

1 **Start codon context controls translation initiation in the fungal kingdom**

2 Edward Wallace^{2*□}, Corinne Maufrais^{1,3□}, Jade Sales-Lee⁴, Laura Tuck², Luciana de Oliveira¹,
3 Frank Feuerbach⁶, Frédérique Moyrand¹, Prashanthi Natarajan⁴, Hiten D. Madhani^{4,5*},
4 Guilhem Janbon^{1*}

5

6

7 1. Institut Pasteur, Unité Biologie des ARN des Pathogènes Fongiques, Département de
8 Mycologie, F-75015, Paris, France

9 2. Institute for Cell Biology and SynthSys, School of Biological Sciences, University of
10 Edinburgh, UK

11 3. Institut Pasteur, HUB Bioinformatique et Biostatistique, C3BI, USR 3756 IP CNRS, F-75015,
12 Paris, France

13 4. Department of Biochemistry and Biophysics, University of California at San Francisco, San
14 Francisco, California 94158, USA

15 5. Chan-Zuckerberg Biohub, San Francisco, CA 94158, USA

16 6 Institut Pasteur, Unité Génétique des Interactions Macromoléculaire, Département
17 Génome et Génétique, F-75015, Paris, France

18 * Corresponding authors

19 □ these authors contribute equality to this work and should be considered as co-first authors

20

1 **Abstract**

2 Eukaryotic protein synthesis initiates at a start codon defined by an AUG and its surrounding
3 Kozak sequence context, but studies of *S. cerevisiae* suggest this context is of little
4 importance in fungi. We tested this concept in two pathogenic *Cryptococcus* species by
5 genome-wide mapping of translation and of mRNA 5' and 3' ends. We observed that
6 upstream open reading frames (uORFs) are a major contributor to translation repression,
7 that uORF use depends on the Kozak sequence context of its start codon, and that uORFs
8 with strong contexts promote nonsense-mediated mRNA decay. Numerous *Cryptococcus*
9 mRNAs encode predicted dual-localized proteins, including many aminoacyl-tRNA
10 synthetases, in which a leaky AUG start codon is followed by a strong Kozak context in-frame
11 AUG, separated by mitochondrial-targeting sequence. Further analysis shows that such dual-
12 localization is also predicted to be common in *Neurospora crassa*. Kozak-controlled
13 regulation is correlated with insertions in translational initiation factors in fidelity-
14 determining regions that contact the initiator tRNA. Thus, start codon context is a signal that
15 programs the expression and structures of proteins in fungi.

1 Introduction

2 Fungi are important in the fields of ecology, medicine, and biotechnology. With 4 million
3 predicted fungal species, this kingdom is the most diverse of the domain Eukarya. Recent
4 initiatives such as the 1000 Fungal Genomes Project at the Joint Genome Institute, or the
5 Global Catalogue of Microorganisms, which aims to produce 2500 complete fungal genomes
6 in the next 5 years, will result in a deluge of genome sequence data (1, 2). Comparative
7 analysis of coding sequences enables the generation of hypotheses on genome biology and
8 evolution (3-6). However, these analyses intrinsically depend on the quality of the coding
9 gene identification and annotation, which have limitations. First, they depend on automatic
10 sequence comparisons, which limit the identification of clade-specific genes. Second, fungal
11 genes generally contain introns whose positions are difficult to predict based on the genome
12 sequence alone (7). An uncertain intron annotation results in a poor annotation of the
13 coding region extremities, which are generally less evolutionary conserved (8). Third,
14 annotation pipelines only predict plausible open reading frames (ORFs), initially for yeast a
15 contiguous stretch of at least 100 codons starting with an AUG codon and ending with a stop
16 codon (9). These approaches do not reveal which ORFs are translated to protein, and are
17 biased against short ORFs (10).

18 The common assumption is that the first AUG of a fungal ORF is used as the translation start
19 codon, yet non-AUG start codons have been observed in every studied eukaryote, including
20 the fungi *S. cerevisiae*, *C. albicans*, *S. pombe*, and *N. crassa* (11-16). In metazoans the rules
21 for selection of AUG start codons were discovered by Kozak: AUGs are efficiently selected by
22 mammalian cells if they are far (>20nt) from the transcription start site, and have a strong
23 sequence context gccRccAUGG, where the R indicates a purine at the -3 position (17, 18).

1 The influence of the motif on the efficiency of translation is organism-dependent. Although
2 in metazoan cells the Kozak context has a major effect on translation initiation (19), in *S.*
3 *cerevisiae* translation usually starts at the first AUG in the mRNA sequence, the strong
4 aaaAUG sequence context only weakly affects endogenous protein output (20), and AUG
5 sequence context somewhat modulates protein output from reporter mRNAs (21, 22). In
6 fact, translation start codon context has a larger effect on non-AUG start codon usage (23).
7 These “*Saccharomyces* rules” have been considered as the paradigm for fungi, but relevant
8 data are lacking in other species of fungi.

9 Weak or inefficient start codons near the 5’ end of mRNA can give rise to translational
10 regulation, explained by the scanning model of eukaryotic translation initiation. Translation
11 starts by the pre-initiation complex binding mRNA at the 5’ cap and then scanning the
12 transcript leader (TL) sequence until it identifies a start codon, at which translation initiates
13 (19). Here we call the 5’ regulatory region of mRNA the TL rather than the 5’ UTR, because
14 short “upstream” ORFs in this region can be translated (24). The pre-initiation complex
15 sediments at 43S, and comprises the small ribosomal subunit, methionyl initiator tRNA, and
16 numerous eukaryotic translation initiation factors (eIFs). Biochemical, genetic, and structural
17 data indicate that eIF1 and eIF1A associate with the 43S pre-initiation complex (25, 26).
18 Recognition of the start codon involves direct interactions of eIF1 and eIF1A with the start
19 codon context and initiator tRNA within a larger eIF2/3/5-containing 48S pre-initiation
20 complex. Start codon selection occurs when eIF1 is replaced by eIF5’s N-terminus (27), then
21 eIF2 is released, the large ribosomal subunit joins catalyzed by eIF5B, and translation begins
22 (26). This work has been largely driven by studies in *S. cerevisiae* and metazoans. Although
23 the core protein and RNA machinery of eukaryotic translation initiation is highly conserved,

1 it is not understood how fungi quantitatively vary in the sequence, structure, and function of
2 their translation initiation machinery.

3 *Cryptococcus* are basidiomycete yeasts with a high density of introns in their coding genes
4 (28). These introns influence gene expression and genome stability (29-31). The current
5 genome annotation of pathogenic *C. neoformans* and *deneoformans* reference strains are
6 based on both automatic and manual curations of gene structures using RNA-Seq data (32,
7 33). Although the high degree of interspecies conservation of intron numbers and positions
8 within coding sequences suggest that these annotations are reliable (33), the regulatory
9 regions (transcript leader and 3' UTRs) at transcript extremities are less well identified. In
10 fact, most fungal genomes lack complete transcript annotations, and thus we do not know
11 how regulatory structure varies across fungi.

12 In this paper, we experimentally determine the beginning and the end of both coding
13 regions and of transcripts in two *Cryptococcus* species, providing an important genomic
14 resource for the field. Furthermore, our joint analysis of TL sequences and translation
15 identifies a Kozak sequence context that regulates start codon selection, affecting upstream
16 ORF regulation and also alternative protein targeting to mitochondria. Comparison with
17 other fungal genomes revealed that these types of regulation are common in this kingdom:
18 the first AUG of an mRNA or an ORF is not always the major start codon in fungi. These
19 studies demonstrate that, in contrast to the situation in *S. cerevisiae*, start codon sequence
20 context is an important gene regulatory signal that programs the levels and structures of
21 proteins in the fungal kingdom.

22

1

2 **RESULTS**

3 **Delineation of transcript ends in *C. neoformans* and *C. deneoformans***

4 To annotate the extremities of the coding genes in *C. neoformans* and *C. deneoformans*, we
5 mapped the 5' ends (Transcription Start Sites; TSS) with TSS-Seq (34), the 3' ends
6 (Polyadenylation sites; PAS) with QuantSeq 3'mRNA-Seq, and sequenced the same samples
7 with stranded mRNA-Seq. These experiments were done in biological triplicate from cells
8 growing at two temperatures (30°C and 37°C) and two stages of growth (exponential and
9 early stationary phases) with external normalization with spike-in controls.

10 We identified 4.7×10^6 unique TSSs and 6.3×10^4 unique PASs in *C. neoformans*. Clustering of
11 these positions revealed between 27,339 and 42,720 TSS clusters and between 9,217 and
12 16,697 PAS clusters depending on the growth conditions (Table S1). We used the clusters
13 associated with the coding genes to produce an initial annotation, using the most distal TSS
14 and PAS clusters for each gene. The predicted positions which changed the extremities of
15 the genes by more than 100 bp were manually curated (n=1131 and n=286 for the TSS and
16 PAS, respectively). We then selected the most prominent clusters that represented at least
17 10% of the normalized reads count per coding gene in at least one condition. Finally, the
18 most distal of these TL-TSS and 3'UTR-PA clusters were labeled as the 5' and 3' ends of the
19 coding genes for our final annotation (Table S1). For the genes for which no TL-TSS cluster or
20 no 3'UTR-PAS cluster could be identified, we maintained the previous annotation. We used
21 the same strategies for *C. deneoformans* and obtained similar results (Table S1).
22 As expected, most of the TSS clusters (62%) were associated with the TL whereas most of
23 the PAS clusters (82%) were associated with the 3'UTR of the coding genes (Table S1). We
24 analyzed the 3'UTR sequences, confirming the ATGHAH motif associated with the PAS (32).

1 In addition, as previously observed in other systems (35) a (C/T)(A/G)-rich motif was
2 associated with the maxima of these transcription start site clusters. Overall, 89% of the
3 coding genes have both their TL and 3'UTR sequences supported by identified TSS and PAS
4 clusters, respectively.

5 The analysis leads to a scheme of a stereotypical *C. neoformans* coding gene (Figure 1A). In
6 average, it is 2,305 bp long (median 2,008 bp) and contains 5.6 short introns (median 5) in its
7 sequence. As previously reported (28), these introns are short (63.4 nt in average) and
8 associated with conserved consensus motifs. The *C. neoformans* TL and 3'UTR have median
9 lengths of 107 nt and 129 nt, respectively (180 nt and 189 nt, mean; Figure 1A,B). Only 887
10 and 429 genes contain one or more introns in their TL and 3'UTR sequence, respectively;
11 these introns are usually larger (118.3 nt) than those that interrupt the CDS. This gene
12 structure is similar in *C. deneoformans* (Table S1) and there are good correlations between
13 the 3'UTR and TL sizes of the orthologous genes in the two species (Figure 1C,D).

14

15 **More than a third of genes have upstream ORFs that affect translation**

16 The analysis of the TL sequences in *C. neoformans* revealed the presence of 10,286 AUG
17 triplets upstream (uAUG) of the annotated translation start codon (aAUG). We include
18 uAUGs that are either out-of-frame from the start codon, or in-frame but with an
19 intervening stop codon, which are very unlikely to encode a continuous polypeptide.
20 Strikingly, 2,942 genes possess at least one uAUG, representing 43% of the genes with an
21 annotated TL in *C. neoformans* (Figure 1F). A similar result was obtained in *C. deneoformans*,
22 in which we found 10,254 uAUGs in 3,057 genes, and uAUG counts are correlated between
23 the species (Fig 1G).

1 Translation initiation at uAUGs results in the translation of uORFs, which can regulate
2 expression of the main ORF (36, 37). To evaluate the functionality of the uAUGs in
3 *Cryptococcus*, we generated riboprofiling data in both species and compared densities of
4 ribosome-protected fragments with those of sample-matched poly(A)+ RNA. Our
5 riboprofiling data passes quality metrics of 3-nucleotide periodicity of reads on ORFs
6 indicating active translation by ribosomes, and appropriate read lengths of 26-30nt (Figure
7 S2.1).

8 Most genes have ribosome occupancy close to that predicted by their RNA abundance, and
9 restricted to the main ORF, for example the most highly translated gene, translation
10 elongation factor eEF1 α /CNAG_06125 (Figure 2A,B). However, we observed dramatic
11 examples of translation repression associated with uORFs in CNAG_06246, CNAG_03140 in
12 *C. neoformans* (Figure 2A,C,D). These patterns are conserved in their homologs in *C.*
13 *deneoformans* (Figure S2.2). Other spectacularly translationally repressed genes,
14 CNAG_07813 and CNAG_07695 and their *C. deneoformans* homologs (Figure 2A, S2.2A)
15 contain conserved uORFs in addition to 5' introns with alternative splicing or intronically
16 expressed non-coding RNAs (Fig S2.3). In all these cases, high ribosomal occupancy on one or
17 more uORFs is associated with low occupancy of the main ORF.

18 The uncharacterized gene CNAG_06246 has two AUG-encoded uORFs that are occupied by
19 ribosomes, and a predicted C-terminal bZIP DNA-binding domain. This gene structure is
20 reminiscent of the multi-uORF-regulated amino-acid responsive transcription factors
21 Gcn4/Atf4 (37), or the *S. pombe* analog Fil1 (38). The sugar transporter homolog
22 CNAG_03140 has six uAUGs, with substantial ribosome occupancy only at the first.
23 Interestingly, *N. crassa* has a sugar transporter in the same major facilitator superfamily

1 regulated by a uORF (*rco-3/sor-4*, (39)), and sugar-responsive translational repression via
2 uORFs has been observed in plants (40).
3 Since these translationally repressed genes have multiple uAUGs, we investigated the
4 relationship between uAUGs and translation efficiency genome-wide. We observed a clear
5 negative relationship between the number of uAUGs and translation efficiency (Figure 2E,
6 S2.2E), suggesting an uAUG-associated negative regulation of gene expression in both
7 species.

8

9 **Position relative to the TSS affects uAUG translation.**

10 Although some uAUGs are recognized and efficiently used as translation start sites, some
11 others are used poorly or not at all, and allow translation of the main ORF. We thus analyzed
12 *Cryptococcus* uAUG position and sequence context to see how translation start codons are
13 specified in these fungi.

14 We compared the translation efficiency of genes containing only uAUGs close to the TSS to
15 those with uAUGs far from the TSS. In *C. neoformans*, 1,627 of the 10,286 uAUGs are
16 positioned within the first 20 nt of the TL, and 816 uAUG-containing genes have no uAUG
17 after this position. The presence of one or several uAUGs close to the TSS (<20 nt) has nearly
18 no effect on translation efficiency, whereas genes containing uAUGs far from the TSS are less
19 efficiently translated (Figure 2F), and similarly in *C. deneoformans* (Figure S2.2F).

20

21 **A Kozak sequence context determines AUG translation initiation.**

22 To analyze the importance of AUG sequence context for translation initiation in
23 *Cryptococcus* we used the 5% most translated genes (hiTrans $n = 330$) to construct a
24 consensus sequence surrounding their annotated translation start codon (Figure 3A). The

1 context contains a purine at the -3 position, a hallmark of the Kozak consensus sequence
2 (19). However, there is very little enrichment for the +1 position, in contrast with the
3 mammalian Kozak context in which a G is present in +1 ((A/G)CCAUGG) (19). Because of the
4 limited sequence bias downstream of the AUG, and its confounding with signals of N-
5 terminal amino acids and codon usage, we do not consider it further. However, we found a
6 slight sequence bias in the positions -10 to -7 that is outside the metazoan Kozak context.
7 We thus calculated “Kozak scores” for all uAUGs against the position weight matrix (pwm) of
8 the Kozak context from -10 from AUG through to AUG (Figure 3A). We compared the Kozak
9 scores of the annotated AUGs (aAUGs) with those of the 5% most highly translated genes,
10 the first upstream AUG (uAUGs) and the first downstream AUG (d₁AUG). Highly translated
11 aAUGs have a higher score than typical aAUGs, and aAUGs have usually a higher score than
12 the uAUGs and d₁AUGs (Figure 3B). This suggests that the sequence context of the AUG
13 codon might be of importance in the selection of the translation start site in *Cryptococcus*.
14 We next asked if the sequence context of uAUGs affected their ability to repress translation
15 of the annotated ORF, focusing on transcripts with only a single uAUG. For uAUGs close to
16 the TSS, there is no correlation between uAUG Kozak score and the translation efficiency of
17 the aORF; there is a weak negative correlation if the uAUG is far from the TSS (Figure 3D).
18 This indicates that the location of an uAUG impacts its activity. The most striking examples of
19 translational repression in Figure 2 tend to have multiple high-score uAUGs (scores
20 CNAG_06246, u₁AUG 0.93, u₂AUG 0.86; CNAG_03140, u₁AUG 0.85, u₂AUG 0.76;
21 CNAG_07813, u₁AUG 0.79; CNAG_07695, u₁AUG 0.97, u₂AUG 0.90).
22
23 We also asked if the AUG score affects the AUG usage transcriptome-wide, by comparing the
24 difference in u₁AUG and aAUG scores with the ratio in A-site ribosome occupancy in a 10-

1 codon neighbourhood downstream of the u₁AUG and aAUG. We considered the relative
2 occupancy to control for transcript-specific differences in abundance and cap-dependent
3 initiation-complex recruitment. A higher score difference is associated with higher relative
4 ribosome occupancy, while the control comparison with RNA-Seq coverage shows a smaller
5 effect (Figure 3E).

6

7 **Nonsense-mediated decay acts on uORF-containing genes.**

8 An mRNA molecule translated using an uAUG can be recognized as a premature stop codon
9 bearing molecule and will be as such degraded by the nonsense-mediated mRNA decay
10 (NMD) (41). To test this concept, we first sequenced RNA from *C. deneoformans* strains with
11 the conserved NMD factor Upf1 deleted (33), finding 370 genes with increased mRNA
12 abundance and 270 with decreased (Figure 4A, Table S2; 2-fold difference in levels at 1%
13 FDR).

14 We next compared the fold-change in abundance of uORF-containing or uORF-free mRNAs.
15 Two genes with extreme increases in *upf1Δ* are also extremely translationally repressed
16 uORF-containing genes we identified above (Figures 2, S2.1, S2.2): CNF00330 (CNAG_07695
17 homolog, 11-fold) and CNG04240 (CNAG_03140 homolog, 8-fold). Another extreme is the
18 carbamoyl-phosphate synthase CND01051 (5-fold up in *upf1Δ*), a homolog of *S. cerevisiae*
19 *CPA1* and *N. crassa arg-2*. These orthologs are regulated by a conserved uORF encoding a
20 arginine attenuator peptide that have all been verified to repress reporter gene synthesis in
21 a *N. crassa* cell-free translation system (42). Consistent with this model, in both *C.*
22 *neoformans* and *C. deneoformans* the native uORF shows strong ribosome occupancy while
23 the aORF is translationally repressed (*CnTE* = 0.47, *CdTE* = 0.38; Figure S4.1).

1 In general, uORF-containing genes are more likely to be up regulated in the *upf1Δ* mutant
2 than uAUG-free genes (Figure 4B), suggesting that uORFs negatively regulate mRNA
3 abundance in *Cryptococcus*, in addition to repressing translation of the main ORF. Similarly,
4 uORF-containing genes are enriched for NMD-sensitivity only when the uAUG is more than
5 20 nt from the TSS (Figure 4C), suggesting that TSS-proximal uAUGs (< 20 nt) are skipped,
6 and generally not used as translation start codons in *Cryptococcus*.

7 Next, we asked if uAUG Kozak score affects mRNA decay via the NMD pathway. Restricting
8 our analysis to genes with a single uAUG (n=1,421), we binned genes according to their
9 Kozak score. We find that mRNAs that contain uORFs which start with a higher Kozak-score
10 uAUG are more likely to increase in abundance in the *upf1Δ* mutant (Figure 4D). Indeed, the
11 abundance increase is monotonically correlated with the mean of the score bins.

12 In conclusion, in *Cryptococcus*, the position and the sequence context of uAUGs determines
13 their usage as translation start codons, and the effect of the uORF on stimulating nonsense-
14 mediated decay of the mRNA.

15

16 **Start codon sequence context and uORF regulation in other fungi**

17 We then examined sequences associated with translation start codons in other fungi, for
18 which both RNA-Seq and Riboprofiling data were available, and for which the annotation
19 was sufficiently complete (i.e. *S. cerevisiae*; *Neurospora crassa*, *Candida albicans* and
20 *Schizosacchomyces pombe*). We analyzed the Kozak context associated with aAUG of all
21 annotated coding genes, of the 5% most translated genes (hiTrans), and for mRNAs
22 encoding cytoplasmic ribosomal proteins (CytoRibo), as a model group of highly expressed
23 and co-regulated genes defined by homology (Table S3). Cytoplasmic ribosomal proteins
24 have informative Kozak contexts, with a strong A-enrichment at the positions -1 to-3 and

1 weak sequence enrichment after the AUG in all these species (Figure 5A). The total
2 information content of the Kozak sequence is higher for CytoRibo genes than HiTrans, and
3 higher for HiTrans than all annotated genes, in all these fungi (Figure 5B). Nevertheless,
4 these contexts have also some species specificity: Kozak sequences for HiTrans and CytoRibo
5 are more informative in *Cryptococcus* and *N. crassa* than in *S. pombe*, *C. albicans*, and *S.*
6 *cerevisiae*. In particular, the C-enrichment at positions -1, -2 and -5 in *Cryptococcus* and *N.*
7 *crassa* is absent in *S. cerevisiae*, and we observed no sequence enrichment upstream of the -
8 4 position for *S. pombe* and very little for *S. cerevisiae*. In contrast, a -8 C enrichment, similar
9 to the *Cryptococcus* pattern, was observed in *N. crassa*. The -10 -6 A/T rich region for *C.*
10 *albicans* is likely to reflect an overall A/T-richness of the TLs in *C. albicans*.

11 The analysis of the TL sequences from these fungi, excluding *C. albicans* for which no TL
12 annotation is available, also shows species specificity. The average TL length in *S. cerevisiae*
13 (84 nt) is less than half that in *Cryptococcus* (Figure 5C). In sharp contrast with *Cryptococcus*,
14 only 985 uAUGs are present in 504 genes, which correspond to 18% of the genes with an
15 annotated TL in *S. cerevisiae*. Moreover, the density of the uAUGs is very low and uAUGs
16 have no global effect on TE in this yeast (Figures 5D,E, Fig S5.1).

17 More broadly, short TLs with very low uAUG density are more the exception than the rule in
18 the fungal kingdom (Figure 5C). Most fungi have large number of uAUGs in their TL and
19 these uAUGs globally down regulate gene expression (Figure S5.1). Our analysis in
20 *Cryptococcus*, together with the strength and quality of the Kozak contexts associated with
21 the aAUG, suggest that fungi in general are able to discriminate between AUGs with
22 different sequence context, and thus are using AUG sequence context to regulate gene
23 expression at the post transcriptional level.

24

1 **Kozak context controls leaky scanning in *Cryptococcus***

2 We earlier calculated the Kozak score of the first downstream AUG (d_1 AUG) within each CDS:
3 these d_1 AUG scores are mostly lower than the score of the aAUGs (Figure 3B), consistent
4 with most annotations correctly identifying a good-context AUG as the start codon. Yet, we
5 identified number of d_1 AUGs with a high score ($n=1109$ above 0.826, the median Kozak score
6 for aAUGs; $n= 131$ above 0.926, the median for hiTrans), which could be efficiently used as
7 translation start codon. The scanning model of translation initiation predicts that the d_1 AUG
8 will be used as the start codon only if pre-initiation complexes leak past the aAUG. Thus if
9 the aAUG has a strong sequence context, most ribosomes will start translation there and the
10 d_1 AUG will be not used as a translation start codon.

11 To identify potential leaky translation initiation events, we compared the aAUG and d_1 AUG
12 scores for each of the 50% most expressed genes (Figure 6A). For above-median aAUG score
13 genes, the score of the d_1 AUGs can be very high or very low. By contrast, for the genes with
14 a low aAUG score, there is a bias toward higher d_1 AUG score, suggesting that for these genes
15 the strong d_1 AUG could be used as alternative translation start site (Figure 6A).

16 To test whether AUG score affects translation initiation, we calculated the ratio of ribosome
17 protected fragment density and RNA-Seq density around each aAUG and d_1 AUG, and the
18 difference in score between these two AUGs (Figure 6B). We found a weak positive
19 correlation between the difference in scores of the two AUGs and RNA-Seq density at these
20 specific loci, raising the possibility that transcription start sites sometimes occur downstream
21 of a weak aAUG. By contrast, the relative ribosome density sharply increases at the d_1 AUGs
22 when the difference in score between the d_1 AUG and the aAUGs increases. This suggests

1 that for these genes, both AUGs can be used as translation start codon, as because a subset
2 of scanning ribosomes leak past the aAUG and initiate at the d₁AUG.

3

4 **Kozak context-controlled scanning specifies alternative N-termini in *Cryptococcus* and** 5 ***Neurospora***

6 We next determined which groups of genes could be affected by potential alternative start
7 codon usage. We focused our analysis on the 50% most expressed genes for which the
8 difference in score between the aAUG and d₁AUG was the highest (difference in wide score
9 d₁AUG – aAUG > 0.1, n = 167 for *C. neoformans*) (Table S4). Strikingly, for 66% of these genes
10 (110/167) the d₁AUG is in frame with the corresponding aAUG, with a median of 69 nt (mean
11 79 nt) between the two AUGs. Thus, alternative usage of in-frame AUGs would result in
12 proteins with different N-terminal ends. Supporting this hypothesis, 37% of these proteins
13 (41/110) possess a predicted mitochondrial targeting sequence located between the two
14 AUGs, far exceeding the 8% genome-wide (560/6788). This suggests that the usage of the
15 annotated start codon would target the isoform to mitochondria, whereas the usage of the
16 d₁AUG would produce a protein specific to the cytoplasm or another organelle. Examples of
17 alternative localization driven by alternative N-termini have been observed across
18 eukaryotes (43).

19 The pattern of predicted dual-localization, i.e. enrichment of high-score d₁AUGs in-frame
20 with predicted mitochondrial localization signal on the longer N-terminal, is conserved in
21 some fungi but not others (Figure 6C). In a null model where coding sequences have random
22 nucleotide content, we would expect roughly 1/3 of d₁AUGs to be in frame. In 6 fungal
23 species we examined, for d₁AUGs whose score is comparable to or less than the aAUG they

1 follow, the proportion in frame is close to (*Cryptococcus*, *N. crassa*) or less than 1/3. These
2 proportions are similar when we considered reasonably expressed (above-median) or low
3 expressed genes. The pattern differs for proteins with a d₁AUGs whose score high relative to
4 the aAUG they follow (d₁AUG score > aAUG score + 0.1). In *Cryptococcus* and *N. crassa*, most
5 reasonably expressed mRNAs are in-frame and over 1/3 of these in-frame high-score d₁AUGs
6 have predicted mitochondrial localization. In *S. cerevisiae* and *C. albicans*, a relative slight
7 enrichment for in-frame d₁AUGs and for protein possessing a mitochondrial targeting
8 sequence for high-scoring d₁AUGs can be also observed. By contrast, in *S. pombe* we see
9 depletion in the in-frame/out-of-frame ratio, even in these proteins with high-scoring
10 d₁AUGs.

11 These results suggest that the extent to which alternate translation start codons regulate
12 proteome diversity is variable in fungi. Accordingly, we identified a number of *Cryptococcus*
13 proteins with potential alternative start codons and N-terminal targeting sequences, whose
14 two homologs in *S. cerevisiae* are known to be necessary in two compartments of the cells.
15 For instance, *CnPUS1/CNAG_06353* is an homolog of both the mitochondrial and
16 cytoplasmic tRNA:pseudouridine synthases encoded by the *PUS1* and *PUS2* paralogs in *S.*
17 *cerevisiae*. In *C. neoformans*, ribosome occupancy at both the aAUG and d₁AUG of
18 CNAG_06353, and the presence of transcription start sites both sides of the aAUG (Figure
19 6D), argues that both AUGs are used as start codons, and transcription and translation
20 regulation could co-operate to set isoform levels. Similarly, *CnGLO1/CNAG_04219* encodes
21 both the cytoplasmic and nuclear isoforms of the glyoxalase I depending of the alternate
22 AUG usage (Figure S6.1B). The next enzyme in this pathway, Glyoxalase II, is likewise
23 encoded by *CnGLO2/CNAG_01128*, which is a homolog of both cytoplasmic (Glo2) and
24 mitochondrial (Glo4) enzymes in *S. cerevisiae*. CNAG_01128 has a weak aAUG, strong d₁AUG,

1 and N-terminal predicted mitochondrial targeting sequence (Figure S6.1C). Finally, we
2 observed that nine members of the amino-acyl tRNA synthetase gene family have predicted
3 alternate localization from alternate AUG start codons.

4

5 **Amino-acyl tRNA synthetases (aaRSs) are frequently single-copy and dual-localized in**

6 ***Cryptococcus***

7 The tRNA charging activity of aaRSs is essential in both cytosol and mitochondria to support
8 translation in each compartment, and examples of alternative localization of two aaRS
9 isoforms of a single gene have been observed in fungi, plants, and animals (44-46). This
10 implies that a eukaryote with a single genomic homolog of an aaRS is likely to make distinct
11 localized isoforms from that locus. Thus, we examined predicted aaRS localization in fungi.
12 We assembled gene lists of aaRSs in diverse fungi from homology databases OrthoDB (47)
13 and PANTHERdb (48), adding a mitochondrial SerRS (CNAG_06763/CNB00380) to the list of
14 *Cryptococcus* aaRSs analysed by Datt and Sharma (49).

15 In *C. neoformans* and *C. deneoformans*, eleven aaRSs are each expressed from a single
16 genomic locus, including the homologs of all five *S. cerevisiae* aaRSs whose dual-localization
17 has been verified (Table S5). Nine of these *Cryptococcus* aaRSs have the same structure of a
18 poor-context annotated AUG followed by a predicted mitochondrial targeting sequence and
19 a strong-context d₁AUG (Figure 7A/B; AlaRS, CysRS, GlyRS, HisRS, ValRS, LysRS, ProRS, ThrRS,
20 TrpRS). The similar annotated AUG contexts, sharing an unfavourable -3U, suggests that the
21 same mechanism could lead to leaky translation initiation at most of these (Table S6). At the
22 downstream AUGs, the strong Kozak context is consistent with efficient translation initiation
23 of the cytoplasmic isoform from this start codon (Table S6).

1 The two remaining single-copy aaRSs have near-AUG translation initiation sites upstream of
2 predicted mitochondrial targeting sequences. Translation of ArgRS starts at an AUU codon
3 with otherwise strong context (cccaccAUU) conserved in both *Cryptococcus* species. This N-
4 terminal extension includes a predicted mitochondrial targeting sequence (mitofates $p >$
5 0.95 for both species). Translation of LeuRS starts at adjacent ACG and AUU codons which
6 collectively provide strong initiation context (gccaccACGAUU in *C. neoformans*, gccACGAUU
7 in *C. deneoformans*). This N-terminal extension also includes a predicted mitochondrial
8 targeting sequence (mitofates $p \approx 0.7$ for both species).

9 In *Cryptococcus*, alternative aaRS isoforms appear to be mostly generated by alternative
10 translation from a single transcript, and sometimes by alternative transcription start sites.
11 On all the predicted dual-localized aaRSs, we observe ribosomal occupancy starting at the
12 earliest start codon (Figure 7C/D, Fig S7.1). LysRS/CNAG_04179 contains only a single cluster
13 of transcription start sites, upstream of the aAUG (Figure 7C). ProRS/CNAG_04082 contains a
14 wider bimodal cluster of TSSs, both upstream of the aAUG. Similarly, most transcription
15 initiation is well upstream of the aAUG in CysRS/CNAG_06713, LeuRS/CNAG_06123,
16 ThrRS/CNAG_06755, and ValRS/CNAG_07473. However, for GlyRS/CNAG_05900, and
17 HisRS/CNAG_01544, we observe alternative transcription start sites closely upstream of the
18 annotated start codon, that are likely to affect the efficiency of start codon usage. In
19 ArgRS/CNAG_03457 there is also an alternative transcription start site, close to the near-
20 AUG start codon for the mitochondrial form. In AlaRS/CNAG_05722 and TrpRS/CNAG_04604
21 we detect some transcription start sites between the alternative start codons, and TrpRS
22 also has an uORF in the transcript leader that is likely to affect translation. These
23 observations suggest that dual-localization of the single-copy aaRSs in *Cryptococcus* is

1 regulated largely by start codon choice. For some genes this regulation is backed up by
2 alternative TSS usage.
3 Some dual-localized genes use an upstream near cognate codon (DualNCC) in all these fungi,
4 but the NCC-initiated aaRS are not the same from one fungus to the other. For instance,
5 both *Cryptococcus* and *N. crassa* AlaRS use DualAUG whereas in *S. pombe*, *S. cerevisiae* and
6 *C. albicans* a DualNCC is used. On the other hand, *S. pombe* GlyRS is regulated by DualNCC
7 whereas the other ones use a DualAUG regulation. Substitution between weak AUG codons
8 and near-cognate codons seems thus to have taken place multiple times in the fungal
9 kingdom.

10

11 **Amino-acyl tRNA synthetases as an evolutionary case study**

12 To understand patterns of dual-localization, we next examined the evolution of aaRSs. The
13 ancestral eukaryote is thought to have had two complete sets of aaRS, one mitochondrial
14 and one cytoplasmic, but all mitochondrial aaRSs have been captured by the nuclear
15 genome and many have been lost (50). Thus we examined aaRS phylogenetic trees in more
16 detail. For some amino acids (Asn, Asp, Glu, Iso, Met, Phe, Ser, Tyr), reference fungi have
17 distinct cytoplasmic and mitochondrial aaRSs that cluster in separate trees (51). We also do
18 not consider Gln, because organellar Gln-tRNA charging in some species is achieved by an
19 indirect pathway (52).

20 Dual-localized AlaRS, CysRS, and HisRS in the 6 fungi we focus on are each monophyletic
21 (51). Even these aaRS can be encoded by two genes in some other fungi: AlaRS is duplicated
22 to one exclusively mitochondrial and another exclusively cytoplasmic gene in the
23 Saccharomycete yeast *Vanderwaltozyma polyspora* (53). For CysRS, *Aspergillus versicolor*
24 (ASPVEDRAFT_141527 and ASPVEDRAFT_46520) and *Coprinus cinerea* (CC1G_03242 and

1 CC1G_14214) have two copies, one of which has a predicted mitochondrial targeting
2 sequence. For HisRS, *Rhizopus delemar* (RO3G_01784 and RO3G_16958) and *Phycomyces*
3 *blakesleanus* (PHYBL_135135 and PHYBL_138952) likewise contain gene duplications.
4 Similarly, *S. cerevisiae* has two ArgRS genes that arose from the whole-genome duplication:
5 *RRS1/YDR341C* is essential, abundant, and inferred to be cytoplasmic (54) while
6 *MSR1/YHR091C* has a mitochondrial localization sequence and MSR1 deletions have a petite
7 phenotype (55), although both have been detected in mitochondria suggesting some
8 residual dual-localization of the cytoplasmic enzyme (56). The second *S. cerevisiae* stress-
9 responsive cytoplasmic copy of GlyRS also arose from the whole-genome duplication (57). *S.*
10 *pombe* cytoplasmic ValRS is monophyletic with dual-localized ValRS in other fungi, and
11 *Schizosaccharomyces* also has a paralogous but diverged mitochondrial ValRS that appears
12 to be descended from an early eukaryotic ValRS of mitochondrial origin (58).

13 LysRS appears to have been duplicated in an ancestor of ascomycetes: ascomycete
14 mitochondrial homologs cluster together, and ascomycete cytoplasmic homologs cluster
15 together, while the single basidiomycete homolog clusters close to the base of this split from
16 other opisthokonts (51). By contrast, LeuRS, ProRS, and TrpRS are each represented by two
17 distinct proteins in ascomycetes, one cytoplasmic and one mitochondrial and of
18 independent descent, but the mitochondrial homolog has been lost in *Cryptococcus* species.
19 In basidiomycetes *Ustilago* and *Puccinia*, homologs of mitochondrial LeuRS and ProRS are
20 not present, but there is a homolog of mitochondrial TrpRS; all these have a single homolog
21 of the cytoplasmic TrpRS (51). Our independent phylogenetic analysis of LysRS and ProRS
22 agrees with the conclusions from PANTHERdb (Figures 7E/F). These analyses show that
23 aaRSs have undergone multiple incidences of at least two processes during fungal evolution:

1 losses associated with the dual-localization of the remaining gene, and duplications followed
2 by specialization.

3

4 **Evolutionary conservation of gene-specific feedback regulation by alternate AUG usage**

5 We also observed striking examples of gene-specific regulation by start codon context in
6 *Cryptococcus*, in translation factors affecting start codon selection, supporting previously
7 proposed models of feedback regulation (59, 60).

8 Translation initiation factor eIF1, which enforces selection of strong context start codons, is
9 encoded by an mRNA with poor start codon context in diverse eukaryotes, driving an
10 autoregulatory program (59, 61). In *C. neoformans*, eIF1 also initiates from a poor-context
11 cuuaguugaAUG start (score 0.75), and ribosome profiling reads are spread across the
12 annotated ORF (Figure 8A). Intriguingly, the next AUG is out-of frame and has strong context
13 cuccaaaaAUG (score 0.98), with a same-frame stop codon 35 codons later, suggesting that
14 this could represent a downstream short ORF that captures ribosomes that have leaked past
15 the poor-context start. To test this hypothesis, we examined the 5' ends of riboprofiling
16 reads, which report on the translation frame of the ribosomes (36). Riboprofiling reads from
17 the 5' and 3' of the eIF1 annotated ORF are roughly 77% in frame 0, 10% in +1, and 13% in
18 +2, as are reads on two other highly expressed genes, eEF1 α and *HSP90*. By contrast, in the
19 hypothesized downstream ORF, reads are only 57% in frame 0, 32% in frame +1, and 11% in
20 frame +2, consistent with translation occurring in both frame 0 and +1. The gene structure is
21 conserved in *C. deneoformans*, with a weak aAUG (score 0.76), a strong d₁AUG (score 0.98)
22 in the +1 frame, followed by an enrichment in +1-frame riboprofiling reads (Figure S8.1A,B).
23 We observe small increases in eIF1 mRNA levels in the *upf1* Δ strain of *C. deneoformans* at
24 both 30°C (1.16x, p=0.04) and 37°C (1.09x), so NMD could regulate this transcript. Overall,

1 our data support the hypothesis that the downstream ORF of eIF1 is translated after read-
2 through of the annotated AUG, and that the downstream ORF contributes to translation
3 regulation of the annotated ORF.

4 Translation initiation factor eIF5 reduces the stringency of strong-context start-codon
5 selection, and is encoded by an mRNA with a repressive uORF initiated from a poor-context
6 uAUG in diverse eukaryotes (60). In *Cryptococcus*, eIF5 (*TIF5/CNAG_01709*) also contains a
7 uAUG with the poor sequence context aaagaguucAUG (score 0.72), while the main ORF of
8 eIF5 is initiated by a strong context cccgcaaaAUG (score 0.94). We detect ribosomal density
9 on the uORF of *TIF5* comparable to that on the main ORF (Figure 8C), suggesting there is
10 substantial translation initiation at the uAUG. There is also clear translation initiation at a
11 further upstream near-cognate CUG codon. The gene structure is conserved in *C.*
12 *deneoformans* *TIF5*, with the same pattern of riboprofiles at upstream poor-context AUG
13 and near-cognate codons (Figure S8). Further, the *C. deneoformans* homolog transcript
14 abundance increases substantially in the *upf1Δ* strain (2.6x, $p < 10^{-50}$). This supports the
15 model that eIF5 translation is repressed by upstream reading frames initiated from poor
16 start codons, leading to nonsense-mediated decay of the transcript.

17 These examples further illustrate that the first start codon is not always used, but rather
18 start codon usage is driven by the sequence context, and that variability in start codon
19 context including the canonical AUG sequence is used for translational regulation.

20

21 **Variable inserts in eTIFs correlate with variation in translation initiation determinants**

22 The conserved proteins eIF1, eIF5, and eIF1A play pivotal roles in start codon selection in *S.*
23 *cerevisiae*; specific mutations in these factors give rise to suppressor of upstream initiation
24 codon (Sui-) phenotypes and their suppressors (Ssu-) (61). To ask if between-species

1 variability in start codon preference is linked to these initiation factors, we generated
2 multiple sequence alignments of their homologs in fungi.
3 Translation initiation factor eIF1 shows striking sequence variation across fungi, notably at
4 multiple *Cryptococcus*-specific sequence insertions that result in a 159-aa protein
5 substantially larger than the 108-aa *S. cerevisiae* homolog (Figure 9A). Variation in eIF1
6 occurs at and around positions known to modulate start codon selection in *S. cerevisiae* (61).
7 For instance, a T15A substitution increases fidelity in ScelF1 (61), and an analogous T15A
8 substitution is present in eIF1s from *Neurospora* and other filamentous fungi, while both
9 *Cryptococcus* homologs have the T15V substitution. The three fungi that tend not to use
10 alternative AUG start codons in the regulation of proteome diversity, *S. cerevisiae*, *C.*
11 *albicans*, and *S. pombe*, all have a threonine residue at position 15. Variation in fungal eIF1
12 extends far beyond this N-terminal region: similar patterns of sequence diversity occur at
13 the positions E48, L51, D61 that have been shown to increase fidelity in ScelF1 (61). By
14 contrast, positions K56, K59, D83, Q84, at which mutations have been shown to reduce
15 fidelity in ScelF1 (61), are highly conserved in fungi.
16 We next tested how the translation pre-initiation complex could be affected by the
17 insertions in *Cryptococcus* eIF1 using published structures of the *S. cerevisiae*/*K. lactis* “Pin”
18 complex engaged in the act of AUG selection (62). We found that the insertions in eIF1 are
19 facing either the methionine initiator tRNA (tRNAⁱ) or the solvent-exposed side (Figure 9B).
20 The N-terminal insertion is not visible in the structure, but could be close to the acceptor
21 arm of tRNAⁱ. The N-proximal loop insertion of CnelF1 extends from the ScelF1 sequence
22 (18-DETATSNY-25) that contacts the acceptor arm of tRNAⁱ. The CnelF1 insertion in loop 2
23 extends the ScelF1 loop 2 (70-KDPEMGE-76) that contacts the D-loop of tRNAⁱ; substitutions
24 D71A/R and M74A/R increase the charge of ScelF1 loop 2 and increase initiation at UUG

1 codons and weak AUG codons (63). *CnelF1* loop 2 has substitutions at both these
2 functionally important sites, and is extended by a further 14 hydrophobic and negative
3 residues. The last insertion in *CnelF1* extends a loop facing the solvent-exposed surface of
4 *ScelF1*. Collectively, this shows that there are likely major differences in the eIF1-tRNA_i
5 interaction surface in *Cryptococcus* relative to other fungi, an interaction critical for start
6 codon selection (63).

7 The N-terminal domain of eIF5 (eIF5-NTD) replaces eIF1 upon start codon recognition, and
8 we found between-species variation in *CnelF5* at tRNA_i interaction surfaces corresponding to
9 variability in *CnelF1* (Figure 9C, S9.1A). *ScelF5* Lys71 and Arg73 in loop 2 make more
10 favourable contacts with the tRNA_i than the corresponding residues of *ScelF1*, so that the
11 shorter loop 2 of *ScelF5* may allow the tRNA_i to tilt more towards the 40S subunit (27).

12 Although Arg73 is conserved across fungi, Lys71 is absent in *CnelF5* loop 2 (67-SMAN-70),
13 which is two amino acids shorter than *ScelF5* loop 2 (66-SISVDK-71). Collectively, the longer
14 loop 2 of *CnelF1* and the shorter loop 2 of *CnelF5* suggest that the conformational changes
15 accompanying start codon recognition may be more exaggerated in *Cryptococcus*, providing
16 a mechanistic hypothesis for stronger genomic patterns of start codon recognition.

17 Fungal eIF1A homologs also diverge from *ScelF1A* at regions that modulate translation
18 initiation fidelity (Figure S9.1B), for example the N-terminal element DSDGP (61). The
19 *Cryptococcus* eIF1A C-terminus is diverged from all other fungi at *ScelF1A* positions 110-120,
20 and along with other basidiomycetes lacks a loop at *ScelF1A* positions 135-149. This C-
21 terminal region of *ScelF1A* contributes to pre-initiation complex assembly and binds eIF5B
22 (64) and eIF5 (65), and domain deletions or local alanine substitutions reduce fidelity of
23 translation start site selection (61, 64, 66).

1 Thus, although structural analysis of the Cryptococcal initiation complex will be required for
2 a detailed mechanistic understanding, our initial analysis suggests that sequence variability
3 in fungal eIFs could plausibly account for differences in start codon selection between
4 different species.

5

1 **DISCUSSION**

2

3 Our annotation of transcript structure and translation in two pathogenic *Cryptococcus*
4 species and our analysis of published data from other species show that start codon context
5 has a major effect on protein production, regulation, diversity, and localization in the fungal
6 kingdom. As such this work represents a useful resource for the field. While the genome-
7 wide effect of start codon context is weak in *S. cerevisiae* (20), we find that other fungi, from
8 *Neurospora* to *Cryptococcus*, use start codon context to regulate translation initiation to a
9 far greater extent. These fungi have long and AUG-rich TLs, and more information-rich and
10 functionally important Kozak sequences. Further, *Cryptococcus* and *Neurospora* display
11 extensive evidence of leaky scanning of weak AUG codons that is used for regulation by
12 upstream ORFs and to generate alternate N-terminal isoforms with different subcellular
13 localization.

14

15 **Widespread leaky scanning controlled by start codon context in *C. neoformans***

16 Translation initiation regulation can be enabled by start codons that are imperfectly used, so
17 that scanning pre-initiation complexes can leak past them. According to the scanning model
18 of translation initiation, a “perfect” strong start codon would prevent this by capturing all
19 the scanning PICs, and leave none for regulatable downstream initiation. For example, the
20 downstream out-of-frame ORF of *Cryptococcus* eIF1 is likely to be translated only by PICs
21 that leak past the annotated AUG. The alternative second in-frame AUG of dual-localized
22 proteins is also initiated only by PICs that have leaked past the initial AUG. Our data show
23 this leakiness-driven dual-localization is common in *Cryptococcus*, in addition to being
24 conserved across eukaryotes in gene classes such as tRNA synthetases. Our data also argue

1 that AUGs that are proximal to the 5' cap, or that have poor sequence context, are
2 commonly leaked past in *Cryptococcus*, as shown previously in studies of yeast (67) and
3 mammals (17, 68). We note that leakiness-driven translation regulation is not the only
4 mechanism regulating alternative translation from a single mRNA and is distinct from those
5 that depend on either blocking scanning, or on recycling of post-termination ribosomes such
6 as in the case of *S. cerevisiae GCN4* (37).

7

8 **Functional role of start codon context varies across the fungal kingdom**

9 *Cryptococcus* and *Neurospora* have long TLs that are AUG-rich, and extended start codon
10 context sequences that suggest a higher ability to discriminate against poor-context AUGs.
11 Several lines of evidence argue that efficiency with which upstream AUGs capture initiation
12 complexes is determined by the AUG sequence context. The most spectacular examples of
13 uORF-associated translation repression in *Cryptococcus* are associated with good-context
14 uAUGs with high ribosome occupancy. However, such strong-context high-occupancy uAUGs
15 are rare. In *Cryptococcus* and *Neurospora*, the leakiness of potential AUG translation start
16 sites is also extensively used to diversity the proteome by alternative N-terminal formation.
17 In comparison, *S. cerevisiae*, *S. pombe* and *C. albicans* appear to be less efficient in
18 discriminating AUGs based on their sequence context. *S. cerevisiae* has minimized the
19 possibility of regulation of gene expression by uORFs: it has unusually short TLs, these TLs
20 are unusually AUG-poor, uAUGs tend to have poor context, and there is no statistical
21 association between uAUG score and translation efficiency of the main ORF. Reporter gene
22 studies (21, 22) and classic examples such as *GCN4* show that uAUGs can repress translation
23 in *S. cerevisiae*, but genome-wide analysis show that this is rare during exponential growth in
24 rich media (Fig S5.1). Recent work on meiosis (69) and stress (70) shows that 5'-extended

1 transcript leaders that contain repressive uAUGs (“long undecoded transcript isoforms”) are
2 more common during alternative growth conditions for this yeast. Moreover, in *S. cerevisiae*,
3 near-cognate codons appear to be more common starts for alternative N-terminal formation
4 (71). This suggests that leaky scanning from near-cognate codons, more than from AUGs,
5 might be an important mode of regulation in *S. cerevisiae*. The situation is different in *S.*
6 *pombe*, which has long AUG-rich TLs but is depleted for downstream in-frame AUGs.
7 Consequently, uAUGs globally repress gene expression, but do not appear to regulate
8 alternative protein production through alternative AUG start codons. We speculate that the
9 comparatively uninformative Kozak context in *S. pombe* might be variable enough to
10 regulate translation initiation rate but not proteome diversity.
11 We found that multiple near-cognate start codons are used for leaky initiation in
12 *Cryptococcus*: ACG for the mitochondrial isoform of LeuRS, AUU for the mitochondrial
13 isoform of ArgRS, and the upstream CUG in eIF5. Further work will be needed to quantify the
14 extent of near-cognate start codon usage in *Cryptococcus* in different growth conditions and
15 compare it to other organisms (14, 72).

16

17 **Leaky scanning through weak AUGs could regulate the mitochondrial proteome**

18 We computationally predicted dozens of dual-localized proteins with alternative start
19 codons that confer an N-terminal mitochondrial targeting sequence in their longest isoform.
20 We did not identify enrichment of proteins with predicted dual-localization in the cytoplasm
21 and in the nucleus, or with a signal peptide followed by an alternative start codon (data not
22 shown). Thus, increasing the efficiency of weak-context to strong-context translation
23 initiation would predominantly upregulate a regulon consisting of the mitochondrial
24 isoforms of dozens of proteins.

1 Mechanisms to control initiation efficiency of a mitochondrial-localized regulon could
2 include intracellular magnesium concentration (73), variations in availability or modification
3 status of shared initiation factors, variations of the ratio of mitochondrial volume to
4 intracellular volume (74), or specialized factors to promote initiation specifically of
5 mitochondrial isoforms with their specialized start codon context. Nakagawa et al (75)
6 previously suggested that distinct Kozak contexts might be recognized by different molecular
7 mechanisms.

8 One candidate mechanism involves the translation initiation factor 3 complex, which has a
9 role in regulating the translation initiation of mitochondrial-localized proteins across
10 eukaryotes. In *S. pombe*, subunits eIF3d/e promote the synthesis of mitochondrial electron
11 transfer chain proteins through a TL-mediated mechanism (76). In *S. cerevisiae* and
12 *Dictyostelium discoideum* the conserved eIF3-associated Clu1/CluA protein affects
13 mitochondrial morphology (77), and the mammalian homolog CLUH binds and regulates
14 mRNAs of nuclear-encoded mitochondrial proteins (78, 79). Metazoans have 12 stably-
15 associated subunits of eIF3, which are conserved in most fungi including *N. crassa* (80),
16 *Cryptococcus*, and the Saccharomycetale yeast *Yarrowia lipolytica* (Table S7). Interestingly,
17 species that tend not to use alternate AUG codons for dual-localization have lost eIF3
18 subunits: eIF3d/e/k/l/m are lost in *C. albicans*, and additionally eIF3f/h in the related *S.*
19 *cerevisiae*; *S. pombe* has independently lost eIF3k/l (Table S7; (51)). Further work will be
20 needed to investigate the role of eIF3 in regulating mitochondrial- and dual-localized
21 proteins in the fungal kingdom.

22

23 **Evolutionary plasticity of translational initiation in the fungal kingdom**

1 Selection on genome compaction in unicellular yeasts, which has independently led to gene
2 loss and high gene density in multiple lineages of yeast, could lead to shorter TLs. However,
3 *Saccharomyces*, *Schizosaccharomyces*, and *Cryptococcus* have all independently evolved
4 yeast lifestyles with compact genomes, yet their average TL lengths differ three-fold.
5 Mutations in gene expression machinery, such as the variation in eIF1 noted above, would
6 alter selective pressure on start codon context, and thus uAUG density. Cells have multiple
7 redundant quality control mechanisms, and flexible protein production through leaky
8 scanning could be buffered by such mechanisms enabling their evolution. Key control
9 mechanisms acting on mRNA, such as RNAi and polyuridylation, have been lost in fungal
10 lineages such as *Saccharomyces*, which might explain their more ‘hard-wired’ mechanism of
11 translation initiation.

12

13 Unexpectedly, highly conserved core translation initiation factors, such as eIF1, have
14 distinctive sequence inserts in *Cryptococcus* that are not shared even by basidiomycetes
15 such as *Puccinia* and *Ustilago*. One possibility is genetic conflict, as genetic parasites hijack
16 the gene expression machinery (81). Thus, the unique aspects of the *Cryptococcus*
17 translation initiation machinery could have arisen from a past genetic conflict in which rapid
18 evolution of initiation factors in an ancestor enabled evasion of a genomic parasite (e.g. a
19 mycovirus) that would otherwise hijack initiation.

20

21 **Data Availability**

22 Raw and summarized sequencing data are available on GEO under accession numbers
23 GSE133695 (RNA-seq, TSS-seq, PAS-seq) and [GSE133125](#) (ribosome profiling and
24 matched RNA-seq).

25 **Acknowledgments**

1 We thank members of the Wallace, Janbon, and Madhani labs for helpful discussions and
2 comments on the manuscript. We thank Juan Mata for sharing intermediate data related
3 to (Duncan and Mata 2017). We are grateful to J. Weissman (UCSF) for advice on
4 ribosome profiling. Work in the Madhani lab is supported by grants from the US
5 National Institutes of Health. H.D.M. is an Investigator of the Chan-Zuckerberg Biohub.
6 E.W.J.W. is a Sir Henry Dale Fellow, supported by a Sir Henry Dale Fellowship jointly
7 funded by the Wellcome Trust and the Royal Society (Grant Number 208779/Z/17/Z). L.T.
8 is supported by a Wellcome-University of Edinburgh ISSF3 award.

9

10 **Material and Methods**

11 **DNA and RNA purification, sequencing library preparation**

12 *C. neoformans* strain H99 and *C. deneoformans* strain JEC21 and were grown in YPD at 30°C
13 or 37°C under agitation up to exponential or early stationary phase as previously described
14 (32). Briefly, early stationary phase was obtained after 18 h of growth (final OD₆₀₀ = 15)
15 starting from at OD₆₀₀ = 0.5. *C. deneoformans* strain NE579 (*upf1Δ*) (33) was grown in YPD at
16 30°C under agitation in exponential phase. Each *Cryptococcus* cell preparation was spiked in
17 with one tenth (OD/OD) of *S. cerevisiae* strain FY834 (82) cells grown in YPD at 30°C in
18 stationary phase. Cells were washed, snap frozen and used to prepare RNA and total DNA
19 samples as previously described (32). Each condition was used to prepare biological triplicate
20 samples.

21 For RNA-Seq, strand-specific, paired-end cDNA libraries were prepared from 10 µg of total
22 RNA by polyA selection using the TruSeq Stranded mRNA kit (Illumina) according to
23 manufacturer's instructions. cDNA fragments of ~400 bp were purified from each library and
24 confirmed for quality by Bioanalyzer (Agilent). DNA-Seq libraries were prepared using the kit
25 TruSeq DNA PCR-Free (Illumina). Then, 100 bases were sequenced from both ends using an
26 Illumina HiSeq2500 instrument according to the manufacturer's instructions (Illumina).

1 TSS-Seq libraries preparations were performed starting with 75 µg of total RNA as previously
2 described (34) replacing the TAP enzyme by the Cap-clip Pyrophosphatase Acid (TebuBio).
3 For each *Cryptococcus* species we also constructed a control “no decap” library.
4 Briefly, for these control libraries, poly A RNAs were purified from 75 µg of RNA from
5 *Cryptococcus* and 75 µg of RNA *S. cerevisiae* before being dephosphorylated using Antarctic
6 phosphatase. Then, *S. cerevisiae* RNAs and one half of the RNAs extracted from *Cryptococcus*
7 were treated with Cap-clip Pyrophosphatase Acid enzyme. The second half of *Cryptococcus*
8 RNAs was mock treated. Each half of Cap-clip Pyrophosphatase Acid *Cryptococcus* RNA
9 samples was mixed with the same quantity of *S. cerevisiae* Cap-clip Pyrophosphatase Acid
10 treated RNAs. The subsequent steps of the library preparation were identical to the
11 published protocol (34). Single 50 bases single end reads were obtained using an an Illumina
12 HiSeq2500 instrument according to the manufacturer’s instructions (Illumina).

13 For QuantSeq 3’ mRNA-Seq preparation we followed the manufacturer instructions (Lexogen
14 GmbH, Austria). 100 base single end reads were obtained using an Illumina HiSeq2000
15 instrument according to the manufacturer’s instructions (Illumina).

16

17 **Sequencing data analyses**

18 For TSS analysis we kept only the reads containing both the oligo 3665
19 (AGATCGGAAGAGCACACGTCTGAAC) and the 11NCGCCGCGNNN tag (34). These sequences
20 were removed and the trimmed reads were mapped to the *Cryptococcus* genome and *S.*
21 *cerevisiae* genomes using Bowtie2 and Tophat2 (83). Their 5’ extremities were considered as
22 potential TSSs. For each condition we kept only the positions that were present in all three
23 replicates. Their coverage was normalized using the normalization factor used for spiked in

1 RNA-Seq. TSS positions were then clustered per condition. As most of the observed TSS sites
2 appeared as clusters, we grouped them into clusters by allowing an optimal maximum intra-
3 cluster distance (at 50 nt) between sites as previously used (34). We then removed the false
4 TSS clusters using the “no-cap” data keeping the clusters i for which

$$5 \quad R = \frac{\text{Weight}_{\text{cluster}_i}}{\sum \text{Weight}_{\text{cluster}}} / \frac{\text{Weight}_{\text{cluster}_{\text{nodecap}_i}}}{\sum \text{Weight}_{\text{cluster}_{\text{nodecap}}}} > 1$$

6

7 Similarly, QuantSeq 3'mRNA-Seq reads containing both the Sequencing and indexing primers
8 (Lexogen) were sorted. The reads were then cleaned using cutadapt/1.18 (84) and trimmed
9 for polyA sequence in their 3' end. PolyA untrimmed and trimmed reads were mapped to the
10 adapted *Cryptococcus* and to the *S. cerevisiae* genomes with Tophat2 (83) with the same
11 setting as for RNA-Seq. To eliminate the polyadenylated reads corresponding to genomic
12 polyA stretches, we considered only the reads that aligned to the genomes after polyA
13 trimming but not before the trimming. The 3' end position of these reads were considered as
14 potential PAS. As for the TSS, for each condition we kept only the positions that were
15 present in all three replicates. Similarly, the PAS dataset was normalized using the spike in
16 normalization factor and the PAS positions were clustered using the same strategies.

17 **Ribosome profiling and matched mRNA-seq**

18 Ribosome profiling was performed on both *C. neoformans* H99 and *C. deneoformans* JEC21,
19 two biological replicates of WT-H99 and one replicate each of H99 *ago1Δ* and H99 *gwo1Δ*
20 strains from (30), and one replicate each of WT-JEC21 and JEC21 *ago1Δ*. There was
21 negligible differential expression detected between these deletions and their background
22 strains, so in our analyses we treat the deletion strains as biological replicates.

1 Cells were grown to exponential phase in 750 mL of YPAD with shaking at 30°C. 100 ug/ml
2 cycloheximide (Sigma) (dissolved in 100% ethanol) was added to the culture and incubated
3 for 2 minutes. 50 mL of the culture was withdrawn for performing RNA-Seq in parallel. Cells
4 were then pelleted, resuspended in 5ml of lysis buffer (50mM Tris-HCl pH. 7.5, 150mM NaCl,
5 10mM MgCl₂, 5mM DTT, 0.5% Triton and 100ug/mL cyclohexamide) and snap frozen. Lysis,
6 clarification, RNaseI digestion, sucrose gradient separation and monosome isolation was
7 performed as previously described (36).

8 Ribosome protected fragments were isolated from the monosome fraction using hot phenol.
9 150ug of the total RNA extracted from the 50 ml of culture in parallel was polyA selected
10 using the Dynabeads mRNA purification kit (Thermo Fisher Scientific) and digested using
11 freshly made fragmentation buffer (100mM NaCO₃ pH. 9.2 and 2mM EDTA) for exactly 20
12 mins.

13 RNA was resolved on a 15% TBU gel. A gel slab corresponding to 28-34 nt was excised for
14 footprint samples and around 50 nt for mRNA samples, then eluted and precipitated.
15 Sequencing libraries were generated from the RNA fragments as described in Dunn et al.
16 with the following modifications (85). cDNA was synthesized using primer oCJ11 (Table S8).
17 Two rounds of subtractive hybridization for rRNA removal was done using oligos ras1-8 listed
18 in Table S8. After circularization Illumina adaptors were added through 9 cycles of PCR.
19 Libraries were sequenced on a HiSeq 2500 (Illumina).

20 **Ribosome profiling data analysis**

21 Ribosome profiling and matched RNA-seq reads were demultiplexed on BaseSpace (Illumina)
22 and then analyzed essentially with the RiboViz pipeline v.1.1.0 (86). In brief, sequencing
23 adapters were removed with cutadapt (84), and then reads aligned to rRNA were removed

1 by alignment with hisat2 (87). Cleaned non-rRNA reads were aligned to (spliced) transcripts
2 with hisat2 (87), sorted and indexed with samtools (88), and then quantified on annotated
3 ORFs with bedtools (89), followed by calculation of transcripts per million (TPM) and quality
4 control with R (90) scripts included in RiboViz. The cleaned non-rRNA reads were also aligned
5 to the genome with hisat2, and processed analogously, then used to generate figures of
6 genome alignments using ggplot2 (91) in R (90).

7 **Data analysis and visualization**

8 Data analysis and visualization were scripted in R (90), making extensive use of dplyr (92),
9 ggplot2 (91), and cowplot (93). Sequence logos were prepared in ggseqlogo (94). Analysis of
10 differential expression for *upf1Δ* data was performed in DeSeq2 (95). Figures were
11 assembled and annotated in Inkscape v0.92 (<https://inkscape.org>).

12 Protein sequences were aligned using muscle (96), with default parameters for protein
13 sequences and 100 iterations. Phylogenetic trees were constructed using ClustalW2 tool v2.1
14 (97) by using the neighbor-joining method with 1000 bootstrap trial replications.

15 Structural figures were prepared in PyMOL (Schrodinger).

16

17 **External datasets**

18 *N. crassa* (strain OR74A) ribosome profiling data from ((14), GEO:GSE97717) was used to
19 generate highly-translated genes, and ribosome profiling and RNA-seq data from ((98),
20 GEO: GSE71032) used to estimated TE. In both cases, we estimated TPMs using the RiboViz
21 pipeline as above, using the NC12 genome annotation downloaded from EnsemblGenomes
22 (99). TL sequences were also obtained from NC12.

1 *S. pombe* (strain 972h) ribosome profiling and RNA-seq data are from (100), and the authors
2 provided us with a table of RPKMs for all replicates as described. Genome sequence and
3 annotation ASM294v2, including TL annotation, were downloaded from EnsemblGenomes
4 (99).

5 *C. albicans* (strain SC5314) ribosome profiling and RNA-seq data are from (101),
6 GEO:GSE52236), processed with the RiboViz pipeline as above using the assembly 22 of the
7 strain SC5414 genome annotation from CGD (102).

8 *S. cerevisiae* (strain S288C/BY4741) highly-translated genes use the RPKM table from ((103),
9 GEO:GSE59573), and highly-expressed genes use (104). For TE estimates, we used matched
10 ribosome profiling and RNA-seq estimates from (105), although we did not use this for the
11 list of highly translated genes because near-duplicate paralogous ribosomal protein genes
12 were not present in the dataset, which thus omits a substantial fraction of highly-translated
13 genes. TL sequences were downloaded from SGD (106)

14 Protein homolog lists were assembled with OrthoDB (47) and PANTHERdb (48), with
15 reference to FungiDB (107). The list of cytoplasmic ribosomal proteins was assembled in *S.*
16 *cerevisiae* based on (108) with help from SGD (106), extended to other fungi with
17 PANTHERdb (47), and manually curated.

References

1. Grigoriev IV, Nikitin R, Haridas S, Kuo A, Ohm R, Otilar R, et al. MycoCosm portal: gearing up for 1000 fungal genomes. *Nucleic Acids Res.* 2014;42(Database issue):D699-D704.
2. Fan G, Sun Q, Li W, Shi W, Li X, Wu L, et al. The global catalogue of microorganisms 10K type strain sequencing project: closing the genomic gaps for the validly published prokaryotic and fungi species. *GigaScience.* 2018;7(5).
3. Shen X-X, Opulente DA, Kominek J, Zhou X, Steenwyk JL, Buh KV, et al. Tempo and Mode of Genome Evolution in the Budding Yeast Subphylum. *Cell.* 2018;175(6):1533-45.e20.
4. Butler G, Rasmussen MD, Lin MF, Santos MAS, Sakthikumar S, Munro CA, et al. Evolution of pathogenicity and sexual reproduction in eight *Candida* genomes. *Nature.* 2009;459:657.
5. Dujon B, Sherman D, Fisher G, Durrens P, Casaregola S, Lafontaine I, et al. Genome evolution in yeasts. *Nature.* 2004;430:35-44.
6. Stajich JE. Fungal Genomes and Insights into the Evolution of the Kingdom. *Microbiology spectrum.* 2017;5(4):10.1128/microbiolspec.FUNK-0055-2016.
7. Stajich JE, Dietrich FS, Roy SW. Comparative genomic analysis of fungal genomes reveals intron-rich ancestors. *Genome Biol.* 2007;8:R223.
8. Coletta A, Pinney JW, Solís DYW, Marsh J, Pettifer SR, Attwood TK. Low-complexity regions within protein sequences have position-dependent roles. *BMC systems biology.* 2010;4:43-.
9. Goffeau A, Barrell BG, Bussey H, Davis RW, Dujon B, Feldmann H, et al. Life with 6000 Genes. *Science.* 1996;274(5287):546.
10. Haas BJ, Zeng Q, Pearson MD, Cuomo CA, Wortman JR. Approaches to Fungal Genome Annotation. *Mycology.* 2011;2(3):118-41.
11. Fervers P, Fervers F, Makołowski W, Jąkowski M. Life cycle adapted upstream open reading frames (uORFs) in *Trypanosoma congolense*: A post-transcriptional approach to accurate gene regulation. *PLOS ONE.* 2018;13(8):e0201461.
12. Duncan CDS, Rodríguez-López M, Ruis P, Bähler J, Mata J. General amino acid control in fission yeast is regulated by a nonconserved transcription factor, with functions analogous to *Gcn4/Atf4*. *Proceedings of the National Academy of Sciences of the United States of America.* 2018;115(8):E1829-E38.
13. Sundaram A, Grant CM. A single inhibitory upstream open reading frame (uORF) is sufficient to regulate *Candida albicans* GCN4 translation in response to amino acid starvation conditions. *RNA (New York, NY).* 2014;20(4):559-67.
14. Ivanov IP, Wei J, Caster SZ, Smith KM, Michel AM, Zhang Y, et al. Translation Initiation from Conserved Non-AUG Codons Provides Additional Layers of Regulation and Coding Capacity. *mBio.* 2017;8(3):e00844-17.
15. von Arnim AG, Jia Q, Vaughn JN. Regulation of plant translation by upstream open reading frames. *Plant Science.* 2014;214:1-12.
16. Barbosa C, Peixeiro I, Romão L. Gene expression regulation by upstream open reading frames and human disease. *PLoS Genet.* 2013;9(8):e1003529-e.
17. Kozak M. Point mutations define a sequence flanking the AUG initiator codon that modulates translation by eukaryotic ribosomes. *Cell.* 1986;44(2):283-92.
18. Kozak M. An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs. *Nucleic acids research.* 1987;15(20):8125-48.
19. Hinnebusch AG, Ivanov IP, Sonenberg N. Translational control by 5'-untranslated regions of eukaryotic mRNAs. *Science.* 2016;352(6292):1413-6.
20. Dever TE, Kinzy TG, Pavitt GD. Mechanism and Regulation of Protein Synthesis in *Saccharomyces cerevisiae*. *Genetics.* 2016;203(1):65-107.
21. Dvir S, Velten L, Sharon E, Zeevi D, Carey LB, Weinberger A, et al. Deciphering the rules by which 5'-UTR sequences affect protein expression in yeast. *Proceedings of the National Academy of Sciences of the United States of America.* 2013;110(30):E2792-E801.

22. Cuperus JT, Groves B, Kuchina A, Rosenberg AB, Jovic N, Fields S, et al. Deep learning of the regulatory grammar of yeast 5' untranslated regions from 500,000 random sequences. *Genome Res.* 2017;27(12):2015-24.
23. Chen S-J, Lin G, Chang K-J, Yeh L-S, Wang C-C. Translational Efficiency of a Non-AUG Initiation Codon Is Significantly Affected by Its Sequence Context in Yeast. *J Biol Chem.* 2008;283(6):3173-80.
24. Wethmar K. The regulatory potential of upstream open reading frames in eukaryotic gene expression. *Wiley Interdisciplinary Reviews: RNA.* 2014;5(6):765-8.
25. Llácer JL, Hussain T, Marler L, Aitken CE, Thakur A, Lorsch JR, et al. Conformational Differences between Open and Closed States of the Eukaryotic Translation Initiation Complex. *Mol Cell.* 2015;59(3):399-412.
26. Hinnebusch AG. Structural Insights into the Mechanism of Scanning and Start Codon Recognition in Eukaryotic Translation Initiation. *Trends Biochem Sciences.* 2017;42(8):589-611.
27. Llácer JL, Hussain T, Saini AK, Nanda JS, Kaur S, Gordiyenko Y, et al. Translational initiation factor eIF5 replaces eIF1 on the 40S ribosomal subunit to promote start-codon recognition. *eLife.* 2018;7:e39273.
28. Janbon G. Introns in *Cryptococcus*. *Memorias do Instituto Oswaldo Cruz.* 2018;113(7):e170519-e.
29. Goebels C, Thonn A, Gonzalez-Hilarion S, Rolland O, Moyrand F, Beilharz TH, et al. Introns regulate gene expression in *Cryptococcus neoformans* in a Pab2p dependent pathway. *PLoS Genet.* 2013;9:e1003686.
30. Dumesic PA, Natarajan P, Chen C, Drinnenberg IA, Schiller BJ, Thompson JD, et al. Stalled spliceosomes are a signal for RNAi-mediated genome defense. *Cell.* 2013;152:957-68.
31. Bonnet A, Grosso AR, Elkaoutari A, Coleno E, Presle A, Sridhara SC, et al. Introns Protect Eukaryotic Genomes from Transcription-Associated Genetic Instability. *Mol Cell.* 2017;67(4):608-21.e6.
32. Janbon G, Ormerod KL, Paulet D, Byrnes III EJ, Chatterjee G, Yadav V, et al. Analysis of the genome and transcriptome of *Cryptococcus neoformans* var. *grubii* reveals complex RNA expression and microevolution leading to virulence attenuation. *PLoS Genet.* 2014;10:e1004261.
33. Gonzalez-Hilarion S, Paulet D, Lee K-T, Hon C-C, Lechat P, Mogensen E, et al. Intron retention-dependent gene regulation in *Cryptococcus neoformans*. *Scientific Reports.* 2016;6:32252.
34. Malabat C, Feuerbach F, Ma L, Saveanu C, Jacquier A. Quality control of transcription start site selection by nonsense-mediated-mRNA decay. *Elife.* 2015;4:e06722.
35. Li H, Hou J, Bai L, Hu C, Tong P, Kang Y, et al. Genome-wide analysis of core promoter structures in *Schizosaccharomyces pombe* with DeepCAGE. *RNA Biology.* 2015;12(5):525-37.
36. Ingolia NT, Ghaemmaghami S, Newman JRS, Weissman JS. Genome-Wide Analysis in Vivo of Translation with Nucleotide Resolution Using Ribosome Profiling. *Science.* 2009;324(5924):218.
37. Hinnebusch AG. TRANSLATIONAL REGULATION OF GCN4 AND THE GENERAL AMINO ACID CONTROL OF YEAST. *Ann Rev Microbiol.* 2005;59(1):407-50.
38. Duncan CDS, Rodríguez-López M, Ruis P, Bähler J, Mata J. General amino acid control in fission yeast is regulated by a nonconserved transcription factor, with functions analogous to Gcn4/Atf4. *Proceedings of the National Academy of Sciences.* 2018;115(8):E1829.
39. Madi L, McBride SA, Bailey LA, Ebbole DJ. rco-3, a gene involved in glucose transport and conidiation in *Neurospora crassa*. *Genetics.* 1997;146(2):499-508.
40. Wiese A, Elzinga N, Wobbles B, Smeekens S. Sucrose-induced translational repression of plant bZIP-type transcription factors. *Biochemical Society Transactions.* 2005;33(1):272.
41. Kervestin S, Jacobson A. NMD: a multifaceted response to premature translational termination. *Nat Rev Mol Cell Biol.* 2012;13:703-12.
42. Hood HM, Spevak CC, Sachs MS. Evolutionary changes in the fungal carbamoyl-phosphate synthetase small subunit gene and its associated upstream open reading frame. *Fungal Genet Biol.* 2007;44(2):93-104.
43. Danpure CJ. How can the products of a single gene be localized to more than one intracellular compartment? *Trends in Cell Biology.* 1995;5(6):230-8.

44. Mireau H, Lancelin D, Small ID. The same Arabidopsis gene encodes both cytosolic and mitochondrial alanyl-tRNA synthetases. *The Plant cell*. 1996;8(6):1027-39.
45. Mudge SJ, Williams JH, Eyre HJ, Sutherland GR, Cowan PJ, Power DA. Complex organisation of the 5'-end of the human glycine tRNA synthetase gene. *Gene*. 1998;209(1):45-50.
46. Natsoulis G, Hilger F, Fink GR. The HTS1 gene encodes both the cytoplasmic and mitochondrial histidine tRNA synthetases of *S. cerevisiae*. *Cell*. 1986;46(2):235-43.
47. Kriventseva EV, Kuznetsov D, Tegenfeldt F, Manni M, Dias R, Simão FA, et al. OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral genomes for evolutionary and functional annotations of orthologs. *Nucleic acids research*. 2019;47(D1):D807-D11.
48. Muruganujan A, Mi H, Thomas PD. PANTHER in 2013: modeling the evolution of gene function, and other gene attributes, in the context of phylogenetic trees. *Nucleic Acids Research*. 2012;41(D1):D377-D86.
49. Datt M, Sharma A. Novel and unique domains in aminoacyl-tRNA synthetases from human fungal pathogens *Aspergillus niger*, *Candida albicans* and *Cryptococcus neoformans*. *BMC Genomics*. 2014;15(1):1069.
50. Duchêne A-M, Pujol C, Maréchal-Drouard L. Import of tRNAs and aminoacyl-tRNA synthetases into mitochondria. *Curr Genet*. 2009;55(1):1-18.
51. Muruganujan A, Ebert D, Mi H, Thomas PD, Huang X. PANTHER version 14: more genomes, a new PANTHER GO-slim and improvements in enrichment analysis tools. *Nucleic Acids Res*. 2018;47(D1):D419-D26.
52. Frechin M, Duchêne A-M, Becker HD. Translating organellar glutamine codons : A case by case scenario? *RNA Biology*. 2009;6(1):31-4.
53. Chang C-P, Tseng Y-K, Ko C-Y, Wang C-C. Alanyl-tRNA synthetase genes of *Vanderwaltozyma polyspora* arose from duplication of a dual-functional predecessor of mitochondrial origin. *Nucleic Acids Res*. 2012;40(1):314-22.
54. Geslain R, Martin F, Delagoutte B, Cavarelli J, Gangloff J, Eriani G. In vivo selection of lethal mutations reveals two functional domains in arginyl-tRNA synthetase. *RNA (New York, NY)*. 2000;6(3):434-48.
55. Merz S, Westermann B. Genome-wide deletion mutant analysis reveals genes required for respiratory growth, mitochondrial genome maintenance and mitochondrial protein synthesis in *Saccharomyces cerevisiae*. *Genome Biol*. 2009;10(9):R95.
56. Sickmann A, Reinders J, Wagner Y, Joppich C, Zahedi R, Meyer HE, et al. The proteome of *Saccharomyces cerevisiae* mitochondria. *Proceedings of the National Academy of Sciences of the United States of America*. 2003;100(23):13207-12.
57. Chen S-J, Wu Y-H, Huang H-Y, Wang C-C. *Saccharomyces cerevisiae* Possesses a Stress-Inducible Glycyl-tRNA Synthetase Gene. *PLOS ONE*. 2012;7(3):e33363.
58. Chiu W-C, Chang C-P, Wen W-L, Wang S-W, Wang C-C. *Schizosaccharomyces pombe* Possesses Two Paralogous Valyl-tRNA Synthetase Genes of Mitochondrial Origin. *Mol Biol Evol*. 2010;27(6):1415-24.
59. Ivanov IP, Loughran G, Sachs MS, Atkins JF. Initiation context modulates autoregulation of eukaryotic translation initiation factor 1 (eIF1). *Proceedings of the National Academy of Sciences of the United States of America*. 2010;107(42):18056-60.
60. Loughran G, Sachs MS, Atkins JF, Ivanov IP. Stringency of start codon selection modulates autoregulation of translation initiation factor eIF5. *Nucleic Acids Res*. 2012;40(7):2898-906.
61. Martin-Marcos P, Cheung Y-N, Hinnebusch AG. Functional elements in initiation factors 1, 1A, and 2 β discriminate against poor AUG context and non-AUG start codons. *Mol Cell Biol*. 2011;31(23):4814-31.
62. Hussain T, Ll acer JL, Fern andez IS, Munoz A, Martin-Marcos P, Savva CG, et al. Structural changes enable start codon recognition by the eukaryotic translation initiation complex. *Cell*. 2014;159(3):597-607.

63. Thakur A, Hinnebusch AG. eIF1 Loop 2 interactions with Met-tRNA(i) control the accuracy of start codon selection by the scanning preinitiation complex. *Proceedings of the National Academy of Sciences of the United States of America*. 2018;115(18):E4159-E68.
64. Olsen DS, Savner EM, Mathew A, Zhang F, Krishnamoorthy T, Phan L, et al. Domains of eIF1A that mediate binding to eIF2, eIF3 and eIF5B and promote ternary complex recruitment in vivo. *The EMBO journal*. 2003;22(2):193-204.
65. Luna RE, Arthanari H, Hiraishi H, Akabayov B, Tang L, Cox C, et al. The interaction between eukaryotic initiation factor 1A and eIF5 retains eIF1 within scanning preinitiation complexes. *Biochemistry*. 2013;52(52):9510-8.
66. Fekete CA, Applefield DJ, Blakely SA, Shirokikh N, Pestova T, Lorsch JR, et al. The eIF1A C-terminal domain promotes initiation complex assembly, scanning and AUG selection in vivo. *The EMBO journal*. 2005;24(20):3588-601.
67. Slusher LB, Gillman EC, Martin NC, Hopper AK. mRNA leader length and initiation codon context determine alternative AUG selection for the yeast gene MOD5. *Proceedings of the National Academy of Sciences*. 1991;88(21):9789.
68. Calvo SE, Pagliarini DJ, Mootha VK. Upstream open reading frames cause widespread reduction of protein expression and are polymorphic among humans. *Proceedings of the National Academy of Sciences of the United States of America*. 2009;106(18):7507-12.
69. Cheng Z, Otto GM, Powers EN, Keskin A, Mertins P, Carr SA, et al. Pervasive, Coordinated Protein-Level Changes Driven by Transcript Isoform Switching during Meiosis. *Cell*. 2018;172(5):910-23.e16.
70. Van Daltsen KM, Hodapp S, Keskin A, Otto GM, Berdan CA, Higdon A, et al. Global Proteome Remodeling during ER Stress Involves Hac1-Driven Expression of Long Undecoded Transcript Isoforms. *Developmental Cell*. 2018;46(2):219-35.e8.
71. Monteuuis G, Miścicka A, Świrski M, Zenad L, Niemitalo O, Wrobel L, et al. Non-canonical translation initiation in yeast generates a cryptic pool of mitochondrial proteins. *Nucleic Acids Res*. 2019:in press.
72. Brar GA. Beyond the Triplet Code: Context Cues Transform Translation. *Cell*. 2016;167(7):1681-92.
73. Feeney KA, Hansen LL, Putker M, Olivares-Yañez C, Day J, Eades LJ, et al. Daily magnesium fluxes regulate cellular timekeeping and energy balance. *Nature*. 2016;532(7599):375-9.
74. Tsuboi T, Viana MP, Xu F, Yu J, Chanchani R, Arceo XG, et al. Mitochondrial volume fraction controls translation of nuclear-encoded mitochondrial proteins. *bioRxiv*. 2019:529289.
75. Nakagawa S, Niimura Y, Gojobori T, Tanaka H, Miura K-i. Diversity of preferred nucleotide sequences around the translation initiation codon in eukaryote genomes. *Nucleic acids research*. 2008;36(3):861-71.
76. Shah M, Su D, Scheliga JS, Pluskal T, Boronat S, Motamedchaboki K, et al. A Transcript-Specific eIF3 Complex Mediates Global Translational Control of Energy Metabolism. *Cell reports*. 2016;16(7):1891-902.
77. Fields SD, Conrad MN, Clarke M. The *S. cerevisiae* CLU1 and *D. discoideum* cluA genes are functional homologues that influence mitochondrial morphology and distribution. *Journal of Cell Science*. 1998;111(12):1717.
78. Gao J, Schatton D, Martinelli P, Hansen H, Pla-Martin D, Barth E, et al. CLUH regulates mitochondrial biogenesis by binding mRNAs of nuclear-encoded mitochondrial proteins. *The Journal of Cell Biology*. 2014;207(2):213.
79. Schatton D, Pla-Martin D, Marx M-C, Hansen H, Mourier A, Nemazanyy I, et al. CLUH regulates mitochondrial metabolism by controlling translation and decay of target mRNAs. *The Journal of cell biology*. 2017;216(3):675-93.
80. Smith MD, Gu Y, Querol-Audí J, Vogan JM, Nitido A, Cate JHD. Human-Like Eukaryotic Translation Initiation Factor 3 from *Neurospora crassa*. *PLOS ONE*. 2013;8(11):e78715.
81. Madhani HD. The frustrated gene: origins of eukaryotic gene expression. *Cell*. 2013;155(4):744-9.

82. Winston F, Dollard C, Ricupero-Hovasse SL. Construction of a set of convenient *saccharomyces cerevisiae* strains that are isogenic to S288C. *Yeast*. 1995;11(1):53-5.
83. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol*. 2013;14(4):R36-R.
84. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnetjournal*; Vol 17, No 1: Next Generation Sequencing Data Analysis. 2011.
85. Dunn JG, Foo CK, Belletier NG, Gavis ER, Weissman JS. Ribosome profiling reveals pervasive and regulated stop codon readthrough in *Drosophila melanogaster*. *eLife*. 2013;2:e01179.
86. Carja O, Xing T, Wallace EWJ, Plotkin JB, Shah P. riboviz: analysis and visualization of ribosome profiling datasets. *BMC bioinformatics*. 2017;18(1):461-.
87. Pertea M, Kim D, Pertea GM, Leek JT, Salzberg SL. Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nature Protocols*. 2016;11:1650.
88. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map format and SAMtools. *Bioinformatics (Oxford, England)*. 2009;25(16):2078-9.
89. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics (Oxford, England)*. 2010;26(6):841-2.
90. Team RC. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. Available online at <https://www.R-project.org/>. 2018.
91. Wickham H, editor. *ggplot2: Elegant Graphics for Data Analysis*: Springer-Verlag New York; 2016.
92. Wickham H, François R, Henry L, Müller K. dplyr: A Grammar of Data Manipulation. R package version 0.7.8. <https://CRAN.R-project.org/package=dplyr>. 2018.
93. Wilke CO. cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'. R package version 0.9.3. <https://CRAN.R-project.org/package=cowplot>. 2018.
94. Wagih O. ggseqlogo: A 'ggplot2' Extension for Drawing Publication-Ready Sequence Logos. R package version 0.1. <https://CRAN.R-project.org/package=ggseqlogo>. 2017.
95. Love MI, Huber W, Anders S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol*. 2014;15:550.
96. Edgar RC. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res*. 2004;32(5):1792-7.
97. Wilm A, Higgins DG, Valentin F, Blackshields G, McWilliam H, Wallace IM, et al. Clustal W and Clustal X version 2.0. *Bioinformatics*. 2007;23(21):2947-8.
98. Yu C-H, Dang Y, Zhou Z, Wu C, Zhao F, Sachs MS, et al. Codon Usage Influences the Local Rate of Translation Elongation to Regulate Co-translational Protein Folding. *Mol Cell*. 2015;59(5):744-54.
99. Kersey PJ, Allen JE, Allot A, Barba M, Boddu S, Bolt BJ, et al. Ensembl Genomes 2018: an integrated omics infrastructure for non-vertebrate species. *Nucleic Acids Res*. 2018;46(D1):D802-D8.
100. Duncan CDS, Mata J. Effects of cycloheximide on the interpretation of ribosome profiling experiments in *Schizosaccharomyces pombe*. *Scientific Reports*. 2017;7(1):10331.
101. Muzzey D, Sherlock G, Weissman JS. Extensive and coordinated control of allele-specific expression by both transcription and translation in *Candida albicans*. *Genome Res*. 2014;24(6):963-73.
102. Skrzypek MS, Binkley J, Binkley G, Miyasato SR, Simison M, Sherlock G. The *Candida* Genome Database (CGD): incorporation of Assembly 22, systematic identifiers and visualization of high throughput sequencing data. *Nucleic Acids Res*. 2017;45(D1):D592-D6.
103. Gerashchenko MV, Gladyshev VN. Translation inhibitors cause abnormalities in ribosome profiling experiments. *Nucleic Acids Res*. 2014;42(17):e134-e.
104. Csárdi G, Franks A, Choi DS, Airoidi EM, Drummond DA. Accounting for Experimental Noise Reveals That mRNA Levels, Amplified by Post-Transcriptional Processes, Largely Determine Steady-State Protein Levels in Yeast. *PLoS Genet*. 2015;11(5):e1005206.

105. Weinberg DE, Shah P, Eichhorn SW, Hussmann JA, Plotkin JB, Bartel DP. Improved Ribosome-Footprint and mRNA Measurements Provide Insights into Dynamics and Regulation of Yeast Translation. *Cell reports*. 2016;14(7):1787-99.
106. Cherry JM, Hong EL, Amundsen C, Balakrishnan R, Binkley G, Chan ET, et al. Saccharomyces Genome Database: the genomics resource of budding yeast. *Nucleic Acids Res*. 2012;40(Database issue):D700-D5.
107. Basenko YE, Pulman AJ, Shanmugasundram A, Harb SO, Crouch K, Starns D, et al. FungiDB: An Integrated Bioinformatic Resource for Fungi and Oomycetes. *Journal of Fungi*. 2018;4(1).
108. Ban N, Beckmann R, Cate JHD, Dinman JD, Dragon F, Ellis SR, et al. A new system for naming ribosomal proteins. *Current opinion in structural biology*. 2014;24:165-9.

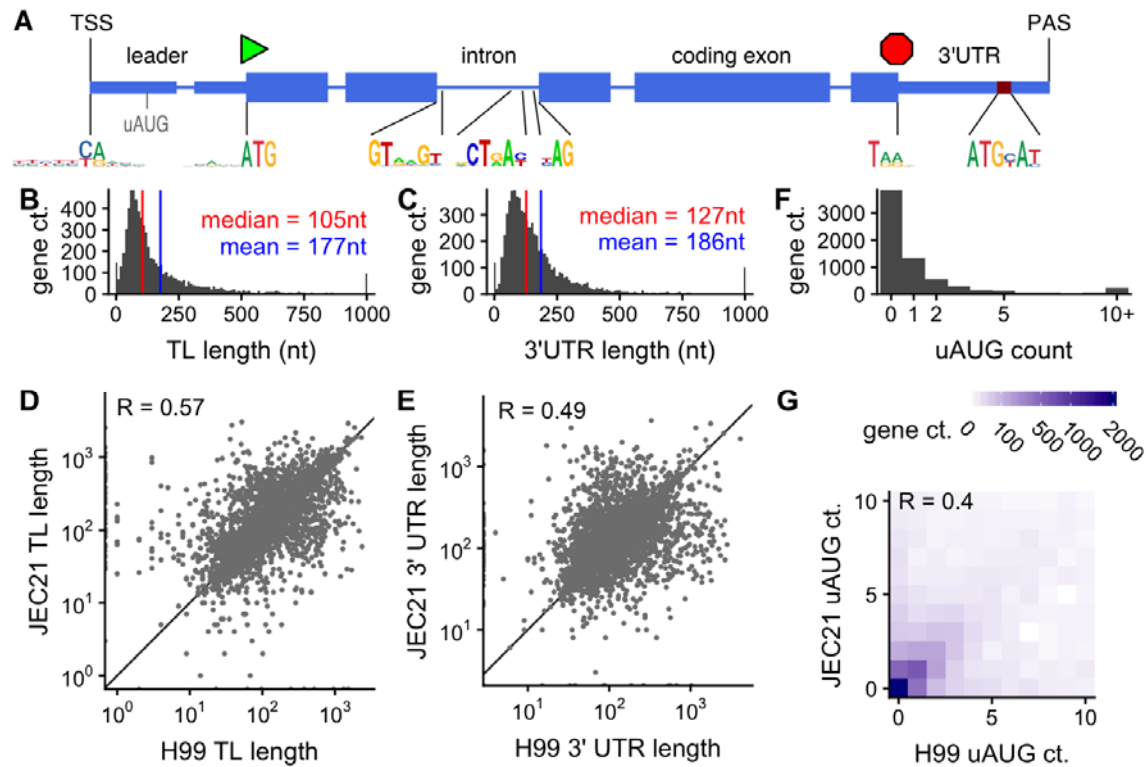


Figure 1: Mapping the coding transcriptome of *Cryptococcus neoformans*. A, Representation of a stereotypical gene of *C. neoformans* H99, showing the sequence logos for the transcription start site (TSS), AUG start codon, intron splicing, stop codon, and polyadenylation site (PAS). B, Distribution of transcript leader (TL) lengths over *C. neoformans* genes, for yeast cells growing exponentially in YPD at 30°C. C, Distribution of 3' untranslated region (3'UTR) lengths over *C. neoformans* genes. C, D : Comparisons of TL and 3'UTR lengths between orthologous genes in *C. neoformans* H99 and *C. deneoformans* JEC21 growing exponentially in YPD at 30°C. F, Distribution of upstream AUG (uAUG) counts over *C. neoformans* genes, and G, comparison of uAUG counts with *C. deneoformans*.

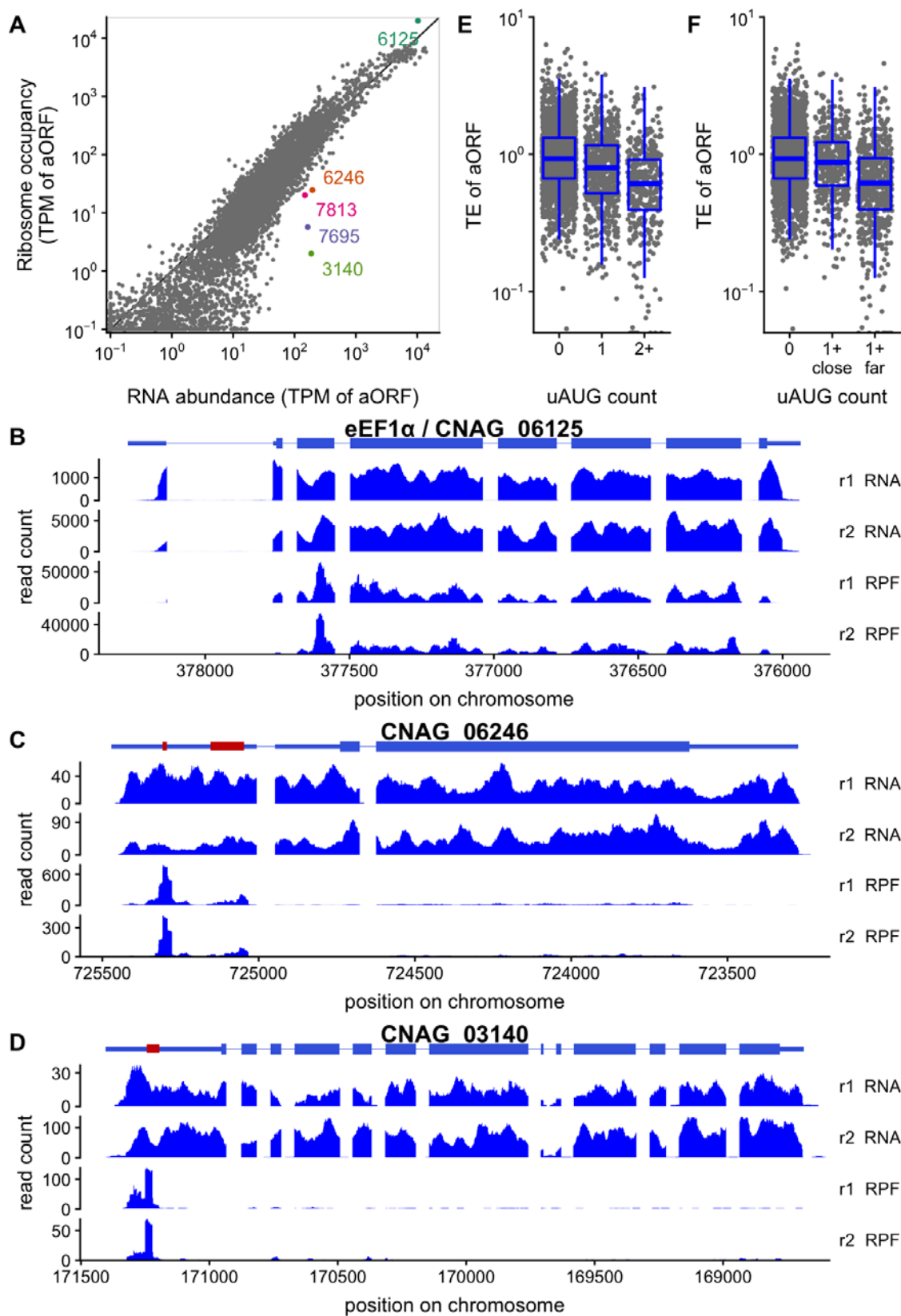


Figure 2: Upstream AUGs repress translation in *C. neoformans*. A, translation regulation of annotated ORFs (aORFs) in *C. neoformans* H99 growing exponentially in YPD at 30°C (equivalent data for *C. deneoformans* shown in Fig S2.1). Ribosome occupancy is plotted against the RNA abundance, both calculated in transcripts per million (TPM) on the aORF. Select genes discussed in the text are highlighted in colour. B, uAUGs are associated with lower translation efficiency (TE) of annotated ORFs, measured as the ratio of ribosome occupancy to RNA-seq reads. C, only uAUGs far from the transcription start site are associated with low TE. A gene is in the “1+ far” category if it has at least one uAUG more than 20nt from the TSS, “1+ close” if all uAUGs are within 20nt of the TSS. D-F, Examples of ribosome occupancy profiles along select RNAs highlighted in 2A (others are shown in Fig S2.2). D, Translation elongation factor eEF2/CNAG_06125 has high ribosome occupancy in the annotated ORF. Translationally repressed mRNAs CNAG_06246 (E) and CNAG_03140 (F) have high ribosome occupancy in uORFs in the transcript leader (red), and low ribosome occupancy in the aORF. Only the first of 5 uORFs in CNAG_03140 is shown. Homologous genes in *C. deneoformans* have similar structure and regulation (Fig S2.1).

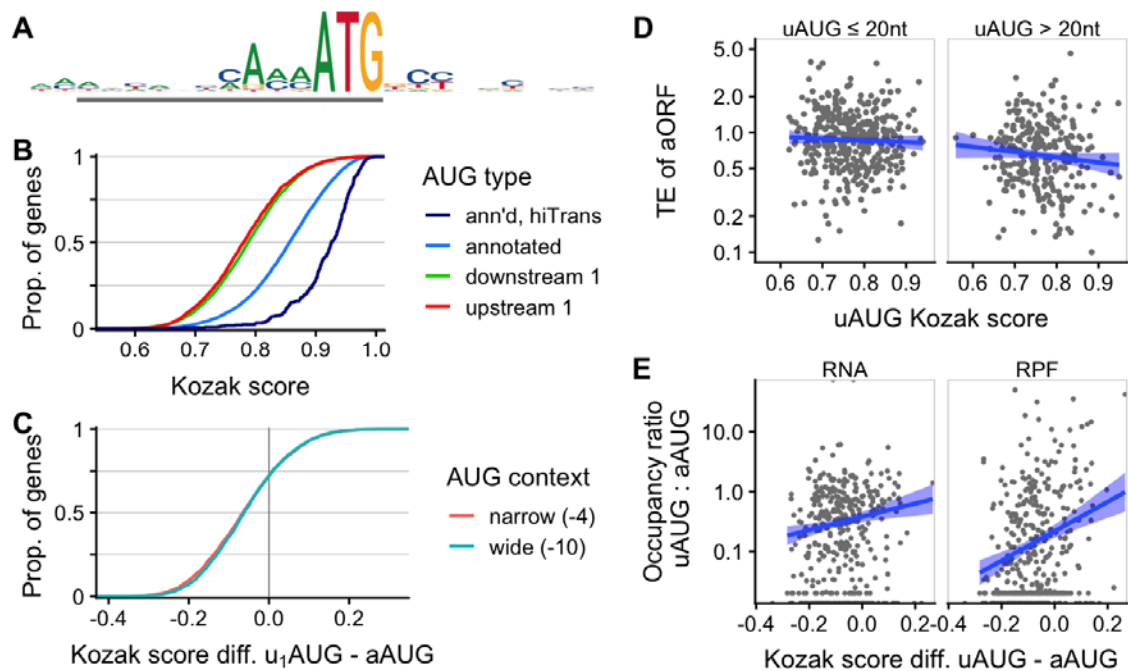


Figure 3: An AUG sequence context is associated with translation in *C. neoformans*. A, Kozak-like sequence context of AUGs, from -12 to +12, for highest-translated 5% of genes (hiTrans). This sequence context is used to create “Kozak scores” of other AUG sequences by their similarity to the consensus from -10 onwards. B, Cumulative density plot of Kozak scores from various categories of AUG, showing that high scores are associated with annotated AUGs of highly translated genes (hiTrans), somewhat with annotated AUGs, and not with the most 5' downstream AUG (downstream 1) or 5' most upstream AUG (upstream 1) in a transcript. C, Cumulative density plot of differences in scores between most 5' upstream (u_1 AUG) and annotated AUG, showing that for 75% of genes the upstream AUG score is less than the annotated AUG, whether we take a wide (-10:AUG) or a narrow (-4:AUG) window to calculate the score. D, High upstream AUG score is weakly associated with translation repression of the annotated ORF, if the uAUG is further than 20nt from the TSS. E, The relative occupancy of ribosomes (RPF) at the upstream AUG and annotated AUG depends on the difference in scores, even when compared to RNA-seq reads; linear model trend fit shown (blue). Panels D and E show data only for genes in the top 50% by RNA abundance, and with only a single upstream AUG.

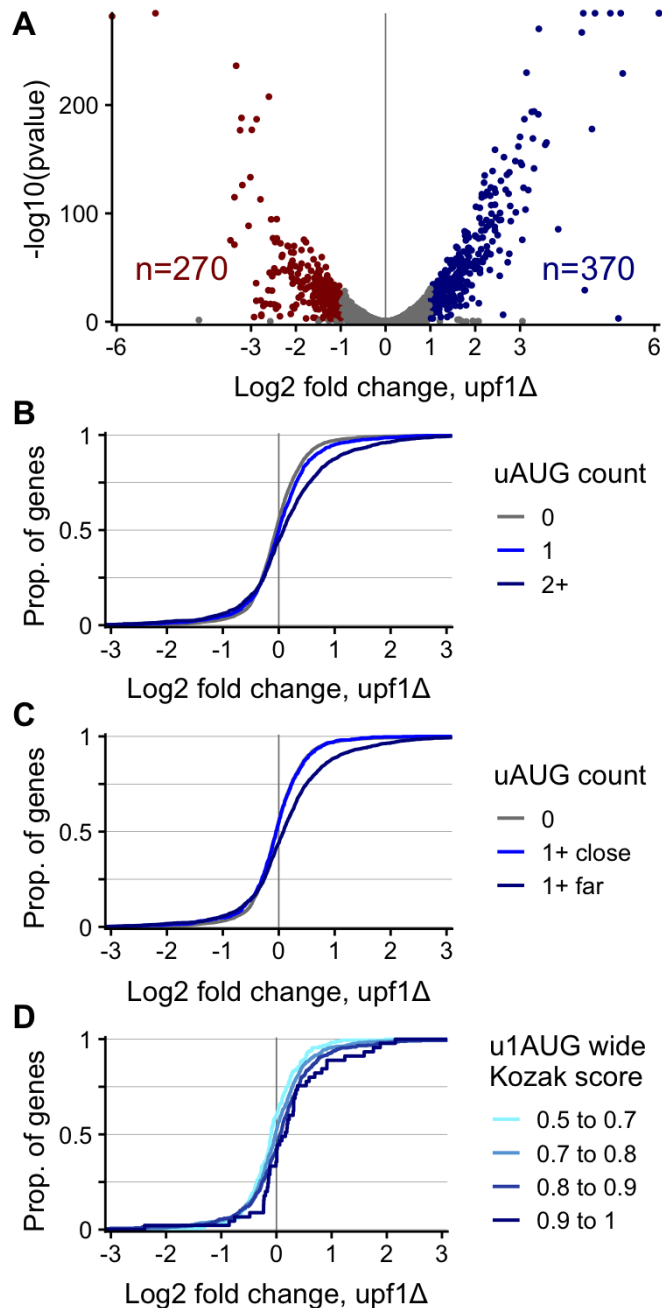


Figure 4: Nonsense-mediated decay (NMD) acts on upstream-ORF-containing mRNAs in *C. deneoformans*. A, Differential expression results from RNA-Seq in *C. deneoformans* JEC21, comparing expression in wild-type cells with a mutant deleted for NMD factor *UPF1/CNC02960*, and using DeSeq2 to identify genes upregulated in the *upf1* Δ mutant. B, uORF containing genes are enriched for NMD-sensitivity. C, uORF-containing genes are enriched for NMD-sensitivity only when the uAUG is more than 20nts from the TSS (1+far), but not when the uAUG is less than 20nts (1+close). D, Start codon sequence context affects NMD sensitivity of genes containing a single upstream AUG: uORFs starting with higher Kozak-score uAUG are more likely to increase in abundance in the *upf1* Δ mutant.

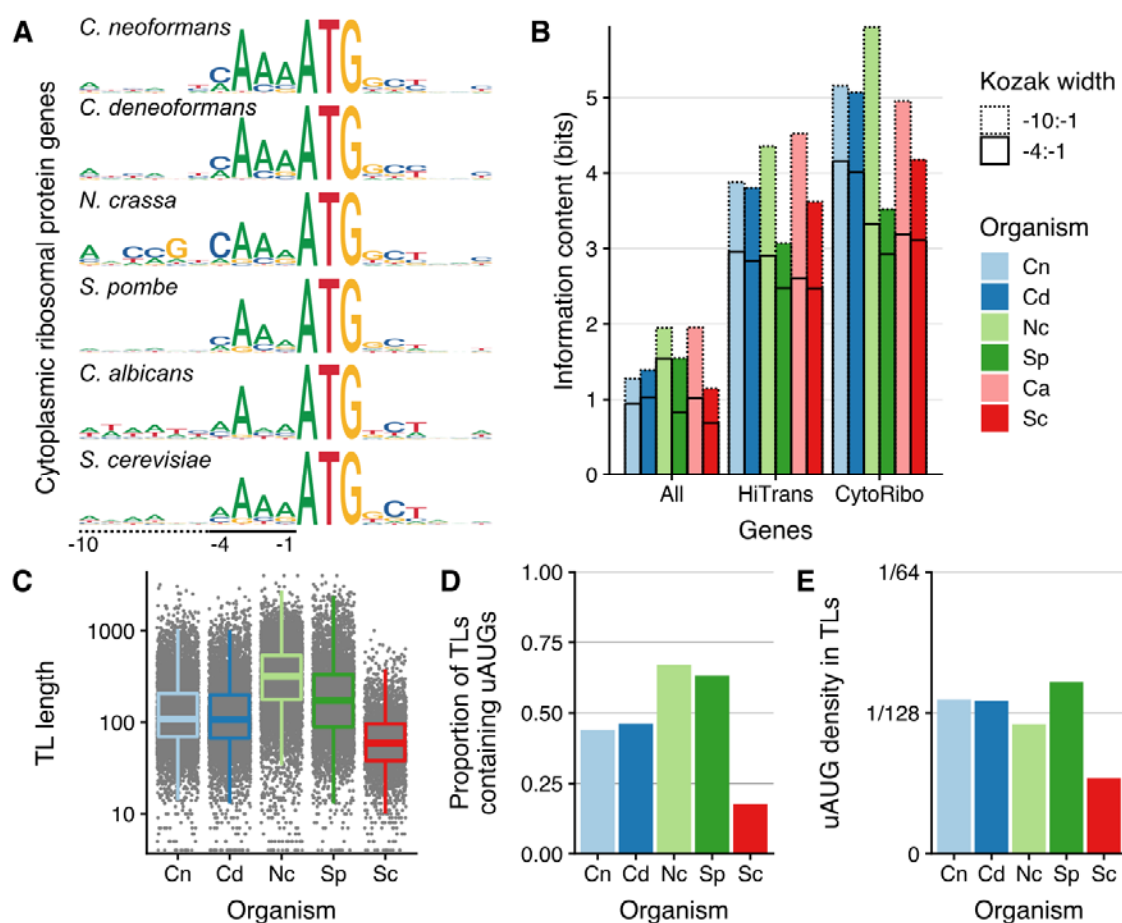


Figure 5: Sequences specifying start codon selection are quantitatively different in different fungi. A, Kozak consensus sequence logo for annotated start codons of cytoplasmic ribosomal protein genes from 6 fungal species. The height of each letter represents the Shannon information content in bits, so that the anchor ATG sequence has height 2 bits. B, Information content at annotated start codons in bits per base (i.e. summed height of stacked letters in sequence logo) for 3 groups of genes, in the 6 fungi from panel A. Solid line indicates information from -1 to -4 of ATG, and dotted line additionally to -10 (see bottom of panel A). Gene groups are all annotated ORFs, highly translated ORFs (HiTrans) and cytoplasmic ribosomal proteins (CytoRibo, as panel A). HiTrans used the highest-translated 5% of genes, or the highest 400 genes for fungi with more than 8000 annotated genes (*C. albicans* and *N. crassa*; see methods). C-E, For 5 fungi for which transcript leader (TL) annotations were available, TL length (C), proportion of annotated TL containing an upstream AUG (D), and proportions of AUGs per nucleotide in the TL (E; a uniform random model would have density 1/64).

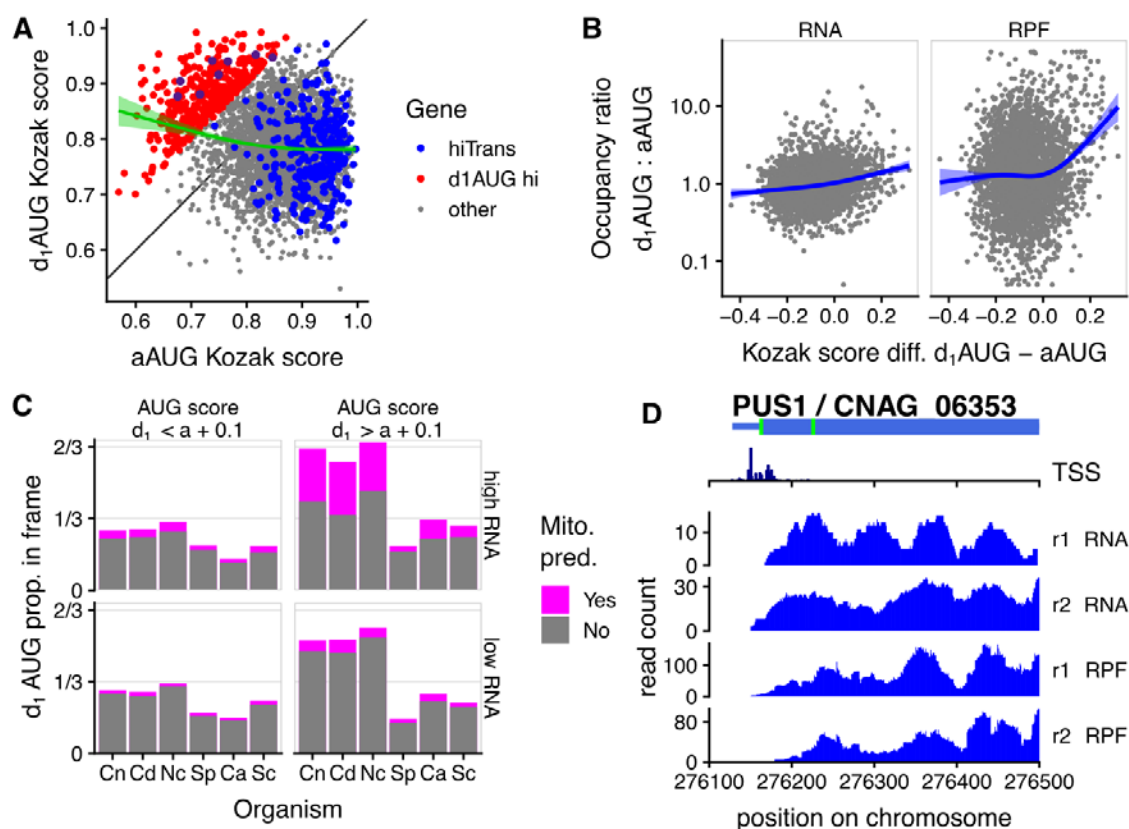


Figure 6: High-scoring downstream AUGs specify alternative N-terminal isoforms in *C. neoformans*.

A, Most genes with reasonable RNA abundance (top 50% by RNA abundance shown), especially very highly-translated genes (blue, top 5%), have lower Kozak score at the 1st downstream AUG than at the annotated AUG. However there are exceptions (red, d₁AUG hi: d₁AUG score > annotated AUG score + 0.1), and there is a trend for genes with low aAUG score to have a higher d₁AUG score (green, generalized additive model fit). **B**, Higher d₁AUG score than aAUG score drives higher ribosome protected fragment (RPF) occupancy at the d₁AUG compared to the aAUG, but much smaller differences in RNA-seq density. Blue line indicates generalized additive model fit. **C**, Downstream AUGs with high Kozak scores (d₁AUG score > annotated AUG score + 0.1) and reasonable RNA abundance (top 50%) are likely to be in-frame and enriched for N-terminal mitochondrial localization signals in *C. neoformans*, *C. deneoformans*, and *N. crassa*, but not in *S. pombe*, *C. albicans*, or *S. cerevisiae*. **D**, The pseudouridine synthase *CnPus1* is a candidate alternate-localized protein with a low-score aAUG and high-score d₁AUG, and transcription start sites on both sides of the aAUG. RNA-Seq and RPF reads on the first exon are shown, and the full length of the gene shown in Fig S6.1.

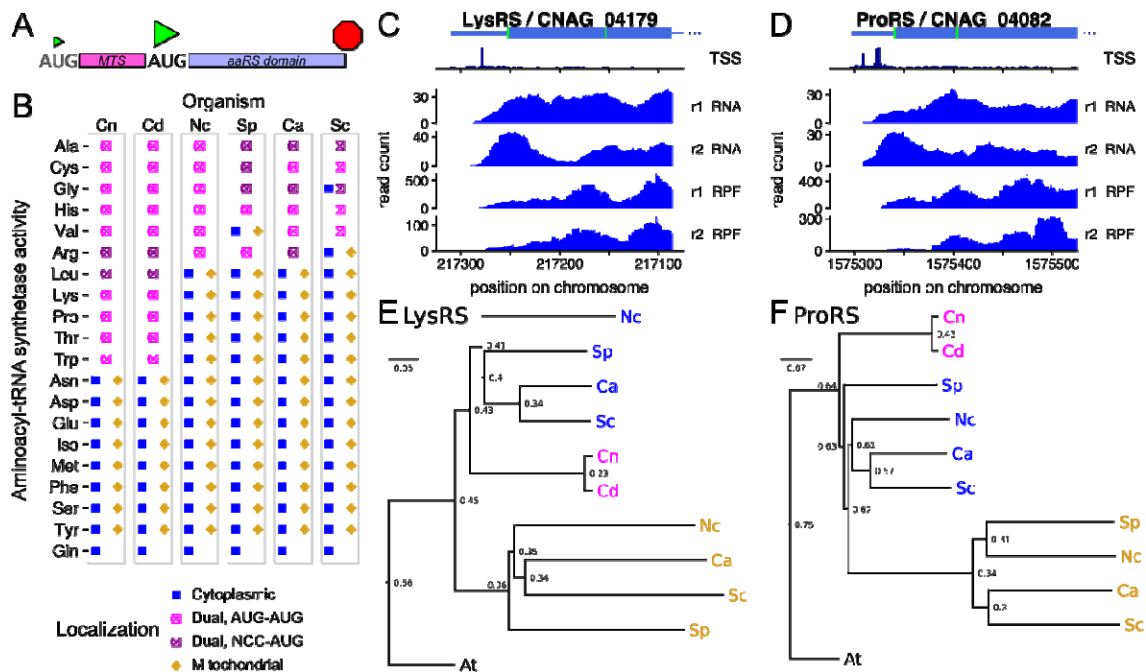


Figure 7: Aminoacyl-tRNA synthetases (aaRSs) are commonly alternatively localized to cytoplasm and mitochondria by use of alternative start codons in fungi. A, Schematic of the structure of a dual-localized aaRS with alternate AUG start codons. B, Predicted localization of all aaRS enzymes in the fungi *C. neoformans* (Cn), *C. deneoformans* (Cd), *N. crassa* (Nc), *S. pombe* (Sp), *C. albicans* (Ca), *S. cerevisiae* (Sc). C/D, Transcription start site reads, RNA-seq, and ribosome profiles of 5'-ends of CnLysRS (C) and CnProRS (D) show that most transcription starts upstream of both AUG start codons (green), and both AUG codons are used for translation initiation. E/F Simplified neighbour-joining phylogenetic trees show that LysRS (E) and ProRS (F) genes were duplicated in ascomycete fungi, and *Cryptococcus* retained a single dual-localized homolog. *Arabidopsis thaliana* (At) was used as an out-group. The scale bar represents the number of amino acid substitutions per residue, and the numbers at nodes are the proportion of substitutions between that node and its parent. See table S5, for details of identifiers for genes (GeneID).

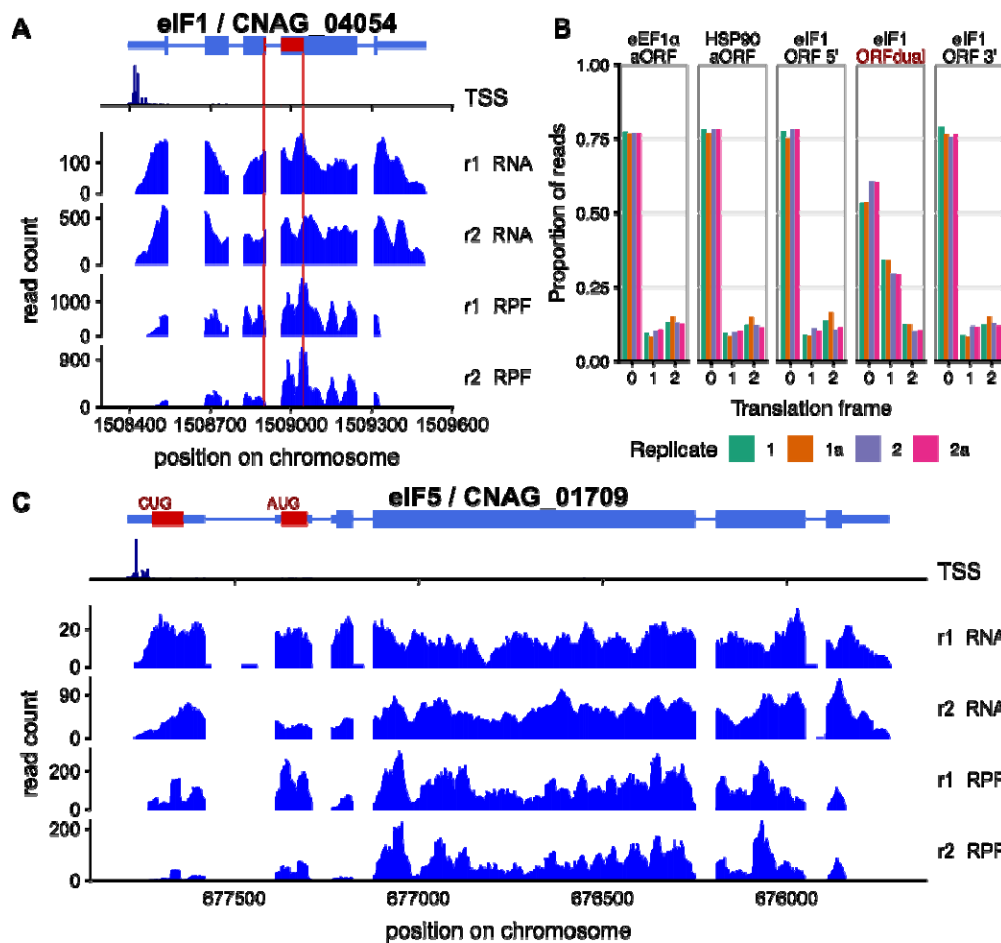


Figure 8. Translation initiation factors eIF1 and eIF5 are regulated by alternate start codon usage in *C. neoformans*. A, Reads on *CneIF1*/CNAG_04054, showing frame +1 “downstream ORF” in dark red, breaking for an intron. B, The downstream ORF of *CneIF1* is dual-translated in two frames. Most ribosome profiling read 5’ ends are in a consistent frame, including in control genes eEF1α/CNAG_06125 and HSP90/CNAG_06125, and in the 5’ and 3’ ends of the *CneIF1* ORF, but there is 2x enrichment of reads in frame+1 in the dual-decoded ORF. C, Reads on *CneIF5*/CNAG_01709 showing substantial ribosomal occupancy over upstream ORFs. The first upstream ORF shown is translated from a CUG start codon and the second from an AUG codon, and other uORFS potentially initiated from near-cognate codons are not shown. *C. deneoformans* homologs have the same structure and regulation (Figure S8).

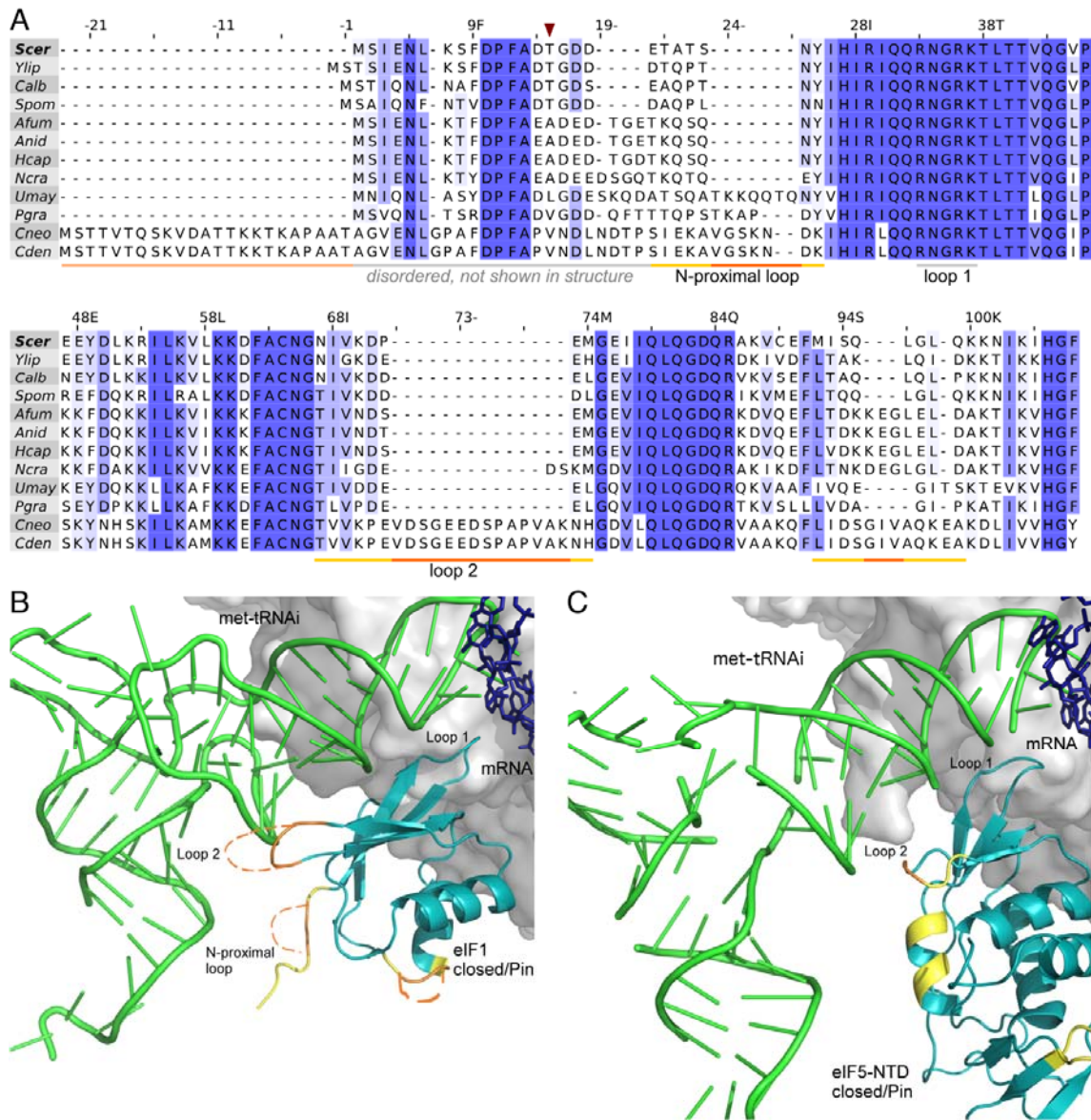


Figure 9. Eukaryotic translation initiation factor 1 is highly variable across fungi. A, Multiple sequence alignment of translation initiation factor eIF1 from 12 fungi, numbered as *S. cerevisiae* (*Scer*, top line). *Cryptococcus* insertions are indicated in orange, and surrounding variable residues in yellow. The N-terminal extension in *Cryptococcus* eIF1, that is predicted disordered, is shown in pale orange, and T15 residue with dark red arrow. B, Structural predictions of insertions (orange) and non-conserved neighborhoods (yellow) in *Cryptococcus* eIF1 mapped onto the closed pre-initiation complex of *S. cerevisiae/K. lactis* (PDB:3J81, Hussein 2015). eIF1 (teal) and Met-tRNAi (green) in closed conformation, shown with synthetic mRNA sequence (pink), and eIF2 (pale pink) and ribosomal subunit surface (greys) in background. Approximate ribosomal contacts are shown as grey background surface and eIF2-alpha subunit is shown as pale pink sticks. B, Structural predictions of variations in *Cryptococcus* eIF5 mapped on to *S. cerevisiae* PIC (PDB:6FYX, (27)). Multiple sequence alignment of eIF5 is shown in figure S9.1A.

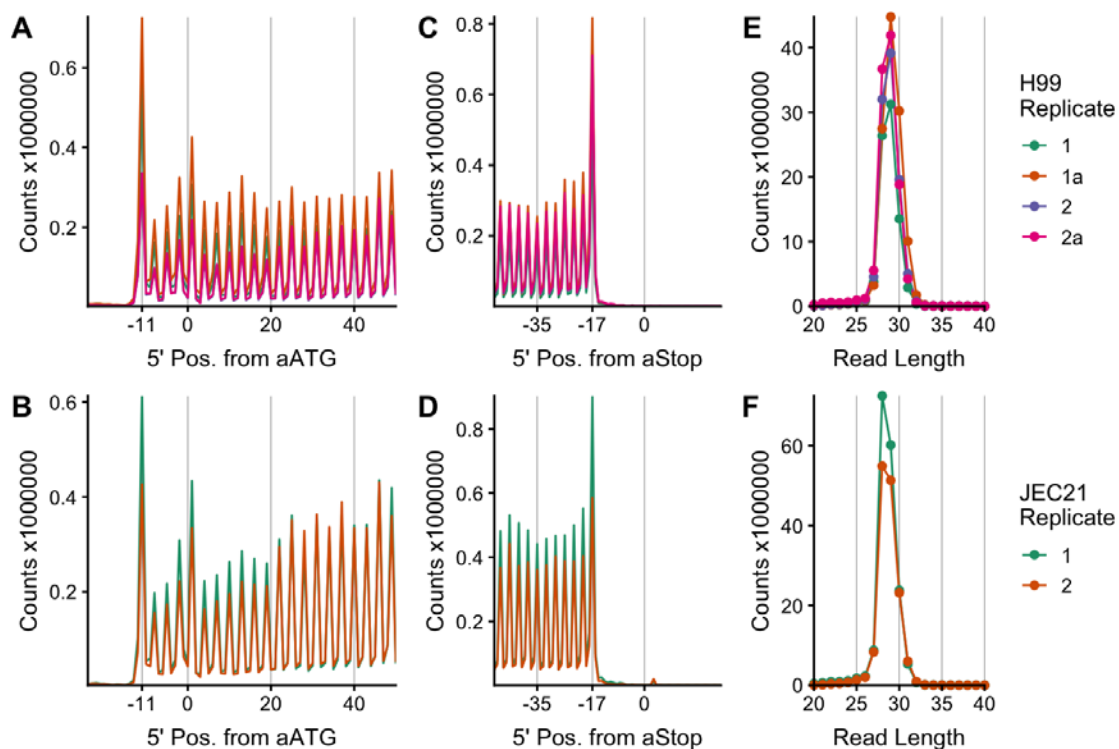


Figure S2.1, related to figure 2: Ribosome profiling data passes quality control metrics. Metagenome profiles of mapped 5' ends of ribosome-protected fragment counts at the 5' end (A,B) and 3' end (C,D) of ORFs, showing 3-nucleotide periodicity indicative of active translation starting at the annotated start codon and ending at the annotated stop codon. Ribosome protected fragment length is of a consistent length with other studies (E,F). Top row is data from 4 replicates of *C. neoformans* H99, bottom row from 2 replicates of *C. deneoformans* JEC21. These figures were made using RiboViz (86).

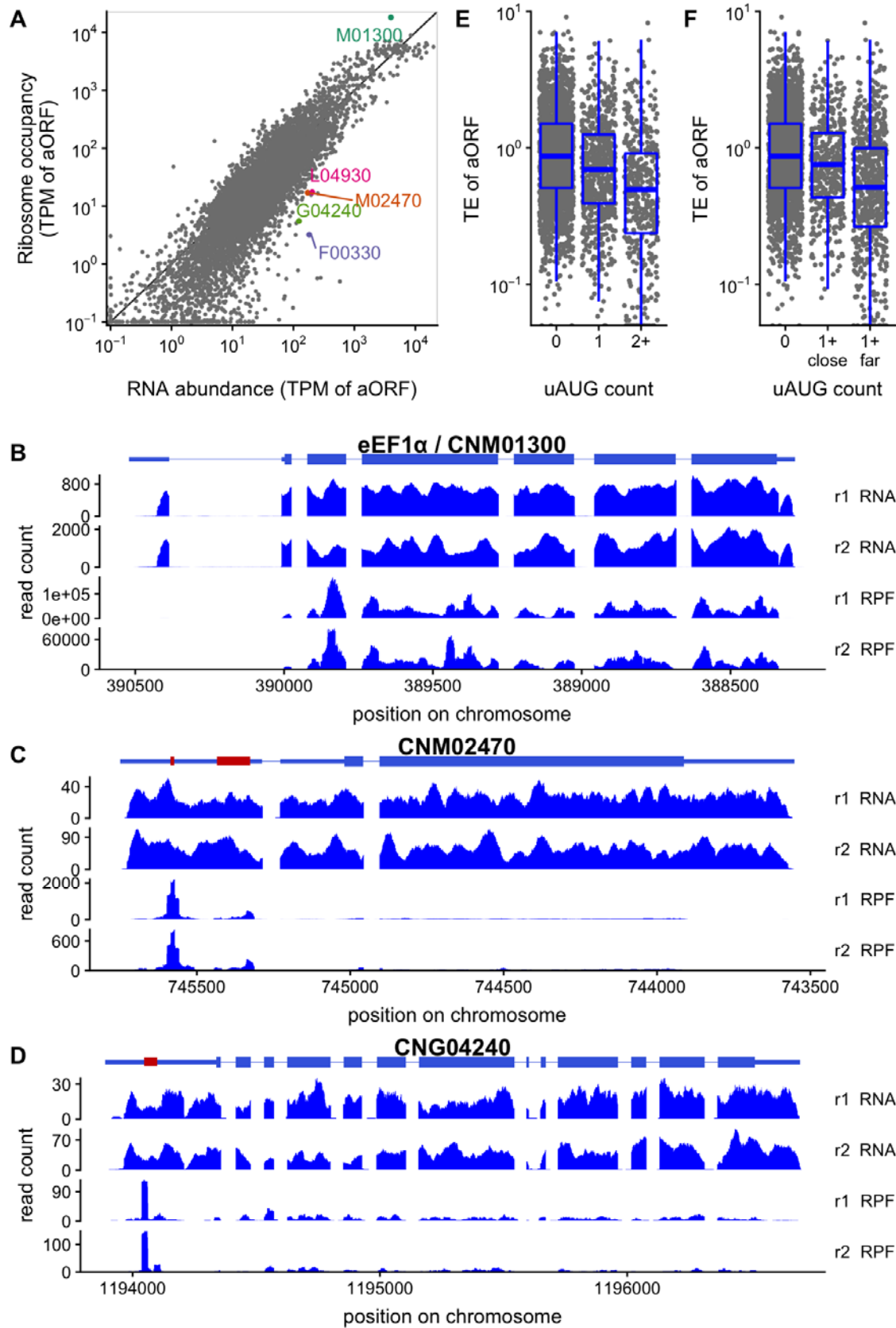


Figure S2.2, related to figure 2: Upstream AUGs repress translation in *C. deneoformans*. A, translation regulation of annotated ORFs (aORFs) in *C. deneoformans* JEC21 growing exponentially in YPD at 30°C. Ribosome occupancy is plotted against the RNA abundance, both calculated in transcripts per million (TPM) on the aORF. B, uAUGs are associated with lower translation efficiency (TE) of annotated ORFs, measured as the ratio of ribosome occupancy to RNA-seq reads. C, only uAUGs far from the transcription start site are associated with low TE. A gene is in the “1+ far” category if it has at least one uAUG more than 20nt from the TSS, “1+ close” if all uAUGs are within 20nt of the TSS. D-F, Examples of ribosome occupancy profiles along select RNAs highlighted in A (others are shown in Fig S2B). D, Translation elongation factor eEF2/CNM01300 (CNAG_06125 homolog) has high ribosome occupancy in the annotated ORF. Translationally repressed mRNAs CNM02470 (CNAG_06246 homolog, E) and CNG04240 (CNAG_03140 homolog, F) have high ribosome occupancy in uORFs in the transcript leader (red), and low ribosome occupancy in the aORF. Only the first of 5 uORFs in CNG04240 is shown, and only transcript isoform t01 is shown, excluding the annotated TL intron in isoform t02.

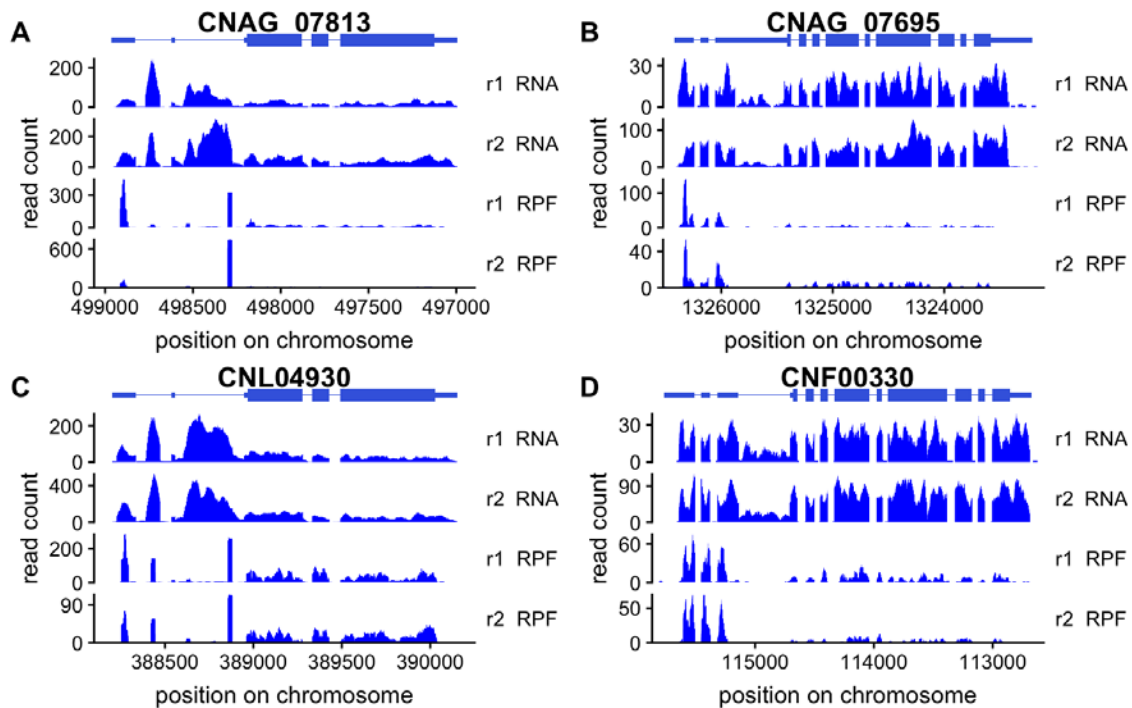


Figure S2.3, related to figure 2: Further examples of upstream AUG and 5'-end regulation in *C. neoformans* and *C. deneoformans*. CNAG_07813 (A) and CNL04930 (C) are paralogs, and in addition to an upstream ORF with ribosome occupancy, they have an intronically encoded non-coding RNA in the TL. CNAG_07695 (B) and CNF00330 (D) are paralogs, and in addition to an upstream ORF with ribosome occupancy, they have an alternatively-spliced intron in the TL that is not occupied by ribosomes.

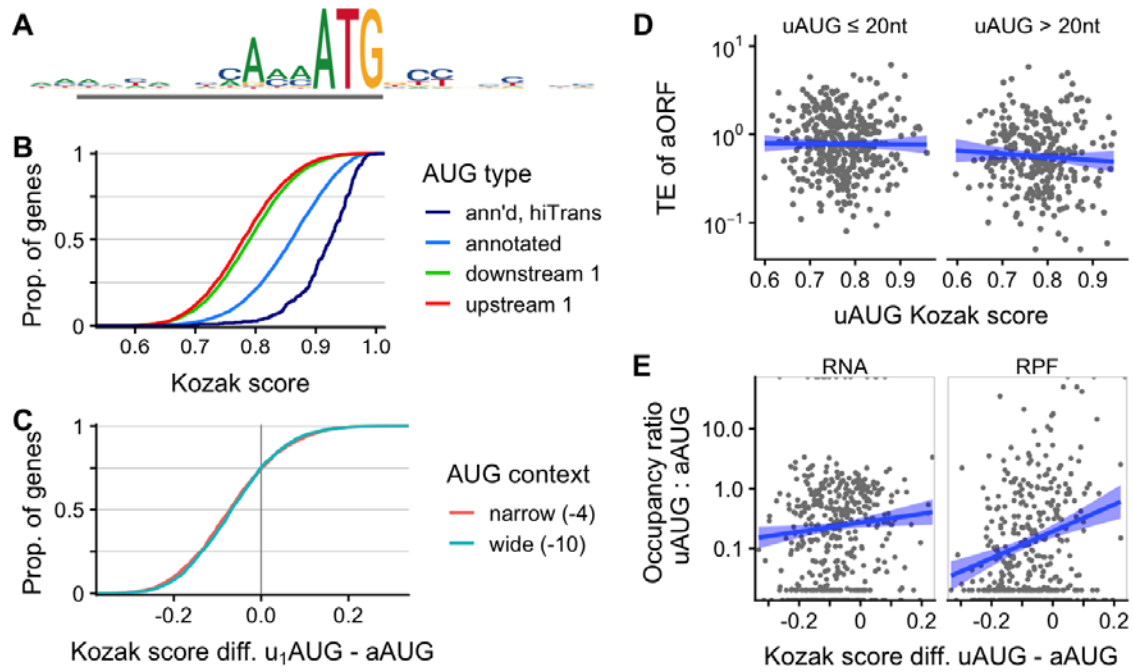


Figure S3.1, related to Figure 3: AUG sequence context is associated with translation in *C. deneoformans*. As for Fig 3, but with data from *C. deneoformans* JEC21.

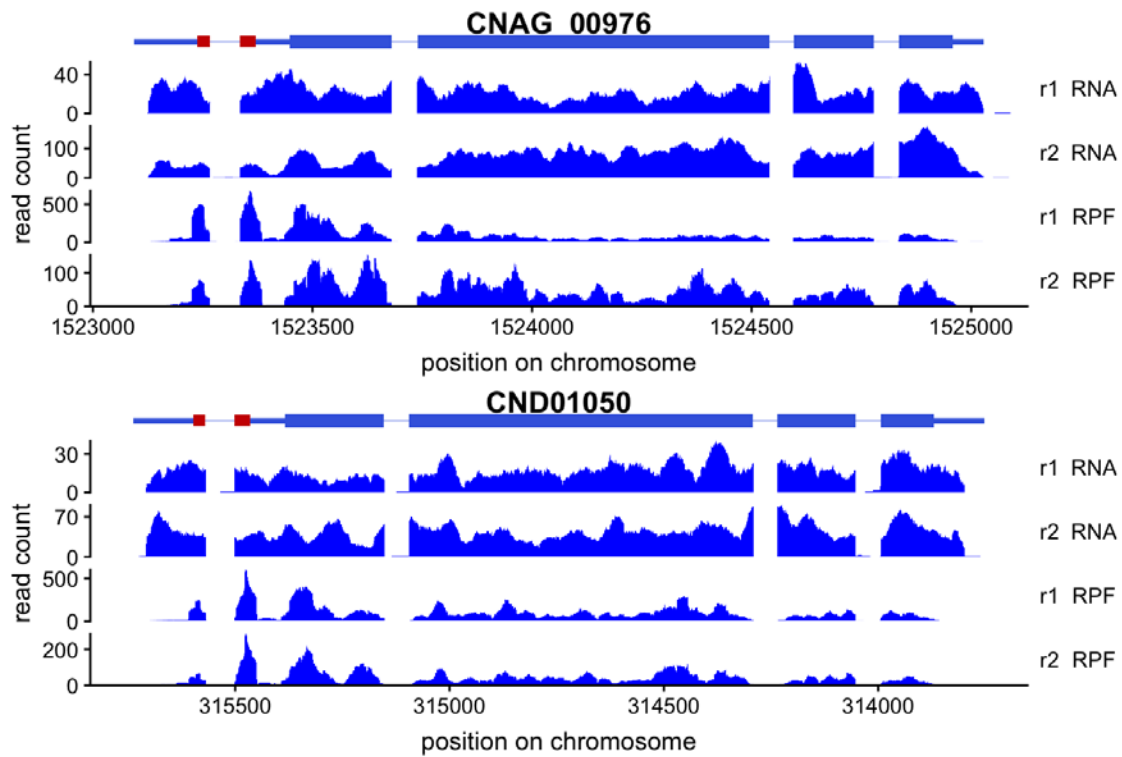


Figure S4.1, related to Figure 4: Carbamoyl-phosphate synthase CPA1 homologs have a conserved uORF that is occupied by ribosomes in *C. neoformans* (A) and *C. deneoformans* (B).

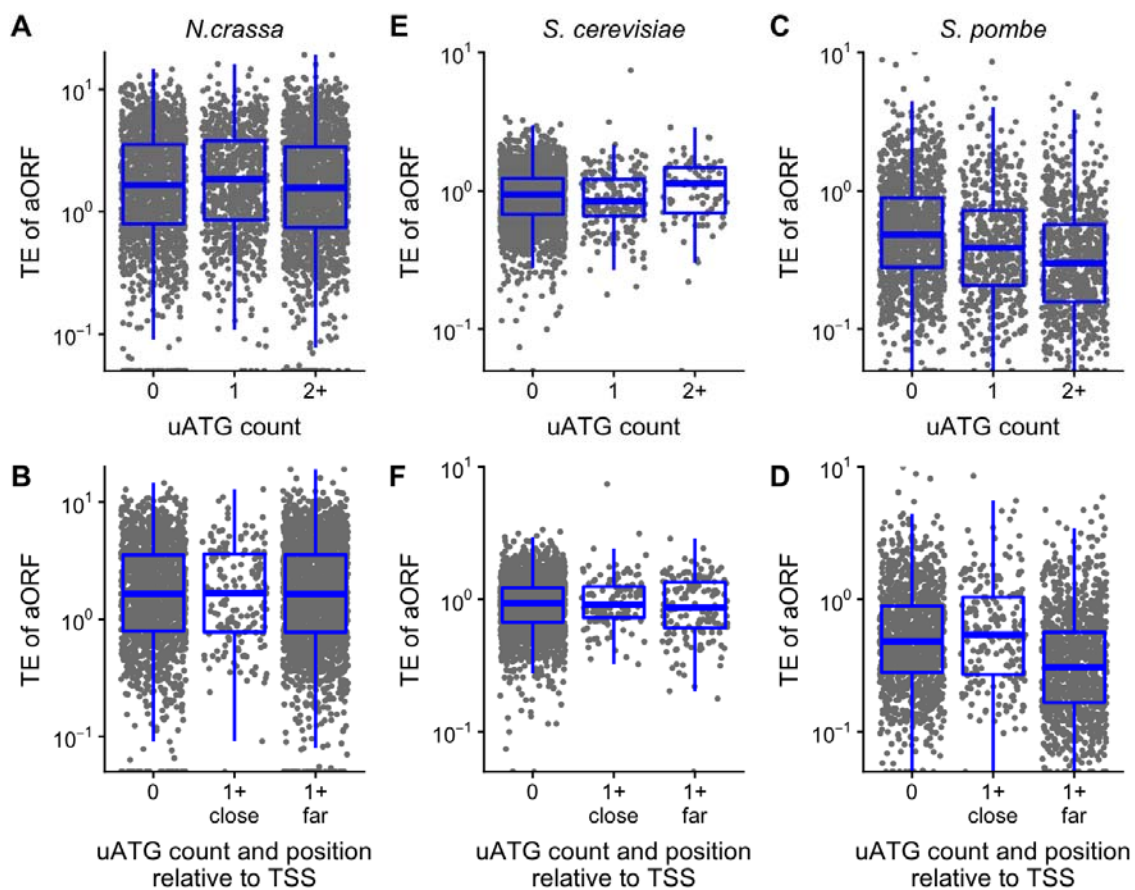


Figure S5.1, related to figure 5: Effect of uATGs on translational efficiency in *N. crassa*, *S. pombe*, and *S. cerevisiae*.

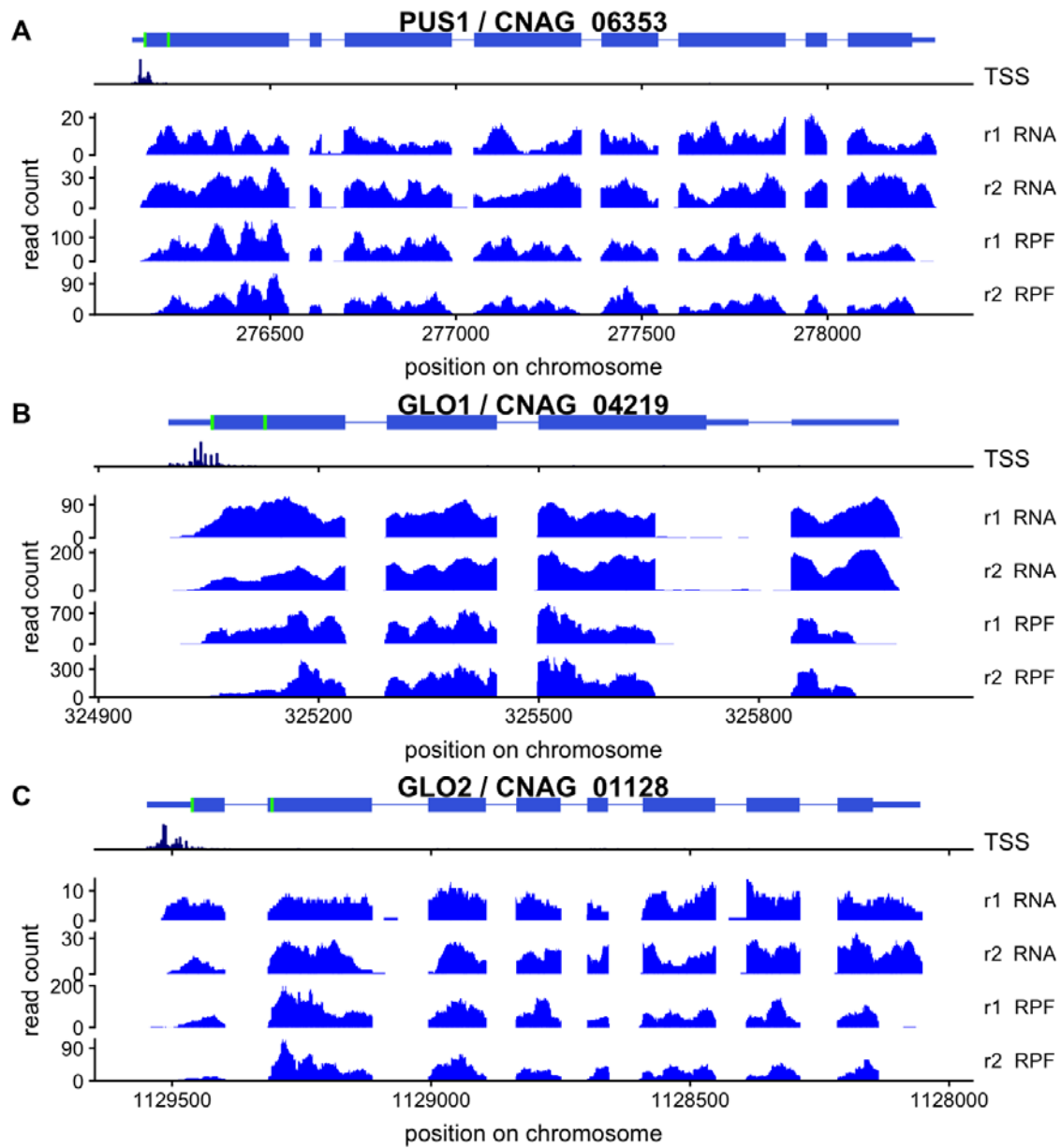


Figure S6.1, related to Figure 6: Ribosome profiles of *C. neoformans* genes with predicted dual-localization specified by alternative N-termini.

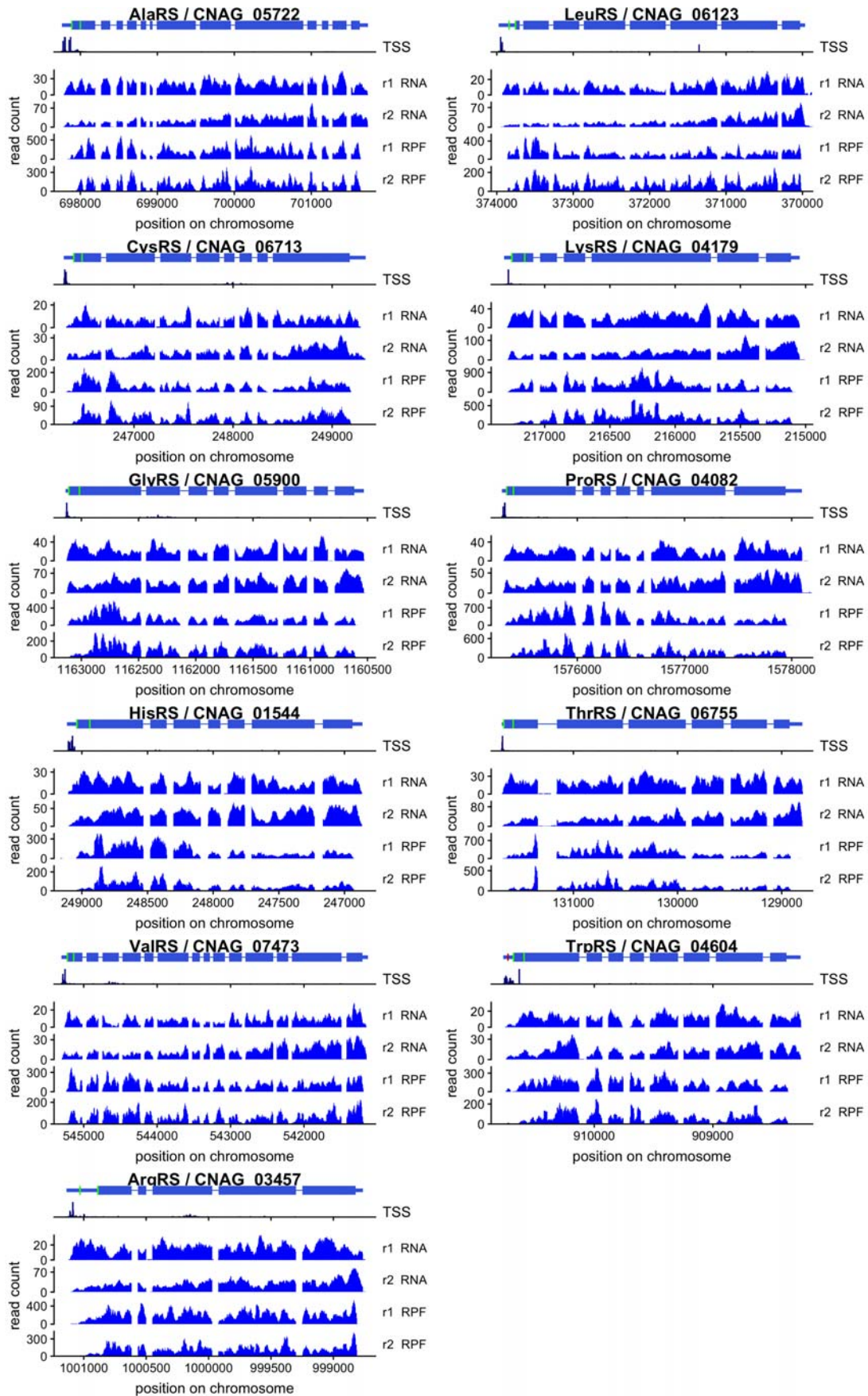


Figure S7.1, related to Figure 7: Ribosome profiles along the 11 *C. neoformans* aaRS genes with predicted dual-localization. Predicted start codons are shown in green, and the uORF of TrpRS in dark red.

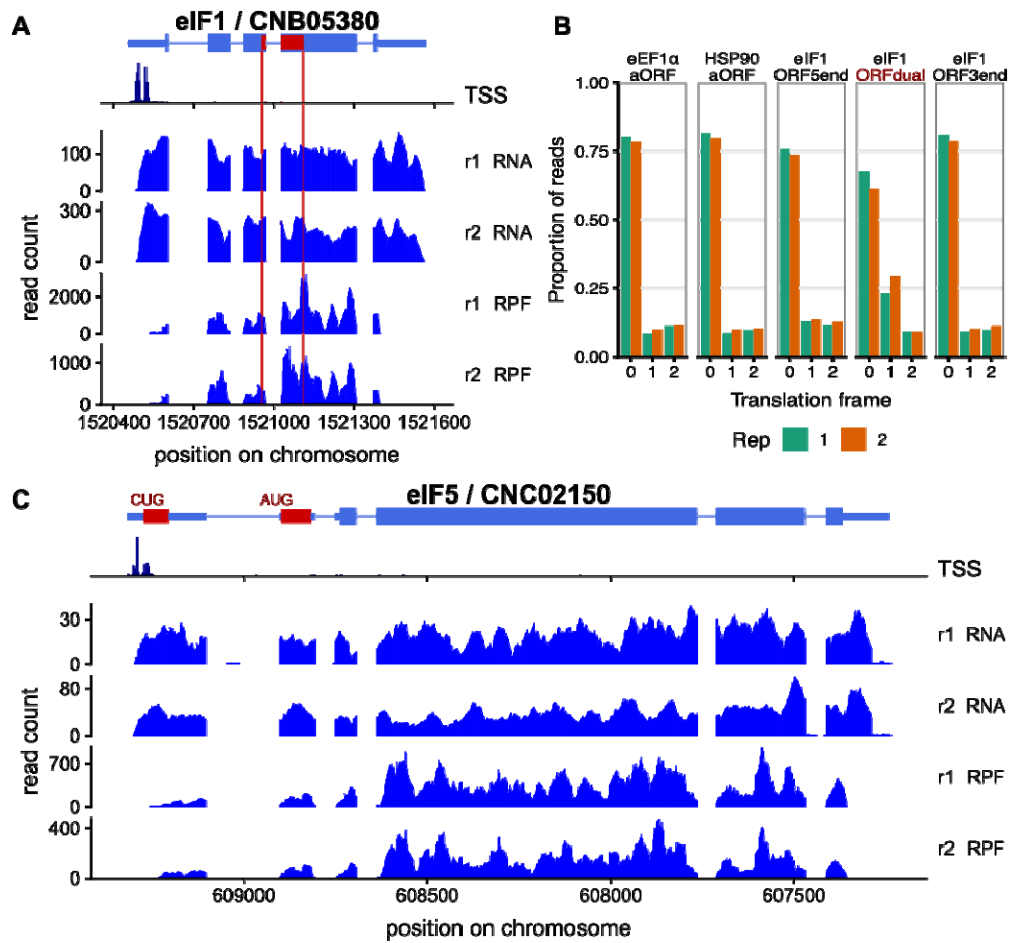


Figure S8, related to Figure 8: Translation initiation factors eIF1 and eIF5 are regulated by alternate start codon usage in *C. deneoformans*. See figure 8 for legend.

Table S1: Sequencing and annotation numbers (Excel table).

Table S2: Differential expression results in *Cryptococcus deneoformans upf1Δ*.

Table S3: Cytoplasmic ribosomal proteins in 6 fungal species.

Table S4: Genes with score d1AUG – aAUG > 0.1, n = 167 (Excel table).

Table S5: List of aaRS in *Cryptococcus* and select fungi (Excel table).

Table S6: Initiation contexts of annotated and downstream AUGs in 9 *Cryptococcus* aaRSs

Table S7: Initiation factor 3 components in 12 fungal species. (Excel table)

We show homologs of all human eIF3 components eIF3a-eIF3m, assembled mostly from PANTHERdb (51). Note that the *C. neoformans* homolog *PRT1* is expressed in two non-identical paralogs at the mating-type locus, so does not have a systematic ORF name.

Table S8: rRNA subtraction oligos used in ribosome profiling.