1 Functional coherence among miRNA targets: a potential metric for assessing biological

2 signal among target prediction methods in non-model species

3

4 Authors:

5

6 Christopher W. Wheat[1*]

7 Rachel A. Steward[1]

8 Yu Okamura[2,3]

9 Heiko Vogel[2]

10 Philipp Lehmann[1,4]

11 Kevin T. Roberts[1]

12

13

14

15 [1]Department of Zoology, Stockholm University, Stockholm, Sweden

16 [2]Department of Insect Symbiosis, Max Planck Institute for Chemical Ecology, Jena,

17 Germany

18 [3]Department of Biological Sciences, Graduate School of Science, University of Tokyo,

19 Tokyo, Japan

20 [4] Zoological Institute and Museum, Greifswald University, Greifswald, Germany

21

22 * Author for Correspondence: Christopher W. Wheat, Department of Zoology,

23 Stockholm University, Stockholm, Sweden. +46 721 958586.

24 chris.wheat@zoologis.su.se

25

26

27

28 **Abstract**

29

30 Although miRNA regulation of protein production is a likely target of adaptive evolution,

31 high false-positive rates in the identification of mRNAs targeted by miRNAs in non-

32 model species' complicates interpretation of recent advances. Here we document the

33 challenges and then outline steps for the community to address these challenges.

34

35

36

37 **Keywords**

38

39 miRNA, target detection, false-positives, functional coherence, gene set enrichment

40 analysis

41

42

43

44

45  One major revelation of the genomics era is that gene regulatory networks (GRNs)

46  exhibit extensive functional coherence, as most transcription factors regulate the

47  transcription of functionally related modules of genes, resulting in co-expressed genes

48  generally comprising coherent developmental and metabolic pathways (Stuart et al.

49  2003; Wolfe et al. 2005). GRNs are at the core of evolutionary biology studies, since it is

50  the modification of GRNs, as well as their co-option into novel developmental contexts,

51  that is the major axis upon with evolutionary adaptations and novelty arise (Bruce &

52  Patel 2020; Erwin 2021). However, mRNA transcription alone does not determine

53  protein concentrations and hence phenotypes, but rather a diverse set of dynamics,

54  including post-transcriptional and post-translational regulation, significantly modify the

55  transcriptome, forming a key feature of the genotype to phenotype map (Liu et al. 2016;

56  Bartel 2018).

57  Here we focus upon post-transcriptional regulation via microRNAs (miRNAs),

58  ~22 nucleotides (nt) long RNAs. In most animals, miRNAs are produced after

59  transcription via a series of processes (hairpin formation, cleavage, export to cytoplasm,

60  cleavage), then bound by the Argonaute protein, creating a silencing complex that

61  selectively binds mRNA based upon a short (6-8 nt) sequence seed matching between

62  the miRNA and mRNA, primarily in the 3' UTR region of mRNA transcripts, which then

63  initiates various forms of translation repression (Bartel 2018). Via this post-

64  transcriptional action regulating the mRNA to protein production relationship, miRNAs

65  play an important role in developmental progression and physiological functioning

66  (Bartel 2018; Gebert & MacRae 2019). Numerous studies over the past decade, across

67  both invertebrates and vertebrates, have found significant differential expression of

68  miRNA genes associated with adaptive phenotypes, suggesting that these "sculptors" of

69  the transcriptome play an important role in adaptive evolution (Bartel 2018; Leung &

70  Sharp 2010; Fruciano et al. 2021). However, investigating how such differential

71  expression of miRNA causally leads to adaptive phenotypes necessitates identifying the

72  mRNAs that are targeted by miRNAs, as only this allows researchers to make causal

73  connections between differential miRNA expression, protein expression changes, and

74  ultimately differential reproductive success. Unfortunately, identifying which mRNAs are

75  targeted by which miRNAs remains a complex problem (Bracken et al. 2016).

76　　　　Based upon insights from model-species (e.g. humans, flies, worms), animals

77　　are expected to have 100's of miRNA families (miRNAs that target the same canonical

78　　motif in mRNA), each of which can effectively reduce the protein production of 100's

79　　genes. In humans these numbers correspond to about 500 miRNAs, 300 of which can

80　　be placed into about 170 gene families, with each family on average

81　　posttranscriptionally repressing roughly 400 genes (Bartel 2018). From the perspective

82　　of a given mRNA sequence, nearly half of fly (~ 40%) and human (> 60%) mRNAs

83　　contain conserved miRNA binding targets, with each mRNA on average containing

84　　multiple miRNA binding sites (of the same and/or different miRNA families). Thus,

85　　across diverse taxa, miRNAs have the potential to sculpt a large faction of the

86　　transcriptome.

87　　　　　Genomic core facilities now routinely provide short RNA sequencing, enabling

88　　quantitative assessments of miRNA abundance in nearly any taxa. However, identifying

89　　the biologically meaningful targets of differentially expressed miRNA remains

90　　challenging, despite technological advances. While direct sequencing of the mRNA pool

91　　bound by the silencing complex is possible (crosslinking-immunoprecipitation-

92　　sequencing, CLIP-seq), a high concentration of cells is required, with results necessarily

93　　averaging over the diverse miRNA regulation dynamics among cells lineages. While a

94　　single cell approach has just been developed (Sekar et al. 2023), neither technique is

95　　able to identify the miRNAs directly involved.

96　　　　　As an initial, or only, foray into miRNA research, many research groups rely upon

97　　bioinformatic prediction of miRNA targets in their focal species, for initial interpretation

98　　of differential miRNA expression. In animals, miRNA binding to mRNA primarily relies

99　　upon 6 to 8 nucleotides of complimentary sequence, referred to as seed pairing. While

100　　legions of such short motifs populate the UTR regions of transcriptome, only a small

101　　fraction are involved in post-transcriptional repression (Agarwal et al. 2015, 2018;

102　　Fridrich et al. 2019). This scenario highlights the inherently challenging nature of target

103　　prediction due to the exceptional potential for statistically significant false positives

104　　(Fridrich et al. 2019), with the challenge of accurate *in silico* prediction spawning yet

105　　another bioinformatics cottage industry (~ 100 different software approaches to date

106　　(Fridrich et al. 2019; Kern et al. 2020; Ritchie et al. 2009).

107    Emerging from diverse efforts in model-species to understand miRNA post-

108    transcriptional regulation comes the robust result that signatures of evolutionary

109    conservation, generated due to consistent purifying selection acting over 10 to 100's of

110    millions of years, provides a powerful means of discriminating functionally important

111    seed regions from other candidates in the dynamically evolving UTR regions of mRNA.

112    Indeed, compared to using only identified motifs in a single species, or in combination

113    with various ways of modeling local thermodynamics, only approaches incorporating

114    evolutionary conservation appear accurate (Friedman et al. 2009; Agarwal et al. 2015),

115    though the field continues to explore additional parameters and approaches (Kern et al.

116    2020). Of direct relevance to this journal's readership, the prediction tools most

117    commonly employed by the ecology and evolution, non-model species community are

118    those using data from only one species without information on evolutionary

119    conservation, which exhibit false-positives rates approaching 50% or fail to identify true-

120    positives in well verified experiments (e.g. miRanda, RNAhybrid; (Agarwal et al. 2015;

121    Pinzón et al. 2017; Fridrich et al. 2019; Krüger & Rehmsmeier 2006)).

122    These observations thereby suggest that our community faces extensive

123    challenges, not only when hypothesizing about the potential range of functional impacts

124    of differentially expressed miRNAs, but when trying to conduct functional validation

125    studies. Currently, it is common to see studies intersecting miRNA expression patterns

126    with RNAseq results, scanning for inverse relationships. Unfortunately, finding

127    meaningful negative correlations between miRNA and mRNA levels is likely to

128    challenging, as the power of such correlations depends upon the number of time points

129    in comparison and the accuracy of identified miRNA-mRNA interactions. Given that

130    each miRNA can have hundreds of predicted targets, we fear that without a

131    substantially large dataset of such comparison across tissues and timepoints, such

132    efforts will always be beset by high false-positive rates. In sum, the aforementioned

133    issues highlight the need for an external means of assessing the accuracy of miRNA

134    target set prediction, especially one that could be used by the non-model species

135    community.

136    Here we present rational for an external means of assessing the accuracy of

137    miRNA target set prediction. We take as our starting point that the regulatory network of

138    miRNAs is non-random, as miRNA targets are significantly higher than expected in

139    genes having positive regulatory motifs and being highly-connected GRN components,

140    such as transcription factors (Cui et al. 2006; Bracken et al. 2016). Co-expressed

141    miRNAs, whether co-localized or not, have also been found to target specific genes and

142    pathways (Lee et al. 2012; Xu & Wong 2008; Bracken et al. 2016). Additionally,

143    individual miRNA gene families have been found to exhibit functional coherence in the

144    genes they target (Tsang et al. 2010). Indeed, the functional coherence of mRNA

145    targets is itself central to resolving the paradox between the small post-transcriptional

146    effect of miRNAs upon individual genes and the larger phenotypic effects of miRNAs, as

147    miRNA action upon multiple steps of a pathway is expected to culminate in larger

148    phenotypic impacts (Bracken et al. 2016). However, currently little is known about the

149    extent of such functional coherence across miRNA gene families as a whole. Specially,

150    we can find no global scale analyses of the functional coherence of individual miRNA

151    targets in species other than humans within a disease context (Bracken et al. 2016;

152    Gusev 2008), highlighting the lack of a general understanding of how such coherence

153    varies among taxa. Nevertheless, identifying a signature of functional coherence,

154    beyond informing on the miRNA GRN and how it evolves, could provide a biologically

155    informative metric for assessing *de novo* target predictions in novel taxa.

156          Our work here began with trying to identify the miRNA targets in a novel species,

157    the Green-veined White butterfly *Pieris napi* (Lepidoptera, Pieridae). Ultimately our goal

158    was to identify the miRNAs involved in the different states of diapause progression, but

159    in order to understand patterns of differentially expressed miRNAs, we needed to

160    identify their potential targets in the transcriptome. We present a comparison of different

161    miRNA prediction approaches, finding that only an approach incorporating information

162    on evolutionary constraint results in a detectable functional coherence among the

163    targets per miRNA. In order to validate this finding, we present evidence using miRNA

164    target predictions across model and non-model species that animals generally exhibit

165    extensive functional coherence across miRNA gene families. Therefore, functional

166    coherence provides a biologically informative metric for assessing *de novo* target

167    predictions in novel taxa that could greatly facilitate ability of the ecology and

168    evolutionary genomics community to make logical connections between miRNA to

169    relevant protein expression changes and their eventual phenotypic impacts.

170

171    **Methods**

172    ***Samples, processing, miRNA identification***

173    Data generation, from collection to sequencing through to miRNA gene and seed

174    identification was performed previously (Roberts et al., in review). Although readers are

175    directed to this other work for methodological details (Roberts et al., in review), here

176    they are briefly presented for clarity. A total of 73 samples were taken throughout pupal

177    progression (12 timepoints (0, 3, 6 days direct development; 0,3,6,24,114,144,155 days

178    diapause development), for each of 2 tissues (head, abdomen), each with 3-4 biological

179    replicates). After library construction using Illumina small RNA library kits they were

180    sequenced using HiSeq 2500 50SR, generating an average of 6.9 M reads / library. The

181    miRTrace pipeline was used to check data quality (v1.0.1; (Kang et al. 2018)),

182    contamination and taxonomic bias, followed by filtering and adapter removal (Roberts et

183    al., in review). Using miRDeep2 processing scripts (Friedlander et al. 2011), reads

184    greater than 17bp were mapped against the chromosomal level assembly for *P. napi*

185    genome GCA_905231885.1 (Lohse, Hayward, et al. 2021), with miRNAs detected using

186    *Bombyx mori* and *Heliconius melpomene* as reference miRNA sets.

187

188    ***Target identification***

189    miRNA targets were identified using two separate approaches, the first relying primarily

190    upon evolutionary conservation and the second using data from a single species. Our

191    first approach aligned genomes of 6 species of Pieridae using the software Progressive

192    Cactus (Armstrong et al. 2020), each increasing evolutionary distance from our focal

193    species *P. napi*, which was used as the reference (*P. napi* (GCA_905231885.1; (Lohse,

194    Hayward, et al. 2021), *P. rapae* (GCA_905147795.1; (Lohse, Ebdon, et al. 2021))*, P.*

195    *brassicae* (GCA_905147105.1; (Lohse, Mackintosh, et al. 2021))*, P. macdunnoughii*

196    (Steward et al. 2021). The last two of these 6 genomes were high-quality draft

197    assemblies, using MaSuRCA (Zimin et al. 2017) for genome assembly of Oxford

198    Nanopore Sequencing data and and Illumina short read data for *P. melete*

199   (PRJEB59056, 376 contigs, 320 Mbp, N50 2.6 Mbp, BUSCO: CS:94.1%, CD:4.4%,

200   F:0.3%, M:1.2% (BUSCO v. 5.5.0 (Manni et al. 2021), n:5286, lepidoptera_odb10), and

201   using Flye ver. 2.7 (Kolmogorov et al. 2019) for *Pontia daplidice* (PRJEB59056, 142

202   contigs, 223 Mbp, N50 3.6 Mbp, BUSCO: CS:97.7%, CD:0.5%, F:0.3%, M:1.4%,

203   (n:5286, lepidoptera_odb10). The last common ancestor of these species was

204   approximately 23 million years ago (Chazot et al. 2019). We next sought to identify

205   3'UTR regions that were expressed in the relevant tissue and developmental stage of

206   our miRNA data.

207        Obtaining accurate 3'UTR annotations is challenging for several reasons. First,

208   the 3'UTR per locus is highly variable, with > 65% of human and *Drosophila* loci

209   producing alternative polyadenelated mRNAs across tissues and development (Derti et

210   al. 2012; Ye et al. 2023; Sanfilippo et al. 2017). This gains relevance as the available

211   genomic annotation of our focal species did not use RNAseq data from diapause

212   relevant tissues for their annotation. Second, methods for predicting 3'UTR regions from

213   DNA alone, or even with RNAseq data, perform with high variability across species and

214   in general, poorly in non-model species (Ye et al. 2023; Bryce-Smith et al. 2023), and

215   though some have tried to directly address this (Huang & Teeling 2017), obtaining

216   meaningful UTR predictions is challenging in novel species. Thus, in order to efficiently

217   move beyond data and bioinformatic limitations, here we deployed a simplified

218   approach for exploring potential 3'UTR regions for our focal species.

219        We assessed the 3'UTR annotation for the *P. napi* genome and found that it had

220   overpredicted UTR regions (GCA_905231885.1; (Lohse, Hayward, et al. 2021), such

221   that UTR regions routinely overlapped with flanking genes. In addition, at the time of our

222   analyses, the annotation of GCA_905231885.1 available from the Darwin Tree of Life

223   Program relied on an early annotation pipeline that was not optimized for Lepidoptera.

224   Accordingly, we chose to rely upon a *de novo* genome annotation we previously

225   generated (Steward et al., in review). This *de novo* annotation was produced using the

226   BRAKER2 pipeline (v.2.1.5, (Brůna et al. 2020; Hoff et al. 2016; Ter-Hovhannisyan et

227   al. 2008; Stanke et al. 2006, 2008; Lomsadze et al. 2005; Hoff et al. 2019), run in

228   protein mode using Arthropoda OrthoDB (v.10) reference proteins. This annotation

229   contained 123,638 exons, 16,449 genes and was found to contain 98.4% complete

230    BUSCOs for Lepidoptera_ODB10. Comparisons between this annotation and two

231    accessed from the Darwin Tree of Life revealed the BRAKER2 annotation to be the

232    most complete (i.e. fewest fragmented BUSCOs, a small proportion of single exon

233    genes, and more total estimated transcripts (see Supplementary methods; Table S1, S2
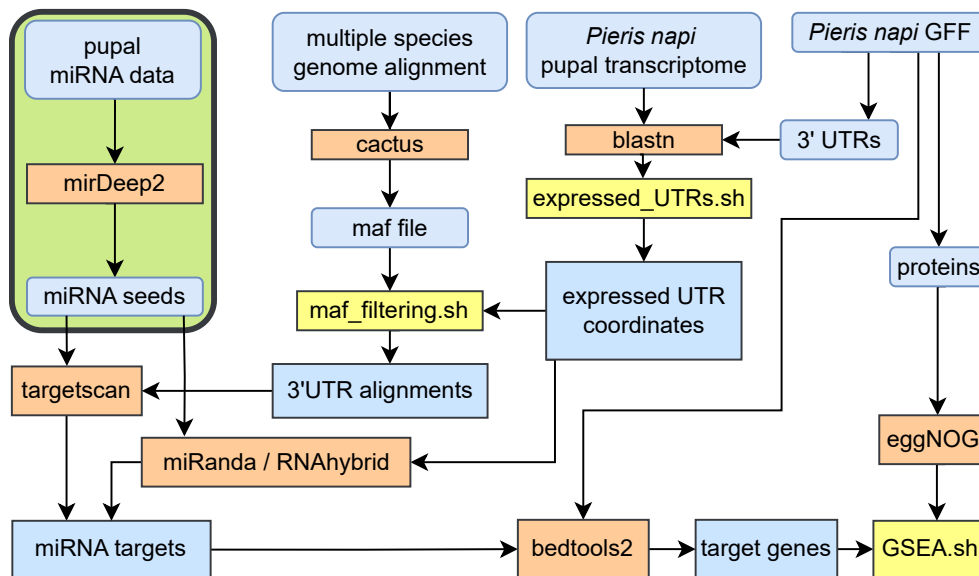
234    in Steward et al., in review).

235         Among moths and flies, the majority of 3'UTR regions are expected to be within 1

236    kb of the stop codon in the terminal coding exon, based upon detailed studies from

237    several *Drosophila* species (Sanfilippo et al. 2017; Wang et al. 2019) and 3'UTR lengths

238    for the an exemplar moth (*Bombyx mori* mean=923 bp, n=27,556) and butterfly

239    (*Heliconius melpomene* mean=600, n=11,770) downloaded from UTRdatabase

240    (Lo Giudice et al. 2023). While alternative UTRs in animals can involve spliced introns,

241    the frequency in 3'UTR regions are lower than 5'UTR, and usually < 10% (Mignone et

242    al. 2002). Based upon these expectations of 3'UTRs, we generated a bed file of likely

243    3'UTR regions, extending 1kb beyond every stop codon (and containing 9 codons (27

244    bp) prior to the terminal codon), of every protein isoform. We then assessed whether

245    any of these candidate 3'UTR regions had a significant match via blastn when searched

246    against the assembled transcriptome of an RNAseq dataset. The assembled

247    transcriptome was generated using Trinity (Haas et al. 2013), default parameters, with

248    RNAseq data comprising all of the same tissues and timepoints of our miRNA samples

249    (Pruisscher et al. 2021). Alignments were filtered to only include candidate 3'UTR

250    regions that had at least 70 bp of 3pUTR (filter settings: DNA identity > 90%, e-value <

251    0.000001, bitscore > 300, alignment length > 100 bp; NCBI BLAST v. 2.2.28+;

252    (Camacho et al. 2009). Coordinates for these post-filtered 3'UTR regions, which we

253    expect to be expressed 3'UTRs, were then used to identify these regions in the *P. napi*

254    genome, then whole genome alignment of all species, followed by the extraction of each

255    expressed 3'UTR region, which were then used as the input for conserved miRNA

256    target identification via targetscan_70.pl, part of TargetScan v.7 (Agarwal et al. 2018).

257    Manipulation of GFF files used bedtools2 (Quinlan & Hall 2010), which was also used to

258    assign nearest coding gene ID to each candidate 3'UTR region, while alignment filtering

259    used maffilter, with default settings unless indicated (remove_duplicates=yes,

260    reference=Pnapi, min_size=6), min_length=50, dist_max=1200; (Dutheil et al. 2014).

261   The other input file for targetscan_70.pl was the seed sequences for each of the

262   identified miRNA genes, predicted from mirDeep2 (Roberts et al., in review).

263        For each identified target region, the resulting output provides information on

264   species depth and seed size, which can be used to filter for differing degrees of

265   evolutionary conservation. Species depth indicates the number of species having the

266   identical target sequence in the alignment, ranging from all of the species down to only

267   2 species. Targets only found in 2 of the 6 species likely identify a region of lower

268   evolutionary constraint compared to targets identical across all species. Seed size of

269   the identified target can vary in size from an 8-mer down to a 6-mer, indicating the

270   length of base pairs of the identified target. Targets shorter in length are more likely to

271   occur by random chance compared to those of longer length. We use this information to

272   explore the quality of targets in later analyses.

273        Our second approach for miRNA target prediction used only two files as the input

274   for miRAnda (Enright et al. 2003) and RNAhybrid (Krüger & Rehmsmeier 2006). These

275   were the expressed 3'UTR coordinates for *P. napi* and seed sequences for *P. napi*, both

276   of which were described above. Both programs were run on default settings. Thresholds

277   for targets were set at e-value < 0.1 for miRanda, and p-value < 0.1 for RNAhybrid.

278



279

280   Fig. 1. Flowchart of miRNA target detection in *Pieris napi*, using two methods that lead

281   to gene set enrichment analyses (GSEA). Shown are the data files (blue), various

282     software programs (orange), and custom bioinformatic scripts (yellow) that were used.

283     Generation of miRNA data through to miRNA seed input file is from previously

284     published work (green enclosed portion of flow chart; Roberts et al., in review). Made

285     using diagrams.net.

286

287

288     ***Functional coherence via gene set enrichment analysis***

289     Target sets predicted per miRNA family were assessed for their functional coherence

290     via gene set enrichment analysis (GSEA) using the r package topGO v2.46 (Alexa &

291     Rahnenfuhrer 2023), with inputs of GO terms assigned to the coding regions of genes

292     having identified 3'UTR targets. For each GSEA of a miRNA target set, we took the -

293     $\log_{10}$ P-values of the top ten most significant categories, and quantified their distribution

294     as a function of the number of aligned species having identical seed sequences, and for

295     different seed pairing lengths, from 6mer to 8mer.

296

297     ***Comparative assessment of functional coherence***

298     In order to gain a robust assessment of miRNA functional coherence, with miRNA target

299     sets independent of our work and for model species having higher quality target

300     prediction, we repeated our analyses on the miRNA targets from 4 additional diverse

301     animals. Three datasets were downloaded from TargetScan databases (*Homo sapiens*:

302     TargetScanHuman release 8.0, Predicted_Targets_Info.default_predictions.txt

303     (McGeary et al. 2019); *Mus musculus*: TargetScanMouse release 8.0,

304     Predicted_Targets_Info.default_predictions.txt (McGeary et al. 2019); *Drosophila*

305     *melanogaster*: TargetScanFly release 7.2,

306     Predicted_Targets_Info.default_predictions.txt, (Agarwal et al. 2018)), while predicted

307     cichlid targets for *Oreochromis niloticus* (Mehta et al. 2022), were provided by Dr. T.

308     Mehta upon request. Note that for each TargetScan species dataset, in order to connect

309     miRNA ID to coding gene ID to GO terms of the latter, for the relevant genome

310     assembly, its GFF annotation was downloaded and protein sequences per ID extracted

311     using gffread from cufflinks-2.2.1 (Trapnell et al. 2010), for which GO annotations were

312     generated using functional annotation via orthology assignment, implemented in the

313   online server eggNOG using default settings (Huerta-Cepas et al. 2019), which was

314   then joined to the miRNA table downloaded from the relevant TargetScan database. An

315   estimate of the evolutionary depth over which 3'UTR alignments were made in order to

316   assess evolutionary constrain was estimated from. Ages for each clades of data upon

317   which miRNA targets were based, i.e. the age of the relevant crown groups (the

318   paraphyletic *Drosophila* genus at 53 MYA (Suvorov et al. 2022); the dataset for *H.*

319   *sapiens* involved using 84 of 100 species of the UCSC multiz alignment (Agarwal et al.

320   2015), including all species sister to, *Latimeria chalumnae*, as well as this coelacanth,

321   with their crown age estimated at roughly 400 MYA (Amemiya et al. 2013); the dataset

322   for *M. musculus* only included 52 species of the 60-way multiz alignment of UCSC, and

323   has a similar crown age as *H. sapiens*; the dataset for target *O. niloticus* has a crown

324   age estimated at 10 MYA (Mehta et al. 2022).

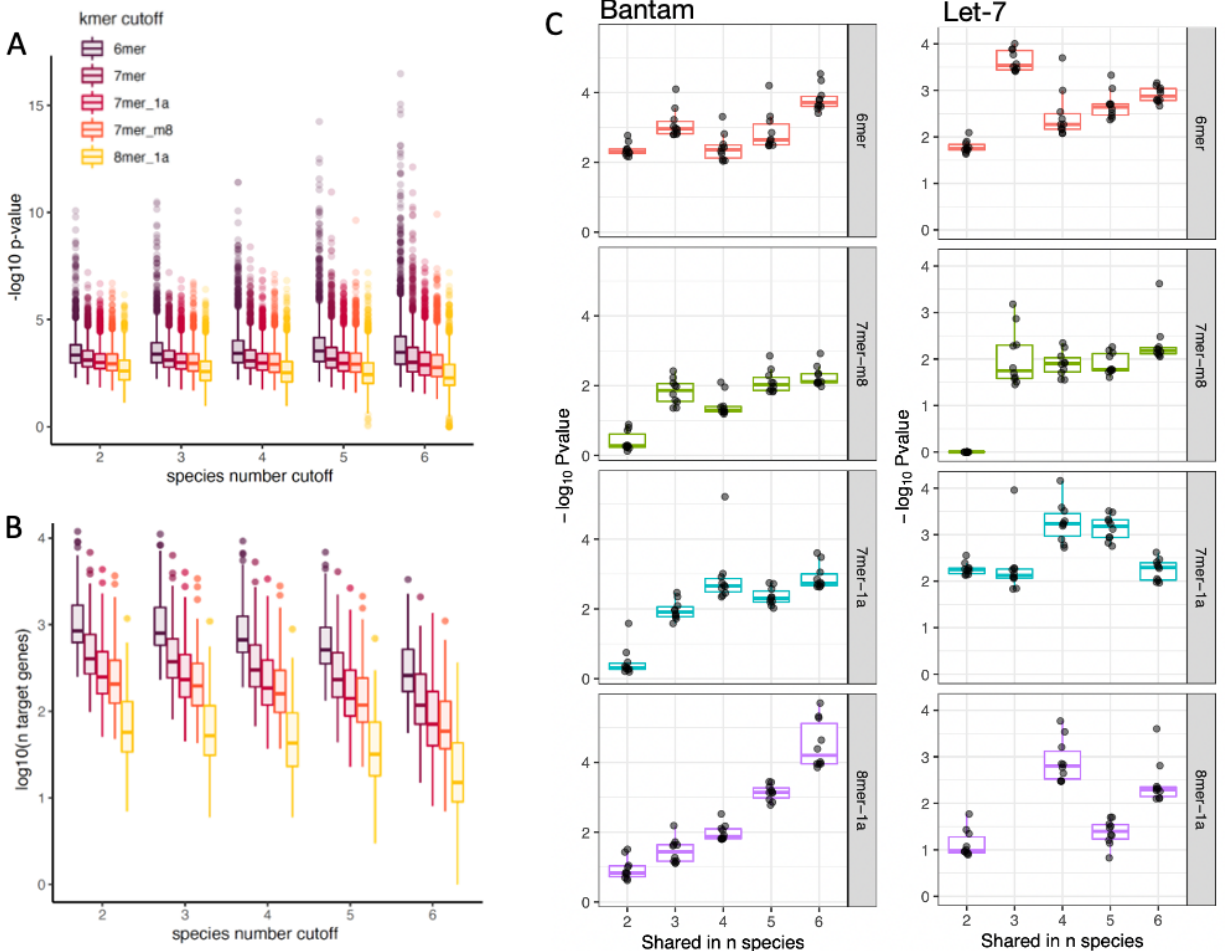325

326   **Results and Discussion**

327   An extensive miRNA sequencing effort has recently identified 257 miRNAs expressed

328   during pupal development of *P. napi* (236 expressed in head tissue, 207 in the

329   abdomen; Roberts et al., in review). Here we use this data to predict mRNA targets of

330   these miRNAs in *P. napi*. We began by identifying which mRNAs, among all candidate

331   3'UTR regions in the genome of *P. napi*, were expressed in a tissue matched RNAseq

332   transcriptome assembly. We then identified these 3'UTR regions if mRNA in a

333   multispecies, whole-genome alignment (n=6 species of Pieridae, Lepidoptera) that span

334   nearly 23 million years of divergence (Chazot et al. 2019). The resulting 3'UTR

335   alignment, together with the seed sequences from the identified miRNA genes of *P.*

336   *napi*, were then used as input for TargetScan v.7, which uses evolutionary conservation

337   in 3'UTRs to predict miRNA targets (Agarwal et al. 2018).

338       Next, we sought an independent means of quantifying whether these predicted

339   target sets per miRNA gene had more biological meaning than random sets, as

340   critiques of target prediction methods suggest that target sets generate from tools such

341   as miRAanda and RNAhybrid may be dominated by false positives (Fridrich et al. 2019;

342   Pinzón et al. 2017; Krüger & Rehmsmeier 2006). We reasoned that since a general

343   feature of gene regulatory networks (GRN) is their extensive functional coherence of

11

344   regulated genes, as most transcription factors regulate related modules of genes (Stuart

345   et al. 2003; Wolfe et al. 2005), the same is likely true for the targets of miRNA (see

346   Methods for additional discussion). Functional coherence was quantified using gene set

347   enrichment analysis (GSEA) upon the predicted set of gene targets for each miRNA,

348   using the average significance of the top ten most enriched GO categories as the

349   representative metric.

350        In order to assess whether there was any functional coherence in our predicted

351   targets, we quantified GSEA of the miRNA target sets using variable levels of

352   evolutionary constraint. TargetScan output provides two axes upon which to vary

353   evolutionary constraint in miRNA target prediction. First, we used differing thresholds of

354   constraint upon the species alignment of the 3'UTR, by varying the number of species

355   for which the seed site was required to be identical. Our lowest evolutionary constraint

356   level required only 2 species to have identical sequences in the alignment for the

357   miRNA seed site (the lowest threshold we could set), while our most stringent required

358   all 6 species to have the same identical sequence for the seed site. Second, there are 5

359   different sizes of target sites for the seed match region of the 3'UTR, ranging from 6 bp

360   (6mer) to 8 bp (8mer) in length. Requiring target sites to be longer in length is a more

361   stringent requirement. In combination, our most relaxed setting was 6mer for only 2

362   species in the alignment, while our most constrained was 8mer for all species. In order

363   to assess the relative tradeoff across these axes of constraint in the prediction of

364   miRNA targets, we explored our results extensively (fig. 2 *A,B*). As the stringency

365   increases, via increasing the number of species having target seed or increasing the

366   size of the seed match category, the predicted number of targets per miRNA gene

367   decreases, suggesting there is a biological signal in our target prediction method. While

368   these results are highly variable across miRNA genes (fig. 2*C*), we concluded that a

369   good balance between over-prediction and power was using a 7mer seed match size

370   and higher (termed 7mer-inclusive, which includes all targets from 7mer variants and

371   8mer) that is present and conserved across all of the aligned species.

372

373

Fig. 2. Assessment of GSEA results across predicted targets per miRNA gene. (*A*)
Significance of the top 10 GO terms per target set per miRNA gene (each dot is one
term) shown as a boxplot of all results, as a function of the number of species for which
seed was identical, for each of 5 different sizes of site type of the seed match (color
scale purple to yellow). As the stringency of predicted targets increases from being
found only in 2 species to all 6 species, the significance values increase for the smaller
seed match sizes (e.g. 6mers increase while 8mers do not). (*B*) Number of targets per
miRNA gene (each dot is count for a miRNA gene), across different prediction
thresholds of species number and miRNA seed match size (as in A). As the stringency
increases, via increasing the number of species having target seed or increasing the
size of the seed match category (color scale purple to yellow), the predicted number of
targets per miRNA gene decreases (6mer in 2 species is largest set, 8mer in 6 species
is the smallest). (*C*) Shown are GSEA results for two miRNA genes (left is Bantam, right

387 is Let-7), displaying effects of stringency increase on significance of the top 10 GO

388 terms per target set per miRNA gene. These exemplify the range of variation between

389 miRNA genes in their GSEA results, with Bantam exhibiting a strong increase in GSEA

390 P-value as evolutionary constraint is maximized (8mer-1a panel) and Let-7 lacking this
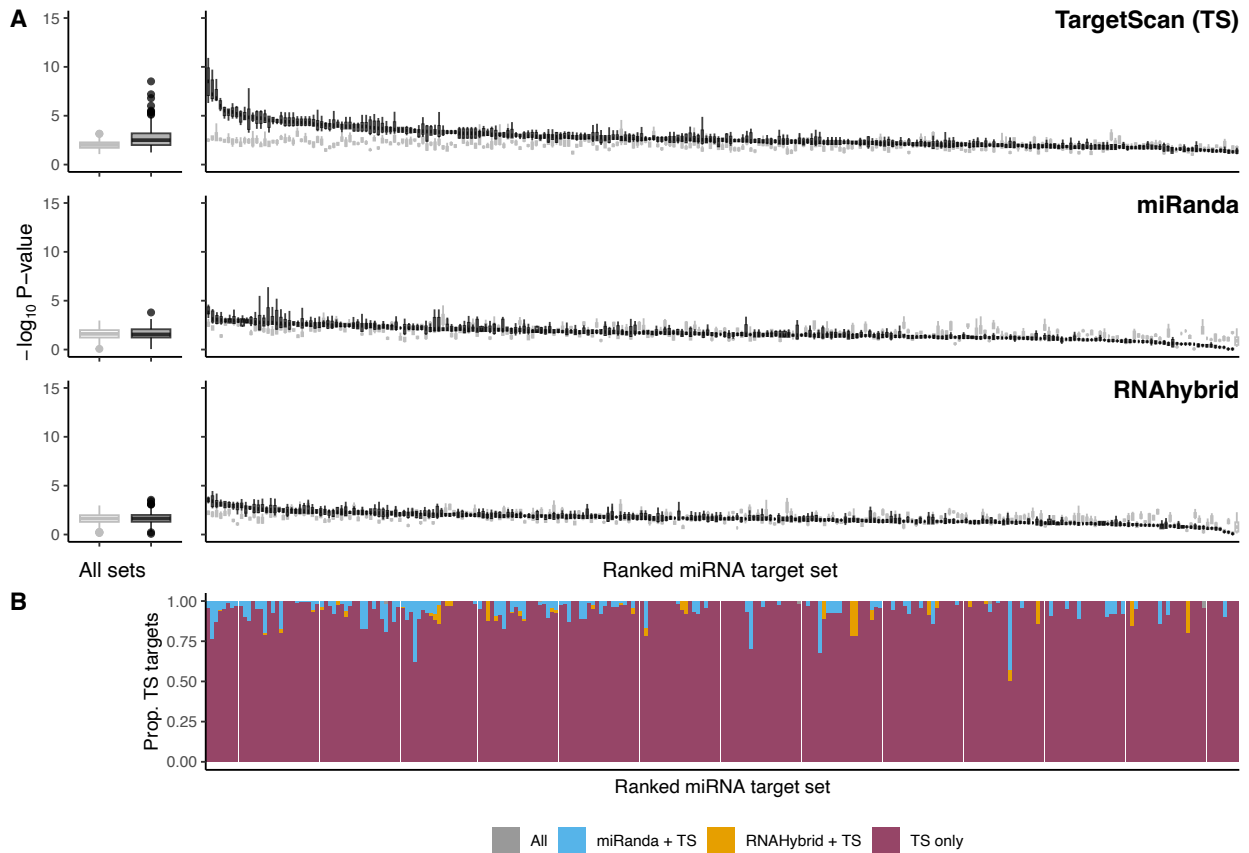
391 trend.

392

393

394      For comparison, we also used single species target prediction methods. Using

395 the 3'UTR regions of *P. napi* and seed sequences of miRNA genes as input, we used

396 the most commonly employed target prediction tool by the ecological and evolutionary

397 genomics community, miRanda (Enright et al. 2003). We additionally employed a

398 second single species tool with the same input data, RNAhybrid (Krüger & Rehmsmeier

399 2006). In order to compare the predicted targets across these tools, we quantified their

400 relative functional coherence via GSEA using the 7mer-inclusive conservation threshold

401 (described above). As a control, a GSEA was conducted on random sets of gene

402 targets conditional on the set size of the observed miRNA targets, which we used as our

403 background expectation of significance given concerns about GSEA significance

404 thresholds when working with miRNA targets (Bleazard et al. 2015).

405      The predicted targets of each miRNA from both methods exhibited significant

406 GSEA results, with average P-values for miRanda of 0.0185 and 0.0420 for RNAhybrid

407 (fig. 3*a*). However, GSEA results on sets of randomly drawn genes had P-value

408 distributions that entirely overlapped with the gene set targets predicted by these

409 methods (fig. 3*a*). Thus, GSEA P-value for targets from miRAnda, RNAhybrid, and

410 random draws were lower than nominal P-value significance thresholds (i.e., alpha =

411 0.05), highlighting two issues. First, these results exemplify previously noted challenges

412 of GSEA when investigating miRNA targets (Bleazard et al. 2015), in that resulting P-

413 values are poorly controlling for diverse many to many relationships, as GSEA were not

414 designed for such relationships. Second, neither miRAnda nor RNAhybrid predicted

415 targets that performed better than random.

416      In stark contrast to the previous results, miRNA targets predicted using

417 evolutionary conservation via TargetScan exhibited extensive functional coherence (fig.

418 3*a*), with GSEA P-values much higher than random draws. This result suggests two

419 mutually exclusive explanations. Either *P. napi* has miRNA targets that lack functional

420 coherence, which could explain the miRanda and RNAhybrid results and therefore

421 justify continued use of such tools by the non-model species community, *or* the miRNAs

422 of this butterfly exhibit functional coherence and only biologically meaningful target sets

423 can reveal this pattern. When facing variable results among target prediction methods,

424 studies in the non-model species community commonly intersect results from various

425 target prediction methods, despite this being explicitly discouraged by experts in the

426 miRNA field (Fridrich et al. 2019; Ritchie et al. 2009). To quantify the performance of

427 such an intersection approach, here we assess the overlap of targets from miRAnda

428 and RNAhybrid with respect to target predictions from TargetScan. We find no

429 substantial overlap across these three methods. Further, the level of overlap among

430 methods does not covary with the degree of functional coherence observed in our
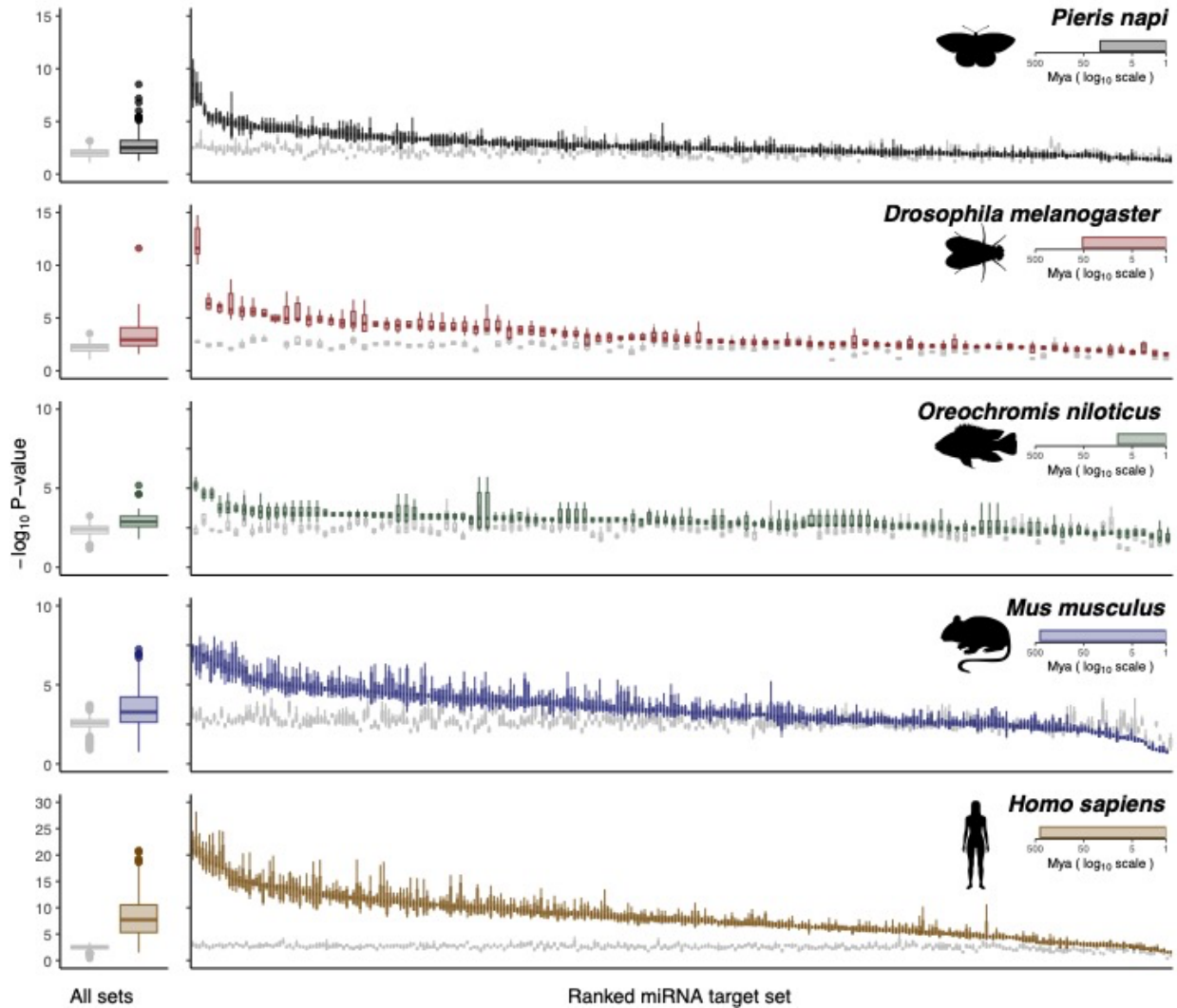
431 TargetScan results (fig. 3*b*).

432

433



Fig. 3. The functional coherence of miRNA targets across animals measured using gene set enrichment analysis (GSEA). (*A*) Comparison of the functional coherence of miRNA target predictions and their relationships, predicted in the butterfly *Pieris napi*. Gene set enrichment analysis P-values for top 10 GO terms for each miRNA (Y-axis) for targets predicted using Targetscan (top panel), miRanda (middle panel), RNAhybird (lower panel). Left-hand panels summarize median P-values for random (light grey) and predicted (black) miRNA target sets, while right-hand panels show results each miRNA target set. (*B*) Intersection of predicted targets from all three methods in relation to TargetScan results, shown as a proportion. Order of miRNAs along X axis are by mean P-value based upon Targetscan GSEA results.

    In order to discriminate between the two aforementioned explanations, we next quantified functional coherence using four published miRNA target sets. Across diverse

449    metazoans, from arthropods to vertebrates, we found extensive functional coherence

450    across many miRNAs (fig. 4). In each species, a large faction of predicted miRNAs

451    exhibited a significantly greater functional coherence than background. Importantly, all

452    of these previously published target sets were generated using the TargetScan

453    framework, using phylogenetic conservatism of miRNA binding sites as a core

454    identification criteria (Friedman et al. 2009; Agarwal et al. 2015). Common to all species

455    is a substantial variation among miRNA gene sets in their functional coherence (the left

456    vs right side of the P-value ranked distribution of miRNA genes). Whether this variation

457    arises due to unequal coherence across miRNAs, variation in the functional annotation

458    of relevant targets, poorly annotated 3'UTRs, or other factors warrants attention.

459    However, the extensive functional coherence seen across nearly all miRNA genes in *H.*

460    *sapiens* suggests such variation likely arises due to factors other than unequal

461    coherence among the target sets of miRNA genes. Among these diverse metazoans,

462    the lower functional coherence observed in these cichlids likely arises due to the young

463    age of the clade analyzed (~ 10 million years), as this necessarily results in a lower

464    power via phylogenetic conservatism. Highlighting the need and challenges of

465    bioinformatic target assessment in young clades, this clade of cichlids is an exemplar of

466    adaptive radiations, having generated > 2000 species in the 10 million years, making

467    observations of their massive reorganization of the miRNA GRN incredibly intriguing for

468    evolutionary study (Mehta et al. 2022).

469

470

471

472

473 Fig. 4. Functional coherence of miRNA targets across animals measured using gene set

474 enrichment analysis (GSEA). Left-hand boxplots summarize median P-values for

475 random and predicted miRNA target sets, while right-hand boxplots show P-values for

476 the top 10 enriched GO terms per per miRNA gene, ordered by median GSEA P-value

477 within each species. Results from predicted targets are colored while results from

478 randomly selected genes are shown in gray. Inset horizontal bars indicate crown age

479 (million years) of the species used to generate miRNA target predictions. Results from

480 *P. napi* (fig. 3a) are presented here, allowing for direct comparison with four divergent

481 taxa whose published datasets were generated using TargetScan.

482

483

**Conclusions**

484        

485        Functional coherence in the targets of miRNA genes appears to be common in

486        the tree of life. Using this observation, together with an in-depth study of miRNA targets

487        in a non-model species, our finding of no biological signal among the miRNA targets

488        produced by miRanda and RNAhybrid predictions is consistent with previous findings

489        and warnings of their low precision (Fridrich et al. 2019; Agarwal et al. 2015, 2018;

490        Pinzón et al. 2017; Ritchie et al. 2009). We conclude that a substantial body of research

491        may benefit from revising hypotheses based upon miRNA expression patterns, when

492        those hypotheses relied upon miRNA target prediction lacking measures of evolutionary

493        conservation.

494        Much remains to be discovered about the role the miRNAs play in adaptive

495        evolution and there has never better time for investigating the role of miRNA

496        posttranslational repression in novel species. An ever-increasing diversity of high-quality

497        genomes provides an unprecedented opportunity for exploiting evolutionary

498        conservation via recent advances in miRNA target prediction (Agarwal et al. 2018). We

499        note however that target predictions are merely another set of hypotheses. Since most

500        miRNAseq studies are also coupled with RNAseq, we further note that correlations

501        between increased miRNA expression and decreases in putative mRNA target

502        expression are also hypotheses fraught with a potential for high false-positives, given

503        the diverse patterns of expression in such datasets coupled with generally few sets of

504        diverse sampling points. Finally, while identified miRNA function in model species can

505        certainly aid hypothesis formulation of miRNA impacts, such relies upon increasingly

506        tenuous assumptions of evolutionarily conserved function (Rusin 2023).

507        Perhaps the most important way forward for the non-model species community

508        seeking to connect miRNA expression changes with adaptive phenotypes will be via

509        harnessing of emerging gene manipulation technologies in the testing of functional

510        hypotheses (Gudmunds et al. 2022). While the diverse many-to-many relationships

511        inherent in the miRNA GRN necessitate careful design and interpretation of such

512        experiments (Bartel 2018), these also offer unique opportunities. For example, consider

513        a scenario where many independent miRNA genes target the same seed sequence

514        within mRNA. While KO of all such miRNA genes could be lethal, knock out of one,

515 several, or many genes within such a gene family could effectively titrate phenotypic

516 effects. Additionally, advances in single cell sequencing of RNA could greatly advance

517 insights (Sekar et al. 2023), especially in the assessment of miRNA interactions with

518 mRNA GRNs across diverse tissues and developmental courses.

519       In conclusion, numerous studies across diverse taxa document differential

520 expression of miRNAs suggestive of a potentially important role in adaptive evolutionary

521 phenotypes. However, much work remains to be conducted in order to establish such

522 genotype to phenotype connections. Here, by drawing attention to the challenges of *de*

523 *novo* miRNA target prediction, we hope that more biologically meaningful hypotheses

524 will emerge that can then be tested by modification of miRNA genes or their target sites,

525 much as mRNA based hypotheses are now routinely explored via CRE and coding

526 region manipulations (Gudmunds et al. 2022).

527

536

537 **Author contributions**

538 C.W.W. performed all the bioinformatic analyses involved in the generation of miRNA

539 targets using TargetScan. R.S. provided R code for generating systematic GSEA for all

540 miRNA gene families and plotting the results. P.E. and K.R. ran the miRanda and

541 RNAhybrid analyses. C.W.W. and K.R. conceived of the study, with input from R.S.

542 C.W.W. wrote the manuscript with feedback from K.R. and the other coauthors. Y.O.

543 and H.V. provided two genomes for analyses. All authors approve of the manuscript.

544

545 **Data Availability Statement**

546    Scripts will be made available upon submission for review and publication.

547

548    **References**

549    Agarwal V, Bell GW, Nam J-W, Bartel DP. 2015. Predicting effective microRNA target sites in
550         mammalian mRNAs Izaurralde, E, editor. eLife. 4:e05005. doi: 10.7554/eLife.05005.
551    Agarwal V, Subtelny AO, Thiru P, Ulitsky I, Bartel DP. 2018. Predicting microRNA targeting
552         efficacy in Drosophila. Genome Biol. 19:152. doi: 10.1186/s13059-018-1504-3.
553    Alexa A, Rahnenfuhrer J. 2023. topGO: Enrichment Analysis for Gene Ontology. doi:
554         10.18129/B9.bioc.topGO.
555    Amemiya CT et al. 2013. The African coelacanth genome provides insights into tetrapod
556         evolution. Nature. 496:311–316. doi: 10.1038/nature12027.
557    Armstrong J et al. 2020. Progressive Cactus is a multiple-genome aligner for the thousand-
558         genome era. Nature. 587:246–251. doi: 10.1038/s41586-020-2871-y.
559    Bartel DP. 2018. Metazoan MicroRNAs. Cell. 173:20–51. doi: 10.1016/j.cell.2018.03.006.
560    Bleazard T, Lamb JA, Griffiths-Jones S. 2015. Bias in microRNA functional enrichment analysis.
561         Bioinformatics. 31:1592–1598. doi: 10.1093/bioinformatics/btv023.
562    Bracken CP, Scott HS, Goodall GJ. 2016. A network-biology perspective of microRNA function
563         and dysfunction in cancer. Nat. Rev. Genet. 17:719–732. doi: 10.1038/nrg.2016.134.
564    Bruce HS, Patel NH. 2020. Knockout of crustacean leg patterning genes suggests that insect
565         wings and body walls evolved from ancient leg segments. Nat. Ecol. Evol. 4:1703–1712.
566         doi: 10.1038/s41559-020-01349-0.
567    Brůna T, Hoff KJ, Lomsadze A, Stanke M, Borodovsky M. 2020. BRAKER2: Automatic
568         Eukaryotic Genome Annotation with GeneMark-EP+ and AUGUSTUS Supported by a
569         Protein Database. bioRxiv. 2020.08.10.245134. doi: 10.1101/2020.08.10.245134.
570    Bryce-Smith S et al. 2023. *Extensible benchmarking of methods that identify and quantify*
571         *polyadenylation sites from RNA-seq data*. Bioinformatics doi:
572         10.1101/2023.06.23.546284.
573    Camacho C et al. 2009. BLAST+: architecture and applications. BMC Bioinformatics. 10:421.
574         doi: 10.1186/1471-2105-10-421.
575    Chazot N et al. 2019. Priors and Posteriors in Bayesian Timing of Divergence Analyses: The
576         Age of Butterflies Revisited. Syst. Biol. 68:797–813. doi: 10.1093/sysbio/syz002.
577    Cui Q, Yu Z, Purisima EO, Wang E. 2006. Principles of microRNA regulation of a human cellular
578         signaling network. Mol. Syst. Biol. 2:46. doi: 10.1038/msb4100089.
579    Derti A et al. 2012. A quantitative atlas of polyadenylation in five mammals. Genome Res.
580         22:1173–1183. doi: 10.1101/gr.132563.111.
581    Dutheil JY, Gaillard S, Stukenbrock EH. 2014. MafFilter: a highly flexible and extensible multiple
582         genome alignment files processor. BMC Genomics. 15:np. doi: 10.1186/1471-2164-15-
583         53.
584    Enright AJ et al. 2003. MicroRNA targets in Drosophila. Genome Biol. 5:R1. doi: 10.1186/gb-
585         2003-5-1-r1.
586    Erwin DH. 2021. A conceptual framework of evolutionary novelty and innovation. Biol. Rev.
587         96:1–15. doi: https://doi.org/10.1111/brv.12643.

588  Fridrich A, Hazan Y, Moran Y. 2019. Too Many False Targets for MicroRNAs: Challenges and
589      Pitfalls in Prediction of miRNA Targets and Their Gene Ontology in Model and Non-
590      model Organisms. BioEssays. 41:1800169. doi: 10.1002/bies.201800169.
591  Friedlander MR, Mackowiak SD, Li N, Chen W, Rajewsky N. 2011. miRDeep2 accurately
592      identifies known and hundreds of novel microRNA genes in seven animal clades.
593      Nucleic Acids Res. 40:37–52. doi: 10.1093/nar/gkr688.
594  Friedman RC, Farh KK-H, Burge CB, Bartel DP. 2009. Most mammalian mRNAs are conserved
595      targets of microRNAs. Genome Res. 19:92–105. doi: 10.1101/gr.082701.108.
596  Fruciano C, Franchini P, Jones JC. 2021. Capturing the rapidly evolving study of adaptation. J.
597      Evol. Biol. 34:856–865. doi: 10.1111/jeb.13871.
598  Gudmunds E, Wheat CW, Khila A, Husby A. 2022. Functional genomic tools for emerging
599      model species. Trends Ecol. Evol. 37:1104–1115. doi: 10.1016/j.tree.2022.07.004.
600  Gusev Y. 2008. Computational methods for analysis of cellular functions and pathways
601      collectively targeted by differentially expressed microRNA. Methods. 44:61–72. doi:
602      10.1016/j.ymeth.2007.10.005.
603  Haas BJ et al. 2013. De novo transcript sequence reconstruction from RNA-seq using the Trinity
604      platform for reference generation and analysis. Nat. Protoc. 8:1494–1512. doi:
605      10.1038/nprot.2013.084.
606  Hoff KJ, Lange S, Lomsadze A, Borodovsky M, Stanke M. 2016. BRAKER1: Unsupervised
607      RNA-Seq-Based Genome Annotation with GeneMark-ET and AUGUSTUS: Table 1.
608      Bioinformatics. 32:767–769. doi: 10.1093/bioinformatics/btv661.
609  Hoff KJ, Lomsadze A, Borodovsky M, Stanke M. 2019. Whole-Genome Annotation with
610      BRAKER. In: Gene Prediction: Methods and Protocols. Kollmar, M, editor. Methods in
611      Molecular Biology Springer: New York, NY pp. 65–95. doi: 10.1007/978-1-4939-9173-
612      0_5.
613  Huang Z, Teeling EC. 2017. ExUTR: a novel pipeline for large-scale prediction of 3′-UTR
614      sequences from NGS data. BMC Genomics. 18:847. doi: 10.1186/s12864-017-4241-1.
615  Huerta-Cepas J et al. 2019. eggNOG 5.0: a hierarchical, functionally and phylogenetically
616      annotated orthology resource based on 5090 organisms and 2502 viruses. Nucleic Acids
617      Res. 47:D309–D314. doi: 10.1093/nar/gky1085.
618  Kang W et al. 2018. miRTrace reveals the organismal origins of microRNA sequencing data.
619      Genome Biol. 19:213. doi: 10.1186/s13059-018-1588-9.
620  Kern F et al. 2020. What's the target: understanding two decades of in silico microRNA-target
621      prediction. Brief. Bioinform. 21:1999–2010. doi: 10.1093/bib/bbz111.
622  Kolmogorov M, Yuan J, Lin Y, Pevzner PA. 2019. Assembly of long, error-prone reads using
623      repeat graphs. Nat. Biotechnol. 37:540–546. doi: 10.1038/s41587-019-0072-8.
624  Krüger J, Rehmsmeier M. 2006. RNAhybrid: microRNA target prediction easy, fast and flexible.
625      Nucleic Acids Res. 34:W451–W454. doi: 10.1093/nar/gkl243.
626  Lee SY, Sohn K-A, Kim JH. 2012. MicroRNA-centric measurement improves functional
627      enrichment analysis of co-expressed and differentially expressed microRNA clusters.
628      BMC Genomics. 13:S17. doi: 10.1186/1471-2164-13-S7-S17.
629  Leung AKL, Sharp PA. 2010. MicroRNA Functions in Stress Responses. Mol. Cell. 40:205–215.
630      doi: 10.1016/j.molcel.2010.09.027.

631 Liu Y, Beyer A, Aebersold R. 2016. On the Dependency of Cellular Protein Levels on mRNA
632     Abundance. Cell. 165:535–550. doi: 10.1016/j.cell.2016.03.014.
633 Lo Giudice C et al. 2023. UTRdb 2.0: a comprehensive, expert curated catalog of eukaryotic
634     mRNAs untranslated regions. Nucleic Acids Res. 51:D337–D344. doi:
635     10.1093/nar/gkac1016.
636 Lohse K, Mackintosh A, et al. 2021. The genome sequence of the large white, Pieris brassicae
637     (Linnaeus, 1758). Wellcome Open Res. 6:262. doi: 10.12688/wellcomeopenres.17274.1.
638 Lohse K, Ebdon S, et al. 2021. The genome sequence of the small white, Pieris rapae
639     (Linnaeus, 1758). Wellcome Open Res. 6:273. doi: 10.12688/wellcomeopenres.17288.1.
640 Lohse K, Hayward A, et al. 2021. The genome sequences of the male and female green-veined
641     white, Pieris napi (Linnaeus, 1758). Wellcome Open Res. 6:288. doi:
642     10.12688/wellcomeopenres.17277.1.
643 Lomsadze A, Ter-Hovhannisyan V, Chernoff YO, Borodovsky M. 2005. Gene identification in
644     novel eukaryotic genomes by self-training algorithm. Nucleic Acids Res. 33:6494–6506.
645     doi: 10.1093/nar/gki937.
646 Manni M, Berkeley MR, Seppey M, Simão FA, Zdobnov EM. 2021. BUSCO Update: Novel and
647     Streamlined Workflows along with Broader and Deeper Phylogenetic Coverage for
648     Scoring of Eukaryotic, Prokaryotic, and Viral Genomes. Mol. Biol. Evol. 38:4647–4654.
649     doi: 10.1093/molbev/msab199.
650 McGeary SE et al. 2019. The biochemical basis of microRNA targeting efficacy. Science.
651     366:eaav1741. doi: 10.1126/science.aav1741.
652 Mehta TK et al. 2022. Evolution of miRNA-Binding Sites and Regulatory Networks in Cichlids
653     Parsch, J, editor. Mol. Biol. Evol. 39:msac146. doi: 10.1093/molbev/msac146.
654 Mignone F, Gissi C, Liuni S, Pesole G. 2002. Untranslated regions of mRNAs. Genome Biol.
655     3:reviews0004.1. doi: 10.1186/gb-2002-3-3-reviews0004.
656 Pinzón N et al. 2017. microRNA target prediction programs predict many false positives.
657     Genome Res. 27:234–245. doi: 10.1101/gr.205146.116.
658 Pruisscher P, Lehmann P, Nylin S, Gotthard K, Wheat CW. 2021. Extensive transcriptomic
659     profiling of pupal diapause in a butterfly reveals a dynamic phenotype. Mol. Ecol. n/a.
660     doi: 10.1111/mec.16304.
661 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
662     features. Bioinformatics. 26:841–842. doi: 10.1093/bioinformatics/btq033.
663 Ritchie W, Flamant S, Rasko JEJ. 2009. Predicting microRNA targets and functions: traps for
664     the unwary. Nat. Methods. 6:397–398. doi: 10.1038/nmeth0609-397.
665 Rusin LY. 2023. Evolution of homology: From archetype towards a holistic concept of cell type.
666     J. Morphol. 284:e21569. doi: 10.1002/jmor.21569.
667 Sanfilippo P, Wen J, Lai EC. 2017. Landscape and evolution of tissue-specific alternative
668     polyadenylation across Drosophila species. Genome Biol. 18:229. doi: 10.1186/s13059-
669     017-1358-0.
670 Sekar V et al. 2023. Detection of transcriptome-wide microRNA–target interactions in single
671     cells with agoTRIBE. Nat. Biotechnol. 1–7. doi: 10.1038/s41587-023-01951-0.
672 Stanke M, Diekhans M, Baertsch R, Haussler D. 2008. Using native and syntenically mapped
673     cDNA alignments to improve de novo gene finding. Bioinformatics. 24:637–644. doi:
674     10.1093/bioinformatics/btn013.

675 Stanke M, Schöffmann O, Morgenstern B, Waack S. 2006. Gene prediction in eukaryotes with a
676           generalized hidden Markov model that uses hints from external sources. BMC
677           Bioinformatics. 7:62. doi: 10.1186/1471-2105-7-62.
678 Steward RA, Okamura Y, Boggs CL, Vogel H, Wheat CW. 2021. The Genome of the Margined
679           White Butterfly (Pieris macdunnoughi): Sex Chromosome Insights and the Power of
680           Polishing with PoolSeq Data Lavrov, D, editor. Genome Biol. Evol. 13:evab053. doi:
681           10.1093/gbe/evab053.
682 Stuart JM, Segal E, Koller D, Kim SK. 2003. A Gene-Coexpression Network for Global
683           Discovery of Conserved Genetic Modules. Science. 302:249–255. doi:
684           10.1126/science.1087447.
685 Suvorov A et al. 2022. Widespread introgression across a phylogeny of 155 Drosophila
686           genomes. Curr. Biol. 32:111-123.e5. doi: 10.1016/j.cub.2021.10.052.
687 Ter-Hovhannisyan V, Lomsadze A, Chernoff YO, Borodovsky M. 2008. Gene prediction in novel
688           fungal genomes using an ab initio algorithm with unsupervised training. Genome Res.
689           18:1979–1990. doi: 10.1101/gr.081612.108.
690 Trapnell C et al. 2010. Transcript assembly and quantification by RNA-Seq reveals unannotated
691           transcripts and isoform switching during cell differentiation. Nat. Biotechnol. 28:511–515.
692           doi: 10.1038/nbt.1621.
693 Tsang JS, Ebert MS, van Oudenaarden A. 2010. Genome-wide Dissection of MicroRNA
694           Functions and Cotargeting Networks Using Gene Set Signatures. Mol. Cell. 38:140–153.
695           doi: 10.1016/j.molcel.2010.03.007.
696 Wang W et al. 2019. Evolutionary and functional implications of 3′ untranslated region length of
697           mRNAs by comprehensive investigation among four taxonomically diverse metazoan
698           species. Genes Genomics. 41:747–755. doi: 10.1007/s13258-019-00808-8.
699 Wolfe CJ, Kohane IS, Butte AJ. 2005. Systematic survey reveals general applicability of 'guilt-
700           by-association' within gene coexpression networks. BMC Bioinformatics. 6:227. doi:
701           10.1186/1471-2105-6-227.
702 Xu J, Wong C. 2008. A computational screen for mouse signaling pathways targeted by
703           microRNA clusters. RNA. 14:1276–1283. doi: 10.1261/rna.997708.
704 Ye W, Lian Q, Ye C, Wu X. 2023. A Survey on Methods for Predicting Polyadenylation Sites
705           from DNA Sequences, Bulk RNA-seq, and Single-cell RNA-seq. Genomics Proteomics
706           Bioinformatics. 21:67–83. doi: 10.1016/j.gpb.2022.09.005.
707 Zimin AV et al. 2017. Hybrid assembly of the large and highly repetitive genome of *Aegilops*
708           *tauschii* , a progenitor of bread wheat, with the MaSuRCA mega-reads algorithm.
709           Genome Res. 27:787–792. doi: 10.1101/gr.213405.116.
710
711