

Bioinformatics, YYYY, 0–0

doi: 10.1093/bioinformatics/xxxxx

Advance Access Publication Date: DD Month YYYY

Original papers

---

*Genome analysis*

# SNP Interaction Pattern Identifier (SIPI): An Intensive Search for SNP-SNP Interaction Patterns

Hui-Yi Lin<sup>1,\*</sup>, Dung-Tsa Chen<sup>2</sup>, Po-Yu Huang<sup>3</sup>, Yung-Hsin Liu<sup>4</sup>, Augusto Ochoa<sup>5</sup>, Jovanny Zabaleta<sup>5</sup>, Donald E. Mercante<sup>1</sup>, Zhide Fang<sup>1</sup>, Thomas A. Sellers<sup>6</sup>, Julio Pow-Sang<sup>7</sup>, Chia-Ho Cheng<sup>2</sup>, Rosalind Eeles<sup>8,9</sup>, Doug Easton<sup>10</sup>, Zsofia Kote-Jarai<sup>8</sup>, Ali Amin Al Olama<sup>10</sup>, Sara Benlloch<sup>10</sup>, Kenneth Muir<sup>11</sup>, Graham G. Giles<sup>12, 13</sup>, Fredrik Wiklund<sup>14</sup>, Henrik Gronberg<sup>14</sup>, Christopher A. Haiman<sup>15</sup>, Johanna Schleutker<sup>16,17</sup>, Børge G. Nordestgaard<sup>18</sup>, Ruth C. Travis<sup>19</sup>, Freddie Hamdy<sup>20</sup>, Nora Pashayan<sup>21, 38</sup>, Kay-Tee Khaw<sup>22</sup>, Janet L. Stanford<sup>23, 24</sup>, William J. Blot<sup>25</sup>, Stephen N. Thibodeau<sup>26</sup>, Christiane Maier<sup>27</sup>, Adam S. Kibel<sup>28, 29</sup>, Cezary Cybulski<sup>30</sup>, Lisa Cannon-Albright<sup>31</sup>, Hermann Brenner<sup>32, 33, 39</sup>, Radka Kaneva<sup>34</sup>, Jyotsna Batra<sup>35</sup>, Manuel R. Teixeira<sup>36</sup>, Hardev Pandha<sup>37</sup>, Yong-Jie Lu<sup>40</sup>, the PRACTICAL consortium<sup>41</sup>, and Jong Y. Park<sup>6</sup>

<sup>1</sup> Biostatistics Program, School of Public Health, Louisiana State University Health Sciences Center, New Orleans, LA 70112, USA. <sup>2</sup> Department of Biostatistics and Bioinformatics, Moffitt Cancer Center & Research Institute, Tampa, FL 33612, USA. <sup>3</sup> Computational Intelligence Technology Center, Industrial Technology Research Institute, Hsinchu City, Taiwan. <sup>4</sup> Department of Biometrics, INC Research, LLC, Raleigh, NC 27609, USA. <sup>5</sup> Stanley S. Scott Cancer Center, Louisiana State University Health Sciences Center, New Orleans, LA 70112, USA. <sup>6</sup> Department of Cancer Epidemiology, Moffitt Cancer Center & Research Institute, Tampa, FL 33612, USA. <sup>7</sup> Department of Genitourinary Oncology, Moffitt Cancer Center & Research Institute, Tampa, FL 33612, USA. <sup>8</sup> The Institute of Cancer Research, London, SM2 5NG, UK, <sup>9</sup> Royal Marsden NHS Foundation Trust, London, SW3 6JJ, UK, <sup>10</sup> Centre for Cancer Genetic Epidemiology, Department of Public Health and Primary Care, University of Cambridge, Strangeways Research Laboratory, Worts Causeway, Cambridge, UK, <sup>11</sup> University of Warwick, Coventry, UK, <sup>12</sup> Cancer Epidemiology Centre, The Cancer Council Victoria, 615 St Kilda Road, Melbourne, Victoria, 3004, Australia, <sup>13</sup> Centre for Epidemiology and Biostatistics, Melbourne School of Population and Global Health, The University of Melbourne, Melbourne, Victoria, Australia, <sup>14</sup> Department of Medical Epidemiology and Biostatistics, Karolinska Institute, Stockholm, Sweden, <sup>15</sup> Department of Preventive Medicine, Keck School of Medicine, University of Southern California/Norris Comprehensive Cancer Center, Los Angeles, California, USA, <sup>16</sup> Department of Medical Biochemistry and Genetics, Institute of Biomedicine, Kiinamylyynkatu 10, FI-20014 University of Turku; and Tyks Microbiology and Genetics, Department of Medical Genetics, Turku University Hospital, <sup>17</sup> BioMediTech, 30014 University of Tampere, Tampere, Finland, <sup>18</sup> Department of Clinical Biochemistry, Herlev Hospital, Copenhagen University Hospital, Herlev Ringvej 75, DK-2730 Herlev, Denmark, <sup>19</sup> Cancer Epidemiology, Nuffield Department of Population Health University of Oxford, Oxford, UK, <sup>20</sup> Nuffield Department of Surgical Sciences, University of Oxford, Oxford, UK, Faculty of Medical Science, University of Ox-

---

ford, John Radcliffe Hospital, Oxford, UK, <sup>21</sup> Centre for Cancer Genetic Epidemiology, Department of Oncology, University of Cambridge, Strangeways Research Laboratory, Worts Causeway, Cambridge, UK, <sup>22</sup> Cambridge Institute of Public Health, University of Cambridge, Forvie Site, Robinson Way, Cambridge CB2 0SR, UK, <sup>23</sup> Division of Public Health Sciences, Fred Hutchinson Cancer Research Center, Seattle, Washington, USA, <sup>24</sup> Department of Epidemiology, School of Public Health, University of Washington, Seattle, Washington, USA, <sup>25</sup> International Epidemiology Institute, 1455 Research Blvd., Suite 550, Rockville, MD 20850, USA, <sup>26</sup> Department of Laboratory Medicine and Pathology, Mayo Clinic, Rochester, MN, USA, <sup>27</sup> Institute of Human Genetics University Hospital Ulm, Germany, <sup>28</sup> Brigham and Women's Hospital/Dana-Farber Cancer Institute, USA, 45 Francis Street- ASB II-3, Boston, MA 02115, USA, <sup>29</sup> Washington University, St Louis, Missouri, USA, <sup>30</sup> International Hereditary Cancer Center, Department of Genetics and Pathology, Pomeranian Medical University, Szczecin, Poland, <sup>31</sup> Division of Genetic Epidemiology, Department of Medicine, University of Utah School of Medicine, USA, <sup>32</sup> Division of Clinical Epidemiology and Aging Research, German Cancer Research Center (DKFZ), Heidelberg, Germany, <sup>33</sup> Division of Preventive Oncology, German Cancer Research Center (DKFZ) and National Center for Tumor Diseases (NCT), Heidelberg, Germany, <sup>34</sup> Molecular Medicine Center and Department of Medical Chemistry and Biochemistry, Medical University - Sofia, 2 Zdrave St, 1431, Sofia, Bulgaria, <sup>35</sup> Australian Prostate Cancer Research Centre-Qld, Institute of Health and Biomedical Innovation and Schools of Life Science and Public Health, Queensland University of Technology, Brisbane, Australia, <sup>36</sup> Department of Genetics, Portuguese Oncology Institute, Porto, Portugal and Biomedical Sciences Institute (ICBAS), Porto University, Porto, Portugal, <sup>37</sup> The University of Surrey, Guildford, Surrey, GU2 7XH, UK, <sup>38</sup> University College London, Department of Applied Health Research, 1-19 Torrington Place, London, WC1E 7HB, UK, <sup>39</sup> German Cancer Consortium (DKTK), German Cancer Research Center (DKFZ), Heidelberg Germany, <sup>40</sup> Centre for Molecular Oncology, Barts Cancer Institute, Queen Mary University of London, John Vane Science Centre, Charterhouse Square, London, EC1M 6BQ, UK, <sup>41</sup> Additional members from the Prostate Cancer Association Group to Investigate Cancer Associated Alterations in the Genome (PRACTICAL) consortium to be provided in the supplement.

\*To whom correspondence should be addressed.

Associate Editor: XXXXXXXX

Received on XXXXXX; revised on XXXXXX; accepted on XXXXXX

## Abstract

**Motivation:** SNP-SNP interactions may be the key for overcoming bottlenecks of genetic association studies. However, related statistical methods for testing SNP-SNP interactions are underdeveloped.

**Results:** We propose the SNP Interaction Pattern Identifier (SIPI), which tests 45 biologically meaningful interaction patterns for a binary outcome. SIPI takes various inheritance modes and model structures (including non-hierarchical models) into consideration. The simulation results show that SIPI has higher power than MDR-LR (Multifactor Dimensionality Reduction-Logistic Regression), AA\_Full, and SNPAssoc in general. Applying SIPI to the prostate cancer PRACTICAL consortium data with approximately 21,000 patients, the two SNP pairs in *EGFR-MMP16* and *EGFR-EGFR* were found to be associated with prostate cancer aggressiveness with the exact pattern in the discovery and validation sets. We demonstrated that SIPI not only searches for more meaningful interaction patterns but can also overcome the unstable nature of interaction patterns.

**Availability:** The SIPI software is freely available at <http://publichealth.lsuhsu.edu/LinSoftware/>.

**Contact:** hlin1@lsuhsu.edu

**Supplementary information:** Supplementary figures and tables are available at Bioinformatics online.

---

## 1 Introduction

During the past decade, the genome-wide association studies (GWAS) have successfully identified many inherited genetic variants (or SNPs) associated with complex diseases, such as cancer or related phenotypes. However, the predictive power of cancer risk for the GWAS-identified SNPs is small by a 1.2 median per-allele odds ratio (Ioannidis, et al., 2010). By combining multiple SNPs in a prediction model, the predictive power of these GWAS SNPs can be improved (Van den Broeck, et al., 2014). We recently reported the polygenic genetic models to estimate their risk for prostate cancer (Al Olama, et al., 2014; Amin Al Olama, et al., 2015; Eeles, et al., 2013). Despite these efforts, major proportion of familial risk of prostate cancer remains unknown. The similar situation applies for using SNPs to predict prostate cancer prognosis (Van den Broeck, et al., 2014). It is well known that biological associations among genes are complicated. The majority of GWAS focus on identification of individual SNP effects, which may not be sufficient to explain the complexity of disease causality. It has been shown that gene-gene/SNP-SNP interactions play an important role in the etiology of complex diseases (Cordell, 2009; Moore, 2003; Moore and Williams, 2002; Onay, et al., 2006). Although SNP-SNP or gene-gene interaction studies have been emerging, the statistical methods for evaluating SNP-SNP interactions are still underdeveloped.

The majority of genetic association studies focus on two-way interactions with two SNPs involved. In the past decade, various statistical methods have been proposed for evaluating two-way SNP-SNP interactions. These methods can be classified either model-based or pattern-based. The most common model-based approach tests an interaction based on a full interaction model with both main effects and their interaction. Examples include PLINK (Purcell, et al., 2007), SNPAssoc (Gonzalez, et al., 2007) and Boolean Operation-based Screening and Testing (BOOST) (Wan, et al., 2010). For the model-based approaches, the impact of an interaction can be distinguished from the main effects, but the number of detectable interaction patterns is limited. In the pattern-based approach, interaction detection is based on risk patterns of the 3x3 genotype combination table of the two SNPs [such as AA+BB/Bb vs. others for SNPA (major/minor allele: A/a) and SNPB (B/b)]. The Hypothesis Free Clinical Cloning (HFCC) tests for 255 patterns for one SNP pair (Gayan, et al., 2008), but some patterns may not be biologically meaningful or are rare. SNPmaxel evaluates 16 interaction patterns and four main effects for a given SNP pair (Boulesteix, et al., 2007). Multifactor dimensionality reduction (MDR) is also a pattern-based approach. MDR generates a binary risk variable (high/low risk) by comparing the case-to-control ratio in each genotype combination to a threshold and classifies each genotype to either a high risk set or low risk set. The K-fold cross-validation is used to relieve over-fitting issue (Ritchie, et al., 2003; Ritchie, et al., 2001). The strength of pattern-based approaches is that they are designed to detect wider range of interaction patterns. The limitation of these approaches is that they search associations that allow for but are not limited to interactions. A significant result detected using the pattern-based approaches may be due to strong main effect without an interaction.

To overcome these weaknesses, we propose SNP Interaction Pattern Identifier (SIPI), which combines the advantages of the model-based and pattern-based approaches. Our approach can examine 45 interaction models that consider biologically meaningful factors. Each model has a straightforward corresponding pattern, and there is a formal statistical test for evaluating the interaction effect only. This approach is powerful, and the identified patterns can be easily applied to assemble risk-prediction models. For evaluating the performance of SIPI, we conducted a simulation study to evaluate power and type I errors of SIPI with other three approaches: MDR-LR (Multifactor dimensionality reduction-Logistic Regression), AA\_full and SNPAssoc. MDR-LR (Multifactor dimensionality reduction- Logistic Regression) is a modified version of MDR for formal testing an interaction (Edwards, et al., 2010).

## 2 Methods

### 2.1 SNP Interaction Pattern Identification (SIPI)

SIPI can intensively and effectively search pairwise SNP-SNP interactions. The conventional approach for identifying SNP-SNP interaction is to search a specific type of interaction using the full interaction model with the additive-additive mode based on the minor allele. The SIPI detects 45 interaction models, which take inheritance mode (both original and reverse), and risk category grouping (model structure) into consideration. The best interaction pattern is selected based on the Bayesian information criterion (BIC), which is used to deal with the trade-off between model fit and complexity of the model. BIC is also shown to be consistent in selecting the true model and tends to select a parsimonious model compared with the Akaike information criterion (AIC), especially in studies with a large sample size (Yang, 2005). Based on these features, we decided use BIC as the selection criteria in SIPI. The concept of SIPI can be applied to different types of outcomes, such as numeric, binary and time-to-event variables. In this study, we focused on the binary outcome using logistic regression models. The two primary components of SIPI are introduced separately below.

#### 2.1a SNP Inheritance Modes

The SNP inheritance modes can impact on power to detect SNP interactions (Lin, et al., 2008). We designate a lowercase letter 'a' to denote the minor (low frequency) allele, and an uppercase 'A' to denote the major (common) allele. Each SNP has three genotype categories: homozygous major type ('AA'), heterozygous type ('Aa') and homozygous minor type ('aa'). For a SNP, the inheritance mode for a disease risk refers to a specific relationship between genotype and phenotype. The inheritance modes include additive, dominant, recessive, genotypic and over-dominant modes. The dominant mode assumes that the impact of having one or two copies of a given allele on the outcome is the same, and the recessive mode implies that the subjects with only homozygous genotypes of a given allele have a higher risk to develop the outcome. Additive mode refers to the impact of each additional copy of a given allele on the outcome being equal. The genotypic mode, treats a SNP as a categorical variable with three groups, and assumes that each genotype has a distinct effect on risk. This genotypic mode needs four degrees of freedom for the interaction term itself, and interpretation of the result is not straightforward. The over-dominant mode, which assumes that heterozygote has a different risk than the other two homozygous genotypes (Aa vs. AA/aa), is a rare case. Therefore, we excluded genotypic and over-dominant mode and evaluated a total of three inheritance modes (dominant, recessive and additive) in this study.

In the majority of genetic association studies, inheritance modes are defined based on the minor (or variant) allele. Under this scenario, the binary inheritance mode (dominant and recessive) is coded as "1" for the group containing the homozygous minor type, and the other group as "0" in modeling. For the AA, Aa and aa genotypes, the additive mode coding is 0, 1 and 2. The reverse coding (=1 - original coding for dominant and recessive mode; and 2-original coding for additive mode) of inheritance mode is seldom to be considered in testing SNP-SNP interactions. The original/reverse coding of inheritance mode does not impact on statistical significance (p-values) for testing individual SNP effects and full interaction model, but dramatically impacts testing SNP-SNP interactions in non-hierarchical interaction models. As shown in Table 1, there are six total possible coding methods for inheritance modes for each SNP. The three inheritance modes with the original coding based on the minor allele are additive (noted as aSNP1 for SNP1), dominant (dSNP1) and recessive (rSNP1). For reverse inheritance modes, the three modes are reverse additive (raSNP1), reverse dominant (rdSNP1), and reverse recessive (rrSNP1).

Table 1. SNP coding scheme by the SNP comparative allele and inheritance mode

SNP1 Maj/Min <sup>1</sup> =A/a	Original mode <sup>2</sup>			Reverse mode <sup>2</sup>		
	Additive (aSNP1)	Dominant (dSNP1)	Recessive (rSNP1)	Reverse Additive (raSNP1)	Reverse Dominant (rdSNP1)	Reverse Recessive (rrSNP1)
AA	0	0	0	2	1	1
Aa	1	1	0	1	0	1
aa	2	1	1	0	0	0
Data type	Continuous	Binary	Binary	Continuous	Binary	Binary

<sup>1</sup>Maj/Min= major/minor allele

<sup>2</sup>Original modes are based on a minor allele 'a'; Reverse modes (1 - original mode) for the dominant and recessive mode and (2 - original mode) for the additive mode

### 2.1b Risk Category Grouping/Model Structure

Both hierarchical and non-hierarchical interaction model were considered in this study. For evaluating 2-way interactions, the hierarchical or full interaction models are the models with two main effects and their interactions. This is the most common model type for testing pairwise SNP-SNP interactions, but this full model tests only one specific interaction pattern. Non-hierarchical models are defined as models with an interaction, and none or one main effects. Studies show pure interactions without main effects are possible in genetic association studies (Lin, et al., 2013; Lin, et al., 2008), so a non-hierarchical model provides a flexible tool to evaluate these interaction patterns. Using non-hierarchical models, a data-driven parsimonious model can be generated; therefore power of detecting these specific interaction patterns increases (Milne, et al., 2008; Piegorsch, et al., 1994). As shown in Equations 1-4, four possible model structures for testing a two-way interaction include models with (1) two main effects plus an interaction (Full-int); (2) the main effect of variable 1 plus an interaction (Main1+int); (3) the main effect of variable 2 plus an interaction (Main2+int); and (4) an interaction only (Int-only).

By considering a binary inheritance mode, there are four inheritance mode combinations (dominant-dominant, dominant-recessive, recessive-dominant and recessive-recessive). When treating SNPs as numeric variables, the additive-additive mode is taken into consideration. Thus, SIPI considers a total of five possible types of inheritance mode combinations. For each inheritance mode combination, there are nine unique interaction models/patterns when taking into consideration different model structures and comparative alleles. Thus, a total of 45 interaction patterns are considered in SIPI for each SNP pair (Table 2).

The best model among the 45 models is based on the lowest value of the Bayesian information criterion (BIC) (Schwarz, 1978). The significance of the interaction effect is tested using the Wald test of the interaction term (H0:  $\beta_3=0$ ). Although the likelihood ratio test (LRT) is usually recommended as the most powerful approach, it requires performing the two models one wishes to compare. The Wald test is similar to LRT in large scale studies and only one model needs to be estimated. In order to ease computation burden for high-dimensional data, the Wald test was primarily used in SIPI. In the SIPI R package, the users can choose to report p-values based on the Wald test or LRT. The Bonferroni method is applied to adjust for multiple comparisons.

Full interaction model (Full-int):

$$\text{logit}[\text{pr}(Y = 1)] = \beta_0 + \beta_1\text{SNP}_1 + \beta_2\text{SNP}_2 + \beta_3\text{SNP}_1 \times \text{SNP}_2 \quad (\text{eq. 1})$$

Main 1+ interaction (Main1+int):

$$\text{logit}[\text{pr}(Y = 1)] = \beta_0 + \beta_1\text{SNP}_1 + \beta_3\text{SNP}_1 \times \text{SNP}_2 \quad (\text{eq. 2})$$

Main 2+ interaction (Main2+int):

$$\text{logit}[\text{pr}(Y = 1)] = \beta_0 + \beta_2\text{SNP}_2 + \beta_3\text{SNP}_1 \times \text{SNP}_2 \quad (\text{eq. 3})$$

Interaction only (Int-only):

$$\text{logit}[\text{pr}(Y = 1)] = \beta_0 + \beta_3\text{SNP}_1 \times \text{SNP}_2 \quad (\text{eq. 4})$$

,where Y is the binary outcome with a value of 0 or 1.

Table 2. List of 45 interaction models by considering the inheritance modes and model structures

SNP1x SNP2 Inheritance mode <sup>1</sup>	Model structure <sup>2</sup>	Model label <sup>3</sup>	Model Details	
Dom-Dom	Full-int	DD_Full	dSNP1 + dSNP2 + dSNP1x dSNP2	
	Main1+int	DD_M1_int_o1	dSNP1 + dSNP1x dSNP2	
		DD_M1_int_r1	rdSNP1 + rdSNP1x dSNP2	
	Main2+int	DD_M2_int_o2	dSNP2 + dSNP1x dSNP2	
		DD_M2_int_r2	rdSNP2 + rdSNP1x dSNP2	
	Int-only	DD_int_oo	dSNP1x dSNP2	
		DD_int_or	dSNP1x rdSNP2	
		DD_int_ro	rdSNP1x dSNP2	
		DD_int_rr	rdSNP1x rdSNP2	
	Dom-Rec	Full-int	DR_Full	dSNP1 + rSNP2 + dSNP1x rSNP2
		Main1+int	DR_M1_int_o1	dSNP1 + dSNP1x rSNP2
			DR_M1_int_r1	rdSNP1 + rdSNP1x rSNP2
Main2+int		DR_M2_int_o2	rSNP2 + dSNP1x rSNP2	
		DR_M2_int_r2	rrSNP2 + dSNP1x rrSNP2	
Int-only		DR_int_oo	dSNP1x rSNP2	
		DR_int_or	dSNP1x rrSNP2	
		DR_int_ro	rdSNP1x rSNP2	
		DR_int_rr	rdSNP1x rrSNP2	
Rec-Dom		Full-int	RD_Full	rSNP1 + dSNP2 + rSNP1x dSNP2
		Main1+int	RD_M1_int_o1	rSNP1 + rSNP1x dSNP2
			RD_M1_int_r1	rrSNP1 + rrSNP1x dSNP2
	Main2+int	RD_M2_int_o2	dSNP2 + rSNP1x dSNP2	
		RD_M2_int_r2	rdSNP2 + rdSNP1x dSNP2	
	Int-only	RD_int_oo	rSNP1x dSNP2	
		RD_int_or	rSNP1x rdSNP2	
		RD_int_ro	rdSNP1x dSNP2	
		RD_int_rr	rdSNP1x rdSNP2	
	Rec-Rec	Full-int	RR_Full	rSNP1 + rSNP2 + rSNP1x rSNP2

Main1+int	RR_M1_int_o1	rSNP1 + rSNP1x rSNP2	
	RR_M1_int_r1	rrSNP1 + rrSNP1x rSNP2	
Main2+int	RR_M2_int_o2	rSNP2 + rSNP1x rSNP2	
	RR_M2_int_r2	rrSNP2 + rrSNP1x rSNP2	
Int-only	RR_int_oo	rSNP1x rSNP2	
	RR_int_or	rSNP1x rrSNP2	
	RR_int_ro	rrSNP1x rSNP2	
	RR_int_rr	rrSNP1x rrSNP2	
Add_Add	Full-int	AA_Full	aSNP1 + aSNP2 + aSNP1x aSNP2
	Main1+int	AA_M1_int_o1	aSNP1 + aSNP1x aSNP2
Main2+int	AA_M1_int_r1	raSNP1 + raSNP1x aSNP2	
	AA_M2_int_o2	aSNP2 + aSNP1x aSNP2	
Int-only	AA_M2_int_r2	raSNP2 + aSNP1x raSNP2	
	AA_int_oo	aSNP1x aSNP2	
	AA_int_or	aSNP1x raSNP2	
	AA_int_ro	raSNP1x aSNP2	
	AA_int_rr	raSNP1x raSNP2	

<sup>1</sup>Dom: dominant, Rec: recessive, Add: additive

<sup>2</sup> Full-int: full interaction model with two main effects plus an interaction; Main1+int: main effect of variable 1 plus an interaction; Main2+int: main effect of variable 2 plus an interaction; and (4) Int-only: an interaction only.

<sup>3</sup> \_o1, \_r1: minor allele (original coding), and reverse coding of SNP1;

\_o2, \_r2: minor allele (original coding), and reverse coding of SNP2;

\_oo, \_or, \_ro, \_rr: based on original-original, original-reverse, reverse-original and reverse-reverse coding for SNP1 and SNP2

### 2.1c Translating Interaction Models to Interaction Patterns

By treating SNPs as binary variables (such as dominant or recessive), we can simplify genotype combinations from a three-by-three panel into a two-by-two panel, resulting in four possible sub-groupings. For the two-by-two panel, we can categorize the genotype combinations to four-, three- and two-risk subgroups. As shown in Figures S1-S2, we can translate the interaction models to the corresponding genotype interaction patterns. Our 45 pattern labels were based on the three-by-three tables with an order of homozygous wild, heterozygous and homozygous variant types (denote as AA, Aa, and aa) and the homozygous major genotypes of the two SNPs as the top left corner.

### 2.2 MDR-LR

MDR-LR, a two-step approach (Edwards, et al., 2010), is formally to test an interaction for the MDR selected interaction models. The first step is to apply the MDR concept to classify each SNP to a binary variable (high/low risk). In the 2<sup>nd</sup> step, a full logistic regression with main effects of these binary SNP variables and their interaction is conducted. The significance of interaction is tested based on the likelihood ratio test of the interaction term.

### 2.3 AA\_Full

The AA\_full [available in PLINK (Purcell, et al., 2007)] approach uses a full logistic regression model with both main effect and interaction. Each SNP is treated as an additive mode based on the minor allele. The significance test is evaluated using the Wald test of the interaction coefficient.

### 2.4 SNPassoc

SNPassoc (Gonzalez, et al., 2007) used the same full logistic regression and allows for five different inheritance modes [additive, dominant, recessive, genotypic, and over-dominant (Aa vs. AA/aa)] based on the minor allele. Two SNPs in the same pair are required to have the same inheritance mode.

### 2.5 Simulation

We conducted a simulation study to compare the power of SIPI with the conventional AA\_Full model, MDR-LR, and SNPassoc approach for detecting two-way SNP-SNP interactions. For simulation settings, one SNP pair was considered. The two candidate SNPs were generated independently based on the Hardy-Weinberg equilibrium. Seven sets of a wide range of minor allele frequencies (MAF=0.05-0.5) for SNP1 and SNP2 were investigated: (0.5, 0.3), (0.5, 0.2), (0.5, 0.05), (0.3, 0.3), (0.3, 0.1), (0.3, 0.05), and (0.1, 0.05). The sample sizes of 1,000 and 5,000 were chosen. All analyses were based on 1,000 simulation runs.

The binary outcome variable (such as case/control) was generated based on outcome prevalence and the proportion of the value of interest (such as disease) in each genotype combination of the two given SNPs using multinomial distribution. We evaluated a total of six designed interaction patterns, including one real-data pattern (Figures 1-2). Most of these simulated models are based on the interaction patterns reported previously (Lin, et al., 2013; Lin, et al., 2012). One null model without an interaction term was also tested. For the effect size of Models 1-4, the outcome prevalence was set to 0.3 or 0.4 in the high-risk subgroups and was 0.2 in the low-risk sub-groups. The corresponding odds ratio (OR) is 1.6

## SNP Interaction Pattern Identifier (SIPI)

and 2.7, respectively. The settings of true interaction models are listed in Figures 1-2.

Models 1-3 were interaction-only models. For Model 1 (RR\_int\_rr pattern), both SNPs are considered as recessive with the reverse coding. The disease prevalence is 0.3 and 0.2 for the high- and low-risk groups, respectively. For Model 2 (DD\_int\_oo pattern), both SNPs are considered as dominant based on the minor alleles. In Model 3 (RD\_int\_rr), SNP1 is considered under a recessive mode, SNP2 is considered as dominant mode, and both SNPs have the reverse coding. Model 4 (DD\_M1\_int\_o1) includes the SNP1 main effect and an interaction, in which both SNPs are considered as dominant based on the minor allele of SNP1. The significance of the interaction term is the same regardless of the inheritance mode coding (original or reverse) for SNP2. Model 5 (AA\_Full) is a full interaction model and both SNPs are treated as an additive mode based on the minor allele. This AA\_full model has the setting of  $\beta_0 = -2.5$  and  $\beta_1 = \beta_2 = \beta_3 = 0.6$ . Model 6 (RD\_int\_oo) was designed based on rs10488141 and rs6994019 from PRACTICAL data (first SNP pair in Figure 4) with an OR of 1.9. For the null model, the disease prevalence of 0.2 was applied for all nine genotype combinations.

### 2.6 Performance Evaluation

Both power and type I error were evaluated in the 1000 simulation runs. Power is defined as the percentage of detecting a significant interaction under the true interaction model. Type I error is defined as percentage of detecting a significant interaction under the null model. The significant tests of the interaction for all four approaches (SIPI, MDR-LR, AA\_full and SNPassoc) were based on testing the coefficient of the interaction term. Statistical significance for SIPI and SNPassoc is defined as a  $p < 0.001$  ( $=0.05/45$ ) and  $p < 0.01$  ( $=0.05/5$ ). For the MDR-LR and AA\_Full approaches, the significance level is 0.05. In addition, we performed the pattern identification rate, which is defined as the percentage of identified correct interaction pattern among the significant simulation runs.

### 2.7 Prostate Cancer Study Application

SIPI was applied in evaluating SNP-SNP interactions in angiogenesis genes associated with prostate aggressiveness using Prostate Cancer Association Group to Investigate Cancer Associated Alterations in the Genome (PRACTICAL) consortium data. The study population includes 21,316 cases of European ancestry (3,812 aggressive and 17,504 non-aggressive) from the 32 study sites. We randomly selected half of the cases as the discovery set and the other half as the validation set in each study site. The sample sizes in the discovery and validation sets are 10,664 and 10,652, respectively. Individuals were excluded from the study based on strict quality control criteria including: overall call rate  $< 95\%$  and extremely high or low heterozygosity ( $p < 1.0 \times 10^{-5}$ ). Aggressive prostate cancer was defined as a Gleason score  $> 8$ , PSA  $> 100$ , disease stage of "distant" (stage IV) or death from PCa. Ethnic groups were defined based on a subset of 37,000 uncorrelated markers that passed quality control (including  $\sim 1,000$  selected as ancestry informative markers). Principal Component Analyses were carried out for the European subgroups. The details of this study population have been published previously (Eeles, et al., 2013).

We evaluated the 148 SNPs in the six angiogenesis genes (*EGFR*, *MMP16*, *ROBO1*, *CSF1*, *FBLN5*, *HSPG2*), which were reported in a genetic interaction network associated with prostate cancer aggressiveness (Lin, et al., 2013). These result in 10,878 SNP pairs. The pairwise interactions among these SNPs associated with prostate cancer aggressiveness (yes/no) were investigated using the SIPI approach in the discovery set first. For the top SNP pairs identified in the discovery set, both SIPI and AA\_Full were conducted in the validation set.

## 3 Results

### 3.1 Simulation

The power of the six simulated models for two SNPs with MAF of (0.5, 0.3), (0.5, 0.2) and (0.5, 0.05) are shown in Figures 1-2. As the sample size increased, power of all four approaches increased. In general, SIPI is more powerful and suffers less negative impact of SNPs with a low MAF

than the other three approaches (MDR-LR, AA\_Full and SNPassoc). In Models 1-4 for a SNP pair with a  $MAF \geq 0.2$  under a sample size of 1,000, SIPI has greater than 49% power while the other three approaches have low power ( $< 25\%$ ). Under a sample size of 1,000 with MAF of (0.5, 0.05), power decreases for all four approaches but SIPI still has the highest power. As the sample size increased to 5,000, SIPI has 100% power in most of the conditions for identifying an interaction with a SNP pair with MAF of (0.5, 0.3) and (0.5, 0.2). The order of power for detecting a SNP-SNP interaction is  $SIPI > MDR-LR > AA\_Full$  (similar with SNPassoc) with a big sample size of 5,000.

With a recessive interaction-only pattern (RR\_int\_rr) in Model 1 for a sample size of 1,000, SIPI has a power of 49-54%, but the other three approaches only have a power  $< 10\%$ . When the sample size increases to 5,000, the power of SIPI is 100% while the other three approaches' power remains low ( $< 35\%$ ). MDR-LR has higher power than AA\_Full and SNPassoc. This demonstrates that MDR-LR, AA\_Full, and SNPassoc have difficulty detecting the 'RR\_int\_rr' pattern.

For Model 2 with a dominant-dominant interaction-only pattern, SIPI has power 58-65%, but the other three approaches only have  $< 20\%$  power in a sample size of 1,000 and MAF of (0.5, 0.3) and (0.5, 0.2). As the sample size increases to 5,000, the power of all methods increase, and SIPI has the highest power compared with the other three approaches.

For Model 3 (RD\_int\_rr), SIPI has the highest power among all testing scenarios in Figure 1. For a sample size of 5,000, the highest power for MDR-LR is 83%, while SIPI's power is 100%. Similarly, the power of Model 4 (DD\_M1\_int\_o1), a dominant-dominant model with SNP1 main effect and an interaction, is 59-73% for SIPI and  $< 25\%$  for the other three approaches when the sample size is 1,000. Power increases to close to 100% for SIPI and 22-78% for others when the sample size becomes 5,000.

For Model 5 (AA\_Full), the AA\_Full method is the most powerful among all testing approaches in most of the conditions, except the condition of low MAF of (0.5, 0.05) in a sample size of 1,000. Under this special condition, SIPI has the highest power and about 70% of the SIPI significant runs selected AA\_M1\_int and AA\_M1\_int\_r pattern. For Model 6 generated according to the first SNP interaction pair in the prostate cancer application (see Figure 4), SIPI is still the most powerful approach in most of the conditions.

SIPI has a smaller negative impact for detecting an interaction of SNPs with a low MAF compared to other statistical approaches. As we expected, power of SIPI decreases when the SNPs' MAF decreases (Figures 1-2). However, this negative impact is much smaller in SIPI. For a SNP pair with (0.5, 0.3) and (0.5, 0.05) in Model 2 with a sample size of 5,000, the power of SIPI only decreases 13% (from 100% to 87%), but MDR-LR, AA\_Full and SNPassoc decreased 50%, 31%, and 39%, respectively.

As we expected, SIPI using the Bonferroni correction is the most conservative method among all testing approaches. As shown in Figure 3, SIPI has the smallest type I errors (0.01-0.02) compared to the other three methods. Some of SNPassoc's type I errors (0.021-0.057) are also less than 0.05. The type I errors for AA\_full and MDR-LR are close to 0.05. As shown in Tables S1-3, the power and type I error comparisons for additional MAF conditions show similar observations.

### 3.2 Pattern Detection Accuracy

The accuracy rate of pattern identification increases (Figures S4-S5) as the sample size increases. For Models 1, 2, 3 and 6 with 1,000 samples, 56-84% of the significant simulation runs identify the correct pattern. For the sample size of 5,000, all models have approximately 100% accuracy in identifying correct interaction patterns. For Models 4-5 and  $MAF = (0.3, 0.3)$  with a sample size of 1,000, the pattern identification rates are low (10% and 2%, respectively). However, these rate becomes 100% for a sample size of 5,000. Although pattern detection accuracy is low for the smaller sample, SIPI's power can still be high due to detection of other similar patterns. Using Model 4 with  $MAF = (0.3, 0.3)$  as an

example, only 10% of the significant runs detect the correct pattern (DD\_M1\_int\_o1) but other three similar patterns (39.9% DD\_int\_oo, 23% DR\_int\_rr, and 12.6% DR\_int\_or) are identified (Figure S5). Thus, its power of detecting any interaction can reach 61.2%.

From the simulation results, we observed an interesting scenario for common variants with a MAF close to 0.5. Under this condition, the minor allele determination is unstable. In the 1,000 simulation runs, the given allele had around a 50% chance to be classified into the major allele and a 50% chance of being classified into the minor allele. This unstable major/minor allele assignment affects SIPI's pattern labels, which are built upon the minor/major allele. As an example shown in Figure S3, a low risk subgroup of a (GG+ GG) combination of SNP1 and SNP2 are classified as the "DD\_int\_rr" pattern when SNP1 is with a major allele of 'G' and a minor allele of 'A' but is classified as "RR\_int\_or" (called a "sister pattern") when SNP1's major allele is 'A'. For an interaction with a SNP with a MAF close to 0.5, the pattern identification rate is the sum of the rates of the designed and sister patterns. We present the pattern identification rates for the significant simulation runs in Figures S4-S5. For Model 1 with a SNP pair with MAF=(0.5, 0.3), a total of 74% runs successfully identified the correct risk pattern (39% designed pattern and 5% sister pattern). A similar observations are presented for other models.

### 3.3 Example of Prostate Cancer Aggressiveness

For the proposed SIPI approach, we considered SNP pairs with a  $p < 1 \times 10^{-7}$  to be statistically significant after the Bonferroni correction for 489,510 tests (=10,878 pairs x 45 models per pair). Although the SNP-SNP interaction results do not appear to be significant after adjusting p values for multiple comparisons, some of them show promising consistent results in both datasets. In the discovery set, 25 SNP pairs had a  $p < 0.001$ . Among these top 25 pairs, four pairs have a  $p$ -value  $< 0.01$  in the validation set. Two pairs (rs10488141+ rs6994019 and rs2058502+ rs4947972) have the exact interaction pattern in both sets. The prevalence of prostate cancer aggressiveness by the nine genotype combinations are shown in Figure 4, and the prediction models are listed in Table 3. The prostate cancer patients with the TT + AC/AA genotype of the SNP pair of EGFR rs10488141 and MMP16 rs6994019 tend to suggest a higher risk of developing aggressive tumors (odds ratio (OR)=1.7,  $p=4.5 \times 10^{-6}$ ). Those with GG+ GG of two SNPs in EGFR (rs2058502 and rs4947972) are less likely to have aggressive prostate cancer tumors (OR=0.8,  $p=5.8 \times 10^{-6}$ ). Those with GG+ AG/AA of two SNPs in EGFR (rs723527 and rs845555) are likely to have aggressive prostate cancer tumors (OR=1.2,  $p=3.1 \times 10^{-4}$ ). The patients with AA/AG and CC in EGFR rs2075110 and CSF1 rs7538029 have a lower chance of developing an aggressive prostate cancer (OR=0.9,  $p=2.6 \times 10^{-5}$ ).

Table 3. Results of the PRACTICAL discovery and validation set for the top 25 SNP-SNP interaction pairs associated with prostate cancer aggressiveness with a  $p < 0.001$  in the discovery set

SNP1	SNP2	Pattern	Discovery <sup>1</sup>		Validation <sup>1</sup>		Pattern Similarity <sup>2</sup>
			$P_d$	$P_v$			
rs10228436	rs723527	DR_int_oo	$1.0 \times 10^{-4}$	DR_int_rr	0.020		
rs13222549	rs16880086	RR_int_oo	$2.0 \times 10^{-4}$	AA_int_ro	0.378		
rs2017000	rs6981717	DR_int_ro	$2.0 \times 10^{-4}$	AA_int_oo	0.043		
rs6956366	rs763317	RD_int_oo	$2.7 \times 10^{-4}$	DR_int_or	0.032		
<b>rs10488141</b>	<b>rs6994019</b>	<b>RD_int_oo</b>	<b><math>2.8 \times 10^{-4}</math></b>	<b>RD_int_oo</b>	<b>0.005</b>		same
rs723527	rs845552	RD_int_oo	$2.9 \times 10^{-4}$	RR_int_oo	0.056		
<b>rs2058502</b>	<b>rs4947972</b>	<b>DD_int_rr</b>	<b><math>8.9 \times 10^{-4}</math></b>	<b>RD_int_or</b>	<b>0.002</b>		Same (sister pattern)
rs6548616	rs7780270	DR_int_ro	$3.2 \times 10^{-4}$	RR_int_ro	0.181		
rs12666347	rs7781264	DR_int_ro	$3.6 \times 10^{-4}$	DD_int_ro	0.082		
rs2017000	rs723527	DR_int_oo	$3.7 \times 10^{-4}$	RR_int_rr	0.079		
<b>rs723527</b>	<b>rs845555</b>	<b>RD_int_oo</b>	<b><math>4.5 \times 10^{-4}</math></b>	<b>RR_int_rr</b>	<b>0.009</b>		similar
rs16880086	rs6954351	AA_int_ro	$4.6 \times 10^{-4}$	RR_int_oo	0.123		
rs10228436	rs7780270	DR_int_oo	$4.7 \times 10^{-4}$	DR_int_rr	0.070		
rs13222549	rs16880099	RD_int_oo	$4.9 \times 10^{-4}$	AA_int_oo	0.424		
rs10225877	rs16880086	AA_int_oo	$5.6 \times 10^{-4}$	RD_int_or	0.053		
rs1519938	rs9842630	DD_int_ro	$5.9 \times 10^{-4}$	DR_int_or	0.040		
rs13224708	rs17290392	DD_int_oo	$6.1 \times 10^{-4}$	DR_int_oo	0.943		
rs10488141	rs1879202	RR_int_oo	$6.4 \times 10^{-4}$	RD_int_oo	0.021		
rs10488141	rs2222294	RD_int_oo	$7.3 \times 10^{-4}$	DR_int_ro	0.063		
<b>rs2075110</b>	<b>rs7538029</b>	<b>RD_int_rr</b>	<b><math>7.7 \times 10^{-4}</math></b>	<b>DD_int_oo</b>	<b>0.007</b>		similar
rs13222549	rs17666091	RD_int_oo	$8.7 \times 10^{-4}$	DR_int_oo	0.021		
rs11986591	rs6954351	AA_int_ro	$9.1 \times 10^{-4}$	DR_int_oo	0.138		
rs11977660	rs9842630	DD_int_ro	$9.2 \times 10^{-4}$	RD_int_ro	0.044		
rs7780270	rs9832396	RD_int_or	$9.6 \times 10^{-4}$	RR_int_oo	0.191		
rs759169	rs9842630	AA_int_rr	$9.8 \times 10^{-4}$	AA_int_rr	0.150		

<sup>1</sup> $P_d$ : p-value in the discovery set,  $P_v$ : p-value in the validation set;  $P_d < 0.001$  and  $P_v < 0.01$  were in bold.

<sup>2</sup>Comparing patterns in the discovery and validation set for the SNP pairs with  $P_d < 0.001$  and  $P_v < 0.01$

Three of the four SNP interaction pairs remain promising (rs10488141+ rs6994019, rs2058502+ rs4947972, and rs2075110+ rs7538029) after including these four SNP pairs and the first five principal components of European ancestry in the model (Table 4). For evaluating whether the SNPs in the top pairs in the discovery are comparable in the validation set, the MAF of these SNPs are calculated. As shown in Table S4, the MAFs for these top SNPs are very similar in these two datasets. The individual effects of these SNPs in the combined dataset are also evaluated, and some SNPs did not have significant main effects. For example, the SNP pairs of rs10488141 and rs6994019 has an interaction with a p-value of  $4.5 \times 10^{-6}$  but without main effects (p-value=0.145 and 0.659, respectively). These show that some pure SNP-SNP interactions (without significant main effects) associated with prostate cancer aggressiveness. In summary, our results demonstrate SNP-SNP interactions in the two gene pairs (EGFR-MMP16 and EGFR-CSF1), and within EGFR. These findings support that EGFR may be the hub of this angiogenesis interaction network, which is consistent with the conclusion of the previous study (Lin, et al., 2013).

## 4 Discussion

SIPI is more powerful than the MDR-LR, AA\_Full and SNPAssoc approach, in general, even after applying stringent Bonferroni correction for multiple comparison justification. The primary strengths of SIPI are (1) taking various non-hierarchical models and inheritance modes into consideration and (2) using BIC to search for a best interaction pattern. In practice, it is challenge to detect a true interaction pattern for studies with a limited sample size. These features ensure that SIPI can search more similar interaction patterns close to the truth for overcoming the unstable nature of detecting SNP-SNP interaction patterns.

Our study demonstrated that SIPI is a more comprehensive and flexible tool for detecting two-way SNP-SNP interactions compared with the other three approaches. AA\_Full in PLINK (Purcell, et al., 2007), SNPAssoc (Gonzalez, et al., 2007), and MDR-LR (Edwards, et al., 2010) are all based on the full interaction model. The difference is how they deal with mode of inheritance. AA\_full treats SNPs as an additive mode. SNPAssoc considers five inheritance modes (additive, dominant, recessive, genotypic and over-dominant) but two SNPs in a pair need to have the same mode. MDR-LR classifies each SNP to a dominant or recessive mode based on risk profile of each individual SNP with the outcome. MDR-LR may not test the exact MDR-selected interaction pattern because the SNP classification is based on main effect. Thus, these three approaches can only detect a limited interaction patterns. For example, AA\_Full, SNPAssoc, and MDR-LR experienced difficulty in detecting the RR\_int\_rr pattern (Model 1, power<35%, Figure 1), but SIPI had 100% power for a large sample size of 5,000.

SIPI also provide advantages compared to other statistical approaches. BOOST (Wan, et al., 2010) uses the log-linear model to test interactions and treats SNPs as the genotypic mode. For BOOST, four degrees of freedom in a model are needed for each interaction pattern, so this worsens high-dimensional data issue. SNPmaxsel (Boulesteix, et al., 2007) evaluates 16 interaction patterns, which are parts of SIPI patterns. These 16 patterns are the interaction-only models for SNPs with a binary mode (dominant or recessive). HFCC (Gayan, et al., 2008) is used to assess 255 patterns, but some are rare or biologically meaningless patterns. Compared with these approaches, SIPI includes 45 biologically meaningful patterns, some of which have been reported previously (Lin, et al., 2013).

For external validation of SNP-SNP interactions, we suggest loosening the validation criteria for evaluating SNP-SNP interactions to allow for similar matches. The optimal goal of a genetic association study is to build prediction models for clinical usage. External validation using an independent dataset is a key in identifying true prediction factors. The majority of previous studies use AA\_Full in the two independent datasets or the exact interaction pattern identified in the discovery set to verify the same pattern in the validation set (Su, et al., 2013). However, this exact match is too stringent for identifying SNP-SNP interactions. Our simulation findings (Figures S4-S5) indicate the unstable nature of interaction patterns due to unsteady risk profiles of the nine genotype sub-

**SNP Interaction Pattern Identifier (SIPI)**

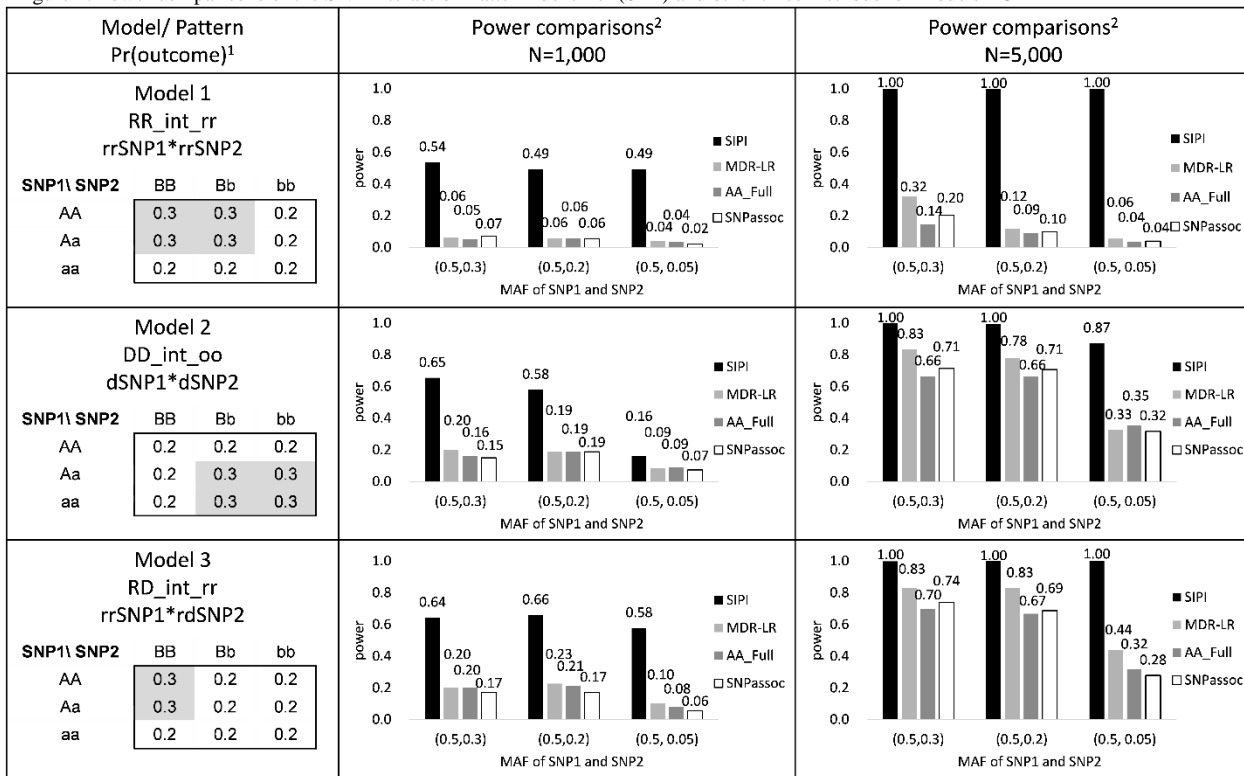
groups. Thus, it should be more effective to allow for similar matches instead of exact matches in SNP-SNP interaction validation, especially in the studies with a small sample size. SIPI provides useful features that work to overcome this unstable pattern nature. SIPI uses the BIC to select the best pattern of 45 patterns so that the true pattern or the most similar pattern can be detected. This provides flexibility in terms of result validation. For a SNP pair with MAF of (0.3, 0.3) in Model 4 with a sample size of 1,000, SIPI can still reach 61% power to detect an interaction of SNP1 and SNP2, even though only 10% of the significant results point to the correct pattern.

The outcome prevalence table stratified using three-by-three genotypes (called the “3x3 outcome table”, available in SIPI software) is a useful way to boost result interpretation for interaction patterns. Using the 3x3 outcome table for real prostate cancer data application, it is easy to observe that two of the top SNP pairs had similar interaction patterns in the discovery set and validation set (Figure 4). Combining the two testing sets with a larger sample size ensures that the interaction pattern is more reliable. In result validation, the sister pattern (one pattern with two different pattern labels) can be easily observed for an interaction with a SNP with a MAF close to 50%. In our prostate cancer application, three out of eight SNPs involved in the top SNP interactions have a

MAF>45%. In practice, the sister pattern issue can be identified by reviewing the 3x3 outcome table. Thus, we cannot purely rely on pattern label to decide whether the two patterns are exactly matched. Due to the sister pattern and similar matching issues, it is suggested that the 3x3 outcome table should be consulted to further review interaction patterns.

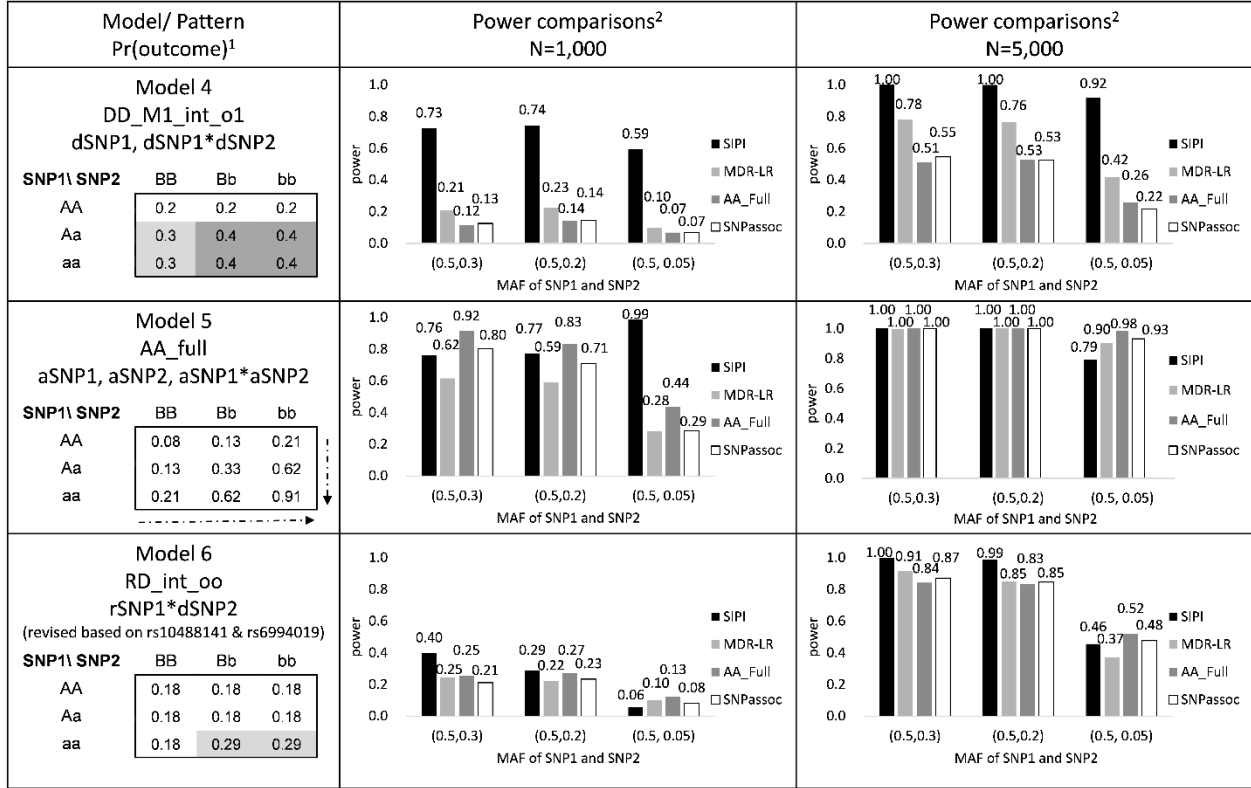
In summary, SIPI is a tool to search for 45 interaction patterns for pairwise SNP interactions. Although only binary outcome models were discussed in this study, it can be extended to various outcome data types, such as numeric and time-to-event data. The promising interaction pairs identified by SIPI can be included in a risk prediction model with other significant individual SNPs, other known clinical risk factors, and biomarkers in order to increase prediction accuracy.

Figure 1. Power comparisons of the SNP Interaction Pattern Identifier (SIPI) and other three methods for Models 1-3



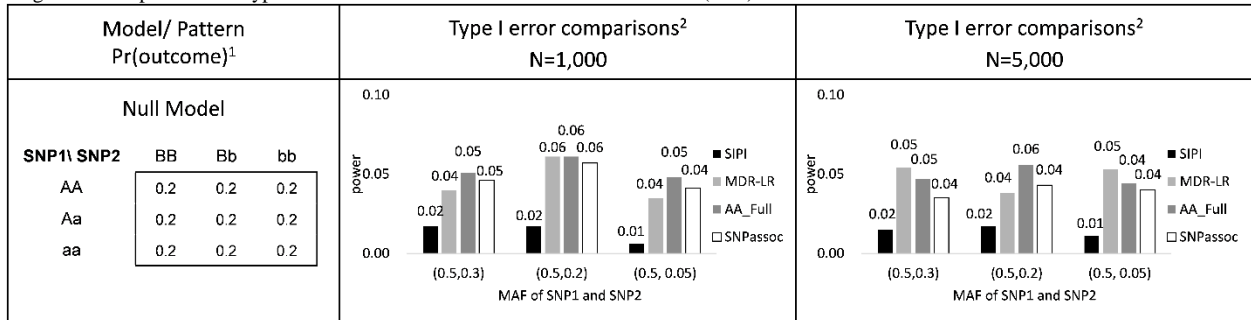
<sup>1</sup>Proportion of the outcome event in the genotype combination of the 3x3 table; a lowercase letter denotes the minor allele, and an uppercase letter denotes the major allele. <sup>2</sup> MDR-LR (Multifactor dimensionality reduction- Logistic Regression), AA\_Full (full interaction model and each SNP is treated as an additive mode), SNPAssoc R package

Figure 2. Power comparisons of the SNP Interaction Pattern Identifier (SIPI) and other three methods for Models 4-6



<sup>1</sup>Proportion of the outcome event in the genotype combination of the 3x3 table; a lowercase letter denotes the minor allele, and an uppercase letter denotes the major allele. <sup>2</sup>MDR-LR (Multifactor dimensionality reduction- Logistic Regression), AA\_Full (full interaction model and each SNP is treated as an additive mode), SNPAssoc R package

Figure 3. Comparisons of Type I errors of the SNP Interaction Pattern Identifier (SIPI) and other three methods



<sup>1</sup>Proportion of the outcome event in the genotype combination of the 3x3 table; a lowercase letter denotes the minor allele, and an uppercase letter denotes the major allele. <sup>2</sup>MDR-LR (Multifactor dimensionality reduction- Logistic Regression), AA\_Full (full interaction model and each SNP is treated as an additive mode), SNPAssoc R package



Figure 4. Proportions of prostate cancer aggressiveness by genotype combinations for the top four SNP-SNP interaction pairs associated with prostate cancer aggressiveness by datasets

Pair	Discovery: Pr(aggr) <sup>1</sup>			Validation: Pr(aggr) <sup>1</sup>			Combined: Pr(aggr) <sup>1</sup>					
	model	p-value		model	p-value		model	p-value				
1	EGFR: rs10488141	MMP16: rs6994019		EGFR: rs10488141	MMP16: rs6994019		EGFR: rs10488141	MMP16: rs6994019				
		CC	AC	AA		CC	AC	AA				
	AA	0.19	0.18	0.17	AA	0.17	0.17	0.20	AA	0.18	0.18	0.18
	AT	0.19	0.18	0.14	AT	0.18	0.16	0.16	AT	0.19	0.17	0.15
TT	0.14	0.28	0.34	TT	0.13	0.25	0.22	TT	0.14	0.27	0.29	
	RD_int_oo	2.8x10 <sup>-4</sup>		RD_int_oo	0.005		RD_int_oo	4.5x10 <sup>-6</sup>				
	AA_Full	0.047		AA_Full	0.842		AA_Full	0.192				
2	EGFR: rs2058502	EGFR: rs4947972		EGFR: rs2058502	EGFR: rs4947972		EGFR: rs2058502	EGFR: rs4947972				
		GG	CG	CC		GG	CG	CC				
	GG	0.16	0.20	0.23	GG	0.14	0.18	0.20	GG	0.15	0.19	0.22
	AG	0.19	0.20	0.21	AG	0.18	0.17	0.18	AG	0.18	0.18	0.19
AA	0.20	0.18	0.17	AA	0.18	0.15	0.21	AA	0.19	0.16	0.19	
	DD_int_rr	8.9x10 <sup>-4</sup>		RD_int_or <sup>2</sup>	0.002		DD_int_rr	5.8x10 <sup>-6</sup>				
	AA_full	2.3x10 <sup>-4</sup>		AA_full	0.023		AA_full	2.4x10 <sup>-5</sup>				
3	EGFR: rs723527	EGFR: rs845555		EGFR: rs723527	EGFR: rs845555		EGFR: rs723527	EGFR: rs845555				
		GG	AG	AA		GG	AG	AA				
	AA	0.19	0.18	0.18	AA	0.17	0.17	0.18	AA	0.18	0.17	0.18
	AG	0.17	0.19	0.18	AG	0.16	0.16	0.18	AG	0.17	0.18	0.18
GG	0.17	0.22	0.22	GG	0.19	0.18	0.20	GG	0.18	0.20	0.21	
	RD_int_oo	4.5x10 <sup>-4</sup>		RR_int_rr	0.009		RD_int_oo	3.1x10 <sup>-4</sup>				
	AA_full	0.034		AA_full	0.776		AA_full	0.186				
4	EGFR: rs2075110	CSF1: rs7538029		EGFR: rs2075110	CSF1: rs7538029		EGFR: rs2075110	CSF1: rs7538029				
		CC	CA	AA		CC	CA	AA				
	AA	0.17	0.20	0.18	AA	0.17	0.18	0.20	AA	0.17	0.19	0.19
	AG	0.17	0.20	0.22	AG	0.16	0.18	0.23	AG	0.17	0.19	0.23
GG	0.21	0.17	0.17	GG	0.17	0.19	0.17	GG	0.19	0.18	0.17	
	RD_int_rr	7.7x10 <sup>-4</sup>		DD_int_oo	0.007		RD_int_rr	2.6x10 <sup>-5</sup>				
	AA_full	0.023		AA_full	0.834		AA_full	0.138				

<sup>1</sup> Pr(aggr): Values in the 3x3 table are proportions of aggressive prostate cancer (=number of aggressive PCa patients/ number of PCa patients)

<sup>2</sup> RD\_int\_or in the validation set indicated the same pattern as DD\_int\_rr in the discovery set. The different pattern name is due to the revise issue of the major and minor allele in the validation set. rs2058502 (minor < major allele): (A<G) in discovery set, and (G<A) in validation set.

Table 4. SNP-SNP interaction models associated with prostate cancer aggressiveness

	Univariate model		Multivariable model <sup>2</sup>	
	Unadjusted	p-value	adjusted	p-value
	OR (95% CI) <sup>1</sup>		OR (95% CI) <sup>1</sup>	
rs10488141+ rs6994019, TT+ AC/AA vs. others	1.7 (1.4-2.1)	4.5x10 <sup>-6</sup>	1.8 (1.4-2.6)	6.3x10 <sup>-7</sup>
rs2058502+ rs4947972, GG+ GG vs. others	0.8 (0.7-0.9)	5.8x10 <sup>-6</sup>	0.8 (0.7-0.9)	5.2x10 <sup>-5</sup>
rs723527+ rs845555, GG+ AG/AA vs. others	1.2 (1.1-1.3)	3.1x10 <sup>-4</sup>	1.1 (1.0-1.3)	1.6x10 <sup>-2</sup>
rs2075110+ rs7538029, AA/AG+ CC vs. others	0.9 (0.8-0.9)	2.6x10 <sup>-5</sup>	0.9 (0.8-0.9)	6.9x10 <sup>-4</sup>

<sup>1</sup> Odds ratio (95% confidence interval)

<sup>2</sup> all four SNP pairs and the first five principal components of European ancestry were included in the multivariable model.

**Acknowledgements**

We thank our anonymous reviewers for their valuable comments, which have led to many improvements to this article. This study was supported by the National Cancer Institute (R01CA128813, PI: Park, JY and R21CA202417, PI: Lin, HY).

Conflict of Interest: none declared.

**References**

Al Olama, A.A., et al. A meta-analysis of 87,040 individuals identifies 23 new susceptibility loci for prostate cancer. *Nat Genet* 2014;46(10):1103-1109.  
 Amin Al Olama, A., et al. Risk Analysis of Prostate Cancer in PRACTICAL, a Multinational Consortium, Using 25 Known Prostate Cancer Susceptibility Loci.

*Cancer Epidemiol Biomarkers Prev* 2015;24(7):1121-1129.

Boulesteix, A.L., et al. Multiple testing for SNP-SNP interactions. *Stat. Appl. Genet. Mol. Biol.* 2007;6:Article37.

Cordell, H.J. Detecting gene-gene interactions that underlie human diseases. *Nat. Rev. Genet.* 2009;10(6):392-404.

Edwards, T.L., et al. A general framework for formal tests of interaction after exhaustive search methods with applications to MDR and MDR-PDT. *PLoS ONE* 2010;5(2):e9363.

Eeles, R.A., et al. Identification of 23 new prostate cancer susceptibility loci using the iCOGS custom

- genotyping array. *Nat Genet* 2013;45(4):385-391, 391e381-382.
- Gayan, J., et al. A method for detecting epistasis in genome-wide studies using case-control multi-locus association analysis. *BMC Genomics* 2008;9:360.
- Gonzalez, J.R., et al. SNPassoc: an R package to perform whole genome association studies. *Bioinformatics* 2007;23(5):644-645.
- Ioannidis, J.P., Castaldi, P. and Evangelou, E. A compendium of genome-wide associations for cancer: critical synopsis and reappraisal. *J. Natl. Cancer Inst.* 2010;102(12):846-858.
- Lin, H.Y., et al. SNP-SNP Interaction Network in Angiogenesis Genes Associated with Prostate Cancer Aggressiveness. *PLoS ONE* 2013;8(4):e59688.
- Lin, H.Y., et al. TRM: a powerful two-stage machine learning approach for identifying SNP-SNP interactions. *Annals of human genetics* 2012;76(1):53-62.
- Lin, H.Y., et al.** Comparison of multivariate adaptive regression splines and logistic regression in detecting SNP-SNP interactions and their application in prostate cancer. *Journal of Human Genetics* 2008;53(9):802-811.
- Milne, R.L., et al. The importance of replication in gene-gene interaction studies: multifactor dimensionality reduction applied to a two-stage breast cancer case-control study. *Carcinogenesis* 2008;29(6):1215-1218.
- Moore, J.H. The ubiquitous nature of epistasis in determining susceptibility to common human diseases. *Hum. Hered.* 2003;56(1-3):73-82.
- Moore, J.H. and Williams, S.M. New strategies for identifying gene-gene interactions in hypertension. *Annals of medicine* 2002;34(2):88-95.
- Onay, V.U., et al. SNP-SNP interactions in breast cancer susceptibility. *BMC Cancer* 2006;6:114.
- Piegorsch, W.W., Weinberg, C.R. and Taylor, J.A. Non-hierarchical logistic models and case-only designs for assessing susceptibility in population-based case-control studies. *Stat. Med.* 1994;13(2):153-162.
- Purcell, S., et al. PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* 2007;81(3):559-575.
- Ritchie, M.D., Hahn, L.W. and Moore, J.H. Power of multifactor dimensionality reduction for detecting gene-gene interactions in the presence of genotyping error, missing data, phenocopy, and genetic heterogeneity. *Genet. Epidemiol.* 2003;24(2):150-157.
- Ritchie, M.D., et al. Multifactor-dimensionality reduction reveals high-order interactions among estrogen-metabolism genes in sporadic breast cancer. *Am. J. Hum. Genet.* 2001;69(1):138-147.
- Schwarz, G. Estimating the Dimension of a Model. *Annals of Statistics* 1978;6(2):461-464.
- Su, W.H., et al. How genome-wide SNP-SNP interactions relate to nasopharyngeal carcinoma susceptibility. *PLoS ONE* 2013;8(12):e83034.
- Van den Broeck, T., et al. The role of single nucleotide polymorphisms in predicting prostate cancer risk and therapeutic decision making. *Biomed Res Int* 2014;2014:627510.
- Wan, X., et al. BOOST: A fast approach to detecting gene-gene interactions in genome-wide case-control studies. *Am. J. Hum. Genet.* 2010;87(3):325-340.
- Yang, Y. Can the strengths of AIC and BIC be shared? A conflict between model identification and regression estimation. *Biometrika* 2005;92(4):937-950.