

Correlation of SNPs with Phylogenetic Footprints

Claudia Fried^{1,2}, Peter Ahnert³, Peter F. Stadler^{1,2}

¹Bioinformatik, Institut für Informatik, Universität Leipzig, Kreuzstraße 7b, D-04103 Leipzig, Germany. Phone: ++49 341 149 5120; Fax: ++49 341 149 5119; Email: claudia@bioinf.uni-leipzig.de.

²Institut für Theoretische Chemie und Molekulare Strukturbiologie Universität Wien, Währingerstraße 17, A-1090 Wien, Austria

³IKIT/BBZ, Medizinische Fakultät, Universität Leipzig, Johannisallee 30, 04103 Leipzig, Germany

We investigate the relative distribution of single nucleotide polymorphisms (SNPs) in exons, putative regulatory sequences identified by phylogenetic footprinting, and surrounding non-functional DNA. Biases in the available SNP databases causes an overrepresentation of SNPs in the exons. On the other hand there is little difference between putative regulatory and non-functional regions.

1. Introduction

Extensive polymorphism in non-coding gene-regulatory sequences were recently reported in particular for the immune system [12]. This type of genetic variation could therefore be functionally and evolutionarily highly significant. A different pattern of polymorphisms between the coding and non-coding regions seem to distinguish “introvert genes” that code for proteins dealing with self molecules, and “extrovert genes” that are targeted towards foreign molecules that enter the body [11]. Non-coding polymorphism appear to dominate in introvert immune genes.

It is likely that systematic patterns in the distribution of polymorphism between coding, non-coding but regulatory functional, and non-functional sequences can be found also in other classes of protein-coding genes. Here we describe a systematic computational approach to address this problem and discuss its limitations based on the currently available data sources. For concreteness we focus on single nucleotide polymorphisms (SNPs) and use a portfolio of immune system related genes that were of interest in another context as illustrative example.

2. Materials and Methods

Location, validation status and other SNP properties were retrieved from NCBI and EBI databases using the powerful data retrieval tool ENSMART¹ [2, 3], which was also used to localize exons and introns consistently with the SNP localization. The SNP data in dbSNP² come from large sequencing efforts and from a variety of laboratories involved in SNP discovery. The quality of the data depends on the SNP detection method and on independent verification. We distinguish validated and non-validated SNPs in our analysis.

¹freely available on the world wide web at <http://www.ensembl.org/Ensmart/>

²URL: <http://www.ncbi.nlm.nih.gov/SNP/index.html>

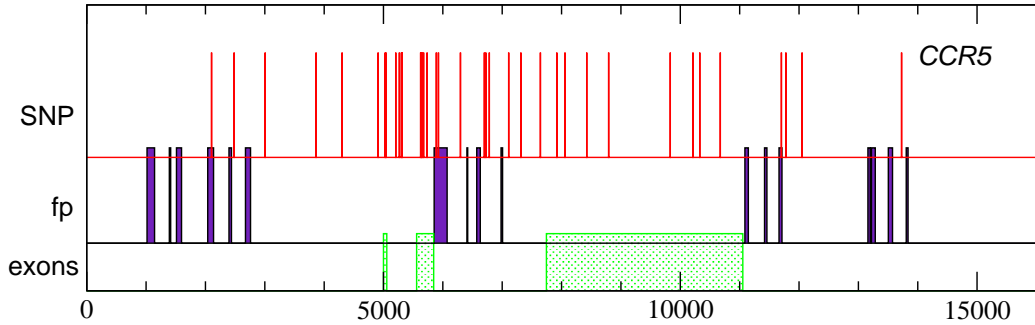


Figure 1. Graphical representation of the SNP distribution of the CCR5 gene. Accession number: NT_0058250, the part of the sequence that we used in our analysis spans the entire gene and surrounding region 5000 basepairs upstream and downstream of the gene

DNA sequences that are regulatory active presumably cover a substantial part of the intergenic regions. On the other hand, only a small number of transcription factor binding sites or promoter sequences are experimentally verified for any given gene. Functional non-coding sequences, however, evolve much slower than the surrounding non-functional DNA because they are subject to stabilizing selection. This is true in particular in vertebrates [4], while invertebrates often show a high rate of binding site turnover. Comparative sequence analysis can therefore be used to detect functional non-coding DNA sequences in the vicinity of the genes of interest [15, 7, 9]. This technique is known as *phylogenetic footprinting*.

Recently we have presented the program *tracker* as an efficient tool for surveying phylogenetic footprints in large datasets [14]. It is based on the initial computation of all pairwise blast alignments from the input sequences with a non-restrictive parameter setting. A hierarchy of filtering steps then removes insignificant matches. The most complicated step is the combination of overlapping alignments to maximal cliques of mutually consistent alignments, thereby producing local multiple alignments of the conserved regions. For details we refer to [14].

For the purpose of the present study we use orthologous sequences retrieved from the genome databases for *homo sapiens*, *mus musculus*, and *rattus norvegicus*. In some cases only the rat or the mouse sequence was available. In all cases the DNA sequence extending 5000nt upstream and downstream of the gene was retrieved.

The genes analyzed here were originally selected for a genotype-phenotype association study of rheumatoid arthritis susceptibility. Most of the selected genes play a role in cytokine balance, which has been suggested to play a significant role in rheumatoid arthritis [10]. Other genes were selected for their role in cartilage metabolism [13] or their known or suggested association with rheumatoid arthritis in other study populations [1, 8].

3. First Results

Figure 1 visualizes a typical data set. We distinguish between SNPs in the exons, in phylogenetic footprints, and in the remaining DNA that presumably is non-functional. Intuitively one would expect that the rate of occurrence of SNPs in the non-functional DNA is largest since there it is not subject to selection, while mutations should be selected against in coding and regulatory sequences.

Table 1 summarizes the distribution of SNPs in 20 genes. We observe a rather large variance in the density of SNPs around the mean of about 2 SNPs per 1000nt.

Table 1. Distribution of SNPs in functional and non-functional sequences of several immune genes. Overrepresentation of SNPs is indicated by \uparrow , underrepresentation by \downarrow . Fisher’s exact test was used to assess the significance level of over- and underrepresentation in the columns labelled FT. Recall that confidence increases with decreasing p -values: $0.05 \leq p < 10^{-3}$ is shown by a single arrow \uparrow , double arrows $\uparrow\uparrow$ indicate $10^{-3} \leq p < 10^{-5}$ and $p < 10^{-5}$ is shown by triple arrows $\uparrow\uparrow\uparrow$ and $\downarrow\downarrow\downarrow$. The column ρ lists the average number of SNPs per 1000nt.

Gene	PF		SNP		Exons			Footprints			other	
	#	#		ρ	length	SNP	FT	length	SNP	FT	length	SNP
CCR5	17	38		2.36	3657	15	\uparrow	1144	4		11253	19
CD8A	74	23		1.40	2202	11	$\uparrow\uparrow$	4435	4		9763	8
CSF2	55	51		4.12	765	6		1528	9		10066	36
IFNA1	4	17		1.56	878	5	$\uparrow\uparrow\uparrow$	467	1		9533	11
IFNG	92	47		3.13	1211	4		8227	33	\uparrow	5535	10
IL13	39	53		4.09	1283	6		8867	14	$\downarrow\downarrow\downarrow$	2788	33
IL18	107	48		1.55	1153	4		549	2		29171	42
IL2	38	31		2.06	791	4		3911	7		10316	20
IL2RA	261	69		1.13	2322	3		27220	14	$\downarrow\downarrow$	31241	52
IL4	105	70		3.74	618	2		8613	14	$\downarrow\downarrow\downarrow$	9462	54
IL5	51	25		2.06	816	5	\uparrow	1616	1		9647	19
IL6	57	69		4.66	1130	8		1630	14	\uparrow	12042	47
IL6R	241	21		0.29	3315	7	$\uparrow\uparrow$	6777	1		61567	13
IRF1	109	29		1.60	2067	2		9231	9	\downarrow	6786	18
LTA	24	0		—	1424	0		916	0		9704	0
LTB	39	2		0.16	900	0		1206	0		9768	2
MIF	32	28		2.58	563	5	\uparrow	615	1		9669	22
MMP3	52	83		4.66	1814	10		6668	21	\downarrow	9319	52
PRG4	163	24		0.86	5011	4		4481	9		18207	11
PRL	69	30		1.48	1359	3		8911	2	$\downarrow\downarrow\downarrow$	9986	25
Total	1629	820		1.92	33279	149	$\uparrow\uparrow\uparrow$	107012	194		285823	477
Frac.					0.0781	0.1817		0.2366	0.2362		0.6708	0.5817
Fisher’s p -value	SNP overrepresented				8.40×10^{-22}			0.1176				
	SNP underrepresented				1			0.845				

The most surprising result is that SNPs are *overrepresented* in exons. We use Fisher’s exact test³ [5, 6] to assess the statistical significance of the differences in SNP density in exons and footprints compared to the remaining non-functional DNA. Significant overrepresentation is found only for some genes (CCR5, CD8A, IFNA1, IL5, IL6R, and MIF). A p -value of $p = 0.05$ or smaller indicates a significant difference in the distribution.

Almost certainly this effect is a bias in the SNP database: If the SNP overrepresentation in the exons of some immune genes were real we would have to postulate an increased mutation rate in the coding sequences compared to non-functional surrounding DNA. Since SNPs are often obtained from cDNA comparison more data are available for the exons than for the surrounding non-coding DNA.

³Computations were performed using the webservice <http://www.matforsk.no/ola/fisher.htm>.

The distribution of SNPs in tracker-predicted phylogenetic footprints does not significantly differ from distribution in non-functional DNA on average. For some genes, however, we find a highly significant underrepresentation of SNPs in the putative regulatory sequences (IL13, IL2RA, IL4, PRL, MMP3, IRF1). As methods for SNP detection do not distinguish between functional and non-functional non-coding DNA we should expect that differences in SNP distribution between phylogenetic footprints and non-functional DNA are biological rather than caused by database biases.

Interestingly, the picture does not change when only the (small) subset of validated SNPs is used instead of all database entries.

Discussion

The human genome program has enabled and spawned many studies attempting to identify the role of genetic variation in the etiology and pathogenesis of many diseases. Of special interest to us are complex diseases like rheumatoid arthritis where a number of genetic factors are thought to be involved.

One approach to identify genes associated with a disease is to select candidate genes according to their known or assumed role in biological processes likely to be involved in the disease. The association of variants of these genes with the disease are then tested by genotyping a number of SNPs distributed along their genomic DNA. These SNPs either modify the function of the gene directly or are in linkage with such a function-modifying polymorphism. Clearly, SNPs located in the coding sequence or in phylogenetic footprints of the gene are most likely to be function-modifying and are therefore the best candidates for further experimental study.

Conversely, the distribution of SNPs in exons, regulatory elements, and non-functional background potentially provides direct information on selection pressures acting on various components of a gene or gene cluster. In order to access this information unbiased subsets of SNP databases will have to be extracted e.g. by removing all entries that are obtained from mRNAs.

In this short contribution we have shown that the statistical analysis of the relationships between polymorphisms in the human genome and functional DNA regions is computationally feasible. Biologically interesting results, however, will have to await more extensive surveys that include larger sets of genes. This will allow the application of clustering techniques to identify protein families whose genes have particular SNP distributions among different types of surrounding non-coding DNA. In particular, a more fine-grained analysis will be of interest that treats GpC islands, proximal promoters, transcription factor binding sites, introns, etc., as separate classes of DNA sequence.

Acknowledgements. This work is supported in part by the *DFG* Bioinformatics Initiative.

References

- [1] C. G. Baerwald, C. Mok, M. Tickly, C. S. Lau, B. P. Wordsworth, B. Ollier, G. S. Panayi, and J. S. Lanchbury. Corticotropin releasing hormone (CRH) promoter polymorphisms in various ethnic groups of patients with rheumatoid arthritis. *Z. Rheumatol.*, 59:29–34, 2000.
- [2] C. Brooksbank, E. Camon, M. A. Harris, M. Magrane, M. J. Martin, N. Mulder, C. O'Donovan, H. Parkinson, M. A. Tuli, R. Apweiler, E. Birney, A. Brazma, K. Henrick, R. Lopez, G. Stoesser, P. Stoehr, and G. Cameron. The european bioinformatics institute's data resources. *Nucl. Acids Res.*, 31:43–50, 2003.
- [3] M. Clamp, D. Andrews, D. Barker, P. Bevan, G. Cameron, Y. Chen, L. Clark, T. Cox, J. Cuff, V. Curwen, T. Down, R. Durbin, E. Eyra, J. Gilbert, M. Hammond, T. Hubbard, A. Kasprzyk, D. Keefe, H. Lehvaslaiho, V. Iyer, C. Melsopp, E. Mongin, R. Pettett, S. Potter, A. Rust, E. Schmidt,

- S. Searle, G. Slater, J. Smith, W. Spooner, A. Stabenau, J. Stalker, E. Stupka, A. Ureta-Vidal, I. Vastrik, and B. E. Ensembl 2002: accommodating comparative genomics. *Nucl. Acids Res.*, 31:38–42, 2003.
- [4] J. W. Fickett and W. W. Wasserman. Discovery and modeling of transcriptional regulatory regions. *Current Opinion in Biotech.*, 11:19–24, 2000.
- [5] R. A. Fisher. The logic of inductive inference. *Journal of the Royal Statistical Society Series A*, 98:39–54, 1935.
- [6] R. A. Fisher. Confidence limits for a cross-product ratio. *Australian Journal of Statistics*, 4:41, 1962.
- [7] D. L. Gumucio, D. A. Shelton, W. Zhu, D. Millinoff, T. Gray, J. H. Bock, J. L. Slightom, and M. Goodman. Evolutionary strategies for the elucidation of cis and trans factors that regulate the developmental switching programs of the beta-like globin genes. *Mol. Phylogenet. Evol.*, 5:18–32, 1996.
- [8] S. John, A. Myerscough, S. Eyre, P. Roby, A. Hajeer, A. J. Silman, W. E. Ollier, and J. Worthington. Linkage of a marker in intron D of the estrogen synthase locus to rheumatoid arthritis. *Arthritis Rheum.*, 42:1617–1620, 1999.
- [9] J. Y. Leung, F. E. McKenzie, A. M. Uglialoro, P. O. Flores-Villanueva, B. C. Sorkin, E. J. Yunis, D. L. Hartl, and A. E. Goldfeld. Identification of phylogenetic footprints in primate tumor necrosis factor- α promoters. *Proc. Natl. Acad. Sci. USA*, 97:6614–6618, 2000.
- [10] P. Miossec. Pro- and antiinflammatory cytokine balance in rheumatoid arthritis. *Clin. Exp. Rheumatol.*, 13:S13–S16, 1995.
- [11] A. Mitchison. Partitioning of genetic variation between regulatory and coding gene segments: the predominance of software variation in genes encoding introvert proteins. *Immunogenetics*, 46:46–52, 1997.
- [12] N. A. Mitchison. Polymorphism in regulatory gene sequences. *Genome Biology*, 2:2001.1–2001.6, 2000.
- [13] G. Murphy, V. Knauper, S. Atkinson, G. Butler, W. English, M. Hutton, J. Stracke, and I. Clark. Matrix metalloproteinases in arthritic disease. *Arthritis Res.*, 4:S39–S49, 2002.
- [14] S. J. Prohaska, C. Fried, C. Flamm, G. P. Wagner, and P. F. Stadler. Surveying phylogenetic footprints in large gene clusters: Applications to *Hox* cluster duplications. *Mol. Phyl. Evol.*, 2003. submitted; SFI preprint #03-02-011.
- [15] D. A. Tagle, B. F. Koop, M. Goodman, J. L. Slightom, D. L. Hess, and R. T. Jones. Embryonic ϵ and γ globin genes of a prosimian primate (*galago crassicaudatus*). Nucleotide and amino acid sequences, developmental regulation and phylogenetic footprints. *J. Mol. Biol.*, 203:439–455, 1988.