

Review

Studying Genomes Through the Aeons: Protein Families, Pseudogenes and Proteome Evolution

Paul M. Harrison and Mark Gerstein*

Department of Molecular
Biophysics and Biochemistry
Yale University, 266 Whitney
Avenue, P.O. Box 208114, New
Haven, CT 06520-8114, USA

Protein families can be used to understand many aspects of genomes, both their “live” and their “dead” parts (i.e. genes and pseudogenes). Surveys of genomes have revealed that, in every organism, there are always a few large families and many small ones, with the overall distribution following a power-law. This commonality is equally true for both genes and pseudogenes, and exists despite the fact that the specific families that are enlarged differ greatly between organisms. Furthermore, because of family structure there is great redundancy in proteomes, a fact linked to the small size of the minimal, indispensable sub-proteome for each organism and the large number of dispensable genes. Pseudogenes in prokaryotes represent families that are in the process of being dispensed with. In particular, the genome sequences of certain pathogenic bacteria (*Mycobacterium leprae*, *Yersinia pestis* and *Rickettsia prowazekii*) show how an organism can undergo reductive evolution on a large-scale (i.e. the dying out of families) as a result of niche change. There appears to be less pressure to delete pseudogenes in eukaryotes. These can be divided into two varieties, duplicated and processed, where the latter involves reverse transcription from an mRNA intermediate. We discuss these collectively in yeast, worm, fly, and human. The fly has few pseudogenes because of its high rate of genomic DNA deletion. In the other three organisms, the distribution of pseudogenes on the chromosome and amongst different families is highly non-uniform. Pseudogenes tend not to occur in the middle of chromosome arms and to be associated with lineage-specific (as opposed to highly conserved) families that have environmental-response functions. This may be because, rather than being dead, they may often form a reservoir of diverse “extra parts” that can be resurrected to help an organism adapt to its surroundings. In yeast, there may be a novel mechanism involving the [PSI⁺] prion that potentially enables this resurrection. In worm, the pseudogenes tend to arise out of families (particularly chemoreceptors) that are greatly expanded in it compared to the fly. The human genome stands out in having many processed pseudogenes. These have a character very different from those of the duplicated variety, essentially just representing random insertions. Thus, their occurrence tends to be roughly in proportion to the amount of mRNA for a particular protein and to reflect the extent of the intergenic sequences. Further information about pseudogenes is available at <http://genecensus.org/pseudogene>

© 2002 Elsevier Science Ltd

*Corresponding author

Keywords: aeons; pseudogenes; proteome evolution

Abbreviations used: LINE, long interspersed nuclear element.

E-mail address of the corresponding author: mark.gerstein@yale.edu

The complete or near-complete sequencing of the genomes of six eukaryotes (at the time of writing) and dozens of prokaryotes is enabling us to examine molecular evolution and diversity of proteins from a “whole-proteome” perspective. In the present review, we discuss various themes and

issues in relation to proteome evolution, examining both the “live” and “dead” proteomes of specific genomes (all the proteins encoded by an organism and all the pseudogenes). We set the stage for discussion of pseudogene populations by surveying different issues relating to protein family redundancy in the live proteome, and how it evolves. In particular, we examine how such redundancy can be viewed in terms of partition into essential and dispensable sub-proteomes. Chiefly, then, we discuss the distribution of proteins and protein families in pseudogene populations for prokaryotes, and specifically for the eukaryotes yeast, worm, fly and human, and the implication of these dead or “dispensed-with” sequences for proteome evolution.

What is a protein family?

A protein family is usually defined as a group of sequences with an obvious evolutionary relationship, judged chiefly by protein sequence comparison, i.e. whose evolution can be studied readily at the sequence level. The definition of the threshold of similarity is arbitrary in practice and different degrees of protein sequence similarity are used depending on the context.^{1–3} Membership of the same protein family is now commonly determined by the occurrence of a sequence motif indicative of sequence, structural and functional similarity, with integrated databases of such motifs used routinely in genome annotation.^{4,5} There are now many databases that cluster protein sequences manually or automatically to varying degrees, at various levels of sequence and structural similarity (e.g. ProtoMap,⁶ SYSTERS,⁷ SCOP⁸ and CATH⁹). As a higher level, a superfamily can then be described in terms of groups of families that have more distant similarity; they may have common evolutionary origin as judged by functional and structural similarities. (This is the definition used in the SCOP database.⁸) Different superfamilies can be grouped together if they have the same protein fold. Sometimes it is more appropriate to group families together into similar functional classes, e.g. the Gene Ontology database,¹⁰ MIPS functional classification¹¹ or GenProtEc for *Escherichia coli*.¹² Although, usually, as for most of the work discussed below, robustness of results is reported for a range of sequence similarity cut-offs, there are a number of caveats in considering assignment of protein families and superfamilies to genomic data.^{13,14} Firstly, such assignment procedures are biased towards larger families and superfamilies, in that sequence-searching procedures, such as the commonly used iterative program PSI-BLAST,¹⁵ operate better for larger known families and are calibrated to search for larger families; secondly, for obvious reasons, gene prediction is more successful for them too.

Surveys of the “live” proteome

There has been extensive recent work on the counting of different levels of proteome parts: protein families, superfamilies and folds.^{5,13,16–24} Initially, this work focused on prokaryotes, but is now shifting emphasis to the recently genomically sequenced eukaryotes. Surveys of protein fold and superfamily occurrence in microbial proteomes shows that a few folds predominate, whereas many folds occur only once; protein fold occurrences tend to rely on the prevalence of a single superfamily, although the rankings for these corresponding folds and superfamilies vary widely.¹⁹ There are similar findings for the eukaryotes (Table 1).

Power-law distribution of protein family size in proteomes

Despite expansion and contraction in the size of individual protein families in proteomes, the redundancy in protein families appears to have a characteristic distribution common to viral, bacterial, archaeal and eukaryotic genomes.^{3,16} An initial analysis of the distribution of the number of sequences in protein families *versus* their occurrence showed that the distribution for protein families in proteomes follows power-law behaviour (i.e. a linear relationship on a log–log plot), with a shallower slope for the relationship in the larger genomes.²⁵ Huynen & Nimwegen³ did a similar analysis for a larger number of microbial genomes and found that the power-law behaviour was maintained over a large range of sequence similarity thresholds used for clustering into families. They argued, using a simple probabilistic formalism, that the power-law distribution implies that gene duplications and deletions within gene families are largely dependent on one another. Other studies have shown that the distribution of the number of protein families and of protein folds in a proteome can be explained by simple evolutionary models that involve only duplication or the creation of new families or folds.^{22,26} An example of this power-law behaviour is illustrated in Figure 1 for families in the yeast proteome, and for protein folds and superfamilies.

Protein family redundancy in proteomes and its evolution in eukaryotes

The total number of protein domain sequence families, or functional diversity, appears to vary much less between organisms than overall proteome size. This is most striking in the eukaryotes.^{2,27,28} For example, despite the wide variation in the number of annotated genes, the yeast, worm, fly and human proteomes seem to contain similarly sized subsets of the InterPro sequence domain database (851 for yeast; 1014 for worm; 1035 for fly; 1262 for human, at the time of writing).^{5,27} The eukaryotic proteomes comprise

Table 1. Top-ranking protein superfamilies and folds in five eukaryote proteomes

Top-ranking superfamilies					Top-ranking folds				
Yeast	Worm	Fly	Mustard weed	Human	Yeast	Worm	Fly	Mustard weed	Human
P-loop NTP hydrolase (438)	P-loop NTP hydrolase (651)	C2H2 Zn finger (823)	P-loop NTP hydrolase (1282)	C2H2 Zn finger, 7.37.1 (3424)	P-loop NTP hydrolase, 3.32 (438)	Ig-like, 2.1 (1044)	<i>Ig-like, 2.1 (999)</i>	<i>α/α Superhelix, 1.111 (1475)</i>	C2H2 Zn finger, 7.37 (3424)
Protein kinase (133)	Ig (571)	P-loop NTP hydrolase (661)	Protein kinase (1070)	Ig, 2.1.1 (1453)	α/α Superhelix, 1.111 (195)	P-loop NTP hydrolase, 3.32 (651)	C2H2 Zn finger, 7.37 (823)	P-loop NTP hydrolase, 3.32 (1282)	Ig-like, 2.1 (3034)
WD-repeat (107)	Protein kinase (500)	<i>Ig, 2.1.1 (548)</i>	<i>Tetratricopeptide repeat, 1.111.8 (787)</i>	P-loop NTP hydrolase, 3.32.1 (1229)	Ferredoxin-like, 4.51 (154)	Protein kinase, 4.130 (500)	P-loop NTP hydrolase, 3.32 (661)	Protein kinase, 4.130 (1070)	P-loop NTP hydrolase, 3.32 (1229)
RNA-binding domain (104)	EGF/laminin (400)	EGF/laminin (330)	RNI-like (709)	EGF/laminin, 7.3.9 (1083)	Protein kinase, 4.130 (133)	Knottin, 7.3 (429)	α/α Superhelix, 1.111 (438)	Leucine-rich repeat, 3.9 (812)	Knottin, 7.3 (1114)
NADP-binding Rossmann fold (99)	C-type lectin (369)	Protein kinase (288)	RING finger (468)	Fibronectin type-III (817)	Seven-bladed β propeller, 2.64 (118)	α/α Superhelix, 1.111 (405)	Ferredoxin-like, 4.51 (357)	Ferredoxin-like, 6.51 (451)	α/α Superhelix, 1.111 (898)
ARM repeat (84)	Glucocorticoid receptor-like (349)	Spectrin repeat (268)	Homeodomain (461)	Protein kinase, 4.130.1 (710)	TIM barrel, 3.1 (114)	C-type lectin, 4.154 (369)	Knottin, 7.3 (345)	DNA/RNA-binding 3-Helical bundle, 1.4 (539)	Protein kinase, 4.130 (710)
DNA/RNA polymerises (59)	Nuclear receptor ligand-binding domain (284)	RNA-binding domain (257)	RNA-binding domain (426)	Cadherin (676)	RNase H, 3.50 (110)	Glucocorticoid receptor-like, 7.39 (349)	Protein kinase, 4.130 (288)	RING finger, 7.44 (468)	Ferredoxin, 4.51 (655)
Actin-like ATPase (56)	Homeodomain (263)	Trypsin-like serine protease (240)	NADP-binding Rossmann-fold domain (366)	RNA-binding domain, 4.51.7 (517)	NADP-binding Rossmann fold, 3.2 (99)	DNA/RNA-binding 3-helical bundle, 1.4 (329)	Spectrin repeat, 1.7 (272)	Seven-bladed β propeller, 2.64 (451)	DNA/RNA-binding 3-helical bundle, 1.4 (510)
Membrane all-α (54)	C2H2 Zn finger (255)	<i>Fibronectin type III, 2.1.2 (219)</i>	α/β Hydrolase (341)	PH domain, 2.52.1 (415)	DNA/RNA-binding 3-helical bundle, 1.4 (59)	Ferredoxin-like, 4.51 (301)	Trypsin-like serine protease, 2.44 (240)	TIM barrel, 3.1 (383)	PH domain, 2.52 (415)
Zn2/Cys6 DNA-binding domain (53)	α/β-Hydrolase (219)	<i>Cadherin, 2.1.6 (213)</i>	<i>ARM repeat, 1.111.1 (284)</i>	Homeodomain, 1.4.1 (339)	DNA/RNA polymerases, 5.8 (59)	Nuclear receptor ligand-binding domain, 1.116 (284)	Seven-bladed β propeller, 2.64 (118)	NADP-binding Rossmann-fold domain, 3.2 (366)	Seven-bladed β propeller, 2.64 (394)

The Table shows the top-ranking folds in eukaryotes from SCOP. There is a pattern similar to that observed in prokaryotes.¹⁹ In particular, for human, the prevalence of a fold tends to be due to a particular superfamily prevalence (superfamilies and folds in bold in the Table. Examples of folds that have multiple prevalent superfamilies are observed; examples for fly and mustard weed (*A. thaliana*) are in italics. Ig, immunoglobulin.

253	316
254	317
255	318
256	319
257	320
258	321
259	322
260	323
261	324
262	325
263	326
264	327
265	328
266	329
267	330
268	331
269	332
270	333
271	334
272	335
273	336
274	337
275	338
276	339
277	340
278	341
279	342
280	343
281	344
282	345
283	346
284	347
285	348
286	349
287	350
288	351
289	352
290	353
291	354
292	355
293	356
294	357
295	358
296	359
297	360
298	361
299	362
300	363
301	364
302	365
303	366
304	367
305	368
306	369
307	370
308	371
309	372
310	373
311	374
312	375
313	376
314	377
315	378

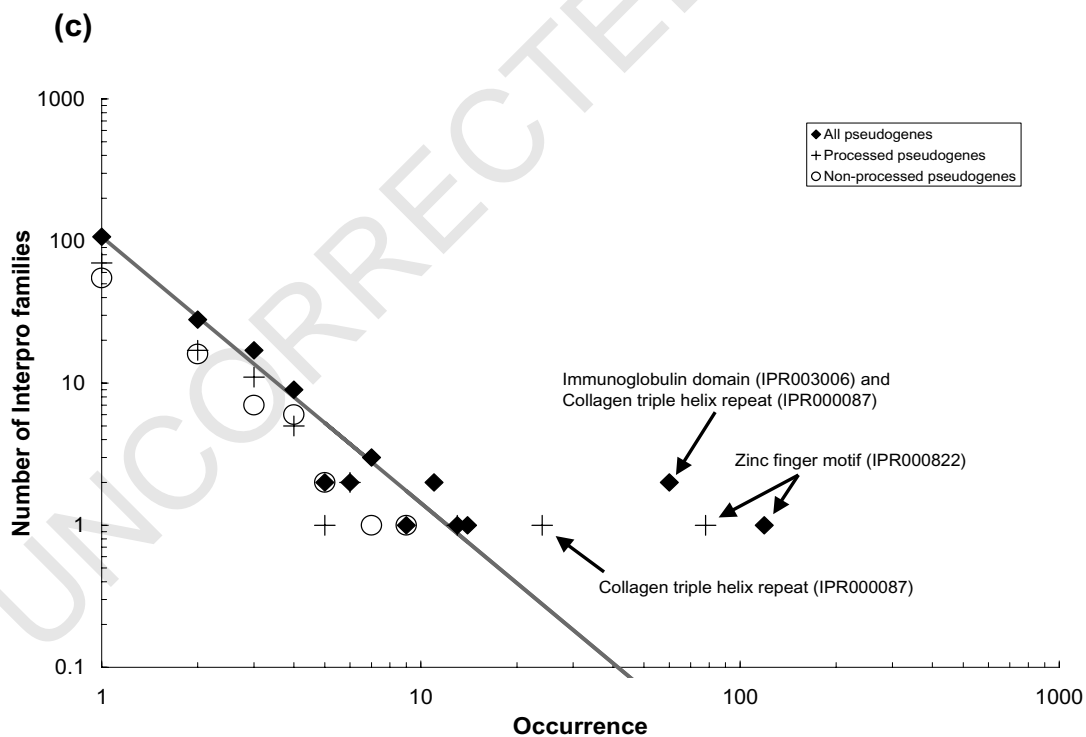
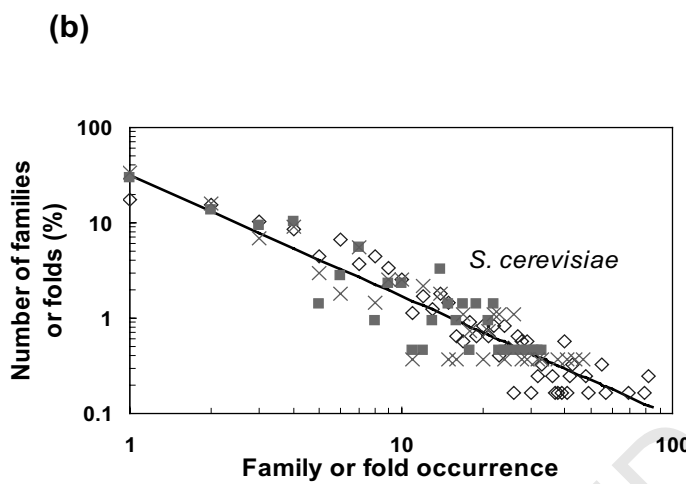
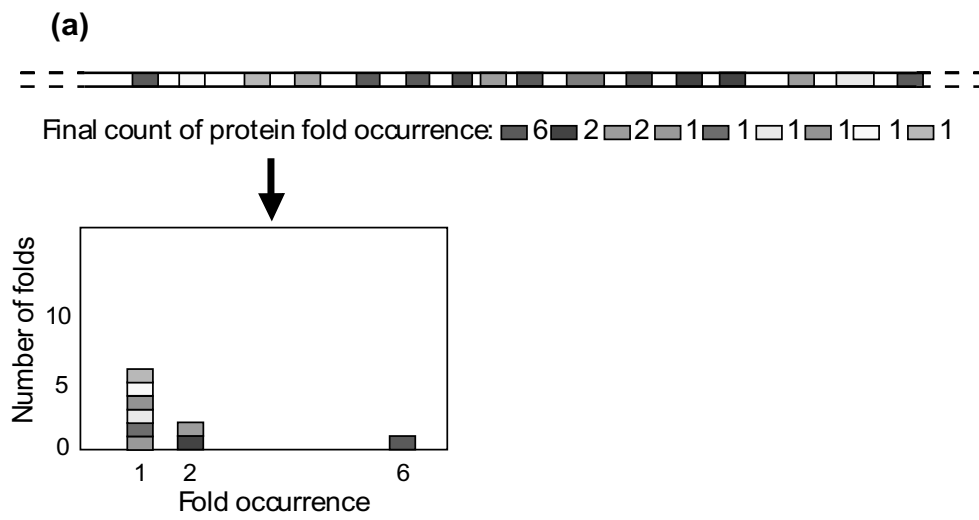
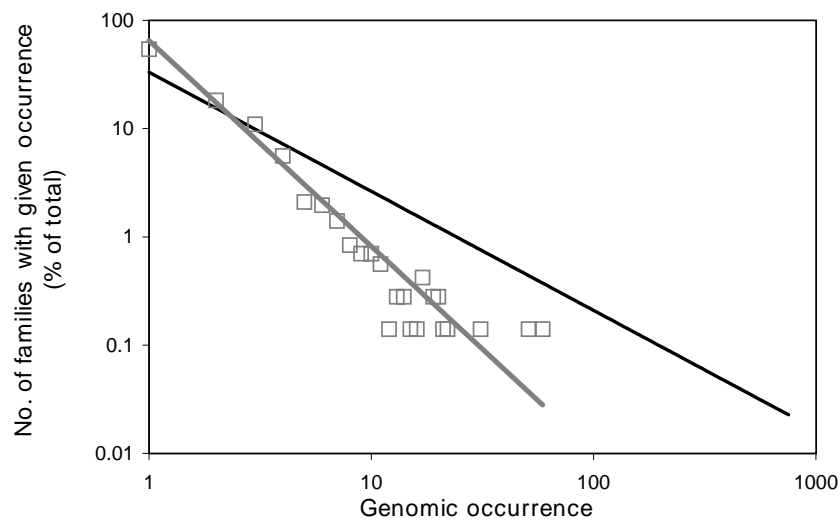


Figure 1 (legend opposite)

379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441

442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504

(d)



law fit to the distribution for pseudogene families (open boxes); the black line is the same fit for the distribution for gene families, clustered as described.⁶⁹ The axes are as for (b).

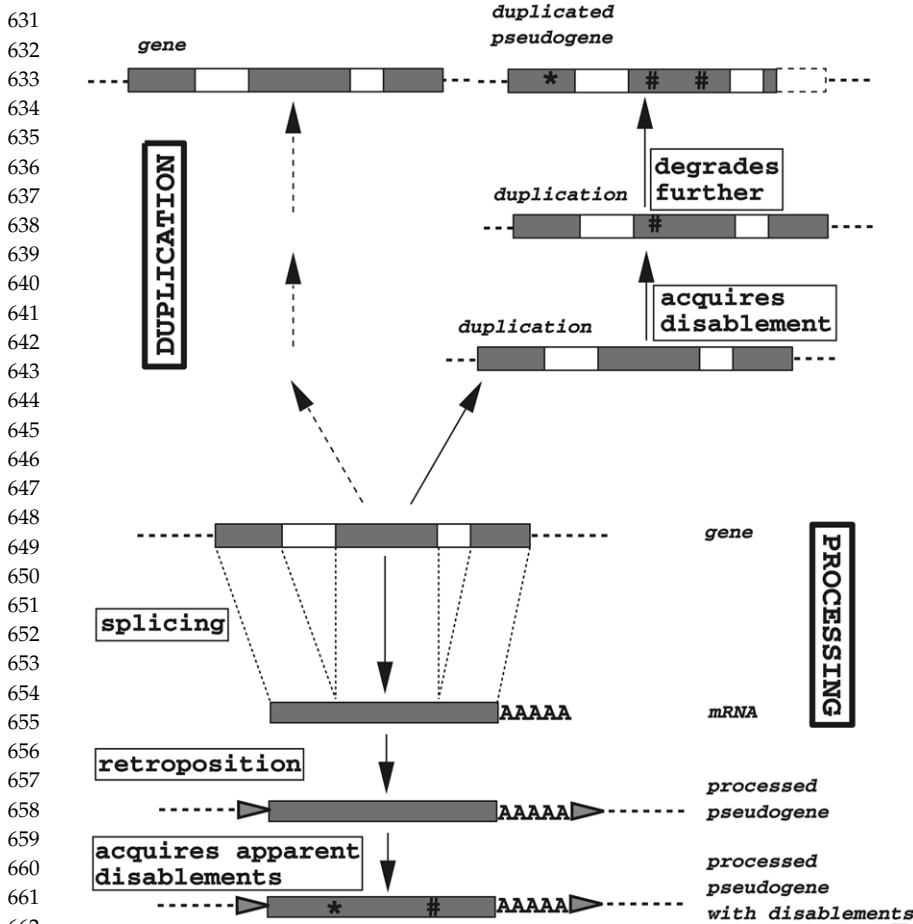
comparable coverage of the SCOP domain database⁸ in terms of superfamilies (between 460 (yeast) and 594 (human)²⁴). Extensive sequence family redundancy is observed at the individual gene level in the eukaryotes, most notably in *Arabidopsis thaliana*, where only 35% of proteins are singletons (i.e. have no paralogs).² (For comparison, the degree of family redundancy is less extensive in the *Saccharomyces cerevisiae* genome, which by the same strict criteria, contains 29% of proteins in families.) In *Arabidopsis*, the extensive redundancy is linked to a large number of segmental chromosomal duplications arising from four distinct large-scale duplication events 100 to 200 million years ago.²⁹ Regardless of the mechanism of formation (whether segmental or local duplication), from an individual gene perspective, new gene duplicates in eukaryotes arise at the rate of about 0.01 per gene per million years, with rates for individual genomes ranging from 0.02 for *Caenorhabditis elegans* to 0.002 for *Drosophila melanogaster*; this is of the same order as the rate of mutation per nucleotide site.³⁰

By what mechanism does the gene family redundancy chiefly arise? For example, Wolfe and colleagues identified homologous arrays of genes on different yeast chromosomes, which they hypothesized had arisen from a single, whole-genome duplication event about 100 million years ago, after separation from the *Saccharomyces kluyveri* yeast branch^{31–33} However, ~90% of the resulting individual duplicated genes arising from this event appear to have been lost. Furthermore, there is no evidence that these duplications occurred at the same time; indeed, many segmental chromoso-

mal duplications may have occurred in yeast at various times over the past 200–300 million years.¹ On the basis of the partial genome sequencing of 13 ascomycete relatives of *S. cerevisiae*, the conservation in yeast of singletons and gene family redundancy was found to arise mostly from local duplication events and did not support the whole-genome duplication hypothesis in yeast evolution.³⁴ Finally, in the human genome, notably, there is much less occurrence of pairs of chromosomal segments where the density of duplicated genes approaches that of *A. thaliana* or *S. cerevisiae*, indicating far less segmental chromosomal duplication.²⁷ Inclusion of detailed pseudogene annotations for the analysis described above would help to pin-point the mechanism of evolution of gene redundancy (see below for a discussion of pseudogene populations).

Indispensable and dispensable sub-proteomes

What is the minimal “indispensable” sub-proteome for the eukaryotic cell? Regardless of how the protein family redundancy in the yeast proteome has arisen, it seems clear from gene disruption experiments that the sub-proteome essential for yeast cell viability contains only ~1000 proteins.^{35,36} This is about three times the number of proteins adjudged essential for a minimal prokaryotic cell.³⁷ Wagner noted, from analysis of gene disruption data for yeast, that there is no strong correlation between gene family redundancy and robustness against gene disruption. This indicates that there is a contribution to the robustness to mutation of a given gene that arises



processed pseudogenes include small direct repeats (grey triangles) at either end of the pseudogene and a polyadenine tail (indicated here by AAAAA). The apparent coding frame of the pseudogene would then acquire obvious disablements, such as premature stops and frameshifts over evolutionary time.

from other genes with no detectable ancestral relationship, which, for instance, could provide alternative routes through pathways.³⁸

From studies in yeast, it seems clear that many proteins have marginal effects on species fitness.³⁹ In a study of 34 *S. cerevisiae* genes that were judged non-essential by gene disruption,³⁵ 70% of them were found to have marginal but significant effects on the fitness of a strain.⁴⁰ This implies that the effective size of the indispensable sub-proteome for yeast can be determined only from study of its behaviour from generation to generation for the reproducing organism. This generation-weighted proteome could perhaps be dubbed the selectome, in analogy to the transcriptome (where the occurrence of different proteins is weighted by their transcription levels at different time-points and under various conditions^{41,42}). The marginality of contribution to fitness in yeast, or protein dispensability, has been shown to be correlated with the molecular rate of evolution, i.e. more dispensable proteins evolve more rapidly.⁴³ It is conceivable that protein families with a higher molecular rate of evolution are more likely to

Figure 2. Two types of pseudo-gene. Pseudogenes are produced chiefly either by duplication or by processing. An example of a gene with three exons (shaded areas) is shown (boxed at the center of the Figure), with no non-coding segment in the exons for simplicity. ATG labels the start of the coding sequence, an asterisk (*) labels a stop codon and hash (#) stands for a frameshift. A non-processed or duplicated pseudogene simply arises when a gene duplication acquires a disablement that leads to: (i) lack of transcription; (ii) degradation *via* nonsense-mediated decay; or (iii) for an unknown subset of pseudogenes that produce messenger RNA transcripts escaping nonsense-mediated decay,¹⁰⁵ to degradation at some later unknown stage, so that a functioning protein chain is not formed. After such an initial disablement, the recently defunct pseudogene will acquire further obvious disablements of its reading frame (such as premature stops arising from point mutation, or truncations and frameshifts arising from deletion or insertion). A processed pseudogene arises when a messenger RNA transcript is reverse transcribed and re-integrated into the genomic DNA. Characteristic signals for these pro-

cessed pseudogenes in the genome. Proteins that have recently been dispensed with from the proteome may remain in the genome as pseudogenes (depending on genome-specific rates of genomic DNA loss and mutation), and this aspect of proteome evolution is discussed below.

The “dead” proteome: pseudogenes and proteome evolution

In the previous sections, we have discussed how the live part of the proteome of an organism is distributed into protein families, and some implications of this sequence redundancy. We now focus on the corresponding dead population of sequences, pseudogenes.

Pseudogenes are disabled copies of genes (or decayed remnants of genes) that do not produce a full-length protein chain. They can generally be divided into two types (Figure 2). Firstly, “processed” pseudogenes arise from reverse transcription from messenger RNA (mRNA) and re-integration into the genomic DNA.⁴⁴ These have

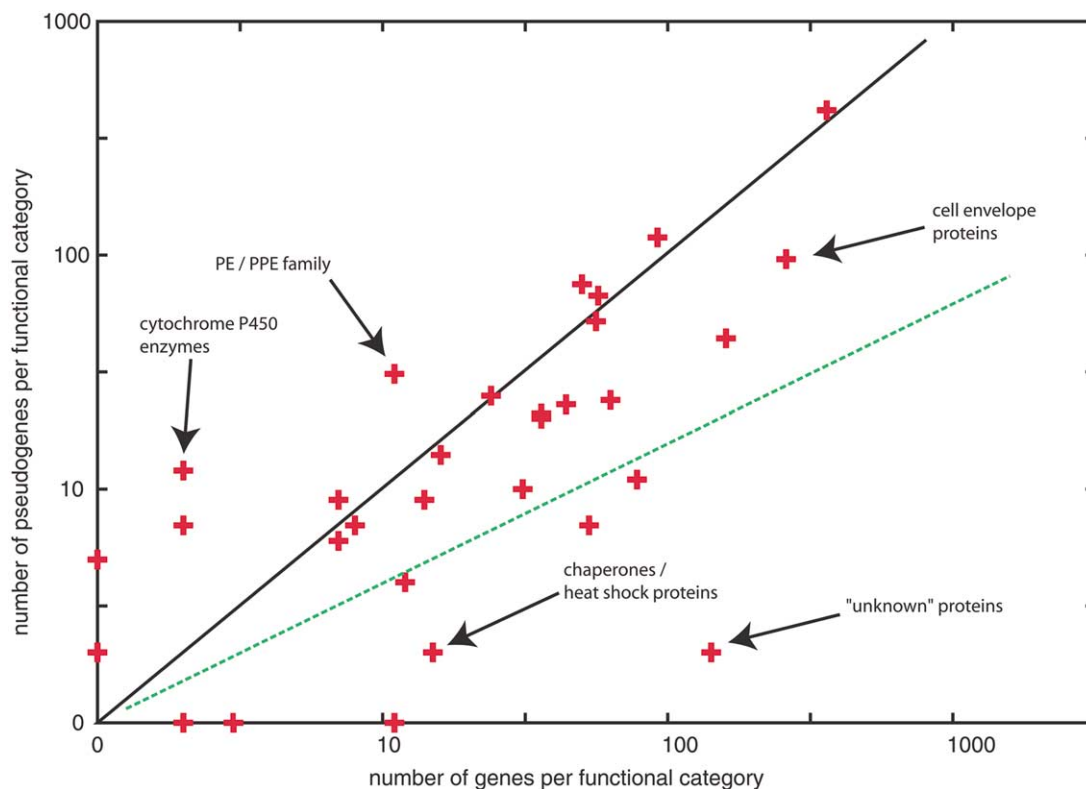


Figure 3. The relationship between the number of pseudogenes and genes for different functional categories in *M. leprae*. Each of the 31 functional categories listed by Cole *et al.*,⁵² (Figure 2 of that paper) is plotted. The continuous line represents the number of pseudogenes being equal to the number of genes. Eleven of the categories are above this line, i.e. are more “dead” than “live”. The dotted line represents the overall ratio of pseudogenes to genes in the proteome. Eight of the categories are below this line, i.e. more live than the overall ratio for live-to-dead.

been observed only in the metazoan animals and flowering plants, and presumably arise from mRNA transcripts in the germ-line cell lineage. In humans, they are probably made as a by-product of long interspersed nuclear element (LINE) retrotransposition.⁴⁵ After integration into the genome, they gradually accumulate disablements (stop codons, frameshifts, inserted repeat elements) of their reading frame. Secondly, “non-processed” or “duplicated” pseudogenes arise from duplication in the genomic DNA and subsequent disablement, most commonly through disruptive frameshift mutation or premature stop codon formation.⁴⁶ Formation of a pseudogene from gene duplication may have effects on the fitness of an organism; for example, if the duplicated gene has diverged very little since the duplication event that formed it (perhaps acquiring a slightly different activity or specificity in its function), the decrease in copy number for the gene family may be mildly deleterious. Conversely, copies of genes may be lost because that particular family is no longer as beneficial for fitness and has become more dispensable.

Pseudogenes, as “molecular fossils”, are important sequences for the study of molecular evolution. Here, we discuss the occurrence of pseudogenes from a whole-proteome perspective, making use,

where appropriate, of comparison of the prevalent families in proteomes and pseudogene populations. Such a perspective, of course, has been possible only recently with the advent of complete genome sequencing. We examine, in turn, the implications for proteome evolution in prokaryotes, and in the eukaryotes yeast, worm, fly and human. In prokaryotes, we see evidence for large-scale reductive evolution that mirrors the expansive evolution arising from horizontal transfer. In eukaryotes, we see that duplicated pseudogenes tend to be associated with environmental and response families. In the yeast, there appears to be a mechanism for conditionally “resurrecting” disabled genes as an evolutionary buffer to environmental fluctuation, perhaps in a concerted fashion. In the worm, the families of sequences that are prevalent in its pseudogene population have corresponding expanded or organism-specific populations in its genome, indicative of recent organism-specific expansions. For the fly, we argue that its apparently very small pseudogene population is linked to the size of its proteome through a very high rate of genomic DNA loss. Finally, for the human, we discuss the substantial number of processed pseudogenes relative to the putative total gene complement.

Prokaryotes: expansive and reductive proteome evolution

Prokaryotes can expand their proteomes by undergoing substantial horizontal transfer of genes from other strains and species.⁴⁷ Comparison of the two complete genomes sequences of *E. coli* strains O157:H7 EDL933 and K-12 MG1655,^{48,49} shows how dramatically dynamic this horizontal transfer can be. Over a quarter (26%, 1387/5416) of the O157:H7 EDL933 genes are specific to that strain compared to K-12 MG1655. Conversely, in the same manner, 528/4405 (12%) of K-12 MG1655 genes are strain-specific. Strain-specific variation such as this has led some to argue that it is perhaps best to compare organisms in terms of a "species genome", with a core sub-proteome, and a variable sub-proteome that comprises the proteins and protein families that vary from strain to strain.^{50,51} It will be interesting to see how closely correspondent such a core sub-proteome is to the indispensable subproteome, as discussed above for yeast.

Conversely, reductive evolution in bacteria may be equally dynamic. The recent sequencing of the genome of the bacterium *Mycobacterium leprae*, the leprosy pathogen, shows that it has undergone massive recent proteome decay.⁵² The *M. leprae* genome contains about ~1100 apparent pseudogenes, and ~1600 genes.⁵² This is a considerable reduction when compared to the ~4000 proteins encoded in the genome of the related bacterium *Mycobacterium tuberculosis* and involves decrease in the redundancy of almost all protein families, with loss of substantial parts of pathways, such as the anaerobic respiratory chain. For example, the repetitive, glycine-rich PE and PPE families comprise 167 genes in the *M. tuberculosis* genome; however, in *M. leprae* it is more dead than live, there are only nine such genes in *M. leprae*, and 30 related pseudogenes. This family is shown on a plot for all of the functional classes reported here with pseudogene number plotted *versus* gene number (Figure 3). On the other hand, the functional class for chaperones and heat-shock proteins has a much smaller dead-to-live ratio than the overall ratio of dead to live proteins. By analogy with the two *E. coli* strains, it would be interesting to see to what extent the observed huge proteome decay is specific for the *M. leprae* strain sequenced, and how this affects the definition of its core sub-proteome.^{50,51}

Proteome decay has been observed for two other pathogenic bacteria. The typhus pathogen *Rickettsia prowazekii* seems to have undergone such reductive evolution recently.^{53,54} Initially, it was thought to harbour only 12 pseudogenes,⁵³ but subsequently this estimate was enlarged. Prokaryote genomes are generally very compact, harbouring little non-coding genomic DNA (generally <10%; *E. coli* K-12 has ~11%⁴⁸), implying that there is rapid deletion of any recently formed pseudogenes. However, the non-

coding DNA in the *R. prowazekii* genome is >24% of the genomic DNA, suggesting that it comprises undetected decayed remnants of genes. Comparison of the *R. prowazekii* genomic sequence to those of other Rickettsias,^{54,55} led to the detection of sequence similarity between (pseudo)genes in one species and the equivalent non-coding DNA in other species. These more fragmentary and disabled pseudogenic sequence homologies were dubbed fossil ORFs⁵⁵ or decayed orthologs.⁵⁴ Inclusion of these more decomposed remnants in *R. prowazekii* raises its total pseudogene population to 241 (compared to 834 live genes). The plague bacterium *Yersinia pestis* has a smaller relative proportion of pseudogenes (160, compared to ~4000 live genes) that appear linked to the loss of an enteropathogenic lifestyle.⁵⁶

Yeast: resurrectable variation between strains

There are very few annotated pseudogenes in the sequenced laboratory strain of *S. cerevisiae*, S288C;⁵⁷ we could find at most 30 such annotations in the SGD and MIPS databases.^{11,58} From the analysis of disabled protein homology matches in the yeast genome, we believe that there may be up to a further 221 un-annotated pseudogenes in the *S. cerevisiae* S288C strain. This number rises further to 241 if we include pairs of existing ORF annotations, termed mORFs, that can be merged into a pseudogene and that could be complete ORFs in a different yeast strain⁵⁹ (Table 2). One of the most important previously documented pseudogenes in the yeast strain S288C is the FLO8 mutation.⁶⁰ This flocculin gene has an intact ORF in other strains but is disrupted by a single stop codon in S288C. This mutation has been shown to be the cause of the lack of diploid pseudohyphal filamentous growth in S288C, and has thus probably been selected in the laboratory so that yeast colonies are round and smooth. Strains that have an active FLO8 gene appear flocculent, having a fluffy colony appearance. The largest sequence families that are relatively prevalent in the S288C strain pseudogene population comprise flocculins like FLO8, the DUP family of double-transmembrane-helix proteins, growth inhibitors, helicases and stress-response proteins, whereas the most populated live families are forms of protein kinase, helicases, a transcriptional regulatory protein domain and the AAA ATPase domain (Table 3). Note how the pseudogenes appear to disproportionately have environmental and stress response functions. They have been found to occur near the ends of the chromosomes, mostly within 20 kb of the telomeres.⁵⁹

Sup35p is part of the surveillance complex in yeast that controls translation termination and nonsense-codon read-through.^{61,62} The [PSI⁺] prion in yeast arises from the propagation of an alternatively folded amyloid-like form of Sup35p.^{61,63} Thus, formation of the alternative

946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008

Table 2. Gene and pseudogene numbers

Organism	No. genes	No. pseudogenes	No. processed pseudogenes	No. duplicated pseudogenes	References
<i>R. prowazekii</i> (B)	834	241	0	241	53,54
<i>M. leprae</i> (B)	1604	1116	0	1116	52
<i>Y. pestis</i> (B)	4061	160	0	160	56
<i>S. cerevisiae</i> strain S288C (E)	6340	221 + 20 = 241 ^a	0	241 ^a	57,59
<i>C. elegans</i> (E)	20,009	1100 (2168) ^b	104 (208)	996 (1962)	66,69
<i>D. melanogaster</i> (E)	14,332	100 +	??	??	28; Harrison <i>et al.</i> , unpublished results
<i>A. thaliana</i> (E)	25,464	785	??	??	2
<i>Homo sapiens</i> (E)	~21,000 to ~39,000	??	~2900	??	27,85
<i>Homo sapiens</i> (E) (just chromosomes 21 + 22)	927	350	178	172	96

^a This total is for dORFs plus mORFs. dORFs are pseudogenic or disabled ORFs that comprise a large fragment of disabled protein sequence homology that is not part of an existing ORF annotation; mORFs (merged ORFs) arise from merging two existing ORF annotations by ignoring their intervening stop codon.⁵⁹

^b This is for a set of disabled protein sequence homologies, supported by protein/cDNA/EST homology evidence. The values in parentheses are upper estimates derived as described.⁶⁹

form of this protein takes Sup35p away from its normal functioning state, and can cause increased levels of nonsense-codon read-through in a particular strain, arguably leading to the full-length resurrection of ORFs that are apparently disabled. This can be seen as an evolutionary “buffering” effect, that enables a small amount of strain-specific variation to be maintained “in store”. Indeed, the ability to form the [PSI +] prion itself may have been selected to enable this buffering effect. Interestingly, a recent study on [PSI +]-engendered phenotypic diversity, showed that one strain is more flocculent when in the [PSI +] state than in the [psi -] state;⁶⁴ this may be due to the resurrection of the complete FLO8 reading frame, or other flocculin genes.

Worm versus fly: comparison in terms of their live and dead proteomes

Despite their comparable genome size (100 Mb for the worm, 120 Mb euchromatic for the fly), and the greater apparent biological complexity of the fly (more cells, longer lifespan, more complicated physiology), the worm (at present) has more genes. The original sequencing projects estimated 19,099 worm and 13,601 fly proteins, although the proteomes comprise comparable functional diversity at the sequence domain level.^{28,65–67} A recent gene prediction study for the fly genome has yielded 1042 additional candidate genes, potentially increasing the *Drosophila* gene total to >14,600 and the total proteome to >15,100.⁶⁸ Furthermore, alternative splicing for the fly may

Table 3. Comparison of the prevalent InterPro sequence motifs in the population of disabled ORFs and in the live proteome of yeast

Disabled ORFs/pseudogenes		Proteins	
No.	Description	No.	Description
12	WD40 (IPR001680) ^a	115	Eukaryotic kinase (IPR000719)
6	DUP membrane protein (IPR001142)	112	Serine–threonine protein kinase (IPR002290)
6	Mitochondrial electron transport (IPR001993)	99	WD40 (IPR001680)
6	Flocculin (IPR001389)	76	Dead-box helicase (IPR001410)
4	Helicase, C-terminal domain (IPR001650)	74	Helicase, C-terminal domain (IPR001650)
4	PIR repeat (IPR000420)	57	Fungal transcriptional regulatory protein (IPR001138)
3	BNR repeat (IPR002860)	57	AAA ATPase superfamily (IPR003593)
3	Zn-containing alcohol dehydrogenase (IPR002085)	55	TyA transposon protein (IPR001042)
3	Dead-box helicase (IPR001410)	54	RNA-binding region RNP-1 (IPR000504)
3	Fungal transcriptional regulatory protein (IPR001138)	53	C2H2 Zn finger (IPR000822)
3	SRP1/TIP1 stress-induced protein (IPR000992)		
3	DNA topoisomerase I DNA-binding domain (IPR003602)		

The name of each InterPro motif is given,⁵ along with its number in parentheses. These counts are for the pseudogenic population derived from combining dORFs and mORFs (see the text and footnote a in Table 2).

^a The 12 of these are all in one protein.

1135 be more extensive than at present documented
1136 (currently about 2% of the documented worm
1137 proteome arises from alternative splicing, and
1138 $\sim 7\%$ for the fly).^{28,65–67}

1139 What about the corresponding sizes of the
1140 pseudogene populations for these two organisms?
1141 Depending on the thresholds used, the worm
1142 genome appears to contain a moderately sized
1143 complement of >1100 pseudogenes.⁶⁹ Only a
1144 small proportion ($< \sim 5\%$) of the pseudogenes
1145 appear to be processed. In general, the number of
1146 pseudogenes associated with each family of pro-
1147 teins is not proportional to the size of the family.⁶⁹
1148 This would be the “default case” if duplicated
1149 pseudogenes were formed randomly from existing
1150 gene families. However, as shown in Table 4, the
1151 largest numbers of pseudogenes are associated
1152 with multiple families of seven-transmembrane
1153 chemoreceptors (these are also a class of “environ-
1154 mental response” proteins, which were observed
1155 above for yeast). Also common are families associ-
1156 ated with a reverse transcriptase and a trans-
1157 posase, which presumably reflects remnants of
1158 decayed transposons (obvious transposons were
1159 screened out before the pseudogene assignment).

1160 There are only 40 annotated pseudogenes for the
1161 fly genome, and a preliminary survey by the
1162 authors suggests at least ~ 60 more (P.M.H. *et al.*,
1163 unpublished results). (One should note, however,
1164 that an unknown number of gene annotations for
1165 either the fly or the worm may be shown to be
1166 pseudogenes, upon further characterization.) The
1167 cohort of olfactory receptors/chemoreceptors and
1168 other seven-transmembrane receptors in the worm
1169 (~ 1100) is almost a scale of magnitude larger than
1170 in the fly (~ 160 seven-transmembrane receptors).
1171 This perhaps indicates a recent evolutionary
1172 organism-specific expansion in these genes for the
1173 worm, or the converse (a contraction in number of
1174 members) for the fly.^{65,66,70} Their predominance
1175 in the worm pseudogene population is presumably
1176 related to this apparent expansion of seven-trans-
1177 membrane receptors in the worm. The substantial
1178 majority of these genes ($\sim 90\%$) appear to be
1179 organism-specific in the worm,⁷¹ although careful
1180 sequence analysis using hidden Markov models
1181 has found mammalian orthologs for ~ 170 of
1182 them.⁷² On a related note, of the estimated ~ 1000
1183 seven-transmembrane olfactory receptor (pseudo)-
1184 genes in the human genome, about two-thirds are
1185 expected to be pseudogenic.^{73,74}

1186 Interestingly, the families that have the largest
1187 number of associated pseudogenes are amongst
1188 the families that are most expanded in the worm
1189 relative to the fly (Table 5). We compared in detail
1190 the list of domain sequence families for the fly
1191 and worm proteomes from the InterPro database.⁵
1192 The families exclusive in this list to either organism
1193 are tabulated, as well as the most expanded large
1194 families (with 30 or more members) relative to the
1195 other organism (Table 5). Three of the largest of
1196 these are for the seven-transmembrane receptor
1197 families (Table 5).

1198 The small number of fly pseudogenes and the
1199 apparently small size of its proteome may be
1200 related to the overall genomic DNA deletion rate.
1201 The larger worm proteome may arise simply
1202 because factors such as genomic DNA deletion
1203 rates and chromosomal rearrangement have
1204 allowed it. It may be that the genomic DNA
1205 deletion rate in the fly (which was previously
1206 evidenced to be very high from the apparent rarity
1207 of true fly pseudogenes^{75–77}) hampers the main-
1208 tenance of recent gene duplications, so that they
1209 have less time to become evolutionarily useful.
1210 Experiments with transposable elements in
1211 *D. melanogaster* and the cricket genus *Laupala* indi-
1212 cate a very rapid loss of genomic DNA in
1213 *Drosophila*.^{78–80} *Drosophila* has an extremely high
1214 rate of chromosomal rearrangement.⁸¹ However,
1215 studies on families of worm chemoreceptor genes
1216 and pseudogenes suggest that the worm has a
1217 rather high genomic DNA deletion rate.^{70,82,83}
1218 Moreover, an analysis looking for small protein
1219 motifs selected from the Prosite database in inter-
1220 genic regions in the fly and the worm suggests
1221 that the fly has as many, if not more, over-
1222 represented motifs (pseudomotifs) than the
1223 worm.⁸⁴ These pseudomotifs may represent frag-
1224 ments of protein fossils. Thus, their prevalence in
1225 the fly in relation to the worm, may indicate that
1226 the fly has much pseudogenic material that has
1227 decayed substantially.

1228 Human: a large processed 1229 pseudogene population

1230 For the human genome, the determination of the
1231 number of pseudogenes is intimately inter-linked
1232 with the determination of the total gene number,
1233 as cDNA/EST coverage for a full range of human
1234 tissues is likely to take many years. The recent
1235 near-complete sequencing of the human genome
1236 has yielded numbers for the human gene total
1237 that seem surprisingly low, of the order of 23,000–
1238 40,000 genes.^{27,85} Efforts to estimate the number of
1239 human genes just prior to the publications of the
1240 sequenced genome, with one notable exception
1241 (which estimated $\sim 120,000$ human genes⁸⁶),
1242 yielded largely similar numbers to these, in the
1243 range $\sim 28,000$ to $\sim 35,500$.^{87–90} A recent compre-
1244 hensive annotation of the draft human genome
1245 estimated about 65,000–75,000 transcriptional
1246 units or genes in the genome.⁹¹

1247 Duplicated pseudogenes are more involved in
1248 the problem of gene prediction than processed
1249 pseudogenes: an exon with a disablement that is
1250 in the region of a gene may or may not be a
1251 part of the extant gene, making it difficult or
1252 impossible to determine if the gene is a pseudo-
1253 gene without cDNA/EST evidence. This is com-
1254 pounded by the prevalence of alternative
1255 splicing in the human genome; three indepen-
1256 dent surveys have shown that $\sim 40\%$ of genes
1257 encode alternatively spliced transcripts.^{92–94}
1258 Estimates for the proportion of gene annotations
1259
1260

1261
1262
1263
1264
1265
1266
1267
1268
1269
1270
1271
1272
1273
1274
1275
1276
1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295
1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323

Table 4. Largest families in terms of proteins and pseudogenes in the worm; adapted from previous family clustering⁶⁹

Pseudogenes		Proteins	
No.	Description	No.	Description
59	Reverse transcriptase (IPR000477)	216	Nuc. hormone receptor ligand-binding domain (IPR000536)
51	7-TM chemoreceptor family #1 (IPR000168, IPR003003)	193	7-TM chemoreceptor family #1 (IPR000168, IPR003003)
31	Unknown domain family #1 ^a	188	7-TM chemoreceptor family #2 (IPR000168)
27	7-TM chemoreceptor family #2 (IPR000168)	124	Eukaryotic kinase (IPR000719)
22	7-TM chemoreceptor family #3 (IPR000168)	93	MATH domain (IPR002083)
21	Major sperm protein (IPR000535)	70	7-TM receptor family #4 (IPR000276)
20	Unknown domain family #3 ^a	70	Guanylyl cyclase recep. tyr kinase (IPR001054)
19	Unknown domain family #4 ^a	70	Cytochrome P450 (IPR001128)
19	TcA transposase (IPR002492)	70	Tyr phosphatase (IPR000242)
17	7-TM receptor family #4 (IPR000276)	68	UDP-glucuronyl transferase (IPR002213)

Corresponding InterPro motifs for some families are indicated in brackets. The thickly outlined boxes are for families that occur in both the top ten pseudogenes and top ten protein families.

^a Those families do not have corresponding InterPro motifs.

1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349
1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386

1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403
1404
1405
1406
1407
1408
1409
1410
1411
1412
1413
1414
1415
1416
1417
1418
1419
1420
1421
1422
1423
1424
1425
1426
1427
1428
1429
1430
1431
1432
1433
1434
1435
1436
1437
1438
1439
1440
1441
1442
1443
1444
1445
1446
1447
1448
1449**Table 5.** Exclusive and expanded large families for assigned INTERPRO domains in the fly and worm proteomes

Largest exclusive to fly ^a		Largest exclusive to worm ^b		Most expanded in worm relative to fly ^b		Most expanded in fly relative to worm	
No.	Description	No.	Description	No. in worm (fly)	Description	No. in fly (worm)	Description
404	Insect cuticle protein (IPR000618)	624	7-TM chemo-receptor family (IPR000168, IPR003003)	60 (1)	DUF23 (IPR002875)	544 (15)	Chymotrypsin serine protease family S1 (IPR001314)
110	Alkaline phosphatase (IPR001952)	322	7-TM chemo-receptor family (IPR000168)	301 (6)	EB module (IPR002899)	950 (35)	Serine protease trypsin family (IPR001254)
99	Glycoside hydrolase family 22 (IPR001916)	276	DUF38 (IPR002900)	44 (1)	ET module (IPR002603)	161 (6)	Lipase (IPR000734)
73	Alpha-tocopherol transport protein (IPR001071)	238	ShK toxin domain (IPR003582)	339 (8)	MATH domain (IPR002083)	48 (2)	Peptidyl di-peptidase A M2 metallo-protease (IPR001548)
54	Hemocyanin (IPR000896)	237	DUF139 (IPR003341)	58 (2)	K + channel (IPR003280)	38 (2)	GMC oxido-reductase (IPR000172)
30	Acylphosphatase (IPR001792)	233	7-TM chemo-receptor family (IPR000168)	167 (6)	Major sperm protein (IPR000535)	37 (2)	NMDA receptor (IPR001508)
29	GYR motif (IPR004011)	184	pol-like reverse transcriptase (IPR003286)	71 (3)	TcA transposase family (IPR002492)	44 (3)	Chaperonin cpn60 60 kDa sub-unit (IPR001844)
26	Mitochondrial brown fat uncoupling protein (IPR002030)	148	SRG family integral membrane protein (IPR000609)	438 (37)	Nuclear hormone receptor ligand-binding domain (IPR000536)	47 (4)	Gamma tubulin (IPR002454)
25	Opsin (IPR001760)	145	Nematode cuticle collagen N-terminal domain (IPR002486)	861 (75)	F box domain (IPR001810)	35 (3)	Neutrophil cytosol factor 2 (IPR000108)
25	NF-κB/Rel/dorsal (IPR000451)	109	WSN (domain of unknown function) (IPR003125)	167 (17)	vWF type A domain (IPR002035)	76 (9)	Insect alcohol dehydrogenase (IPR002424)

These data are taken from the lists provided on the InterPro proteome analysis Website (<http://www.ebi.ac.uk/interpro>). The symbols and abbreviations are explained in Table 4. The boxed families occur also in the top ten pseudogene family list for worm.

^a The four lists are sorted in decreasing order of the degree of expansion. The degree of expansion in a family is simply the size of the family in one organism divided by its size in the other. Only families with 30 or more members in either organism are considered for this analysis.

^b The family numberings differ here from those in Table 2, as these are derived by motif scanning in individual sequences, whereas the Table 2 families are derived by our own sequence clustering procedure (see Table 2).

1450
1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511
1512

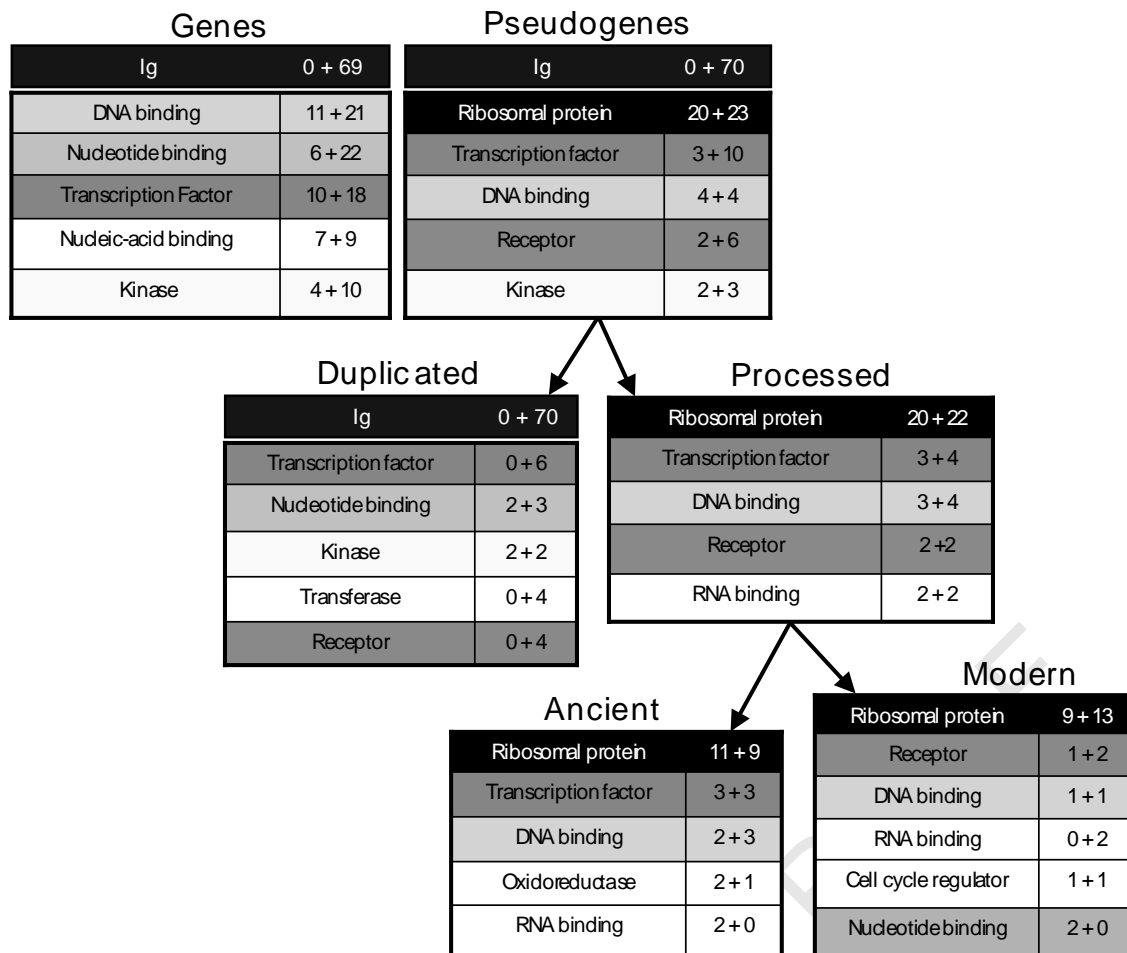


Figure 4. Functional categories of genes and pseudogenes in chromosomes 21 and 22: adapted from data given by Harrison *et al.*⁹⁶ Gene Ontology (GO) functional classes were assigned to predicted genes and pseudogenes for chromosomes 21 and 22 in combination. Those for pseudogenes are separated into processed and duplicated, with processed pseudogenes further separated into ancient and modern processed pseudogenes on the basis of their degree of sequence identity with the closest-matching human gene from the Ensembl data set (<http://www.ensembl.org>)

that may actually be pseudogenes lie in the range 4–22%.^{27,87,95}

Processed pseudogenes will be less likely to interfere with the accuracy of gene predictions; they will, on average, tend to be longer than the average human exon size, and comprise characteristic signals, including a C-terminal polyadenine tail.^{44,46} If they occur in relatively large numbers, they are also, in a sense, evidence that their parent gene is transcribed and most likely functional. Estimated numbers of processed pseudogenes in the human genome are substantial compared to those estimated for the gene total. In the completed chromosome 22 sequence, Dunham *et al.* initially predicted at least 545 genes and 134 pseudogenes (one for every ~4.1 genes).⁸⁷ They surmised that 82% of these pseudogenes were processed, as they contained single blocks of homology and lacked the characteristic exonic structure of the closest matching gene. This gives a predicted proportion of one processed pseudogene for every ~5.0 genes. Venter *et al.*, observed

evidence for at least ~2900 processed pseudogenes arising from their human gene set.⁸⁵ These were identified by searching for continuous spans of homology of >70% sequence identity over >70% of the length of the matching coding sequences from their gene annotations. No effort was made to look for the other characteristics of processed pseudogenes, such as evidence for polyadenylation. This data set of processed pseudogenes gives a smaller proportion of processed pseudogenes, in the region of about one for every ten genes. A survey by the authors of pseudogenes on chromosomes 21 and 22 that included searching for polyadenylation yielded an estimate of about one processed pseudogene for every four genes.⁹⁶ In this survey, we found that about half of all detected pseudogenes are processed (Table 2). The large amount of processing in the human genome may simply reflect its large amount of intergenic sequence and perhaps, the genomic mobility of transposable elements such as LINE-1.⁴⁵

The prevalence of the encoded proteins in the processed pseudogene population appears to be related to expression. Goncalves *et al.* analysed 181 genes that were reported to have one or more processed pseudogenes.⁹⁷ They found that such genes tend to be short, highly conserved and widely expressed. In the survey of ~2900 potential processed pseudogenes by Venter *et al.*⁸⁵ (noted above), by far the most prevalent class of transcripts (>60%) were for ribosomal proteins, which are very highly (and, of course, widely) expressed. The possibility of a large number of processed pseudogenes for ribosomal proteins was first noted during cloning of the mouse ribosomal protein rpL32⁹⁸. As shown in Figure 4, data by the authors from a survey of chromosomes 21 and 22 for processed and duplicated pseudogenes⁹⁶ also indicate that ribosomal proteins predominate in the processed pseudogene population, albeit, to less of an extent than in the survey by Venter *et al.*,⁸⁵ we found that ~20% of processed pseudogenes were ribosomal, and that there was little difference in this prevalence for either modern or ancient processed pseudogenes.

Figure 4 shows that the duplicated pseudogenes found in the survey of chromosomes 21 and 22 tend to be immunoglobulin gene fragments, reflecting their prevalence on chromosome 22. This preference continues the environmental-response theme discussed above for the worm and the yeast.

Concluding remarks

Comparing and contrasting the distribution of protein families in proteomes and in pseudogene populations gives us new perspectives on how proteomes evolve. A number of over-arching themes and implications are apparent. There are three distinct populations of pseudogenes.

Three types of pseudogenes

Prokaryotic pseudogene: dying genes resulting from a niche change

Prokaryotic pseudogenes appear to be genes that are dying and disappearing from the genome, in response to a fundamental niche change for an organism. In particular, there are now three bacterial pathogenic genomes (*M. leprae*, *Y. pestis* and *R. prowazekii*) that exhibit large-scale degradation of the proteome, with the lost or depleted families evidencing apparent niche change. In the most extreme case, *M. leprae* has large-scale patterning in its pseudogene population that indicates modular loss of metabolic pathways and branches of pathways, such as part of the anaerobic respiratory chain, when compared with *M. tuberculosis*, its closest sequenced relative. It is interesting, however, that this organism has lost *dnaQ*-mediated proofreading activities of DNA

polymerase III.⁵² Perhaps, this loss of function may actually have been selected so that removal of redundant genes could be accelerated. Although selection for deletion of pseudogenic DNA may not be sufficiently strong in eukaryote genomes,⁷⁹ there may be strong selection pressures for such deletion in small prokaryotic genomes that are undergoing niche change, and discarding many genes.

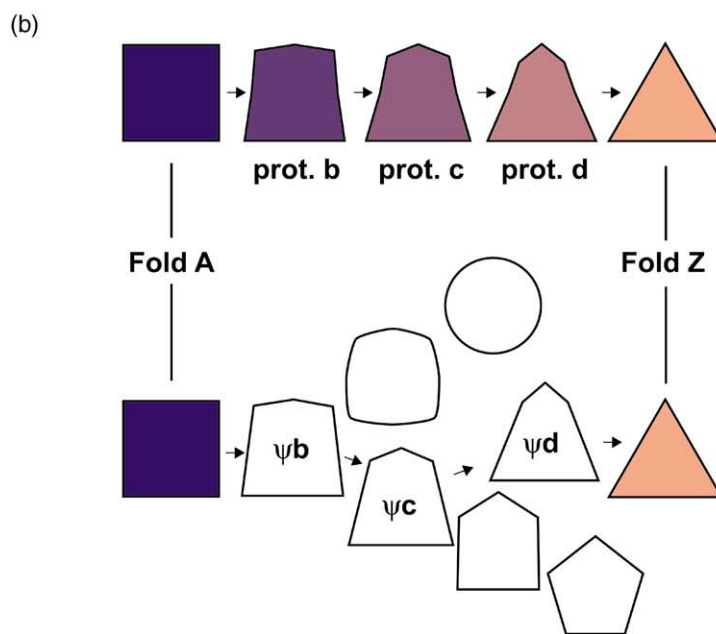
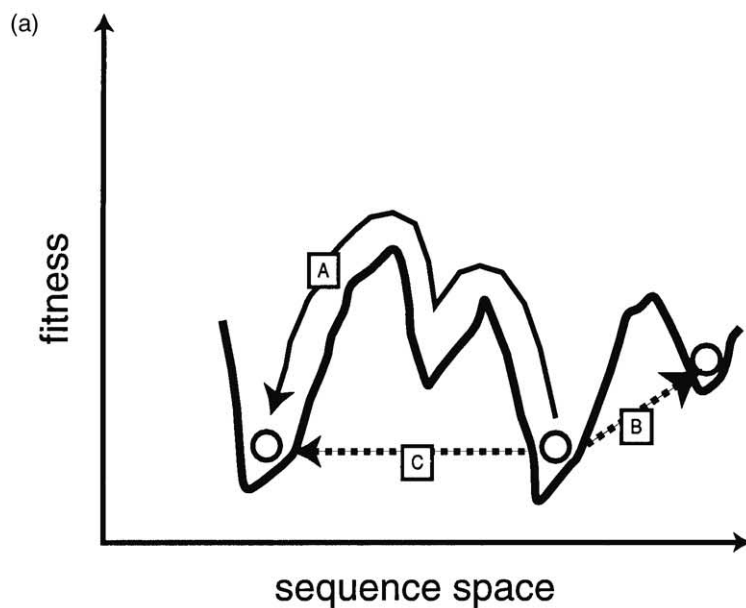
Eukaryotic processed pseudogenes: random insertion events

Processed pseudogenes arise from reverse-transcription of mRNA and re-integration into the genome. In humans, they are probably made as a by-product of LINE retrotransposition^{45,99}. That is, the processed pseudogene is formed from reverse transcribing a spliced mRNA into a cDNA using the reverse transcriptase from the LINE and re-integrating into the genome.^{45,99} Initial surveys suggest that their occurrence is largely based on simply random insertions, with their prevalence based on (1) the amount of mRNA to be inserted (expression levels) and (2) the amount of intergenic DNA available for insertion. The first factor accounts for the large numbers of ribosomal protein families found in processed pseudogenes.^{85,96} The second factor explains the large number of processed pseudogenes in the human genome, relative to the worm. It appears that the number of processed pseudogenes per 10⁶ bases of non-coding DNA is almost the same for both organisms. For human (chromosomes 21 and 22) the ratio is 2.6, which is 178 processed pseudogenes per 67 Mb of non-coding DNA. For the worm, the comparable number is 3.0, which is 208 per 70 Mb. (This ratio uses the high estimate for numbers of pseudogenes in the worm. It would decrease by 50% if one used the lower estimate (see Table 2).)

Eukaryotic duplicated pseudogenes: a resurrectable reservoir of extra parts for environmental response?

Duplicated eukaryotic pseudogenes appear to be most intriguing. They tend to arise for organism-specific environmental response functions. This tendency may reflect genomic mechanisms that an organism uses to generate proteins that deal with changes in its environment. We suggest below that pseudogenes or pseudogenic parts for such classes of gene may occasionally be resurrected and used to enable larger random leaps in sequence space (see below).

Eukaryotic pseudogenes tend to occur for organism-specific families. Pseudogenes in yeast are about twice as likely as a live protein to be yeast-specific.⁵⁹ Similarly, in the worm, the vast majority of the most prominent pseudogene families (those for the 7-TM chemoreceptors, major sperm protein and some unknown domains) are worm-specific or represent families vastly



expanded in the worm relative to the fly (Tables 4 and 5).

Pseudogenicity in eukaryotes appears to be linked to protein functions that are needed for environmental response, needing functional "breadth". In the worm, pseudogenicity is linked to 7-TM chemoreceptor families.⁶⁹ In the yeast, flocculins (which perform a variety of functions involving cell adhesion), growth-inhibitors, and stress-response proteins have the highest numbers of pseudogenes.⁵⁹ Finally, in the human, immunoglobulins have a high degree of pseudogenicity. For example, the immunoglobulin locus containing lambda variable-region gene segments on chromosome 22 is about 50% pseudogenic.⁹⁶ Also, a recent survey shows that there are ~1000 olfactory recep-

tors in the human genome, with 60% of these pseudogenic.⁷³

Pseudogene resurrection as a general evolutionary mechanism

In certain cases, as a rare or occasional evolutionary event, the resurrection of duplicated pseudogenic DNA to an expressed protein may enable sampling of more sequence space for a protein or protein family (Figure 5(a)). In particular, pseudogenes or parts of pseudogenes may be re-used, after having drifted randomly without selection for a period of evolution. The idea of such "untranslatable intermediates" in the evolution of a protein was first postulated about 30 years ago

Figure 5. Aspects of pseudogene resurrection as an evolutionary mechanism. (a) A schematic evolutionary landscape showing a sequence (represented by an open circle) in a favourable fitness minimum, with three evolutionary routes A, B and C. Route A (continuous line) arises from mutation under the pressures of natural selection. Route B (dotted line) represents what happens when a sequence undergoes random drift as a pseudogene, but which, when "resurrected" as a genic sequence, is unfit. Route C represents what happens when a sequence undergoes random drift as a pseudogene, but reaches another favourable fitness minimum in a shorter span of time than would be possible under continuous natural selection. (b) The top panel shows the conventional view of protein fold evolution where every intermediate along the pathway has to be transcribed and translated. The bottom panels shows a pathway that involves pseudogenic fragments.

by Koch.¹⁰⁰ Although generally one would expect this mechanism to produce unviable or unfavourable leaps in sequence space, occasionally it may provide a shorter evolutionary route to another favourable evolutionary energetic minimum (Figure 5(a)).

There are number of cases that one can point to as evidence of such resurrection. A pseudogene of bovine seminal ribonuclease that lay dormant for ~20 million years, appears to have been resurrected to form a functioning gene, probably *via* a gene conversion event.¹⁰¹ As discussed above, the presence of the [PSI +] prion in yeast strains may enable resurrection or extension of ORFs from the yeast genome that have been able to drift without selection pressures since the occurrence of their disrupting mutations.⁶⁴ The large cohort of pseudogenes for chemo- or olfactory receptors (ORs) in metazoans (60% of the ORs in the human genome are pseudogenetic) may be resurrectable by gene conversion events. There appears to have been a large number of gene conversion events (>20) in a cluster of olfactory receptors on chromosome 17 over the course of primate evolution.¹⁰² This cluster contains 16 OR genes and 6 OR pseudogenes in the human genomic DNA. Gene conversion events in OR gene clusters may help to generate diversity at the odorant binding site.¹⁰² Occasional resurrection of OR pseudogenes by gene conversion may contribute to this generation of diversity in binding capability. Finally, in the chicken, diversity of immunoglobulin heavy chain variable-region gene segments appears to be generated by gene conversion of a single functional gene with >80 pseudogenetic gene segments.¹⁰³

Resurrectable pseudogenes may help resolve a paradox about protein fold evolution

Considering duplicated pseudogenes as a resurrectable reservoir of diversity may help to resolve an evolutionary paradox presented by structural biology. How do new folds evolve? An early observation from structural genomics analyses was that there appear to be folds unique to certain phylogenetic groups.^{16,25} For instance, an initial analysis showed that of 275 folds, 46 were present only in eubacteria and 73 only in eukaryotes, and of the 229 total folds in eukaryotes, 20 were only in plants and 90 only in animals.¹⁶ How does one get new unique folds in certain phylogenetic groups? As shown in Figure 5(b), in some cases it may be difficult to imagine a scenario for this where each intermediate form has to be a functioning protein that is transcribed and translated. (This is in contrast to other evolutionary pathways, where functioning and selected intermediates are more plausible.) One can speculate that resurrectable pseudogenes could eliminate this paradox to some degree. A sequence comprising a particular domain fold or (more likely) part of a domain could become pseudogenetic. It could

then drift freely as a pseudogene, and evolve to a new domain fold upon or after resurrection. In this scheme, each intermediate does not have the constraint that it be a folded functional protein.

Elimination of pseudogenes

Pseudogenes can be eliminated from the genome due to deletion events. There is obviously greater pressure to do this for prokaryotes than for eukaryotes. Thus, it is important to point out that the lack of a large pseudogene population for prokaryotes does not imply that an organism has not undergone gene loss as drastic as that seen in *M. leprae*, over a similar evolutionary period. An organism with a higher rate of genomic DNA deletion would delete pseudogenetic DNA more efficiently, and we would therefore not see such a large pseudogene population at present. For *M. leprae*, it may be that the rate of disablement of ORFs is raised, without there being a concomitant increase in the rate of deletion of intergenic DNA. Rates of intergenic DNA deletion vary widely from organism to organism.⁸⁰ For the eukaryote *Drosophila*, although the overall genomic deletion rate is very high, the observed spectrum of deletion sizes in transposable elements implies that it has not been selected for to aid genome compaction.⁷⁹ The *Drosophila* genomic DNA deletion rate seems to explain the dearth of pseudogenes in the fly that are detectable by sequence homology.^{78,80} To find very decayed remnants of proteins in the genome not amenable to sequence alignment, we are currently developing a probabilistic approach based on scanning the genome for decayed protein motifs (termed pseudomotifs).⁸⁴ Over even longer evolutionary periods, gene loss can be inferred from careful comparative proteome analysis. For example, comparison of the *S. cerevisiae* proteome with the near-complete proteome of the fission yeast *Schizosaccharomyces pombe*, indicates the possible loss of about 300 proteins in *S. cerevisiae*, and provides an explanation for the small degree of gene splicing in *S. cerevisiae*, involving deletion of signalosome and spliceosome components.¹⁰⁴ (The fission yeast has extensive gene splicing.)

Power-law behaviour and the size of duplicated pseudogene populations

We noted above that the size of protein families in the live proteomes is governed by a power-law distribution (Figure 1). This behaviour is observed for the distribution of protein families in the pseudogene population (the dead proteome) of chromosomes 21 and 22, and of the worm genome^{69,96} (Figure 1). (It is observed even for the distribution of pseudomotifs in the fly and worm genomes.⁸⁴) This may imply that conservation pressures do not cause such power-law behaviour, but rather the flow of change from old to new families over evolution. Qian *et al.*,²² found that the power-law distribution of protein families and

2017 folds is well described by a simple model in which
2018 existing gene sequences can be duplicated, but
2019 with the occasional creation or addition of a novel
2020 gene.

2021 Thus, despite the great differences in specific
2022 protein families prevalent in various organisms in
2023 both the living and the dead proteomes, we can
2024 see a clear commonality in their occurrence: one
2025 has a few families occurring many times and most
2026 occurring just a few times. In all aspects of geno-
2027 mic biology, one never gets a uniform distribution
2028 of occurrence over families.

2032 Acknowledgments

2033 Thanks to Julian Gough (MRC) for providing data on
2034 protein folds in eukaryotic proteomes, and to Nick
2035 Luscombe and Jiang Qian for part of [Figure 1](#). Thanks
2036 also to Hedi Hegyi, Suganthi Balasubramaniam, Paul
2037 Bertone, Nathaniel Echols, Nick Luscombe and Ted
2038 Johnson for help, and to Alan Weiner for comments on
2039 pseudogene formation. M.G. acknowledges support
2040 from the Keck foundation and the NIH structural
2041 genomics initiative (P50 GM62413-01).

2044 References

- 2045 1. Friedman, R. & Hughes, A. L. (2001). Gene dupli-
2046 cation and the structure of eukaryotic genomes.
2047 *Genome*, **11**, 373–381.
- 2048 2. Arabidopsis Genome Initiative, T. (2000). Analysis
2049 of the genome sequence of the flowering plant
2050 *Arabidopsis thaliana*. *Nature*, **408**, 796–815.
- 2051 3. Huynen, M. A. & van Nimwegen, E. (1998). The fre-
2052 quency distribution of gene family sizes in com-
2053 plete genomes. *Mol. Biol. Evol.* **15**, 583–589.
- 2054 4. Nevill-Manning, C. G., Wu, T. D. & Brutlag, D. L.
2055 (1998). Highly specific protein sequence motifs for
2056 genome analysis. *Proc. Natl Acad. Sci. USA*, **95**,
2057 5865–5871.
- 2058 5. Apweiler, R., Attwood, T. K., Bairoch, A., Bateman,
2059 A., Birney, E., Biswas, M. *et al.* (2000). InterPro—an
2060 integrated documentation resource for protein
2061 families, domains and functional sites. *Bio-*
2062 *informatics*, **16**(12), 1145–1150.
- 2063 6. Yona, G., Linial, N. & Linial, M. (2000). ProtoMap:
2064 automatic classification of protein sequences and
2065 hierarchy of protein families. *Nucl. Acids Res.* **28**(1),
2066 49–55.
- 2067 7. Krause, A., Stoye, J. & Vingron, M. (2000). The
2068 SYSTERS protein sequence cluster set. *Nucl. Acids*
2069 *Res.* **28**(1), 270–272.
- 2070 8. Murzin, A. G., Brenner, S. E., Hubbard, T. &
2071 Chothia, C. (1995). SCOP: a structural classification
2072 of proteins database for the investigation of
2073 sequences and structures. *J. Mol. Biol.* **247**(4),
2074 536–540.
- 2075 9. Pearl, F., Todd, A. E., Bray, J. E., Martin, A. C.,
2076 Salamov, A. A., Suwa, M. *et al.* (2000). Using the
2077 CATH domain database to assign structures and
2078 functions to the genome sequences. *Biochem. Soc.*
2079 *Trans.* **28**(2), 269–275.
- 2080 10. Ashburner, M., Ball, C. A., Blake, J. A., Botstein, D.,
2081 Butler, H., Cherry, J. M. *et al.* (2000). Gene ontology:
2082 tool for the unification of biology. The Gene
2083 Ontology Consortium. *Nature Genet.* **25**(1), 25–29.
- 2084 11. Mewes, H. W., Frishman, D., Gruber, C., Geier, B.,
2085 Haase, D., Kaps, A. *et al.* (2000). MIPS: a database
2086 for genomes and protein sequences. *Nucl. Acids*
2087 *Res.* **28**(1), 37–40.
- 2088 12. Riley, M. & Space, D. B. (1996). Genes and proteins
2089 of *Escherichia coli* (GenProtEc). *Nucl. Acids Res.*
2090 **24**(1), 40.
- 2091 13. Gerstein, M. & Hegyi, H. (1998). Comparing gen-
2092 omes in terms of protein structure: surveys of a
2093 finite parts list. *FEMS Microbiol. Rev.* **24**, 1–28.
- 2094 14. Gerstein, M. (1998). How representative are the
2095 known structures of the proteins in a complete ge-
2096 nome? A comprehensive structural census. *Fold. Des.*
2097 **3**, 497–512.
- 2098 15. Altschul, S. F., Madden, T. L., Schaffer, A. A.,
2099 Zhang, J., Zhang, Z., Miller, W. & Lipman, D. J.
2100 (1997). Gapped BLAST and PSI-BLAST: a new gen-
2101 eration of protein database search programs. *Nucl.*
2102 *Acids Res.* **25**(17), 3389–3402.
- 2103 16. Gerstein, M. & Levitt, M. (1997). A structural census
2104 of the current population of protein sequences. *Proc.*
2105 *Natl Acad. Sci. USA*, **94**(22), 11911–11916.
- 2106 17. Sonnhammer, E. L. & Durbin, R. (1997). Analysis of
2107 protein domain families in *Caenorhabditis elegans*.
2108 *Genomics*, **46**(2), 200–216.
- 2109 18. Salamov, A. A., Suwa, M., Orengo, C. A. &
2110 Swindells, M. B. (1999). Genome analysis: assigning
2111 protein coding regions to three-dimensional struc-
2112 tures. *Protein Sci.* **8**(4), 771–777.
- 2113 19. Hegyi, H., Lin, J. & Gerstein, M. (2002). Structural
2114 genomics analysis for twenty proteomes: common,
2115 shared and unique protein folds. *Proteins: Struct.*
2116 *Funct. Genet.* in press.
- 2117 20. Wolf, Y. I., Brenner, S. E., Bash, P. A. & Koonin, E. V.
2118 (1999). Distribution of protein folds in the three
2119 superkingdoms of life. *Genome Res.* **9**, 17–26.
- 2120 21. Wolf, Y. I., Grishin, N. V. & Koonin, E. V. (2000).
2121 Estimating the number of protein folds and families
2122 from complete genome data. *J. Mol. Biol.* **299**,
2123 897–905.
- 2124 22. Qian, J., Luscombe, N. M. & Gerstein, M. (2001).
2125 Protein family and fold occurrence in genomes:
2126 power-law behaviour and evolutionary model.
2127 *J. Mol. Biol.* **313**(4), 673–681.
- 2128 23. Lin, J. & Gerstein, M. (2000). Whole-genome trees
2129 based on the occurrence of folds and orthologs:
2130 implications for comparing genomes on different
2131 levels. *Genome Res.* **10**, 808–818.
- 2132 24. Gough, J. & Chothia, C. (2002). SUPERFAMILY:
2133 HMMs representing all proteins of known struc-
2134 ture. SCOP sequence searches, alignments and ge-
2135 nome assignments. *Nucl. Acids Res.* **30**(1), 268–272.
- 2136 25. Gerstein, M. (1997). A structural census of genomes:
2137 comparing bacterial, eukaryotic, and archaeal ge-
2138 nomes in terms of protein structure. *J. Mol. Biol.*
2139 **274**(4), 562–576.
- 2140 26. Yanai, I., Camacho, C. J. & DeLisi, C. (2000). Predic-
2141 tions of gene family distributions in microbial ge-
2142 nomes: evolution by gene duplication and
2143 modification. *Phys. Rev. Letters*, **85**, 2641–2644.
- 2144 27. Lander, E. S., Linton, L. M., Birren, B., Nusbaum, C.,
2145 Zody, M. C., Baldwin, J. *et al.* (2001). Initial sequen-
2146 cing and analysis of the human genome. Inter-
2147 national Human Genome Sequencing Consortium.
2148 *Nature*, **409**(6822), 860–921.
- 2149 28. Adams, M. D., Celniker, S. E., Holt, R. A., Evans,
2150 C. A., Gocayne, J. D., Amanatides, P. G. *et al.*

- 2143 (2000). The genome sequence of *Drosophila*
2144 *melanogaster*. *Science*, **287**, 2185–2195.
- 2145 29. Vision, T. J., Brown, D. G. & Tanksley, S. D. (2000).
2146 The origins of genomic duplications in *Arabidopsis*.
2147 *Science*, **290**, 2114–2117.
- 2148 30. Lynch, M. & Conery, J. S. (2000). The evolutionary
2149 fate and consequences of duplicate genes. *Science*,
2150 **290**(5494), 1151–1155.
- 2151 31. Seoighe, C. & Wolfe, K. H. (1999). Yeast genome
2152 evolution in the post-genome era. *Curr. Opin.*
2153 *Microbiol.* **2**(5), 548–554.
- 2154 32. Seoighe, C. & Wolfe, K. H. (1999). Updated map of
2155 duplicated regions in the yeast genome. *Gene*,
2156 **238**(1), 253–261.
- 2157 33. Wolfe, K. H. & Shields, D. C. (1997). Molecular
2158 evidence for an ancient duplication of the entire
2159 yeast genome. *Nature*, **387**(6634), 708–713.
- 2160 34. Llorente, B., Durrens, P., Malpertuy, A., Aigle, M.,
2161 Artiguenave, F., Blandin, G. *et al.* (2000). Genomic
2162 exploration of the hemiascomycetous yeasts: 20.
2163 Evolution of gene redundancy compared to
2164 *Saccharomyces cerevisiae*. *FEBS Letters*, **487**(1),
2165 122–133.
- 2166 35. Winzeler, E. A., Shoemaker, D. D., Astromoff, A.,
2167 Liang, H., Anderson, K., Andre, B. *et al.* (1999).
2168 Functional characterization of the *S. cerevisiae* ge-
2169 nome by gene deletion and parallel analysis. *Science*,
2170 **285**, 901–906.
- 2171 36. Delneri, D., Brancia, F. L. & Oliver, S. G. (2001).
2172 Towards a truly integrative biology through the
2173 functional genomics of yeast. *Curr. Opin. Biotech.*
2174 **12**, 87–91.
- 2175 37. Mushegian, A. (1999). The minimal genome con-
2176 cept. *Curr. Opin. Genet. Dev.* **9**(6), 709–714.
- 2177 38. Wagner, A. (2000). Robustness against mutations in
2178 genetic networks of yeast. *Nature Genet.* **24**,
2179 355–361.
- 2180 39. Tautz, D. (2000). A genetic uncertainty problem.
2181 *Trends Genet.* **16**, 475–477.
- 2182 40. Thatcher, J. W., Shaw, J. M. & Dickinson, W. J.
2183 (1998). Marginal fitness contributions of non-
2184 essential genes in yeast. *Proc. Natl Acad. Sci. USA*,
2185 **95**, 253–257.
- 2186 41. Velculescu, V. E., Zhang, L., Zhou, W., Vogelstein, J.,
2187 Basrai, M. A., Bassett, D. E. *et al.* (1997). Characteriz-
2188 ation of the yeast transcriptome. *Cell*, **88**(2),
2189 243–251.
- 2190 42. Jansen, R. & Gerstein, M. (2000). Analysis of the
2191 yeast transcriptome with structural and functional
2192 categories: characterizing highly expressed
2193 proteins. *Nucl. Acids Res.* **28**(6), 1481–1488.
- 2194 43. Hirsh, A. E. & Fraser, H. B. (2001). Protein dispensa-
2195 bility and rate of evolution. *Nature*, **411**, 1046–1049.
- 2196 44. Vanin, E. F. (1985). Processed pseudogenes: charac-
2197 teristics and evolution. *Annu. Rev. Genet.* **19**,
2198 253–272.
- 2199 45. Esnault, C., Maestre, J. & Heidmann, T. (2000).
2200 Human LINE retrotransposons generate processed
2201 pseudogenes. *Nature Genet.* **24**, 363–367.
- 2202 46. Mighell, A. J., Smith, N. R., Robinson, P. A. &
2203 Markham, A. F. (2000). Vertebrate pseudogenes.
2204 *FEBS Letters*, **468**, 109–114.
- 2205 47. Eisen, J. A. (2000). Horizontal gene transfer among
2206 microbial genomes: new insights from complete
2207 genome analysis. *Curr. Opin. Genet. Dev.* **10**,
2208 606–611.
- 2209 48. Blattner, F. R., Plunkett, G., III, Bloch, C. A., Perna,
2210 N. T., Burland, V., Riley, M. *et al.* (1997). The com-
2211 plete genome sequence of *Escherichia coli* K-12.
2212 *Science*, **277**(5331), 1453–1474.
- 2213 49. Perna, N. T., Plunkett, G., 3rd, Burland, V., Mau, B.,
2214 Glasner, J. D., Rose, D. J. *et al.* (2001). Genome
2215 sequence of enterohaemorrhagic *Escherichia coli*
2216 O157:H7. *Nature*, **409**(6819), 529–533.
- 2217 50. Lan, R. & Reeves, P. R. (2000). Intraspecies variation
2218 in bacterial genomes: the need for a species genome
2219 concept. *Trends Microbiol.* **8**, 396–401.
- 2220 51. Boucher, Y., Nesbo, C. L. & Doolittle, W. F. (2001).
2221 Microbial genomes: dealing with diversity. *Curr.*
2222 *Opin. Microbiol.* **4**, 285–289.
- 2223 52. Cole, S. T., Eigimeier, K., Parkhill, J., James, K. D.,
2224 Thomson, N. R., Wheeler, P. R. *et al.* (2001). Massive
2225 gene decay in the leprosy bacillus. *Nature*, **409**,
2226 1007–1011.
- 2227 53. Andersson, S. G., Zomorodipour, A., Andersson,
2228 J. O., Sicheritz-Ponten, T., Alsmark, U. C.,
2229 Podowski, R. M. *et al.* (1998). The genome sequence
2230 of *Rickettsia prowazekii* and the origin of mito-
2231 chondria. *Nature*, **396**(6707), 133–140.
- 2232 54. Ogata, H., Audic, S., Renesto-Audiffren, P.,
2233 Fournier, P. E., Barbe, V., Samson, D. *et al.* (2001).
2234 Mechanisms of evolution in *Rickettsia conorii* and
2235 *R. prowazekii*. *Science*, **293**(5537), 2093–2098.
- 2236 55. Andersson, J. O. & Andersson, S. G. (2001). Pseudo-
2237 genes, junk DNA and the dynamics of rickettsia
2238 genomes. *Mol. Biol. Evol.* **18**, 829–839.
- 2239 56. Parkhill, J., Wren, B. W., Thomson, N. R., Titball,
2240 R. W., Holden, M. T., Prentice, M. B. *et al.* (2001).
2241 Genome sequence of *Yersinia pestis*, the causative
2242 agent of plague. *Nature*, **413**(6855), 523–527.
- 2243 57. Goffeau, A., Barrell, B. G., Bussey, H., Davis, R. W.,
2244 Dujon, B., Feldmann, H. *et al.* (1996). Life with 6000
2245 genes. *Science*, **274**, 546, 563–567.
- 2246 58. Cherry, J. M., Adler, C., Ball, C., Chervitz, S. A.,
2247 Dwight, S. S., Hester, E. T. *et al.* (1998). SGD:
2248 *saccharomyces* genome database. *Nucl. Acids Res.*
2249 **26**(1), 73–79.
- 2250 59. Harrison, P. M., Kumar, A., Lan, N., Echols, N.,
2251 Snyder, M. & Gerstein, M. (2002). A small reservoir
2252 of disabled ORFs in the sequenced yeast genome
2253 and its implications for the dynamics of proteome
2254 evolution. *J. Mol. Biol.* in press.
- 2255 60. Liu, H., Styles, C. A. & Fink, G. R. (1996).
2256 *S. cerevisiae* S288C has a mutation in FLO8, a gene
2257 required for filamentous growth. *Genetics*, **144**,
2258 967–978.
- 2259 61. Serio, T. R. & Lindquist, S. L. (2000). Protein-only
2260 inheritance in yeast: something to get [PSI +]-ched
2261 about. *Trends Cell Biol.* **10**, 98–105.
- 2262 62. Eaglestone, S. S., Cox, B. S. & Tuite, M. F. (1999).
2263 Translation termination efficiency can be regulated
2264 in *S. cerevisiae* by environmental stress through a
2265 prion-mediated mechanism. *EMBO J.* **18**,
2266 1974–1981.
- 2267 63. Tuite, M. F. (2000). Yeast prions and their prion-
2268 forming domain. *Cell*, **100**, 289–292.
- 2269 64. True, H. L. & Lindquist, S. L. (2000). A yeast prion
2270 provides a mechanism for genetic variation and
2271 phenotypic diversity. *Nature*, **407**(6803), 477–483.
- 2272 65. Chervitz, S. A., Aravind, L., Sherlock, G., Ball, C. A.,
2273 Koonin, E. V., Dwight, S. S. *et al.* (1998). Comparison
2274 of the complete protein sets of worm and yeast:
2275 orthology and divergence. *Science*, **282**, 2022–2028.
- 2276 66. *C. elegans* Sequencing Consortium, T. (1998).
2277 Genome sequence of the nematode *C. elegans*: a
2278 platform for investigating biology. *Science*, **282**,
2279 2012–2018.

- 2269 67. Rubin, G. M., Yandell, M. D., Wortman, J. R., Gabor
2270 Miklos, G. L., Nelson, C. R., Hariharan, I. K. *et al.*
2271 (2000). Comparative genomics of the eukaryotes.
2272 *Science*, **287**, 2204–2215. 2332
- 2273 68. Gopal, S., Schroeder, M., Pieper, U., Sczyrba, A.,
2274 Aytekin-Kurban, G., Bekiranov, S. *et al.* (2001).
2275 Homology-based annotation yields 1042 new candi-
2276 date genes in the *Drosophila melanogaster* genome.
2277 *Nature Genet.* **27**(3), 337–340. 2333
- 2278 69. Harrison, P. M., Echols, N. & Gerstein, M. (2001).
2279 Digging for dead genes: an analysis of the charac-
2280 teristics and distribution of the pseudogene popu-
2281 lation in the *C. elegans* genome. *Nucl. Acids Res.* **29**,
2282 818–830. 2334
- 2283 70. Robertson, H. M. (2000). The large srh family of
2284 chemoreceptor genes in *Caenorhabditis nematodes*
2285 reveals processes of genome evolution involving
2286 large duplications and deletions and intron gains
2287 and losses. *Genome Res.* **10**(2), 192–203. 2337
- 2288 71. Bargmann, C. I. (1998). Neurobiology of the
2289 *Caenorhabditis elegans* genome. *Science*, **282**,
2290 2028–2033. 2338
- 2291 72. Remm, M. & Sonnhammer, E. (2000). Classification
2292 of transmembrane protein families in the *Caenorhabditis elegans* genome and identification of human
2293 orthologs. *Genome Res.* **10**, 1679–1689. 2339
- 2294 73. Glusman, G., Yanai, I., Rubin, I. & Lancet, D. (2001).
2295 The complete human olfactory subgenome. *Genome*
2296 *Res.* **11**, 685–702. 2340
- 2297 74. Zozulya, S., Echeverri, F. & Nguyen, T. (2001). The
2298 human olfactory receptor repertoire. *Genome Biol.* **2**,
2299 research0018.0011–research0018.0012. 2341
- 2300 75. Robin, G. C. Q., Russell, R. J., Cutler, D. J. &
2301 Oakeshott, J. G. (2000). The evolution of an alpha-
2302 esterase pseudogene inactivated in the *Drosophila*
2303 *melanogaster* lineage. *Mol. Biol. Evol.* **17**, 563–575. 2342
- 2304 76. Currie, P. D. & Sullivan, D. T. (1994). Structure,
2305 expression and duplication of genes which encode
2306 phosphoglyceromutase of *Drosophila melanogaster*.
2307 *Genetics*, **138**, 353–363. 2343
- 2308 77. Sullivan, D. T., Starmer, W. H., Curtiss, S. W.,
2309 Menotti-Raymond, M. & Yum, J. (1994). Unusual
2310 molecular evolution of an Adh pseudogene in
2311 *Drosophila*. *Mol. Biol. Evol.* **11**, 443–458. 2344
- 2312 78. Petrov, D. A., Lozovskaya, E. R. & Hartl, D. L.
2313 (1996). High intrinsic rate of DNA loss in *Drosophila*.
2314 *Nature*, **384**, 346–349. 2345
- 2315 79. Petrov, D. A. & Hartl, D. L. (2000). Pseudogene
2316 evolution and natural selection for a compact ge-
2317 nome. *J. Heredit.* **91**, 221–227. 2346
- 2318 80. Petrov, D. A. (2001). Evolution of genome size: new
2319 approaches to an old problem. *Trends Genet.* **17**,
2320 23–28. 2347
- 2321 81. Ranz, J. M., Casals, F. & Ruiz, A. (2001). How malle-
2322 able is the eukaryotic genome? Extreme rate of
2323 chromosomal rearrangement in the genus *Drosophila*.
2324 *Genome Res.* **11**, 230–239. 2348
- 2325 82. Robertson, H. M. (1998). Two large families of
2326 chemoreceptor genes in the nematodes *Caenorhabditis elegans* and *Caenorhabditis briggsae* reveal
2327 extensive gene duplication, diversification, move-
2328 ment and intron loss. *Genome Res.* **8**, 449–463. 2349
- 2329 83. Robertson, H. M. (2001). Updating the str and srj
2330 (stl) families of chemoreceptors in *Caenorhabditis nematodes* reveals frequent gene movement within
2331 and between chromosomes. *Chem. Senses*, **26**(2),
2332 151–159. 2350
- 2333 84. Zhang, Z., Harrison, P.M. & Gerstein, M. (2002).
2334 Digging deep for ancient relics: a survey of protein
2335 motifs in the intergenic DNA of four eukaryotic
2336 genomes. submitted 2351
- 2337 85. Venter, J. C., Adams, M. D., Myers, E. W., Li, P. W.,
2338 Mural, R. J., Sutton, G. G. *et al.* (2001). The sequence
2339 of the human genome. *Science*, **291**(5507),
2340 1304–1351. 2352
- 2341 86. Liang, F., Holt, I., Pertea, G., Karamychea, S.,
2342 Salzberg, S. & Quackenbush, J. (2000). GENE index
2343 analysis of the human genome estimates approxi-
2344 mately 120,000 genes. *Nature Genet.* **24**, 239–240. 2353
- 2345 87. Dunham, I., Shimizu, N., Roe, B. A., Chissole, S.,
2346 Hunt, A. R., Collins, J. E. *et al.* (1999). The DNA
2347 sequence of human chromosome 22. *Nature*,
2348 **402**(6761), 489–495. 2354
- 2349 88. Hattori, M., Fujiyama, A., Taylor, T. D., Watanabe,
2350 H., Yada, T., Park, H. S. *et al.* (2000). The DNA
2351 sequence of human chromosome 21. The chromo-
2352 some 21 mapping and sequencing consortium.
2353 *Nature*, **405**(6784), 311–319. 2355
- 2354 89. Crollius, H. R., Jaillon, O., Bernot, A., Dasilva, C.,
2355 Bouneau, L., Fischer, C. *et al.* (2000). Estimate of
2356 human gene number provided by genome-wide
2357 analysis using *Tetraodon nigroviridis* DNA sequence.
2358 *Nature Genet.* **25**, 235–238. 2359
- 2359 90. Ewing, B. & Green, P. (2000). Analysis of expressed
2360 sequence tags indicates 35,000 human genes. *Nature*
2361 *Genet.* **232**, 232–233. 2362
- 2362 91. Wright, F. A., Lemon, W. J., Zhao, W. D., Sears, R.,
2363 Zhuo, D., Wang, J. P. *et al.* (2001). A draft annotation
2364 and overview of the human genome. *Genome Biol.*
2365 **2**(7). 2366
- 2366 92. Mironov, A. A., Fickett, J. W. & Gelfand, M. S.
2367 (1999). Frequent alternative splicing of human
2368 genes. *Genome Res.* **9**, 1288–1293. 2369
- 2369 93. Brett, D., Hanke, J., Lehmann, G., Haase, S.,
2370 Delbruck, S., Krueger, S. *et al.* (2000). EST compari-
2371 son indicates 38% of human mRNAs contain
2372 possible alternative splice forms. *FEBS Letters*, **474**,
2373 83–86. 2374
- 2374 94. Modrek, B., Resch, A., Grasso, C. & Lee, C. (2001).
2375 Genome-wide detection of alternative splicing in
2376 expressed sequences of human genes. *Nucl. Acids*
2377 *Res.* **29**(13), 2850–2859. 2378
- 2378 95. Yeh, R.-F., Lim, L. P. & Burge, C. (2001). Compu-
2379 tational inference of homologous gene structures in
2380 the human genome. *Genome Res.* **11**, 803–816. 2381
- 2381 96. Harrison, P. M., Hegyi, H., Balasubramaniam, S.,
2382 Luscombe, N. M., Bertone, P., Echols, N. *et al.*
2383 (2002). Molecular fossils in the human genome:
2384 Identification and analysis of the pseudogenes on
2385 chromosomes 21 and 22. *Genome Res.* in press. 2386
- 2386 97. Goncalves, I., Duret, L. & Mouchiroud, D. (2000).
2387 Nature and structure of human genes that generate
2388 retropseudogenes. *Genome Res.* **10**(5), 672–678. 2389
- 2389 98. Dudov, K. P. & Perry, R. P. (1984). The gene family
2390 encoding the mouse ribosomal protein L32 contains
2391 a uniquely expressed intron-containing gene and an
2392 unmutated processed gene. *Cell*, **37**(2), 457–468. 2393
- 2393 99. Weiner, A. M. (2000). Do all SINES lead to LINES?
2394 *Nature Genet.* **24**, 332–333. 2394
- 2394 100. Koch, A. L. (1972). Enzyme evolution I. The
2395 importance of untranslatable intermediates.
2396 *Genetics*, **72**, 297–316. 2397
- 2397 101. Trabesinger-Ruef, N., Jermann, T., Zankel, T.,
2398 Durrant, B., Frank, G. & Benner, S. A. (1996).
2399 Pseudogenes in ribonuclease evolution: a source of
2400 new biomacromolecular function? *FEBS Letters*,
2401 **382**(3), 319–322. 2399

- 2395 102. Sharon, D., Glusman, G., Pilpel, Y., Khen, M.,
2396 Gruetzner, F., Haaf, T. & Lancet, D. (1999). Primate
2397 evolution of an olfactory receptor cluster: diversifi-
2398 cation by gene conversion and recent emergence of
2399 pseudogenes. *Genomics*, **61**(1), 24–36.
- 2400 103. Ota, T. & Nei, M. (1995). Evolution of immuno-
2401 globulin VH pseudogenes in chickens. *Mol. Biol.*
2402 *Evol.* **12**, 94–102.
- 2403
- 2404
- 2405
- 2406
- 2407
- 2408
- 2409
- 2410
- 2411
- 2412
- 2413
- 2414
- 2415
- 2416
- 2417
- 2418
- 2419
- 2420
- 2421
- 2422
- 2423
- 2424
- 2425
- 2426
- 2427
- 2428
- 2429
- 2430
- 2431
- 2432
- 2433
- 2434
- 2435
- 2436
- 2437
- 2438
- 2439
- 2440
- 2441
- 2442
- 2443
- 2444
- 2445
- 2446
- 2447
- 2448
- 2449
- 2450
- 2451
- 2452
- 2453
- 2454
- 2455
- 2456
- 2457
104. Aravind, L., Watanabe, H., Lipman, D. J. & Koonin,
E. V. (2000). Lineage-specific loss and divergence of
functionally linked genes in eukaryotes. *Proc. Natl*
Acad. Sci. USA, **97**, 11319–11324.
105. Lykke-Andersen, J. (2001). mRNA quality control:
marking the message for life or death. *Curr. Biol.*
11(3), R88–R91.
- 2458
- 2459
- 2460
- 2461
- 2462
- 2463
- 2464
- 2465
- 2466
- 2467
- 2468
- 2469
- 2470
- 2471
- 2472
- 2473
- 2474
- 2475
- 2476
- 2477
- 2478
- 2479
- 2480
- 2481
- 2482
- 2483
- 2484
- 2485
- 2486
- 2487
- 2488
- 2489
- 2490
- 2491
- 2492
- 2493
- 2494
- 2495
- 2496
- 2497
- 2498
- 2499
- 2500
- 2501
- 2502
- 2503
- 2504
- 2505
- 2506
- 2507
- 2508
- 2509
- 2510
- 2511
- 2512
- 2513
- 2514
- 2515
- 2516
- 2517
- 2518
- 2519
- 2520

Edited by F.E. Cohen

(Received 14 September 2001; received in revised form 1 February 2002; accepted 2 February 2002)