

## RESEARCH

# Evolutionary model for the statistical divergence of paralogous and orthologous gene pairs generated by whole genome duplication and speciation

Yue Zhang, Chunfang Zheng and David Sankoff\*

\*Correspondence:

sankoff@uottawa.ca

Department of Mathematics and Statistics, University of Ottawa, 585 King Edward Avenue, Ottawa, Ontario, Canada, K1N 6N5

Full list of author information is available at the end of the article

## Abstract

We outline a principled approach to the analysis of duplicate gene similarity distributions, based on a model integrating sequence divergence and the process of fractionation of duplicate genes resulting from whole genome duplication (WGD). This model allows us predict duplicate gene similarity distributions for series of two or three WGD, for whole genome triplication followed by a WGD, and for triplication, followed by speciation, followed by WGD. We calculate the probabilities of all possible fates of a gene pair as its two members proliferate or are lost, predicting the number of surviving pairs from each event. We discuss how to calculate maximum likelihood estimators for the parameters of these models, illustrating with an analysis of the distribution of paralog similarities in the poplar genome.

**Keywords:** mixture of distributions; fractionation; probability model

## Introduction

Speciation creates a set of orthologous gene pairs involving all or almost all genes in the two daughter genomes, and these pairs all evolve according to a dynamic of decaying similarity as gene sequence and amino acid sequences inexorably diverge through random single nucleotide mutation. Whole genome duplication (WGD) creates a set of paralogous pairs involving all genes in the affected genome, and these pairs also diverge through the same processes of random mutation. In addition, paralogous pairs may disappear through the process of fractionation, whereby one of the two genes is excised, pseudogenized or otherwise removed as a recognizable coding gene.

A widespread practice in comparative genomics is to infer the nature and timing of evolutionary events through the examination of the distribution of similarities between orthologous or paralogous gene pairs. This is done by identifying local modes or peaks in the distribution, and inferring that duplications around these points were generated by speciation or WGD events. The identification of the peaks may be accomplished by visual inspection or, if the data seem noisy, by software available for the analysis of mixture of normal distributions, such as EMMIX [1].

There is, however, no rigorous methodology for interpreting the volume of the individual normal distributions inferred by such general methods. Indeed, many

possible outputs may not be conceivable as produced by genomic events. Moreover, there is a general tendency to overfit – to infer components of the mixture that are really just relict statistical fluctuation in the data.

In this paper, we outline a principled approach to the analysis of duplicate gene similarity distributions. It is based on the simplest, one-parameter model of sequence divergence, as well as an equally simple, one-parameter model of the fractionation process. We extend this to build models of a series of two or three WGD, of whole genome triplication followed by a WGD, characteristic of the core eudicots, and of a triplication, followed by a speciation and then a WGD in one of the two daughter species. In all these cases, we calculate the probabilities of all possible fates of a gene pair as its two members proliferate or are lost, to predict the number of surviving pairs from each event. To our knowledge, this is the first method to account for the volume of the component normals of a distribution of similarities, preliminary to an evolutionarily meaningful inference procedure.

We also outline how to infer the parameters of a model, using maximum likelihood methods. We illustrate with an analysis of the distribution of paralog similarities in the genome of the poplar, *Populus trichocarpa*.

We will conclude with a detailed discussion of the advantages and difficulties of our approach and detailed proposals for further research.

## The building blocks

We model gene pair divergence in terms of a probability  $p$  reflecting *similarity* – the proportion of nucleotide positions that are occupied by the same base in the two orthologs (or paralogs), although the same principles hold for *synonymous distance*  $K_s$  – the proportion of synonymous changes (not affecting translation to an amino acid) over all eligible positions, or *fourfold degenerate synonymous distance*  $4dTv$  – the transversion rate at fourfold degenerate third codon positions [2].

We represent by  $G$  the gene length, in terms of the number of nucleotides in the genes' coding region, setting aside for the moment that this varies greatly from gene to gene. We assume  $p$  follows the normal approximation to the sum of  $G$  binomial distributions, divided by  $G$ , and is related to the time  $t \in [0, \infty)$  elapsed since the event that gave rise to the pair:

$$\begin{aligned} \text{mean : } E[p] &= \frac{1}{4} + \frac{3}{4}e^{-\lambda t} \in [0, 1] \\ \text{variance : } E(p - E[p])^2 &= \frac{3}{16} \frac{(1 + 3e^{-\lambda t})(1 - e^{-\lambda t})}{G}, \end{aligned} \quad (1)$$

where  $\lambda > 0$  is a divergence rate parameter.

In practice,  $p$  for duplicate gene pairs is generally much greater than 0.25, so we base our analysis on those pairs with similarity greater than, say, 0.5.

Fractionation, the loss of one gene (and only one) from a pair, is represented by a parameter  $u \in [0, 1]$ , representing the probability, for a pair of genes, that neither gene is lost over a time interval of length  $t$ . The assumption that any gene pair has a constant probability (over time) of being fractionated entails

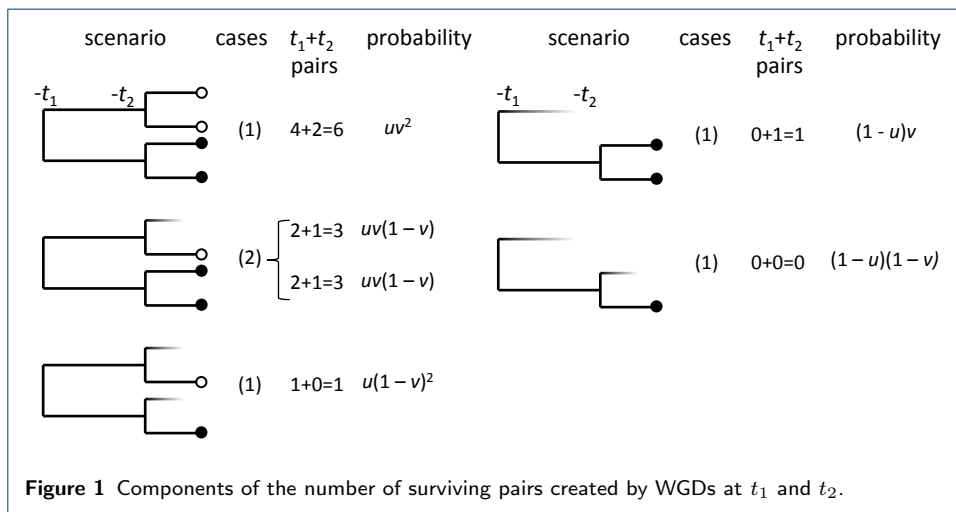
$$u = e^{-\rho t}, \quad (2)$$

where  $\rho$  is the fractionation parameter.

Thus, in the case of a single WGD, the mean of the distribution of duplicate gene pair similarities is an estimate of  $p$  (and also leads to an estimate of  $t$ ), and the number of pairs compared to the number of unpaired genes provides an estimate of  $u$  (and of  $\rho$ ).

### Two WGD

Consider a genome that has undergone two successive WGD. We denote by



“ $t_1$ -pairs” and “ $t_2$ -pairs” those duplicated gene pairs created at  $t_1$  and  $t_2$  respectively, with expected similarities  $p_1$  and  $p_2$ . For fixed  $\rho$ ,  $u$  and  $v$  are functions of  $t_1$  and  $t_2$  only, representing the probabilities  $e^{-\rho(t_1-t_2)}$  and  $e^{-\rho(t_2-0)} = e^{-\rho t_2}$ , respectively, for a pair of genes present at the start of the time interval, that neither gene is lost by the end of the interval. Note that in this and later models, we assume, for simplicity, that a fractionation regime from one WGD is supplanted by that set into operation by the next WGD. That is, fractionation involving older pairs is no longer operative.

In Figure 1, let

$$\begin{aligned}
 A &= \mathbf{E}(t_1 \text{ pairs}) \\
 &= 4uv^2 + 4uv(1-v) + u(1-v)^2 \\
 &= u(1+v)^2
 \end{aligned} \tag{3}$$

$$\begin{aligned}
 B &= \mathbf{E}(t_2 \text{ pairs}) \\
 &= 2uv^2 + 2uv(1-v) + (1-u)v \\
 &= v(1+u)
 \end{aligned} \tag{4}$$

$$\begin{aligned}
 C &= \mathbf{E}(\text{unpaired genes}) \\
 &= (1-u)(1-v)
 \end{aligned} \tag{5}$$

$$P(A) = \text{proportion of } t_1 \text{ pairs}$$

$$= \frac{A}{A + B + C} \quad (6)$$

$P(B)$  = proportion of  $t_2$  pairs

$$= \frac{B}{A + B + C} \quad (7)$$

$P(C)$  = proportion unpaired

$$= \frac{C}{A + B + C} \quad (8)$$

For a fixed gene length  $G$  and  $\lambda$ , let  $\mathbf{N}_p(s)$  be the density at point  $s$  of a normal distribution with mean  $p$  and variance  $\frac{p(1-p)}{G}$ . The probability that gene pair will be observed to be with similarity  $s \in [0, 1]$  is

$$Q(s) = P(A)\mathbf{N}_{p_1}(s) + P(B)\mathbf{N}_{p_2}(s). \quad (9)$$

and the probability of an unpaired gene is

$$Q^* = P(C). \quad (10)$$

The likelihood of a data set with gene pairs at  $s_1, \dots, s_l$  and  $k$  unpaired genes is

$$\mathcal{L} = \prod_{i=1}^l Q(s_i) Q^{*k}. \quad (11)$$

The log likelihood  $L = \log \mathcal{L}$  is

$$\begin{aligned} L &= \sum_{i=1}^l \log Q(s_i) + k \log Q^* \\ &= \sum_{i=1}^l [\log(P(A)\mathbf{N}_{p_1}(s_i) + P(B)\mathbf{N}_{p_2}(s_i))] + k \log Q^* \end{aligned} \quad (12)$$

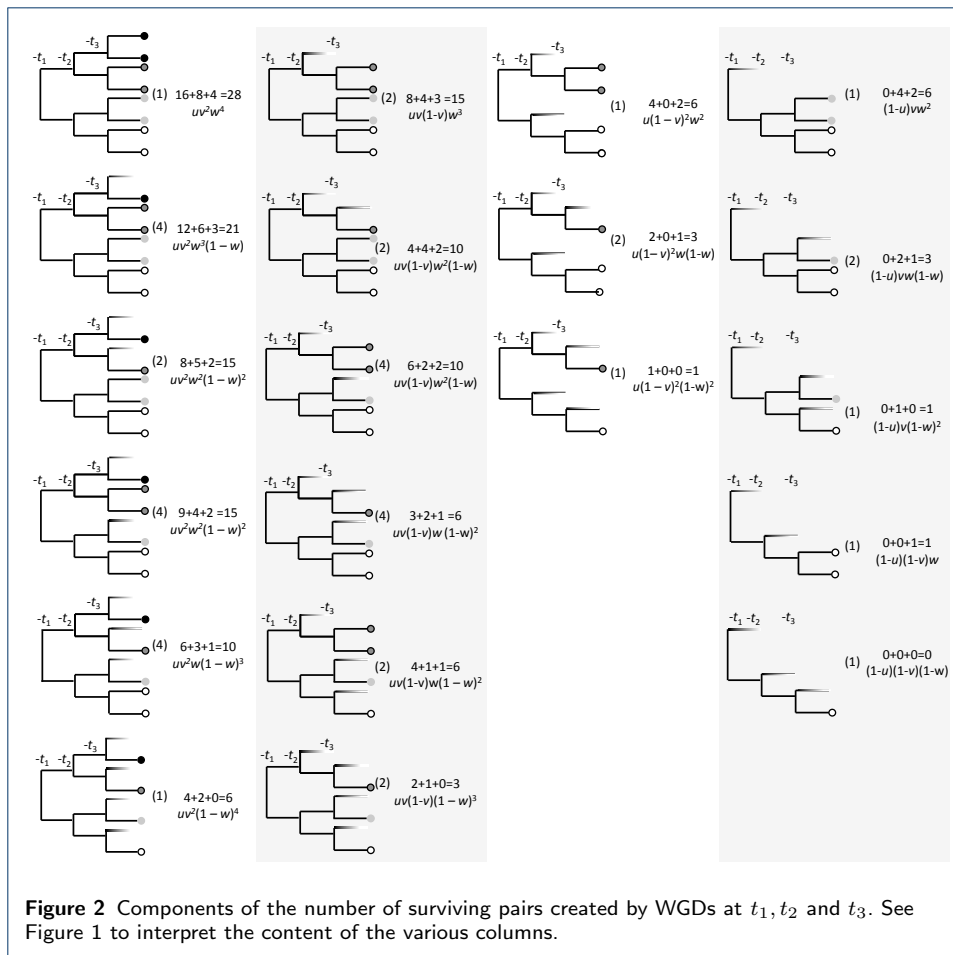
There is no closed form for the maximum likelihood of a mixture of normals, so in practice we use numerical means such as Newton-Raphson or an EM algorithm to derive the MLE.

### Three WGD

Consider now three successive WGD affecting a genome (for example the  $\tau, \sigma$  and  $\rho$  WGD that occurred in the common ancestor of the cereals [3]). The scenarios producing various numbers of gene pairs of various ages are depicted in Figure 2, where  $u, v$  and  $w$  are the retention probabilities for pairs produced at  $t_1, t_2$  and  $t_3$ .

$$\begin{aligned} E(t_1 \text{ pairs}) &= (1 - 3w^2 + 2w)uv^2 + (2 + 6w^2 + 4w)uv + (1 + w^2 + 2w)u \\ E(t_2 \text{ pairs}) &= ((1 + w^2 + 2w)u + 1 + w^2 + 2w)v \\ E(t_3 \text{ pairs}) &= -2wv^2w^2 + ((2w^2 - w)u + w)v + uv + w \\ E(\text{unpaired}) &= (1 - u)(1 - v)(1 - w) \end{aligned} \quad (13)$$

From this analysis, we can predict the number of pairs remaining from the each of the three events, and perform MLE calculations to determine the parameters.



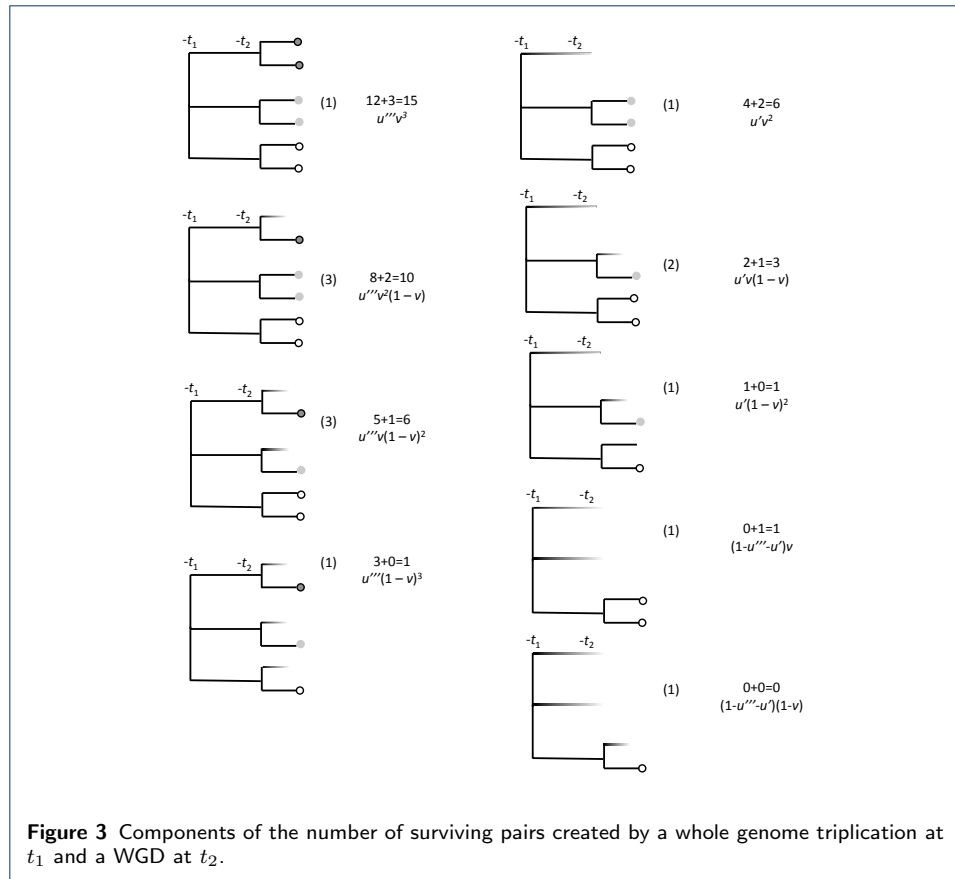
### Whole genome triplication followed by WGD

The core eudicots contain more species than all the other groups of flowering plants combined. A whole genome triplication occurred in the eudicot lineage just before the emergence of the core eudicots, and a large proportion of these have undergone further WGD. The model analyzed in Figure 4 is appropriate for this case. Here

$$\begin{aligned}
 E(t_1 \text{ pairs}) &= (u' + 3u''')v2 + (2u' + 6u''')v + b + 3u''' \\
 E(t_2 \text{ pairs}) &= -3u'''v3 + 3u'''v2 + (1 + 2u''' - u')v \\
 E(\text{unpaired}) &= (1 - u''' - u')(1 - v)
 \end{aligned}
 \tag{14}$$

### The effect of speciation

Up to now we have considered only WGD events, including triplications. In comparing two species, there are peaks at times corresponding to their shared WGD, followed by a single peak dating from their speciation event, but no further peaks. Figure 4 contains the analysis of a whole genome triplication, followed by a speciation event, and a further WGD in one of the daughter genomes. Here the distribution



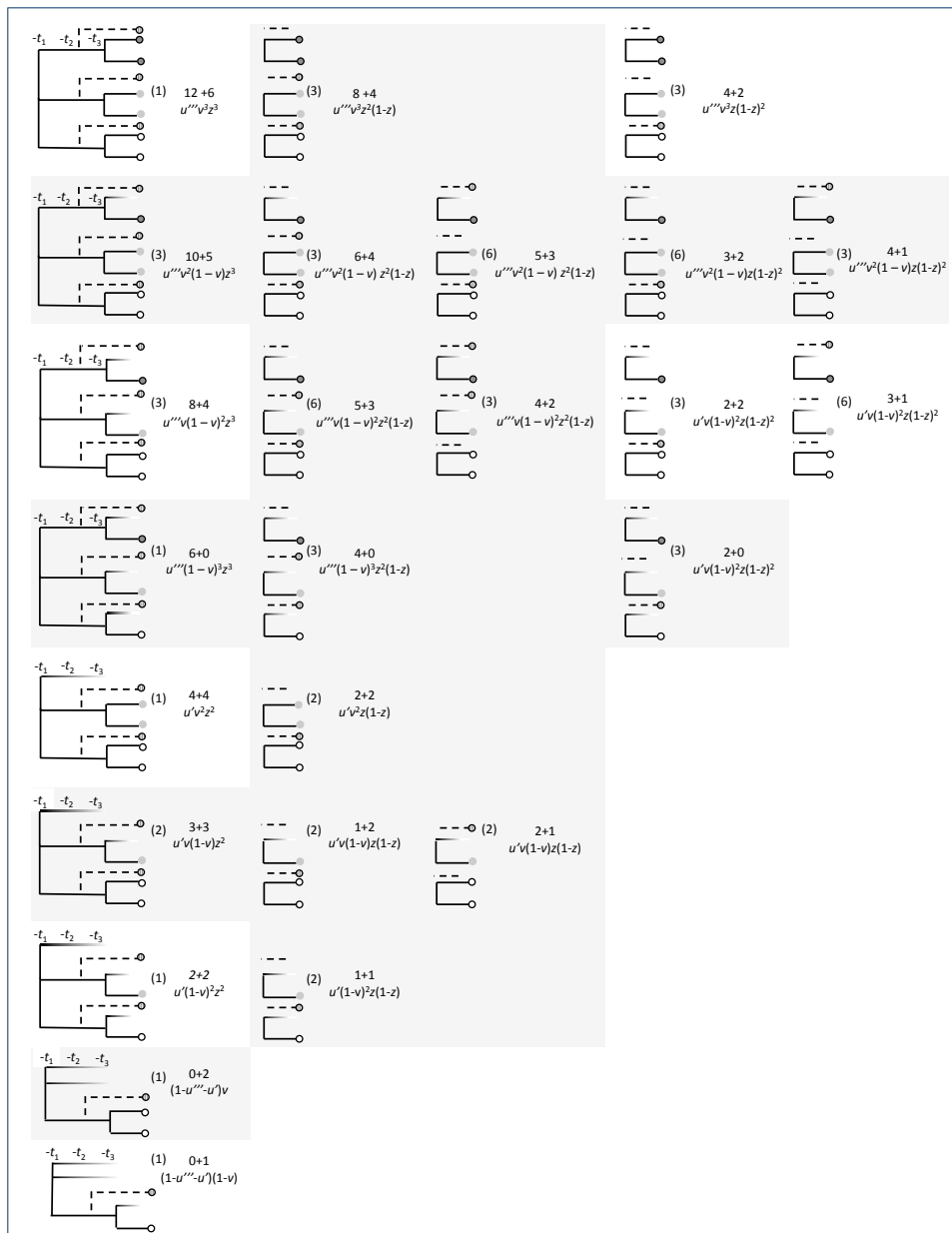
of orthologous gene pair similarities is predicted by:

$$\begin{aligned}
 \mathbf{E}[t_1 \text{ pairs}] &= (6u'''z^2(1-z) - 30u'''z(1-z)^2 + 12u'''z(1-z) \\
 &\quad + 30u'z(1-z)^2)v^3 \\
 &\quad + (30u'''z(1-z)^2 - 60u'z(1-z)^2)v^2 \\
 &\quad + (2u'z(1-z) + 2u'z^2 + 6u'''z^2(1-z) + 6u'''z^3 \\
 &\quad + 30u'z(1-z)^2)v \\
 &\quad + 2u'z^2 + 12u'''z^2(1-z) + 6u'''z^3 + 2u'z(1-z) \quad (15)
 \end{aligned}$$

$$\begin{aligned}
 \mathbf{E}[t_2 \text{ pairs}] &= (3u'''z^3 + 12u'''z^2(1-z) - 9u'''z(1-z)^2 + 6u'z(1-z)^2)v^3 \\
 &\quad + (-9u'''z^3 + 15u'''z(1-z)^2 - 18u'''z^2(1-z) \\
 &\quad - 24u'z(1-z)^2)v^2 \\
 &\quad + (1 - u''' - u' + 2u'z(1-z) + 12u'z(1-z)^2 \\
 &\quad + 2u'z^2 + 12u'''z^3 + 24u'''z^2(1-z))v \\
 &\quad 1 - u''' - u' + 2u'z^2 + 2u'z(1-z) \quad (16)
 \end{aligned}$$

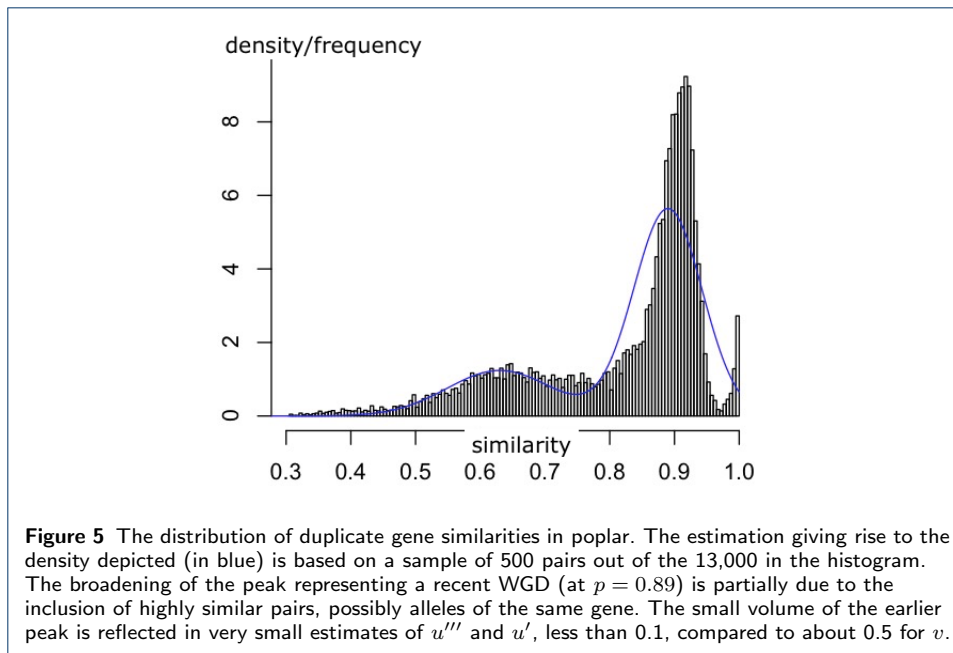
### The case of *Populus trichocarpa*

Though the work we have presented consists of combinatorial models, and the inference procedures are not implemented in a user-friendly package, we did analyze



**Figure 4** Components of the set of surviving orthologous pairs in two species diverging at time  $t_2$  after a (shared) WGD at  $t_1$ , where one species undergoes a second WGD at  $t_3$ . Components ordered vertically according to increased fractionation in the genome with the additional WGD, and ordered horizontally according to increased fractionation in the genome with no additional WGD (dashed lines). Number of cases of same component with different labelling in parentheses. “ $x + y$ ” indicates  $x$  pairs dating from the common WGD and  $y$  pairs created by the speciation event.

one data set using functions on the R platform according to the model in the previous section. We extracted data on the poplar genome [4] from the CoGe platform [5, 6], and calculated gene pairs, producing the distribution of similarities in Figure 5. In estimating the parameters using our preliminary code, we were confined to sampling 500 out of the 13,000 pairs. Running the program repeatedly, the results were quite reproducible.



## Conclusions and directions for further work

We have presented the first model of the simultaneous processes of duplicate gene divergence and fractionation of in the evolution of one or more species affected by WGD. This allows the prediction of both the location, shape and amplitude of the evolutionary signals, both speciation and WGD, contained in pairwise genome comparisons.

The parameter  $G$  affects the spread of the normally distributed contribution by an individual event to the overall distribution of gene pair similarities. It reflects the length of a gene, in terms of the number of nucleotides in the coding sequence, or of the number of synonymous sites, or of the number of four-fold degenerate sites. Length, however, is variable, from gene to gene, and from genome to genome, degrading the signal in the empirical distribution of similarities. One important direction for further work would be to incorporate the known properties of gene length distribution into the theory. Conceptually, this should pose little problem since since  $G$  is approximately log-normally distributed [7], so that  $\sigma^2$  in equation (1) can be adjusted directly.

Duplicate genes are also produced by mechanisms other than WGD. These can be largely avoided by requiring pairs to be corresponding syntenic contexts when extracting them from the data. This eliminates most of the problems due to tandem gene duplicates.

The assumption of a constant rates of gene divergence is another first order simplification made for analytical tractability, reducing as far as possible the number of parameters to be estimated. Rates in fact are variable among genes, between lineages, and over time [8]. Though small differences in rates may not be an overriding concern in the study of a sequence of events in a single genome, neglect of these differences may lead to serious errors in speciation-WGD-based phylogenetics [9]. Our approach allows for the introduction of rate variation, while controlling the number of parameters to be estimated.



While differences in divergence rates of gene pairs within and among genomes is relatively well understood, the same is not true of fractionation rates. There are a few quantitative studies of fractionation in the short term [10] and long term [11], but little coherent comparative literature at the whole genome level. There may be great variability of rates consequent to different events, due to the presence or absence of subgenome dominance in allopolyploids versus autopolyploids [12], the combination and timing of composite events giving rise to hexaploidy, e.g., in the Solanaceae ancestor [13], and other factors, making comparisons difficult. Our approach offers a new way of estimating fractionation rates, allowing different rates after different events or in different lineages.

#### Competing interests

The authors declare that they have no competing interests.

#### Author's contributions

The study was planned by DS, CZ and YZ. The mathematical research was carried out by YZ and DS, who also wrote the paper. All authors read and approved the paper.

#### Acknowledgements

Research supported in part by grants to DS from the Natural Sciences and Engineering Research Council of Canada. DS holds the Canada Research Chair in Mathematical Genomics.

#### References

1. McLachlan, G.J., Peel, D., Basford, K.E. and Adams, P. (1999) The Emmix software for the fitting of mixtures of normal and *t*-components. *Journal of Statistical Software* **4** 1–14. 1999.
2. Kumar, S. and Subramanian, S. (2002) Mutation rates in mammalian genomes. *Proceedings of the National Academy of Sciences of the United States of America* **99** 803–808.
3. McKain, M.R., Tang, H., McNeal, J.R., Ayyampalayam, S., Davis, J.I., DePamphilis, C.W., Givnish, T.J., Pires, J.C., Stevenson, D. W. and Leebens-Mack, J.H. (2016) A phylogenomic assessment of ancient polyploidy and genome evolution across the Poales. *Genome Biology and Evolution* **8** 1150–1164.
4. Tuskan, G.A. et al. (2006) The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science* **313** 1596–604.
5. Lyons, E. and Freeling, M. (2008) How to usefully compare homologous plant genes and chromosomes as DNA sequences. *The Plant Journal* **53** 661–673.
6. Lyons, E., Pedersen, B., Kane, J., Alam, M., Ming, R., Tang, H., Wang, X., Bowers, J., Paterson, A., Lisch, D. and Freeling, M. (2008) Finding and comparing syntenic regions among *Arabidopsis* and the outgroups papaya, poplar and grape: CoGe with rosids. *Plant Physiology* **148** 1772–1781.
7. Lipman, D.J., Souvorov, A., Koonin, E.V., Panchenko, A.R. and Tatusova, T.A. (2002) The relationship of protein conservation and sequence length. *BMC Evolutionary Biology* **2**:20.
8. Wolf, Y.I., Novichkov, P.S., Karev, G. P., Koonin, E. V. and Lipman, D. J., (2009) The universal distribution of evolutionary rates of genes and distinct characteristics of eukaryotic genes of different apparent ages. *Proceedings of the National Academy of Sciences* **106**, 7273–7280.
9. Sankoff, D., Zheng, C., Lyons, E., and Tang, H. (2016) The trees in the peaks. *Proceedings of the Third International Conference on Algorithms for Computational Biology (AICoB)*, Botón-Fernández, M., Martín-Vide, C., Santander-Jiménez, S., and Vega-Rodríguez, M.A. (eds.), Lecture Notes in Computer Science 9702, 3–16.
10. Buggs, R.J.A., Chamala, S., Wu, W., Tate, J.A., Schnable, P.S., Soltis, D.E., Soltis, P.S., and Barbazuk, W.B. (2012) Rapid, repeated, and clustered loss of duplicate genes in allopolyploid plant populations of independent origin. *Current biology* **22** 248–52.
11. Sankoff, D., Zheng, C., and Zhu, Q. (2010). The collapse of gene complement following whole genome duplication. *BMC genomics* **11**:1.
12. Garsmeur, O., Schnable, J.C., Almeida, A., Jourda, C., D'Hont, A. and Freeling, M., (2014) Two evolutionarily distinct classes of paleopolyploidy. *Molecular Biology and Evolution* **31**, 448–454.
13. Denoeud, F., Carretero-Paulet, L., Dereeper, A., Droc, G., Guyot, R., Pietrella, M., Zheng, C., Alberti, A., Anthony, F., Aprea, G., Aury, J.M., et al. (2014) The coffee genome provides insight into the convergent evolution of caffeine biosynthesis. *Science* **345** 1181–1184 (Supplementary materials).