidly scan through markers genome-wide
- Implements the CHP method for linkage analysis with sequence data of pedigrees in VCF format
- Supports output of CHP coded markers to formats compatible with other linkage programs
 - FASTLINK, MEGA2, Merlin, PLINK
- Performs two-point linkage analysis involving multiple families, maximizing linkage signals across families to allow for locus heterogeneity (HLOD score calculation)
- Provides results in text, graphical and table format output, organized with a user-friendly webpage interface
- Can be used in conjunction with variant filtering method to analyze sequence data of human pedigrees

SEQLinkage is freely available at

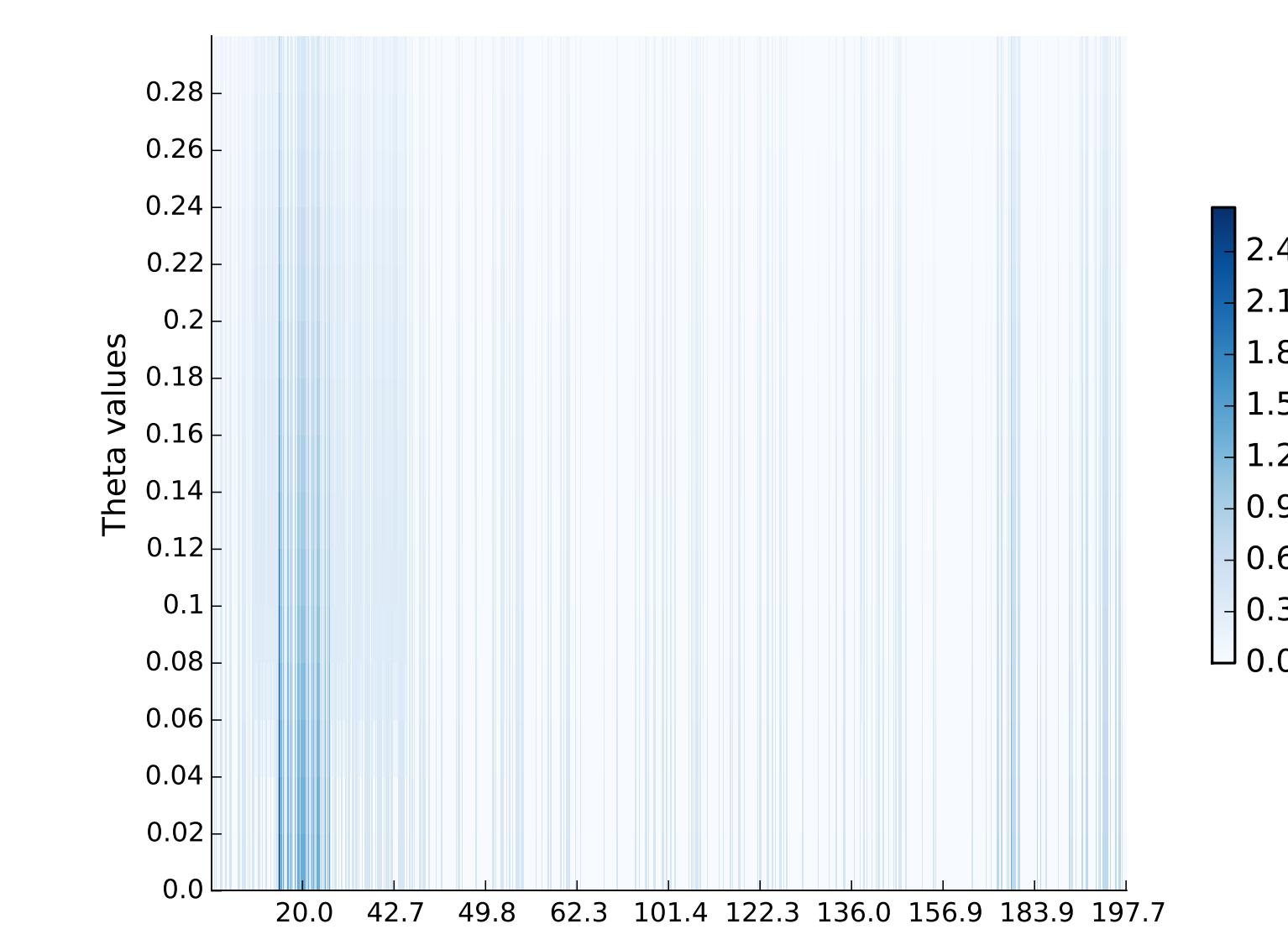
<http://bioinformatics.org/seqlink>

SEQLinkage analysis output

Whole exome sequences of 2 nuclear families each with 4 offspring (3 affected) are analyzed using **SEQLinkage**

0=0.0	0=0.02	0=0.04	0=0.06	0=0.08
Marker name chr:start-end	Marker name chr:start-end	Marker name chr:start-end	Marker name chr:start-end	Marker name chr:start-end
2.658 GRIP2 3:14530618-14583598	2.536 GRIP2 3:14530618-14583598	2.412 GRIP2 3:14530618-14583598	2.286 GRIP2 3:14530618-14583598	2.158 GRIP2 3:14530618-14583598
1.329UBE2Q1 1:154521050-154531120	1.276UBE2Q1 1:154521050-154531120	1.222UBE2Q1 1:154521050-154531120	1.166UBE2Q1 1:154521050-154531120	1.11UBE2Q1 1:154521050-154531120
1.329FANCF 11:22644078-22647387	1.276FANCF 11:22644078-22647387	1.222FANCF 11:22644078-22647387	1.166FANCF 11:22644078-22647387	1.11FANCF 11:22644078-22647387
1.329GRIN1 1:140033608-140063214	1.276GRIN1 1:140033608-140063214	1.222GRIN1 1:140033608-140063214	1.166GRIN1 1:140033608-140063214	1.11GRIN1 1:140033608-140063214
1.329FOXQ1 6:1312674-1314993	1.276FOXQ1 6:1312674-1314993	1.222FOXQ1 6:1312674-1314993	1.166FOXQ1 6:1312674-1314993	1.11FOXQ1 6:1312674-1314993
1.329SP3 5:52084135-52249485	1.276SP3 5:52084135-52249485	1.222SP3 5:52084135-52249485	1.166SP3 5:52084135-52249485	1.11SP3 5:52084135-52249485

Chromosome 3



Acknowledgments

We would like to thank Regie Lyn Santos-Cortez, Daniel Weeks, Alejandro Schaffer, Jeffrey O'Connell and Jürg Ott for helpful discussions and support. This work was funded by National Institute of Health grants DC003594, DC011651 and HG006493.