

The majority of human genes have regions repeated in other human genes

Roy J. Britten*

California Institute of Technology, 101 Dahlia Avenue, Corona del Mar, CA 92625

Contributed by Roy J. Britten, February 17, 2005

Amino acid sequence comparisons have been made between all of 25,193 human proteins with each of the others by using BLAST software (National Center for Biotechnology Information) and recording the results for regions that are significantly related in sequence, that is, have an expectation of $<1 \times 10^{-3}$. The results are presented for each amino acid as the number of identical or similar amino acids matched in these aligned regions. This approach avoids summing or dealing directly with the different regions of any one protein that are often related to different numbers and types of other proteins. The results are presented graphically for a sample of 140 proteins. Relationships are not observed for 26.5% of the 12,728,866 amino acids. The average number of related amino acids is 36.5 for the majority (73.5%) that show relationships. The median number of recognized relationships is ≈ 3 for all of the amino acids, and the maximum number is 718. The results demonstrate the overwhelming importance of gene regional duplication forming families of proteins with related domains and show the variety of the resulting patterns of relationship. The magnitude of the set of relationships leads to the conclusion that the principal process by which new gene functions arise has been by making use of preexisting genes.

domains | protein | relationships

It has long been clear that genes come in families, with selection for functionally useful domains. Genes can be thought to be the result of past events of domain shuffling, splicing, fusion, deletion, and duplication during the evolution of protein families. This article reports the results of an attempt to count the amino acid sequence relationships among all of the presently identified human genes. It amounts to an overview. For this purpose, the relationships of individual amino acids are counted when they are part of regions of proteins that have significant relationships to other proteins. We find that the vast majority of the proteins are involved in such relationships and on average are related to ≈ 26 other proteins.

Methods

The 25,193 human gene amino acid sequences were downloaded from the National Center for Biotechnology Information. Using BLAST, each sequence as probe was compared with all of the other sequences (called targets), rejecting any alignments with an expectation value $>1 \times 10^{-3}$. In comparing amino acid sequence pairs, BLAST often reports alignments of short regions of particular probes and targets that may multiply overlap. Each amino acid of the probe that was matched to either an identical or similar amino acid within one of these regions was given a score of 1, regardless of how many times it was reported in overlapping regions with the same target. Then, the scores of all of the matches to different targets were added up to determine for each amino acid in each gene the number of matches present within regions of significant relationships. A table for the 25,193 genes was constructed in which a name index and the length are listed for each, followed by a listing of the number of similar or identical amino acids recognized for each of the amino acids of the protein. This 30-megabyte table is available on CD by e-mail request. Also available on CD is FORTRAN software which allows

the graphical presentation for each gene of the patterns of frequency of matches for all amino acids. The identification number in Table 2 and Table 3, which is published as supporting information on the PNAS web site, is derived from the National Center for Biotechnology Information code NM_XXXXXX for the mRNA for that protein, which is converted to 1XXXXX. For the CD disk, the identification method is to search the PubMed nucleotides database for “*Homo sapiens*” and a protein name or identifier. Among the results will be the NM_XXXXXX, which is converted as above to find the protein on the disk, or XM_XXXXXX, which is converted to 2XXXXX.

Results

Numbers of Matches. The approach counts all of the individual amino acid matches (identical or similar) that are part of a region of a gene that has significant sequence relationships with other genes. In Table 1, the two columns show percent matches with an expectation value of $<1 \times 10^{-3}$ to the left and $<1 \times 10^{-6}$ to the right. There are no important differences between the two columns, and an expectation $<1 \times 10^{-3}$ was considered significant.

Genes often include significant matches to different sets of other genes in different regions of the gene, making it difficult to count the number of relationships of a gene as a whole. Therefore, a direct and simple method for counting the number of individual amino acid relationships is preferable, avoiding counting the sum or average of all of the relationships of different regions. This analysis was done with 25,193 genes that on average each contain 505 aa, adding to a total length of 12,728,886 aa. There were 339 million matches between amino acids in regions with an expectation $\leq 1 \times 10^{-3}$, giving an average of 36.5 matches per amino acid, counting just those amino acids with any matches. Table 1 summarizes the number of matches for the individual amino acids of each of the human genes in this fairly complete set with the other human genes in this set. The data are restricted to significantly matching regions with expectations of 1×10^{-3} .

Of interest is the number of genes that do or do not contain matching sequences. By count, 4,286 genes, or only 17.0% of the total of 25,193 genes, do not have any matching regions with expectation $\leq 1 \times 10^{-3}$. By whatever measure is chosen, the majority of genes or amino acids of gene sequences are part of sets of related sequences. Seventy-three percent of the amino acids match others in regions of significant relationship, with an average of 36.5 matches for amino acids in genes with regions that match significantly. The number of matches ranges from 0 to 718 for given amino acids. The highest-frequency amino acids are part of C2H2 zinc finger proteins.

Although the average number of matches for all of the amino acids of all genes on this list is ≈ 26 , the median is nearly 3. One view of this observation is that the majority of amino acids are in regions of proteins significantly related to few other proteins

Freely available online through the PNAS open access option.

*E-mail: r.britten@comcast.net or rbritten@caltech.edu.

© 2005 by The National Academy of Sciences of the USA

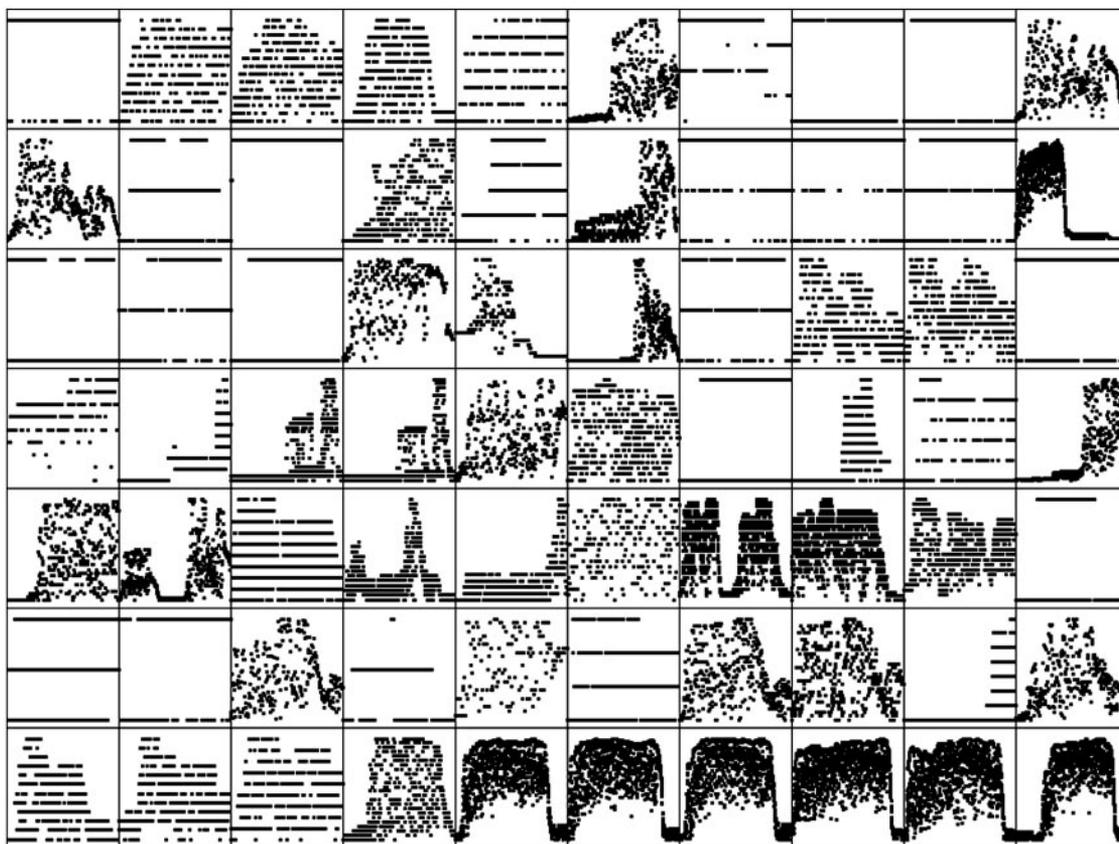


Fig. 2. Profiles for 70 human proteins. Shown are the first 70 cases in which there are matches. Table 2 identifies the individual proteins, and their locations are identified in columns 4 and 5 by numbers 1–7 (Y) for the rows (vertical) and 0–9 (X) for the columns (horizontal). All examples are normalized by frequency and length, as in Fig. 1, to fit their individual boxes on the figure. For example, case 1,0 (row, column) (*Upper Left*) reports a duplicated gene where one copy has diverged from the other at a modest number of amino acids scattered along its length, as shown by the dotted line at zero. Table 2 also lists the length of the protein and maximum frequency of any amino acid in the protein.

sequences. The maximum number of matching genes with expectation <10 was 252 for the original and 31 for the randomized set. Curiously, all of the randomized set had at least one match, whereas 190 of the original set found no matches at this open criterion.

To estimate how many of this set of genes can be considered to have significant matches with expectation $>1 \times 10^{-3}$, note that 85 of them have more matches than any of the randomized examples. This small number is hardly worth considering and suggests that the limit of 1×10^{-3} is reasonable. Because of evolutionary relationships, many of the so-called single-copy proteins may, in fact, have domains that are very divergently related to the well recognized functional domains. It will be difficult to determine the number in this set because of the well recognized problem of seeing distant relationships. Each member of the set of 26,193 was retested against all of the set by using a new version of BLAST, and the results were consistent with those shown in Figs. 2 and 4. Graphs show the same general shape in every case but appear to include more “hits,” indicating an increase in sensitivity to more divergent relationships, even though the limit was set to the same value of 1×10^{-3} . The percentage showing no relationships was 15.5%, compared with the previous 17%. Results with either BLAST version would be satisfactory, but the distributed CD contains the newer, more sensitive version results. This paper can be considered just the beginning of this sort of study of relationships with the tools at hand. Future studies, for example, by setting up hidden Markov model analysis, will likely show many additional relationships.

A Translated Fragment of an Alu Repeat in a Zinc Finger Gene. When REPEATMASKER (Institute for Systems Biology, Seattle) software examined the coding sequence of the zinc finger gene 91 described in Fig. 1, it reported that the last 104 nucleotides (or ≈ 34 aa) are part of an Alu Ye sequence matching at 94%. It turns out that the last exon of this gene is short and is made up only of the fragment of Alu sequence, similar to Alu Ye. Thus, this is probably an example of a gene variation caused by the presence of an exon derived from a mobile element (1–3). A search turned up five other cases of zinc finger genes that have similar short terminal Alu sequences, as well as a few more zinc finger genes with fragments of other mobile elements inserted into the coding sequences.

Regarding the interpretation of Fig. 1, the smoothly down-sloping region at the right-hand end is mostly the Alu sequence translation. This region indicates that there are hundreds of related amino acid sequence regions that include these amino acids, but it does not show that they are linked to zinc finger genes. When REPEATMASKER is used to search all of the 25,193 gene coding sequences for mobile elements, REPEATMASKER finds many hundreds but adding up to $<1\%$ of the length of these coding regions. Among them are ≈ 74 examples of fragments of Alu sequences in genes that produce known proteins. Presumably, some of these Alu sequences are part of the sequences that are matched by the amino acids of the terminal short region of Fig. 1, and they occur among many different gene coding sequences. However, the short terminal region that descends like the tail of an elephant at the right edge of Fig. 1 indicates that there are >400 similar amino acids in significantly related

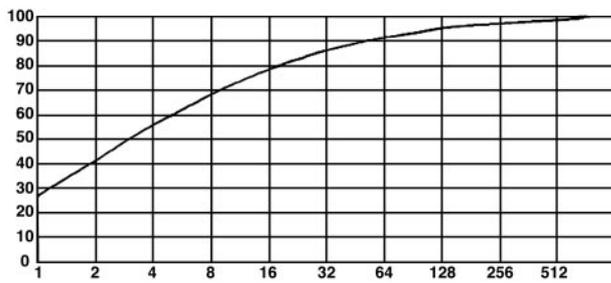


Fig. 3. Cumulative distribution of numbers of matches. The cumulated percent of amino acids plotted against the log of the number of identical or similar amino acids recognized, always restricted to significantly related aligned regions (expectation $< 1 \times 10^{-3}$). The curve starts with the 27% of amino acids that have no matches, passes through the median at ≈ 3 matches, and passes through 90% with fewer than ≈ 50 matches. Finally, there is a step representing the zinc finger proteins and a few other proteins with between 500 and 718 matches.

terminal region segments of various lengths (unpublished data). In each case, fewer amino acid similarities were found than in Fig. 1, and more were found the longer the segment that was used. The 35-aa translated Alu sequence alone found very few matches.

Although many years ago it was a surprise to find Alu sequences as part of the coding sequences of genes (4, 5), there are now well recognized examples (1–3, 6–9). However, in this example, Alu sequences appear to be present in a small number zinc finger genes.

Discussion

By whatever measure is chosen, the majority of genes or amino acids are part of sets of related sequences. Seventy-three percent

of the amino acids match others in regions of significant relationship, with an average of 36.5 matches for these amino acids. Eighty-three percent of the 25,193 proteins contain regions that significantly match other proteins on this list, which probably contains the great majority of human genes. Thus, the majority of genes are members of families of genes in the sense that the sequences of significant functional regions are related. The point must be made that a large fraction of the genes on this list are hypothetical or computer-derived. However, that may not much change the statistics, because the first 165 genes on the list are not hypothetical, and the fraction without significant shared domains is about the same as the total.

The cases shown in Figs. 2 and 4 show that the relationships often are restricted to local regions. For additional examples, FORTRAN is available by e-mail request. The data will be distributed on a CD for all of the 25,193 cases.

The basic issue is the mechanism that underlies these massive sets of relationships. These relationships are presumably the result of past events of regional shuffling, splicing, fusion, deletion, and duplication during evolution of protein families. During this set of processes, the individual proteins have been selected to carry out their function whether it be to produce a regulatory molecule, an enzyme, or a structural molecule. The copying, duplication, or multiplication processes have long been recognized as important (10). An overview of the resulting recognized domains is provided in ref. 11. These data show the true scale in the human genome for the 25,000 proteins recognized so far and leave no doubt that the principal process by which new gene functions arise is by making use of preexisting genes.

John Williams (California Institute of Technology) was responsible for much of the computer analysis and construction of PERL programs. Eric H. Davidson's laboratory at the California Institute of Technology supported this work.

1. Nekrutenko, A. & Li, W. H. (2001) *Trends Genet.* **17**, 619–621.
2. Lorenc, A. & Makalowski, W. (2003) *Genetics* **118**, 183–191.
3. Sorek, R. R., Ast, G. & Graur, D. (2002) *Genome Res.* **12**, 1060–1067.
4. Brownell, E., Mittereder, N. & Rice, N. R. (1989) *Oncogene* **4**, 935–942.
5. Caras, I. W., Davitz, M. A., Rhee, L., Weddell, G., Martin, D. W., Jr., & Nussenzweig, V. (1987) *Nature* **325**, 545–549.
6. Smit, A. F. (1999) *Curr. Opin. Genet. Dev.* **9**, 657–663.
7. Brosius, J. (1999) *Gene* **228**, 115–134.
8. Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C., Zody, M. C., Baldwin, J., Devon, K., Dewar, K., Doyle, M., FitzHugh, W., et al. (2001) *Nature* **409**, 860–921.
9. Dagan, T., Sorek, R., Sharon, E., Ast, G. & Graur, D. (2004) *Nucleic Acids Res.* **32**, D489–D492.
10. Ohno, S. (1970) *Evolution by Gene Duplication* (Springer, Berlin).
11. Apweiler, R., Bairoch, A., Wu, C. H., Barker, W. C., Boeckmann, B., Ferro, S., Gasteiger, E., Huang, H., Lopez, R., Magrane, M., et al. (2004) *Nucleic Acids Res.* **32**, D115–D119.