

Gene finding using multiple related species: a classification approach

Manolis Kellis

MIT Broad Institute for Biomedical Research, Cambridge, MA, USA

1. Introduction

Ideally, we should be able to systematically discover all the functional genes in a newly sequenced genome from its sequence alone. Computational discovery methods rely both on the direct signals used by the cell to guide transcription, splicing, and translation, and also on indirect signals such as evolutionary conservation. In this paper, we summarize the principles of a classification-based approach for systematic gene identification, on the basis of comparative sequence information from multiple, closely related species. We first frame gene identification as a classification problem of distinguishing real genes from spurious gene predictions. We then present the Reading Frame Conservation (RFC) test, a new computational method implementing such a classification approach, on the basis of the patterns of nucleotide changes in the alignment of orthologous regions. We finally summarize our results of applying this method to reannotate the yeast genome, and the challenges of using related methods to discover all functional genes in the human genome.

2. Defining real genes

What is a real gene? This question is relatively easy to answer for those genes that are abundantly expressed, encode well-characterized proteins, and whose disruption affects a specific function in the organism. Beyond these, the distinction is much more subtle between functional genes and spurious gene predictions.

Experimentally, the definition of a functional gene comes largely as accumulating evidence of its usage, including gene expression (Velculescu, 1997), protein fragments (Rezaul *et al.*, 2005), biochemical function (Jackman *et al.*, 2003), protein interactions (Bai and Elledge, 1997), or the effect of its disruption (McAlister and Holland, 1982). It is important to note, however, that absence of experimental evidence does not imply that a gene annotation is spurious: a real

2 Gene Finding and Gene Structure

gene may be missing experimental evidence because it is not used in the particular conditions surveyed. Conversely, any individual report of gene usage could be due to experimental noise, cross-hybridization, or chance transcription events due to a basal level of intergenic transcription. Thus, experimental evidence alone is insufficient to distinguish between real genes and spurious gene predictions.

Computationally, the processes of transcription, splicing, and translation can be thought of as a series of decisions taken by the cell, based on signals in the genome. These signals include distant enhancer elements, regulatory elements surrounding the transcription start site, splicing enhancer and repressor signals in the message, and translation signals in the sequence or structure of mature mRNAs (Fairbrother *et al.*, 2002; Wang, 2004; Thanaraj and Robinson, 2000; Yeo and Burge, 2004). The subset of these signals that we currently understand is insufficient to specify the set of known genes. Thus, in addition to the *direct* signals used by the cell, gene identification methods (Burge and Karlin, 1997; Majoros *et al.*, 2004; Kulp *et al.*, 1996; Krogh, 1997; Henderson *et al.*, 1997; Stanke and Waack, 2003; Salzberg *et al.*, 1998) routinely identify genes using additional *indirect* signals (reviewed in Fickett and Tung, 1992), which are not available to the cell, although they are generally good indicators of protein-coding genes. These include the frequency of each codon in protein-coding regions, the overall length of the translated protein product, and importantly, the evolutionary conservation of protein sequences across related species.

Evolutionary conservation is perhaps the strongest indicator that a predicted gene is functional. A gene that confers even a slight evolutionary advantage can be conserved across millions of years, regardless of the rarity of its usage. Hence, even if experimental methods fail to detect the usage of a gene in a given set of experimental conditions, evolutionary methods are able to detect *indirect* evidence of its usage, by observing the pressure to preserve its function over millions of years.

3. Gene finding as a classification problem

The challenge of using evolutionary conservation to define genes lies in the ability to reject spurious genes. Although it is generally well understood that genes conserved across large evolutionary distances are functional, lack of conservation is generally attributed to evolutionary divergence, thus providing no evidence toward either accepting or rejecting a gene prediction.

Yet, the comparison of related genomes contains information much richer than simply the presence or absence of a protein in a genome. By working with closely related genomes, we can define regions of conserved gene order, or *synteny blocks*, which span several genes, and potentially entire chromosomes. Defining conserved synteny blocks allows us to construct global alignments, spanning both well-conserved regions and regions of low conservation. In particular, these give us access to full nucleotide alignments for all predicted genes and for all intergenic regions within orthologous segments.

With the availability of orthologous alignments for both protein-coding and non-coding regions, we can study their distinct properties and build a classifier between the two types of region. The simplest and most commonly used such classifier

observes the overall level of nucleotide conservation in the alignment, and selects regions of high conservation that are likely to be functional. Other classifiers may observe more subtle signals, such as the number of amino acid changes per nucleotide substitution (Hurst, 2002), the frequency of insertions and deletions, the periodicity of mutations, and so on. By working with properties of the nucleotide alignment, rather than protein alignment, we are able to apply it uniformly to evaluate both genes and intergenic regions in the same test.

One particular conservation property unique to protein-coding segments is the pressure to preserve the reading frame of translation. Since protein sequences are translated every three nucleotides, the length of insertions and deletions (indels) is largely constrained to remain a multiple of three, thus preserving the frame of translation. Within coding regions, indels that disrupt the frame of translation are excluded, or compensated with nearby indels that restore the reading frame. In noncoding regions, the length of indels does not have this constraint, and short spacing changes are tolerated.

To evaluate this property quantitatively, we developed the RFC test. This test evaluates the pressure to preserve the reading frame in a fully aligned interval, by measuring the portion of nucleotides in this interval for which the reading frame has been locally conserved (Kellis *et al.*, 2004). The RFC test provides a classifier between coding and noncoding regions that is completely independent of the traditional signals used in gene identification. It does not rely on start, stop, or splicing signals, nor does it rely on the conservation of protein sequence. It can, therefore, be combined with existing gene-finding tools and provide a highly informative score for any interval considered.

4. Reannotation of yeast

A classification approach is particularly well suited for the yeast genome. The general scarcity of introns makes it hard to rely on splicing signals to discover genes. Thus, the annotation of *Saccharomyces cerevisiae* has traditionally relied solely on the length of predicted proteins to annotate genes, resulting in 6062 annotated open reading frames (ORFs), which potentially encode proteins of at least 100 amino acids (aa). Additionally, a tentative functional annotation has been inferred for as many as 3966 of these ORFs, on the basis of classical genetic experiments and systematic genome-wide studies of gene expression, deletion phenotype, and protein–protein interaction. Together, the interval-based annotation and the large set of well-known genes make it possible to apply the RFC test systematically to evaluate the functional significance of the remaining ORFs.

To apply the test, we constructed multiple sequence alignments for every ORF and every intergenic region across four closely related species. We sequenced and assembled the genomes of *S. paradoxus*, *S. mikatae*, and *S. bayanus* (Kellis *et al.*, 2003), and defined genome-wide synteny blocks with *S. cerevisiae*, on the basis of discrete anchors provided by unique protein blast hits (Kellis *et al.*, 2004). Within these synteny blocks, we constructed global alignments of orthologous genes and intergenic regions across the four species using CLUSTALW (Higgins and Sharp, 1988), and systematically evaluated each alignment. We compared *S. cerevisiae* to

each species in turn, and every comparison cast a vote based on its overall RFC and a species-specific cutoff. A decision was then reached for each gene by tallying the votes from all comparisons.

We evaluated the sensitivity and specificity of the approach based on the 3966 genes with functional annotations (“known genes”), and 340 randomly chosen intergenic sequences with lengths similar to the annotated ORFs. The RFC test correctly accepted 3951 known genes (99.6%) and rejected only 15 known genes; upon manual inspection, these 15 are indeed likely to be spurious (most lack experimental evidence, and deletion phenotypes of the rest are likely due to their overlap with the promoter of other known genes). The method also correctly rejected 326 intergenic regions (96%), accepting only 14 intergenic regions (of which 10 appear to define short ORFs or extend annotated ORFs, suggesting that at most 1% of true intergenic regions failed to be rejected by the RFC test). In summary, the RFC test shows a very strong discrimination between genes and intergenic regions, with sensitivity and specificity values greater than 99%.

We then applied the RFC test to all previously annotated genes, leading to a major revisiting of the yeast gene set. For ORFs with lengths greater than 100 aa, our analysis accepted 5538 ORFs and rejected 503 (of which 376 were immediately rejected, 105 were rejected with additional criteria, and 32 were merged with neighboring ORFs); the classifier abstained from making a decision in 20 cases. The rejected ORFs show an abundance of frame-shifting indels across their entire length, in-frame stop codons, and low conservation of protein sequence, in addition to the low RFC score; their length distribution and atypical codon usage additionally suggest that they are likely occurring by chance (Goffeau, 1996; Dujon, 1994; Sharp and Li, 1987); furthermore, previous systematic experimentation showed no compelling evidence that these may encode a functional gene. Thus, it appears that more than 500 previously annotated ORFs in the yeast genome are spurious predictions.

Below the 100 amino acid cutoff, no previous systematic annotation or experimentation was available. Thus, to validate our method, we compared the results of the RFC test with an independent metric: the proportion of the *S. cerevisiae* ORF that was also free of stop codons in the other three species. Between 50 and 99 aa, we found that the method is still reliable and reports 43 candidate new genes at that length. As ORFs smaller than 50 aa were tested, we found that the specificity of the test decreased, since small intervals tend to be devoid of indels by chance rather than by presence of selective pressure. Thus, additional constraints, and additional species, will be needed to discover genes reliably at such short lengths.

In addition to the discovery of genes themselves, the comparative analysis allowed us to refine gene boundaries. Once an interval was determined to be under selective pressure for RFC, the boundaries of that interval were adjusted on the basis of the conservation of start/stop/splice signals and the boundaries over which the reading frame is conserved. This led to a large-scale reannotation of gene structure, which affected hundreds of genes (146 start codon changes, 67 intron changes, 32 merges of consecutive ORFs, and 45 changes of ORF ends). It is worth noting that in 134 cases, the inferred boundary changes pinpointed sequencing errors in the primary sequence of *S. cerevisiae*, ~50 of which were tested and

corrected by resequencing. These boundary changes reveal the true location of promoter regions, new protein domains in elongated ORFs, and previously overlooked functional relationships in the case of merged ORFs

In summary, the RFC test was able to reliably distinguish between genes and intergenic regions, leading to a systematic reannotation of the yeast genome. The comparative analysis led to the rejection of more than 500 previously annotated genes, and to the discovery of many novel genes. The results agree with similar comparative analyses carried out from a number of yeast species (Cliften *et al.*, 2003; Blandin, 2000; Wood *et al.*, 2001; Brachat, 2003). In addition, it allowed us to refine the gene structure of hundreds of genes, adjusting start and stop boundaries, merging consecutive ORFs, and discovering many new introns. Moreover, by using multiple species, the signals leading to gene identification were powerful enough to pinpoint sequencing errors in any individual species, including *S. cerevisiae*.

5. Implications for the human genome

The challenge of gene finding in the human genome is far greater than in yeast, due to the vastly larger intergenic regions, numerous and large introns, small exons, and alternatively spliced genes. Most approaches to gene finding in higher eukaryotes have relied on Hidden Markov Models (Burge and Karlin, 1997; Majoros *et al.*, 2004; Kulp *et al.*, 1996; Krogh, 1997; Henderson *et al.*, 1997; Stanke and Waack, 2003), which inherently emphasize the importance of exon chaining and rely on knowledge of the expected length distributions for both exons and introns. These have been recently extended to use sequences from multiple species (Yeh *et al.*, 2001; Siepel and Haussler, 2004a; Siepel and Haussler, 2004b; Batzoglou *et al.*, 2000; Korf *et al.*, 2001; Dewey *et al.*, 2004; Rinner and Morgenstern, 2002; Parra, 2003; Meyer and Durbin, 2002; Novichkov *et al.*, 2001). These approaches have limitations however, in cases of alternative splicing (Brett, 2000; Mironov *et al.*, 1999; Cawley and Pachter, 2003), differences in splicing between species (Nurtdinov *et al.*, 2003), widely varying exon and intron lengths, and noncanonical splice sites.

An exon-based classification approach to gene finding can help overcome these limitations. By exhaustively enumerating and testing all candidate protein-coding intervals, classification approaches make exon detection independent of chaining. Splice site models can be applied to each species independently and then compared with each other, which is much stronger than the evidence in any one species alone. The intervals can be then tested on the basis of discriminating variables that distinguish genes from noncoding regions (Fickett and Tung, 1992; Goldman and Yang, 1994). New discriminating variables can be defined as alignments of each type of region are systematically compared, and these can be combined and weighted, leveraging traditional machine-learning techniques for feature selection and classification (Moore and Lake, 2003). Once relevant intervals have been identified, their boundaries can be adjusted for optimal chaining into complete genes. In particular, chaining can also leverage an inferred frame of translation for each exon, based on the higher mutation rate observed in largely degenerate third codon positions (Hurst, 2002). The exon-chaining step produces full gene models,

and is able to cope with alternative splicing and missing exons, since it is not constrained to a single optimal path through the exons.

By systematically observing alignment properties of large sequence regions, we can build new rigorous approaches for sequence analysis. These are widely applicable beyond coding exons, and similar classification-based approaches can be used to distinguish CpG islands, 5'- and 3'-untranslated regions, promoter regions, regulatory islands, and other functional elements. Through the lens of evolutionary selection, our ability to directly interpret genomes is revolutionized. Coupled with systematic experimentation and validation, these analyses can lead to a systematic catalog of functional elements in the human genome, forming the future foundations of biomedical research.

References

- Bai C and Elledge SJ (1997) Gene identification using the yeast two-hybrid system. *Methods in Enzymology*, **283**, 141–156.
- Batzoglou S, Pachter L, Mesirov JP, Berger B and Lander ES (2000) Human and mouse gene structure: Comparative analysis and application to exon prediction. *Genome Research*, **10**, 950–958.
- Blandin G, Durrens P, Tekaiia F, Aigle M, Bolotin-Fukuhara M, Bon E, Casaregola S, de Montigny J, Gaillardin C, Lepingle A, *et al.* (2000) Genomic exploration of the hemiascomycetous yeasts: 4. The genome of *Saccharomyces cerevisiae* revisited. *FEBS Letters*, **487**, 31–36.
- Brachat S, Dietrich FS, Voegeli S, Zhang Z, Stuart L, Lerch A, Gates K, Gaffney T and Philippsen P (2003) Reinvestigation of the *Saccharomyces cerevisiae* genome annotation by comparison to the genome of a related fungus: *Ashbya gossypii*. *Genome Biology*, **4**, R45.
- Brett D, Hanke J, Lehmann G, Haase S, Delbruck S, Krueger S, Reich J and Bork P (2000) EST comparison indicates 38% of human mRNAs contain possible alternative splice forms. *FEBS Letters*, **474**, 83–86.
- Burge C and Karlin S (1997) Prediction of complete gene structures in human genomic DNA. *Journal of Molecular Biology*, **268**, 78–94.
- Cawley SL and Pachter L (2003) HMM sampling and applications to gene finding and alternative splicing. *Bioinformatics*, **19**(Suppl 2), II36–II41.
- Cliften P, Sudarsanam P, Desikan A, Fulton L, Fulton B, Majors J, Waterston R, Cohen BA and Johnston M (2003) Finding functional features in *Saccharomyces* genomes by phylogenetic footprinting. *Science*, **301**, 71–76.
- Dewey C, Wu JQ, Cawley S, Alexandersson M, Gibbs R and Pachter L (2004) Accurate identification of novel human genes through simultaneous gene prediction in human, mouse, and rat. *Genome Research*, **14**, 661–664.
- Dujon B, Alexandraki D, Andre B, Ansorge W, Baladron V, Ballesta JP, Banrevi A, Bolle PA, Bolotin-Fukuhara M and Bossier P (1994) Complete DNA sequence of yeast chromosome XI. *Nature*, **369**, 371–378.
- Fairbrother WG, Yeh RF, Sharp PA and Burge CB (2002) Predictive identification of exonic splicing enhancers in human genes. *Science*, **297**, 1007–1013.
- Fickett JW and Tung CS (1992) Assessment of protein coding measures. *Nucleic Acids Research*, **20**, 6441–6450.
- Goffeau A, Barrell BG, Bussey H, Davis RW, Dujon B, Feldmann H, Galibert F, Hoheisel JD, Jacq C, Johnston M, *et al.* (1996) Life with 6000 genes. *Science*, **274**, 546, 563–567.
- Goldman N and Yang Z (1994) A codon-based model of nucleotide substitution for protein-coding DNA sequences. *Molecular Biology and Evolution*, **11**, 725–736.
- Henderson J, Salzberg S and Fasman KH (1997) Finding genes in DNA with a hidden Markov model. *Journal of Computational Biology*, **4**, 127–141.

- Higgins DG and Sharp PM (1988) CLUSTAL: A package for performing multiple sequence alignment on a microcomputer. *Gene*, **73**, 237–244.
- Hurst LD (2002) The Ka/Ks ratio: Diagnosing the form of sequence evolution. *Trends in Genetics*, **18**, 486.
- Jackman JE, Montange RK, Malik HS and Phizicky EM (2003) Identification of the yeast gene encoding the tRNA m1G methyltransferase responsible for modification at position 9. *RNA*, **9**, 574–585.
- Kellis M, Patterson N, Birren B, Berger B and Lander ES (2004) Methods in comparative genomics: Genome correspondence, gene identification and regulatory motif discovery. *Journal of Computational Biology*, **11**, 319–355.
- Kellis M, Patterson N, Endrizzi M, Birren B and Lander ES (2003) Sequencing and comparison of yeast species to identify genes and regulatory elements. *Nature*, **423**, 241–254.
- Korf I, Flicek P, Duan D and Brent MR (2001) Integrating genomic homology into gene structure prediction. *Bioinformatics*, **17**(Suppl 1), S140–S148.
- Krogh A (1997) Two methods for improving performance of an HMM and their application for gene finding. *Proceedings/International Conference on Intelligent Systems for Molecular Biology*, **5**, 179–186.
- Kulp D, Haussler D, Reese MG and Eeckman FH (1996) A generalized hidden Markov model for the recognition of human genes in DNA. *Proceedings/International Conference on Intelligent Systems for Molecular Biology*, **4**, 134–142.
- Majoros WH, Pertea M and Salzberg SL (2004) TigrScan and GlimmerHMM: Two open source ab initio eukaryotic gene-finders. *Bioinformatics*, **20**, 2878–2879.
- McAlister L and Holland MJ (1982) Targeted deletion of a yeast enolase structural gene. Identification and isolation of yeast enolase isozymes. *The Journal of Biological Chemistry*, **257**, 7181–7188.
- Meyer IM and Durbin R (2002) Comparative *ab initio* prediction of gene structures using pair HMMs. *Bioinformatics*, **18**, 1309–1318.
- Mironov AA, Fickett JW and Gelfand MS (1999) Frequent alternative splicing of human genes. *Genome Research*, **9**, 1288–1293.
- Moore JE and Lake JA (2003) Gene structure prediction in syntenic DNA segments. *Nucleic Acids Research*, **31**, 7271–7279.
- Novichkov PS, Gelfand MS and Mironov AA (2001) Gene recognition in eukaryotic DNA by comparison of genomic sequences. *Bioinformatics*, **17**, 1011–1018.
- Nurtdinov RN, Artamonova II, Mironov AA and Gelfand MS (2003) Low conservation of alternative splicing patterns in the human and mouse genomes. *Human Molecular Genetics*, **12**, 1313–1320.
- Parra G, Agarwal P, Abril JF, Wiehe T, Fickett JW and Guigo R (2003) Comparative gene prediction in human and mouse. *Genome Research*, **13**, 108–117.
- Rezaul K, Wu L, Mayya V, Hwang SI and Han DK (2005) A systematic characterization of mitochondrial proteome from human T leukemia cells. *Molecular and Cellular Proteomics*, **4**, 169–181.
- Rinner O and Morgenstern B (2002) AGenDA: Gene prediction by comparative sequence analysis. *In Silico Biology*, **2**, 195–205.
- Salzberg S, Delcher AL, Fasman KH and Henderson J (1998) A decision tree system for finding genes in DNA. *Journal of Computational Biology*, **5**, 667–680.
- Sharp PM and Li WH (1987) The codon Adaptation Index—a measure of directional synonymous codon usage bias, and its potential applications. *Nucleic Acids Research*, **15**, 1281–1295.
- Siepel A and Haussler D (2004a) Combining phylogenetic and hidden Markov models in biosequence analysis. *Journal of Computational Biology*, **11**, 413–428.
- Siepel A and Haussler D (2004b) Phylogenetic estimation of context-dependent substitution rates by maximum likelihood. *Molecular Biology and Evolution*, **21**, 468–488.
- Stanke M and Waack S (2003) Gene prediction with a hidden Markov model and a new intron submodel. *Bioinformatics*, **19**(Suppl 2), II215–II225.

8 Gene Finding and Gene Structure

- Thanaraj TA and Robinson AJ (2000) Prediction of exact boundaries of exons. *Briefings in Bioinformatic*, **1**, 343–356.
- Velculescu VE, Zhang L, Zhou W, Vogelstein J, Basrai MA, Bassett DE Jr, Hieter P, Vogelstein B and Kinzler KW (1997) Characterization of the yeast transcriptome. *Cell*, **88**, 243–251.
- Wang Z, Rolish ME, Yeo G, Tung V, Mawson M and Burge CB (2004) Systematic identification and analysis of exonic splicing silencers. *Cell*, **119**, 831–845.
- Wood V, Rutherford KM, Ivens A, Rajandream M-A and Barrell B (2001) A re-annotation of the *Saccharomyces cerevisiae* genome. *Comparative and Functional Genomics*, **2**, 143–154.
- Yeh RF, Lim LP and Burge CB (2001) Computational inference of homologous gene structures in the human genome. *Genome Research*, **11**, 803–816.
- Yeo G and Burge CB (2004) Maximum entropy modeling of short sequence motifs with applications to RNA splicing signals. *Journal of Computational Biology*, **11**, 377–394.