# HAPLOTYPE BLOCK PARTITIONING AND TAGSNP Selection using Mutual Information<sup>†</sup>

1 0

0 1 0

**1 0** 0

0 1 0

 $0 \ 0 \ 0$ 

0 0 1

 $0 \ 0 \ 1$ 

0 0 0

1 0 1

0 1 0

**0 1** 0

**1 1** 0

0 0 0

1 0

1 1 0

0 1 0

1 0 1

0 0 0

 $0 \ 0 \ 1$ 

 $0 \ 0 \ 1$ 

0 0 0

1 0  $0 \ 0 \ 0$ 

1 0

0 1 0

0 1 0

0 0

0 0

0 1 0

0 1 0

1 1 0

0 0 0

1

0 0 1

0 1 0

0

1 1 0 0 0

 $\mathbf{0}$ 1 0 **1 0** 0

 $0 \ 0 \ 0$ 

1 0 1

0 0 0 **1 0** 1

0 1 0

0 0 0

1 1 0

0 1 0

1 0 1

 $0 \ 0 \ 0$ 

0 1 0

0 1 0

 $0 \ 0 \ 0$ 

**1 0** 1

**1 0** 0

0 0 0

 $\mathbf{0}_{-}\mathbf{0}_{-}\mathbf{0}$ 

1 0 1

1 0

0 0

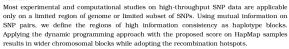
0 1 0

1

1

Katanforoush A. (1,2), Sadeghi M. (1,2), Pezeshk H. (3)

### Abstract



Keywords: Haplotype Blocks, HapMap, TagSNP selection



0

0 4

100110

1 0 0 1

 $0 \quad 0$ 

1 0 0

1 0 0

0 0 0

0 0 1

 $0 \ 0 \ 1$ 

 $0 \ 0 \ 0$ 

1 0 0

1 0 0

 $0 \ 0 \ 0$ 

0 0 0

0 0

 $0 \quad 0$ 

 $0 \quad 0$ 

 $0 \ 0 \ 0$ 

 $0 \ 0 \ 0$ 

0 1 0 0

1

1 0 0

1 1 0

1 0 1

0

0 0 0

1 1 1

0 0 1

1 0 0

1 **0 1** 

0 1 0 1

0 0 0

1 0 0

1 0

1  $0 \quad 0$ 

 $0 \quad 0$ 

0 1

1 1



## INTRODUCTION

Current high-throughput techniques in molecular genetics can determine millions of single nucleotide polymorphism on human genome. These data provide informative materials for disease sociation studies and analyzing models of population genetic Most practical and computational studies on these genome-wide data are applicable only on a limited region of genome or limited subset of SNPs. Using the mutual information of SNP pairs, we define the chromosomal regions of high information consistency as haplotype blocks. We have incorporated dynamic programming to find the block partitioning maximizing the sum of block scores. Also using the linear programming approaches for the problem of dominating set, we develop an efficient algorithm for tagSNP selection. The power of the case-control association test obtained by this tagging method also compared with some other common methods. It is shown that the proposed tagging SNP's improves levels of performance while holding the desired accuracy.

## **M**ETHOD

Essentially the mutual information reflects the genetic association between two SNPs. By doing some algebra, it would be shown that mutual information asymptotically tends to the logarithm of the Fisher exact test value.

$$I(i,j) = \sum_{\substack{a \in \{0,1\} \\ b \in \{0,1\}}} P(i=a \land j=b) . \log_2 \left( \frac{P(i=a \land j=b)}{P(i=a) . P(j=b)} \right)$$

$$F_{ex} = \frac{\binom{n_{1X}}{n_{11}} \binom{n_{0X}}{n_{01}}}{\binom{n}{n}}$$

evolutionary events such as selections, migrations and genetic drift. Such messy bunch of factors behaves like a noisy channel communicates words from one SNP to the other. Assuming that there is no other operation affecting SNPs then the capacity of such a channel is determined by the mutual information. Therefore high mutual information shows a strong associated pair of SNPs. The maximum value for the mutual information defined for bi-allelic SNPs is one bit per message and its minimum is zero.

## TagSNP selection

The general approach to define tagging SNPs aims at finding smallest subset of SNPs that any difference between every two distinct haplotypes can be captured by at least one tagSNP. These methods rely on some parsimonious assumption on haplotype diversity. Instead, we devise a method based on minimal dominating set from graph theory. By this approach every SNP either is a tagSNP or is in strong association with at least one tagging SNP.

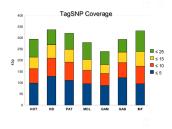




Among methods for haplotypes block partitioning, we have selected some of the most cited (Table 1). Also, recombination hotspots can realize boundaries of haplotype blocks on genome. In hotspot regions the recombination rate is high and the probability of observing less diverse haplotype is low. Hotspots for all individuals are not necessarily the same. However, there are some regions on genome which could be recognized as the consensus

# RESULT

It is observed that the blocks produced by LD-based methods like Gabriel. Infres also the blocks defined by hotspots are broader than what determined by the greedy algorithm and 4-gamete test.



In general, the methods producing broader blocks entail less number of tagSNPs. It reveals that there is not any essential relation between the minimal number of tagSNPs and grouping high associated SNPs within blocks. The higher numbers corresponded to the wider blocks indicate that a usual tagSNP is capable of capturing much longer haplotypes than those broken by diversity criteria. The conventional methods based on block diversity usually generate more



# The sum of all mutual information of SNP pairs in a block is introduced as score after subtracting the sum by a constant factor which bounds the type I error.

$$S(i;d) = \sum_{s,d \in S(d)} f(s,t) - \theta$$

For each interval of length d ending at SNP i, S(i;d) denotes the ore of the block and shows how strong a certain set of SNPs tends to be integrated into a block.

The free parameter,  $\theta \in [01]$  incorporated in the formula allows to alter average length of blocks by the following optimization step. For θ's close to zero, the algorithm tends to produce wider blocks while for  $\theta$ 's close to one, it produces smaller blocks same as the 4-gameter test approach.

The optimum solution is obtained using a dynamic programming approach. The number of fundamental operations to calculate blocks scores and to find optimal partitioning is O(nwL) where n is the number of haplotype samples,  $\boldsymbol{w}$  is the window size, and  $\boldsymbol{L}$  is the number of SNPs. In general, we refer to the method described above

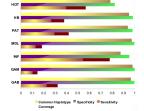
 $x_j = 1$  if SNP j is selected as a tag SNP  $Min \sum_{j=1}^{\infty} x_{j}$ s.t.  $\sum_{i=1}^{n} f_{ij} x_{j} \ge k$ , i=1,...,m

Method	Block criteria	Implement at ior	Literature
Patil's	diversity / greedy algorithm	HapBlock v.3	Patil et al. 2001
HapBlock	diversity / dynamic programming	HapBlock v.3	K. Zhang and L. Jin 2003
Gabriel's	pairwise LD	Haploview v.4	Gabriel et al. 2002
4 gamete test	perfect phylogeny	Haploview v.4	Wang et al. 2002
MDBlock	Minimum Description Length	MDBlock v.1	Anderson and Novembre 200
Hotspots	Fine-scale LD map		Myres et al. 2005

Table 1: Haplotype Block partitioning methods

(i) Department of Bioinformatics, Institute of Biochemistry and Biophysis, University of Biotun, Refran, Fors.
(2) Institute for Studies in Theoretical Physics and Madesman; (1974), Naturan Spatra, Patric, Prince, Prince,

# wer of Haplotype Block Partitioning Algorithms



# CONCLUSION

Finding minimum number of tagSNPs can pose a two fold object. The first is the genome partitioning in a way that reported blocks could be captured by as limited number of tagSNPs as possible. The second object is to get fewer numbers of tagSNPs to be associated with all SNPs. Results obtained by available methods show that in practice there is a gap between two extremes. The algorithms aimed to recognize evolutionary conserved haplotype blocks such as greedy diversity based algorithms, four-gamete-test, and Gabriel's LD based method obtain the blocks which are tagged by just few number of tagSNPs, but they assign too many tagSNPs for the whole chromosome. On the other hand, the blocks being introduced by hotspots neglect the diversity considerations while giving less number of tagSNPs. The information theoretic method, Infres, seems to cover this gap by satisfying both the objects on average.

Blefantane and November L (2001) Finding Hashington Bleck Brandines by Using the Maintann Description Assessment of the State Control of Hashington Contr Wing, N. et al. (2020) Distribution of recombination crossovers and the origin of haplotype blocks: the interplay long, N. et al. (2020) Distribution of recombination crossovers and the origin of haplotype blocks: the interplay population bistory, recombination, and mutation. Ann. J. Hant. Georg. 71, 1227–34.
Zhang, K. and Jini, L. (2000) Hiptelfolds/Finder: Haplotype block analyses. Bioinformaries, 19, 1300-1301.