

# The Genotype-Tissue Expression (GTEx) project

The GTEx Consortium\*

**Genome-wide association studies have identified thousands of loci for common diseases, but, for the majority of these, the mechanisms underlying disease susceptibility remain unknown. Most associated variants are not correlated with protein-coding changes, suggesting that polymorphisms in regulatory regions probably contribute to many disease phenotypes. Here we describe the Genotype-Tissue Expression (GTEx) project, which will establish a resource database and associated tissue bank for the scientific community to study the relationship between genetic variation and gene expression in human tissues.**

In the past decade, genome-wide association studies (GWAS) have documented a strong statistical association between common genetic variation at thousands of loci and more than 250 human traits<sup>1</sup>. Yet, the functional effects of most GWAS-implicated variants remain largely unexplained. The finding that nearly 90% of these sites occur outside of protein-coding sequences<sup>2</sup> suggests that many associated variants may instead have a role in gene regulation. The careful examination of gene expression and its relationship to genetic variation has thus become a critical next step in the elucidation of the genetic basis of common disease. Cell context is a key determinant of gene regulation; but, to date, the challenge of collecting large numbers of diverse tissues in humans has largely precluded such studies outside of a few easily sampled cell types.

Expression quantitative trait locus (eQTL) mapping offers a powerful approach to elucidate the genetic component underlying altered gene expression<sup>3</sup>. Studies primarily in blood, skin, liver, adipose and brain indicate that eQTLs are common in humans<sup>4–6</sup>. Genetic variation can also influence gene expression through alterations in splicing, noncoding RNA expression and RNA stability<sup>7–9</sup>.

eQTLs regulating nearby or distant genes are commonly referred to as *cis* eQTLs and *trans* eQTLs, respectively<sup>3</sup>. Gene expression is differentially regulated across tissues, and many human transcripts are expressed in a limited set of cell types or during a limited developmental stage. Several studies have reported tissue-specific eQTLs<sup>10,11</sup>, and combining eQTL studies with network analyses across multiple tissues has helped to define complex networks of gene interaction<sup>12,13</sup>.

Complementing eQTL data with information on other molecular phenotypes, for example, from epigenomic assays<sup>14</sup>, on the same tissues and linking to resources such as the Encyclopedia of DNA Elements (ENCODE)<sup>15</sup> will provide a powerful means of dissecting gene-regulatory and higher-order networks across multiple tissues. Analyzing multiple tissues will be important because evaluation of the functional consequences of a disease-associated SNP is ideally performed in a disease-relevant cell context. However, for most tissue types, human biospecimens are very difficult to obtain from living donors (for example, brain, heart and pancreas), and most eQTL studies so far have been performed with RNA isolated from immortalized lymphoblasts or lymphocytes<sup>6</sup> and a few additional readily sampled tissues.

To fully enable this critical next step in the study of the genetic basis of common disease, it will be of enormous value to have a resource of blood samples from individuals who have been comprehensively genotyped

(and eventually completely sequenced), with genotyping data linked to genome-wide gene expression patterns across a wide range of tissue types. Initially, this resource would enable the research community to perform a comprehensive search for eQTLs (both tissue-type specific and across tissue types) and establish their association with disease-associated variants from GWAS or sequencing studies. Eventually, as other molecular phenotypes are added, the relationship between genetic variation and gene expression could expand to include correlations with epigenetics and proteomics data as well as other molecular characteristics. Although such a catalog would have been unthinkable a few years ago, new genomic technologies are now making the problem approachable.

This convergence of unmet scientific need and new technologies prompted a US National Institutes of Health (NIH) workshop held in June 2008 to discuss the advisability and feasibility of a large-scale public resource for human genetic variation and gene expression across tissues. On the basis of the output from this workshop and ongoing consultation, the NIH developed the concept of the GTEx project (**Box 1**). Many of the specifics of the pilot project described here were contributed by funded investigators and were influenced by early, experimental biospecimen collections.

## Design of the GTEx project

The GTEx project of the NIH Common Fund aims to establish a resource database and asso-

\*A full list of members and affiliations appears at the end of the paper. Correspondence should be addressed to G. Getz ([gadgetz@broadinstitute.org](mailto:gadgetz@broadinstitute.org)), J.P.S. ([struewij@mail.nih.gov](mailto:struewij@mail.nih.gov)) or W.W. ([winckler@broadinstitute.org](mailto:winckler@broadinstitute.org)).

## BOX 1 GOALS OF THE GTEx PROJECT

- To create a data resource to enable the systematic study of genetic variation and the regulation of gene expression in multiple reference human tissues
- To provide the scientific community with a biospecimen resource including tissues, nucleic acids and cell lines upon which to determine other molecular phenotypes
- To support and disseminate the results of a study of the ethical, legal and social issues related to donor recruitment and consent
- To support the development of novel statistical methods for the analysis of human eQTLs, alone and in the context of other molecular phenotypes
- To make data available to the research community as rapidly as possible
- To support the dissemination of knowledge, standards and protocols related to biospecimen collection and analysis methods developed during the project

ciated tissue bank in which to study the relationship between genetic variation and gene expression and other molecular phenotypes in multiple reference tissues (**Supplementary Fig. 1**). The GTEx project began with a 2.5-year pilot phase to test the feasibility of establishing a rapid autopsy program that would yield high-quality nucleic acids and robust gene expression measurements. Having met milestones of donor enrollment, RNA quality and eQTL findings, the project is scaling up to include approximately 900 post-mortem donors by the end of 2015. The power to detect eQTLs is dependent on multiple factors that are difficult to quantify precisely, but power estimates over a range of effect sizes and allele frequencies are described (**Fig. 1**).

GTEx donors are identified through low-PMI (post-mortem-interval) autopsy or organ and tissue transplantation settings. To compare the quality of results for tissues derived from autopsy and surgery, a small subset of tissue types routinely discarded during surgical amputation, such as skin, fat and muscle, are also collected. In addition, peripheral blood samples are collected and used as both a source of DNA for whole-genome SNP and copy number variant (CNV) genotyping and to establish lymphoblastoid cell lines. Skin samples are collected from the same region of the lower leg, both for measurement of gene expression and to establish fibroblast cultures. Quantification of gene expression is performed primarily through massively parallel sequencing of RNA, but some pilot-phase tissues were analyzed both by sequencing and by gene expression array to enable a comparison of the results with different technologies. eQTLs are identified and will be made accessible to the scientific community through the National Center for Biotechnology (NCBI) GTEx database and a GTEx data portal. In addition, GTEx raw data will be made available through the database of Genotypes and Phenotypes (dbGaP) on a periodic basis.

GTEx project structure during the pilot phase (**Supplementary Fig. 2**) included

entities for biospecimen acquisition, processing, storage and verification; a study on ethical, legal and social issues (ELSI study); the Laboratory, Data Analysis and Coordinating Center (LDACC); the GTEx-eQTL browser; novel statistical methods development grants; and a brain bank. The scale-up is organized similarly to the pilot; the current structure of the project and information on funding opportunities are available from the NIH Common Fund website.

### Biospecimen acquisition

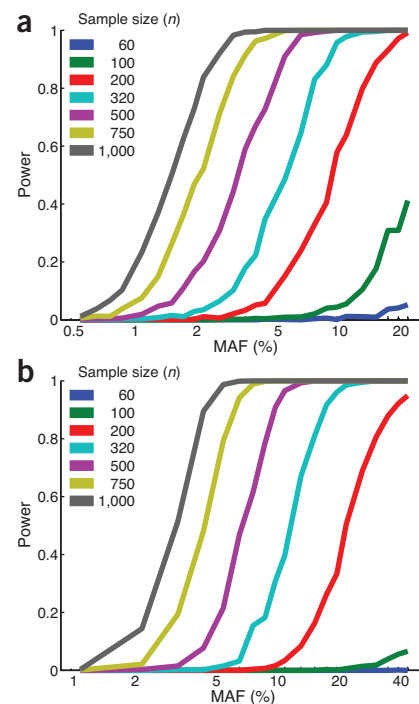
These functions are designed and organized under the Cancer Human Biobank (caHUB) of the National Cancer Institute. caHUB has enrolled under contract several Biospecimen Source Sites (BSSs), a Comprehensive Biospecimen Resource (CBR), a Comprehensive Data Resource (CDR) and pathology and quality management teams to perform acquisition of biospecimens and associated data. Details on all standard operating procedures for donor enrollment and sample collection are available from the caHUB website.

Donors of either sex from any ancestry group are eligible if they are aged 21–70 and if biospecimen collection can start within 24 h of death. There are few medical exclusionary criteria: human immunodeficiency virus (HIV) infection or high-risk behaviors, viral hepatitis, metastatic cancer, chemotherapy or radiation therapy for any condition within the past 2 years, whole-blood transfusion in the past 48 h or body mass index of >35 or <18.5. Each BSS collects, where feasible, aliquots from many predesignated tissue sites and organs (**Supplementary Table 1**), including the brain of deceased donors who were not on a ventilator for the 24 h before death. Immediately after excision, most aliquots are stabilized in a solution containing alcohols (ethanol and methanol), acetic acid and a soluble organic compound that fixes primarily by protein precipitation (PAXgene Tissue, Qiagen) and shipped to the CBR. Only blood samples and

full-thickness skin biopsies are sent unfixed to the LDACC for cell line initiation. The majority of the brain and brainstem are also left unfixed and shipped overnight on wet ice to a brain bank. Further details of donor recruitment and sample collection, including standard operating procedures, are available through caHUB.

### Pathology review and clinical annotation

At the CBR, an aliquot from each sampled tissue is paraffin embedded, sectioned and stained for histological analysis. A dedicated team of pathologists reviews slides from all tissue specimens to verify the organ source and to characterize both general pathological characteristics, such as autolysis, as well as organ-specific pathological states and inflammation. Of course, not all organs will be entirely normal, but donor eligibility is broad and is not restricted to specific diseases or conditions, and it is expected that many organs will be free of major disease processes. An aliquot of each tissue, fixed and stabilized in PAXgene Tissue solution but without paraffin embedding, is sent to the LDACC for molecular analysis. Policies and systems for



**Figure 1** Effect of sample size and MAF on power to detect eQTLs. (a) Power for *cis*-eQTL analysis in which we assume  $\alpha = 0.05/200,000$ , reflecting Bonferroni correction for 200,000 hypotheses based on 20,000 genes and an average of 10 non-redundant SNPs in the region  $\pm 100$  kb of each gene. (b) Power for *trans*-eQTL analysis in which we test 20,000 genes against 5 million SNPs in a total of  $1 \times 10^{11}$  tests with  $\alpha = 5 \times 10^{-13}$ .

accessing stored tissue samples are currently being developed. The CBR's histological sections are viewed as digitally scanned images, allowing precise annotations to be made to indicate where downstream studies, for example, tissue microarray and laser-capture micro-extraction, on selected portions of a specimen can focus (for example, lymphoid nodules in the ileal mucosa or squamous epithelium in the esophageal mucosa).

The clinical data collected for each GTEx donor belong to one of two categories: donor-level data or sample-level data. Donor-level data encompass all clinical measures of the donor, which include basic demographics, medication use, medical history, laboratory test results and the circumstances surrounding the donor's death. These data are collected from the donor (surgical biospecimens) or next of kin (post-mortem biospecimens) and verified against the donor's medical record, when readily available. Summary frequency distributions for clinical variables are available in dbGaP. Sample-level data are attributes belonging to each sample collected and include the tissue type, ischemic time, comments from the prosector and pathology reviewer, and process metadata such as batch ID and the amount of time spent in PAXgene Tissue fixative. Both donor- and sample-level data are examined for quality and completeness before being released.

### Brain bank

Aliquots from single regions of the cortex and cerebellum are sampled and preserved in PAXgene Tissue at the BSS, and the remaining whole brain, with attached brain stem and cervical spinal cord, when possible, is shipped on wet ice to an NIH-funded brain bank. After sectioning of the brain at the brain bank, frozen samples from additional anatomical regions of the brain are analyzed at the LDACC, and the remaining brain is banked for future uses.

### LDACC

The LDACC performs nucleic acid extractions and quality assessment, DNA genotyping and RNA expression analysis. The LDACC integrates results from the molecular analysis with phenotype data, performs eQTL analysis, deposits data into dbGaP and provides a portal for open-access data, standard operating procedures for sample processing and data generation, and results.

DNA is genotyped using the Illumina HumanOmni5M-Quad BeadChip to collect whole-genome SNP and CNV information from each donor's blood sample (or an alternate tissue, if blood is unavailable). The Illumina assay contains over 4 million probes,

with robust coverage of both SNPs and CNVs. DNA is also characterized using the Illumina Infinium HumanExome BeadChip to obtain high-quality SNP calls in coding regions.

A portion of each tissue is processed for RNA and DNA extraction, quantification and quality assessment. Extraction protocols that preserve both mRNA and microRNA are being used and are available from the data portal. For measurement of gene expression during the pilot phase, the LDACC analyzed approximately 1,000 samples using both microarrays (Affymetrix Human Gene 1.1 ST Array) and next-generation RNA sequencing (Illumina HiSeq 2000) to establish the comparability of these methods using post-mortem tissue. RNA sequencing (RNA-seq) uses a 76-base, paired-end Illumina TruSeq RNA protocol, averaging ~50 million aligned reads per sample. This read depth was selected to maximize sequencing value with the budget available and should make it possible to accurately measure moderately expressed transcripts, as well as some with low-level expression, but will have limited ability to accurately quantify rare transcripts and splice isoforms. It should provide gene expression measurements that have equivalent or better accuracy than those obtained with expression arrays and should include a higher dynamic range (with coefficient of variation < 0.1 for at least 12,000 genes; **Supplementary Fig. 3**). RNA-seq allows one to evaluate allele-specific expression in heterozygous individuals, improving the power to identify *cis* regulatory variants. With the target depth of 50 million aligned reads, we expect to have sufficient power to detect allele-specific expression in the top tertile of expressed genes (**Supplementary Fig. 4** and **Supplementary Note**). As the cost of RNA-seq drops, greater read depth will be possible, but, with current resources, the strategy is to maximize the number of samples analyzed.

The fresh-blood and full-thickness skin samples are used to establish Epstein-Barr virus (EBV)-transformed lymphoblastoid cell lines and primary fibroblast cell lines. Because many existing human eQTL studies have used EBV-immortalized cell lines, having these lines in addition to all the other peripheral tissues will allow researchers to evaluate the limitations of using only a lymphoblastoid cell line.

### GTEx-eQTL browser

eQTLs are available and can be queried in browsers hosted both at the LDACC GTEx portal and at NCBI, who will verify the eQTL results provided by the project and both display them and make them available to other genome browsers and the scientific community.

### Statistical analysis development

To promote the analysis of eQTL results across a wide range of human tissues, the NIH funded five centers to develop improved methods for statistical analysis. Investigators funded through this request for applications (RFA) form an analysis consortium that will provide innovative approaches for the analyses of GTEx data and other similar data sets. Investigators also collaborate with the LDACC to perform data quality assessment and quality control before release of the data into dbGaP. The initial GTEx Consortium publications, anticipated in 2013, will include genome-wide analysis of *cis* and *trans* eQTLs, allele-specific expression and splicing quantitative trait loci and a comparison of gene expression results obtained by array and RNA-seq.

### Sample access and molecular analyses

The NIH is interested in making maximal use of this unique biospecimen resource, rich with clinical and genomic information. An access system, including mechanisms for requesting samples, is under development. Except for the fibroblast and lymphoblastoid cell lines, biospecimens are of limited quantity and are non-renewable. Potential uses that are comprehensive (for example, genomic versus single gene or small gene network and proteomic versus single protein or small protein network) and complementary to existing gene expression and variation data are preferred. Scientific questions that are equally well addressed using other sample sets will probably not be suitable, whereas those that take full advantage of the unique aspects of GTEx data, such as the multiple tissues from each donor and the gene expression information, are particularly sought. All data resulting from the analysis of GTEx samples must be made widely available to the scientific community. In addition to scientific review, all proposals to use GTEx samples would also go through a Biospecimen Access Committee (currently being formed).

### Power analysis

To set expectations and guide the design of the full GTEx project, we built a framework to evaluate the statistical power to detect eQTLs. Statistical power depends on various parameters, some known more accurately than others. These parameters include the number of donors, the eQTL effect size and the presence of noise, as well as the significance threshold selected, which is chosen on the basis of the number of hypotheses tested. Assuming we are testing the *cis* eQTL effects of the 10 non-redundant SNPs (on average) in the vicinity ( $\pm 100$  kb of the start site) of each of 20,000 genes, the overall number of hypotheses is

200,000. Therefore, using Bonferroni correction, we set the significance threshold  $\alpha$  to 0.05/200,000. For *trans*-eQTL analysis, a conservative estimate of  $\alpha$  is  $\sim 5 \times 10^{-13}$  (20,000 transcripts tested against 5 million loci). We model the expression data as having a log normal distribution with a log standard deviation of 0.13 within each genotype class (AA, AB, BB). This level of noise is based on estimates from initial GTEx data. The effect size depends both on the minor allele frequency of the SNP (MAF) and the actual log expression change between genotype classes ( $\Delta$ ). **Figure 1a** shows the statistical power of *cis*-eQTL analysis, and **Figure 1b** shows *trans*-eQTL analysis, with each analysis using an ANOVA statistical test as a function of the number of subjects and MAF and assuming  $\Delta = 0.13$  (equivalent to detecting a log expression change similar to the standard deviation within a single genotype class). A final GTEx resource of 900 or more donors would realistically yield  $\sim 750$  samples of any given tissue, as not all organs are available for collection from each donor. At an effective sample size of 750, we would have 80% power to detect *cis* eQTLs with MAF as low as 2% and *trans* eQTLs with MAF as low as 4%. Statistical power may be higher using methods that leverage the fact that multiple tissues are collected and analyzed for each donor. Because the underlying parameters were merely rough estimates, we repeated power analysis with different values (10–20 SNPs and 20,000–100,000 transcripts) and showed that 80% power is achieved for MAFs between 3 and 4% for *cis* eQTLs. For *trans* eQTLs, this range in transcript numbers gives sufficient power with MAFs between 4 and 5% (**Supplementary Fig. 5**).

### Data access and publication policy

GTEx is designated by NIH as a community resource and, as such, aims to share as much of the data (some of which will be unique and identifiable) as rapidly as possible, according to NIH guidelines. It is recognized that quantifying the risk of identifying a donor on the basis of genetic and other information lies on a continuum and is a complex issue dependent on many factors, such as the availability of other sources of data and the evolution of analytical methods<sup>16,17</sup>. Sharing of any information unique to an individual carries a small but difficult-to-define risk of allowing identification of the donor, but this risk must be balanced with the benefits of data sharing to the advancement of science.

Some data from the GTEx project is openly available, meaning that it can be accessed directly through the Internet. However, to reduce risks of sharing potentially identify-

ing data, some data elements are available to the scientific community only through a controlled-access system, dbGaP. Standard operating procedures, details of data collection instruments, histopathological interpretations, molecular data that do not provide direct genetic variation information (for example, data from expression arrays, summary sequence-based gene expression estimates stripped of variant information and eQTL results), laboratory processing variables (for example, cDNA library preparation methods) and a very limited set of medical and sociodemographic variables (for example, sex and age at death in intervals) will be openly available. The LDACC will host an open-access data portal, and specimen acquisition standard operating procedures and information on associated data collection instruments will be available through caHUB. Medical and other epidemiological information, molecular results that contain direct genetic variation information (for example, SNP genotyping files and RNA-seq reads) and summary results that allow accurate inference of allele frequencies<sup>18</sup> will be available only through controlled access. Direct HIPAA (Health Insurance Portability and Accountability Act) identifiers, such as dates that include the month and day, will not be available through either open or controlled access.

Implementation of these data release policies and processes is a topic of ongoing discussion and may need to be modified as risks of identifiability are better quantified for various data types and as the size of the study increases. After initial processing of raw data (such as sequence reads and genotyping files), basic data quality checks are completed by the LDACC and statistical methods investigators, and data are then released immediately through dbGaP. The first dbGaP data release, consisting of data from 62 individuals, occurred in June 2012. For the pilot phase of the project, which concluded in January 2013, the data set comprised genotype data from 190 individuals from whom 1,814 total tissues (from 47 separate tissue sites) were profiled by RNA-seq to a median depth of 80 million aligned reads. These data are in the process of being released to dbGaP, and we anticipate releasing data two to four times per year until the project is completed. We expect total enrollment to increase to over 400 by 2013, to over 700 in 2014 and to approximately 900 by the end of 2015.

The GTEx project falls under the Ft. Lauderdale meeting principles of rapid, prepublication data release. These principles involve publication of a manuscript near the outset to describe the scope and vision of the project and plans to make data available. The

continued success of rapid prepublication data release relies on the scientific community to respect the data producer's interest to publish a full analysis of the data first. Although others are free to analyze GTEx data immediately upon release, the GTEx Consortium envisions publication of both a comprehensive description of the sample acquisition and processing system and a series of genome-wide analyses of genetic variation and gene expression, as described for statistical analysis and the development of methods.

### Ethical, legal and social issues

The GTEx project involves potentially sensitive recruitment, institutional review board (IRB) and consent issues, particularly for deceased donors and their families. The collection of biospecimens from deceased individuals is not legally classified as human subjects research under 45 CFR 46; nonetheless, the depth of the genetic information obtained from the specimens of deceased donors has direct implications for the families of the donors. In recognition of this understanding, sites were required to obtain written or recorded verbal authorization from next of kin for the participation of deceased donors in GTEx, typically through an addendum or modification to an existing authorization form for donation of tissues and organs for research. This authorization included statements common in consent forms, such as the intention to perform genetic analyses, establish cell lines and share data with the scientific community. Work under way is more closely identifying familial concerns and may result in modifications to authorization procedures. Living surgery donors participate only after full, written informed consent is obtained.

In addition, an ELSI study of the consent and authorization process is being carried out at one BSS to assess both the effectiveness of the process in informing participants of the risks and benefits of the study and its potential psychosocial effects on donors and/or their families. The ELSI study is fully integrated with biospecimen collection efforts and will be expanded during the scale-up of the GTEx program.

### Conclusions

A large-scale GTEx resource will be a powerful tool in unraveling the complex patterns of genetic variation and gene regulation across diverse human tissue types. The GTEx project will aid in the interpretation of GWAS findings for translational research by providing data and resources on eQTLs in a wide range of tissues of relevance to many diseases. But the value of a large GTEx resource, especially one that includes other molecular phenotypes,

goes well beyond GWAS follow-up, by providing a deeper understanding of the functional elements of the genome and their underlying biological mechanisms.

**URLs.** Catalog of published GWAS, <http://www.genome.gov/gwastudies>; GTEx LDACC data portal, <http://www.broadinstitute.org/gtex/>; caHUB, <http://cahub.cancer.gov/>; caHUB standard operating procedures, <http://biospecimens.cancer.gov/resources/sops/default.asp>; GTEx project on dbGaP, [http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study\\_id=phs000424](http://www.ncbi.nlm.nih.gov/projects/gap/cgi-bin/study.cgi?study_id=phs000424); US GTEx project on NIH Common Fund, <http://commonfund.nih.gov/GTEX/>; GTEx on National Human Genome Research Institute, <http://genome.gov/gtex/>; NCBI GTEx eQTL Browser, <http://www.ncbi.nlm.nih.gov/gtex/test/GTEX2/gtex.cgi/>; Request for information for the GTEx project, <http://grants.nih.gov/grants/guide/notice-files/NOT-RM-12-028.html>; seeQTL, [http://www.bios.unc.edu/research/genomic\\_software/seeQTL/](http://www.bios.unc.edu/research/genomic_software/seeQTL/); SCAN, <http://www.scandb.org/newinterface/about.html>; US NIH community resource policy for GWAS, <http://gwas.nih.gov/03policy2.html>; Sharing Data from Large-Scale Biological Research Projects: a System of Tripartite Responsibility (Wellcome Trust), <http://www.wellcome.ac.uk/About-us/Publications/Reports/Biomedical-science/WTD003208.htm>; US NIH GTEx working group members, <http://commonfund.nih.gov/GTEX/members.aspx>.

#### ACKNOWLEDGMENTS

The authors would like to acknowledge and thank the donors and their families for making organ and tissue donations, both for transplantation and for the GTEx research study. The authors acknowledge the following funding sources: contracts X10S170, X10S171 and X10172, SAIC-Frederick, Inc., National Cancer Institute and NIH Common Fund, US NIH

to the National Disease Research Interchange, the Roswell Park Cancer Institute and Science Care, Inc.; contract HHSN268201000029C, National Heart, Lung, and Blood Institute and NIH Common Fund, US NIH to the Broad Institute of Harvard and MIT (W.W., contact principal investigator); R01 DA006227-17, National Institute of Drug Abuse, National Institute of Mental Health and National Institute of Neurological Disorders and Stroke, US NIH to the University of Miami School of Medicine (D. Mash, principal investigator); contract 10ST1035, SAIC-Frederick, Inc., National Cancer Institute and NIH Common Fund, US NIH to the Van Andel Institute; prime contract HHSN261200800001E, National Cancer Institute and NIH Common Fund, US NIH to SAIC-Frederick, Inc.; R01 MH090941, National Institute of Mental Health and NIH Common Fund, US NIH to the University of Geneva (E.T.D., contact principal investigator); R01 MH090951, National Institute of Mental Health and NIH Common Fund, US NIH to the University of Chicago (J. Pritchard, principal investigator); R01 MH090937, National Institute of Mental Health, National Human Genome Research Institute, National Heart, Lung, and Blood Institute and NIH Common Fund, US NIH to the University of Chicago (N.J.C., contact principal investigator); R01 MH090936, National Institute of Mental Health and NIH Common Fund, US NIH to the University of North Carolina at Chapel Hill (I.R., contact principal investigator); and R01 MH090948, National Institute of Mental Health, National Human Genome Research Institute and NIH Common Fund, US NIH to Harvard University (J. Liu, contact principal investigator). This research was supported in part by the Intramural Research Program of the National Library of Medicine at the US NIH. The views presented in this article do not necessarily reflect those of the US NIH.

#### AUTHOR CONTRIBUTIONS

**Biospecimen and data collection, processing, quality control, storage and pathological review.** caHUB-BSS: J. Lonsdale, J.T., M. Salvatore, R.P., E.L., S. Shad, R.H., G.W., F.G., N.Y., B.F., M. Moser, E.K., B.G., K. Ramsey, S. Sullivan, J.B., H.M., J.S. and J.F. caHUB-ELSI study: L. Siminoff, H.T., M. Mosavel and L.B. caHUB-CBR: S.J., D.R., D. Maxim, D.F., P.H., E.C., B.B., L.T., E.H. and K.F. caHUB-PRC: L. Sobin, J.R. and P.B. caHUB-CDR: G.K., C.S., D.T., L.Q., K.G. and S.N. caHUB-Operations Management: S.B., A.Z., A. Smith, R.B., K. Robinson, K.V., D.B., M.C., N.D.-M., M. Kennedy, T.E., P.W. and K.E.

**Laboratory analysis, data analysis and study coordination.** K.A., W.W., G. Getz, D.D., D. MacArthur, M. Kellis, A.T., T.Y., E. Gelfand, M.D., Y. Meng and G. Grant. **Brain bank operations.** D. Mash, Y. Marcus and M.B. **Statistical methods development and data analysis.** J. Liu, J.Z., Z.T., E.T.D., T. Lappalainen, R.G., J. Monlong, M. Sammeth, D.K., A.B., S.M., M. McCarthy, M.R., J. Maller, I.R., A.N., F.W., A. Shabalina, N.J.C., D.L.N., E.R.G., H.K.L., A.K., J. Pritchard, M. Stevens, T.F. and X.W. **Database.** M.F., N.S., A. Sturcke and J. Paschal. **Program management.** J.M.A., E.L.W., L.K.D., E.D.G., J.P.S., G.T., S.V., J.T.B., E.J.T., M.S.G., C.N., A.A., D.C., T.R.I., S.E.K., A.R.L., P.K.B., T. Lehner, Y.Y., C.C.C., J.B.V., S. Sawyer, N.C.L., J.D. and H.F.M.

#### COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.



This work is licensed under a Creative Commons Attribution-NonCommercial-Share Alike 3.0

Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-sa/3.0/>.

1. Altshuler, D., Daly, M.J. & Lander, E.S. *Science* **322**, 881–888 (2008).
2. Hindorf, L.A. *et al. Proc. Natl. Acad. Sci. USA* **106**, 9362–9367 (2009).
3. Gilad, Y., Rifkin, S.A. & Pritchard, J.K. *Trends Genet.* **24**, 408–415 (2008).
4. Emilsson, V. *et al. Nature* **452**, 423–428 (2008).
5. Schadt, E.E. *et al. PLoS Biol.* **6**, e107 (2008).
6. Stranger, B.E. *et al. Nat. Genet.* **39**, 1217–1224 (2007).
7. Pickrell, J.K. *et al. Nature* **464**, 768–772 (2010).
8. Pickrell, J.K., Pai, A.A., Gilad, Y. & Pritchard, J.K. *PLoS Genet.* **6**, e1001236 (2010).
9. Borel, C. *et al. Genome Res.* **21**, 68–73 (2011).
10. Petretto, E. *et al. PLoS Comput. Biol.* **6**, e1000737 (2010).
11. Grundberg, E. *et al. Nat. Genet.* **44**, 1084–1089 (2012).
12. Zhong, H. *et al. PLoS Genet.* **6**, e1000932 (2010).
13. Zhao, E. *et al. Mamm. Genome* **20**, 476–485 (2009).
14. Bernstein, B.E. *et al. Nat. Biotechnol.* **28**, 1045–1048 (2010).
15. ENCODE Project Consortium. *Nature* **489**, 57–74 (2012).
16. Craig, D.W. *et al. Nat. Rev. Genet.* **12**, 730–736 (2011).
17. Schadt, E.E., Woo, S. & Hao, K. *Nat. Genet.* **44**, 603–608 (2012).
18. Jacobs, K.B. *et al. Nat. Genet.* **41**, 1253–1257 (2009).

John Lonsdale<sup>1</sup>, Jeffrey Thomas<sup>1</sup>, Mike Salvatore<sup>1</sup>, Rebecca Phillips<sup>1</sup>, Edmund Lo<sup>1</sup>, Saboor Shad<sup>1</sup>, Richard Hasz<sup>2</sup>, Gary Walters<sup>3</sup>, Fernando Garcia<sup>4</sup>, Nancy Young<sup>5</sup>, Barbara Foster<sup>6</sup>, Mike Moser<sup>6</sup>, Ellen Karasik<sup>6</sup>, Bryan Gillard<sup>6</sup>, Kimberley Ramsey<sup>6</sup>, Susan Sullivan<sup>7</sup>, Jason Bridge<sup>7</sup>, Harold Magazine<sup>8</sup>, John Syron<sup>8</sup>, Johnelle Fleming<sup>8</sup>, Laura Siminoff<sup>9</sup>, Heather Traino<sup>9</sup>, Maghboeba Mosavel<sup>9</sup>, Laura Barker<sup>9</sup>, Scott Jewell<sup>10</sup>, Dan Rohrer<sup>10</sup>, Dan Maxim<sup>10</sup>, Dana Filkins<sup>10</sup>, Philip Harbach<sup>10</sup>, Eddie Cortadillo<sup>10</sup>, Bree Berghuis<sup>10</sup>, Lisa Turner<sup>10</sup>, Eric Hudson<sup>10</sup>, Kristin Feenstra<sup>10</sup>, Leslie Sobin<sup>11</sup>, James Robb<sup>11</sup>, Phillip Branton<sup>12</sup>, Greg Korzeniewski<sup>11</sup>, Charles Shive<sup>11</sup>, David Tabor<sup>11</sup>, Liquan Qi<sup>11</sup>, Kevin Groch<sup>11</sup>, Sreenath Nampally<sup>11</sup>, Steve Buia<sup>11</sup>, Angela Zimmerman<sup>11</sup>, Anna Smith<sup>11</sup>, Robin Burges<sup>11</sup>, Karna Robinson<sup>11</sup>, Kim Valentino<sup>11</sup>, Deborah Bradbury<sup>11</sup>, Mark Cosentino<sup>11</sup>, Norma Diaz-Mayoral<sup>11</sup>, Mary Kennedy<sup>11</sup>, Theresa Engel<sup>11</sup>, Penelope Williams<sup>11</sup>, Kenyon Erickson<sup>12</sup>, Kristin Ardlie<sup>13</sup>, Wendy Winckler<sup>13</sup>, Gad Getz<sup>13,14</sup>, David DeLuca<sup>13</sup>, Daniel MacArthur<sup>13,14</sup>, Manolis Kellis<sup>13,15</sup>, Alexander Thomson<sup>13</sup>, Taylor Young<sup>13</sup>, Ellen Gelfand<sup>13</sup>, Molly Donovan<sup>13</sup>, Yan Meng<sup>13</sup>, George Grant<sup>13</sup>, Deborah Mash<sup>16</sup>, Yvonne Marcus<sup>16</sup>, Margaret Basile<sup>16</sup>, Jun Liu<sup>17</sup>, Jun Zhu<sup>18</sup>, Zhidong Tu<sup>18</sup>, Nancy J Cox<sup>19</sup>, Dan L Nicolae<sup>19</sup>, Eric R Gamazon<sup>19</sup>, Hae Kyung Im<sup>19</sup>, Anuar Konkashbaev<sup>19</sup>, Jonathan Pritchard<sup>19,20</sup>, Matthew Stevens<sup>19</sup>, Timothée Flutre<sup>19</sup>, Xiaoquan Wen<sup>19</sup>, Emmanouil T Dermitzakis<sup>21</sup>, Tuuli Lappalainen<sup>21</sup>, Roderic Guigo<sup>22</sup>, Jean Monlong<sup>22</sup>, Michael Sammeth<sup>22</sup>, Daphne Koller<sup>23</sup>, Alexis Battle<sup>23</sup>, Sara Mostafavi<sup>23</sup>, Mark McCarthy<sup>24</sup>, Manual Rivas<sup>24</sup>, Julian Maller<sup>24</sup>, Ivan Rusyn<sup>25</sup>, Andrew Nobel<sup>25</sup>, Fred Wright<sup>25</sup>, Andrey Shabalina<sup>25</sup>, Mike Feolo<sup>26</sup>, Nataliya Sharopova<sup>26</sup>, Anne Sturcke<sup>26</sup>, Justin Paschal<sup>26</sup>, James M Anderson<sup>26</sup>, Elizabeth L Wilder<sup>26</sup>, Leslie K Derr<sup>27</sup>, Eric D Green<sup>28</sup>, Jeffery P Struewing<sup>28</sup>, Gary Temple<sup>28</sup>, Simona Volpi<sup>28</sup>, Joy T Boyer<sup>28</sup>, Elizabeth J Thomson<sup>28</sup>, Mark S Guyer<sup>28</sup>, Cathy Ng<sup>28</sup>, Assya Abdallah<sup>28</sup>, Deborah Colantuoni<sup>28</sup>, Thomas R Insel<sup>29</sup>, Susan E Koester<sup>29</sup>, A Roger Little<sup>29</sup>, Patrick K Bender<sup>29</sup>, Thomas Lehner<sup>29</sup>, Yin Yao<sup>29</sup>, Carolyn C Compton<sup>30</sup>, Jimmie B Vaught<sup>30</sup>, Sherilyn Sawyer<sup>30</sup>, Nicole C Lockhart<sup>30</sup>, Joanne Demchok<sup>30</sup> & Helen F Moore<sup>30</sup>

<sup>1</sup>National Disease Research Interchange, Philadelphia, Pennsylvania, USA. <sup>2</sup>Gift of Life Donor Program, Philadelphia, Pennsylvania, USA.

<sup>3</sup>LifeNet Health, Virginia Beach, Virginia, USA. <sup>4</sup>Drexel University College of Medicine, Philadelphia, Pennsylvania, USA. <sup>5</sup>Albert Einstein Medical Center, Philadelphia, Pennsylvania, USA. <sup>6</sup>Roswell Park Cancer Institute, Buffalo, New York, USA. <sup>7</sup>Upstate New York Transplant Service, Buffalo, New York, USA. <sup>8</sup>Science Care, Inc., Phoenix, Arizona, USA. <sup>9</sup>Virginia Commonwealth University, Richmond, Virginia, USA.

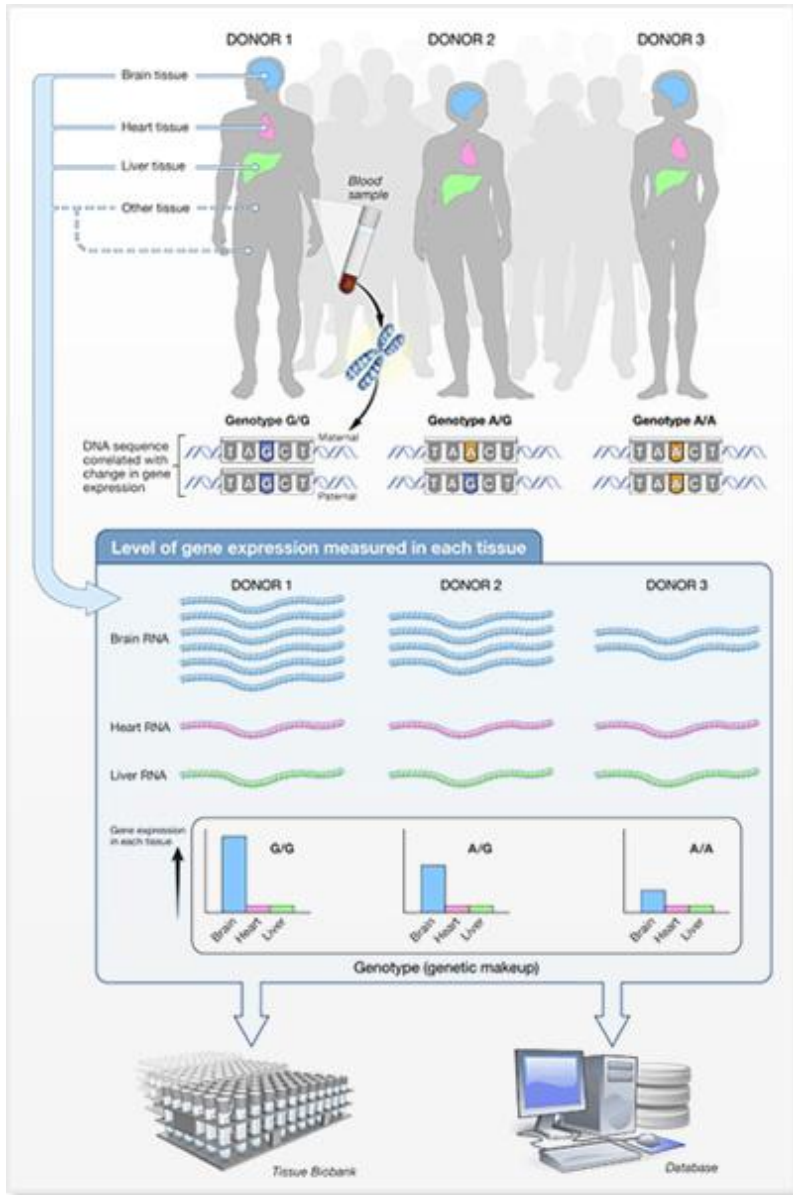
<sup>10</sup>Van Andel Institute, Grand Rapids, Michigan, USA. <sup>11</sup>SAIC-Frederick, Inc., Frederick, Maryland, USA. <sup>12</sup>Sapient Government Services, Arlington, Virginia, USA. <sup>13</sup>The Broad Institute of Harvard and MIT, Cambridge, Massachusetts, USA. <sup>14</sup>Massachusetts General Hospital Cancer Center, Boston, Massachusetts, USA. <sup>15</sup>Massachusetts Institute of Technology, Cambridge, Massachusetts, USA. <sup>16</sup>University of Miami School of Medicine, Miami, Florida, USA. <sup>17</sup>Harvard University, Boston, Massachusetts, USA. <sup>18</sup>Mount Sinai School of Medicine, New York, New York, USA. <sup>19</sup>University of Chicago, Chicago, Illinois, USA. <sup>20</sup>Howard Hughes Medical Institute, Chicago, Illinois, USA. <sup>21</sup>University of Geneva, Geneva, Switzerland. <sup>22</sup>Center for Genomic Regulation, Barcelona, Spain. <sup>23</sup>Stanford University, Palo Alto, California, USA. <sup>24</sup>Oxford University, Oxford, UK. <sup>25</sup>University of North Carolina at Chapel Hill, Chapel Hill, North Carolina, USA. <sup>26</sup>National Center for Biotechnology Information, National Library of Medicine, US National Institutes of Health, Bethesda, Maryland, USA. <sup>27</sup>Division of Program Coordination, Planning and Strategic Initiatives, Office of Strategic Coordination (Common Fund), Office of the Director, US National Institutes of Health, Bethesda, Maryland, USA.

<sup>28</sup>National Human Genome Research Institute, Bethesda, Maryland, USA. <sup>29</sup>National Institute of Mental Health, Bethesda, Maryland, USA.

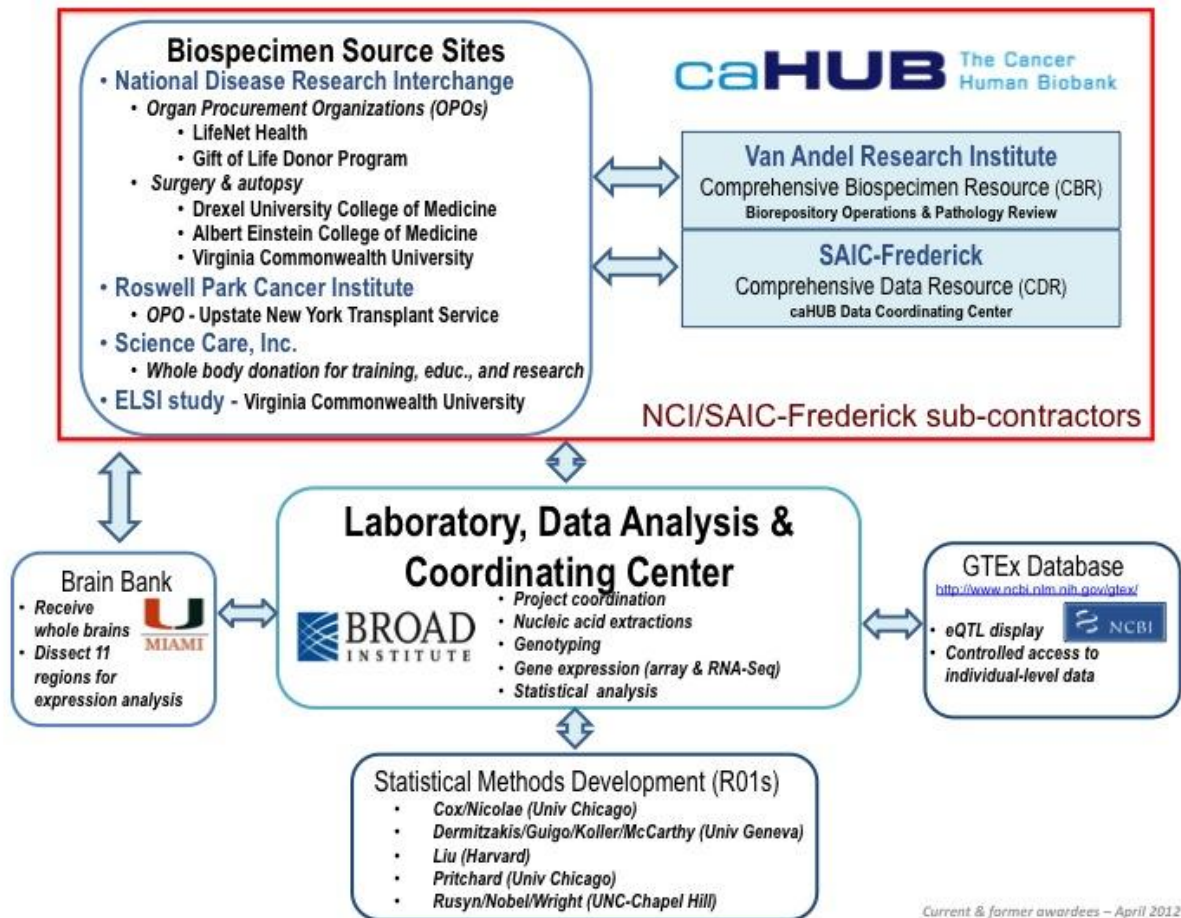
<sup>30</sup>US National Cancer Institute, Bethesda, Maryland, USA.

## Supplementary Figure 1: Design of the Genotype-Tissue Expression (GTEx) Project

Correlations between genotype and tissue-specific gene expression levels will help identify regions of the genome that influence whether and how much a gene is expressed. GTEx will help researchers to understand inherited susceptibility to disease and will be a resource database and tissue bank for many studies in the future.



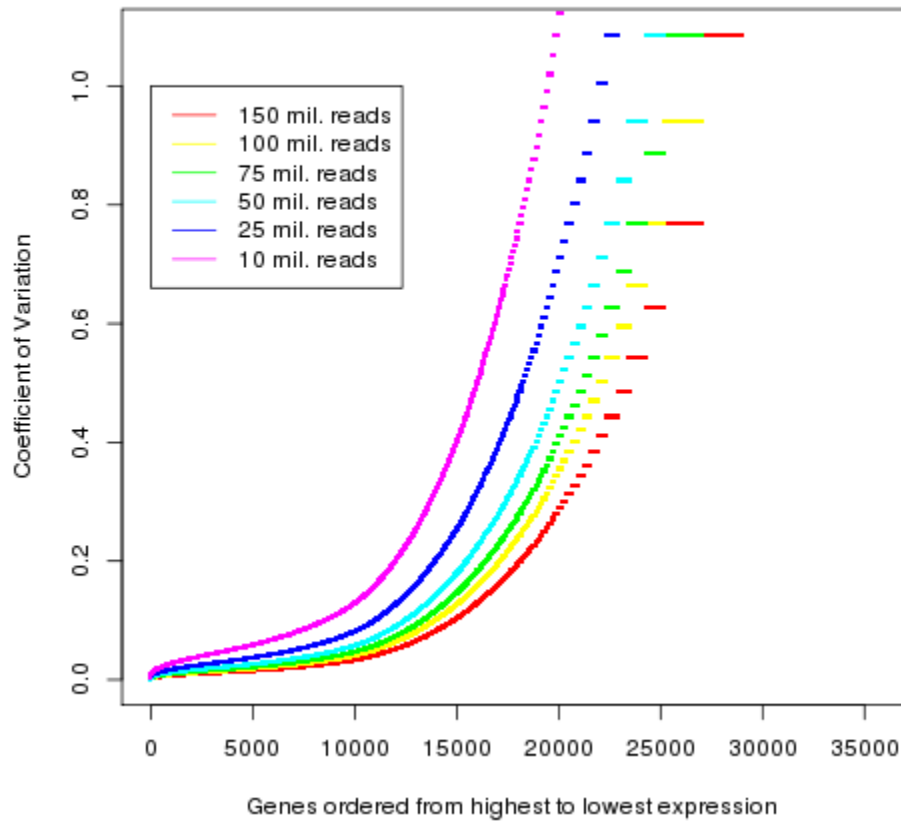
**Supplementary Figure 2: The GTEx Project Consortium Organization.**





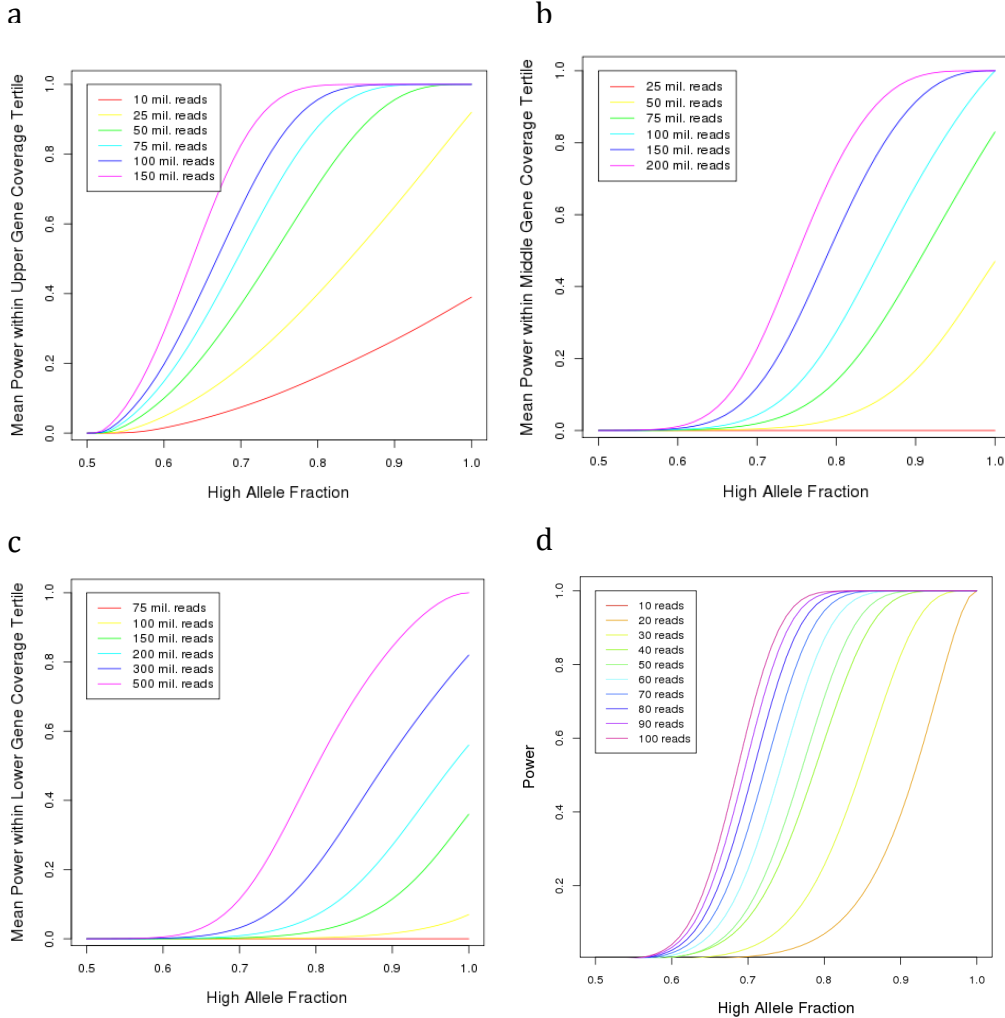
### Supplementary Figure 3 – Sampling variation of gene expression.

Variation in observed expression levels due to the sampling process of RNA-seq. The coefficient of variation (CV) modeled using the observed read distribution in a deep RNA-seq run (176 Million reads) and scaled down to yields between 10 and 150 million reads. The CV (=standard deviation/mean) is estimated as  $\sqrt{1/T}$  following a Poisson distribution. The GTEx target of 50 million mapped reads provides a CV of less than 0.1 for >12,000 genes.



### Supplementary Figure 4 - Power Analysis for Allele Specific Expression (ASE).

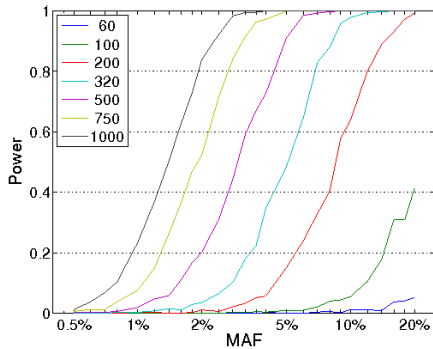
Power analysis for allele specific expression, segregated by tertiles of sequencing coverage per gene (a: upper coverage tertile; b: middle coverage tertile; c: lower coverage tertile). The mean power for genes within a tertile is reported as a function of effect size (expressed as high allele fraction). A range of sequencing yields (from 10 million to 500 million reads) was assessed. 50 million reads are necessary to discover ASE in the upper tertile (a) of genes. For the lower tertile (c), as much as 500 million reads would be necessary to power ASE analysis. (d) Power analysis for allele specific expression based on coverage at heterozygous sites. Power is given as a function of effect size (x axis, expressed as high allele fraction) as well as the number of reads covering the putative ASE site (from 20 to 100).



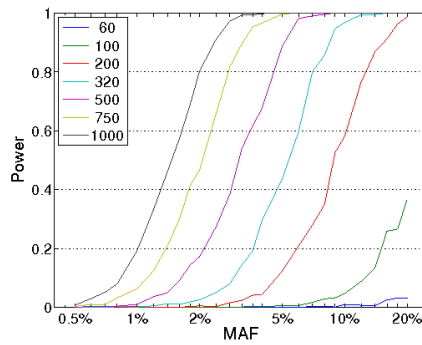
### Supplementary Figure 5 - Data for eQTL Power Analysis.

Increasing the number of *cis* eQTL tests performed has a modest effect on power. The power to detect eQTLs is given as a function of the minor allele frequency. The results for a range of sample sizes (60 to 1000) are given in each plot. Comparing 200,000 tests (a) with two million tests (d) reveals an increase in detectable MAF by 1%. In the case of 750 samples, increasing the tests from 200,000 to 2 million increases the detectable MAF from 3 to 4%. This relatively modest effect indicates that this power analysis is reasonably stable with regard to  $\alpha$ . For *trans* eQTLs (e and f), the number of SNPs is constant at 5 million, but we can vary the number of transcripts from 20,000 to 100,000. Accordingly, increasing the number of tests from 1 billion to 5 billion results in a 1% shift in detectable MAF from 4 to 5%.

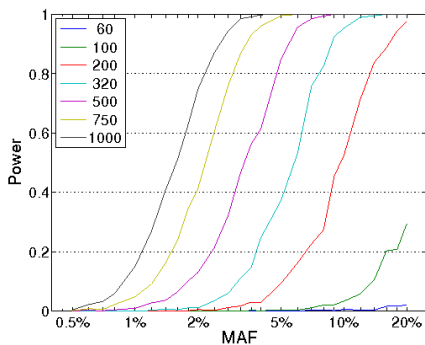
a) 200,000 Tests



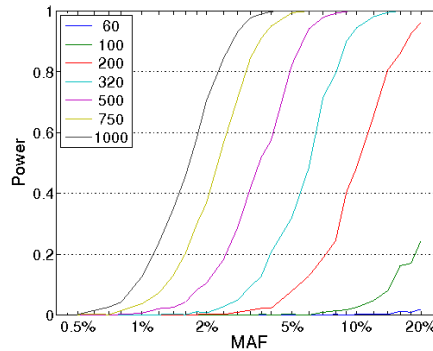
b) 400,000 Tests



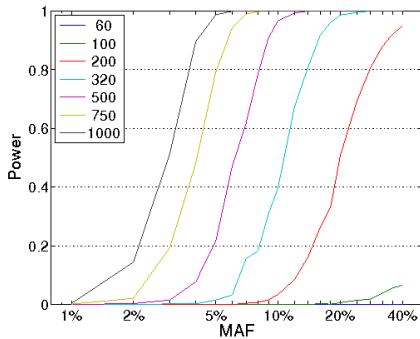
c) 1 Million Tests



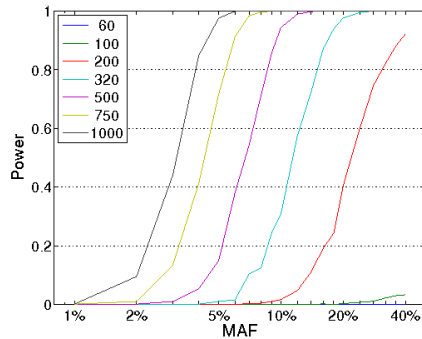
d) 2 Million Tests



e) 100 Billion Tests



f) 500 Billion Tests



**Supplementary Table 1: Complete list of tissues collected in at least one GTEx donor.**

Adipose, Subcutaneous #  
Adipose, Visceral (Omentum)  
Adrenal Gland  
Artery - Aorta (Ascending or thoracic)  
Artery - Coronary  
Artery - Tibial #  
Bladder - Urinary \*  
Blood #  
Brain, Cerebellar Hemisphere  
Brain, Frontal Cortex  
Brain, Hippocampus  
Brain, Substantia nigra  
Brain, Anterior cingulate cortex (BA24)  
Brain, Amygdala  
Brain, Caudate (basal ganglia)  
Brain, Nucleus accumbens (basal ganglia)  
Brain, Putamen (basal ganglia)  
Brain, Hypothalamus  
Brain, Spinal cord (cervical c-1)  
Breast, Mammary Tissue  
Cervix - Ectocervix \*  
Cervix - Endocervix \*  
Colon - Sigmoid  
Colon - Transverse  
Esophagus - Mucosa  
Esophagus - Muscularis  
Fallopian Tube \*  
Gastroesophageal junction  
Heart - Right Atrium  
Heart - Left Ventricle  
Ilium - Terminal (Peyer's patch)  
Kidney - Cortex  
Kidney - Medulla \*  
Liver  
Lung  
Muscle, Skeletal #  
Nerve - Tibial #  
Ovary  
Pancreas  
Pituitary Gland  
Prostate  
Salivary Gland, Minor  
Skin - Sun exposed (Leg) #  
Skin - Not sun exposed (Suprapubic area)  
Spleen  
Stomach  
Testis  
Thyroid  
Uterus  
Vagina

# These tissue samples are also being collected from surgical donors.

\* Tissues collected only during the early pilot phase for which there will likely be fewer than 100 in total.

## Supplementary Note - Power Analysis for Allele Specific Expression (ASE).

The power to detect ASE is primarily influenced by sequencing depth. The following power analysis is aimed at providing a sense of which ASE effect sizes will be detectable given the target sequencing depth of 50 million mapped reads in the GTEx project. We apply the binomial test as a means of detecting whether an observed allelic ratio in heterozygous site-spanning reads is significantly different from 0.5 (after correcting for multiple hypotheses). The coverage at heterozygous sites is not only dependent on the sequencing yield, but is also a function of the level of gene expression. We estimate the coverage distribution of expressed heterozygous sites using preliminary GTEx deep RNA-seq data.

One of the most interesting applications of GTEx data will be allele specific expression (ASE) analysis within eQTL genes. As such, we correct for 500 hypotheses in this analysis representing a rough estimate for the number of eQTL genes containing heterozygous sites within the coding region. Therefore we use as a significant threshold,  $\alpha=0.05/500=10^{-4}$ . We calculated the power within the lower, middle and upper tertiles of the observed coverage distribution. Within each tertile we report the mean power as a function of effect size. This analysis was repeated for a range of sequencing yields between 10 and 150 million reads (**Supplementary Figure 4a-c**, below), scaling the coverage distribution accordingly.

For a yield of 150 millions reads there were 14,933 genes with an average depth of at least 5 reads, or 3393 genes per tertile. For the upper tertile of genes, the targeted 50 million reads will provide sufficient power (>80%) to detect an expected allele fraction of 0.83, i.e. an odds-ratio of ~4.9. This power is achieved when having at least 40 reads covering the heterozygous site (**Supplementary Figure 4d**). Comparable power within the middle tertile would require 150 million reads. Based on current variations in sequencing yield, we expect several dozen GTEx samples (~3%) to achieve this depth providing extra power for the middle tertile. It appears that detecting ASE within the lower tertile would require over 500 million reads which could be achieved in a more targeted approach.