OXFORD
UNIVERSITY PRESS | DNA Research

# New insights into the interplay between codon bias determinants in plants

SCHOLARONE™
Manuscripts

1       **New insights into the interplay between codon bias determinants in plants**

2       S. Camiolo, S. Melito, A. Porceddu*

3

4       Università degli Studi di Sassari, Dipartimento di Agraria, SACEG, Sassari

5

6

7       **Corresponding author:**

8       A. Porceddu

9       e-mail: aporceddu@uniss.it

10

11      **Running title:** Codon bias in plants

12

13

14

15

16

17

18

19

20

21      Supplemental material is available online

1    **Abstract**

2    Codon bias is the non-random use of synonymous codons, a phenomenon has been observed in

3    species as diverse as bacteria, plants and mammals. The preferential use of particular synonymous

4    codons may reflect neutral mechanisms (e.g. mutational bias, G|C-based gene conversion, genetic

5    drift) and/or selection for mRNA stability, translational efficiency and accuracy. The extent to

6    which these different factors influence codon usage is unknown, so we dissected the contribution of

7    mutational bias and selection towards codon bias in genes from 17 eudicots and 4 monocots. We

8    analysed the frequency of mononucleotides, dinucleotides and trinucleotides, and investigated

9    whether the compositional genomic background could account for the observed codon usage

10    profiles. Neutral forces such as mutational pressure and G|C-based gene conversion appeared to

11    underlie most of the observed codon bias, although there was also evidence for the selection of

12    optimal translational efficiency and mRNA folding. Our data confirmed the compositional

13    differences between monocots and dicots, with the former featuring in general a lower background

14    compositional bias but a higher overall codon bias.

15

16

17    **Keywords:** Codon bias, mutational bias, translational selection, plant genetics.

18

19

20

21

22

23

24

1    **Introduction**

2    The genetic code is redundant, with most amino acids encoded by two or more synonymous

3    codons [1–4]. The non-random use of synonymous codons is known as codon bias, and it may reflect

4    several underlying factors including mutational bias in the genome and translational selection. The

5    possibility that mutational bias affects codon usage has led to the neutralist model, in which codon

6    identity is mainly determined by nucleotide substitution patterns in the genome. In contrast, the

7    possibility of translational selection has led to the selective model, in which the choice of

8    synonymous codons reflects tRNA abundance [5] to optimize the efficiency [4] and accuracy [6] of

9    translation. These models are not mutually exclusive, i.e. the choice among synonymous codons

10   may reflect a balance between selective and mutational pressures [7].

11   Species-dependent differences in codon usage are well known [8], but recent studies have identified

12   variations within species that must also be addressed by the neutralist and selective models. For

13   example, the abundance of tRNA may vary during development and in response to external

14   stimuli [9], suggesting that codon bias may represent an adaptive response to tRNA levels that differ

15   among plant tissues [10]. Nucleotide substitution patterns are also unequally distributed in the

16   genome [11], e.g. there are large homogeneous blocks G|C-rich sequences known as isochores in the

17   genomes of warm-blooded vertebrates [12,13]. Likewise, differences in nucleotide substitution patterns

18   have been observed in rice (*Oryza sativa*), with genes expressed in the roots being predominantly

19   G|C-rich and genes expressed in seeds and leaves being predominantly A|T-rich [14].

20   The compositional context can also influence synonymous codon selection, a phenomenon known

21   as context-dependent codon bias (CCDB). In mammals, bacteria and plants, the first nucleotide

22   after each codon drives synonymous codon choice, because several dinucleotide sequences such as

23   CG, GA and TA are underrepresented [15–17]. In plant genomes, there is also a general bias in the use

24   of specific dinucleotides and trinucleotides in different genomic regions [18].

1   The composition of coding sequences is determined by a complex series of interacting factors, so it

2   is difficult to identify the relative impact of different components. A model to determine the

3   influence of nucleotide substitution patterns on codon bias has been developed by building a new

4   set of sequences in which the third codon position in the coding sequence is replaced with a random

5   nucleotide from the neighbouring intergenic region [19]. Such intergenic corrected coding sequences

6   (ICCSs) retain the same amino acid sequence while mirroring the background nucleotide

7   substitution pattern of the genome. Comparing the codon bias between the original coding sequence

8   (CS) and ICCS datasets can therefore highlight the influence of the background nucleotide

9   substitution pattern on the coding sequence composition [19].

10  Although the selective model has been studied in several plant species, the impact of background

11  composition on gene structure has been largely overlooked. Codon bias in *Arabidopsis thaliana*

12  (Arabidopsis) tends to be associated with the composition of the 3' flanking region in both strongly

13  and weakly expressed genes [20], although the impact of selection on both classes of genes has also

14  been recognized. Here we used the Hershberg and Petrov approach [19] in order to determine the

15  effect of background composition on the codon bias of 21 plant species while also accounting for

16  bias in the frequencies of dinucleotides and trinucleotides. We discuss in detail the impact of these

17  multiple factors and others influencing codon bias in plants.

18

19  **Materials and methods**

20  *Sequence data*

21  The genomic sequences and annotation data for 21 plant species were downloaded from Phytozome

22  (http://www.phytozome.net)[21] . We analysed the sequences of 4 monocots (*Brachypodium*

23  *distachyon* (BD), *Oryza sativa* (OS), *Sorgum bicolour* (SB) and *Zea mays* (ZM)), 15 dicots

24  (*Arabidopsis lyrata* (AL)*, A. thaliana* (AT)*, Brassica rapa* (BR)*, Citrus clementine* (CC)*, C.*

25  *sinensis* (CS)*, Eucalyptus grandis* (EG)*, Phaseolus vulgaris* (FV)*, Glycine max* (GM)*, Linum*

1    *usitatissimum* (LU)*, Medicago truncatula* (MT)*, Populus trichocarpa* (PopT)*, Solanum*

2    *lycopersicum* (SL)*, Solanum tuberosum* (ST)*, Theluginella halophile* (TH) and *Vitis vinifer*a (VV))

3    and two mosses (*Selaginella moellendorffii* (SM)*, Physcomitrella patens* (PP)) .

4    We used gff2sequence [22] to identify coding sequences, proteins, introns and intergenic sequences.

5    Coding sequences featuring non-canonical bases (other than A, C, G or T), missing stop codons or

6    incomplete triplets were excluded. Finally the longest splicing variant was chosen when multiple

7    transcripts representing the same gene were annotated. Monocot coding sequences were divided

8    into three subsets for analysis: (1) the entire genome, (2) high-G|C sequences (GC content > 60%,

9    HGC) and (3) low-G|C sequences (GC content ≤ 60%, LGC).

10

11   *Expression data*

12   Gene expression data for Arabidopsis and rice were downloaded from the Plexdb database

13   ([http://www.plexdb.org/](http://www.plexdb.org/))[23]. We chose the expression atlases representing Arabidopsis dataset AT40

14   ([http://www.plexdb.org/modules/PD_browse/experiment_browser.php?experiment=AT40](http://www.plexdb.org/modules/PD_browse/experiment_browser.php?experiment=AT40)) and rice

15   dataset                                                                                    OS5

16   ([http://www.plexdb.org/modules/PD_browse/experiment_browser.php?experiment=OS5](http://www.plexdb.org/modules/PD_browse/experiment_browser.php?experiment=OS5)).  All  the

17   expression data were RMA normalized. An expression value was calculated for each gene by

18   averaging the replicates within each experiment and then computing a mean value over all the

19   experiments in which the corresponding gene was expressed [24].

20

21   *Intergenic controlled coding sequences.*

22   For each gene, the first 100 two-fold and four-fold degenerate codons were used to create the

23   coding sequence for analysis. Transcripts with fewer degenerate codons were excluded. The ICCSs

24   were generated to estimate the influence of the background composition on coding sequence codon

1  bias. Upstream and downstream sequences were extracted from the leading strand and joined

2  together to form a set of concatenated intergenic sequences (CISs) and those shorter than 50

3  nucleotides were excluded at this stage.

4  Four different background controls were used to generate ICCS datasets, beginning with the

5  mononucleotide composition as originally used in the Hershberg and Petrov method [19]. Briefly, a

6  subsequence of 100 consecutive base pairs was randomly selected from the CIS, and the third base

7  of each codon in the coding sequence was replaced with a nucleotide from this CIS subsequence

8  (Figure 1a). This yielded a new dataset called monoICCS. The intergenic dinucleotide composition

9  was used to generate a second class of ICCS (dinuICCS) by choosing a random subsequence of 200

10  consecutive base pairs from the CIS and picking the second base of each coding sequence codon

11  randomly from within that subsequence. The adjacent base was then selected as the third codon

12  positon in the ICCS (Figure 1b). If the CIS was shorter than 100 bp it was excluded from the

13  monoICCS and if it was shorter than 200 bp it was excluded from the dinuICCS. Genes were also

14  excluded from further analysis if any base in the CIS occurred fewer than four times.

15  Intergenic trinucleotide controlled coding sequences (trinuICCS) were produced by randomly

16  picking the first dinucleotide for each coding sequence codon from within the CIS and selecting the

17  adjacent nucleotide as the third base for the ICCS (Figure 1c). Finally, CDCB was estimated by

18  randomly selecting an interrupted dinucleotide comprising the second base of each coding sequence

19  codon and the first base of the subsequent triplet. The intervening nucleotide was then selected as

20  the third codon base in the context-dependent codon bias intergenic controlled coding sequence

21  (cdcbICCS). Genes in the trinuICCS and cdICCS datasets were excluded from further analysis if

22  the corresponding dinucleotide appeared fewer than four times in the CIS.

23  In order to maintain the coding sequence amino acid structure, two-fold degenerate codons ending

24  in A|G were used to create an ICCS ending in A when the corresponding intergenic position was A

25  or T, and otherwise the ICCS ended in G. Similarly, two-fold degenerate codons ending in T|C were

1    used to create an ICCS featuring T at the third codon position when the corresponding CIS

2    nucleotide was A or T, and otherwise the ICCS ended in C.

3

4    *Codon bias measurements*

5    The effective number of codons (Nc)[25] was used to estimate the overall codon bias for each gene in

6    the coding sequence and ICCS datasets. The Nc index generated higher values for less-biased genes

7    (theoretical values between 21 and 64). To explore the contribution of each individual codon,

8    relative synonymous codon usage (RSCU) values were computed for informative two-fold and

9    four-fold degenerate codons (i.e. excluding methionine, tryptophan and stop codons). RSCU values

10   were calculated as the ratio of the observed and expected codon frequencies, i.e. the random use of

11   all codons within a specific degenerate family[26] .

12

13   *Signature of selection*

14   For the 44 codons with two-fold or four-fold degeneracy, differences between the average RSCU

15   values in the CS and ICCS datasets were calculated to highlight overrepresentation and

16   underrepresentation in the coding sequences. For each codon *cod*, the deviation from background

17   was calculated as follows:

18
$$\Delta RSCU^{cod} = \frac{RSCU_{CS}^{cod} - RSCU_{ICCS}^{cod}}{Deg^{cod}}$$

19   where $Deg^{cod}$ is the degeneracy of the codon. The significance of highlighted differences was tested

20   using a paired Wilcoxon test. The same statistical analysis was applied to highlight differences in

21   the effective number of codons between the CS and ICCS datasets.

22   The association between the RSCU values of the CS and ICCS datasets and the gene expression

23   levels was investigated in Arabidopsis and rice by sorting genes on the basis of their expression

1 levels into 20 bins containing the same number of genes. The bin rank was then plotted against the

2 average CS|ICCS RSCU value within the bin using either the intergenic or intron portion for the

3 construction of the ICCS dataset. In the latter case, all introns within the same gene were

4 concatenated. In rice, this analysis was also carried out separately on the HGC and LGC gene sets.

5

6 **Results**

7 **Codon bias in the intergenic corrected coding sequences**

8 Mutational bias is one of the main forces affecting synonymous codon choice in bacteria [19],

9 plants [27] and humans [28]. In theory, if no additional forces act on the coding sequences, the third base

10 of each codon should reflect the background nucleotide frequency in the genome. However, the

11 direction and strength of the codon bias should be investigated in the local genomic context to

12 correct for compositional pattern heterogeneity. Previous studies have revealed the non-random

13 distribution of the four nucleotides in several eukaryotic genomes [29] including the presence of

14 compositionally homogenous isochores in mammals and birds [12,13]. Plant genomes also comprise a

15 mosaic of compositionally homogenous segments although the overall compositional heterogeneity

16 is much less extreme than the sequences found in mammals [12]. The non-random usage of specific

17 dinucleotides [30] and trinucleotides [18] should also be considered during the analysis of local

18 mutational bias to determine the impact on the overall compositional pattern.

19 We used the composition of intergenic regions as a proxy for background bias according to the

20 Hershberg and Petrov method [19]. The portion of codon bias caused by background composition was

21 measured by generating several ICCS datasets taking into account the mononucleotide,

22 dinucleotide, trinucleotide and context-dependent intergenic composition of each plant species. This

23 allowed us to compare the actual structure of the coding sequences with the hypothetical structure

24 based on compositional features of the flanking intergenic regions. We found that codon bias in the

25 coding sequence, and to a lesser extent in the ICCSs, fluctuated along the chromosomes of several

1    species (Figure S1) highlighting the presence of position-specific compositional patterns. We

2    calculated spatial autocorrelations between the chromosome location and the RSCU values in

3    Arabidopsis and rice to determine whether codon bias was conserved among clustered genes.

4    Significant spatial autocorrelation was observed for several codons in the CS datasets but for only a

5    few codons in the ICCS datasets, although the trinuICCS dataset was an exception (Figure S2). This

6    result mirrors the more uniform composition of isochores in plants [29] and contrasts with equivalent

7    results generated by analysing human genes, where the isochore structure is more heterogeneous

8    (Figure S2).

9    Next we investigated the direction of mutational bias by analysing the codon bias among all the

10   ICCS datasets. If the four nucleotides are distributed randomly in the background, then there should

11   be no codon bias in any of the gene sequences. However, we observed the opposite trend, i.e. Nc

12   values compatible with a multilevel background compositional bias. Interestingly, the lowest Nc

13   values were found in the dinuICCS dataset, revealing the strongest bias for all species (Figure 2).

14   Whereas the legumes (*P. vulgaris*, *M. truncatula* and *G. max*) showed the highest ICCS codon bias

15   at all levels, the opposite trend was observed for the monocots (rice, *Z. mays*, *S. bicolor* and *B.*

16   *distachyon*).

17   Nc values provide a snapshot of overall codon bias within genes but do not reveal the specific

18   contributions of each codon. ICCS RSCU values were therefore calculated for all codons with two-

19   fold or four-fold degeneracy (Figure 3). The frequency of specific mononucleotides, dinucleotides

20   and trinucleotides in the intergenic regions was associated, by construction, with the RSCU values

21   of the codons containing them (e.g. overrepresentation of the dinucleotide CA would be associated

22   with the higher RSCU values for codon GCA in the dinuICCS dataset). We observed a significant

23   overrepresentation of codons ending in A|T in the monoICCS dataset, particularly in *V. vinifera*,

24   *Solanum* spp. and the legumes. This highlighted the high degree of A|T-enrichment within the

1    intergenic regions in all the eudicot species we analysed, and the lower degree of enrichment in the

2    monocots and *S. moellendorffii.*

3    A similar picture emerged from the analysis of the dinuICCS and trinuICCS RSCU values although

4    several compositional signatures also emerged. The codons ending in A|T generally showed higher

5    RSCU values at the expense of those ending in G|C, but the extent of the bias was variable. Indeed,

6    codons ending in CG, AG and AC tended to be suppressed more strongly in the dinuICCS dataset,

7    e.g. these dinucleotides were underrepresented in the intergenic regions of all the plants. But codons

8    ending in CA showed the opposite trend. As previously observed for the monoICCS dataset, the

9    codon bias in the monocot dinuICCS and trinuICCS datasets was weaker than the corresponding

10   eudicot datasets. Only marginal differences were observed between the dinuICCS and trinuICCS

11   datasets, indicating that the intergenic trinucleotide bias is predominantly caused by bias in the

12   frequency of the underlying dinucleotides.

13

14   **Codon bias in the coding sequences**

15   If synonymous codon choice solely reflects the distribution of nucleotides in the genomic

16   background, there should be no differences in Nc value between the CS and ICCS datasets.

17   However, our results revealed significant differences (Wilkoxon paired sample test) between the

18   datasets, although the direction and extent of divergence differed among the species we investigated

19   and divergence was particularly evident between eudicots and monocots. Although the coding

20   sequences were more biased than the monoICCS dataset in monocots, only small differences were

21   observed in dicots (Figure 2). Interestingly, the coding sequences were even less biased than the

22   corresponding ICCSs in some species, with the trend most noticeable when comparing the CS and

23   dinuICCS datasets. These data suggest that additional forces shape the coding sequences and

24   oppose the mutational bias, e.g. resulting in G|C enrichment despite background A|T enrichment.

1  Differences in RSCU values between the CS and ICCS datasets allowed us to focus on codons

2  whose frequency cannot be explained by background bias alone. Comparisons between the CS

3  dataset and the four ICCS datasets constructed using alternative approaches led to similar results

4  (Figures 4 and S3). All codons ending in A|T were less frequent in the coding sequence, with GTA

5  suffering the most suppression. However, several codons ending in G|C were overrepresented in the

6  coding sequence, with certain species-dependent exceptions. Codons CCC and GGG were

7  suppressed in many species, together with codons CCG, GCG and ACG. The underrepresentation

8  of GTA was more striking when comparing the CS and monoICCS datasets and less pronounced

9  when comparing the CS and dinuICCS datasets, suggesting that bias against GTA in the coding

10  sequences in part reflects the general suppression of the TA dinucleotide in the intergenic regions

11  (Figure 3). Similarly, the underrepresentation of CCG, GCG and ACG observed when comparing

12  the CS and monoICCS datasets was not so apparent when comparing the CS and dinuICCS

13  datasets, suggesting it reflects the general suppression of the CG dinucleotide in the intergenic

14  regions (Figure 3).

15  Some differences between the CS and the ICCS datasets were also taxon-dependent, e.g. the

16  preference for G|C at the third codon position was more apparent in monocots than eudicots,

17  particularly *M. truncatula, V. vinifera* and *Solanum* spp., where there was little evidence for

18  preference.

19

20  **Correlation between codon bias and gene expression in Arabidopsis and rice**

21  Codon bias in the species we investigated was not fully explained by background compositional

22  differences so we investigated the influence of gene expression on synonymous codon usage in the

23  CS and ICCS datasets. Codons that are translated more rapidly or accurately should be

24  preferentially found in the coding sequences of genes, particularly those expressed at high levels.

25  We chose Arabidopsis and rice as representative dicot and monocot species and assigned genes to

1 20 bins based on expression levels, and then mapped the RSCU values of the CS and the ICCS

2 datasets onto these bins (Figure 5). This strategy should not only reveal selection but also the

3 impact of selection on the coding and noncoding regions. For example, a codon whose frequency is

4 positively associated with expression level may be under selection to optimize translation efficiency

5 and/or accuracy. However, the same positive association should also be observed in the ICCS

6 datasets, e.g. genes with comparable expression levels should share the sequence composition of the

7 intergenic sequences. Both scenarios rely on events that are not mutually exclusive, but this method

8 can highlight the participation of additional forces in shaping the compositional pattern of the

9 coding sequences. Indeed if RSCU values in the CS and ICCS databases differ regardless of the

10 expression level, then forces other than the translational selection are likely to be responsible.

11 We found that two-fold degenerate codons ending in G|C are more frequent in the coding sequences

12 of rice and Arabidopsis, although in the latter case the aspartic acid codon GAC was an exception.

13 Although previous studies have shown that optimal codons tend to end in G|C [31,32] this is not solely

14 dependent on translational selection because a positive association between the RSCU value and

15 expression was evident only at extreme expression levels. There was no association between

16 expression level and the RSCU values of the ICCSs, suggesting that expression-dependent variation

17 in codon bias does not reflect differences in the background composition of the Arabidopsis

18 genome.

19 A more complex picture emerged from the analysis of rice genes. The correlation between RSCU

20 values and gene expression was non-monotonic for codons with two-fold degeneracy, although as

21 in Arabidopsis the codons ending in G|C were still used on average more frequently in the coding

22 sequences. However an initial reduction in the frequency of such codons was complemented by an

23 increase in RSCU values at higher expression levels. If we exclude the effect of the background

24 composition, which again does not vary with the expression level, such non-monotonic trends may

25 reflect the coexistence of contrasting forces whose strength could be dependent on the expression

1    level. Alternatively, the different expression bins in rice may be populated with genes that do not

2    experience the same selective pressure. Indeed, it is well known that monocots have two classes of

3    genes that differ in GC content. For this reason, the above analysis was repeated focusing

4    specifically on LGC and HGC genes. The LGC genes accounted for most of the dataset and

5    mirrored the overall trends discussed above, but the HGC genes showed a general positive

6    association between the expression level and the RSCU values of two-fold degenerate codons

7    ending in G|C.

8    Weaker trends were observed in Arabidopsis when the four-fold degenerate codons were analysed

9    both in terms of deviation from the background and variation with expression. Indeed, the

10   frequency of many codons mirrored the background composition of the entire set of bins, so that the

11   curves generated by the CS and ICCS datasets could be superimposed. However, there was a

12   positive association between the RSCU values and expression levels of codons ending in C and the

13   opposite trend was observed for codons ending in A. Interestingly, the frequency of several codons

14   ending in T was positively associated with the expression level, particularly the alanine codon CGT.

15   In rice, four-fold degenerate codons ending in G|C were always used more frequently in the coding

16   sequence than the ICCS whereas the opposite trend was observed for codons ending in A|T.

17   Furthermore, codons ending in G|C were generally used less frequently in strongly expressed genes,

18   whereas there was a positive association between codons ending in T and expression level. As

19   previously reported in monocots, such trends are often non-monotonic, either disappearing or

20   changing direction when HGC and LGC genes are analysed separately (Figure S4).

21   The positive association between the RSCU value and expression level of codons ending in T is

22   reminiscent of transcription-associated mutational bias (TAMB), a well-known repair system that

23   increases bias towards G|T rather than C|A in the pre-mRNA sequence [33]. For this reason, we

24   repeated our analysis by constructing a new set of ICCSs using intron sequences as a proxy for the

25   background composition. If codons ending in T become more frequent in strongly expressed genes

1    due to TAMB, the same trend should be observed in intron-corrected coding sequences. Our results

2    revealed no such association (Figure S6), indicating that TAMB is not responsible for the positive

3    association between gene expression level and the frequency of codons ending in T.

4

5    **Discussion**

6    **Background composition**

7    Several ICCS datasets were generated to investigate whether codon bias in the coding sequence

8    reflected the genomic background nucleotide composition of the plants included in this study. Four

9    ICCS datasets constructed using different strategies were analysed to determine the effective

10    number of codons and relative synonymous codon usage. A strong dinucleotide compositional bias

11    was evident in all the species, as shown by the higher levels of codon bias (i.e. lower Nc values) in

12    the dinuICCS dataset compared to the others (Figure 2). The RSCU values of the monoICCS

13    dataset revealed the overrepresentation of A|T compared to G|C although to a different extent in

14    each species. The legumes and solanaceous species showed the greatest difference in the

15    representation of A|T and G|C, whereas there was a weaker distinction in the monocot species and

16    the brassicas showed intermediate values. G|C to A|T mutations may be more frequent in

17    Arabidopsis, leading to AT enrichment in the genome. The deamination of methylated cytosine

18    residues at CG dinucleotide motifs, and the UV-induced mutagenesis of dipyrimidines (CC and TC)

19    may explain this phenomenon[34]. Our data support this model in all the plants included in this

20    investigation. We found that codons ending in CG were among the most underrepresented in the

21    dinuICCS dataset, whereas codons ending in CA (the reverse complement of the C deamination

22    product in CG) were among the most overrepresented, as recently reported in Arabidopsis and rice

23    [30]. Codons ending in CC and TC were also suppressed in the dinuICCS dataset, confirming the

24    underrepresentation of these dinucleotides in the non-coding sequences. Several factors may

25    account for the observed differences in the background nucleotide composition of plants. In

1    monocots, the vertical leaf orientation, protective basal sheath, and concealed apical meristem make

2    the interception of solar radiation less efficient [35] thus reducing the prevalence of AT-enrichment

3    induced by UV light. Furthermore, A|T pairs contain seven nitrogen atoms compared to the eight

4    present in G|C pairs, which may drive A|T-enrichment in non-cultivated plants [36]. However, the

5    cluster analysis of ICCS RSCU values revealed a pattern that is consistent with the genomic

6    composition of plants solely depending on their phylogenesis (Figure S5).

7

8    **Differences between the CS and ICCS datasets**

9    Paired Wilcoxon tests generally revealed significant differences in codon usage between the CS and

10   ICCS datasets in terms of both the Nc and RSCU values. The observed differences were similar in

11   magnitude when comparing the CS dataset with all four methods for the construction of ICCS

12   datasets, despite the clear dinucleotide signature of the background composition. This validates the

13   Hershberg and Petrov method for the identification of optimal codons even though it only considers

14   the mononucleotide composition of the background sequences. Indeed, the optimal codon datasets

15   calculated using this method could be almost precisely superimposed over datasets generated using

16   the other three methods (Figure S7).

17   The CS|ICCS comparisons suggested there was enrichment for codons entirely composed of G and

18   C (hereafter described as GC3 codons) particularly in monocots, although among codons ending in

19   G|C the frequency of the complementary codons GGG and CCC was close to the genomic

20   background. This suggests there may be selection against codons that promote the formation of

21   complex mRNA tertiary structures, and the prevalence of this phenomenon in monocots with their

22   higher overall G|C content may emphasize such an underlying mechanism.

23   Codons ACG, CCG and GCG were marginally overrepresented in the coding sequences of some

24   species and underrepresented in others. This supports the observation that the CG dinucleotide is

25   suppressed, in part due to deleterious methylation/deamination events. However, this is actually a

genomic tendency and the effect is not seen when comparing the CS and dinuICCS datasets (Figure

4).

Finally, the general suppression of codon GTA was one of the most conserved features among the

plant species we investigated. This is not surprising because the codon ends with dinucleotide TA,

which is known to be suppressed in the coding sequences of several plant species [18,30] possibly to

discourage insertion events that target the TA dinucleotide [37], to reduce the likelihood of mutations

leading to stop codons, and to prevent attacks by TA-specific RNases [38]. The less striking

divergence between the CS and dinuICCS datasets in terms of GTA preference suggests that TA

suppression is also a general genomic signature (Figure 3).


**Mutational bias in coding sequences**

The background nucleotide pattern alone cannot explain the observed codon bias in the coding

sequences, so additional forces must be involved (although the signature of mutational pressure may

still be observed at sites that are more loosely constrained in some plant species). Indeed, modestly

frequent amino acids (or underrepresented amino acids in the proteins encoded by strongly

expressed genes) may be under weaker selection, and for this reason their codons may better

tolerate changes driven by mutational bias. For example, histidine and cysteine are among the less

abundant amino acids in the Arabidopsis proteome and there is a significant negative correlation

between their frequency and the expression level of the corresponding genes (Table S1). The

frequency of codons CAC (His) and TGC (Cys) in Arabidopsis is similar in the CS and ICCS

datasets (Figure 5). The codons GAC (Asp), GGC (Gly) and GCC (Ala) also revealed biases that

mirrored the background composition. Interestingly, these codons feature the generic sequence

CNG which may represent the core of the primitive genetic code [39]. Such a trend may therefore

reflect the more prolonged effect of mutational bias on the most ancient codons. In contrast, there

1    appeared to be little mutational bias in monocots, where other factors increase the frequency of CG3

2    codons.

3

4    **Evidence for biased gene conversion in plants**

5    The increasing G|C content of plant genomes may have been driven by G|C-biased gene conversion

6    (gBGC) during double-strand break repair followed by a recombination event. This phenomenon is

7    accompanied by the correction of eventual mismatches between two paired DNA strands featuring

8    high sequence similarity, with such a correction being biased toward the placement of either G or C.

9    By definition, gBGC is associated with the recombination rate along the genome and has

10   contributed to isochore formation in mammals and birds. For this reason, a positive correlation

11   between the G|C content and the recombination rate provides evidence that supports gBGC, and

12   such an association has been observed in grasses but not Arabidopsis. Nevertheless, the association

13   between recombination rate and G|C content may be perturbed (or possibly lost) when the

14   evolutionary history of a species features a large number of genomic rearrangements, as is the case

15   for Arabidopsis and the eudicots in general [40]. Although gBGC should not be restricted to genes,

16   and differences in codon bias between the CS and ICCS datasets would therefore not highlight

17   gBGC events, mutations caused by gBGC are more likely to be fixed in the coding sequence

18   because codons ending in G|C are known to mirror the most abundant tRNAs in most plant species.

19   There is also evidence that recombination in plant genomes occurs mainly in genes [41].

20   Our data provided several lines of evidence supporting the occurrence of gBGC in plants. First,

21   higher standardized ΔRSCU values were observed for codons with two-fold degeneracy ending in

22   G|C compared to those with four-fold degeneracy, which indicates the occurrence of gBGC events

23   featuring a more diluted effect on codons with four-fold degeneracy. Moreover, if gBGC rather than

24   selection is considered to be the main source of GC3 enrichment, such an effect should be evident

25   in all the transcripts regardless of the expression level. As shown in Figure 5, this was the case for

1    the majority of the codons and expression bins in both Arabidopsis and rice, with the latter featuring

2    wider divergence from the background in accordance with previous reports of gBGC in grasses [42].

3

4    **Translational selection and mRNA stability**

5    Genes were assigned to 20 expression bins whose average RSCU values were plotted for the CS

6    dataset and all four ICCS datasets for Arabidopsis and rice, representing the eudicots and monocots

7    respectively (Figure 5). ICCS RSCU values were not associated with the expression level

8    suggesting that co-expression clustering within the genome of these two species cannot be detected

9    using this analytical approach.

10   In Arabidopsis, we observed a positive correlation between expression level and the frequency of

11   optimal codons [31,32] ending in G|C mainly for genes assigned to the last few expression bins,

12   underlining the marginal effect of translational selection in this species. A more cumbersome

13   scenario emerged in rice, where non-monotonic trends were observed for codons ending in G|C, i.e.

14   a reduction in frequency in the first expression bins changing to an increase in frequency in the last

15   few. This behaviour may be typical of G|C-rich monocot genomes and may indicate a compromise

16   between the advantage of using optimal codons and the avoidance of tightly-packed mRNA tertiary

17   structures. Indeed a significant negative correlation was found between the expression level and

18   both the coding sequence length ($r = –0.05$, $p < 0.0001$) and the G|C content ($r = –0.15$, $p < 0.0001$)

19   of rice genes. Taken together, these data suggest that weakly expressed genes are longer and have a

20   lower content of CG3 codons than strongly expressed genes. Longer transcripts are more likely to

21   form strongly-packed mRNA tertiary structures if they are enriched in optimal codons ending in

22   G|C, so the accumulation of such codons may be counterselected in genes expressed at low or

23   moderate levels. Such a trend would diminish in shorter genes expressed at high levels. This

24   hypothesis was confirmed by analysing HGC and LGC rice genes separately. We found that short

1    HGC genes were generally characterized by a monotonic positive association between the

2    expression level and the frequency of optimal codons ending in G|C.

3

4    **Additional factors**

5    Factors other than mutational bias, selection and gBGC that shape codon bias include TAMB, a

6    well-known DNA repair system that influences the composition of primary transcripts (both exons

7    and introns) by increasing the prevalence of G|T over C|A. Despite this general GC3 enrichment,

8    we observed the enrichment of four-fold degenerate codons ending in T, suggesting that TAMB

9    may affect codon bias in Arabidopsis and rice, so we repeated our analysis using introns rather than

10   intergenic DNA as an index of background composition. As shown in Figure S7, we found no

11   differences in the frequency of codons ending in T when we used introns rather than intergenic

12   DNA, suggesting the overrepresentation of such codons is a general genomic trend that probably

13   reflects non-localized effects such as the preference for nitrogen saving base pairs in non-cultivated

14   species (the A|T and A|U base pairs contain only seven nitrogen atoms, compared to eight in the

15   G|C base pair). In contrast, crop species are provided with ample nitrogen and phosphorous which

16   removes such constraints. However, the convergence of these trends between the brassicas and the

17   nitrogen-fixing legumes, and between the monocots and the ancestral species *S. moellendorffii*

18   (Figure S5) suggests that the overrepresentation of codons ending in T may reflect the phylogenesis

19   of the species rather than the efficiency of nitrogen utilization.

20

1     **Figure 1:** Method used for the construction of the mononucleotide ICCS (monoICCS), dinucleotide ICCS (dinuICCS)

2     and trinucleotide ICCS (trinuICCS) datasets.

3     **Figure 2**: Heat map showing the average Nc values for the CS and ICCS datasets.

4     **Figure 3:** RSCU values calculated for the monoICCS, dinuICCS and trinuICCS datasets (only codons ending in G|C

5     with two-fold degeneracy are shown).

6     **Figure 4**: Calculation of ΔRSCU values as the standardized differences between the CS RSCU and (a) the monoICCS

7     and (b) the dinuICCS RSCUs. Black crosses in white squares show insignificant differences between the datasets.

8     **Figure 5:** RSCU values of (a) Arabidopsis and (b) rice genes after splitting the datasets into 20 expression bins.

9     Intergenic sequences were used for the construction of the ICCS datasets. Colour codes: blue = CS, red = monoICCS,

10     green = dinuICCS, yellow = trinuICCS, brown = cdcbICCS.

11

12

1      **References**

2    1.    Chamary, J. V, Parmley, J. L., and Hurst, L. D. 2006, Hearing silence: non-neutral evolution
3          at synonymous sites in mammals. *Nat. Rev. Genet.*, **7**, 98–108.

4    2.    Duret, L. 2002, Evolution of synonymous codon usage in metazoans. *Curr. Opin. Genet.*
5          *Dev.*, **12**, 640–9.

6    3.    Plotkin, J. B., Dushoff, J., Desai, M. M., and Fraser, H. B. 2006, Codon usage and selection
7          on proteins. *J. Mol. Evol.*, **63**, 635–53.

8    4.    Plotkin, J. B., and Kudla, G. 2011, Synonymous but not the same: the causes and
9          consequences of codon bias. *Nat. Rev. Genet.*, **12**, 32–42.

10   5.    Ikemura, T. 1981, Correlation between the abundance of Escherichia coli transfer RNAs and
11         the occurrence of the respective codons in its protein genes: a proposal for a synonymous
12         codon choice that is optimal for the E. coli translational system. *J. Mol. Biol.*, **151**, 389–409.

13   6.    Akashi, H. 1994, Synonymous codon usage in Drosophila melanogaster: natural selection
14         and translational accuracy. *Genetics*, **136**, 927–35.

15   7.    Bulmer, M. 1991, The Selection-Mutation-Drift Theory of Synonymous Codon Usage.
16         *Genetics*, **129**, 897–907.

17   8.    Sharp, P. M., Cowe, E., Higgins, D. G., Shields, D. C., Wolfe, K. H., and Wright, F. 1988,
18         Codon usage patterns in Escherichia coli, Bacillus subtilis, Saccharomyces cerevisiae,
19         Schizosaccharomyces pombe, Drosophila melanogaster and Homo sapiens; a review of the
20         considerable within-species diversity. *Nucleic Acids Res.*, **16**, 8207–11.

21   9.    Najafabadi, H. S., Goodarzi, H., and Salavati, R. 2009, Universal function-specificity of
22         codon usage. *Nucleic Acids Res.*, **37**, 7014–23.

23   10.   Camiolo, S., Farina, L., and Porceddu, A. 2012, The relation of codon bias to tissue-specific
24         gene expression in Arabidopsis thaliana. *Genetics*, **192**, 641–9.

25   11.   Wong, G. K.-S., Wang, J., Tao, L., et al. 2002, Compositional gradients in Gramineae genes.
26         *Genome Res.*, **12**, 851–6.

27   12.   Bernardi, G., Olofsson, B., Filipski, J., et al. 1985, The mosaic genome of warm-blooded
28         vertebrates. *Science*, **228**, 953–8.

29   13.   Carels, N., and Bernardi, G. 2000, Two classes of genes in plants. *Genetics*, **154**, 1819–25.

30   14.   Liu, Q. 2012, Mutational bias and translational selection shaping the codon usage pattern of
31         tissue-specific genes in rice. Robinson-Rechavi, M., (ed.), . *PLoS One*, **7**, e48295.

32   15.   Berg, O. G., and Silva, P. J. 1997, Codon bias in Escherichia coli: the influence of codon
33         context on mutation and selection. *Nucleic Acids Res.*, **25**, 1397–404.

16. De Amicis, F., and Marchetti, S. 2000, Intercodon dinucleotides affect codon choice in plant genes. *Nucleic Acids Res.*, **28**, 3339–45.

17. Karlin, S., and Mrázek, J. 1996, What drives codon choices in human genes? *J. Mol. Biol.*, **262**, 459–72.

18. Porceddu, A., and Camiolo, S. 2011, Spatial analyses of mono, di and trinucleotide trends in plant genes. *PLoS One*, **6**, e22855.

19. Hershberg, R., and Petrov, D. A. 2009, General rules for optimal codon choice. *PLoS Genet.*, **5**, e1000556.

20. Morton, B. R., and Wright, S. I. 2007, Selective constraints on codon usage of nuclear genes from Arabidopsis thaliana. *Mol. Biol. Evol.*, **24**, 122–9.

21. Goodstein, D. M., Shu, S., Howson, R., et al. 2012, Phytozome: a comparative platform for green plant genomics. *Nucleic Acids Res.*, **40**, D1178–86.

22. Camiolo, S., and Porceddu, A. 2013, gff2sequence, a new user friendly tool for the generation of genomic sequences. *BioData Min.*, **6**, 15.

23. Dash, S., Van Hemert, J., Hong, L., Wise, R. P., and Dickerson, J. A. 2012, PLEXdb: gene expression resources for plants and plant pathogens. *Nucleic Acids Res.*, **40**, D1194–201.

24. Camiolo, S., Rau, D., and Porceddu, A. 2009, Mutational biases and selective forces shaping the structure of Arabidopsis genes. *PLoS One*, **4**, e6356.

25. Wright, F. 1990, The "effective number of codons" used in a gene. *Gene*, **87**, 23–9.

26. Sharp, P. M., and Li, W. H. 1986, An evolutionary perspective on synonymous codon usage in unicellular organisms. *J. Mol. Evol.*, **24**, 28–38.

27. Palidwor, G. A., Perkins, T. J., and Xia, X. 2010, A general model of codon bias due to GC mutational bias. *PLoS One*, **5**, e13431.

28. Sémon, M., Lobry, J. R., and Duret, L. 2006, No evidence for tissue-specific adaptation of synonymous codon usage in humans. *Mol. Biol. Evol.*, **23**, 523–9.

29. Nekrutenko, A., and Li, W. H. 2000, Assessment of compositional heterogeneity within and between eukaryotic genomes. *Genome Res.*, **10**, 1986–95.

30. Karlin, S. 1998, Global dinucleotide signatures and analysis of genomic heterogeneity. *Curr. Opin. Microbiol.*, **1**, 598–610.

31. Wright, S. I., Yau, C. B. K., Looseley, M., and Meyers, B. C. 2004, Effects of gene expression on molecular evolution in Arabidopsis thaliana and Arabidopsis lyrata. *Mol. Biol. Evol.*, **21**, 1719–26.

32. Porceddu, A., Zenoni, S., and Camiolo, S. 2013, The signatures of selection for translational accuracy in plant genes. *Genome Biol. Evol.*, **5**, 1117–26.

1   33.   Green, P., Ewing, B., Miller, W., Thomas, P. J., and Green, E. D. 2003, Transcription-
2         associated mutational asymmetry in mammalian evolution. *Nat. Genet.*, **33**, 514–7.

3   34.   Ossowski, S., Schneeberger, K., Lucas-Lledó, J. I., et al. 2010, The rate and molecular
4         spectrum of spontaneous mutations in Arabidopsis thaliana. *Science*, **327**, 92–4.

5   35.   Cline, M. G., and Salisbury, F. B. 1966, Effects of ultra-violet alone and simulated solar
6         ultra-violet radiation on the leaves of higher plants. *Nature*, **211**, 484–6.

7   36.   Günther, T., Lampei, C., and Schmid, K. J. 2013, Mutational bias and gene conversion affect
8         the intraspecific nitrogen stoichiometry of the Arabidopsis thaliana transcriptome. *Mol. Biol.*
9         *Evol.*, **30**, 561–8.

10  37.   Mashkova, T. D., Oparina, N. Y., Lacroix, M. H., et al. 2001, Structural rearrangements and
11        insertions of dispersed elements in pericentromeric alpha satellites occur preferably at
12        kinkable DNA sites. *J. Mol. Biol.*, **305**, 33–48.

13  38.   Beutler, E., Gelbart, T., Han, J. H., Koziol, J. A., and Beutler, B. 1989, Evolution of the
14        genome and the genetic code: selection at the dinucleotide level by methylation and
15        polyribonucleotide cleavage. *Proc. Natl. Acad. Sci. U. S. A.*, **86**, 192–6.

16  39.   Ikehara, K., Omori, Y., Arai, R., and Hirose, A. 2002, A novel theory on the origin of the
17        genetic code: a GNC-SNS hypothesis. *J. Mol. Evol.*, **54**, 530–8.

18  40.   Salse, J., Abrouk, M., Bolot, S., et al. 2009, Reconstruction of monocotelydoneous proto-
19        chromosomes reveals faster evolution in plants than in animals. *Proc. Natl. Acad. Sci. U. S.*
20        *A.*, **106**, 14908–13.

21  41.   Rafalski, A., and Morgante, M. 2004, Corn and humans: recombination and linkage
22        disequilibrium in two genomes of similar size. *Trends Genet.*, **20**, 103–11.

23  42.   Serres-Giardi, L., Belkhir, K., David, J., and Glémin, S. 2012, Patterns and evolution of
24        nucleotide landscapes in seed plants. *Plant Cell*, **24**, 1379–97.
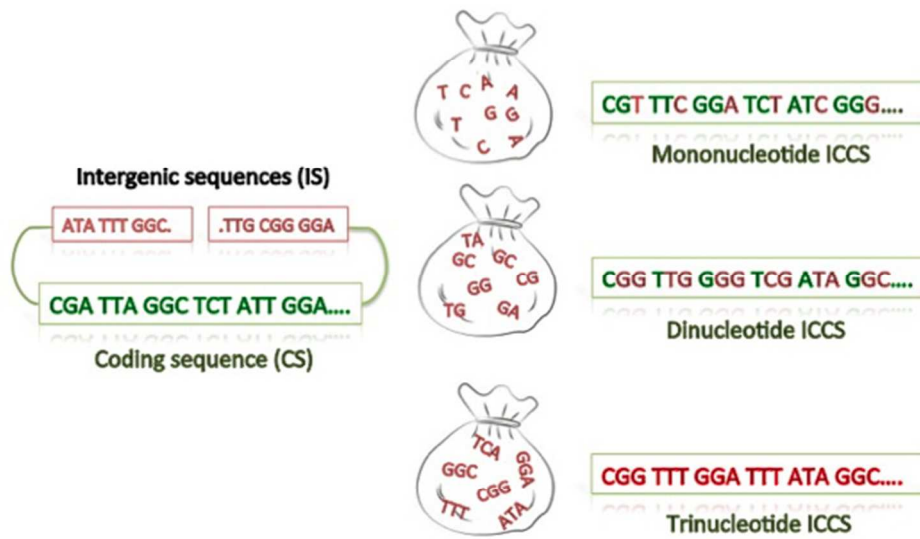
25

Figure 1: Method used for the construction of the mononucleotide ICCS (monoICCS), dinucleotide ICCS (dinuICCS) and trinucleotide ICCS (trinuICCS) datasets.
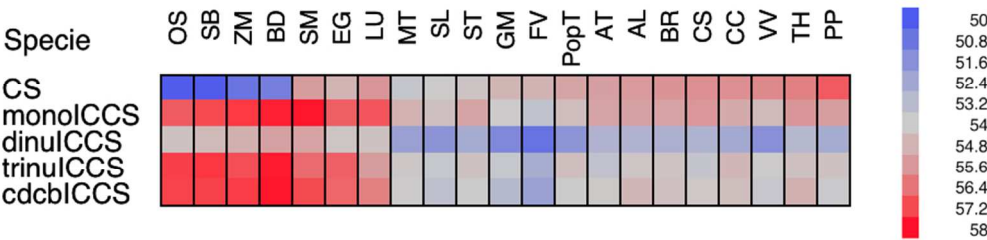243x182mm (72 x 72 DPI)

Figure 2: Heat map showing the average Nc values for the CS and ICCS datasets.
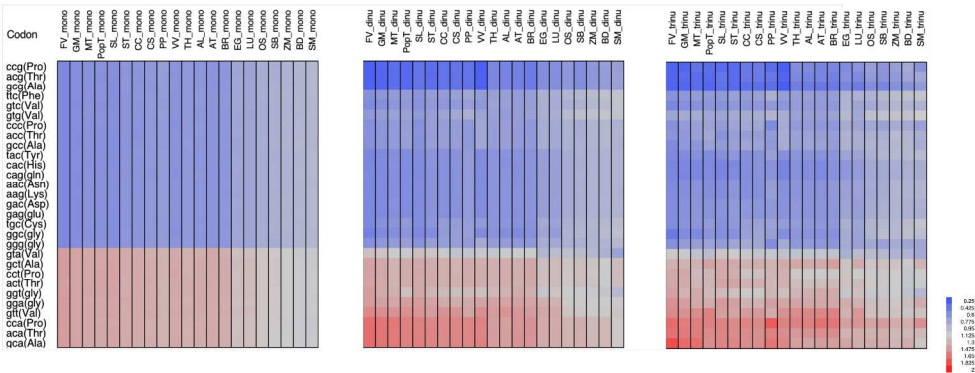323x82mm (72 x 72 DPI)

Figure 3: RSCU values calculated for the monoICCS, dinuICCS and trinuICCS datasets (only codons ending in G|C with two-fold degeneracy are shown).
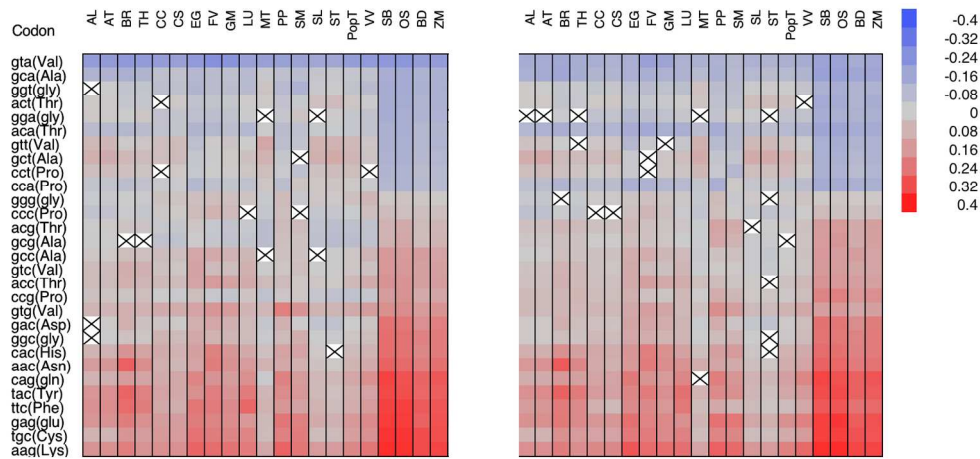799x310mm (150 x 150 DPI)

Figure 4: Calculation of ΔRSCU values as the standardized differences between the CS RSCU and (a) the monoICCS and (b) the dinuICCS RSCUs. Black crosses in white squares show insignificant differences between the datasets.
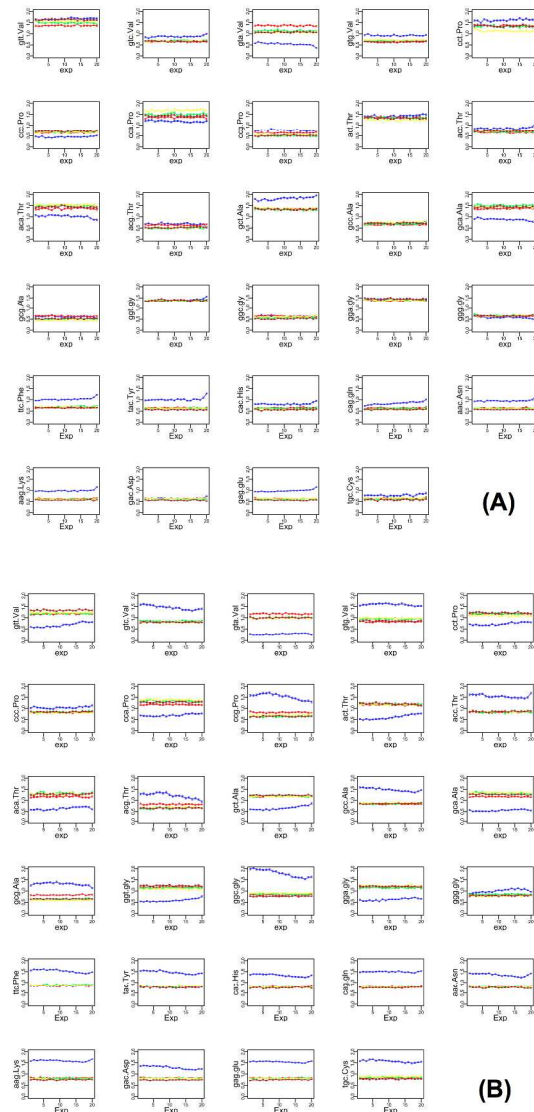610x278mm (72 x 72 DPI)

Figure 5: RSCU values of (a) Arabidopsis and (b) rice genes after splitting the datasets into 20 expression bins. Intergenic sequences were used for the construction of the ICCS datasets. Colour codes: blue = CS, red = monoICCS, green = dinuICCS, yellow = trinuICCS, brown = cdcbICCS.
199x399mm (300 x 300 DPI)