# SUPPLEMENTARY MATERIAL

# Different sequence signatures in the upstream regions of plant and animal tRNA genes shape distinct modes of regulation

Gong Zhang[1], Radoslaw Lukoszek[1,2], Bernd Mueller-Roeber[2], Zoya Ignatova[1]

[1]Biochemistry, Institute of Biochemistry and Biology, University of Potsdam, Potsdam, Germany

[2]Molecular Biology, Institute of Biochemistry and Biology, University of Potsdam, Potsdam, Germany

**Corresponding authors:** Zoya Ignatova, Biochemistry, University of Potsdam, Karl-Liebknecht-Str. 24-25, 14467 Potsdam, Germany; Tel: +49 331 977 5130, Fax: +49 331 977 5128; e-mail: ignatova@uni-potsdam.de

Bernd Mueller-Roeber, Molecular Biology, University of Potsdam, Karl-Liebknecht-Str. 24-25, 14467 Potsdam, Germany; Tel: +49 331 977 2810, Fax: +49 331 977 2512; e-mail: bmr@uni-potsdam.de

**Table S1**. Animal tRNAs encoding the same amino acid share common motif in their upstream regions. tRNAs from each organism were grouped dependent on the amino acid they decode and the tRNAs with at least one conserved motif in their upstream region are marked with +. Upstream segments without any conserved motif are marked with −.

| Amino acid | *Caenorhabditis elegans* | *Drosophila melanogaster* | *Anopheles gambiae* | *Gallus gallus* | *Mus musculus* | *Rattus norvegicus* | *Homo sapiens* |
|---|---|---|---|---|---|---|---|
| Ala | + | − | − | + | + | + | + |
| Arg | + | + | + | − | − | − | − |
| Asn | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. | + |
| Asp | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. |
| Cys | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. |
| Gln | + | − | + | + | + | + | + |
| Glu | + | + | + | + | + | + | + |
| Gly | + | + | - | + | + | − | + |
| His | + | s.a. | s.a. | + | s.a. | s.a. | s.a. |
| Ile | + | − | + | + | + | + | + |
| Leu | + | − | − | − | + | + | + |
| Lys | − | + | + | − | + | − | − |
| Met | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. |
| Phe | s.a. | s.a. | s.a. | + | s.a. | + | s.a. |
| Pro | − | − | − | + | − | − | + |
| Ser | − | + | − | + | − | + | − |
| Thr | − | − | − | − | − | − | − |
| Trp | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. |
| Tyr | s.a. | s.a. | s.a. | s.a. | s.a. | s.a. | + |
| Val | − | − | − | − | + | + | + |

s.a. denotes that all tRNAs decoding this amino acid bear the same anticodon.

**Table S2.** Human Brf1 and Bdp1 splicing variants. Data were retrieved from the ECgene alternative splicing database [http://genome.ewha.ac.kr/ECgene/].

| Gene | ID | mRNA length (bp) | CDS length (bp) |
|---|---|---|---|
| **Brf1** | H14C11627.1 | 3511 | 579 |
| | H14C11627.2 | 2874 | 609 |
| | H14C11627.3 | 913 | 402 |
| | H14C11627.4 | 496 | 259 |
| | H14C11627.5 | 967 | 483 |
| | H14C11627.6 | 1673 | 402 |
| | H14C11627.7 | 4169 | 1419 |
| | H14C11627.8 | 2861 | 1317 |
| | H14C11627.9 | 3538 | 2031 * |
| | H14C11627.10 | 3603 | 1950 |
| | H14C11627.11 | 3948 | 2031 * |
| | H14C11627.12 | 1906 | 450 |
| | H14C11627.13 | 1492 | 339 |
| | H14C11627.14 | 574 | 297 |
| | H14C11627.15 | 1452 | 624 |
| **Bdp1** | H5C7152.1 | 740 | 476 |
| | H5C7152.2 | 4447 | 4116 |
| | H5C7152.3 | 7265 | 6762 |
| | H5C7152.4 | 10600 | 7872 |
| | H5C7152.5 | 5031 | 1392 |
| | H5C7152.6 | 5013 | 1797 |

* mRNA splice variants H14C11627.9 and H14C11627.11 are identical in protein sequence but differ in the non-coding flanking regions.

A

mismatch = 0

**Figure S1.** Best "word" frequency search within the 2000 nucleotide-long upstream regions of the mature tRNA sequences revealed no additional significant motif that could serve as a putative recognition sites for TFIIIB. The occurrence of the most frequent six-letter "word" at each position (x-axis) is plotted on the y-axis allowing a search with none (A) or one mismatch (B). Position 0 is determined by the first 5'-nucleotide of the mature tRNA. The significant peaks within the first 100 nucleotides upstream of the 5'-initial nucleotide of the mature tRNA are identical to those identified in Figure 1.

**Figure S2.** The best frequency "word" search in two yeast genomes, *Saccharomyces cerevisiae* (SacCer1, Oct. 2003) and *Schizosaccharomyces pombe* (build 1.1., Apr. 2007), revealed different conserved motifs. The search was performed in the near upstream region (up to -100 nucleotides) of the tRNAs genes allowing none (A) or one mismatch (B) (for details see the legend to Figure 1). Position 0 is determined by the first nucleotide of the mature tRNA. The dashed line indicates the average frequency value obtained from randomized sequences. Maximum and minimum frequency values determine the boundaries of the shadowed area and were also derived from the randomized dataset.

**Figure S3**. Motifs identified within the 100 nucleotides upstream of the 5'-start of the mature tRNAs. The top five most significant motifs for each genome are shown. TA-rich motifs are detected in every plant genome, but are completely lacking in the animal genomes. The motif search was performed using MEME version 4.4.0 (http://meme.sdsc.edu/meme4_4_0/cgi-bin/meme.cgi).

A

For anticodon ATA, no motif detected.

B

Cys

| Label | Combined p-value |
|---|---|
| GCA, chr17, #30 | 1.48e-09 |
| GCA, chr7, #19 | 1.48e-03 |
| GCA, chr7, #8 | 2.76e-04 |
| GCA, chr7, #14 | 3.49e-05 |
| GCA, chr7, #11 | 1.92e-04 |
| GCA, chr7, #23 | 6.50e-04 |
| GCA, chr7, #16 | 3.71e-03 |
| GCA, chr17, #29 | 3.28e-02 |
| GCA, chr7, #25 | 8.79e-04 |
| GCA, chr7, #22 | 2.07e-04 |
| GCA, chr3, #6 | 4.19e-05 |
| GCA, chr7, #13 | 5.11e-02 |
| GCA, chr7, #6 | 2.91e-04 |
| GCA, chr14, #8 | 3.81e-09 |
| GCA, chr1, #127 | 1.48e-09 |
| GCA, chr3, #7 | 2.07e-02 |
| GCA, chr15, #3 | 1.48e-09 |

His

| Label | Combined p-value |
|---|---|
| GTG, chr3, #4 | 2.44e-03 |
| GTG, chr1, #106 | 2.61e-10 |
| GTG, chr1, #111 | 4.76e-10 |
| GTG, chr1, #118 | 3.34e-06 |
| GTG, chr1, #16 | 4.76e-10 |
| GTG, chr1, #21 | 4.76e-10 |
| GTG, chr15, #1 | 4.37e-04 |
| GTG, chr15, #8 | 2.30e-02 |
| GTG, chr15, #9 | 2.80e-05 |
| GTG, chr6, #33 | 4.30e-05 |
| GTG, chr9, #7 | 5.13e-02 |

Met

| Label | Combined p-value |
|---|---|
| CAT, chr9, #1 | 1.24e-01 |
| CAT, chr6, #92 | 6.04e-05 |
| CAT, chr17, #20 | 3.99e-03 |
| CAT, chr6, #129 | 3.14e-02 |
| CAT, chr6, #142 | 1.82e-02 |
| CAT, chr6, #169 | 5.33e-02 |
| CAT, chr6, #171 | 2.64e-02 |
| CAT, chr6, #21 | 5.04e-05 |
| CAT, chr6, #162 | 5.09e-09 |
| CAT, chr6, #164 | 4.49e-10 |
| CAT, chr6, #27 | 4.49e-10 |
| CAT, chr6, #75 | 5.15e-02 |
| CAT, chr6, #97 | 1.97e-03 |
| CAT, chr16, #20 | 3.75e-04 |

11

**Figure S4.** Majority of the human tRNA genes with the same anticodon bear conserved motifs in their upstream regions as identified by MEME 4.4. tRNA families are grouped based on the amino acid they bear and subdivided into amino acids encoded by more than one codon (A) and by one single codon (B). In panel A, the square brackets on the left group the tRNAs with the same anticodon. Each tRNA is specified by its anticodon, chromosome location and serial chromosome number (data were retrieved from the genomic tRNA database). The horizontal lines represent the 100-nucleotide-long region upstream of the 5'-end of the mature tRNA genes; position 0 determines the first nucleotide of the mature tRNA. The positions of the three most significant motifs are colored (cyan, blue, red) and the corresponding sequence logos are presented on the right side of each group. The combined *p*-value is the product of the *p*-values of all motifs detected within the upstream sequence; they are inversely proportional to the significance of the motifs. Motifs with $p<0.05$ were considered as significant.

**Figure S5**. Distribution of the Bdp1 and Brf1 splice isoforms among different tissues. The read hits for specific splice variant were normalized to the total read count of the sequencing data sets (for details see the description in the Results section). Numbering of the splicing variants is according to the ECgene alternative splicing database (Table S2). Some other theoretically possible splicing variants of Bdp1 or Brf1 (Table S2) were also detected in the deep-sequencing datasets however at a very low level (less than 3 reads per 100 million reads) and were therefore not considered in the analysis.

```
1   1  MVW--CNHCVKNVPGIR-PY-DGALACNLCGRILENFHFSTEVTFVKNAAGQSQASGNIVRSVQSGIT------------  64
2   1  MVW--CKHCGKNVPGIR-PY-DAALSCDLCGRILENFNFSTEVTFVKNAAGQSQASGNILKSVQSGMS------------  64
3   1  MFV--CKNCHGTEFERDLSNANNDLVCKACGVVSEDNPIVSEVTFGETSAGAAVVQGSFI---GAGQSHAAFGG---SSA  72
4   1  MTGRVCRGCGGTDIELD-A-ARGDAVCTACGSVLEDNIIVSEVQFVESSGGGSSAVGQFVSLDGAGKTPTLGGGFHVNLG  78

1  65  -SSRERRFRIARDELMNLKDALGIGDERDDVIVIAAKFFEMAVEQNFTKGRRTELVQASCLYLTCRELNIALLLIDFSSY  143
2  65  -SSRERIIRKATDELMNLRDALGIGDDRDDVIVMASNFFRIALDHNFTKGRSKELVFSSCLYLTCRQFKLAVLLIDFSSY  143
3  73  LESREATLNNARRKLRAVSYALHIPE---YITDAAFQWYKLALANNFVQGRRSQNVIASCLYVACRKEKTHHMLIDFSSR  149
4  79  KESRAQTLQNGRRHIHHLGNQLQLNQ---HCLDTAFNFFKMAVSRHLTRGRKMAHVIAACLYLVCRTEGTPHMLLDLSDL  155

1 144  LRVSVYELGSVYLQLCEMLYLVENRNYEKLVDPSIFMDRFSNSLLKGKNNKDVVATARDIIASMKRDWIQTGRKPSGICG  223
2 144  LRVSVYDLGSVYLQLCDMLYITENHNYEKLVDPSIFIPRFSNMLLKGAHNNKLVLTATHIIASMKRDWMQTGRKPSGICG  223
3 150  LQVSVYSIGATFLKMVKKLHITEL----PLADPSLFIQHFAEKLDLADKKIKVVKDAVKLAQRMSKDWMFEGRRPAGIAG  225
4 156  LQVNVYVLGKTFLLLARELCINA-----PAIDPCLYIPRFAHLLEFGEKNHEVSMTALRLLQRMKRDWMHTGRRPSGLCG  230

1 224  AALYTAALSHGIKCSKTDIVNIVHICEATLTKRLIEFGDTDSGNLNVNELRERESHK-----RSFTM----KPTSNKEAV  294
2 224  AALYTAALSHGIKCSKTDIVNIVHICEATLTKRLIEFGDTEAASLTADELSKTEREK-----ETAALRSKRKPNFYKEGV  298
3 226  ACILLACRMNNLRRTHTEIVAVSHVAEETLQQRLNEFKNTKAAKLSVQKFRENDVEDGEARPPSFVK-NRKKERKIKDSL  304
4 231  AALLVAARMHDFRRTVKEVISVVKVCESTLRKRLTEFDTPTSQLTIDEFMKIDLEE-ECDPPSYTA-GQRKLR------  302

1 295  ------------LCMHQ-----------DSKPFGYGLCEDCYK-----DFINVSGGL-------VGGSNPPAFQRA----  335
2 299  -----------VLCMHQ-----------DCKPVDYGLCESCYD-----EFMTVSGGL-------EGGSDPPAFQRA----  340
3 305  DKEEMFQTSEEALNKNPILTQVLGEQELSSKEVLFYLKQFSERRARVVERIKATNGIDGENIYHEGSENETRKRKLSEVS  384
4 303  MKQ-----LEQVLSKK--LEEVEGEISSYQDAIEIELEN---------SRPKAKGGL--ASLAKDGSTEDTASSLCGEED  364

1 336  -EKERME-KAARE----------------ENEGGISSLNHDEQLYSDYCSMSKRGKQCSEKGEKDKDGAEEHADTSDESD  397
2 341  -EKERMEEKASSE----------------ENDKQVNLDGH--------------------------------SDESS  368
3 385  IQNEHVE-GEDKE-TEGTEEKVKKVKTKTSEEKKENESGHFQDAIDGY-SLETDPYCPRNLHLLPTTDTYLSKVSDDPDN  461
4 365  TEDEELE-AAASHLNKDLYRELLGGAPGSSEAAGSPEWGGRPPAL-GS-LLDPLPTAA-SLGISDSIRECISSQSSDPKD  440

1 398  -------NFSDISDDEVNGYINNEEETHYKTITWTEMNKDYLEEQAAKEAALKAASEALKASNSNCPEDARKAFEAAKAD  470
2 369  -------TLSDVDDRELDCYFRTPEEVRLVKIFFDHENPGYDEKEAAK-----------KAAGLNACNNASNIFEASKAA  430
3 462  --------LEDVDDEELNAHLLNEEASKLKERIWIGLNADFLLEQESKRLKQE------------------ADIATGNTS  515
4 441  ASGDGELDLSGIDDLEIDRYILNESEARVKAELWMRENAEYLREQREKEARIA-----------------KEKELGIYK  502

1 471  AAKSRKEKQQKKAEEA------------------------------KNAAPPATAVEAVRRTLDKKRLSSVINYDVLESLF  521
2 431  AAKSRKEKRQQRAEEE------------------------------KNAPPPATGIEAVDSMVKRKKFRD-INCDYLEELF  480
3 516  VKKKRTRRRNNTRSDEPTKTVDAAAAIGLMSDLQDKSGLHAALKAAEESGDFTTADSVKNMLQKASFSKKINYDAIDGLF  595
4 503  EHKPKKSCKRR------------------------------EPIQASTAREAIEKMLEQKKISSKINYSVLRGLS  547

1 522  DTSA-----PEKSPKRSKTETDIEKKK---------------EENKEMKSNE-------------------------  553
2 481  DASV------EKSPKRSKTETVMEKKK--------------KEEHEIVENE-------------------------  511
3 596  R------------------------------------------------------------------------  596
4 548  SAGGGSPHREDAQPEHSASARKLSRRRTPASRSGADPVTSVGKRLRPLVSTQPAKKVATGEALLPSSPTLGAEPARPQAV  627

1 554  --------HENGENEDEDEEDEEEGNVESYDMKTDFQNGEKFYEEDEEEEEDGNDFGLY  604
2 512  --------QEEEDYAAPYEQDEED-YAAPYEMNTD----KKFYESEVEEEEDGYDFGLY  557
3        ---------------------------------------------------------
4 628  LVESGPVSYHADEEADEEEPDEEDGEPCVSALQMMGSN----DYGCDGDEDDGY-----  677
```

**Figure S6**. Sequence alignment of the two putative Brf1-homologues from *Arabidopsis thaliana*.  1-AT3G09360 2 - AT2G45100, 3 – Brf1 from *Saccharomyces cerevisiae*, and 4 – Brf1 from *Homo sapiens*. Conserved regions are highlighted using the same color code as in Figure 4: black – represents the TFIIB zinc-binding domain, gray - cyclin fold boxes and white box - Brf1-like TBP-binding domain.