

Computational challenges in genome wide association studies: data processing, variant annotation and epistasis

Pablo Cingolani

PhD.

School of Computer Science

McGill University

Montreal, Quebec

March 2015

A thesis submitted to McGill University in partial fulfillment of the
requirements of the degree of Doctor of Philosophy

Pablo Cingolani 2015

TABLE OF CONTENTS

1	Introduction	1
1.1	Genomes and genetic variants	2
1.2	DNA and disease risk	6
1.2.1	Heritability and Missing heritability	7
1.2.2	Conclusions	9
1.3	Identification of genetic variants	9
1.3.1	Sequencing data	10
1.3.2	Read mapping	11
1.3.3	Variant calling	14
1.4	Functional annotations of genomic variants	17
1.4.1	Variant types	19
1.4.2	Types of genetic annotations	19
1.4.3	Coding variant annotation	21
1.4.4	Loss of function variants	23
1.4.5	Variants with low or moderate impact	25
1.4.6	Non-coding variant annotation	27
1.4.7	Impact assessment of non-coding variants	31
1.4.8	Clinical effect of variants	34
1.4.9	Data structures and computational efficiency	36
1.4.10	Conclusions	37
1.5	Genome wide association studies	37
1.5.1	Single variant tests and models	38
1.5.2	Multiple variant tests	39
1.5.3	Continuous traits and correcting for co-factors	40
1.5.4	Population structure	40
1.5.5	Population as confounding variable	41
1.5.6	Common and Rare variants	42
1.5.7	Rare variants test	43
1.5.8	Epistasis	43
1.5.9	Epistatic GWAS	45
1.6	Thesis roadmap and Contributions	47
2	BigDataScript: A scripting language for data pipelines	53
2.1	Preface	53
2.2	Introduction	54
2.3	Methods	58

2.3.1	Language overview	59
2.3.2	Abstraction from resources	59
2.3.3	Robustness	61
2.3.4	Other features	63
2.3.5	BDS implementation	65
2.4	Results	69
2.5	Discussion	73
3	A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of <i>Drosophila melanogaster</i> strain $w^{1118}; iso - 2; iso - 3$	75
3.1	Preface	75
3.2	Abstract	76
3.3	Introduction	77
3.4	Results	79
3.5	Discussion	94
3.6	Methods	99
3.7	Acknowledgements	102
3.8	Epilogue	102
3.8.1	Data structures for annotations	103
4	Epistatic GWAS analysis	104
4.1	Preface	104
4.2	Abstract	105
4.3	Introduction	106
4.4	Methods	110
4.4.1	Substitution model for pairs of interacting amino acids	110
4.4.2	Calculating likelihood of individual and pairs of alignment columns	112
4.4.3	GWAS model	113
4.4.4	Putting it all together	117
4.5	Results	118
4.5.1	Co-evolutionary substitution models	119
4.5.2	Co-evolutionary model validation	120
4.5.3	Epistatic GWAS analysis	123
4.6	Discussion	127
5	Conclusions	130
5.1	Contributions	130
5.2	Future work	132
5.3	Perspectives	133
References	135

CHAPTER 1

Introduction

How does one's DNA influence their risk of getting a disease? Contrary to popular belief, your future health is not "hard wired" in our DNA. Only in a few diseases, referred as "Mendelian diseases", are there well known, almost certain, links between genetic mutations and disease susceptibility. For the majority of what are known as "complex traits", such as cancer or diabetes, genomic predisposition is subtle and, so far, not fully understood.

With the rapid decrease in the cost of DNA sequencing, the complete genome sequence of large cohorts of individuals can now be routinely obtained. This wealth of sequencing information is expected to ease the identification of genetic variations linked to complex traits. In this work, I investigate the analysis of genomic data in relation to complex diseases, which offers a number of important computational and statistical challenges. We tackle several steps necessary for the analysis of sequencing data and the identification of links to disease. Each step, which corresponds to a chapter in my thesis, is characterized by very different problems that need to be addressed.

- i) The first step is to analyze large amounts of information generated by sequencers to obtain a set of "genomic variants" present in each individual. To address these big data processing problems, Chapter 2 shows how we designed a programming language (BigDataScript), that simplifies the creation robust, scalable data pipelines.
- ii) Once genomic variants are obtained, we need to prioritize and filter them to discern which variants should be considered "important" and which ones are likely to be less relevant. We created the SnpEff & SnpSift [30],

- 31] packages that, using optimized algorithms, solve several annotation problems: a) standardize the annotation process, b) calculate putative genetic effects, c) estimate genetic impact, d) add several sources of genetic information, and e) facilitate variants filtering.
- iii) Finally, we address the problem of finding associations between interacting genetic loci and disease. One of the main problems in GWAS, known as “missing heritability”, is that most of the phenotypic variance attributed to genetic causes remains unexplained. Since interacting genetic loci have been pointed out as one of the possible causes of missing heritability, finding links between such interactions and disease has great significance in the field. We propose a methodology to increase the statistical power of this type of approaches by combining population-level genetic information with evolutionary information.

In a nutshell, this thesis addresses computational, analytical, algorithmic and methodological problems of transforming raw sequencing data into biological insight in the aetiology of complex disease. In the rest of this introduction we give the background required to understand the material shown in Chapters 2, 3 & 4 while providing motivations for our research.

1.1 Genomes and genetic variants

DNA is composed of four basic building blocks, called “bases” or “nucleotides”. These four nucleotides, usually abbreviated $\{A, C, G, T\}$, are Adenine, Cytosine, Guanine, and Thymine. Bases form pairs, either as $A - T$ or $C - G$, that pile-up forming two long polymers, with backbones that run in opposite directions giving rise to a double-helix structure. Arbitrarily, one of the polymers is called the positive strand and the other is called the negative strand.

Proteins are composed by chains of amino acids and, as explained by the central dogma of biology, DNA is the template that instructs cellular machinery how to produce proteins. There are 20 amino acids, which are the building blocks of all proteins. Each of the twenty amino acids is encoded by a group of three DNA bases called “codon”. More than one codon can code for the same amino acid (i.e. $4^3 = 64$ codons > 20 amino acids) allowing for code redundancy. Additionally, there are codons that mark the end of the protein, these are called “STOP” and signal molecular machinery to end the translation process.

Proteins are the most important part of the cell composing up to 50% of a cell’s dry weight compared to 3% of the DNA [5]. Proteins perform their functions mainly by interacting with other proteins, forming complex pathways that lead to a vast array of cellular functions including catalysis of chemical reactions, cell signaling, and structural conformation of the cell. The 3-dimensional structure of the protein, also called “tertiary structure”, is tailored to bind to other proteins in a specific manner to accomplish a functionality.

The human genome has a total of 3 Giga-base-pairs (Gb), and those bases are divided into 22 “autosomal” chromosome pairs (in each pair one chromosome is maternally inherited and the other paternally inherited) and one “sex” chromosome pair. The longest of the autosomal chromosomes is roughly 250 Mega-bases (Mb) and the shortest one is 50 Mb.

In order to compare DNA from different individuals (or samples), we need a “reference genome”. Having a standard reference sequence facilitates comparisons and analysis. For most well studied organisms, “reference genome”

sequences are available and current large scale sequencing projects are extending significantly the number of genomes known, e.g. one project seeks to sequence 10,000 mammalian genomes [71], another is targeting all microbes that live within humans guts [152]. The human reference genome (e.g. GRCh37) does not correspond to the DNA of any particular person, but to a “mosaic” of the genomes of thirteen anonymous volunteers from Buffalo, New York [140].

When the genome of an individual is sequenced, the DNA is compared to the “reference genome”. Most of the DNA is the same, but there are differences. These differences, generically known as “genomic variants” (or “variants”, for short), describe the particular genetic make-up of each individual. There are several different ways a sample can differ from a reference genome. Each variant is the result of a mutations that happened at some point in the evolutionary history of the individual (or that of the reference genome). Variant types can be roughly categorized in the following way:

Single nucleotide variants (SNV) or Single nucleotide polymorphism (SNP) are the simplest and more common variants produced by single base difference (e.g. a base in the reference genome, at a given coordinate, is an ‘A’, whereas the sample is ‘C’). Depending on whether the variant was identified in an individual or in a population, it is called a Single Nucleotide Variant (SNV) or Single Nucleotide Polymorphism (SNP). It is estimated that there are roughly $3.6M$ SNPs per individual [36]. There are several biological mechanisms responsible for this type of variants: i) replication errors, ii) errors introduced by DNA repair mechanism, iii) deamination (a base is changed by hydrolysis which may not be corrected by DNA repair mechanisms), iv) tautomerism (and alteration on the hydrogen bond that results in an incorrect pairing).

Multiple nucleotide polymorphism (MNP) are sequence differences affecting several consecutive nucleotides and are typically treated as a single variant locus if they are in perfect linkage disequilibrium (e.g. reference is ACG whereas the sample is TGC). .

Insertions (INS) refer to a sample having one or more extra base(s) compared to the reference genome (e.g. reference is AT and sample is ACT). Short insertions and deletions (indels) of a chromosome region range from 1 to 20 bases in length are approximately 30 times less frequent than SNV [36]. Small insertion are usually attributed to DNA polymerase slipping and replicating the same base/s (this produces a type of insertion known as duplication). Large insertions can be caused by unequal cross-over event (during meiosis) or transposable elements.

Deletions (DEL) are the opposite of insertions, the sample has some base(s) removed with respect to the reference genome (e.g. reference is ACT and sample is AT). As in the case of insertions, deletions can also be caused by ribosomal slippage, cross-over events during meiosis. Those include large deletions, which can result in the loss of an exon or one or more whole genes.

Copy number variations (CNVs) arise when the sample has two or more copies of the same genomic region (e.g. a whole gene that has been duplicated or triplicated) or conversely, when the sample has fewer copies than the reference genome. Copy number variations are often attributed to homologous recombination events.

Rearrangements such as inversions and translocations are events that involve two or more genomic breakpoints and a reorganization of genomic

segments, possibly resulting in gene fusions or loss of critical regulatory elements. Inversions, a type of rearrangement, result from a whole genomic region being inverted.

As humans have two copies of each autosome, variants could affect zero, one or two of the chromosomes and are called “homozygous reference”, “heterozygous”, and “homozygous alternative” respectively. Variants are also classified based on how common they are within the population: common ($\geq 5\%$), low frequency ($\leq 5\%$), or rare ($\leq 0.1\%$). How these types of genetic variants influence traits or disease risk is a topic of intense research that is discussed throughout this thesis.

1.2 DNA and disease risk

It would be fair to say that the Garrod family was fascinated by urine. As a physician at Kings College, Alfred Baring Garrod, discovered gout related abnormalities in uric acid [83]. His son, Sir Archibald Garrod, was interested in a condition known as alkaptonuria, in which children are mostly asymptomatic except for producing brown or black urine, but by the age of 30 individuals develop pain in joints of the spine, hips and knees. In 1902, Archibald observed that the family inheritance pattern of alkaptonuria resembled Mendel's recessive pattern and postulated that a mutation in a metabolic gene was responsible for the disease. Publishing his finding he gave birth to a new field of study known as “Human biochemical genetics” [83].

Diseases having simple inheritance patterns, such as alkaptonuria, cystic fibrosis, phenylketonuria and Huntington's are also known as Mendelian diseases [83]. The genetic components of several Mendelian diseases have been discovered since the mechanism was first elucidated by Garrod in 1902 and the process has been accelerated in recent years, thanks to the application of DNA sequencing techniques [13].

In complex diseases (or complex traits), such as diabetes or Alzheimers, affected individuals cannot be segregated within pedigrees (i.e. no simple pattern of inheritance can be identified). As opposed to Mendelian diseases the aetiological of complex traits is complicated due to factors such as: incomplete penetrance (symptoms are not always present in individuals who have the disease-causing mutation) and genetic heterogeneity (caused by any of a large number of alleles). This makes it difficult to pinpoint the genetic variants that increase risk of complex disease.

1.2.1 Heritability and Missing heritability

We all know that “tall parents tend to have tall children”, which is an informal way to say that height is a highly heritable trait. It is said that there are 30 cm from the tallest 5% to the shortest 5% of the population and genetics account for 80% to 90% of this variation [165], which means that 27cm of variance are assumed to be “carried” by DNA variants from parents to offspring. Since 2010 the GIANT consortia has been investigating the genetic component of complex traits like height, body mass index (BMI) and waist to hip ratio (WHR). Even though they found many variants associated those traits, their findings only explain 10% of the phenotypic variance which corresponds to only a few centimeters in height [165].

In order to measure heritability we need a formal definition. Heritability is defined as the proportion of phenotypic variance that is attributed to genetic variations. The total phenotypic variation is assumed to be caused by a combination of “environmental” and genetic variations $Var[P] = Var[G] + Var[E] + 2Cov[G, E]$.

The environmental variance $Var[E]$ is the phenotypic variance attributable only to environment, that is the variance for individuals having the same

genome $Var[E] = Var[P|G]$. This can be estimated by studying monozygotic and dizygotic twins.

If the covariance factor $Cov[G, E]$ is assumed to be zero, we can define heritability as $H^2 = \frac{Var[G]}{Var[P]}$. This is called “broad sense heritability” because $Var[G]$ takes into account all possible forms of genetic variance: $Var[G] = Var[G_A] + Var[G_D] + Var[G_I]$, where $Var[G_A]$ is the additive variance, $Var[G_D]$ is the variance from dominant alleles, and $Var[G_I]$ is the variance from interacting alleles (epistasis). Non-additive terms are difficult to estimate, so a simpler form of heritability called “narrow sense heritability” that only takes into account additive variance is defined as $h^2 = \frac{Var[G_A]}{Var[P]}$ [174].

Focusing on narrow sense heritability, the concept of “explained heritability” is defined as the part of heritability due to known variants with respect to phenotypic variation ($\pi_{explained} = h_{known}^2/h_{all}^2$). Similarly, missing heritability is defined as $\pi_{missing} = 1 - \pi_{explained} = 1 - h_{known}^2/h_{all}^2$. When all variants associated with traits are known, then $\pi_{missing} = 0$.

Until recently, it was widely assumed by the research community that the problem of missing heritability lied in finding the appropriate genetic variants to account for the numerator of the equation (h_{known}^2) [174]. However, in a series of theorems published recently, it has been proposed that there is a problem in the way the denominator is estimated [174]. The authors created a limiting pathway model ($LP(k)$) that accounts for epistasis (gene-gene interactions) in k biological pathways. They showed that a severe inflation of h_{all}^2 estimators occurs even for small values of k (e.g. $k \in [2, 10]$). As a result, genetic variants estimated to account only for 20% of heritability, could actually account for as much as 80% using an appropriate model [174].

Even though this result is encouraging, the problem is now shifted to detecting epistatic interactions, a problem that we discuss in section 1.5.8 and

Chapter 4. In the same work [174], the authors show an example of power calculation assuming relatively large genetic effect that would require sequencing roughly 5,000 individuals to detect links to genetic variants, which is a large but nowadays not uncommon, sample size. Nevertheless other estimates place the sample size requirements as high as 500,000 individuals [174]. Even though this sounds as an extremely large number of samples, it is quickly becoming possible thanks to large technological advances and cost reductions in sequencing and genotyping technologies.

1.2.2 Conclusions

Although some genetic causes of complex traits, such as type II diabetes, have been found, only a small portion of the phenotypic variance can be explained. This might indicate that many risk variants are yet to be discovered. Recent studies on the topic of missing heritability report that these “difficult to find genetic variants” might be in epistatic interaction (analyzed in section 1.5.9) or rare variants (see section 1.5.6). Analysis of either them requires more complex statistical models and larger sample sizes with the corresponding increase in computational requirements. In Chapter 4 of this thesis, we focus on methods for finding epistatic interactions related to complex disease and develop computationally tractable algorithms that can process data from sequencing experiments involving large number of samples in a reasonable amount of time.

1.3 Identification of genetic variants

Two of the main milestones in genetics were the discovery of the DNA structure in 1953 [160], followed by the first draft of the human genome in 2004 [34]. The cost of sequencing the first human reference genome was around \$3 billion (unadjusted US dollars) and it was an endeavor that took around 10 years. Since that time, sequencing technology has evolved substantially so that

a human genome can now be sequenced in three days for a price of less than \$1,000, according to prices estimated by Illumina, one of the main genome sequencer manufacturers.

The amount of information delivered by sequencing devices is growing faster than computer speed (Moore’s law) and data storage capacity. Just as a crude example, a leading edge sequencing system is advertized to be capable of delivering 18,000 human genomes at $30\times$ coverage per year, yielding over 3.2 PB of information. Having to process huge amounts of sequencing information poses several challenges, a problem informally known as “data deluge”. In this section, we explain how sequencing data is generated and how the huge amount of information delivered by a sequencer can be handled in order to make the problem tractable. We want to transform this raw data into knowledge of genomic variants that contribute to disease risk with the ultimate goal to translate these risk variants into biological knowledge. As expected, processing huge datasets consisting of thousands of sample is a complex problem. In Chapter 2 we show how mitigate or solve some of these issues, by designing a computer language specially tailored to tackle what are know as “Big data” problems.

1.3.1 Sequencing data

DNA sequencing machines (or sequencers) are based on different technologies. In a nutshell, a sequencer detects a set of polymers (or chains) of DNA nucleotides and outputs a set of strings of A, C, G, and Ts. Unfortunately, current technological limitations make it impossible to “read” a full chromosome as one long DNA sequence. Instead, modern sequencers produce a large number of “short reads”, which range from 100 bases to 20 Kilo-bases (Kb) in length, depending on the technology. Since sequencers are unable to read long DNA chains, preparing the DNA for sequencing involves fragmenting

it into small pieces. These DNA fragments are a random sub-samples of the original chromosomes. Reading each part of the genome several times allows to increase accuracy and ensure that the sequencer reads as much as possible of the original chromosomes. The coverage of a sequencing experiment is defined as the number of times each base of the genome is read on average. For instance, if the sequencing experiment is designed to produce one billion reads, and each read is 150 bases long, then the total number of bases read is 150Gb. Since the human genome is 3Gb, the coverage is said to be $50\times$.

After sequencing a sample, we have millions of reads but we do not know where these reads originate from in the genome. This is solved by aligning (also called mapping) reads to the reference genome, which is assumed to be very similar to the genome being sequenced. Once the reads are mapped, we can infer if the samples DNA has any differences with respect to the reference genome, a problem is known as “variant calling”.

Although sequencing costs are dropping fast, it is still expensive to sequence thousands of samples and in some cases it makes sense to focus on specific areas of the genome. A popular experimental setup is to focus on coding regions (exons). A technique called “exome sequencing” consists of capturing exons using a DNA chip and then sequencing the captured DNA fragments only. Exons are roughly 1.2% of the genome, thus this technique reduces sequencing costs significantly, for which it has been widely used by many research groups although it has the disadvantage of only analysing coding genomic variation.

1.3.2 Read mapping

Once the samples have been sequenced, we have a set of reads from the sequencer. The first step in the analysis is finding the location in the reference genome where each read is supposed to originate from, a process that is

complicated by a several factors: i) there are differences between the reference genome and the sample genome, ii) sequencing reads may contain errors, iii) several parts of the reference genome are quite similar making reads from those regions indistinguishable, and iv) a typical sequencing experiment generates millions of reads.

Local sequence alignment. We introduce a problem known as *local sequence alignment*: Given two sequences s_1 and s_2 from an alphabet (e.g. $\Sigma = \{A, C, G, T\}$), the alignment problem is to add gap characters ('-') to both sequences, so that a distance, such as Levenshtein distance, $d(s_1, s_2)$ is minimized. This problem has a well known solution, the Smith-Waterman algorithm [144], which is a variation of the global sequence alignment solution from Needleman-Wunsch [116], having an algorithm complexity $O(l_1 \cdot l_2)$ where l_1 and l_2 are the length of the sequences. So, Smith-Waterman algorithm is slow since in this case one of the sequences is the entire genome.

In order to speed up sequence alignments, several heuristic approaches emerged. Most notably, BLAST [6], which could be for mapping sequences to a reference genome. BLAST uses an index of the genome to map parts of the query sequence, called seeds, to the reference genome. Once these seeds have been positioned against the reference, BLAST joins the seeds performing an alignment only using a small part of the reference.

Read mapping. Sequence alignment has an exact algorithm solution and several faster heuristic solutions. But even the fastest solutions are too slow to be used with the millions of reads generated in a typical sequencing experiment. Faster algorithms can be used if we relax our requirements in two ways: i) we allow for sub-optimal results, and ii) instead of requiring the output to be a complete local alignment between a read and the genome, we just want to know the region in the reference genome where the read sequence is from.

This relaxed version of the alignment algorithm is called “read mapping” and the reduced complexity is enough to speed up the computations significantly. In a nutshell, a read mapping is regarded as correct if it overlaps the true reference genome region where the read originated. Once the mapping is performed, the read is locally aligned, a strategy similar to BLAST algorithm [97, 91].

Reformulating this as a *mapping* problem allows us to use data structures such as suffix trees to index the reference genome. Using suffix trees we can query for a substring (read) [53] of the indexed string in $O(m)$ time, where m is the length of the query. Alternatively, we can use suffix arrays which are a space optimized alternative to suffix trees [53]. An implicit assumption in this solution, is that the read will be very similar to the reference and that there will be no big gaps. Suffix arrays algorithms are fast but, even though they are memory optimized versions of suffix trees, memory requirements are still high ($O[n \log(n)]$), where n is the length of the indexed sequence -the genome-) and this becomes the limiting factor. In order to reduce memory footprint of suffix arrays, Ferragina and Manzini [56] created a data structure based on the Burrows-Wheeler transform. This structure, known as an FM-Index, is memory efficient yet fast enough to allow mapping high number of reads. An FM-index for the human genome can be built in only 1Gb of memory, compared to 12Gb required for an equivalent suffix array [97]. Given a genome G and a read R , an FM-index search can find the N_{occ} exact occurrences of R in G in $O(|R| + N_{occ})$ time, where $|R|$ is the length of R [97].

We should keep in mind that suffix trees, suffix arrays and FM-indexes are guaranteed to find all matching substring occurrences, nevertheless a sequencing read may not be an exact substring of the reference genome (due to sample’s genome differences with the reference genome, read errors, etc.). So,

even if efficient indexing and heuristic algorithms can decrease mapping time considerably, these algorithms are not guaranteed to find an optimal mapping. Several parameters, such as read length, sequencing error profile, and genome complexity profile can affect performance. The most commonly used implementation of the FM-index mapping algorithms are BWA [97, 98] and Bowtie [91, 90]. Each of them provide optimized versions for the two most common sequencing types: i) short reads with high accuracy [97, 91] or ii) longer reads with lower accuracy [98, 90].

Mapping quality. Sequencers not only provide sequence information, but also provide an error estimate for each base [96]. This is often referred as a quality (Q) value, which is the probability of an error, measured in negative decibels $Q = -10 \log_{10}(\epsilon)$, where ϵ is the error probability. Mapping quality is an estimation of the probability that a read is incorrectly mapped to the reference genome. Mapping algorithms provide estimates of mapping errors. In the MAQ model [99], which is one of the earliest models for calculating mapping quality, three main sources of error are explored: i) the probability that a read does not originate from the reference genome (e.g. sample contamination); ii) the probability that the true position is missed by the algorithm (e.g. mapping error); and iii) the probability that the mapping position is not the true one (e.g. if we have several possible mapping positions). It is assumed that the total error probability can be approximated as $\epsilon \approx \max(\epsilon_1, \epsilon_2, \epsilon_3)$.

1.3.3 Variant calling

Genome-wide variant calling has until recently largely been done using genotyping arrays (for SNVs) or Comparative Genomic Hybridization arrays (for CNVs). The inherent limitations of these technologies, particularly their

ability to only assay genotypes at sites that are known in advance to be polymorphic, combined with the declining cost of sequencing, have now made approaches based on high-throughput resequencing the tool of choice for variant calling in clinical studies.

Once the sequencing reads have been mapped to the reference genome, we can try to find the differences between a sequenced sample and the reference genome. This process is called “variant calling” [117]. Several factors complicate this task, the two main ones being sequencing errors and mapping errors, described in 1.3.2. Based on sequencing data and mapping error estimates, tools such as GATK [110] and SamTools/BcfTools [99] use maximum likelihood models can infer when there is a mismatch between a sample and the reference genome and whether the sample is homozygous or heterozygous for the variant. This method works best for differences of a single base (SNV), but it can also work with different degrees of success for short insertions or deletions (InDels) usually consisting of less than 10 bases.

Aligning sequences that contain InDels (gaps) is more difficult than ungapped alignments since finding optimal gap boundary depends on the scoring method being used. This biases variant calling algorithms towards detecting false SNVs near InDels [51]. An approach to reduce this problem is to look for candidate InDels and perform a local realignment in those regions. This local re-alignment process reduces significantly the number of false positive SNVs [51]. Another approach to reduce the number of false positive SNVs calls near InDels involves the “Base Alignment Quality” (BAQ) [95], which is the probability of misalignment for each base. It can be shown that replacing the original base quality with the minimum between base quality and BAQ produces an improvement in SNV calling accuracy. The BAQ can be calculated using a special type of “Hidden Markov Model” (HMM) designed for

sequence alignment [95, 53]. A more sophisticated option for reducing errors consist of performing a local genome re-assembly on each polymorphic region (e.g. HaplotypeCaller algorithm [149]).

Finally, one should note that the error probabilities inferred by the sequencers are far from perfect. Once the variants have been called, empirical error probabilities can be easily calculated [110] by comparing sequenced variants to a set of “gold standard variants” (i.e. variants that have been extensively validated). This allows to re-calibrate or re-estimate the error profile of the reads. This is known as a re-calibration step, and usually improves the number of false positives calls [51].

Due to the nature of short reads, this family of methods does not work for structural genomic variants, such as large insertions, deletions, copy number variations, inversions, or translocations. A different family of algorithms are used to identify structural variants generally making use of pair end reads or split reads, but their accuracy so far has been low compared to SNV calling algorithm [118].

Caveats

- i Using current technologies and computational methods for variant calling, detection accuracy varies significantly for different variant types. SNVs are by far the most accurately detected. Insertions and deletions, collectively referred as InDels, can be detected less efficiently depending on their sizes. Small InDels consisting of ten bases or less are easier to detect than large InDels consisting of 200 bases or more. The reason being that the most commonly used sequencers reads DNA in stretches roughly 200 bases long. Due to this technological limitations, detection is less reliable for more complex variant types.

1.4 Functional annotations of genomic variants

The development of cost-effective, high-throughput next generation sequencing (NGS) technologies is poised to have a profound impact on our ability to study the effects of individual genetic variants on the pathogenesis and progression of both monogenic and common polygenic diseases. As sequencing costs decrease and throughput increases, it has now become possible to quickly identify a large number of sequence polymorphisms (SNVs, indels, structural) using samples from affected and unaffected subjects and investigate these in epidemiologic studies to identify genomic regions where mutations increase disease risk. However, translating this information into biological or clinical insights is challenging as it is often difficult to determine which specific polymorphisms are the main pathogenetic drivers of disease across a population; and more importantly, how they affect the activity of disease-related molecular pathways in tissues and organism a specific patient. In part, this difficulty results from the large number of genetic variants that are observed in individual genomes (the human population is believed to contain approximately 3.5 million polymorphic sites with minor allele frequency above 5%) combined with the limited ability of computational approaches to distinguish variants with no impact on genome function (the vast majority) from variants affecting gene function or expression that may be associated with disease risk or drug response (the minority). The development of algorithms for automated variant annotation, which link each variant with information that may help predict its molecular and phenotypic impact, is a critical step towards prioritizing variants that may have a functional impact from those that are harmless or have irrelevant functional effects. In this section we review the key concepts and existing approaches in this important field. In Chapter 3 we introduce an approach to collect relevant information that will help answer questions about

genetic variants discovered in next-generation sequencing studies, including: (i) will a given coding variant affect the ability of a protein to carry its functions; (ii) will a given non-coding variant affect the expression or processing of a given gene; and ultimately (iii) will a given coding or non-coding variant have any impact on phenotypes of interest?

Answering these questions is essential for many types of analyses that use large-scale genomics datasets to study quantitative traits and diseases, particularly when only a small number of individuals is studied comprehensively at a genome-wide level. For example, most genome-wide association studies (GWAS) or exome sequencing studies lack the statistical power to identify rare variants or variants with small effects associated with a disease, in part due to the large number of variants assayed. This limitation can be addressed by directing both statistical analysis and subsequent experimental steps to focus on smaller sets of genetic variants that have been prioritized based on external evidence of their putative impact. The common impairment of DNA repair mechanisms and chromatin stability in malignant cells leads to a similar challenge in cancer genomics, where the hundreds or thousands of mutations that distinguish an individual's tumor and germline genomes need to be classified on the basis of their putative phenotypic effects and potential roles in carcinogenesis.

The large number of databases containing potentially helpful information about a given variant make the process of gathering and presenting relevant data challenging, despite excellent tools that already exist to analyze large genomics datasets (including GATK [110] and Galaxy [64]) and visualize the results (such as the UCSC [81] or Ensembl [58] genome browsers). Each of these databases uses its own format and is updated asynchronously, which makes it difficult for any analysis to remain up to date. In addition, the lack

of comprehensive and computationally efficient models that allow integrative analyses using these resources, makes the task of comprehensive variant annotation overwhelming. By efficiently combining information from tens or hundreds of genome-wide databases, the tools described here are designed to greatly facilitate the process of variant annotation, and make it accessible to groups with limited bioinformatics expertise or resources.

1.4.1 Variant types

Although variant calling is a challenging task and remains an important area of research, many high-quality tools exist for calling SNVs and indels. We discuss here the problem of annotating the variants identified by some of these tools. The most common type of variant identified by current technologies and analysis approaches is a single base difference with respect to the reference genome (SNV) followed by multiple base differences (MNP), as well as small insertions and deletions (InDels). Here, we focus on annotating those variants (or combinations of them, called "Mixed" variants), which comprise most of the variants in a typical sequencing experiment. We do not address the annotation of large rearrangements due to the challenges involved in their identification and functional characterization and their relative rarity in the germ line.

1.4.2 Types of genetic annotations

The process of genetic variant annotation consists of the collection, integration, and presentation of experimental and computational evidence that may shed light on the impact of each variant on gene or protein activity and ultimately on disease risk or other phenotypes. Variant annotation has traditionally been divided in two apparently independent but actually interrelated tasks based on the variant's location with respect to known protein-coding genes (see Table 1 for a list of commonly used variant annotations). Coding

variant annotation focuses on variants that are located within coding regions of annotated protein-coding genes and attempts to assess their impact on the function of the encoded protein. In contrast, non-coding variant annotation focuses on variants located outside the coding portion of genes (i.e. in intergenic regions, UTRs, introns, or non-protein-coding genes) and aims to assess their potential impact on transcriptional and post-transcriptional gene regulation. These two categories of variant annotations are not mutually exclusive, as variants located within exons can often have an impact on the gene transcripts processing (splicing). In addition, some transcripts can have both protein-coding and non-coding functions. Despite the intermingling of the notion of coding and non-coding variants, we will consider each type of annotation separately as assessing their impact requires different sources of data and algorithms.

The ultimate goal of variant annotation is to predict the impact of a sequence variant, although this is an ill-defined term. On the one hand, one may be interested in the molecular impact of a variant on the activity of a protein. On the other, others may be interested in a variant's impact on much higher-level phenotypes such as disease risk. Mutations that are predicted to completely abrogate a gene's activity are called loss-of-function (LOF) mutations. Those that are predicted to have less severe consequences are called moderate or low impact mutations. In practice, a variant will be predicted to cause LOF if it has two properties: (i) its molecular impact is reliably predictable by existing computational approaches (e.g. gain of stop-codon); and (ii) its functional impact, reflected by altered protein activity or expression levels, is expected to be large. Many types of variants, including most non-coding variants, may have a large functional impact but lack predictability, and as a consequence are typically not predicted to be LOF variants.

1.4.3 Coding variant annotation

Coding variants occur in translated exons. When a reliable gene annotation is available, their main impact can be classified by determining their effect on the translated amino acid sequence (if any). A synonymous coding variant (also called silent) does not change the sequence of amino acids encoded by the gene, although it may impact aspects of post-transcriptional regulation such as splicing and translation efficiency and can affect the total protein activity through changes in the amount of translated protein that is made in the cell. In contrast, a non-synonymous coding variant changes one or more amino acids encoded by the gene and can directly alter the protein's activity, localization or stability. Non-synonymous variants include missense substitutions that change a single amino acid, nonsense substitutions that lead to the gain of a stop codon, frame-preserving indels that insert or delete one or more amino acids, and frame-shifting indels that may completely alter the protein's amino acid sequence. Primary annotation and assessment of impact, determines whether a variant falls in any of these categories.

Caveats

- i *Gene misannotation.* Genomic variants that have a significant effect on a protein's expression or function represent a very small fraction of all variants. Assembly and gene annotation errors or genomic oddities that break classical computational models are also rare, but lead to false positives. This implies that one is likely to find a non-negligible fraction of false-positive high-impact variants among the list of what appear to be the strongest candidates for variants with severe effects. Tools such as SnpEff can anticipate some of the most common causes of misannotation, but the number and diversity of the type of events that can lead to false-positives makes the task very challenging. As a consequence, one should

always manually inspect the top candidates to ensure that they have been assigned to the correct genes and transcripts.

- ii *Gene isoforms.* In higher eukaryotes, most genes have more than one transcript (or isoform), due to alternative promoters, splicing, or polyadenylation sites. For example, a human gene has an average of 8.8 annotated messenger RNA (mRNA) isoforms and some genes are believed to have over 4,000 isoforms resulting from complex splicing programs. For these genes, a variant may be coding with respect to one mRNA isoform and non-coding with respect to another. There are two frequent approaches to address this situation: (i) annotate a variant using the most severe functional effect predicted for at least one mRNA isoform; or (ii) use only a single canonical transcript per gene to perform primary annotation.
- iii *Variant calling for indels.* Variant annotation relies on knowing the exact genomic coordinates of the variant: this is rarely a problem for isolated SNVs; however, insertions and deletions often cannot be located unambiguously. Consider for example the variant $AA \rightarrow A$. This mutation results in the loss of a single base, but was it the first or second A that was deleted? From the standpoint of the cell, this question is irrelevant and deletion of any A will have the same effect. In contrast, from the standpoint of most variant annotation software, deleting the first A is different from deleting the second. Consider the scenario of a previously annotated transcript where the first A is part of the 5' UTR and the second is the first base of a start codon. If the missing base is assigned to the leftmost position in the motif (as is the current convention), the deletion would be annotated as a low impact 5'UTR variant. However, assigning it to the rightmost A would make it appear (incorrectly) to be

a high-impact start-codon deletion. Similar issues may arise when considering conservation scores or transcription factor binding site (TFBS) predictions.

1.4.4 Loss of function variants

True LOF variants are difficult to predict computationally, but specific types of genetic changes will frequently lead to severely impaired protein activity. These include (i) stop-gains (nonsense mutations) and start-loss; (ii) indels causing frameshifts; (iii) large deletions that remove either the first exon or at least 50% of the protein coding sequence; and (iv) loss of splice acceptor or donor sites that alter the protein-coding sequence. Variants that introduce premature in-frame stop codons (nonsense mutations and most frameshift indels) are expected to abolish protein function, unless the variant is very near the C-terminus of the coding region [167] (effectively, downstream of the last functional domain in the protein). Such mutations may have severe consequences in affected cells, tissues or organism, as is seen for mutations that cause monogenic diseases [139]. In addition, a new stop codon that lies upstream of the last exon will likely trigger nonsense mediated decay (NMD), a process that degrades mRNA before protein synthesis occurs [114]. NMD predictions are not exact and many factors can affect mRNA degradation, including the variant's distance from the last exon-exon junction or poly-A tail, and the possibility that transcription may re-initiate downstream of the LOF variant [20].

A variant that leads to the loss of a stop codon, sometimes called read-through mutation, will result in an elongated protein-coding transcript that terminates at the next in-frame stop codon. While there are no general models that predict how deleterious this may be, such variants can also result in aberrant folding and degradation of the nascent proteins, leading to activation of

cellular stress response pathways, in addition to their direct effects on protein activity and expression levels [139].

The effect of the loss of a start codon depends on the location of a replacement start codon with respect to the translation start site and reading frame of the native protein. If the new start codon maintains the reading frame, the only consequence may be the loss of a few amino acids in the protein transcript; however, in many cases, the new start codon will not be in-frame, thus producing a frame-shifted protein that is later degraded. In addition, the new start codon may lack an appropriate regulatory context (for example, if there is no Kozak sequence nearby or if it disrupts 5' UTR folding) leading to reduced expression of an N-terminally truncated protein. Consequently, losing a start codon is thought to be highly deleterious in most cases, due to the potential that it may reduce both protein production and activity.

Caveats

- i *Rare amino acids.* Through a process called translational recoding, a UGA “Stop” codon located in the appropriate mRNA context (determined by both primary mRNA sequence and secondary structure) may be translated to incorporate a selenocysteine amino acid (Sec / U) [5]. In humans, it is known to occur 100 codons located in mRNAs whose 3' UTR contains a Selenocysteine insertion sequence element (SECIS). Since the translation machinery goes so far to encode these special rare amino acids, the expectation is that mutations at those sites would be highly deleterious. This is supported by evidence that reduced efficiency of selenocysteine incorporation is linked to severe clinical outcomes, such as early onset myopathy [105] and progressive cerebral atrophy [4].
- ii *False-positives in LOF predictions.* Variants predicted to result in a LOF sometimes actually produce proteins that are partially functional

[104]. In fact, an apparently healthy individual is typically heterozygous for around 100 predicted LOF variants, and homozygous for roughly 10, but many of those are unlikely to completely abolish the protein function. Indeed, these variants are enriched toward the 3' end of the gene, where they are likely to be less deleterious.

1.4.5 Variants with low or moderate impact

Compared to the high impact variants discussed above, where extensive prior biological evidence strongly suggests that a specific type of variant will severely impair protein activity, there are few guidelines that can reliably predict how the majority of nonsynonymous (missense) variants will alter protein function or expression. As a result, the primary annotation performed by SnpEff and most related software packages will broadly categorize missense substitutions and their accompanying amino acid changes (e.g. $K154 \rightarrow L154$) as moderate impact variants. Short indels whose length is a multiple of three are treated similarly, unless they introduce a stop codon, as their effect will usually be localized.

Once missense and frame-preserving indel variants are identified, a more detailed estimation of their impact on protein function can be performed using heuristic and statistical models. The most common approaches are based on conservation, either amongst orthologous or homologous proteins, or protein domains, sometimes adding information of the physio-chemical properties of the reference and variant amino acids (e.g. differences in side chain charge, hydrophobicity, or size). The SIFT algorithm [87] assesses the degree of selection against specific amino acid changes at a given position of a protein sequence by analyzing the substitution process at that site throughout a collection of predicted homologous proteins identified by PSI-BLAST [7]. Based on this multiple sequence alignment and the highly conserved regions it contains,

SIFT calculates a normalized probability of amino acid replacement (called the SIFT score), which estimates the mutation’s effect on protein function. Polyphen [3], another commonly used tool, takes the process one step further by searching UniProtKB/Swiss-Prot [39] and the DSSP database of secondary structure assignments [79] to determine if the variant is located in a known active site in the protein. In contrast to other methods that categorize each variant individually, VAAST [130], a commercially available package, computes scores for groups of variants located within a given gene and “collapses” them into a single category, a concept similar to burden testing performed for rare variants identified in exome sequencing studies. For human proteins, SnpEff makes use of the Database for Nonsynonymous SNVs’ Functional Predictions [102] (dbNSFP), which collects scores produced by several impact assessment algorithms in a single database. Individually, impact assessment methods usually have an estimated accuracy of 60% to 80% when compared to manually curated databases of human variants, but predictions from several algorithms can be combined to provide a stringent, but more accurate estimate of impact [28].

In most cases these algorithms apply best to SNVs since these are common in populations and there is more genomic sequence and experimental data available to refine the statistical methods. However, some recently developed algorithms are capable of assessing variants other than SNVs, including PROVEAN [28], which extends SIFT to assess the functional impact of indels.

Caveats

- i *Imprecise models of protein function.* Accurate impact assessment of coding variants remains an open problem and most computational predictions are riddled with both false positives and false negatives. While both missense variants and frame-preserving indels are broadly cataloged

as having moderate effects, this is mostly due to lack of a comprehensive model and the extremely complex computations that would be required for an in-depth analysis (such as protein structure predictions). In these cases, proteomic information can be revealing. SnpEff adds annotations from curated proteomic databases, such as NextProt [89], which can help to elucidate if a mutation alters a critical protein amino acid or domain (such as amino acids that are post-translationally modified as part of a signaling cascade or that are form the active site of an enzyme) resulting in a protein may no longer function.

- ii *Gain of deleterious function.* Computational variant annotation may eventually be able to fairly accurately predict the molecular impact of a variant in terms of the degree to which it translates in a loss of function for the encoded protein. However, gains of function, including the acquired ability to interact with new partners and disrupt their function, remain vastly more difficult to tackle, although several such variants have been linked to disease [163].
- iii *Unanticipated effects of synonymous variants.* In most cases, synonymous variants are regarded as non-deleterious (or low impact); however, one needs to seriously consider the possibility that they may have greater functional effects by altering mRNA splicing [43] or secondary structure [133]. Synonymous SNVs may also alter translation efficiency, by changing a frequently used to a rarely used codon and have been linked to changes in protein expression [136].

1.4.6 Non-coding variant annotation

Although coding variants represent less than 2% of variants in the human genome, they make up the vast majority of confirmed disease-related variants

that have been validated at a functional level. This may result from ascertainment bias (since variants in coding regions are straightforward to discover and characterize at a basic level and many studies have largely ignored non-coding variants); or may be explained by the increased complexity of computational approaches and lab assays required to predict and validate the impact of non-coding variants; or by their potentially more subtle impact on gene expression or cell function. Nonetheless, in a compendium of current GWAS studies, roughly 40% of the variants are intergenic and 30% intronic. Functional studies of these variants are increasingly emphasizing the importance of non-coding genetic variation at risk loci for complex genetic diseases and traits [73].

Functional non-coding regions of the genome encompass a wide variety of regulatory elements contained in DNA and RNA molecules that are involved in transcriptional and post-transcriptional regulation. Cis-regulatory elements include (i) binding sites for DNA-binding proteins such as transcription factors and chromatin remodelers; (ii) binding sites for RNA-binding proteins involved in splicing, mRNA localization, or translational regulation; (iii) micro RNA (miRNA) target sites; and (iv) long non-coding RNA (lncRNA) targets on DNA, RNA and proteins. Non-coding transcripts include well-characterized regulatory RNAs (e.g. miRNA, snoRNA, snRNA, piRNA and lncRNAs) as well as RNAs involved directly in protein synthesis (e.g. tRNA and rRNA). The annotation and impact assessment of non-coding variants presents a significant challenge for several reasons: (i) reliable technologies to study transcriptional regulatory regions on a genome-wide basis are only just reaching maturity and provide limited resolution of binding sites for individual transcription factors and regulatory RNA molecules; (ii) non-coding functional regions of most genomes remain incompletely mapped as they vary widely among different cell types and cell states (for example, in diseased and

healthy tissues); (iii) non-coding regulatory elements often are part of complex transcriptional programs that are time-dependent, contain many redundant linkages or reciprocal connections between genes and respond to a wide range of intra- and extracellular signals; and (iv) genomic regulatory elements rarely have a strict consensus sequence (for example, compare the position weight matrices used to identify transcription factor or miRNA binding sites with the amino acid triplet code) making the effect of a mutation on gene regulatory programs difficult to predict. As a result, high-quality annotation of non-coding variants relies more heavily on experimental data than is the case for coding variants: since many of these experimental techniques did not study the effects of SNVs on gene regulatory programs, they can only be used to annotate variants and not to predict their effects on gene transcription. In the few cases where the effects of SNVs have been studied (for example, the effects of SNVs that are common in a population and located in genetic loci associated with complex diseases), experimental approaches provide highly accurate functional assessment at a cost of reduced coverage compared to computational approaches.

Large-scale projects such as ENCODE [38] and modENCODE [24] have made major steps toward mapping gene transcription and transcriptional regulatory regions in many tissues and cell types, but similar studies in diseased tissues remain at an early stage (for example, the growing collection of disease-related epigenomes from the Epigenome Roadmap [16]). The base-by-base resolution and number of cell states studied for different types of regulatory elements and non-coding transcripts varies widely among datasets; in part due to the lack of sensitive, comprehensive and high-resolution technologies to study the different molecular species and modes of interaction that can be

altered by non-coding variants. Efficient technologies for genome-wide, high-throughput mapping of binding sites for RNA-binding proteins (PAR-CLIP [10]), miRNAs (PAR-CLIP [69] and CLASH [72]) are starting to be applied on a broad scale as are protocols to map transcription factor binding sites (TFBS) which can improve resolution to a single base (ChIP-exo [128]). However, in most cases, DNA and RNA binding sites are only imprecisely located within ChIP-Seq peaks that span genomic regions hundreds of base pairs in length, with computational approaches being used to pinpoint the bases most likely mediating the interaction. In the absence of more precise localization data, de novo computational prediction of binding sites for DNA and RNA binding proteins remains insufficiently accurate to be of much use in annotating single noncoding variants.

This limitation is particularly critical for functional predictions of putative target sites for microRNAs and other regulatory RNA species. MicroRNAs are short RNA molecules that regulate gene expression post-transcriptionally by binding the messenger RNA of a gene through complementary, usually in the 3' region of the transcript, which leads to mRNA degradation or inhibits translation. Sequence variants that cause the loss or gain of a miRNA target site would lead to dysregulation of the gene, with likely deleterious effects. Although miRNAs are relatively well documented in most model organisms including human, their binding sites are only starting to be mapped experimentally, and computational predictions have very low specificity. Meaningful information regarding the possible role of a variant in disrupting a miRNA target site is starting to emerge [100], although variants that create new miRNA binding sites remain under the radar.

Even if the position of a functional element could be perfectly determined, predicting a variant's impact on chromatin conformation, promoter activity,

gene expression, or transcript processing remains challenging. For transcription factors, this involves predicting whether the protein will still be able to recognize its mutated site (and with what affinity), as well as predicting the impact of these changes on gene expression levels. The latter is particularly hard to predict as a result of interactions, competition, and redundancy contained in regulatory networks of transcription factors or RNA binding proteins. As a consequence, computational prediction of the functional impact of non-coding variants remains a very active area of research and there is no broad consensus on the best methodology to use [159]. One significant exception is the identification of variants affecting canonical splice sites, defined as two bases on the 3' end on the intron (splice site acceptor) and 5' end of the intron (splice site donor). Variants that affect canonical splice sites are easily detected and typically lead to abnormal mRNA processing, involving exon loss or extension that leads to loss of function of the encoded protein.

1.4.7 Impact assessment of non-coding variants

Two broad classes of publicly available genome-wide datasets are commonly combined to assess the functional impact of non-coding genetic variants: (i) computational predictions of sequence conservation and sites involved in molecular interactions such as transcription factor and RBP binding, as well as miRNA-mRNA target interactions; and (ii) experimental genome-wide localization assays for DNA binding proteins, histone modifications, and chromatin accessibility.

Computational sources of evidence:. Interspecies sequence conservation plays a key role in scoring and prioritizing non-coding variants. This is based on the assumption is that sites or regions that have been more conserved across species than expected under a neutral model of evolution are likely to be functional; suggesting that mutations contained in them are likely

to be deleterious. In the absence of strong experimental data, sequence conservation measures calculated from whole genome multiple alignments, (for example using PhastCons [143], SciPhy [60], PhyloP [124] , and GERP [48]), have been developed to provide a generic indicator of function for non-coding variants. Although high conservation scores generally mean that a genomic region may be functional, the converse is not true and many experimentally proven functional noncoding regions show only modest sequence conservation (for example due to binding site turnover events). Finally, some regulatory regions (e.g. specific elements regulating immune response [126]) are under positive selection and may thus show less conservation than surrounding neutral regions.

In human, genome-wide computational predictions of transcription factor binding sites based on matching to publicly available position weight matrices are available from variety of sources, including Ensembl [58] and Jaspar [22]. Because of the low information content of most binding affinity profiles, the specificity of the predictions is generally very low. Related approaches exist to predict splicing regulatory regions [54] and miRNA target sites [173], some of which are precomputed for whole genomes and available from the UCSC or Ensembl genome browsers. Recent efforts to determine RNA-binding protein sequence affinities can also be used to identify putative binding sites for these proteins in mRNA [127].

Experimental sources of evidence: To investigate the potential impact of variants on transcriptional regulation, many published experimental data sets produced by large-scale projects such as ENCODE [38], modENCODE [24] and Roadmap Epigenomics [16], can be used directly by annotation packages. These include: (i) ChIP-seq or ChIP-exo experiments that identify

TFBSs on a genome-wide basis; (ii) DNaseI hypersensitivity or Formaldehyde-Assisted Isolation of Regulatory Elements (FAIRE) assays that identify regions with open chromatin; and (iii) ChIP-seq studies to identify the presence of specific promoter or enhancer-associated histone post-translational modifications, which can be combined to identify active, poised, and inactive enhancers and promoters [127]. Most of these data sets are easily available through Galaxy [64] (as tracks from the UCSC Genome Browser) or through SnpEff (as downloadable pre-computed datasets). In parallel with the types of studies described above, expression quantitative trait loci (eQTLs) represent an agnostic way to map putative regulatory regions. An increasing number of such loci are available through the GTEx database [103]. Experimental data that may support assessment of the impact of variants on post-transcriptional regulation remain sparser, although databases such as doRiNa [8] or starBase [168] contain genome-wide datasets obtained by CLIP-Seq and degradome sequencing. To our knowledge, these data have yet to be used in the context of variant annotation studies.

Combining sources of evidence:. Despite the variety of computational and experimental sources of evidence available, impact assessment for non-coding variants remains relatively crude, due to the fact that biological models of gene regulation remain fairly simple. Nonetheless, significant steps forward have been made recently, and two web-based tools, HaploReg [158] and RegulomeDb [18], perform SNV and indel impact assessment for variants from dbSNP on the basis of a broad body of computational and experimental evidence. Both use pre-computed scores for variants from dbSnp and therefore cannot be used for rare variants, but they are extremely valuable for exploration by associating the variant of interest with a variant in dbSnp via linkage disequilibrium.

Caveats

- i *Sparseness of functional sites within ChIP-seq peaks.* Even if a noncoding variant is located in a region that contains a ChIP-seq peak for a given TF and has all the hallmark signatures of regulatory chromatin, the likelihood that it is deleterious remains low, because most DNA bases contained within a peak are non-functional.
- ii *Gain of function mutations.* While this section has focused on variants causing the loss of a functional regulatory element, genetic variants may also create new or more effective transcription factor binding sites. These are substantially harder to detect as they can occur in regions that show no evidence of function in individuals possessing the reference allele, and show little conservation across species. Furthermore, computational methods to predict gain of affinity for a given TF caused by a variant have insufficient specificity to be of much use on their own.

1.4.8 Clinical effect of variants

One of the most revealing types of annotation of both coding and non-coding variants reports whether the variant has previously been implicated in a phenotype or disease. Although such information is available for only a small minority of all deleterious variants, their number is growing and should be the first type of annotation one seeks out. Clinical annotations, until recently, have been scattered in a large number of specialized databases of medical conditions with a genetic basis, including the comprehensive, manually curated collection of genetic loci, variants and phenotypes in the Online Mendelian Inheritance in Man database (OMIM, www.omim.org); web pages containing detailed clinical and genetic information about uncommon disorders in the Swedish National Board of Health and Welfare Database for Rare Diseases (www.socialstyrelsen.se/rarediseases) and the peer-reviewed

NIH GeneReviews collection [22] (www.ncbi.nlm.nih.gov/books/NBK1116); and a curated collection of over 140,000 mutations associated with common and rare genetic disorders in the commercial Human Gene Mutation Database (HGMD, www.hgmd.org/). In most cases, these datasets do not use standardized data collection or reporting formats; are designed to primarily provide information to patients and health professionals through a web interface; and rely on heterogeneous criteria to describe disease phenotypes and clinical outcomes; pathological and other clinical laboratory data; as well as the genetic and biologic experiments that have been used to demonstrate disease mechanisms at a molecular or cellular level. These shortcomings are being addressed by initiatives that provide centralized, evidence-based, comprehensive collections of known relationships between human genetic variants and their phenotype that are suitable for computational analysis, such as the NIH effort to aggregate records from OMIM, GeneReviews and locus-specific databases in ClinVar (www.ncbi.nlm.nih.gov/clinvar).

Another important application of variant detection and annotation is in the study of cancer genomes, which is occurring increasingly in clinical settings to support treatment decisions for advanced tumors. Annotation of variants detected in tumor sequences can be analyzed for clinical cohorts, using similar techniques as other complex traits, as well as for individual patients, using techniques to identify differences between somatic (tumor) and germline (healthy) tissues. In the latter case, one looks for cancer-associated mutations that distinguish the somatic genome of cancer cells of an individual from the germline genome in order to find the driving mutations that pinpoint the specific mechanisms underlying tumorigenesis or metastasis. Ideally, these mutations can be used to select a treatment for the patient, establish prognosis, or to identify causative mutations that have led to the cancer's progression.

In such a setting, given that sequence differences between the cancer and germline genomes are of greater interest than the background genetic changes between the germline and a reference genome, variant calling is performed using specialized algorithms, such as MuTect [29] and SomaticSniper [92].

Caveats

- i *Annotation accuracy.* Biological knowledge, as well as molecular and phenotypic evidence supports the identification of certain groups of high impact variants based on simple criteria (such as premature stops, frameshifts, start lost and rare amino acid mutations); however, it is often hard to predict whether non-synonymous variants will have equally large effects on an organism’s health. Even when the accepted “rules of thumb” used in the primary annotation indicate that protein function is impaired, we should consider that these predictions may be based on a small number of model genes and will require appropriate wet-lab validation or confirmatory studies in cohorts. In addition, as more human genomes are sequenced, it is likely that some genetic variants that have been linked to Mendelian diseases will be found in healthy individuals [129]; and in many cases, may not actually be disease-causing mutations [15].

1.4.9 Data structures and computational efficiency

Most computational pipelines for genomic variant annotation and primary impact assessment are relatively efficient and can annotate variants obtained from large resequencing projects involving thousands of samples within a few minutes or hours even using a moderately powered laptop. This is typically achieved through two key optimizations: (i) creation of reference annotation databases and (ii) implementation of efficient search algorithms. Reference database creation refers to the process of creating and storing precomputed genomic data from the reference genome, which can be searched quickly to

extract information relevant to each variant. This process needs to be performed only once per reference genome and most annotation tools have pre-computed databases for many organisms available for users to download. Since these databases are typically quite large, efficient search algorithms are used together with appropriate data structures to optimize the search process. In ANNOVAR [157], each chromosome is subdivided in a set of intervals of size k and genomic features for a given chromosome are stored in a hash table of size L/k , where L is the length of the chromosome. Another approach, used by SnpEff, is to use an “interval forest”, which is a hash of interval trees [42] indexed by chromosome. Querying an interval tree requires $O[\log(n) + m]$ time, where n is the number of features in the tree and m is the number of features in the result.

1.4.10 Conclusions

In Chapter 3 we introduce SnpEff & SnpSift, two approaches we designed for efficiently performing functional annotations of sequencing variants. These packages allow annotating, prioritizing, filtering and manipulating variant annotations as well as combining several public or custom-created databases. It should be noted SnpEff was one of the first annotation packages and has become one of the most widely used annotation software in both research and clinical environments.

1.5 Genome wide association studies

A genome wide association study aims at identifying genetic variants associated to a particular phenotype. First, the genomes (or exome, depending on the study design) of affected individuals (cases) and healthy individuals (controls) need to be sequenced, variants called, annotated and filtered. Then, the goal is to find variants that exhibit some statistical association with the trait or phenotype of interest, which could be a disease status (e.g. diabetic

vs healthy), a biomedical measurement (e.g. cholesterol level), or any measurable characteristic (e.g. height). Since the genome is so large, patterns of mutations that suggest correlation may be encountered by chance, so we need to establish statistical significance in order to distinguish true association from spurious ones. Like most studies, we will focus our discussion on SNVs, but most methods can be extended to other genomic variants.

1.5.1 Single variant tests and models

Consider a simple situation where there is only one variant in the whole genome for the cohort we are analysing. Since each individual has two sets of chromosomes, the variant can be present in one, both, or neither chromosomes, so the number of times a non-reference allele is present in an individual, is $N_{nr} = \{0, 1, 2\}$.

When the trait of interest is binary (e.g healthy vs disease), a cohort can be divided into cases and controls and we can build a 3 by 2 contingency table:

	<i>Homozygous Reference</i> $(N_{variant} = 0)$	<i>Heterozygous</i> $(N_{nr} = 1)$	<i>Homozygous non - reference</i> $(N_{nr} = 2)$
<i>Cases</i>	$N_{ca,ref}$	$N_{ca,het}$	$N_{ca,hom}$
<i>Controls</i>	$N_{co,ref}$	$N_{co,het}$	$N_{co,hom}$

Further assumptions about how many variants are required to increase disease risk can reduce this 3×2 table to a 2×2 table. In the “dominant model”, the effect of a mutated gene dominates over the healthy one, so one variant is enough to increase risk. The opposite, called “recessive model”, is when both chromosomes have to be mutated in order to increase risk [12, 33]. In these models, we can count how many cases and controls have at least one variant (dominant model) or two variants (recessive model). This simplifies

the previous table, yielding a 2×2 contingency table, than can be tested using either a χ^2 test or a Fisher exact test [12].

Two other commonly used models, are the “multiplicative” and the “additive” models [12, 33]. In these models, a disease risk is assumed to be multiplied (or increased) by a factor γ with every variant present. We cannot simplify the contingency table, so we assess significance using a Cochran-Armitage test [33].

1.5.2 Multiple variant tests

In a real case scenario there are thousands or millions of variants. We can extend the concept shown in the previous section by performing individual tests for each variant present in the cohort. Multiple testing can be addressed either by performing a correction, such as False Discovery Rate [12, 33], or using a stricter genome wide significance level. There are 3×10^9 bases in the genome, but taking into account the correlation between nearby variants (linkage disequilibrium), the genome wide significance level is generally accepted to be $p_{value} \leq 10^{-8}$.

In order to check if the null hypothesis of a significance tests is adequate, a QQ-plot is used (i.e. plotting the $y = -\log(p_{value})$ vs $x = -\log[\text{rank}(p_{value})/(N+1)]$, where N is the total number of variants). Adherence of the p-values to a 45 degree line on most of the range implies few systematic sources of association [12, 33]. If the p-values have a higher slope than the $y = x$ line, there might be “inflation”, possibly due to co-factors, such as population structure (see section 1.5.4). If the inflation is not too high (e.g. less than 5%), this bias can be corrected by shifting the p-values towards the 45 degree slope. More sophisticated methods are explained in section 1.5.4.

1.5.3 Continuous traits and correcting for co-factors

Methods described so far are suitable for binary “traits” or “phenotypes”. Statistical methods that link genetic information to traits can also be used on continuous or “quantitative” traits (e.g. weight, height, cholesterol level, etc.). A linear regression can be used assuming the traits are approximately normally distributed [12, 33]. A significance test (p_{value}) for linear models can be calculated using an F statistic, but more sophisticated methods are also available [12, 33].

Using linear models, it is easy to include known co-factors to correct for biases or inflation. For instance, if it is known that a risk increases with age or that males are more susceptible than females, age and sex can be added to the linear equation in order to correct for these effects [12, 33]. In a similar manner, we can add co-factors to binary traits using logistic regression.

1.5.4 Population structure

It is widely accepted that humans started in Africa and migrated to Europe, then to Asia and later to America [70]. Out of an initial population, a few individuals migrate and colonize a new territory. This implies that the genetic variety of the new colony is significantly reduced, compared to the previous population, since the genetic pool is only a small “founder population”. The “Out of Africa” hypothesis implies that each new migration produced a reduction in genetic variety, also known as a “population bottleneck” [70].

As we previously mentioned, each individual inherits two chromosome sets, a maternal and a paternal one. Through recombination a chromosome is formed by a crossover combining maternal and paternal chromosomes and then passed down, thus the offspring has two sets of chromosomes (one from each parent) that on average have half chromosome from each grandparent.

This breaking and shuffling of chromosomes every generation, increases genetic diversity. Nevertheless if variants are located nearby in the chromosome, the chances that they are broken apart by recombination event are smaller than if they are further away from each other. This produces a correlation of close variants or “linkage disequilibrium” (LD). Nearby highly correlated variants are said to be in the same “LD-block” [70]. If a population has low genetic variety, the LD-blocks are large. So African population has more variety (smaller LD-blocks) and conversely, European, Asian and Amerindian populations have less variety (larger LD-blocks) [70].

1.5.5 Population as confounding variable

Imagine that we have a cohort of individuals drawn from two populations (P_A and P_B) and that individuals in P_A have much higher risk of diabetes than individuals from P_B . Now imagine that individuals from P_A have a variant v_A more often, but v_A is actually neutral and has no health effects whatsoever. If we do not take into account population factors, our study would conclude that v_A is the cause of diabetes, just because we see v_A more often in affected individuals. In this case it is clear that population structure is a confounding variable. We could avoid this problem by analyzing each population separately [120], but this would cause a loss of statistical power since we have fewer samples.

A population that is a mixture of two or more population is known as an “admixed population”. For instance the “African-American” population is a mixture of, roughly, 80% African and 20% European genomes [70, 12]. This means that analyzing a cohort of African-American individuals, we would get population structure as a confounding variable because of population admixture [70]. Obviously, in this case we cannot analyze each population separately, because each individual in the sample is a mixture of two populations.

The admixed population problem can be studied by performing a correction using the eigen-structure of the sample covariance matrix [120]. Samples can be arranged as a matrix C where each row is a sample and each column represents a position in the genome where there is a variant. The numbers $C_{i,j}$ in the matrix indicate the number of non-reference alleles in a sample (row) at a genomic position (column j). Since the allele can be present in zero, one, or two chromosomes in each individual, the possible values for $C_{i,j}$ are $\{0, 1, 2\}$. The covariance matrix is calculated as $M = \hat{C}^T \cdot \hat{C}$, where \hat{C} is the matrix C corrected to have zero mean columns. Usually, the first two to ten principal components of M are used as factors in linear models (see section 1.5.3) to correct for population structure [120].

Whether a cohort has any population structure and needs correction or not can be tested using two methods: a) plotting the projections of the first two principal components and empirically observing the number of clusters in the chart, or b) using a statistic of the eigenvalues of M based on Tracy-Widom's distribution [120].

1.5.6 Common and Rare variants

The “allele frequency” (AF) is defined as the frequency a variant appears in a population. Variants are usually categorized according to AF into three groups: i) Common variants ($AF \geq 5\%$), “low frequency” ($1\% < AF < 5\%$), and iii) “rare variants” ($AF < 1\%$). Common variants originated earlier in the population while rare variants are either relatively recent or selected against.

There are three main models for disease susceptibility [70, 63]:i) the Common-Disease-Common-Variant hypothesis (CDCV) assumes that if disease is common, it must be caused by a common variant; ii) the “infinitesimal hypothesis” proposes that there are many common variants each having small

risk effects; and iii) the Common-Disease-Rare-Variant hypothesis proposes that there exists many rare variants, each one having large risk effects.

1.5.7 Rare variants test

The “rare variant model” assumes that multiple rare variants have large effects on a trait. The problem is that, since these variants are rare, huge sample sizes are required for tests to identify statistically significant associations. To overcome this problem, methods known as “burden tests” collapse groups of rare variants that are hypothesised to have similar effect and perform statistical significance tests on the group [94]. An example of collapsing technique is to count the number of rare variant in a given window and apply a Fisher exact test, as shown in section 1.5.1. A limitation of some burden tests is that they implicitly assume that all rare variants have the same direction of effect, although in reality they might have no effect, be deleterious, or protective [94, 166].

Several techniques allow weighting rare variants by collapsing them using a kernel matrix. This allows to incorporate other information, such as allele frequency and functional annotations. It can be shown that the statistic induced by kernel weighting functions follows a mixture of χ^2 distributions and there is an efficient way to approximate it [94, 166], avoiding computationally expensive permutations tests.

1.5.8 Epistasis

In this section we introduced the basic concepts and methodologies used in GWAS. Although fairly mature, there is still heavy research and continuous improvement on GWAS statistical methods. Not only it is well known that traditional (i.e. single marker) GWAS methods fail under non-additive models [45], but also variants so far discovered using these methods do not account for all the expected phenotypic variance attributed to genetic causes (i.e. missing

heritability). As other authors pointed out [41, 174, 175], this might be because we need to look for epistatic variants which are not taken into account using these methods. In the next section, and in Chapter 4, we cover the topic of epistatic GWAS analysis.

Genome wide association studies have traditionally focused on single variants or nearby groups of variants. An often cited reason for the lack of discovery of high impact risk factors in complex disease is that these models ignore loci interactions [41] which have recently been pointed out as a potential solution for the “missing heritability” problem [174, 175]. With interactions being so ubiquitous in cell function, one may wonder why they have been so neglected by GWAS. There are several reasons: i) models using interactions are much more complex [59] and by definition non-linear, ii) information on which proteins interacts with which other proteins is incomplete [154], iii) in the cases where there protein-protein interaction information is available, precise interacting sites are rarely known [154]. Taking into account the last two items, we need to explore all possible loci combinations, thus the number of N order interactions grows as $O(M^N)$ where M is the number of variants [49]. This requires exponentially more computational power than single loci models. This also severely reduces statistical power, which translates into requiring larger cohort, thus increasing sample collection and sequencing costs [49].

In Chapter 4 we develop a computationally tractable model for analysing putative interaction of pairs of variants from GWAS involving large case / control cohorts of complex disease. Our model is based on analysing cross-species multiple sequence alignments using a co-evolutionary model in order to obtain informative interaction prior probabilities that can be combined

to perform GWAS analysis of pairs of non-synonymous variants that may interact.

1.5.9 Epistatic GWAS

The definition of epistasis from a statistical perspective is a “departure from a linear model” [41]. This means that in a logistic regression model the input for sample s includes terms with each of the genotypes at loci i and j), as well as an “interaction term” $g_{s,i} \cdot g_{s,j}$ [40].

$$\begin{aligned} P(d_s | g_{s,i}, g_{s,j}) &= \phi[\theta_0 + \theta_1 g_{s,i} + \theta_2 g_{s,j} + \theta_3(g_{s,i}g_{s,j}) \\ &\quad \dots + \theta_4 c_{s,1} + \dots + \theta_m c_{s,N_{cov}}] \end{aligned}$$

where d_s is disease status, $\phi(\cdot)$ is the sigmoid function, $c_{s,1}, c_{s,2}, \dots$ are covariates for sample s .

Models involving interactions between more than two variants can be defined similarly, but require more parameters and extremely large samples are required to accurately fit them.

Several families of approaches for epistatic GWAS exist. Here we mention a few:

- Allele frequency: In [2], an analysis of imbalanced allele pair frequencies is performed under the assumption that an implicit test for fitness can be achieved looking for over/under-represented allele pairs in a given population. In another study [172] the authors infer that interactions can create LD in disease population under two-loci model, then they show how LD-based p-values can uncover interaction and sometimes (in their simulations) outperform logistic regression tests.

- Bayesian model: In [171], a “Bayesian partitioning model” is used by providing Dirichlet prior distributions for each partition and computing posterior probabilities using Markov chain Monte Carlo (MCMC) algorithms. The methodology first test individual makers and picks only the top 10% to further investigate for epistasis, because it is prohibitive to test all loci.
- Machine learning: From a machine learning point of view, finding interacting variants is simply an “*optimisation procedure is to find a set of parameters that allows the machine-learning model to most accurately predict class membership (e.g. affected vs unaffected)*” [111]. Several approaches have emerged to tackle the “interaction problem” and used a variety of different techniques [85, 111] , such as neural networks, cellular automata, random forests, multifactor dimensionality reduction, support vector machines, etc.

Although all these models have advantages under some assumptions, none of them seems to be a “clear winner” over the rest [41]. All of these models suffer from the increase in number of tests that need to be performed, which raises two issues: i) multiple testing, which is often resolved by stringent significance threshold, and ii) computational feasibility, which is solved by efficient algorithms, parallelization, and heuristic approaches to quickly discard uninformative loci combinations. So far, no method for epistatic GWAS has been widely adopted and there is need of different approaches to be explored. In Chapter 4 we propose an approach to combine co-evolutionary models and GWAS epistasis of pairs of putatively interacting loci.

1.6 Thesis roadmap and Contributions

The original research presented in this thesis covers topics related to the computational and statistical methodologies related to the analysis of sequencing variants to unveil genetic links to complex disease. Broadly speaking, we address three types of problems: (i) Data processing of large datasets from high throughput biological experiments such as resequencing in the context of a GWAS (Chapter 2); (ii) functional annotations, i.e. calculating variant's impact at the molecular, cellular or even clinical level (Chapter 3); (iii) identification of genetic risk factors for complex disease using models that combine population-level and evolutionary-level data to detect putative epistatic interactions (Chapter 4). When applicable, background material specific to each chapter is presented in a preface, together with an explanation of how that chapter ties in with the rest of the thesis.

This thesis comprises text and figures of articles that have either been published, submitted for publication, or ready to be submitted (waiting upon data embargo restrictions):

Chapter 2

1. P. Cingolani, R. Sladek, and M. Blanchette. “BigDataScript: a scripting language for data pipelines.” *Bioinformatics* 31.1 (2015): 10-16.

For this paper, PC conceptualized the idea and performed the language design and implementation. RS & MB helped in designing robustness testing procedures. PC, RS & MB wrote the manuscript.

Chapter 3

2. **P. Cingolani**, A. Platts, M. Coon, T. Nguyen, L. Wang, S.J. Land, X. Lu, D.M. Ruden, et al. “A program for annotating and predicting the effects of single nucleotide polymorphisms, snpeff: Snps in the genome of drosophila melanogaster strain $w^{1118}; iso-2; iso-3$ ”. Fly, 6(2), 2012.

For this paper, PC conceptualized the idea, implemented the program and performed testing. AP contributed several feature ideas, software testing and suggested improvements. XL, DR, SL, LW, TN, MC, LW performed mutagenesis and sequencing experiments. XL and DR performed the biological interpretation of the data. All authors contributed to the manuscript.

SnpEff’s accompanying publication (SnpSift):

3. **P. Cingolani**, V. M. Patel, M. Coon, T. Nguyen, S. Land, D. M. Ruden, and X. Lu. “Using drosophila melanogaster as a model for genotoxic chemical mutational studies with a new program, snpsift”. Toxicogenomics in non-mammalian species, page 92, 2012.

We used SnpEff & SnpSift and developed a number of new functionalities in the context of two collaborative GWAS projects on type II diabetes:

4. M. McCarthy, T2D Genes Consortia. “Variation in protein-coding sequence and predisposition to type 2 diabetes”, Ready for submission.

5. A. Mahajan, X. Sim, H. Ng, A. Manning, M. Rivas, H. Heather, A. Locke, N. Grarup, H. K. Im, **P. Cingolani**, et. al. “Identification and Functional Characterization of G6PC2 Coding Variants Influencing Glycemic Traits Define an Effector Transcript at the G6PC2-ABCB11 Locus.” PLoS genetics 11.1 (2015): e1004876-e1004876.

Chapter 4

6. **P. Cingolani**, R. Sladek, and M. Blanchette. “A co-evolutionary approach for detecting epistatic interactions in genome-wide association studies”. Ready for submission (data embargo restrictions).

For this paper, PC designed the methodology under the supervision of MB and RS. PC implemented the algorithms. PC, RS & MB wrote the manuscript. This work uses data from the T2D consortia, thus it cannot be published until the main T2D paper is accepted for publication (according to T2D data embargo).

Other contributions

During my thesis I have co-authored several other scientific articles (grouped by topic) published, submitted for publication, or ready to be submitted, not mentioned in this thesis:

Epigenetics

7. **P. Cingolani**, X. Cao, R. Khetani, C.C. Chen, M. Coon, A. Bollig-Fischer, S. Land, Y. Huang, M. Hudson, M. Garfinkel, and others. “Intronic Non-CG DNA hydroxymethylation and alternative mRNA splicing in honey bees.” *BMC genomics* 14.1 (2013): 666.
8. M. Senut, A. Sen, **P. Cingolani**, A. Shaik, S. Land, Susan J and D. M. Ruden. “Lead exposure disrupts global DNA methylation in human embryonic stem cells and alters their neuronal differentiation.” *Toxicological Sciences* (2014).
9. D. M. Ruden, **P. Cingolani**, A. Sen, W. Qu, L. Wang, M. Senut, M. Garfinkel, V. Dollars, X. Lu, “Epigenetics as an answer to Darwin’s ‘special difficulty’ Part 2: Natural selection of metastable epialleles in honeybee castes”, *Frontiers in Genetics* (2015).
10. M. Senut, A. Sen, **P. Cingolani**, A. Shaik, S. Land, Susan J and D. M. Ruden. “Lead exposure induces changes in 5-hydroxymethylcytosine clusters in CpG islands in human embryonic stem cells and umbilical cord blood”, Submitted to ‘Epigenomics’.
11. M. Senut, **P. Cingolani**, A. Sen, Arko, A. Kruger, A. Shaik, H. Hirsch, S. Suhr, D. Ruden. “Epigenetics of early-life lead exposure and effects on brain development.” *Epigenomics* 4.6 (2012): 665-674.

GWAS & Disease

12. K. Oualkacha, Z. Dastani, R. Li, **P. Cingolani**, T. Spector, C. Hammond, J. Richards, A. Ciampi, C. Greenwood. “Adjusted sequence kernel association test for rare variants controlling for cryptic and family relatedness.” *Genetic epidemiology* 37.4 (2013): 366-376.
13. S. Bongfen, I. Rodrigue-Gervais, J. Berghout, S. Torre, **P. Cingolani**, S. Wiltshire, G. Leiva-Torres, L. Letourneau, R. Sladek, M. Blanchette, and others. “An N-ethyl-N-nitrosourea (ENU)-induced dominant negative mutation in the JAK3 kinase protects against cerebral malaria.” *PloS one* 7.2 (2012): e31012.
14. C. Meunier, L. Van Der Kraak, C. Turbide, N. Groulx, I. Labouba, Ingrid, **P. Cingolani**, M. Blanchette, G. Yeretssian, A. Mes-Masson, M. Saleh, and others. “Positional mapping and candidate gene analysis of the mouse Ccs3 locus that regulates differential susceptibility to carcinogen-induced colorectal cancer.” *PloS one* 8.3 (2013): e58733.
15. G. Caillard, G. Leiva-Torres, M. Leney-Greene, B. Charbonneau, A. Dumaine, N. Fodil-Cornu, M. Pyzik, **P. Cingolani**, J. Schwartzentruber, J. Dupaul-Chicoine, and others. “Genome-wide mouse mutagenesis reveals CD45-mediated T cell function as critical in protective immunity to HSV-1.” *PLoS pathogens* 9.9 (2013): e1003637.
16. M. Bouttier, D. Laperriere, M. Babak Memari, M. Verway, E. Mitchell, **P. Cingolani**, T. Wang, M. Behr, R. Sladek, M. Blanchette, S. Mader and J. White. “Genomics analysis reveals elevated LXR signaling reduces *M. tuberculosis* viability”, Submitted to *Journal of Clinical Investigation*.

17. M. Bouttier, D. Laperriere, M. Babak Memari, M. Verway, E. Mitchell, **P. Cingolani**, T. Wang, M. Behr, R. Sladek, M. Blanchette, S. Mader and J. White. “Genomic analysis of enhancers engaged in *M. tuberculosis*-infected macrophages reveals that LXR signaling reduces mycobacterial burden”, Submitted to PLOS Pathogens.

Fuzzy logic

18. **P. Cingolani** and Jesus Alcala-Fdez. “jFuzzyLogic: a robust and flexible Fuzzy-Logic inference system language implementation.” FUZZ-IEEE. 2012.
19. **P. Cingolani** and Jesus Alcala-Fdez. “jFuzzyLogic: a java library to design fuzzy logic controllers according to the standard for fuzzy control programming.” International Journal of Computational Intelligence Systems (2013), vol 6, pages 65-75.

CHAPTER 2

BigDataScript: A scripting language for data pipelines

2.1 Preface

The overall goal in this thesis is to find genetic loci related to complex disease. In order to have enough statistical power to find these risk loci, we need to sequence thousands of cases and controls (i.e. patients and healthy individuals). Obviously the first step is to find all these patients, obtain patients consent, take samples and keep track of clinically relevant variables (such as age, sex, BMI, and glycemic traits). Just by the sheer number of patients involved, its easy to see that the logistics are challenging, to say the least.

Once the sequencing of each patients DNA is performed, we need to process the raw sequencing information by performing what is known as “primary sequencing analysis”, which involves mapping reads to the reference genome, calling variants, as well as performing several types of quality controls. The term “primary analysis” makes it sound as if this step is simple, but it is not. Managing such volume of information is a huge task that requires large computational resources, and coordinating the process involved at every stage of the analysis is not trivial, even if the jobs are relatively easy to parallelize.

As an example of the complexity and data volumes involved in these analysis pipelines, mapping the raw reads to the reference genome (i.e. the first stage of the primary analysis) for our T2D sequencing data is estimated to take over 12,000 CPU hours, that is over 32 CPU/years, under the most optimistic assumptions. At this magnitude hardware and failures become a significant issue since the probability of one or more nodes malfunction while the data is being processed is quite high.

We designed and implemented a simple script-like programming language called BigDataScript (BDS), with a clean and minimalist syntax to develop and manage pipeline execution and provide robustness to various types of software and hardware failures as well as portability. This programming language specifically tailored for data processing pipelines, improves abstraction from hardware resources and assists with robustness. Hardware abstraction allows BDS pipelines to run without modification on a wide range of computer architectures, from a small laptop to multi-core servers, server farms, clusters, clouds or even whole datacenters. BDS achieves robustness by incorporating the concepts of absolute serialization and lazy processing, thus allowing pipelines to recover from errors. By abstracting pipeline concepts at programming language level, BDS simplifies implementation, execution and management of complex bioinformatics pipelines, resulting in reduced development and debugging cycles as well as cleaner code. BDS was used to create data analysis pipelines required for our research, including the ones described throughout this thesis, and is currently used by other research groups and sequencing facilities in both academic and private environments.

The rest of the chapter is published in: Cingolani, Pablo, Rob Sladek, and Mathieu Blanchette. “BigDataScript: a scripting language for data pipelines.” *Bioinformatics* 31.1 (2015): 10-16.

2.2 Introduction

Processing large amounts of data is becoming increasingly important and common in research environments as a consequence of technology improvements and reduced costs of high-throughput experiments. This is particularly the case for genomics research programs, where massive parallelization of microarray and sequencing-based assays can support complex genome-wide experiments involving tens or hundreds of thousands of patient samples [175].

With the democratization of high-throughput approaches and simplified access to processing resources (e.g. cloud computing), researchers must now routinely analyze large datasets. This paradigm shift with respect to the access and manipulation of information creates new challenges by requiring highly specialized skill, such as implementing data-processing pipelines, to be accessible to a much wider audience.

A data-processing pipeline, referred as “pipeline” for short, is a set of partially ordered computing tasks coordinated to process large amounts of data. Each of these tasks is designed to solve specific parts of a larger problem, and their coordinated outcomes are required to solve the problem as a whole. Many of the software tools used in pipelines that solve big data genomics problems are CPU, memory or I/O intensive and commonly run for several hours or even days. Creating and executing such pipelines require running and coordinating several of these tools to ensure proper data flow and error control from one analysis step to the next. For instance, a processing pipeline for a sequencing-based genome-wide association study may involve the following steps [11]: (i) mapping DNA sequence reads obtained from thousands of patients to a reference genome; (ii) identifying genetic changes present in each patient genome (known as “calling” variants); (iii) annotating these variants with respect to known gene transcripts or other genome landmarks; (iv) applying statistical analyses to identify genetic variants that are associated with differences in the patient phenotypes; and (v) quality control on each of the previous steps. Even though efficient tools exist to perform each of these steps, coordinating these processes in a scalable, robust and flexible pipeline is challenging because creating pipelines using general-purpose computer languages (e.g. Java, Python or Shell scripting) involves handling many low-level process synchronization and scheduling details. As a result, process coordination usually depends on

specific features of the underlying systems architecture, making pipelines difficult to migrate. For example, a processing pipeline designed for a “multi-core server” cannot directly be used on a cluster because running tasks on a cluster requires queuing them using cluster-specific commands (e.g. `qsub`). Therefore, if using such a language, programmers and researchers must spend significant efforts to deal with architecture-specific details that are not germane to the problem of interest, and pipelines have to be reprogrammed or adapted to run on other computer architectures. This is aggravated by the fact that the requirements change often and the software tools are constantly evolving.

In the context of bioinformatics, there are several frameworks to help implement data-processing pipelines; although a full comparison is beyond the scope of this article, we mention a few that relate to our work: (i) Snakemake (Koster and Rahmann, 2012) written as a Python domain-specific language (DSL), which has a strong influence from `make` command. Just as in `make`, the workflow is specified by rules, and dependencies are implied between one rules input files and another rules output files. (ii) Ruffus (Goodstadt, 2010), a Python library, uses a syntactic mechanism based on decorations. This approach tends to spread the pipeline structure throughout the code, making maintenance cumbersome [134]. (iii) Leaf [115], which is also written as a Python library, expresses pipelines as graphs drawn using ASCII characters. Although visually rich, the authors acknowledge that this representation is harder to maintain than the traditional code. (iv) Bpipe [134] is implemented as a DSL on top of Groovy, a Java Virtual Machine (JVM)-based language. Bpipe facilitates reordering, removing or adding pipeline stages, and thus, it is easy for running many variations of a pipeline. (v) NextFlow (www.nextflow.io), another Groovy-based DSL, is based on data flow programming paradigm. This paradigm simplifies parallelism and lets the programmer

focus on the coordination and synchronization of the processes by simply specifying their inputs and outputs.

Each of these systems creates either a framework or a DSL on a pre-existing general-purpose programming language. This has the obvious benefit of leveraging the languages power, expressiveness and speed, but it also means that the programmer may have to learn the new general-purpose programming language, which can be taxing and take time to master. Some of these pipeline tools use new syntactic structures or concepts (e.g. NextFlows data-flow programming model or Leafs pipeline drawings) that can be powerful, but require programming outside the traditional imperative model, and thus might create a steep learning curve.

In this article, we introduce a new pipeline programming language called BigDataScript (BDS), which is a scripting language designed for working with big data pipelines in system architectures of different sizes and capabilities. In contrast to existing frameworks, which extend general-purpose languages through libraries or DSLs, our approach helps to solve the typical challenges in pipeline programming by creating a simple yet powerful and flexible programming language. BDS tackles common problems in pipeline programming by transparently managing infrastructure and resources without requiring explicit code from the programmer, although allowing the programmer to remain in tight control of resources. It can be used to create robust pipelines by introducing mechanisms of lazy processing and absolute serialization, a concept similar to continuations (Reynolds, 1993) that helps to recover from several types of failures, thus improving robustness. BDS runs on any Unix-like environment (we currently provide Linux and OS.X pre-compiled binaries) and can be ported to other operating systems where a Java runtime and a GO compiler are available.

Unlike other efforts, BDS consists of a dedicated grammar with its own parser and interpreter, rather than being implemented on top of an existing language. Our language is similar to commonly used syntax and avoids inventing new syntactic structures or concepts. This results in a quick-to-learn, clean and minimalistic language. Furthermore, creating our own interpreter gives better control of pipeline execution and allows us to create features unavailable in general-purpose language (most notably, absolute serialization). This comes at the expense of expressiveness and speed. BDS is not as powerful as Java or Python, and our simple interpreter cannot be compared with sophisticated just-in-time execution or JVM-optimized byte-code execution provided by other languages. Nonetheless, in our experience, most bioinformatics pipelines rely on simple programmatic constructs. Furthermore, in typical pipelines, the vast majority of the running time is spent executing external programs, making the executing time of the pipeline code itself a negligible factor. For these reasons, we argue that BDS offers a good trade-off between simplicity and expressiveness or speed.

2.3 Methods

In our experience, using general-purpose programming languages to develop pipelines is notably slow owing to many architecture-specific details the programmer has to deal with. Using an architecture agnostic language means that the pipeline can be developed and debugged on a regular desktop or laptop using a small sample dataset and deployed to a cluster to process large datasets without any code changes. This significantly reduces the time and effort required for development cycles. As BDS is intended to solve or simplify the main challenges in implementing, testing and programming data processing pipelines without introducing a steep learning curve, our main design goals

are (i) simple programming language; (ii) abstraction from systems architecture; and (iii) robustness to hardware and software failure during computationally intensive data analysis tasks. In the next sections, we explore how these concepts are implemented in BDS.

2.3.1 Language overview

BDS is a scripting language whose syntax is similar to well-known imperative languages. BDS supports basic programming constructs (`if/ else, for, while, etc.`) and modularity constructs such as functions and `include` statements, which are complemented with architecture-independent mechanisms for basic pipeline runtime control (such as `task`, `sys`, `wait` and `checkpoint`). At runtime, the BDS backend engine translates these high-level commands into the appropriate architecture-dependent instructions. At the moment, BDS does not support object-oriented programming, which is indeed supported by other pipeline tools based on libraries/DSL extending general-purpose programming languages. The complete language specification and documentation is available online at <http://pcingola.github.io/BigDataScript>.

Unlike most scripting languages, BDS is strongly typed, allowing detection of common type conversion errors at the initial parsing stage (pseudo-compilation) rather than at runtime (which can happen after several hours of execution). As the syntax of strict typing languages tends to be more verbose owing to longer variable declaration statements, we provide a type inference mechanism (operator `:=`) that improves code readability. For example (Listing 1), the variables `in` and `out` are automatically assigned the types the first time they are used (in this case, the type is assigned to be `string`).

2.3.2 Abstraction from resources

One of the key features of BDS is that it provides abstraction from most architecture-specific details. In the same way that high-level programming

languages such as C or Java allow abstraction of the CPU type and other hardware features, BDS supports system-level abstraction, including the number and the type of computing-nodes or CPU-cores that are available to the pipeline and its component tasks, whether firing another process may saturate the servers memory or whether a process is executed immediately or queued.

Pipeline programming requires effective task management, particularly the ability to launch processes and wait for processes to finish execution before starting others. Task management can be performed using a single BDS statement, independently of whether this is running on a local computer or a cluster. Processes are executed using the task statement, which accepts an optional list of resources required by the task (for example, see Listing 1). The task consists of running a fictitious system command myProcess and diverting the output to `output.file`. BDS currently supports the following architectures: (i) local, single or multi-core computer; (ii) cluster, using GridEngine, Torque and Moab; (iii) server farm, using ssh access; and (iv) cloud, using EC2 and StarCluster. Depending on the type of architecture on which the script is run, the task will be executed by calling the appropriate queuing command (for a cluster) or by launching it directly (for a multi-core server).

Listing 2.1: `pipeline.bds` program. A simple pipeline example featuring and a maximum of 6 h of execution time (Line 5).

```
1 #!/usr/bin/env bds
2 in := "input.file"
3 out := "output.file"
4 task ( out <- in, cpus=2, timeout=6*HOUR ) {
5     sys myProcess $in > $out # Invoke command
6 }
```

BDS performs process monitoring or cluster queue monitoring to make sure all tasks end with a successful exit status and within required time limits.

This is implemented using the `wait` command, which acts as a barrier to ensure that no statement is executed until all tasks finished successfully. Listing 2 shows a two-step pipeline with task dependencies using a `wait` statement (Line 13). If one or more of the `task` executions fail, BDS will wait until all remaining tasks finish and stop script execution at the `wait` statement. An implicit `wait` statement is added at the end of the main execution thread, which means that a BDS script does not finish execution until all tasks have finished running. It is common for pipelines to need multiple levels of parallel execution; this can be achieved using the `parallel` statement (or `par` for short). Wait statements accept a list of task IDs/parallel IDs in the current execution thread.

In addition to supporting explicitly defined task dependencies, BDS also automatically models implicit dependencies using a directed acyclic graph (DAG) that is inferred from information provided in the dependency operators (`<`) contained in `task` statements (see Listing 2, line 8). Finally, the `dep` expression defines a task whose conditions are not evaluated immediately (as it happens in `task` expressions) but only executed if required to satisfy a `goal`. Using `dep` and `goal` makes it easier to define pipelines in a **declarative** manner that is similar to other pipeline tools, as tasks are executed only if the output needs to be updated with respect to the inputs, independent of the intermediate results file, which might have been deleted.

2.3.3 Robustness

BDS provides two different mechanisms that help create robust pipelines: lazy processing and absolute serialization. When a processing pipeline fails, BDS automatically cleans up all stale output files to ensure that rerunning the pipeline will produce a correct output. If a BDS program is interrupted, typically by pressing Ctrl-C on the console, all scheduled tasks and running jobs are terminated or deallocated from the cluster. In addition to immediately

releasing computing resources, a clean stop means that users do not have to manually dequeue tasks, which allows them to focus on the problem at hand without having to worry about restoring a clean state.

Lazy processing. Complex processing pipelines are bound to fail owing to unexpected reasons that range from data format problems to hardware failures. Rerunning a pipeline from scratch means wasting days on recalculating results that have already been processed. One common approach, when using general-purpose scripting languages, is to edit the script and comment out some steps to save processing time, which is inelegant and error prone. A better approach is to develop pipelines that incorporate the concept of lazy processing [115], a concept popularized by the `make` command (Feldman, 1979) used to compile programs, and which simply means the work is not done a task invoking a fictitious command `myProcess` defined to require 2 CPUs twice. This concept is at the core of many of the pipeline programming tools, such as SnakeMake, Ruffus, Leaf and Bpipe. By design, when lazy processing pipelines are rerun using the same dataset, they avoid unnecessary work. In the extreme case, if a lazy processing pipeline is run on an already successfully processed dataset, it should not perform any processing at all.

BDS facilitates the creation of lazy processing pipelines by means of the dependency operator (`<-`) and conditional task execution (see Listing 1, line 5 for an example). The task is defined as `task (out < in)`, meaning that it is executed only if `out` file needs to be updated with respect to `in` file: for example, if `output.file` file does not exist, has zero length, is an empty directory or has been modified before `input.file`.

Absolute serialization. This refers to the ability to save and recover a snapshot of the current execution state, compiled program, variables, scopes and program counter, a concept similar to continuations (Reynolds, 1993).

BDS can perform an absolute serialization of the current running state and environment, producing checkpoint files from which the program can be re-executed, either on the same computer or on any other computer, exactly from the point where execution terminated. Checkpoint files (or `checkpoints` for short) also allow all variables and the execution stack to be inspected for debugging purposes (`bds -i checkpoint.chp`). The most common use of checkpoints is when a task execution fails. On reaching a `wait` statement, if one or more tasks have failed, BDS creates a checkpoint, reports the reasons for task execution failure and terminates. Using the checkpoint, pipeline execution can be resumed from the point where it terminated (in this case, at the most recently executed `wait` statement) and can properly re-execute pending tasks (i.e. the tasks that previously failed execution).

Limitations. BDS is designed to afford robustness to the most common types of pipeline execution failures. However, events such as full cluster failures, emergency shutdowns, head node hardware failures or network problems isolating a subset of nodes may result in BDS being unable to exit cleanly, leading to an inconsistent pipeline state. These problems can be mitigated by a special purpose `checkpoint` statement that, as the name suggests, allows the programmer to explicitly create checkpoints. Given that the overhead of creating checkpoints is minimal (a few milliseconds compared with hours of processing time for a typical pipeline), carefully crafted checkpoint statements within the pipeline code can be useful to prevent losing processed data, mitigate damage and minimize the overhead when rerun, which can be critical for long running pipelines.

2.3.4 Other features

Here we mention some selected features that are useful in pipeline programming. Extensive documentation is available at <http://pcingola.github.io/BigDataScript>.

Automatic logging. Logging all actions performed in pipelines is important for three reasons: (i) it helps debugging; (ii) it improves repeatability; and (iii) it performs audits in cases where detailed documentation and logging are required by regulatory authorities (such as clinical trials).

Listing 2.2: `pipeline_2.bds` program. A two-step pipeline with task dependencies. The first step (line 9) requires to run `myProcess` command on a hundred input files, which can be executed in parallel. The second step (line 19) processes the output of those hundred files and creates a single output file (using fictitious `myProcessAll` command). It should be noted that we never explicitly state which hardware we are using: (i) if the pipeline is run on a dual-core computer, as each process requires 2 CPUs, one `myProcess` instance will be executed at the time until the 100 tasks are completed; (ii) if it is run on a 64-core server, then 32 `myProcess` instances will be executed in parallel; (iii) if it is run on a cluster, then 100 `myProcess` instances will be scheduled and the cluster resource management system will decide how to execute them; and (iv) if it is run on a single-core computer, execution will fail owing to lack of resources. Thus, the pipeline runs independent of the underlying architecture. The task defined in line 18 depends on all the outputs from tasks in line 8 (`mainOut <- outs`).

```

1  #!/usr/bin/env bds
2  // Step 1: Parallel processing of input files
3  string[] outs // Define a list of strings
4  for( int i=0 ; i < 100 ; i++ ) {
5      in := "input_${i}.file"
6      out := "output_${i}.file"
7      task ( out <- in, cpus=2, timeout=6*HOUR ) {
8          sys myProcess $in > $out
9      }
10     outs.add( out )           // Add all output files here
11 }
12 wait // Optional: Wait for all tasks to finish
13

```

```

14 // Step 2: Process all outputs from previous step
15 mainOut := "main.txt"
16 mainIn := outs.join(      ) // Create a string with all names
    (space-separated)
17 task ( mainOut <- outs, mem=10*G ) {
18     sys myProcessAll $mainIn > $mainOut
19 }
```

Creating log files is simple, but it adds boilerplate code and increases the complexity of the pipeline. BDS performs automatic logging in three different ways. First, it directs all process StdOut/StdErr output to the console. Second, as having a single output can be confusing when dealing with thousands of processes running in parallel, BDS automatically logs each processs outputs (StdOut and StdErr) and exit codes in separate clearly identified files. Third, BDS creates a report showing both an overview and details of pipeline execution (Fig. 2–3).

Automatic command line parsing. Programming flexible data pipelines often involves parsing command-line inputsa relatively simple but tedious task. BDS simplifies this task by automatically assigning values to variables specified through the command line. As an example, if the program in Listing 1 is called `pipeline.bds`, then invoking the program as `pipeline.bds -in another.file` will automatically replace the value of variable `in` with `another.file`.

Task re-execution. Tasks can be re-executed automatically on failure. The number of retries can be configured globally (as a command-line argument) or by a task (using the `retry` variable). Only after failing `retry+1` times will a task will be considered to have failed.

2.3.5 BDS implementation

BDS is programmed using Java and GO programming languages. Java is used for high-level actions, such as performing lexical analysis, parsing,

creating abstract syntax trees (AST), controlling AST execution, serializing processes, queuing tasks, etc. Low-level details, such as process execution control, are programmed in GO. As BDS is intended to be used by programmers, it does not rely on graphical interfaces and does not require installation of complex dependencies or Web servers.

Figure 2–2 shows the cascade of events triggered when a BDS program is invoked. First the script pipeline.bds (Fig. 2–2A) is compiled to an AST structure (Fig. 2–2B) using ANTLR (Parr, 2007). After creating the AST, a runnable-AST (RAST) is created. RAST nodes are objects representing statements, expressions and blocks from our BDS implementation. These nodes can execute BDS code, serializing their state, and recover from a serialized file, thus achieving absolute serialization. The script is run by first creating a scope and then properly traversing the RAST (Fig. 2–2C). We note that if needed, this approach could be tuned to perform efficiently, as demonstrated by modern languages, such as Dart.

When recovering from a checkpoint, the scopes and RAST are deserialized (i.e. reconstructed from the file) and then traversed in recovery mode, meaning that the nodes do not execute BDS code. When the node that was executed at the time of serialization event is reached, BDS switches to run mode and the execution continues. This achieves execution recovery from the exact state at serialization time. Checkpoints are the full state of a programs instance and are intended as a recovery mechanism from a failed execution. This includes failures owing to corrupted or missing files, as BDS will re-execute all failed tasks when recovering, thus correcting outputs from those tasks. However, checkpoints are not intended to recover from programming errors, where the user modifies the program to fix a bug, as a previously generated checkpoint is no longer valid respect to the new source code.

When a task statement is invoked, process requirements, such as memory, CPUs and timeouts, can optionally be specified. Depending on the architecture, BDS either checks that the underlying system has appropriate resources (CPUs and memory) to run the process (e.g. local computer or ssh-farm) or relies on the cluster management system to appropriately allocate the task. If all task requirements are met, a script file is created (Fig. 2–2D), and the task is executed by running an instance of bds-exec, a program that controls execution (Fig. 2–2E). This indirection is necessary for five reasons, which are described in detail below: (i) process identification, (ii) timeout enforcement, (iii) logging, (iv) exit status report and (v) signal handling.

Process identification means that bds-exec reports its process ID (PID), so that BDS can kill all child processes if the BDS script execution is terminated for some reason (e.g. the Ctrl-C key is pressed at the console).

Timeout enforcement has to be performed by bds-exec as many underlying systems do not have this capability (e.g. a process running on a server). When a timeout occurs, bds-exec sends a kill signal to all child processes and reports a timeout error exit status that propagates to the user terminal and log files.

Logging a process means that bds-exec redirects stdout and stderr to separate log files. These files are also monitored by the main BDS process, which shows the output on the console. As there might be thousands of processes running at the same time and operating systems have hard limits on the number of simultaneous file descriptors available for each user, opening all log files is not an option. To overcome this limit, BDS polls log file sizes, only opening and reading the ones that change.

Exit status has to be collected to make sure a process finished successfully. Unfortunately, there is no unified way to do this, and some cluster

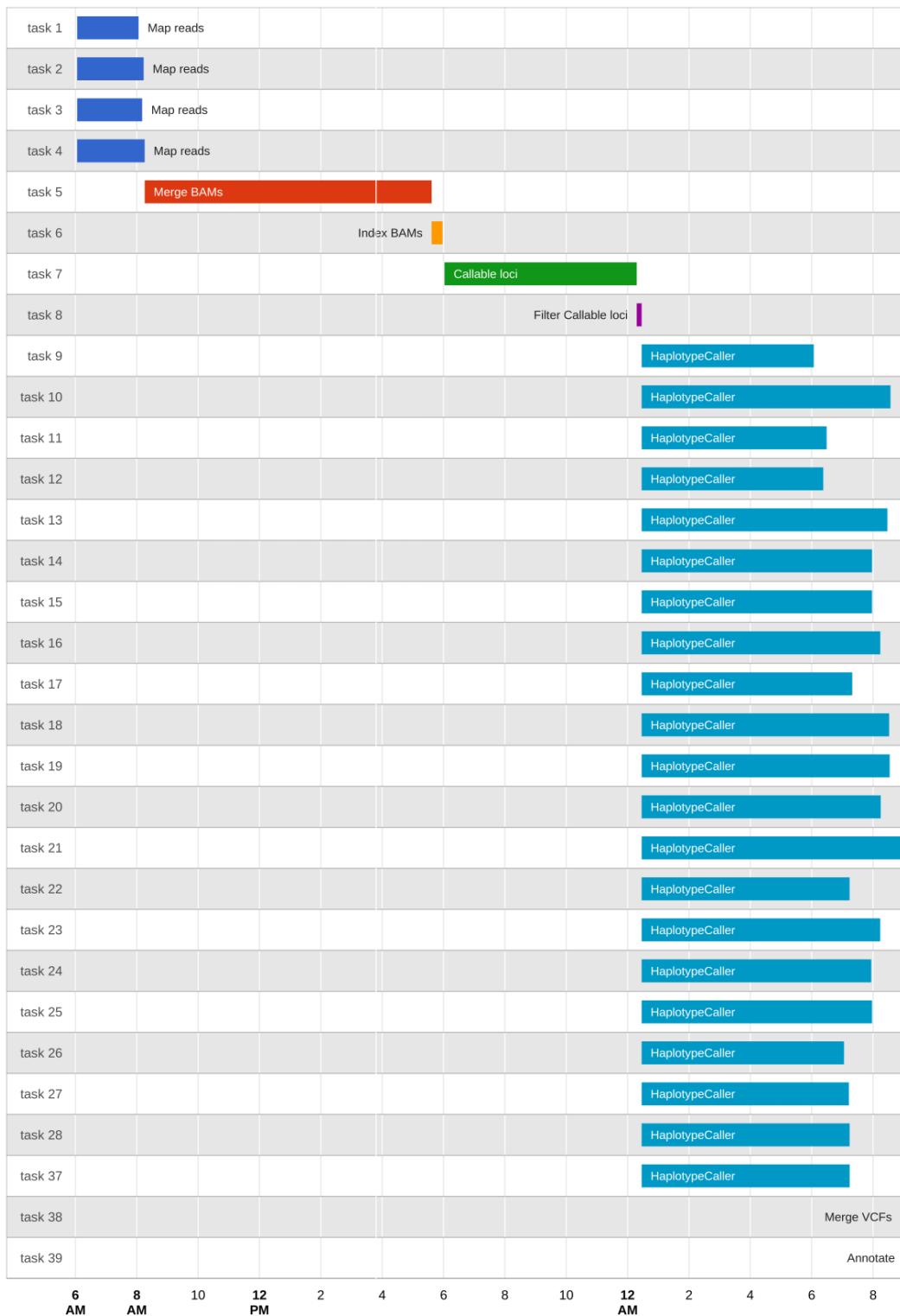


Figure 2–1: BDS report showing pipelines task execution timeline

systems do not provide this information directly. By saving the exit status to a file, bds-exec achieves two goals: (i) unified exit status collection and (ii) exit status logging.

Signal handling is also enforced by bds-exec making sure that a kill signal correctly propagated to all subprocesses, but not to parent processes. This is necessary because there is no limit on the number of indirect processes that a task can run, and Unix/Posix systems do not provide a unified way to obtain all nested child processes. To be able to keep track of all subprocesses, bds-exec creates a process group and spawns the subprocess in it. When receiving a signal from the operating system, bds-exec traps the signal and propagates a kill signal to the process group.

2.4 Results

To illustrate the use of BDS in a real-life scenario, we present an implementation of a sequencing data analysis pipeline. This example illustrates three key BDS properties: architecture independence, robustness and scalability. The data we analyzed in this example consist of high-quality short-read sequences (200 coverage) of a human genome corresponding to a person of European ancestry from Utah (NA12877), downloaded from Illumina platinum genomes (<http://www.illumina.com/platinumgenomes>).

The example pipeline we created follows current best practices in sequencing data analysis [110], which involves the following steps: (i) map reads to a reference genome using BWA (Li and Durbin, 2009), (ii) call variants using GATKs HaplotypeCaller and (iii) annotate variants using SnpEff [30] and SnpSift [31]. The pipeline makes efficient use of computational resources by making sure tasks are parallelized whenever possible. Figure 2–3 shows a flowchart of our implementation, while the pipelines source code is available at

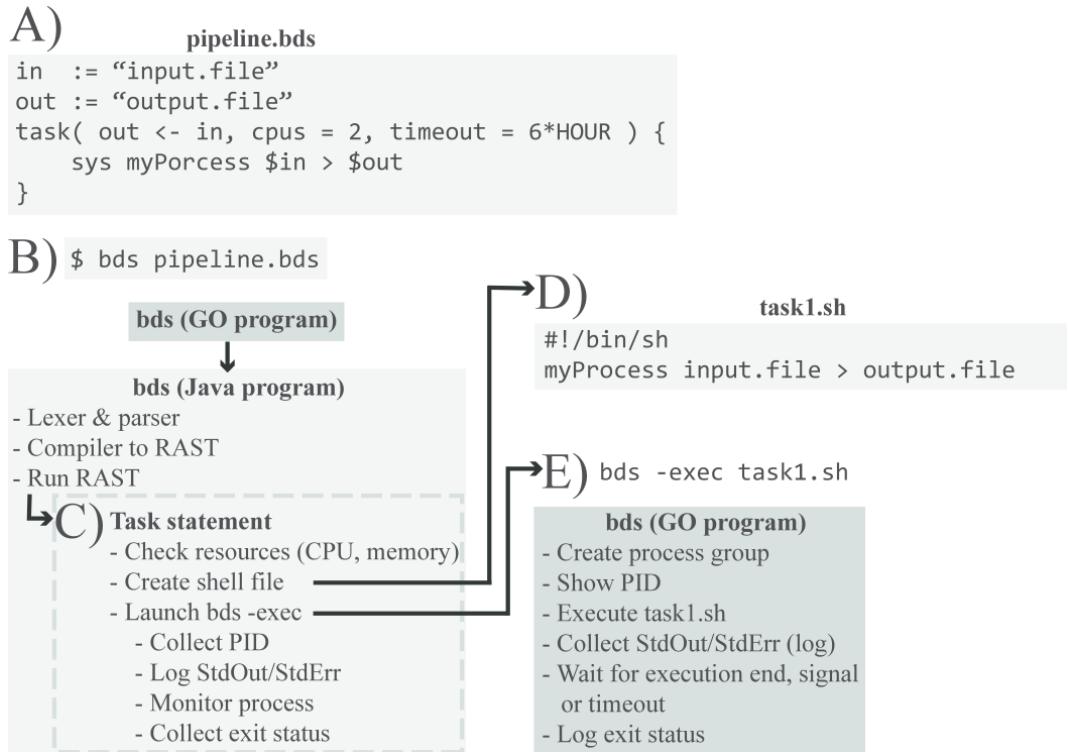


Figure 2–2: Execution example. (A) Script `pipeline.bds`. (B) The script is executed from a terminal. The GO executable invokes main BDS, written in JAVA, performs lexing, parsing, compilation to AST and runs AST. (C) When the task statement is run, appropriate checks are performed. (D) A shell script `task1.sh` is created, and a `bds-exec` process is fired. (E) `bds-exec` reports PID, executed the script `task1.sh` while capturing stdout and stderr as well as monitoring timeouts and OS signals. When a process finishes execution, the exit status is logged

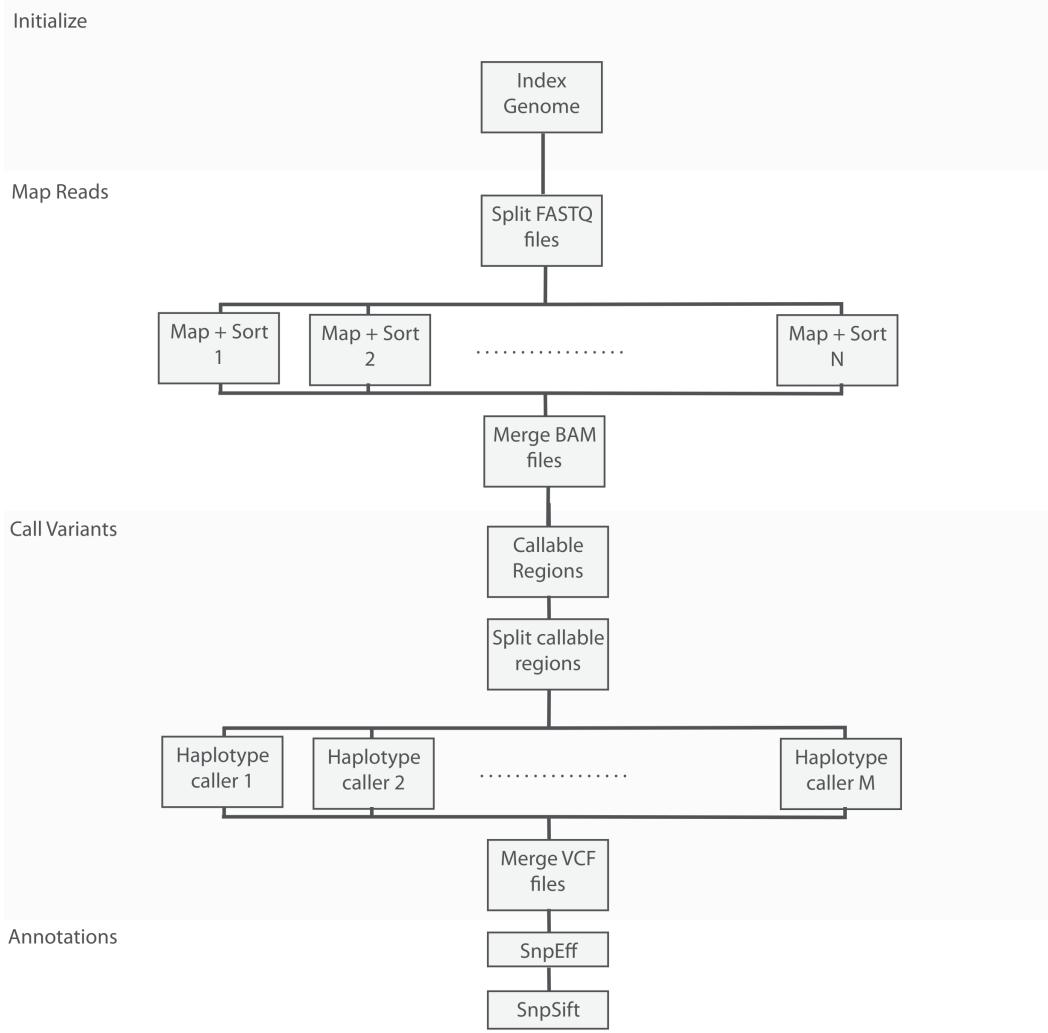


Figure 2–3: Whole-genome sequencing analysis pipeline’s flow chart, showing how computations are split across many nodes

`include/bio/seq` directory of our projects source code (<https://github.com/pcingola/BigDataScript>).

Architecture independence. We ran the exact same BDS pipeline on (i) a laptop computer; (ii) a multi-core server (24 cores, 256 GB shared RAM); (iii) a server farm (5 servers, 2 cores each); (iv) a 1200-core cluster; and (v) the Amazon AWS Cloud computing infrastructure (Table 2–4). For the purpose of this example and to accommodate the fact that running the pipeline on a laptop using the entire dataset would be prohibitive, we limited our experiment

System	CPUs	RAM	Notes
Laptop (OS.X)	4	8 GB	
Server (Linux)	24	256 GB	
Server farm (ssh)	16	8 Gb	Server farm using 8 nodes, 2 cores each.
Cluster (PBS Torque)	1200	12 TB	High load cluster (over 95%).
Cluster (MOAB) (Random failures)	1200	12 TB	High load cluster (over 95%). Hardware induced failures.
Cloud (AWS + SGE)	Inf.	Inf.	StarCluster, 8 m1.large instances.

Table 2–4: Architecture independence example

to reads that map to chromosome 20. The architectures involved were based on different operating systems and spanned about three orders of magnitude in terms of the number of CPUs (from 4 to 1200) and RAM (from 8GB to 12TB). BDS can also create a cluster from a server farm by coordinating raw SSH connections to a set of computers. This minimalistic setup only requires that the computers have access to a shared disk, typically using NFS, which is a common practice in companies and university networks.

In all cases, the overhead required to run the BDS script itself accounted for 52 ms per task, which is negligible compared with typical pipeline runtimes of several hours.

Robustness. To assess BDSs robustness, we ran the pipeline on a cluster where 10% of the nodes have induced hardware failures. As opposed to software failures, which are usually detected by cluster management systems, hardware node failures are typically more difficult to detect and recover from. In addition, we elevated the cluster load to 495% to make sure the pipeline was running on less than ideal conditions. As shown in Table 2–4, the pipeline finished successfully without any human intervention and required only 30% more time than in the ideal case scenario because BDS had to rerun several failed tasks. This shows how BDS pipelines can be robust and recover from multiple failures by using lazy processing and absolute serialization mechanisms.

Dataset	Dataset size	System	CPUs	RAM
chr20	2 GB	Laptop (OS.X)	4	8 GB
Whole genome	1.5 TB	Cluster (MOAB)	22000	80 TB

Table 2–5: Scaling dataset sized by a factor of 1000

Scalability. To assess BDSs scalability, we ran exactly the same pipeline on two datasets that vary in size by several orders of magnitude (Table 2–5): (i) a relatively small dataset (chromosome 20 subset, 2GB) that would typically be used for development, testing and debugging and (ii) a high-depth whole-genome sequencing dataset (over 200 coverage, roughly 1.5 TB).

2.5 Discussion

We introduced BDS, a programming language that simplifies implementing, testing and debugging complex data analysis pipelines. BDS is intended to be used by programmers in a similar way to shell scripts, by providing glue for several tools to ensure that they execute in a coordinated way. Shell scripting was popularized when most personal computers had a single CPU and clusters or clouds did not exist. One can thus see BDS as extending the hardware abstraction concept to data-center level while retaining the simplicity of shell scripting.

BDS tackles common problems in pipeline programming by abstracting task management details at the programming language level. Task management is handled by two statements (`task` and `wait`) that hide system architecture details, leading to cleaner and more compact code than general-purpose languages. BDS also provides two complementary robustness mechanisms: lazy processing and absolute serialization.

A key feature is that being architecture agnostic, BDS allows users to code, test and debug big data analysis pipelines on different systems than the ones intended for full-scale data processing. One can thus develop a pipeline on a laptop and then run exactly the same code on a large cluster. BDS also provides mechanisms that eliminate many boilerplate programming tasks, which in our experience significantly reduce pipeline development times. BDS can also reduce CPU usage, by allowing the generation of code with fewer errors and by allowing more efficient recovery from both software and hardware failures. These benefits generally far outweigh the minimal overhead incurred in typical pipelines.

CHAPTER 3

A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain $w^{1118}; iso - 2; iso - 3$

3.1 Preface

As this thesis is focused on extracting biological insight from sequencing data, in this chapter we examine algorithms we created for calculating “functional annotations” of genomic variants. In essence, functional variant annotations are bits of biological knowledge that allow us to make prioritize variants that are assumed to be more relevant to the phenotypic trait under study and to filter out variants assumed irrelevant. The spectrum of functional annotations for a genomic variant is wide and may involve information on which genes are affected by the variant, how the protein product is affected, how conserved is the genomic region the variant lies onto, and which clinically relevant information is associated with the loci; just to mention a few typical use cases.

When trying to find variants that affect risk of complex disease, statistical power is paramount. We need to be able to “separate wheat from chaff”. In our context this means two different but closely related tasks: i) performing functional annotations, and ii) using that information for prioritizing variants (and filtering out the ones we suspect are not related to the particular trait under study). Failing to efficiently filter out irrelevant variants would reduce our statistical power as more statistical tests are calculated, thus would decrease our chances of finding the associations we are looking for. In order to efficiently annotate and filter variants, we created two software packages

called SnpEff and SnpSift that deal with the annotation and filtering aspects respectively.

The rest of the chapter is published in: **P. Cingolani**, A. Platts, M. Coon, T. Nguyen, L. Wang, S.J. Land, X. Lu, D.M. Ruden, et al. “A program for annotating and predicting the effects of single nucleotide polymorphisms, snpeff: Snps in the genome of drosophila melanogaster strain $w^{1118}; iso - 2; iso - 3$ ”. Fly, 6(2), 2012.

3.2 Abstract

We describe a new computer program, SnpEff, for rapidly categorizing the effects of single nucleotide polymorphisms (SNPs) and other variants such as multiple nucleotide polymorphism (MNPs) and insertion-deletions (InDels), in whole genome sequences. Once a genome is sequenced, the SnpEff program can be used to annotate and classify genetic polymorphisms based on their effects on annotated genes, such as synonymous or non-synonymous SNPs, start codon gains or losses, stop codon gains or losses; or based on their genomic locations, such as intronic, 5' untranslated region (5' UTR), 3' UTR, upstream, downstream or intergenic regions. Here the use of SnpEff is illustrated by annotating $\tilde{3}56,660$ candidate SNPs in $\tilde{1}17$ Mb unique sequences, representing a substitution rate of $\tilde{1}/305$ nucleotides, between the Drosophila melanogaster $w^{1118}; iso - 2; iso - 3$ strain and the reference $y^1; cn^1bw^1sp^1$ strain [123]. We show that $\tilde{1}5,842$ SNPs are synonymous and $\tilde{4},467$ SNPs are non-synonymous ($N/S \tilde{0}.28$) and the remainder are in other categories, such as stop codon gains (38 SNPs), stop codon losses (8 SNPs) and start codon gains (297 SNPs) in the 5' UTR. We found, as expected, that the SNP frequency is proportional to the recombination frequency (i.e., highest in the middle of chromosome arms). We also found that start-gained and stop-lost SNPs in Drosophila melanogaster often encode N-terminal and C-terminal amino acids

that are conserved in other *Drosophila* species. This suggests that the 5' and 3' UTRs are reservoirs of cryptic genetic variation that can be used multiple times during the evolution of the *Drosophila* genus. At this time, SnpEff has been set up for annotating DNA polymorphisms of over 320 genome versions of multiple species including the human genome. It has already been used by over 50 institutions and universities in the bioinformatics community. Tools such as SnpEff are valuable because, as sequencing becomes cheaper and more available, whole genome sequencing is becoming more important in model organism genetics.

3.3 Introduction

When we re-sequenced the $w^{1118}; iso - 2; iso - 3$ genome in 2009, [123] bioinformatics tools available then were unable to rapidly categorize the $\tilde{3}56,660$ SNPs as comparing to the $y^1; cn^1bw^1sp^1$ reference strain. At the time, other available tools such as ENSEMBLs variant web application (<http://ensembl.org>) could only analyze a few hundred to a few thousand SNPs per batch. Therefore, over the past couple of years, we have been developing a new program, SnpEff, which is able to analyze and annotate thousands of variants per second. In addition to SnpEff, other programs to annotate genomic variants are currently now available, such as Annotate Variation (ANNOVAR) [157] and Variant Annotation, Analysis and Search Tool (VAAST) [130]. However, SnpEff supports more genome versions, is open source for any user, supports variant call format (VCF) files and it is marginally faster (although the speeds of SnpEff, ANNOVAR and VAASST are comparable). Table S1 shows a feature comparison of some currently available software packages.

SnpEff, an abbreviation of “SNP effect,” is a multi-platform open source variant effect predictor program. SnpEff annotates variants and predicts the

# SNP	Gene_name	Effect	Old_AA/new_AA	Old_codon/New_codon	Codon_Num (CDS)	CDS_size
chr2L:10006682_C/T	CG31755	UPSTREAM: 541 bases				
chr2L:10006758_G/A	CG31755	UPSTREAM: 465 bases				
chr2L:10007289_G/A	CG4747	SYNONYMOUS_CODING	L/L	TTG/TTA	489	1809
chr2L:10007319_G/C	CG4747	SYNONYMOUS_CODING	G/G	GGG/GGC	499	1809
chr2L:10007356_A/T	CG4747	INTRON				1809
chr2L:10007363_T/A	CG4747	INTRON				1809

Table 3–1: Output of SnpEff.

coding effects of genetic variations, such as SNPs, insertions and deletions (INDELs) and multiple nucleotide polymorphisms (MNPs) (<http://SnpEff.sourceforge.net/>). The main features of SnpEff include: (1) speedthe ability to make thousands of predictions per second; (2) flexibilitythe ability to add custom genomes and annotations; (3) the ability to integrate with Galaxy, an open access and web-based platform for computational bioinformatic research (<http://gmod.org/wiki/Galaxy>); (4) compatibility with multiple species and multiple codon usage tables (e.g., mitochondrial genomes); (5) integration with Genome Analysis Toolkit (GATK) [110]; and (6) ability to perform non-coding annotations. When SnpEff was integrated into the GATK, it replaced the ANNOVAR program for variant analyses.

A simple walk-through example on how to analyze sequencing data to calculate variants and their effects is shown in Listing SL1. This example is intended for illustration purposes only since many additional steps are routinely used in re-sequencing data analysis pipelines, but design of a fully featured pipeline is beyond the scope of this paper.

Here, we report the results of SnpEff (version 1.9.6) analyses of the 356,660 candidate SNPs that we identified in *w¹¹¹⁸;iso-2;iso-3* with respect to the *y¹;cn¹bw¹sp¹* reference strain as reported in our previous paper. 1 This is of great interest to the Drosophila community because thousands of transposon

insertion stocks [150] and hundreds of deficiency stocks [119],[119] were generated in the $w^{1118}; iso - 2; iso - 3$ genetic background. The large number and potential severity of many SNPs in the two laboratory strains was a surprising finding, and the possible evolutionary implications of this finding are discussed.

3.4 Results

Formats used in SnpEff. To understand the potential effects of large numbers of SNPs in genome sequence comparisons, we developed an open-source tool, SNPeff, to classify SNPs based on gene annotations. Table 3–1 shows the beginning portion of the output generated by SnpEff when the SNPs in $w^{1118}; iso - 2; iso - 3$ were compared with the reference genome, $y^1; cn^1bw^1sp^1$ that is represented in *Drosophila melanogaster* release 5.3. A more complete SnpEff effect list is shown in Table 3–2. Before using SnpEff, an input file must be generated that lists all of the SNPs and INDELs in a genome. We published the input file for $w^{1118}; iso - 2; iso - 3$ in our previous paper, 1 and it was derived by comparing hundreds of millions of short sequence reads ($\tilde{20}$ -fold genome coverage) and identifying SNPs based on a Sequence Alignment/Map tools (SAMtools) quality score for each nucleotide in the genome [95].

Input formats supported by SnpEff are variant call format (VCF) [46], tab separated TXT format; and and the SAMtools

Pileup format [95]. VCF was created by the 1,000 Genomes project and it is currently the de facto standard for variants in sequencing applications. The TXT and Pileup formats are currently deprecated and being phased out.

SnpEff also supports two output formats, TXT and VCF. The information provided in both of them includes four main groups: (i) variant information (genomic position, the reference and variant sequences, change type, heterozygosity, quality and coverage); (ii) genetic information (gene Id, gene

name, gene biotype, transcript ID, exon ID, exon rank); and (iii) effect information (effect type, amino acid changes, codon changes, codon number in CDS, codon degeneracy, etc.).

Whenever multiple transcripts for a gene exist, the effect and annotations on each transcript are reported, so one variant can have multiple output lines. Table 3–3 shows the information provided by each column in TXT format and Table 3–4 shows the information provided in VCF format. When using VCF format, the effect information is added to the information (INFO) fields using an effect (EFF) tag. As in the case of TXT output, if multiple alternative splicing products are annotated for a particular gene, SnpEff provides this information for each annotated version (see Sup. Data File 1 for the complete SnpEff output for *w¹¹¹⁸;iso – 2;iso – 3*).

Predicted effects are with respect to protein coding genes. Variants affecting non-coding genes are annotated and the corresponding biotype is identified, whenever the information is available. A “biotype” is a group of organisms having the same specific genotype.

According to SnpEff (version 1.9.6), the largest number of SNPs in *w¹¹¹⁸;iso – 2;iso – 3* are in introns (130,126) followed by those in upstream (76,155), downstream (71,645) and intergenic (51,783) regions (Fig. 3–5). “Upstream” is defined as 5 kilobase (kb) upstream of the most distal transcription start site and “downstream” is defined as 5 kb downstream of the most distal polyA addition site, but these default variables can be easily adjusted. SnpEff also found thousands of SNPs within the exons. For example, there are 3,718 SNPs in the 3’ untranslated regions (3’ UTR) and 2,508 SNPs in the 5’ untranslated regions (5’ UTR). The SNPs in the upstream, downstream, 5’ and 3’ UTR regions might affect transcription or translation, but the actual effects have to

Effect	Note
INTERGENIC	The variant is in an intergenic region
UPSTREAM	Upstream of a gene (default length: 5K bases)
UTR_5_PRIME	Variant hits 5'UTR region
UTR_5_DELETED	The variant deletes and exon which is in the 5'UTR of the transcript
START_GAINED	A variant in 5'UTR region produces a three base sequence that can be a START codon
SPLICE_SITE_ACCEPTOR	The variant hits a splice acceptor site (defined as two bases before exon start, except for the first exon)
SPLICE_SITE_DONOR	The variant hits a Splice donor site (defined as two bases after coding exon end, except for the last exon)
START_LOST	Variant causes start codon to be mutated into a non-start codon
SYNONYMOUS_START	Variant causes start codon to be mutated into another start codon
CDS	The variant hits a CDS
GENE	The variant hits a gene
TRANSCRIPT	The variant hits a transcript
EXON	The variant hits an exon
EXON_DELETED	A deletion removes the whole exon
NON_SYNONYMOUS_CODING	Variant causes a codon that produces a different amino acid
SYNONYMOUS_CODING	Variant causes a codon that produces the same amino acid
FRAME_SHIFT	Insertion or deletion causes a frame shift
CODON_CHANGE	One or many codons are changed
CODON_INSERTION	One or many codons are inserted
CODON_CHANGE_PLUS_CODON_INSERTION	One codon is changed and one or many codons are inserted
CODON_DELETION	One or many codons are deleted
CODON_CHANGE_PLUS_CODON_DELETION	One codon is changed and one or more codons are deleted
STOP_GAINED	Variant causes a STOP codon
SYNONYMOUS_STOP	Variant causes stop codon to be mutated into another stop codon
STOP_LOST	Variant causes stop codon to be mutated into a non-stop codon
INTRON	Variant hits intron. Technically, hits no exon in the transcript
UTR_3_PRIME	Variant hits 3'UTR region
UTR_3_DELETED	The variant deletes and exon which is in the 3'UTR of the transcript
DOWNTSTREAM	Downstream of a gene (default length: 5K bases)
INTRON_CONSERVED	The variant is in a highly conserved intronic region
INTERGENIC_CONSERVED	The variant is in a highly conserved intergenic region

Table 3–2: Detailed effect list from SnpEff

be confirmed case-by-case. In the next few sections, we present examples of several types of SNPs that might affect the protein function.

Heterozygosity is not considered in the *w¹¹¹⁸; iso-2; iso-3* sequence because the stock was isogenized and only high quality (i.e., homozygous SNPs) were used for this analysis. 1

The SnpEff website (<http://snpeff.sourceforge.net/SnpSift.html>) has a frequently asked questions (FAQ) section that addresses most issues that a user might have in operating this program.

SNPs that generate new start codons. There are 297 SNPs that potentially generate a new translation initiation codon in the 5' UTR (start-gained SNPs). The most common translation initiation codon is AUG, which is coded by ATG in the genome. To be thorough, we also included CUG and UUG codons, which code for leucine, as these codons can also be used to initiate translation in rare genes in Drosophila and mammals [148],[76]. There are 60 genes with ATG start-gained SNPs (Table 3–6), 99 genes with CTG start-gained SNPs and 120 genes with TTG startgained SNPs in *w¹¹¹⁸; iso – 2; iso – 3*, all by definition in 5' UTR regions, compared with the reference genome (the reading frame is indicated on the SnpEff table). Most of the ATG start-gained SNPs are within 1 kb of the annotated translation start (Table 3–6), but this probably reflects the fact that most 5' UTR sequences are less than 1 kb long. Less than expected by chance, only $\tilde{2}5\%$ of the ATG start-gain SNPs are in the same reading frame as the annotated translation start point (Table 3–6). Since 33% of in frame ATG start-gained SNPs are expected by chance, this suggests that there might be weak selection against this class of SNPs. Of the 60 genes with ATG start-gained SNPs, five genes have two ATG start-gained SNPs and one gene has three startgained SNPs; the remaining 54 genes have a single start-gained SNP. Since SnpEff does not take into account the Kozak consensus sequence flanking the AUG site, 5'-ACC AUG G-3', that is generally required for efficient translation [86], and thus further assessment is required to determine whether a start-gained SNP is actually used.

Gene ontology (GO) pathway analysis of the genes affected by the 297 start-gain SNPs in *w¹¹¹⁸; iso – 2; iso – 3* was done using DAVID (Database for Annotation, Visualization and Integrated Discovery) [50, 74]. We found that

Column	Notes
Chromosome	Chromosome name (usually without any leading 'chr' string)
Position	One based position
Reference	Reference
Change	Sequence change
Change type	Type of change (SNP, MNP, INS, DEL)
Homozygous	Is this homozygous or heterozygous (Hom, Het)
Quality	Quality score (from input file)
Coverage	Coverage (from input file)
Warnings	Any warnings or errors.
Gene_ID	Gene ID (usually ENSEMBL)
Gene_name	Gene name
Bio_type	BioType, as reported by ENSEMBL
Trancript_ID	Transcript ID (usually ENSEMBL)
Exon_ID	Exon ID (usually ENSEMBL)
Exon_Rank	Exon number on a transcript
Effect	Effect of this variant. See details below
old_AA/new_AA	Amino acid change
old_codon/new_codon	Codon change
Codon_Num(CDS)	Codon number in CDS
Codon_degeneracy	Codon degeneracy
CDS_size	CDS size in bases
Custom_interval_ID	If any custom interval was used, add the IDs here (may be more than one)

Table 3–3: Information provided by SnpEff in tab separated output format (TXT)

the GO categories “tissue morphogenesis,” “immunoglobulin like,” “developmental protein,” and “alternative splicing” are significantly enriched after multiplecomparisons correction by false-discovery rate (FDR \downarrow 0.001; Table 3–7). These categories are interesting because they predominantly contain proteins that show a wide degree of intra- and interspecies variability. For example, the immunoglobulin loci, which are highly divergent among humans and in other vertebrates, are used for antigen recognition [93]. Also, developmental proteins and proteins involved in tissue morphogenesis often have both conserved domains, such as the Hox domain, and highly divergent domains that maintain morphological diversity within a species, such as the trans-activation domains [131, 67].

Our previous analyses suggest that most of the SNPs that we identified in *w*¹¹¹⁸; *iso* – 2; *iso* – 3 are probably genuine and can be validated by capillary sequencing. A common worry about nextgeneration sequencing data in general is that SNPs are vastly over estimated. One might think that if a large fraction of the identified SNPs had the predicted “effects”, the organism would not be viable. However, since short-read next-generation sequencing has a high error rate, such as the short-read sequences we obtained with the Illumina platform, further validation of specific SNPs is needed to be absolutely certain. Further validation of SNPs is best done with long-range DNA sequencing, such as with traditional capillary sequencing, or sequencing with the Roche [162], and many other DNA sequencing instruments that are now available [138] (see [123] for validation examples with capillary sequencing).

An example of a start-gained SNP is found in the 5’ UTR of Ecdysone inducible protein 63E (Eip63E) gene, which is predicted to be a cyclin J dependent kinase required for oogenesis and embryonic development (Fig. 3–8) [101]. The potential start-gain SNP (A \downarrow G) in Eip63E changes 5'-ATA-3'

to 5'-ATG-3' in the same reading frame with no in-frame intervening stop codons (Fig. 3–8A). If translation occurs at the new start-gained SNP, it would produce a protein with 57 additional N-terminal amino acids compared with the reference gene (Fig. 3–8B). However, the three bases prior to the new 5'-ATG-3' sequence, 5'-AAT-3', is a poor match to the Kozak consensus sequence, 5'-ACC-3', discussed above in reference 12. Therefore, it is unclear whether the startgain SNP in Eip63E is recognized by the ribosomal machinery.

It is interesting that a BLASTp search of the protein database reveals that the N-terminal 57 amino acids in Eip63E are 63% identical (36/57) to the 58 N-terminal amino acids of the orthologous gene in *Drosophila yakuba*, but not to any other *Drosophila* species. *D. yakuba* is very close to *D. melanogaster* in the phylogeny. This suggests that the 5' UTR of Eip63E might be a source for cryptic genetic variation encoding novel N-terminal protein sequences that potentially modulates protein function (see Discussion).

SNPs that generate new stop codons. Another surprise in our SnpEff analysis was the identification of 28 stop-gained SNPs and 5 stop-lost SNPs in *w¹¹¹⁸; iso – 2; iso – 3* (Table 3–9). A stop-gained SNP, classically called a nonsense SNP, has a coding codon changed to a stop codon, UAA, UAG, UGA [19]. Three genes, *oc/ otd*, LRP1 and *trol9*, have two stop-gained SNPs. Surprisingly at least 8 of the stop-gained SNPs are in genes that encode essential proteins, and these are *Dif*, *dp*, *ex*, *MESR4*, *mew*, *oc/ otd*, *tai* and *trol*. It is not known whether the other stop-gained SNPs also affect essential protein-coding genes because their functions have not yet been characterized (according to www.flybase.org). We note that what would be a stop-gained SNP in *w¹¹¹⁸; iso – 2; iso – 3* would be a stop-lost SNP in the reference strain,

and vice versa, because the sequence of the ancestral *Drosophila melanogaster* strain that gave rise to both of these strains is not known.

An important consideration with stop-gained and stop-lost SNPs is whether the C-terminal amino acids in the longest version of the protein that are not present in the shortest version of the protein are conserved in other *Drosophila* species. If the additional C-terminal amino acids are not conserved, then these amino acids might not affect the essential function of the protein but they might exert modulatory effects. If the additional C-terminal amino acids are conserved in multiple *Drosophila* species, then their loss might adversely affect the function of the protein. Therefore, in Table 3–9, we further classify the stop-gained and stop-lost SNPs into four categories: Category 1, including 23 genes, with both the N-terminal and novel C-terminal regions conserved among *Drosophila* species and other organisms; Category 2, including only one gene, with the entire gene sequence not conserved even among other *Drosophila* species; Category 3, with two genes, with the novel C-termini not conserved among other *Drosophila* species. In this category, the N-termini are conserved among *Drosophila* species, but this conservation is not maintained beyond the *Drosophila* genus (this class is likely a novel gene that arose in the *Drosophila* genus); and Category 4, including seven genes, with the novel C-terminal regions conserved among other *Drosophila* species but not beyond the *Drosophila* genus. In this category, the N-terminus is conserved beyond the *Drosophila* genus (this class probably has a C-terminal domain with a modulatory role in the *Drosophila* genus but not beyond the genus).

An example of an essential protein-coding gene in Category 4, where the novel C-terminus is not conserved outside the *Drosophila* genus, is *oceliless* (*oc*), also known as *orthodenticle* (*otd*) (Fig. 3–10). The *oc/otd* gene has two in-frame stop-gained SNPs in *w¹¹¹⁸; iso – 2; iso – 3*. The *oc/otd* gene is a

Sub-field	Notes
Effect	Effect of this variant. See details below
Codon_Change	Codon change: old_codon/new_codon
Amino_Acid_change	Amino acid change: old_AA/new_AA
Warnings	Any warnings or errors
Gene_name	Gene name
Gene_BioType	BioType, as reported by ENSEMBL
Coding	[CODING NON_CODING]. If information reported by ENSEMBL (e.g., has 'protein_id' information in GTF file)
Transcript	Transcript ID (usually ENSEMBL)
Exon	Exon ID (usually ENSEMBL)
Warnings	Any warnings or errors (not shown if empty)

Table 3–4: Information provided by SnpEff in variant call format (VCF)

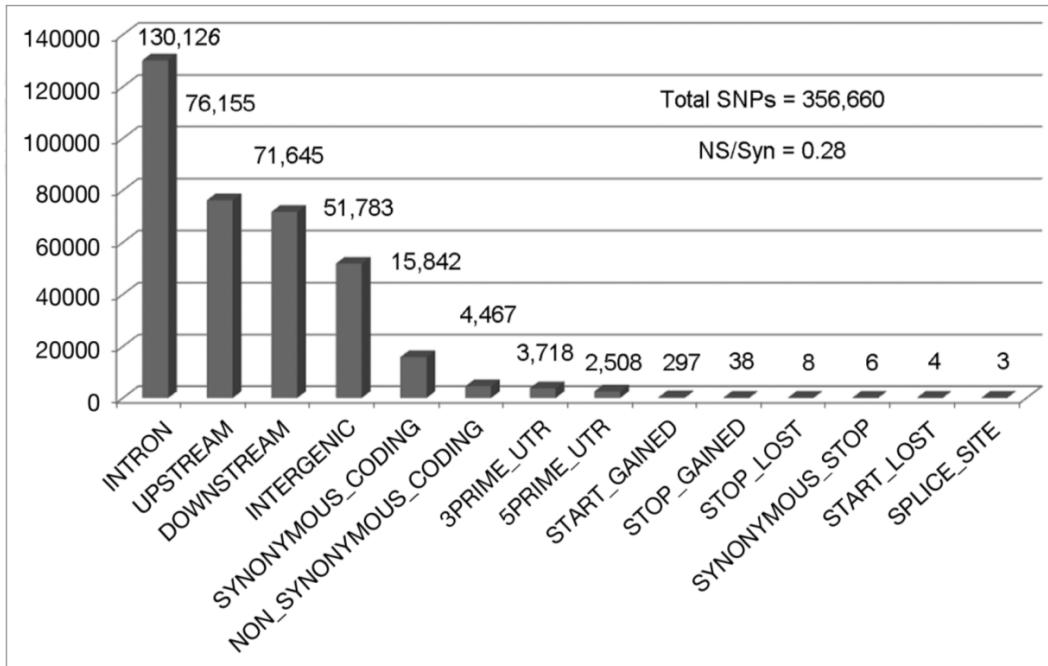


Figure 3–5: Classification of SNPs in $w^{1118}; iso - 2; iso - 3$. The number of NSPs in each class is shown above the bar. The quality score was arbitrarily set at 70 and above for this graph.

Gene_name	bases from TSS	Gene_name	bases from TSS	Gene_name	bases from TSS
a	386 (-)	CG4766	367 (-)	MESR3	454 (-)
Ace	652 (-)	CG4839	293 (-)	Mipp2	67 (-)
Axn	107 (-)	CG5103 (2)	104/17 (-/-)	osp	358 (-)
btsz	228 (+)	CG6024	269 (-)	p120ctn	119 (-)
Calx	582 (+)	CG7985	60 (+)	Pld	144 (+)
CAP	1224 (+)	CG8026	612 (+)	Pli	196 (-)
CG10186	402 (+)	CG8176	128 (-)	Pvr (2)	472/915 (-/+)
sesn	147 (+)	cpo	168 (+)	pxb (2)	50/76 (-/-)
CG12355	151 (-)	dac	103 (-)	rib	2 (-)
CG13802 (3)	490/575/635 (-/-/-)	dpr15	433 (-)	rn	142 (-)
haf	89 (-)	EcR	160 (-)	Samuel	517 (-)
CG15086	114	Eip63E	171 (+)	sli	307 (-)
CG15878	52 (-)	fdl (2)	307/437 (-/-)	so	5252 (-)
CG18522	40 (-)	frtz	196 (-)	Sobp	24 (+)
CG30419	253 (-)	GC	76 (-)	sppt	358 (-)
CG31163	998 (-)	Gug	70 (-)	Strn-Mlck (2)	210/228 (+/+)
Dscam3	269 (-)	inv	771 (+)	tai	203 (-)
CG31688	430 (-)	lpk1	376 (-)	vn	1793 (-)
CG32048	63 (+)	klu	576 (+)	wg	231
CG32150	747 (+)	Mbs	10 (-)	Wnt4	680 (-)

Table 3–6: 60 Genes with start-gained SNPs with ATGs

Hox-family transcription factor required for photoreceptor development in the compound eye and the light-sensing ocellus, embryonic development and brain segmentation [1, 170]. The Hox domain is 60 amino acids, 59 of which are identical with the human Otd protein. The Hox domains, which arose before invertebrates and vertebrates split several hundred million years ago, are among the most conserved protein domains in bilaterally-symmetric organisms in evolution [78]. The two stop-gained SNPs are in the non-conserved C-terminal region of Oc/Otd, which is thought to have a transcriptional-regulatory function. Since both strains are viable, both oc/otd genes are apparently functional although they encode a protein with 489 amino acids in *w*¹¹¹⁸; *iso* – 2; *iso* – 3, and a protein with 543 amino acids in the reference genome (Table 3–7).

An example of a stop-lost gene in class c, where the C-terminus is not conserved even among the Drosophila genera, is CG13958 that encodes a protein of unknown function (Fig. 3–11). In *w*¹¹¹⁸; *iso* – 2; *iso* – 3, CG13958 encodes

Term	Count	%	pvalue	List Total	Pop Hits	Pop Total	Fold Enrichment	Bonferroni	Benjamini	FDR
tissue morphogenesis	21	8.898305	2.07E-08	147	247	7937	4.590515	2.37E-05	2.37E-05	3.33E-05
immunoglobulin-like	16	6.779661	3.40E-08	198	132	10196	6.241812	1.42E-05	7.08E-06	4.77E-05
developmental protein	29	12.28814	2.75E-07	229	540	12980	3.043992	3.99E-05	3.99E-05	3.27E-04
alternative splicing	31	13.13559	3.82E-07	229	616	12980	2.852464	5.53E-05	2.77E-05	4.53E-04
tissue morphogenesis	FRTZ, NRX-IV, ESG, WG, PBL, SFL, MBS, RIB, TOW, WNT4, FORM3, SLI, EIP63E, PHL, YRT, FAS, SRC64B, TWI, DLG1, BTSZ, HS6ST									
immunoglublin-ilke	CG31814, DPR15, PVR, DPR16, CG14521, KLG, VN, CG12484, BEAT-IB, CG10186, DPR2, STRN-MLCK, CG34371, KEK5, FAS, CG15630									
developmental protein	VN, ESG, DEI, INV, DAB, AWH, SCRIB, BICC, MST87F, WNT4, RIG, SLI, NUMB, PIP, INE, TWI, DLG1, FOXO, PTP10D, WG, AXN, EIP74EF, BUN, SO, FZ2, FDL, SCYL, SRC64B, POXN									
alternative splicing	CPO, CPN, ECR, VN, CG11299, RN, DAB, AWH, SCRIB, INX7, SLI, PIP, NRV2, INE, DLG1, L(1)G0196, CG32048, FOXO, PTP10D, CYCT, WG, EIP74EF, BUN, CG13624, GLUT1, OSP, FDL, SSP4, PHL, SCYL, RDGC									

Table 3–7: Genes with start-gained SNP GO categories in *w*¹¹¹⁸; *iso*–2; *iso*–3

a protein of 48 amino acids but in the reference genome it encodes a protein with 84 amino acids. When BLASTp was done with the non-redundant (nr) data set, there was not much homology beyond the 38th amino acid within the *Drosophila* genus. However, there was a near perfect (37/38) identity of the first 38 amino acids in four other *Drosophila* species: *Drosophila grimshawi*, *Drosophila yakuba*, *Drosophila erecta* and *Drosophila virilis* (Fig.3–11). This protein likely arose in the *Drosophila* genus since it has no known homologs outside of this genus.

There are also five stop-lost SNPs in *w*¹¹¹⁸; *iso* – 2; *iso* – 3 (Table 3–7). All of these SNPs are in predicted protein-coding genes, metabotropic GABA-B receptor subtype 1 (GABA-B-R1), CG13958, CG4975, brown (bw), and POU domain motif 3 (pdm3). It is not known whether any of the these genes are essential in *Drosophila* besides bw, which is not required for viability. However, the metabotropic GABA-B receptor subtype 1 (GABAB-R1) gene is required for normal behavior in mice [78] and the ortholog is therefore likely also essential in *Drosophila*, although no phenotypic data are available (www.flybase.org). The bw gene is classic gene first described in 1921 by Waaler, [155] which causes the eyes to be brown rather than red and encodes an

ATPase binding cassette (ABC) transporter [137]. The bw 1 mutation in the reference strain is a spontaneous allele with a 412-transposon repeat insertion [52], which would have been missed in our nextgeneration sequencing data because the input sequence we analyzed contained only short-read sequences that mapped uniquely to the reference genome.

Not much is known about the functions of several genes with in-frame stop-gained SNPs. The pdm3 gene is expressed in the larval and adult nervous system, and it encodes a highlyconserved Hox domain, but no phenotypic data are available (www.flybase.org). No phenotypic data are available for either CG13958 or CG4975. The protein encoding CG13958 has no known conserved domain, and its peak expression is observed within 0624 h of embryogenes, during early larval stages, at stages throughout the pupal period, and in the adult male (www.flybase.org). The protein encoded by CG4975 has an Armadillolike helical domain and an Ataxin-10 domain and has expression in the hind gut during the late larval and periods (www.flybase.org) [27].

Some of the stop-lost SNPs have interesting consequences. For example, a stop-lost SNP in *w¹¹¹⁸; iso – 2; iso – 3* is in the CG13958 gene and causes an extension of eight amino acids before the next stop codon in 3' UTR sequence is reached (Fig. 3–13). Since the C-termini of CG13958 vary in *w¹¹¹⁸; iso – 2; iso – 3* and the reference strains of *Drosophila melanogaster*, it is conceivable that the C-terminus might also fluctuate in other *Drosophila* species. To test this idea, we investigated the C-terminal regions of CG13958 homologs in other *Drosophila* species.

We found that CG13958 homologs have variable C-terminal amino acids in different species of *Drosophila*. When the CG13958 protein is analyzed by protein Basic Local Alignment Search Tool (BLASTp) with the non-redundant (nr) protein database (<http://www.ncbi.nlm.nih.gov/>), at least two *Drosophila*

A	
G	start gain SNP in w ¹¹¹⁸
<p>↑</p> <p>Stop His Cys Ser Arg Ser Leu Ser Ala Ala Gly Ala Val Glu Ala Thr Thr Thr TGA CTA ATA CAC TGC TCG CGA AGT TTG TCA GCT GCT GGC GCC GTG GAA GCA ACA ACA ACA</p> <p>Lys Leu Thr Thr Ser Thr Ser Ala Thr Thr Ser Ala Phe Tyr Arg Ala Ala Thr Ser AAA CTA ACC ACA TCC ACA TCG GCA ACA ACA TCG GCT TTC TAC AGA GCA GCG ACG TCG</p> <p>Ala Ser Ala Glu Ala Ser Ala Cys Thr Thr Pro Ala Thr Iso Lys Ser Lys Thr Lys Thr GCG TCG GCA GAG GCC TCT GCC TGC ACA ACA CCA GCA ATA AAA TCA AAA ACT AAA ACT</p> <p>→ translation start site in <i>Drosophila melanogaster</i> reference strain</p> <p>Met Ala Thr Thr Thr Thr Thr Thr Gln Ala Thr Asn Ala Lys Asp Gln Val ATG GCC ACC ACC ACA ACA ACA ACG GGG GCA ACA AAT GCT AAA GAT GGC GTC</p>	
B	
<p>→</p> <p>D mel 1 MHCSRSLSAAGAVEATTKLTTSTSATTTS-AFYRAATSASAEASACTTPATIKSKTKTM 60 MHCSRSLSAAGAV+ATTT T++++ TT++ FYRAATSASAEAS CTT</p> <p>D yak 1 MHCSRSLSAAGAVDATTTLKLTTSATTTSASFYRAATSASAEASVCTT-PATTKSKTGM 61</p> <p>D mel 61 TTTTTTTTGATNAKDGVTMREKKGGALQKLKRLSHSFGRLTISREDGDESTHHHHHH 120 TTTTTTTGATNAKDGVTMREKKGGALQKLKRLSHSFGRLTISREDGDESTHHHHHH</p> <p>D yak 62 ATTTTTTTGATNAKDGVTMREKKGGALQKLKRLSHSFGRLTISREDGDESTHHHHHH 121</p>	

Figure 3–8: Analysis of Eip63E start-gained SNP in *w*¹¹¹⁸; *iso* – 2; *iso* – 3. (A), Location of the start-gained SNP at the Eip63E locus. Notice that the reading frame is the same as the normal translation start site (TSS). (B), Conservation of 60 amino acid N-terminal region of Eip63E in *w*¹¹¹⁸; *iso* – 2; *iso* – 3 with *Drosophila yakuba* orthologous gene. The other sequenced *Drosophila* species do not have this N-terminal sequence (not shown).

stop gained	location	length	phenotype	stop gained	location	length	phenotype
ade3	255K/*	435	ND ^a	ex	693Q/*	1428	Lethal ^d
CG10126	11W/*	228	ND ^a	lbk	1130Y/*	1174	ND ^d
CG15394	120Q/*	186	ND ^a	MESR4	1509E/*	2072	lethal ^d
CG31145	27L/*	764	ND ^a	mew	752Q/*	1050	Lethal ^a
CG31784	1049Q/*	1078	ND ^a	NFAT	12G/*	1420	ND ^d
CG32115	468W/*	476	ND ^a	oc/otd	389Y/*, 453Y/*	543	Lethal ^d
LRP1	2917Y/*, 2918E/*	4700	ND ^a	Pde9	255C/*	1527	ND ^a
CG34006	121R/*	202	ND ^b	rho-4	140W/*	418	ND ^a
CG34326	49Y/*	84	ND ^c	Synd	375S/*	495	ND ^a
CG3493	1419E/*	1490	ND ^a	tai	1420Q/*	2048	Lethal ^d
CG3964	509Y/*	983	ND ^a	trol	811Y/*, 808E/*	4180	Lethal ^a
CG4068	379Q/*	623	ND ^d	stop lost			
CG7236	70E/*	502	ND ^a	GABA-B-R1	*/L (+9 aa)	837	ND ^a
Cht6	4175L/*	4542	ND ^a	CG13958	*/G (+8 aa)	539	ND ^a
Cyp4s3	260W/*	496	ND ^a	CG4975	*/Q (+1aa)	353	ND ^a
Dif	263C/*	668	lethal ^a	bw	*/Q (+71 aa)	417	eye color ^c
Dp	17353L/*	22972	Lethal ^a	CG14755/pdm3	*/Q (+5 aa)	285	ND ^a

Table 3–9: Stop gained and stop lost in *w*¹¹¹⁸; *iso* – 2; *iso* – 3

species have extended C-terminal amino acids and at least three Drosophila species have missing amino acids at the C-termini (Fig. 3–13). For example, *Drosophila pseudoobscura* has three of the extended amino acids found in *w¹¹¹⁸; iso – 2; iso – 3* and *Drosophila mojavensis* has four of them. In contrast, *Drosophila simulans* is missing the last terminal amino acid, *Drosophila erecta* is missing the last two terminal amino acids, and *Drosophila yakuba* is missing the last three amino acids found in the reference strain (Fig. 3–13). The large number of stop-gain and stop-lost SNPs in *Drosophila* likely has important implications on the evolution of protein function (see Discussion).

Synonymous and non-synonymous SNPs in *w¹¹¹⁸; iso – 2; iso – 3*. There are 15,842 synonymous SNPs and 4,467 nonsynonymous SNPs in annotated coding regions in *w¹¹¹⁸; iso – 2; iso – 3* (Fig. 3–5). A synonymous SNP (silent SNP) is defined as a SNP that does not change the amino acid in the protein, whereas a nonsynonymous SNP does. The genome-wide normalized *N/S* ratio (*dN/dS*), also called ω (i.e., $\omega = dN/dS$), is by definition normalized to 1 in most evolutionary studies [146]. The non-normalized *N/S* ratio is $\tilde{0.28}$ in *w¹¹¹⁸; iso – 2; iso – 3* compared with the reference genome, *y¹; cn¹bw¹sp¹* (i.e., *N/S* = 4,467/15,842; Table 3–1).

We examined the distribution of synonymous and nonsynonymous SNPs genome-wide for *w¹¹¹⁸; iso – 2; iso – 3* and saw higher levels of both classes of SNPs in the middle of the chromosome arms and lower levels near the centromeres and telomeres (Fig. 3–12 and left). This was expected because the number of SNPs is proportional to the recombination frequencies in the different regions of the chromosomes [14, 26]. Also, our previous analyses of the distribution of total SNPs revealed a similar pattern. We observed higher *N/S* ratios near the telomeres and centromeres and lower *N/S* ratios in the middle of the chromosome arms (Fig. 3–12 and right).

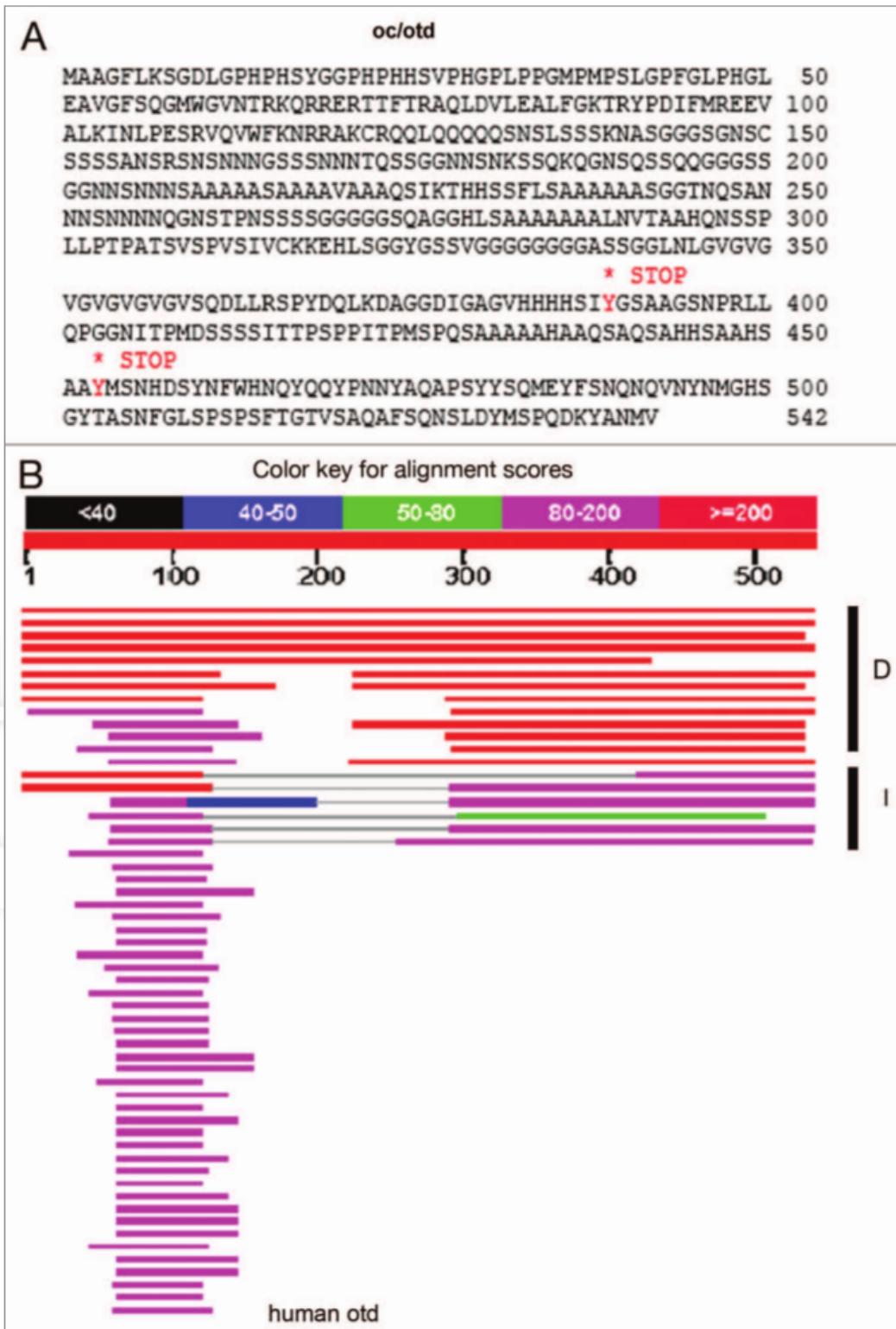


Figure 3–10: Oc/Otd has two stop-gained SNPs in *w¹¹¹⁸*; *iso* – 2; *iso* – 3. (A) Location of the two stop gained SNPs in *oc/otd*. (B) Protein BLAST of Oc/Otd against the non-redundant (nr) protein database shows that only the 60 amino Hox domain flanking amino acid 100 is conserved from *Drosophila* to humans. The color coding shows the alignment scores.

3.5 Discussion

In this paper, we used SnpEff to categorize the ~356,660 SNPs in *w¹¹¹⁸; iso-2; iso-3* and place them into 14 different classes based on their predicted effects on protein function. In order of prevalence, these 14 classes are intron, upstream, downstream, intergenic, synonymous, non-synonymous, 3' UTR, 5' UTR, start-gained, stop-gained, stop-lost, synonymous-stop, start-lost and splice-site SNPs (Fig. 3-5). The reason for cataloging the SNPs in *w¹¹¹⁸; iso2; iso-3* is to get a better appreciation of evolution of genome sequences and genome organization in this common laboratory strain. We appreciate the fact that both *w¹¹¹⁸; iso-2; iso-3* and *y¹; cn¹bw¹sp¹* are derived and highly manipulated laboratory strains and do not represent natural populations. Therefore, we do not mean to imply that the analyses in this paper are significant but rather just observational. To be meaningful, these observations need to be followed up with natural populations. Hundreds of Drosophila natural populations have already been or are in the process of being sequenced, so this should be feasible in the near future with a program such as SnpEff [9].

Many of the stop-gained and stoplost SNPs in *w¹¹¹⁸; iso-2; iso-3* occur in essential genes that apparently still function after amino acid truncations caused by the stop-gained SNPs (Table 3-7). These non-critical effects of the stop-gained SNPs are worth noting because nonsense codons in the transcribed mRNAs generally result in nonfunctional protein products. For example, some genetic disorders, such as thalassemia and Duchenne muscular dystrophy (DMD), result from nonsense SNPs [57, 151, 25]. Also, nonsense SNP-mediated RNA decay exists in yeast, Drosophila and humans, and usually ensures that mRNAs with premature stop codons are degraded [61].

The stop-gained and stop-lost SNPs in essential genes, if they are validated, could have profound evolutionary implications and suggest the involvement of prions, analogous to [PSI^+], in the retention and selection of these SNPs. Brian Cox, a geneticist working with the yeast *Saccharomyces cerevisiae*, discovered [PSI^+] in 1965 as a non-genetically transmissible trait with a cytoplasmic pattern of inheritance similar to mitochondria [44]. He isolated a yeast strain auxotrophic for adenine due to a nonsense mutation is able to survive in media lacking adenine when [PSI^+] is present [44]. Reed Wickner showed in 1994 that [PSI^+] resulted from a prion form of the translation termination factor, Sup35 [164]. Lindquist and colleagues showed in 2008 that the [PSI^+] prion provides survival advantages in several stressful environments, such as high salt conditions [153]. They have speculated that Sup35 is an evolutionary capacitor that, when inactivated in the PSI^+ form, releases cryptic genetic variation that allow expression of novel C-terminal amino acids in hundreds of proteins, some of which are beneficial in stressful environments.

How might prions be involved in revealing cryptic genetic variation in the 5' and 3' UTRs? While most prions are thought to not directly mutate DNA sequences, they could provide an environment that would make the retention and selection of beneficial SNPs more likely. For example, a stop-lost SNP would allow a modified protein with the new C-terminal tail to be always expressed, even when the prion is lost [153]. Therefore, a stop-lost SNP would more likely occur in a strain with beneficial codons in the 3' UTR because the cryptic C-terminal amino acids encoded by these nucleotides would provide a selective advantage in stressful (i.e., [PSI^+]) environments when they are translated.

It is attractive to speculate that a similar prion-mediated evolutionary mechanism might occur in *Drosophila*, for both stoploss and stop-gained SNPs,

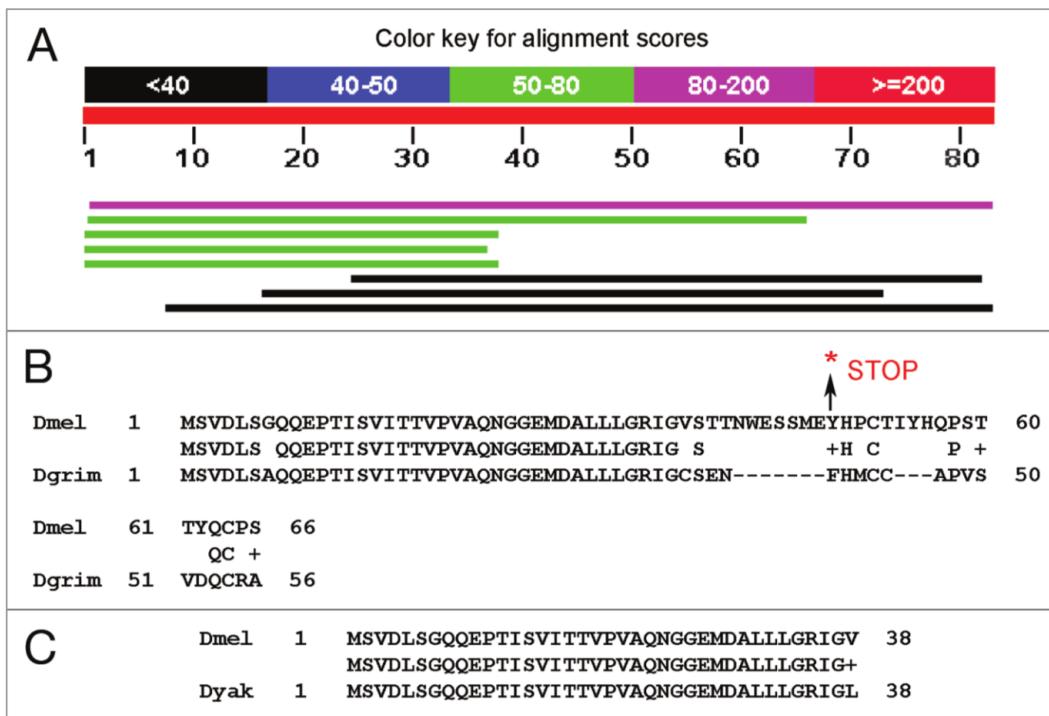


Figure 3–11: CG34326 has one stop-gained SNP in *w¹¹¹⁸*; *iso* – 2; *iso* – 3 in the non-conserved C-terminal region. (A) Protein BLAST of CG34326 against the non-redundant (nr) protein database shows that only the 38 N-terminal amino acids are conserved among Drosophila species and not beyond Drosophila. The colored lines represent the homologs from the following organisms: Drosophila melanogaster, Drosophila grimshawi, Drosophila yakuba, Drosophila erecta, Drosophila virilis, Ixodes scapularis, Ixodes scapularis and Nycticebus coucang. (B) Alignment of Drosophila melanogaster CG34326 with orthologous gene from Drosophila grimshawi. (C) Alignment of Drosophila melanogaster CG34326 with orthologous gene from Drosophila yakuba.

and that this might help explain the large number of SNPs that we see in these categories. We note that *Drosophila* has several Sup35 orthologs, some of which have N-terminal repeats that are known to be potentially prion-forming domains [153]. We acknowledge that this is a highly speculative explanation for the high numbers of start-gained and stop-lost SNPs, but we believe that it is worthy of further investigation.

The many potential start-gained SNPs in *Drosophila* might also have evolutionary implications. Similar to the cryptic genetic variation that is revealed by stop-lost mutations in the 3' UTR, start-gained SNPs reveal cryptic genetic variation in the 5' UTR. Uncovering the cryptic genetic variation in times of environmental stress, such as by inducing transcription initiation at start sites upstream of the normally-used transcription start sites, could be one mechanism to facilitate the use of potential start-gained SNPs. Further mutations and selection of the potential start-gained SNPs, such as by introducing better Kozak consensus sequences or more commonly used 5'-AUG-3' translation initiation codons, can stabilize the cryptic genetic variation further if it leads to improved survival or reproductive fitness in a stressful environment. While amino acid extensions and deletions in known essential genes occur only 8 times in *w¹¹¹⁸; iso-2; iso-3* compared with the reference strain (Table 3-9), as laboratories begin to sequence hundreds or even thousands of individuals in a population, extensions and deletions are likely to be found in a large proportion of functional genes.

Finally, we recently upgraded SnpEff further by including over 320 databases for different reference genome versions that can be analyzed (<http://snpeff.sourceforge.net/SnpSift.html>). Sources of information for creating these databases are ENSEMBL, UCSC Genome Bioinformatics website as well as organism specific databases, such as FlyBase (*Drosophila melanogaster*), WormBase

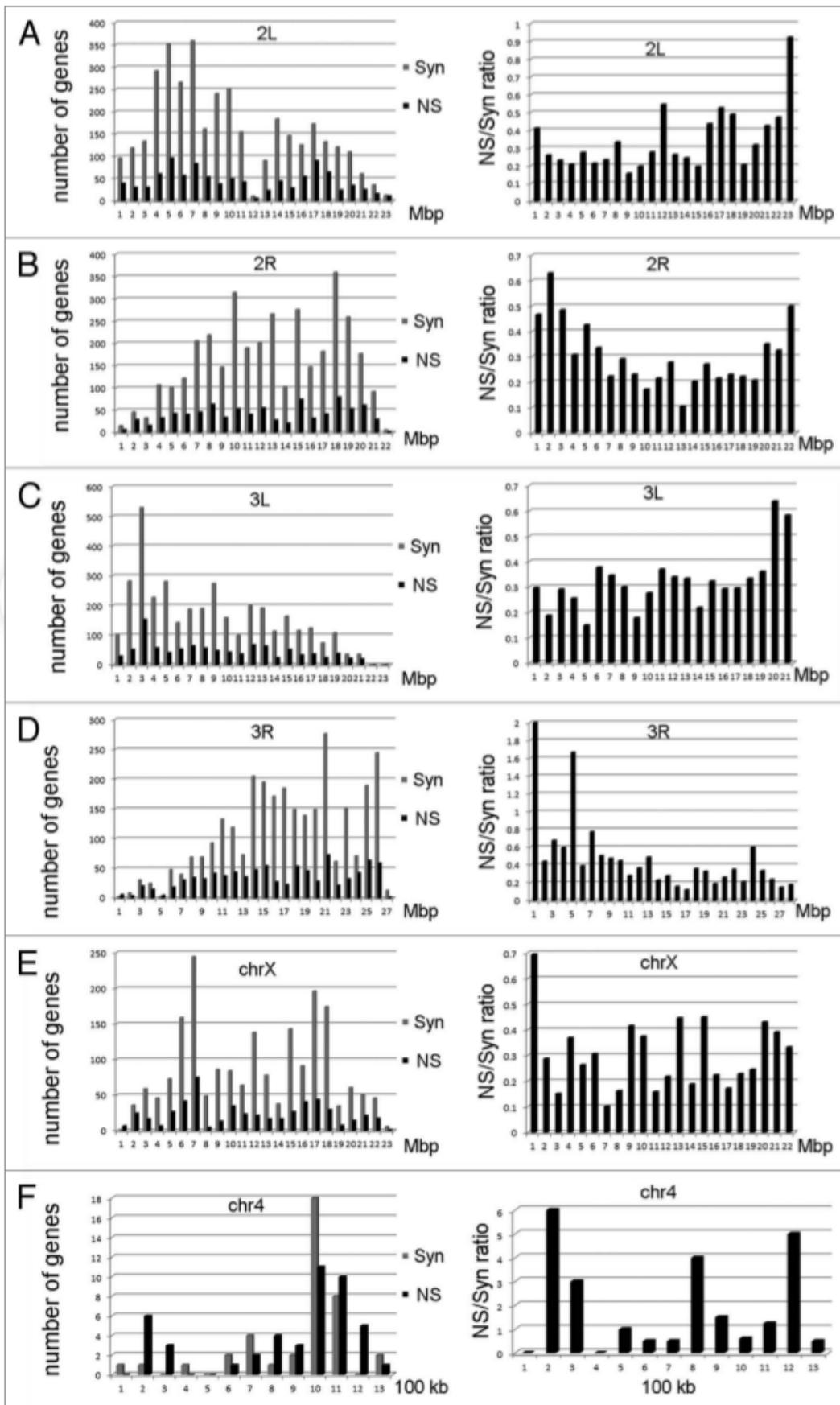


Figure 3–12: Nonsynonymous to synonymous ratios along the chromosome arms in *w*¹¹¹⁸; *iso* – 2; *iso* – 3. (A) Left, Nonsynonymous SNPs at 1 Mbp intervals along the 2L chromosome arm (black) and synonymous SNPs (gray).

(*C. elegans*) and TAIR (*Arabidopsis thaliana*), to name a few. The program SnpEff is open access and additional genomes can be added and assistance in using SnpEff can be provided upon request. Rapid analyses of whole-genome sequencing data should now be feasible to perform by any laboratory

3.6 Methods

SnpEff overview. The program is divided in two main parts (i) database build and (ii) effect calculation. Part (i) Database build is usually not run by the user, because many databases containing genomic annotations are available. Databases are build using a reference genome, a FASTA file, and an annotation file, usually GTF, GFF or RefSeq table, provided by ENSEMBL, UCSC Genome Bioinformatics website or other specific websites, such as FlyBase, WormBase and TAIR. SnpEff databases are gzip serialized objects that represent genomic annotations.

Part (ii) Effect calculations can be performed once the user has downloaded, or built, the database. The program loads the binary database and builds a data structure called “interval forest,” used to perform an efficient interval search (see next section). Input files, usually in VCF format, are parsed and each variant queries the data structures to find intersecting genomic annotations. All intersecting genomic regions are reported and whenever these regions include an exon, the coding effect of the variant is calculated (hence the name of the program). A list of the reported effects and annotations is shown in Table 3–2, additional information produced by the program, is shown in Table 3–3 and Table 3–4, for different output formats.

SnpEff algorithms. In order to be able to process thousands of variants per second, we implemented an efficient data structure that allows for arbitrary interval overlaps. We created an interval forest, which is a hash of interval trees indexed by chromosome. Each interval tree [42] is composed of nodes. Each

node has five elements (i) a center point, (ii) a pointer to a node having all intervals to the left of the center, (iii) a pointer to a node having all intervals to the right of the center, (iv) all intervals overlapping the center point sorted by start position and (v) all intervals overlapping the center point, sorted by end position.

Querying an interval tree requires $O(\log n + m)$ time, where n is the number of intervals in the tree and m is the number of intervals in the result. Having a hash of trees, optimizes the search by reducing the number of intervals per tree.

In order to create this the interval forest, genomic information can be parsed from three main annotation formats: GTF (version 2.2), GFF (versions 3 and 2), UCSC Genome Bioinformatics website RefSeqTables and tab separated text files (TXT). Once the interval forest is created, the structure is serialized and compressed (GZIP) into a binary database. There are over 250 genomic binary databases that are currently distributed with SnpEff, which include all genomes from ENSEMBL.

SnpEff accuracy. As part of our standard development cycle, we perform accuracy testing by comparing SnpEff to ENSEMBL “Variant effect predictor,” which we consider it is the “gold standard.” Current unity testing includes over a hundred test cases with thousands of variants each to ensure predictions are accurate.

SnpEff integration. SnpEff provides integration with third party tools, such as Galaxy [62], which creates a web based interface for bioinformatic analysis pipelines. Integration with Genome analysis tool kit 4 (GATK) was provided by the GATK team. Detailed information on how to download, install and run, as well as usage examples of the program, can be found at <http://snpEff.sourceforge.net>.

Drosophila melanogaster

Dm-ref	522	LVLQQCDSVQGYMEVSL*	538	
		LVLQQCDSVQGYMEVSL*		+8
Dm-w ¹¹¹⁸	522	LVLQQCDSVQGYMEVSLQIFNNINI*	546	

Drosophila simulans

Dm-ref	522	LVLQQCDSVQGYMEVS	537	
		LVLQQCDSVQGYMEVS		-1
Sbjct	522	LVLQQCDSVQGYMEVS	537	

Drosophila erecta

Dm-ref	522	LVLQQCDSVQGYMEV	536	
		LVLQQCDSVQGYMEV		-2
Sbjct	522	LVLQQCDSVQGYMEV	536	

Drosophila yakuba

Dm-ref	481	LVLQQCDSVQGYME	535	
		LVLQQCDSVQGYME		-3
Sbjct	481	LVLQQCDSVQGYME	535	

Drosophila mojavensis

Query	522	LVLQQCDSVQGYMEVS-LQIF	541	
		LVLQQCDSVQGY+EV L+IF		+3
Sbjct	517	LVLQQCDSVQGYIEVRYLKIF	537	

Drosophila pseudoobscura pseudoobscura

Query	522	LVLQQCDSVQGYMEVSLQIFN	542	
		LVLQQCDSVQGY+EV +F+		+4
Sbjct	571	LVLQQCDSVQGYIEVFCALFH	591	

Figure 3–13: CG13958 has a stop lost SNP in w^{1118} ; iso – 2; iso – 3. The top comparison shows the alignment of the Drosophila melanogaster reference genome with w^{1118} ; iso – 2; iso – 3. Notice that the stop lost causes an extension of 9 amino acids. The second through sixth comparisons shows the alignment of Drosophila simulans, Drosophila erecta, Drosophila yakuba, Drosophila mojavensis and Drosophila pseudoobscura pseudoobscura (Sbjct) with the Drosophila melanogaster reference genome (Dm-ref). The number of terminal amino acids missing or gained is shown (-1 to +3).

Data access. SnpEff Data can be accessed from the Supplemental data file for w^{1118} ; $iso - 2$; $iso - 3$ or by contacting D.M.R.

Disclosure of Potential Conflicts of Interest No potential conflicts of interest were disclosed.

3.7 Acknowledgements

This work was supported by a Michigan Core Technology grant from the State of Michigan's 21st Century Fund Program to the Wayne State University Applied Genomics Technology Center. This work was also supported by the Environmental Health Sciences Center in Molecular and Cellular Toxicology with Human Applications Grant P30 ES06639 at Wayne

State University, NIH R01 grants (ES012933) to D.M.R. and DK071073 to X.L. We thank David Roazen, Eric Banks and Mark DePristo in the GATK team at the Broad Institute who integrated SnpEff with the Genome Analysis Toolkit (GATK).

Note Supplemental material can be found at: <http://www.landesbioscience.com/journals/fly/article/19695>

3.8 Epilogue

At the beginning of my Ph.D., functional annotation of genomic variants was an unsolved problem with many research labs creating in-house custom solutions that oftentimes were inefficient and lacking of rigorous testing. As a consequence, shortly after SnpEff & SnpSift were released they quickly became widely adopted by the research community as well as many private organizations. Currently SnpEff & SnpSift has over 250 downloads per week (as reported by SourceForge, where the tools are hosted). So far SnpEff & SnpSift have been cited over 400 times.

3.8.1 Data structures for annotations

A very simple approach used by ANNOVAR [157] is to create an index by dividing each chromosome into N bins of equal size. All genomic features are stored in a hash table indexed by chromosome name and bin number. This approach has running time of $O(n)$ where n is the number of features, but it can be easily tuned by creating small bins, at the cost of increased memory requirements.

Another approach [30] is to use an “interval forest”, which is a hash of “interval trees” indexed by chromosome. Each interval tree is composed of nodes. Each node has five elements i) a center point, ii) a pointer to a node having all intervals to the left of the center, iii) a pointer to a node having all intervals to the right of the center, iv) all intervals overlapping the center point sorted by start position, and v) all interval overlapping the center point sorted by end position. Querying an interval tree requires $O[\log(n) + m]$ time, where n is the number of features in the tree and m is the number of features in the result. Having a hash of trees optimizes the search by reducing the number of intervals per tree.

CHAPTER 4

Epistatic GWAS analysis

4.1 Preface

Although this thesis focusses on the development of computational approaches that could be applied to the study of a number of complex diseases, our focus has been on type II diabetes mellitus (T2D), a complex disease first described by the Egyptians in 1500 BCE. Later the Greeks in 230 BCE used the term “diabetes” meaning “pass through” (or “siphon”) denoting the constant thirst and frequent urination of the patients. In the 1700s the term “mellitus” (from honey) was added to denote that the urine was sweet and would “attracts ants”.

Diabetes symptoms include frequent urination, thirst, and constant hunger, high blood sugar (hyperglycemia) and insulin resistance. Long term complication from T2D may include eyesight problems, heart disease, strokes and kidney failure. Type II diabetes, is highly correlated with obesity and disease rate has increased dramatically during the last 50 years. According to the World Health Organisation the prevalence of diabetes is 9% in adults and an estimated 1.5 millions deaths were caused by diabetes in 2012 [68], which is predicted to be the 7th leading cause of death by 2030. The costs associated to treating diabetes patients only in the U.S. are estimated around \$245 billion dollars.

In recent years, over 80 genetic loci related to T2D have been identified [113, 37]. Nevertheless, the overall effect sizes of these loci account for less than 10% of the overall disease predisposition [107]. This poses the question of why, given that so much efforts has been directed at finding the genetic components

of this disease, the loci found so far have such modest effects. This lack of large genetic effects, known as the “missing heritability” problem, does not only arise in T2D but also in almost all complex traits. In recent studies about missing heritability [174, 175] it was proposed that this effect might be partly explained by taking into account epistasis (i.e. gene interactions).

In this chapter, we propose a novel framework that takes into account putative epistatic interactions into genome wide association studies (GWAS). We detect putative interactions by means of a co-evolutionary model that we develop based on Markov evolutionary theory. We fit the co-evolutionary model’s parameters using protein structural data (Pdb) an genome wide multiple sequence alignments (UCSC’s 100-way). Once this model is parametrized, we perform genome wide association using of pairs of variants from a sequencing study by incorporating the co-evolutionary model as priors in a Bayesian Framework.

The rest of the chapter is published in: **P. Cingolani**, R. Sladek, M. Blanchette, “A co-evolutionary approach for detecting epistatic interactions in genome-wide association studies”

4.2 Abstract

Motivation. Since proteins perform their functions mainly by interacting with other proteins, epistasis genome wide association studies are assumed to be key in unveiling genetic risk of disease. Due to their high complexity and sometimes prohibitive computational requirements, epistatic GWAS have often been neglected. In this paper, we propose a novel methodology for analysing putative epistatic interactions by combining multiple genome alignments and sequencing information.

Results. We propose a methodology to perform genome wide association on pairs of variants whose putative epistatic interaction might have a

phenotypic effect. Using Pdb information and genome wide multiple species alignment we create a statistical coevolutionary model that increases priors of putative interacting sites, these priors are then used as basis for genome wide association in a Bayesian framework. Our optimized algorithms can be applied to genome wide scale sequencing studies for tens of thousands of samples, that typically yield millions of variants.

4.3 Introduction

Genetic studies aim to discover how a phenotype of interest, such as disease risk or height, is affected by genetic background. Genome wide association studies (GWAS) are powerful techniques aimed to find statistical association between phenotype and genetic variants (mutations or polymorphisms) [33]. Although several genetic variants related to different phenotypes have been found, variants discovered in GWAS so far can only explain a small part for the phenotypic inheritance. For instance, all genetic variants associated to height collectively account for few centimetres in the offspring’s height [165], similarly all combined variants related to type 2 diabetes risk only explain 5% to 10% of the overall variance in disease predisposition [113, 37]. This problem is known as “missing heritability” [107] and recent theories suggest that genetic interactions (epistasis) might play an important role in it [174, 175].

The foundations for epistasis, broadly defined as “genetic interaction” [59], have been proposed almost a hundred years ago by Bateson (1909) and Fisher (1918). It was the latter who coined the term to denote a “statistical deviation of multilocus genotype values from an additive linear model for the value of a phenotype” [59]. There is evidence of such interactions being involved in complex diseases. For instance an interaction between BACE1 and APOE4 having a significant association with Alzheimer’s disease has consistently been replicated in different studies [35]. One of the main problems in

finding association between interactions and disease is that out of the whole set of molecular interactions (the interactome) only a small part of it has been characterized [154].

Interacting proteins can be identified experimentally through several types of approaches (Yeast-Two-Hybrid, Protein fragment complementation assay, Glutathione-s-transferase, Affinity purification coupled to mass spectrometry, Tandem affinity purification, etc. [141]) and large databases of protein-protein interactions are now available for human [145, 141]. Nevertheless, these methods predict the interaction between proteins and may not discern the exact residues mediating such interactions Furthermore, it is estimated that up to 80% large of the human protein-protein interactions remains unknown [154].

These issues can be addressed using computational predictions of either pairs of interacting proteins or interacting residues [142]. A type of approaches that has been gaining popularity recently is one that makes use of the plethora of genomic sequences available for species other than human in order to discover evolutionary evidence of selective pressure on individual residues to identify interacting sites and interfaces [108]. Interacting residues and their neighbors may then be subject to compensating epistasis, where a mutation at a residue in one protein may be compensated by another mutation at a residue in the second protein [121]. Assuming that evolutionary pressure acts on both interaction sites simultaneously, co-occurring compensatory mutations can become fixed in the population with higher probability than non-compensatory ones. In light of this hypothesis, we can use statistical methods on multiple sequence alignments (M_{SA}) of proteins from different organisms to find coevolving sites. This types of approaches has been used to identify coevolving sites both within a protein (e.g. N-terminal and C-terminal domains in PKG protein [65] or the GroES-L chaperoning system [132], α and β haemoglobin

subunits [121]), and between interacting proteins (e.g. G-protein coupled receptors and protein ligands [65]).

Many methods exist to find putative interaction loci, both within and across proteins, based on evolutionary evidence (see [49] for a review). One of the simplest methods for inferring co-evolution uses mutual information (MI) between two loci [108] in a multiple sequence alignment (M_{sa}). However, methods based on correlation or mutual information are known to have systematic biases due to phylogeny [49], or sequence heterogeneity problems [161]. More sophisticated methods, such as DCA [112], PSICOV [77] or mdMI [32] try to overcome these biases, however they are usually not suitable for GWAS-scale analysis for two main reasons: i) they require multiple alignments of a very large number of sequences (from 400 to $25L$, where L is the length of the protein [32]), such depth is not usually available at whole genome scale; and ii) they are computationally demanding (e.g. running for minutes or even days for each interacting pair of proteins being considered), making them unsuitable for analyses involving millions of variants spanning over thousands of proteins. Furthermore, a recent study shows that overall agreement between methods is not high (65% or less) and predictive power is quite low (only 6% of the “top scoring pairs” are real interactions) [32].

Applying epistatic interaction models to GWAS studies is a challenging problem for several reasons: i) interaction models are by definition non-linear [59]; ii) analyzing all order N variant combinations in a sequencing experiment requires great computational power and efficient algorithms because the number tests grows exponentially with N [122]; iii) multiple testing correction can render association tests underpowered for all but very large cohorts [59, 122]; and iv) there is no consensus of what genetic interaction means [106], which is reflected in the difficulty to find a unified model [122, 106]. For all these

reasons and due to the lack of sequencing cohorts large enough to detect these interactions, the application of epistatic models to sequencing studies has not been widespread. Furthermore, there is no clear consensus on the required sample size to detect epistatic interactions. Depending on phenotypic effect size and variant’s allele frequency some estimates assume in the order of 10,000 to 500,000 cases [80] to be required. Such cohorts are now becoming feasible due to improvements and cost reductions in sequencing technology.

Approaches for epistatic GWAS do exist and they apply a wide array of methodologies, although we don’t intend write a comprehensive review, we mention a few examples. In [172] the authors infer probabilities by noting that interactions create linkage disequilibrium patterns in the disease population. A “Bayesian partitioning model” is applied in [171]. In [2] the authors look for over / under-represented allele pairs in a given population by performing an analysis of imbalanced allele pair frequencies. Finally, finding interacting variants can be viewed as an optimisation / attribute selection procedure, thus many machine learning methodologies have been proposed [111].

In this work we propose methodologies to prioritize pairs of variants from sequencing experiments by combining genome wide association with epistatic interaction models. In a nutshell, our method uses recently computed 100-way vertebrate genome alignments to calculate interaction posterior probabilities for any given pair of residues in human proteins. This is achieved by contrasting the likelihood of the observed pair of alignment columns under a joint substitution model that models dependencies between interacting sites, and a null model of independent evolution. These posterior probabilities are then used as priors to modulate the evidence of epistatic interaction derived from GWAS data. We engineer efficient algorithms that can be applied to GWAS-scale datasets of tens of thousands of samples.

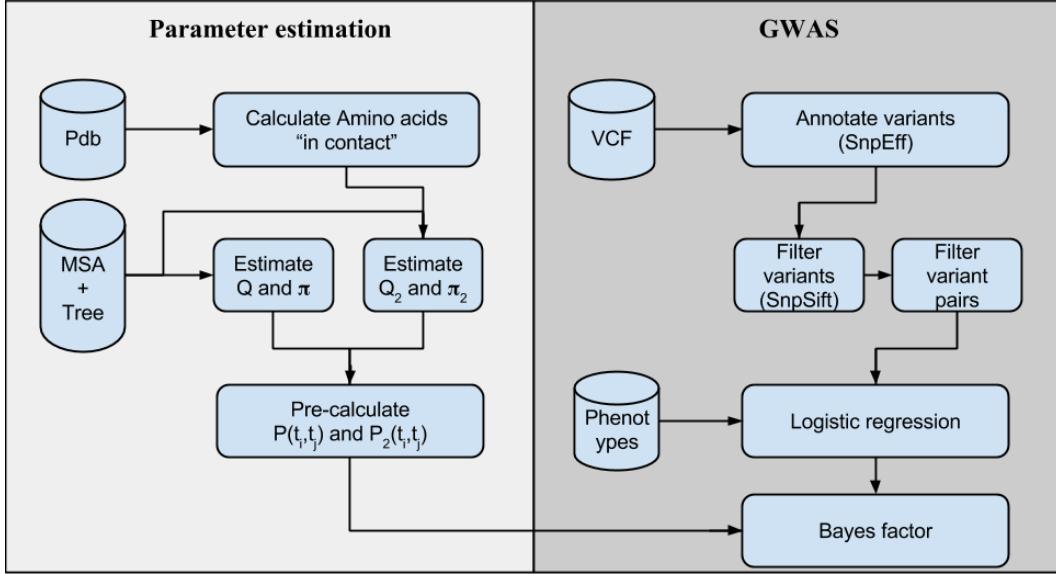


Figure 4–1: Complete pipeline example

Finally we apply this methods to a cohort of $\sim 13,000$ individuals in a case-control study of type II diabetes (for study details, see [109]). We are able to find suggestive association of putatively epistatic interactions.

4.4 Methods

Our epistatic GWAS analysis pipeline involves three key steps, as shown in Figure 4–1. First, we learn a co-evolutionary substitution rate matrix for pairs of amino acids that are in contact in proteins. Second, we analyze a GWAS data set to identify pairs of non-synonymous SNPs that show (possibly weak) evidence of epistasis. Third, for each pair of SNP identified in step 2, we measure the evidence of co-evolution of the pair of encoded amino acids, and combine it with the GWAS evidence by adding up the corresponding Bayes factors.

4.4.1 Substitution model for pairs of interacting amino acids

In this section, we describe how we estimate two substitution rate matrices. The first is the usual 20×20 substitution rate matrix Q describing the

evolution of individual amino acids. The second, Q_2 is a 400×400 substitution rate matrix for pairs of interacting residues.

We used the 100-way vertebrate multiple sequence alignment and accompanying phylogenetic tree T available from the UCSC Genome Browser [81]. This alignment includes the DNA sequences of 100 species whose genome is completely or nearly completely sequenced, with 6 Afrotheria, 14 Avians, 14 Euarchontoglires, 16 Fish, 25 Laurasiatheria, 5 Mammilians, 12 Primates and 8 Sarcopterygii. The multiple alignment is performed using “multiz” algorithm [17, 84].

From the $\sim 21,000$ human protein structures (resolution less than 3 \AA) available in Protein Data Bank, we extracted a set of $\sim 770,000$ pairs of “within protein amino acids interactions” residues, defined as pairs of residues from the same protein where at least one pair of atoms is within 3 \AA or less. Similarly, from the set of $\sim 5,700$ models of co-crystallized complexes in PDB, we extracted a set of $\sim 12,000$ pairs of “protein-protein interacting” residues, defined as amino acids from different proteins that satisfy the same distance criterion.

To derive rate matrix Q , we consider the complete set of $n \sim 22M$ protein coding sites present in the alignment (that unambiguously map to the genome), irrespective of the presence or absence of contacts. Q is obtained following classical sequence evolution theory ([169] [55]). First, for each pair of species s_i and s_j , we obtain $c_i(a)$ defined as the count of amino acid a in species s_i , and $c_{i,j}(a, b)$ defined as the number of sites that have had a transition from amino acid a in s_i to b in s_j . Stationary probability of amino acid a in genome s_i is then defined as $\pi_i = c_i(a)/n$. Assuming a time reversible model, we get the frequency of change from a to b : $f_{i,j}(a, b) = f_{j,i}(a, b) = (c_{i,j}(a, b) + c_{j,i}(a, b))/(2n)$. Let $P_{i,j}$ be the amino acid transition probability matrix from

s_i to s_j , i.e. $P_{i,j}(a, b)$ is the probability that s_j has amino acid b given than s_i has amino acid a . Then P is obtained through the relation $f_{i,j}(a, b) = \pi_i(a) \cdot P_{i,j}(a, b)$, or $P_{i,j}(a, b) = f_{i,j}(a, b)/\pi_i(a)$. Let $t_{i,j}$ be the total branch length between s_i and s_j (obtained from the phylogenetic tree). Assuming time reversibility, we have $P_{i,j} = e^{Q \cdot t_{i,j}}$, and thus $Q = \log[P_{i,j}/t_{i,j}]$ [169]. Taking into account the estimation error, the equation becomes $\hat{Q}(t_i + t_j) = Q = \log[P_{i,j}/t_{i,j}] + \epsilon_{i,j}$, where $\epsilon_{i,j}$ is an error matrix. Under the assumption that the mean error is zero, we can approximate the rate matrix by the calculating an average of all estimates:

$$\hat{Q} = \frac{1}{N(N-1)/2} \sum_{i < j} \hat{Q}(t_i + t_j) \quad (4.1)$$

$$= \frac{2}{N(N-1)} \sum_{i < j} \frac{1}{t_i + t_j} \log[\hat{P}(t_i + t_j)] \quad (4.2)$$

The much larger substitution matrix Q_2 describes the substitution rate from any pair of amino acid (a, b) to any other pair (c, d) . It is derived similarly to Q , but considering only pairs of amino acids from the set of within protein interacting pairs of amino acids. We only take into account amino acids pairs within the same chain, that are separated by 20 amino acids or more from high resolution PDB structures of human proteins.

4.4.2 Calculating likelihood of individual and pairs of alignment columns

Given a substitution rate matrix Q , the likelihood $L_1(i)$ of an alignment column i assigning an amino acid to each leaf in the tree T is calculated using the well known Felsenstein algorithm [55]. This is achieved in time $O(N \cdot |\Sigma|^2)$, where $|\Sigma| = 20$. Given matrix Q_2 , the same algorithm can be used to compute the likelihood $L_2(i, j)$ of a pair of alignment columns (i, j) , but now in time $O(N \cdot |\Sigma|^4)$.

A test for co-evolution of two positions i, j of the same or different proteins is obtained using the likelihood ratio under the two models:

$$LR[M_{SA}(i, j)] = \frac{P(i, j|M_{SA}, Q_2)}{P(a_i|M_{SA}, Q) \cdot P(a_j|M_{SA}, Q)}$$

where the denominator assumes that the amino acids i and j evolve independently. Similarly, the log-likelihood is defined as

$$LL[M_{SA}(i, j)] = \log \left[\frac{P(i, j|M_{SA}, Q_2)}{P(a_i|M_{SA}, Q) \cdot P(a_j|M_{SA}, Q)} \right] \quad (4.3)$$

Given a pair of columns in the multiple sequence alignment M_{sa} , the corresponding phylogenetic tree \mathcal{T} , we calculate the likelihood ratio by propagating the probabilities from the leaves up the root of the phylogenetic tree using Felsenstein's algorithm [55] using transition matrices Q and Q_2 for the null and alternative models respectively.

Because the calculations described in this section will need to be performed on a very large number of pairs of sites, optimizations we required to ensure manageable running time. First, pre-calculation of matrix exponentials $P(t) = e^{Qt}$ is necessary for all values of t corresponding to individual branch lengths. Another optimization ("constant-tree caching") is to cache likelihood values for subtrees of the phylogenetic tree where all nodes have the same amino acid values. This optimization results in speedup only if the phylogenetic tree remains constant throughout the genome.

4.4.3 GWAS model

Given N_S samples (individuals), we use the standard notation for phenotypes and code them as $d_s = 1$ meaning that the individual s is affected by disease and $d_s = 0$ if the individual is "healthy". Let $\bar{d} = [d_1, \dots, d_{N_s}]$ be a

phenotype vector and $g_{s,i} \in \{0, 1, 2\}$ a genomic variant for sample s , loci i (we assume there are N_V variants). A logistic model of disease risk [12] is

$$\begin{aligned} p_{s,i} &= P(d_s | g_{s,i}) \\ &= \phi(\theta_0 + \theta_1 g_{s,i} + \theta_2 c_{s,1} + \theta_4 c_{s,2} + \dots) \\ &= \frac{1}{1 + e^{\theta_0 + \theta_1 g_{s,i} + \theta_2 c_{s,1} + \theta_4 c_{s,2} + \dots}} \\ &= \phi(\bar{\theta}^T \bar{g}_{s,i}) \end{aligned}$$

where $\phi(\cdot)$ is the sigmoid function, $c_{s,1}, c_{s,2}, \dots$ are covariates for each individual s (these covariates usually include sex, age and eigenvalues from population structure analysis [125]), $\bar{g}_{s,i} = [1, g_{s,i}, c_{s,1}, c_{s,2}, \dots, c_{s,N_C}]$, and $\bar{\theta} = [\theta_1, \theta_2, \dots, \theta_m]$. The parameter estimates $\bar{\theta}$ are obtained by solving the maximum likelihood equation

$$\begin{aligned} L(\bar{\theta}) &= \prod_{s=1}^{N_S} P(d_s | \bar{\theta}, g_{s,i}) \\ &= \prod_{s=1}^{N_S} p_{s,i}^{d_s} (1 - p_{s,i})^{1-d_s} \\ &= \prod_{s=1}^{N_S} \phi(\bar{\theta}^T \bar{g}_{s,i})^{d_s} (1 - \phi(\bar{\theta}^T \bar{g}_{s,i}))^{1-d_s} \end{aligned}$$

where $p_{s,i} = P(d_s | \bar{\theta}, g_{s,i})$ is the probability of individual s disease outcome, given a genomic variant.

Using this model, we have two hypotheses: i) the null hypothesis, H_0 , assumes that genotype does not influence disease probability (i.e. $\theta_1 = 0$). ii) the alternate hypothesis, H_1 , assumes that the genotype does influences disease probability. We can compare these two hypothesis using a likelihood ratio test, we define

$$LR = \frac{L(\bar{\theta}_{alt}|H_1)}{L(\bar{\theta}_{null}|H_0)} \quad (4.4)$$

$$LL = \log [LR] = \log \left[\frac{L(\bar{\theta}_{alt}|H_1)}{L(\bar{\theta}_{null}|H_0)} \right] \quad (4.5)$$

where $\bar{\theta}_{null}$ and $\bar{\theta}_{alt}$ are the maximum likelihood estimates for null and alt model respectively. According to Wilk's theorem, the log likelihood ratio has a χ^2_1 distribution, so we can easily calculate a p-value.

Next, we extend the logistic model to accommodate interacting loci. For an individual (sample s), we model interactions between two genetic loci i and j , having genotypes $g_{s,i}$ and $g_{s,j}$, by extending the logistic model with N_{cov} covariates $c_{s,j}$ as

$$P(d_s|g_{s,i}, g_{s,j}, H_1) = \phi[\theta_0 + \theta_1 g_{s,i} + \theta_2 g_{s,j} + \theta_3(g_{s,i}g_{s,j}) \quad (4.6)$$

$$\dots + \theta_4 c_{s,1} + \dots + \theta_m c_{s,N_{cov}}] \quad (4.7)$$

$$= \phi(\bar{\theta}^T \bar{g}_{s,i,j})) \quad (4.8)$$

where $\bar{g}_{s,i,j} = [1, g_{s,i}, g_{s,j}, (g_{s,i}g_{s,j}), c_{s,1}, c_{s,2}, \dots, c_{s,N_{cov}}]^T$. An implicit assumption in this equation is that $g_{s,i}$ and $g_{s,j}$ are not correlated (e.g. they don't sit in the same LD-Block). This can be enforced either by using haplotype structure information (e.g. from HapMap) or by limiting the application of the model to variants either in different chromosomes or over 1M bases apart. The null hypothesis H_0 assumes that variants act independently

$$P(d_s|g_{s,i}, g_{s,j}, H_0) = \phi[\theta'_0 + \theta_1 g_{s,i} + \theta_2 g_{s,j} + \theta_3 c_{s,1} + \dots] \quad (4.9)$$

$$= \phi(\bar{\theta}^T \bar{g}'_{s,i,j}) \quad (4.10)$$

where $\bar{g}'_{s,i,j} = [1, g_{s,i}, g_{s,j}, c_{s,1}, c_{s,2}, \dots, c_{s,N_{cov}}]^T$. Since these models are not nested (i.e. H_0 is not included in H_1), we cannot directly apply Wilks theorem.

Logistic regression parameter fitting was implemented using several different algorithms to reach the desired performance. The fastest convergence is obtained using Iterative Reweighted Least Squares (IRWLS [47]) and BroydenFletcherGoldfarbShanno algorithm (BFGS [21]) with some code optimizations. In most cases IRWLS converges faster, so it was selected as the default implementation in our analysis.

Another way to compare the null hypothesis to the alternative hypothesis, is using a Bayesian formulation [82, 156]

$$\begin{aligned} P(H_1|\mathcal{D}) &= \frac{P(\mathcal{D}|H_1)P(H_1)}{P(\mathcal{D})} = \frac{\int P(\mathcal{D}|\theta, H_1)P(\theta|H_1)d\theta}{P(\mathcal{D})} \\ \Rightarrow \frac{P(H_1|\mathcal{D})}{P(H_0|\mathcal{D})} &= \frac{\int P(\mathcal{D}|\theta, H_1)P(\theta|H_1)d\theta}{\int P(\mathcal{D}|\theta, H_0)P(\theta|H_0)d\theta} \frac{P(H_1)}{P(H_0)} = B_F \frac{P(H_1)}{P(H_0)} \end{aligned}$$

where B_F , the ratio of the two integrals, is the Bayes factor. Using a Bayesian formulation has two main advantages: i) the hypothesis are automatically corrected for model complexity since Bayes factor asymptotically converge to Bayesian Information Criteria (BIC) is implicitly within Bayes Factors [82], and ii) we can compare non-nested models. The Bayes factor for the epistatic model becomes:

$$B_F = \frac{\int \prod_{s=1}^{N_S} \phi(\bar{\theta}^T \bar{g}_s)^{d_s} [1 - \phi(\bar{\theta}^T \bar{g}_s)]^{1-d_s} P(\bar{\theta}|H_1) d\theta}{\int \int \prod_{s=1}^{N_S} \phi(\bar{\theta}^T \bar{g}_{s,i}) \phi(\bar{\theta}^T \bar{g}_{s,j})^{d_s} [1 - \phi(\bar{\theta}^T \bar{g}_{s,i}) \phi(\bar{\theta}^T \bar{g}_{s,j})]^{1-d_s} P(\bar{\theta}|H_0) P(\bar{\theta}'|H_0) d\theta d\theta'} \quad (4.11)$$

Unfortunately, calculating Bayes factors is not trivial and most of the times there are no closed form equations. Calculating the integrals using numerical algorithms is possible, but imposes a significant computational burden thus making it impractical for large datasets, such as GWAS data, even using

large computing clusters. We can approximate the integrals using Laplace's method [82]. If $g(x)$ has a minimum at x_0 , it can be shown that

$$\int e^{\lambda g(x)} h(x) dx \simeq h(x_0) e^{\lambda g(x_0)} \sqrt{\frac{2\pi}{\lambda g''(x_0)}}$$

The multivariate case, for $\bar{x} \in \Re^d$, is analogous, we just need a Hessian matrix instead of a second derivate of $g(\cdot)$

$$\int e^{\lambda g(\bar{x})} h(\bar{x}) d\bar{x} \simeq h(\bar{x}_0) e^{\lambda g(\bar{x}_0)} \left(\frac{2\pi}{\lambda} \right)^{d/2} \left[\frac{\partial^2 g(\bar{x})}{\partial \bar{x} \partial \bar{x}^T} \right]^{-1/2} \quad (4.12)$$

Using equation 4.12 we can try to approximate the complex integrals in equation 4.11 by the transformation $L(\bar{\theta}) = e^{\ell(\bar{\theta})}$, where $\ell(\cdot)$ is the log-likelihood of the data. So, we can use Laplace approximation by using Eq.4.12, at the point of the maximum likelihood. In order to do so, we need to calculate the Hessian matrix in Eq.4.12. Fortunately, for logistic models, we can make a few simplifications. Considering that $L(\bar{\theta}) = \prod_{s=1}^{N_s} \phi(\bar{\theta}^T \bar{g}_s)^{d_s} [1 - \phi(\bar{\theta}^T \bar{g}_s)]^{1-d_s}$, it can be shown that for genotype terms

$$\frac{\partial^2 \ell(\bar{\theta})}{\partial \theta_i \partial \theta_j} = \sum_s g_{s,i} g_{s,j} p_s (1 - p_s)$$

Using analogous derivation for the covariates, we can find an analytic form of the Hessian, which completes the Laplace approximation formula.

4.4.4 Putting it all together

In summary, we first calculate the transitions matrices for the Markov models (Q and Q_2) based on observations from protein structures (Pdb) and multiple sequence alignments (UCSC's 100-way). We analyze variants from

genome sequencing data first by filtering only for non-synonymous variants, then analyzing all possible pairs of variants and filtering out those that are unsuitable for further analysis (e.g. in linear dependence, Hardy Weinberg p-value less than 10^{-6} , etc.). From the pairs of variants that pass filtering, we fit two logistic regression models (null and alternative hypothesis), then calculate a p-value using the log-likelihood ratio, and keeping pairs of variants having p-values below a predefined threshold (10^{-6}). These pairs of variants are then analyzed under our co-evolutionary model, we find the corresponding columns in the multiple sequence alignment and calculate the likelihoods for the null and alternative models by means of Felsenstein’s algorithm (using matrices Q and Q_2 in respectively). Finally, likelihoods from co-evolutionary model and likelihoods from logistic regression models are incorporated into a Bayes Factor equation, which is calculated using Laplace’s approximation.

4.5 Results

Our approach, which is summarized in Figure ??, involves three main components. First we estimate evolutionary substitution rates for individual amino acids in a protein as well as for pairs of amino acids (either from the same protein or not) that are physically interacting. Given a set of multiple sequence alignment of protein sequences, these evolutionary models can be used to calculate the likelihood of interaction between any two given amino acids, without the need for any structural information. Second, a statistical test for epistasis is developed to identify pairs of non-synonymous SNPs that show (often weak) evidence of interaction in the way they associate to a given trait. Finally, information from the co-evolution component is combined with that from the epistasis component to give more power to the epistasis test.

4.5.1 Co-evolutionary substitution models

The approach described in Methods was used to obtain substitution rate matrix Q for individual amino acids and Q_2 for pairs of physically interacting residues within the same protein. Unsurprisingly, Q (or more precisely a transition matrix $P(t)$ obtained from Q) is very similar to well known transitions matrixes such as PAM [] (Supplementary Figure ?? and Table ??).

Estimating Q_2 requires information about amino acids that are known to be “interacting”. A pair of amino acids is considered to be “interacting” if any pair of atoms (one from each amino acid) has a distance of 3 Å or less [23].

The structure of Q_2 , which describes substitution rates between one pair of interacting amino acids to another, is richer (Supplementary Figure ?? and supplementary file ??). Of particular interest are the pairs of pairs of amino acids for which the ratio $R(ab, cd) = Q_2(ab, cd)/(Q(a, c) \cdot Q(b, d))$ is large. Those substitution pairs are the ones that are most strongly indicative of an interaction. Figure ?? shows that the number of pairs for which R deviates significantly from 1 is quite large, arguing that interacting sites have co-evolutionary rates that differ from the bulk of non-interacting sites.

For example, the case with the highest rate ratio is [V.I → W.W]’ (i.e. amino acid ‘V’ switched to ‘W’ in the one sequence, and amino acid ‘I’ changed to ‘W’ in the other). In fact, the top 10 pairs are transitions to ’W.W’ amino acid pairs. This makes sense considering (i) individual amino acid substitution rates to tryptophan are generally very low, but that (ii) tryptophan pairs are well known β -hairpin stabilizers and are considered as a paradigm for designing stable β -hairpins [135].

Another type of pair transitions with large ratio is the double transitions to a pair of phenylalanine amino acids from a pairs of hydrophobic amino acids (Lysine, Asparagine, Glutamine, Arginine, Aspartic acid and Glutamic acid).

Phenylalanine-Phenylalanine interaction pairs are assumed to conform $\pi - \pi$ interactions which are predicted and experimentally observed to be favorable [75].

4.5.2 Co-evolutionary model validation

We first the ability of our co-evolutionary model to detect interacting sites located within the same protein by computing the likelihood ratio of the evolutionary history of a candidate pair of sites under an co-evolutionary model (Q_2) versus under independence (Q). Although such pairs of sites are unlikely to exhibit evidence of epistasis in GWAS studies (due to linkage), accurate prediction of interacting sites in a given protein are useful for many other purposes, such as protein structure prediction and prediction of the impact of individual mutations. Figure ?? shows interacting sites tend to have higher likelihood ratio scores than non-interacting ones. Although the likelihood ratio score it itself cannot perfectly discriminate between the two classes, only 25.9% of non-interacting pairs have a likelihood ratio above the median likelihood ratio of interacting pairs.

To confirm that an evolutionary model estimated based on pairs of interacting sites from the same protein is useful at predicting pairs of interacting sites between proteins, we repeated the same type of analysis on $\sim 3,000$ pairs of interacting ($< 3 \text{ \AA}$) and $\sim 3,000$ pairs of non-interacting ($> 30 \text{ \AA}$) residues from distinct proteins, obtained from co-crystal structures in PDB (see Methods). As seen on Figure ??, the two classes of sites have substantially different likelihood ratio distributions (Mann-Whitney one sided test: $p - value < 2.2E^{-16}$), although slightly less so than for sites from the same protein. Only 29% of non-interacting sites have a likelihood ratio larger than the median for interacting sites. These empirical distributions, allow us to approximate of the log odds of the “interacting” vs “non-interacting” amino acids

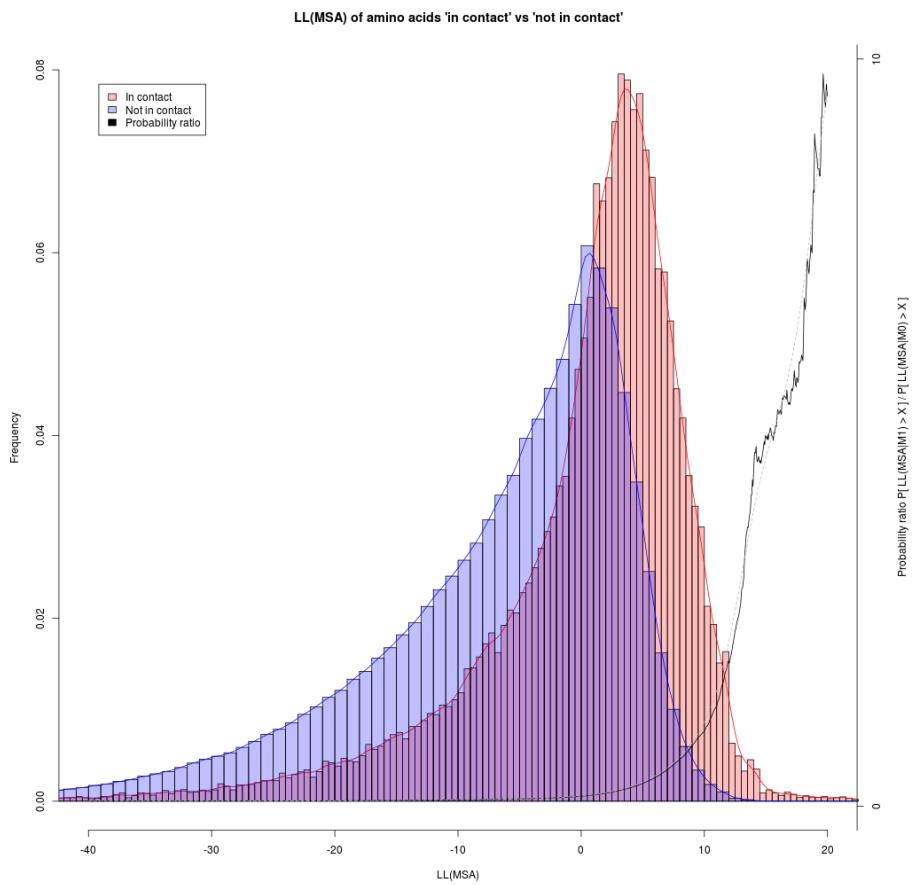


Figure 4–2: Red & Blue: Histogram of $LL(M_{SA})$ “in contact” vs ‘null’ within amino acids within the protein (Pdb). Black: Cummulative probability ratios $R_{M_{SA}} = P[LL(M_{SA}|M1) > X]/P[LL(M_{SA}|M0) > X]$. Dotted grey: Smoothed $R_{M_{SA}}$

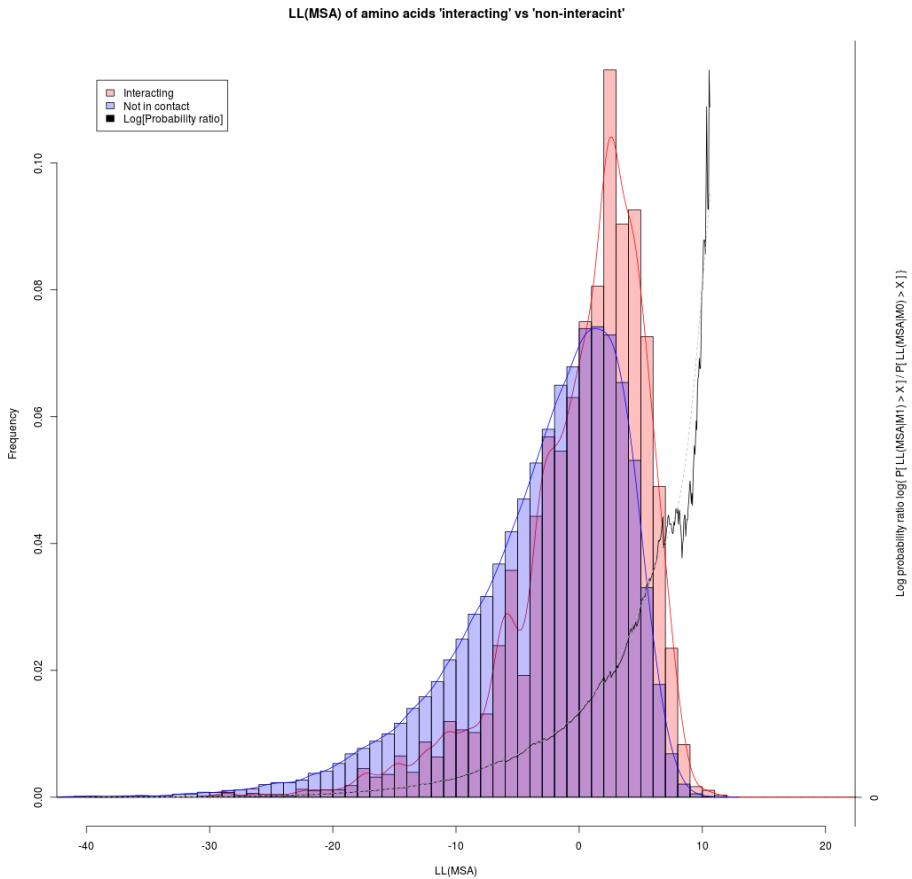


Figure 4–3: A) $LL(M_{SA})$ in contact vs 'null' interacting amino acids in co-crystallized proteins (from Pdb), roughly 40% of interacting records have $LL(M_{SA}) > 1$. B) Log odds ratio of cummulative $LL(M_{SA})$ probability (interacting / non-interacting)

distributions as $\log_{odds}(x) = \log[P(LL(M_{SA}|Q2) \geq x)/P(LL(M_{SA}|Q) \geq x)]$, which fits well an exponential function $e^{0.195x} - 1.018$ where $\alpha = 0.195$ and $\beta = 1.018$. The log odds value is capped to 4.0 to avoid artificially increasing Bayes Factors.

Figure 4 shows the example of a predicted contact $LL(M_{SA}) = 7.7$ between SENP1 and SUMO1 proteins detected by our method. The co-crystallized structure from PDB highlights the interacting amino acids (less than 3 Å apart) and the corresponding M_{SA} matched with the phylogenetic tree suggests co-evolutionary evidence.

Although our approach aims at identifying contacting residues from different proteins, it can also be used to predict interactions between proteins as a whole. We extracted from BioGrid [145] a set of $\sim 3,000$ pairs of human proteins with evidence of interaction, and further required that both proteins belong to the same pathway (MsigDb, C2 groups [147]), and their corresponding genes are expressed in the same tissue (GTex [103], expression of 1 FPKM or more, tissues $\in \{\text{skeletal muscle, adipose tissue, pancreatic Islets}\}$). We define “non-interacting” pairs of proteins as those not fulfilling any of the three conditions, and chose $\sim 3,000$ such pairs randomly.

Let the two proteins considered have amino acid sequences $a_1 \dots a_m$ and $b_1 \dots b_n$. To obtain the prediction score for this pair of proteins, we identify the pair of length- k substrings $a_i, a_{i+1}, \dots, a_{i+k-1}$ and $b_j, b_{j+1}, \dots, b_{j+k-1}$ that exhibit the strongest support for parallel or antiparallel interactions, i.e. for which $\max\{\sum_{l=0}^{k-1} LL(a_{i+l}, b_{j+l}), \sum_{l=0}^{k-1} LL(a_{i+l}, b_{j+k-1-l})\}$ is maximized. Empirically, the value of k that seems to provide the best predictive power is $k = 3$. We calculate the average likelihood ratio of three consecutive pairs of amino acids ($\text{avg}_3[LL(M_{SA})]$), either in the forward or reverse directions. For a given pair of genes, we calculate the highest $\text{avg}_3[LL(M_{SA})]$ and pick the highest number as representative for that pair. As shown in Figure ??), prediction accuracy is quite good (p-value $< 2 \cdot 10^{-42}$), considering the modest amount of information considered.

4.5.3 Epistatic GWAS analysis

We apply our methods to a cohort of $\sim 13,000$ individuals in a case-control study of type II diabetes [109]. This multi-ethnic study covers exons of unrelated individuals from five major ancestral groups (European descent, South Asian, East Asian, Hispanic and African American descent) using an average sequencing coverage over $80\times$ yielding $1.7M$ coding variants.

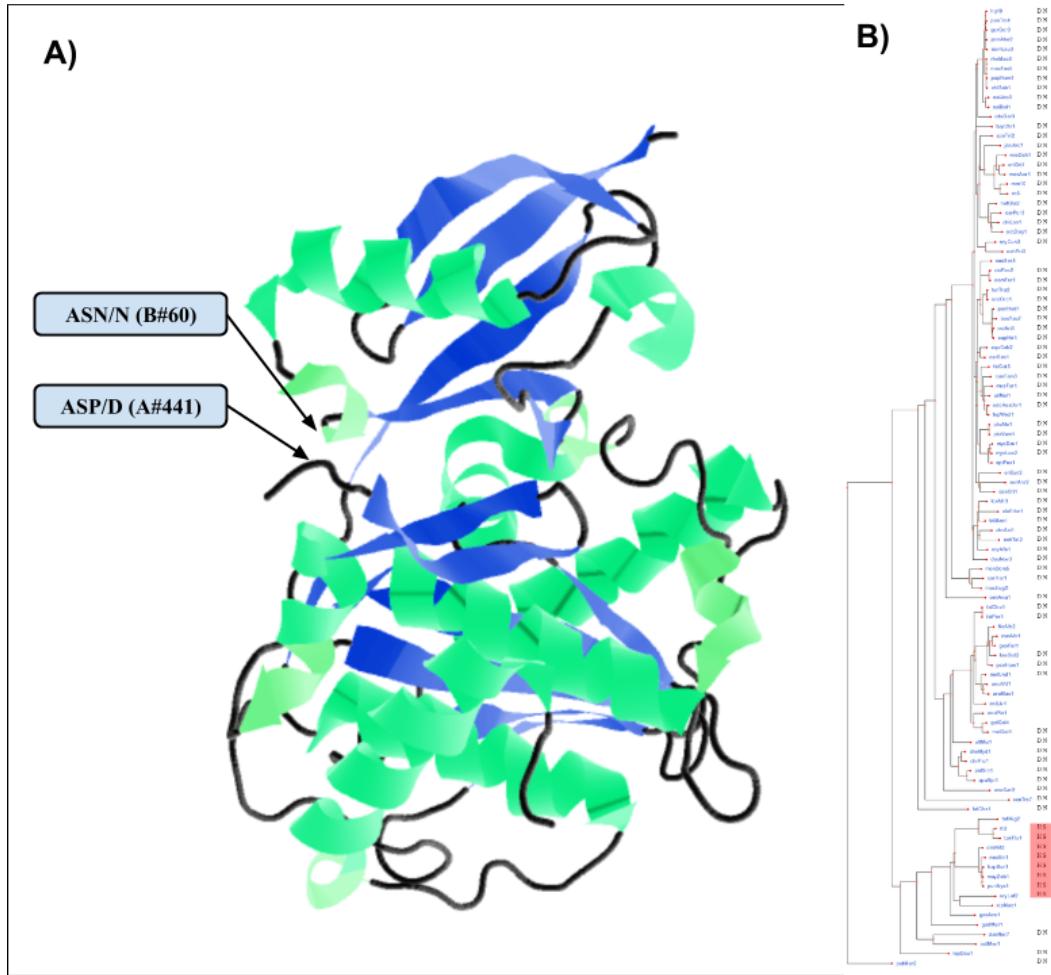


Figure 4–4: Example of amino acid interaction between SENP1 and SUMO1 proteins detected by our method $LL(M_{SA}) = 7.7$. A) PDB structure 2G4D, shows that the amino acids are in close proximity: amino acid #441 in SENP1 (ASP/D) interacts with amino acid #60 in SUMO1 (ASN/N). B) Multiple sequence alignment and phylogenetic tree shows putative compensatory amino acid substitution pair “D-N” replaced by “H-S” (colors added for visual reference).

Computational and statistical issues. It is easy to see that the computational burden for detection of pairs of interacting genetic loci affecting disease risk is significantly larger than in a standard (single variant) GWAS study. A priori all pairs of variants should be analyzed thus significantly increasing the number of statistical tests. This also reduces statistical power since the required p-value significance level would be orders of magnitude smaller. A naive approach would estimate that if a typical genetic sequencing study has 10^6 variants, a GWAS on epistatic variants would square that number of statistical tests, thus p-values required for significance would be in the order of $0.05/(10^6)^2 = 5 \cdot 10^{-14}$.

Fortunately these numbers can be reduced significantly. First, in this study, we only concentrate on non-synonymous coding variants. Second, as required by our co-evolutionary model, only variants overlapping a multiple sequence alignment are taken into account (when several multiple sequence alignments overlapped a region, the alignment with the longest number of proteins was selected). Third, if two variants for which the interaction term $(g_{s,i}g_{s,j})$ is zero in all samples, which usually happens for pairs of rare variants, there is no need to calculate the logistic regression since in the epistatic model has a trivial solution by setting θ_3 to zero, thus we can skip these variant pairs. Fourth, if the variants and the epistatic term $[g_{s,i}, g_{s,j}, g_{s,i}g_{s,j}]$ are linearly dependent, the logistic regression result will be meaningless, so we can safely skip such variant pairs. Fifth, if one of the variants has high allele frequency respect to the other, all non-zero epistatic terms may lie in the same positions as non-zero genotypes from the low frequency variant, causing logistic regression estimates to artificially inflate the coefficients of the low frequency variant and the epistatic term thus creating an artificially high association (low p-value). So we filter out these variant pairs as well. Finally, we filter

out all variants having Hardy-Weinberg p-value of less than 10^{-6} , since these variants also artificially inflate the logistic regression coefficients. Once the results are obtained, we can focus on interactions by further filtering results and keeping variant pairs whose alternative logistic model (see equation 4.6) has small absolute values for θ_1 and θ_2 while having large absolute values for θ_3 , specifically we keep results if $|\theta_3| > K(|\theta_1| + |\theta_2|)$ (based on empirical data, we set $K = 3$).

After all filters have been applied, using a model GWAS data which involves $\sim 13,000$ individuals and over $1.7M$ variants, the number of variant pairs to be analyzed is less than 10^8 variants, as opposed to the naive estimation of $\sim 10^{12}$ variant pairs. By means of the z-score relationship [66], we set the GWAS significance threshold for 50 million pairs at $\log_{10}[BF] = 8.0$.

Calculating Bayes factors involves using prior parameter distributions. In order to estimate the distributions, we run the logistic regression fitting analysis and plot the parameter distributions for different levels of significance. As expected most parameters have unimodal distribution, except for θ_3 , which has a multimodal distribution (Figure ??). For all parameters, except θ_3 , we use a normal distribution centered at the mean and variance was set to at least one ($\sigma = 1$) even though most times the variance is much smaller (this is done to avoid penalizing outliers too heavily and to have smooth derivatives near the maximum likelihood estimates). For θ_3 , which has a multimodal distribution, we fit mixed model parameters using an EM algorithm on high significant term ($BF_{raw} \geq 3$), as shown in Supplementary Figure / Table ??.

Results. Variant annotated and filtered according to the previous paragraphs lead to $\sim 32M$ pairs of variants having high log likelihood in our logistic regression model ($LL_{LogReg} > 6$, in equation 4.4) that were further analysed under co-evolutionary and Bayesian models. The complete analysis takes less

Combined model	Logistic regression			Variant 1				Variant 2			
	log10(BF)	p-value	log10(BF)	Coordinate	Gene	Functional annotation	Coordinate	Gene	Functional annotation		
8.21	7.80E-07	6.94	20:49129705_G/T,A	PTGIS	STOP_GAINED	16:81348733_G/T	GAN	NON_SYNONYMOUS_CODING			
7.12	2.73E-06	6.47	4:90743415_T/C	SNCA	NON_SYNONYMOUS_CODING	2:179659911_G/A	TTN	NON_SYNONYMOUS_CODING			
7.35	2.26E-06	6.43	3:53213690_G/C	PRKCD	NON_SYNONYMOUS_CODING	14:75746689_C/T	FOS	NON_SYNONYMOUS_CODING			
6.80	2.82E-06	6.36	2:242795125_G/A	PDCD1	NON_SYNONYMOUS_CODING	8:13356801_G/A	DLC1	NON_SYNONYMOUS_CODING			
6.35	3.85E-06	6.22	11:57582923_G/A	CTNND1	AA_modification:Phosphoserine	2:220420784_G/A	OBSL1	NON_SYNONYMOUS_CODING			
8.08	6.60E-06	6.07	1:112298763_T/C	DDX20	NON_SYNONYMOUS_CODING	3:125826058_T/C	ALDH1L1	NON_SYNONYMOUS_CODING			
7.22	4.77E-06	6.04	2:179457146_G/A	TTN	NON_SYNONYMOUS_CODING	22:35695930_C/A	TOM1	NON_SYNONYMOUS_CODING			
7.20	9.62E-06	6.00	1:45797504_C/G	MUTYH	NON_SYNONYMOUS_CODING	8:30982425_T/G	WRN	NON_SYNONYMOUS_CODING			
7.45	1.58E-05	5.90	16:4855278_A/G	GLYR1	AA_modification:Phosphoserine	8:13259100_G/A	DLC1	NON_SYNONYMOUS_CODING			
6.81	1.55E-05	5.81	11:45975129_C/T	PHF21A	NON_SYNONYMOUS_CODING	19:49458190_C/A	BAX	NON_SYNONYMOUS_CODING			
8.67	1.02E-05	5.81	7:128490102_G/A	FLNC	NON_SYNONYMOUS_CODING	8:13357339_G/C	DLC1	NON_SYNONYMOUS_CODING			
6.42	9.91E-06	5.79	11:236090_G/A	SIRT3	NON_SYNONYMOUS_CODING	4:110615838_C/T	CASP6	NON_SYNONYMOUS_CODING			
7.87	1.08E-05	5.79	8:144993930_C/G	PLEC	NON_SYNONYMOUS_CODING	14:73422258_C/G	DCAF4	NON_SYNONYMOUS_CODING			
7.86	2.30E-05	5.70	1:45224997_A/C	KIF2C	NON_SYNONYMOUS_CODING	2:108921032_C/T	SULT1C2	NON_SYNONYMOUS_CODING			
6.47	1.56E-05	5.68	11:134252895_C/T	B3GAT1	NON_SYNONYMOUS_CODING	15:75012984_T/C	CYP1A1	NON_SYNONYMOUS_CODING			
5.67	1.24E-05	5.63	6:116441645_C/G	COL10A1	NON_SYNONYMOUS_CODING	20:30072135_G/A	REM1	NON_SYNONYMOUS_CODING			
7.40	2.89E-05	5.61	1:201016295_G/A	CACNA1S	NON_SYNONYMOUS_CODING	19:17000695_G/A	F2RL3	NON_SYNONYMOUS_CODING			
5.87	8.39E-07	5.59	7:43351409_T/G	HECW1	NON_SYNONYMOUS_CODING	2:219294200_A/G	VIL1	NON_SYNONYMOUS_CODING			
5.84	2.34E-05	5.58	15:67457334_A/G	SMAD3	NON_SYNONYMOUS_CODING	1:201052381_A/G	CACNA1S	NON_SYNONYMOUS_CODING			
7.00	1.46E-05	5.56	10:53822300_A/G	PRKG1	NON_SYNONYMOUS_CODING	9:140007465_G/A	DPP7	NON_SYNONYMOUS_CODING			
6.93	2.96E-05	5.56	2:225362477_C/T	CUL3	NON_SYNONYMOUS_CODING	9:120476787_C/G	TLR4	STOP_GAINED			

Table 4–5: Results from epistatic GWAS analysis of type II diabetes sequencing data

than 2 days using a 1,000 CPU-cluster, thus that an epistatic GWAS analysis of large cohort sequencing data is feasible using current computational resources.

. Table 4–5 shows the main results from our GWAS epistatic analysis. Selected genes are shown in red in table 4–5 belong to a hand curated set of genes either associated with diabetes or known to be in diabetes related pathway. does this look like something interesting?

It should be noted that some of the top results include known amino acid modification sites such as Phosphoserine (or Glycosylation, not shown), these sites any biological meaning? it looks like these site should indeed be interacting....

From the top results what story can we tell?

4.6 Discussion

In this paper, we propose a novel methodology for genome wide association studies of pairs of variants under putative epistatic interaction. Due

to the large number of statistical tests required in epistatic analysis, and the corresponding reduction of statistical power, this type of analysis is meant to be applied to datasets consisting of large number of samples, but our highly optimized algorithms are suitable for large scale sequencing genomic studies.

The co-evolutionary model we propose in section 4.4.1 requires multiple sequence alignment and the corresponding phylogenetic tree. Both the tree and the number of sequences in the M_{SA} should remain constant throughout the genome in order to take advantage of computational optimizations (matrix exponential precalculation and “constant tree caching”) that allow the algorithm to be applied at genome-wide scale. Some multiple sequence alignments (such as Pfam) usually have different number of sequences for each protein (thus different phylogenetic trees). This poses two main disadvantages for our methodology: i) we cannot benefit from the previously mentioned optimizations, since they require a constant phylogenetic tree throughout the whole genome; and ii) we would add the problem of reconciling different phylogenetic trees from two proteins, which may lead to inconsistencies. For this reasons, we selected UCSC’s multi-100way [81], a genome wide multiple sequence alignment of 100 organisms, which has single genome wide phylogenetic tree.

We tested whether our co-evolutionary model can separate clinically relevant variants from ClinVar database [88] according to their clinical significance attribute (CLNSIG). Interestingly, variants categorized as “benign” or “druggable” have higher scores (mean $LL(M_{SA})$ within protein) than variants categorized as pathogenic (Supplementary Tables ??, ?? and Figure ??). We speculate that this might be because amino acids that can be compensated would be characterized as “benign” whereas deleterious amino acids changes cannot be compensated by mutation.

Finally, we show the application of our methods to a large scale exome sequencing study for type II diabetes consisting of $\sim 13,000$ samples and $\sim 1,7M$ variants. First, this shows that it is feasible to apply our methods GWAS-scale datasets. Second we show that some association of pairs of putatively interacting variants with type II diabetes have Bayes Factors showing suggestive association. Although larger cohorts are needed in order to find risk alleles that may be hidden by lower frequencies, thus not captured by this study.

As future work, we plan to extend our method to analyze domain specific transition pairs, this would allow to obtain better estimates for known interaction domains. Another line of work is to perform GWAS using kernel based statistics of multiple variants [166] thus allowing simultaneous analysis of nearby variants in a putative interaction hotspot, in this case the epistatic information would be used as a function modifying the kernel, instead of a bayesian prior.

CHAPTER 5

Conclusions

5.1 Contributions

In this report we show three steps involved in the analysis of sequencing data and identifying the links to disease. Each step is characterized by very different problems that need to be addressed.

- i) The first step is to reduce large amounts of information generated by high throughput experiments into a manageable subset. In our case, it involves reducing the raw sequencing information to a variant call set, but it could be any other features to be analyzed (RNA expression, transcript structure, enrichment peaks, genome reference assembly, etc.). This is mainly done by mapping reads into a reference genome and then using variant call algorithms. This step is characterized by requiring fast parallel algorithms and usually, due to the amount of data involved, I/O can be one of the bottlenecks. Algorithm that work on “chunks of data” instead of the whole data-set are preferred, and in many cases exist, because it makes the problem trivial to parallelize. Usually several stages of these highly specialized algorithms are combined into a “data analysis pipeline”. Programming data analysis pipelines is not trivial since it requires process coordinations, robustness, scalability and flexibility (data pipelines, particularly in research environments, tend to change often). Although many solutions are available (usually in the form of libraries), these tend to make pipeline programming cumbersome or create new programming paradigms thus introducing steep learning curves. In Chapter 2, we solve problems related to pipeline programming in a

novel way by creating a new programming language, BDS, that simplifies the creation of robust, scalable and flexible data pipelines. Although the main goal was managing our sequencing data pipelines, BDS is a flexible datacenter-scale programming language that can be applied to many large data pipelines (a.k.a. Big Data problems).

- ii) The second step in our data analysis, consists of functional annotations, prioritization and filtering. The main concern in the annotation step performing an adequate filtering of what should be considered relevant variants for our experiment from irrelevant ones. Functional annotation of genomic variants was until not long ago an unsolved problem and shortly after created SnpEff & SnpSift, they quickly became widely adopted by the research community. In Chapter 3 we describe the challenges of variant annotations and some of the solutions we implemented in our algorithms.
- iii) Finally, in Chapter 4, we analyse the problem of finding genetic links to complex disease. This is known to be a difficult problem affected by several hidden co-factors that bias the results (e.g. population structure). Furthermore there are unsolved problems, such as missing heritability, implying that genomic links to complex disease may not be found using traditional GWAS methodologies. We believe that alternative models that combine higher level information, may help to boost statistical significance.
- iii.a) We were involved in two major projects on GWAS of type II diabetes using: a) cohorts of multi-ethnic unrelated individuals and b) family pedigrees. Results uncovered new genes linked to diabetes. Also, the studies indicate that one of the main hypothesis in the field, the “Rare variant hypothesis”, might not hold strong.

iii.b) We propose a new methodology for addressing a difficult problem:
detection of two interacting genomic loci that affect disease risk.

Our models combine genotype information and co-evolutionary methods. We show that efficient algorithms make these studies computationally feasible, albeit using large computational resources, and we apply them to real data from type II diabetes sequencing study of over 13,000 individuals.

These three Chapters (three steps) complete our journey from “raw data” to “biological insight” trying to find the genetic causes of complex disease.

5.2 Future work

Here we propose several improvements, extensions and future lines of work for each of the methods developed in this thesis (some of them are currently being developed / explored):

BigDataScript. We are adding native support for new clusters and frameworks, such as LSF, Mesos, Kubernetes as well as a “*Generic cluster*” API which allows the user to customize BigDataScript for any cluster or framework by encapsulating task management via user defined scripts. On the language specification side, we are exploring ways to add functional constructs such as `map`, `apply`, `filter` as well as support for *map/reduce* and *scatter/gather* which are convenient ways to define some problems in data pipeline programming. Finally, will be incorporating user defined data structures or a basic class mechanism (BDS currently supports maps and list).

Annotations. We are creating new annotation standard for VCF files, in an effort coordinated with the developer of other annotations tools (such as ANNOVAR, ENSEMBLs Variant effect predictor -VEP-, JAnnovar, etc.). We are actively contributing with the “*Global Alliance for genomic and Health*”

(GA4GH) in the variant annotation specification & API definitions. We plan to extend SnpEff’s variant annotations capabilities to *haplotype-based* annotations, which means taking into account phasing information (either by phasing trios or “read phasing”) to calculate compound variant effects (e.g. phased SNPs affecting the same codon or compensating frame shifts within the same DNA strand). Finally, we are using information theoretic analysis of splice sites from several species in order to improve splicing effect predictions.

GWAS Epistasis. As future work, we’d like to evaluate the possibility of incorporating contextual information, such as protein domain, in order to build more specific co-evolutionary models. Other improvements include further optimization of logistic regression and Bayes factor algorithms since any improvement greatly reduces computational times. We also plan to use our methods on even larger type II diabetes cohorts that are currently being sequenced. Finally, we are evaluating the possibility of incorporating higher order interactions by clustering genes from our variant-pairs analysis and then evaluate them in a joint analysis.

5.3 Perspectives

Genomic research for complex disease is trending towards larger and larger cohorts in order to improve statistical power. Some years ago, projects involving hundreds to a thousand individuals were common. To put this in perspective, that’s the population of a village, or a small town. Nowadays, projects like the T2D consortia, sequence in the order of 20,000 people (i.e. the population of a large town). I am aware, through personal communications with other researchers, that projects being drafted for sequencing over 100,000 individuals (i.e. the population of a whole city). This quest for ever bigger sample sizes shows how elusive the genetic causes of complex diseases are.

The methods developed here aim to help in the processing of these huge datasets (BDS), annotate and prioritize the variants (SnpEff) before testing for significance. But also help in looking at these variants from another perspective (epistatic GWAS) than the traditional “single variant association” approach. It might be true that huge sample sizes are needed to uncover risk loci, but perhaps one of the reasons why traditional GWAS studies are not finding as many associations as expected is just that we they are looking in the wrong place. In science we must explore all possibilities.

References

- [1] Dario Acampora, Virginia Avantaggiato, Francesca Tuorto, Paolo Barone, Heinrich Reichert, Robert Finkelstein, and Antonio Simeone. Murine otx1 and drosophila otd genes share conserved genetic functions required in invertebrate and vertebrate brain development. *Development*, 125(9):1691–1702, 1998.
- [2] Marit Ackermann and Andreas Beyer. Systematic detection of epistatic interactions based on allele pair frequencies. *PLoS genetics*, 8(2):e1002463, 2012.
- [3] I.A. Adzhubei, S. Schmidt, L. Peshkin, V.E. Ramensky, A. Gerasimova, P. Bork, A.S. Kondrashov, and S.R. Sunyaev. A method and server for predicting damaging missense mutations. *Nature methods*, 7(4):248–249, 2010.
- [4] Orly Agamy, Bruria Ben Zeev, Dorit Lev, Barak Marcus, Dina Fine, Dan Su, Ginat Narkis, Rivka Ofir, Chen Hoffmann, Esther Leshinsky-Silver, et al. Mutations disrupting selenocysteine formation cause progressive cerebello-cerebral atrophy. *The American Journal of Human Genetics*, 87(4):538–544, 2010.
- [5] Bruce Alberts, Dennis Bray, Julian Lewis, Martin Raff, Keith Roberts, James D Watson, and AV Grimstone. Molecular biology of the cell (3rd edn). *Trends in Biochemical Sciences*, 1995.
- [6] S.F. Altschul, W. Gish, W. Miller, E.W. Myers, D.J. Lipman, et al. Basic local alignment search tool. *Journal of molecular biology*, 215(3):403–410, 1990.
- [7] Stephen F Altschul, Thomas L Madden, Alejandro A Schäffer, Jinghui Zhang, Zheng Zhang, Webb Miller, and David J Lipman. Gapped blast and psi-blast: a new generation of protein database search programs. *Nucleic acids research*, 25(17):3389–3402, 1997.
- [8] Gerd Anders, Sebastian D Mackowiak, Marvin Jens, Jonas Maaskola, Andreas Kuntzagk, Nikolaus Rajewsky, Markus Landthaler, and Christoph Dieterich. dorina: a database of rna interactions in post-transcriptional regulation. *Nucleic acids research*, page gkr1007, 2011.

- [9] Jennifer A Anderson, William D Gilliland, and Charles H Langley. Molecular population genetics and evolution of drosophila meiosis genes. *Genetics*, 181(1):177–185, 2009.
- [10] Manuel Ascano, Markus Hafner, Pavol Cekan, Stefanie Gerstberger, and Thomas Tuschl. Identification of rna–protein interaction networks using par-clip. *Wiley Interdisciplinary Reviews: RNA*, 3(2):159–177, 2012.
- [11] Geraldine A Auwera, Mauricio O Carneiro, Christopher Hartl, Ryan Poplin, Guillermo del Angel, Ami Levy-Moonshine, Tadeusz Jordan, Khalid Shakir, David Roazen, Joel Thibault, et al. From fastq data to high-confidence variant calls: The genome analysis toolkit best practices pipeline. *Current protocols in bioinformatics*, pages 11–10, 2013.
- [12] D.J. Balding. A tutorial on statistical methods for population association studies. *Nature Reviews Genetics*, 7(10):781–791, 2006.
- [13] Michael J Bamshad, Sarah B Ng, Abigail W Bigham, Holly K Tabor, Mary J Emond, Deborah A Nickerson, and Jay Shendure. Exome sequencing as a tool for mendelian disease gene discovery. *Nature Reviews Genetics*, 12(11):745–755, 2011.
- [14] David J Begun and Charles F Aquadro. Levels of naturally occurring dna polymorphism correlate with recombination rates in *d. melanogaster*. 1992.
- [15] Callum J Bell, Darrell L Dinwiddie, Neil A Miller, Shannon L Hateley, Elena E Ganusova, Joann Mudge, Ray J Langley, Lu Zhang, Clarence C Lee, Faye D Schilkey, et al. Carrier testing for severe childhood recessive diseases by next-generation sequencing. *Science translational medicine*, 3(65):65ra4–65ra4, 2011.
- [16] Bradley E Bernstein, John A Stamatoyannopoulos, Joseph F Costello, Bing Ren, Aleksandar Milosavljevic, Alexander Meissner, Manolis Kellis, Marco A Marra, Arthur L Beaudet, Joseph R Ecker, et al. The nih roadmap epigenomics mapping consortium. *Nature biotechnology*, 28(10):1045–1048, 2010.
- [17] Mathieu Blanchette, W James Kent, Cathy Riemer, Laura Elnitski, Arrian FA Smit, Krishna M Roskin, Robert Baertsch, Kate Rosenbloom, Hiram Clawson, Eric D Green, et al. Aligning multiple genomic sequences with the threaded blockset aligner. *Genome research*, 14(4):708–715, 2004.
- [18] Alan P Boyle, Eurie L Hong, Manoj Hariharan, Yong Cheng, Marc A Schaub, Maya Kasowski, Konrad J Karczewski, Julie Park, Benjamin C Hitz, Shuai Weng, et al. Annotation of functional variation in personal genomes using regulomedb. *Genome research*, 22(9):1790–1797, 2012.

- [19] Sydney Brenner, AOW Stretton, and S Kaplan. Genetic code: the non-sensetriplets for chain termination and their suppression. 1965.
- [20] Saverio Brogna and Jikai Wen. Nonsense-mediated mrna decay (nmd) mechanisms. *Nature structural & molecular biology*, 16(2):107–113, 2009.
- [21] Charles George Broyden. The convergence of a class of double-rank minimization algorithms 1. general considerations. *IMA Journal of Applied Mathematics*, 6(1):76–90, 1970.
- [22] Jan Christian Bryne, Eivind Valen, Man-Hung Eric Tang, Troels Marstrand, Ole Winther, Isabelle da Piedade, Anders Krogh, Boris Lenhard, and Albin Sandelin. Jaspar, the open access database of transcription factor-binding profiles: new content and tools in the 2008 update. *Nucleic acids research*, 36(suppl 1):D102–D106, 2008.
- [23] Lukas Burger and Erik van Nimwegen. Disentangling direct from indirect co-evolution of residues in protein alignments. *PLoS computational biology*, 6(1):e1000633, 2010.
- [24] Susan E Celniker, Laura AL Dillon, Mark B Gerstein, Kristin C Gunsalus, Steven Henikoff, Gary H Karpen, Manolis Kellis, Eric C Lai, Jason D Lieb, David M MacAlpine, et al. Unlocking the secrets of the genome. *Nature*, 459(7249):927–930, 2009.
- [25] Judy C Chang and Yuet Wai Kan. beta 0 thalassemia, a nonsense mutation in man. *Proceedings of the National Academy of Sciences*, 76(6):2886–2889, 1979.
- [26] B Charlesworth, JA Coyne, and NH Barton. The relative rates of evolution of sex chromosomes and autosomes. *American Naturalist*, pages 113–146, 1987.
- [27] Venkateswara R Chintapalli, Jing Wang, and Julian AT Dow. Using flyatlas to identify better drosophila melanogaster models of human disease. *Nature genetics*, 39(6):715–720, 2007.
- [28] Yongwook Choi, Gregory E Sims, Sean Murphy, Jason R Miller, and Agnes P Chan. Predicting the functional effect of amino acid substitutions and indels. *PloS one*, 7(10):e46688, 2012.
- [29] Kristian Cibulskis, Michael S Lawrence, Scott L Carter, Andrey Sivachenko, David Jaffe, Carrie Sougnez, Stacey Gabriel, Matthew Meyerson, Eric S Lander, and Gad Getz. Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nature biotechnology*, 31(3):213–219, 2013.

- [30] P. Cingolani, A. Platts, M. Coon, T. Nguyen, L. Wang, S.J. Land, X. Lu, D.M. Ruden, et al. A program for annotating and predicting the effects of single nucleotide polymorphisms, snpeff: Snps in the genome of drosophila melanogaster strain w1118; iso-2; iso-3. *Fly*, 6(2):0–1, 2012.
- [31] Pablo Cingolani, Viral M Patel, Melissa Coon, Tung Nguyen, Susan J Land, Douglas M Ruden, and Xiangyi Lu. Using drosophila melanogaster as a model for genotoxic chemical mutational studies with a new program, snpsift. *Toxicogenomics in non-mammalian species*, page 92, 2012.
- [32] Greg W Clark, Sharon H Ackerman, Elisabeth R Tillier, and Domenico L Gatti. Multidimensional mutual information methods for the analysis of covariation in multiple sequence alignments. *BMC bioinformatics*, 15(1):157, 2014.
- [33] G.M. Clarke, C.A. Anderson, F.H. Pettersson, L.R. Cardon, A.P. Morris, and K.T. Zondervan. Basic statistical analysis in genetic case-control studies. *Nature protocols*, 6(2):121–133, 2011.
- [34] FS Collins, ES Lander, J. Rogers, RH Waterston, and I. Conso. Finishing the euchromatic sequence of the human genome. *Nature*, 431(7011):931–945, 2004.
- [35] Onofre Combarros, Mario Cortina-Borja, A David Smith, and Donald J Lehmann. Epistasis in sporadic alzheimer’s disease. *Neurobiology of aging*, 30(9):1333–1349, 2009.
- [36] 1000 Genomes Project Consortium et al. An integrated map of genetic variation from 1,092 human genomes. *Nature*, 491(7422):56–65, 2012.
- [37] Diabetes SAT2D Consortium, Diabetes MAT2D Consortium, Anubha Mahajan, Min Jin Go, Weihua Zhang, Jennifer E Below, Kyle J Gaulton, Teresa Ferreira, Momoko Horikoshi, Andrew D Johnson, et al. Genome-wide trans-ancestry meta-analysis provides insight into the genetic architecture of type 2 diabetes susceptibility. *Nature genetics*, 46(3):234–244, 2014.
- [38] ENCODE Project Consortium et al. An integrated encyclopedia of dna elements in the human genome. *Nature*, 489(7414):57–74, 2012.
- [39] UniProt Consortium et al. Update on activities at the universal protein resource (uniprot) in 2013. *Nucleic acids research*, 41(D1):D43–D47, 2013.
- [40] Heather J Cordell. Epistasis: what it means, what it doesn’t mean, and statistical methods to detect it in humans. *Human molecular genetics*, 11(20):2463–2468, 2002.

- [41] Heather J Cordell. Detecting gene–gene interactions that underlie human diseases. *Nature Reviews Genetics*, 10(6):392–404, 2009.
- [42] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, Clifford Stein, et al. *Introduction to algorithms*, volume 2. MIT press Cambridge, 2001.
- [43] Jasmin Coulombe-Huntington, Kevin CL Lam, Christel Dias, and Jacek Majewski. Fine-scale variation and genetic determinants of alternative splicing across individuals. *PLoS genetics*, 5(12):e1000766, 2009.
- [44] BS Cox, MF Tuite, and CS McLaughlin. The ψ factor of yeast: a problem in inheritance. *Yeast*, 4(3):159–178, 1988.
- [45] Robert Culverhouse, Brian K Suarez, Jennifer Lin, and Theodore Reich. A perspective on epistasis: limits of models displaying no main effect. *The American Journal of Human Genetics*, 70(2):461–471, 2002.
- [46] Petr Danecek, Adam Auton, Goncalo Abecasis, Cornelis A Albers, Eric Banks, Mark A DePristo, Robert E Handsaker, Gerton Lunter, Gabor T Marth, Stephen T Sherry, et al. The variant call format and vcftools. *Bioinformatics*, 27(15):2156–2158, 2011.
- [47] Ingrid Daubechies, Ronald DeVore, Massimo Fornasier, and C Sinan Güntürk. Iteratively reweighted least squares minimization for sparse recovery. *Communications on Pure and Applied Mathematics*, 63(1):1–38, 2010.
- [48] E.V. Davydov, D.L. Goode, M. Sirota, G.M. Cooper, A. Sidow, and S. Batzoglou. Identifying a high fraction of the human genome to be under selective constraint using gerp++. *PLoS computational biology*, 6(12):e1001025, 2010.
- [49] David de Juan, Florencio Pazos, and Alfonso Valencia. Emerging methods in protein co-evolution. *Nature Reviews Genetics*, 14(4):249–261, 2013.
- [50] Glynn Dennis Jr, Brad T Sherman, Douglas A Hosack, Jun Yang, Wei Gao, H Clifford Lane, Richard A Lempicki, et al. David: database for annotation, visualization, and integrated discovery. *Genome biol*, 4(5):P3, 2003.
- [51] M.A. DePristo, E. Banks, R. Poplin, K.V. Garimella, J.R. Maguire, C. Hartl, A.A. Philippakis, G. Del Angel, M.A. Rivas, M. Hanna, et al. A framework for variation discovery and genotyping using next-generation dna sequencing data. *Nature genetics*, 43(5):491–498, 2011.

- [52] TD Dreesen, DH Johnson, and S Henikoff. The brown protein of *drosophila melanogaster* is similar to the white protein and to components of active transport complexes. *Molecular and cellular biology*, 8(12):5206–5215, 1988.
- [53] R. Durbin. *Biological sequence analysis: probabilistic models of proteins and nucleic acids*. Cambridge Univ Pr, 1998.
- [54] William G Fairbrother, Ru-Fang Yeh, Phillip A Sharp, and Christopher B Burge. Predictive identification of exonic splicing enhancers in human genes. *Science*, 297(5583):1007–1013, 2002.
- [55] Joseph Felsenstein and Joseph Felenstein. *Inferring phylogenies*, volume 2. Sinauer Associates Sunderland, 2004.
- [56] P. Ferragina and G. Manzini. Opportunistic data structures with applications. In *Foundations of Computer Science, 2000. Proceedings. 41st Annual Symposium on*, pages 390–398. IEEE, 2000.
- [57] Kevin M Flanigan, Diane M Dunn, Andrew von Niederhausern, Michael T Howard, Jerry Mendell, Anne Connolly, Carol Saunders, Ann Modrcin, Majed Dasouki, Giacomo P Comi, et al. Dmd trp3x nonsense mutation associated with a founder effect in north american families with mild becker muscular dystrophy. *Neuromuscular Disorders*, 19(11):743–748, 2009.
- [58] Paul Flicek, Ikhlaq Ahmed, M Ridwan Amode, Daniel Barrell, Kathryn Beal, Simon Brent, Denise Carvalho-Silva, Peter Clapham, Guy Coates, Susan Fairley, et al. Ensembl 2013. *Nucleic acids research*, page gks1236, 2012.
- [59] Hong Gao, Julie M Granka, and Marcus W Feldman. On the classification of epistatic interactions. *Genetics*, 184(3):827–837, 2010.
- [60] M. Garber, M. Guttman, M. Clamp, M.C. Zody, N. Friedman, and X. Xie. Identifying novel constrained elements by exploiting biased substitution patterns. *Bioinformatics*, 25(12):i54–i62, 2009.
- [61] David Gatfield, Leonie Unterholzner, Francesca D Ciccarelli, Peer Bork, and Elisa Izaurralde. Nonsense-mediated mrna decay in *drosophila*: at the intersection of the yeast and mammalian pathways. *The EMBO journal*, 22(15):3960–3970, 2003.
- [62] Belinda Giardine, Cathy Riemer, Ross C Hardison, Richard Burhans, Laura Elnitski, Prachi Shah, Yi Zhang, Daniel Blankenberg, Istvan Albert, James Taylor, et al. Galaxy: a platform for interactive large-scale genome analysis. *Genome research*, 15(10):1451–1455, 2005.

- [63] G. Gibson. Rare and common variants: twenty arguments. *Nature Reviews Genetics*, 13(2):135–145, 2012.
- [64] Jeremy Goecks, Anton Nekrutenko, James Taylor, et al. Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biol*, 11(8):R86, 2010.
- [65] Chern-Sing Goh, Andrew A Bogan, Marcin Joachimiak, Dirk Walther, and Fred E Cohen. Co-evolution of proteins with their interaction partners. *Journal of molecular biology*, 299(2):283–293, 2000.
- [66] Steven N Goodman. Toward evidence-based medical statistics. 2: The bayes factor. *Annals of internal medicine*, 130(12):1005–1013, 1999.
- [67] L Guarente, YS Lin, M Carey, L Breeden, MA Osley, J Gould, S Kim, M Kane, and L Hereford. Generating yeast transcriptional activators containing no yeast protein sequences. *Nature*, 350, 1991.
- [68] L Guariguata, DR Whiting, I Hambleton, J Beagley, U Linnenkamp, and JE Shaw. Global estimates of diabetes prevalence for 2013 and projections for 2035. *Diabetes research and clinical practice*, 103(2):137–149, 2014.
- [69] Markus Hafner, Steve Lianoglou, Thomas Tuschl, and Doron Belenky. Genome-wide identification of mirna targets by par-clip. *Methods*, 58(2):94–105, 2012.
- [70] D.L. Hartl and A.G. Clark. *Principles of population genetics*. Sinauer associates Sunderland, Massachusetts, 2007.
- [71] David Haussler, Stephen J O’Brien, Oliver A Ryder, F Keith Barker, Michele Clamp, Andrew J Crawford, Robert Hanner, Olivier Hanotte, Warren E Johnson, Jimmy A McGuire, et al. Genome 10k: a proposal to obtain whole-genome sequence for 10 000 vertebrate species. *Journal of Heredity*, 100(6):659–674, 2009.
- [72] Aleksandra Helwak, Grzegorz Kudla, Tatiana Dudnakova, and David Tollervey. Mapping the human mirna interactome by clash reveals frequent noncanonical binding. *Cell*, 153(3):654–665, 2013.
- [73] Lucia A Hindorff, Praveen Sethupathy, Heather A Junkins, Erin M Ramos, Jayashri P Mehta, Francis S Collins, and Teri A Manolio. Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proceedings of the National Academy of Sciences*, 106(23):9362–9367, 2009.

- [74] Douglas A Hosack, Glynn Dennis Jr, Brad T Sherman, H Clifford Lane, Richard A Lempicki, et al. Identifying biological themes within lists of genes with ease. *Genome Biol*, 4(10):R70, 2003.
- [75] Christopher A Hunter, Juswinder Singh, and Janet M Thornton. π - π interactions: The geometry and energetics of phenylalanine-phenylalanine interactions in proteins. *Journal of molecular biology*, 218(4):837–846, 1991.
- [76] Ivaylo P Ivanov, Andrew E Firth, Audrey M Michel, John F Atkins, and Pavel V Baranov. Identification of evolutionarily conserved non-aug-initiated n-terminal extensions in human coding sequences. *Nucleic acids research*, 39(10):4220–4234, 2011.
- [77] David T Jones, Daniel WA Buchan, Domenico Cozzetto, and Massimiliano Pontil. Psicov: precise structural contact prediction using sparse inverse covariance estimation on large multiple sequence alignments. *Bioinformatics*, 28(2):184–190, 2012.
- [78] Kenneth A Jones, Beth Borowsky, Joe A Tamm, Douglas A Craig, Margaret M Durkin, Meng Dai, Wen-Jeng Yao, Mary Johnson, Caryn Gunwaldsen, Ling-Yan Huang, et al. Gabab receptors function as a heteromeric assembly of the subunits gababr1 and gababr2. *Nature*, 396(6712):674–679, 1998.
- [79] Robbie P Joosten, Tim AH Te Beek, Elmar Krieger, Maarten L Hekkelman, Rob WW Hooft, Reinhard Schneider, Chris Sander, and Gert Vriend. A series of pdb related databases for everyday needs. *Nucleic acids research*, 39(suppl 1):D411–D419, 2011.
- [80] Luke Jostins, Adam P Levine, and Jeffrey C Barrett. Using genetic prediction from known complex disease loci to guide the design of next-generation sequencing experiments. *PloS one*, 8(10):e76328, 2013.
- [81] Donna Karolchik, Galt P Barber, Jonathan Casper, Hiram Clawson, Melissa S Cline, Mark Diekhans, Timothy R Dreszer, Pauline A Fujita, Luvina Guruvadoo, Maximilian Haeussler, et al. The ucsc genome browser database: 2014 update. *Nucleic acids research*, 42(D1):D764–D770, 2014.
- [82] Robert E Kass and Adrian E Raftery. Bayes factors. *Journal of the american statistical association*, 90(430):773–795, 1995.
- [83] Martin Alexander Kennedy. Mendelian genetic disorders. *Encyclopedia of Life Sciences*, 2001.

- [84] Szymon M Kiełbasa, Raymond Wan, Kengo Sato, Paul Horton, and Martin C Frith. Adaptive seeds tame genomic sequence comparison. *Genome research*, 21(3):487–493, 2011.
- [85] Ching Lee Koo, Mei Jing Liew, Mohd Saberi Mohamad, and Abdul Hakim Mohamed Salleh. A review for detecting gene-gene interactions using machine learning methods in genetic epidemiology. *BioMed research international*, 2013, 2013.
- [86] Marilyn Kozak. An analysis of 5'-noncoding sequences from 699 vertebrate messenger rnas. *Nucleic acids research*, 15(20):8125–8148, 1987.
- [87] P. Kumar, S. Henikoff, and P.C. Ng. Predicting the effects of coding non-synonymous variants on protein function using the sift algorithm. *Nature protocols*, 4(7):1073–1081, 2009.
- [88] Melissa J Landrum, Jennifer M Lee, George R Riley, Wonhee Jang, Wendy S Rubinstein, Deanna M Church, and Donna R Maglott. Clinvar: public archive of relationships among sequence variation and human phenotype. *Nucleic acids research*, page gkt1113, 2013.
- [89] Lydie Lane, Ghislaine Argoud-Puy, Aurore Britan, Isabelle Cusin, Paula D Duek, Olivier Evalet, Alain Gateau, Pascale Gaudet, Anne Gleizes, Alexandre Masselot, et al. nextprot: a knowledge platform for human proteins. *Nucleic acids research*, 40(D1):D76–D83, 2012.
- [90] B. Langmead and S.L. Salzberg. Fast gapped-read alignment with bowtie 2. *Nature Methods*, 2012.
- [91] B. Langmead, C. Trapnell, M. Pop, and S.L. Salzberg. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol*, 10(3):R25, 2009.
- [92] David E Larson, Christopher C Harris, Ken Chen, Daniel C Koboldt, Travis E Abbott, David J Dooling, Timothy J Ley, Elaine R Mardis, Richard K Wilson, and Li Ding. Somaticsniper: identification of somatic point mutations in whole genome sequencing data. *Bioinformatics*, 28(3):311–317, 2012.
- [93] C Lazure, WT Hum, and DM Gibson. Sequence diversity within a subgroup of mouse immunoglobulin kappa chains controlled by the igkef2 locus. *The Journal of experimental medicine*, 154(1):146–155, 1981.
- [94] B. Li and S.M. Leal. Methods for detecting associations with rare variants for common diseases: application to analysis of sequence data. *The American Journal of Human Genetics*, 83(3):311–321, 2008.

- [95] H. Li. Improving snp discovery by base alignment quality. *Bioinformatics*, 27(8):1157–1158, 2011.
- [96] H. Li. A statistical framework for snp calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*, 27(21):2987–2993, 2011.
- [97] H. Li and R. Durbin. Fast and accurate short-read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25(5), 2009.
- [98] H. Li and R. Durbin. Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinformatics*, 26(5):589, 2010.
- [99] H. Li, J. Ruan, and R. Durbin. Mapping short DNA sequencing reads and calling variants using mapping quality scores. *Genome research*, 18(11):1851, 2008.
- [100] Chenxing Liu, Fuquan Zhang, Tingting Li, Ming Lu, Lifang Wang, Weihua Yue, and Dai Zhang. Mirsnp, a database of polymorphisms altering mirna target sites, identifies mirna-related snps in gwas snps and eqtls. *BMC genomics*, 13(1):661, 2012.
- [101] Dongmei Liu and Russell L Finley. Cyclin y is a novel conserved cyclin essential for development in drosophila. *Genetics*, 184(4):1025–1035, 2010.
- [102] Xiaoming Liu, Xueqiu Jian, and Eric Boerwinkle. dbnsfp: a lightweight database of human nonsynonymous snps and their functional predictions. *Human mutation*, 32(8):894–899, 2011.
- [103] John Lonsdale, Jeffrey Thomas, Mike Salvatore, Rebecca Phillips, Edmund Lo, Saboor Shad, Richard Hasz, Gary Walters, Fernando Garcia, Nancy Young, et al. The genotype-tissue expression (gtex) project. *Nature genetics*, 45(6):580–585, 2013.
- [104] D.G. MacArthur, S. Balasubramanian, A. Frankish, N. Huang, J. Morris, K. Walter, L. Jostins, L. Habegger, J.K. Pickrell, S.B. Montgomery, et al. A systematic survey of loss-of-function variants in human protein-coding genes. *Science*, 335(6070):823–828, 2012.
- [105] Baijayanta Maiti, Sandrine Arbogast, Valérie Allamand, Mark W Moyle, Christine B Anderson, Pascale Richard, Pascale Guicheney, Ana Ferreiro, Kevin M Flanigan, and Michael T Howard. A mutation in the sepn1 selenocysteine redefinition element (sre) reduces selenocysteine incorporation and leads to sepn1-related myopathy. *Human mutation*, 30(3):411–416, 2009.

- [106] Ramamurthy Mani, Robert P St Onge, John L Hartman, Guri Giaever, and Frederick P Roth. Defining genetic interaction. *Proceedings of the National Academy of Sciences*, 105(9):3461–3466, 2008.
- [107] Teri A Manolio, Francis S Collins, Nancy J Cox, David B Goldstein, Lucia A Hindorff, David J Hunter, Mark I McCarthy, Erin M Ramos, Lon R Cardon, Aravinda Chakravarti, et al. Finding the missing heritability of complex diseases. *Nature*, 461(7265):747–753, 2009.
- [108] Debora S Marks, Thomas A Hopf, and Chris Sander. Protein structure prediction from sequence variation. *Nature biotechnology*, 30(11):1072–1080, 2012.
- [109] et. al. McCarthy, Mark. Variation in protein-coding sequence and predisposition to type 2 diabetes. *To be published*, 2015.
- [110] A. McKenna, M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytsky, K. Garimella, D. Altshuler, S. Gabriel, M. Daly, et al. The genome analysis toolkit: a mapreduce framework for analyzing next-generation dna sequencing data. *Genome research*, 20(9):1297–1303, 2010.
- [111] Brett A McKinney, David M Reif, Marylyn D Ritchie, and Jason H Moore. Machine learning for detecting gene-gene interactions. *Applied bioinformatics*, 5(2):77–88, 2006.
- [112] Faruck Morcos, Andrea Pagnani, Bryan Lunt, Arianna Bertolino, Debora S Marks, Chris Sander, Riccardo Zecchina, José N Onuchic, Terence Hwa, and Martin Weigt. Direct-coupling analysis of residue coevolution captures native contacts across many protein families. *Proceedings of the National Academy of Sciences*, 108(49):E1293–E1301, 2011.
- [113] Andrew P Morris, Benjamin F Voight, Tanya M Teslovich, Teresa Ferreira, Ayellet V Segré, Valgerdur Steinthorsdottir, Rona J Strawbridge, Hassan Khan, Harald Grallert, Anubha Mahajan, et al. Large-scale association analysis provides insights into the genetic architecture and pathophysiology of type 2 diabetes. *Nature genetics*, 44(9):981, 2012.
- [114] Eszter Nagy and Lynne E Maquat. A rule for termination-codon position within intron-containing genes: when nonsense affects rna abundance. *Trends in biochemical sciences*, 23(6):198–199, 1998.
- [115] Francesco Napolitano, Renato Mariani-Costantini, and Roberto Tagliaferri. Bioinformatic pipelines in python with leaf. *BMC bioinformatics*, 14(1):201, 2013.

- [116] Saul B Needleman and Christian D Wunsch. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of molecular biology*, 48(3):443–453, 1970.
- [117] Rasmus Nielsen, Joshua S Paul, Anders Albrechtsen, and Yun S Song. Genotype and snp calling from next-generation sequencing data. *Nature Reviews Genetics*, 12(6):443–451, 2011.
- [118] Jason ORawe, Tao Jiang, Guangqing Sun, Yiyang Wu, Wei Wang, Jingchu Hu, Paul Bodily, Lifeng Tian, Hakon Hakonarson, W Evan Johnson, et al. Low concordance of multiple variant-calling pipelines: practical implications for exome and genome sequencing. *Genome med*, 5(3):28, 2013.
- [119] Annette L Parks, Kevin R Cook, Marcia Belvin, Nicholas A Dompe, Robert Fawcett, Kari Huppert, Lory R Tan, Christopher G Winter, Kevin P Bogart, Jennifer E Deal, et al. Systematic generation of high-resolution deletion coverage of the drosophila melanogaster genome. *Nature genetics*, 36(3):288–292, 2004.
- [120] N. Patterson, A.L. Price, and D. Reich. Population structure and eigenanalysis. *PLoS genetics*, 2(12):e190, 2006.
- [121] Florencio Pazos, Manuela Helmer-Citterich, Gabriele Ausiello, and Alfonso Valencia. Correlated mutations contain information about protein-protein interaction. *Journal of molecular biology*, 271(4):511–523, 1997.
- [122] Patrick C Phillips. Epistasisthe essential role of gene interactions in the structure and evolution of genetic systems. *Nature Reviews Genetics*, 9(11):855–867, 2008.
- [123] Adrian Platts, Susan J Land, Lang Chen, Grier Page, Parsa Rasouli, Luan Wang, Xiangyi Lu, and Douglas M Ruden. Massively parallel resequencing of the isogenic drosophila melanogaster strain w1118; iso-2; iso-3 identifies hotspots for mutations in sensory perception genes. *Fly*, 3(3):192–204, 2009.
- [124] Katherine S Pollard, Melissa J Hubisz, Kate R Rosenbloom, and Adam Siepel. Detection of nonneutral substitution rates on mammalian phylogenies. *Genome research*, 20(1):110–121, 2010.
- [125] Alkes L Price, Nick J Patterson, Robert M Plenge, Michael E Weinblatt, Nancy A Shadick, and David Reich. Principal components analysis corrects for stratification in genome-wide association studies. *Nature genetics*, 38(8):904–909, 2006.
- [126] Towfique Raj, Manik Kuchroo, Joseph M Replogle, Soumya Raychaudhuri, Barbara E Stranger, and Philip L De Jager. Common risk alleles

- for inflammatory diseases are targets of recent positive selection. *The American Journal of Human Genetics*, 92(4):517–529, 2013.
- [127] Debashish Ray, Hilal Kazan, Kate B Cook, Matthew T Weirauch, Hamed S Najafabadi, Xiao Li, Serge Guerousov, Mihai Albu, Hong Zheng, Ally Yang, et al. A compendium of rna-binding motifs for decoding gene regulation. *Nature*, 499(7457):172–177, 2013.
 - [128] Ho Sung Rhee and B Franklin Pugh. Chip-exo method for identifying genomic location of dna-binding proteins with near-single-nucleotide accuracy. *Current Protocols in Molecular Biology*, pages 21–24, 2012.
 - [129] Erin Rooney Riggs, Karen E Wain, Darlene Riethmaier, Melissa Savage, Bethanny Smith-Packard, Erin B Kaminsky, Heidi L Rehm, Christa Lese Martin, David H Ledbetter, and W Andrew Fauci. Towards a universal clinical genomics database: the 2012 international standards for cytogenomic arrays consortium meeting. *Human mutation*, 34(6):915–919, 2013.
 - [130] A.F. Rope, K. Wang, R. Evjenth, J. Xing, J.J. Johnston, J.J. Swensen, B. Moore, C.D. Huff, L.M. Bird, J.C. Carey, et al. Using vaast to identify an x-linked disorder resulting in lethality in male infants due to n-terminal acetyltransferase deficiency. *The American Journal of Human Genetics*, 2011.
 - [131] Douglas M Ruden, D Curtis Jamison, Barry R Zeeberg, Mark D Garfinkel, John N Weinstein, Parsa Rasouli, and Xiangyi Lu. The edge hypothesis: Epigenetically directed genetic errors in repeat-containing proteins (rcps) involved in evolution, neuroendocrine signaling, and cancer. *Frontiers in neuroendocrinology*, 29(3):428–444, 2008.
 - [132] Mario X Ruiz-González and Mario A Fares. Coevolution analyses illuminate the dependencies between amino acid sites in the chaperonin system groes-l. *BMC evolutionary biology*, 13(1):156, 2013.
 - [133] Radhakrishnan Sabarinathan, Hakim Tafer, Stefan E Seemann, Ivo L Hofacker, Peter F Stadler, and Jan Gorodkin. The rnasnp web server: predicting snp effects on local rna secondary structure. *Nucleic acids research*, 41(W1):W475–W479, 2013.
 - [134] Simon P Sadedin, Bernard Pope, and Alicia Oshlack. Bpipe: a tool for running and managing bioinformatics pipelines. *Bioinformatics*, 28(11):1525–1526, 2012.
 - [135] Clara M Santiveri and M Jiménez. Tryptophan residues: Scarce in proteins but strong stabilizers of β -hairpin peptides. *Peptide Science*, 94(6):779–790, 2010.

- [136] Zuben E Sauna and Chava Kimchi-Sarfaty. Understanding the contribution of synonymous mutations to human disease. *Nature Reviews Genetics*, 12(10):683–691, 2011.
- [137] William Saurin, Maurice Hofnung, and Elie Dassa. Getting in or out: early segregation between importers and exporters in the evolution of atp-binding cassette (abc) transporters. *Journal of molecular evolution*, 48(1):22–41, 1999.
- [138] Eric Schadt. First steps on a long road. *Science*, 331(6018):691–691, 2011.
- [139] Gert C Scheper, Marjo S van der Knaap, and Christopher G Proud. Translation matters: protein synthesis defects in inherited disease. *Nature Reviews Genetics*, 8(9):711–723, 2007.
- [140] Valerie Schneider and Deanna Church. Genome reference consortium. 2013.
- [141] Benjamin A Shoemaker and Anna R Panchenko. Deciphering protein–protein interactions. part i. experimental techniques and databases. *PLoS computational biology*, 3(3):e42, 2007.
- [142] Benjamin A Shoemaker and Anna R Panchenko. Deciphering protein–protein interactions. part ii. computational methods to predict protein and domain interaction partners. *PLoS computational biology*, 3(4):e43, 2007.
- [143] Adam Siepel, Gill Bejerano, Jakob S Pedersen, Angie S Hinrichs, Minmei Hou, Kate Rosenbloom, Hiram Clawson, John Spieth, LaDeana W Hillier, Stephen Richards, et al. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome research*, 15(8):1034–1050, 2005.
- [144] Temple F Smith and Michael S Waterman. Identification of common molecular subsequences. *Journal of molecular biology*, 147(1):195–197, 1981.
- [145] Chris Stark, Bobby-Joe Breitkreutz, Teresa Reguly, Lorrie Boucher, Ashton Breitkreutz, and Mike Tyers. Biogrid: a general repository for interaction datasets. *Nucleic acids research*, 34(suppl 1):D535–D539, 2006.
- [146] Nina Stoletzki and Adam Eyre-Walker. The positive correlation between dn/ds and ds in mammals is due to runs of adjacent substitutions. *Molecular biology and evolution*, 28(4):1371–1380, 2011.

- [147] Aravind Subramanian, Pablo Tamayo, Vamsi K Mootha, Sayan Mukherjee, Benjamin L Ebert, Michael A Gillette, Amanda Paulovich, Scott L Pomeroy, Todd R Golub, Eric S Lander, et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proceedings of the National Academy of Sciences of the United States of America*, 102(43):15545–15550, 2005.
- [148] H Sugihara, V Andrisani, and PM Salvaterra. Drosophila choline acetyltransferase uses a non-aug initiation codon and full length rna is inefficiently translated. *Journal of Biological Chemistry*, 265(35):21714–21719, 1990.
- [149] GATK team. The genome analysis toolkit. Accessed: 2015.
- [150] Stephen T Thibault, Matthew A Singer, Wesley Y Miyazaki, Brett Miliash, Nicholas A Dompe, Carol M Singh, Ross Buchholz, Madelyn Demsky, Robert Fawcett, Helen L Francis-Lang, et al. A complementary transposon tool kit for drosophila melanogaster using p and piggybac. *Nature genetics*, 36(3):283–287, 2004.
- [151] Van Khanh Tran, Yasuhiro Takeshima, Zhujun Zhang, Yasuaki Habara, Kazuhiro Hagiwara, Atsushi Nishiyama, Mariko Yagi, and Masafumi Matsuo. A nonsense mutation-created intraexonic splice site is active in the lymphocytes, but not in the skeletal muscle of a dmd patient. *Human genetics*, 120(5):737–742, 2007.
- [152] Peter J Turnbaugh, Ruth E Ley, Micah Hamady, Claire M Fraser-Liggett, Rob Knight, and Jeffrey I Gordon. The human microbiome project. *Nature*, 449(7164):804–810, 2007.
- [153] Jens Tyedmers, Maria Lucia Madariaga, and Susan Lindquist. Prion switching in response to environmental stress. *PLoS biology*, 6(11):e294, 2008.
- [154] Kavitha Venkatesan, Jean-Francois Rual, Alexei Vazquez, Ulrich Stelzl, Irma Lemmens, Tomoko Hirozane-Kishikawa, Tong Hao, Martina Zenkner, Xiaofeng Xin, Kwang-Il Goh, et al. An empirical framework for binary interactome mapping. *Nature methods*, 6(1):83–90, 2009.
- [155] GEORG HM WAALER. The location of a new second chromosome eye colour gene in drosophila melanogaster. *Hereditas*, 2(3):391–394, 1921.
- [156] Jon Wakefield. Bayes factors for genome-wide association studies: comparison with p-values. *Genetic epidemiology*, 33(1):79–86, 2009.
- [157] K. Wang, M. Li, and H. Hakonarson. Annovar: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic acids research*, 38(16):e164–e164, 2010.

- [158] Lucas D Ward and Manolis Kellis. Haploreg: a resource for exploring chromatin states, conservation, and regulatory motif alterations within sets of genetically linked variants. *Nucleic acids research*, 40(D1):D930–D934, 2012.
- [159] Lucas D Ward and Manolis Kellis. Interpreting noncoding genetic variation in complex traits and human disease. *Nature biotechnology*, 30(11):1095–1106, 2012.
- [160] J.D. Watson and F.H.C. Crick. Molecular structure of nucleic acids. *Nature*, 171(4356):737–738, 1953.
- [161] Martin Weigt, Robert A White, Hendrik Szurmant, James A Hoch, and Terence Hwa. Identification of direct residue contacts in protein–protein interaction by message passing. *Proceedings of the National Academy of Sciences*, 106(1):67–72, 2009.
- [162] David A Wheeler, Maithreyan Srinivasan, Michael Egholm, Yufeng Shen, Lei Chen, Amy McGuire, Wen He, Yi-Ju Chen, Vinod Makijani, G Thomas Roth, et al. The complete genome of an individual by massively parallel dna sequencing. *nature*, 452(7189):872–876, 2008.
- [163] David C Whitcomb, Michael C Gorry, Robert A Preston, William Furey, Michael J Sossenheimer, Charles D Ulrich, Stephen P Martin, Lawrence K Gates, Stephen T Amann, Phillip P Toskes, et al. Hereditary pancreatitis is caused by a mutation in the cationic trypsinogen gene. *Nature genetics*, 14(2):141–145, 1996.
- [164] Reed B Wickner. [ure3] as an altered ure2 protein: evidence for a prion analog in *saccharomyces cerevisiae*. *Science*, 264(5158):566–569, 1994.
- [165] Andrew R Wood, Tonu Esko, Jian Yang, Sailaja Vedantam, Tune H Pers, Stefan Gustafsson, Audrey Y Chu, Karol Estrada, Jian'an Luan, Zoltán Kutalik, et al. Defining the role of common variation in the genomic and biological architecture of adult human height. *Nature genetics*, 46(11):1173–1186, 2014.
- [166] M.C. Wu, S. Lee, T. Cai, Y. Li, M. Boehnke, and X. Lin. Rare-variant association testing for sequencing data with the sequence kernel association test. *The American Journal of Human Genetics*, 2011.
- [167] Yumi Yamaguchi-Kabata, Makoto K Shimada, Yosuke Hayakawa, Shinsei Minoshima, Ranajit Chakraborty, Takashi Gojobori, and Tadashi Imanishi. Distribution and effects of nonsense polymorphisms in human genes. *PloS one*, 3(10):e3393, 2008.

- [168] Jian-Hua Yang, Jun-Hao Li, Peng Shao, Hui Zhou, Yue-Qin Chen, and Liang-Hu Qu. starbase: a database for exploring microrna–mrna interaction maps from argonaute clip-seq and degradome-seq data. *Nucleic acids research*, 39(suppl 1):D202–D209, 2011.
- [169] Ziheng Yang. *Computational molecular evolution*, volume 284. Oxford University Press Oxford, 2006.
- [170] Amelia Younossi-Hartenstein, Patricia Green, Gwo-Jen Liaw, Karen Rudolph, Judith Lengyel, and Volker Hartenstein. Control of early neurogenesis of the drosophilabrain by the head gap genes *stll*, *otd*, *ems*, and *btd*. *Developmental biology*, 182(2):270–283, 1997.
- [171] Yu Zhang and Jun S Liu. Bayesian inference of epistatic interactions in case-control studies. *Nature genetics*, 39(9):1167–1173, 2007.
- [172] Jinying Zhao, Li Jin, and Momiao Xiong. Test for interaction between two unlinked loci. *The American Journal of Human Genetics*, 79(5):831–845, 2006.
- [173] Jesse D Ziebarth, Anindya Bhattacharya, Anlong Chen, and Yan Cui. Polymirts database 2.0: linking polymorphisms in microrna target sites with human diseases and complex traits. *Nucleic acids research*, page gkr1026, 2011.
- [174] O. Zuk, E. Hechter, S.R. Sunyaev, and E.S. Lander. The mystery of missing heritability: Genetic interactions create phantom heritability. *Proceedings of the National Academy of Sciences*, 109(4):1193–1198, 2012.
- [175] Or Zuk, Stephen F Schaffner, Kaitlin Samocha, Ron Do, Eliana Hechter, Sekar Kathiresan, Mark J Daly, Benjamin M Neale, Shamil R Sunyaev, and Eric S Lander. Searching for missing heritability: Designing rare variant association studies. *Proceedings of the National Academy of Sciences*, 111(4):E455–E464, 2014.