



**HAL**  
open science

## Quantitative global studies reveal differential translational control by start codon context across the fungal kingdom

Edward W J Wallace, Corinne Maufrais, Jade Sales-Lee, Laura R Tuck, Luciana de Oliveira, Frank Feuerbach, Frederique Moyrand, Prashanthi Natarajan, Hiten D Madhani, Guilhem Janbon

### ► To cite this version:

Edward W J Wallace, Corinne Maufrais, Jade Sales-Lee, Laura R Tuck, Luciana de Oliveira, et al.. Quantitative global studies reveal differential translational control by start codon context across the fungal kingdom. *Nucleic Acids Research*, 2020, 48 (5), pp.2312-2331. 10.1093/nar/gkaa060 . pasteur-02651962

**HAL Id: pasteur-02651962**

**<https://pasteur.hal.science/pasteur-02651962>**

Submitted on 29 May 2020

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Quantitative global studies reveal differential translational control by start codon context across the fungal kingdom

Edward W. J. Wallace<sup>1,\*</sup>, Corinne Maufrais<sup>2,3,†</sup>, Jade Sales-Lee<sup>4</sup>, Laura R. Tuck<sup>1</sup>,  
Luciana de Oliveira<sup>2</sup>, Frank Feuerbach<sup>5</sup>, Frédérique Moyrand<sup>2</sup>, Prashanthi Natarajan<sup>4</sup>,  
Hiten D. Madhani<sup>4,6,\*</sup> and Guilhem Janbon<sup>2,\*</sup>

<sup>1</sup>Institute for Cell Biology and SynthSys, School of Biological Sciences, University of Edinburgh, UK, <sup>2</sup>Institut Pasteur, Unité Biologie des ARN des Pathogènes Fongiques, Département de Mycologie, F-75015 Paris, France, <sup>3</sup>Institut Pasteur, HUB Bioinformatique et Biostatistique, C3BI, USR 3756 IP CNRS, F-75015 Paris, France, <sup>4</sup>Department of Biochemistry and Biophysics, University of California at San Francisco, San Francisco, CA 94158, USA, <sup>5</sup>Institut Pasteur, Unité Génétique des Interactions Macromoléculaire, Département Génome et Génétique, F-75015 Paris, France and <sup>6</sup>Chan-Zuckerberg Biohub, San Francisco, CA 94158, USA

Received July 12, 2019; Revised January 13, 2020; Editorial Decision January 19, 2020; Accepted January 20, 2020

## ABSTRACT

Eukaryotic protein synthesis generally initiates at a start codon defined by an AUG and its surrounding Kozak sequence context, but the quantitative importance of this context in different species is unclear. We tested this concept in two pathogenic *Cryptococcus* yeast species by genome-wide mapping of translation and of mRNA 5' and 3' ends. We observed thousands of AUG-initiated upstream open reading frames (uORFs) that are a major contributor to translation repression. uORF use depends on the Kozak sequence context of its start codon, and uORFs with strong contexts promote nonsense-mediated mRNA decay. Transcript leaders in *Cryptococcus* and other fungi are substantially longer and more AUG-dense than in *Saccharomyces*. Numerous *Cryptococcus* mRNAs encode predicted dual-localized proteins, including many aminoacyl-tRNA synthetases, in which a leaky AUG start codon is followed by a strong Kozak context in-frame AUG, separated by mitochondrial-targeting sequence. Analysis of other fungal species shows that such dual-localization is also predicted to be common in the ascomycete mould, *Neurospora crassa*. Kozak-controlled regulation is correlated with insertions in translational initiation factors in fidelity-determining regions that contact the initiator tRNA. Thus, start codon context

is a signal that quantitatively programs both the expression and the structures of proteins in diverse fungi.

## INTRODUCTION

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

\*To whom correspondence should be addressed. Tel: +44 131 6513348; Email: edward.wallace@ed.ac.uk

Correspondence may also be addressed to Hiten D. Madhani. Email: hiten.madhani@ucsf.edu

Correspondence may also be addressed to Guilhem Janbon. Email: guilhem.janbon@pasteur.fr

†The authors wish it to be known that, in their opinion, the first two authors should be regarded as joint First Authors.

ORFs are translated to protein, and are biased against short ORFs (11).

The rule of thumb that first AUG of an ORF is used as the start codon can be wrong in both directions: poor-context AUGs may be skipped, and non-AUGs may be used. The rules for selection of start codons in eukaryotes were discovered by Kozak: the most 5' AUG is generally selected if it is far (>20nt) from the transcription start site, and in metazoans efficient selection is associated with a sequence context gccRccAUGG, where the R indicates a purine at the -3 position (12,13). In *S. cerevisiae*, translation likewise generally starts at the 5' AUG in the mRNA sequence, and early studies suggested that aaaAUG was the preferred sequence context despite 'only modest impacts of flanking nucleotides on AUG start codon selection', as reviewed in (14). By contrast, it is now clear that AUG sequence context strongly modulates protein output from reporter mRNAs (15–17), and AUG-proximal context alone predicts over 15% of the genome-wide variation in translation efficiency (18). Usage of non-AUG start codons has also been observed in diverse eukaryotes, including the fungi *Saccharomyces cerevisiae*, *Candida albicans*, *Schizosaccharomyces pombe* and *Neurospora crassa* (19–24), and in *S. cerevisiae*, sequence context has a large effect on the usage of non-AUG start codons (25).

Weak or inefficient start codons near the 5' end of mRNA can give rise to translational regulation, explained by the scanning model of eukaryotic translation initiation. Translation starts by the pre-initiation complex binding mRNA at the 5' cap and then scanning the transcript leader (TL) sequence in a 3' direction until it identifies a start codon, at which translation initiates (26). Here, we call the 5' regulatory region of mRNA the TL rather than the 5' UTR, because short 'upstream' ORFs in this region can be translated (27). The pre-initiation complex sediments at 43S, and comprises the small ribosomal subunit, methionyl initiator tRNA and numerous eukaryotic translation initiation factors (eIFs). Biochemical, genetic and structural data indicate that eIF1 and eIF1A associate with the 43S pre-initiation complex (28,29). Recognition of the start codon involves direct interactions of eIF1 and eIF1A with the start codon context and initiator tRNA within a larger 48S pre-initiation complex. Start codon selection occurs when eIF1 is replaced by eIF5's N-terminus (30), then eIF2 is released, the large ribosomal subunit joins catalyzed by eIF5B and translation begins (29). This work has been largely driven by studies in *S. cerevisiae* and metazoans. Although the core protein and RNA machinery of eukaryotic translation initiation is highly conserved, it is not understood how fungi quantitatively vary in the sequence, structure, and function of their translation initiation machinery.

*Cryptococcus* are basidiomycete yeasts with a high density of introns in their coding genes (31). These introns influence gene expression and genome stability (32–34). The current genome annotation of pathogenic *C. neoformans* and *C. deneoformans* reference strains are based on both automatic and manual curations of gene structures using RNA-Seq data (35,36). Although the high degree of interspecies conservation of intron numbers and positions within coding sequences suggest that these annotations are reliable (36), the regulatory regions (transcript leader and 3'

UTRs) at transcript extremities are less well identified. In fact, most fungal genomes lack complete transcript annotations, thus we do not know how regulatory structure varies across fungi.

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

## MATERIALS AND METHODS

### DNA and RNA purification, sequencing library preparation

*Cryptococcus neoformans* strain H99 and *C. deneoformans* strain JEC21 were grown in 100 ml YPD at 30°C or 37°C under agitation up to exponential or early stationary phase as previously described (35). Briefly, early stationary phase was obtained after 18 h of growth (final OD<sub>600</sub> = 15) starting from at OD<sub>600</sub> = 0.5. *Cryptococcus deneoformans* strain NE579 (*upf1* Δ) (36) was grown in YPD at 30°C under agitation in exponential phase. Each *Cryptococcus* cell preparation was spiked in with one tenth (OD/OD) of *S. cerevisiae* strain FY834 (37) cells grown in YPD at 30°C in stationary phase. Cells were washed, snap frozen and used to prepare RNA and total DNA samples as previously described (38,39). Briefly, total DNA was extracted by bead-beating and phenol:chloroform extraction, and RNA was extracted from lyophilized cells using Trizol. Each condition was used to prepare biological triplicate samples.

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

TSS-Seq libraries preparations were performed starting with 75 μg of total RNA as previously described (40) replacing the TAP enzyme by the Cap-clip Pyrophosphatase Acid (TebuBio). For each *Cryptococcus* species we also constructed a control 'no decap' library.

Briefly, for these control libraries, poly A RNAs were purified from 75 μg of RNA from *Cryptococcus* and 75 μg of RNA from *S. cerevisiae* before being dephosphorylated using Antarctic phosphatase. Then, *S. cerevisiae* RNAs and one half of the RNAs extracted from *Cryptococcus* were treated with Cap-clip Pyrophosphatase Acid enzyme. The second half of *Cryptococcus* RNAs was mock treated. Each half of Cap-clip Pyrophosphatase Acid *Cryptococcus* RNA

samples was mixed with the same quantity of *S. cerevisiae* Cap-clip Pyrophosphatase Acid treated RNAs. The subsequent steps of the library preparation were identical to the published protocol (40). Fifty base single end reads were obtained using an Illumina HiSeq2500 instrument according to the manufacturer's instructions (Illumina).

For QuantSeq 3'mRNA-Seq preparation we followed the manufacturer's instructions for the QuantSeq fwd kit (Lexogen GmbH, Austria). One hundred base single end reads were obtained using an Illumina HiSeq2000 instrument according to the manufacturer's instructions (Illumina).

### Sequencing data analyses

For TSS analysis, we kept only the reads containing both the oligo 3665 (AGATCGGAAGAGCACACGTCTGAA C) and the 11NCGCCGCGNNN tag (40). These sequences were removed and the trimmed reads were mapped to the *Cryptococcus* genome and *S. cerevisiae* genomes using Bowtie2 and Tophat2 (41). Their 5' extremities were considered as potential TSSs. For each condition, we kept only the positions that were present in all three replicates. Their coverage was normalized using the normalization factor used for spiked in RNA-Seq. TSS positions were then clustered per condition. As most of the observed TSS sites appeared as clusters, we grouped them into clusters by allowing an optimal maximum intra-cluster distance (at 50 nt) between sites as previously used (40). We then removed the false TSS clusters using the 'no-cap' data keeping the clusters  $i$  for which

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

Similarly, QuantSeq 3'mRNA-Seq reads containing both the Sequencing and indexing primers (Lexogen) were sorted. The reads were then cleaned using cutadapt/1.18 (42) and trimmed for polyA sequence in their 3' end. PolyA untrimmed and trimmed reads were mapped to the adapted *Cryptococcus* and to the *S. cerevisiae* genomes with Tophat2 (41) with the same setting as for RNA-Seq. To eliminate the polyadenylated reads corresponding to genomic polyA stretches, we considered only the reads that aligned to the genomes after polyA trimming but not before the trimming. The 3' end position of these reads were considered as potential PAS. As for the TSS, for each condition we kept only the positions that were present in all three replicates. Similarly, the PAS dataset was normalized using the spike in normalization factor and the PAS positions were clustered using the same strategies.

### Ribosome profiling and matched mRNA-seq

Ribosome profiling (riboprofiling) was performed on both *C. neoformans* H99 and *C. deneoformans* JEC21, two biological replicates of WT-H99 and one replicate each of H99 *ago1*  $\Delta$  and H99 *gwo1*  $\Delta$  strains from (33), and one replicate each of WT-JEC21 and JEC21 *ago1*  $\Delta$ . We detected negligible differential abundance between these deletions and their background strains, so in our analyses we treat the deletion strains as biological replicates.

Cells were grown to exponential phase in 750 ml of YPAD with shaking at 30°C. 100  $\mu\text{g/ml}$  cycloheximide (Sigma) (dissolved in 100% ethanol) was added to the culture and incubated for 2 min. 50 ml of the culture was withdrawn for performing RNA-Seq in parallel. Cells were then pelleted, resuspended in 5 ml of lysis buffer (50 mM Tris-HCl pH. 7.5, 150 mM NaCl, 10 mM MgCl<sub>2</sub>, 5 mM DTT, 0.5% Triton and 100  $\mu\text{g/ml}$  cycloheximide) and snap frozen. Lysis, clarification, RNaseI digestion, sucrose gradient separation and monosome isolation was performed as previously described (43).

Ribosome protected fragments were isolated from the monosome fraction using hot phenol. 150  $\mu\text{g}$  of the total RNA extracted from the 50 ml of culture in parallel was polyA selected using the Dynabeads mRNA purification kit (Thermo Fisher Scientific) and digested using freshly made fragmentation buffer (100 mM NaCO<sub>3</sub> pH. 9.2 and 2 mM EDTA) for exactly 20 min.

RNA was resolved on a 15% TBE-urea gel. A gel slab corresponding to 28–34 nt was excised for footprint samples and ~50 nt for mRNA samples, then eluted and precipitated. Sequencing libraries were generated from the RNA fragments as described in Dunn *et al.* (44) with the following modifications. cDNA was synthesized using primer oCJ11 (Supplementary Table S8). Two rounds of subtractive hybridization for rRNA removal was done using oligos asDNA1-8 (Supplementary Table S8). After circularization Illumina adaptors were added through 9 cycles of PCR. Libraries were sequenced on a HiSeq 2500 (Illumina).

### Ribosome profiling data analysis

Riboprofiling and matched RNA-seq reads were demultiplexed on BaseSpace (Illumina) and then analyzed essentially with the RiboViz pipeline v.1.1.0 (45). In brief, sequencing adapters were removed with cutadapt (42), and then reads aligned to rRNA were removed by alignment with hisat2 (46). Cleaned non-rRNA reads were aligned to (spliced) transcripts with hisat2 (46), sorted and indexed with samtools (47), and then quantified on annotated ORFs with bedtools (48), followed by calculation of transcripts per million (TPM) and quality control with R (49) scripts included in RiboViz. The cleaned non-rRNA reads were also aligned to the genome with hisat2, and processed analogously, then used to generate figures of genome alignments using ggplot2 (50) in R (49).

### Data analysis and visualization

Data analysis and visualization were scripted in R (49), making extensive use of dplyr (51), ggplot2 (50) and cowplot (52). Sequence logos were prepared in gseqlogo (53). Analysis of differential mRNA abundance for *upf1*  $\Delta$  data was performed in DeSeq2 (54). Some figures were assembled and annotated in Inkscape v0.92 (<https://inkscape.org>).

Protein sequences were aligned using muscle (55), with default parameters for protein sequences and 100 iterations. Phylogenetic trees were constructed using ClustalW2 tool v2.1 (56) by using the neighbor-joining method with 1000 bootstrap trial replications.

Structural figures were prepared in PyMOL (Schrödinger).

## External datasets

*Neurospora crassa* (strain OR74A) riboprofiling data from ((22), GEO:GSE97717) was used to generate highly-translated genes, and riboprofiling and RNA-seq data from ((57), GEO: GSE71032) used to estimate TE. In both cases, we estimated TPMs using the RiboViz pipeline as above, using the NC12 genome annotation downloaded from EnsemblGenomes (58). TL sequences were also obtained from NC12.

*Schizosaccharomyces pombe* (strain 972h) riboprofiling and RNA-seq data are from (59), and the authors provided us with a table of RPKMs for all replicates as described. Genome sequence and annotation ASM294v2, including TL annotation, were downloaded from EnsemblGenomes (58).

*Candida albicans* (strain SC5314) riboprofiling and RNA-seq data are from (60), GEO:GSE52236), processed with the RiboViz pipeline as above using the assembly 22 of the strain SC5414 genome annotation from CGD (61).

*Saccharomyces cerevisiae* (strain S288C/BY4741) highly-translated mRNAs use the RPKM table from ((62), GEO:GSE59573), and highly-abundant mRNAs use (63). For TE estimates, we used matched riboprofiling and RNA-seq estimates from (64), although we did not use this for the list of highly translated genes because near-duplicate paralogous ribosomal protein genes were not present in the dataset, which thus omits a substantial fraction of highly-translated genes. TL sequences were downloaded from SGD (65).

Protein homolog lists were assembled with OrthoDB (66) and PANTHERdb (67), with reference to FungiDB (68). The list of cytoplasmic ribosomal proteins was assembled in *S. cerevisiae* based on (69) with help from SGD (65), extended to other fungi with PANTHERdb (47), and manually curated.

## RESULTS

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

To annotate the extremities of the coding genes in *C. neoformans* and *C. deneoformans*, we mapped the 5' ends (Transcription Start Sites; TSS) with TSS-Seq (40), the 3' ends (Polyadenylation sites; PAS) with QuantSeq 3'mRNA-Seq, and sequenced the same samples with stranded mRNA-Seq. These experiments were done in biological triplicate from cells growing at two temperatures (30°C and 37°C) and two stages of growth (exponential and early stationary phases) with external normalization with spike-in controls.

We identified  $4.7 \times 10^6$  unique TSSs and  $6.3 \times 10^4$  unique PASs in *C. neoformans*. Clustering of these positions revealed between 27 339 and 42 720 TSS clusters and between 9217 and 16 697 PAS clusters depending on the growth conditions (Supplementary Table S1). We used the clusters associated with the coding genes to produce an initial annotation, using the most distal TSS and PAS clusters for each gene. The predicted positions which changed the extremities of the genes by >100 bp were manually curated ( $n = 1131$  and  $n = 286$  for the TSS and PAS, respectively). We

then selected the most prominent clusters that represented at least 10% of the normalized reads count per coding gene in at least one condition (i.e. sum across three normalized replicate samples), for wild-type strains. Finally, the most distal of these TL-TSS and 3'UTR-PAS clusters were labeled as the 5' and 3' ends of the coding genes for our final annotation (Supplementary Table S1). For the genes for which no TL-TSS cluster or no 3'UTR-PAS cluster could be identified, we maintained the previous annotation. We used the same strategies for *C. deneoformans* and obtained similar results (Supplementary Table S1).

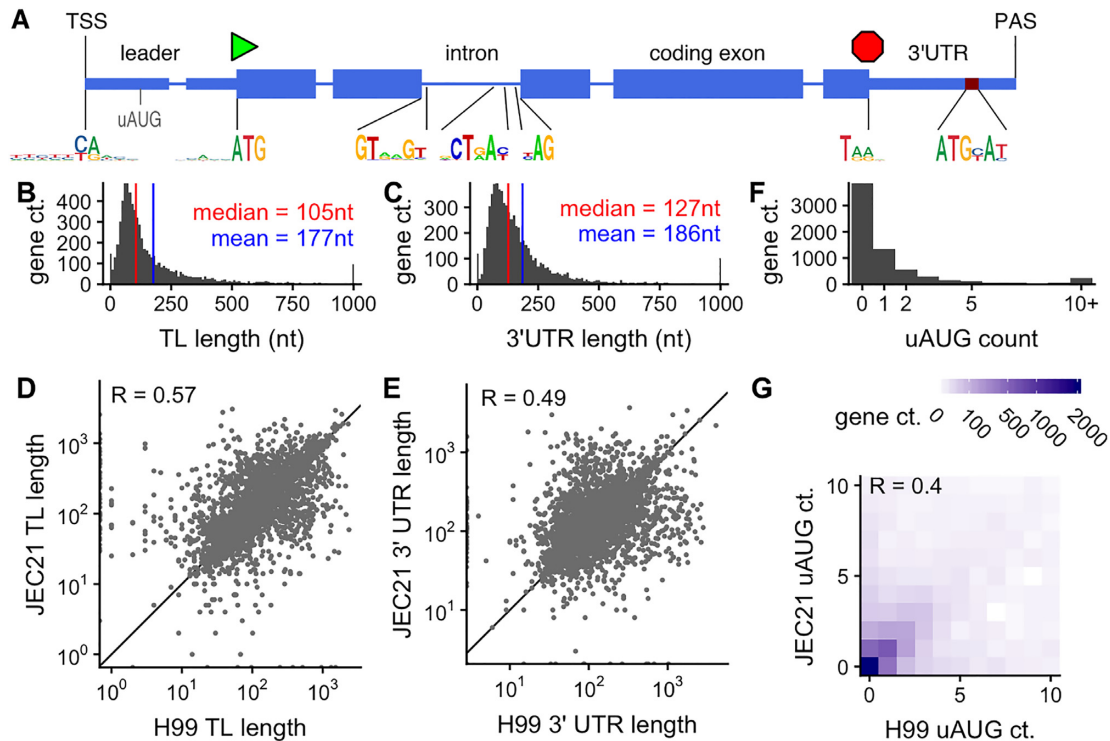
As expected, most of the TSS clusters (62%) were associated with the TL whereas most of the PAS clusters (82%) were associated with the 3'UTR of the coding genes (Supplementary Table S1). We analyzed the 3'UTR sequences, confirming the ATGHAH motif associated with the PAS (35). In addition, as previously observed in other systems (70) a (C/T)(A/G)-rich motif was associated with the maxima of these transcription start site clusters. Overall, 89% of the coding genes have both their TL and 3'UTR sequences supported by identified TSS and PAS clusters, respectively.

The analysis leads to a scheme of a stereotypical *C. neoformans* coding gene (Figure 1A). In average, it is 2305 bp long (median 2008 bp) and contains 5.6 short introns (median 5) in its sequence. As previously reported (31), these introns are short (63.4 nt in average) and associated with conserved consensus motifs. The *C. neoformans* TL and 3'UTR have median lengths of 105 nt and 127 nt, respectively (177 nt and 186 nt, mean; Figure 1B, C). Only 887 and 429 genes contain one or more introns in their TL and 3'UTR sequence, respectively; these introns are usually larger (118.3 nt) than those that interrupt the CDS. This gene structure is similar in *C. deneoformans* (Supplementary Table S1) and there are good correlations between the 3'UTR and TL sizes of the orthologous genes in the two species (Figure 1D, E).

### More than a third of genes have upstream AUGs that affect translation

The analysis of the TL sequences in *C. neoformans* revealed the presence of 10 286 AUG triplets upstream (uAUG) of the annotated translation start codon (aAUG). We include uAUGs that are either out-of-frame from the start codon, or in-frame but with an intervening stop codon, which are very unlikely to encode a continuous polypeptide. Strikingly, 2942 genes possess at least one uAUG, representing 43% of the genes with an annotated TL in *C. neoformans* (Figure 1F). A similar result was obtained in *C. deneoformans*, in which we found 10 254 uAUGs in 3057 genes, and uAUG counts are correlated between orthologous mRNAs in the two species (Figure 1G). This is consistent with previous findings of conserved uAUG-initiated ORFs in *Cryptococcus* species (71).

Translation initiation at uAUGs results in the translation of uORFs, which can regulate translation of the main ORF (43,72). To evaluate the functionality of the uAUGs in *Cryptococcus*, we generated riboprofiling data in both species and compared densities of ribosome-protected fragments with those of sample-matched poly(A)+ RNA. Our riboprofiling data passes quality metrics of 3-nucleotide periodicity of reads on ORFs indicating active translation by



**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. (D, E) 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*.

ribosomes, and appropriate read lengths of 26–30 nt (Supplementary Figure S1).

Most genes have ribosome occupancy close to that predicted by their RNA abundance, and restricted to the main ORF, for example the most highly translated gene, translation elongation factor eEF1 $\alpha$ /CNAG\_06125 (Figure 2A, B). However, we observed dramatic examples of translation repression associated with uORFs in CNAG\_06246 and CNAG\_03140 in *C. neoformans* (Figure 2A, C, D). These patterns are conserved in their homologs in *C. deneoformans* (Supplementary Figure S2). Other spectacularly translationally repressed genes, CNAG\_07813 and CNAG\_07695 and their *C. deneoformans* homologs (Figure 2A, Supplementary Figure S2A) contain conserved uORFs in addition to 5' introns with alternative splicing or intronically expressed non-coding RNAs (Supplementary Figure S3). In all these cases, high ribosomal occupancy on one or more uORFs is associated with low occupancy of the main ORF.

The uncharacterized gene CNAG\_06246 has two AUG-encoded uORFs that are occupied by ribosomes, and a predicted C-terminal bZIP DNA-binding domain. This gene structure is reminiscent of the multi-uORF-regulated amino-acid responsive transcription factors Gcn4/Atf4 (72), or the *S. pombe* analog Fill1 (73). The sugar transporter homolog CNAG\_03140 has six uAUGs, with substantial ribosome occupancy only at the first. Interestingly, *N. crassa* has a sugar transporter in the same major facilitator super-

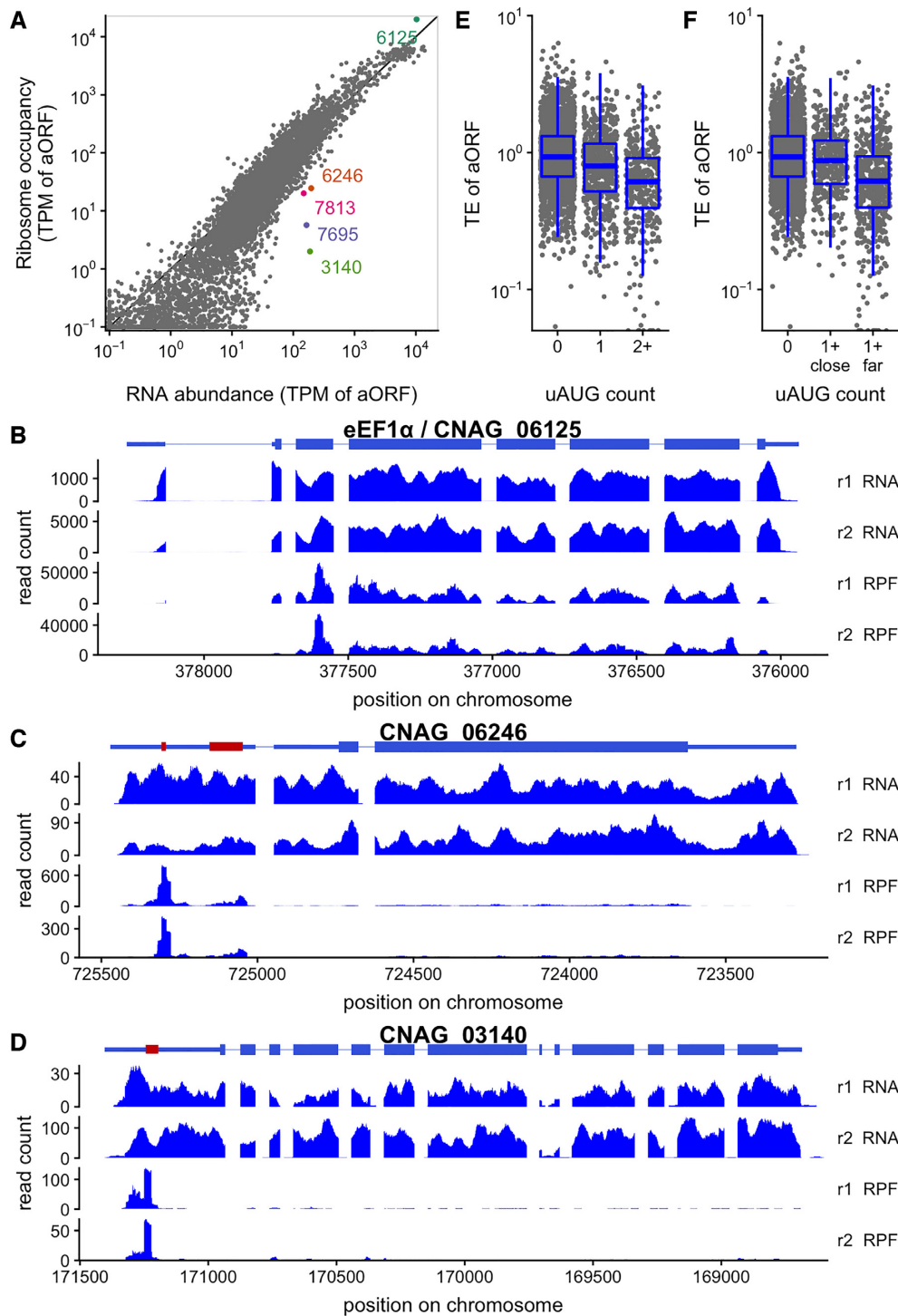
family regulated by a uORF (rco-3/sor-4, (74)) and sugar-responsive translational repression via uORFs has been observed in plants (75).

Since these translationally repressed genes have multiple uAUGs, we investigated the relationship between uAUGs and translation efficiency genome-wide. We observed a clear negative relationship between the number of uAUGs and translation efficiency (Figure 2E, Supplementary Figure S2E), suggesting an uAUG-associated negative regulation of translation in both species.

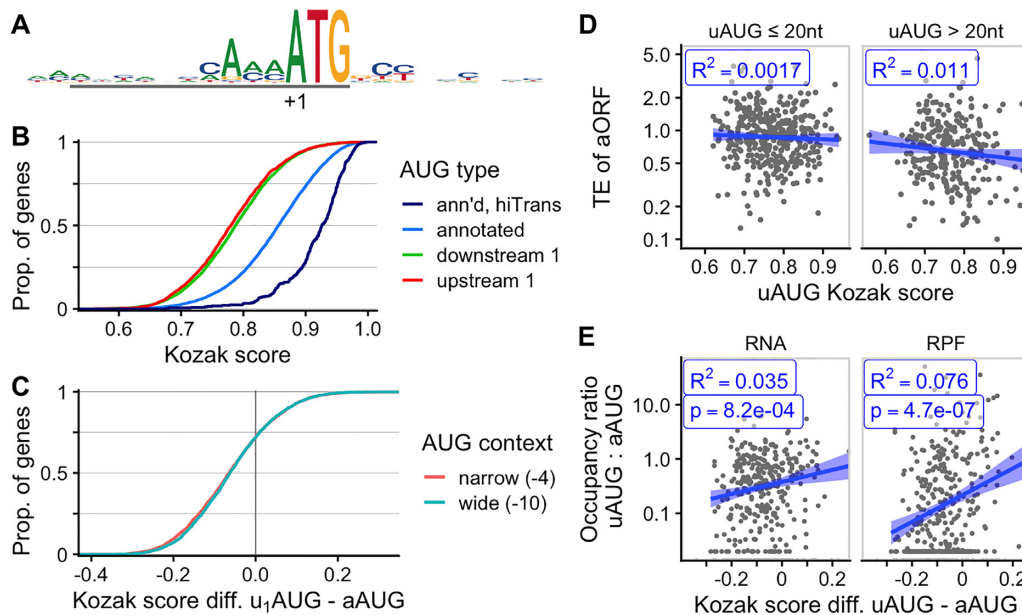
### Position relative to the TSS affects uAUG translation

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

We compared the translation efficiency of genes containing only uAUGs close to the TSS to those with uAUGs far from the TSS. In *C. neoformans*, 1627 of the 10 286 uAUGs are positioned within the first 20nt of the TL, and 816 uAUG-containing genes have no uAUG after this position. The presence of one or several uAUGs close to the TSS (<20nt) has nearly no effect on translation efficiency, whereas genes containing uAUGs far from the TSS are less efficiently translated (Figure 2F), and similarly in *C. deneoformans* (Supplementary Figure S2F).



**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 Supplementary Figure S2). 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 color. (B) Translation elongation factor eEF1 $\alpha$ /CNAG\_06125 has high ribosome occupancy in the annotated ORF. Translationally repressed mRNAs CNAG\_06246 (C) and CNAG\_03140 (D) have high ribosome occupancy in uORFs in the transcript leader (red), and low ribosome occupancy in the aORF. Only the first of five uORFs in CNAG\_03140 is shown. Other genes highlighted in panel A are shown in Supplementary Figure S3. Homologous genes in *C. deneoformans* have similar structure and regulation (Supplementary Figures S2 and S3). (E) uAUGs are associated with lower translation efficiency (TE) of annotated ORFs, measured as the ratio of ribosome occupancy to RNA-seq reads. (F) 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.



**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  $-10$  onward. (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 and not significantly associated with translation repression of the annotated ORF. (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) with  $R^2$ , and  $P$ -value of associated  $t$ -test. Panels D and E show data only for genes in the top 50% by RNA abundance, and with only a single upstream AUG. Supplementary Figure S4 shows homologous data for *C. deneoformans*.

### A Kozak sequence context determines AUG translation initiation

To analyze the importance of AUG sequence context for translation initiation in *C. neoformans*, we used the 5% most translated genes (hiTrans,  $n = 330$ ) to construct a consensus sequence surrounding their annotated translation start codon (Figure 3A). The context contains a purine at the  $-3$  position, a hallmark of the Kozak consensus sequence (26). However, there is very little enrichment for the  $+4$  position, in contrast with the mammalian Kozak context in which a G is present in  $+4$  ((A/G)CCAUGG) (26). Because of the limited sequence bias downstream of the AUG, and its confounding with signals of N-terminal amino acids and codon usage, we do not consider it further. However, we found a slight sequence bias in the positions  $-10$  to  $-7$  that is outside the metazoan Kozak context.

We thus calculated ‘Kozak scores’ for all uAUGs against the position weight matrix of the Kozak context from  $-10$  from AUG through to AUG (Figure 3A). We compared the Kozak scores of the annotated AUGs (aAUGs) with those of the 5% most highly translated genes, the first upstream AUG (uAUGs) and the first downstream AUG ( $d_1$ AUG). Highly translated aAUGs have a higher score than typical aAUGs, and aAUGs have usually a higher score than the uAUGs and  $d_1$ AUGs (Figure 3B). On a given transcript, the uAUG score is usually lower than the aAUG score (Figure 3C).

We next asked if the sequence context of uAUGs affected their ability to repress translation of the annotated ORF, focusing on transcripts with only a single uAUG. Surprisingly, there is a weak and insignificant correlation between uAUG Kozak score and the translation efficiency of the aORF, whether the uAUG is close to or far from the TSS (Figure 3D). However, the most striking examples of translational repression in Figure 2 tend to have multiple high-score uAUGs (scores CNAG.06246,  $u_1$ AUG 0.93,  $u_2$ AUG 0.86; CNAG.03140,  $u_1$ AUG 0.85,  $u_2$ AUG 0.76; CNAG.07813,  $u_1$ AUG 0.79; CNAG.07695,  $u_1$ AUG 0.97,  $u_2$ AUG 0.90). This is consistent with direct biochemical evidence that AUG context determines translation repression by uORFs in *N. crassa* and *S. cerevisiae* (76).

We also asked if the AUG score affects the AUG usage transcriptome-wide, by comparing the difference in  $u_1$ AUG and aAUG scores with the ratio in A-site ribosome occupancy in a 10-codon neighbourhood downstream of the  $u_1$ AUG and aAUG. We considered the relative occupancy to control for transcript-specific differences in abundance and cap-dependent initiation-complex recruitment. We restrict our analysis to a short neighborhood to control for start-codon specific biases in ribosome occupancy caused by addition of cycloheximide prior to cell lysis (59,62). A higher score difference is associated with higher relative ribosome occupancy, while the control comparison with RNA-Seq coverage shows a smaller effect (Fig-



ure 3E). We find the same patterns of AUG consensus, scores, and occupancy in *C. deneoformans* (Supplementary Figure S4).

### Nonsense-mediated decay acts on uORF-containing genes

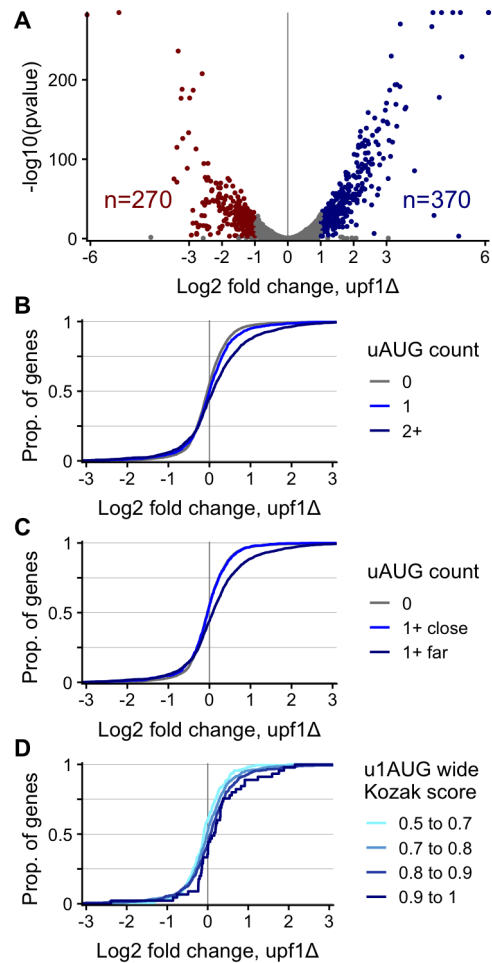
An mRNA molecule translated using an uAUG can be recognized as a premature stop codon bearing molecule and may be as such degraded by the nonsense-mediated mRNA decay (NMD) (77). In *S. cerevisiae*, uAUGs are associated with NMD genome-wide (78). To test this concept in *Cryptococcus*, we first sequenced RNA from *C. deneoformans* strains with the conserved NMD factor Upf1 deleted (36), finding 370 genes with increased mRNA abundance and 270 with decreased (Figure 4A, Supplementary Table S2; 2-fold difference in levels at 1% FDR).

We next compared the fold-change in abundance of uAUG-containing or uAUG-free mRNAs. Two genes with extreme increases in *upf1*Δ are also extremely translationally repressed uORF-containing genes we identified above (Figures 2, Supplementary Figures S2 and S3): CNF00330 (CNAG.07695 homolog, 11-fold) and CNG04240 (CNAG.03140 homolog, 8-fold). Another extreme is the carbamoyl-phosphate synthase CND01050 (5-fold up in *upf1*Δ), a homolog of *S. cerevisiae* CPA1 and *N. crassa* arg-2. These orthologs are regulated by a conserved uORF encoding a arginine attenuator peptide that have all been verified to repress reporter gene synthesis in a *N. crassa* cell-free translation system (79); both *S. cerevisiae* and *N. crassa* orthologs are NMD substrates, which for *Sc*CPA1 depends on the uORF (80,81). Consistent with this model, in both *C. neoformans* and *C. deneoformans* the native uORF shows strong ribosome occupancy while the aORF is translationally repressed (*Cn*TE = 0.47, *Cd*TE = 0.38; Supplementary Figure S5).

In general, uAUG-containing genes are more likely to be upregulated in the *upf1*Δ mutant than uAUG-free genes (Figure 4B), suggesting that uORFs negatively regulate mRNA abundance in *Cryptococcus*, in addition to repressing translation of the main ORF. Similarly, uAUG-containing genes are enriched for NMD-sensitivity only when the uAUG is >20nt from the TSS (Figure 4C), suggesting that TSS-proximal uAUGs (<20nt) are skipped, and generally not used as translation start codons in *Cryptococcus*.

Next, we asked if uAUG Kozak score affects mRNA decay via the NMD pathway. Restricting our analysis to genes with a single uAUG ( $n = 1421$ ), we binned genes according to their Kozak score. We find that mRNAs that contain higher Kozak-score uAUG are more likely to increase in abundance in the *upf1*Δ mutant (Figure 4D). Indeed, the abundance increase is monotonically correlated with the mean of the score bins. This could explain the weak effect size of uAUG score on translation efficiency (Figure 3D), as higher-scoring uAUGs repress the RNA abundance (denominator of TE) in addition to repressing translation of the main ORF (numerator).

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



**Figure 4.** Nonsense-mediated decay (NMD) acts on upstream-AUG-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) uAUG containing genes are enriched for NMD-sensitivity. A one-sided Kolmogorov–Smirnov test shows that these differences are significant comparing 1 uAUG to 0 ( $P = 5.7 \times 10^{-5}$ ) and 2 or more AUGs to 1 ( $P = 7.3 \times 10^{-11}$ ). (C) uAUG-containing genes are enriched for NMD-sensitivity only when the uAUG is more than 20nts from the TSS (1+ far;  $P < 2.2 \times 10^{-16}$ ), but not when the uAUG is less than 20nts (1+ close;  $P = 0.73$ ). (D) Start codon sequence context affects NMD sensitivity of genes containing a single upstream AUG: RNAs starting with higher Kozak-score uAUG are more likely to increase in abundance in the *upf1*Δ mutant ( $P = 1.1 \times 10^{-4}$ , comparing score < 0.8 with score > 0.8).

### Start codon sequence context and uORF regulation in other fungi

We then examined sequences associated with translation start codons in other fungi, for which both RNA-Seq and riboprofiling data were available, and for which the annotation was sufficiently complete (i.e. *S. cerevisiae*; *N. crassa*, *C. albicans* and *S. pombe*). We analyzed the Kozak context associated with aAUG of all annotated coding genes, of the 5% most translated genes (hiTrans), and for mRNAs encoding cytoplasmic ribosomal proteins (CytoRibo), as a model group of highly expressed and co-regulated genes defined by homology (Supplementary Table S3). Cytoplasmic riboso-

mal proteins have informative Kozak contexts, with a strong A-enrichment at the positions  $-1$  to  $-3$  and weak sequence enrichment after the AUG in all these species (Figure 5A). The total information content of the Kozak sequence is higher for CytoRibo genes than HiTrans, and higher for HiTrans than all annotated genes, in all these fungi (Figure 5B). Nevertheless, these contexts have also some species specificity: Kozak sequences for HiTrans and CytoRibo are more informative in *Cryptococcus* and *N. crassa* than in *S. pombe*, *C. albicans* and *S. cerevisiae*. In particular, the C-enrichment at positions  $-1$ ,  $-2$  and  $-5$  in *Cryptococcus* and *N. crassa* is absent in *S. cerevisiae*, and we observed no sequence enrichment upstream of the  $-4$  position for *S. pombe* and very little for *S. cerevisiae*. In contrast, a  $-8$  C enrichment, similar to the *Cryptococcus* and mammalian pattern, was observed in *N. crassa*, confirming previous results (82). The  $-10$ – $-6$  A/T rich region for *C. albicans* is likely to reflect an overall A/T-richness of the TLs in *C. albicans*.

The analysis of the TL sequences from these fungi, excluding *C. albicans* for which no TL annotation is available, also shows species specificity. The average TL length in *S. cerevisiae* (84nt) is less than half that in *Cryptococcus* (Figure 5C). In sharp contrast with *Cryptococcus*, only 985 uAUGs are present in 504 genes, which correspond to 18% of the genes with an annotated TL in *S. cerevisiae*. Moreover, the density of the uAUGs is very low and uAUGs have no global effect on TE in this yeast (Figures 5D, E, Supplementary Figure S6). The short uAUG-depleted TLs observed in the SGD annotations of *S. cerevisiae* are conserved in a recent annotation of other *Saccharomyces* species (83) (Supplementary Figure S7).

More broadly, short TLs with very low uAUG density are more the exception than the rule in the fungal kingdom (Figure 5C). However, fungi vary in how much these uAUGs globally down regulate gene translation (Supplementary Figure S6). Our analysis shows that fungi quantitatively vary in the sequence context of the AUGs that they use, and in the distribution of AUGs in their TLs. Thus, distinct fungi may differ in how much they use AUG sequence context to regulate gene expression at the post transcriptional level.

### Kozak context programs leaky scanning in *Cryptococcus*

We earlier calculated the Kozak score of the first downstream AUG ( $d_1$ AUG) within each CDS: these  $d_1$ AUG scores are mostly lower than the score of the aAUGs (Figure 3B), consistent with most annotations correctly identifying a good-context AUG as the start codon. Yet, we identified a number of  $d_1$ AUGs with a high score ( $n = 1109$  above 0.826, the median Kozak score for aAUGs;  $n = 131$  above 0.926, the median for hiTrans), which could be efficiently used as a translation start codon. The scanning model of translation initiation predicts that the  $d_1$  AUG will be used as the start codon only by pre-initiation complexes that leak past the aAUG, which is unlikely if the aAUG has a strong sequence context.

To identify probable leaky translation initiation events, we thus compared the aAUG and  $d_1$ AUG scores within each of the 50% most abundant mRNAs (Figure 6A). For above-median aAUG score genes, the score of the  $d_1$ AUGs

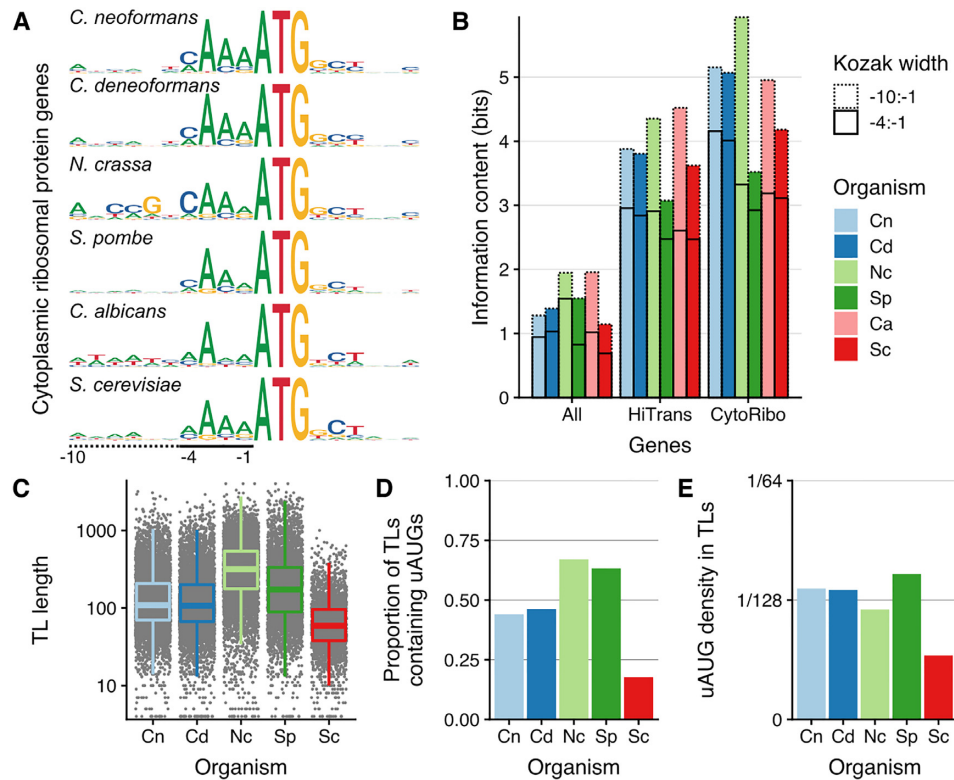
can be very high or very low. By contrast, for the genes with a low aAUG score, there is a bias toward higher  $d_1$ AUG score, suggesting that for these genes the strong  $d_1$ AUG could be used as alternative translation start site (Figure 6A).

To test whether AUG score affects translation initiation, we calculated the ratio of ribosome protected fragment density and RNA-Seq density around each aAUG and  $d_1$ AUG on the same mRNA, and the difference in score between these two AUGs (Figure 6B), using the same 10-codon neighborhood as for our earlier uAUG-aAUG comparison. We found a weak positive correlation between the difference in scores of the two AUGs and RNA-Seq density at these specific loci, raising the possibility that transcription start sites sometimes occur downstream of a weak aAUG. The relative ribosome density is equal on average when the  $d_1$ AUG score is less than the aAUG score. However, a non-linear generalized additive model shows that the relative density sharply increases at  $d_1$ AUGs when their score increases above that of the aAUG. This suggests that for these genes, both AUGs can be used as translation start codon, because a subset of scanning ribosomes leak past a lower-score aAUG and then initiate at the higher-score  $d_1$ AUG.

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

We next determined which groups of genes could be affected by potential alternative start codon usage. We focused our analysis on the 50% most abundant RNAs for which the difference in score between the aAUG and  $d_1$ AUG was the highest (difference in wide score  $d_1$ AUG – aAUG > 0.1,  $n = 167$  for *C. neoformans*) (Supplementary Table S4). Strikingly, for 66% of these genes (110/167) the  $d_1$ AUG is in frame with the corresponding aAUG, with a median of 69nt (mean 79nt) between the two AUGs. Thus, alternative usage of in-frame AUGs would result in proteins with different N-terminal ends. Supporting this hypothesis, 37% of these proteins (41/110) possess a predicted mitochondrial targeting sequence located between the two AUGs, far exceeding the 8% genome-wide (560/6788). This suggests that the usage of the annotated start codon would target the isoform to mitochondria, whereas the usage of the  $d_1$ AUG would produce a protein specific to the cytoplasm or another organelle. Examples of alternative localization driven by alternative N-termini have been observed across eukaryotes (84,85).

The pattern of predicted dual-localization, i.e. enrichment of high-score  $d_1$ AUGs in-frame with predicted mitochondrial localization signal on the longer N-terminal, is conserved in some fungi but not others (Figure 6C). In a null model where coding sequences have random nucleotide content, we would expect roughly one third of  $d_1$ AUGs to be in frame. In six fungal species we examined, for  $d_1$ AUGs whose score is comparable to or less than the aAUG they follow, the proportion in frame is close to (*Cryptococcus*, *N. crassa*) or less than 1/3. These proportions are similar when we considered high abundance (top 50%) or low abundance (bottom 50%) mRNAs. The pattern differs for mRNAs with a  $d_1$ AUG whose score is high relative to the aAUG they follow ( $d_1$ AUG score > aAUG score + 0.1). In *Cryptococcus*



**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 five 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  $1/64$ ).

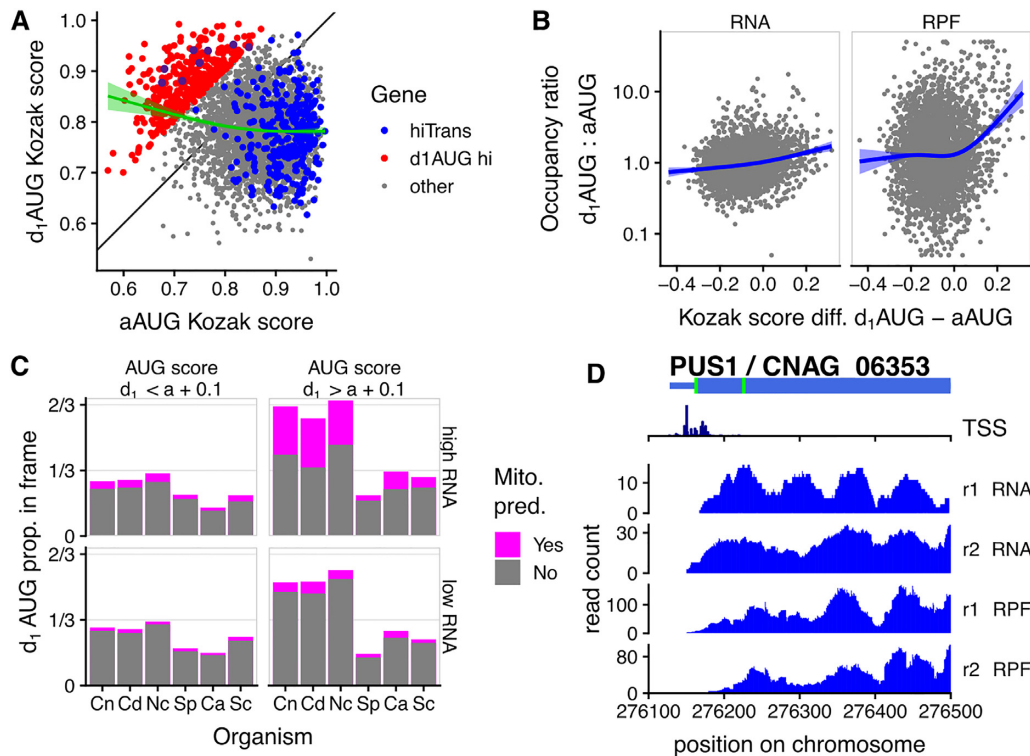
and *N. crassa*, most high abundance mRNAs are in-frame and over one third of these in-frame high-score  $d_1$ AUGs have predicted mitochondrial localization. In *S. cerevisiae* and *C. albicans*, we observe only a slight relative enrichment for high-scoring  $d_1$ AUGs to be in-frame and to follow a mitochondrial targeting sequence. By contrast, in *S. pombe* we see depletion in the in-frame/out-of-frame ratio, even in these proteins with high-scoring  $d_1$ AUGs.

These results suggest that the extent to which alternate translation start codons regulate proteome diversity is variable in fungi. Accordingly, we identified a number of *Cryptococcus* proteins with potential alternative start codons and N-terminal targeting sequences, whose two homologs in *S. cerevisiae* are known to be necessary in two compartments of the cells. For instance, *CnPUS1/CNAG\_06353* is an homolog of both the mitochondrial and cytoplasmic tRNA:pseudouridine synthases encoded by the *PUS1* and *PUS2* paralogs in *S. cerevisiae*. In *C. neoformans*, ribosome occupancy at both the aAUG and  $d_1$ AUG of CNAG\_06353, and the presence of transcription start sites both sides of the aAUG (Figure 6D, Supplementary Figure S8A), argues that both AUGs are used as start codons, and transcription and translation regulation could co-operate to set isoform levels. Similarly, *CnGLO1/CNAG\_04219* encodes both the cytoplasmic and nuclear isoforms of the glyoxalase I depending

on the alternate AUG usage (Supplementary Figure S8B). The next enzyme in this pathway, Glyoxalase II, is likewise encoded by *CnGLO2/CNAG\_01128*, which is a homolog of both cytoplasmic (Glo2) and mitochondrial (Glo4) enzymes in *S. cerevisiae*. CNAG\_01128 has a weak aAUG, strong  $d_1$ AUG, and N-terminal predicted mitochondrial targeting sequence (Supplementary Figure S8C). Finally, we observed that nine members of the amino-acyl tRNA synthetase gene family have predicted alternate localization from alternate AUG start codons (Figure 7A/B).

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

The tRNA charging activity of aaRSs is essential in both cytosol and mitochondria to support translation in each compartment, and examples of alternative localization of two aaRS isoforms of a single gene have been observed in fungi, plants, and animals (86–88). This implies that a eukaryote with a single genomic homolog of an aaRS activity is likely to make distinct localized isoforms from that locus. Thus, we examined predicted aaRS localization in fungi. We assembled gene lists of aaRSs in diverse fungi from homology databases OrthoDB (66) and PANTHERdb (67), adding a mitochondrial SerRS (CNAG\_06763/CNB00380) to the



**Figure 6.** High-scoring downstream AUGs specify alternative N-terminal isoforms in *C. neoformans*. (A) Most genes with high 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, d<sub>1</sub>AUG hi: d<sub>1</sub>AUG score > annotated AUG score + 0.1), and there is a trend for genes with low aAUG score to have a higher d<sub>1</sub>AUG score (green, generalized additive model fit). (B) Higher d<sub>1</sub>AUG score than aAUG score drives higher ribosome protected fragment (RPF) occupancy at the d<sub>1</sub>AUG 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 (d<sub>1</sub>AUG score > annotated AUG score + 0.1) and high 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 *Cn*Pus1 is a candidate alternate-localized protein with a low-score aAUG and high-score d<sub>1</sub>AUG, 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 Supplementary Figure S8.

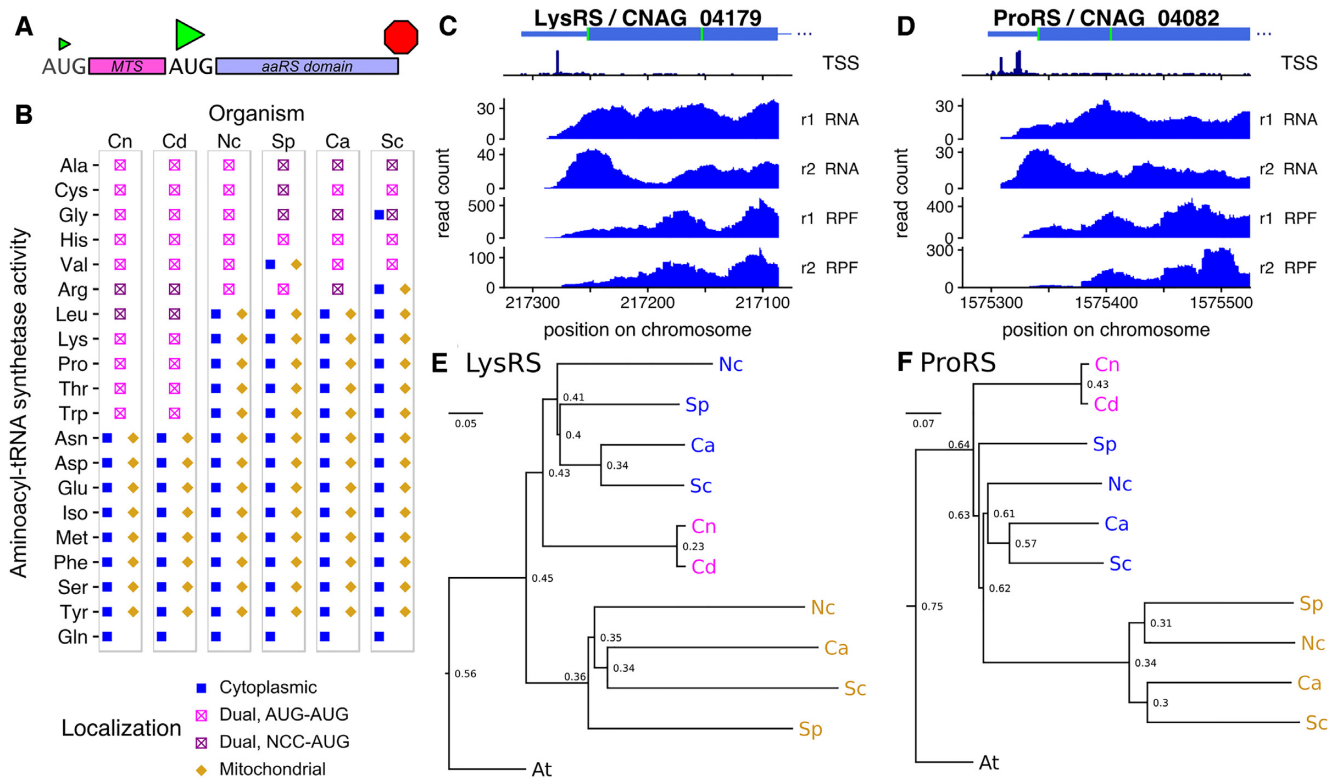
list of *Cryptococcus* aaRSs analysed by Datt and Sharma (89).

In *C. neoformans* and *C. deneoformans*, 11 aaRSs are each expressed from a single genomic locus, including the homologs of all five *S. cerevisiae* aaRSs whose dual-localization has been verified (Supplementary Table S5). Nine of these *Cryptococcus* aaRSs have the same structure of a poor-context annotated AUG followed by a predicted mitochondrial targeting sequence and a strong-context d<sub>1</sub>AUG (Figure 7A, B; AlaRS, CysRS, GlyRS, HisRS, ValRS, LysRS, ProRS, ThrRS, TrpRS). The similar annotated AUG contexts, sharing an unfavourable -3U, suggests that the same mechanism could lead to leaky translation initiation at most of these (Supplementary Table S6). At the downstream AUGs, the strong Kozak context is consistent with efficient translation initiation of the cytoplasmic isoform from this start codon (Supplementary Table S6).

The two remaining single-copy aaRSs have near-AUG translation initiation sites upstream of predicted mitochondrial targeting sequences. Translation of ArgRS starts at an AUU codon with otherwise strong context (ccaccAAU) conserved in both *Cryptococcus* species. This N-terminal extension includes a predicted mitochondrial targeting se-

quence (mitofates  $P > 0.95$  for both species). Translation of LeuRS starts at adjacent ACG and AUU codons which collectively provide strong initiation context (gccaccACGAUU in *C. neoformans*, gccACGAUU in *C. deneoformans*). This N-terminal extension also includes a predicted mitochondrial targeting sequence (mitofates  $P \approx 0.7$  for both species).

In *Cryptococcus*, alternative aaRS isoforms appear to be mostly generated by alternative translation from a single transcript, and sometimes by alternative transcription start sites. On all the predicted dual-localized aaRSs, we observe ribosomal occupancy starting at the earliest start codon (Figure 7C, D and Supplementary Figure S9). LysRS/CNAG\_04179 contains only a single cluster of transcription start sites, upstream of the aAUG (Figure 7C). ProRS/CNAG\_04082 contains a wider bimodal cluster of TSSs, both upstream of the aAUG. Similarly, most transcription initiation is well upstream of the aAUG in CysRS/CNAG\_06713, LeuRS/CNAG\_06123, ThrRS/CNAG\_06755 and ValRS/CNAG\_07473. However, for GlyRS/CNAG\_05900, and HisRS/CNAG\_01544, we observe alternative transcription start sites closely upstream of the annotated start codon, that are likely to affect the efficiency of start codon usage. In ArgRS/CNAG\_03457 there is also an alternative transcription start site, close to



**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 Supplementary Table S5, for details of identifiers for genes (GeneID).

the near-AUG start codon for the mitochondrial form. In AlaRS/CNAG\_05722 and TrpRS/CNAG\_04604 we detect some transcription start sites between the alternative start codons, and TrpRS also has an uORF in the transcript leader that is likely to affect translation. These observations suggest that dual-localization of the single-copy aaRSs in *Cryptococcus* is regulated largely by start codon choice. For some genes, this regulation is backed up by alternative TSS usage.

Some dual-localized genes use an upstream near cognate codon (DualNCC) in all these fungi, but the NCC-initiated aaRS are not the same from one fungus to the other. For instance, both *Cryptococcus* and *N. crassa* AlaRS use DualAUG whereas in *S. pombe*, *S. cerevisiae* and *C. albicans* a DualNCC is used. On the other hand, *S. pombe* GlyRS is regulated by DualNCC whereas the other ones use a DualAUG regulation. Substitution between weak AUG codons and near-cognate codons seems thus to have taken place multiple times in the fungal kingdom.

#### Amino-acyl tRNA synthetases as an evolutionary case study

To understand patterns of dual-localization, we next examined the evolution of aaRSs. The ancestral eukaryote is thought to have had two complete sets of aaRS, one

mitochondrial and one cytoplasmic, but all mitochondrial aaRSs have been captured by the nuclear genome and many have been lost (90). Thus we examined aaRS phylogenetic trees in more detail. For some amino acids (Asn, Asp, Glu, Iso, Met, Phe, Ser, Tyr), reference fungi have distinct cytoplasmic and mitochondrial aaRSs that cluster in separate trees (91). We also do not consider Gln, because organellar Gln-tRNA charging in some species is achieved by an indirect pathway (92).

Dual-localized AlaRS, CysRS and HisRS in the six fungi we focus on are each monophyletic (91). Even these aaRS can be encoded by two genes in some other fungi: AlaRS is duplicated to one exclusively mitochondrial and another exclusively cytoplasmic gene in the Saccharomycete yeast *Vanderwaltozyma polyspora* (93). For CysRS, *Aspergillus versicolor* (ASPVEDRAFT\_141527 and ASPVEDRAFT\_46520) and *Coprinus cinerea* (CC1G\_03242 and CC1G\_14214) have two copies, one of which has a predicted mitochondrial targeting sequence. For HisRS, *Rhizopus delemar* (RO3G\_01784 and RO3G\_16958) and *Phycomyces blakesleeanus* (PHYBL\_135135 and PHYBL\_138952) likewise contain gene duplications. Similarly, *S. cerevisiae* has two ArgRS genes that arose from the whole-genome duplication: *RRS1/YDR341C* is essential, abundant, and inferred to be cytoplasmic (94) while *MSR1/YHR091C* has

a mitochondrial localization sequence and MSR1 deletions have a petite phenotype (95), although both have been detected in mitochondria suggesting some residual dual-localization of the cytoplasmic enzyme (96). The second *S. cerevisiae* stress-responsive cytoplasmic copy of GlyRS also arose from the whole-genome duplication (97). *S. pombe* cytoplasmic ValRS is monophyletic with dual-localized ValRS in other fungi, and *Schizosaccharomyces* also has a paralogous but diverged mitochondrial ValRS that appears to be descended from an early eukaryotic ValRS of mitochondrial origin (98).

LysRS appears to have been duplicated in an ancestor of ascomycetes: ascomycete mitochondrial homologs cluster together, and ascomycete cytoplasmic homologs cluster together, while the single basidiomycete homolog clusters close to the base of this split from other opisthokonts (91). By contrast, LeuRS, ProRS, and TrpRS are each represented by two distinct proteins in ascomycetes, one cytoplasmic and one mitochondrial and of independent descent, but the mitochondrial homolog has been lost in *Cryptococcus* species. In basidiomycetes *Ustilago* and *Puccinia*, homologs of mitochondrial LeuRS and ProRS are not present, but there is a homolog of mitochondrial TrpRS; all these have a single homolog of the cytoplasmic TrpRS (91). Our independent phylogenetic analysis of LysRS and ProRS agrees with the conclusions from PANTHERdb (Figures 7E, F). These analyses show that aaRSs have undergone multiple incidences of at least two processes during fungal evolution: losses associated with the dual-localization of the remaining gene, and duplications followed by specialization.

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

We also observed striking examples of gene-specific regulation by start codon context in *Cryptococcus*, in translation factors affecting start codon selection, supporting previously proposed models of feedback regulation (99,100).

Translation initiation factor eIF1, which enforces the accurate selection of start codons, is encoded by an mRNA with poor start codon context in diverse eukaryotes, driving an autoregulatory program (99,101). In *C. neoformans*, eIF1 (SUI1/CNAG\_04054) also initiates from a poor-context cuuaguugaAUG start (score 0.75), and riboprofiling reads are spread across the annotated ORF (Figure 8A). Intriguingly, the next AUG is out-of frame and has strong context cucaaaaAUG (score 0.98), with a same-frame stop codon 35 codons later, suggesting that this could represent a downstream short ORF that captures ribosomes that have leaked past the poor-context start. To test this hypothesis, we examined the 5' ends of riboprofiling reads, which report on the translation frame of the ribosomes (43). Riboprofiling reads from the 5' and 3' of the eIF1 annotated ORF are ~77% in frame 0, 10% in +1 and 13% in +2, as are reads on two other highly expressed genes, eEF1 $\alpha$  and *HSP90*. By contrast, in the hypothesized downstream ORF, reads are only 57% in frame 0, 32% in frame +1, and 11% in frame +2, consistent with translation occurring in both frame 0 and +1. The gene structure is conserved in *C. deneoformans* eIF1 (CNB05380), with a weak aAUG (score 0.76), a strong

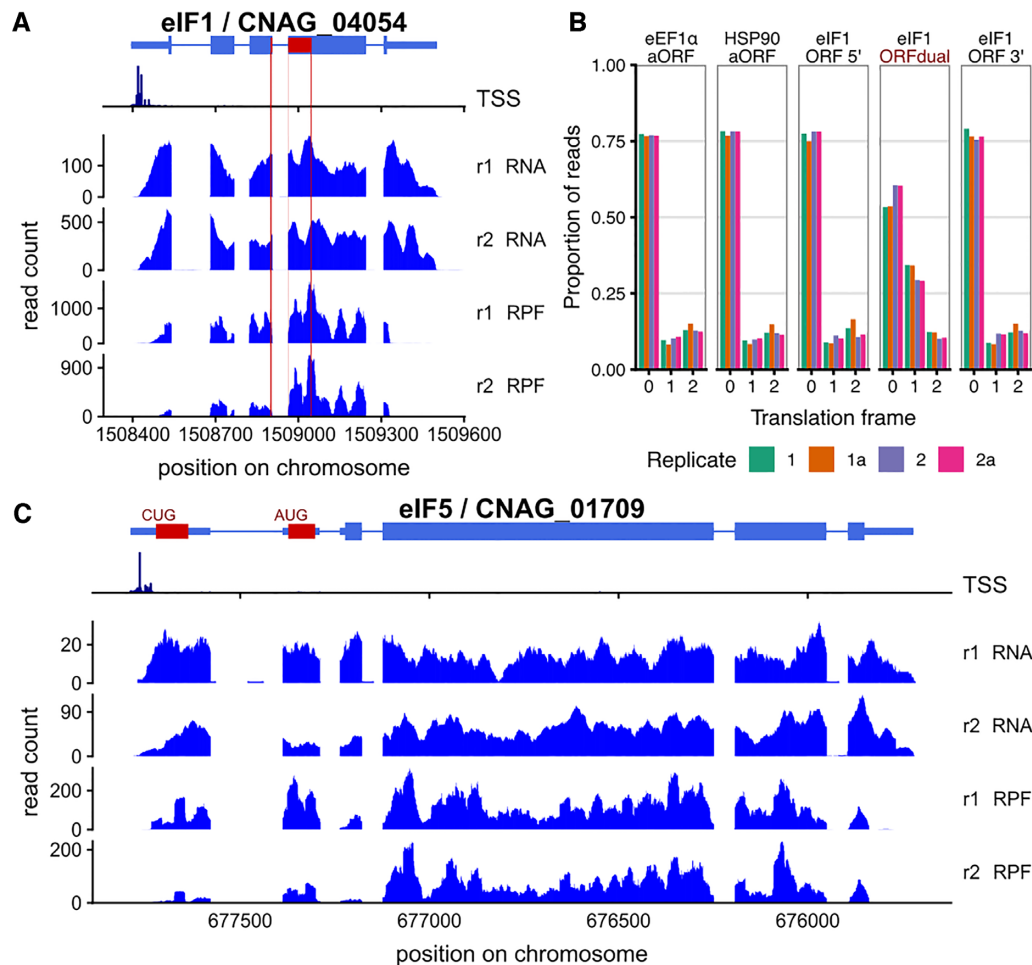
d<sub>1</sub>AUG (score 0.98) in the +1 frame, followed by an enrichment in +1-frame riboprofiling reads (Supplementary Figure S10A,B). We observe small increases in eIF1 mRNA levels in the *upf1*  $\Delta$  strain of *C. deneoformans* at both 30°C (1.16 $\times$ ,  $P = 0.04$ ) and 37°C (1.09 $\times$ ), so NMD could regulate this transcript. Overall, our data support the hypothesis that the downstream ORF of eIF1 is translated after leaky scanning past the annotated AUG, and that the downstream ORF contributes to translation regulation of the annotated ORF.

Translation initiation factor eIF5 reduces the stringency of start codon selection, and is encoded by an mRNA with a repressive uORF initiated from a poor-context uAUG in diverse eukaryotes (100). In *C. neoformans*, eIF5 (*TIF5*/CNAG\_01709) also contains a uAUG with the poor sequence context aaagagucaAUG (score 0.72), while the main ORF of eIF5 is initiated by a strong context ccgcaaaaAUG (score 0.94). We detect ribosomal density on the uORF of *TIF5* comparable to that on the main ORF (Figure 8C), suggesting substantial translation initiation at the uAUG, while there is also clear translation initiation at a further upstream CUG codon. The gene structure is conserved in *C. deneoformans* eIF5 (CNC02150), with the same pattern of riboprofiles at upstream poor-context AUG and near-cognate codons (Supplementary Figure S10C). Further, the *C. deneoformans* homolog transcript abundance increases substantially in the *upf1*  $\Delta$  strain (2.6 $\times$ ,  $P < 10^{-50}$ ). In *N. crassa*, eIF5 has two uORFs and direct analysis of mRNA stability indicated that its transcript is a NMD target (81). The present data support the model that eIF5 translation in *Cryptococcus* is also repressed by upstream reading frames initiated from poor start codons, leading to nonsense-mediated decay of the transcript.

### Variable inserts in eTIFs correlate with variation in translation initiation determinants

The conserved proteins eIF1, eIF5 and eIF1A play pivotal roles in start codon selection, and specific mutations in these factors give rise to suppressor of upstream initiation codon (Sui-) phenotypes and their suppressors (Ssu-) (101). To ask if between-species variability in start codon preference is linked to these initiation factors, we generated multiple sequence alignments of their homologs in fungi.

Translation initiation factor eIF1 shows striking sequence variation across fungi, notably at multiple *Cryptococcus*-specific sequence insertions that result in a 159-aa protein substantially larger than the 108-aa *S. cerevisiae* homolog (Figure 9A). Variation in eIF1 occurs at and around positions known to modulate start codon selection in *S. cerevisiae* (101). For instance, a T15A substitution increases fidelity in *SceIF1* (101), and an analogous T15A substitution is present in eIF1s from *Neurospora* and other filamentous fungi, while both *Cryptococcus* homologs have the T15V substitution. The three fungi that tend not to use alternative AUG start codons in the regulation of proteome diversity, *S. cerevisiae*, *C. albicans* and *S. pombe*, all have a threonine residue at position 15. Variation in fungal eIF1 extends far beyond this N-terminal region: similar patterns of sequence diversity occur at the positions E48, L51, D61 that have been shown to increase fidelity in



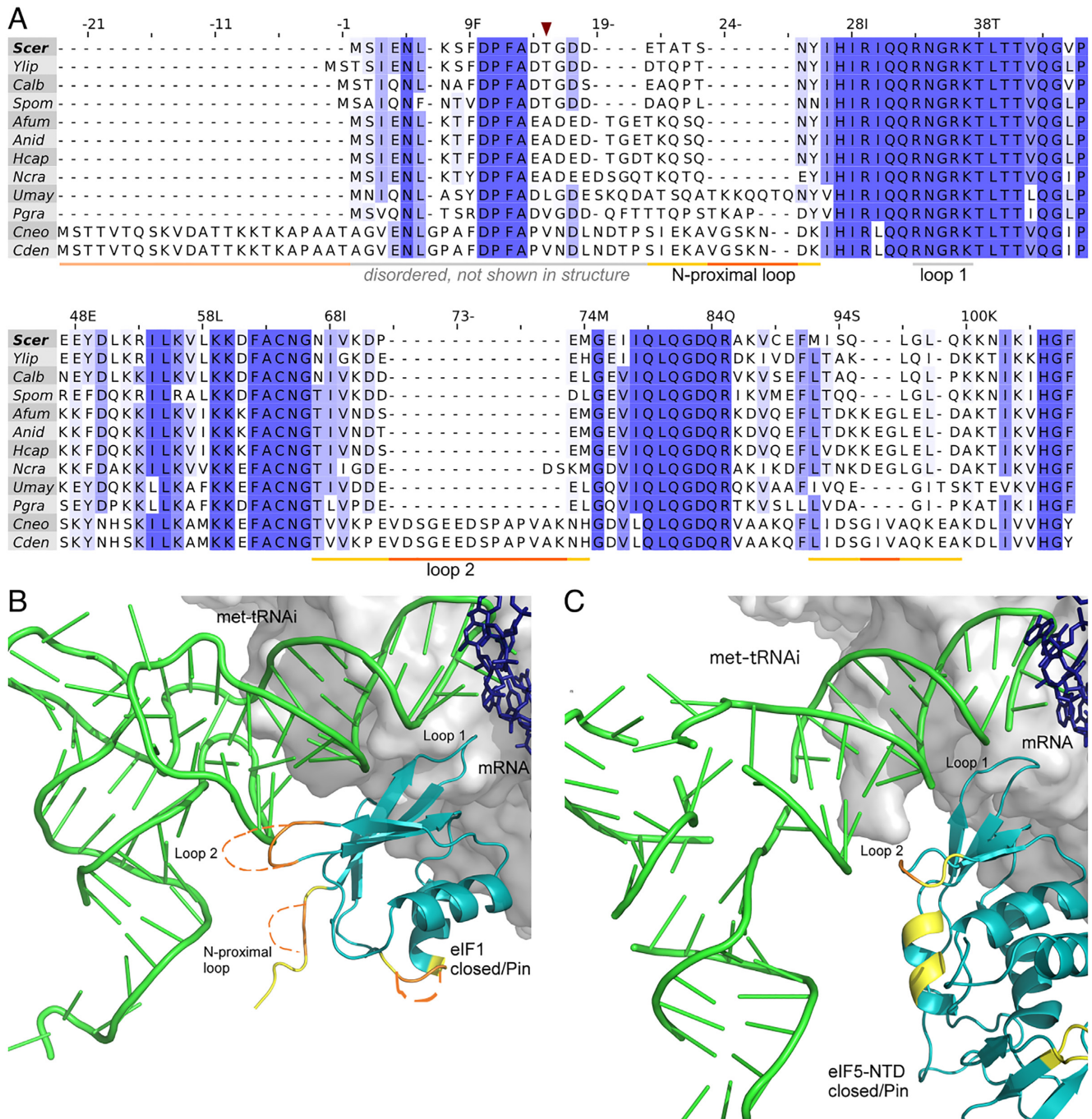
**Figure 8.** Translation initiation factors eIF1 and eIF5 are regulated by alternate start codon usage in *C. neoformans*. (A) Reads on *CneIF1*/CNAG\_04054, showing frame +1 'downstream ORF' in dark red, breaking for an intron. (B) The downstream ORF of *CneIF1* is dual-translated in two frames. Most riboprofiling read 5' ends are in a consistent frame, including in control genes eEF1 $\alpha$ /CNAG\_06125 and HSP90/CNAG\_06125, and in the 5' and 3' ends of the *CneIF1* ORF, but there is 2 $\times$  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 (Supplementary Figure S10).

*SceIF1* (101). By contrast, positions K56, K59, D83, Q84, at which mutations have been shown to reduce fidelity in *SceIF1* (101), are highly conserved in fungi.

We next tested how the translation pre-initiation complex could be affected by the insertions in *Cryptococcus* eIF1 using published structures of the *S. cerevisiae*/*K. lactis* 'Pin' complex engaged in the act of AUG selection (102). We found that the insertions in eIF1 are facing either the methionine initiator tRNA (tRNA<sub>i</sub>) or the solvent-exposed side (Figure 9B). The N-terminal insertion is not visible in the structure, but could be close to the acceptor arm of tRNA<sub>i</sub>. The N-proximal loop insertion of *CneIF1* extends from the *SceIF1* sequence (18-DETATSNY-25) that contacts the acceptor arm of tRNA<sub>i</sub>. The *CneIF1* insertion in loop 2 extends the *SceIF1* loop 2 (70-KDPEMGE-76) that contacts the D-loop of tRNA<sub>i</sub>; substitutions D71A/R and M74A/R increase the charge of *SceIF1* loop 2 and increase initiation at UUG codons and weak AUG codons (103). *CneIF1* loop 2 has substitutions at both these functionally

important sites, and is extended by a further 14 hydrophobic and negative residues. The last insertion in *CneIF1* extends a loop facing the solvent-exposed surface of *SceIF1*. Collectively, this shows that there are likely major differences in the eIF1-tRNA<sub>i</sub> interaction surface in *Cryptococcus* relative to other fungi, an interaction critical for start codon selection (103).

The N-terminal domain of eIF5 (eIF5-NTD) replaces eIF1 upon start codon recognition, and we found between-species variation in *CneIF5* at tRNA<sub>i</sub> interaction surfaces corresponding to variability in *CneIF1* (Figure 9C, Supplementary Figure S11A). *SceIF5* Lys71 and Arg73 in loop 2 make more favourable contacts with the tRNA<sub>i</sub> than the corresponding residues of *SceIF1*, so that the shorter loop 2 of *SceIF5* may allow the tRNA<sub>i</sub> to tilt more towards the 40S subunit (30). Although Arg73 is conserved across fungi, Lys71 is absent in *CneIF5* loop 2 (67-SMAN-70), which is two amino acids shorter than *SceIF5* loop 2 (66-SISVDK-71). Collectively, the longer loop 2 of *CneIF1* and



**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, (102)). 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. (C) Structural predictions of variations in *Cryptococcus* eIF5 mapped on to *S. cerevisiae* PIC (PDB:6FYX, (28)). Multiple sequence alignment of eIF5 is shown in Supplementary Figure S11A.



the shorter loop 2 of *CneIF5* suggest that the conformational changes accompanying start codon recognition may be more exaggerated in *Cryptococcus*, providing a mechanistic hypothesis for stronger genomic patterns of start codon recognition.

Fungal eIF1A homologs also diverge from *SceIF1A* at regions that modulate translation initiation fidelity (Supplementary Figure S11B), for example the N-terminal element DSDGP (101). The *Cryptococcus* eIF1A C-terminus is diverged from all other fungi at *SceIF1A* positions 110–120, and along with other basidiomycetes lacks a loop at *SceIF1A* positions 135–149. This C-terminal region of *SceIF1A* contributes to pre-initiation complex assembly and binds eIF5B (104) and eIF5 (105), and domain deletions or local alanine substitutions reduce fidelity of translation start site selection (101,104,106).

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

## DISCUSSION

Our annotation of transcript structure and translation in two pathogenic *Cryptococcus* species and our analysis of published data from other species show that start codon context has a major effect on protein production, regulation, diversity, and localization in diverse fungi. As such this work represents a useful resource for the field. We find that the use of start codon context to regulate translation initiation varies quantitatively between fungal species. Compared to the model *Saccharomyces*, both *Cryptococcus* and *Neurospora* have long and AUG-rich TLs, and more information-rich and functionally important Kozak sequences. Further, *Cryptococcus* and *Neurospora* display extensive evidence of leaky scanning of weak AUG codons that is used for regulation by upstream ORFs and to generate alternate N-terminal isoforms with different subcellular localization.

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

Translation initiation regulation can be enabled by start codons that are imperfectly used, so that scanning pre-initiation complexes can leak past them. According to the scanning model of translation initiation, a ‘perfect’ strong start codon would prevent this by capturing all the scanning PICs, and leave few for downstream initiation. For example, the downstream out-of-frame ORF of *Cryptococcus* eIF1 is likely to be translated only by PICs that leak past the annotated AUG. The alternative second in-frame AUG of dual-localized proteins is also initiated only by PICs that have leaked past the initial AUG. Our data show this leakiness-driven dual-localization is common in *Cryptococcus*, in addition to being conserved across eukaryotes in gene classes such as tRNA synthetases. Our data also argue that AUGs that are proximal to the 5' cap, or that have poor sequence context, are commonly leaked past in *Cryptococcus*, as shown previously in studies of yeast (107) and

mammals (12,108). We note that leakiness-driven translation regulation is not the only mechanism regulating alternative translation from a single mRNA and is distinct from those that depend on either blocking scanning, or on recycling of post-termination ribosomes such as in the case of *S. cerevisiae* *GCN4* (72).

### Functional role of start codon context varies across the fungal kingdom

*Cryptococcus* and *Neurospora* have long TLs that are AUG-rich, and extended start codon context sequences that suggest a higher ability to discriminate against poor-context AUGs. Several lines of evidence argue that the efficiency with which upstream AUGs capture initiation complexes is determined by the AUG sequence context, notably *in vitro* translation studies in *N. crassa* and *S. cerevisiae* from the Hinnebusch and Sachs labs (76). The most spectacular examples of uORF-associated translation repression in *Cryptococcus* are associated with good-context uAUGs with high ribosome occupancy. However, such strong-context high-occupancy uAUGs are rare. In *Cryptococcus* and *Neurospora*, the leakiness of potential AUG translation start sites is also extensively used to diversify the proteome by alternative N-terminal formation.

In comparison, *S. cerevisiae*, *S. pombe* and *C. albicans* appear to be less efficient in discriminating AUGs based on their sequence context. *S. cerevisiae* has minimized the possibility of regulation of translation by uORFs: it has unusually short TLs, these TLs are unusually AUG-poor, uAUGs tend to have poor context, and there is no statistical association between uAUG score and translation efficiency of the main ORF. Reporter gene studies (15,16) and classic examples such as *GCN4* show that uAUGs can repress translation in *S. cerevisiae*, but genome-wide analysis shows that this is rare during exponential growth in rich media (Supplementary Figure S5.1). Recent work on meiosis (109) and stress (110) shows that 5'-extended transcript leaders that contain repressive uAUGs (‘long undecoded transcript isoforms’) are more common during alternative growth conditions for this yeast. Moreover, in *S. cerevisiae*, near-cognate codons appear to be more common starts for alternative N-terminal formation (111). This suggests that leaky scanning from near-cognate codons, more than from AUGs, might be an important mode of regulation in *S. cerevisiae*. The situation is different in *S. pombe*, which has long AUG-rich TLs but is depleted for downstream in-frame AUGs. Consequently, uAUGs globally repress aORF translation, but do not appear to regulate alternative protein production through alternative AUG start codons. We speculate that the comparatively uninformative Kozak context in *S. pombe* might be variable enough to regulate translation initiation rate but not proteome diversity.

We found that multiple near-cognate start codons are used for leaky initiation in *Cryptococcus*: ACG for the mitochondrial isoform of LeuRS, AUU for the mitochondrial isoform of ArgRS, and the upstream CUG in eIF5. Further work will be needed to quantify the extent of near-cognate start codon usage in *Cryptococcus* in different growth conditions and to compare it to other organisms (22,112).

### Leaky scanning through weak AUGs could regulate the mitochondrial proteome

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

Mechanisms to control initiation efficiency of a mitochondrial-localized regulon could include intracellular magnesium concentration (113), variations in availability or modification status of shared initiation factors, variations of the ratio of mitochondrial volume to intracellular volume (114), or specialized factors to promote initiation specifically of mitochondrial isoforms with their specialized start codon context. Nakagawa *et al.* (115) previously suggested that distinct Kozak contexts might be recognized by different molecular mechanisms.

One candidate mechanism involves the translation initiation factor 3 complex, which has a role in regulating the translation initiation of mitochondrial-localized proteins across eukaryotes. In *S. pombe*, subunits eIF3d/e promote the synthesis of mitochondrial electron transfer chain proteins through a TL-mediated mechanism (116). In *S. cerevisiae* and *Dictyostelium discoideum*, the conserved eIF3-associated Clu1/CluA protein affects mitochondrial morphology (117), and the mammalian homolog CLUH binds and regulates mRNAs of nuclear-encoded mitochondrial proteins (118,119). Metazoans have 12 stably-associated subunits of eIF3, which are conserved in most fungi including *N. crassa* (120), *Cryptococcus* and the Saccharomycetale yeast *Yarrowia lipolytica* (Supplementary Table S7). Interestingly, species that tend not to use alternate AUG codons for dual-localization have lost eIF3 subunits: eIF3d/e/k/l/m are lost in *C. albicans*, and additionally eIF3f/h in the related *S. cerevisiae*; *S. pombe* has independently lost eIF3k/l (Supplementary Table S7; (91)). Further work will be needed to investigate the role of eIF3 in regulating mitochondrial- and dual-localized proteins in the fungal kingdom.

### How could evolutionary plasticity of translational initiation in the fungal kingdom have arisen?

Selection on genome compaction in unicellular yeasts, which has independently led to gene loss and high gene density in multiple lineages of yeast, could lead to shorter TLs. However, *Saccharomyces*, *Schizosaccharomyces* and *Cryptococcus* have all independently evolved yeast lifestyles with compact genomes, yet their average TL lengths differ three-fold. Mutations in gene expression machinery, such as the variation in eIF1 noted above, would alter selective pressure on start codon context, and thus uAUG density. Cells have multiple redundant quality control mechanisms, and flexible protein production through leaky scanning could be buffered by such mechanisms enabling their evolution. Key control mechanisms acting on mRNA, such as RNAi

and polyuridylation, have been lost in fungal lineages such as *Saccharomyces*, which might explain their more 'hard-wired' mechanism of translation initiation.

Unexpectedly, highly conserved core translation initiation factors, such as eIF1, have distinctive sequence inserts in *Cryptococcus* that are not shared even by basidiomycetes such as *Puccinia* and *Ustilago*. One possibility is genetic conflict, as genetic parasites hijack the gene expression machinery (121). Thus, the unique aspects of the *Cryptococcus* translation initiation machinery could have arisen from a past genetic conflict in which rapid evolution of initiation factors in an ancestor enabled evasion of a genomic parasite (e.g. a mycovirus) that would otherwise hijack initiation.

### DATA AVAILABILITY

Raw and summarized sequencing data are available on GEO under accession numbers GSE133695 (RNA-seq, TSS-seq, PAS-seq, DNA-seq) and GSE133125 (ribosome profiling and matched RNA-seq). Custom analysis code in R, and intermediate data files are available at <https://github.com/ewallace/CryptoTranscriptome2018>, doi:10.5281/zenodo.3627874.

### SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

### ACKNOWLEDGEMENTS

We thank members of the Wallace, Janbon, and Madhani labs for helpful discussions and comments on the manuscript. We thank Juan Mata for sharing intermediate data related to (59). We are grateful to J. Weissman (UCSF) for advice on ribosome profiling. We thank three anonymous reviewers for their helpful comments.

### FUNDING

This work in the Madhani lab was supported by grants from the US National Institutes of Health [R01AI120464, R01GM71801 to H.D.M.]; H.D.M. is an Investigator of the Chan-Zuckerberg Biohub; E.W.J.W. is a Sir Henry Dale Fellow, supported by a Sir Henry Dale Fellowship jointly funded by the Wellcome Trust and the Royal Society [208779/Z/17/Z]; L.T. is supported by a Wellcome-University of Edinburgh ISSF3 award; This work in the Janbon lab was supported by an Infect-ERA grant (project Cryptoview). Funding for open access charge: University of Edinburgh.

*Conflict of interest statement.* None declared.

### REFERENCES

1. Hawksworth, D.L. and Lücking, R. (2017) Fungal diversity revisited: 2.2 to 3.8 million species. *Microbiol. Spectrum*, **5**, doi:10.1128/microbiolspec.FUNK-0052-2016.
2. Grigoriev, I.V., Nikitin, R., Haridas, S., Kuo, A., Ohm, R., Otilar, R., Riley, R., Salamov, A., Zhao, X., Korzeniewski, F. *et al.* (2014) MycoCosm portal: gearing up for 1000 fungal genomes. *Nucleic Acids Res.*, **42**, D699–D704.

3. Fan, G., Sun, Q., Li, W., Shi, W., Li, X., Wu, L., Ma, J., Kim, C. Y., Lee, J.-S., Zhou, Y. *et al.* (2018) The global catalogue of microorganisms 10K type strain sequencing project: closing the genomic gaps for the validly published prokaryotic and fungi species. *GigaScience*, **7**, doi:10.1093/gigascience/giy026.
4. Shen, X.-X., Oplulent, D.A., Kominek, J., Zhou, X., Steenwyk, J.L., Buh, K.V., Haase, M.A.B., Wisecaver, J.H., Wang, M., Doering, D.T. *et al.* (2018) Tempo and mode of genome evolution in the budding yeast subphylum. *Cell*, **175**, 1533–1545.
5. Butler, G., Rasmussen, M.D., Lin, M.F., Santos, M.A.S., Sakthikumar, S., Munro, C.A., Rheinbay, E., Grabherr, M., Forche, A., Reedy, J.L. *et al.* (2009) Evolution of pathogenicity and sexual reproduction in eight *Candida* genomes. *Nature*, **459**, 657.
6. Dujon, B., Sherman, D., Fisher, G., Durrens, P., Casaregola, S., Lafontaine, I., De Montigny, J., Marck, C., Neugeglise, C., Talla, E. *et al.* (2004) Genome evolution in yeasts. *Nature*, **430**, 35–44.
7. Stajich, J.E. (2017) Fungal genomes and insights into the evolution of the kingdom. *Microbiol. Spectrum*, **5**, doi:10.1128/microbiolspec.FUNK-0055-2016.
8. Stajich, J.E., Dietrich, F.S. and Roy, S.W. (2007) Comparative genomic analysis of fungal genomes reveals intron-rich ancestors. *Genome Biol.*, **8**, R223.
9. Coletta, A., Pinney, J.W., Solis, D.Y.W., Marsh, J., Pettifer, S.R. and Attwood, T.K. (2010) Low-complexity regions within protein sequences have position-dependent roles. *BMC Syst. Biol.*, **4**, 43–43.
10. Goffeau, A., Barrell, B.G., Bussey, H., Davis, R.W., Dujon, B., Feldmann, H., Galibert, F., Hoheisel, J.D., Jacq, C., Johnston, M. *et al.* (1996) Life with 6000 Genes. *Science*, **274**, 546.
11. Haas, B.J., Zeng, Q., Pearson, M.D., Cuomo, C.A. and Wortman, J.R. (2011) Approaches to fungal genome annotation. *Mycology*, **2**, 118–141.
12. Kozak, M. (1986) Point mutations define a sequence flanking the AUG initiator codon that modulates translation by eukaryotic ribosomes. *Cell*, **44**, 283–292.
13. Kozak, M. (1987) An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs. *Nucleic Acids Res.*, **15**, 8125–8148.
14. Dever, T.E., Kinzy, T.G. and Pavitt, G.D. (2016) Mechanism and regulation of protein synthesis in *Saccharomyces cerevisiae*. *Genetics*, **203**, 65–107.
15. Dvir, S., Velten, L., Sharon, E., Zeevi, D., Carey, L.B., Weinberger, A. and Segal, E. (2013) Deciphering the rules by which 5'-UTR sequences affect protein expression in yeast. *Proc. Natl Acad. Sci. U.S.A.*, **110**, E2792–E2801.
16. Cuperus, J.T., Groves, B., Kuchina, A., Rosenberg, A.B., Jojic, N., Fields, S. and Seelig, G. (2017) Deep learning of the regulatory grammar of yeast 5' untranslated regions from 500,000 random sequences. *Genome Res.*, **27**, 2015–2024.
17. Park, H. and Subramaniam, A.R. (2019) Inverted translational control of eukaryotic gene expression by ribosome collisions. *PLoS Biol.*, **17**, e3000396.
18. Li, J.J., Chew, G.-L. and Biggin, M.D. (2019) Quantitative principles of cis-translational control by general mRNA sequence features in eukaryotes. *Genome Biol.*, **20**, 162.
19. Fervers, P., Fervers, F., Makalowski, W. and Jankalski, M. (2018) Life cycle adapted upstream open reading frames (uORFs) in *Trypanosoma congolense*: A post-transcriptional approach to accurate gene regulation. *PLoS One*, **13**, e0201461.
20. Duncan, C.D.S., Rodríguez-López, M., Ruis, P., Bähler, J. and Mata, J. (2018) General amino acid control in fission yeast is regulated by a nonconserved transcription factor, with functions analogous to Gcn4/Atf4. *Proc. Natl Acad. Sci. U.S.A.*, **115**, E1829–E1838.
21. Sundaram, A. and Grant, C.M. (2014) A single inhibitory upstream open reading frame (uORF) is sufficient to regulate *Candida albicans* GCN4 translation in response to amino acid starvation conditions. *RNA*, **20**, 559–567.
22. Ivanov, I.P., Wei, J., Caster, S.Z., Smith, K.M., Michel, A.M., Zhang, Y., Firth, A.E., Freitag, M., Dunlap, J.C., Bell-Pedersen, D. *et al.* (2017) Translation initiation from conserved Non-AUG codons provides additional layers of regulation and coding capacity. *mBio*, **8**, e00844–00817.
23. von Arnim, A.G., Jia, Q. and Vaughn, J.N. (2014) Regulation of plant translation by upstream open reading frames. *Plant Sci.*, **214**, 1–12.
24. Barbosa, C., Peixeiro, I. and Romão, L. (2013) Gene expression regulation by upstream open reading frames and human disease. *PLoS Genet.*, **9**, e1003529–e1003529.
25. Chen, S.-J., Lin, G., Chang, K.-J., Yeh, L.-S. and Wang, C.-C. (2008) Translational efficiency of a Non-AUG initiation codon is significantly affected by its sequence context in yeast. *J. Biol. Chem.*, **283**, 3173–3180.
26. Hinnebusch, A.G., Ivanov, I.P. and Sonenberg, N. (2016) Translational control by 5'-untranslated regions of eukaryotic mRNAs. *Science*, **352**, 1413–1416.
27. Wethmar, K. (2014) The regulatory potential of upstream open reading frames in eukaryotic gene expression. *Wiley Interdiscip. Rev.: RNA*, **5**, 765–768.
28. Llácer, J.L., Hussain, T., Marler, L., Aitken, C.E., Thakur, A., Lorsch, J.R., Hinnebusch, A.G. and Ramakrishnan, V. (2015) Conformational differences between open and closed states of the eukaryotic translation initiation complex. *Mol. Cell*, **59**, 399–412.
29. Hinnebusch, A.G. (2017) Structural insights into the mechanism of scanning and start codon recognition in eukaryotic translation initiation. *Trends Biochem. Sci.*, **42**, 589–611.
30. Llácer, J.L., Hussain, T., Saini, A.K., Nanda, J.S., Kaur, S., Gordiyenko, Y., Kumar, R., Hinnebusch, A.G., Lorsch, J.R. and Ramakrishnan, V. (2018) Translational initiation factor eIF5 replaces eIF1 on the 40S ribosomal subunit to promote start-codon recognition. *eLife*, **7**, e39273.
31. Janbon, G. (2018) Introns in *Cryptococcus*. *Mem. Inst. Oswaldo Cruz*, **113**, e170519.
32. Goebels, C., Thonn, A., Gonzalez-Hilarion, S., Rolland, O., Moyrand, F., Beilharz, T.H. and Janbon, G. (2013) Introns regulate gene expression in *Cryptococcus neoformans* in a Pab2p dependent pathway. *PLoS Genet.*, **9**, e1003686.
33. Dumesic, P.A., Natarajan, P., Chen, C., Drinnenberg, I.A., Schiller, B.J., Thompson, J.D., Moresco, J.J., Yates III, J.R., Bartel, D.P. and Madhani, H.D. (2013) Stalled spliceosomes are a signal for RNAi-mediated genome defense. *Cell*, **152**, 957–968.
34. Bonnet, A., Grosso, A.R., Elkaoutari, A., Coleno, E., Presle, A., Sridhara, S.C., Janbon, G., Géli, V., de Almeida, S.F. and Palancade, B. (2017) Introns protect eukaryotic genomes from transcription-Associated genetic instability. *Mol. Cell*, **67**, 608–621.
35. Janbon, G., Ormerod, K.L., Paulet, D., Byrnes, E.J. III, Chatterjee, G., Yadav, V., Mullanpudi, N., Hon, C.C., Billmyre, R.B., Brunel, F. *et al.* (2014) Analysis of the genome and transcriptome of *Cryptococcus neoformans* var. *grubii* reveals complex RNA expression and microevolution leading to virulence attenuation. *PLoS Genet.*, **10**, e1004261.
36. Gonzalez-Hilarion, S., Paulet, D., Lee, K.-T., Hon, C.-C., Lechat, P., Mogensen, E., Moyrand, F., Proux, C., Barboux, R., Bussotti, G. *et al.* (2016) Intron retention-dependent gene regulation in *Cryptococcus neoformans*. *Sci. Rep.*, **6**, 32252.
37. Winston, F., Dollard, C. and Ricupero-Hovasse, S.L. (1995) Construction of a set of convenient *saccharomyces cerevisiae* strains that are isogenic to S288C. *Yeast*, **11**, 53–55.
38. Lee, N. and Janbon, G. (2006) In: Kavanagh, K (ed). *Med Mycol*, pp. 275–304.
39. Moyrand, F., Lafontaine, I., Fontaine, T. and Janbon, G. (2008) *UGE1* and *UGE2* regulate the UDP-glucose/UDP-galactose equilibrium in *Cryptococcus neoformans*. *Eukaryot. Cell*, **7**, 2069–2077.
40. Malabat, C., Feuerbach, F., Ma, L., Saveanu, C. and Jacquier, A. (2015) Quality control of transcription start site selection by nonsense-mediated-mRNA decay. *Elife*, **4**, e06722.
41. Kim, D., Pertea, G., Trapnell, C., Pimentel, H., Kelley, R. and Salzberg, S.L. (2013) TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol.*, **14**, R36–R36.
42. Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal*, **17**, doi:10.14806/ej.17.1.200.
43. Ingolia, N.T., Ghaemmaghami, S., Newman, J.R.S. and Weissman, J.S. (2009) Genome-wide analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science*, **324**, 218.
44. Dunn, J.G., Foo, C.K., Belletier, N.G., Gavis, E.R. and Weissman, J.S. (2013) Ribosome profiling reveals pervasive and regulated stop codon readthrough in *Drosophila melanogaster*. *eLife*, **2**, e01179.

45. Carja, O., Xing, T., Wallace, E.W.J., Plotkin, J.B. and Shah, P. (2017) riboviz: analysis and visualization of ribosome profiling datasets. *BMC Bioinformatics*, **18**, 461–461.
46. Pertea, M., Kim, D., Pertea, G.M., Leek, J.T. and Salzberg, S.L. (2016) Transcript-level expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat. Protoc.*, **11**, 1650.
47. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R. and Genome Project Data Processing, S. (2009) The sequence alignment/map format and SAMtools. *Bioinformatics*, **25**, 2078–2079.
48. Quinlan, A.R. and Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, **26**, 841–842.
49. R.C.Team (2018) *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna.
50. Wickham, H. (2016) In: *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag, NY.
51. Wickham, H., François, R., Henry, L. and Müller, K. (2018) dplyr: A Grammar of Data Manipulation. R package version 0.7.8.
52. Wilke, C.O. (2018) cowplot: Streamlined Plot Theme and Plot Annotations for 'ggplot2'. R package version 0.9.3.
53. Wagih, O. (2017) gseqlogo: A 'ggplot2' Extension for Drawing Publication-Ready Sequence Logos. R package version 0.1.
54. Love, M.I., Huber, W. and Anders, S. (2014) Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.*, **15**, 550.
55. Edgar, R.C. (2004) MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.*, **32**, 1792–1797.
56. Wilm, A., Higgins, D.G., Valentin, F., Blackshields, G., McWilliam, H., Wallace, I.M., Thompson, J.D., Larkin, M.A., Brown, N.P., McGettigan, P.A. et al. (2007) Clustal W and Clustal X version 2.0. *Bioinformatics*, **23**, 2947–2948.
57. Yu, C.-H., Dang, Y., Zhou, Z., Wu, C., Zhao, F., Sachs, M.S. and Liu, Y. (2015) Codon usage influences the local rate of translation elongation to regulate Co-translational protein folding. *Mol. Cell*, **59**, 744–754.
58. Kersey, P.J., Allen, J.E., Allot, A., Barba, M., Boddu, S., Bolt, B.J., Carvalho-Silva, D., Christensen, M., Davis, P., Grabmueller, C. et al. (2018) Ensembl genomes 2018: an integrated omics infrastructure for non-vertebrate species. *Nucleic Acids Res.*, **46**, D802–D808.
59. Duncan, C.D.S. and Mata, J. (2017) Effects of cycloheximide on the interpretation of ribosome profiling experiments in *Schizosaccharomyces pombe*. *Sci. Rep.*, **7**, 10331.
60. Muzzey, D., Sherlock, G. and Weissman, J.S. (2014) Extensive and coordinated control of allele-specific expression by both transcription and translation in *Candida albicans*. *Genome Res.*, **24**, 963–973.
61. Skrzypek, M.S., Binkley, J., Binkley, G., Miyasato, S.R., Simison, M. and Sherlock, G. (2017) The candida genome database (CGD): incorporation of Assembly 22, systematic identifiers and visualization of high throughput sequencing data. *Nucleic Acids Res.*, **45**, D592–D596.
62. Gerashchenko, M.V. and Gladyshev, V.N. (2014) Translation inhibitors cause abnormalities in ribosome profiling experiments. *Nucleic Acids Res.*, **42**, e134.
63. Csárdi, G., Franks, A., Choi, D.S., Airoidi, E.M. and Drummond, D.A. (2015) Accounting for experimental noise reveals that mRNA Levels, amplified by post-transcriptional processes, largely determine Steady-State protein levels in yeast. *PLoS Genet.*, **11**, e1005206.
64. Weinberg, D.E., Shah, P., Eichhorn, S.W., Hussmann, J.A., Plotkin, J.B. and Bartel, D.P. (2016) Improved ribosome-footprint and mRNA measurements provide insights into dynamics and regulation of yeast translation. *Cell Rep.*, **14**, 1787–1799.
65. Cherry, J.M., Hong, E.L., Amundsen, C., Balakrishnan, R., Binkley, G., Chan, E.T., Christie, K.R., Costanzo, M.C., Dwight, S.S., Engel, S.R. et al. (2012) *Saccharomyces genome database: the genomics resource of budding yeast*. *Nucleic Acids Res.*, **40**, D700–D705.
66. Kriventseva, E.V., Kuznetsov, D., Tegenfeldt, F., Manni, M., Dias, R., Simão, F.A. and Zdobnov, E.M. (2019) OrthoDB v10: sampling the diversity of animal, plant, fungal, protist, bacterial and viral genomes for evolutionary and functional annotations of orthologs. *Nucleic Acids Res.*, **47**, D807–D811.
67. Mi, H., Muruganujan, A. and Thomas, P.D. (2013) PANTHER in 2013: modeling the evolution of gene function, and other gene attributes, in the context of phylogenetic trees. *Nucleic Acids Res.*, **41**, D377–D386.
68. Basenko, Y.E., Pulman, A.J., Shanmugasundram, A., Harb, S.O., Crouch, K., Starns, D., Warrenfeltz, S., Aurrecochea, C., Stoeckert, J.C., Kissinger, C.J. et al. (2018) FungiDB: an integrated bioinformatic resource for fungi and oomycetes. *J. Fungi*, **4**, E39.
69. Ban, N., Beckmann, R., Cate, J.H.D., Dinman, J.D., Dragon, F., Ellis, S.R., Lafontaine, D.L.J., Lindahl, L., Liljas, A., Lipton, J.M. et al. (2014) A new system for naming ribosomal proteins. *Curr. Opin. Struct. Biol.*, **24**, 165–169.
70. Li, H., Hou, J., Bai, L., Hu, C., Tong, P., Kang, Y., Zhao, X. and Shao, Z. (2015) Genome-wide analysis of core promoter structures in *Schizosaccharomyces pombe* with DeepCAGE. *RNA Biol.*, **12**, 525–537.
71. Neafsey, D.E. and Galagan, J.E. (2007) Dual modes of natural selection on upstream open reading frames. *Mol. Biol. Evol.*, **24**, 1744–1751.
72. Hinnebusch, A.G. (2005) Translational regulation of *GCN4* and the general amino acid control of yeast. *Annu. Rev. Microbiol.*, **59**, 407–450.
73. Duncan, C.D.S., Rodríguez-López, M., Ruis, P., Bähler, J. and Mata, J. (2018) General amino acid control in fission yeast is regulated by a nonconserved transcription factor, with functions analogous to *Gcn4/Atf4*. *Proc. Natl. Acad. Sci. U.S.A.*, **115**, E1829.
74. Madi, L., McBride, S.A., Bailey, L.A. and Ebole, D.J. (1997) *rco-3*, a gene involved in glucose transport and conidiation in *Neurospora crassa*. *Genetics*, **146**, 499–508.
75. Wiese, A., Elzinga, N., Wobbes, B. and Smeeckens, S. (2005) Sucrose-induced translational repression of plant bZIP-type transcription factors. *Biochem. Soc. Trans.*, **33**, 272.
76. Gaba, A., Wang, Z., Krishnamoorthy, T., Hinnebusch, A.G. and Sachs, M.S. (2001) Physical evidence for distinct mechanisms of translational control by upstream open reading frames. *EMBO J.*, **20**, 6453–6463.
77. Kervestin, S. and Jacobson, A. (2012) NMD: a multifaceted response to premature translational termination. *Nat. Rev. Mol. Cell Biol.*, **13**, 703–712.
78. Arribere, J.A. and Gilbert, W.V. (2013) Roles for transcript leaders in translation and mRNA decay revealed by transcript leader sequencing. *Genome Res.*, **23**, 977–987.
79. Hood, H.M., Spevak, C.C. and Sachs, M.S. (2007) Evolutionary changes in the fungal carbamoyl-phosphate synthetase small subunit gene and its associated upstream open reading frame. *Fungal Genet. Biol.*, **44**, 93–104.
80. Gaba, A., Jacobson, A. and Sachs, M.S. (2005) Ribosome occupancy of the yeast CPA1 upstream open reading frame termination codon modulates Nonsense-Mediated mRNA Decay. *Mol. Cell*, **20**, 449–460.
81. Zhang, Y. and Sachs, M.S. (2015) Control of mRNA stability in fungi by NMD, EJC and CBC factors through 3'UTR introns. *Genetics*, **200**, 1133.
82. Wei, J., Zhang, Y., Ivanov, I.P. and Sachs, M.S. (2013) The stringency of start codon selection in the filamentous fungus *Neurospora crassa*. *J. Biol. Chem.*, **288**, 9549–9562.
83. Spealman, P., Naik, A.W., May, G.E., Kuersten, S., Freeberg, L., Murphy, R.F. and McManus, J. (2018) Conserved non-AUG uORFs revealed by a novel regression analysis of ribosome profiling data. *Genome Res.*, **28**, 214–222.
84. Danpure, C.J. (1995) How can the products of a single gene be localized to more than one intracellular compartment? *Trends Cell Biol.*, **5**, 230–238.
85. Silva-Filho, M.C. (2003) One ticket for multiple destinations: dual targeting of proteins to distinct subcellular locations. *Curr. Opin. Plant Biol.*, **6**, 589–595.
86. Mireau, H., Lancelin, D. and Small, I.D. (1996) The same Arabidopsis gene encodes both cytosolic and mitochondrial alanyl-tRNA synthetases. *Plant Cell*, **8**, 1027–1039.
87. Mudge, S.J., Williams, J.H., Eyre, H.J., Sutherland, G.R., Cowan, P.J. and Power, D.A. (1998) Complex organisation of the 5'-end of the human glycine tRNA synthetase gene. *Gene*, **209**, 45–50.

88. Natsoulis,G., Hilger,F. and Fink,G.R. (1986) The HTS1 gene encodes both the cytoplasmic and mitochondrial histidine tRNA synthetases of *S. cerevisiae*. *Cell*, **46**, 235–243.
89. Datt,M. and Sharma,A. (2014) Novel and unique domains in aminoacyl-tRNA synthetases from human fungal pathogens *Aspergillus niger*, *Candida albicans* and *Cryptococcus neoformans*. *BMC Genomics*, **15**, 1069.
90. Duchêne,A.-M., Pujol,C. and Maréchal-Drouard,L. (2009) Import of tRNAs and aminoacyl-tRNA synthetases into mitochondria. *Curr. Genet.*, **55**, 1–18.
91. Muruganujan,A., Ebert,D., Mi,H., Thomas,P.D. and Huang,X. (2018) PANTHER version 14: more genomes, a new PANTHER GO-slim and improvements in enrichment analysis tools. *Nucleic Acids Res.*, **47**, D419–D426.
92. Frechin,M., Duchêne,A.-M. and Becker,H.D. (2009) Translating organellar glutamine codons: a case by case scenario? *RNA Biol.*, **6**, 31–34.
93. Chang,C.-P., Tseng,Y.-K., Ko,C.-Y. and Wang,C.-C. (2012) Alanyl-tRNA synthetase genes of *Vanderwaltozyma polyspora* arose from duplication of a dual-functional predecessor of mitochondrial origin. *Nucleic Acids Res.*, **40**, 314–322.
94. Geslain,R., Martin,F., Delagoutte,B., Cavarelli,J., Gangloff,J. and Eriani,G. (2000) In vivo selection of lethal mutations reveals two functional domains in arginyl-tRNA synthetase. *RNA*, **6**, 434–448.
95. Merz,S. and Westermann,B. (2009) Genome-wide deletion mutant analysis reveals genes required for respiratory growth, mitochondrial genome maintenance and mitochondrial protein synthesis in *Saccharomyces cerevisiae*. *Genome Biol.*, **10**, R95.
96. Sickmann,A., Reinders,J., Wagner,Y., Joppich,C., Zahedi,R., Meyer,H.E., Schönfisch,B., Perschil,I., Chacinska,A., Guiard,B. *et al.* (2003) The proteome of *Saccharomyces cerevisiae* mitochondria. *Proc. Natl Acad. Sci. U.S.A.*, **100**, 13207–13212.
97. Chen,S.-J., Wu,Y.-H., Huang,H.-Y. and Wang,C.-C. (2012) *Saccharomyces cerevisiae* possesses a stress-inducible glycyl-tRNA synthetase gene. *PLoS One*, **7**, e33363.
98. Chiu,W.-C., Chang,C.-P., Wen,W.-L., Wang,S.-W. and Wang,C.-C. (2010) *Schizosaccharomyces pombe* possesses two paralogous Valyl-tRNA synthetase genes of mitochondrial origin. *Mol. Biol. Evol.*, **27**, 1415–1424.
99. Ivanov,I.P., Loughran,G., Sachs,M.S. and Atkins,J.F. (2010) Initiation context modulates autoregulation of eukaryotic translation initiation factor 1 (eIF1). *Proc. Natl Acad. Sci. U.S.A.*, **107**, 18056–18060.
100. Loughran,G., Sachs,M.S., Atkins,J.F. and Ivanov,I.P. (2012) Stringency of start codon selection modulates autoregulation of translation initiation factor eIF5. *Nucleic Acids Res.*, **40**, 2898–2906.
101. Martin-Marcos,P., Cheung,Y.-N. and Hinnebusch,A.G. (2011) Functional elements in initiation factors 1, 1A, and 2 $\beta$  discriminate against poor AUG context and non-AUG start codons. *Mol. Cell Biol.*, **31**, 4814–4831.
102. Hussain,T., Llácer,J.L., Fernández,I.S., Munoz,A., Martin-Marcos,P., Savva,C.G., Lorsch,J.R., Hinnebusch,A.G. and Ramakrishnan,V. (2014) Structural changes enable start codon recognition by the eukaryotic translation initiation complex. *Cell*, **159**, 597–607.
103. Thakur,A. and Hinnebusch,A.G. (2018) eIF1 Loop 2 interactions with Met-tRNA(i) control the accuracy of start codon selection by the scanning preinitiation complex. *Proc. Natl Acad. Sci. U.S.A.*, **115**, E4159–E4168.
104. Olsen,D.S., Savner,E.M., Mathew,A., Zhang,F., Krishnamoorthy,T., Phan,L. and Hinnebusch,A.G. (2003) Domains of eIF1A that mediate binding to eIF2, eIF3 and eIF5B and promote ternary complex recruitment in vivo. *EMBO J.*, **22**, 193–204.
105. Luna,R.E., Arthanari,H., Hiraishi,H., Akabayov,B., Tang,L., Cox,C., Markus,M.A., Luna,L.E., Ikeda,Y., Watanabe,R. *et al.* (2013) The interaction between eukaryotic initiation factor 1A and eIF5 retains eIF1 within scanning preinitiation complexes. *Biochemistry*, **52**, 9510–9518.
106. Fekete,C.A., Applefield,D.J., Blakely,S.A., Shirokikh,N., Pestova,T., Lorsch,J.R. and Hinnebusch,A.G. (2005) The eIF1A C-terminal domain promotes initiation complex assembly, scanning and AUG selection in vivo. *EMBO J.*, **24**, 3588–3601.
107. Slusher,L.B., Gillman,E.C., Martin,N.C. and Hopper,A.K. (1991) mRNA leader length and initiation codon context determine alternative AUG selection for the yeast gene MOD5. *Proc. Natl Acad. Sci. U.S.A.*, **88**, 9789.
108. Calvo,S.E., Pagliarini,D.J. and Mootha,V.K. (2009) Upstream open reading frames cause widespread reduction of protein expression and are polymorphic among humans. *Proc. Natl Acad. Sci. U.S.A.*, **106**, 7507–7512.
109. Cheng,Z., Otto,G.M., Powers,E.N., Keskin,A., Mertins,P., Carr,S.A., Jovanovic,M. and Brar,G.A. (2018) Pervasive, coordinated protein-level changes driven by transcript isoform switching during meiosis. *Cell*, **172**, 910–923.
110. Van Dalen,K.M., Hodapp,S., Keskin,A., Otto,G.M., Berdan,C.A., Higdon,A., Cheunkarndee,T., Nomura,D.K., Jovanovic,M. and Brar,G.A. (2018) Global proteome remodeling during ER stress involves Hac1-Driven expression of long undecoded transcript isoforms. *Dev. Cell*, **46**, 219–235.
111. Monteuiis,G., Mišcicka,A., Świrski,M., Zenad,L., Niemitalo,O., Wrobel,L., Alam,J., Chacinska,A., Kastaniotis,A.J., Kufel,J. *et al.* (2019) Non-canonical translation initiation in yeast generates a cryptic pool of mitochondrial proteins. *Nucleic Acids Res.*, **47**, 5777–5791.
112. Brar,G.A. (2016) Beyond the triplet code: Context cues transform translation. *Cell*, **167**, 1681–1692.
113. Feeney,K.A., Hansen,L.L., Putker,M., Olivares-Yañez,C., Day,J., Eades,L.J., Larrondo,L.F., Hoyle,N.P., O'Neill,J.S. and van Ooijen,G. (2016) Daily magnesium fluxes regulate cellular timekeeping and energy balance. *Nature*, **532**, 375–379.
114. Tsuboi,T., Viana,M.P., Xu,F., Yu,J., Chanchani,R., Arceo,X.G., Tutucci,E., Choi,J., Chen,Y.S., Singer,R.H. *et al.* (2019) Mitochondrial volume fraction controls translation of nuclear-encoded mitochondrial proteins. bioRxiv doi: <https://doi.org/10.1101/529289>, 25 January 2019, preprint: not peer reviewed.
115. Nakagawa,S., Niimura,Y., Gojobori,T., Tanaka,H. and Miura,K.-i. (2008) Diversity of preferred nucleotide sequences around the translation initiation codon in eukaryote genomes. *Nucleic Acids Res.*, **36**, 861–871.
116. Shah,M., Su,D., Scheliga,J.S., Pluskal,T., Boronat,S., Motamedchaboki,K., Campos,A.R., Qi,F., Hidalgo,E., Yanagida,M. *et al.* (2016) A transcript-specific eIF3 complex mediates global translational control of energy metabolism. *Cell Rep.*, **16**, 1891–1902.
117. Fields,S.D., Conrad,M.N. and Clarke,M. (1998) The *S. cerevisiae* CLU1 and *D. discoideum* cluA genes are functional homologues that influence mitochondrial morphology and distribution. *J. Cell Sci.*, **111**, 1717.
118. Gao,J., Schatton,D., Martinelli,P., Hansen,H., Pla-Martin,D., Barth,E., Becker,C., Altmueller,J., Frommolt,P., Sardiello,M. *et al.* (2014) CLUH regulates mitochondrial biogenesis by binding mRNAs of nuclear-encoded mitochondrial proteins. *J. Cell Biol.*, **207**, 213–223.
119. Schatton,D., Pla-Martin,D., Marx,M.-C., Hansen,H., Mourier,A., Nemazany,I., Pessia,A., Zentis,P., Corona,T., Kondylis,V. *et al.* (2017) CLUH regulates mitochondrial metabolism by controlling translation and decay of target mRNAs. *J. Cell Biol.*, **216**, 675–693.
120. Smith,M.D., Gu,Y., Querol-Audi,J., Vogan,J.M., Nitido,A. and Cate,J.H.D. (2013) Human-Like eukaryotic translation initiation factor 3 from *Neurospora crassa*. *PLoS One*, **8**, e78715.
121. Madhani,H.D. (2013) The frustrated gene: origins of eukaryotic gene expression. *Cell*, **155**, 744–749.