

A Genome-wide Study of Dual Coding Regions in Human Alternatively Spliced Genes

Han Liang¹ and Laura F. Landweber²
Department of Chemistry¹ and Ecology & Evolutionary Biology²
Princeton University, Princeton, NJ 08544, USA

Supplementary File 2

Fig. S1

Characterization of alternatively spliced genes with multiple coding regions.

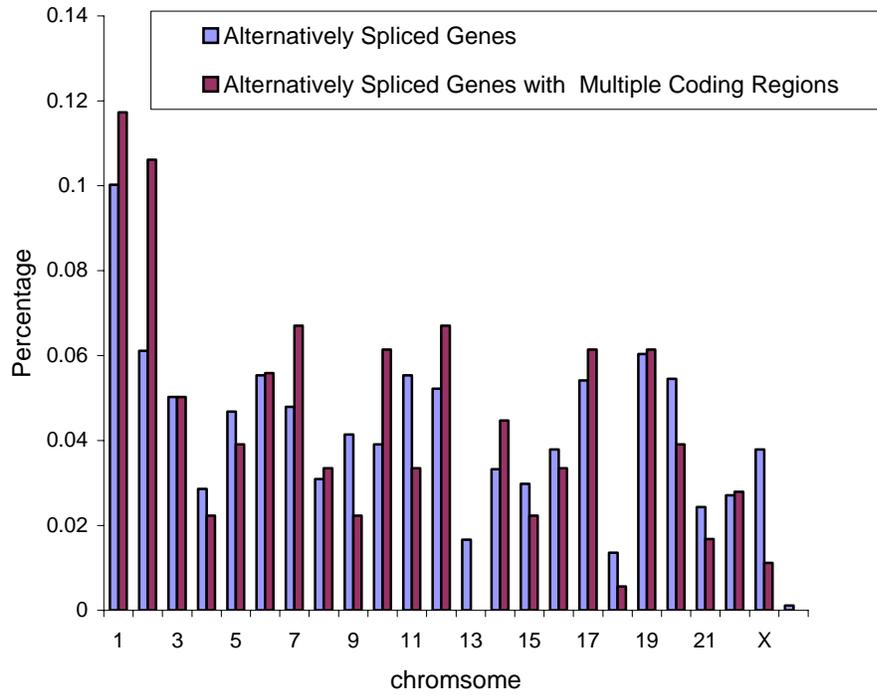
(a) The distribution over different chromosomes. The blue bars represent the percentage of alternatively spliced genes in each chromosome; the red bars represent the percentage of alternatively spliced genes with multiple coding regions in each chromosome.

(b) The distribution in different biological processes. The blue bars represent the percentage of alternatively spliced genes in each process; the red bars represent the percentage of alternatively spliced genes with multiple coding regions in each process.

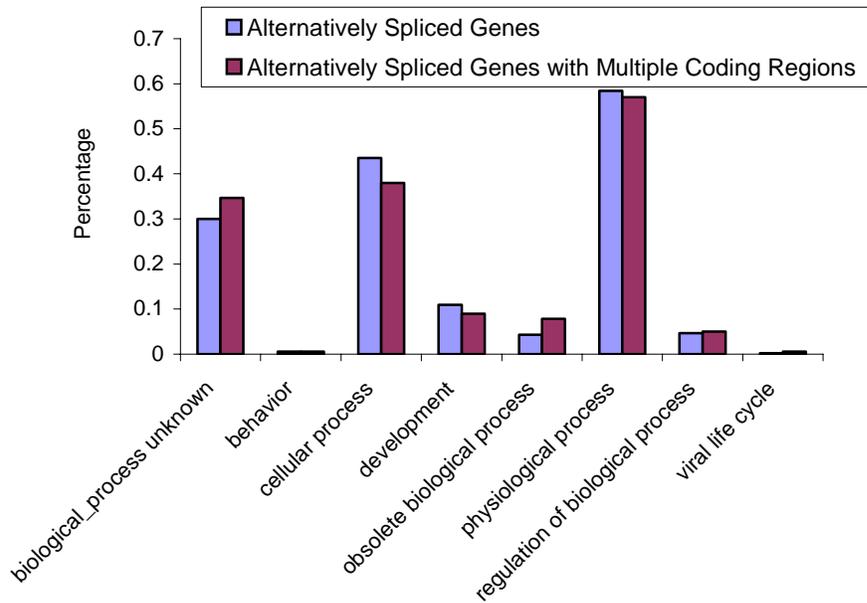
(c) The distribution in different biological function categories. The blue bars represent the percentage of alternatively spliced genes in each category; the red bars represent the percentage of alternatively spliced genes with multiple coding regions in each category.

(d) The distribution in different biological components. The blue bars represent the percentage of alternatively spliced genes in each component; the red bars represent the percentage of alternatively spliced genes with multiple coding regions in each component.

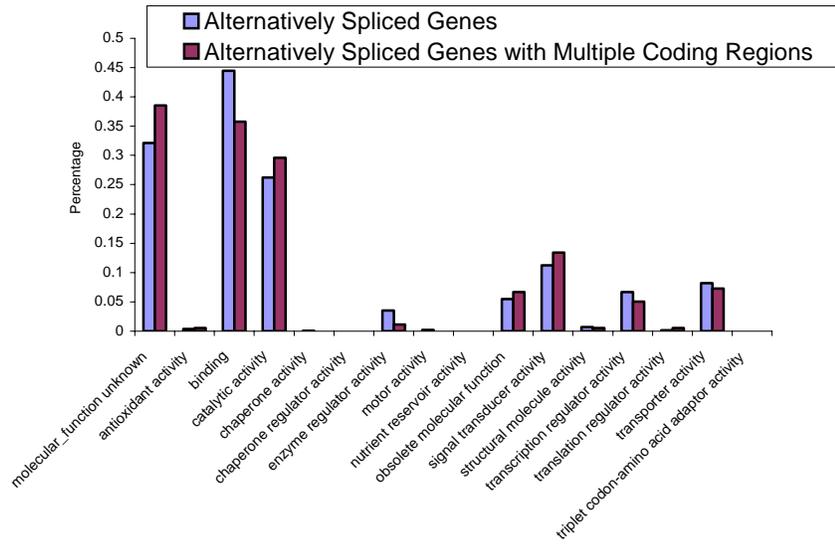
(a)



(b)



(c)



(d)

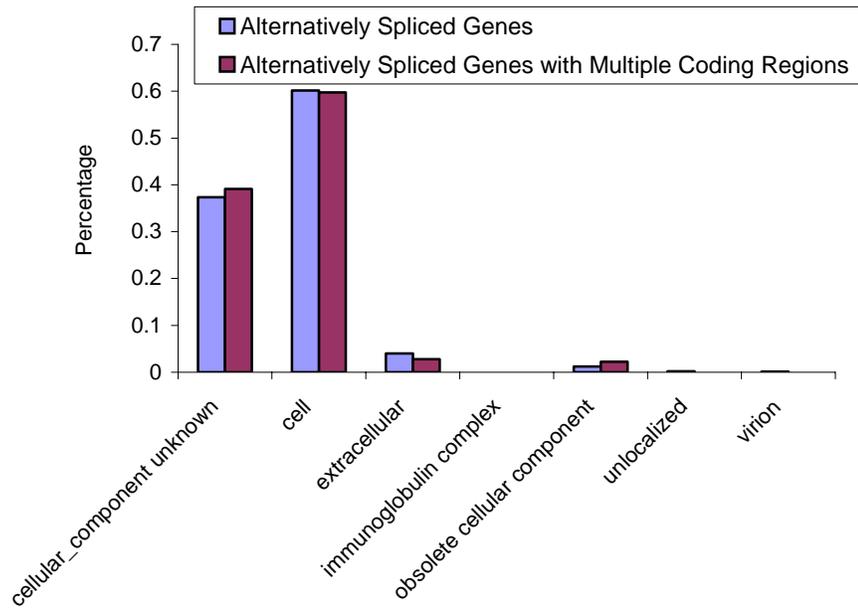


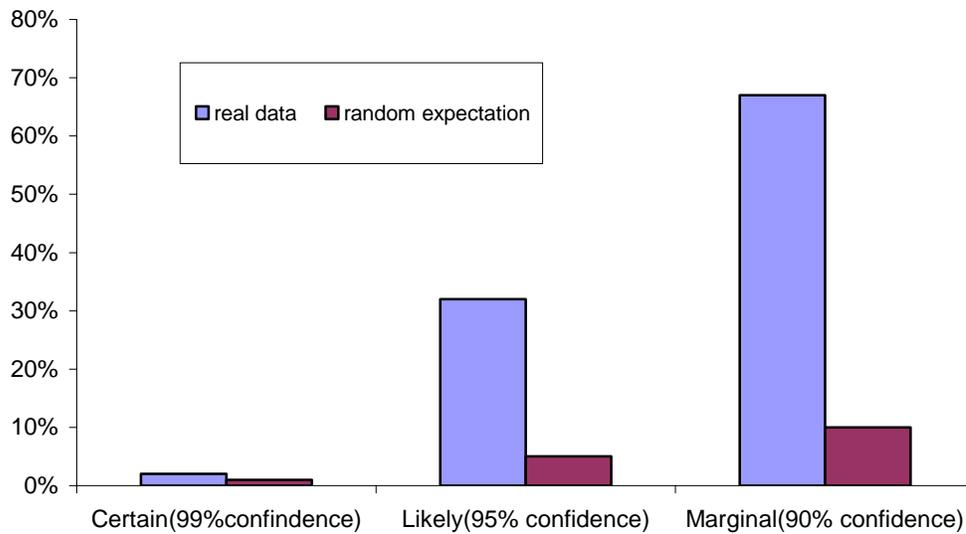
Fig. S2

(a) Statistical significance of sequence similarity search (by BLASTP) in the peptide sequences in DRFs in dual coding regions.

	Hits to PDB	No Hits
Peptide sequences in DRFs	12	90
Randomly shuffled control	4	98

$\chi^2 = 4.3$, P-value < 0.04

(b) Statistical significance of threading prediction (by FUGUE) in the peptide sequences in DRFs in dual coding regions.



$\chi^2 = 70$, P-value < 10^{-14}

Fig. S3

Schematic representation of a dual coding region in the human *UPK3B* gene.

Exons are represented by boxes, and introns are represented by connecting lines.

Numbers inside the boxes refer to base pairs. Roman numerals indicate intron phases.

The dual coding region is marked by a red arrow. NM_030570 and NM_182684 are RefSeq accession numbers.

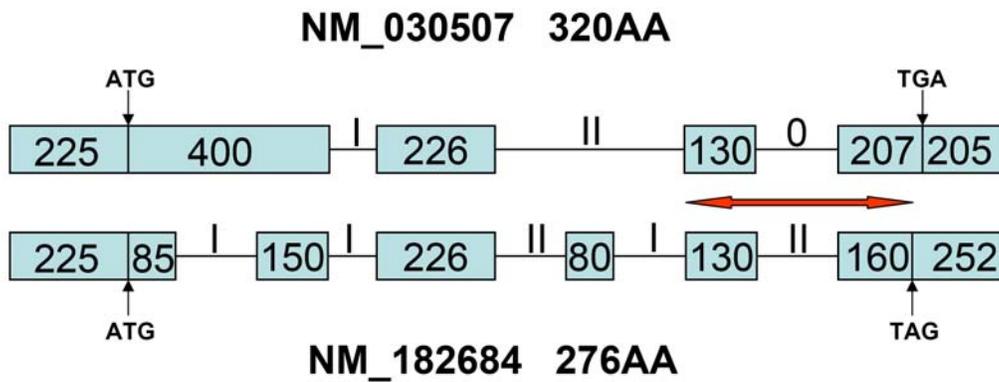


Fig. S4. The distribution of amino acids in Human DRFs corresponding to stop codons in dog DRFs.

