1    **Asymmetric evolution of the transcription profiles and *cis*-regulatory sites contributes to**

2    **the retention of transcription factor duplicates**

3    Nicholas L. Panchy[1], Christina B. Azodi[2], Eamon F. Winship[3†], Ronan C. O'Malley[5], Shin-Han

4    Shiu[1,2,4*]

5

6    [1]Genetics Program, [2]Department of Plant Biology, [3]Department of Biochemistry and Molecular

7    Biology, and [4]Ecology, Evolutionary Biology, and Behavior Program, Michigan State University,

8    East Lansing, MI 48824, [5]DOE Joint Genome Institute, Walnut Creek, CA 94598

9

10

11    *Corresponding Authors:

12    Shin-Han Shiu

13    Michigan State University

14    Plant Biology Laboratories

15    612 Wilson Road, Room 166

16    East Lansing, MI 48824-1312

17    517-353-7196

18    shius@msu.edu

19

20    [†]Current address: MYcroarray 5692 Plymouth Rd Ann Arbor MI, 48105

21

1  **Abstract**

2    Transcription factors (TFs) play a key role in regulating plant development and response

3    to environmental stimuli. While most genes revert to single copy after a duplication event,

4    transcription factors are retained at a significantly higher rate. However, it is unclear why TF

5    duplicates have higher rates of retention relative to other genes. In this study, we compared

6    three types of features (expression, sequence, and conservation) of retained TFs following

7    whole genome duplication (WGD) events to genes with other functions, using *Arabidopsis*

8    *thaliana* as a model. We found that gene function groups with higher maximum expression but

9    lower mean expression tended to have higher duplicate retention rate post WGD, though TFs in

10   particular are retained more often than would be expected based on the features examined.

11   Conversely, expression of individual genes was not associated with duplication, but sequence

12   conservation was. Furthermore, we found that the evolution of TF expression patterns and cis-

13   regulatory cites favors the partitioning of ancestral states among the resulting duplicates. In

14   particular, we found that one duplicate retains the majority of ancestral expression and cis-

15   regulatory sites, while the "non-ancestral" duplicate was enriched for novel regulatory sites. To

16   investigate how this pattern of partitioning pattern evolved, we modeled the retention of

17   ancestral states in duplicate pairs using a system of differential equations. Our findings indicate

18   that duplicate pairs evolve to a partitioned state more often than away from it, which in

19   combination with accumulation of new regulatory sites in non-ancestral duplicates, suggest that

20   selection favors partitioning via neofunctionalization.

1    **Author Summary**

2    Gene expression is controlled by regulatory proteins known as transcription factors. These

3    factors control how an organism develops and responds to its environment. The evolution of

4    transcription factor functions also contributes to the emergence of new species and crop

5    domestication. In plants, new transcription factors mainly arise due to polyploidy, multiplication

6    of the genome. Although most duplicated copies are lost following a genome duplication event,

7    transcription factors are exceptional because they are often kept. Furthermore, we found that

8    transcription factor duplicates that tend to diverge in how they are expressed and regulated in

9    an unusual way where one copy mirrors the original, pre-duplication functional states of the

10   ancestral gene, while the other loses the ancestral status and instead accumulates novel

11   regulatory sites. Our results suggest these duplicate transcription factors may have been kept

12   because one copy preserve ancestral function while the other has evolved new ones.

13

1    **Introduction**

2         Plant genomes are replete with paralogous genes derived from a variety of duplication

3    events and mechanisms (Panchy et al., 2016). Among them, whole genome duplication (WGD)

4    events are responsible for the majority of extant duplicate genes (Panchy et al. 2016). Analysis

5    of sequenced plant genomes has revealed evidence for two ancient WGD events prior to the

6    divergence of angiosperms (Jiao et al. 2011) and, since then, more than a dozen WGD events

7    have occurred across a variety of angiosperm lineages (Lyons et al. 2008; Lee et al. 2013;

8    Myburg et al. 2014; Renny-Byfield et al. 2014; Soltis et al. 2014; Wang et al. 2014), including

9    three in the lineage leading to *Arabidopsis thaliana* (Bowers et al. 2003). This suggests that

10   WGD occurs more frequently in plants relative to other lineages, as the last known WGD event

11   in the *Saccharomyces cerevisiae* (Wolfe and Shields 1997; Kellis et al. 2004) and human

12   (Panopoulou et al. 2003; Dehal and Boore 2005) lineages occurred prior to the radiation of

13   angiosperms. Due to frequent WGD events and ensuing gene losses, the number of genes in

14   extant angiosperms vary widely but certain gene families, particularly transcription factor (TF)

15   families, have expanded dramatically in plant genomes (Lespinet et al. 2002; Shiu et al. 2005).

16        The duplication of TFs is of particular interest because duplicate TFs contribute

17   significantly to plant adaption (Lehti-Shiu et al. 2016), agricultural traits (Zhang et al. 2011), and

18   domestication (Liu et al. 2015). Not only does the expansion of several TF families coincide with

19   major events in the evolution of plants (i.e. migration to land and expansion of flowering plants,

20   Weirauch and Hughes 2011), but TF duplication has specifically influenced the evolution of

21   flowering time (Schranz et al. 2002), floral structures (Theissen and Melzer 2007) and the co-

22   option of floral regulators to fruit development (Litt and Irish 2003, McCarthy et al. 2015). WGD

23   accounts for ~90% of the expansion of TF families across plants lineages (Maere et al. 2005).

24   Compared to other eukaryotes, plant TF duplicates derived from WGD are retained at higher

25   rates than most plant genes with other functions (Seoighe and Gehring 2004; Shiu et al. 2005).

1    Furthermore, TFs are consistently enriched among WGD duplicates across divergent plant

2    species (Carretero-Paulet and Fares 2012), although the overall retention rates in these species

3    lineages vary greatly (Lynch and Conery 2000; Lynch and Conery 2003; Moghe and Shiu 2014).

4    Because WGD results in duplication of all genes in a genome, the observed differences in the

5    expansion of different gene families (e.g. Blanc and Wolfe 2004; Seoighe and Gehring 2004;

6    Hanada et al. 2008; Li et al. 2016) must result from differential rates of gene retention.

7    Previously, sequence features (e.g. gene length), biochemical activities (e.g. expression level),

8    evolutionary characteristics (e.g. substitution rates), and annotated function have been used to

9    assess the general properties of retained duplicates in plant genomes (Jiang et al. 2013; Moghe

10    et al. 2014). These properties of extant TFs, however, tell us little about the evolutionary

11    trajectories of TF duplicates. In order to determine what properties are associated with retained

12    TFs and how well those properties explain the differences in retention rates between TFs and

13    other function groups, it is critical to include information about how retained pairs have changed

14    since duplication.

15         Modeling approaches to infer ancestral expression levels based on extant gene

16    properties have been used to determine the rate of gene activation and repression in duplicate

17    genes in *Drosophila melanogaster* (Oakley et al. 2006) and analyze the evolution of stress

18    response in *A. thaliana* (Zou et al. 2009b). This approach allows for the explicit characterization

19    of how duplicate TFs may have deviated (or not) from their ancestral state over the course of

20    evolution, which in turn may provide information about the mechanism(s) contributing to TF

21    retention. Among theories proposed to explain how duplicated genes may be preserved,

22    subfunctionalization (Force et al. 1999) and gene balance (Birchler and Veitia 2007; Birchler

23    and Veitia 2010; Baker et al. 2013) have been hypothesized to explain the retention of TFs in

24    particular, as well as other genes with high interactivity (Seoighe and Gehring 2004; Maere et al.

25    2005; Shiu et al. 2005). In addition, neofunctionalization (Ohno,1970) and escape from adaptive

26    conflict (Des Marais and Rausher 2008) can lead to retention of both duplicate copies through

1    selection on new function or novel optimization of existing functions respectively. Nevertheless,

2    to assess the contribution of these different mechanisms, knowledge of the ancestral functional

3    states of a TF duplicate is essential.

4         In this study, we compared the retention rate of TF WGD duplicates against all other

5    genes in *A. thaliana* and against genes in each of 19 other "function groups" of similar size to

6    TFs. First, to understand the differences in retention rates amongst different function groups of

7    genes, we modeled retention rate as a function of the expression, structural features, and

8    conservation genes in each function group. Using the same feature set, we also classified

9    genes within each function group as either duplicated or not. In particular, we were interested to

10    see what features of duplicate genes distinguish them from non-duplicates. In addition, we

11    examined whether the correlations between a feature and retention status was consistent

12    across function groups or if some groups, like TFs, deviated from the norm. Second, we asked

13    to what degree the ancestral and extant functions of duplicate pairs were related to their

14    retention or loss. To determine how gene expression and *cis*-regulatory sites of TF duplicates

15    have likely evolved post WGD we inferred the ancestral expression and *cis*-regulatory states of

16    existing TF duplicates. We then asked whether the evolution of TF functions as defined by cis-

17    regulatory sites and expression patterns were correlated, which would be expected if regulatory

18    sites are controlling expression. Finally, to test if natural selection is contributing to the observed

19    distribution of ancestral states amongst extant duplicates, we modeled the evolution of TF WGD

20    duplicates as a system of differential equations which tracks the change in frequency of

21    duplicate pairs retaining the ancestral state in both, one, or neither.

1    **Results**

2    **Retention rates of duplicate genes in different function groups following WGD**

3    To assess the factors contributing to the differential retention of TF duplicates from WGD

4    events and duplicates from WGD events involved in other functions, we first quantified the

5    retention rates of *A. thaliana* WGD duplicates in 20 different function groups. These function

6    groups include TFs annotated by the Plant Transcription Factor Database (Jin et al. 2014) and

7    19 other groups defined based on Gene Ontology (GO, Ashburner et al. 2000) molecular

8    functions (see **Methods, Table S1**). Within each function group, genes were classified as

9    "WGD-duplicates" (both duplicate copies retained) or "WGD-singletons" (only one copy

10   retained) depending on whether there were paralogs in corresponding duplicate blocks (Bowers

11   et al. 2003). Because duplicate retention rates are expected to differ across different WGD

12   events, duplicate pairs derived from the α, β, and γ WGD events (Bowers et al. 2003) were

13   analyzed separately. To represent the duplicate retention rate of genes in a function group, we

14   calculated the log odds ratio of genes having a retained WGD-duplicate derived from a specific

15   WGD event for each function group relative to all *A. thaliana* genes (see **Methods**). Thus a

16   more positive retention rate indicates that a function group contains a proportionally higher

17   number of retained WGD-duplicates compared to the genome average, while a more negative

18   rate indicates a proportionally lower number of duplicates. Among the 20 function groups

19   examined, the retention rates were highly heterogeneous and only TFs and protein kinases had

20   retention rate significantly higher than the genome average for all three WGD events (**Fig 1**).

21   Compared to protein kinases, the retention rate of TFs are even higher for the older (β

22   and γ) duplication events, indicating that on average, the longevity of TF WGD-duplicates is

23   higher than that of protein kinases. Furthermore, the retention rates for the γ TF duplicates

24   (log2(odds) 1.72) is significantly higher than α (log2(odds) 0.95) and β (log2(odds) 1.07) TF

25   WGD-duplicates. This is in spite of the fact that the γ WGD occurred prior to the divergence of

7

1   monocots and dicots, making it at least twice as old as the α event (Bowers et al. 2003).

2   Importantly, the higher rate of retention for TF WGD-duplicates from the γ event cannot be

3   explained by the fact that these WGD-duplicates went through subsequent duplication events

4   (i.e. α and β). Rather, either more TFs were retained immediately following the γ WGD (i.e.

5   before subsequent duplication) than either the α or β event, or WGD-duplicate retention in

6   subsequent WGD events favored genes with retained WGD-duplicates following the γ WGD

7   (see **Supplementary File S1**). In summary, TFs were retained more frequently post WGD than

8   most other function groups. Additionally, the retention rate of γ WGD-duplicate TFs is

9   significantly higher than other events, potentially reflecting the contribution of γ WGD-duplicate

10  TFs to the early radiation of flowering plant species.

11  **Linear model of WGD-duplicate retention rates across function groups**

12          Amongst function groups, TFs stand out as being one of only two which are retained

13  more often than the genome average consistently across all WGD events (the other being

14  protein kinases, one of the largest gene families in plants, Lehti-Shiu and Shiu, 2012). For the

15  rest of the function groups, the retention rate varies above and below the genome average

16  across WGD events without any clear association to either more recent or ancient WGDs.

17  Additionally, although a few highly retained functions also have a higher gene count, there is no

18  clear relationship between the size of a function group and its retention rate for any WGD event

19  (**Fig 2A**). Therefore, the reason for retention rates differences remains unclear but must involve

20  factors beyond gene function and timing of duplication. To address why retention rates differ, we

21  examined sequence, expression, conservation, and interactivity related features (**Fig 2B, Table**

22  **S2**) of WGD-duplicate and WGD-singleton genes. We also asked how well the retention rate

23  differences between function groups can be explained by these different features.

24          To see how well the features we considered could explain the retention rate differences

25  among function groups, we constructed a linear model of retention rates for each WGD event. A

8

1    subset of the features examined here has previously been shown to be significantly associated

2    with retention of WGD-duplicates as a whole (Jiang et al. 2013), however we choose to

3    separate events in this case both because there was a large variance in retention across events

4    within our function groups (**Fig 1**). Furthermore, WGD events were separated because the

5    correlations between retention rates and feature values have different signs (see black arrows,

6    **Fig 2B**) and magnitudes (see white arrows, **Fig 2B**) depending on the WGD events. Hence, a

7    model which describes the relationship between the average features of function groups and

8    duplicate retention rates for a particular WGD event may not be generalizable to other events.

9    Beginning with the full set of 34 features, for each WGD event we determined the subset of

10   features (between 5 and 6 in each case) which maximized the F-statistic of the model (see

11   **Methods, Table 1**). Full model equations can be found in **Supplementary File S2**. Our models

12   explained 87%, 83%, and 65% of the variance in retention rates for the α, β and γ events

13   respectively. Applying the F-test to the maximum F-statistic for each model, we found that each

14   model performs significantly better at explaining the retention rates of function groups than the

15   null model (i.e. fitting retention rates to their mean) at p-value threshold of 0.05.

16

17   **Table 1.** Statistics for the best fitting model of duplicate retention rate for each WGD-event.

| WGD Event | # Features[1] | CoD [2] | F-statistic[3] | *p*-value[4] |
|---|---|---|---|---|
| α | 6 | 0.87 | 13.8 | 5.6E-05 |
| β | 5 | 0.83 | 13.2 | 7.1E-05 |
| γ | 5 | 0.65 | 5.1 | 7.2E-03 |

18   1: The number of explanatory variables (features) used in the best fitting model

19   2: Coefficient of Determination ($R^2$)

20   3: the F-statistic is a measure of the goodness of fit of the model to the observed retention rate.

21   4: p-value of fit based on the F-statistic. A significant p-value (p<0.05) indicates that the model

22   performs better than fitting the mean value to the data, after accounting for the number of

23   features in the model.

24

1    To determine the importance of individual features in explaining the retention rate

2    differences among function groups, we determined the change in explained variance caused by

3    independently removing each feature from the model (**Table 2**). As expected, features which

4    are used in more models have a greater impact on variance explained when removed. Note,

5    maximum expression (RNA-seq) with a positive sign and expression mean (AtGenExpress

6    microarray) with a negative sign were important to models for all three WGD events. This would

7    suggest that more specific expression (i.e. lower average across all conditions, but higher

8    maximum expression under a few specific conditions) increases the likelihood of duplicate

9    retention. Additionally, fewer domains, lower expression correlation with paralogs, and lower

10   nucleotide diversity, were all significantly features of the β and γ models, indicating that these

11   features had a greater impact on older duplication events. These findings suggest long term

12   retention of duplicates favors genes experiencing stronger purifying selection at the primary

13   sequence level (low nucleotide diversity) than at the level of gene expression patterns (lower

14   expression correlation reflecting higher degrees of expression divergence). The remaining

15   feature were found in only one of the models and had significant but much smaller impacts on

16   the variance explained than the features discussed above (**Table 2**).

17

18   **Table 2.** The importance of all features used in the linear models of duplicate retention in

19   function groups across each WGD event

| Feature | Sign[1] | $\alpha^2$ | $\beta^2$ | $\gamma^2$ |
|---|---|---|---|---|
| Expression Mean (AtGenExpress) | - | -0.29 | -0.09 | -0.49 |
| Expression Maximum (RNASeq) | + | -0.56 | -0.59 | -0.14 |
| Number of Domains | - | -0.06 | -0.36 | n/a |
| Nucleotide Diversity (Pi) | - | -0.06 | n/a | -0.32 |
| Expression Correlation (AtGenExpress) | - | n/a | -0.24 | -0.21 |
| Expression MAD/Median (AtGenExpress) | - | -0.09 | n/a | n/a |
| Protein Length (in Amino Acids) | + | -0.07 | n/a | n/a |
| Paralog Dn | + | n/a | -0.07 | n/a |

| Maximum Percent Identity | + | n/a | n/a | -0.2 |

1: The sign of the association between the feature and duplicate retention

2: Importance of features measured as the decrease in $R^2$ when the feature is removed from the

model, with more negative values indicating greater impact and therefore greater importance.

An n/a indicates the feature was not used in the model for that event.

To assess how well the models explained the retention rates of WGD duplicates in

different functional groups, we compared the actual and the predicted retention rates. In

general, the retention rates predicted by the models closely align with the actual values for each

function groups across each event (**Fig 2C**). However, TF retention rates were consistently

underestimated (**Fig S1**), particularly in the γ model the TF retention rate is predicted to be only

76% of the actual value (black arrows, **Fig 2C**). As such, while our models provide a general

explanation for the differences in duplicate retention between function groups, the behavior of

TFs departs from the norm. To conclude, we demonstrated that retention rates of genes in

different function groups are related to expression level/pattern and sequence divergence.

However, while these features are useful for predicting retention rates for some function groups,

they systematically underestimated TF retention rates.

**Machine learning classification of duplication status of individual genes**

One explanation for the underestimation of TF retention rates is that the relationships

between retention rate and feature values of TFs are significantly different than those of other

genes, such that a general model gives poor predictions. To explore this possibility, we used a

machine learning approach, Random Forest, to classify each individual gene as either having or

lacking a duplicate from a WGD event based on the gene's individual properties (see **Methods**).

We constructed separate classifiers for TFs, protein kinases (another function group with high

retention but better linear model fit, see white arrows **Fig 2C**), and all genes regardless of the

function groups. For this analysis, duplicates from all three WGD events were combined into

11

1    one classifier as small numbers of β and γ made the difficult to correctly classify on their own.

2    To evaluate the performance of our classifiers, we determined receiver operating characteristic

3    curves (ROCs) for each model (**Fig 3**) and calculated the Area Under Curve (AUC-ROC), a

4    metric that summarizes the ability of the classifier to recover true positive WGD-duplicate genes

5    at different false positive rates. An AUC-ROC of 0.5 indicates that the classifier is no better than

6    randomly labeling genes as having a retained duplicate or not, while an AUC-ROC of 1.0

7    indicates that the classifier can make predictions without error.

8         Among the classifiers, the one characterizing the full genome performed best (AUC-

9    ROC = 0.86), followed closely by protein kinases (AUC-ROC = 0.82), while the classifier for

10   TFs, while much better than random, did not perform as well (AUC-ROC = 0.74). To investigate

11   the source of the difference, we determined the importance of each feature to the classifier by

12   calculating the Mean Decrease in Accuracy (MDA) which is the average number of genes

13   misclassified across multiple runs as a result of removing a feature (**Table 3**). Given TFs are the

14   least well predicted, we suspected the informative features for predicting retention in TFs would

15   differ greatly from those for the genome at large and the protein kinases. Contrary to this

16   expectation, the ranking of importance for TF WGD-duplicate prediction was more similar to the

17   ranking of features for the whole genome prediction (Spearman's rank, $\rho = 0.86$) than the

18   ranking of features protein kinases to the whole genome prediction (Spearman's rank, $\rho = 0.51$).

19   This finding suggests that the feature value distributions of TF WGD-duplicate and WGD-

20   singletons are more similar to the genome at large. Therefore, the reason that TF duplicate

21   prediction model had lower performance was not simply because their feature values were

22   substantially different from other duplicate genes. Instead, the features examined simply have

23   lower importance in general for predicting TF retention (average MDA=11.3) than for other

24   genes (average MDA=47.9), suggesting there are additional features important for TF retention

25   that were not considered. For example, we might expect the number of DNA binding sites to be

26   predictive of duplication status as an indication of the breadth of function of the TF which is

12

1    related to the probability that a duplicate copy has been retained through subfunctionalization or

2    gene balance.

3

4    **Table 3.** The importance (rank) of all features used in the classification of individual duplicate

5    genes

| Feature | Genome[1] | Kinases[1] | TFs[1] |
|---|---|---|---|
| Maximum Percent Identity (Paralog) | 171.6 (1) | 57.9 (1) | 29.4 (2) |
| Sequence Conservation (Viridiplantae) | 109.3 (2) | 27.2 (2) | 31.8 (1) |
| Gene Family Size (OrthoMCL) | 81.6 (3) | 2.0 (19) | 27.7 (3) |
| Protein Length (in Amino Acids) | 52.5 (4) | 14.2 (3) | 10.5 (11) |
| Expression Breadth (AtGenExpression) | 46.2 (5) | 9.4 (10) | 11.2 (8) |
| Expression MAD/Median (AtGenExpress) | 41.0 (6) | 5.4 (11) | 11.8 (5) |
| Expression Mean (LightDev Data) | 40.8 (7) | 10.4 (7) | 10.7 (9) |
| Expression Breadth (RNASeq) | 40.0 (8) | 3.3 (15) | 11.5 (6) |
| Expression Mean (Control Data) | 39.2 (9) | 10.7 (6) | 10.7 (10) |
| Expression Median (AtGenExpress) | 37.5 (10) | 10.2 (8) | 10.2 (12) |
| Expression Mean (Stress Data) | 37.0 (11) | 11.5 (4) | 12.0 (4) |
| Expression Mean (AtGenExpress) | 36.9 (12) | 11.3 (5) | 11.3 (7) |
| Expression Median (RNASeq) | 34.8 (13) | 4.7 (12) | 4.9 (16) |
| Sequence Conservation (Metazoa) | 34.5 (14) | 3.6 (13) | 4.4 (18) |
| Nucleotide Diversity (Pi) | 32.4 (15) | 9.7 (9) | 6.2 (14) |
| Expression Mean (Diff Data) | 31.6 (16) | 1.7 (2) | 7.9 (13) |
| Sequence Conservation (Fungi) | 30.6 (17) | 2.4 (17) | 4.6 (17) |
| Number of Protein Domains | 28.4 (18) | 2.3 (18) | 5.6 (15) |
| Expression Mean (RNASeq) | 18.6 (19) | 3.0 (16) | 2.9 (19) |
| Expression Maximum (RNASeq) | 12.9 (20) | 3.3 (14) | 0.4 (20) |

6    1: The importance of the feature as defined by the mean decrease in accuracy of the

7    classification when the feature is removed. Features are ordered according to the rank of their

8    importance in the whole genome model and the rank of each value for each model is indicated

9    by (),

10

13

1    Furthermore, the most informative feature for classifying kinases and the whole genome,

2    the percent identity to the best matching paralog in *A. thaliana,* was less important when applied

3    to TFs (**Table 3**). Although the maximum percent identity of WGD-duplicates compared to

4    WGD-singletons is significantly higher in full genome ($p$ = 1e-320), protein kinases ($p$ = 1.1e-

5    36), and TFs ($p$ = 6.2e-12), the magnitude of the difference was greater for protein kinases

6    (11.2%) and the whole genome (11.3%) than TFs (4.4%). This is due to WGD duplicate TFs

7    having lower maximum percent identity (71.3%) than either kinases (75.2%, $p$=4.1e-24, t-test)

8    or all genes (72.5%, $p$=5.9e-83 , t-test), while WGD-singletons TF had higher identity (66.9%)

9    than kinases (64.0%, $p$=4.2E-35, t-test) and all genes (61.3%, $p$=1.9e-223, t-test). This

10   observation may related to non-duplicate TF genes having apparent paralogs more often than

11   non-duplicate genes do on average across the *A. thaliana genome* (**Fig S2**). The variance in the

12   importance of maximum percent identity accounts for most of the performance difference across

13   the classifiers as removing this feature yields similar results from all three (**Fig S3**). Similarly,

14   inflating the difference in the percent identity of TF WGD-duplicates and WGD-singletons from

15   4.4% to 11.2% (the difference for protein kinases) would raise the predicted retention of TF from

16   the γ WGD from 2.50 to 2.94, making up for more than half of the original error.

17   We would expect that other features used in our linear models (**Table 2**) would also be

18   useful for classifying genes within function groups. However, the average importance rank of

19   features found in more than one linear model was low (13.9 of 20), with the maximum

20   expression value in RNA-seq being the worst feature in both the whole genome and TF

21   classifiers. Of the four linear model features, mean expression in AtGenExpress had the highest

22   rank in the whole genome (12th), TF (7th), and kinase classifiers (5th). However, the difference

23   in mean expression between WGD-duplicates and WGD-singletons was not consistent: WGD

24   duplicates genes were more highly expressed across the whole genome (+0.32, $p$=4.0e-23),

25   and TFs (+0.37, $p$=1.0e-4), but in protein kinases WGD-singletons were more highly expressed,

26   though not at a significant level (-1.1, $p$=0.77). Hence, not only does relationship between gene

14

1    features and retention depend on the gene function, but the relationship within individual

2    function groups can be the opposite direction of the relationship across function groups. For

3    example, the high retention of the TF function group is in part due to relatively low average

4    expression in AtGenExpress, but within TFs, genes with higher average expression are more

5    often WGD duplicates. This suggests that selection for duplicate retention is dependent not only

6    on function and features, but their interaction as well, though the exact nature of these

7    interactions is beyond the scope of this study.

8    **Partitioning of ancestral expression states following TF duplication**

9         While the gene features (**Table 2**) were generally useful predictors of WGD-duplicates,

10    they were less useful for predicting TF duplicates specifically. To further explore what

11    characteristics retained TF WGD-duplicates possess, we examined how the functions of

12    retained TF WGD-duplicates have evolved following WGD events. To do this, we first used

13    expression patterns as a proxy of TF function(s) and inferred the likely expression states of the

14    ancestral TFs prior to WGD (see **Methods**). Ancestral expression values were inferred from

15    extant gene expression values that had been discretized into quartiles (expression state = 0, 1,

16    2, or 3) based on the distribution of expression levels for each array experiment. Additionally,

17    expression data were grouped into four subsets, including control conditions (Ctrl), light and

18    development sets (LightDev), abiotic and biotic stress treatments (Stress), and differential

19    expression between stress treatments and controls (Diff), and analyzed separately. This

20    grouping was used to distinguish between trends that were universal or specific to certain

21    datasets. We were able to infer 165,385 ancestral expression states across 474 TF WGD-

22    duplicate pairs (a detailed breakdown of inferred states can be found in **Table S3**).

23         First, to test how often the expression states of TFs are retained post-duplication, we

24    compared the expression states of individual, extant TF WGD-duplicate to its inferred ancestral

25    states (**Fig 4A**). Although all possible changes in expression state were observed between

15

1    ancestral and extant TFs in each expression data subset, the most common ancestral-extant

2    expression state combination was that the ancestral and extant TFs had the same expression

3    quartiles (diagonal red boxes, **Fig 4A**). This is true across all expression quartiles, though the

4    deviation from expectation was greatest for expression values in the lowest (0) and highest (3)

5    quartiles. This general pattern holds across all four data subsets (**Fig S4**), suggesting that most

6    TFs WGD-duplicates retain their original expression irrespective of what that expression is.

7    However, when considering a pair of duplicates, we found that when the ancestral state was

8    retained in one duplicate, it was lost more often in the other duplicate than expected by random

9    chance (**Fig 4B**). This "partitioned" state of TF WGD-duplicates pairs is most over-represented

10    in duplicates from the α event compared to duplicates from the older β and γ events (**Fig 4B**). In

11    these older WGD events, having neither duplicate inherit the ancestral expression state is more

12    common than the partitioned state where only one copy inherits the ancestral state. Using

13    ANOVA, we confirmed that there is indeed significant interaction between the expression state

14    of a TF WGD-duplicate pair and the timing of the WGD event ($p$<2e-16), which indicated that

15    partitioning occurred relatively quickly after the most recent WGD, but that these partitioned

16    patterns were not necessarily retained as the duplicates age.

17          Next we asked if TF duplicates expression tends to increase or decrease when they

18    deviate away from the ancestral state. Because we found a significant interaction between the

19    expressions state evolution of TF WGD-duplicate pairs and the subset of the expression data

20    used ($p$=2.5e-05), we asked this question for each subset of expression datasets individually.

21    For the LightDev (**Fig 4B**), Ctrl, and Stress expression subsets (**Fig S5**), partitioning of ancestral

22    expression states among duplicates favors small, negative changes from the ancestral states.

23    Based on an earlier study showing that *A. thaliana* duplicates tend to be expressed at a lower

24    level compared to the ancestral state (Zou et al. 2009b), we anticipated TFs would lose their

25    ancestral stress response. However, when we looked at the Diff subset, which looks at

26    differential expression as opposed to raw expression, we found that TFs were equally likely to

1   increase or decrease differential expression in response to stress compared to the ancestral

2   state.

3           To test the significance of this difference in non-ancestral expression amongst

4   expression subsets, we modeled the evolution of ancestral expression (O) to higher (+) and

5   lower (-) expression states following a WGD duplication event was modeled using ordinary

6   differential equations (see **Methods**). We compared a one-parameter model where the rates of

7   transition from (O) to (+) and (-) were set to be equal to a two-parameter model where the rates

8   to (+) and (-) were allowed to differ (**Fig S6**). The two parameter model was only significantly

9   better than the one parameter model for the LightDev (likelihood ratio test, $p$=2.2e-11), Ctrl

10  ($p$=2.7e-3), and Stress ($p$=2.9e-3) subsets. For these subsets, the rate of evolution from (O) to (-

11  ) was 1.9~3.1 times more frequent than that from (O) to (+). For the Diff subset, (O) to (-) was

12  only 1.2 times more frequent, which was not significant ($p$=0.43). In summary, these results

13  suggest that the evolution of TF duplicates favors decreasing expression levels relative to the

14  ancestral expression state (Control, LightDev, and Stress). However, when looking at differential

15  expression in response to stress, TF duplicates can evolve in either direction with approximately

16  equal likelihood. Thus, following duplication, TF duplicates may have increased or decreased

17  responses to stress, rather than losing the response altogether.

18  **Asymmetry in the partitioning of ancestral expression and regulatory sites**

19          Thus far we show that an ancestral expression state tends to be retained by only one

20  copy of a TF WGD-duplicate pair. One outstanding question is whether each copy would retain

21  different parts of the ancestral expression state, as would be expected if the TF duplicates were

22  retained due to subfunctionalization (Force et al. 1999). To address this, we considered all of

23  the partitioned expression states (i.e. all expression series showing partitioning) across a pair of

24  TF WGD-duplicates. If partitioning were random, we would expect that the number of ancestral

25  states retained by a single WGD-duplicate to follow a binomial distribution for the given number

17

1    of partitioned expression states (n) and a retention probability of 0.5 such that each copy is

2    equally likely to retain ancestral express states. Under this scenario, the expected asymmetry of

3    a duplicate pair (the difference in the fraction of ancestral states inherited between duplicates) is

4    0.18 given the observed distribution of partitioned states. However, the actual mean asymmetry

5    between TF WGD-duplicates was 0.67, which is unlikely to have been generated by random

6    partitioning ($p<$1e323) expected under the subfunctionalization model. In fact, in 35.1% of

7    cases, one WGD-duplicate retains all the ancestral expression and the distribution is highly

8    biased towards higher asymmetry (**Fig 5A**). As with the mean asymmetry, the skewed

9    distribution of asymmetry values is also significantly different from what was expected from

10   random partitioning (Kolmogorov–Smirnov test, $p<$2.2e-16). This biased partitioning was also

11   found within the Ctrl (mean=0.84), LightDev (mean=0.67), Stress (mean=0.69), and Diff

12   (mean=0.56) subsets. Given these results, for TF WGD-duplicate pairs we can generally define

13   one duplicate as being "ancestral" and the other as being "non-ancestral". Why then is the non-

14   ancestral copy being retained? One hypothesis is that the non-ancestral copy is retained

15   because it has acquired a novel function.

16          To test whether the non-ancestral copies tend to have novel functions, we first applied

17   our model of ancestral-state partitioning to *cis*-regulatory sites. We used *cis*-regulatory sites

18   here because the discretized expression levels used above allowed us to determine the

19   direction of changes away from the ancestral expression state, but not whether an expression

20   state was novel. The *cis*-regulatory sites used here are from putative binding sites of 345 *A.*

21   *thaliana* TFs (O'Malley et al. 2016). We applied the same methodology used to infer ancestral

22   gene expression to infer ancestral *cis*-regulatory sites of ancestral TFs (see **Methods**). Among

23   16,415 ancestral-extant site comparisons, the majority (58.9%) involved the loss of an ancestral

24   site in only one WGD-duplicate gene, which is significantly different compared what would be

25   expected if WGD-duplicate and ancestral genes were randomly associated (42.3%; t-test, $p<$1e-

26   323). In contrast, retention (15.8%, $p<$1e-323) and loss (10.2%, $p<$1e-323) of ancestral *cis*-

18

1    regulatory sites in both WGD-duplicates were significantly less frequent than randomly

2    expected. Similar to ancestral expression state evolution, the partitioning patterns of ancestral

3    *cis*-regulatory sites were highly asymmetric (**Fig 5B**), resulting in a distribution that is

4    significantly different from the distribution generated by random partitioning (Kolmogorov–

5    Smirnov test, $p$< 2.2e-16). Thus, much like what we observed for expression, these results

6    suggest that, with regard to *cis*-regulatory sites, TF WGD-duplicates can be classified into

7    ancestral and non-ancestral copies. Most importantly, amongst the 249 duplicate pairs with at

8    least one novel regulatory site, in 71.0% of cases the non-ancestral copy had more novel *cis*-

9    regulatory sites (**Fig 5C**). This is significantly different from what would be expected if ancestral

10   site retention and gain of novel regulatory sites associated independently (49.8%, $p < 3.8e-12$).

11   Furthermore, in 61.8% of pairs only the non-ancestral copy had novel sites, while the ancestral

12   copy contained all of the novel sites in only 14.0% of pairs. These patterns suggested that, the

13   gains of novel *cis*-regulatory sites likely contribute to the retention of the non-ancestral TF

14   duplicate copies.

15        Note that we can divide each pair of WGD-duplicates into an ancestral copy and a non-

16   ancestral copy based on either expression or *cis*-regulatory site information. Given the important

17   roles of *cis*-regulatory sequences in regulating gene expression, we expected that the ancestral

18   and non-ancestral designation defined according to expression data should be similar to that

19   defined based on *cis*-regulatory sites. Among the 179 TF WGD-duplicate pairs with data on

20   gene expression and regulatory sites, 59.8% follow this expected pattern. This percentage was

21   higher than expected by random association (24.6%, $p = 1.8e-20$). After examining the evolution

22   of expression states of TFs and *cis*-regulatory sites controlling TF expression, the next question

23   is how the regulatory targets differed between TF WGD-duplicates. This is currently not feasible

24   because the *A. thaliana* TF-target data set was too sparse for clear inference, however, we

25   have demonstrated that evolution of WGD-duplicates clearly favors the partitioning of ancestral

26   states into distinct ancestral and non-ancestral duplicates.

19

**Patterns of WGD-duplicate divergences and partitioning results from evolutionary bias**

We have demonstrated that partitioning of ancestral expression and regulation into ancestral and non-ancestral duplicates is favored following WGD-duplication of TFs. It remains an open question if this preference is due to bias towards partitioning or simply results from the progressive loss of ancestral function from duplicate genes. One possible explanation for the observed frequency of partitioning is that it results from the timing of WGD duplications in the *A. thaliana* lineage. Assuming that mutations occur and are fixed by drift, the time required for mutations that alter expression patterns to accumulate in two WGD-duplicated loci is expected to be longer than that for a single locus. Under the scenario where all mutations occur and are fixed at a constant rate, there would be a time when we would expect to find partitioning of ancestral state enriched simply because duplicate pairs are more likely to have lost the ancestral state at a single locus than both loci or neither loci. In contrast, if there is bias for the partitioning of ancestral states, we would expect the loss of the ancestral state at the first locus of a TF WGD-duplicate pair to occur must faster than at the second locus.

To determine whether the observed frequency of partitioning of ancestral states requires biased evolution or not, we modeled loss of ancestral state at pairs of loci in TF WGD-duplicates using a system of ordinary differential equations (see **Methods**). Using the synonymous substitution rate of TF WGD-duplicate pairs derived from the α, β, or γ events as a proxy for time, the rate of transition between WGD-duplicate pairs where both (state II), only one (state I) and neither (state O) duplicate had lost ancestral expression was modeled (**Fig 6A**). We compared a model where the rates of transitions between all states were equivalent (one-parameter model, **Fig 6B**) with a model where the transition rates between state II and I were allowed to vary from those between state I and O (two-parameter model) (**Fig 6C**). These models were applied to all expression subsets, the results for the LightDev dataset are shown in **Fig 6** and the remainder can be found in the supplement (**Fig S7**).

1    We found the two-rate model to be significantly better at explaining the observed

2    difference in WGD-duplicate states over time (Likelihood Ratio Test, p-value < 2e-14).

3    Regardless of the expression data set, the transition rates between state O (ancestral

4    expression in both duplicates) and I (ancestral expression in on duplicate) were 7-13 times

5    higher than the rates between state I and II (ancestral expression in neither duplicate). Given

6    that a pair of TF WGD duplicate would have the same expression patterns as their ancestral

7    gene initially (state O), our finding suggests that the number of partitioned WGD-duplicates

8    accumulated relatively rapidly post WGD, followed by slow accumulation of WGD-duplicates

9    pairs where ancestral expression had been lost entirely. Applying this same approach to model

10   regulatory site evolution revealed an even more extreme difference, as the rates governing the

11   transition between state O and I are two orders of magnitude higher than between state I and II

12   (**Fig 6E-G**). Additionally, since the best fit model for the regulatory data involved allowing all four

13   rate parameters to vary (p-values 4.8e-13 and 1.2e-11 vs. 1 and 2 parameter models

14   respectively), we can state that rates for transition to state I are higher than the accompanying

15   transition rates away (**Fig 6H**). This pattern does not hold when the four parameter model is fit

16   to ancestral expression state (**Fig 6D**), however, the model overall was not significantly better

17   than the two parameter model at a threshold of p-value < 0.05, implying that faster transition

18   between O and I compared to I and II is a better at explaining the evolution of expression states.

19   The non-equivalency in the rate at which the ancestral state is lost in the first and

20   second duplicates indicates that the frequency of partitioning is not trivial, but rather results from

21   a bias against losing the ancestral state from the second TF WGD-duplicate relative to the first.

22   Furthermore, in the case of ancestral regulatory sites, we found that evolving towards a

23   partitioned state was always favored over the corresponding path away from partitioning.  This

24   suggests that partitioned TF WGD-duplicates pairs do not accumulate simply because of

25   selection against having both duplicates ancestrally expressed, but that maintaining the original

26   number of ancestral states is favored. In summary, these results combined with the finding that

21

1    novel *cis*-regulatory sites tend to accumulate in the non-ancestral duplicate, suggest that

2    partitioned WGD-duplicate TFs result from selection on one copy that has neofunctionalized

3    while the other maintains the ancestral gene function.


4    **Discussion**

5        WGD is unique amongst the mechanisms of gene duplication in that all loci are

6    duplicated and thus affected equally. As such, any bias in duplicate retention must occur in the

7    aftermath of the WGD event as the genome experiences further rearrangement and reduction.

8    In allopolyploids, where WGD is the consequence of merging two related parental genomes, the

9    process of fractionation (i.e. gene losses after WGD) is biased with regards to the parent of

10   origin and the distribution of deleterious mutations (Thomas et al. 2006; van Hoek and Hogeweg

11   2007; Schnable et al. 2011). However, we would not expect this to result in different retention

12   rates amongst function groups unless genes with certain functions have different predispositions

13   to loss of function following WGD. Therefore, the significant differences in the frequency of

14   duplicate retention between function groups are likely the result of selective pressures on these

15   gene functions relative to their duplicate status.

16       It is well established that genes with certain molecular functions are enriched among

17   retained WGD-duplicates (Blanc and Wolfe 2004; Carretero-Paulet and Fares 2012), including

18   protein kinases and TFs (Maere et al. 2005; Shiu et al. 2005). However, the reason these genes

19   are enriched amongst WGD-duplicate was unknown. In this study, we have shown that

20   duplicates are retained at different rates across function groups depending on their expression,

21   conservation, and structure, suggesting these features are related to the selection for or against

22   retaining a duplicate pair. Importantly, the trends that apply to function groups in general do not

23   necessarily hold for individual duplicates, suggesting unaccounted for interaction between

24   function and features. In addition, the retention rate of TF WGD duplicates, once normalized

25   against the retention rates of other genes in the genome, is higher for older duplication events.

22

1    We propose two explanations for this observation: (1) TF duplicates produced by the γ event

2    were retained more frequently than duplicates produced by later events, (2) TF duplicates

3    retained from the γ event were more frequently retained following subsequent duplication by

4    later events (see **Supplemental File S1**). It has been shown that duplicate retention in *A.*

5    *thaliana* is affected by whether the ancestral copy had been transposed prior to the WGD event

6    (Woodhouse et al. 2011), so the notion of prior genome duplicate/rearrange events effecting

7    subsequent retention is not unprecedented. Distinguishing these two scenarios would provide

8    insight into how subsequent duplication events interact with existing WGD duplicates, but in

9    either case, our results suggest non-equivalence either between WGD events or between

10    WGD-duplicates based on their duplication state prior to a WGD event.

11    Why is there an apparent bias in favor of retaining TF WGD-duplicates? Surprisingly,

12    despite the fact that the TF duplicates were generated ~50 and ~140 million years ago for the α

13    and β WGD events, respectively; retention of ancestral expression is still the most common

14    outcome and even amongst our 1,239 partitioned duplicate pairs, 83.1% have retained

15    ancestral expression in both copies under ≥5 other conditions. The frequency with which

16    ancestral expression is conserved suggests that the retention of a significant number of TF

17    duplicates is not due to selection acting on their divergent functions. Furthermore, the fact that a

18    single duplicate pair can both retain ancestral expression under some conditions and be non-

19    randomly partitioned under other conditions, raises the possibility that both subfunctionalization

20    (Force et al. 1999) and gene balance played a role in retention of these WGD-duplicates

21    (Birchler and Veitia 2007; Birchler and Veitia 2010). The question then becomes whether initial

22    retention by gene balance creates the opportunity for later subfunctionalization or if gene

23    balance restrains the divergence of ancestral expression in duplicates pairs that have already

24    been selected for.

25    The observed evolution of extant gene expression post duplication demonstrates a

26    preference for maintaining partitioned functional states amongst retained TF WGD-duplicate

23

1    pairs where one duplicate mirrors the ancestral gene's functional states while the other copy

2    inherits significantly fewer ancestral states. On an experiment-by-experiment basis, the

3    partitioning of ancestral expression states appears to support the notion of WGD-duplicate

4    retention by subfunctionalization (Force et al., 1999). However, when examining the ancestral

5    state partitioning patterns across multiple experiments, we find an extreme bias where one TF

6    retains most of the ancestral states and other, non-ancestral, copy retains few. Most

7    importantly, the non-ancestral copy tends to gain novel *cis*-regulatory sites. This pattern harkens

8    back to the notion of there being an ancestral copy and a neofunctionalized copy after

9    duplication, contributing to the retention of both duplicates (Ohno, 1970). Thus, for duplicate

10    copies with significant expression and *cis*-regulatory site differences, both neofunctionalization

11    and subfunctionalization are likely important for duplicate retention, with the former playing a

12    more important role.

13          Rather than suggesting a singular explanation of why TF WGD-duplicates are retained,

14    our findings suggest a nuanced pattern of expression and *cis*-regulatory state evolution between

15    duplicates. This raises the question of what are the relative contributions of theorized models of

16    duplicate retention, such as subfunctionalization, gene balance, neofunctionalization, and

17    escape from adaptive conflict. One way to assess contributions of different retention models

18    would be to extend our differential equation models to include more recent WGD events. The α,

19    β, and γ events are relatively ancient taking place ~50-300 million years ago (Bowers et al.

20    2003). As such, any part of our model prior to the α event, including the initial rapid drop of

21    WGD-duplicates in state O, are extrapolated from what best fits the data from the α, β, and γ

22    events. In addition, for the most ancient γ event, there are so few duplicates that the model

23    quality is significantly impacted. Thus, instead of focusing on even more ancient events,

24    incorporating earlier events could reliably resolve the initial behavior of WGD-duplicate and

25    elucidate the roles of different models of retention. For example, were data from more recent

26    WGD events to indicate initial lag period prior to the accumulation of WGD-duplicates with

24

1  partitioned ancestral functions, it would support the hypothesis that future sub- and/or neo-

2  functionalization is enabled by initial retention because of gene balance (Veitia et al. 2013).

3      Additionally, while WGD-Duplicates TFs are known to be preferentially retained across

4  many plant species (Carretero-Paulet and Fares 2012), the patterns of ancestral expression and

5  regulatory site partitioning we have uncovered is in *A. thaliana*. It remains unclear if these

6  patterns are specific to *A. thaliana*, shared by other plant lineages sharing the α, β, and γ

7  events, or is common to any species with similar patterns of WGD events. Even if the overall

8  pattern of TF evolution is consistent across multiple species, it is possible that TF families may

9  evolve differently from each other. In future studies, it will be important to directly compare the

10 size, rate of retention, and rate of partitioning both within and across species in individual

11 families. A multi-species comparison will also be informative for further our understanding of the

12 interplay between WGD and TF retention. For example, we found that γ duplicate TFs have a

13 relatively higher rate of retention compared to the α and β events. Comparing retention of

14 duplicate TFs across multiple species which share the γ duplication event, but have

15 independent subsequent WGDs would allow us to assess how older and younger WGD events

16 may jointly influence TF retention.

17     Although there are many questions yet to be answered about the factors affecting

18 whether or not a duplicate pair is retained following WGD, we have found that the frequency

19 with which duplicate genes are retained following a WGD events is strongly correlated with the

20 expression, conservation, and structural features of that gene. And although our general model

21 does not give perfect predictions for all function groups, it can serve as a basis for exploring

22 more complicated interactions underlying duplicate retention, namely, the potential interaction

23 between gene features and annotated gene function suggested by our results. Furthermore, the

24 partitioning of ancestral expression and the non-random distribution of ancestral and new *cis*-

25 regulatory sites suggest that selection for both existing and novel functions plays a major role in

26 the retention of TF duplicates. In contrast, most duplicates retain ancestral expression levels

25

1    under some conditions in both duplicates, making it likely that the complete retention of

2    ancestral expression is preferred under certain conditions. Further investigation of which genes

3    and, more specifically, which expression/regulatory states are preferentially retained,

4    partitioned, or neofunctionalized will benefit from the modeling approaches developed here and

5    increase our knowledge of how duplication interacts with gene function


6    **Methods**


7    **Genome sequence, gene annotation, and Expression Data**

8    Genome sequences, protein sequences, and gene annotation information for *A. thaliana*

9    was obtained from Phytozome v10 (https://phytozome.jgi.doe.gov/pz/portal.html). WGDs were

10    defined according to Bowers et al. (2003) and tandem genes in *A. thaliana* were defined as

11    pairs of reciprocal best BLAST hits with an e-value < 1e-10 and ≤ 5 intervening genes.

12    Expression microarray data for this study was taken from AtGenExpress (Schmid et al. 2005;

13    Kilian et al. 2007; Goda et al. 2008), normalized using RMA (Irizarry et al. 2003) in R as

14    performed previously (Zou et al. 2009a), and divided into four groups: control conditions (in

15    environmental condition experiments, Ctrl), light and development set (LightDev), abiotic and

16    biotic stress treatments (Stress), and differential expression between stress treatments and

17    controls (Diff). The individual experiments in the Controls, Development and Light, and Stress

18    Treatment groups are described in **Table S4**. The Diff data contains the log2 normalized

19    difference between data sets for each stress condition/treatment/duration and its corresponding

20    controls. In addition to microarray data, we have included a set of 214 RNA-sequencing

21    samples (**Table S5**) from *A. thaliana* Col1 wildtype from the Sequence Read Archive

22    (https://www.ncbi.nlm.nih.gov/sra). Raw sequence reads were processed using Trimmomatic

23    (Bolger et al. 2014), with a quality threshold of 20, window size of 4, and hard-clipping length of

24    3 for leading and trailing bases. Processed reads were then mapped to the *A. thaliana* genome

1    using Tophat2 (Kim et al. 2013) and expression levels calculated with Cufflinks (Trapnell et al.

2    2010), both with a maximum intron length of 5,000bp.

3    **Defining TFs and other groups of genes in *A. thaliana***

4         TFs were defined according to the criteria used by the Plant Transcription Factor

5    Database (Jin et al. 2014), which has annotated 1,717 unique TF loci in *A. thaliana*. Additionally,

6    19 additional function groups were defined using Gene Ontology (GO) Terms in the molecular

7    function and biological process categories (The *Arabidopsis* Information Resource,

8    https://www.arabidopsis.org/), each containing 100-2,000 genes and ≥20 WGD-duplicate pairs.

9    We excluded GO:0006355 (regulation_of_transcription, DNA-templated) due to its substantial

10   overlap with the TF group we have defined above and because this category include other types

11   of regulators in addition to DNA-binding TFs. The GO terms for the other function groups

12   include: ATP Binding (GO:0005524), catalytic activity (GO:0003824), defense response

13   (GO:0006952), DNA endoreduplication (GO:0042023), hydrolase activity hydrolyzing O-glycosyl

14   compounds (GO:0004553), kinase activity (GO:0016301), lipid binding (GO:0008289),

15   oxidoreductase activity (GO:001649), oxygen binding (GO:0019825), protein binding

16   (GO:0005515), proteolysis (GO:0006508), response to auxin (GO:0009733), response to chitin

17   (GO:0010200), RNA binding (GO:0003723), transferase activity, transferring glycosyl groups

18   (GO:0016757), translation (GO:0006412), transporter activity (GO:0005215), ubiquitin-protein

19   transferase activity (GO:0004842), zinc ion binding (GO:0008270). A list of genes in each group

20   can be found in **Table S6**.

21   **Fitting duplicate retention rate within each group of genes for each WGD event using**

22   **linear models**

23        A gene was designated as a WGD-duplicate if its paralog derived from a particular WGD

24   event is present. For a gene without its paralog from WGD, it was designated as a WGD-

27

1    singleton gene. The retention rate for each function group, $g$, after a specific WGD event, $w$, is

2    defined as:

$$R_{g,w} = \frac{(D_{g,w}/S_{g,w})}{(D_{\neg g,w}/S_{\neg g,w})}$$

3    Where $D_{\neg g,w}$ and $D_{g,w}$ are the numbers of WGD-duplicate genes in group $g$ and those not in

4    group g (¬g), respectively. $S_{\neg g,w}$ and $S_{g,w}$ are the numbers of WGD-singleton genes in group $g$

5    and those not in group g (¬g), respectively. For each gene group/WGD event, we established a

6    general linear model with the glm function in the R environment which relates the $R_{g,w}$ to a set of

7    features of each gene group. The 34 features (predictor variables) (**Table S2**) were filtered to

8    prevent over-fitting because the number of observed retention rates was 20. We calculated the

9    correlation between all features to find all cases where the absolute value of correlation was >

10   0.7. The considerations for which features to keep included: (1) how well each feature

11   correlated with $R_{g,w}$ on its own, (2) whether the feature was derived from a subset of another

12   feature, and (3) the number of other features with a correlation > 0.7 (favored the elimination of

13   more features). In addition to the above criteria, one data set (protein-protein interactions) was

14   eliminated because of a high frequency of missing values (88%). The synonymous substitution

15   rate ($d_S$) feature and any feature using $d_S$ in their calculation were also excluded because they

16   would be highly correlated with WGD timing and confound our analyses comparing the three

17   WGD events. The filtering step left 11 features for building the general linear model in the

18   following, iterative fashion. Following fitting the glm function, features were ranked according to

19   their $p$ values from the least to the greatest and the feature with the largest $p$ value was

20   dropped. The model was then fit to the reduced feature set and features were once again

21   ranked. This process was repeated until the F-statistic (a measure of goodness of fit of the

22   given model against a null model where all coefficients are set to zero) of the model was

23   maximized and the final $p$ value was calculated based on the maximal F-stat. The final model

24   for each event can be found in **Supplementary File S2**.

**1** **Predicting WGD-duplicate retention status of individual genes using machine learning**

**2** In addition to the linear model, machine learning models for each group of genes (TFs,

**3** kinases, or all genes in the genome) were generated to predict whether a gene in a particular

**4** group had a retained WGD paralog from either the α, β, and γ event. The machine learning was

**5** performed using the Random Forest algorithm implement in the R package randomForest

**6** (https://cran.r-project.org/web/packages/randomForest/index.html). We filtered the gene level

**7** feature set from a previous study (Lloyd et al. 2015) by removing those with missing values for

**8** ≥5% of genes. For the remaining features, missing values were imputed with the rfImpute

**9** algorithm in randomForest using 10 iterations of 500 trees. The final matrix of genes and

**10** features for TFs, kinases, and the whole genome can be found in **Tables S7, S8,** and **S9**,

**11** respectively. Using the imputed data set for each group of genes and for each WGD event, we

**12** ran the Random Forest algorithm 10 times with 500 trees (each time with 10 fold cross

**13** validation) and collected the resulting votes (retained or not) for constructing Receiver Operating

**14** Characteristic curves (ROCs). The importance of each individual feature was assessed using

**15** Mean Decrease in Accuracy (MDA), the average number of genes misclassified across multiple

**16** runs as a result of removing the feature in question. The statistical significance of the difference

**17** in values of a feature between WGD-duplicates and WGD-singletons was determined using

**18** Welch's t-test.

**19** **Inferring ancestral expression levels and *cis*-regulatory sites**

**20** DNA-binding domains were identified in TF protein coding sequences using hmmscan

**21** via HMMER3 (Mistry et al. 2013) based on the Pfam-A HMMs (version 29.0, Finn et al. 2016)

**22** with a threshold e-value of 1e-5. TFs were classified into families according to their DNA-binding

**23** domains and 44 of 59 TF families with ≥4 members were used for further analysis (**Table S10**).

**24** For each TF family, full-length protein sequences were aligned using MAFFT (Katoh and

**25** Standley 2013) with default parameters. The phylogeny of each TF family was obtained using

29

1 RAxML (Stamatakis 2014) with the following approach: rapid Bootstrapping algorithm, 100 runs,

2 GAMMA rate heterogeneity, and the JTT amino-acid substitution model. These trees were then

3 mid-point rooted with retree in PHYLIP (Felsenstein, 1989). These trees were then used to infer

4 the ancestral gene expression states and the *cis*-regulatory sites of WGD-duplicate TF pairs

5 with BayesTrait (Pagel et al. 2004) as was done in our earlier study (Zou et al. 2009a). The

6 expression data sets used are described in **Table S4**. The discretized gene expression state

7 (0,1,2,3) was based on the quartiles of gene expression levels within each experiment. Thus the

8 inferred, ancestral expression state was also discretized. For *cis*-regulatory sites the binding

9 targets of 345 *A. thaliana* TFs were defined based DAP-Seq data generated by O'Malley et al.

10 (2016). All in-vitro binding targets from the Plant Cistrome Database

11 (http://neomorph.salk.edu/dap_web/pages/index.php) where 5% of the read associated with a

12 site were found to be in the 200bp peak region. The inference was whether a particular site was

13 present or absent (0,1). For both expression and DAP-Seq data, in cases where there was a

14 missing value, it was explicitly included as an ambiguous state. To call the ancestral state from

15 the expression or *cis*-regulatory site data, we required a posterior probability >0.5. Cases where

16 the called state was ambiguous or no majority existed were excluded from further analysis.

17 **Asymmetry of the retention of ancestral expression and regulatory sites**

18  For determining expression state asymmetry, only TF WGD-duplicates with ≥5

19 partitioned ancestral expression states in one of the four expression datasets (Ctrl, LightDev,

20 Stress, and Diff) were considered. For a WGD-duplicate pair with genes A and B, if the number

21 of inherited ancestral expression states in A was larger or equal to that in B, then A and B were

22 defined as the ancestral and the non-ancestral duplicate copies, respectively. The degree of

23 asymmetry ($Y_{A,B}$) of expression states between two duplicates was defined as:

$$Y_{A,B} = \max(F_A, F_B) - (1 - \max(F_A, F_B))$$

1    Where $F_A$ and $F_B$ are the frequency with which ancestral expression was retained across the

2    partitioned states for duplicates A and B respectively. By definition, $F_A + F_B = 1$, such that $Y_{A,B}$ is

3    value between 0 and 1 which equals 0 when $F_A = F_B$ (no asymmetry) and 1 when either $F_A$ or $F_B$

4    = 1 (maximum asymmetry).

5        With the asymmetry values for each TF pair, an average asymmetry value was

6    calculated for each expression dataset, as well as for the union of all WGD-duplicates from all

7    datasets (1,239 values total) to assess how the observed degree of asymmetry compared to

8    what would be expected from random partitioning. The expected distribution of asymmetry

9    values for the expression states of TF WGD-duplicates was determined by randomly partitioning

10   a data set with the same number of pairs and the same distribution of expression states

11   amongst pairs 1,000 times.

12       For *cis*-regulatory site asymmetry, only TF WGD-duplicates with ≥5 inferred ancestral

13   *cis*-regulatory sites we considered (402 WGD-duplicate pairs total). Similar to expression state

14   asymmetry, in each duplicate pair the ancestral and non-ancestral duplicates were defined

15   according to the number of inherited ancestral sites. For each WGD-duplicate pair, the degree

16   of asymmetry of *cis*-regulatory site among a TF pair was defined analogous to what was done

17   for expression. The expected distribution of asymmetry values for the *cis*-regulatory sites of TF

18   WGD-duplicates was determined by randomly partitioning a data set with the same number of

19   pairs and the same distribution of expression states amongst pairs 1,000 times

20

21   **Ordinary differential equation models of TF function evolution**

22       The change in expression states from the ancestral expression quartile to either a higher

23   or lower quartile in an extant TF was modeled as a system of ordinary differential equations

24   such that:

$$\frac{d}{dt}\begin{pmatrix} 0 \\ + \\ - \end{pmatrix} = \begin{pmatrix} 0 \\ + \\ - \end{pmatrix}\begin{pmatrix} -(a+b) & c & d \\ -c & b & 0 \\ -d & 0 & a \end{pmatrix}$$

31

1    Where *O*, +, and - are the frequency of TF WGD duplicate genes retaining the ancestral

2    expression states, a higher than ancestral expression, and a lower than ancestral expression,

3    respectively. The parameters a, b, c, and d, define the transition rate between these states. This

4    system of equations was solved in Maxima (http://maxima.sourceforge.net/index.html) and best

5    parameters for the observed distribution of duplicates pairs were determined using maximum

6    likelihood estimates calculated with the bbmle package in R (https://cran.r-project.org/web/

7    packages/bbmle/index.html). Non-linear minimization was used to approximate an initial guess,

8    although the actual initial parameters often needed to be adjusted to reach a convergent

9    solution. The best fit parameters for this single duplicate expression state evolution model can

10    be found in **Table S11**.

11        The loss of ancestral expression states in a pair of duplicated TFs was modeled as a

12    system of ordinary differential equations such that:

$$\frac{d}{dt}\begin{pmatrix} O \\ I \\ II \end{pmatrix} = \begin{pmatrix} O \\ I \\ II \end{pmatrix}\begin{pmatrix} -a & b & 0 \\ b & -(a+c) & d \\ 0 & c & -d \end{pmatrix}$$

13    Where *O, I,* and *II* are the frequency of TF WGD duplicate pairs where both, one, or neither

14    duplicate retained the ancestral expression state. The parameters a, b, c, and d, define the

15    transition rate between these states. This system of equations was solved and the initial and

16    best parameters were estimated in the same fashion as above. The best fit parameters for this

17    pairwise expression state evolution model can be found in **Table S12**. The same model was

18    also applied to ancestral regulatory sites with *O, I,* and *II* representing the frequency of TF WGD

19    duplicate pairs where both, one, or neither duplicate retained the ancestral regulatory site.

20    **Acknowledgements**

3    **References**

4    Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, Cherry JM, Davis AP, et al. Gene

5        Ontology: Tool for the Unification of Biology. The Gene Ontology Consortium. Nature

6        Genet. 2000; 25: 25–29.

7    Baker CR, Hanson-Smith V, Johnson AS. Following Gene Duplication, Paralog Interference

8        Constrains Transcriptional Circuit Evolution. Science. 2013; 342: 104–8.

9    Birchler JA, Veitia RA. The Gene Balance Hypothesis: From Classical Genetics to Modern

10       Genomics. Plant Cell. 2007; 19: 395–402.

11   Birchler JA, Veitia RA. The Gene Balance Hypothesis: Implications for Gene Regulation,

12       Quantitative Traits and Evolution. New Phytol. 2010; 186: 54–62.

13   Blanc G, Wolfe KH. Functional Divergence of Duplicated Genes Formed by Polyploidy during

14       Arabidopsis Evolution. Plant Cell. 2004; 16: 1679–91.

15   Bolger AM, Lohse M, Usadel B. Trimmomatic: A Flexible Trimmer for Illumina Sequence Data.

16       Bioinformatics. 2014; 30: 2114–20.

17   Bowers JE, Chapman BA, Rong J, Paterson AH. Unravelling Angiosperm Genome Evolution by

18       Phylogenetic Analysis of Chromosomal Duplication Events. Nature. 2003; 422: 433–38.

19   Carretero-Paulet L, Fares MA. Evolutionary Dynamics and Functional Specialization of Plant

20       Paralogs Formed by Whole and Small-Scale Genome Duplications. Mol Bio Evol. 2012; 29:

21       3541–51.

22   Dehal P, Boore JL. Two Rounds of Whole Genome Duplication in the Ancestral Vertebrate.

23       PLoS Biol. 2005; 3: e314.

24   Des Marais DL, Rausher DM. Escape from adaptive conflict after duplication in an anthocyanin

pathway gene. Nature. 2008; 454: 762-5.

Felsenstein, J. PHYLIP - Phylogeny Inference Package (Version 3.2). Cladistics.1989; 5: 164-166.

Finn RD, Coggill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, et al. The Pfam Protein Families Database: Towards a More Sustainable Future. Nucleic Acids Res. 2016; 44: D279–85.

Force A, Lynch M, Pickett FB, Amores A, Yan YL, Postlethwait J. "Preservation of Duplicate Genes by Complementary, Degenerative Mutations." Genetics. 1999; 151: 1531–45.

Hideki G, Sasaki E, Akiyama K, Maruyama-Nakashita A, Nakabayashi K, Li W, et al. The AtGenExpress Hormone and Chemical Treatment Data Set: Experimental Design, Data Evaluation, Model Data Analysis and Data Access. Plant J. 2008; 55: 526–42.

Hanada K, Zou C, Lehti-Shiu MD, Shinozaki K, and Shiu SH. Importance of Lineage-Specific Expansion of Plant Tandem Duplicates in the Adaptive Response to Environmental Stimuli. Plant Physiol. 2008; 148: 993–1003.

van Hoek MJ, Hogeweg P. The Role of Mutational Dynamics in Genome Shrinkage. Mol Bio Evol. 2007; 24: 2485–94.

Irizarry RA, Hobbs B, Collin F, Beazer-Barclay YD, Antonellis KJ, Scherf U, et al. Exploration, normalization, and summaries of high density oligonucleotide array probe level data. Biostatistics. 2003; 4: 249-64.

Jiang WK, Liu YL, Xia EH, Gao LZ. Prevalent Role of Gene Features in Determining Evolutionary Fates of Whole-Genome Duplication Duplicated Genes in Flowering Plants. Plant Physiol. 2013; 161: 1844–61.

Jiao Y, Wickett NJ, Ayyampalayam S, Chanderbali AS, Landherr L, Ralph PE, et al. Ancestral Polyploidy in Seed Plants and Angiosperms. Nature. 2011; 473: 97–100.

Jin J, Zhang H, Kong L, Gao G, Luo J. PlantTFDB 3.0: A Portal for the Functional and Evolutionary Study of Plant Transcription Factors. Nucleic Acids Res. 2014; 42: D1182–87.

34

1    Katoh K, Standley DM. MAFFT Multiple Sequence Alignment Software Version 7:

2         Improvements in Performance and Usability. Mol Bio Evol. 2013; 30: 772–80.

3    Kellis M, Birren BW, Lander ES. Proof and Evolutionary Analysis of Ancient Genome

4         Duplication in the Yeast Saccharomyces Cerevisiae. Nature. 2004; 428: 617–24.

5    Kilian J, Whitehead D, Horak J, Wanke D, Weinl S, Batistic O, et al. The AtGenExpress Global

6         Stress Expression Data Set: Protocols, Evaluation and Model Data Analysis of UV-B Light,

7         Drought and Cold Stress Responses. Plant J. 2007; 50: 347–63.

8    Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: Accurate Alignment

9         of Transcriptomes in the Presence of Insertions, Deletions and Gene Fusions. Genome

10        Biology. 2013; 14: R36.

11   Lee, TH, Tang H, Wang X, Paterson AH. PGDD: A Database of Gene and Genome Duplication

12        in Plants. Nucleic Acids Res. 2013; 41: D1152–58.

13   Lehti-Shiu MD, Panchy N, Wang P, Uygun S, Shiu SH. Diversity, expansion, and evolutionary

14        novelty of plant DNA-binding transcription factor families. BBA. 2016; 1860: 3-20.

15   Lespinet O, Wolf YI, Eugene V. Koonin EV, Aravind L. The Role of Lineage-Specific Gene

16        Family Expansion in the Evolution of Eukaryotes. Genome Res. 2002; 12: 1048–59.

17   Li Z, Defoort J, Tasdighian S, Maere S, Van de Peer Y, De Smet R. Gene Duplicability of Core

18        Genes Is Highly Consistent across All Angiosperms. Plant Cell. 2016; 28: 326–44.

19   Litt A, Irish VF. Duplication and Diversification in the APETALA1/FRUITFULL Floral Homeotic

20        Gene Lineage: Implications for the Evolution of Floral Development. Genetics. 2003; 165:

21        821-33.

22   Liu H, Liu H, Zhou L, Zhang Z, Zhang X, Wang M, et al. Parallel Domestication of the Heading

23        Date 1 Gene in Cereals. Mol Biol Evol. 2015; 32: 2726-37.

24   Lloyd JP, Seddon AE, Moghe GD, Simenc MC, Shiu SH. Characteristics of Plant Essential

25        Genes Allow for within- and between-Species Prediction of Lethal Mutant Phenotypes.

26        Plant Cell. 2015; 27: 2133–47.

1   Lynch M, Conery JS. The Evolutionary Fate and Consequences of Duplicate Genes. Science.

2       2000; 290: 1151–55.

3   Lynch M, Conery JS. The Evolutionary Demography of Duplicate Genes. J Struct Funct

4       Genomics. 2003; 3:35–44.

5   Lyons E, Pedersen B, Kane J, Alam M, Ming R, Tang H, et al. Finding and Comparing Syntenic

6       Regions among Arabidopsis and the Outgroups Papaya, Poplar, and Grape: CoGe with

7       Rosids. Plant Physiol. 2008; 148: 1772–81.

8   Maere S, De Bodt S, Raes J, Casneuf T, Van Montagu M, Kuiper M, et al. Modeling Gene and

9       Genome Duplications in Eukaryotes. Proc Natl Acad Sci USA. 2005; 102: 5454–59.

10  McCarthy EW, Mohamed A, Litt A. Functional Divergence of APETALA1 and FRUITFULL is due

11      to Changes in both Regulation and Coding Sequence. Front Plant Sci. 2015; 6: 1076.

12  Mistry J, Finn RD, Eddy SR, Bateman A, Punta M. Challenges in Homology Search: HMMER3

13      and Convergent Evolution of Coiled-Coil Regions. Nucleic Acids Res. 2013; 41: e121.

14  Moghe GD, Hufnagel DE, Tang H, Xiao Y, Dworkin I, Town CD, et al. Consequences of Whole-

15      Genome Triplication as Revealed by Comparative Genomic Analyses of the Wild Radish

16      Raphanus Raphanistrum and Three Other Brassicaceae Species. Plant Cell. 2014; 26:

17      1925–37.

18  Moghe GD, Shiu SH. The Causes and Molecular Consequences of Polyploidy in Flowering

19      Plants. Ann NY Acad Sci. 2014; 1320: 16–34.

20  Myburg AA, Grattapaglia D, Tuskan GA, Hellsten U, Hayes RD, Grimwood J, et al. The Genome

21      of Eucalyptus Grandis. Nature. 2014; 510: 356–62.

22  O'Malley RC, Huang SS, Song L, Lewsey MG, Bartlett A, Nery JR, et al. Cistrome and

23      Epicistrome Features Shape the Regulatory DNA Landscape. Cell. 2016; 165: 1280–92.

24  Oakley TH, Østman B, Wilson ACV. Repression and loss of gene expression outpaces

25      activation and gain in recently duplicated fly genes. Proc Natl Acad Sci USA. 2006; 103:

26      11637-41.

1  Ohno S. Evolution by Gene Duplication. New York: Springer-Verlag; 1970.

2  Pagel M, Meade A, Barker D. Bayesian Estimation of Ancestral Character States on

3     Phylogenies. Syst Biol. 2004; 53: 673–84.

4  Panchy N, Lehti-Shiu M, Shiu SH. Evolution of Gene Duplication in Plants. Plant Physiol. 2016;

5     171: 2294–2316.

6  Panopoulou G, Hennig S, Groth D, Krause A, Poustka AJ, Herwig R, et al. New Evidence for

7     Genome-Wide Duplications at the Origin of Vertebrates Using an Amphioxus Gene Set and

8     Completed Animal Genomes. Genome Research. 2003; 13: 1056–66.

9  Renny-Byfield S, Gallagher JP, Grover CE, Szadkowski E, Page JT, Udall JA, et al. Ancient

10    Gene Duplicates in Gossypium (cotton) Exhibit near-Complete Expression Divergence.

11    Genome Biol Evol. 2014; 6: 559–71.

12  Schmid M, Davison TS, Henz SR, Pape UJ, Demar M, Vingron M, et al. A Gene Expression

13    Map of Arabidopsis Thaliana Development. Nature Genet. 2005; 37: 501–6.

14  Schnable JC, Springer NM, Freeling M. Differentiation of the maize subgenomes by genome

15    dominance and both ancient and ongoing gene loss. Proc Natl Acad Sci USA. 2011; 108:

16    4069-74.

17  Schranz ME, Quijada P, Sung SB, Lukens L, Amasino R, Osborn TC. Characterization and

18    Effects of the Replicated Flowering Time Gene *FLC* in *Brassica rapa*. Genetics. 2002; 3:

19    1457-68.

20  Seoighe C, Gehring C. 2004. Genome Duplication Led to Highly Selective Expansion of the

21    Arabidopsis Thaliana Proteome. Trends Genet. 2004; 20: 461–64.

22  Shiu SH, Shih MC, Li WH. Transcription Factor Families Have Much Higher Expansion Rates in

23    Plants than in Animals. Plant Physiol. 2005; 139: 18–26.

24  Soltis DE, Visger CJ, Soltis PS. The Polyploidy Revolution Then…and Now: Stebbins Revisited.

25    Am J Bot. 2014; 101: 1057–78.

26  Stamatakis,A. 2014. RAxML Version 8: A Tool for Phylogenetic Analysis and Post-Analysis of

1      Large Phylogenies. Bioinformatics. 2014; 30: 1312–13.

2    Theissen G, Melzer R. Molecular Mechanisms Underlying Origin and Diversification of the

3      Angiosperm Flower. Ann Bot. 2007; 100: 603-619.

4    Thomas BC, Pedersen B, Freeling M. Following Tetraploidy in an Arabidopsis Ancestor, Genes

5      Were Removed Preferentially from One Homeolog Leaving Clusters Enriched in Dose-

6      Sensitive Genes. Genome Res. 2006; 16: 934–46.

7    Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, et al. Transcript

8      Assembly and Quantification by RNA-Seq Reveals Unannotated Transcripts and Isoform

9      Switching during Cell Differentiation. Nature Biotechnol. 2010; 28: 511–15.

10   Veitia RA, Bottani S, Birchler JA. Gene Dosage Effects: Nonlinearities, Genetic Interactions, and

11     Dosage Compensation. Trends Genet. 2013; 29: 385–93.

12   Wang K, Wang Z, Li F, Ye W, Wang J, Song G, et al. The Draft Genome of a Diploid Cotton

13     Gossypium Raimondii. Nat Genet. 2012; 44: 1098–1103.

14   Wang W, Haberer G, Gundlach H, Gläßer C, Nussbaumer T, Luo MC, et al. The Spirodela

15     Polyrhiza Genome Reveals Insights into Its Neotenous Reduction Fast Growth and Aquatic

16     Lifestyle. Nat Commun. 2014; 5: 3311.

17   Weirauch MT, Hughes TR. A catalogue of eukaryotic transcription factor types, their

18     evolutionary origin, and species distribution. Subcell Biochem. 2011; 52: 25-73.

19   Wolfe KH, Shields DC. Molecular Evidence for an Ancient Duplication of the Entire Yeast

20     Genome. Nature. 1997; 387: 708–13.

21   Woodhouse MR, Tang H, Freeling M. Different Gene Families in Arabidopsis thaliana

22     Transposed in Different Epochs and at Different Frequencies throughout the Rosids. Plant

23     Cell. 2011; 23: 4241-53.

24   Zhang Z, Belcram H, Gornicki P, Charles M, Just J, Huneau C, et al. Duplication and partitioning

25     in evolution and function of homoeologous Q loci governing domestication characters in

26     polyploid wheat. Proc Natl Acad Sci USA. 2011; 108: 18737-42.

1    Zou C, Lehti-Shiu MD, Thibaud-Nissen F, Prakash T, Buell CR, Shiu SH. Evolutionary and

2        Expression Signatures of Pseudogenes in Arabidopsis and Rice. Plant Physiol. 2009; 151:

3        3–15.

4    Zou C, Lehti-Shiu MD, Thomashow M, Shiu SH. Evolution of Stress-Regulated Gene

5        Expression in Duplicate Genes of Arabidopsis Thaliana. PLoS Genet. 2009; 5: e1000581.

6    Zou C, Sun K, Mackaluso JD, Seddon AE, Jin R, Thomashow MF, et al. Cis-Regulatory Code of

7        Stress-Responsive Transcription in Arabidopsis Thaliana. Proc Natl Acad Sci USA. 2001;

8        108: 14992–97.

9

1  **Figure Legends**

2  **Figure 1.** Retention of WGD-duplicate genes in *A. thaliana*

3  The odds of duplicate gene retention within 20 function groups relative to whole genome.

4  Groups are ordered by the odds of duplicate retention in the alpha event. Colors represent

5  different WGD duplication events (α = orange, ß = green, γ = blue). Error bars indicate the 95%

6  confidence interval of the odds of retention as determined by the Fisher's Exact Test

7  implemented in R.

8

9  **Figure 2. (A)** Relationships between gene counts and retention rate of WGD duplicates across

10  functional groups (α = orange, ß = green, γ = blue). **(B)** A heatmap of the Pearson product-

11  moment correlation coefficient between the values of a feature across different function groups

12  (y-axis) and the retention rate of functions groups from a particular WGD event (x-axis, indicated

13  by the symbols α, ß, and γ). Darker red: stronger positive correlation. Darker blue: stronger

14  negative correlation. Features with different sign of correlation across WGD events are indicated

15  by black arrows. Features with a large (~0.20) difference in magnitude with the same sign are

16  indicated by white arrows. **(C)** The observed odds of duplicate retention (x-axis) for each group

17  plotted against the predicted odds of retention (y-axis) from the best model for each event (α =

18  orange, ß = green, γ = blue). The dotted line represents equality between predicted and

19  observed odds of retention. Values from TFs are indicated by a black arrow while values from

20  kinases are indicated by a white arrow.

21

22  **Figure 3.** Predicting WGD-duplicate status of *A. thaliana* genes using Random Forest.

23  Receiver-operator characteristic (ROC) curves for Random Forest models predicting duplicate

24  status of all genes in *A. thaliana* (orange), kinases (green), and TFs (blue). The dotted line

40

1    indicates the curve for random guessing with an area under the ROC curve (AUC-ROC) of 0.5.

2    The AUC-ROC of each model is indicated in parentheses.

3

4    **Figure 4.** Evolution of expression in TF WGD-duplicates. **(A)** Difference in expression quartile of

5    individual TF WGD-duplicates compared to their ancestral state. Heatmaps show the z-scores

6    of the observed frequency of each difference compared to the expected frequency for LightDev

7    (left column) and Diff (right column) data across all three duplication events (α = top, ß = middle,

8    γ = bottom). Color correlates with the magnitude of the z-score, with darker red values indicating

9    counts further above random expectation and darker blue values indicating counts further below

10   random expectation. **(B)** Deviation of pairs of TF WGD-duplicates from their ancestral state,

11   defined as the difference value that each duplicate in a pair has from its ancestral state.

12   Heatmaps show the z-scores of the observed frequency of WGD-duplicate pair deviation

13   compared to the expected frequency for LightDev (left column) and Diff (right column) data

14   across all three duplication events (α = top, ß = middle, γ = bottom). Color correlates with the

15   magnitude of the z-score as in (A)

16

17   **Figure 5.** Asymmetry of ancestral state retention in TF WGD-duplicates. **(A)** The asymmetry

18   value (see Methods) of ancestral expression partitioning between TF WGD-duplicates. **(B)** The

19   asymmetry values of ancestral *cis*-regulatory site partitioning between TF WGD-duplicates. **(C)**

20   The frequency distribution of the difference in number of novel *cis*-regulatory sites between

21   ancestral and non-ancestral WGD duplicates. Values on the x-axis to the left of zero indicate the

22   magnitude of differences in the number of sites favoring the ancestral duplicate while values to

23   the right of zero indicate the magnitude of differences favoring the non-ancestral duplicate.

24

25   **Figure 6.** ODE models of TF WGD-duplicate expression and *cis*-regulatory site evolution

26   relative to the ancestral state **(A)** In this model, we consider the transition of WGD-duplicate pair

41

1    expression between three possible states relative to their ancestral state (O = both retained, I =

2    one retained, II = neither retained) using four variables representing the rate of transition

3    between state (x,y,w,z). **(B)** Results for the one parameter version of the model (z=y=w=z)

4    showing the change in time (x-axis) of the frequency (y-axis) of each WGD-duplicate-pair state

5    (O = orange, I = blue, II = green). Curves represent the continuous output of the model with

6    symbols indicate the observed values on which the models were built (O = circle, I = square, II =

7    triangle). **(C)** Results for the two parameter version of the model (z=y|w=z). Axis, color, lines

8    and symbols are used the same as in (B). **(D)** Bar graph of the parameter values (x,y,w,z) for

9    the one (orange), two (green), and four (blue) parameter versions for the expression ODE

10   model. **(E)** Results for the one parameter version of the model (z=y=w=z) showing the change

11   in time (x-axis) of the frequency (y-axis) of each WGD-duplicate-pair state (O = orange, I = blue,

12   II = green). Curves represent the continuous output of the model with symbols indicate the

13   observed values on which the models were built (O = circle, I = square, II = triangle). **(F)** Results

14   for the two parameter version of the model (z=y|w=z). Axis, color, lines and symbols are used

15   the same as in (E). **(G)** Results for the two parameter version of the model (z|y|w|z). Axis, color,

16   lines and symbols are used the same as in (E). **(H)** Bar graph of the parameter values (x,y,w,z)

17   for the one (orange), two (green), and four (blue) parameter versions for the *cis*-regulatory site

18   ODE model.

**A.**



**C.**



**B.**

**WGD event retention rate**

| | | α | β | γ | |
|---|---|---|---|---|---|
| | Paralog | 0.08 | 0.16 | -0.20 | ← |
| | A. lyrata | 0.12 | 0.20 | -0.04 | ← |
| $d_n/d_s$ | O. Sativa | -0.19 | -0.13 | -0.31 | |
| | P. Patens | -0.37 | -0.29 | -0.46 | |
| | P. trichocarpa | -0.21 | -0.26 | -0.20 | |
| | V. vinfera | -0.29 | -0.37 | -0.24 | |
| | Mean (All) | -0.21 | -0.19 | -0.20 | |
| | Mean (Stress) | -0.24 | -0.21 | -0.21 | |
| | Mean (LightDev) | -0.16 | -0.15 | -0.19 | |
| | Mean (Control) | -0.16 | -0.15 | -0.19 | |
| | Mean (Diff) | -0.40 | -0.34 | -0.04 | |
| Microarray | Median | 0.02 | 0.00 | -0.20 | ← |
| | MAD/Median | -0.43 | -0.35 | -0.19 | |
| | Breadth | -0.36 | -0.34 | -0.14 | |
| | Correlation | 0.17 | -0.14 | -0.14 | ← |
| | Correlation (Ds<2) | 0.13 | -0.01 | -0.23 | ← |
| | Corr. Module Size | 0.04 | -0.19 | -0.10 | ← |
| | Breadth | -0.03 | -0.42 | -0.06 | |
| RNASeq | Mean | 0.30 | 0.16 | -0.06 | ← |
| | Median | -0.16 | -0.26 | -0.13 | |
| | Maximum | 0.47 | 0.27 | 0.12 | ⇦ |
| Seq. Cons. | Viridiplantae | 0.05 | 0.02 | 0.14 | |
| | Fungi | 0.04 | -0.22 | -0.02 | ← |
| | Metazoa | 0.00 | -0.24 | -0.08 | |
| | Group Size | 0.21 | 0.20 | 0.48 | ⇦ |
| | Family Size | 0.22 | -0.10 | 0.01 | ← |
| | Protein Length | -0.23 | -0.36 | -0.16 | |
| | Protein Domain # | -0.15 | -0.28 | -0.05 | |
| Other | Max % ID | 0.40 | 0.21 | 0.18 | ⇦ |
| | Paralog Ds | -0.03 | 0.06 | 0.22 | ← |
| | Paralog Dn | 0.11 | 0.34 | 0.13 | ⇦ |
| | Pi | -0.04 | 0.00 | -0.23 | |
| | AraNet | 0.15 | 0.03 | -0.18 | ← |
| | PPI | -0.20 | -0.29 | 0.11 | ← |

**A.**

**Extant Quartile**

| Z-Score | LightDev 0 | 1 | 2 | 3 | Diff 0 | 1 | 2 | 3 |
|---|---|---|---|---|---|---|---|---|
| α 0 | 136.7 | -38.8 | -58.5 | -45.6 | 167.9 | -34.3 | -51.0 | -92.0 |
| α 1 | -39.4 | 114.6 | -44.8 | -42.0 | -40.4 | 139.2 | -19.6 | -53.4 |
| α 2 | -54.8 | -34.5 | 119.2 | -30.4 | -49.1 | -18.4 | 142.1 | -52.9 |
| α 3 | -46.8 | -42.6 | -17.7 | 130.4 | -90.8 | -58.1 | -49.2 | 168.1 |
| β 0 | 72.4 | -20.0 | -32.9 | -28.3 | 80.1 | -11.8 | -20.6 | -52.1 |
| β 1 | -17.9 | 60.8 | -25.2 | -27.2 | -16.5 | 63.8 | -5.5 | -25.2 |
| β 2 | -29.7 | -15.6 | 65.7 | -15.5 | -20.0 | -3.7 | 64.1 | -23.7 |
| β 3 | -29.0 | -27.2 | -4.6 | 76.2 | -51.9 | -26.9 | -20.6 | 82.4 |
| γ 0 | 50.0 | -11.4 | -23.2 | -22.5 | 39.7 | -3.2 | -13.3 | -26.8 |
| γ 1 | -14.2 | 27.1 | -1.0 | -11.2 | -7.9 | 26.3 | -1.1 | -10.8 |
| γ 2 | -20.1 | -0.3 | 32.8 | -10.2 | -16.1 | 1.1 | 30.9 | -8.2 |
| γ 3 | -23.3 | -7.6 | -6.8 | 43.4 | -23.4 | -12.9 | -8.2 | 39.4 |

*(Ancestral Quartile, rows; Z-Score columns)*

**B.**

**Duplicate 1 Deviation**

LightDev:

| Dup2 Dev | -3 | -2 | -1 | 0 | 1 | 2 | 3 |
|---|---|---|---|---|---|---|---|
| α -3 | -1 | -2 | 1 | 3 | n/a | n/a | n/a |
| α -2 | | -1 | -6 | 12 | -9 | n/a | n/a |
| α -1 | | | -13 | 24 | -16 | -5 | n/a |
| α 0 | | | | -30 | 19 | 9 | 4 |
| α 1 | | | | | -5 | -5 | -2 |
| α 2 | | | | | | -1 | -1 |
| α 3 | | | | | | | 0 |
| β -3 | 1 | 2 | 1 | 1 | n/a | n/a | n/a |
| β -2 | | 0 | -5 | 7 | -6 | n/a | n/a |
| β -1 | | | -2 | 10 | -11 | -2 | n/a |
| β 0 | | | | -13 | 8 | 3 | 1 |
| β 1 | | | | | 4 | -2 | -1 |
| β 2 | | | | | | 0 | 0 |
| β 3 | | | | | | | 0 |
| γ -3 | 1 | 0 | -1 | 2 | n/a | n/a | n/a |
| γ -2 | | 3 | 4 | 0 | -5 | n/a | n/a |
| γ -1 | | | 3 | -1 | -4 | -2 | n/a |
| γ 0 | | | | -1 | 2 | 1 | -1 |
| γ 1 | | | | | 4 | -2 | -1 |
| γ 2 | | | | | | -1 | 0 |
| γ 3 | | | | | | | 0 |

Diff:

| Dup2 Dev | -3 | -2 | -1 | 0 | 1 | 2 | 3 |
|---|---|---|---|---|---|---|---|
| α -3 | 0 | -5 | -7 | 14 | n/a | n/a | n/a |
| α -2 | | -2 | -9 | 16 | -8 | n/a | n/a |
| α -1 | | | -8 | 23 | -13 | -9 | n/a |
| α 0 | | | | -38 | 23 | 16 | 15 |
| α 1 | | | | | -8 | -9 | -8 |
| α 2 | | | | | | -4 | -5 |
| α 3 | | | | | | | -2 |
| β -3 | 1 | -3 | -5 | 8 | n/a | n/a | n/a |
| β -2 | | 4 | -3 | 5 | -4 | n/a | n/a |
| β -1 | | | 4 | 6 | -5 | -4 | n/a |
| β 0 | | | | -12 | 6 | 6 | 5 |
| β 1 | | | | | 2 | -5 | -5 |
| β 2 | | | | | | -3 | -4 |
| β 3 | | | | | | | -3 |
| γ -3 | 8 | 1 | -2 | 0 | n/a | n/a | n/a |
| γ -2 | | 4 | -1 | 1 | -2 | n/a | n/a |
| γ -1 | | | 5 | 2 | -4 | -3 | n/a |
| γ 0 | | | | -1 | -2 | 1 | 3 |
| γ 1 | | | | | 9 | -3 | -4 |
| γ 2 | | | | | | -1 | -2 |
| γ 3 | | | | | | | 0 |

*(Duplicate 2 Deviation, rows; Z-Score columns)*

**A.**



**B.**



**C.**