

Title:

Additional layer of regulation via convergent gene orientation in yeasts

Authors: Jules Gilet^{1,2†}, Romain Conte^{1,2}, Claire Torchet², Lionel Benard^{2*} and Ingrid Lafontaine^{1*}

Author Affiliations:

¹Institut de Biologie Physico-Chimique, UMR7141 Laboratoire de Biologie du Chloroplaste et Perception de la Lumière chez les Microalgues, CNRS, Sorbonne Université, , F-75005, Paris, France.

²Institut de Biologie Physico-Chimique, UMR8226, CNRS, Sorbonne Université, Laboratoire de Biologie Moléculaire et Cellulaire des Eucaryotes, F-75005, Paris, France.

†Present address: Institut Curie Immunologie - INSERM U932, PSL, F-75005 Paris, France.

* *To whom correspondence should be addressed:* ingrid.lafontaine@ibpc.fr (I.L.), lionel.benard@ibpc.fr (L.B.),

Keywords: yeasts, evolution, mRNA duplexes, post-transcriptional RNA processing

Abstract

Convergent gene pairs can produce transcripts with complementary sequences. We had shown that mRNA duplexes form *in vivo* in *Saccharomyces cerevisiae* via interactions of their 3'-ends and can lead to post-transcriptional regulatory events. Here we show that mRNA duplex formation is restricted to convergent genes separated by short intergenic distance, independently of their 3'-UTR length. We disclose an enrichment in genes involved in biological processes related to stress among these convergent genes. They are markedly conserved in convergent orientation in budding yeasts, meaning that this mode of post-transcriptional regulation could be shared in these organisms, conferring an additional level for modulating stress response. We thus investigated the mechanistic advantages potentially conferred by 3'-UTR mRNA interactions. Analysis of genome-wide transcriptome data revealed that Pat1 and Lsm1 factors, having 3'-UTR binding preference and participating to the remodeling of messenger ribonucleoprotein particles, bind differently these mRNA duplexes in comparison to canonical mRNAs. Functionally, the translational repression upon stress also appears limited on mRNA duplexes. We thus propose that mRNA duplex formation modulates the regulation of mRNA expression by limiting their access to translational repressors. Our results thus show that post-transcriptional regulation is an additional factor that determine the order of coding genes.

Introduction

The transcriptional orientation of genes relative to their adjacent gene neighbors along the chromosome can be either co-orientation ($\rightarrow\rightarrow$), divergence ($\leftarrow\rightarrow$) or convergence ($\rightarrow\leftarrow$) (fig. 1A). This genomic neighborhood may reveal functional constraints. In eukaryotes, neighboring genes are likely to be co-expressed, independently of their relative orientation (Cohen et al. 2000; Hurst et al. 2004; Michalak 2008). To date, most attention has been devoted to the link between genomic neighborhood and co-transcriptional regulation. Co-orientation can allow co-regulation of transcription of the two genes by a single promoter in an operon-like fashion (Osborn and Field 2009) and divergence can allow co-regulation of transcription by means of a bi-directional promoter (Wei et al. 2011). In the case of convergent gene pairs that do not share any promoter region, co-transcription could be mediated by chromatin effects rather than by direct interactions (Chen et al. 2010). Most importantly, there is increasing evidence that convergent gene orientation can also mediate regulation at the post-transcriptional

level. Transcriptome analyses have shown that convergent gene pairs can produce tail-to-tail 3'-overlapping mRNA pairs that can theoretically form mRNA duplexes in *S. cerevisiae* (Pelechano and Steinmetz 2013; Wilkening et al. 2013) and in other eukaryotes (Jen et al. 2005; Makalowska et al. 2005; Sanna et al. 2008). We previously demonstrated that such mRNAs duplexes exist extensively and can interact in the cytoplasm in *Saccharomyces cerevisiae* (Sinturel et al. 2015). RNA duplexes can promote post-transcriptional regulatory events including no-go decay, precisely when the major cytoplasmic 5'-3' exoribonuclease Xrn1 is regulated or deficient (Doma and Parker 2006; Passos et al. 2009; Sinturel et al. 2015).

We thus hypothesize that mRNA duplex formation can modulate interactions with RNA binding proteins that preferentially bind at mRNA 3'-end, like Pat1 and Lsm1. Pat1 and Lsm1 are considered as main translation repressors activated during stress and also as key players in 5' to 3' mRNA decay, linking deadenylation to decapping (Tharun and Parker 2001; Tharun et al. 2000; Chowdhury and Tharun 2009). Pat1 and Lsm1 are components of messenger ribonucleoprotein particles (mRNP), named P-bodies, which contain translational repressors and the mRNA decay machinery (Mitchell and Parker 2014).

If genomic neighborhood plays some critical role in gene expression, it should be conserved when under selection, and some favorable new genomic neighborhoods should appear during evolution. Although there are many evidences for cis-regulatory constraints on gene order, our understanding of the determinants of the evolution of gene order in eukaryotes is still limited. Globally, gene pair conservation decreases as intergenic distance increases (Hurst et al. 2002; Poyatos and Hurst 2007). In yeasts, gene pairs that are highly co-expressed are more conserved than gene pairs that are not co-expressed and it has been reported that only divergent gene pairs are under selection for high co-expression (Kensche et al. 2008; Wang et al. 2011; Yan et al. 2016). However, the co-expression of linked genes persists long after their separation by chromosomal rearrangements whatever their original relative orientation, and natural selection often favors chromosomal rearrangements in which co-expressed genes become neighbors. Thus, selectively favorable co-expression appears not to be restricted to bi-directional promoters (Wang et al. 2011).

In order to determine the possible role of mRNA duplex formation in gene regulation, we performed a genomic analysis on an evolutionary perspective to determine i) the extent to which mRNA duplexes could form in 45 budding yeasts (*Saccharomycotina* subphylum), covering an evolutionary distance of *ca.* 300 MYA (Marcet-Houben and Gabaldón 2015), ii)

their functional properties by a Gene Ontology (GO) enrichment analysis, iii) the conservation of convergent orientation among yeasts as a proxy for their functional importance. We also compared the properties of mRNA duplexes and canonical mRNAs that do not form duplexes. First by analyzing cross-linking immunoprecipitation (CLIP) data used to map the interaction sites of Pat1 and Lsm1 on mRNAs (Mitchell et al. 2013) in condition of glucose deprivation, a condition triggering the formations of P-bodies, particularly requiring Pat1 and Lsm1. Secondly by analyzing ribosome loading data used to determine the ability of Pat1 and Lsm1 factors to repress translation of mRNAs in stress conditions (Garre et al. 2018).

Our results show that mRNA duplexes form between genes that are less than 200 bp apart, independently of the length of their 3'-UTR. They are functionally enriched in biological processes occurring during response to stress in *S. cerevisiae* and provided orthology-function relationships are preserved, it is also the case in many yeasts for convergent genes less than 200 bp, theoretically able to form mRNA duplexes. We propose that mRNA-mRNA interactions can interfere with canonical mRNP remodelers such as Pat1 and Lsm1, contributing to limit the translational repression on mRNA duplexes and thus participating in modulating gene expression upon stress. Furthermore, convergent orientation between neighboring genes is in general more conserved at short intergenic distances than co-oriented or divergent orientation in all 45 studied genomes, which suggests that convergent orientation allowing post-translation regulation of mRNA of genes involved in stress response is widely shared.

Results

RNA duplexes in *S. cerevisiae* occur between genes separated by short intergenic distances

The distribution of 3'-UTR length of the 365 mRNA duplexes determined experimentally are not statistically different than the distributions of the other mRNAs that do not form any duplex, previously named solo mRNAs (Sinturel et al. 2015) but hereafter named canonical mRNAs (Mann-Whitney tests, p-value < 0.05, fig. 1B). Conversely, the intergenic distance, defined as the distance between CDS regions of adjacent genes (fig. 1A), between convergent genes forming mRNA duplexes are shorter (median 155 bp) than between convergent genes producing canonical mRNAs (median 236 bp) (Mann-Whitney test, p-value < 10^{-13} , fig. 1C). This suggests that the short intergenic distances between convergent genes are the major determinant for mRNA duplex formation. We note that the bimodal distribution of

intergenic lengths among divergent gene pairs is in agreement with previous reports (Hermsen et al. 2008).

The reconstructed phylogenetic tree of the 45 yeasts studied, congruent with the backbone of the Saccharomycotina phylogeny (Shen et al. 2016), is presented in fig. 2. Within these genomes, convergent gene pairs are separated by the smallest intergenic distance, with a median of 158 bp, compared to co-oriented gene pairs (median of 405 bp) or divergent gene pairs (median of 517 bp) (supplementary table S1), a trend previously observed (Chen et al. 2011). In addition, a comparative transcriptomic analysis in different yeasts revealed that 3'-UTR lengths are also broadly similar (Moqtaderi et al. 2013), demonstrating that the majority of convergent transcripts overlap and are theoretically able to form mRNA duplexes in yeasts. In 36 out of 45 genomes, the proportion of co-oriented gene pairs are significantly smaller than expected under a neutral model of gene order evolution, where genes would be equally distributed among the two DNA strands (50% of co-oriented, 25% of divergent and convergent) (supplementary table S2). Neighboring genes are then more often encoded on opposite strands, probably due to a greater impact of bidirectional promoters and of chromatin context for transcriptional regulation, or a greater impact of mRNA duplex formation for post-transcriptional regulation. Interestingly, in 27 of these 36 genomes, the proportion of convergent pairs is higher than those of divergent pairs.

GO analysis of convergent genes

According to the Yeast GO Slims annotation (see Methods), the 365 validated mRNA duplexes in *Saccharomyces cerevisiae* are significantly enriched (more than 2-fold enrichment, hypergeometric test, adjusted p-values < 0.01) into cellular response to DNA damage stimulus, DNA metabolism (repair, recombination and replication) together with mRNA processing and RNA splicing (Table 1). We found the same results by estimating the probability of occurrence of each GO Slim term among 365 randomly selected genes (permutation p-values < 0.01) (supplementary table S3).

For each of the other Saccharomycotina species, we determined the functional distribution of GO Slim terms among the total number of convergent genes with an identified ortholog in *S. cerevisiae* (N_{conv}) -assuming they share the same function as their *S. cerevisiae* ortholog- that are less than 200 bp apart and thus theoretically able to form mRNA duplexes

(see previous section). We next calculated the probability of occurrence of each GO Slim term among N_{conv} randomly selected genes with an identified *S. cerevisiae* ortholog (fig. 3 and supplementary table S4). As in *S. cerevisiae*, in at least 18 genomes, there is a significant enrichment (permutation p-value < 0.05) for cellular response to DNA damage stimulus, DNA repair and RNA splicing. In addition, mRNA modification, tRNA processing, chromosome segregation and protein complex biogenesis are also enriched in more than one third of the yeast species. Such functional enrichment of terms that could be linked to stress response among convergent genes theoretically able to form RNA duplexes indicates that their mode of post-transcriptional regulation could be shared in most of these species.

Table 1. Goslim terms of biological processes significantly enriched (p-value < 0.01) more than 2 fold in *S. cerevisiae* mRNA duplexes validated experimentally.

GO Slim term	# all genes ^a	# conv duplex ^b	adj. p-value ^c
DNA repair	245	32	3.03E-04
cellular response to DNA damage stimulus	299	33	2.12E-03
DNA replication	152	18	8.73E-03
DNA recombination	174	20	8.73E-03
mRNA processing	167	19	8.73E-03
RNA splicing	134	17	8.73E-03

Note - ^a number of *S. cerevisiae* genes annotated with GOslim term. ^c p-value of the hypergeometric test adjusted with the Benjamini-Hochberg procedure for multiple testing (see Methods).

Convergent relative orientation is more conserved than divergent and co-oriented relative orientations at short intergenic distances

A further insight at the potential importance of convergent orientation is its conservation during evolution. We defined the orthologs between each pair of the 45 yeast genomes (see Methods), and determined their relative orientation to their gene neighbors in the two considered species.

We first considered a conserved gene orientation when two orthologs share the same relative orientation with respect to their 3' neighbor in both genomes, independently of the orthology relationship of the neighboring genes (genomic context). It allows to estimate the extant at which the relative orientation is functionally important in itself. On average, co-orientation is the most conserved gene orientation (78%) followed by convergence (75%) and divergence (72%) (supplementary table S5). Indeed, conservation decreases as the evolutionary

distance between species increases (fig. 4). At four large pairwise evolutionary distances (1, 1.1, 1.4 and 1.5) co-orientation is the less conserved orientation and convergence the most conserved one. These distances always concern the species that are the most isolated from all the other studied species: *Cyberlindnera fabianii*, the only species from the *Phaffomycetaceae* clade studied (Kurtzman et al. 2008), and *Y. Lipolytica* and *B. adenivorans* belonging to the most distant clade from all the others (fig. 2). These three species also share the lowest numbers of ortholog with all other species: 804, 1273 and 1722 orthologs shared with another species for *Y. Lipolytica*, *A. adenivorans* and *C. fabiani* respectively, for an average of 2874 for the 45 species studied. This particular trend could reflect either a sampling bias or that the most conserved gene orientation between different yeast clades is the convergent one.

As convergent genes are separated by smaller intergenic regions than divergent and co-oriented genes, we looked whether the physical proximity of convergent genes was the main explanation to their preferential conservation. To that end, we estimated the proportion of conserved gene relative orientation in windows of non-overlapping intergenic distances ranging from 100 to 1000 bp. Strikingly, for intergenic distances lower than 200 bp, convergent gene orientation, thus theoretically able to form mRNA duplexes, is more conserved than divergent and co-oriented gene orientations (fig. 5). In addition, conservation of convergent orientation decreases more rapidly (exponential decrease) as the intergenic distance increases compared to the conservation of divergent and co-oriented pairs (polynomial decrease) (fig. 5). Above 400 bp, there is an opposite trend, the convergent gene orientation being less conserved, while co-oriented gene orientation is the most conserved. These trends hold when considering species pairs at different evolutionary distances (supplementary fig. S1A). We also performed these calculations by counting conserved orientation at a given intergenic distance imposed only in one of the two genomes, thus considering the possibility that the ancestor intergenic distance could correspond either to one or the other distance in the two genomes. As expected, the conservation proportions increases in this case (supplementary fig. S1B).

We next estimated the conservation of gene pairs as pairs of adjacent protein-coding genes with adjacent orthologs with the same relative orientation in the two genomes (microsynteny). Like convergent orientation, convergent gene pairs are also the most conserved ones compared to divergent and co-oriented pairs at small intergenic distances but to a lower extent in terms of proportion and in terms of evolutionary distance, as this is a more stringent criterion (supplementary fig. S2 A and B). Thus, the conservation of the genomic context is more important than the conservation of microsynteny. As expected under a neutral model of

evolution of gene order, the probability of gene pair conservation decreases with the length of the intergenic region between the genes for all types of pairs because the probability of recombination between two genes increases with the distance separating them. It has been shown indeed that intergenic distance is the major determinant of gene pairs conservation in yeasts (Poyatos and Hurst 2007). However, at small intergenic distances, the convergent pairs are less prone to recombination than co-oriented and divergent ones. This either reflects a selective pressure to maintain convergent pairs at small intergenic distances allowing RNA duplexes formation, and/or a counterselection of co-oriented and divergent pairs.

In summary, at small intergenic distances which allow for mRNA duplex formation, the conservation of convergent gene pairs appears not neutral, and the conservation of the convergent orientation of a given gene is even more important, suggesting that convergent gene pairs are either conserved or recruited by chromosomal rearrangements for functional constraints.

mRNA duplexes limit Lsm1 and Pat1 interactions in 3'-UTRs

The enrichment of convergent genes in functions related to stress and their conservation in Saccharomycotina species encouraged us to question how these mRNA-mRNA interactions could confer an advantage along evolution. One hypothesis is that 3'-end RNA interactions affect mRNA access to mRNP remodeling proteins known to preferentially bind to the 3'-ends of mRNAs, such as the Pat1 and Lsm1 translational repressors comprised in P-bodies that participate in stress response (Chowdhury et al. 2007; He and Parker 2001). We analyzed CLIP data (see Methods; (Mitchell et al. 2013)) used to map the interaction sites of different P-body components, including Pat1, Lsm1, Dhh1 -that has no clear RNA sequence binding preferences- and Sbp1 that is involved in enhancing the decapping of mRNA that binds preferentially to 5'-UTR presumably resulting from its affinity to eIF4G (Rajyaguru et al. 2012; Sheth and Parker 2003; Mitchell et al. 2013).

A metagene representation of specific protein interactions has been computed for mRNA duplexes and canonical mRNAs from normalized reads (in RPKM), corresponding to protein interaction sites. For canonical mRNAs, as previously observed (Mitchell et al. 2013), Pat1 and Lsm1 preferentially bind the 3'-end of mRNA, Dhh1 has no positional bias and Sbp1 positions are biased towards the 5'-UTR region (fig. 6). In contrast, a significant shift in binding peaks of Pat1 and Lsm1 is observed in the 3'-UTR region of mRNA duplexes (fig. 6) suggesting

that 3'-RNA interactions might limit Pat1 and Lsm1 access. Spb1 also shows a decrease in the preferential binding on the 5'-UTR of mRNAs duplexes in addition to a mild decrease in their 3'-UTRs. The peak distribution associated to Dhh1 is not significantly different for canonical and mRNA duplexes, in accordance with the fact that Dhh1 has no clear RNA sequence binding preferences (Mitchell et al., 2013). Taken together these observations argue for significant altered associations of mRNA duplexes, with Pat1 and Lsm1 leading to mRNP differing from those assembled from canonical RNAs. Therefore, the fate of mRNA duplexes should differ from the fate of canonical mRNAs upon stress.

mRNA duplexes escape the ribosome access control governed by Pat1 upon stress

In order to further examine how the decrease in interactions of Lsm1 and Pat1 on 3'-UTR regions of mRNA duplexes can affect ribosome dynamics *-i.e* translation initiation-, we took advantage of a published genome-wide analysis performed in condition of osmotic stress, during which P-bodies, involving Pat1 and Lsm1, are formed. These were used to determine ribosome mRNA associations in wild type (WT), *lsm1* and *pat1* mutant strains (Garre et al. 2018). We thus compared the ribosome accumulation at 5'-UTR of mRNA duplexes and canonical mRNAs in WT and *pat1* and *lsm1* mutants. The ribosome loading (*i.e.* log₂ ratio of 5' sequencing reads obtained upstream versus downstream of the mRNA translation start site (Garre et al. 2018)) in WT, *lsm1* and *pat1* mutants for each mRNA category in normal growth condition and osmotic stress is presented in fig. 7. A positive shift of ribosome loading for a given mRNA between two genetic backgrounds will reflect an increase in ribosome access, thus revealing a decreased translational repression (Garre et al. 2018). A positive shift of ribosome loading in *pat1* and *lsm1* mutants was observed when all mRNAs are globally computed, as previously reported (Garre et al. 2018). We observed a similar shift for canonical mRNAs and mRNA duplexes, which confirms the general role of Lsm1 and Pat1 in limiting ribosome access on mRNAs in normal growth conditions, (left panels, fig. 7). In stress conditions, this positive shift in ribosome loading in mutants versus WT was lost for mRNA duplexes in *pat1* mutants only (right panels, fig. 7). This difference between mRNA duplexes for *pat1* and *lsm1* mutants is also visible when considering the ratios of ribosome loading in mutant relative to WT in stress condition (supplementary fig. S3), whereas these differences are no longer visible for canonical mRNAs. Taken together, these results suggest that the ribosome accumulation on mRNA duplexes is not dependent on the presence of Pat1 upon stress. Pat1 being considered a main translation repressor, this suggests that mRNPs formed by

mRNA duplexes differ from canonical mRNPs and thus might be differently controlled at the translational level.

Discussion

In this study, we showed that convergent genes separated by short intergenic spaces are likely to produce mRNAs duplexes with 3'-end overlapping and complementary sequences, as we previously described (Sinturel et al. 2015), independently of their 3'-UTR length. Given that the median length of intergenes separating convergent gene pairs in Saccharmycotina genomes is of 158 bp, we propose that mRNA duplexes can form in most of these yeasts. Indeed, intergenes between convergent pairs are the smallest ones, while those between divergent pairs are the longest ones, as previously observed among fungi (Kensche et al. 2008).

As intergenic distances is the major determinant of gene pair conservation (Poyatos and Hurst 2007), one could argue that convergent pairs, having smaller intergenic regions will inherently be more conserved, independently of selection. However, we have shown that at short intergenic distances (< 200 bp), convergent gene pairs are more conserved than divergent and convergent ones. Thus, the close proximity between convergent genes can also be considered as strongly beneficial, because of tightening their linkage. A trend already observed in higher plants: at small intergenic distances (< 250 bp) there is a higher conservation of convergent gene pairs than divergent ones between *Arabidopsis*, *Populus* and Rice genomes (Krom and Ramakrishna 2008). One could also argue that the low conservation of divergent pairs is counter-selected at the smallest intergenic distances because they can barely contain a canonical promoter region that helps the anchoring of the transcription machinery that ranges from *ca.* 115 ± 50 bp in yeasts (Venters and Pugh 2009; Chen et al. 2011; Lubliner et al. 2013). The bimodal distribution of intergenic distances between divergent pairs most probably reflect additional cis-regulatory constraints, as previously reported (Hermsen et al. 2008). In line with this view, among recently formed gene pairs in yeasts, divergent ones are counter-selected and are separated by very long intergenic regions (978 bp on average) (Chen et al. 2011; Sugino and Innan 2012). However, when only considering the conservation of a gene relative orientation with respect to its neighbor, the same trend holds, i.e. conservation of gene orientation is higher for convergent genes than for the other orientations at small intergenic distances and we showed that the decreased conservation as the intergenic distance increases

has not the same behavior for genes in convergent orientation (exponential decrease) than genes in the two other orientations (polynomial decrease). This could reflect a functional advantage of convergent genes with small intergenic spacers, related to their ability to form RNA duplexes and its possible influence on the post-transcriptional regulation of their expression. This is in agreement with previous analyses posing that selectively favorable co-expression appears not to be restricted to bi-directional promoters (Wang et al. 2011). This is further supported by our observation that genes in convergent orientation present an enrichment in functions related to stress in all studied genomes. Thus, the selective pressure would rather be exerted on the conservation (or creation) of convergent relative orientation of genes, rather than on the conservation of microsynteny.

To investigate the structure of mRNPs produced by mRNA duplexes, we reconsidered CLIP data previously used to map the distribution of different mRNA binding proteins, Lsm1, Pat1, Dhh1, and Sbp1 on mRNAs in conditions of stress (Mitchell et al., 2013). Lsm1, Pat1, Dhh1, and Sbp1 are components of P-bodies, foci formed by stress. Interestingly, Lsm1 and Pat1 were found less frequently associated with the 3'-UTR of mRNAs forming mRNA duplexes than with the 3'-UTR of others canonical mRNAs. Previous analysis did not determine a particular consensus explaining why these factors bind preferentially 3'-UTR regions of canonical mRNAs (Mitchell et al., 2013) but we found that 3'-end mRNA-mRNA interactions significantly counteract Lsm1 and Pat1 associations. Here we demonstrated that Dhh1-mRNA association is not affected by mRNA-mRNA interactions, confirming that Dhh1 interaction with mRNA is not region specific. Then, the less frequent associations of both factors and the moderate altered association with Sbp1 reflect a particular assembly of mRNPs. We cannot exclude that the limited Lsm1/Pat1 association also reflects a preference for other mRNA binding proteins whose access will be facilitated by the existence of double-stranded RNA sequences. In this regards, mRNP structures are complex and a multitude of other mRNA binding proteins might participate in structure assemblies of mRNA duplexes (Mitchell et al., 2013; Garre et al., 2018).

It was thus critical to assess the role of an apparent decrease in 3'-UTR associations for Lsm1 or Pat1 in the functionality of mRNA duplexes. From analysis of a genome-wide functional assay investigating the impact of Lsm1 and Pat1 on ribosome access on mRNAs (Garre et al., 2018), we found that ribosome access on mRNA duplexes is not modulated by Pat1, in contrast to that observed for canonical mRNAs. However, we found that Lsm1 still

modulates the ribosome access on mRNA duplexes, suggesting that Lsm1 and Pat1 have different roles for this mRNA category although their association deduced from CLIP data are similar. We thus propose that Pat1 and Lsm1 protein networks may not completely overlap and thus differently impact mRNP assemblies. In this regard, Pat1 has been proposed as a key component in promoting the formation of P-bodies (Sachdev et al. 2019). We thus propose that mRNAs forming mRNA duplexes escape to the Pat1-dependent translation repression upon stress.

In conclusion, we showed that the conservation of the convergent orientation of genes separated by short intergenic distances is important in budding yeasts and that those convergent genes are functionally associated with stress response. Convergent genes can produce mRNA forming duplexes *in vivo* and our results argue for a remodeling of mRNP by those mRNA-mRNA interactions, thus providing a selective advantage for modulating gene expression upon stress (fig. 8). Such a post-translational regulation process is most probably conserved among budding yeasts and should be considered as a possible part of the stress response in other living cells. Thus, it should be considered as an additional factor that determine the order of coding genes.

Materials and Methods

Data collection

The genome annotation of the *S. cerevisiae*, as well as GO Slim terms for *S. cerevisiae* genes (version 05/18/2013) were retrieved at SGD (www.yeastgenome.org), Accession numbers and address retrieval for the other 44 genomes are given in supplementary table S6.

Orthology relationships

Pairwise orthology relationships among all 45 genomes were defined between syntenic homologs retrieved with the SynCHro algorithm (Drillon et al. 2014), with the version available in June 2015.

Conservation of relative orientation

Between two species, we estimated the proportion of orthologs that are in the same relative orientation with respect to their gene neighbor in the two genomes. We also determined the proportion of the pairs of adjacent genes whose orthologs in a given genome also form a pair of adjacent genes in that genome.

In a first approximation, without knowledge of the ancestral intergenic distance separating the ortholog and his gene neighbor, we define the smallest intergenic distance between an ortholog and its neighbor of the two contemporary genomes as the intergenic distance of the conservation.

Phylogenetic analyses

By transitivity, we inferred 224 groups of syntenic homologs composed of only one gene per species in the 45 yeasts studied. A multiple alignment of each group of orthologs was generated at the amino acid level with the MAFFT algorithm (v7.310, auto implementation, default parameters) (Kato and Toh 2008). Concatenation of the 224 alignments was used to estimate a concatenation tree with IQtree v1.6.7 (Nguyen et al. 2015; Kalyanamorthy et al. 2017). The best-fit estimated model is LG+F+I+G4. Maximum likelihood distances between each species pair were estimated from the concatenated alignment and used as the evolutionary distance between species.

Mapping of protein RNA binding sites

Analysis are based on the CLIP sequencing datasets for the RNA binding proteins (RBP) Dhh1, Lsm1, Pat1 and Sbp1 from (Mitchell et al. 2013), downloaded at www.ncbi.nlm.nih.gov/geo/. Adapter sequences were excluded from the reads with Cutadapt v1.1 (Martin 2011) and sequences less than 22 nucleotides long were removed. Bowtie2 v2.2.3 (end-to-end mode) (Langmead and Salzberg 2012) was used to align CLIP sequencing data (22-40 bp long) against 5'-UTR, ORF and 3'-UTR from 4415 coding transcripts of the reference genome (version R57-1-1, downloaded from <http://www.yeastgenome.org>) and (Nagalakshmi et al. 2008) for UTR coordinates. Aligned reads received a penalty score of -6 per mismatch, -5+(-3*gap length) per gap and were excluded if penalty score was less than the default threshold (between -13.8 and -24.6 for 22 and 40 bp reads respectively). Thus, aligned reads were allowed for less than a mismatch per 10bp, (1 mismatch per 9.52 to 9.75 bp, respectively) dynamically taking into account UV-light induced mutations consecutive to the sample processing. Duplicated PCR reads, as well as reads mapping to non-coding RNAs were excluded with samtools v1.2 (Li et

al. 2009) and uniquely mapped sequenced reads with a MAPQ score > 20 were kept (average MAPQ=33.26).

For each RBP-associated data, mapping of the peak interactions was constructed from a metagene aggregation procedure. The alignment depth for each gene at each nucleotide position has been determined with *samtools* (Li et al. 2009). In order to compensate control values without any sequencing signal, enrichment in the depth of sequencing signal per nucleotide coordinate S_n was defined as:

$$S_n = \frac{S_{nRPKMCLIPseq} + 1}{S_{nCTRLseq} + 1}$$

where $S_{nRPKMCLIPseq}$ is the sequencing signal at position n in the CLIP-seq experiment and $S_{nCTRLseq}$ is the sequencing signal at position n in the control experiment (RNA library without crosslinking and immunoprecipitation) of (Mitchell et al. 2013). Signal S_n has been normalized to the total number of aligned reads per experiment (*i.e.* per CLIP file), and to the length of each gene (reads per kilobase of transcript per millions of aligned reads, aka. RPKM).

In order to compare the relative binding sites of each protein along the transcripts, we performed a metagene analysis. Each single gene nucleotide coordinates n' have been adjusted to the longest sequence of either 5'-UTR, ORF or 3'-UTR regions to prevent any loss of information according the formula:

$$n' = \frac{n}{N} \times L$$

With n the position within the gene, N the gene length and L the longest nucleotide sequence in a defined region (5'-UTR, ORF and 3'-UTR) per CLIP experiment. Information computed for each region of a single gene were then concatenated to construct the final metagene of final length:

$$M = L_{5'UTR} + L_{ORF} + L_{3'UTR}$$

A direct interpolation was then conducted to compute the normalized depth of sequencing at each nucleotide position for the whole metagene length. Finally, the summed signal from each

$S_{n'}$ has been normalized according to the number of transcripts interacting with the RBP, in sort that in the final metagene representation:

$$\sum_{n=1}^M S_{n'} = 1$$

Analysis of ribosome 5'-UTR protection

Analysis of co-translational mRNA decay by global 5'P sequencing which allows the determination of ribosome mRNA protection was previously described (Garre et al. 2018). Data were analyzed using python in-house scripts.

Analysis of functional annotations

GO enrichment analyses among the genes forming RNA duplexes in *S. cerevisiae* with respect to the entire gene set of *S. cerevisiae* were performed with the hypergeometric test. The probability of occurrence of GO Slim terms at random in the subset of convergent genes in a given species have been computed over 1000 simulation trials. GO Slim terms for *S. cerevisiae* (go_slim_mapping.20130518.tab.gz) were retrieved at SGD's downloads site (<https://www.yeastgenome.org/search?category=download>). For species others than *S. cerevisiae*, the simulations were performed by considering all genes that have orthologs in *S. cerevisiae* and that are annotated with a GO Slim term.

Statistical analysis

Mann-Whitney U tests were performed to compare distributions of 3'-UTR length, intergenic distance, normalized CLIP seq signal along metagenes and ribosome accumulation. Hypergeometric tests were performed to determine the enrichment of GO Slim terms for genes forming RNA duplexes in *S. cerevisiae*, compared to the entire *S. cerevisiae* gene set. The expected probability of observing GO Slim terms for convergent genes in the other yeast genomes were performed as described in the paragraph above.

A false positive risk of $\alpha = 0.05$ was chosen as a significance threshold for all tests. P-values were adjusted with the Benjamini-Hochberg false discovery rate (Benjamini and Hochberg 1995) for GOSlim enrichment and with the Holm correction in the other cases (Holm 1979). All statistical calculations were performed with R functions and with functions from the Python *scipy* module.

Bibliography

- Chen W-H, Meaux J de, Lercher MJ. 2010. Co-expression of neighbouring genes in Arabidopsis: separating chromatin effects from direct interactions. *BMC Genomics* **11**: 178.
- Chen W-H, Wei W, Lercher MJ. 2011. Minimal regulatory spaces in yeast genomes. *BMC Genomics* **12**: 320.
- Chowdhury A, Mukhopadhyay J, Tharun S. 2007. The decapping activator Lsm1p-7p-Pat1p complex has the intrinsic ability to distinguish between oligoadenylated and polyadenylated RNAs. *Rna* **13**: 998–1016.
- Chowdhury A, Tharun S. 2009. Activation of decapping involves binding of the mRNA and facilitation of the post-binding steps by the Lsm1-7-Pat1 complex. *Rna* **15**: 1837–48.
- Cohen BA, Mitra RD, Hughes JD, Church GM. 2000. A computational analysis of whole-genome expression data reveals chromosomal domains of gene expression. *Nature Genetics* **26**: 183–186.
- Doma MK, Parker R. 2006. Endonucleolytic cleavage of eukaryotic mRNAs with stalls in translation elongation. *Nature* **440**: 561–4.
- Drillon G, Carbone A, Fischer G. 2014. SynChro: A Fast and Easy Tool to Reconstruct and Visualize Synteny Blocks along Eukaryotic Chromosomes. *PLoS ONE* **9**: e92621.
- Garre E, Pelechano V, Pino MS del, Alepuz P, Sunnerhagen P. 2018. The Lsm1-7/Pat1 complex binds to stress-activated mRNAs and modulates the response to hyperosmotic shock. *PLOS Genetics* **14**: e1007563.
- He W, Parker R. 2001. The yeast cytoplasmic Lsm1/Pat1p complex protects mRNA 3' termini from partial degradation. *Genetics* **158**: 1445–55.
- Hermesen R, ten Wolde PR, Teichmann S. 2008. Chance and necessity in chromosomal gene distributions. *Trends in Genetics* **24**: 216–219.
- Holm S. 1979. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics* **6**: 65–70.
- Hurst LD, Williams EJB, Pál C. 2002. Natural selection promotes the conservation of linkage of co-expressed genes. *Trends in Genetics* **18**: 604–606.
- Hurst LD, Pal C, Lercher MJ. 2004. The evolutionary dynamics of eukaryotic gene order. *Nature Reviews Genetics Nat Rev Genet* **5**: 299–310.
- Jen CH, Michalopoulos I, Westhead DR, Meyer P. 2005. Natural antisense transcripts with coding capacity in Arabidopsis may have a regulatory role that is not linked to double-stranded RNA degradation. *Genome biology* **6**: R51.
- Kalyanamoorthy S, Minh BQ, Wong TKF, Haeseler A von, Jermini LS. 2017. ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature Methods* **14**: 587–589.
- Katoh K, Toh H. 2008. Recent developments in the MAFFT multiple sequence alignment program. *Brief Bioinform* **9**: 286–98.

- Kensche PR, Oti M, Dutilh BE, Huynen MA. 2008. Conservation of divergent transcription in fungi. *Trends in Genetics* **24**: 207–211.
- Krom N, Ramakrishna W. 2008. Comparative Analysis of Divergent and Convergent Gene Pairs and Their Expression Patterns in Rice, Arabidopsis, and Populus. *PLANT PHYSIOLOGY* **147**: 1763–1773.
- Kurtzman CP, Robnett CJ, Basehoar-Powers E. 2008. Phylogenetic relationships among species of *Pichia*, *Issatchenkia* and *Williopsis* determined from multigene sequence analysis, and the proposal of *Barnettozyma* gen. nov., *Lindnera* gen. nov. and *Wickerhamomyces* gen. nov. *FEMS Yeast Research* **8**: 939–954.
- Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nature Methods* **9**: 357–359.
- Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079.
- Lublinter S, Keren L, Segal E. 2013. Sequence features of yeast and human core promoters that are predictive of maximal promoter activity. *Nucleic Acids Res* **41**: 5569–5581.
- Makalowska I, Lin CF, Makalowski W. 2005. Overlapping genes in vertebrate genomes. *Computational biology and chemistry* **29**: 1–12.
- Marcet-Houben M, Gabaldón T. 2015. Beyond the Whole-Genome Duplication: Phylogenetic Evidence for an Ancient Interspecies Hybridization in the Baker's Yeast Lineage. *PLoS Biology* **13**. <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4529251/> (Accessed June 27, 2017).
- Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* **17**: 10–12.
- Michalak P. 2008. Coexpression, coregulation, and cofunctionality of neighboring genes in eukaryotic genomes. *Genomics* **91**: 243–248.
- Mitchell SF, Jain S, She M, Parker R. 2013. Global analysis of yeast mRNPs. *Nat Struct Mol Biol* **20**: 127–133.
- Mitchell SF, Parker R. 2014. Principles and properties of eukaryotic mRNPs. *Molecular cell* **54**: 547–58.
- Moqtaderi Z, Geisberg JV, Jin Y, Fan X, Struhl K. 2013. Species-specific factors mediate extensive heterogeneity of mRNA 3' ends in yeasts. *PNAS* **110**: 11073–11078.
- Nagalakshmi U, Wang Z, Waern K, Shou C, Raha D, Gerstein M, Snyder M. 2008. The Transcriptional Landscape of the Yeast Genome Defined by RNA Sequencing. *Science* **320**: 1344–1349.
- Nguyen L-T, Schmidt HA, von Haeseler A, Minh BQ. 2015. IQ-TREE: A Fast and Effective Stochastic Algorithm for Estimating Maximum-Likelihood Phylogenies. *Mol Biol Evol* **32**: 268–274.
- Osborn AE, Field B. 2009. Operons. *Cell Mol Life Sci* **66**: 3755–3775.
- Passos DO, Doma MK, Shoemaker CJ, Muhlrads D, Green R, Weissman J, Hollien J, Parker R. 2009. Analysis of Dom34 and its function in no-go decay. *Mol Biol Cell* **20**: 3025–32.
- Pelechano V, Steinmetz LM. 2013. Gene regulation by antisense transcription. *Nature reviews Genetics*

14: 880–93.

Poyatos JF, Hurst LD. 2007. The determinants of gene order conservation in yeasts. *Genome Biol* **8**: R233.

Rajyaguru P, She M, Parker R. 2012. Scd6 targets eIF4G to repress translation: RGG motif proteins as a class of eIF4G-binding proteins. *Molecular cell* **45**: 244–54.

Sachdev R, Hondele M, Linsenmeier M, Vallotton P, Mugler CF, Arosio P, Weis K. 2019. Pat1 promotes processing body assembly by enhancing the phase separation of the DEAD-box ATPase Dhh1 and RNA. *eLife* **8**.

Sanna CR, Li WH, Zhang L. 2008. Overlapping genes in the human and mouse genomes. *BMC genomics* **9**: 169.

Shen X-X, Zhou X, Kominek J, Kurtzman CP, Hittinger CT, Rokas A. 2016. Reconstructing the Backbone of the Saccharomycotina Yeast Phylogeny Using Genome-Scale Data. *G3* **6**: 3927–3939.

Sheth U, Parker R. 2003. Decapping and decay of messenger RNA occur in cytoplasmic processing bodies. *Science* **300**: 805–8.

Sinturel F, Navickas A, Wery M, Descrimes M, Morillon A, Torchet C, Benard L. 2015. Cytoplasmic Control of Sense-Antisense mRNA Pairs. *Cell Rep* **12**: 1853–1864.

Sugino RP, Innan H. 2012. Natural Selection on Gene Order in the Genome Reorganization Process After Whole-Genome Duplication of Yeast. *Mol Biol Evol* **29**: 71–79.

Tharun S, He W, Mayes AE, Lennertz P, Beggs JD, Parker R. 2000. Yeast Sm-like proteins function in mRNA decapping and decay. *Nature* **404**: 515–8.

Tharun S, Parker R. 2001. Targeting an mRNA for decapping: displacement of translation factors and association of the Lsm1p-7p complex on deadenylated yeast mRNAs. *Molecular cell* **8**: 1075–83.

Venters BJ, Pugh BF. 2009. A canonical promoter organization of the transcription machinery and its regulators in the *Saccharomyces* genome. *Genome Res* **19**: 360–371.

Wang G-Z, Chen W-H, Lercher MJ. 2011. Coexpression of Linked Gene Pairs Persists Long after Their Separation. *Genome Biol Evol* **3**: 565–570.

Wei W, Pelechano V, Järvelin AI, Steinmetz LM. 2011. Functional consequences of bidirectional promoters. *Trends Genet* **27**: 267–276.

Wilkening S, Pelechano V, Jarvelin AI, Tekkedil MM, Anders S, Benes V, Steinmetz LM. 2013. An efficient method for genome-wide polyadenylation site mapping and RNA quantification. *Nucleic acids research* **41**: e65.

Yan C, Wu S, Pocetti C, Bai L. 2016. Regulation of cell-to-cell variability in divergent gene expression. *Nature Communications* **7**: 11099.

Authors contribution

IL and LB conceived the project; IL devised the work plan; JG, RC and IL performed research; IL and LB wrote the manuscript. All authors provided critical reading of the manuscript.

Acknowledgments

This study was supported by basic funding from CNRS and Sorbonne Université, by the “Initiative d’Excellence” program from the French State (Grant ‘DYNAMO’, ANR-11-LABX-0011- 01) and the AAP Emergence Sorbonne Université, SU-16-R-EMR-03. J.G. was supported by a fellowship from the Edmond de Rothschild Foundation.

We thank Benoist Laurent from the IBPC bioinformatic platform for technical support, M. Cavaiulo for critical reading of the manuscript and F.A. Wollman for fruitful discussions and critical reading of the manuscript.

Figure legends

Figure 1. mRNA duplexes form at small intergenic distances, independently of their lengths. (A) Schematic representation of relative orientation of adjacent genes. Intergenes are delimited by dashed blue lines and their distances indicated by double arrows. (B) mRNA 3’-UTR lengths (logarithmic scale) for different gene groups in *S. cerevisiae* taken from (Nagalakshmi et al. 2008). Median values are indicated for each group. (C) Intergenic lengths (logarithmic scale) for the same groups of genes. (D) Intergenic lengths for convergent (grey), co-oriented (lightgrey) and divergent (white) gene pairs in the 45 species studied. Species are named with a 4 letters code available in supplementary table S1 and ordered according to their evolutionary distance from *S. cerevisiae*. The dashed horizontal line at 155 bp indicates the median of 3’-UTR length of mRNA duplexes in *S. cerevisiae*. *conv_duplex*: convergent pairs forming experimentally validated mRNA duplexes, *conv_canonical*: convergent pairs with no experimentally validated RNA duplexes, *divergent*: genes in divergent orientation. *co-oriented*: genes in co-orientation.

Figure 2. Phylogenetic relationships of the 45 Saccharomycotina yeasts species studied.

Phylogeny of 45 Saccharomycotina species inferred from a maximum likelihood analysis of a concatenated alignment of 224 groups of syntenic homologs present in every genome (See Methods).

Figure 3. Functional enrichment in convergent genes. Biological process GO Slim terms with a significant enrichment (> 2 fold with a p -value < 0.05) among convergent genes compared to the whole gene population. *S. cerevisiae* GO slim terms have been attributed to orthologs in the 44 other species. Species are named with a 4 letters code available in supplementary table S6 and ordered according to their evolutionary distance from *S. cerevisiae*.

Figure 4. The conservation of relative orientation of orthologs decreases with the evolutionary distance between the species. Distribution of the frequency of orthologs that are in the same orientation relative to their adjacent gene in two species in function of the evolutionary distance between the two species (see Methods). Values correspond to the median of the frequencies for all pairs of species. Green dots: genes in convergent orientation (conv); black dots: co-oriented genes (coor) and blue dots: divergent genes (div).

Figure 5 The conservation of relative orientation is higher for convergent pairs less than 200 bp apart. Distribution of the frequency of orthologs that are in the same relative orientation relative to their adjacent gene in the two species in function of the length of the intergenic distance separating them. Values correspond to the median of the frequencies for all pairs of species for intergenic distances at most the value along the X axis. Same color legend as in fig. 4.

Figure 6. mRNA in duplexes have a marked loss of Pat1 and Lsm1 binding in their 3'-UTR. Metagene representation of sequence reads enriched using CLIP over the control sequence data for individual mRNAs. Normalized reads for canonical mRNAs and mRNAs forming duplexes are represented for Pat1, Lsm1, Dhh1 and Spb1 proteins (see methods, Mitchell et al., 2013). 5'-UTR, ORF and 3'-UTR regions are indicated. Lengths are scaled to the average 5'-UTR, ORF and 3'-UTR lengths over the entire genome. Green: mRNA duplexes, gray: canonical mRNAs, dark green: overlay observed in regions equivalently bent within the two mRNA classes. Mann-Whitney tests p -values for comparison of the distributions between the two mRNA classes are indicated. ***: $p < 1.0e-3$; n.s.: $p \geq 0.05$.

Figure 7. Impact of *lsm1* or *pat1* deletion on the over-accumulation of ribosomes in the 5'-UTR regions of canonical mRNAs or mRNA duplexes in both control and stress conditions. Ribosome loadings are calculated as a \log_2 ratio between ribosome profiling reads upstream of the start codon versus downstream of the start codon for each mRNA in WT, *pat1*

and *lsm1* strains, as previously determined by (Garre et al. 2018). Each mRNA is ranked along the Y axis according to its ribosome loading (X axis) calculated in different genetic backgrounds: WT (blue triangles), *pat1* (green squares), *lsm1* (red line). Individual panel represents ribosome loading of all mRNAs, canonical mRNAs or mRNA duplexes in both control (no stress, left) and stress (right) conditions. Mann-Whitney tests p-values for comparison of the distributions between WT and mutants are indicated. ****: $p < 1.0e-4$, ***: $p < 1.0e-3$; **: $p < 1.0e-3$; *: $p < 0.05$. n.s.: $p \geq 0.05$.

Figure 8. Model of the post-transcriptional regulation mediated by mRNA duplex formation. Upon stress, the translational repressor Pat1 binds preferentially to the 3'-UTR of canonical mRNAs, limits ribosome access on mRNA 5'-UTRs and promotes their aggregation into P-bodies, composed by a variety of mRNA-processing factors and translational repressors. mRNA duplexes escape Pat1 repression by masking 3'-UTR access and then fully participate in stress response.

Figure 1.

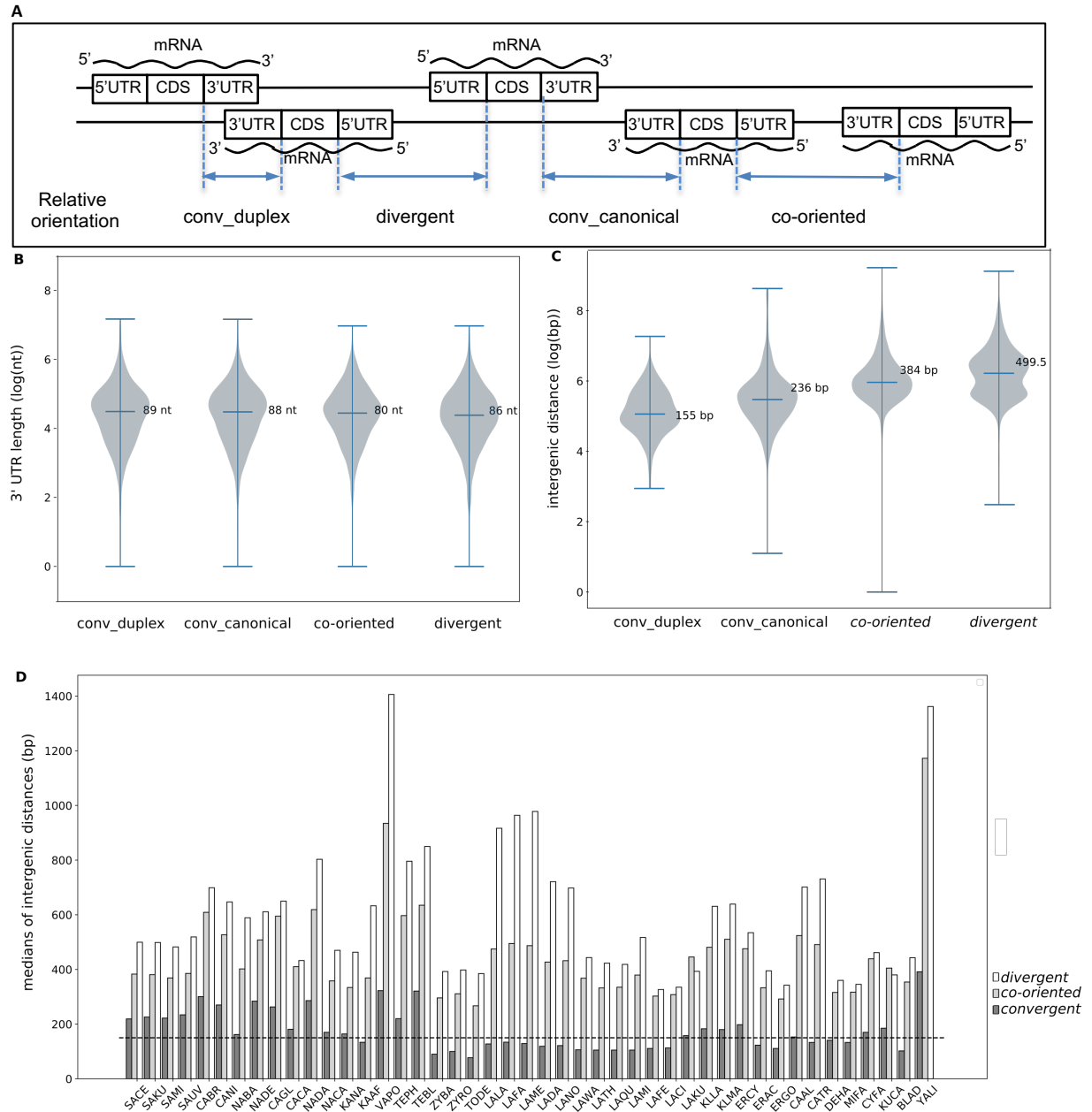


Figure 2.

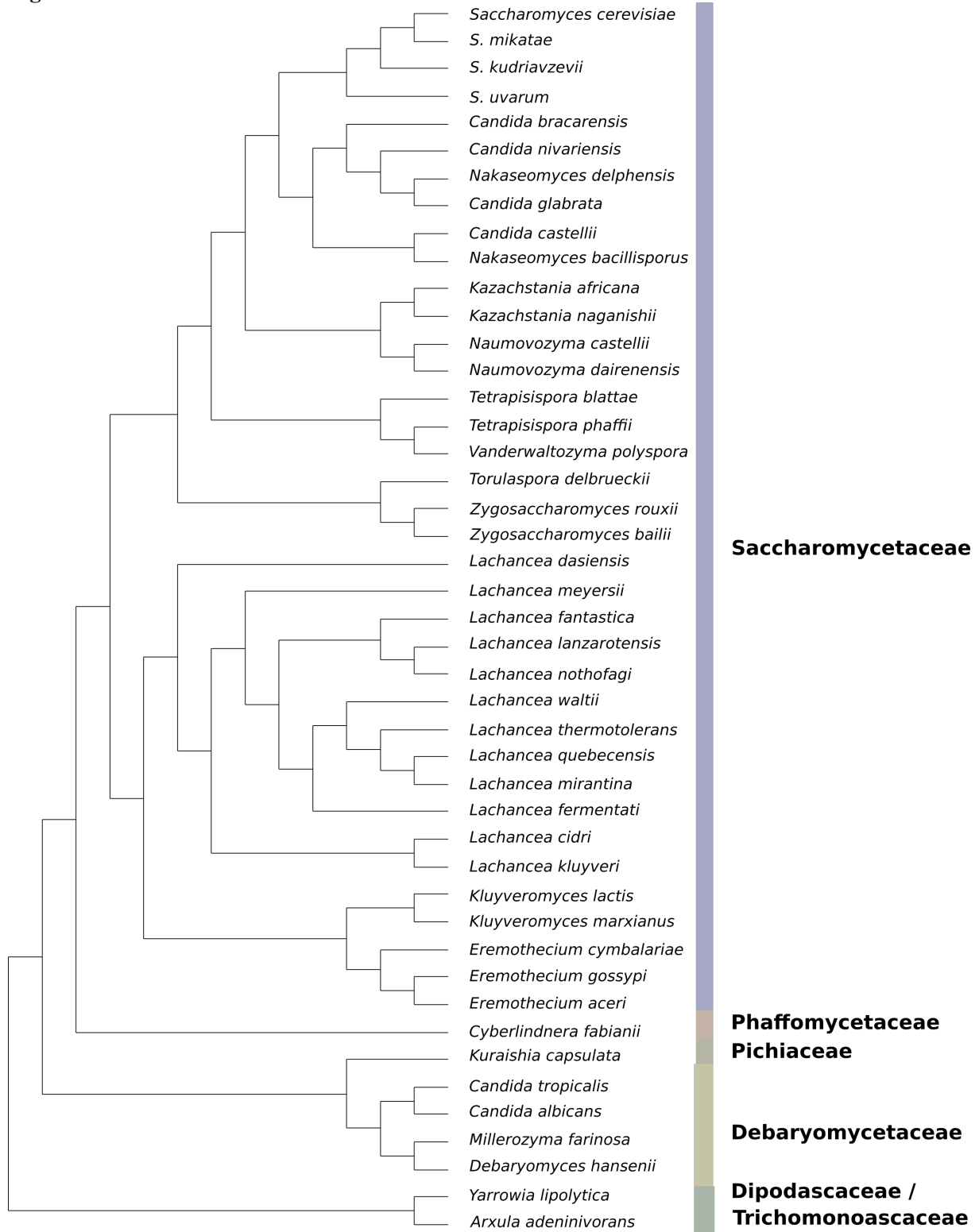


Figure 3

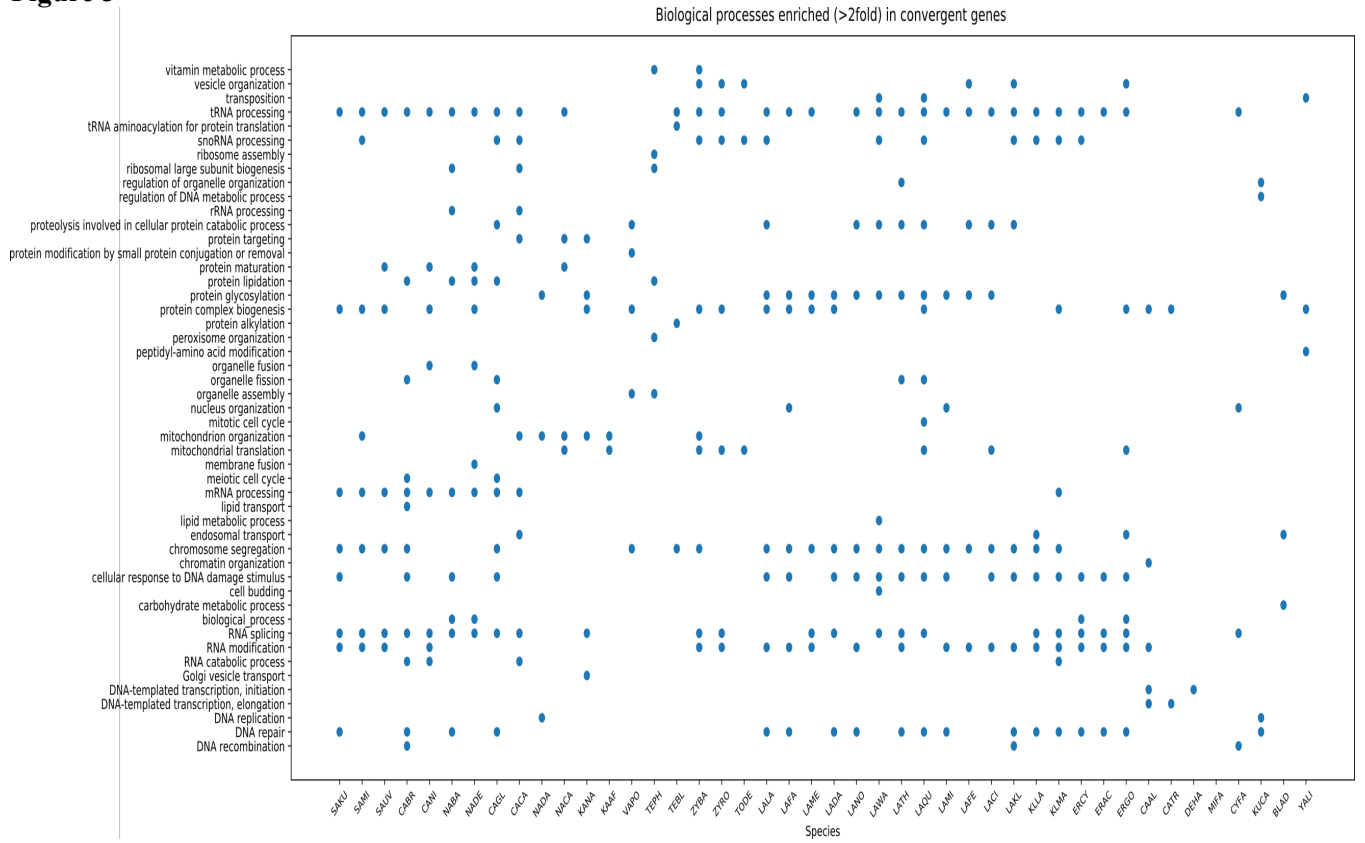


Figure 4

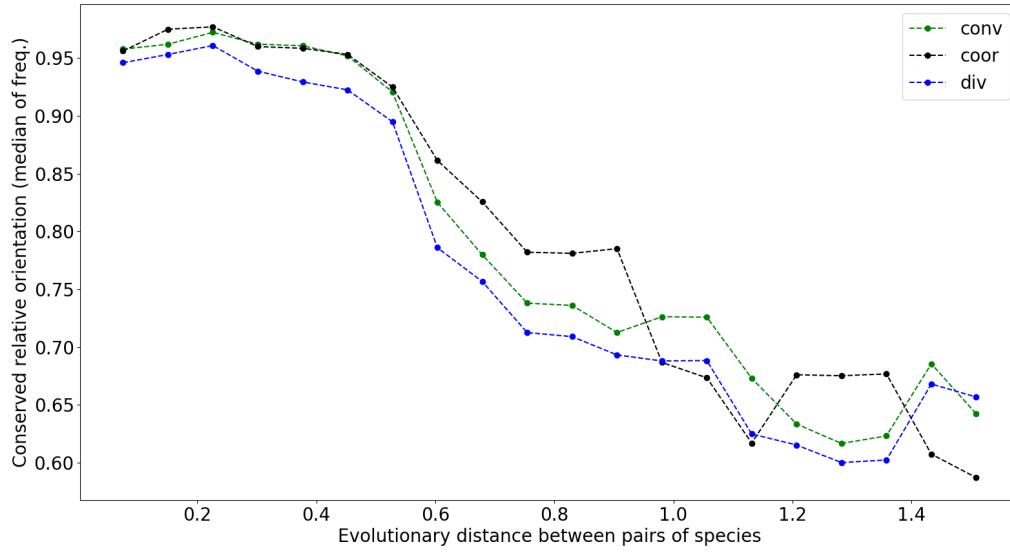


Figure 5.

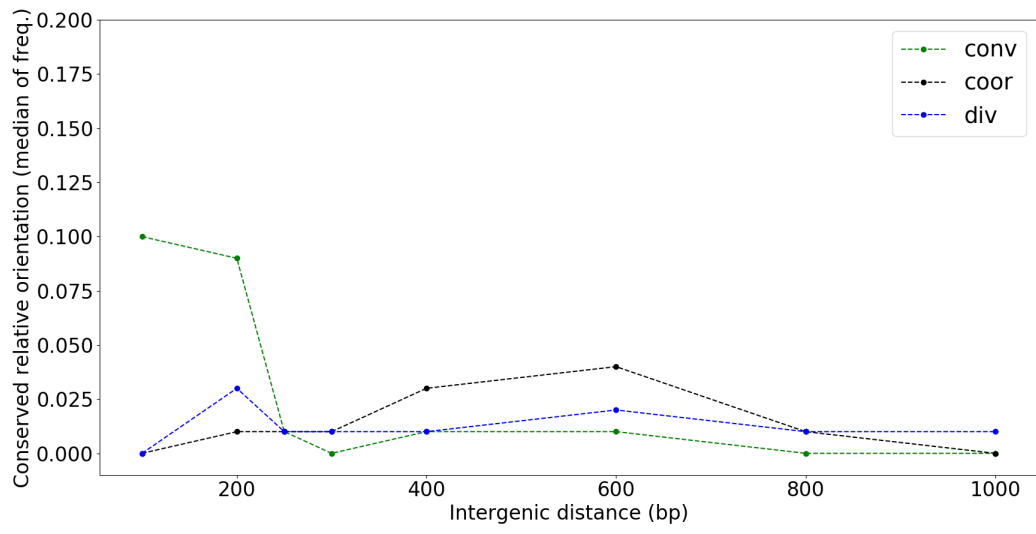


Figure 6

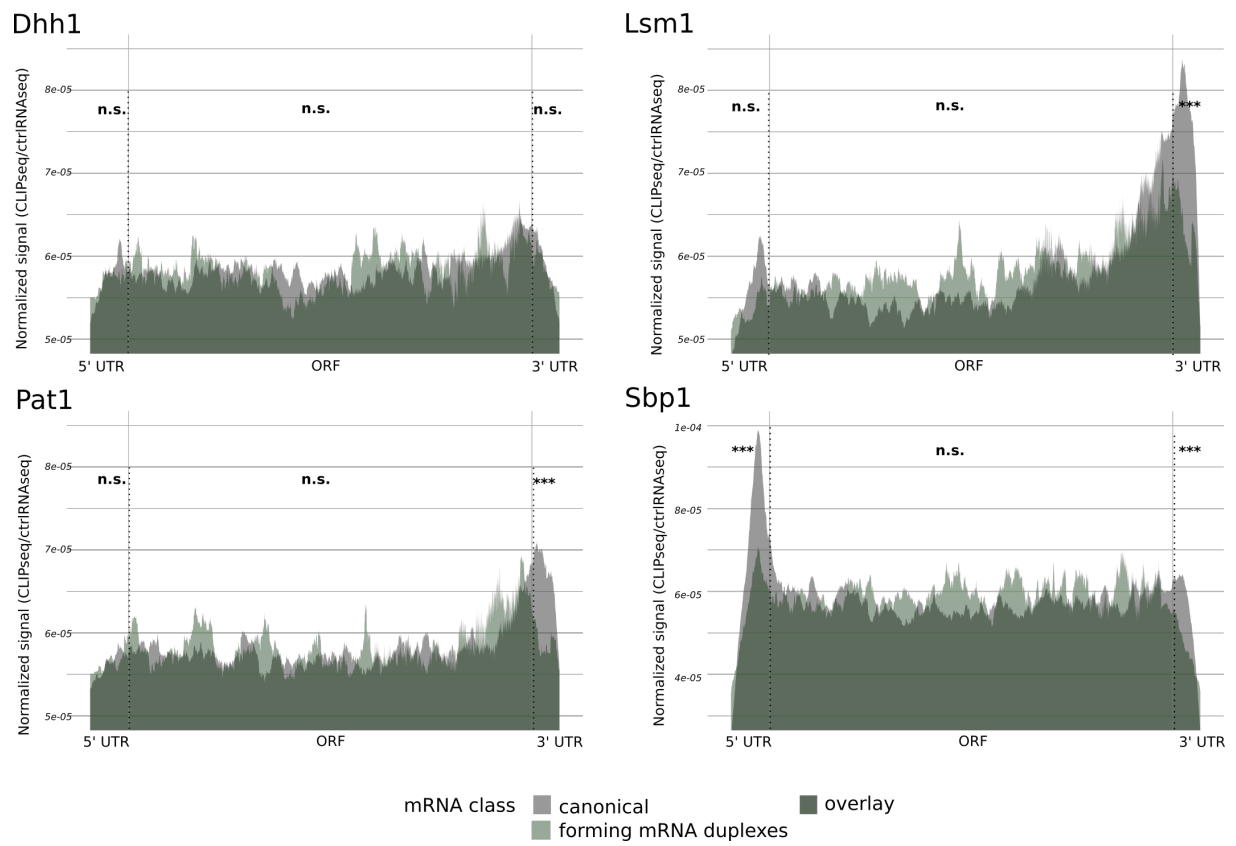


Figure 7

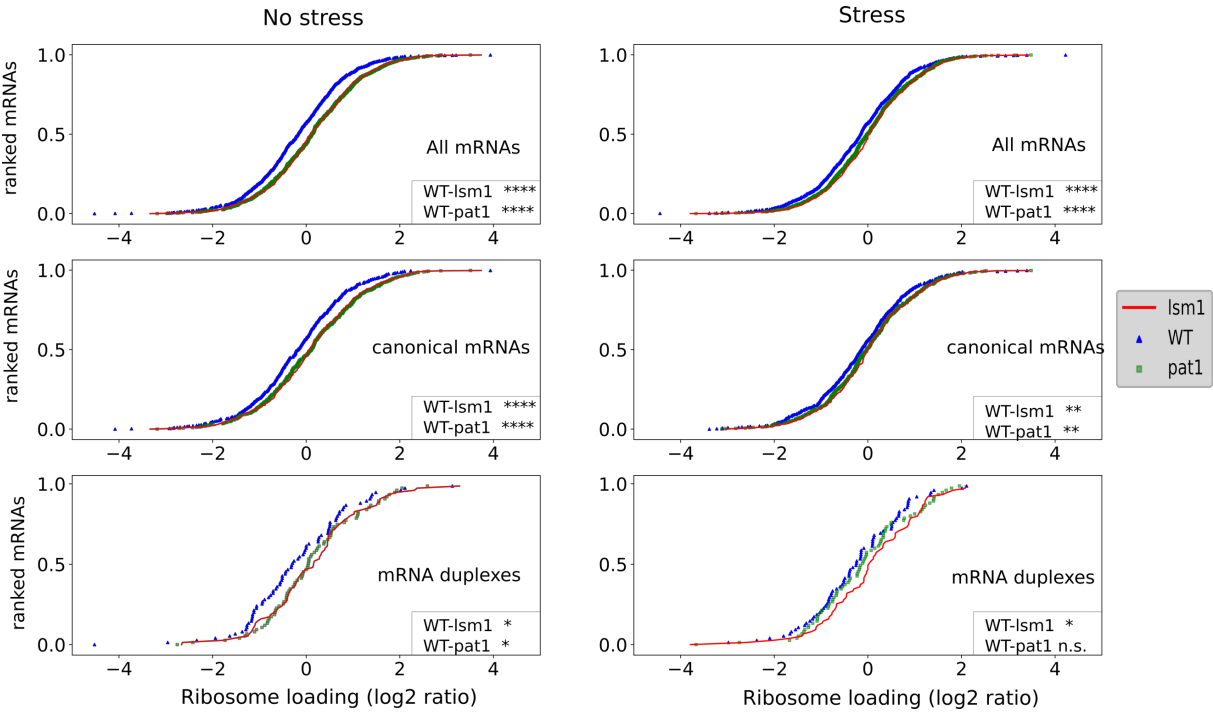


Figure 8.

