1

# On the origin and evolution of RNA editing in metazoans

2  Qiye Li[1,2,19*], Pei Zhang[1,3,19], Ji Li[2,4], Hao Yu[1], Xiaoyu Zhan[1], Yuanzhen Zhu[1,5], Qunfei Guo[1,6],

3  Huishuang Tan[1,7], Nina Lundholm[8], Lydia Garcia[8], Michael D. Martin[9,10], Meritxell Antó

4  Subirats[11], Yi-Hsien Su[12], Iñaki Ruiz-Trillo[11,13,14], Mark Q. Martindale[15], Jr-Kai Yu[12,16], M.

5  Thomas P. Gilbert[9,17], Guojie Zhang[1,2,3,18*]

6

7  [1] BGI-Shenzhen, Shenzhen 518083, China

8  [2] State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology,

9  Chinese Academy of Sciences, Kunming 650223, China

10  [3] Section for Ecology and Evolution, Department of Biology, University of Copenhagen, DK-

11  2100 Copenhagen, Denmark

12  [4] China National Genebank, BGI-Shenzhen, Shenzhen 518120, China

13  [5] School of Basic Medicine, Qingdao University, Qingdao 266071, China

14  [6] College of Life Science and Technology, Huazhong University of Science and Technology,

15  Wuhan 430074, China

16  [7] Center for Informational Biology, University of Electronic Science and Technology of China,

17  Chengdu 611731, China

18  [8] Natural History Museum of Denmark, University of Copenhagen, Copenhagen 1350,

19  Denmark

20  [9] Department of Natural History, NTNU University Museum, Norwegian University of Science

21  and Technology (NTNU), NO-7491 Trondheim, Norway

22  [10] Center for Theoretical Evolutionary Genomics, Dept. of Integrative Biology, University of

23  California Berkeley, Berkeley, California 94720, USA

24  [11] Institute of Evolutionary Biology, UPF-CSIC Barcelona, 08003 Barcelona, Spain

25  [12] Institute of Cellular and Organismic Biology, Academia Sinica, 11529 Taipei, Taiwan

26  [13] ICREA, Passeig Lluís Companys 23, 08010 Barcelona, Catalonia, Spain

27  [14] Departament de Genètica, Microbiologia i Estadística, Facultat de Bilogia, Universitat de

28  Barcelona (UB), Barcelona 08028, Spain

29  [15] The Whitney Laboratory for Marine Bioscience, University of Florida, St. Augustine, FL

30  32080, USA

31  [16] Marine Research Station, Institute of Cellular and Organismic Biology, Academia Sinica,

32  26242 Yilan, Taiwan

33    [17] Section for Evolutionary Genomics, The GLOBE Institute, University of Copenhagen,

34    Copenhagen 1352, Denmark

35    [18] Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences,

36    650223, Kunming, China

37    [19] These authors contributed equally

38    * Correspondence: liqiye@genomics.cn (Q.L.) and guojie.zhang@bio.ku.dk (G.Z.).

39

40

## Abstract

42    Extensive adenosine-to-inosine (A-to-I) editing of nuclear-transcribed RNAs is the hallmark

43    of metazoan transcriptional regulation, and is fundamental to numerous biochemical processes.

44    Here we explore the origin and evolution of this regulatory innovation, by quantifying its

45    prevalence in 22 species that represent all major transitions in metazoan evolution. We provide

46    substantial evidence that extensive RNA editing emerged in the common ancestor of extant

47    metazoans. We find the frequency of RNA editing varies across taxa in a manner independent

48    of metazoan complexity. Nevertheless, cis-acting features that guide A-to-I editing are under

49    strong constraint across all metazoans. RNA editing seems to preserve an ancient mechanism

50    for suppressing the more recently evolved repetitive elements, and is generally nonadaptive in

51    protein-coding regions across metazoans, except for *Drosophila* and cephalopods. Interestingly,

52    RNA editing preferentially target genes involved in neurotransmission, cellular

53    communication and cytoskeleton, and recodes identical amino acid positions in several

54    conserved genes across diverse taxa, emphasizing broad roles of RNA editing in cellular

55    functions during metazoan evolution that have been previously underappreciated.

## Introduction

The central dogma of molecular biology emphasizes how genetic information passes faithfully from DNA, to RNA, to proteins. However, this dogma has been challenged by the phenomenon of RNA editing — a post/co-transcriptional-processing mechanism that can alter RNA sequences by insertion, deletion or substitution of specific nucleotides, thus producing transcripts that are not directly encoded in the genome [1]. In metazoans, the most prevalent form of RNA editing is the deamination of adenosine (A) to inosine (I), which is catalyzed by a family of adenosine deaminases acting on RNA (ADARs) [2,3]. As inosine is recognized *in vivo* as guanosine by ribosomes and other molecular machinery, RNA editing can affect almost all aspects of cellular RNA functions, from changing mRNA coding potential by altering codons or splicing patterns, to regulating the cellular fate of mRNA by editing its microRNA (miRNA) binding sites [4-6]. RNA editing is particularly pervasive in neural systems, where it has been shown to modulate neural development processes [7,8], neural network plasticity [9,10] and organismal adaptation to environmental changes [11-13]. Defects in RNA editing machinery have been linked to a variety of neurological diseases, autoimmune disorders and cancers [14-18].

Although recent high-throughput sequencing-based analyses have identified a surprisingly large number of RNA-editing sites in different metazoans, including humans [19-23], mice [24,25], *Caenorhabditis elegans* [26], fruit flies [27-30], ants [31], bumblebees [32] and cephalopods [33,34], conclusions about the evolutionary patterns of this phenomenon are inconsistent. For example, while almost all human RNA-editing sites occur in Alu repeat elements [20,21], editing in *Drosophila* primarily targets exonic (particularly coding) regions [27,28]. Additionally, while recoding RNA editing, which leads to nonsynonymous substitutions in protein-coding sequences, is abundant and affects around half of the protein-coding genes in coleoid cephalopods [33,34], it is relatively rare in mammals and insects [21,24,28,31,32]. Furthermore, while recoding editing in humans is generally nonadaptive [35], it is typically adaptive in *Drosophila* and cephalopods [28,34]. More importantly, although the *ADAR* gene family is considered to have originated in the common ancestor of extant metazoans [36], the functional activity of ADARs in catalyzing RNA editing in most metazoan lineages actually remains unknown, especially in those earliest branching lineages like Ctenophora and Porifera.

In summary, many fundamental questions about the nature of metazoan RNA editing remain to be investigated, including: When did RNA editing emerge during metazoan evolution? Are there conserved sequence features that underly RNA editing in all metazoans? What genes and genomic elements are the primary targets of metazoan RNA editing? How does the prevalence

89  of recoding editing vary by lineage, and does it generally provide adaptive amino acid changes
90  in metazoans? Addressing these questions requires the characterization of RNA editomes
91  across the diversity of metazoans and their closest unicellular relatives, thus we systematically
92  investigated the prevalence and characteristics of RNA editing in 22 lineages that encompass
93  the key transitions in metazoan evolution.

94

95  **Results**

96  **Profiling the RNA editomes across the phylogeny of metazoans.**

97  We performed both DNA-seq and strand-specific RNA-seq for 18 species, including 14
98  metazoans and 4 unicellular eukaryotes closely related to animals. 14 out of these 18 species
99  have not been subjected to transcriptome-wide RNA editing investigation previously (Fig. 1a).
100  For each species, two to three (mostly three) biological replicates were sequenced, yielding
101  3.27 Tbp (tera base pairs) sequencing data in total, with the average DNA and RNA coverage
102  achieving 75X (ranging 15-345X) and 45X (ranging 10-162X) respectively for each biological
103  replicate after alignment (Supplementary Table 1). Together with published sequencing data
104  from *C. elegans* [26], ant [31], octopus [37] and human [22] (Supplementary Table 1), we were able to
105  profile and compare the RNA editomes of 22 species, which represent nearly all the major
106  phyla of extant metazoans, including the earliest-branching lineages Ctenophora, Porifera and
107  Placozoa, as well as their closest unicellular relatives Choanoflagellatea, Filasterea and
108  Ichthyosporea (Fig. 1a). These data thus provide the first opportunity to phylogenetically
109  investigate the prevalence of RNA editing within Holozoa, the clade that includes animals and
110  their closest single-celled relatives [38].

111  Given that some RNA-editing sites tend to appear in clusters, while others remain isolated, we
112  adopted two complementary methods to identify the RNA editomes for each species. Briefly,
113  we first employed RES-Scanner [39] to identify RNA-editing sites by comparing the matching
114  DNA- and RNA-seq data from the same sample. This method has high accuracy when
115  searching for RNA-editing sites that are isolated or not heavily clustered. We next performed
116  hyper-editing detection [40], using the RNA reads that failed to align by RES-Scanner, in order
117  to capture the hyper-edited reads and the clusters of editing sites they harbored. The results of
118  RES-Scanner and hyper-editing detection were combined to yield the RNA editome of each
119  sample (Supplementary Table 2). We have compiled the whole pipeline as an easy-to-use
120  software package named RES-Scanner2, which is applicable to transcriptome-wide

121     identification of RNA-editing sites in any species with matching DNA- and RNA-seq data (see

122     Methods for details).

123

124     **Extensive RNA editing emerged in the last common ancestor of modern metazoans**

125     **accompanied by the origin of *ADARs*.**

126     We detected very few putative RNA-editing sites (ranging 23-304) in the four unicellular

127     holozoans (Fig. 1b and Supplementary Table 2). No dominant type of nucleotide substitution

128     was observed (Fig. 1c), and the frequency of each type of nucleotide substitution was close to

129     that of genetic polymorphism (Supplementary Fig. 1a), implying that RNA-editing sites

130     detected in these species likely represent noise. In contrast thousands, to hundreds of thousands,

131     of RNA-editing sites were identified in almost all the sampled metazoans, including the

132     earliest-branching Ctenophora and Porifera, with the vast majority (>90%) consisting of A-to-

133     G substitutions (i.e. A-to-I editing; Fig. 1b,c). The only exception was *Trichoplax adhaerens*,

134     a morphology-simplified metazoan belonging to Placozoa (a sister group to Cnidaria and

135     Bilateria) [41]. Concordantly, we confirmed the existence of *ADAR-like* genes in all the sampled

136     species except *T. adhaerens* and the unicellular taxa (Fig. 1a and Supplementary Table 3; See

137     Methods). Our results thus provide direct evidence that extensive editing of nuclear-transcribed

138     RNAs first emerged in the last common ancestor of modern metazoans, alongside the

139     appearance of ADAR-mediated A-to-I editing, which is pervasively preserved in most extant

140     animal lineages. We also highlight that our detection methods do not depend on any prior

141     knowledge about the dominate type of RNA editing in any species studied (see Methods), thus,

142     our results also imply that RNA editing in any manner other than A-to-I, is either extremely

143     rare, or non-existent, in the animal kingdom.

144     We next calculated the occurrence rate of RNA editing per genome by counting the number of

145     RNA-editing sites per million transcribed genomic sites (i.e. sites with RNA depth $\geq$ 2X). Our

146     results indicate that the octopus exhibits the highest, and *Drosophila* the lowest, number and

147     occurrence rate among the sampled taxa that have the RNA-editing machinery. Surprisingly,

148     the occurrence rates in the ctenophore *Mnemiopsis leidyi* and sponge *Amphimedon*

149     *queenslandica* are higher than that of all sampled cnidarians and many bilaterians (Fig. 1b),

150     while humans are among the species with lowest rates (Fig. 1b). Similar patterns were obtained

151     if we weighted each editing site with its editing level, or if we only considered A-to-I editing

152     (Supplementary Fig. 1b-e). These results suggest that the global level of RNA editing has

153    changed considerably during the diversification of metazoan, and does not increase directly

154    alongside organismal complexity.

155

156    **The A-to-I editing associated sequence features are under strong constraint in metazoans.**

157    Consistent with the double-stranded RNA (dsRNA) binding property of ADAR enzymes [2,6],

158    we observed that A-to-I editing sites in all the sampled metazoans with *ADARs* were

159    preferentially located in potential dsRNA regions that could form by intramolecular folding of

160    pre-mRNA. Specifically, we found on average that 37% (ranging 6% to 86%) of the editing

161    sites target regions that show a reverse-complement alignment within their upstream or

162    downstream sequences, which is significantly higher than the expected levels of ~1%

163    calculated from randomly selected transcribed adenosines (Fig. 2a; See methods). These results

164    confirm that a stable dsRNA structure is critical for establishing A-to-I editing *in vivo* across

165    metazoans [42], and further reveal that intramolecular folding of pre-mRNA is a major way to

166    form dsRNA substrates for A-to-I editing in most species.

167    Intermolecular hybridization of sense and antisense transcripts is another potential mechanism

168    to form dsRNA [43], but its role in inducing A-to-I editing is thought to be negligible in mammals

169    [44]. Taking advantage of the strand information provided by strand-specific RNA-seq, we found

170    that the proportions of editing sites that were located in regions containing transcription signals

171    on both strands (mean 17%, ranging 3% to 64%) were significantly higher than the expected

172    levels (mean 8%, ranging 3% to 32%) in 8 out of the 17 metazoans with *ADARs* (Fig. 2b; See

173    methods). In particular, while for most species there are generally many more editing sites

174    found in potential dsRNA regions formed by intramolecular folding, the ctenophore *M. leidyi*

175    and sea squirt *Ciona savignyi* showed a reverse tendency, with higher proportions of editing

176    sites found in regions with transcription signals in both strands (Fig. 2c). This implies that

177    intermolecular hybridization of sense and antisense transcripts likely represents an important

178    means for forming dsRNA substrates for A-to-I editing, in at least some taxa. This conclusion

179    is further supported by the significantly higher-than-expected proportion of A-to-I editing sites

180    locating in regions targeted by RNA editing on both strands in many species (Supplementary

181    Fig. 2a,c).

182    With regards to the genomic distribution of A-to-I editing, we found on average 81% (ranging

183    41% to 97%) of the metazoan editing sites were clustered, which is significantly higher than

184    the expected levels of less than l% (Fig. 2d). The median distances between any two adjacent

185    editing sites were mostly around 5 nt (ranging 4 to 81 nt; Supplementary Table 4). Furthermore,

186  editing levels of the clustered editing sites were generally higher than those of isolated sites,
187  except in *Hydra vulgaris*, *Drosophila*, *C. savignyi* and humans (Supplementary Fig. 2b). A
188  typical metazoan editing cluster (i.e. a region with ≥ 3 A-to-I editing sites and the distance of
189  two adjacent sites ≤ 30 nt) was ~50 nt in length, and harbored 9 A-to-I editing sites, and we
190  estimated that up to 52% of the adenosines within a cluster were targeted by RNA editing
191  (Supplementary Table 4). Taken together, our results indicate that the majority of metazoan A-
192  to-I editing sites are organized in dense clusters, within RNA regions that can form stable
193  dsRNA structures.

194  Since ADARs recognize dsRNA when exerting A-to-I editing, we then asked what primary
195  sequence motifs guide ADARs to preferentially edit certain adenosines rather than others in
196  their dsRNA substrates. By comparing the surrounding sequence context of edited adenosine
197  sites to neighboring unedited adenosine sites (i.e. unedited adenosines with RNA depth ≥ 2X
198  and within ± 50 nt of the edited adenosines), we observed clear and conserved nucleotide
199  preferences for the positions that are directly 5' and 3' adjacent to the edited adenosines (i.e.
200  the -1 and +1 positions). Specifically, the 5' adjacent position strongly favored uridine and
201  adenosine, but disfavored guanosine across all metazoans, and to a lesser extent, cytosine was
202  also disfavored (Fig. 2e and Supplementary Fig. 3). In contrast, the nucleotide preference for
203  the 3' adjacent position is relatively weaker, and less conserved, with guanosine being favored
204  and uridine being disfavored in most species (Fig. 2e and Supplementary Fig. 3). This implies
205  that the 5' adjacent position has the most influential and a conserved role on determining
206  whether an adenosine will be edited. Concordantly, we found the nucleotide triplets of U<u>A</u>G
207  and A<u>A</u>G, with the edited adenosines in the center, to be the most likely edited triplets, while
208  G<u>A</u>U was the least likely edited triplet in metazoans (Fig. 2f).

209  Interestingly, *C. elegans* also displayed a strong sequence preference for the 5' second nearest
210  (-2) position of the edited adenosines that is not observed in other metazoans, with uridine
211  being strongly favored and adenosine being strongly disfavored (Fig. 2e and Supplementary
212  Fig. 3). We speculate that this *C. elegans* specific motif adjustment is associated with the high
213  sequence divergence of the *C. elegans* ADARs against other metazoan ADARs, as
214  phylogenetic analyses separate both the *C. elegans* ADR-1 and ADR-2 from other metazoans
215  (Supplementary Fig. 2d), and both *C. elegans* ADR-1 and ADR-2 show high nonsynonymous
216  substitution rates ($d_N$) against ADARs from other metazoans (Supplementary Fig. 2e).

217

218  **Evolutionarily young repetitive elements are the primary targets of metazoan RNA**
219  **editing.**

220  In all metazoans sampled except for the two fruit flies and sea squirt, repetitive elements
221  including transposons and tandem repeats were the major targets of A-to-I editing, and
222  harbored on average 83% (ranging 73% to 95%) of the editing sites (Fig. 3a). This suggests
223  that extensive editing of repeat-containing transcripts is the ancestral and predominant feature
224  for metazoan RNA editing, probably because these regions are more likely to hybridize with
225  nearby oppositely oriented repeats, creating the dsRNA structures suitable for ADARs binding
226  (Supplementary Fig. 4c). It is noteworthy that, even for those sites on pre-mRNA (i.e. exon +
227  intron) of protein-coding genes, especially those outside coding regions, the majority (>70%)
228  were also associated with repetitive elements (Supplementary Fig. 4d). This implies that most
229  editable sites on protein-coding genes were actually introduced by the invasion of repetitive
230  elements into gene regions.

231  Given that the total lengths of the different genomic elements vary greatly within each genome,
232  we next calculated the A-to-I editing density for each type of genomic element, by counting
233  the number of editing sites per million of transcribed adenosine sites (i.e. RNA depth $\geq$ 2X).
234  After this normalization, we observed that the editing densities of protein-coding gene-related
235  elements (i.e. 5'-UTR, CDS, intron and 3'-UTR) were close to the whole genome average level
236  in all metazoans (Fig. 3b). However, editing densities generally increased from 5' to 3' of
237  mRNA transcripts, with 3' UTR being relatively more favored by A-to-I editing than 5'-UTR
238  and CDS (Fig. 3c), consistent with previous observation in *Drosophila* [27,28]. In contrast, the
239  editing densities of repetitive elements, especially DNA transposons, short interspersed nuclear
240  elements (SINEs), long interspersed nuclear elements (LINEs) or Helitrons depending on
241  species, were significantly higher than the whole genome average. This further supports the
242  hypothesis that repetitive elements are the most favorable targets of A-to-I editing in metazoans.
243  Similar results were obtained even if we weighted each editing site with its editing level
244  (Supplementary Fig. 4e,f). Moreover, we observed negative correlations between the
245  divergence rates and the editing densities of repetitive elements in most species (Fig. 3d and
246  Supplementary Fig. 4g), suggesting that A-to-I editing preferentially targets evolutionarily
247  young repetitive elements that likely only relatively recently invaded the genome of each
248  species. Given that hyper-edited dsRNAs can be degraded by endonuclease V [45], RNA editing
249  may therefore serve as a guardian mechanism to avoid the overactivation of repetitive elements
250  in metazoans.

251

252 **Recoding RNA editing is rare and generally nonadaptive in metazoans.**

253 The phenomenon of recoding editing has gained considerable research interest, as it can result

254 in nonsynonymous substitutions in protein-coding sequences, and thus has the potential to

255 increase proteome diversity by introducing novel protein isoforms [3,6]. We observed that the

256 number of recoding sites varied greatly across species, with the octopus having an

257 overwhelming higher number (29,464) than all other species (median 850). In general, the

258 proportion of recoding sites among all A-to-I editing sites was low, ranging from less than 1%

259 to 7% in the majority of metazoans, with only 1% to 5% of all expressed protein-coding genes

260 being recoded. However, the proportions of recoding sites in the fruit flies and the sea squirt

261 were prominently high, reaching 33% (711/2,149), 30% (641/2,165) and 14% (850/6,254) in

262 *D. melanogaster*, *D. simulans* and *C. savignyi*, respectively (Fig. 4a). This may possibly be due

263 to the reduced proportion of editing sites in repetitive elements for these species (see Fig. 3a

264 and Supplementary Fig. 4c,d).

265 We next examined the effect of natural selection on recoding sites. It has been previously

266 reported that nonsynonymous editing is generally adaptive in fruit flies and cephalopods [28-30,34].

267 If this is so, one would expect that, in relation to synonymous editing, which is expected to be

268 neutral, the frequency of nonsynonymous editing ($f_n$) calculated as the number of A-to-I

269 editing sites causing nonsynonymous changes against all potential nonsynonymous adenosine

270 sites if A is replaced with G, is higher than that of synonymous editing ($f_s$) (see Methods).

271 When considering all recoding sites together, we observed that the frequencies of

272 nonsynonymous editing were either close to, or significantly lower than, synonymous editing

273 in all species (Fig. 4b). This therefore argues against the adaptive hypothesis, and suggests that

274 the recoding editing events observed in coding regions of most metazoans are generally neutral

275 or deleterious, consistent with previous reports in humans [35]. Consistently, editing levels of A-

276 to-I sites in coding regions were generally lower than the genome average and other types of

277 genomic elements (Fig. 4c), implying that editing of coding regions tends to be suppressed.

278 However, when we divided the recoding sites of each species into lowly (editing level < 0.2)

279 and highly (editing level ≥ 0.2) edited groups, we found that the frequencies of nonsynonymous

280 editing in fruit flies and octopus became significantly higher than synonymous editing in the

281 highly edited group (Fig. 4b). This demonstrates a relatively larger portion of adaptive recoding

282 sites exists in these two lineages than in other metazoans.

283    If recoding editing is generally nonadaptive, one would also expect that nonsynonymous
284    editing is depleted from evolutionarily conserved genes which are less tolerant to mutations.
285    We thus divided the genes of each species into three groups according to the degree of
286    evolutionary conservation (see Methods). Group I and II comprise genes that have orthologs
287    in closely-related species, but with relatively low and high $d_N/d_S$ ratios, representing the most
288    and moderately conserved groups, respectively. Group III comprises all the remaining genes,
289    that cannot find orthologs, and represents the least conserved group. As expected, the genes
290    subjected to recoding editing were generally enriched in the least conserved groups in most
291    metazoans (Fig. 4d), suggesting that recoding editing tends to be purged from the
292    evolutionarily conserved genes in most metazoans. Nevertheless, an inverse tendency can be
293    observed in the fruit flies and octopus, probably due to the relatively larger portions of adaptive
294    recoding sites in these species. This also implies that adaptive recoding editing more likely
295    emerged in the evolutionarily conserved genes, which benefit from increasing protein diversity
296    without introducing DNA mutation in these genes.

297

298    **RNA editing preferentially affects cellular communication and cytoskeleton related genes**
299    To uncover the functional preference of genes targeted by A-to-I editing in metazoans, we
300    conducted gene ontology (GO) based functional enrichment analysis for the RNA-recoded
301    genes (i.e. genes with at least one recoding site of which the average editing level across
302    samples > 0.1 or shared by at least two samples) in each species. Consistent with previous
303    observations in mammals [7,9], insects [28,29,31] and cephalopods [33,34], we found that
304    neurotransmission-related functions such as ion transmembrane transport, synaptic
305    transmission and gated channel activity were significantly enriched in diverse species including
306    human, zebrafish, acorn worm, *Drosophila*, ant and octopus (Fig. 5a), confirming the important
307    role of RNA editing in modulating neural function in bilaterians. Representative examples are
308    the voltage-gated $K^+$ channels, that show the same recoding events on two highly conserved
309    amino acid residues within the ion transport domain among *Drosophila*, ant, octopus and even
310    human (Fig. 5b and Supplementary Table 5), and the glutamate ionotropic receptors in
311    vertebrates (Supplementary Fig. 5a and Supplementary Table 5). Interestingly, although a
312    nervous system is absent in the sponge [46], functional categories related to cellular
313    communication, signal transduction and response to stimulus were significantly enriched in
314    this early-branching and morphologically simple metazoan. Given that neurotransmission is
315    also part of the cell communication and signal transduction processes which mediate cellular
316    response to internal and external stimulus [47], these results imply that RNA editing might have

317    been adopted to modulate the molecular pathways of stimulus response during the early stage

318    of metazoan evolution.

319    However, it is unexpected that significant enrichment of cytoskeleton-related functions such

320    as cytoskeletal protein binding, actin cytoskeleton organization and motor activity, was

321    frequently observed in diverse bilaterians (Fig. 5a). Recoding editing of genes involved in

322    cytoskeleton-related functions has been only rarely reported previously [33,34]. The only well-

323    documented cases so far are the actin crosslinking proteins filamin α (*FLNA*) and filamin β

324    (*FLNB*), of which a conserved Q-to-R recoding event occur at the same position of both

325    proteins in mammals [48]. Our data not only confirms that recoding editing of *FLNA* or *FLNB*

326    occurs in humans, but also detects recoding events in sea urchin (*FLNB*), *Drosophila* (*FLNA*),

327    and acorn worm (*FLNA*). Other representative examples comprise the *cilia and flagella*

328    *associated protein 100* (*CFAP100*) which contains a S-to-G recoding event shared by oyster

329    and acorn worm (Fig. 5c and Supplementary Table 5), and *fascin* (an actin filament-bundling

330    protein) which has a Q-to-R recoding event shared by octopus, sea urchin and lancelet

331    (Supplementary Fig. 5b and Supplementary Table 5). The repeated emergence of same

332    recoding editing in the cytoskeleton-related genes in different lineages emphasizes an

333    important, but previously unappreciated, role of RNA editing in regulating cytoskeleton-

334    related functions in metazoans.

335

## Discussion

337    The phenomenon of RNA editing has been reported previously across a diverse range of

338    eukaryotes including metazoans, protists, fungi and plants, and to affect different types of RNA

339    [1,49]. However, while in most eukaryotes it is exclusively limited to mitochondrial or chloroplast

340    RNA, the extensive editing of nuclear-transcribed mRNA is phylogenetically rare, and

341    restricted to metazoans and some filamentous ascomycetes in which it originated through

342    independent mechanisms [3,50]. In this study, we present the first direct evidence that this method

343    for extensive alteration of nuclear DNA-encoded genetic information was adopted alongside

344    the origin of *ADARs* by the last common ancestor of extant metazoans ca 800 million years ago

345    [51], following its divergence from unicellular choanoflagellates. This raises the possibility that

346    ADAR-meditated RNA editing is an ancient regulatory process that was fundamental for initial

347    metazoan evolution. The evolutionary maintenance of ADAR-meditated RNA editing in

348    almost all extant metazoan lineages also emphasizes its essentiality in metazoan biology.

349 Consistent with the evolutionary constraint of ADARs, we show that the nucleotide sequence
350 and structural features surrounding A-to-I editing sites, including the strong favor of
351 uridine/adenosine and disfavor of guanosine in the adjacent 5' positions, and the tendency of
352 the underlying sequences to form dsRNA structures, are under strong constraint across the
353 animal kingdom, from the earliest branching ctenophore and sponge to human. These findings
354 might be valuable for ADAR-based RNA engineering, such as the recently reported RESTORE
355 and LEAPER approaches, which can recruit endogenous ADAR to specific transcripts for site-
356 directed RNA editing in human cells [52,53], as these conserved features imply that the approaches
357 developed based on one species (usually human) may well be easily applicable to other
358 metazoan species with *ADARs*.

359 It is now generally acknowledged that the complexity of transcriptional regulation coincides
360 with organismal complexity [54]. RNA editing and alternative splicing have long been proposed
361 to serve as important co/post-transcriptional regulatory mechanisms for increasing
362 transcriptome diversity [3,55]. However, while alternative splicing has been demonstrated to be
363 strongly associated with organismal complexity [56], we do not observe such a relationship
364 between the extent of global RNA-editing and organismal complexity in metazoans. Together
365 with our observations that in metazoans A-to-I editing preferentially targets evolutionary
366 young repetitive elements, and that recoding events in protein-coding sequences are generally
367 neutral or slightly deleterious, these findings question the ancestral role of A-to-I RNA editing
368 as a transcriptome or proteome diversifier in metazoans. Recent *ADAR1*-knockout studies in
369 human cells and mice indicated that ADAR1-mediated A-to-I editing of endogenous dsRNAs
370 formed by inverted repeats, plays a key role in preventing cellular sensing of endogenous
371 dsRNA as nonself (e.g. viral RNA), thus avoiding autoinflammation [18,57]. This suggests that
372 the avoidance of aberrant immune responses triggered by the accumulation of endogenous
373 dsRNA represents the primary driving force for preserving the extensive A-to-I editing in most
374 metazoan lineages. Alternatively, given that most editing sites are only edited at low to
375 moderate levels in all the species examined, and thus might not be sufficient to unwind dsRNAs
376 to avoid immune response, we hypothesize that metazoans may benefit from the maintenance
377 of mild single-nucleotide mutations in the RNA pool, as these mutations can provide plentiful
378 transcript variants that might help metazoans cope with unpredictable future conditions in their
379 life.

380 Our extensive survey across the phylogeny of metazoans also emphasizes that *Drosophila* and
381 cephalopods, whose RNA editomes harbor relatively high proportions of adaptive recoding

382   events subject to positive selection, are actually evolutionary exceptions in the animal kingdom.
383   The abundant recoding editing in cephalopods has been demonstrated to emerge in the ancestor
384   of coleoids after splitting from nautiloids, with the expansion of the cephalopod RNA editomes
385   [34]. In contrast, we find that the *Drosophila* RNA editomes have been greatly contracted in
386   comparison to most metazoans, while a considerable portion of recoding events is maintained
387   by natural selection. When this *Drosophila* pattern emerged during the evolution of insects
388   remains unknown. At least, the fact that more 'classic' RNA editomes, in which the majority
389   of sites targeting repetitive elements and rare recoding editing, are observed in ants and recently
390   in bumble bees [32], indicates that this *Drosophila* pattern must emerge after the divergence of
391   Diptera and Hymenoptera ca 345 million years ago [58].

392   RNA editing has been long acknowledged to regulate neural functions, affecting genes
393   encoding ion channels and neuroreceptors [7,9,10], consistent with the results of our functional
394   enrichment analysis of recoded genes in diverse species. Thus what is most surprising about
395   our observations is the over-representative of recoded genes encoding cytoskeleton-related
396   functions in diverse species, implying that post-transcriptional diversification of cytoskeleton-
397   related genes via RNA editing might be an important way through which to increase cellular
398   complexity during the evolution of metazoans. In particular, in some cases, we find exactly the
399   same positions are edited and cause the same amino acid changes in evolutionarily conserved
400   residues in distantly related species. The cytoskeleton is an interconnected network of
401   filamentous polymers and regulatory proteins, which carries out broad functions including
402   spatially organizing the contents of the cell, connecting the cell physically and biochemically
403   to the external environment and generating coordinated forces that enable the cell to move and
404   change shape [59]. It will be necessary for future studies to ascertain which aspect of
405   cytoskeleton-related functions RNA editing regulate.

406   In summary, our study provides the first large-scale and unbiased transcriptome-wide
407   investigation of RNA editing across the phylogeny of metazoans. These resources are valuable
408   for our understanding of the biological role and evolutionary principle of RNA editing in the
409   animal kingdom.
410

411   **Methods**
412   **Sample collection**
413   To rule out that false positives resulted from genetic variation during RNA-editing site
414   identification,   matching   DNA   and   RNA   sequences   generated   from   the   same

415    individual/specimen are the ideal data for use in RNA editing studies [39,60]. Thus, for the
416    metazoan species with sufficient body mass, both genomic DNA and total RNA were extracted
417    from the same individual, after grinding of the tissue/whole organism in liquid nitrogen. Two
418    to three individuals were collected as biological replicates. These species included the comb
419    jelly *Mnemiopsis leidyi* (three whole adults), the sponge *Amphimedon queenslandica* (three
420    biopsies from three adults), the sea anemone *Nematostella vectensis* (three whole adults), the
421    sea hare *Aplysia californica* (three whole juveniles), the oyster *Crassostrea gigas* (three whole
422    adults after removing shells), the sea urchin *Strongylocentrotus purpuratus* (three pairs of
423    gonad and non-gonad tissues dissected from one female and two male adults; non-gonad tissues
424    comprised the digestive, water vascular, and nervous systems), the acorn worm *Ptychodera*
425    *flava* (three whole adults), the lancelet *Branchiostoma belcheri* (three whole adults), the sea
426    squirt *Ciona savignyi* (two whole adults) and the zebrafish *Danio rerio* (three whole adults).

427    For metazoan species from which a single individual is not sufficient to allow the simultaneous
428    extraction of sufficient DNA and RNA for sequencing library construction, 10-15 individuals
429    with similar genetic background were pooled together, then both genomic DNA and total RNA
430    were extracted from the same pool of organisms after the whole pool was ground in liquid
431    nitrogen. These included the hydra *Hydra vulgaris* (10 adults per pool, two pools to serve as
432    biological replicates), the fruit fly *Drosophila melanogaster* (15 male adults per pool, two
433    pools), and *Drosophila simulans* (15 male adults per pool, two pools).

434    For the unicellular species and tiny metazoan species, biomass was first increased by the
435    propagation of a single colony with the same genetic background, then both genomic DNA and
436    total RNA were extracted from the same culture of organisms. These included the
437    ichthyosporean *Sphaeroforma arctica* (three cultures to serve as biological replicates), the
438    filasterean *Capsaspora owczarzaki* (three cultures), the choanoflagellate *Salpingoeca rosetta*
439    (three cultures) and *Monosiga brevicollis* (three cultures), and the metazoan *Trichoplax*
440    *adhaerens* (three cultures).

441    All the species were either collected from conventionally grown lab conditions, or obtained
442    from the wild. With the exception of the sea hare samples which were purchased from the
443    National Resource for Aplysia, University of Miami, 4600 Rickenbacker Causeway, Miami,
444    FL 33149, samples of all the other species were kindly provided by researchers who have
445    worked on corresponding species for years. The strain identifier (if applicable), geographical
446    origin and providers of each species were listed in Supplementary Table 1.

447    Genomic DNA of all species was extracted with the phenol/chloroform/isopentanol (25:24:1)

448    protocol. The integrity of the DNA samples was assayed by agarose gel electrophoresis

449    (concentration: 1 %; voltage: 150 V; Time: 40 min) before DNA-seq library construction. Total

450    RNA of all species except the choanoflagellates was extracted using TRIzol Reagent according

451    to manufacturer's protocol (Invitrogen, CA, USA). Total RNA of the choanoflagellates *S.*

452    *rosetta* and *M. brevicollis* was extracted using the RNAqueous Kit (Ambion, CA, USA). The

453    quality of the RNA samples was assayed by the Agilent 2100 Bioanalyzer (Thermo Fisher

454    Scientific, MA, USA) before RNA-seq library construction. In summary, a total of 53 DNA

455    and 53 RNA samples were obtained in this study. After quality control before library

456    construction, two out of the three RNA samples of *M. brevicollis* and one out of the three RNA

457    samples of *N. vectensis* were discarded due to poor RNA integrity (RIN < 6).

458

459    **Library construction and sequencing**

460    The strand-specific RNA-seq libraries for all the RNA samples were prepared using the TruSeq

461    Stranded mRNA LT Sample Prep kit (RS-122-2101, Illumina) with 1 μg total RNA as input,

462    then sequenced on the Illumina HiSeq 4000 platform using the PE100 chemistry, according to

463    the manufacturer's instructions (Illumina, San Diego, CA, USA).

464    The genomic DNA samples were either sequenced on an Illumina HiSeq 4000 or a BGISEQ-

465    500RS platform. The Illumina HiSeq and BGISEQ-500 platforms have been proved to generate

466    data with comparable quality and show high concordance for calling single nucleotide variants

467    by multiple independent studies [61-63]. For the Illumina DNA libraries, 1 μg genomic DNA per

468    sample was fragmented by a Covaris ultrasonicator, followed by end repair, 3′-end addition of

469    dATP and adapter ligation. The ligated fragments were then size selected at 300 bp on an

470    agarose gel and amplified by 10 cycles of PCR. The amplified libraries were purified using the

471    AxyPrep Mag PCR Clean-Up Kit (Axygen, MA, USA) then sequenced on the Illumina HiSeq

472    4000 platform using the PE100 chemistry according to the manufacturer's instructions

473    (Illumina, San Diego, CA, USA). The BGISEQ DNA sequencing libraries were prepared using

474    the MGIEasy DNA Library Prep Kit (V1.1, MGI Tech) with 1 μg genomic DNA as input, and

475    sequenced on the BGISEQ-500RS platform using the PE100 chemistry according to the

476    manufacturer's instructions (MGI Tech Co., Ltd., Shenzhen, China). Details about the

477    sequencing platform and data production for each sample were presented in Supplementary

478    Table 1.

479

**Identification of RNA-editing sites**

*(i) Quality control for raw sequencing data*

All the DNA- and RNA-seq reads were first submitted to SOAPnuke (v1.5.6) [64] for quality control by removal of adapter-contaminated reads and low-quality reads before subsequent analyses with parameters *-G -l 20 -q 0.2 -E 60 -5 1 -Q 2*.

*(ii) Adjustment of reference genome with DNA-seq data*

Given that many samples were collected from wild animals, which have high levels of heterozygosity, or were from strains which are genetically different from those used for assembling the reference genomes, we employed Pilon (v1.21) [65] to adjust the reference genome of each species using the DNA-seq data from different samples separately, generating sample-specific reference genomes for each species before RNA-editing site identification. Specifically, DNA sequence reads from each sample of a species were first aligned to the published reference genome using BWA-MEM (v0.7.15) [66] with default parameters. Then, genome adjustment was performed by Pilon with default parameters except that *--fix snps* was set, using the original reference genome FASTA and the DNA BAM files as input. It is noteworthy that we only adjusted SNPs in the reference genomes in order to ensure that the adjusted genomes from different samples of the same species have the same length and the same coordinate system. The version and source of the original reference genome for each species were listed in Supplementary Table 1.

*(iii) Identification of RNA-editing sites with RES-Scanner*

RNA-editing sites from each sample were first identified by RES-Scanner (v20160713), a software package that was designed to identify transcriptome-wide RNA-editing sites with matching DNA- and RNA-seq data from the same individual or specimen [39]. Briefly, RES-Scanner invoked BWA-ALN (v0.7.15) [67] to align the DNA and RNA reads that passed quality control to the adjusted reference genome of each species, followed by filtering low-quality alignments, calling homozygous genotype from DNA data, and identifying candidate RNA-editing sites from RNA data by ruling out false-positives resulted from genetic variants and sequencing or alignment errors. In general, default parameters were used for the whole pipeline, except that the mapping quality cutoff was set to 5 for DNA alignment (default 20) and the numbers of bases masked at the 5'- and 3'-end of a DNA read was set to 0 (default 6). This was done as we found that lowering these requirements for the DNA data could yield RNA-editing sites with higher accuracy in many species, manifesting as the higher proportions of A-to-I editing sites out of all identified editing sites.

513     ***(vi) Identification of hyper-editing sites***

514     Given that A-to-I editing sites tend to occur in clusters, the heavily edited RNA reads

515     (commonly called hyper-edited reads) which contain many of the same type of substitutions in

516     relation to the reference genome, often fail to be aligned during normal alignment process. In

517     order to capture these hyper-edited reads and the clusters of editing sites they harbor, we next

518     performed hyper-editing detection for each sample following a scheme originally proposed by

519     Porath *et al* [40].

520     We first collected the RNA read pairs that could not be aligned to the adjusted reference

521     genome or that had mapping quality < 20 from the RNA BAM files generated by the RES-

522     Scanner pipeline as described above. We then removed the read pairs for which one or both

523     reads contained more than 10% of Ns along their lengths, or had particularly large (>60%) or

524     small (<10%) percentage of a single-type nucleotide as recommended by Porath *et al* [40]. Next,

525     we adopted a "three-letter" alignment strategy to align these potential hyper-edited reads, in

526     order to overcome the excess mismatches in relation to the reference genome. For example, to

527     align the RNA reads with many A-to-I editing sites (i.e. many A-to-G mismatches), all Ts in

528     the first read of a read pair were transformed to Cs, and all the As in the second read of a read

529     pair were transformed to Gs. This is because, for read pairs generated from the dUTP-based

530     strand-specific RNA-seq libraries, the second read is from the original RNA strand/template

531     while the first read is from the opposite strand [68]. In the meantime, two versions of the reference

532     genome were created, of which the first version was named the *positive* reference, with all As

533     transformed to Gs, and the second version was named the *negative* reference, with all Ts

534     transformed to Cs.

535     Next, the transformed read pairs were aligned to both the *positive* and *negative* references by

536     BWA-ALN with parameters *-n 0.02 -o 0*, yielding the *positive* and *negative* alignments,

537     respectively. Then, we filtered both alignments by removing read pairs that were not aligned

538     to the reference genome concordantly, and the reads within concordantly aligned pairs that had

539     mapping score < 20. In addition, for *positive* alignment, we further required that the first read

540     in a pair was the reverse complement of the reference genome, while the second read was

541     aligned to reference genome directly; for *negative* alignment, we required that the first read in

542     a read pair was directly aligned to reference genome, while the second read was the reverse

543     complement of the reference genome.

544     After the strict quality control for the BWA alignments, we converted the transformed reads to

545     their original sequences, followed by trimming the first and last 10 bases of each read in the

546     alignments. Then we identified hyper-edited reads by requiring the mismatch rate of a trimmed

547     read to be > 5%, and the proportion of the expected mismatches (i.e. A-to-G substitution in this

548     example) against all mismatches to be > 60% as recommended by Porath *et al* [40]. Finally, BAM

549     files of hyper-edited RNA reads were submitted to RES-Scanner to extract potential editing

550     sites together with the matching DNA BMA files generated in the previous step. RES-Scanner

551     was run with default parameters in general, except that the mapping quality cutoff was set to 5

552     for DNA alignment, the numbers of bases masked at the 5'- and 3'-end of a read were set to 0

553     for both DNA and RNA reads, the minimum number of RNA reads supporting editing was set

554     to 2 (default 3), and the minimum editing level was set to 0 (default 0.05).

555     The above hyper-editing detection method was undertaken for all of the 12 possible

556     substitution types of RNA editing in each sample of a species, and the results from all the 12

557     substitution types were combined together by discarding those sites that presented different

558     editing types in any single genomic position.

### (v) Combing the results of RES-Scanner and hyper-editing detection

560     To generate the representative RNA-editing sites for a species, and to improve the

561     identification of editing sites in each sample, we combined the editing sites identified by RES-

562     Scanner (step iii) and hyper-editing detection (step vi) in each sample, to obtain a

563     comprehensive map of potentially editable positions in the reference genome of each species.

564     If a genomic position was identified as an editing site in both methods, we respectively added

565     the numbers of RNA reads supporting editing, and the number supporting non-editing as

566     generated by these two methods. We then retrieved the missed editing sites in each sample in

567     these editable positions using the criteria of at least one RNA read supporting editing and the

568     false discovery rate (FDR) [69] adjusted $p$ value for this site to be resulted from sequencing error

569     < 0.01. Specifically, statistical tests were performed based on the binomial distribution B($k$, $n$,

570     $p$), where $p$ was set to be the maximal probability of an RNA base to be a sequencing error (i.e.

571     0.1% here as we only used RNA bases with Phred quality score $\geq$ 30), $n$ was equal to the total

572     read depth of a given candidate editing site, and $k$ denoted the number of reads supporting

573     editing. We also used the DNA-seq data from multiple samples to further remove false-

574     positives resulted from genetic variants, by discarding those editing sites for which the genomic

575     DNA showing the same type of substitution as RNA editing (i.e. the frequency of edited base

576     versus the total number of bases covering this position > 0.1) in any one of the multiple DNA

577     samples. RNA-editing sites that displayed different editing types in different samples of a

578     species were also discarded.

579      We have updated the software package RES-Scanner we previously established for RNA-

580      editing site scanning by compiling above steps (step *i* to *v*). This RES-Scanner2 now can also

581      identify hyper-editing sites. It works from raw sequencing reads and is applicable to RNA-

582      editing site detection in any species with matching DNA- and RNA-seq data.

583

584      **RNA-editing sites for additional metazoan species**

585      To increase the phylogenetic coverage of the investigated species, we collected the matching

586      DNA-seq and strand-specific RNA-seq data from the nematode *Caenorhabditis elegans*

587      (pooled whole organisms collected from three larval stages and two adult stages) [26], the leaf-

588      cutting ant *Acromyrmex echinatior* (three pooled head samples of the small worker caste

589      collected from three colonies, respectively) [31], the octopus *Octopus bimaculoides* (four neural

590      tissue samples including faxial nerve cord, optic lobe, subesophageal ganglia and

591      supraesophageal ganglia) [37] and human (three brain samples from three male adults,

592      respectively) [22]. The SRA accession numbers and statistics of the downloaded sequencing data

593      were presented in Supplementary Table 1. RNA-editing sites in each of the four species were

594      identified using the same procedure (step *i* to *v*) as described above.

595

596      **Refining the ORFs and annotating UTRs for protein-coding genes**

597      Protein-coding genes (GFF/GTF and corresponding cds/pep FASTA files) were downloaded

598      from public databases along with the reference genomes, of which the sources were presented

599      in Supplementary Table 1. The correctness of the open-reading frames (ORFs) in the GFF/GTF

600      files were checked for all the protein-coding genes, with the defective ORFs such as those that

601      were not the integer multiple of 3 in length or not exactly matching the protein sequences

602      presented in the downloaded pep FASTA files being carefully corrected by in-house scripts.

603      Then the transcript model with the longest ORF was chosen as the representative model for a

604      locus if multiple transcript models were annotated in this locus.

605      5'- and 3'-UTRs for the representative ORFs were annotated using the RNA-seq data used in

606      this study, for all the species except for human. Briefly, RNA-seq reads that passed quality

607      control as described above were first aligned to the reference genome of each species by

608      HISAT2 (v2.1.0) [70], with default parameters except setting *--rf*, followed by removing those

609      reads that could be mapped to multiple positions of the genome. Then, transcribed regions with

610      continual RNA depth $\geq$ 5X were extended from the 5'- and 3'-end of each representative ORF

611      to serve as initial 5'- and 3'-UTRs, respectively. Next, an iterative process was used to further

612   recruit the upstream or downstream transcribed regions that were apart from, but linked by $\geq 5$

613   junction reads to previously defined UTRs. If a gene had different 5'- or 3'-UTRs annotated in

614   different samples, the longest one was chosen as the representative 5'- or 3'-UTR for this gene.

615

616   **Gene expression quantification and transcript assembly with RNA-seq data**

617   HISAT2 alignments generated in the above analysis were used to quantify gene expression

618   levels for the refined representative gene models in each species. Only the RNA-seq reads that

619   were aligned to one position of the reference genome, and that overlapped with annotated exons

620   were kept for expression quantification. Gene expression levels were measured by RPKM

621   (reads per kilobase per million mapped exonic reads), and the RPKM values in all the

622   sequenced samples from the same species were adjusted by a scaling normalization method

623   based on TMM (trimmed mean of M values) to normalize the sequencing bias among samples

624   [71]. We also assembled transcripts for each species with StringTie (v1.3.4d) [72] with default

625   parameters using the HISAT2 alignments as input. These transcript models were regarded as

626   one kind of reference models during the manual annotation of *ADAR* genes as described below.

627

628   **Annotation of *ADAR* genes in each species**

629   ADAR protein sequences of *Nematostella vectensis* (XP_001642062.1 and XP_001629615.1),

630   *Drosophila melanogaster* (NP_569940.2), *Caenorhabditis elegans* (NP_492153.2 and

631   NP_498594.1), *Crassostrea gigas* (EKC20855.1 and EKC32699.1), *Strongylocentrotus*

632   *purpuratus* (XP_011680614.1 and XP_781832.1), *Ciona intestinalis* (XP_002128212.1),

633   *Danio rerio* (NP_571671.2, NP_571685.2, XP_021334693.1 and XP_686426.5) and *Homo*

634   *sapiens* (XP_024305442.1, NP_056648.1 and NP_061172.1) collected from NCBI were used

635   as queries to search for *ADAR* genes in reference genomes of all the 22 species by TBLASTN

636   (blast-2.2.23) [73] with parameters *-F F -e 1e-5*, followed by the determination of gene structure

637   and protein sequences in the target species with GeneWise (wise2.2.0) [74]. The predicted

638   proteins were then aligned to the NCBI nr database to confirm whether they were ADARs.

639   Next, we manually compared the gene models in the putative *ADAR* loci resulted from

640   homologous predictions, transcript assemblies by StringTie and the published gene set of each

641   species, and we chose the models with the longest ORFs as the representative models. Domain

642   organizations of the manually confirmed ADAR proteins were predicted using the CD-Search

643   tool in NCBI (CDD v3.17; https://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi) and Pfam

644   (release-32.0; https://pfam.xfam.org) with default settings, and only ADARs with at least one

645    dsRNA binding domain (dsRBDs) and one adenosine-deaminase domain (AD) were regarded

646    as potential *ADAR* genes. Of note, *ADAD* genes, which usually contain one or more dsRBDs

647    and one AD, were also identified as potential *ADARs* by our criteria, but they could be

648    distinguished from *ADARs* according to phylogenetic analysis (see below). The information of

649    *ADAR* genes annotated in each species, including the coding nucleotide sequences, protein

650    sequences, domain annotations and editing sites are presented in Supplementary Table 3.

651    Phylogenetic analysis of all the potential ADARs identified above, were performed with the

652    AD peptide sequences (ca 324 amino acids in length) using MEGA7 with the neighbor-joining

653    method [75]. We did not perform phylogenetic analysis whit the dsRBDs, as the lengths of the

654    dsRBDs were generally very short (ca 40 to 60 amino acids) and the copy number of dsRBDs

655    varied among ADARs both within and between species. The peptide sequences of ADs used

656    for phylogenetic analysis were aligned using ClustalW as implemented in MEGA7. Reliability

657    of the trees was estimated using 1,000 bootstrap replications (Supplementary Fig. 2d). To

658    further estimate the divergence between any two potential ADARs, we calculated the

659    nonsynonymous substitution rates ($d_N$) for any pair of potential ADARs using PAML (v4.9i)

660    [76] with the Yang & Nielsen (2000) method [77], according to the codon alignment of the ADs

661    (Supplementary Fig. 2e).

662

663    **Identification of editing sites locating in potential dsRNA regions**

664    The dsRNA regions formed by two potential mechanisms, intramolecular folding of pre-

665    mRNA and intermolecular hybridization of sense-antisense transcripts, were tested for the

666    enrichment of A-to-I editing sites.

667    For the mechanism of intramolecular folding, we extracted a 401 nt sequence centered on each

668    A-to-I editing site, then searched this query sequence against a 4001 nt sequence centered on

669    corresponding A-to-I editing site using BLASTN (v2.2.26) with parameters *-F F -e 1e-2*. Then

670    an A-to-I editing site was identified as locating in a dsRNA region formed by intramolecular

671    folding, if a reverse-complement alignment was detected with identity $\geq$ 80%, the aligned

672    length was $\geq$ 50 nt, and the aligned region of the query sequence spanned the edited adenosine.

673    For the mechanism of intermolecular hybridization of sense-antisense transcripts, we examined

674    the RNA coverage of a 101 nt region centered on each A-to-I editing site, and searched for the

675    regions with RNA depth $\geq$ 2X along >50% of the region length, on both strands.

676    To estimate the expected ratio of A-to-I editing sites that occurred in dsRNA regions formed

677    by the above two different mechanisms in each sample, we randomly selected an adenosine

678    site with comparable RNA depth (i.e. within ± 20% of the editing site) for each editing site in
679    a sample, and performed the same analyses for these control adenosine sites. The significance
680    levels for the difference between the observed and expected ratios were examined by two-tailed
681    paired t-tests in each species.
682

683    **Definition of clustered and isolated editing sites**
684    For each sample of a species, we considered a genomic region containing ≥ 3 A-to-I editing
685    sites, of which the distance for two adjacent sites was ≤ 30 nt, as an RNA-editing cluster. The
686    genomic locations of the first and last editing sites in a cluster were assigned as the start and
687    end genomic positions of this cluster. A-to-I editing sites located in the defined editing clusters
688    were regarded as clustered editing sites, and those outside editing clusters were regarded as
689    isolated editing sites. To estimate the expected ratio of A-to-I editing sites occurring in clusters
690    in each sample, we randomly selected an adenosine site with comparable RNA depth (i.e.
691    within ± 20% of the editing site) for each editing site in a sample, and calculated the ratio of
692    these control adenosine sites occurring in clusters. The significance levels for the difference
693    between the observed and expected ratios were examined by two-tailed paired t-tests in each
694    species.
695

696    **Analysis of the neighboring nucleotide preference for A-to-I editing**
697    The Two Sample Logo software (v1.21) [78] was used to analyze the neighboring nucleotide
698    preference of A-to-I editing sites with parameters *-K N -T binomial -C nucleo_weblogo -y*.
699    Specifically, for each species, the eleven-nucleotide sequences with the edited adenosines in
700    the center were used as the foreground dataset, while the eleven-nucleotide sequences centered
701    by the transcribed (RNA depth ≥ 2X) but unedited adenosines locating within ± 50 nt of the
702    edited adenosines, were used as the background dataset for Two Sample Logo analysis.
703    Nucleotides were plotted using the size of the nucleotide that was proportional to the difference
704    between the foreground and background datasets.
705

706    **Annotation of repetitive elements**
707    Considering that the repetitive elements of many species investigated in this study are either
708    not well annotated and/or not publicly available, we re-annotated the repetitive elements of all
709    the sampled species except human using the same strategy. Repetitive elements of the human
710    genome (GRCh38/hg38) have been well annotated and thus were downloaded from UCSC
711    directly.

712    Repetitive elements in the genome assembly of other sampled species were identified by

713    homology searches against known repeat databases and *de novo* predictions as previously

714    described [79]. Briefly, we carried out homology searches for known repetitive elements in each

715    genome assembly by screening the Repbase-derived RepeatMasker libraries with

716    RepeatMasker (v4.0.6; setting *-nolow -no_is -norna -engine ncbi*) [80] and the transposable

717    element protein database with RepeatProteinMask (an application within the RepeatMasker

718    package; setting *-noLowSimple -pvalue 0.0001 -engine ncbi*). For *de novo* prediction,

719    RepeatModeler (v1.0.8) [81] was executed on the genome assembly to build a *de novo* repeat

720    library for each species, respectively. Then RepeatMasker was employed to align the genome

721    sequences to the *de novo* library for identifying repetitive elements. We also searched each

722    genome assembly for tandem repeats using Tandem Repeats Finder (v4.07) [82] with parameters

723    *Match=2 Mismatch=7 Delta=7 PM=80 PI=10 Minscore=50 MaxPeriod=2000*. To confirm

724    the reliability of our annotations, we compared our repeat annotation results of the fruit fly

725    *Drosophila melanogaster* and the zebrafish *Danio rerio* with those downloaded from UCSC

726    and observed good consistency (Supplementary Fig. 3a,b).

727

728    **Calculation of RNA-editing density for different genomic elements**

729    To compare the probability of different genomic elements targeted by A-to-I editing, including

730    the protein-coding genes related elements (5'-UTR, CDS, intron and 3'-UTR) and the repeat-

731    associated elements (SINE, LINE, LTR, DNA transposon, Helitron, tandem repeat and other

732    unclassified repeat loci), we calculated the A-to-I editing density for each type of genomic

733    element by counting the number of A-to-I editing sites located in this element type, out of the

734    total number of transcribed adenosines (RNA depth $\geq$ 2X) from this element type. The editing

735    density of each element type was first calculated for each sample of a species separately, then

736    the mean editing density across samples was calculated as the representative value for a species.

737    When calculating the editing-level-weighted editing densities for each element type, an editing

738    site with for example an editing level of 0.1, would be regarded as 0.1 editing site instead of 1

739    editing site, when counting the number of editing sites for an element type. Only editing sites

740    and transcribed adenosines with RNA depth $\geq$ 10X were used in the weighted analysis.

741

742    **Analysis of relationship between repeat divergence and editing density**

743    The divergence rates of repetitive elements in each species were estimated by RepeatMasker,

744    by comparing the repeat sequences to the ancestral consensus sequences identified by

745     RepeatModeler during the repeat annotation process as described above. Only the transcribed

746     repeat loci with no less than 50 nucleotides covered by $\geq 2$ RNA reads were used for this

747     analysis. The transcribed repeat loci were first sorted according to divergence rate from the

748     lowest to the highest (i.e. the youngest to oldest), then divided into 10 equal bins with the same

749     transcribed repeat loci in each bin. Next the editing density for each bin was calculated, as the

750     number of A-to-I editing sites located in repeat loci belonging to this bin, divided by the total

751     number of transcribed adenosines (RNA depth $\geq 2X$) from the repeat loci in this bin. The

752     editing density of each bin was first calculated for each sample of a species separately, then the

753     mean editing density across samples was calculated as the representative value for a species.

754     The relationships between repeat divergence rate and editing density in all species were

755     displayed by a heatmap as presented in Fig. 3d.

756

757     **Estimating the potentials of repeat and non-repeat regions to form dsRNA**

758     The potential of repeat and non-repeat genomic regions to form dsRNA was approximatively

759     measured as the ratios of repeat and non-repeat derived genomic sites locating in regions that

760     could find a reverse-complement alignment in nearby regions. Briefly, we randomly selected

761     100,000 sites from the genomic regions annotated as repeat and non-repeat, respectively. Then,

762     we extracted a 401 nt sequence centered on each randomly selected site and searched this query

763     sequence against a 4001 nt sequence centered on the corresponding repeat or non-repeat

764     genomic site using BLASTN (v2.2.26) with parameters *-F F -e 1e-2*. Then a repeat or non-

765     repeat derived genomic site was regarded as locating in a potential dsRNA region formed by

766     intramolecular folding, if a reverse-complement alignment was detected with identity $\geq 80\%$,

767     aligned length $\geq 50$ nt, and the aligned region of the query sequence spanned this randomly

768     selected site. The ratio of such sites against all randomly selected sites was calculated to

769     represent the potential of repeat or non-repeat regions to form dsRNA in a species, and the

770     same process was iterated for 100 times to estimate the distribution (see Supplementary Fig.

771     4c).

772

773     **Analyzing the adaptive potential of recoding editing**

774     Recoding editing sites were identified as the sites where the editing events could cause

775     nonsynonymous changes in protein-coding regions. Given that the numbers of recoding sites

776     were generally small in most species, for the evolutionary analysis of recoding editing,

777     recoding sites from different samples of a species were first combined according to their

778   genomic locations. The editing level of a combined recoding site was measured as the mean

779   editing level across samples with RNA coverage $\geq$ 10X in this position.

780   To examine the adaptive potential of recoding editing in a species, we compared the frequency

781   of nonsynonymous editing ($f_n$) to the frequency of nonsynonymous editing ($f_s$) as previously

782   described [35]. Specifically, $f_n$ was calculated as the number of A-to-I editing sites causing

783   nonsynonymous changes ($n$), divided by the number of potential nonsynonymous adenosine

784   sites (RNA depth $\geq$ 2X in at least one sample) if A is replaced with G ($N$) from the genes with

785   $\geq$ 1 editing site in their coding regions. $f_s$ was calculated as the number of A-to-I editing sites

786   causing synonymous changes ($s$), divided by the number of potential synonymous adenosine

787   sites (RNA depth $\geq$ 2X in at least one sample) if A is replaced with G ($S$) from the same set of

788   genes. If recoding editing is generally adaptive in a species, one would expect that $f_n$ is

789   significantly larger than $f_s$ in this species. The significance level for the difference between $f_n$

790   and $f_s$ in a species was assessed by a two-tailed Fisher's exact test using the values of $n$, $N$, $s$

791   and $S$ from this species.

792   To compare the adaptive potential for recoding sites with different editing levels, the same

793   analyses were performed for recoding sites with relatively high ($\geq$ 0.2) and low ($<$ 0.2) editing

794   levels separately, using the sites with RNA depth $\geq$ 10X and the genes with one or more editing

795   sites achieving this RNA depth in their coding regions.

796

797   **Analyzing the evolutionary conservation of recoded genes**

798   Recoded genes were defined as the protein-coding genes with at least one recoding site. To

799   evaluate the evolutionary conservation of the recoded genes in the seventeen species with

800   reliable A-to-I editing (the target species), we identified the orthologous gene of each recoded

801   gene in a closely-related species with a publicly available reference genome (the related

802   species), and calculated the $d_N/d_S$ ratio (i.e. the ratio of the number of nonsynonymous

803   substitutions per nonsynonymous site ($d_N$) to the number of synonymous substitutions per

804   synonymous site ($d_S$)) for each orthologous pair. The closely-related species chosen for each

805   target species is presented in Supplementary Table 6.

806   Briefly, all the protein sequences from each target species were first aligned to its related

807   species genome using TBLASTN (blast-2.2.26) with parameters *-F F -e 1e-5*, followed by

808   chaining the syntenic blocks and picking one candidate locus for each target-species protein

809   with the highest TBLASTN bit score by in-house scripts. Then the genomic sequences of

810    these candidate loci together with 2 kb flanking sequences, were extracted from the related-

811    species genome and submitted to GeneWise (wise-2.4.1) to determine the protein sequences

812    by aligning the target-species proteins to these related-species genomic sequences. The

813    related-species proteins were then aligned back to all the protein sequences of the target

814    species using BLASTP (blast-2.2.26) with parameters *-F F -e 1e-5*, and only those hitting the

815    expected proteins in the target species with the highest BLASTP bit score were identified as

816    orthologous proteins in related species. Next, the protein sequences of each orthologous pair

817    were aligned using MAFFT (v6.923) [83] with parameters *--maxiterate 1000 –localpair*,

818    followed by the replacement of the amino acids by their corresponding codons for each species.

819    The orthologous pairs of which the MAFFT alignments with invalid sites (i.e. presented as "-"

820    in one of the two aligned sequences) exceeding 50% of the alignment length were discarded.

821    Then the $d_N/d_S$ ratio for each qualified orthologous pair was calculated using PAML (v4.9i) [76]

822    with the Yang & Nielsen (2000) method [77].

823    Finally, the genes of each target species were divided into three groups according to the degree

824    of evolutionary conservation, and the observed/expected number of recoded genes among

825    different groups was calculated. Specifically, group I was comprised of genes with orthologs

826    in closely-related species and $d_N/d_S$ ratios lower than the median value among all orthologous

827    pairs, representing the most conserved group; Group II was comprised of genes with orthologs

828    in closely-related species with $d_N/d_S$ ratios higher than the median value among all orthologous

829    pairs, representing the moderately conserved group; Group III was comprised of all the

830    remaining genes that cannot find orthologs in closely-related species, representing the least

831    conserved group. The expected probability of a gene being recoded in a species was estimated

832    as the number of recoded genes out of all transcribed protein-coding genes (RPKM > 1 in at

833    least one sample) in this species, and the expected number of recoded genes in each

834    conservation group was calculated as the number of genes in this group multiplied by the

835    expected probability of a gene being recoded. The significance level for the difference between

836    observed and expected numbers in each conservation group was estimated by a two-tailed

837    binomial test.

838

839    **Functional annotation and enrichment analysis of recoded genes**

840    GO annotations for the protein-coding genes were downloaded from Ensembl (*Caenorhabditis*

841    *elegans*, *Ciona savignyi*, *Danio rerio* and *Homo sapiens*) or Ensembl Metazoa (*Mnemiopsis*

842    *leidyi*, *Amphimedon queenslandica*, *Drosophila melanogaster*, *Drosophila simulans*,

843    *Crassostrea gigas*, *Octopus bimaculoides*, *Nematostella vectensis* and *Strongylocentrotus*

844    *purpuratus*) via the BioMart function. For *Hydra vulgaris*, *Aplysia californica*, *Acromyrmex*

845    *echinatior*, *Ptychodera flava* and *Branchiostoma belcheri* that do not have publicly available

846    GO annotations, we first aligned all the proteins of these species to the UniProt database

847    (release-2019_04) using BLASTP (blast-2.2.26) with parameters *-F F -e 1e-5*. Then the best

848    hit of each query gene was retained based on its BLASTP bit score, and the GO annotations

849    of this best hit was assigned to the query gene.

850    GO enrichment analysis was conducted for genes with at least one recoding site of which the

851    mean editing level across samples > 0.1, or the editing event shared by at least two samples, in

852    order to reduce the influence of nonadaptive recoding sites that are likely the by-products of

853    promiscuous ADAR activity. Hypergeometric tests were employed to examine whether the

854    recoded genes of a species was enriched in a specific GO term in relation to background genes

855    as previously described [31], by comparing the number of recoded genes annotated to this GO

856    term, the number of recoded genes not annotated to this GO term, the number of background

857    genes (i.e. the protein-coding genes with RPKM > 1 in at least one sample after excluding the

858    recoded genes in the species) annotated to this GO term, and the number of background genes

859    not annotated to this GO term. *P*-values were adjusted for multiple testing by applying FDR [69],

860    and the GO terms with adjusted *p*-values < 0.05 in at least three species (Note: GO terms shared

861    by *D. melanogaster* and *D. simulans* were only counted once here) were considered as the

862    general functional categories preferred by metazoan recoding editing.

863

864    **Identification of conserved recoding events shared by multiple species**

865    To identify recoding events shared by two or more species, we first identified the orthologous

866    groups of genes (i.e. gene families) from the seventeen metazoan species with reliable RNA

867    editing using OrthoFinder (v2.2.7) [84] with default parameters. For the gene families that

868    contained recoded genes from multiple species, we aligned the protein sequences of the

869    recoded genes using MUSCLE (v3.8.31) [85] with parameter *-maxiters 1000* and filtered poorly

870    aligned positions using Gblocks (v0.91b) [86]. Next recoding events occurring in the same

871    position in the alignments and causing the same amino acid changes among at least two species

872    were identified as conserved recoding events. Recoding events only shared by *D. melanogaster*

873    and *D. simulans* were removed. Only recoding sites in which the mean editing levels were no

874    less than 0.1 across samples of a species, or were shared by at least two samples, were used in

875    this analysis. The complete list of recoding events shared by multiple species was presented in
876    Supplementary Table 5.
877

## Data and code availability

879    The raw sequencing reads generated in this study are deposited in NCBI under the BioProject
880    accession PRJNA557895 and are also deposited in the CNGB Nucleotide Sequence Archive
881    (CNSA) with accession number CNP0000504 (https://db.cngb.org/cnsa/). RNA-editing sites,
882    refined gene annotations and repeat annotations used in this study are deposited in the figshare
883    repository under the link https://doi.org/10.6084/m9.figshare.10050437. Codes are available
884    upon request.

885

## Acknowledgements

904

## Author Contributions

906    Q.L. and G.Z. conceived the study; M.T.P.G. and Q.L. coordinated the sample collection from
907    different labs around the world; N.L. and L.G. conducted lab work for the culture and collection

908    of *T. adhaerens* samples; M.D.M. conducted lab work for the culture and DNA/RNA extraction

909    of *S. rosetta* and *M. brevicollis*; M.A.S. and I.R.-T. conducted lab work for the culture and

910    DNA/RNA extraction of *S. arctica and C. owczarzaki*; M.Q.M. collected the *M. leidyi* samples

911    and performed DNA/RNA extraction; Y.-H.S. and J.-K.Y. collected the *S. purpuratus*, *P. flava*

912    and *B. belcheri* samples and performed dissection for *S. purpuratus*; X.Z. managed library

913    construction and sequencing of all species; P.Z., J.L., H.Y., Y.Z., Q.G., H.T. and X.Z.

914    performed bioinformatic analyses under the supervision of Q.L.; G.Z., N.L., I.R.-T., M.Q.M.

915    and J.-K.Y. contributed reagents and materials; Q.L. and G.Z. wrote the manuscript with the

916    inputs from all authors. All authors read and approved the final manuscript.

917

918    **Competing interests**

919    The authors declare no competing interests.

920

921    **References**

922   1    Gott, J. M. & Emeson, R. B. Functions and mechanisms of RNA editing. *Annual review*
923      *of genetics* **34**, 499-531, doi:10.1146/annurev.genet.34.1.499 (2000).

924   2    Nishikura, K. Functions and regulation of RNA editing by ADAR deaminases. *Annual*
925      *review of biochemistry* **79**, 321-349, doi:10.1146/annurev-biochem-060208-105251
926      (2010).

927   3    Eisenberg, E. & Levanon, E. Y. A-to-I RNA editing - immune protector and
928      transcriptome diversifier. *Nature reviews. Genetics* **19**, 473-490, doi:10.1038/s41576-
929      018-0006-1 (2018).

930   4    Nishikura, K. Editor meets silencer: crosstalk between RNA editing and RNA
931      interference. *Nat Rev Mol Cell Biol* **7**, 919-931, doi:10.1038/nrm2061 (2006).

932   5    Rieder, L. E. & Reenan, R. A. The intricate relationship between RNA structure, editing,
933      and splicing. *Semin Cell Dev Biol* **23**, 281-288, doi:10.1016/j.semcdb.2011.11.004
934      (2012).

935   6    Nishikura, K. A-to-I editing of coding and non-coding RNAs by ADARs. *Nat Rev Mol*
936      *Cell Biol* **17**, 83-96, doi:10.1038/nrm.2015.4 (2016).

937   7    Behm, M. & Ohman, M. RNA Editing: A Contributor to Neuronal Dynamics in the
938      Mammalian Brain. *Trends in genetics : TIG* **32**, 165-175, doi:10.1016/j.tig.2015.12.005
939      (2016).

940   8    Hwang, T. *et al.* Dynamic regulation of RNA editing in human brain development and
941      disease. *Nature neuroscience* **19**, 1093-1099, doi:10.1038/nn.4337 (2016).

9    Jepson, J. E. & Reenan, R. A. RNA editing in regulating gene expression in the brain. *Biochimica et biophysica acta* **1779**, 459-470, doi:10.1016/j.bbagrm.2007.11.009 (2008).

10   Li, J. B. & Church, G. M. Deciphering the functions and regulation of brain-enriched A-to-I RNA editing. *Nature neuroscience* **16**, 1518-1522, doi:10.1038/nn.3539 (2013).

11   Garrett, S. & Rosenthal, J. J. RNA editing underlies temperature adaptation in K+ channels from polar octopuses. *Science* **335**, 848-851, doi:10.1126/science.1212795 (2012).

12   Rieder, L. E. *et al.* Dynamic response of RNA editing to temperature in *Drosophila*. *BMC biology* **13**, 1, doi:10.1186/s12915-014-0111-3 (2015).

13   Buchumenski, I. *et al.* Dynamic hyper-editing underlies temperature adaptation in *Drosophila*. *PLoS genetics* **13**, e1006931, doi:10.1371/journal.pgen.1006931 (2017).

14   Zipeto, M. A., Jiang, Q., Melese, E. & Jamieson, C. H. RNA rewriting, recoding, and rewiring in human disease. *Trends Mol Med* **21**, 549-559, doi:10.1016/j.molmed.2015.07.001 (2015).

15   Ben-Aroya, S. & Levanon, E. Y. A-to-I RNA Editing: An Overlooked Source of Cancer Mutations. *Cancer Cell* **33**, 789-790, doi:10.1016/j.ccell.2018.04.006 (2018).

16   Maas, S., Kawahara, Y., Tamburro, K. M. & Nishikura, K. A-to-I RNA editing and human disease. *RNA biology* **3**, 1-9 (2006).

17   Rice, G. I. *et al.* Mutations in ADAR1 cause Aicardi-Goutieres syndrome associated with a type I interferon signature. *Nat Genet* **44**, 1243-1248, doi:10.1038/ng.2414 (2012).

18   Chung, H. *et al.* Human ADAR1 Prevents Endogenous RNA from Triggering Translational Shutdown. *Cell* **172**, 811-824 e814, doi:10.1016/j.cell.2017.12.038 (2018).

19   Peng, Z. *et al.* Comprehensive analysis of RNA-Seq data reveals extensive RNA editing in a human transcriptome. *Nature biotechnology* **30**, 253-260, doi:10.1038/nbt.2122 (2012).

20   Bazak, L. *et al.* A-to-I RNA editing occurs at over a hundred million genomic sites, located in a majority of human genes. *Genome research*, gr. 164749.164113 (2013).

21   Ramaswami, G. *et al.* Identifying RNA editing sites using RNA sequencing data alone. *Nature methods* **10**, 128-132 (2013).

22   Picardi, E. *et al.* Profiling RNA editing in human tissues: towards the inosinome Atlas. *Sci Rep* **5**, 14941, doi:10.1038/srep14941 (2015).

23   Bahn, J. H. *et al.* Accurate identification of A-to-I RNA editing in human by transcriptome sequencing. *Genome Res* **22**, 142-150, doi:10.1101/gr.124107.111 (2012).

979   24    Danecek, P. *et al.* High levels of RNA-editing site conservation amongst 15 laboratory
980         mouse strains. *Genome Biol* **13**, 26, doi:10.1186/gb-2012-13-4-r26 (2012).

981   25    Tan, M. H. *et al.* Dynamic landscape and regulation of RNA editing in mammals.
982         *Nature* **550**, 249-254, doi:10.1038/nature24041 (2017).

983   26    Zhao, H. Q. *et al.* Profiling the RNA editomes of wild-type *C. elegans* and ADAR
984         mutants. *Genome Res* **25**, 66-75, doi:10.1101/gr.176107.114 (2015).

985   27    St Laurent, G. *et al.* Genome-wide analysis of A-to-I RNA editing by single-molecule
986         sequencing in *Drosophila. Nature structural & molecular biology* **20**, 1333-1339,
987         doi:10.1038/nsmb.2675 (2013).

988   28    Duan, Y., Dou, S., Luo, S., Zhang, H. & Lu, J. Adaptation of A-to-I RNA editing in
989         *Drosophila. PLoS genetics* **13**, e1006648, doi:10.1371/journal.pgen.1006648 (2017).

990   29    Yu, Y. *et al.* The Landscape of A-to-I RNA Editome Is Shaped by Both Positive and
991         Purifying Selection. *PLoS genetics* **12**, e1006191, doi:10.1371/journal.pgen.1006191
992         (2016).

993   30    Zhang, R., Deng, P., Jacobson, D. & Li, J. B. Evolutionary analysis reveals regulatory
994         and functional landscape of coding and non-coding RNA editing. *PLoS genetics* **13**,
995         e1006563, doi:10.1371/journal.pgen.1006563 (2017).

996   31    Li, Q. *et al.* Caste-specific RNA editomes in the leaf-cutting ant *Acromyrmex echinatior.
997         Nat Commun* **5**, 4943, doi:10.1038/ncomms5943 (2014).

998   32    Porath, H. T. *et al.* RNA editing is abundant and correlates with task performance in a
999         social bumblebee. *Nature communications* **10**, 1605 (2019).

1000  33    Alon, S. *et al.* The majority of transcripts in the squid nervous system are extensively
1001        recoded by A-to-I RNA editing. *eLife* **4**, doi:10.7554/eLife.05198 (2015).

1002  34    Liscovitch-Brauer, N. *et al.* Trade-off between Transcriptome Plasticity and Genome
1003        Evolution in Cephalopods. *Cell* **169**, 191-202 e111, doi:10.1016/j.cell.2017.03.025
1004        (2017).

1005  35    Xu, G. & Zhang, J. Human coding RNA editing is generally nonadaptive. *Proceedings
1006        of the National Academy of Sciences of the United States of America* **111**, 3769-3774,
1007        doi:10.1073/pnas.1321745111 (2014).

1008  36    Grice, L. F. & Degnan, B. M. The origin of the ADAR gene family and animal RNA
1009        editing. *BMC evolutionary biology* **15**, 4, doi:10.1186/s12862-015-0279-3 (2015).

1010  37    Albertin, C. B. *et al.* The octopus genome and the evolution of cephalopod neural and
1011        morphological novelties. *Nature* **524**, 220-224, doi:10.1038/nature14668 (2015).

1012  38    Lang, B. F., O'Kelly, C., Nerad, T., Gray, M. W. & Burger, G. The closest unicellular
1013        relatives of animals. *Curr Biol* **12**, 1773-1778 (2002).

1014  39    Wang, Z. *et al.* RES-Scanner: a software package for genome-wide identification of
1015        RNA-editing sites. *Gigascience* **5**, 37, doi:10.1186/s13742-016-0143-4 (2016).

1016   40   Porath, H. T., Carmi, S. & Levanon, E. Y. A genome-wide map of hyper-edited RNA
1017        reveals numerous new sites. *Nat Commun* **5**, 4726, doi:10.1038/ncomms5726 (2014).

1018   41   Srivastava, M. *et al.* The *Trichoplax* genome and the nature of placozoans. *Nature* **454**,
1019        955-960, doi:10.1038/nature07191 (2008).

1020   42   Porath, H. T., Knisbacher, B. A., Eisenberg, E. & Levanon, E. Y. Massive A-to-I RNA
1021        editing is common across the Metazoa and correlates with dsRNA abundance. *Genome*
1022        *Biol* **18**, 185, doi:10.1186/s13059-017-1315-y (2017).

1023   43   Carmichael, G. G. Antisense starts making more sense. *Nature biotechnology* **21**, 371-
1024        372, doi:10.1038/nbt0403-371 (2003).

1025   44   Neeman, Y., Dahary, D., Levanon, E. Y., Sorek, R. & Eisenberg, E. Is there any sense
1026        in    antisense    editing?    *Trends    in    genetics    :    TIG*    **21**,    544-547,
1027        doi:10.1016/j.tig.2005.08.005 (2005).

1028   45   Morita, Y. *et al.* Human endonuclease V is a ribonuclease specific for inosine-
1029        containing RNA. *Nat Commun* **4**, 2273, doi:10.1038/ncomms3273 (2013).

1030   46   Mah, J. L. & Leys, S. P. Think like a sponge: The genetic signal of sensory cells in
1031        sponges. *Dev Biol* **431**, 93-100, doi:10.1016/j.ydbio.2017.06.012 (2017).

1032   47   Marks, F., Klingmüller, U. & Müller-Decker, K. *Cellular signal processing: an*
1033        *introduction to the molecular mechanisms of signal transduction*. (Garland Science,
1034        2008).

1035   48   Stulic, M. & Jantsch, M. F. Spatio-temporal profiling of Filamin A RNA-editing reveals
1036        ADAR preferences and high editing levels outside neuronal tissues. *RNA biology* **10**,
1037        1611-1617, doi:10.4161/rna.26216 (2013).

1038   49   Gray, M. W. Evolutionary origin of RNA editing. *Biochemistry* **51**, 5235-5242,
1039        doi:10.1021/bi300419r (2012).

1040   50   Bian, Z., Ni, Y., Xu, J. R. & Liu, H. A-to-I mRNA editing in fungi: occurrence, function,
1041        and evolution. *Cell Mol Life Sci* **76**, 329-340, doi:10.1007/s00018-018-2936-3 (2019).

1042   51   Erwin, D. H. *et al.* The Cambrian conundrum: early divergence and later ecological
1043        success    in    the    early    history    of    animals.    *Science*    **334**,    1091-1097,
1044        doi:10.1126/science.1206375 (2011).

1045   52   Merkle, T. *et al.* Precise RNA editing by recruiting endogenous ADARs with antisense
1046        oligonucleotides. *Nature biotechnology* **37**, 133-138, doi:10.1038/s41587-019-0013-6
1047        (2019).

1048   53   Qu, L. *et al.* Programmable RNA editing by recruiting endogenous ADAR using
1049        engineered RNAs. *Nature biotechnology*, doi:10.1038/s41587-019-0178-z (2019).

1050   54   Levine, M. & Tjian, R. Transcription regulation and animal diversity. *Nature* **424**, 147-
1051        151, doi:10.1038/nature01763 (2003).

1052  55  Chen, M. & Manley, J. L. Mechanisms of alternative splicing regulation: insights from
1053      molecular and genomics approaches. *Nat Rev Mol Cell Biol* **10**, 741-754,
1054      doi:10.1038/nrm2777 (2009).

1055  56  Chen, L., Bush, S. J., Tovar-Corona, J. M., Castillo-Morales, A. & Urrutia, A. O.
1056      Correcting for differential transcript coverage reveals a strong relationship between
1057      alternative splicing and organism complexity. *Mol Biol Evol* **31**, 1402-1413,
1058      doi:10.1093/molbev/msu083 (2014).

1059  57  Liddicoat, B. J. *et al.* RNA editing by ADAR1 prevents MDA5 sensing of endogenous
1060      dsRNA as nonself. *Science* **349**, 1115-1120, doi:10.1126/science.aac7049 (2015).

1061  58  Misof, B. *et al.* Phylogenomics resolves the timing and pattern of insect evolution.
1062      *Science* **346**, 763-767, doi:10.1126/science.1257570 (2014).

1063  59  Fletcher, D. A. & Mullins, R. D. Cell mechanics and the cytoskeleton. *Nature* **463**, 485-
1064      492, doi:10.1038/nature08908 (2010).

1065  60  Ramaswami, G. *et al.* Accurate identification of human Alu and non-Alu RNA editing
1066      sites. *Nat Methods* **9**, 579-581, doi:10.1038/nmeth.1982 (2012).

1067  61  Chen, J., Li, X., Zhong, H., Meng, Y. & Du, H. Systematic comparison of germline
1068      variant calling pipelines cross multiple next-generation sequencers. *Sci Rep* **9**, 9345,
1069      doi:10.1038/s41598-019-45835-3 (2019).

1070  62  Patch, A. M. *et al.* Germline and somatic variant identification using BGISEQ-500 and
1071      HiSeq X Ten whole genome sequencing. *PloS one* **13**, e0190264,
1072      doi:10.1371/journal.pone.0190264 (2018).

1073  63  Mak, S. S. T. *et al.* Comparative performance of the BGISEQ-500 vs Illumina
1074      HiSeq2500 sequencing platforms for palaeogenomic sequencing. *Gigascience* **6**, 1-13,
1075      doi:10.1093/gigascience/gix049 (2017).

1076  64  Chen, Y. *et al.* SOAPnuke: a MapReduce acceleration-supported software for
1077      integrated quality control and preprocessing of high-throughput sequencing data.
1078      *Gigascience* **7**, 1-6, doi:10.1093/gigascience/gix120 (2018).

1079  65  Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant
1080      detection and genome assembly improvement. *PloS one* **9**, e112963,
1081      doi:10.1371/journal.pone.0112963 (2014).

1082  66  Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-
1083      MEM. *arXiv preprint arXiv:1303.3997* (2013).

1084  67  Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler
1085      transform. *Bioinformatics* **25**, 1754-1760, doi:10.1093/bioinformatics/btp324 (2009).

1086  68  Parkhomchuk, D. *et al.* Transcriptome analysis by strand-specific sequencing of
1087      complementary DNA. *Nucleic acids research* **37**, e123, doi:10.1093/nar/gkp596 (2009).

1088  69  Benjamini, Y., Drai, D., Elmer, G., Kafkafi, N. & Golani, I. Controlling the false
1089      discovery rate in behavior genetics research. *Behavioural brain research* **125**, 279-284
1090      (2001).

1091  70  Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low
1092      memory requirements. *Nat Methods* **12**, 357-360, doi:10.1038/nmeth.3317 (2015).

1093  71  Robinson, M. D. & Oshlack, A. A scaling normalization method for differential
1094      expression analysis of RNA-seq data. *Genome Biol* **11**, R25, doi:10.1186/gb-2010-11-
1095      3-r25 (2010).

1096  72  Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from
1097      RNA-seq reads. *Nature biotechnology* **33**, 290-295, doi:10.1038/nbt.3122 (2015).

1098  73  Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
1099      alignment search tool. *Journal of molecular biology* **215**, 403-410, doi:10.1016/S0022-
1100      2836(05)80360-2 (1990).

1101  74  Birney, E., Clamp, M. & Durbin, R. GeneWise and Genomewise. *Genome Res* **14**, 988-
1102      995, doi:10.1101/gr.1865504 (2004).

1103  75  Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular Evolutionary Genetics
1104      Analysis Version 7.0 for Bigger Datasets. *Mol Biol Evol* **33**, 1870-1874,
1105      doi:10.1093/molbev/msw054 (2016).

1106  76  Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**,
1107      1586-1591, doi:10.1093/molbev/msm088 (2007).

1108  77  Yang, Z. & Nielsen, R. Estimating synonymous and nonsynonymous substitution rates
1109      under realistic evolutionary models. *Mol Biol Evol* **17**, 32-43,
1110      doi:10.1093/oxfordjournals.molbev.a026236 (2000).

1111  78  Vacic, V., Iakoucheva, L. M. & Radivojac, P. Two Sample Logo: a graphical
1112      representation of the differences between two sets of sequence alignments.
1113      *Bioinformatics* **22**, 1536-1537, doi:10.1093/bioinformatics/btl151 (2006).

1114  79  Cai, H. *et al.* A draft genome assembly of the solar-powered sea slug *Elysia chlorotica*.
1115      *Sci Data* **6**, 190022, doi:10.1038/sdata.2019.22 (2019).

1116  80  Smit, A. F., Hubley, R. & Green, P. Available fom http://www.repeatmasker.org. (26
1117      July 2018 date last accessed).

1118  81  Smit, A. & Hubley, R. Available fom http://www.repeatmasker.org/RepeatModeler/.
1119      (26 July 2018 date last accessed).

1120  82  Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic acids*
1121      *research* **27**, 573-580 (1999).

1122  83  Katoh, K., Misawa, K., Kuma, K. & Miyata, T. MAFFT: a novel method for rapid
1123      multiple sequence alignment based on fast Fourier transform. *Nucleic acids research*
1124      **30**, 3059-3066, doi:10.1093/nar/gkf436 (2002).

1125   84   Emms, D. M. & Kelly, S. OrthoFinder: solving fundamental biases in whole genome
1126        comparisons dramatically improves orthogroup inference accuracy. *Genome Biol* **16**,
1127        157, doi:10.1186/s13059-015-0721-2 (2015).

1128   85   Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high
1129        throughput. *Nucleic acids research* **32**, 1792-1797, doi:10.1093/nar/gkh340 (2004).

1130   86   Castresana, J. Selection of conserved blocks from multiple alignments for their use in
1131        phylogenetic     analysis.     *Mol     Biol     Evol*     **17**,     540-552,
1132        doi:10.1093/oxfordjournals.molbev.a026334 (2000).

1133   87   Laumer, C. E. *et al.* Revisiting metazoan phylogeny with genomic sampling of all phyla.
1134        *Proceedings.  Biological  sciences  /  The  Royal  Society* **286**, 20190831,
1135        doi:10.1098/rspb.2019.0831 (2019).

1136   88   Ryan, J. F. *et al.* The genome of the ctenophore Mnemiopsis leidyi and its implications
1137        for cell type evolution. *Science* **342**, 1242592, doi:10.1126/science.1242592 (2013).

1138   89   Simion, P. *et al.* A Large and Consistent Phylogenomic Dataset Supports Sponges as
1139        the   Sister   Group   to   All   Other   Animals.   *Curr   Biol* **27**,   958-967,
1140        doi:10.1016/j.cub.2017.02.031 (2017).
1141
1142

## Figure Legends

### Figure 1



**Figure 1 | The origin and variation of RNA editing in metazoans.**

(**a**) The phylogeny of the 22 species examined in this study, with the inferred origin and lineage-specific loss of *ADARs* indicated. The topology of the phylogenetic tree is derived according to previous reports [87-89]. The copy number of *ADAR-like* genes identified in the genome of each species is present in parenthesis after the Latin name. Asterisks (*) indicate species that have not previously been subject to transcriptome-wide RNA editing analyses. Full names for the 22 species from top to bottom are *Sphaeroforma arctica* (ichthyosporean), *Capsaspora owczarzaki* (filasterean), *Salpingoeca rosetta* (choanoflagellate), *Monosiga brevicollis* (choanoflagellate), *Mnemiopsis leidyi* (comb jelly), *Amphimedon queenslandica* (sponge), *Trichoplax adhaerens* (placozoan), *Hydra vulgaris* (hydra), *Nematostella vectensis*

(sea anemone), *Aplysia californica* (sea hare), *Crassostrea gigas* (oyster), *Octopus bimaculoides* (octopus), *Caenorhabditis elegans* (roundworm), *Acromyrmex echinatior* (ant), *Drosophila melanogaster* (fruit fly), *Drosophila simulans* (fruit fly), *Strongylocentrotus purpuratus* (sea urchin), *Ptychodera flava* (acorn worm), *Branchiostoma belcheri* (lancelet), *Ciona savignyi* (sea squirt), *Danio rerio* (zebrafish) and *Homo sapiens* (human). (**b**) The occurrence rate of RNA editing in each species, which was measured as the number of RNA-editing sites per million transcribed genomic sites (RNA depth $\geq$ 2X). The mean number of editing sites identified in each species is presented on the right of each dot. (**c**) Percentage of editing sites across the 12 possible types of nucleotide substitutions. Error bars in panels **b** and **c** represent s.d. across samples.

## Figure 2



**Figure 2 | The structure and sequence preferences of A-to-I editing in metazoans.**

(**a**) Percentage of A-to-I editing sites locating in dsRNA regions potentially formed by intramolecular folding of pre-mRNA, measured as the proportion of sites locating in a region (± 200 nt centered on the edited adenosine) that shows a reverse-complement alignment (identity ≥ 80% and aligned length ≥ 50 nt) within its flanking sequences (± 2 knt centered on the edited adenosine). Control sites were randomly selected transcribed adenosines with the same number and comparable RNA depth of the A-to-I editing sites in each sample of each species (see Methods for details). (**b**) Percentage of A-to-I editing sites locating in dsRNA regions potentially formed by intermolecular hybridization of sense-antisense transcripts, measured as the proportion of sites locating in a region (± 50 nt centered on the edited adenosine) with transcription signal (RNA depth ≥ 2X along >50% of the region) in the both strands. Control sites were the same set of adenosine sites used in panel **a**. (**c**) Percentage of A-

to-I editing sites locating in dsRNA regions, either potentially formed by intramolecular folding of pre-mRNA, or by intermolecular hybridization of sense-antisense transcripts or by both mechanisms. (**d**) Percentage of A-to-I editing sites occurring in clusters. A cluster contains $\geq 3$ A-to-I editing sites of which the distance for two adjacent sites $\leq 30$ nt. Control sites were the same set of adenosine sites used in panel **a**. (**e**) Neighboring nucleotide preference of the edited adenosines in comparison to the unedited transcribed adenosines within $\pm 50$ nt of the edited adenosines. For the lineages with more than one representative species, the same number of editing sites were randomly selected according to the species with the lowest number of editing sites in this lineage. Nucleotides were plotted using the size of the nucleotide that was proportional to the difference between the edited and unedited datasets, with the upper part presenting enriched nucleotides in the edited dataset and lower part presenting depleted nucleotides. (**f**) The relative frequency of the 16 nucleotide triplets with adenosine in the center subject to A-to-I editing, measured as the frequency of a triplet observed among all edited adenosines against the frequency of this triplet observed among all neighboring unedited adenosines within $\pm 50$ nt of the edited adenosines. Boxplot shows the distribution across the 17 metazoans with *ADARs*, and nucleotide triplets are ordered according to highest median value across species to the lowest. Error bars in panels **a, b** and **d** represent s.d. across samples, and asterisks indicate significance levels estimated by two-tailed paired t-tests with "∗" representing $p < 0.05$, "∗∗" $< 0.01$ and "∗∗∗" $< 0.001$.
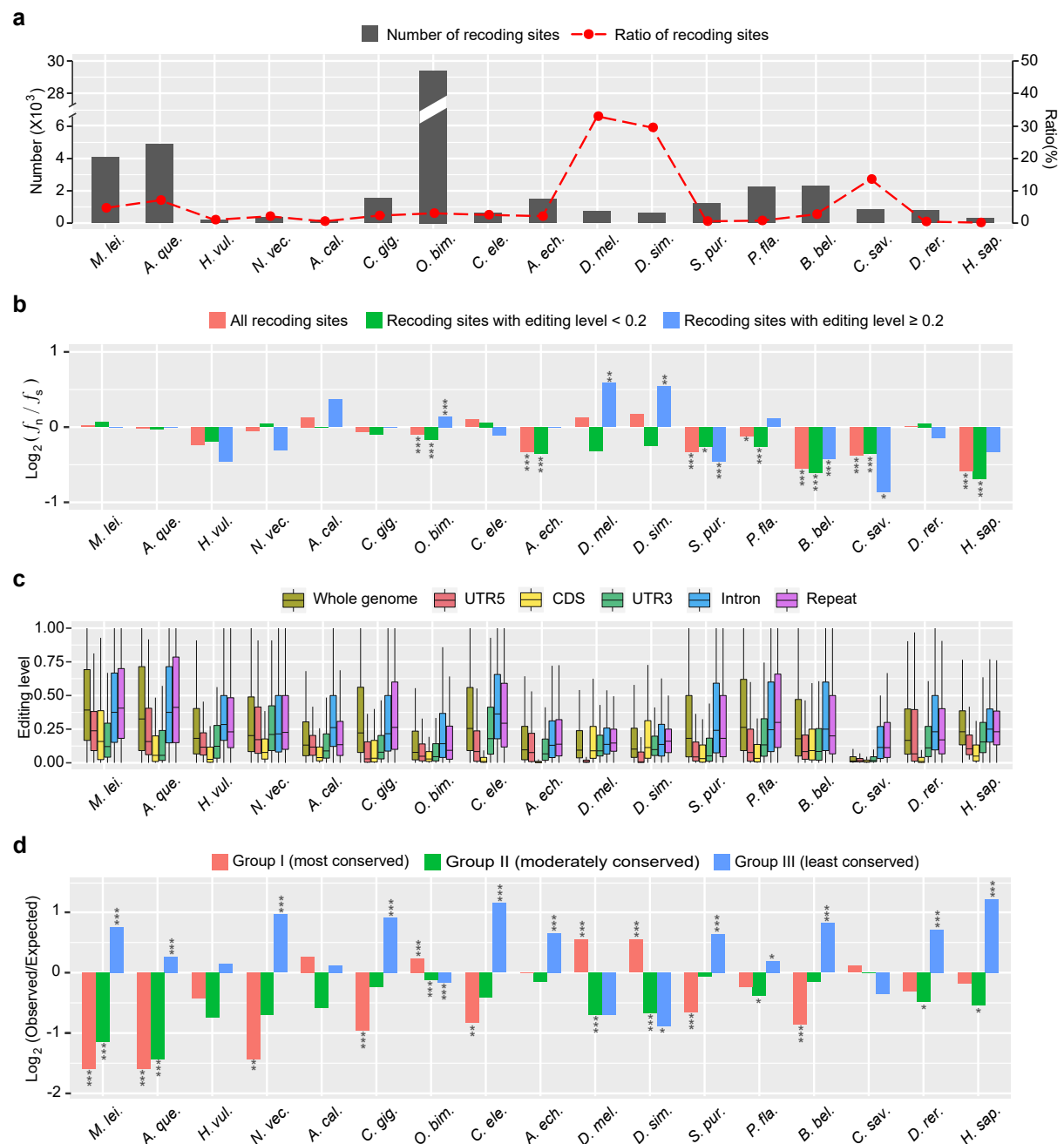
**Figure 3**



**Figure 3 | The primary genomic targets of metazoan A-to-I editing.**

(**a**) The proportion of A-to-I editing sites versus total A-to-I editing sites in different genomic regions. Genic regions include untranslated regions (5'-UTR and 3'UTR), CDS and intron of all protein-coding genes. Repeats include transposons and tandem repeats annotated for each species in this study (see Methods). (**b**) Comparison of A-to-I editing density across different genomic elements in each species. Editing density of an element was calculated as the number of A-to-I editing sites locating in this element divided by the number of transcribed adenosines (RNA depth ≥ 2X) in this element. (**c**) Comparison of editing density across 5'-UTR, CDS and 3'-UTR. Note the different scales used between panel **b** and **c** in order to show the difference of editing density among genic elements. (**d**) The negative correlation between the sequence divergence and editing density of repetitive elements.
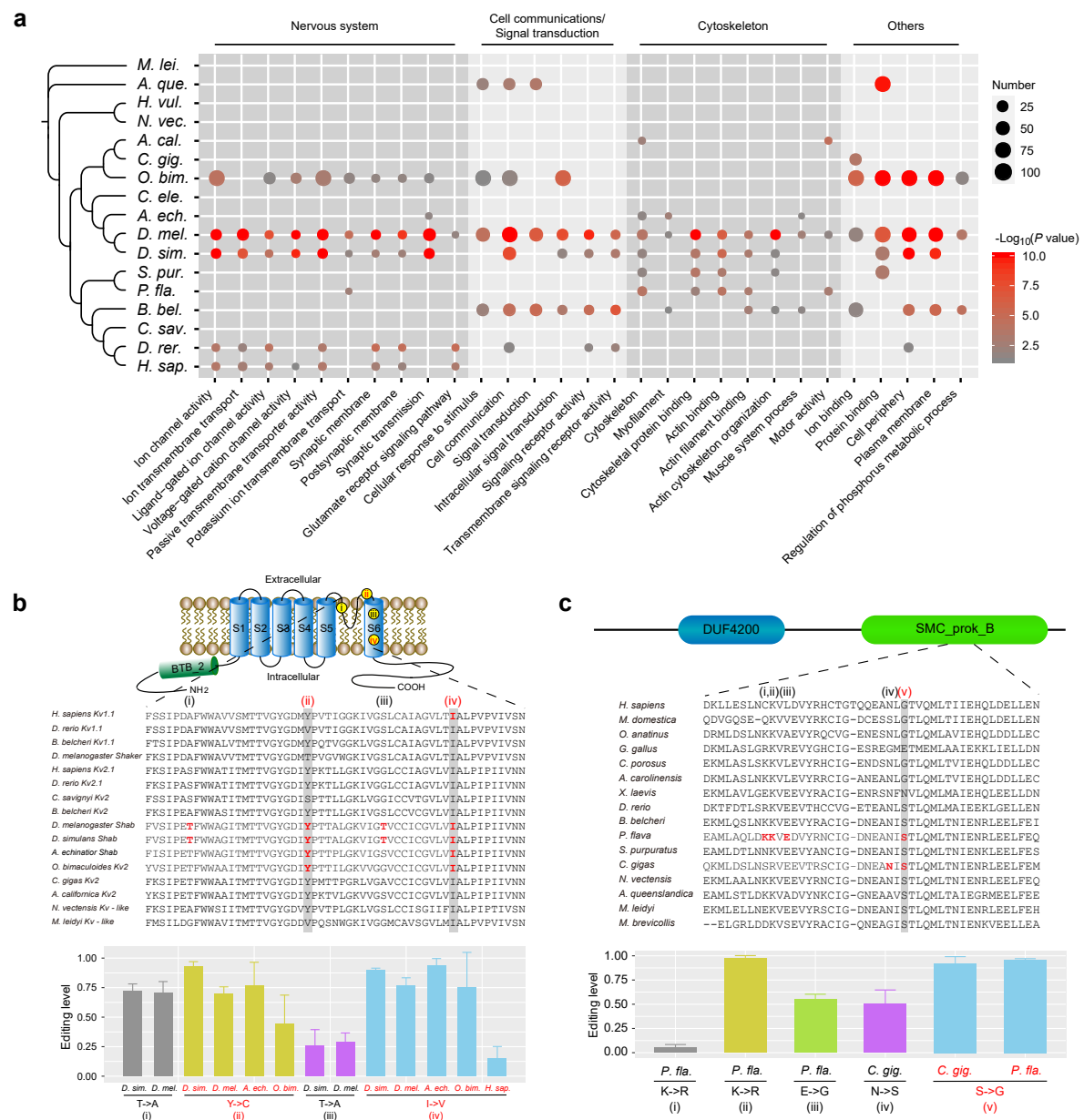
**Figure 4**



**Figure 4 | The prevalence and adaptive potential of recoding editing in metazoans.**

(**a**) The number (left y-axis) and proportion (right y-axis) of recoding sites versus total A-to-I editing sites in each species. (**b**) The frequency of nonsynonymous editing ($f_n$) versus synonymous editing ($f_s$) for all A-to-I editing sites, and A-to-I editing sites with low ($< 0.2$) and high ($\geq 0.2$) editing levels. $f_n$ (or $f_s$) was calculated as the number of A-to-I editing sites causing nonsynonymous (or synonymous) changes against the number of potential nonsynonymous (or synonymous) adenosine sites if A is replaced with G from the genes with $\geq 1$ editing site in their coding regions. Significance levels were estimated by two-tailed

Fisher's exact tests. (**c**) The comparison of editing levels of A-to-I editing sites in different genomic elements. For each editing site in a species, the mean editing level across samples with RNA depth $\geq$ 10X was calculated and used in this analysis. (**d**) The observed/expected number of genes subject to recoding editing among genes with different conservation levels. The expected probability of a gene being recoded in a species was estimated as the number of recoded genes (i.e. genes with $\geq$ 1 recoding site) out of all transcribed protein-coding genes (i.e. RPKM > 1 in at least one sample) in this species, and the expected number of recoded genes in each conservation group was calculated as the number of genes in this group multiplied by the expected probability of a gene being recoded. Significance levels were estimated by two-tailed binomial tests. Asterisks in panels **b** and **d** indicate the significance levels with "∗" representing $p < 0.05$, "∗∗" < 0.01 and "∗∗∗" < 0.001.
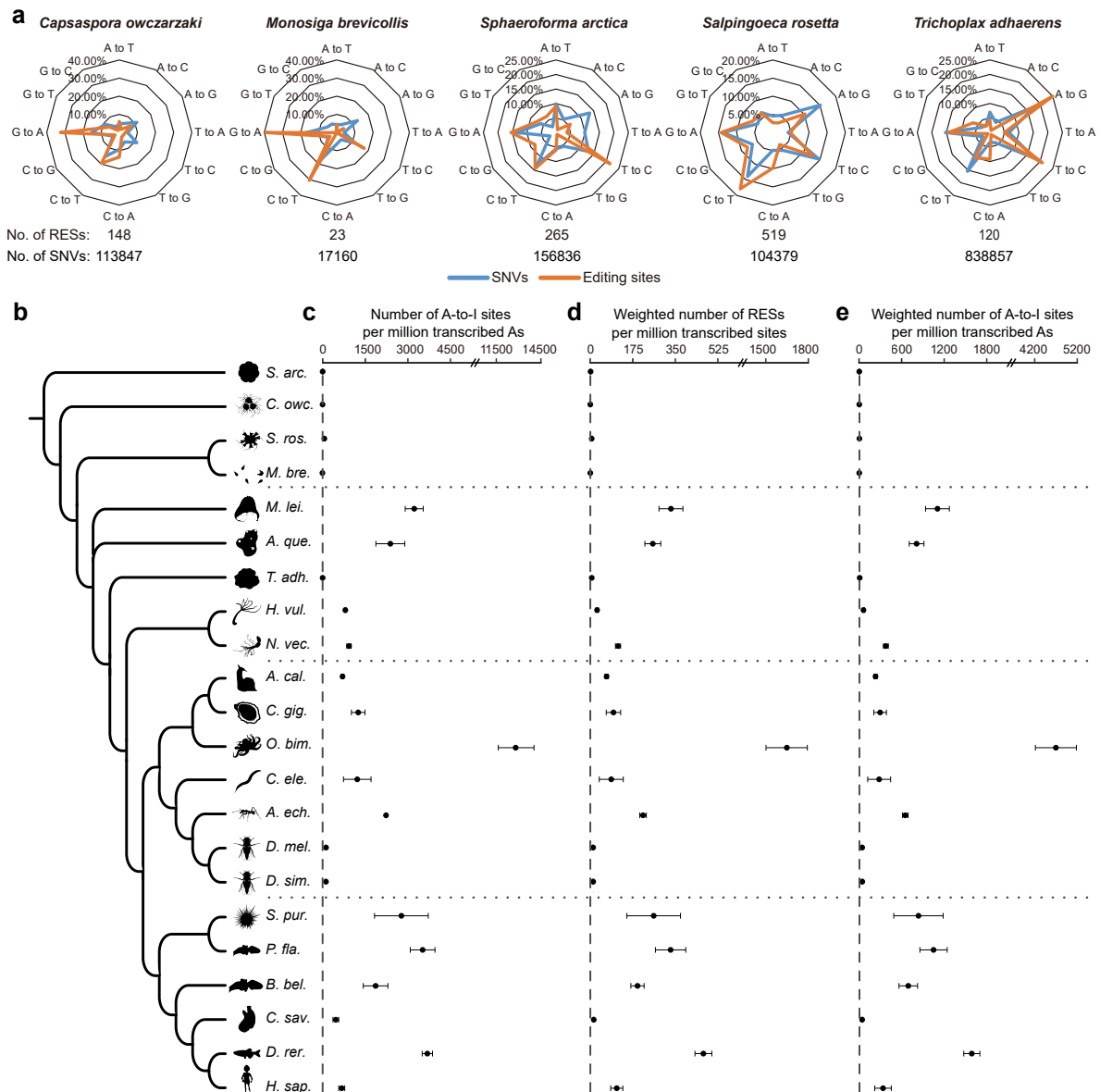
**Figure 5**



**Figure 5 | Functional preference of recoding editing in metazoans.**

(**a**) Functional categories that are enriched by recoded genes in no less than three species (Hypergeometric test adjusted $p < 0.05$). Only genes with at least one recoding site of which the average editing level across samples > 0.1 or shared by at least two samples in each species were used for the functional enrichment analysis (see Methods for details). (**b**) An example showing the convergent evolution of recoding editing in the voltage-gated $K^+$ (Kv) channels among *Drosophila*, ant, octopus and human. The upper part shows the classic structure of the Kv channel within the cell membrane that contains six membrane-spanning domains (S1–S6). Yellow dots indicate the locations of recoding sites. The middle part shows the multiple sequence alignment of the region containing the four recoding sites, with recoding sites
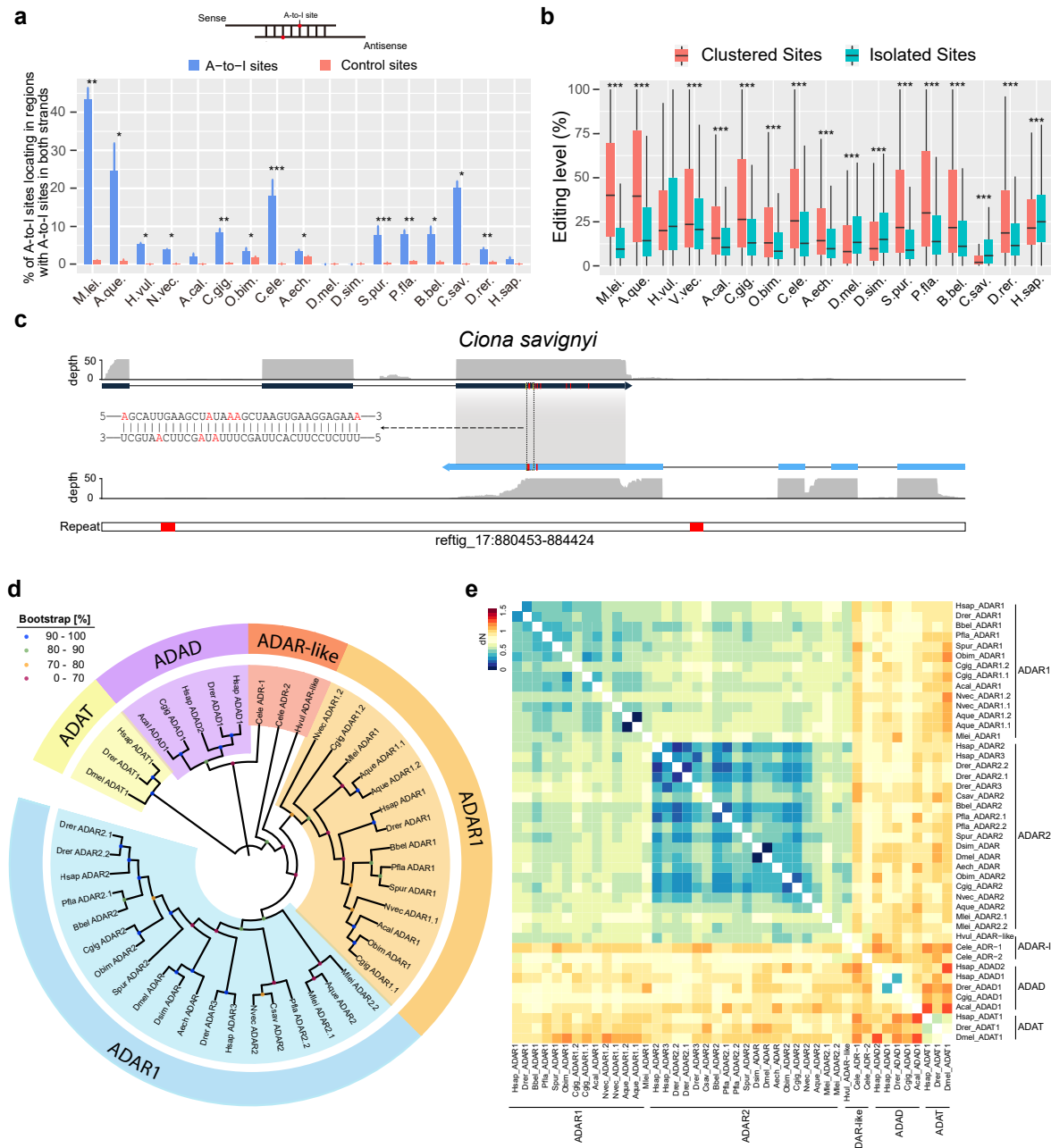
identified in each species highlighted by red color and the conserved recoding events shared across phyla highlighted by gray shadow. The lower part shows the editing levels of the recoding events in each species. (**c**) An example showing the convergent evolution of a recoding editing event in the same amino acid residue of the cytoskeleton-related gene *CFAP100* between the oyster *C. gigas* and acorn worm *P. flava*. The upper part shows the domain organization of the oyster CFAP100 protein (XP_011420958.1) annotated by the Conserved Domain Database (CDD) of NCBI. The middle and lower part are similar with panel **b**. Error bars in panels **b** and **c** represent s.d. across samples.

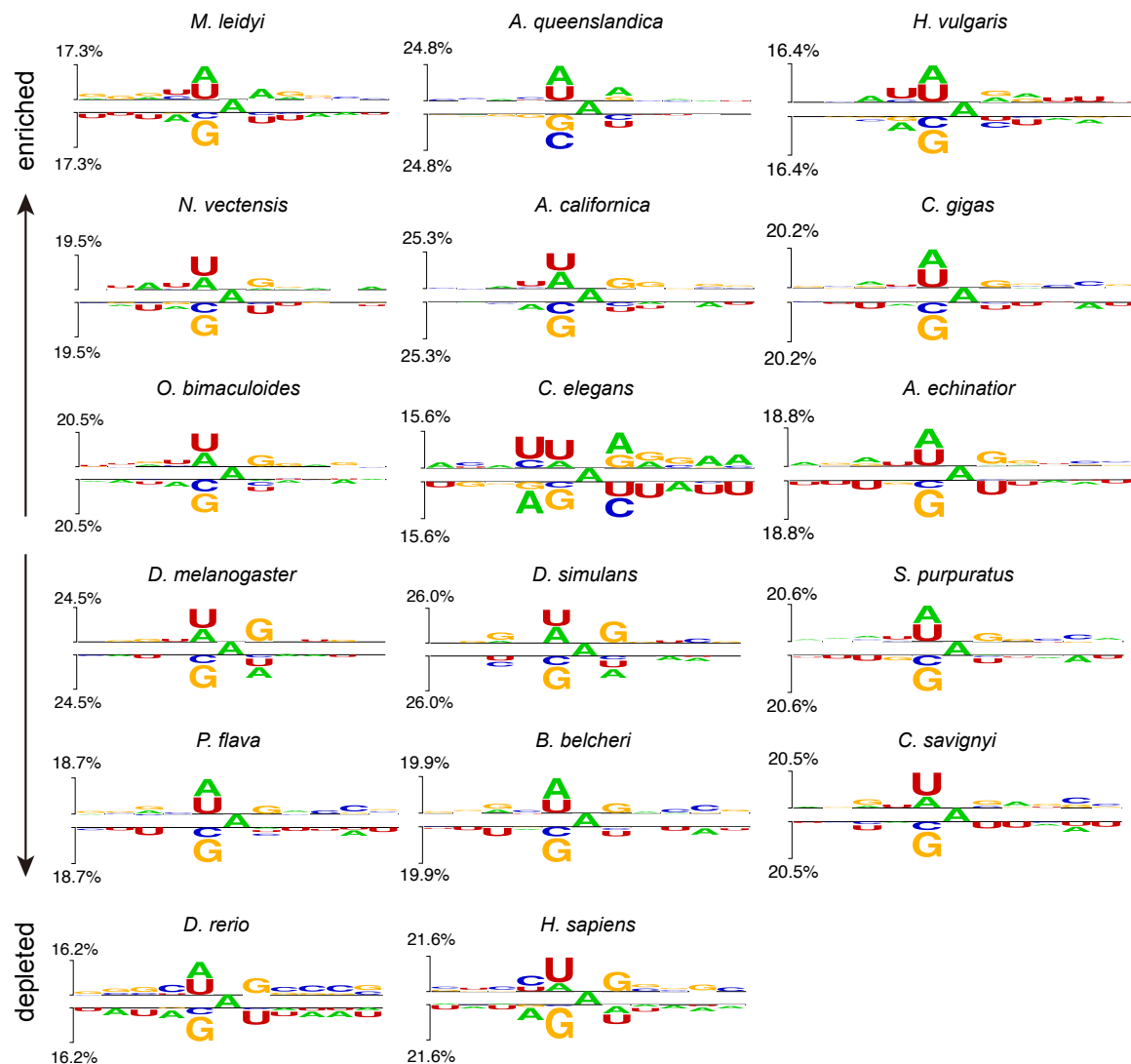**Supplementary Figure 1 | The origin and variation of RNA editing in metazoans.**
(**a**) The frequency of each type of nucleotide substitution among all potential RNA-editing sites (RESs) and genetic single nucleotide variants (SNVs) identified in the five species without *ADARs*. The types of nucleotide substitutions for RESs and SNVs were both inferred according to the genotypes present in the plus strand of the reference genome in this analysis. That is, an A-to-G RES from the minus strand of the reference genome was regarded as a T-to-C substitution, while substitution types of RESs from the plus strand remained unchanged. The RESs and SNVs from different samples of the same species were first combined according to their genomic locations, respectively, before the frequency calculation. (**b**) The phylogeny of the 22 species examined in this study. The topology of the phylogenetic tree is derived according to previous reports [1,2]. (**c**) The occurrence rate of A-to-I editing in each species, which was calculated as the number of A-to-I sites divided by the number of transcribed adenosine sites (RNA depth ≥ 2X) then multiplied by one million. (**d**) The editing-level-weighted occurrence rate of RNA editing in each species, which was calculated as the summed editing level for all RNA-editing sites with RNA depth ≥ 10X divided by the number of transcribed genomic sites with RNA depth ≥ 10X then multiplied by one million. (**e**) The editing-level-weighted occurrence rate of A-to-I editing in each species, which was calculated as the summed editing level for all A-to-I sites with RNA depth ≥ 10X divided by the number of transcribed adenosine sites with RNA depth ≥ 10X then multiplied by one million. Error bars in panels **c**, **d** and **e** represent s.d. across samples.

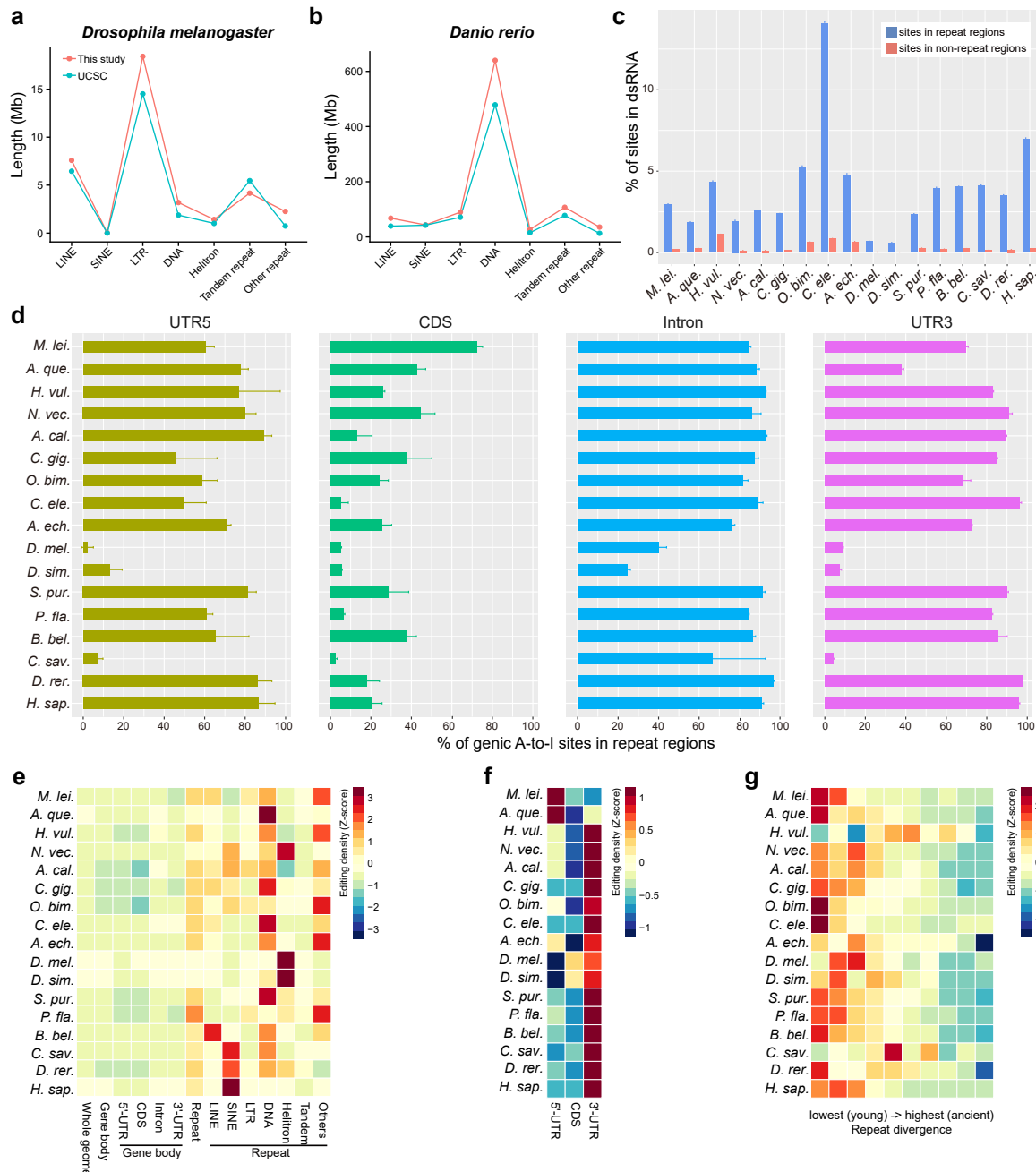**Supplementary Figure 2 | The structure and sequence preferences of A-to-I editing in metazoans.**
(**a**) The proportion of A-to-I editing sites locating in regions targeted by RNA editing on both strands, measured as the proportion of sites locating in a region ($\pm$ 50 nt centered on the edited adenosine) with at least one A-to-I editing site found on the opposite strand. Control sites were randomly selected transcribed adenosines with the same number and comparable RNA depth of the A-to-I editing sites in each sample of each species (see Methods). The asterisks indicate significance levels estimated by two-tailed paired t-tests with "*" representing $p < 0.05$, "**" < 0.01 and "***" < 0.001. (**b**) The comparison of editing levels for the clustered and isolated A-to-I editing sites in each species. The asterisks indicate significance levels estimated by two-tailed Wilcoxon signed-rank tests with "*" representing $p < 0.05$, "**" < 0.01 and "***" < 0.001. (**c**) An example of a sense-antisense transcript pairing in *Ciona savignyi*, showing the RNA coverage of both transcript models, the pairing region (grey shadow), the A-to-I editing sites on both transcripts (red vertical bars), the detail view of a 34 nt regions (dashed box) with edited adenosines highlighted in red, and the distribution of repeats in this region (red boxes in the bottom track).
(**d**) The phylogenetic trees of the deamination domains of ADARs using the neighbor-joining (NJ) method with MEGA7. Protein sequences of the deamination domains were aligned with ClustalW implemented in MEGA7.

The deamination domains of ADAT1 from *D. melanogaster*, *D. rerio* and *H. sapiens* were selected as the outgroups. All protein sequences used for the phylogenetic analyses are present in Supplementary Table 3. (**e**) Pairwise nonsynonymous substitution rates (*dN*) for any pair of ADs in panel **d** using PAML (v4.9i) with the Yang & Nielsen (2000) method.



**Supplementary Figure 3 | Neighboring nucleotide preference of edited adenosines.**
The neighboring nucleotide preference (± 5 nt) of the edited adenosines in each species was estimated in comparison to the unedited transcribed adenosines within ± 50 nt of the edited adenosines, by the Two Sample Logo software [3]. Nucleotides were plotted using the size of the nucleotide that was proportional to the difference between the edited and unedited datasets, with the upper part presenting enriched nucleotides in the edited dataset and lower part presenting depleted nucleotides.

**Supplementary Figure 4 | The primary genomic targets of metazoan A-to-I editing.**
(**a-b**) The non-redundant lengths of different repeat families in *D. melanogaster* (**a**) and *D. rerio* (**b**), according to the annotations generated in this study (see Methods) and those from UCSC, showing the good consistency between these two annotation results. (**c**) The potentials of repeat and non-repeat regions to form dsRNA in each species, measured as the ratios of repeat and non-repeat derived genomic sites locating in regions that could find a reverse-complement alignment in nearby regions (see Methods). *P*-values were estimated by Monte Carlo simulations (100 times) and < 0.01 for all species. (**d**) The percentage of genic A-to-I editing sites locating regions annotated as repeats. Genic editing sites were defined as editing sites locating in protein-coding gene related elements including 5'-UTR, CDS, intron and 3'-UTR. (**e**) Comparison of editing-level-weighted editing density across different genomic elements in each species. The weighted editing density of an element was calculated as the summed editing level of A-to-I editing sites (RNA depth ≥ 10X) locating in this element divided by the number of transcribed adenosines (RNA depth ≥ 10X) in this element. (**f**) Comparison of editing-level-weighted editing density across 5'-UTR, CDS and 3'-UTR. Note the different scales used between panel **e** and **f** in order to show the difference of editing density among genic elements. (**g**) The negative correlation between the sequence divergence and editing-level-weighted editing density of repetitive elements.

**Supplementary Figure 5 | Functional preference of recoding editing in metazoans.**
(**a**) The convergent evolution of two recoding editing events in the glutamate ionotropic receptors in human and zebrafish. The upper part shows the domain organization of the human GRIA2 protein (NP_001077088.1) annotated by the Conserved Domain Database (CDD) of NCBI. The middle part shows the multiple sequence alignments of the two regions containing the shared recoding sites. Recoding sites were highlighted by red color and the positions of conserved recoding events were highlighted by gray shadow. The lower part shows the editing levels of the recoding events in each species. (**b**) The convergent evolution of a recoding editing event in fascin (an actin filament-bundling protein), which was shared by the octopus *O. bimaculoides*, the sea urchin *S. purpuratus* and the lancelet *B. belcheri*. The upper part shows the domain organization of lancelet fascin-like protein (XP_019642030.1) annotated by the CDD of NCBI. The middle and lower part are the same as panel **a**. Error bars in panels **a** and **b** represent s.d. across samples.

## Supplementary References

1    Laumer, C. E. *et al.* Revisiting metazoan phylogeny with genomic sampling of all phyla. *Proceedings. Biological sciences / The Royal Society* **286**, 20190831, doi:10.1098/rspb.2019.0831 (2019).

2    Ryan, J. F. *et al.* The genome of the ctenophore Mnemiopsis leidyi and its implications for cell type evolution. *Science* **342**, 1242592, doi:10.1126/science.1242592 (2013).

3    Vacic, V., Iakoucheva, L. M. & Radivojac, P. Two Sample Logo: a graphical representation of the differences between two sets of sequence alignments. *Bioinformatics* **22**, 1536-1537, doi:10.1093/bioinformatics/btl151 (2006).