

1 The complex relationship of 2 gene duplication and 3 essentiality

4 **Takashi Makino, Karsten Hokamp and Aoife McLysaght**

5 Smurfit Institute of Genetics, University of Dublin, Trinity College, Dublin 2, Ireland

6 *Corresponding author:* McLysaght, A. (aoife.mclysaght@tcd.ie).

7 In yeast and worm, duplicate genes overlap in
8 function so that deleting one of a pair from the
9 genome is less likely to be lethal than deleting a
10 singleton gene. By contrast, previous analyses
11 showed that mouse duplicate genes were as
12 essential as singletons. We show that the
13 relationship between gene duplication and
14 essentiality is complex in multicellular organisms,
15 with developmental genes and genes that were
16 duplicated by whole genome duplication being
17 more essential than other duplicated genes.

18 The 'essentiality' of duplicated genes

19 A gene is considered 'essential' if its removal
20 results in a lethal or sterile phenotype. Gene
21 duplication is frequent in eukaryotic genomes and
22 is the primary source of new genes [1-3].
23 Duplicate genes can have a backup role and can
24 functionally compensate for the loss of their
25 duplicated copies [4]. This concept was verified by
26 genome-wide gene knockout or knockdown
27 experiments in yeast and worm demonstrating
28 that the essentiality of duplicate genes is
29 significantly lower than that of singletons [4,5].
30 In addition, double knockout experiments in
31 yeast of paralogs derived from whole genome
32 duplication (WGD) strongly support functional
33 compensation by duplicated genes [6,7]. By
34 contrast, recent studies in mouse reported no
35 significant difference in essentiality between
36 duplicated genes and singletons [8,9]. This
37 surprising result indicated that duplicate genes
38 in mammals do not carry out a backup role and
39 indicated that the factors governing the evolution
40 and retention of duplicate genes differ between
41 mammals and less complex eukaryotes.

42 Mouse gene knockout dataset is enriched for 43 developmental genes

44 The data leading to the conclusions on essential
45 genes in yeast and worm were based on whole-
46 genome studies; however, the mouse studies [8,9]
47 relied on data from <4000 genes available from
48 Mouse Genome Informatics (MGI;
49 <http://www.informatics.jax.org/>) collected from
50 many individual studies. The patchiness of the

51 dataset makes it susceptible to potential data
52 biases because individual researchers might
53 preferentially report a gene with a discernable
54 phenotype in the knockout experiment.
55 Therefore, reports of gene knockouts with no
56 phenotypic change are likely to be dramatically
57 under-represented even in cases in which the
58 requisite experiment has actually been carried
59 out. By contrast, the stronger the knockout
60 phenotype, the more likely it is that the
61 observations are reported.

62 Liao and Zhang [9] investigated potential data
63 bias by comparison of their estimate of the
64 proportion of embryonic lethal genes from the
65 knockout dataset (14.0%) with an estimate from a
66 random mutagenesis study (13.7%) [10]. The
67 consistency of these two estimates led them to
68 conclude there was no significant data bias.
69 However, we found that 1523 out of 5078
70 knockout genes (30.0%) cause prenatal-perinatal
71 lethality in the most recent knockout dataset (see
72 methods in the supplementary material online),
73 strongly indicating that the knockout dataset is
74 not a representative sample. We considered the
75 possibility that there might be a functional bias
76 in the genes selected for knockout experiments,
77 and in particular genes involved in development
78 are likely to have a prenatal-perinatal lethal
79 knockout phenotype.

80 We tested the hypothesis of a functional bias in
81 knockout gene datasets for mouse and fly (see
82 methods in supplementary material online). Out
83 of 5078 knockout genes in mouse, 4609 genes
84 were annotated with at least one Gene Ontology
85 (GO) ID. We found that 18 GO terms are over-
86 represented in the knockout dataset with respect
87 to their frequency in the entire genome
88 (Table S1). Notably, GO terms related to early
89 development, such as GO:0007525 (multicellular
90 organismal development) and GO:0030154 (cell
91 differentiation), were highly over-represented in
92 reported knockout genes in mouse (genes with
93 either of these GO terms are hereafter referred to
94 as 'developmental genes'). Even though only 11%
95 of genes in the genome are annotated as

1 developmental (2682/23727), they constitute 37%
2 of the knockout dataset (1863/5078). We also
3 found a similar bias in fly (Table S2 and methods
4 in the supplementary material online). Thus,
5 there is a large bias in the reported knockout set
6 towards genes that function in development.

7 **Are developmental genes essential in mouse and** 8 **fly?**

9 If there is a large difference in essentiality
10 between developmental genes and others, then
11 this knockout dataset might give a misleading
12 impression of the genome-wide trend. To
13 investigate whether developmental genes are
14 more essential than other genes, we compared the
15 essentiality of developmental genes with non-
16 developmental genes. Using the same approach
17 as Liang and Li [8], and Liao and Zhang [9], we
18 defined an essential gene in mouse as one with
19 the knockout phenotype of sterility or lethality
20 before maturity [8,9]. The proportion of essential
21 genes (P_E) of developmental genes was
22 significantly higher than that of non-
23 developmental genes (mouse, $P < 2.2 \times 10^{-16}$; fly, P
24 $< 2.2 \times 10^{-16}$; c^2 test; Table 1). These results are
25 consistent with a recent report that showed
26 greater essentiality of genes highly expressed in
27 early development [11]. The greater likelihood of
28 fly and mouse developmental genes being
29 essential is understandable given the importance
30 of the developmental process.

31 **The essentiality of developmental and non-** 32 **developmental duplicates and singletons**

33 Given their overall high essentiality, we
34 wondered whether developmental genes were
35 subject to less functional compensation by
36 duplicate copies and whether the abundance of
37 developmental genes in the knockout dataset had
38 the potential to mask functional compensation in
39 other genes. Therefore, we subdivided the
40 developmental and non-developmental genes into
41 duplicates and singletons (see methods in the
42 supplementary material online). We found that
43 the essentiality of non-developmental duplicated
44 genes was significantly lower than that of non-
45 developmental singletons in mouse and fly
46 (mouse, $P = 0.00051$; fly, $P = 2.7 \times 10^{-8}$; c^2 test;
47 Table 1), following the trend observed in yeast
48 and worm [4,5]. Interestingly, the essentiality of
49 developmental duplicated genes was significantly
50 higher than that of developmental singletons in
51 mouse ($P = 0.0086$, χ^2 test; Table 1), and there
52 was no difference in essentiality between
53 developmental duplicated genes and singletons in
54 0.0051 fly $P = 0.98$, c^2 test; Table 1. Thus,
55 developmental genes are likely to be essential
56 irrespective of gene duplication.

57 **The influence of whole genome duplication on the** 58 **essentiality of duplicate genes**

59 Two rounds of WGD occurred early in the
60 vertebrate lineage [12–18] and duplicate
61 developmental genes created by these events
62 were preferentially retained in vertebrate
63 genomes [19–21]. Interestingly, developmental
64 genes were also preferentially retained after
65 WGD in plants [22], thus indicating particular
66 evolutionary dynamics after WGD in
67 multicellular organisms. Recent analysis of yeast
68 WGD duplicated genes indicated that they are
69 less essential than small-scale duplication (SSD)
70 duplicated genes [23,24]. We investigated the
71 essentiality of WGD and SSD duplicated genes in
72 mouse. We identified 1669 WGD duplicated genes
73 [17] and 2039 SSD duplicated genes with GO ID
74 and knockout data (see methods in the
75 supplementary material online). We confirm that
76 duplicate developmental genes are preferentially
77 generated by WGD rather than SSD, even when
78 we consider only genes from the knockout dataset
79 ($P = 3.0 \times 10^{-10}$, c^2 test; Figure 1a). Furthermore,
80 the P_E of WGD duplicated genes (45.4%) was
81 significantly greater than SSD duplicated genes
82 (38.1%; $P = 3.1 \times 10^{-6}$, c^2 test; Figure 1a). This
83 result is true even when we control for age
84 differences between WGD and SSD duplicates
85 (see methods in the supplementary material
86 online). We found there was no difference in
87 essentiality between WGD duplicated genes
88 (45.4%) and singletons (42.2%; $P = 0.10$, c^2 test)
89 in the entire mouse gene knockout set, but that
90 the P_E of SSD duplicated genes (38.1%) was
91 significantly lower than that of singletons (42.2%;
92 $P = 0.0027$, c^2 test). This is contrary to the
93 findings in yeast [23,24].

94 **Correlation between sequence divergence from** 95 **closest paralog and essentiality of duplicated** 96 **genes**

97 Previous studies reported that there is a positive
98 correlation between sequence divergence from the
99 closest paralog (most similar protein sequence)
100 and essentiality of duplicated genes in yeast and
101 worm [4,5]; that is, the greater the sequence
102 similarity between duplicated genes, the greater
103 the propensity for mutual functional
104 compensation. By contrast, in mouse there is a
105 negative correlation between sequence divergence
106 from the closest paralog and essentiality of
107 duplicated genes [9], or no correlation [25].

108 We examined the relationship between
109 sequence divergence from the closest paralog and
110 essentiality of duplicated genes used in above
111 analyses (see methods in the supplementary
112 material online). We found that the lower the
113 divergence from the closest paralog (i.e. lower

Comment [J1]: Author: correct?

1 K_A), the lower the P_E for SSD duplicated genes in
2 mouse (Pearson's product-moment correlation
3 coefficient $R = 0.94$, $P = 0.017$), but this trend was
4 not observed in other groups of duplicated genes
5 (Figure 1b). However, when we focused on genes
6 with $K_A > 0.2$, because highly constrained genes
7 might have unusual properties (e.g. ribosomal
8 proteins) [4,9], we observed a positive correlation
9 for non-developmental duplicated genes in mouse
10 ($R = 0.90$, $P = 0.039$; Figure 1b) and fly ($R = 0.92$,
11 $P = 0.027$; Figure 1c).

12 Concluding remarks

13 The relationship between gene essentiality and
14 gene duplication is complex in mouse owing to the
15 constraints on the developmental process and the
16 history of genome duplications in the vertebrate
17 lineage. Many transcription factors, members of
18 protein complexes and developmental genes are
19 sensitive to their relative dosage to other genes
20 (i.e. they are dosage-balanced) [26–28]. Dosage-
21 balanced genes are not robust to gene loss and
22 gene duplication [27,28]. WGD duplicates all
23 genes simultaneously and therefore does not
24 perturb relative dosages. Whereas SSD of dosage-
25 balanced genes is likely to be deleterious, WGD
26 should be neutral. Furthermore, subsequent loss
27 of dosage-balanced genes after WGD will be
28 deleterious unless contemporaneous loss is
29 somehow achieved. Therefore, the only
30 opportunity to duplicate dosage-balanced genes
31 might be when WGD occurs [27,28].

32 Our finding that developmental genes and
33 genes duplicated by WGD are more essential than
34 expected could be explained by dosage-balance
35 constraints. Subunits of a protein complex are
36 particularly likely to be dosage-balanced [27]. We
37 found significant enrichment for protein complex
38 membership for both WGD duplicated genes
39 (21.8%; 388/1781) and developmental genes
40 (20.0%; 372/1863) compared with the total
41 dataset (17.9%; 906/5068; see methods in the
42 supplementary material online). In addition, the
43 WGD-duplicated genes and developmental genes
44 in our dataset are significantly enriched for the
45 functional category GO:0030528 'transcription
46 regulator activity' (data not shown), which are
47 likely to be dosage-balanced [27,28].

48 In yeast, genes duplicated by WGD are less
49 essential than those duplicated by SSD [23,24].
50 The contrast with observations in mouse can be
51 explained by the comparatively simple
52 development process of this unicellular organism.
53 Similarly, worm, with only ~1000 cells, has less
54 complex development than fly or mammals [29]
55 and has not experienced WGD.

56 We suggest that the constraints inherent in
57 development of complex organisms (especially

58 dosage constraints) combined with the unique
59 evolutionary opportunities granted by the
60 simultaneous duplication by WGD of all
61 components of a pathway or complex explains the
62 high essentiality of these genes [30,31]. Because
63 WGD-duplicated genes and developmental genes
64 together constitute 26% of the mouse genome, but
65 57% of the knockout dataset, we expect that when
66 the data become available the genome-wide trend
67 in mouse will show that with these notable
68 exceptions, singletons are more essential than
69 duplicates, as is predicted by functional
70 compensation models.

71 Acknowledgements

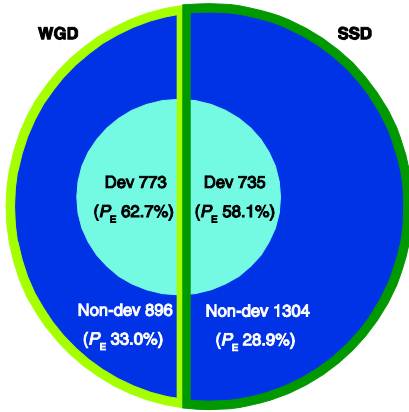
72 We would like to thank Yoichiro Nakatani for supplying lists of the WGD
73 duplicated genes and all the members of the McLysaght laboratory for
74 helpful discussions. This work is supported by Science Foundation
75 Ireland.

76 References

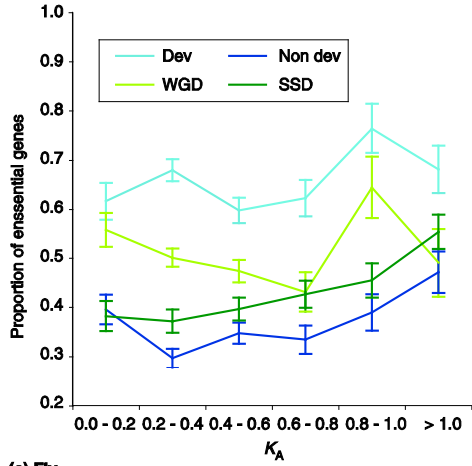
- 77 1 Ohno, S. (1970) Evolution by gene duplication.
78 Springer-Verlag
- 79 2 Long, M. *et al.* (2003) The origin of new genes:
80 glimpses from the young and old. *Nat. Rev. Genet.* 4,
81 865–875
- 82 3 Lynch, M. and Conery, J.S. (2003) The origins of
83 genome complexity. *Science* 302, 1401–1404
- 84 4 Gu, Z. *et al.* (2003) Role of duplicate genes in
85 genetic robustness against null mutations. *Nature* 421,
86 63–66
- 87 5 Conant, G.C. and Wagner, A. (2004) Duplicate
88 genes and robustness to transient gene knock-downs in
89 *Caenorhabditis elegans*. *Proc Biol Sci* 271, 89–96
- 90 6 DeLuna, A. *et al.* (2008) Exposing the fitness
91 contribution of duplicated genes. *Nat. Genet.* 40, 676–
92 681
- 93 7 Musso, G. *et al.* (2008) The extensive and
94 condition-dependent nature of epistasis among whole-
95 genome duplicates in yeast. *Genome Res.* 18, 1092–
96 1099
- 97 8 Liang, H. and Li, W.H. (2007) Gene essentiality,
98 gene duplicability and protein connectivity in human
99 and mouse. *Trends Genet.* 23, 375–378
- 100 9 Liao, B.Y. and Zhang, J. (2007) Mouse duplicate
101 genes are as essential as singletons. *Trends Genet.* 23,
102 378–381
- 103 10 Wilson, L. *et al.* (2005) Random mutagenesis of
104 proximal mouse chromosome 5 uncovers predominantly
105 embryonic lethal mutations. *Genome Res.* 15, 1095–
106 1105
- 107 11 Roux, J. and Robinson-Rechavi, M. (2008)
108 Developmental constraints on vertebrate genome
109 evolution. *PLoS Genet.* 4, e1000311
- 110 12 McLysaght, A. *et al.* (2002) Extensive genomic
111 duplication during early chordate evolution. *Nat.*
112 *Genet.* 31, 200–204
- 113 13 Hokamp, K. *et al.* (2003) The 2R hypothesis and
114 the human genome sequence. *J. Struct. Funct.*
115 *Genomics* 3, 95–110
- 116 14 Panopoulou, G. *et al.* (2003) New evidence for
117 genome-wide duplications at the origin of vertebrates
118 using an amphioxus gene set and completed animal
119 genomes. *Genome Res.* 13, 1056–1066
- 120 15 Vandepoel, K. *et al.* (2004) Major events in the
121 genome evolution of vertebrates: paranome age and
122 size differ considerably between ray-finned fishes and

1	land vertebrates. <i>Proc. Natl. Acad. Sci. U. S. A.</i> 101,	26	23	Guan, Y. <i>et al.</i> (2007) Functional analysis of gene
2	1638–1643	27	24	duplications in <i>Saccharomyces cerevisiae</i> . <i>Genetics</i> 175,
3	16	28	25	933–943
4	Dehal, P. and Boore, J.L. (2005) Two rounds of	29	26	Hakes, L. <i>et al.</i> (2007) All duplicates are not
5	whole genome duplication in the ancestral vertebrate.	30	27	equal: the difference between small-scale and genome
6	<i>PLoS Biol.</i> 3, e314	31	28	duplicate. <i>Genome Biol.</i> 8, R209
7	17	32	29	Su, Z. and Gu, X. (2008) Predicting the proportion
8	Nakatani, Y. <i>et al.</i> (2007) Reconstruction of the	33	30	of essential genes in mouse duplicates based on biased
9	vertebrate ancestral genome reveals dynamic genome	34	31	mouse knockout genes. <i>J Mol Evol</i> (in press)
10	reorganization in early vertebrates. <i>Genome Res.</i> 17,	35	32	26
11	1254–1265	36	33	Veitia, R.A. (2002) Exploring the etiology of
12	18	37	34	haploinsufficiency. <i>Bioessays</i> 24, 175–184
13	Putnam, N.H. <i>et al.</i> (2008) The amphioxus	38	35	27
14	genome and the evolution of the chordate karyotype.	39	36	Papp, B. <i>et al.</i> (2003) Dosage sensitivity and the
15	<i>Nature</i> 453, 1064–1071	40	37	evolution of gene families in yeast. <i>Nature</i> 424, 194–
16	19	41	38	197
17	Blomme, T. <i>et al.</i> (2006) The gain and loss of	42	39	28
18	genes during 600 million years of vertebrate evolution.	43	40	Wapinski, I. <i>et al.</i> (2007) Natural history and
19	<i>Genome Biol.</i> 7, R43	44	41	evolutionary principles of gene duplication in fungi.
20	20	45	42	<i>Nature</i> 449, 54–61
21	Brunet, F.G. <i>et al.</i> (2006) Gene loss and	46	43	29
22	evolutionary rates following whole-genome duplication	47	44	Nelson, C.E. <i>et al.</i> (2004) The regulatory content
23	in teleost fishes. <i>Mol. Biol. Evol.</i> 23, 1808–1816	48	45	of intergenic DNA shapes genome architecture.
24	21	49	46	<i>Genome Biol.</i> 5, R25
25	Huften, A.L. <i>et al.</i> (2008) Early vertebrate whole	50	47	30
	genome duplications were predated by a period of	51	48	Freeling, M. and Thomas, B.C. (2006) Gene-
	intense genome rearrangement. <i>Genome Res.</i> 18, 1582–		49	balanced duplications, like tetraploidy, provide
	1591		50	predictable drive to increase morphological complexity.
	22		51	<i>Genome Res.</i> 16, 805–814
	22			31
	Maere, S. <i>et al.</i> (2005) Modeling gene and genome			Otto, S.P. (2007) The evolutionary consequences
	duplications in eukaryotes. <i>Proc. Natl. Acad. Sci. U. S.</i>			of polyploidy. <i>Cell</i> 131, 452–462
	<i>A.</i> 102, 5454–5459			

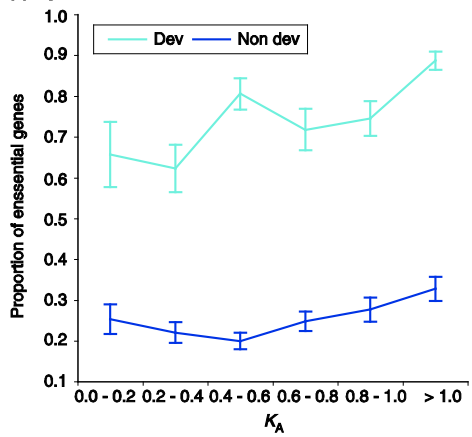
(a) Mouse duplicated genes



(b) Mouse



(c) Fly



1 **Figure 1.** The relationship of proportion of essential genes (P_E) and function, divergence, and origin of duplicated genes. **(a)** Venn diagram of P_E of
2 developmental, non-developmental, WGD and SSD duplicated genes in the mouse gene knockout dataset. **(b, c)** Relationship of sequence divergence and
3 proportion of essential genes for mouse (b) and fly (c) duplicate genes. The x-axis indicates the non-synonymous substitution rate (K_a) between a duplicated
4 gene and its closest paralog. The y-axis indicates the P_E in each K_a category. Error bars indicate standard error. Color code: Light blue, developmental
5 genes; dark blue, non-developmental genes; light green, WGD genes; and dark green, SSD duplicated genes in the mouse gene knockout dataset.

6 **Table 1. Proportion of essential genes for mouse and fly genes**

Species		Developmental genes	Non-developmental genes	Total
Mouse	Singletons	52.7% (187/355)	38.5% (210/546)	44.1% (397/901)
	Duplicated genes	60.5% (912/1508)	30.6% (673/2200)	42.7% (1585/3708)
	Total	59.0% (1099/1863)	32.2% (883/2746)	43.0% (1982/4609)
Fly	Singletons	79.1% (474/599)	34.3% (522/1520)	47.0% (996/2119)
	Duplicated genes	78.9% (607/769)	25.6% (487/1905)	41.1% (1094/2674)
	Total	79.0% (1081/1368)	29.5% (1009/3425)	43.6% (2090/4793)

7