

# UCLA

## UCLA Previously Published Works

### Title

Current Challenges and New Opportunities for Gene-Environment Interaction Studies of Complex Diseases.

### Permalink

<https://escholarship.org/uc/item/9731f4gk>

### Journal

American journal of epidemiology, 186(7)

### ISSN

0002-9262

### Authors

McAllister, Kimberly  
Mechanic, Leah E  
Amos, Christopher  
[et al.](#)

### Publication Date

2017-10-01

### DOI

10.1093/aje/kwx227

Peer reviewed

# Current Challenges and New Opportunities for Gene-Environment Interaction Studies of Complex Diseases

Kimberly McAllister, Leah E. Mechanic, Christopher Amos, Hugues Aschard, Ian A. Blair, Nilanjan Chatterjee, David Conti, W. James Gauderman, Li Hsu, Carolyn M. Hutter, Marta M. Jankowska, Jacqueline Kerr, Peter Kraft, Stephen B. Montgomery, Bhramar Mukherjee, George J. Papanicolaou, Chirag J. Patel, Marylyn D. Ritchie, Beate R. Ritz, Duncan C. Thomas, Peng Wei, John S. Witte on behalf of GxE meeting participants

Running Title: Gene-Environment Interaction Studies of Complex Diseases

Corresponding Author:

Leah E. Mechanic, Ph.D., M.P.H.

Genomic Epidemiology Branch

Epidemiology and Genomics Research Program

Division of Cancer Control and Population Sciences

National Cancer Institute

9609 Medical Center Drive, Rm. 4E104, MSC 9763

Bethesda, MD 20892

(For express delivery, use Rockville, MD 20850)

Phone: 240-276-6847

Email: [mechanil@mail.nih.gov](mailto:mechanil@mail.nih.gov)

## **Abbreviations:**

CHARGE, Cohorts for Heart and Aging Research in Genomic Epidemiology; ENCODE, Encyclopedia of DNA Elements; GxE, gene-environment; GWAS, genome-wide association study; GTEx, Genotype-Tissue Expression; TCGA, the Cancer Genome Atlas

## **Abstract**

Recently, many new approaches, study designs, and statistical and analytical methods have emerged for studying gene-environment interactions (GxE) in large-scale human population studies. There are currently opportunities in this field, particularly with respect to the incorporation of -omics and next-generation sequencing data and continual improvement in measures of environmental exposures implicated in complex disease outcomes. A workshop held on October 17-18, 2014 by the National Institute of Environmental Health Sciences and the National Cancer Institute in conjunction with the annual American Society of Human Genetics meeting explored new approaches and tools developed in recent years for GxE interaction discovery. This paper will highlight current and critical issues and themes in GxE research discussed that need additional consideration including the topics of improved data analytical methods, environmental exposure assessment, and incorporation of functional data and annotations.

Keywords: gene-environment , genome-wide association study , environmental exposure

## **Introduction**

Genetic and environmental factors are thought to contribute to the etiology of most complex diseases. Through genome-wide association studies (GWAS), thousands of common loci associated with complex diseases have been identified (1-3). Researchers have been motivated to discover and describe how the interplay of these factors influence disease risk and outcomes. Several reasons for studying gene-environment (GxE) interaction include: providing insights into the biology of disease (e.g. developing new models for disease etiology based on observed GxE findings); building better prognostic models (e.g. using genotype to inform treatment and prognosis); identifying possible high-penetrance subgroups (e.g. increased genotype-specific risk in pre-menopausal women); or conversely, identifying genetic subgroups with higher exposure-specific disease risk for prevention efforts (e.g. increase environmental-specific risk for individuals with a particular genotype) (4-7). Furthermore, in the search for novel genes via GWAS, the modifying effects of environmental risk factors are not often taken into account; therefore, leveraging GxE may result in discovery of additional disease susceptibility loci (5, 8, 9). Despite interest in GxE, there are few agreed upon successes where the effect of exposure differs across genotypes (and vice versa). Numerous reasons have been suggested to contribute to the small number of successes including: the inherent low power of tests for GxE, complexity of measurement of environmental exposures and difficulty of incorporating temporality of environmental exposures, measurement error, limited range of genetic and/or environmental variation, scale dependence in the definition of statistical interaction, and lack of data on the biological consequences of most genetic variants (10-13).

The past few years have seen an emergence of new approaches, study designs, and statistical and analytical methods for exploring gene-environment interactions (GxE) in large-

scale human population studies. Further, new opportunities in this field, with respect to the incorporation of -omics and next-generation sequencing data and improvements in measures of environmental exposures implicated in complex disease outcomes, continue to be developed. Therefore, on October 17-18, 2014, National Institute of Environmental Health Sciences and National Cancer Institute held a workshop at the 64<sup>th</sup> Annual Meeting of the American Society of Human Genetics to explore these new approaches and tools for GxE interaction discovery. Based on the discussions, we prepared four papers that provide an update on: 1) the state of the science in analytical methods (14); 2) opportunities for incorporation of biological knowledge into GxE analyses (15); 3) advances in environmental exposure assessment in human populations (Chirag J. Patel, Department of Biomedical Informatics, Harvard Medical School, unpublished manuscript); and 4) lessons learned from GxE successes (Beate R. Ritz, Department of Epidemiology, Fielding School of Public Health, University of California Los Angeles, unpublished manuscript). In addition, this current paper develops some overarching themes and sets the stage for this series. As environmental factors may be modifiable, defining subpopulations of individuals most susceptible to environmental factors through GxE analysis may provide targets to improve public health. This idea is consistent with the goal for President Obama's recently launched Precision Medicine initiatives at the National Institutes of Health (NIH) --to better understand how individual variability contributes to differences in response to treatment or prevention (16, 17).

### **Analytical Methods**

Studies of GxE interaction require much larger sample sizes than studies targeting either genetic or environmental main effects alone (18). Further, when performing GxE on a genome-

wide scale, sometimes referred to as genome-wide interaction studies, sample size requirements are substantially further inflated to account for the multiple comparisons performed (5, 19).

Therefore, a goal of GxE methods development has been to improve power to detect associations. As detailed in the accompanying manuscript (14), many different methods have been explored in the context of a case-control studies as alternatives to traditional GxE tests, including case-only (20), empirical Bayes (21), Bayes Model Averaging (22), joint tests (9, 23, 24), case parent approaches (25-27), and 2-step approaches (19, 23, 28-33). Other approaches include set-based methods, which combine multiple variants or GxEs and which may be particularly appropriate for studies of rare variants (34-38). In addition, several methods have been developed to analyze GxE for quantitative outcomes (39-46).

The large number of available methods, as well as novel software tools to support the application of these methods (29, 47, 48), create opportunities to better study GxE interactions in genome-wide settings. Researchers may therefore wonder which method to use for their studies. Several previous simulation studies suggest that none of these GxE methods is universally the most powerful approach (29, 30, 49-52). Therefore, decisions about the most appropriate approach depend on several considerations including the hypotheses to be tested, likely genetic architecture, study design attributes, and characteristics of the population being studied. Investigators should be cautious about applying multiple methods to their data without an *a priori* basis for choosing among the results, as simply picking those with the most "significant" findings to report would clearly be a biased strategy that could contribute to spurious associations and to what has been referred to as a "vibration of effects" (53, 54). Some of the new methods, however, provide flexible frameworks for combining multiple tests with an appropriate permutation procedure to evaluate the significance of the overall results (29, 30).

The collection of methods allows investigators to address specific scientific questions and offers new opportunities for studies of GxE in large populations.

### **Functional Validation and Discovery**

Despite the recent success of GWAS at identifying risk loci, by design variants identified are not usually the causal variants, defined as the functional genetic variant that influences risk of disease and explains the association. Currently, the underlying biological mechanism contributing to disease risk is only known for a small proportion of these loci. Therefore, more research to functionally characterize risk loci is now being performed, providing opportunities by which GxE analyses may shed new insights into disease development (55). An understanding of the biological consequences of particular genetic differences could lead to specific mechanistic hypotheses, identifying relevant exposures to test and specifying relevant statistical models. As described in the accompanying manuscript (15), these approaches include utilizing functional annotations for discovery and validation, studying molecular phenotypes (e.g. epigenetics or gene expression) to improve GxE discovery, and leveraging *in vitro* and *in vivo* models for these studies.

Several large public databases [such as Encyclopedia of DNA Elements (ENCODE), Epigenomics Roadmap, Genotype-Tissue Expression (GTEx), and the Cancer Genome Atlas (TCGA)] have facilitated the functional annotation and interpretation of many genomic regions, which can be used to prioritize candidate GxE markers (30). Many disease-associated GWAS SNPs appear to be located in non-coding or regulatory regions which are often affected by environmental exposures (56-58). The ENCODE and Roadmap Epigenomics programs have helped to define many of the regulatory regions, and new tools developed by these programs and

others now allow functional annotation information, such as the genomic location of histone modification states, methylation patterns, transcription factor binding sites and DNase hypersensitivity sites or other higher order chromosomal structural information, to be overlaid with GWAS results and could be integrated into GxE analyses (59-63). Projects like GTEx have greatly increased the compendium of putative biological functions of genetic variants. However, neither GTEx nor large-scale epigenomics projects provide information on effects of genetic and genomic functions across a range of environmental conditions. To explore genetic effects in response to environment, *in vitro* studies have now perturbed cells and recorded responses to various drugs, infections, and other exposures. Through use of intermediate molecular phenotypes such as gene expression, these efforts have demonstrated success in illustrating how an exposure may impact gene function, suggesting potential candidate genes or variants for GxE studies (64-68).

In addition to data resources, the use of population-based mouse resources (such as the Collaborative Cross, Diversity Outbred, and Hybrid Mouse Diversity Panel) and other appropriate mouse models have also been leveraged to assist in the discovery or replication of GxE interactions. These population-based variant enriched mouse resources have been designed to mimic the genetic diversity of human populations and can be used to replicate or inform GxE hypotheses by utilizing carefully controlled exposures in the mouse studies. Several recent examples have exemplified the power of these resources to map genetic variants related to susceptibility to environmental exposures (69, 70). Although both *in vitro* and model systems have led to potential mechanistic insights, linking of these to human populations remains challenging.

There are many approaches for incorporating biological knowledge to improve analytical methods (71) for GxE interaction in both the discovery and the validation phase. Incorporating functional annotation data and *a priori* biological information (such as metabolomics or gene expression data collected on individuals or knowledge on biological pathways) to inform data analytical GxE methods have aided in the discovery of new GxE findings in recent years (72). For example, Bayesian Variable Selection (73, 74), the Algorithm for Learning Pathways (75) and PEAK (72) are all methods that incorporate external biological information and properties of the dataset itself to increase power over agnostic approaches to detect interactions. Another approach is to use 2-stage modeling where functional annotations are used to prioritize variants (76, 77) for GxE studies. As one example, Biofilter was designed to build biologically plausible models of gene-gene and GxE interactions to test for associations based on biological features using biological knowledge from the public domain (76, 78). These types of filtering approaches are also being explored to prioritize environmental exposures by using databases such as the Comparative Toxicogenomics Database, which links exposures to genes (79). However, challenges still exist in linking environmental exposures into currently available ontological knowledge resources, though some investigators are beginning to navigate these challenges (80). Furthermore, all these databases and functional annotations depend on the quality and extent of existing biological knowledge (71).

### **Environmental Exposures**

The complex realities of environmental exposures have long made measurement of exposures substantially more complicated than inherited genetic measurements (e.g. genotypes) and single nucleotide variants in particular; the technologies and approaches to incorporate

exposures into human population studies have therefore lagged behind genomics capabilities (11, 81). Assessing exposure impact must take into context not just the variety of exposures themselves (physical, often complex chemical mixtures, biological, and psychosocial) but also the source and place of exposure, the timing during a person's life trajectory, the route of contact (skin, lung, diet), metabolism/excretion, and distribution in target tissues. All of these factors may impact the ultimate disease risk associated with environmental exposures. In addition, in the classic environmental exposure paradigm, studies may focus on measurements to capture internal versus external exposure, early markers of disease, or an ultimate biological response, which further adds to the complexity of exploring the impact of environmental exposures.

In recent years, however, exciting new opportunities have become available for environmental exposure assessment. The potential importance of examining the totality of internal and external exposures, referred to as the 'exposome', has been recognized (81, 82). Several recent commentaries described considerations for measurements of the exposome (83-86). Innovative technologies including activity monitors, improved sensors, global positioning systems, and Geographic Information Systems, which enable new and more detailed exposure measurements. Although issues of the timing of exposure measures persist and should be considered. Moreover, development of biological response markers for assessment of exposure, such as changes in gene expression, transcriptomic signatures, and DNA methylation profiles, has been useful for GxE discovery (87-90). Another opportunity is the exploration of environmental exposures in a more agnostic discovery-based fashion, similar to GWAS. These studies, termed environment wide association studies (EWAS) led to new discoveries of environmental factors associated with disease (91-94).

Key challenges and considerations remain associated with assessing environmental exposures in GxE studies, as detailed in the accompanying manuscript (Chirag J. Patel, Department of Biomedical Informatics, Harvard Medical School, unpublished manuscript) including how to: select most appropriate study designs, incorporate high throughput -omic measures (e.g. metagenome, metabolome) and sensor technologies into human population-based studies, assess long term exposure, integrate a variety of divergent external exposure and internal response data, and further advance statistical approaches to handle the dynamic nature of exposure data. We are now at the early stages of exploring what novel exposure assessment technologies can be appropriately applied to larger population studies most effectively. To this end, some two-stage study designs have been investigated (24, 95-98). Given the extreme cost of incorporating some sophisticated environmental measures in a large scale human population study, the question of what can be accomplished with dense (i.e. repeated measures of a marker or measurement of multiple analytes using an -omic platform) environmental measures on a subsample and extrapolating to a larger sample size (and whether simulations can demonstrate that this approach increases power to detect GxE) is currently being explored (49, 95, 99).

Several analytical methods have been developed for the unique considerations of exposure assessment. New statistical methods can adjust for exposure misclassification (which has been shown to lead to inflated type I errors and substantially reduced power) much better; these approaches should allow for obtaining greater power with smaller sample sizes. In addition, novel statistical methods have been developed to detect gene by longitudinal exposure interactions by taking into account long term time-varying exposures (100). Importantly, as researchers begin to combine exposure data to obtain larger sample sizes required for GxE research across studies, they have to address that exposures may have been measured using

different approaches or have very different distributions in and between populations such that exposure misclassification could produce spurious associations (14). There is also the challenge of exposure-related population stratification for studies relating to GxE interactions (101). Meanwhile, multiple measures can sometimes increase power for detecting associations. For example, in a recent study, continuous monitoring was shown to reduce the sample size required in a clinical trial context (102).

### **GxE Examples from Human Population Studies**

By examining GxE successes, it may be possible to improve the design of GxE studies for the future. Examples of GxE successes span from Mendelian-like traits (e.g. phenylketonuria) to complex diseases (*NAT2* variants, smoking and bladder cancer) and response to therapies (*HLA-B\*1502* variant and carbamazepine induced Stevens Johnsons Syndrome) (Beate R. Ritz, Department of Epidemiology, Fielding School of Public Health, University of California Los Angeles, unpublished manuscript). In addition, several recent studies examined the use of polygenic risk scores, generated from common genetic variation, to assess the impact of environmental factors on individuals with low compared with higher genetic risk (103-106). In the accompanying manuscript highlighting some of the most successful GxE interactions identified to date (Beate R. Ritz, Department of Epidemiology, Fielding School of Public Health, University of California Los Angeles, unpublished manuscript), several common themes have emerged including: the strength of focusing on metabolic pathways for a specific exposure, the utility of studying unique, highly or diversely exposed populations, the necessity of using high-quality exposure assessment methods, the need for large sample sizes, and the utility of model

systems to demonstrate genetic function when replication is challenging in population-based studies. These suggest important avenues for undertaking successful future research in GxE.

### **Themes and Future Directions**

Inclusion of diverse populations may facilitate GxE research by improving power for discovery of casual genetic variants and environmental factors associated with disease. Trans-ethnic differences in the distribution of linkage disequilibrium can be leveraged to improve fine mapping to identify potential causal alleles (107-110). Combining admixture mapping with conventional GWAS may also facilitate discovery of novel loci (111). Using this later approach, novel loci were identified associated with total IgE levels (112) and asthma (113). Lastly, using geographically diverse populations might expand the distribution of the environmental exposure and thus increase power to detect interactions (13). Performing genetic studies on populations of diverse ancestry may improve our understanding of disease mechanisms and such studies are required to ensure all populations benefit equally from this research (114).

Replication is an essential component to genetic association studies, and the requirement for independent replication contributed to the success of GWAS (115, 116). However, replication and meta-analysis becomes challenging as GxE studies become sophisticated in analytical methods, exposure assessment, and incorporation of functional information.

Differences in the underlying distribution of environmental exposures, genetic linkage disequilibrium (LD) structure, and genetic modifiers can reduce the power to detect the same level of interaction in independent studies. Moreover, an appropriate human replication study may not (yet) exist: in studies of a rare disease, genetic variant, or environmental exposure;

where exposures are unique to particular populations; or where the initial finding was obtained within a large consortium comprising all known studies of a specific outcome (12). As illustrated in the manuscripts describing GxE successes and incorporation of biological knowledge, in some situations functional studies could serve to provide support for initial GxE observations in absence of a suitable replication population. Moreover, as the field considers gene and pathway based approaches to study GxE, replication may become further complicated as different combinations of genes in different datasets may be observed in the interaction. Some have argued that replication requirements might be met if the underlying biological pathway is the same even if replication was not observed with the individual SNP or gene (15). More consideration of standards for replication, definitions of replication and alternative approaches for replication and verification of GxE results is needed.

Many exciting opportunities exist for studies of GxE. There is the emerging recognition that developmental exposures may lead to disease throughout life and efforts have focused on beginning to address how much of environmental exposure risk for many disease outcomes may be attributable to *in utero* exposures or other particularly vulnerable windows of susceptibility (childhood, adolescence, etc.). Successful integration of large volumes of diverse data types (including Geographic Information Systems, sensor, metabolomics and other omics data) will create the opportunity for generating unique insights. Epigenetics tools open up new opportunities to directly link environmental exposure to the genome and generate new exposure biomarkers (e.g. methylation of cancer specific genes associated with dietary folate and alcohol in colorectal cancer (117) or smoking exposure in lung cancer (118)) [for review (119, 120)]. Moreover, epigenomics, as well as other -omic technologies, may elucidate mechanisms by which exposures contribute to disease. The role of the microbiome as a key environmental risk

factor for many complex disease phenotypes is starting to be appreciated and extensively studied. In addition, molecular phenotype data creates opportunities to examine disease subtypes or more precisely classify disease. This may eventually reduce heterogeneity in studies and improve power to study GxE associations, assuming molecular characterizations are performed with the correct cell type, tissue, or appropriate surrogate tissue for the hypothesis being tested.

Additional areas of research may allow further advances in GxE discovery and replication. The field needs to determine how to best leverage experimental studies in animals or human cell lines to aid in discovering and functionally validating GxE interactions. Moreover, it is unclear how to best leverage existing family and twin based studies for examining GxE. In incorporating functional information into GxE studies, questions remain about the appropriate balance between using prior or external information versus the characteristics of the dataset being studied in building analytical models and appropriate methods for linking environmental exposures information into available biological knowledge databases that are usually focused on genes and pathways. In addition, since many GxE findings to date have modest effect sizes or have not been extensively replicated (11, 121), exploring the general question of when the effort of attempting to identify these complex types of interactions is worth it. Though even with modest effect sizes, if a GxE finding is sufficiently replicated in human populations and supported by other experimental data, this information could provide insights into possible disease mechanisms. Finally, given the reduced power to detect GxE combinations with present methods, approaches that examine higher order interactions should be taken on cautiously.

Despite many recent advances in analytical methods for GxE discovery and some validation in recent years, additional statistical methods are needed for studies of copy number and rare genetic variation, survival traits, analysis of trios, and meta-analysis and pooling in large

consortia. In addition, many of the assumptions about expected GxE findings are based on results from genetic simulation studies, but these expectations have not always directly correlated to GxE observations in real population studies. Therefore, the question remains whether simulation studies have been designed with realistic assumptions about the underlying genetic architecture of the traits and whether better simulation approaches are needed (122).

Extended collaboration and data sharing will also advance GxE research. Large epidemiological consortia, such as the Cohorts for Heart and Aging Research in Genomic Epidemiology (CHARGE) consortium, that have longitudinal measures of environmental exposures have been heavily leveraged in recent years as a way to examine repeated environmental exposures over time and attempt to incorporate cumulative and time-varying exposures into assessments of complex disease risk (123). There is also a need for further collaboration to allow validation of biomarkers in larger cohorts. Meta-analysis and pooling methodology and efforts will likely need to be advanced to have the power to detect GxE in rarer diseases. Standards are needed to describe the adequate criteria for identifying, reproducing, and reporting a GxE finding; a place to publish negative findings would allow researchers to avoid repeating failed experiments (11). There is also a need for greater integration and education with other fields to better design studies of GxE. Specifically, toxicology expertise will be needed to allow validation in experimental models of GxE discoveries. Lastly, compared to genomic data sharing, the sharing of environmental and epidemiological data has lagged behind. Some have suggested that an environmental data sharing policy mirroring the National Institutes of Health genomic data sharing policy could advance data sharing in the environmental health science fields. However, there are unique sensitivities and ethical issues related to the sharing of environmental data that must be considered., including participant confidentiality and privacy

issues (i.e. environmental exposure data with global positioning systems information can allow specific identification of the sources of exposure) and legal/ regulatory matters (i.e. regulatory reporting, remediation, and reform).

Researchers are exploring how to apply GxE findings to risk prediction studies as a possibility for targeted screening or intervention. Questions remain about the optimal approaches for risk prediction models, including how to integrate biomarkers and external exposures and how best to model the joint effects of genetic markers, biomarkers, and lifestyle and environmental exposures (124). Although most statistical methods for detecting GxE focus on identifying departures from a multiplicative relative risk model, the absence of multiplicative interactions will typically imply the presence of additive interaction (i.e., when there are marginal genetic and environmental effects). Additive interactions may have public health implications, as they suggest the difference in absolute risks between exposed and unexposed groups differs across genetically defined subgroups (103-106, 124). If an exposure causes disease, then an intervention to remove the exposure will prevent more cases in a genetically sensitive population than in an equivalently sized genetically insensitive population. Important challenges that remain include: determining whether the exposure in fact causes disease, developing effective interventions to change exposures, and evaluating whether targeted or population-level interventions optimize the risk: benefit trade-off. As with main effects, where it is well understood that observational findings of associations across individuals do not necessarily imply that an intervention to change exposure will change any individual's outcomes, so an additive interaction does not necessarily imply that a genetically-targeted intervention would be a more effective prevention strategy. Modern methods of causal inference (125, 126) may be useful for estimating the causal difference in disease rates between genetically-targeted

and population-wide exposure interventions. Finally, the lessons and approaches for research into how the combination of genes and environment contribute to disease relates broadly to the studies of precision medicine and precision prevention. These types of studies may lead to insights for targeting prevention, intervention, or treatment in the future.

### **Acknowledgements**

#### **Author Affiliations:**

Genes, Environment, and Health Branch, National Institute of Environmental Health Sciences (NIEHS), National Institutes of Health (NIH), Research Triangle Park, North Carolina (Kimberly McAllister); Epidemiology and Genomics Research Program, Division of Cancer Control and Population Sciences, National Cancer Institute (NCI), NIH, Bethesda, Maryland (Leah E. Mechanic); Department of Biomedical Data Science, Dartmouth College, Lebanon, New Hampshire (Christopher Amos); Department of Epidemiology, Harvard T.H. Chan School of Public Health, Boston, Massachusetts; Centre de Bioinformatique, Biostatistique et Biologie Intégrative (C3BI), Institut Pasteur, Paris, France (Hugues Aschard); Center of Excellence in Environmental Toxicology and Penn SRP Center, Perelman School of Medicine, University of Pennsylvania Philadelphia, Pennsylvania; Department of Systems Pharmacology and Translational Therapeutics, Perelman School of Medicine, University of Pennsylvania Philadelphia, Pennsylvania (Ian A. Blair); Department of Biostatistics, Bloomberg School of Public Health, Department of Oncology, School of Medicine, Johns Hopkins University, Baltimore, Maryland (Nilanjan Chatterjee); Department of Preventive Medicine, University of Southern California, Los Angeles, California (David Conti); Department of Preventive Medicine, University of Southern California, Los Angeles, California (W. James Gauderman); Biostatistics and Biomathematics Program, Division of Public Health Sciences, Fred Hutchinson Cancer

Research Center, Seattle, Washington (Li Hsu); Division of Genome Sciences, National Human Genome Research Institute, NIH, Bethesda, Maryland (Carolyn M. Hutter); California Institute for Telecommunications and Information Technology, Qualcomm Institute, University of California San Diego, La Jolla California (Marta M. Jankowska); Department of Family Medicine and Public Health, University of California San Diego, La Jolla, California (Jacqueline Kerr); Department of Epidemiology, Harvard T.H. School of Public Health, Boston, Massachusetts (Peter Kraft); Departments of Genetics and Pathology, Stanford University School of Medicine, Stanford, California (Stephen B. Montgomery); Department of Biostatistics, University of Michigan School of Public Health, Ann Arbor, Michigan (Bhramar Mukherjee); Division of Cardiovascular Sciences, Prevention and Population Sciences Program, National Heart, Lung, and Blood Institute, National Institutes of Health, Bethesda, Maryland (George J. Papanicolaou); Department of Biomedical Informatics, Harvard Medical School, Boston, Massachusetts (Chirag J. Patel); Department of Biochemistry and Molecular Biology, Center for Systems Genomics, The Pennsylvania State University, University Park, Pennsylvania; Biomedical and Translational Informatics, Geisinger Health System, Danville, Pennsylvania (Marylyn D. Ritchie); Department of Epidemiology, Fielding School of Public Health, University of California Los Angeles, Los Angeles, California (Beate R. Ritz); Department of Preventive Medicine, University of Southern California, Los Angeles, California (Duncan C. Thomas); Department of Biostatistics, The University of Texas MD Anderson Cancer Center, Houston, Texas (Peng Wei); Department of Epidemiology and Biostatistics, University of California, San Francisco, California (John S. Witte)

We acknowledge NIEHS for providing funds to support this meeting. Research reported in this publication was supported by the National Cancer Institute, National Heart Lung and Blood Institute, National Human Genome Research Institute, and National Institute of Environmental Health Sciences of the National Institutes of Health, and the National Science Foundation under award numbers R21HG007687 to H.A.; R01CA140561, R01CA201407, and P01CA196569 to D.C.; R01CA189532, R01CA195789, and P01CA53996 to L.H.; R21CA169535 and R01CA179977 to J.K.; R21ES020811 and NSF DMS 1406712 to B.M.; R00ES023504 and R21ES025052 to C.J.P.; R01CA169122, R01HL116720 and R21HL126032 to P.W.; and R01CA201358 to J.S.W. S.B.M. is supported by the National Institutes of Health through R01HG008150, R01MH101814, U01HG007436, and U01HG00908001. This work is funded, in part, under a grant with the Pennsylvania Department of Health (SAP 4100070267) to M.D.R. The authors thank the participants in the workshop “Current Challenges and New Opportunities for Gene-Environment Interaction Studies of Complex Diseases”. The Pennsylvania Department of Health specifically disclaims responsibility for any analyses, interpretations or conclusions. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health.

The authors have no conflicts of interest to report.

## References

1. Hindorff LA, Gillanders EM, Manolio TA. Genetic architecture of cancer and other complex diseases: lessons learned and future directions. *Carcinogenesis* 2011;32(7):945-954.
2. Hindorff LA, Sethupathy P, Junkins HA, et al. Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc Natl Acad Sci U S A* 2009;106(23):9362-9367.
3. Stadler ZK, Thom P, Robson ME, et al. Genome-wide association studies of cancer. *J Clin Oncol* 2010;28(27):4255-4267.
4. Hunter DJ. Gene-environment interactions in human diseases. *Nature reviews Genetics* 2005;6(4):287-298.
5. Thomas D. Gene-environment-wide association studies: emerging approaches. *Nature reviews Genetics* 2010;11(4):259-272.
6. Thompson WD. Effect modification and the limits of biological inference from epidemiologic data. *J Clin Epidemiol* 1991;44(3):221-232.
7. Le Marchand L, Wilkens LR. Design considerations for genomic association studies: importance of gene-environment interactions. *Cancer Epidemiol Biomarkers Prev* 2008;17(2):263-267.
8. Boffetta P, Winn DM, Ioannidis JP, et al. Recommendations and proposed guidelines for assessing the cumulative evidence on joint effects of genes and environments on cancer occurrence in humans. *Int J Epidemiol* 2012;41(3):686-704.
9. Kraft P, Yen YC, Stram DO, et al. Exploiting gene-environment interaction to detect genetic associations. *Hum Hered* 2007;63(2):111-119.
10. Bookman EB, McAllister K, Gillanders E, et al. Gene-environment interplay in common complex diseases: forging an integrative model-recommendations from an NIH workshop. *Genet Epidemiol* 2011;35(4):217-225.
11. Hutter CM, Mechanic LE, Chatterjee N, et al. Gene-environment interactions in cancer epidemiology: a National Cancer Institute Think Tank report. *Genet Epidemiol* 2013;37(7):643-657.
12. Mechanic LE, Chen HS, Amos CI, et al. Next generation analytic tools for large scale genetic epidemiology studies of complex diseases. *Genet Epidemiol* 2012;36(1):22-35.
13. Kraft P, Aschard H. Finding the missing gene-environment interactions. *Eur J Epidemiol* 2015;30(5):353-355.
14. Gauderman WJ, Mukheerjee B, Aschard H, et al. Update on the State of the Science for Analytical Methods. *Am J Epidemiol* 2017;in press.
15. Ritchie MD, Davis JR, Aschard H, et al. Incorporation of Biological Knowledge into the Study of GxE. *Am J Epidemiol* 2017;in press.
16. Collins FS, Varmus H. A new initiative on precision medicine. *N Engl J Med* 2015;372(9):793-795.
17. Khoury MJ, Gwinn ML, Glasgow RE, et al. A population approach to precision medicine. *Am J Prev Med* 2012;42(6):639-645.
18. Aschard H. A perspective on interaction effects in genetic association studies. *Genet Epidemiol* 2016;40(8):678-688.
19. Murcray CE, Lewinger JP, Conti DV, et al. Sample size requirements to detect gene-environment interactions in genome-wide association studies. *Genet Epidemiol* 2011;35(3):201-210.
20. Piegorsch WW, Weinberg CR, Taylor JA. Non-hierarchical logistic models and case-only designs for assessing susceptibility in population-based case-control studies. *Stat Med* 1994;13(2):153-162.

21. Mukherjee B, Chatterjee N. Exploiting gene-environment independence for analysis of case-control studies: an empirical Bayes-type shrinkage estimator to trade-off between bias and efficiency. *Biometrics* 2008;64(3):685-694.
22. Li D, Conti DV. Detecting gene-environment interactions using a combined case-only and case-control approach. *Am J Epidemiol* 2009;169(4):497-504.
23. Dai JY, Logsdon BA, Huang Y, et al. Simultaneously testing for marginal genetic association and gene-environment interaction. *Am J Epidemiol* 2012;176(2):164-173.
24. Han SS, Rosenberg PS, Ghosh A, et al. An exposure-weighted score test for genetic associations integrating environmental risk factors. *Biometrics* 2015;71(3):596-605.
25. Kistner EO, Shi M, Weinberg CR. Using cases and parents to study multiplicative gene-by-environment interaction. *Am J Epidemiol* 2009;170(3):393-400.
26. Umbach DM, Weinberg CR. Designing and analysing case-control studies to exploit independence of genotype and exposure. *Stat Med* 1997;16(15):1731-1743.
27. Weinberg CR, Umbach DM. A hybrid design for studying genetic influences on risk of diseases with onset early in life. *Am J Hum Genet* 2005;77(4):627-636.
28. Dai JY, Kooperberg C, Leblanc M, et al. Two-stage testing procedures with independent filtering for genome-wide gene-environment interaction. *Biometrika* 2012;99(4):929-944.
29. Gauderman WJ, Zhang P, Morrison JL, et al. Finding novel genes by testing G x E interactions in a genome-wide association study. *Genet Epidemiol* 2013;37(6):603-613.
30. Hsu L, Shuo J, Dai Y, et al. Powerful cocktail methods for detecting genome-wide gene-environment interaction. *Genet Epidemiol* 2012;36(3):183-194.
31. Kooperberg C, Leblanc M. Increasing the power of identifying gene x gene interactions in genome-wide association studies. *Genet Epidemiol* 2008;32(3):255-263.
32. Murcray CE, Lewinger JP, Gauderman WJ. Gene-environment interaction in genome-wide association studies. *Am J Epidemiol* 2009;169(2):219-226.
33. Gauderman WJ, Thomas DC, Murcray CE, et al. Efficient genome-wide association testing of gene-environment interaction in case-parent trios. *Am J Epidemiol* 2010;172(1):116-122.
34. Chen H, Meigs JB, Dupuis J. Incorporating gene-environment interaction in testing for association with rare genetic variants. *Hum Hered* 2014;78(2):81-90.
35. Jiao S, Hsu L, Bezieau S, et al. SBERIA: set-based gene-environment interaction test for rare and common variants in complex diseases. *Genet Epidemiol* 2013;37(5):452-464.
36. Lin X, Lee S, Christiani DC, et al. Test for interactions between a genetic marker set and environment in generalized linear models. *Biostatistics* 2013;14(4):667-681.
37. Lin X, Lee S, Wu MC, et al. Test for rare variants by environment interactions in sequencing association studies. *Biometrics* 2016;72(1):156-164.
38. Tzeng JY, Zhang D, Pongpanich M, et al. Studying gene and gene-environment effects of uncommon and common variants on continuous traits: a marker-set approach using gene-trait similarity regression. *Am J Hum Genet* 2011;89(2):277-288.
39. Aschard H, Zaitlen N, Tamimi RM, et al. A nonparametric test to detect quantitative trait loci where the phenotypic distribution differs by genotypes. *Genet Epidemiol* 2013;37(4):323-333.
40. Brown AA, Buil A, Vinuela A, et al. Genetic interactions affecting human gene expression identified by variance association mapping. *Elife* 2014;3:e01381.
41. Levene H. Robust tests for equality of variances. In: Olkin I, ed. *Contributions to Probability and Statistics: Essays in Honor of Harold Hotelling*. Stanford: Stanford University Press, 1960:278-292.
42. O'Brien PC. Procedures for comparing samples with multiple endpoints. *Biometrics* 1984;40(4):1079-1087.

43. Pare G, Cook NR, Ridker PM, et al. On the use of variance per genotype as a tool to identify quantitative trait interaction effects: a report from the Women's Genome Health Study. *PLoS genetics* 2010;6(6):e1000981.
44. Wang G, Yang E, Brinkmeyer-Langford CL, et al. Additive, epistatic, and environmental effects through the lens of expression variability QTL in a twin cohort. *Genetics* 2014;196(2):413-425.
45. Yang J, Loos RJ, Powell JE, et al. FTO genotype is associated with phenotypic variability of body mass index. *Nature* 2012;490(7419):267-272.
46. Zhang P, Lewinger JP, Conti D, et al. Detecting gene-environment interactions for a quantitative trait in a genome-wide association study. *Genet Epidemiol* 2016;40(5):394-403.
47. Bhattacharjee S, Chatterjee N, Han S, et al. An R package for analysis of case-control studies in genetic epidemiology. R package version 3.10.0. Bethesda, MD; 2012. (<http://bioconductor.org/packages/release/bioc/html/CGEN.html>).
48. Su Y-R, Di C, Hsu L, et al. A Unified Powerful Set-based Test for Sequencing Data Analysis of GxE Interactions. *Biostatistics* 2016;in press.
49. Boonstra PS, Mukherjee B, Gruber SB, et al. Tests for gene-environment interactions and joint effects with exposure misclassification. *Am J Epidemiol* 2016;183(3):237-247.
50. Cornelis MC, Tchetgen Tchetgen EJ, Liang L, et al. Gene-environment interactions in genome-wide association studies: a comparative study of tests applied to empirical studies of type 2 diabetes. *Am J Epidemiol* 2012;175(3):191-202.
51. Mukherjee B, Ahn J, Gruber SB, et al. Testing gene-environment interaction in large-scale case-control association studies: possible choices and comparisons. *Am J Epidemiol* 2012;175(3):177-190.
52. Thomas DC, Lewinger JP, Murcray CE, et al. Invited commentary: GE-whiz! Ratcheting gene-environment studies up to the whole genome and the whole exposome. *Am J Epidemiol* 2012;175(3):203-207.
53. Ioannidis JP. Why most discovered true associations are inflated. *Epidemiology* 2008;19(5):640-648.
54. Patel CJ, Burford B, Ioannidis JP. Assessment of vibration of effects due to model specification can demonstrate the instability of observational associations. *J Clin Epidemiol* 2015;68(9):1046-1058.
55. Freedman ML, Monteiro AN, Gayther SA, et al. Principles for the post-GWAS functional characterization of cancer risk loci. *Nat Genet* 2011;43(6):513-518.
56. Maurano MT, Humbert R, Rynes E, et al. Systematic localization of common disease-associated variation in regulatory DNA. *Science* 2012;337(6099):1190-1195.
57. John S, Sabo PJ, Thurman RE, et al. Chromatin accessibility pre-determines glucocorticoid receptor binding patterns. *Nat Genet* 2011;43(3):264-268.
58. Nicolae DL, Gamazon E, Zhang W, et al. Trait-associated SNPs are more likely to be eQTLs: annotation to enhance discovery from GWAS. *PLoS genetics* 2010;6(4):e1000888.
59. Boyle AP, Hong EL, Hariharan M, et al. Annotation of functional variation in personal genomes using RegulomeDB. *Genome Res* 2012;22(9):1790-1797.
60. Ernst J, Kheradpour P, Mikkelsen TS, et al. Mapping and analysis of chromatin state dynamics in nine human cell types. *Nature* 2011;473(7345):43-49.
61. Guo Y, Conti DV, Wang K. Enlight: web-based integration of GWAS results with biological annotations. *Bioinformatics* 2015;31(2):275-276.
62. Ward LD, Kellis M. HaploReg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic Acids Res* 2012;40(Database issue):D930-934.

63. Yao L, Tak YG, Berman BP, et al. Functional annotation of colon cancer risk SNPs. *Nat Commun* 2014;5:5114.
64. Barreiro LB, Tailleux L, Pai AA, et al. Deciphering the genetic architecture of variation in the immune response to Mycobacterium tuberculosis infection. *Proc Natl Acad Sci U S A* 2012;109(4):1204-1209.
65. Fairfax BP, Humburg P, Makino S, et al. Innate immune activity conditions the effect of regulatory variants upon monocyte gene expression. *Science* 2014;343(6175):1246949.
66. Grundberg E, Adoue V, Kwan T, et al. Global analysis of the impact of environmental perturbation on cis-regulation of gene expression. *PLoS genetics* 2011;7(1):e1001279.
67. Qiu W, Rogers AJ, Damask A, et al. Pharmacogenomics: novel loci identification via integrating gene differential analysis and eQTL analysis. *Hum Mol Genet* 2014;23(18):5017-5024.
68. Wei P, Yang Y, Guo X, et al. Identification of an Association of TNFAIP3 Polymorphisms With Matrix Metalloproteinase Expression in Fibroblasts in an Integrative Study of Systemic Sclerosis-Associated Genetic and Environmental Factors. *Arthritis & rheumatology (Hoboken, NJ)* 2016;68(3):749-760.
69. French JE, Gatti DM, Morgan DL, et al. Diversity Outbred Mice Identify Population-Based Exposure Thresholds and Genetic Factors that Influence Benzene-Induced Genotoxicity. *Environ Health Perspect* 2015;123(3):237-245.
70. Rasmussen AL, Okumura A, Ferris MT, et al. Host genetic diversity enables Ebola hemorrhagic fever pathogenesis and resistance. *Science* 2014;346(6212):987-991.
71. Ritchie MD, Holzinger ER, Li R, et al. Methods of integrating data to uncover genotype-phenotype interactions. *Nature reviews Genetics* 2015;16(2):85-97.
72. Baurley JW, Conti DV. A scalable, knowledge-based analysis framework for genetic association studies. *BMC Bioinformatics* 2013;14:312.
73. Quintana MA, Conti DV. Integrative variable selection via Bayesian model uncertainty. *Stat Med* 2013;32(28):4938-4953.
74. Quintana MA, Schumacher FR, Casey G, et al. Incorporating prior biologic information for high-dimensional rare variant association studies. *Hum Hered* 2012;74(3-4):184-195.
75. Baurley JW, Conti DV, Gauderman WJ, et al. Discovery of complex pathways from observational data. *Stat Med* 2010;29(19):1998-2011.
76. Pendergrass SA, Frase A, Wallace J, et al. Genomic analyses with biofilter 2.0: knowledge driven filtering, annotation, and model development. *BioData Min* 2013;6(1):25.
77. Sun X, Lu Q, Mukherjee S, et al. Analysis pipeline for the epistasis search - statistical versus biological filtering. *Front Genet* 2014;5:106.
78. Biofilter. 2016. (<https://ritchielab.psu.edu/research/research-areas/expert-knowledge-bioinformatics/methods/biofilter>). (Accessed 9/1/2016 2016).
79. Davis AP, Grondin CJ, Lennon-Hopkins K, et al. The Comparative Toxicogenomics Database's 10th year anniversary: update 2015. *Nucleic Acids Res* 2015;43(Database issue):D914-920.
80. Audouze K, Brunak S, Grandjean P. A computational approach to chemical etiologies of diabetes. *Sci Rep* 2013;3:2712.
81. Wild CP. Complementing the genome with an "exposome": the outstanding challenge of environmental exposure measurement in molecular epidemiology. *Cancer Epidemiol Biomarkers Prev* 2005;14(8):1847-1850.
82. Wild CP. The exposome: from concept to utility. *Int J Epidemiol* 2012;41(1):24-32.
83. Cui Y, DM B, Kwok R, et al. The exposome - embracing the complexity for discovery in environmental health. *Environ Health Perspect* 2016;127(8):A137-A140.

84. Dennis KK, Auerbach SS, Balshaw DM, et al. The importance of the biological impact of exposure to the concept of the exposome. *Environ Health Perspect* 2016;124(10):1504-1510.
85. Dennis KK, Marder ME, Balshaw DM, et al. Biomonitoring in the era of the exposome [available online ahead of print July 6, 2016]. *Environ Health Perspect* 2016;DOI: 10.1289/EHP474.
86. Turner MC, Nieuwenhuijsen M, Anderson K, et al. Assessing the exposome with external measures: commentary on the state of the science and research recommendations. *Annu Rev Public Health* 2017;38, in press, DOI: 10.1146/annurev-publhealth-082516-012802.
87. Gibson G. The environmental contribution to gene expression profiles. *Nature reviews Genetics* 2008;9(8):575-581.
88. van Breda SG, Wilms LC, Gaj S, et al. The exposome concept in a human nutrigenomics study: evaluating the impact of exposure to a complex mixture of phytochemicals using transcriptomics signatures. *Mutagenesis* 2015;30(6):723-731.
89. Shaw JG, Vaughan A, Dent AG, et al. Biomarkers of progression of chronic obstructive pulmonary disease (COPD). *J Thorac Dis* 2014;6(11):1532-1547.
90. Alexander N, Wankerl M, Hennig J, et al. DNA methylation profiles within the serotonin transporter gene moderate the association of 5-HTTLPR and cortisol stress reactivity. *Translational psychiatry* 2014;4:e443.
91. Patel CJ, Chen R, Kodama K, et al. Systematic identification of interaction effects between genome- and environment-wide associations in type 2 diabetes mellitus. *Hum Genet* 2013;132(5):495-508.
92. Patel CJ, Bhattacharya J, Butte AJ. An Environment-Wide Association Study (EWAS) on type 2 diabetes mellitus. *PLoS One* 2010;5(5):e10746.
93. Hall MA, Dudek SM, Goodloe R, et al. Environment-wide association study (EWAS) for type 2 diabetes in the Marshfield Personalized Medicine Research Project Biobank. *Pac Symp Biocomput* 2014:200-211.
94. McGinnis DP, Brownstein JS, Patel CJ. Environment-Wide Association Study of Blood Pressure in the National Health and Nutrition Examination Survey (1999–2012). *Sci Rep* 2016;6:30373.
95. Ahn J, Mukherjee B, Gruber SB, et al. Bayesian semiparametric analysis for two-phase studies of gene-environment interaction. *Ann Appl Stat* 2013;7(1):543-569.
96. Breslow NE, Chatterjee N. Design and analysis of two-phase studies with binary outcome applied to Wilms tumour prognosis. *Journal of the Royal Statistical Society: Series C (Applied Statistics)* 1999;48(4):457-468.
97. Chatterjee N, Chen Y-H. Maximum likelihood inference on a mixed conditionally and marginally specified regression model for genetic epidemiologic studies with two-phase sampling. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* 2007;69(2):123-142.
98. Wacholder S, Weinberg CR. Flexible maximum likelihood methods for assessing joint effects in case-control studies with complex sampling. *Biometrics* 1994;50(2):350-357.
99. Stenzel SL, Ahn J, Boonstra PS, et al. The impact of exposure-biased sampling designs on detection of gene-environment interactions in case-control studies with potential exposure misclassification. *Eur J Epidemiol* 2015;30(5):413-423.
100. Wei P, Tang H, Li D. Functional Logistic Regression Approach to Detecting Gene by Longitudinal Environmental Exposure Interaction in a Case-Control Study. *Genet Epidemiol* 2014;38(7):638-651.
101. Shi M, Umbach DM, Weinberg CR. Family-based gene-by-environment interaction studies: revelations and remedies. *Epidemiology* 2011;22(3):400-407.
102. Dodge HH, Zhu J, Mattek NC, et al. Use of High-Frequency In-Home Monitoring Data May Reduce Sample Sizes Needed in Clinical Trials. *PLoS One* 2015;10(9):e0138095.

103. Garcia-Closas M, Gunsoy NB, Chatterjee N. Combined associations of genetic and environmental risk factors: implications for prevention of breast cancer. *J Natl Cancer Inst* 2014;106(11):dju305.
104. Garcia-Closas M, Rothman N, Figueroa JD, et al. Common genetic polymorphisms modify the effect of smoking on absolute risk of bladder cancer. *Cancer Res* 2013;73(7):2211-2220.
105. Joshi AD, Lindstrom S, Husing A, et al. Additive interactions between susceptibility single-nucleotide polymorphisms identified in genome-wide association studies and breast cancer risk factors in the Breast and Prostate Cancer Cohort Consortium. *Am J Epidemiol* 2014;180(10):1018-1027.
106. Maas P, Barrdahl M, Joshi AD, et al. Breast cancer risk from modifiable and nonmodifiable risk factors among white women in the United States. *JAMA Oncol* 2016.
107. Franceschini N, van Rooij FJ, Prins BP, et al. Discovery and fine mapping of serum protein loci through transethnic meta-analysis. *Am J Hum Genet* 2012;91(4):744-753.
108. Auton A, Brooks LD, Durbin RM, et al. A global reference for human genetic variation. *Nature* 2015;526(7571):68-74.
109. Liu CT, Buchkovich ML, Winkler TW, et al. Multi-ethnic fine-mapping of 14 central adiposity loci. *Hum Mol Genet* 2014;23(17):4738-4744.
110. Wu Y, Waite LL, Jackson AU, et al. Trans-ethnic fine-mapping of lipid loci identifies population-specific signals and allelic heterogeneity that increases the trait variance explained. *PLoS genetics* 2013;9(3):e1003379.
111. Seldin MF, Pasaniuc B, Price AL. New approaches to disease mapping in admixed populations. *Nature reviews Genetics* 2011;12(8):523-528.
112. Pino-Yanes M, Gignoux CR, Galanter JM, et al. Genome-wide association study and admixture mapping reveal new loci associated with total IgE levels in Latinos. *J Allergy Clin Immunol* 2015;135(6):1502-1510.
113. Galanter JM, Gignoux CR, Torgerson DG, et al. Genome-wide association study and admixture mapping identify different asthma-associated loci in Latinos: the Genes-environments & Admixture in Latino Americans study. *J Allergy Clin Immunol* 2014;134(2):295-305.
114. Bustamante CD, Burchard EG, De la Vega FM. Genomics for the world. *Nature* 2011;475(7355):163-165.
115. Chanock SJ, Manolio T, Boehnke M, et al. Replicating genotype-phenotype associations. *Nature* 2007;447(7145):655-660.
116. Kraft P, Zeggini E, Ioannidis JP. Replication in genome-wide association studies. *Stat Sci* 2009;24(4):561-573.
117. van Engeland M, Weijenberg MP, Roemen GM, et al. Effects of dietary folate and alcohol intake on promoter methylation in sporadic colorectal cancer: the Netherlands cohort study on diet and cancer. *Cancer Res* 2003;63(12):3133-3137.
118. Zochbauer-Muller S, Lam S, Toyooka S, et al. Aberrant methylation of multiple genes in the upper aerodigestive tract epithelium of heavy smokers. *Int J Cancer* 2003;107(4):612-616.
119. Cortessis VK, Thomas DC, Levine AJ, et al. Environmental epigenetics: prospects for studying epigenetic mediation of exposure-response relationships. *Hum Genet* 2012;131(10):1565-1589.
120. Bakulski KM, Fallin MD. Epigenetic epidemiology: promises for public health research. *Environ Mol Mutagen* 2014;55(3):171-183.
121. Simonds NI, Ghazarian AA, Pimentel CB, et al. Review of the Gene-Environment Interaction Literature in Cancer: What Do We Know? *Genet Epidemiol* 2016;40(5):356-365.
122. Chen HS, Hutter CM, Mechanic LE, et al. Genetic simulation tools for post-genome wide association studies of complex diseases. *Genet Epidemiol* 2015;39(1):11-19.

123. Psaty BM, O'Donnell CJ, Gudnason V, et al. Cohorts for Heart and Aging Research in Genomic Epidemiology (CHARGE) Consortium: Design of Prospective Meta-Analyses of Genome-Wide Association Studies From 5 Cohorts. *Circ Cardiovasc Genet* 2009;2(1):73-80.
124. Chatterjee N, Shi J, Garcia-Closas M. Developing and evaluating polygenic risk prediction models for stratified disease prevention. *Nature reviews Genetics* 2016;17(7):392-406.
125. VanderWeele TJ, Robins JM. The identification of synergism in the sufficient-component-cause framework. *Epidemiology* 2007;18(3):329-339.
126. VanderWeele TJ. Sufficient cause interactions and statistical interactions. *Epidemiology* 2009;20(1):6-13.

ORIGINAL UNEDITED MANUSCRIPT