



Evolutionary dynamics of polyadenylation signals and their recognition strategies in protists

Marcin P Sajek, Danielle Y Bilodeau, Michael A Beer, et al.

Genome Res. published online September 26, 2024
Access the most recent version at doi:[10.1101/gr.279526.124](https://doi.org/10.1101/gr.279526.124)

P<P	Published online September 26, 2024 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see https://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

Evolutionary dynamics of polyadenylation signals and their recognition strategies in protists

Marcin P Sajek^{1,2,3}, Danielle Y Bilodeau^{1,2}, Michael A Beer^{4,5}, Emma Horton^{1,2}, Yukiko Miyamoto⁶, Katrina B Velle⁷, Lars Eckmann⁶, Lillian Fritz-Laylin⁷, Olivia S Rissland^{1,2*}, and Neelanjan Mukherjee^{1,2*}

¹Department of Biochemistry and Molecular Genetics, University of Colorado School of Medicine, Aurora CO 80045, USA

²RNA Bioscience Initiative, University of Colorado School of Medicine, Aurora CO 80045, USA

³Institute of Human Genetics, Polish Academy of Sciences, 60-479 Poznan, Poland

⁴Department of Biomedical Engineering, Johns Hopkins University School of Medicine, Baltimore, MD 21205, USA

⁵McKusick-Nathans Department of Genetic Medicine, Johns Hopkins University School of Medicine, Baltimore, MD 21205, USA

⁶Department of Medicine, University of California San Diego, La Jolla, CA 92093, USA

⁷Department of Biology, University of Massachusetts, Amherst, MA 01003, USA

*Corresponding Author: olivia.rissland@gmail.com,
neelanjan.mukherjee@cuanschultz.edu

1 **ABSTRACT**

2 The poly(A) signal, together with auxiliary elements, directs cleavage of a pre-mRNA
3 and thus determines the 3' end of the mature transcript. In many species, including
4 humans, the poly(A) signal is an AAUAAA hexamer, but we recently found that the
5 deeply branching eukaryote *Giardia lamblia* uses a distinct hexamer (AGURAA) and
6 lacks any known auxiliary elements. Our discovery prompted us to explore the
7 evolutionary dynamics of poly(A) signals and auxiliary elements in the eukaryotic
8 kingdom. We used direct RNA sequencing to determine poly(A) signals for four
9 protists within the Metamonada clade (which also contains *Giardia lamblia*) and two
10 outgroup protists. These experiments revealed that the AAUAAA hexamer serves as
11 the poly(A) signal in at least four different eukaryotic clades, indicating that it is likely
12 the ancestral signal, whereas the unusual *Giardia* version is derived. We found that
13 the use and relative strengths of auxiliary elements are also plastic; in fact, within
14 Metamonada, species like *Giardia lamblia* make use of a previously unrecognized
15 auxiliary element where nucleotides flanking the poly(A) signal itself specify genuine
16 cleavage sites. Thus, despite the fundamental nature of pre-mRNA cleavage for the
17 expression of all protein-coding genes, the motifs controlling this process are
18 dynamic on evolutionary timescales, providing motivation for future biochemical and
19 structural studies as well as new therapeutic angles to target eukaryotic pathogens.

20

21 **Keywords:** polyadenylation, poly(A) signal, protists, machine learning, long read
22 sequencing, evolution

23

24

1 INTRODUCTION

2 Cleavage and polyadenylation (CPA) are key steps in eukaryotic mRNA maturation,
3 specifying the 3' end of the transcript and addition of the poly(A) tail. Required for the
4 proper expression of nearly all mRNAs, the site of pre-mRNA cleavage is determined
5 by *cis*-regulatory RNA motifs, including the poly(A) [for polyadenylation] signal and
6 auxiliary elements found upstream and downstream of the poly(A) signal. Decades
7 of research have been spent defining these signals and their corresponding trans-
8 acting factors in model systems like humans and yeast, but we know much less
9 about the diversity and evolutionary dynamics throughout the eukaryotic tree of life.

10 The poly(A) signal is required for 3' end processing. The most well-known
11 poly(A) signal is the AAUAAA hexamer used in humans and other metazoans,
12 originally discovered in the mid-1970s (Proudfoot and Brownlee 1976; Fitzgerald and
13 Shenk 1981; Higgs et al. 1983; Montell et al. 1983). Elegant biochemical and
14 structural studies have shown how this hexamer is recognized by a multiprotein
15 complex known as the cleavage and polyadenylation specificity factor (CPSF) and
16 directly bound by two components of the CPSF complex (CPSF30 and WDR33) with
17 the remaining components facilitating cleavage (Mandel et al. 2008; Shi et al. 2009;
18 Mandel et al. 2006). Outside of metazoans, AAUAAA has been reported to be used
19 in *Arabidopsis* and *S. pombe* (Graber et al. 1999b; Liu et al. 2017), while *S.*
20 *cerevisiae* prefer an A-rich poly(A) signal (Graber et al. 1999b; Liu et al. 2017). Some
21 unicellular eukaryotes have been reported to use different and shorter poly(A) signal
22 sequences; for instance, a UGUAA is used in the green algae *Chlamydomonas*
23 *reinhardtii* and *Oscrococcus lucimarinus* (Shen et al. 2008; Zhao et al. 2019), while
24 UAAA/GUAA tetramers or UAA trimers are used in red algae and diatoms,
25 respectively (Zhao et al. 2019). We recently discovered that an assemblage A strain

1 of the protist *Giardia lamblia* uses an unusual AGURAA poly(A) signal (Bilodeau et
2 al. 2022). Thus, although the poly(A) signal is necessary for the correct expression of
3 nearly all coding genes, there appears to be plasticity in the sequence itself, which is
4 likely reflected in corresponding changes to the CSPF complex.

5 In most species studied to date, the poly(A) signal is rarely sufficient for
6 cleavage and proper cleavage requires additional auxiliary elements (Sheets et al.
7 1990; Birse 1997). For instance, metazoans have two major auxiliary elements: an
8 upstream U-rich motif and downstream U- and GU-rich motifs. The most canonical
9 U-rich motif is a UGUA tetramer recognized by proteins in the Cleavage factor Im
10 (CFIm) family (Brown and Gilmartin 2003; Venkataraman et al. 2005). U- and GU-
11 rich sequences downstream from the cleavage site are recognized by cleavage
12 stimulation factor proteins (CstF) (Zarudnaya et al. 2003; Beyer et al. 1997; Takagaki
13 and Manley 1997). Auxiliary elements boost the assembly of the cleavage and
14 polyadenylation machinery on the poly(A) signal and direct the endonuclease
15 CPSF73 for cleavage of the nascent RNA (Takagaki and Manley 1997; Hu et al.
16 2005; Mandel et al. 2006; Sullivan et al. 2009). In yeast, UA-rich elements located
17 ~40 nt upstream are bound by the cleavage and polyadenylation factor 1B (CF1B),
18 while U-rich elements surrounding the cleavage site are bound by the
19 polyadenylation factor complex CPF (Guo and Sherman 1996; Graber et al. 1999a).
20 Although less studied than the Amorphea clade (which contains humans and yeast;
21 Figure 1A), auxiliary elements have also been defined in plants where far upstream
22 elements are used although they lack a highly conserved consensus sequence (Wu
23 et al. 1995; Rothnie 1996). GU-rich downstream sequences, similar to human ones,
24 were also found in the parasitic protist *Blastocystis hominis* within the SAR clade (Li
25 and Du 2014). We previously found that *G. lamblia* assemblage A does not use any

1 known auxiliary element (Bilodeau et al. 2022), highlighting the potential diversity of
2 cleavage site recognition within the eukaryotic tree.

3 A major limitation to our understanding of pre-mRNA cleavage is that
4 knowledge of poly(A) signals, auxiliary motifs, and even 3'UTR annotations are
5 limited to an extremely small subset of eukaryotic species (the majority of which lie in
6 the Amorphea clade) and to only a few of the ~200,000 protist species. Thus, we
7 know very little about the evolution and diversity of the motifs specifying pre-mRNA
8 cleavage. For instance, a major gap is how poly(A) signals were identified in the
9 eukaryotic common ancestor. One explanation is that the examples of AAUAAA in
10 humans and possibly plants could reflect convergent evolution, possibly due to other
11 sequence constraints. In such a case, perhaps the ancestral Metamonada species
12 went down a different trajectory to use AGURAA. An alternative model is that
13 AAUAAA is in fact the ancestral sequence, and the AGURAA hexamer emerged
14 within the Metamonada clade, which would then have led to potentially fascinating –
15 and unexplored – genomic adaptations. Similar questions surround the use of
16 auxiliary elements, especially given the lack of any known ones in *Giardia lamblia*.

17 To shed light on these fundamental evolutionary questions, we set out to
18 characterize the diversity and evolutionary dynamics of strategies used to recognize
19 poly(A) sites in six protists, including four within the Metamonada clade, one within
20 Amorphea, and one within the Discoba clade. We also aimed to determine how
21 changes to poly(A) signals and auxiliary elements underlie genomic adaptations
22 allowing Metamonada species, and *Giardia lamblia* in particular, to specify genuine,
23 as opposed to premature, poly(A) signals.

24

25

1 RESULTS

2 Annotation of 3'UTRs across diverse eukaryotic species

3 We previously found that *G. lamblia* uses the distinct, but well-defined, AGURAA
4 poly(A) signal (Bilodeau et al. 2022) rather than the AAUAAA signal used in many
5 organisms, including humans (Chan et al. 2011). To understand how the poly(A)
6 signal evolved in eukaryotes, we explored poly(A) signals in other eukaryotic
7 species. Although our initial experiments with *G. lamblia* assemblage A used both
8 short and long-read sequences to identify poly(A) signals (Bilodeau et al. 2022), we
9 reasoned that direct-RNA Oxford Nanopore Technology would be sufficient for
10 identification of poly(A) signals because the method primes from the 3' end of a
11 transcript and does not require high quality genome annotations. Direct RNA
12 sequencing also enables determination of other features of interest, including 3'UTR
13 annotations and poly(A)-tail length.

14 We focused on several criteria to identify possible species to include in our
15 survey: evolutionary distance to *G. lamblia*; availability of RNA samples; and existing
16 annotations and other biological insights. Based on these criteria, we selected four
17 protists from the Metamonada clade, including two related *Giardia* species, *G.*
18 *lamblia* assemblage B (which is genetically diverse from the previously examined
19 assemblage A strain (Adam et al. 2013; Zajackowski et al. 2021)), and *G. muris*,
20 and two *Trichomonadida* species (*Trichomonas vaginalis* and *Tritrichomonas foetus*)
21 from the separate *Parabasalia* phylum (Figure 1A). We also selected two outgroup
22 protists, *Naegleria gruberi* from the Discoba clade and *Entamoeba histolytica* from
23 the Amorphea clade, which also contains well-studied organisms like *Homo sapiens*
24 and *S. cerevisiae*.

1 We analyzed at least two biological replicates using direct RNA sequencing
2 for each organism (Supplemental Fig S1A). Results from the replicates were highly
3 correlated, and the vast majority of reads contained untemplated adenosines, as
4 expected (Supplemental Fig S1C). In general, the median tail length was similar to
5 that seen in all eukaryotic species, with the shortest occurring in *N. gruberi* (median
6 length = 33). Nonetheless, even in this species, the poly(A) tail is sufficiently long to
7 accommodate the predicted 26-nucleotide footprint of poly(A) binding protein (Baer
8 and Kornberg 1983) (Supplemental Fig S1B). From a practical standpoint, tail
9 lengths in these protists are long enough to introduce few biases during the standard
10 oligo(dT)-enrichment step of RNA sequencing (except possibly for *N. gruberi* where
11 care is warranted). For the rest of the analysis, we discarded reads with untemplated
12 poly(A) tails shorter than 15 nucleotides for *N. gruberi* and 30 nucleotides for other
13 protists, reasoning that these reads could represent decay intermediates rather than
14 mature transcripts. Consistent with this interpretation, overall expression was lower
15 for the discarded transcripts in comparison to the retained ones (Supplemental Fig
16 S1D).

17 Based on these criteria and minimum of 10 reads combined from all
18 replicates, we were able to determine 3'UTRs for 1,409–5,867 genes across the six
19 species (Supplemental Fig S1E). In many cases, our annotations are an
20 improvement over previous annotations: for instance, prior to our study, no 3'UTRs
21 were annotated in *T. foetus* or in the assemblage B strain of *G. lamblia*. Our analysis
22 identified 5,867 and 2,630 3'UTRs, respectively, corresponding to 23% and 59% of
23 all annotated genes for these organisms. We were also able to annotate 101–914
24 novel isoforms of previously annotated 3'UTRs (Supplemental Table S1). In general,
25 protist 3'UTRs were substantially shorter than human 3'UTRs (Supplemental Fig

1 S1F), and in some cases even shorter than those of *G. lamblia*, which has been
2 previously noted for its short untranslated regions (Franzén et al. 2013; Bilodeau et
3 al. 2022). For instance, the median length in *E. histolytica* and *N. gruberi* was 22 and
4 39 nucleotides, respectively – meaning that the “functional” 3’UTR sequence space
5 for RNA-binding proteins seems to be minimal for many transcripts in several
6 protists, especially once the ribosome shadow of ~28 nucleotides is considered
7 (Takyar et al. 2005; Yusupova et al. 2001). Thus, it may be that ultrashort 3’UTRs
8 are more widespread than previously thought.

9

10 **WGURAA polyadenylation signal is a derived trait in the *Giardia* genus**

11 To determine the poly(A) signals for all seven protists, we next analyzed the
12 sequence surrounding the cleavage sites of annotated 3’UTRs. Because poly(A)
13 signals are upstream of cleavage sites, we performed *de novo* motif discovery on the
14 40 nucleotides upstream of the transcript end. Consistent with our previous work,
15 both *G. lamblia* assemblage A and *G. lamblia* assemblage B showed enrichment of
16 AGURAA-like sequences (Figure 1B), as did *G. muris*. This extended analysis also
17 revealed that UGUAAA was used to some extent in all *Giardia* species, leading us to
18 redefine the consensus poly(A) signal as WGURAA (Figure 1C, blue sequences). In
19 other words, all *Giardia* species analyzed to date prefer “G”, not “A”, in the second
20 position of the poly(A) signal.

21 A different picture emerged when we analyzed the other species. For the
22 other Metamonada organisms (*T. vaginalis* and *T. foetus*), the most frequent
23 hexamer was AAUAAA, followed by AUUAAA. Similarly, for *N. gruberi* in the Discoba
24 clade, AAUAAA was again the top hexamer. This result is consistent with previous
25 reports based on limited short-read data (Huang et al. 2013; Espinosa et al. 2002;

1 Fuentes et al. 2012). Although *E. histolytica* showed more diverse poly(A) signal
2 hexamer use, its most highly used sequences align with the previously identified
3 AAWUDA sequence (Hon et al. 2013). Thus, together with previous reports in *A.*
4 *thaliana* and humans (Graber et al. 1999b), the AAUAAA hexamer has now been
5 found in four clades (Archaeplastida, Amorphea, Discoba, and Metamonada), while
6 the WGURAA hexamer appears to be restricted to Metamonada.

7 In most multicellular eukaryotes, canonical replication-dependent histone
8 transcripts lack polyadenylation signals and poly(A) tails. In our analysis, we
9 observed that most of the annotated histone transcripts were polyadenylated, which
10 is consistent with the lack of a *SLBP* homolog in *N. gruberi* and *G. lamblia* (Dávila
11 López and Samuelsson 2008; Yee et al. 2007). For all the protist species in our
12 analysis, the poly(A) signal usage patterns for histone transcripts were similar to
13 other transcripts (Supplemental Fig S1G, Figure 1C), which suggests common
14 mechanisms of 3' end formation.

15 To ensure that patterns of poly(A) signals usage we observed were not biased
16 by transcripts that are polyadenylated but not transcribed by RNA polymerase II or
17 processed by the nuclear 3' end processing machinery, we repeated the analysis
18 focusing on mRNAs encoding ribosomal proteins. These transcripts were highly
19 conserved as well as highly expressed (Supplemental Fig S1H). We observed the
20 poly(A) signal usage pattern for transcripts encoding ribosomal proteins reflected the
21 pattern for annotated 3'UTRs (Supplemental Fig S1I, Figure 1C). Thus, the observed
22 differences in poly(A) signal hexamers are unlikely to be explained by non-pol II
23 transcripts with alternative 3' end formation.

24 Additionally, we checked poly(A) signal usage in the transcripts removed from
25 the analysis due to short poly(A) tail, and it was consistent with the pattern observed

1 in transcripts kept in the analysis (Supplemental Figure 1J, Figure 1C), except in the
2 case of *G. muris*, when only ~half of 3'UTRs contain WGURAA or AWUAAA signals.
3 However, this result can be explained because 29/61 of the transcripts removed in
4 *G. muris* were rRNAs. This finding is in line with our previous observation that rRNA
5 transcripts have short poly(A) tails likely processed by the TRAMP complex
6 (Bilodeau et al. 2022).

7 We also investigated the positional preference of the poly(A) signals relative
8 to the cleavage site because this distance was shown to be an important
9 determinant of cleavage position (Chen et al. 1995). The poly(A) signals of *E.*
10 *histolytica*, *N. gruberi*, *T. foetus*, and *T. vaginalis* were typically ~15-20 nucleotides
11 upstream of the cleavage sites, which is similar to the distances reported for
12 Metazoans (Kumar et al. 2019) (Figure 1D). The poly(A) signals in all three *Giardia*
13 strains were much closer to the cleavage site and only eight to nine nucleotides
14 upstream.

15 We conclude that AWUAAA is the ancestral poly(A) signal, and that the use of
16 “G” at the second position is a derived trait within, at least, the *Giardia* genus. This
17 sequence preference coincides with positional preferences for cleavage sites, which
18 may reflect changes in the underlying cleavage and polyadenylation machinery or
19 other evolutionary constraints on the surrounding sequence elements.

20

21 **“Dual-use” stop codons are most frequent for WGURAA hexamers in *Giardia***
22 **species**

23 One hallmark of organisms with highly compact genomes, like *G. lamblia* and
24 *Spironucleus salmonicida*, are so-called “dual use” stop codons, where the stop
25 codon itself also provides part of the poly(A) signal (Xu et al. 2014, 2020; Bilodeau et

1 al. 2022) (Figure 2A). The shift in poly(A) signal increases the number of possible
2 dual use stop codons from five to nine (Figure 2A). To examine the potential extent
3 of dual-use signals, we calculated the percentage of annotated protein coding genes
4 in which AAUAAA or WGURAA overlaps with a stop codon. We found that 7.4, 12.8,
5 and 18.3% of them had potential dual-use signals for the WGURAA hexamer in *G.*
6 *lamblia* A, *G. lamblia* B, and *G. muris*, respectively, but only 1–2.1% for the AAUAAA
7 hexamer (Figure 2B).

8 Given the short 3'UTR lengths in non-*Giardia* protists such as *N. gruberi*, we
9 examined the possibility that dual-use stop codons might play a role in these protists.
10 (Due to the apparent degeneracy in poly(A) signal specification in *E. histolytica*, we
11 did not include this species in this or further analysis.) Consistent with WGURAA not
12 being a poly(A) signal in *N. gruberi*, and *T. foetus*, there was little overlap of this
13 sequence with stop codons. In contrast, for AAUAAA, ~11–13% of stop codons
14 overlapped with this signal in *N. gruberi* and *T. foetus*. Stop codons overlapped with
15 AAUAAA and WGURAA to similar extent in *T. vaginalis* (Figure 2B). To assess, the
16 potential of these putative dual-use stop codons to act as poly(A) signals, we
17 examined the 3'UTR lengths of the corresponding transcripts (Figure 2C). As
18 expected, and previously reported (Bilodeau et al. 2022), the median length of
19 WGURAA dual-use sites in *G. lamblia* assemblage A was 11 nucleotides, which is in
20 line with the distance between poly(A) signals and cleavage sites in this species
21 (Figure 1D); for AAUAAA-containing sites, the median length was 24. The opposite
22 picture emerged for *N. gruberi*, where the AAUAAA-containing sites had a median
23 length of 22 nucleotides compared to 39 nucleotides for the WGURAA sites. These
24 results indicate that *N. gruberi* also has dual-use stop codons, although they may be
25 less common than in *Giardia*. In contrast, based on 3'UTR lengths we saw limited

1 evidence for dual-use sites in *T. vaginalis* and *T. foetus*. Thus, dual-use stop codons
2 appear to be a common strategy to allow the production of ultrashort 3'UTRs within
3 some protist species.

4

5 **Known auxiliary elements are absent in Metamonada**

6 Poly(A) signals are necessary, but typically not sufficient to specify cleavage sites,
7 and auxiliary elements are often used to help discriminate between “genuine” and
8 “premature” sites. One important motif is the upstream UGUA motif, which is
9 recognized in humans by CFIm. We had previously found little evidence for its use in
10 *G. lamblia* assemblage A (Bilodeau et al. 2022) and wondered how evolution has
11 shaped the use of the UGUA motif across eukaryotes. To look for evidence of its
12 use, we compared the occurrence of the UGUA motif 20-50 nucleotides upstream
13 from the cleavage site to the shuffled versions (Figure 3A; Supplemental Fig S2). As
14 expected, there was a clear signal in humans for this motif, and we also found
15 evidence in *N. gruberi* for the presence of this motif. In contrast, we found no
16 evidence for the UGUA motif in any Metamonada species, irrespective of poly(A)
17 signal identity. Given the evidence of an analogous “far upstream element” in plants
18 (Li and Hunt 1997), these data are consistent with a loss of the UGUA element early
19 in Metamonada evolution.

20 Because downstream GU- and U-rich elements also help specify poly(A) sites
21 in metazoans, we next examined these motifs among the protists, comparing their
22 occurrence 40 nucleotides downstream of the cleavage site with those 40
23 nucleotides upstream (as a control set). Using the ratio in humans as a benchmark,
24 we found that all protist species in our analysis exhibited significantly less
25 downstream GU- and U-rich positional bias (Figure 3B). Taken together, our results

1 highlight a substantial plasticity in the role of auxiliary elements within eukarya; most
2 notably, those in the Metamonada clade lack known auxiliary elements, which
3 suggests there may be alternative mechanism(s) of poly(A) signals recognition and
4 discrimination in this clade.

5

6 **Nucleotides flanking the poly(A) signal are crucial for proper discrimination in**
7 ***G. lamblia***

8 Premature cleavage has the potential to result in reduced or absent gene expression
9 due to truncated open reading frames and subsequent repression by surveillance
10 pathways. Our results thus far posed a riddle: how did a new poly(A) signal evolve
11 without triggering premature cleavage at sites that previously would not have been
12 recognized, especially given the apparent absence of other *cis*-elements to
13 discriminate between true and premature sites? We hypothesized that one
14 mechanism might be reduced occurrence of the signal itself within open reading
15 frames. Because coding region sequences are also shaped by amino acid biases,
16 we first compared the occurrences of the AGU–AAA dicodons (which encode serine
17 – lysine and could also serve as a poly(A) signal) with the eleven other synonymous
18 codon combinations (Figure 4A). The AGU-AAA dicodon was depleted relative to the
19 other combinations in the *Giardia* species but not in *T. vaginalis* or *T. foetus*, which
20 use the ancestral signal. We observed similar results with the other minor poly(A)
21 signals (Supplemental Fig S3A). Such depletions were not observed when we
22 reversed the codons (*e.g.* AAA-AGU; Supplemental Fig S3A). Thus, the WGURAA
23 poly(A) signal is depleted specifically in the coding regions of *Giardia species* in a
24 manner that cannot be explained by amino acid bias. This depletion is especially
25 pronounced in *G. muris*, where there are only 199 occurrences of the WGURAA

1 hexamer in any frame (Figure 4B), roughly 15x less than in *Giardia lamblia* A and
2 *Giardia lamblia* B.

3 Nonetheless, in *G. lamblia* A and *G. lamblia* B, we still found thousands of
4 examples of the poly(A) signal within coding regions (Figure 4B), suggesting that
5 mechanisms exist that allow discrimination of the correct poly(A) signal to avoid
6 premature cleavage and polyadenylation in these situations. Consequently, we
7 hypothesized that unique but heretofore unknown *cis*-elements must surround the
8 true poly(A) signals. To identify such elements, we focused first on *G. lamblia* and
9 used a gapped *k*-mer support vector machine (gkmSVM) (Ghandi et al. 2014, 2016)
10 to classify the 80 nucleotide regions surrounding WGURAA signals in the coding
11 sequence (*i.e.*, “premature”) and 3’UTRs (*i.e.*, “genuine”) (Figure 4C).

12 Testing several *k*-mer lengths revealed minimal increase in classifier
13 performance beyond eight nucleotides with the gkmSVM model performing
14 exceedingly well (as judged by an F1 score exceeding 0.9; Supplemental Table S2).
15 Given that the 8-nucleotide model only minimally exceeds the poly(A) signal length,
16 we wondered whether nucleotides directly surrounding the poly(A) signal might be
17 important for its recognition. We extracted gkmSVM scores for every 8-mer
18 containing WGURAA sequences and applied a linear model to determine the extent
19 to which the upstream nucleotide, downstream nucleotide, and poly(A) signal
20 sequence could explain the variation in 8-mer scores. Specifically, pyrimidines were
21 enriched downstream of the 3’UTR classified sites, and AGUAAA was more
22 associated with 3’UTRs, consistent with a model that this hexamer is “stronger” than
23 AGUGAA (Figure 4D, Supplemental Fig S3B).

24 These three features – the upstream and downstream nucleotides, and
25 poly(A) signal – explained at least 90% of the variance in 8-mer scores for both *G.*

1 *lamblia* A and *G. lamblia* B. When the contribution of each feature was determined
2 individually, the downstream nucleotide had the largest contribution, explaining more
3 than 70% of the variance, followed by the poly(A) signal identity (e.g., AGUAAA vs
4 AGUGAA, 11% of the variance) and then the upstream nucleotide (5% of the
5 variance; Figure 4E, Supplemental Fig. 3C). Thus, the flanking nucleotides,
6 especially the nucleotides directly downstream of the poly(A) signal, are important for
7 distinguishing genuine poly(A) signals from premature ones.

8 To explore the possibility that additional, but as-yet unknown, auxiliary
9 elements might support poly(A) signal discrimination and detect putative regulatory
10 elements shorter than 5 nucleotides (the minimal *k*-mer length for gkmSVM), we
11 generated *de novo* position weight matrices (PWMs) by systematically merging the
12 most predictive *k*-mers from trained gkmSVM models (Ghandi et al. 2014). The top
13 putative regulatory elements still overlapped with WGURAA and therefore contained
14 at least a partial poly(A) signal and flanking nucleotide (Supplemental Fig S3D, E),
15 but outside these features, the sequences were highly degenerate. Thus, it appears
16 that little beyond the eight-nucleotide element containing the poly(A) signal helps
17 differentiate between premature and true cleavage sites.

18 Despite the high performance of gkmSVM models, some coding sequences
19 were misclassified as 3'UTRs. These exceptions were especially interesting because
20 they contained 3'UTR-like features. We wondered if these might be examples of
21 premature cleavage and visually inspected them. Of the 23 *G. lamblia* potentially
22 false positive genes, 14 genes, including *GL50803_1890* (Figure 4F), contained
23 reads consistent with premature cleavage events. The remaining sites were due to
24 either incorrect gene annotations (seven cases) or classifier errors (two cases).
25 Similarly, of the eight *G. lamblia* B false positives, five showed evidence of

1 premature cleavage events (Supplemental Fig S3F). Given that such premature
2 cleavage is likely to reduce protein expression, it was interesting to note that three of
3 these sites were within duplicated genes (two in *G. lamblia* A; one in *G. lamblia* B),
4 and at least one copy of the duplicated genes lacked WGURAA hexamer
5 (Supplemental Fig S3G and H), providing a potential explanation for how the
6 consequences of a functional premature cleavage site may be minimized within *G.*
7 *lamblia*.

8

9 **Majority of coding region WGURAA hexamers are recognized as poly(A)** 10 **signals in *G. muris***

11 Given that *G. muris* contains very few WGURAA hexamers, we wondered about the
12 extent to which flanking nucleotides helped discriminate between poly(A) signals in
13 this organism. We again used the gkmSVM classifier and found the *G. muris* model
14 also had high performance with an F1 score of 0.99. In contrast to our results with *G.*
15 *lamblia* however, the linear model only explained 51% of the variance in 8-mer
16 scores (Fig. 5A, Supplemental Fig S4A). Given the low performance of the linear
17 model, we asked if dependencies between flanking nucleotides and poly(A) signal
18 sequences may be important. Indeed, by including interactions between terms, the
19 linear model now explained 85% of the variance (Figure 5A).

20 The *G. muris* model differed in several important ways from *G. lamblia*. First,
21 although the poly(A) signal and flanking nucleotides each contributed to the model,
22 the poly(A) signal was now the dominant feature with AGUGAA indicating a coding
23 region site, but AGUAAA and UGUAAA occurring very infrequently in this region. As
24 in *G. lamblia*, downstream purines were associated with coding region sites in *G.*
25 *muris*, but these effects were strongest when coupled with upstream pyrimidines

1 (Figure 5B, Supplemental Fig S4B, C). These findings indicate that in *G. muris*, the
2 full octomeric context is required to distinguish between premature and true cleavage
3 sites.

4 Although the classifier was able to correctly recognize all coding sequences
5 containing WGURAA hexamers, we wondered about the extent of premature
6 cleavage in *G. muris*. We had sufficient read support for 106 genes and observed
7 premature cleavage in 86% cases (Figure 5C). We conclude that despite flanking
8 sequence differences distinguishing WGURAA hexamers within coding sequences
9 from those within 3'UTRs, a substantial fraction has remained as functional poly(A)
10 signals.

11

12 **Flanking nucleotides play a role throughout the Metamonada clade**

13 Having found evidence for a role of flanking nucleotides within *Giardia* species, we
14 explored if the flanking nucleotides might have co-evolved with the derived poly(A)
15 signal, perhaps enabling adaptation to the change in the recognized hexamer. As
16 before, we trained gkmSVM classifiers for *T. foetus*, *T. vaginalis*, and the outgroup
17 species of *N. gruberi*. The classifier showed high performance for all three with the
18 highest being for *N. gruberi* (F1 = 0.9), and lower for *T. foetus* and *T. vaginalis* (F1 =
19 0.78 for both).

20 The linear models for *T. foetus* and *T. vaginalis* performed less well than for
21 *Giardia* (57% and 55%, respectively), but they still revealed the relative importance
22 of poly(A) signals and flanking nucleotides. Indeed, contrary to our original
23 hypothesis, the upstream and downstream flanking nucleotides contributed to model
24 performance for both *T. vaginalis* and *T. foetus*. In other words, despite these
25 species using the ancestral poly(A) signal, flanking nucleotides still play a role in

1 poly(A) site discrimination (Figure 6A, B). Incorporating dependencies between
2 poly(A) signals and flanking nucleotides in *T. foetus* and *T. vaginalis* did not improve
3 model performance (Supplemental Table S2). Thus, the three major elements of
4 poly(A) signal and flanking nucleotides contribute substantially to poly(A) site
5 recognition in all five Metamonada species analyzed. In the case of *T. foetus* and *T.*
6 *vaginalis*, we suspect they utilize additional unknown mechanism(s) given the
7 frequent occurrence of AAUAAA hexamers in the coding sequences.

8 In *N. gruberi*, the linear model for 8-mer scores explained 93% of the
9 observed variance with most of the contribution coming from the poly(A) signal, while
10 the flanking nucleotides explained none of the variance. The strongest differentiator
11 appeared to be AAUAAA in 3'UTR sites, while AUUAAA occurred in coding regions,
12 an observation which is consistent with AAUAAA acting as a stronger poly(A) signal.
13 However, we still found 1,238 instances of AAUAAA hexamers in coding sequences
14 (Figure 6C). We were curious to understand if and how proper poly(A) signal
15 discrimination occurred for transcripts containing the AAUAAA signal in both coding
16 sequences and 3'UTRs and re-trained the classifier on these sequences. This model
17 also had very high performance (0.96), which suggested the presence of other
18 auxiliary elements since the poly(A) signal for this restricted set was the same
19 between coding sequences and 3'UTRs. As with the full set, flanking nucleotides still
20 made almost no contribution. To figure out the identity of the auxiliary elements, we
21 clustered the top 50 8-mers with the highest gkmSVM scores. One of the sequences
22 we identified contained a UGUA tetranucleotide and was preferentially located
23 upstream of the poly(A) signal (Figure 6D). The other sequence was U-rich and
24 preferentially located downstream of the poly(A) signal in 3'UTRs but not in coding
25 sequences (Figure 6E). Thus, *N. gruberi* uses UGUA as an upstream element, which

1 is consistent with its enrichment (Figure 3A). Our earlier analysis showed that, en
2 masse, U- and GU-rich sequences had a lower downstream to upstream ratio in *N.*
3 *gruberi* compared to humans (Figure 3B). Together, these data indicate that *N.*
4 *gruberi* chiefly use the identity of the poly(A) signal for proper discrimination, which is
5 supplemented by the same primary upstream and downstream auxiliary elements
6 used in metazoans.

7

8

1 **DISCUSSION**

2 In this study, we used long-read direct RNA sequencing to map cleavage and
3 poly(A) sites in six protist species to understand the evolution of poly(A) signals and
4 auxiliary regulatory elements. We generated high-resolution annotation for 18,852
5 3'UTRs across six protist species, which will be a valuable resource for the scientific
6 community. We previously found that *G. lamblia* A uses the unusual poly(A) signal,
7 WGURAA (Bilodeau et al. 2022). Here, we confirmed that finding and extended it by
8 identifying the same poly(A) signal in other species within the *Giardia* genus.
9 However, other Metamonada protists (*T. vaginalis*, *T. foetus*) as well as those in
10 Discoba (*N. gruberi*) use the AAUAAA poly(A) signal, now providing evidence of this
11 signal in four eukaryotic clades. Based on parsimony, we conclude that AAUAAA
12 represents the ancestral poly(A) signal and that WGURAA is a derived trait within the
13 *Giardia* genus. This result is in agreement with previous study suggesting that
14 AAUAAA or its 3' part UAAA is an ancestral poly(A) signal (Zhao et al. 2019).
15 Because the fish pathogen *Spironucleus salmonicida* has been reported to use an
16 AGUGA poly(A) signal (Xu et al. 2014), we propose that the shift from the ancestral
17 signal occurred in the Diplomonada order (Figure 7), although more precise
18 evolutionary placement will require analysis of additional species.

19 Previous work in several species had focused on the known auxiliary
20 elements of the upstream UGUA and downstream U- and GU-rich elements. In the
21 case of *N. gruberi*, UGUA motifs are enriched upstream of poly(A) signals, and U-
22 rich sequences downstream. Both seem to be used to discriminate AAUAAA poly(A)
23 signals from ~1,200 such hexamers in coding sequences. Nonetheless, in
24 Metamonada, irrespective of the poly(A) signal used, known auxiliary elements are
25 neither enriched nor contribute to recognition based on our modeling. Given that

1 these elements play a role in *T. brucei* (in the Discoba clade), *A. thaliana*, and *H.*
2 *sapiens*, we propose that they were present in the eukaryotic ancestral recognition
3 and then lost early in Metamonada evolution.

4 In fact, one surprise of our analysis is the diversity of strategies to distinguish
5 premature poly(A) sites from mature ones aside from the previously described
6 upstream and downstream motifs. Our analysis revealed at least three other
7 strategies. One strategy is the identity of the poly(A) signal itself. The best example
8 of this strategy is *N. gruberi*, although this specification also occurs in *G. lamblia*.
9 Another strategy is widespread depletion of poly(A) signals from coding sequences.
10 This strategy seems to be very prevalent in *G. muris*, but not *G. lamblia*, indicating a
11 recent, widespread genomic adaptation.

12 Finally, all studied Metamonada species rely, to some extent, on the
13 nucleotides flanking the poly(A) signal. To our knowledge, a role for these flanking
14 nucleotides – to make a functional octameric poly(A) signal – has not been observed
15 in any species to date. The role of the flanking nucleotides, especially the nucleotide
16 directly downstream of the WGURAA hexamer, is especially strong for *G. lamblia*
17 species, raising the intriguing possibility that increasing reliance on these poly(A)
18 site-adjacent nucleotides may have been an alternative “strategy” to adapt to the
19 change in poly(A) signal. These flanking nucleotides contribute to recognition in all
20 Metamonada species, including *T. vaginalis* and *T. foetus*, which use the ancestral
21 signal. However, weaker model performance for these two species together with the
22 frequent occurrence of AAUAAA hexamers in coding sequences suggest that
23 unknown additional *cis*-elements are likely involved in *T. vaginalis* and *T. foetus*
24 poly(A) signal recognition.

1 Given that all of these RNA motifs are recognized by RNA-binding proteins, it
2 is reasonable to predict that these proteins have also changed to recognize the
3 derived poly(A) signals and auxiliary elements. In the case of *G. lamblia*, the
4 streamlined polyadenylation machinery has only six to seven homologs to the human
5 core complex (Bilodeau et al. 2022; Ospina-Villa et al. 2020). The CPSF30 homolog,
6 which is directly involved in poly(A) signal recognition, contains *Giardia*-specific
7 amino acid changes within two critical zinc finger motifs that are responsible for
8 recognition of the second and fourth nucleotide in the poly(A) signal (Bilodeau et al.
9 2022; Ospina-Villa et al. 2020). It is tempting to hypothesize that these mutations
10 may provide a molecular explanation for the difference in poly(A) signals. Similarly,
11 we previously found a putative WDR33 homolog (Bilodeau et al. 2022); this protein
12 also contains amino acid substitutions specific to the *Giardia* genus, which may have
13 augmented the role of the downstream nucleotide in poly(A) signal recognition.
14 Understanding the structural basis of the so-called octomeric poly(A) signal will be
15 an important next step.

16 Our study thus highlights a hitherto unknown plasticity associated with pre-
17 mRNA cleavage and polyadenylation and shows the power of exploring fundamental
18 molecular processes beyond the standard model systems. We highlight how new
19 methods like long-read sequencing can improve our understanding of RNA biology
20 across the eukaryotic tree of life and motivate future studies into non-traditional
21 organisms, especially eukaryotic pathogens. In addition to enabling molecular
22 biologic studies (e.g., through the creation of robust transgenic gene expression
23 cassettes), our results set the stage for future structural and molecular studies to
24 better understand interactions between *cis*-elements and *trans* factors within protists.
25 Our discoveries highlight a potential avenue for pharmacologically targeting

- 1 eukaryotic pathogens by virtue of the differences in poly(A) site recognition
- 2 compared to metazoans.
- 3

1 **METHODS**

2 *G. lamblia* (assemblage A, strain WB clone C6) (ATCC 50803), and *G. lamblia*
3 (assemblage B, strain GS/M, clone H7) (ATCC 50581) were grown in Keister's
4 modified TYI-S-33 medium (Keister 1983). *E. histolytica* strain HM1-IMSS (ATCC
5 30459) was grown in TYI-S-33 medium (Diamond et al. 1978). *T. foetus* strain D1
6 and *T. vaginalis* strain F1623 were grown in trypticase-yeast extract maltose (TYM)
7 medium, supplemented with 10% horse serum (DH-05, Omega, Tarzana, CA),
8 0.74% ammonium iron (II) sulfate hexahydrate (215406, Sigma-Aldrich, St. Louis,
9 MO) and 1% penicillin-streptomycin (P4333, Sigma-Aldrich, St. Louis, MO), adjusted
10 to pH 6.2. *G. muris* Roberts-Thomson isolate was obtained from Waterborne Inc
11 (P105). Cysts were incubated in 15 ml of the induction medium (0.1 M potassium
12 phosphate buffer pH 7, 0.3 M sodium bicarbonate) in capped centrifuge tubes at
13 37°C for 10 minutes ($0.5-1.0 \times 10^6$ cysts/tube). Samples were then centrifuged at
14 600 x g for 5 minutes, washed in 15 ml of 0.1 M potassium phosphate buffer (pH
15 7.0), concentrated by centrifugation for 5 minutes at 600 x g and suspended in 1-2
16 ml of TYI-S-33 medium. All protists were cultured at 37°C. *N. gruberi* strain NEG-M
17 (ATCC 30224) was obtained from Chandler Fulton at Brandeis University and were
18 grown in M7 medium (0.362 g/L KH_2PO_4 , 0.5 g/L Na_2HPO_4 , 5.4 g/L glucose, 5 g/L
19 yeast extract (Difco), 45 mg/L L-methionine, 10% fetal bovine serum) at 26°C without
20 shaking in 25 cm² plug-seal tissue culture flasks (CellTreat Cat#229330).

21

22 **RNA extraction**

23 *G. lamblia*, *E. histolytica*, *T. vaginalis* and *T. foetus* were grown in log phase to no
24 more than 80% confluence, detached by cooling on ice (*E. histolytica* and *G. lamblia*)

1 or by vigorous shaking (*T. vaginalis*), and washed twice by centrifugation and
2 resuspension in ice-cold PBS.

3 *G. lamblia* assemblage A RNA was isolated using hot acid phenol as
4 previously described (Collart and Oliviero 2001). For remaining protists, the washed
5 pellets were lysed in TRI Reagent (Zymo research) and RNA was extracted and
6 purified using Direct-zol RNA Miniprep Kit (Zymo research) following the
7 manufacturer's instructions. *G. muris* samples were centrifuged (1,000 x g for 5
8 minutes), washed twice with PBS, resuspended in 1 ml of TRIzol and frozen at -
9 20°C. The following day, RNA was isolated following the manufacturer's protocol.
10 Total RNA concentration and purity (A260/A280 ratio) were quantified using a
11 NanoDrop™ 1000 instrument (Thermo Fisher Scientific). RNA concentrations were
12 >100 ng/μl and A260/A280 ratios > 1.75 for all samples. *N. gruberi* samples were
13 washed once in 2 mM Tris and centrifuged at 1500 x g at room temperature. The
14 pellet was suspended in 1 ml TRIzol, vortexed, and stored at -80°C until RNA
15 extraction. Samples were lysed using FastPrep homogenizer with bead beating in
16 TRIzol. Lysate was cleaned up using a Zymo kit with on-column DNase treatment.
17 The RNA concentration and integrity were measured by Qubit (Thermo Fisher
18 Scientific). Concentration was > 350 ng/μl and RIN value was between 5.8-6.6. Cell
19 harvesting and RNA extraction were repeated for each parasite strain on three
20 different days.

21

22 **RNA sequencing and analysis**

23 Nanopore libraries were prepared as described in (Bilodeau et al. 2022). Libraries
24 were sequenced on a FLO-MIN106 flow cell and minION sequencing device. For *G.*
25 *lamblia* assemblage A, we prepared one new library and used it as a third replicate.

1 Raw data for the other two replicates were from our previous study (Bilodeau et al.
 2 2022). Base-calling was performed using guppy v. 5.0.11+2b6dbffa5. Reads were
 3 aligned using minimap2 (Li 2018) v. 2.17-r941 with the following parameters: “-a -x
 4 splice -uf -k 14”. The following genomic annotations were used for the analyses:
 5

organism	genome	source
<i>Entamoeba histolytica</i>	Entamoeba_histolytica.JCVI-ESG2-1.0.52	Ensembl
<i>Giardia lamblia</i> A	GiardiaDB-56_GintestinalisAssemblageAWB	VEuPathDB
<i>Giardia lamblia</i> B	GiardiaDB-56_GintestinalisAssemblageBGS	VEuPathDB
<i>Giardia muris</i>	GiardiaDB-56_GmurisRobertsThomson	VEuPathDB
<i>Naegleria gruberi</i>	Naegleria_gruberi_gca_000004985.V1.0.52	Ensembl
<i>Trichomonas vaginalis</i>	TrichDB-56_TvaginalisG3	VEuPathDB
<i>Tritrichomonas foetus</i>	TrichDB-56_TfoetusK	VEuPathDB

6

7 Poly(A) tail length was estimated using nanopolish (Loman et al. 2015) v. 0.13.2.
 8 Mapped nanopore reads were assigned to their corresponding genes using

1 featureCounts (Liao et al. 2014) v. 2.8.2 from Rsubread (Liao et al. 2019) using R (R
2 Core Team 2021) .

3

4 **Cleavage sites identification and 3'UTR annotation**

5 To annotate the cleavage sites, first we select the reads with poly(A) tail length ≥ 15
6 (*N. gruberi*) or ≥ 30 (other protists) and for which the 5' end of the read fell within the
7 open reading frame of the associated gene. Then we selected genes with at least 10
8 reads total and grouped them by a read start position (3' end). Reads within the 20 nt
9 window were grouped together as a putative 3' end if the group contained at least
10 10% of total reads for a particular gene. Cleavage site was annotated at the 3' end of
11 the group. If more than one group was present for the gene they were considered
12 3'UTR isoforms with alternative cleavage and polyadenylation sites. The analysis
13 was performed using custom R scripts.

14

15 **Poly(A) signals, “dual use” signals and auxiliary elements identification and 16 quantification**

17 To identify putative poly(A) signals we performed de novo motif discovery on the 40
18 nucleotides upstream of the annotated cleavage sites using MEME (Bailey and
19 Elkan 1994) v. 5.0.5 with the following parameters -rna -mod zoops -minw 3 -maxw
20 10, and we selected the top motif based on the number of occurrences. MEME
21 PWMs were used to select top 6-mers for all organisms and calculate their
22 frequencies within 40 nucleotides upstream of the annotated cleavage sites.

23 To identify putative “dual use” signals where stop codon overlaps with poly(A)
24 signal we extracted stop codons coordinates for all annotated protein coding genes
25 in each organism. To capture all possible overlaps, we extended these coordinates 2

1 nucleotides upstream and 5 nucleotides downstream and searched for the
2 sequences listed in Figure 2A.

3 Occurrences of UGUA upstream auxiliary element and shuffled versions of
4 this tetramer were calculated in the region 20-50 nucleotides upstream from the
5 cleavage site. The statistical significance of the differences between humans and
6 protists was checked using the chi square test. Occurrences of GU and U-rich
7 auxiliary elements (CUGCCU, CUGGGG, CUGUGU, GUCUGU, GUGUCU,
8 GUGUGU, UGUCUC, UGUCUG, UGUGUG, UGUUUU, UUAUUU, UUUCUU,
9 UUUUUU) were compared 40 nt downstream from the cleavage vs 40 nt upstream
10 to the cleavage site region. Statistical significance was calculated in comparison to
11 human enrichment using a one-sided Wilcoxon rank sum test. The analyses were
12 performed using custom R scripts.

13

14 **Di-codons occurrence quantification**

15 All in-frame codon combinations encoding Ser-Lys, Lys-Ser, Ser-Cys, Cys-Ser, Ser-
16 Glu and Glu-Ser were counted within annotated coding sequences for every
17 organism using oligonucleotideFrequency function from Biostrings v. 2.62 R library
18 with the following settings: width = 6, step = 3.

19

20 **Gapped *k*-mer support vector machine analysis**

21 The positive and negative sets were prepared as follows: every
22 AGUAAA|UGUAAA|AGUGAA (*G. lamblia*, *G. muris*) or AAUAAA|AUUAAA (*N.*
23 *gruberi*, *T. vaginalis*, *T. foetus*) hexamer was extracted from either 3'UTR (positive
24 set) or coding sequence (negative set) and extended 37 nt upstream and 37 nt
25 downstream (total length 80 nt). 80-mers were used to train gkmSVM classifier as

1 described in <https://www.beerlab.org/gkmsvm/gkmsvm-tutorial.htm>. We used
2 gkmSVM (Ghandi et al. 2014, 2016) v. 0.81 R library with the following settings: *k*-
3 mer length: 8 (8, 10 and 12 were tested for *G. lamblia* and *G. muris* - Supplemental
4 Table S2), number of informative columns: 6, max number of mismatches: 3, add
5 reverse complement: FALSE. For *G. lamblia* and *G. muris* datasets all positive and
6 negative 80-mers were used for training. For *N. gruberi*, *T. vaginalis* and *T. foetus*,
7 negative sets were split into 5 subsets because of sample imbalance and training
8 was repeated for each subset. Model performance was measured using caret (Kuhn
9 2008) v. 6.0.92 R library and F1 metrics (Supplemental Table S2). *N. gruberi*, *T.*
10 *vaginalis* and *T. foetus* models trained on negative subsets 1 were used for further
11 analysis. The gkmSVM models were used to score all possible 8-mers containing
12 either AGUAAA|UGUAAA|AGUGAA or AAUAAA|AUUAAA. Scores were then used
13 to build linear models with three independent variables: upstream nucleotide,
14 hexamer sequence and downstream nucleotide (lm R function). To predict putative
15 new auxiliary elements de novo PWMs were built by systematically merging the most
16 predictive *k*-mers from trained gkmSVM models, as described in (Ghandi et al.
17 2014). *K*-mers clustering for *N. gruberi* AAUAAA-containing sequences was
18 performed using kmer v. 1.1.2 R library.

19

20 **Premature cleavage identification**

21 Coding sequences containing WGURAA hexamers misclassified as 3'UTRs (*G.*
22 *lamblia* A, *n* = 23, *G. lamblia* B, *n* = 14) or all coding sequences containing WGURAA
23 hexamer (*G. muris*, *n* = 199) were visually inspected in Integrative Genomics Viewer
24 (Robinson et al. 2011) v. 2.14. If according to the investigator assessment gene was
25 correctly annotated with a substantial coverage drop within 20 nt downstream from

1 WGURAA hexamer it was defined as premature cleavage event. Coverage data for
2 genes selected to figure preparation was converted from BAM to bigWig format
3 using deepTools (Ramírez et al. 2014) v. 3.5.1 and visualized using ggcoverage
4 (Song and Wang 2023) v. 1.2.0 R library. Duplication events for genes undergoing
5 premature cleavage were identified using OrthoFinder (Emms and Kelly 2019) v.
6 2.5.4.

7

8 **Data analysis, statistical tests, and visualization**

9 Statistical details of specific experiments can be found in the Results, Methods,
10 and/or Figure Legends. Raw plots were generated using the R ggplot2 (Wickham
11 2016) or ggpubr (Kassambara 2019) packages. Final figures were generated using
12 Adobe Illustrator.

13

14 **Data access**

15 All raw and processed sequencing data generated in this study have been submitted
16 to the NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>)
17 under accession number GSE260731. All code used for data analysis and figure
18 generation is accessible at https://github.com/mukherjeelab/2024_protists_polyA and
19 in Supplemental Code.

20

21 **Competing interest statement**

22 OSR is a member of the *Cell Reports* and *Molecular Cell* Scientific Advisory Boards.

23

24 **Acknowledgments**

1 We thank J. Mathew Taliaferro, David Bentley, and Srinivas Ramachandran for the
2 critical comments. This work was supported by the Polish National Agency for
3 Academic Exchange Bekker Program PPN/BEK/2019/1/00173 (M.P.S.), the
4 University of Colorado Anschutz Medical Campus RNA Bioscience Initiative (M.P.S.,
5 N.M., O.S.R.), National Institute of Health grants R01 AI158612 (L.E.), P30
6 DK120515 (L.E.), R35 GM128680 (O.S.R.), R35 GM147025 (N.M.) and a Pew
7 Biomedical Scholar Award to L.F.-L. who is a fellow in the Canadian Institute for
8 Advanced Research Fungal Kingdom program.

9

10 *Author contributions:* Conceptualization: O.S.R., N.M., M.P.S.; Methodology: M.P.S.,
11 M.A.B., N.M., O.S.R.; Software: M.P.S., M.A.B.; Formal Analysis: M.P.S., M.A.B.;
12 Investigation: M.P.S., D.Y.B., N.M., O.S.R.; Resources: Y.M., K.B.V., L.E., L.F-L.
13 N.M., O.S.R.; Data Curation: M.P.S.; Writing – Original Draft: M.P.S.; Writing –
14 Review & Editing: M.P.S., N.M., O.S.R.; Visualization: M.P.S., E.H., N.M., O.S.R.;
15 Supervision: N.M., O.S.R.; Project Administration: N.M., O.S.R.; Funding
16 Acquisition: M.P.S., L.E., L.F-L., N.M., O.S.R.

17

18

1 **REFERENCES**

- 2 Adam RD, Dahlstrom EW, Martens CA, Bruno DP, Barbian KD, Ricklefs SM,
3 Hernandez MM, Narla NP, Patel RB, Porcella SF, et al. 2013. Genome
4 sequencing of *Giardia lamblia* genotypes A2 and B isolates (DH and GS) and
5 comparative analysis with the genomes of genotypes A1 and E (WB and Pig).
6 *Genome Biol Evol* **5**: 2498–2511. <http://dx.doi.org/10.1093/gbe/evt197>.
- 7 Baer BW, Kornberg RD. 1983. The protein responsible for the repeating structure of
8 cytoplasmic poly(A)-ribonucleoprotein. *J Cell Biol* **96**: 717–721.
9 <http://dx.doi.org/10.1083/jcb.96.3.717>.
- 10 Bailey TL, Elkan C. 1994. Fitting a mixture model by expectation maximization to
11 discover motifs in biopolymers. *Proc Int Conf Intell Syst Mol Biol* **2**: 28–36.
12 <https://www.ncbi.nlm.nih.gov/pubmed/7584402>.
- 13 Beyer K, Dandekar T, Keller W. 1997. RNA ligands selected by cleavage stimulation
14 factor contain distinct sequence motifs that function as downstream elements
15 in 3'-end processing of pre-mRNA. *J Biol Chem* **272**: 26769–26779.
16 <http://dx.doi.org/10.1074/jbc.272.42.26769>.
- 17 Bilodeau DY, Sheridan RM, Balan B, Jex AR, Rissland OS. 2022. Precise gene
18 models using long-read sequencing reveal a unique poly(A) signal in *Giardia*
19 *lamblia*. *RNA* **28**: 668–682. <http://dx.doi.org/10.1261/rna.078793.121>.
- 20 Birse CE. 1997. Transcriptional termination signals for RNA polymerase II in fission
21 yeast. *EMBO J* **16**: 3633–3643. <http://dx.doi.org/10.1093/emboj/16.12.3633>.

- 1 Brown KM, Gilmartin GM. 2003. A mechanism for the regulation of pre-mRNA 3'
2 processing by human cleavage factor im. *Mol Cell* **12**: 1467–1476.
3 [http://dx.doi.org/10.1016/s1097-2765\(03\)00453-2](http://dx.doi.org/10.1016/s1097-2765(03)00453-2).
- 4 Chan S, Choi E-A, Shi Y. 2011. Pre-mRNA 3'-end processing complex assembly
5 and function. *Wiley Interdiscip Rev RNA* **2**: 321–335.
6 <http://dx.doi.org/10.1002/wrna.54>.
- 7 Chen F, MacDonald CC, Wilusf J. 1995. Cleavage site determinants in the
8 mammalian polydenylation signal. *Nucleic Acids Res* **23**: 2614–2620.
9 <http://dx.doi.org/10.1093/nar/23.14.2614>.
- 10 Collart MA, Oliviero S. 2001. Preparation of yeast RNA. *Curr Protoc Mol Biol*
11 **Chapter 13**: Unit13.12. <http://dx.doi.org/10.1002/0471142727.mb1312s23>.
- 12 Dávila López M, Samuelsson T. 2008. Early evolution of histone mRNA 3' end
13 processing. *RNA* **14**: 1–10. <http://dx.doi.org/10.1261/rna.782308>.
- 14 Diamond LS, Harlow DR, Cunnick CC. 1978. A new medium for the axenic
15 cultivation of *Entamoeba histolytica* and other *Entamoeba*. *Trans R Soc Trop*
16 *Med Hyg* **72**: 431–432. [http://dx.doi.org/10.1016/0035-9203\(78\)90144-x](http://dx.doi.org/10.1016/0035-9203(78)90144-x).
- 17 Emms DM, Kelly S. 2019. OrthoFinder: phylogenetic orthology inference for
18 comparative genomics. *Genome Biol* **20**. [http://dx.doi.org/10.1186/s13059-](http://dx.doi.org/10.1186/s13059-019-1832-y)
19 [019-1832-y](http://dx.doi.org/10.1186/s13059-019-1832-y).
- 20 Espinosa N, Hernández R, López-Griego L, López-Villaseñor I. 2002. Separable
21 putative polyadenylation and cleavage motifs in *Trichomonas vaginalis*
22 mRNAs. *Gene* **289**: 81–86. [http://dx.doi.org/10.1016/s0378-1119\(02\)00476-6](http://dx.doi.org/10.1016/s0378-1119(02)00476-6).

- 1 Fitzgerald M, Shenk T. 1981. The sequence 5'-AAUAAA-3' forms parts of the
2 recognition site for polyadenylation of late SV40 mRNAs. *Cell* **24**: 251–260.
3 [http://dx.doi.org/10.1016/0092-8674\(81\)90521-3](http://dx.doi.org/10.1016/0092-8674(81)90521-3).
- 4 Franzén O, Jerlström-Hultqvist J, Einarsson E, Ankarklev J, Ferella M, Andersson B,
5 Svärd SG. 2013. Transcriptome profiling of *Giardia intestinalis* using strand-
6 specific RNA-seq. *PLoS Comput Biol* **9**: e1003000.
7 <http://dx.doi.org/10.1371/journal.pcbi.1003000>.
- 8 Fuentes V, Barrera G, Sánchez J, Hernández R, López-Villaseñor I. 2012.
9 Functional analysis of sequence motifs involved in the polyadenylation of
10 *Trichomonas vaginalis* mRNAs. *Eukaryot Cell* **11**: 725–734.
11 <http://dx.doi.org/10.1128/ec.05322-11>.
- 12 Ghandi M, Lee D, Mohammad-Noori M, Beer MA. 2014. Enhanced regulatory
13 sequence prediction using gapped k-mer features. *PLoS Comput Biol* **10**:
14 e1003711. <http://dx.doi.org/10.1371/journal.pcbi.1003711>.
- 15 Ghandi M, Mohammad-Noori M, Ghareghani N, Lee D, Garraway L, Beer MA. 2016.
16 gkmSVM: an R package for gapped-kmer SVM. *Bioinformatics* **32**: 2205–
17 2207. <http://dx.doi.org/10.1093/bioinformatics/btw203>.
- 18 Graber JH, Cantor CR, Mohr SC, Smith TF. 1999a. Genomic detection of new yeast
19 pre-mRNA 3'-end-processing signals. *Nucleic Acids Res* **27**: 888–894.
20 <http://dx.doi.org/10.1093/nar/27.3.888>.
- 21 Graber JH, Cantor CR, Mohr SC, Smith TF. 1999b. In silico detection of control
22 signals: mRNA 3'-end-processing sequences in diverse species. *Proc Natl*

- 1 *Acad Sci U S A* **96**: 14055–14060.
2 <http://dx.doi.org/10.1073/pnas.96.24.14055>.
- 3 Guo Z, Sherman F. 1996. Signals sufficient for 3'-end formation of yeast mRNA. *Mol*
4 *Cell Biol* **16**: 2772–2776. <http://dx.doi.org/10.1128/MCB.16.6.2772>.
- 5 Higgs DR, Goodbourn SE, Lamb J, Clegg JB, Weatherall DJ, Proudfoot NJ. 1983.
6 Alpha-thalassaemia caused by a polyadenylation signal mutation. *Nature* **306**:
7 398–400. <http://dx.doi.org/10.1038/306398a0>.
- 8 Hon C-C, Weber C, Sismeiro O, Proux C, Koutero M, Deloger M, Das S, Agrahari M,
9 Dillies M-A, Jagla B, et al. 2013. Quantification of stochastic noise of splicing
10 and polyadenylation in *Entamoeba histolytica*. *Nucleic Acids Res* **41**: 1936–
11 1952. <http://dx.doi.org/10.1093/nar/gks1271>.
- 12 Hu J, Lutz CS, Wilusz J, Tian B. 2005. Bioinformatic identification of candidate *cis*-
13 regulatory elements involved in human mRNA polyadenylation. *RNA* **11**:
14 1485–1493. <http://dx.doi.org/10.1261/rna.2107305>.
- 15 Huang K-Y, Shin J-W, Huang P-J, Ku F-M, Lin W-C, Lin R, Hsu W-M, Tang P. 2013.
16 Functional profiling of the *Trichomonas foetus* transcriptome and proteome.
17 *Mol Biochem Parasitol* **187**: 60–71.
18 <http://dx.doi.org/10.1016/j.molbiopara.2012.12.001>.
- 19 Kassambara A. 2019. *GGPlot2 essentials*. Independently Published.
- 20 Keister DB. 1983. Axenic culture of *Giardia lamblia* in TYI-S-33 medium
21 supplemented with bile. *Trans R Soc Trop Med Hyg* **77**: 487–488.
22 [http://dx.doi.org/10.1016/0035-9203\(83\)90120-7](http://dx.doi.org/10.1016/0035-9203(83)90120-7).

- 1 Kuhn M. 2008. Building predictive models in R using the caret Package. *J Stat Softw*
2 **28**. <http://dx.doi.org/10.18637/jss.v028.i05>.
- 3 Kumar A, Clerici M, Muckenfuss LM, Passmore LA, Jinek M. 2019. Mechanistic
4 insights into mRNA 3'-end processing. *Curr Opin Struct Biol* **59**: 143–150.
5 <http://dx.doi.org/10.1016/j.sbi.2019.08.001>.
- 6 Li H. 2018. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*
7 **34**: 3094–3100. <http://dx.doi.org/10.1093/bioinformatics/bty191>.
- 8 Li Q, Hunt AG. 1997. The polyadenylation of RNA in plants. *Plant Physiol* **115**: 321–
9 325. <http://dx.doi.org/10.1104/pp.115.2.321>.
- 10 Li X-Q, Du D. 2014. Motif types, motif locations and base composition patterns
11 around the RNA polyadenylation site in microorganisms, plants and animals.
12 *BMC Evol Biol* **14**. <http://dx.doi.org/10.1186/s12862-014-0162-7>.
- 13 Liao Y, Smyth GK, Shi W. 2014. featureCounts: an efficient general purpose
14 program for assigning sequence reads to genomic features. *Bioinformatics*
15 **30**: 923–930. <http://dx.doi.org/10.1093/bioinformatics/btt656>.
- 16 Liao Y, Smyth GK, Shi W. 2019. The R package Rsubread is easier, faster, cheaper
17 and better for alignment and quantification of RNA sequencing reads. *Nucleic*
18 *Acids Res* **47**: e47–e47. <http://dx.doi.org/10.1093/nar/gkz114>.
- 19 Liu X, Hoque M, Larochelle M, Lemay J-F, Yurko N, Manley JL, Bachand F, Tian B.
20 2017. Comparative analysis of alternative polyadenylation in *S. cerevisiae* and
21 *S. pombe*. *Genome Res* **27**: 1685–1695.
22 <http://dx.doi.org/10.1101/gr.222331.117>.

- 1 Loman NJ, Quick J, Simpson JT. 2015. A complete bacterial genome assembled de
2 novo using only nanopore sequencing data. *Nat Methods* **12**: 733–735.
3 <http://dx.doi.org/10.1038/nmeth.3444>.
- 4 Mandel CR, Bai Y, Tong L. 2008. Protein factors in pre-mRNA 3'-end processing.
5 *Cell Mol Life Sci* **65**: 1099–1122. [http://dx.doi.org/10.1007/s00018-007-7474-](http://dx.doi.org/10.1007/s00018-007-7474-3)
6 [3](http://dx.doi.org/10.1007/s00018-007-7474-3).
- 7 Mandel CR, Kaneko S, Zhang H, Gebauer D, Vethantham V, Manley JL, Tong L.
8 2006. Polyadenylation factor CPSF-73 is the pre-mRNA 3'-end-processing
9 endonuclease. *Nature* **444**: 953–956. <http://dx.doi.org/10.1038/nature05363>.
- 10 Montell C, Fisher EF, Caruthers MH, Berk AJ. 1983. Inhibition of RNA cleavage but
11 not polyadenylation by a point mutation in mRNA 3' consensus sequence
12 AAUAAA. *Nature* **305**: 600–605. <http://dx.doi.org/10.1038/305600a0>.
- 13 Ospina-Villa JD, Tovar-Ayona BJ, López-Camarillo C, Soto-Sánchez J, Ramírez-
14 Moreno E, Castañón-Sánchez CA, Marchat LA. 2020. mRNA polyadenylation
15 machineries in intestinal protozoan parasites. *J Eukaryot Microbiol* **67**: 306–
16 320. <http://dx.doi.org/10.1111/jeu.12781>.
- 17 Proudfoot NJ, Brownlee GG. 1976. 3' non-coding region sequences in eukaryotic
18 messenger RNA. *Nature* **263**: 211–214. <http://dx.doi.org/10.1038/263211a0>.
- 19 Ramírez F, Dünder F, Diehl S, Grüning BA, Manke T. 2014. deepTools: a flexible
20 platform for exploring deep-sequencing data. *Nucleic Acids Res* **42**: W187–
21 W191. <http://dx.doi.org/10.1093/nar/gku365>.

- 1 Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G,
2 Mesirov JP. 2011. Integrative genomics viewer. *Nat Biotechnol* **29**: 24–26.
3 <http://dx.doi.org/10.1038/nbt.1754>.
- 4 Rothnie HM. 1996. Plant mRNA 3'-end formation. *Plant Mol Biol* **32**: 43–61.
5 <https://www.ncbi.nlm.nih.gov/pubmed/8980473>.
- 6 Sheets MD, Ogg SC, Wickens MP. 1990. Point mutations in AAUAAA and the poly
7 (A) addition site: effects on the accuracy and efficiency of cleavage and
8 polyadenylation *in vitro*. *Nucleic Acids Res* **18**: 5799–5805.
9 <http://dx.doi.org/10.1093/nar/18.19.5799>.
- 10 Shen Y, Liu Y, Liu L, Liang C, Li QQ. 2008. Unique features of nuclear mRNA
11 poly(A) signals and alternative polyadenylation in *Chlamydomonas reinhardtii*.
12 *Genetics* **179**: 167–176. <http://dx.doi.org/10.1534/genetics.108.088971>.
- 13 Shi Y, Di Giammartino DC, Taylor D, Sarkeshik A, Rice WJ, Yates JR 3rd, Frank J,
14 Manley JL. 2009. Molecular architecture of the human pre-mRNA 3'
15 processing complex. *Mol Cell* **33**: 365–376.
16 <http://dx.doi.org/10.1016/j.molcel.2008.12.028>.
- 17 Song Y, Wang J. 2023. ggcoverage: an R package to visualize and annotate
18 genome coverage for various NGS data. *BMC Bioinformatics* **24**.
19 <http://dx.doi.org/10.1186/s12859-023-05438-2>.
- 20 Sullivan KD, Steiniger M, Marzluff WF. 2009. A core complex of CPSF73, CPSF100,
21 and symplekin may form two different cleavage factors for processing of
22 poly(A) and histone mRNAs. *Mol Cell* **34**: 322–332.
23 <http://dx.doi.org/10.1016/j.molcel.2009.04.024>.

- 1 Takagaki Y, Manley JL. 1997. RNA recognition by the human polyadenylation factor
2 CstF. *Mol Cell Biol* **17**: 3907–3914. <http://dx.doi.org/10.1128/MCB.17.7.3907>.
- 3 Takyar S, Hickerson RP, Noller HF. 2005. mRNA helicase activity of the ribosome.
4 *Cell* **120**: 49–58. <http://dx.doi.org/10.1016/j.cell.2004.11.042>.
- 5 Venkataraman K, Brown KM, Gilmartin GM. 2005. Analysis of a noncanonical
6 poly(A) site reveals a tripartite mechanism for vertebrate poly(A) site
7 recognition. *Genes Dev* **19**: 1315–1327.
8 <http://dx.doi.org/10.1101/gad.1298605>.
- 9 Wickham H. 2016. *Ggplot2*. 2nd ed. Springer International Publishing, Basel,
10 Switzerland.
- 11 Wu L, Ueda T, Messing J. 1995. The formation of mRNA 3'-ends in plants. *Plant J* **8**:
12 323–329. <http://dx.doi.org/10.1046/j.1365-313x.1995.08030323.x>.
- 13 Xu F, Jerlström-Hultqvist J, Einarsson E, Ástvaldsson Á, Svärd SG, Andersson JO.
14 2014. The genome of *Spironucleus salmonicida* highlights a fish pathogen
15 adapted to fluctuating environments. *PLoS Genet* **10**: e1004053.
16 <http://dx.doi.org/10.1371/journal.pgen.1004053>.
- 17 Xu F, Jiménez-González A, Einarsson E, Ástvaldsson Á, Peirasmaki D, Eckmann L,
18 Andersson JO, Svärd SG, Jerlström-Hultqvist J. 2020. The compact genome
19 of *Giardia muris* reveals important steps in the evolution of intestinal
20 protozoan parasites. *Microb Genom* **6**.
21 <http://dx.doi.org/10.1099/mgen.0.000402>.

- 1 Yee J, Tang A, Lau W-L, Ritter H, Delport D, Page M, Adam RD, Müller M, Wu G.
2 2007. Core histone genes of *Giardia intestinalis*: genomic organization,
3 promoter structure, and expression. *BMC Mol Biol* **8**: 26.
4 <http://dx.doi.org/10.1186/1471-2199-8-26>.
- 5 Yusupova GZ, Yusupov MM, Cate JH, Noller HF. 2001. The path of messenger RNA
6 through the ribosome. *Cell* **106**: 233–241. [http://dx.doi.org/10.1016/s0092-](http://dx.doi.org/10.1016/s0092-8674(01)00435-4)
7 [8674\(01\)00435-4](http://dx.doi.org/10.1016/s0092-8674(01)00435-4).
- 8 Zajaczkowski P, Lee R, Fletcher-Lartey SM, Alexander K, Mahimbo A, Stark D, Ellis
9 JT. 2021. The controversies surrounding *Giardia intestinalis* assemblages A
10 and B. *Curr Res Parasitol Vector Borne Dis* **1**: 100055.
11 <http://dx.doi.org/10.1016/j.crpvbd.2021.100055>.
- 12 Zarudnaya MI, Kolomiets IM, Potyahaylo AL, Hovorun DM. 2003. Downstream
13 elements of mammalian pre-mRNA polyadenylation signals: primary,
14 secondary and higher-order structures. *Nucleic Acids Res* **31**: 1375–1386.
15 <http://dx.doi.org/10.1093/nar/gkg241>.
- 16 Zhao Z, Wu X, Ji G, Liang C, Li QQ. 2019. Genome-wide comparative analyses of
17 polyadenylation signals in eukaryotes suggest a possible origin of the
18 AAUAAA signal. *Int J Mol Sci* **20**: 958.
19 <http://dx.doi.org/10.3390/ijms20040958>.
- 20 R Core Team. 2021. R: A language and environment for statistical computing.
21 <https://www.R-project.org/>.
- 22

1 FIGURE LEGENDS

2 **Figure 1. Characterization of poly(A) signals in diverse protists.** **A.** Simplified
3 view of phylogenetic tree indicating evolutionary relationships between sequenced
4 protist species (highlighted red) and other organisms. Known poly(A) signal
5 sequences are shown next to the organism's name. **B.** MEME analysis of sequence
6 motifs enriched in the regions 40 nt upstream from the cleavage sites. Analysis was
7 performed for all sequenced protist species and human annotated 3'UTR
8 sequences. Sequences representing previously described unusual AGURAA *G.*
9 *lamblia* poly(A) signals were enriched within the *Giardia* genus. Sequences with
10 enriched AWUAAA-like motifs are highlighted green, AGURAA-like sequences in
11 blue, and others in black. **C.** Quantification of hexamers frequencies in annotated
12 3'UTRs in the regions 40 nt upstream from the cleavage sites, based on the
13 enrichment results from panel B. Color schemes the same as in B. **D.** Distribution of
14 the distances between poly(A) signals and cleavage sites. Poly(A) signals in the
15 *Giardia* genus are positioned closer to the cleavage sites than in other organisms.
16 Most common distance is shown above the plot. Color scheme the same as in B.

17

18 **Figure 2. 'Dual use' signals in protists.** **A.** Switching the poly(A) signal from
19 AAUAAA to WGURAA extends the number of possible "dual use" signals from 5 (left
20 panel) to 9 (right panel). Poly(A) signal sequence is colored green (AAUAAA) or blue
21 (WGURAA), stop codon sequence is highlighted gray. **B.** Stop codons of 8-20% of
22 protein coding genes in *Giardia* overlap with poly(A) signals. Almost all overlapping
23 poly(A) signals are WGURAA. In contrast, AAUAAA makes up the majority
24 overlapping poly(A) signals for *N. gruberi*. **C.** Length comparison between 3'UTRs
25 with and without "dual use" signals for *G. lamblia* A and *N. gruberi*. 3'UTRs

1 containing these signals are significantly shorter. Median length is shown below the
2 plot.

3

4 **Figure 3. Absence of metazoan-like auxiliary elements in protists. A.**

5 Comparison of the UGUA vs its shuffled versions 20-50 nt upstream from the
6 cleavage site in human vs protists. Each data point represents total number of either
7 UGUA tetramers (red) or its shuffled versions (gray) and bar shows mean count
8 value for all compared combinations. We observed human-like enrichment only in *N.*
9 *gruberi*. Statistical analysis is presented in Supplemental Fig S2. **B.** Enrichment of
10 GU and U-rich elements calculated in the region 40 nt downstream from the
11 cleavage site in comparison to the region 40 nt upstream to the cleavage site. Each
12 point on box and whisker plot represents \log_2 ratio from downstream vs upstream
13 quantity of one of the CUGCCU, CUGGGG, CUGUGU, GUCUGU, GUGUCU,
14 GUGUGU, UGUCUC, UGUCUG, UGUGUG, UGUUUU, UUAUUU, UUUCUU,
15 UUUUUU sequences. Statistical significance was calculated in comparison to human
16 enrichment using a one-sided Wilcoxon rank sum test.

17

18 **Figure 4. Importance of 3'adjacent nucleotide of the WGURAA hexamer for**
19 **distinguishing poly(A) signals from open reading frame hexamers. A.**

20 Occurrence of AGU AAA di-codon in comparison to other Ser-Lys encoding di-
21 codons in coding sequences within Metamonada clade. Di-codons occurrence was
22 calculated only in frame. Each point represents one di-codon combination and bar
23 shows mean count value for all compared combinations. AGU AAA was shown to be
24 depleted in the *Giardia* genus, but not in other species. Data for AAA AGU, UGU
25 AAA and AGU GAA di-codons are shown in Supplemental Fig S3A. **B.** WGURAA

1 hexamer occurrence in the coding sequences of *Giardia* species. Occurrence was
2 calculated independent of the reading frame. **C.** Schematic representation of the
3 machine learning approach to distinguish WGURAA sequences in 3'UTRs vs coding
4 sequences. WGURAA sequences were extracted from coding sequences and
5 3'UTRs and together with 37 flanking nucleotides from both sides (80-mer) put into a
6 gapped *k*-mer support vector machine classifier, which performed sequence
7 classification and *k*-mers scoring. **D.** Variance explained by the linear model applied
8 to WGURAA-containing *k*-mers scores from gkmSVM classifier by the full model,
9 upstream nucleotide, poly(A) signal and downstream nucleotide in *G. lamblia*.
10 Explained variance was measured as adjusted R^2 value. Data for *G. lamblia* B are
11 shown in Supplemental Fig S3C. **E.** Beta-coefficient values from the linear model
12 applied to WGURAA-containing *k*-mers scores from gkmSVM classifier,
13 corresponding to the upstream nucleotide, poly(A) signal and downstream nucleotide
14 in *G. lamblia*. Data for *G. lamblia* B are shown in Supplemental Fig S3B. **F.** Example
15 of *G. lamblia* A gene GL50803_1890 where a hexamer in the coding sequence was
16 misclassified by gkmSVM. Premature cleavage after the hexamer inside the coding
17 sequences indicated as coverage drop was observed. Similar example from *G.*
18 *lamblia* B is shown in Supplemental Fig S3D.

19

20 **Figure 5. Functional poly(A) signals in majority of WGURAA hexamers in**
21 **coding sequences in *G. muris* .** **A.** Variance explained by the linear model applied
22 to WGURAA-containing *k*-mers scores from gkmSVM classifier by the full model,
23 upstream nucleotide, poly(A) signal and downstream nucleotide in *G. muris*.
24 Explained variance was measured as adjusted R^2 value. Green bar represents a full
25 linear model without interactions, blue - with interactions, gray - components of the

1 model with interactions. Beta-coefficient values are found in Supplemental Fig S4A.
2 **B.** Influence of upstream and downstream nucleotide interaction to AGUGAA
3 hexamers classification. Upstream nucleotide is color-coded at the top of each
4 barplot. Data for AGUAAA and UGUAAA hexamers are shown in Supplemental Fig
5 S4B and C, respectively. **C.** Quantification of premature cleavage events in *G. muris*
6 coding sequences by hexamer identity. Among 199 WGURAA hexamers in *G. muris*
7 coding sequences, 106 had sufficient read support and 91 were cleaved.

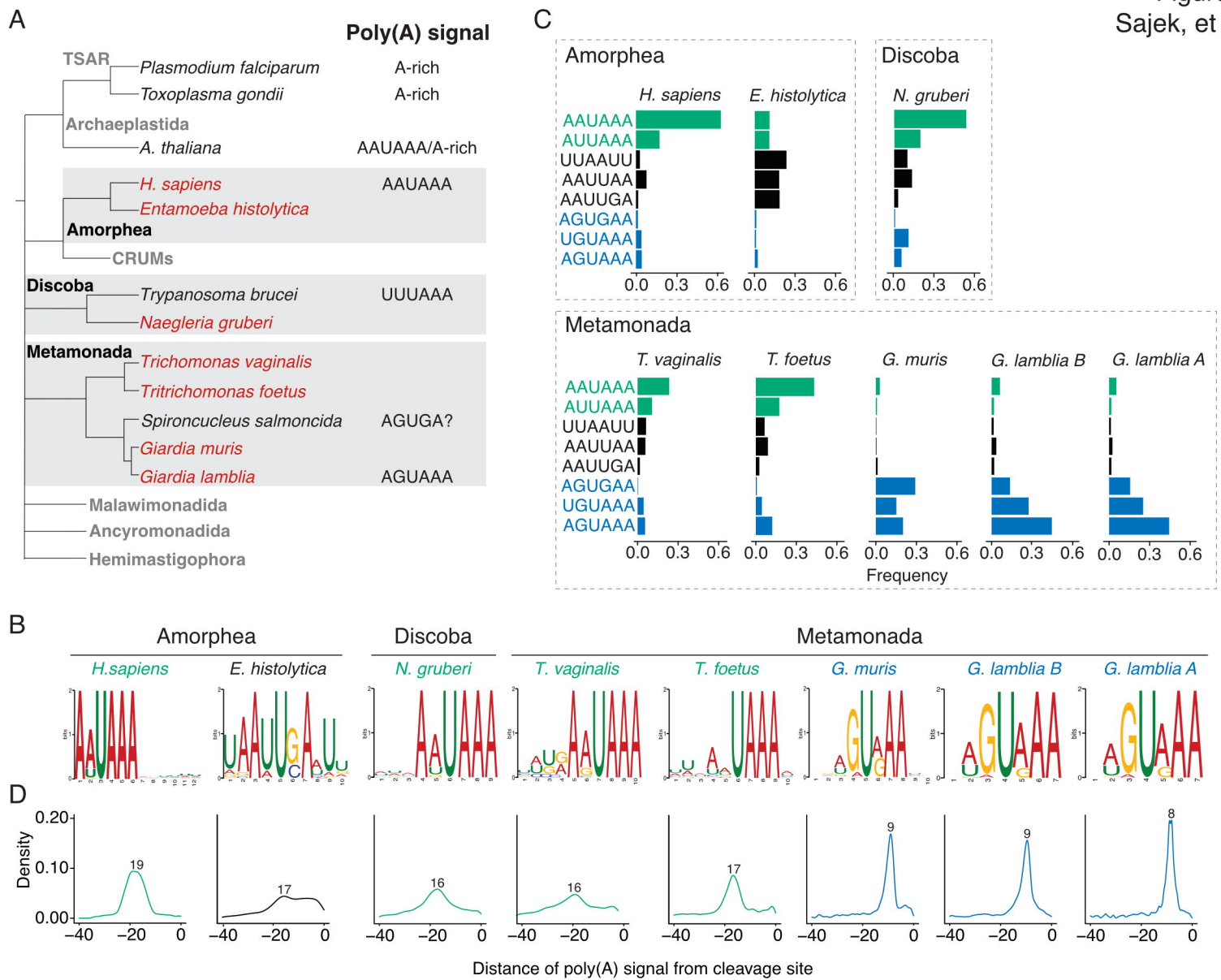
8

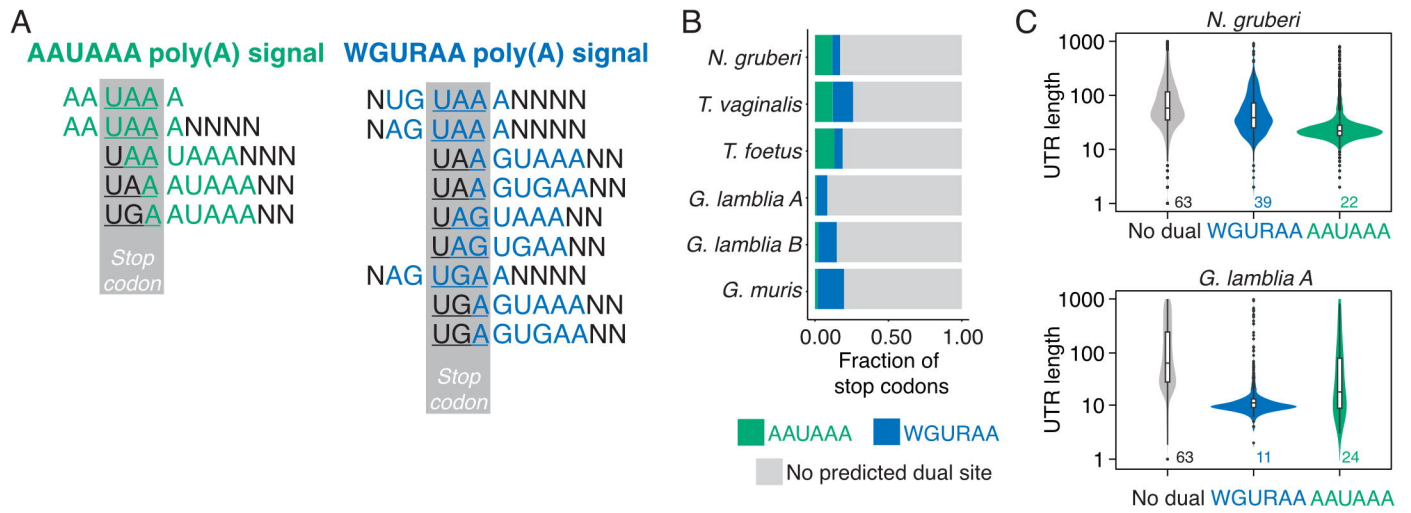
9 **Figure 6. Flanking nucleotides play a role in poly(A) signal recognition in other**
10 **Metamonada, but not *N. gruberi*.** **A.** Variance explained by the linear model
11 applied to AWUAAA-containing *k*-mers scores from gkmSVM classifier. In the
12 outgroup organism *N. gruberi*, signal sequence is the main determinant, whereas in
13 *T. vaginalis* and *T. foetus* from Metamonada upstream nucleotides have substantial
14 contributions. **B.** Beta-coefficient values from the linear model form A. **C.** AAUAAA
15 and AUUAAA hexamer occurrence in the coding sequences of *N. gruberi*.
16 Occurrence was calculated independent of the reading frame. **D, E.** Auxiliary
17 elements help specify genuine AAUAAA poly(A) signals in *N. gruberi*, putative
18 upstream (D) and downstream (E) motifs and their quantifications around AAUAAA
19 poly(A) signals.

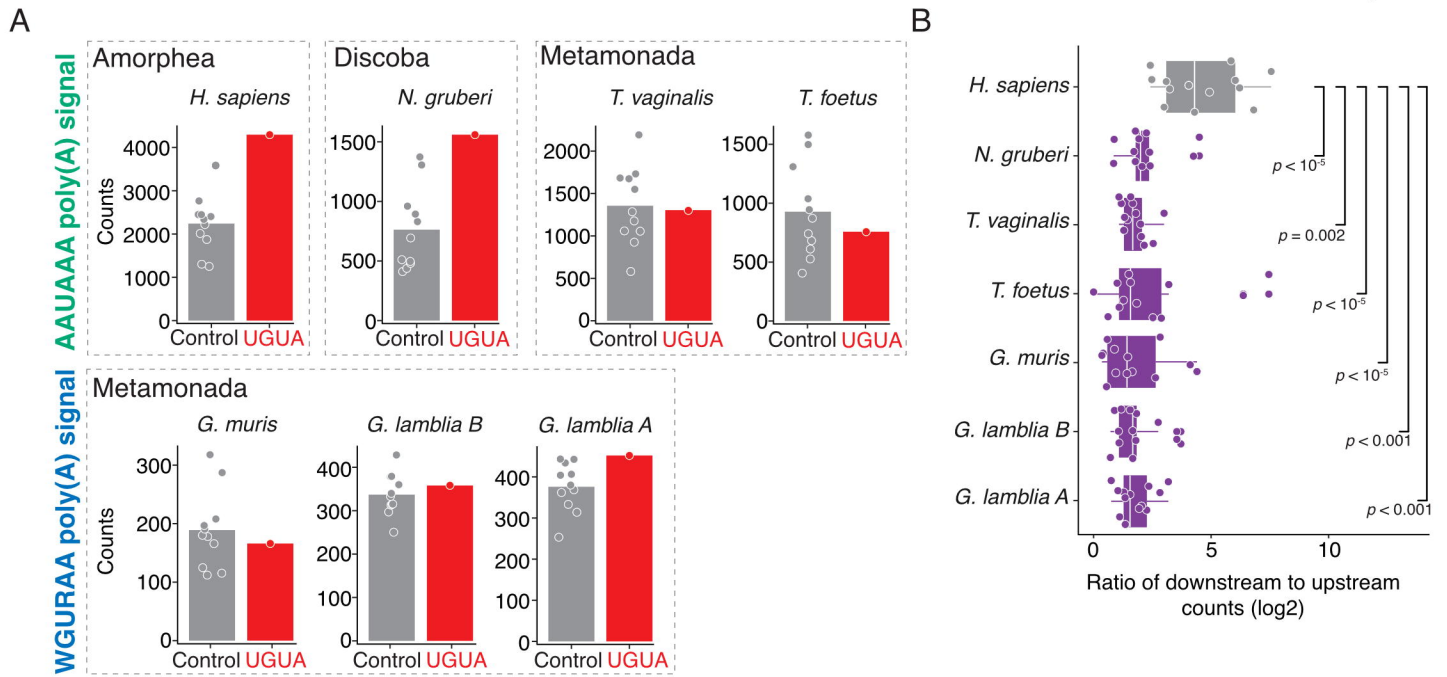
20

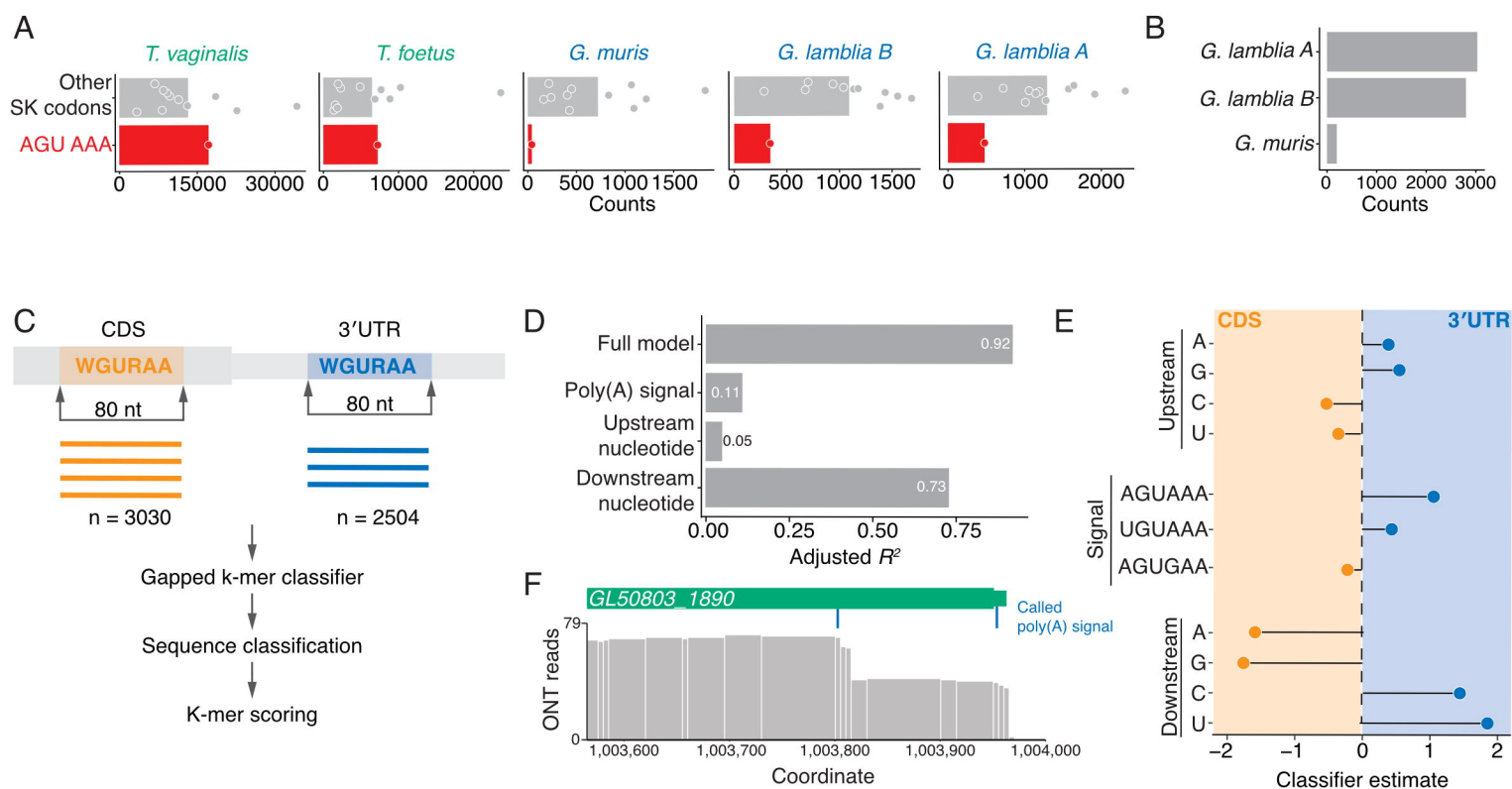
21 **Figure 7. Putative evolutionary events that shaped the evolutionary dynamics**
22 **of poly(A) signal recognition.** Multiple evolutionary events have shaped poly(A)
23 signal recognition throughout eukaryotes, including involvement of flanking
24 nucleotides, loss of importance for auxiliary elements, and change in the sequence

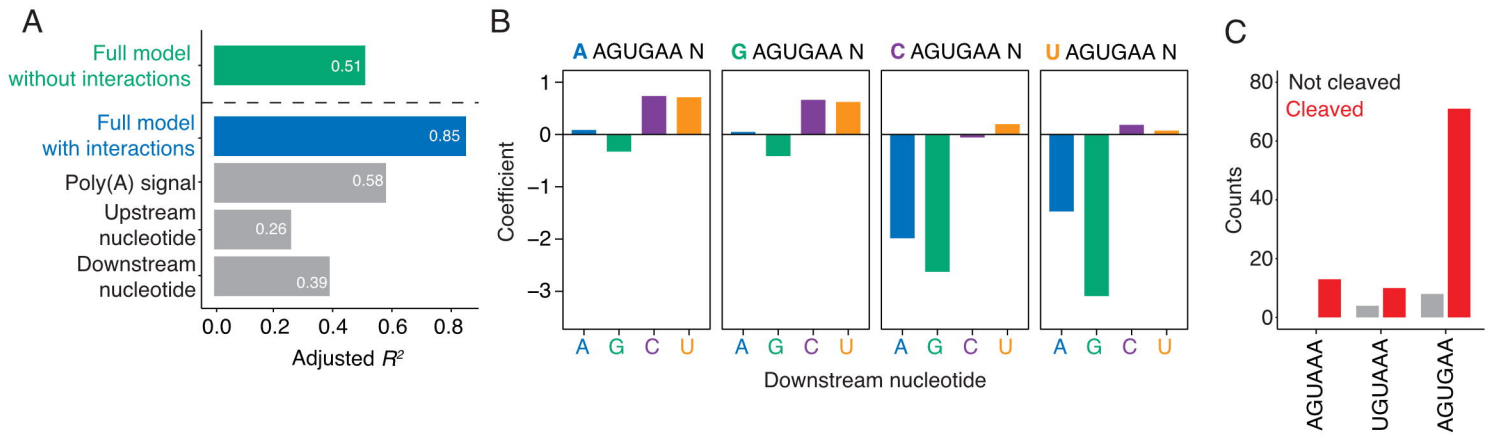
- 1 of the poly(A) signal itself. To identify more precise placement of the hexamer
- 2 change, additional annotations from more organisms are required.

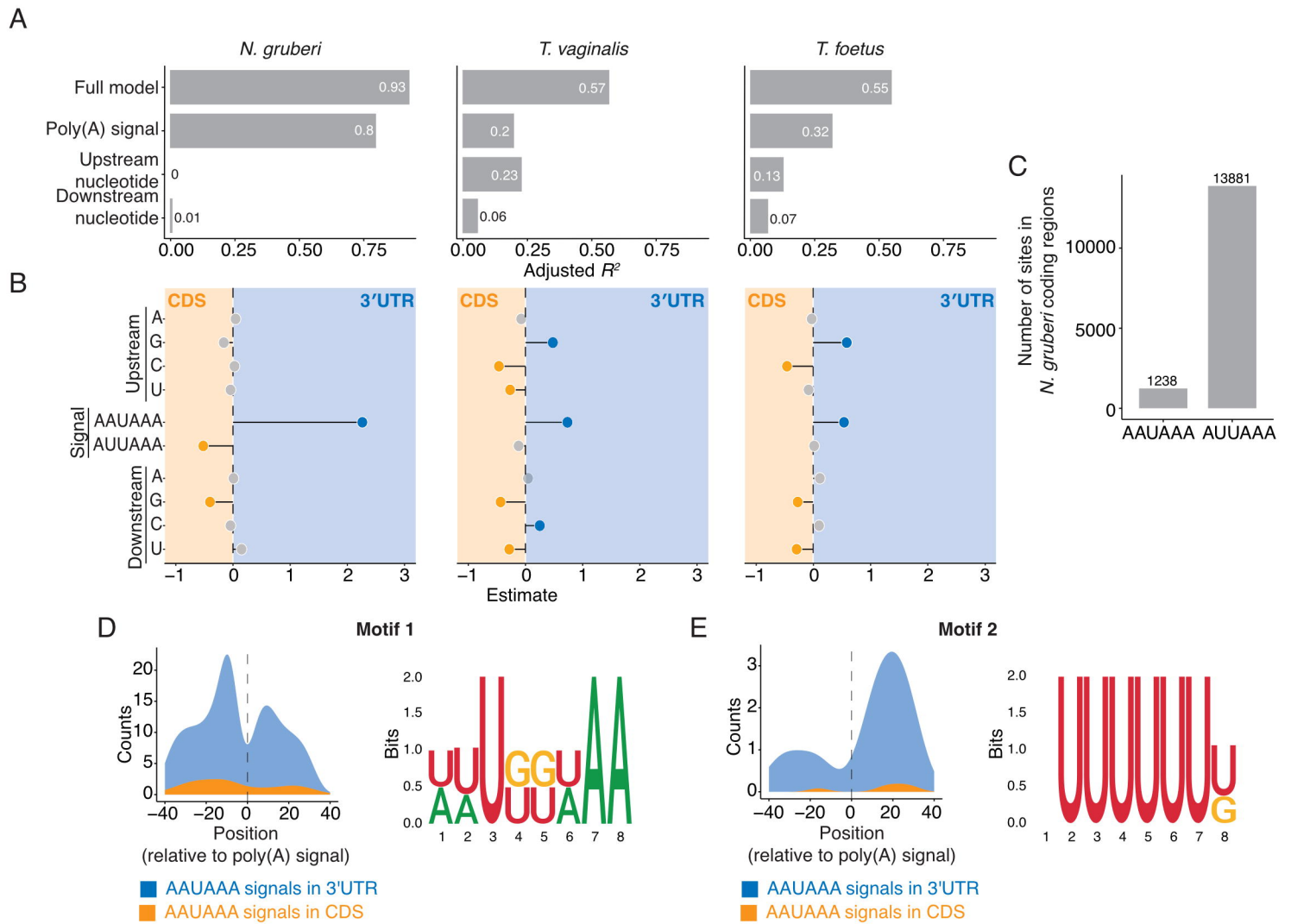


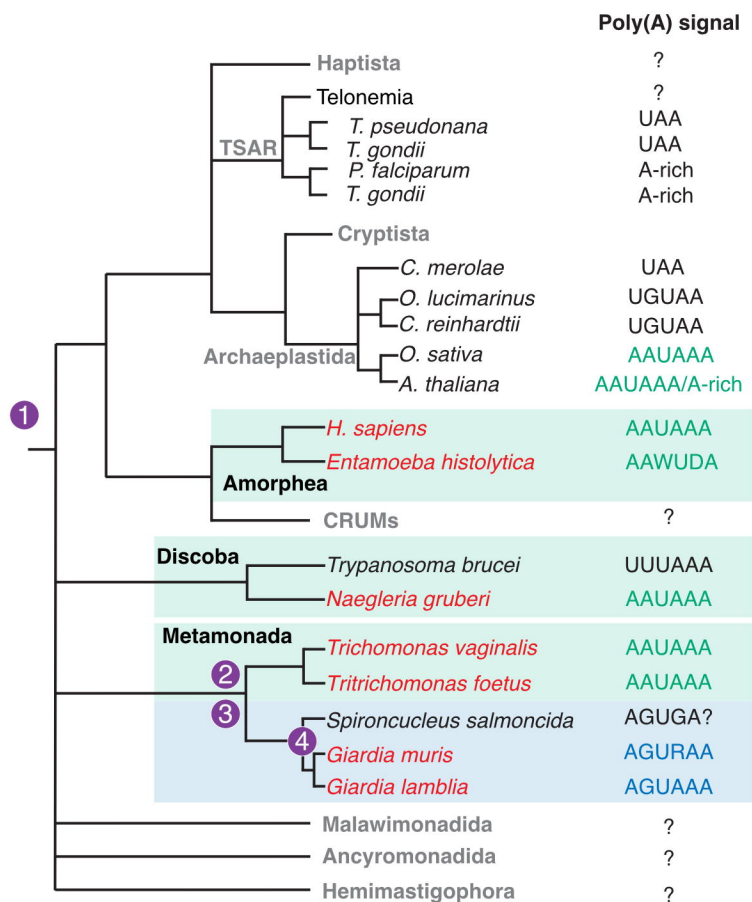












- 1 Ancestral poly(A) signal: AAUAA + auxiliary motifs (?)
- 2 Loss of importance of auxiliary motifs
- 3 Use of flanking nucleotides
- 4 Change in hexamer to AGURAA and genome adaptation