

# **Statistical Methods for High Dimensional Biology**

STAT/BIOF/GSAT 540

Lecture I – course introduction

Paul Pavlidis

January 5 2015

## **Today's topics**

- What the course is about
- Course mechanics
- Introduction to high-dimensional biology

## Your instructors

- Dr. Gaby Cohen-Freue – Assistant Professor of Statistics  
– [gcohen@stat.ubc.ca](mailto:gcohen@stat.ubc.ca)
- Dr. Paul Pavlidis – Professor of Psychiatry/CHiBi  
– [paul@chibi.ubc.ca](mailto:paul@chibi.ubc.ca)
- Dr. Sara Mostafavi – Assistant Professor of Statistics / Medical Genetics – [saram@stat.ubc.ca](mailto:saram@stat.ubc.ca)
- TAs: Evan Durno ([wdurno@gmail.com](mailto:wdurno@gmail.com)), Alice Zhu ([jingyunalice@gmail.com](mailto:jingyunalice@gmail.com))

## Course audience

- Researchers who want to know how to analyze large data sets from biological studies
- Genomics-focused, but information is broadly applicable
- Statistics students might find the math parts easy
- Biology students might find the biology easy
- We are counting on you to help make it work: help your peers!

## Prerequisites

Officially, none. But:

- **Statistics** – You should have already taken university level “Statistics 101”. You’ll get a refresher, but you should be prepared to get comfortable thinking about things like “probabilities” and “specificity”.
- **Biology – No requirements**, but you are expected to learn things like the difference between a DNA and RNA and a gene and a genome. We assume you are here because you are interested in biology and will pick it up.
- No **R** experience required but you must be prepared to do a lot of self-guided learning.
- You’ll use your own computer to run R. If you can’t install R on your computer, ask us for options.

## What you can expect to learn

- Conceptual and practical knowledge you need to handle large biological data sets
  - Less about specific types of data, more about generally applicable approaches and principles
- You will be able to critically evaluate analyses in the literature
- Implementation of analyses using the R/Bioconductor computing environment

## Not about:

- Formal mathematical theory underpinning the approaches
- Gory details of how to analyze any particular type of data at a low level

## Topics covered

Probability foundations  
Exploratory data analysis  
Data QC and preprocessing  
Basic statistical inference (“one gene at a time”)  
Large-scale inference (“genome-wide”)  
Count-based data (e.g. RNA-seq) analysis  
DNA methylation analysis  
Principal Component Analysis  
Clustering  
Classification  
Resampling and bootstrap  
Model selection and regularization  
Gene sets and gene networks

## Course mechanics

## Course web site

<http://stat540-ubc.github.io/>

- Lecture notes
- Lab notes
- Assignments

Much interaction via Github (discussions, submission)

<https://github.com/STAT540-UBC>

## Lectures

- ESB 4192
- Lectures shared among three professors
- Notes provided on web before class

## Sections/Labs

- Wednesdays in room ESB 1042
- Officially from 12-1, but we will start at 11
  - 11-12: Lab available
  - 12-1: TA Office hour (this week: Mol. Bio. Primer)
- Self-guided exercises to help you learn to use R for analysis.
- Using your own computer (other options possible)
- Exercise material will be made available ahead of time
- Towards end of course, more time devoted to working on group projects.

## Readings

- No textbook, but we can give suggestions
- Lectures often come with suggested background papers (reviews or primary literature)
- Helpful to access journals online (e.g. via the UBC VPN)

– <http://it.ubc.ca/services/email-voice-internet/myvpn/setup-documents>

## Evaluation

- **Homeworks**
  - Two assignments worth 20 points each
  - +5 points each for peer evaluation
- **Group project**
  - Planning + project + poster session – 40 points
- 10 Points for “other”
  - e.g. Preparedness, participation.

## Homework assignment

- One for February, one for March.
- Involve detailed analysis of real data
- Deliverables include a short report and R code
- Two weeks from assignment to due date
- Lateness penalties

## Group projects

- Starts today – start thinking about it
- A few minutes for group project pitches on Jan 19 and Jan 21.
- Form groups by Fri Jan 23 (3-4 people)
- Friday Jan 30: initial project proposals
- Feedback to groups Feb 13 – Proposals finalized by Feb ~15
- Work on projects over rest of term
- Final session of the course is the poster session



## Group projects: where do they come from?

- Historically, almost all projects have been based on a data set provided by a student (i.e., collected in their lab).
- Occasionally, instead based on an idea from a student, using published data.
- If you need help thinking up an idea for a project let us know. But this has never been needed before (beyond refinement). If you are unsure of where you are going to get a project from, wait until you hear the project pitches.

## Examples of past group projects

- Genomic copy number alterations for prognosis of prostate cancer
- Learning about proteins from other proteins: Protein Database Prediction
- Conditional epistasis profiling in yeast
- Epigenetic biomarkers for cancer diagnosis
- Comparative metagenomics : metabolic potential
- Epigenome and transcriptome in rice strains
- Analysis of HPV E2 protein on host gene expression
- Effects of Mutations in Histone Modifying Enzymes on Gene Expression Profiles
- Methodological considerations in analysis of Illumina Infinium methylation data
- Gene expression in invasive ragweeds
- Modeling time-course expression of SET domain-containing genes in mouse embryos
- Gene expression in blood of humans with asthma challenged with allergen

2011 and 2012 project titles, paraphrased

## High-dimensional biology

1. What is it
2. What kinds of methods are used to analyze it
3. Some examples

## Collecting data the low-dimensional way

- Pick one variable (e.g. “activity of a protein”) and study it under various conditions.
- Repeat this for another variable
- Usually “hypothesis-driven”
- Powerful, but knowledge accumulates slowly and synthesis is difficult

## Biology is complicated

- Thousands of “parts”
- Limitations of the “one thing at a time approach” – how do the parts work together?
- Technology enabling increasingly detailed analyses – measure many things in parallel
- Drawback: Fishing expeditions?

## Defining “high dimensional”

- Large number of features measured in each sample/subject/individual (“high content”) – Genes, proteins, DNA sites, brain regions, etc.
- Not *usually* talking about huge numbers of samples (e.g. individuals studied) – often 10s, but can be 1000s (some genetics studies)
- Studies can sometimes be “non-hypothesis driven”

## Example of a question answered with a high-dimensional approach

- Tumor type A is deadly and type B is more easily treatable (but still bad)
- Telling A from B is difficult
  - Cells look the same, etc. – we only find out by seeing what happens to the patients.
- We know that cancer is a “gene” disease

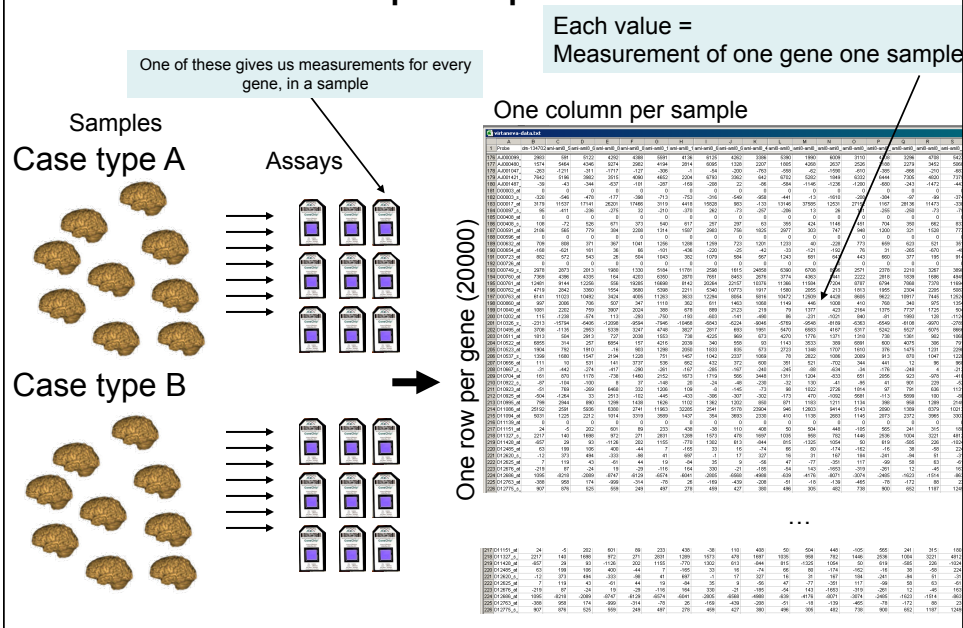
### Questions:

- Where is the difference?
- Can we find new targets for drugs or for diagnosis?
  - (Drug targets are usually proteins, encoded by genes)

## Looking for insight from genomics

- Since cancer is a disease of genes, let's look at the genes - not just one, but all of them
- We are hypothesizing that there is *some* difference in genes between the two types, if only we could find it
- But we're not starting with a *specific* hypothesis. We're going to test thousands of hypotheses
- In this example, we're going to look at “gene expression levels” – a measure of “how active” is each gene.

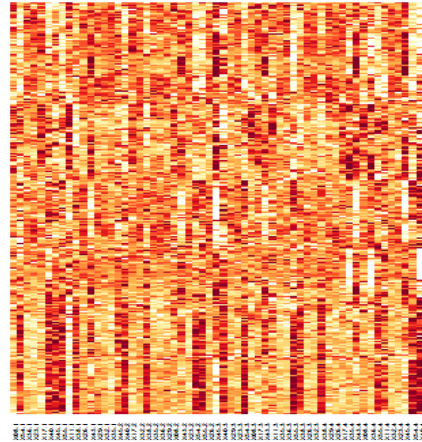
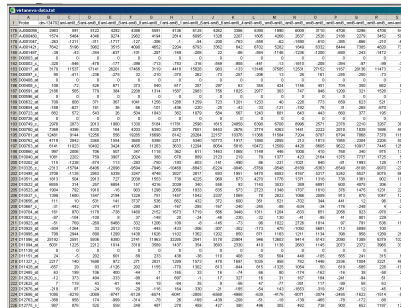
## Example experiment



## A partial list of things to assay

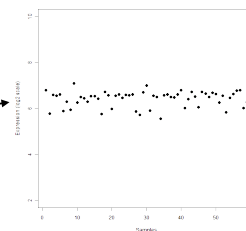
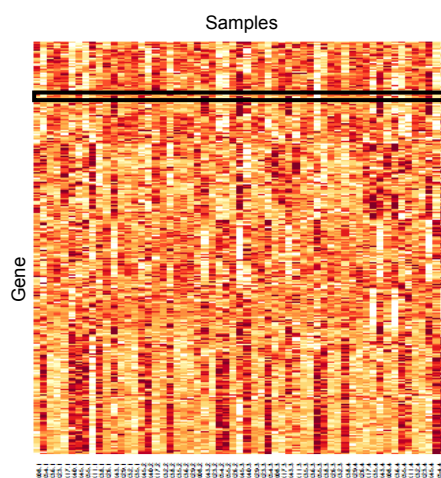
- DNA/Chromatin
  - Genotypes, copy numbers ("mutations" and variants)
  - DNA methylation
  - Chromatin state (histone marks, transcription factors ...)
- RNA
  - Quantification of transcripts (protein coding, non-coding)
  - Transcript variants (splicing, editing)
- Proteins
  - Detection, Quantification
  - Binding and complexes
- Metabolites and other small molecules
- Phenotypic screens
  - RNAi (etc.)
  - Genetic interactions
- Cellular composition of a sample (cell types)

## Alternative representation



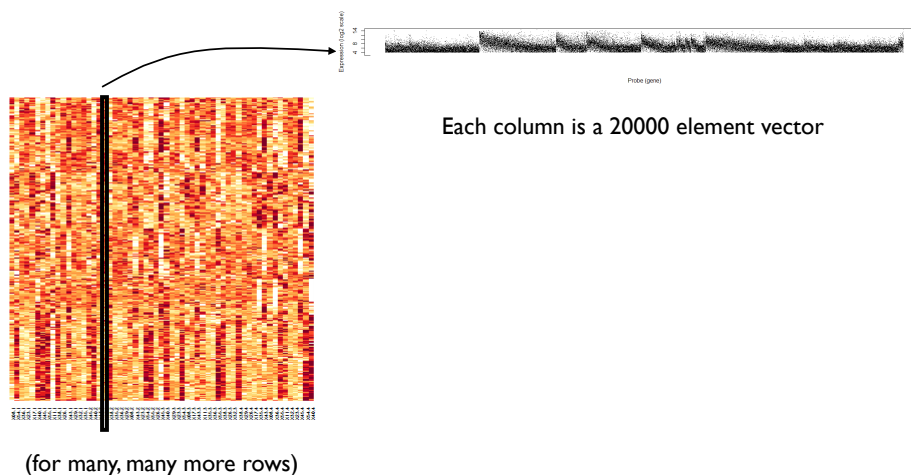
Lighter colours mean higher levels  
of gene expression ("activity")  
Only show part of the data!

## Profile for a gene



Profile for a gene. This is a  
59-element vector

## Profile for a sample

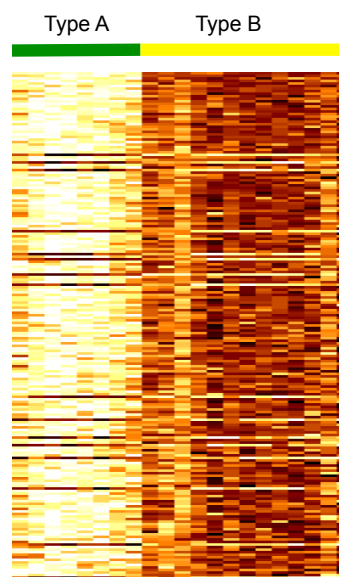


\*This is a schematic. The graph and color map don't match

## One type of analysis

- I've ranked the genes by how different they are between types A and B (t-statistic)
- Only the first few genes are shown
- Though it can be a lot more complicated, most "high-dimensional" studies boil down to something like this, at least in part

**What's the big deal?**



## Pitfalls and challenges

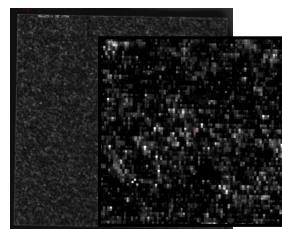
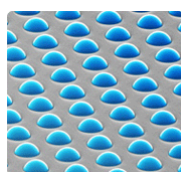
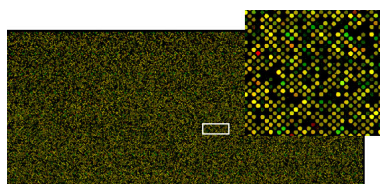
- Signals can be small and buried in lots of non-signals; False positives are a danger.
- Need to detect outliers, batch effects and other confounds
- Can we make better use of the fact that we're testing 20,000 genes than just doing a t-test on each one?
- Data sets (and questions) are often much more complex
- Getting just a list of “hits” isn't enough – can we understand something more about the “system”

## High-dimensional technologies

- DNA & RNA sequencing
  - Transcriptomes, exomes, full genomes
- Complex gene library construction
  - Expression vectors, protein tags, knockdowns
- Microarrays and other robotic/parallel tech.
  - Screens, high-content assays ...
- Mass spectroscopy
- Flow cytometry
- Imaging



## Microarrays



Agilent SurePrint

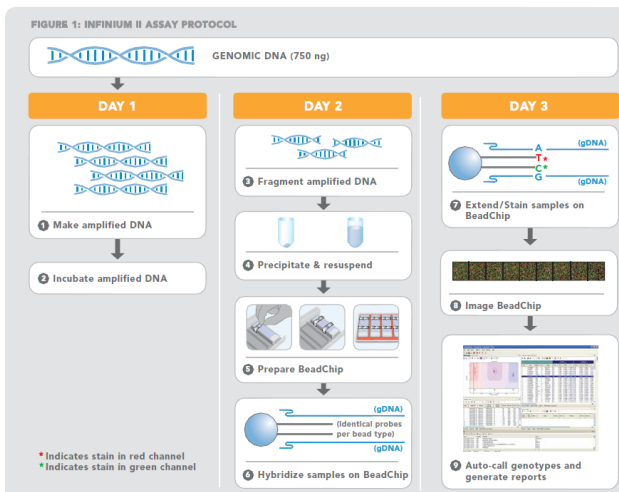
Illumina Beadarray

Affymetrix Genechip

## SNP arrays

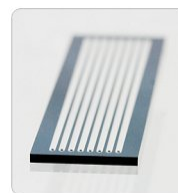
- Similar idea to the RNA arrays, but hybridize genomic DNA, and probe is designed to be sensitive to the allele
- Intensities are converted into a “call”, with a quality score.
- Low-quality calls are usually simply treated as missing data.
- DNA methylation arrays involves specialized versions of these (+bisulphite conversion)

Illumina



## Sequencing-based assays

- Instead of using hybridization to a designed probe, determine the DNA sequence of the sample
- Several competing platforms
- Genotyping: Compare to a reference
- RNA: quantify how many times you see a sequence



Up to eight samples can be loaded onto the flow cell for simultaneous analysis on the Illumina Genome Analyzer.

Illumina HiSeq

## Analysis modes

- What is the general toolkit available for the analysis of data?
- How are these specialized for high-dimensional data?

## Exploratory analysis

- The first thing you do with your data
- Graphs and other visualizations, often combined with data reduction
- Use to spot problems, formulate hypotheses
- Often rely on power of human brain
- Data reduction essential to make exploration tractable for large data sets, even then it can be a challenge
- Follow up with more formal analysis

## Model fitting and hypothesis testing

- Formally test a specific question about the data
- Is what I see “statistically significant”?
- False positives are a major risk in large data sets
- Can exploit repeating structure of the data to improve ability to find true positives

## Unsupervised learning

- “Learn” undiscovered groupings in the data
- Clustering -- how do my samples or features group together?
- Useful as an exploratory technique as well as “data mining” when backed with quantitative analyses
- Example: Finding previously unknown groups of subjects based on a gene profile

## Supervised learning

- Can I predict an unmeasured feature of a sample from a measured one?
- Less common than unsupervised learning, most commonly used in clinically-oriented settings – development of biomarkers
- Example: predicting tumour drug response based on gene profiles

## Other methods

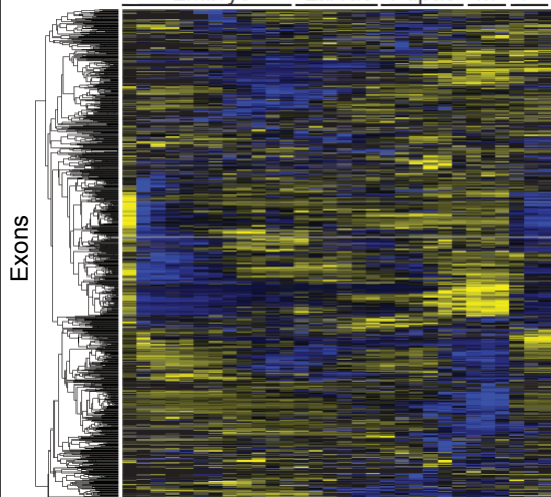
- Many analyses just give a list of genes
- “Downstream” analysis needed to make sense of it - “biological interpretation”
  - Overlay/combine/compare with other data
  - Transform one data set into another type of data at a different granularity
    - Genes → pathways
- Usually these end up returning to exploratory etc. modes

## More examples

- Illustrate some real-life cases of high-dimensional data
- We hope to teach you enough in the course to do at least primitive versions of these analyses
- ... or at least be able to read the papers
- ... even if it's a type of experiment we don't teach in detail.

Developmental stage

Embryo	Larvae	Pupae	M	F
--------	--------	-------	---	---



Colour indicates use pattern of the exon

<http://www.nature.com/nature/journal/vaop/ncurrent/full/nature09715.html>

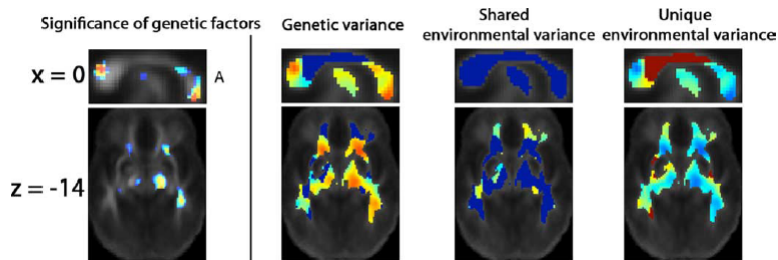
- 30 developmental stages
- Analysis at the exon level
- Heat maps
- Clustering
- GO enrichment
- Generates many new hypotheses

# The developmental transcriptome of *Drosophila melanogaster*

Freeman R. Grave<sup>1</sup>, Angela N. Boyce<sup>2</sup>, Joseph W. Carlen<sup>3</sup>, Michael G. Duff<sup>4</sup>, Jane M. Landislin<sup>5</sup>, LiYang<sup>6</sup>, Carlo G. Ariotti<sup>7</sup>, Marjorie L. van Rantwijk<sup>8</sup>, J. Robert Smith<sup>9</sup>, Justin B. Smith<sup>10</sup>, Lucy Chertan<sup>11</sup>, Carrie A. Davis<sup>12</sup>, David A. D'Amico<sup>13</sup>, Jennifer L. Davis<sup>14</sup>, John H. Matlock<sup>15</sup>, Nicholas A. Matlock<sup>16</sup>, David M. Miller<sup>17</sup>, Joseph M. Miller<sup>18</sup>, Chris Zakaria<sup>19</sup>, Deyang Zhang<sup>20</sup>, Marco Balarska<sup>21,22</sup>, Sandrine Pardin<sup>23</sup>, Brian E. Bach<sup>24</sup>, Richard E. Green<sup>25</sup>, Ann Hammonds<sup>26</sup>, Lichun Jiang<sup>27</sup>, Mark Kapranov<sup>28</sup>, Laura Langdon<sup>29</sup>, Norbert Perle<sup>30</sup>, Jeremy S. Sandler<sup>31</sup>, Kenneth H. Wain<sup>32</sup>, Aaron Willingham<sup>33</sup>, Yu Zhang<sup>34</sup>, Yi Zhou<sup>35</sup>, James Andrews<sup>36</sup>, Peter J. Bickel<sup>37</sup>, Steven E. Brenner<sup>38</sup>, Michael R. Brent<sup>39</sup>, Peter Chertan<sup>40</sup>, Thomas R. Gingeras<sup>41</sup>, Roger A. Hinkins<sup>42</sup>, Thomas C. Kaufman<sup>43</sup>, Brian Oliver<sup>44</sup> & Susan E. Goldstein<sup>45</sup>



- Relatedness (twins)
- IQ



\* Fractional anisotropy / White matter integrity

22712 • The Journal of Neuroscience, February 18, 2009 • 29(7):22712–22724

Behavioral/Systems/Cognitive

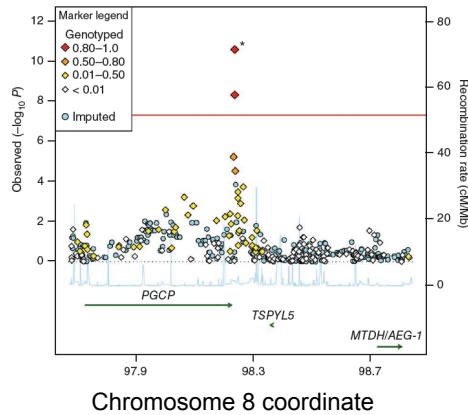
Genetics of Brain Fiber Architecture and Intellectual Performance

Ming-Chang Chiang,<sup>1</sup> Marina Barysheva,<sup>1</sup> David W. Shattuck,<sup>1</sup> Agatha D. Lee,<sup>1</sup> Sarah K. Madsen,<sup>1</sup> Christina Avedissian,<sup>2</sup> Andrea D. Klunder,<sup>1</sup> Arthur W. Toga,<sup>1</sup> Katie L. McMahon,<sup>2</sup> Greig I. de Zubicaray,<sup>1</sup> Margaret J. Wright,<sup>3</sup> Anuj Srivastava,<sup>4</sup> Nikolay Balov,<sup>4</sup> and Paul M. Thompson<sup>1</sup>

<http://www.ncbi.nlm.nih.gov/pubmed/19228974>

- 92 identical or fraternal twins
- 1.5 million voxels per subject
- Linear models, factor analysis
- Multiple test correction
- Heat maps
- Genetics explains 80% of the variance
- Brain structures correlated with IQ  $\sim 0.3$

## Example 3: Genetics of migraine



- 13,500 individuals
- 429,912 DNA sites tested
- Analysis\* with multiple test correction to identify markers associated with migraine
- One site is "A" in 0.267 of the migraine-affecteds but only 0.216 of the controls
- 40% higher risk of migraine if you have "A"

Genome-wide association study of migraine implicates a common susceptibility variant on 8q22.1

Verner Anttila<sup>1,2,3,4</sup>, Heinn Stefansson<sup>5</sup>, Mikko Kallela<sup>6</sup>, Unda Todt<sup>6,8</sup>, Gisela M Terwindt<sup>7</sup>, M Stella Calafato<sup>1,8</sup>, Dale R Nyholt<sup>9</sup>, Antigone S Dimas<sup>10,11</sup>, Tobias Freilinger<sup>12,13</sup>, Bertram Müller-Myhsok<sup>14</sup>, Ville Arto<sup>15</sup>, Michael Innes<sup>1,12</sup>, Kirsi Alakurki<sup>12</sup>, Mari A Kaunisto<sup>1,16</sup>, Eija Hämmäläinen<sup>12</sup>, Rosalie de Vries<sup>17</sup>, Antine H Stam<sup>7</sup>, Claudia M Weller<sup>18</sup>, Axel Heinze<sup>17</sup>, Katja Heinze-Kuhn<sup>17</sup>, Ingrid Goebel<sup>18</sup>, Guntram Brock<sup>18</sup>, Hartmut Göbel<sup>17</sup>, Stacy Steinberg<sup>19</sup>, Christiane Wolf<sup>14</sup>, Asgeir Björnsson<sup>5</sup>, Gretar Gudmundsson<sup>19</sup>, Malene Kirchmann<sup>19</sup>, Anne Haug<sup>19</sup>, Thomas Werge<sup>20</sup>, Jean Schoen<sup>21</sup>, Johan G Eriksson<sup>22,23</sup>, Knut Hagen<sup>24</sup>, Lars Stovner<sup>25</sup>, H-Erich Wichmann<sup>26-28</sup>, Thomas Meitinger<sup>29,30</sup>, Michael Alexander<sup>31,32</sup>, Susanne Moebus<sup>33</sup>, Stefan Schreiber<sup>34,35</sup>, Yuri S Aulchenko<sup>36</sup>, Monique M B Breteler<sup>36</sup>, Andre G Uitterlinden<sup>37</sup>, Albert Hofman<sup>38</sup>, Cornelia M van Duijn<sup>38</sup>, Pieter Tikkas-Klemola<sup>39</sup>, Salli Vepäläinen<sup>40</sup>, Susanne Lucae<sup>41</sup>, Federica Tuzzi<sup>42</sup>, Pierandrea Muglia<sup>39,40</sup>, Jeffrey Barrett<sup>43</sup>, Jaakko Kaprio<sup>34,44</sup>, Markus Fiekkilä<sup>45</sup>, Leena Palomaa<sup>1,2,46,47</sup>, Kari Stefansson<sup>5</sup>, John-Anker Zwart<sup>48,49</sup>, Michel D Ferrari<sup>50</sup>, Jes Olesen<sup>51</sup>, Mark Daly<sup>42</sup>, Majja Wessman<sup>2,16</sup>, Arn M J M van den Maagdenberg<sup>7,15</sup>, Martin Dichgans<sup>12,15</sup>, Christian Kubisch<sup>52,53,54,55</sup>, Emmanouil T Dermizakis<sup>41</sup>, Rune R Frants<sup>56</sup> & Aurora Pakiet<sup>1,2,45,46,47</sup> for the International Headache Genetics Consortium

VOLUME 42 | NUMBER 10 | OCTOBER 2010 NATURE GENETICS

<http://www.nature.com/ng/journal/v42/n10/abs/ng.652.html>

\* Cochran-Mantel-Haenszel, like Fisher's exact test