1    Title

2    Repeated duplication of Argonaute2 is associated with strong selection and testis

3    specialization in *Drosophila*

4

5    Authors

6    Samuel H. Lewis*,†, Claire L. Webster*,‡, Heli Salmela§ & Darren J. Obbard*,**

7

8    Affiliations

9    *Institute of Evolutionary Biology, University of Edinburgh, Ashworth Laboratories,

10    EH9 3FL, United Kingdom

11    †Present Address: Department of Genetics, University of Cambridge, Downing

12    Street, Cambridge, CB2 3EH, United Kingdom

13    ‡Present Address: Life Sciences, University of Sussex, United Kingdom

14    §Department of Biosciences, Centre of Excellence in Biological Interactions,

15    University of Helsinki, Helsinki, Finland

16    **Centre for Immunity, Infection and Evolution, University of Edinburgh, Ashworth

17    Laboratories, EH9 3FL, United Kingdom

18

19    Supporting Data

20    All new sequences produced in this study have been submitted to Genbank as

21    KX016642-KX016771.

22

23

1    Running Title

2    Adaptive specialization of Drosophila Argonaute2 duplicates

3

4    Keywords

5    Argonaute, RNAi, *Drosophila*, duplication, testis

6

7    Corresponding Author

8    Name: Samuel H. Lewis

9    Mailing address: Department of Genetics, University of Cambridge, Downing Street,

10    Cambridge, CB2 3EH, United Kingdom

11    Telephone number: +441223 332584

12    Email address: sam.lewis@gen.cam.ac.uk

13

1  Abstract

2  Argonaute2 (Ago2) is a rapidly evolving nuclease in the *Drosophila melanogaster*

3  RNA interference (RNAi) pathway that targets viruses and transposable elements in

4  somatic tissues. Here we reconstruct the history of Ago2 duplications across the

5  *Drosophila obscura* group, and use patterns of gene expression to infer new

6  functional specialization. We show that some duplications are old, shared by the

7  entire species group, and that losses may be common, including previously

8  undetected losses in the lineage leading to *D. pseudoobscura*. We find that while the

9  original (syntenic) gene copy has generally retained the ancestral ubiquitous

10  expression pattern, most of the novel Ago2 paralogues have independently

11  specialized to testis-specific expression. Using population genetic analyses, we

12  show that most testis-specific paralogues have significantly lower genetic diversity

13  than the genome-wide average. This suggests recent positive selection in three

14  different species, and model-based analyses provide strong evidence of recent hard

15  selective sweeps in or near four of the six *D. pseudoobscura* Ago2 paralogues. We

16  speculate that the repeated evolution of testis-specificity in *obscura* group Ago2

17  genes, combined with their dynamic turnover and strong signatures of adaptive

18  evolution, may be associated with highly derived roles in the suppression of

19  transposable elements or meiotic drive. Our study highlights the lability of RNAi

20  pathways, even within well-studied groups such as *Drosophila*, and suggests that

21  strong selection may act quickly after duplication in RNAi pathways, potentially giving

22  rise to new and unknown RNAi functions in non-model species.

23

24

1    Introduction

2    Argonaute genes are found in almost all eukaryotes, where they play a key role in

3    antiviral immune defence, gene regulation and genome stability. They perform this

4    diverse range of functions through their role in RNA interference (RNAi)

5    mechanisms, an ancient system of nucleic acid manipulation in which small RNA

6    (sRNA) molecules guide Argonaute proteins to nucleic acid targets through base

7    complementarity (reviewed in Meister 2013). Gene duplication has occurred

8    throughout the evolution of the Argonaute gene family, with ancient duplication

9    events characteristic of some lineages – such as three duplications early in plant

10   evolution (Singh et al. 2015), and multiple expansions and losses throughout the

11   evolution of nematodes (reviewed in Buck and Blaxter 2013) and the Diptera (Lewis

12   et al. 2016). After duplication, Argonautes have often undergone functional

13   divergence, involving changes in expression patterns and altered sRNA binding

14   partners (Lu et al. 2011; Leebonoi et al. 2015; Miesen et al. 2015). Duplication early

15   in eukaryotic evolution produced two distinct Argonaute subfamilies, Ago and Piwi,

16   which have since been retained in the vast majority of Metazoa (Cerutti and Casas-

17   Mollano 2006). Members of the Ago subfamily are expressed in both somatic and

18   germline tissue, and variously bind sRNAs derived from host transcripts (miRNAs,

19   endo-siRNAs) or transposable elements (TE endo-siRNAs) and viruses (viRNAs). In

20   contrast, in most vertebrates and arthropods, the Piwi subfamily members are

21   expressed primarily in association with the germline (reviewed in Ross et al. 2014),

22   and bind sRNAs from TEs and host loci (piRNAs), suggesting that the Piwi subfamily

23   specialised to a germline-specific role on the lineages leading to vertebrates and

24   arthropods.

After the early divergence of the Ago and Piwi subfamilies, subsequent duplications gave rise to three Piwi subfamily members (Ago3, Aubergine (Aub) and Piwi) and two Ago subfamily members (Ago1 & Ago2) in *Drosophila melanogaster*. All three Piwi subfamily genes are associated with the germline and bind Piwi-interacting RNAs (piRNAs) derived from TEs and other repetitive genomic elements: Ago3 and Aub amplify the piRNA signal through the "Ping-Pong" cycle (reviewed in Luteijn and Ketting 2013), and Piwi suppresses transposition by directing heterochromatin formation (Sienski et al. 2012). These functional differences are associated with contrasting selective regimes, with Aub evolving under positive selection (Kolaczkowski et al. 2011) and more rapidly than Ago3 and Piwi (Obbard, Gordon, et al. 2009). In contrast, Ago1 binds microRNAs (miRNAs), and regulates gene expression by inhibiting translation and marking transcripts for degradation (reviewed in Eulalio et al. 2008). This function imposes strong selective constraint on Ago1, resulting in slow evolution and very few adaptive substitutions (Obbard et al. 2006; Obbard, Gordon, et al. 2009; Kolaczkowski et al. 2011). Finally, Ago2 binds small interfering RNAs (siRNAs) from viruses (viRNAs) and TEs (endo-siRNAs), and functions in gene regulation (Wen et al. 2015), dosage compensation (Menon and Meller 2012), and the ubiquitous suppression of viruses (Li et al. 2002; van Rij et al. 2006) and TEs (Chung et al. 2008; Czech et al. 2008). Ago2 also evolves under strong positive selection, with frequent selective sweeps (Obbard et al. 2006; Obbard, Gordon, et al. 2009; Obbard, Welch, et al. 2009; Kolaczkowski et al. 2011; Obbard et al. 2011), possibly driven by an arms race with virus-encoded suppressors of RNAi (VSRs) (Obbard et al. 2006; Marques and Carthew 2007; van Mierlo et al. 2014).

1  In contrast to *D. melanogaster*, from which most functional knowledge of Ago2 in

2  arthropods is derived, an expansion of Ago2 has been reported in *D. pseudoobscura*

3  (Hain et al. 2010), providing an opportunity to study how the RNAi pathway evolves

4  after duplication. Given the roles of *D. melanogaster* Ago2 in antiviral defence (Li et

5  al. 2002; van Rij et al. 2006), TE suppression (Chung et al. 2008; Czech et al. 2008),

6  dosage compensation (Menon and Meller 2012), and gene regulation (Wen et al.

7  2015), we hypothesized that these *D. pseudoobscura* Ago2 paralogues may have

8  diverged in function. To elucidate the evolution and function of Ago2 paralogues in

9  *D. pseudoobscura* and its relatives, we identified and dated Ago2 duplication events

10  across available *Drosophila* genomes and transcriptomes, tested for divergence in

11  expression patterns between the Ago2 paralogues in *D. subobscura*, *D. obscura* and

12  *D. pseudoobscura*, and quantified the evolutionary rate and positive selection acting

13  on each of these paralogues. We find that testis-specificity of Ago2 paralogues has

14  evolved repeatedly in the *obscura* group, and that the majority of paralogues show

15  evidence of recent positive selection.

16

17  Materials and Methods

18  Identification of Ago2 homologues in the Drosophilidae

19  We used tBLASTx to identify Ago2 homologues in transcriptomes and genomes of

20  39 species of the Drosophilidae, using previously-characterised Ago2 from the

21  closest possible relative to provide the query for each species. If blast returned

22  partial hits, we aligned all hits from the target species to all Argonautes from the

23  query species, and assigned hits to the appropriate Ago lineage based on a

24  neighbour-joining tree. For each query sequence, we then manually curated partial

1    blast hits into complete genes using Geneious v5.6.2 (http://www.geneious.com,

2    Kearse et al. 2012) (see Supplementary Materials for sequence accessions).

3    Additionally, we used degenerate PCR to identify Ago2 paralogues in *D. azteca* and

4    *D. affinis*, and paralogue-specific PCR with a touchdown amplification cycle to

5    validate the Ago2 paralogues identified in *D. subobscura*, *D. obscura* and *D.*

6    *pseudoobscura*. For each reaction, unincorporated primers were removed with

7    ExonucleaseI (New England Biolabs) and 5' phosphates were removed with

8    Antarctic Phosphatase (NEB). The PCR products were sequenced by Edinburgh

9    Genomics using BigDye V3 reagents on a capillary sequencer (Applied Biosystems),

10    and Sanger sequence reads were trimmed and assembled using Geneious v.5.6.2

11    (http://www.geneious.com, Kearse et al. 2012). We also used a combination of PCR

12    and blast searches to locate *D. pseudoobscura* Ago2a1 & Ago2a3, which lie on the

13    unplaced "Unknown_contig_265" in release 3.03 of the *D. pseudoobscura* genome

14    (all PCR primers are detailed in Table S4)*.*

15    <u>Phylogenetic analysis of drosophilid Ago2 paralogues</u>

16    To characterise the evolutionary relationships between Ago2 homologues in the

17    Drosophilidae, we aligned sequences using translational MAFFT (Katoh et al. 2002)

18    with default parameters. We noted that there is a high degree of codon usage bias

19    (CUB) in *D. pseudoobscura* Ago2e (effective number of codons (ENC)=34.24) and

20    *D. obscura* Ago2e (ENC=40.36), and a lesser degree in *D. subobscura* Ago2f

21    (ENC=45.63) and *D. obscura* Ago2f (ENC=48.39), and comparison with genome-

22    wide patterns of codon usage bias placed these genes in the lower half of the

23    distribution of ENC (Figure S5). To reduce the impact of CUB, which

24    disproportionately affects synonymous sites, we stripped all third position sites in this

25    analysis (Behura and Severson 2013). We then inferred a gene tree using the

1 Bayesian approach implemented in BEAST v1.8.1 (Drummond et al. 2012) under a

2 nucleotide model, assuming a GTR substitution model, variation between sites

3 modelled by a gamma distribution with four categories, and base frequencies

4 estimated from the data. We used the default priors for all parameters, except tree

5 shape (for which we specified a birth-death speciation model) and the date of the

6 *Drosophila-Sophophora* split. To estimate a timescale for the tree, we specified a

7 normal distribution for the date of this node using values based on mutation rate

8 estimates in Obbard et al. 2012, with a mean value of 32mya, standard deviation of

9 7mya, and lower and upper bounds of 15mya and 50mya respectively. We ran the

10 analysis for 50 million steps, recording samples from the posterior every 1,000 steps,

11 and inferred a maximum clade credibility tree with TreeAnnotator v1.8.1 (Drummond

12 et al. 2012). Note that precise date estimates are not a primary focus of this study,

13 but that other calibrations (Russo et al. 1995; Tamura 2004) would lead to more

14 ancient estimates of divergence, and thus stronger evidence for selective

15 maintenance.

16 Domain architecture and structural modelling of Ago2 paralogues in the *obscura*

17 group

18 To infer the location of each domain in each paralogue identified in *D. subobscura*,

19 *D. obscura* and *D. pseudoobscura*, we searched the Pfam database (Finn et al.

20 2009). To test for structural differences between the *D. pseudoobscura* paralogues,

21 we built structural models of each paralogue based on the published X-ray

22 crystallographic structure of human Ago2 (Schirle and Macrae 2012). We used the

23 MODELER software in the Discovery Studio 4.0 Modeling Environment (Accelrys

24 Software Inc., San Diego, 2013) to calculate ten models, selected the most

25 energetically favourable for each protein, and assessed model quality with the 3D-

1 profile option in the software. To assess variation in selective pressure across the

2 structure of each paralogue, we mapped variable residues onto each structure

3 (Figure S7) using PyMol v.1.7.4.1 (Schrödinger, LLC).

4 <u>Quantification of virus-induced expression of Ago2 paralogues</u>

5 We exposed 48-96hr post-eclosion virgin males and females of *D. melanogaster*, *D.*

6 *subobscura*, *D. obscura* and *D. pseudoobscura* to *Drosophila C virus* (DCV), by

7 puncturing the thorax with a pin contaminated with DCV at a dose of approximately

8 $4 \times 10^7$ TCID$^{50}$ per ml. Infection with DCV using this method has previously been

9 shown to lead to a rapid and ultimately fatal increase in DCV titre in *D. melanogaster*

10 and *obscura* group species (Longdon et al. 2015). All flies were incubated at 18C

11 under a 12L:12D light cycle, with *D. melanogaster* on Lewis medium and *D.*

12 *subobscura*, *D. obscura* and *D. pseudoobscura* on banana medium. We sampled 4-7

13 individuals per species at 0, 8, 16, 24, 48 and 72 hours post infection. At each time-

14 point we extracted RNA using TRIzol reagent (Ambion) and a chloroform/isopropanol

15 extraction, treated twice with TURBO DNase (Ambion), and reverse-transcribed

16 using M-MLV reverse transcriptase (Promega) primed with random hexamers. We

17 then quantified the expression of Ago2 paralogues in these samples by qPCR, using

18 Fast Sybr Green (Applied Biosystems) and custom-designed paralogue-specific

19 qPCR primer pairs (see Table S6 for primer sequences). Due to their high level of

20 sequence similarity (99.9% identity), no primer pair could distinguish between *D.*

21 *pseudoobscura* Ago2a1 and Ago2a3, so combined expression of these two genes is

22 presented as "Ago2a". All qPCR reactions for each sample were run in duplicate,

23 and scaled to the internal reference gene Ribosomal Protein L32 (RpL32). To

24 capture the widest possible biological variation, the three biological replicates for

1   each species each used a different wild-type genetic background (see Table S3 for

2   backgrounds used).

3   <u>Quantification of Ago2 paralogue expression in different tissues and life stages</u>

4   For *D. subobscura*, *D. obscura* and *D. pseudoobscura*, we extracted RNA from the

5   head, testis/ovaries and carcass of 48-96hr post-eclosion virgin adults, with males

6   and females extracted separately. Each sample consisted of 8-15 individuals in *D.*

7   *subobscura*, 10 individuals in *D. obscura* and 15 individuals in *D. pseudoobscura*.

8   We then used qPCR to quantify the expression of each Ago2 paralogue in each

9   tissue, with two technical replicates per sample (reagents, primers and cycling

10  conditions as above). We carried out five replicates per species, each using a

11  different wild-type background (see Table S3 for details of backgrounds used). To

12  provide an informal comparison with the expression pattern of Ago2 before

13  duplication (an "ancestral" expression pattern), we used the BPKM (bases per

14  kilobase of gene model per million mapped bases) values for Ago2 calculated from

15  RNA-seq data from the body (carcass and digestive system), head, ovary and testis

16  of 4 day old *D. melanogaster* adults by Brown et al. 2014, scaling each BPKM value

17  to the value for RpL32 in each tissue. Due to the design of that experiment, the body

18  data are derived from pooled samples of males and females (Brown et al. 2014).

19  To quantify expression of Ago2 paralogues in *D. pseudoobscura* embryos, we

20  collected eggs within 30 minutes of laying, and used qPCR to measure the

21  expression of each Ago2 paralogue (reagents and primers as above) in two separate

22  wild-type genetic backgrounds (MV8 and MV10). As above, we estimated an

23  ancestral expression pattern of Ago2 before duplication from the BPKM values for

24  Ago2 in 0-2hr old *D. melanogaster* embryos according to Brown et al. 2014, scaled

25  to the BPKM value for RpL32 in embryos. To determine any changes in the

10

1    expression of other *D. pseudoobscura* Argonautes (Ago1, Ago3, Aub & Piwi) that are

2    associated with Ago2 duplication, we measured their expression in adult tissues and

3    embryos as detailed above, and compared this with the expression of the

4    Argonautes in *D. melanogaster* as measured by Brown et al. 2014.

5    <u>Testing for evolutionary rate changes associated with tissue-specificity of Ago2</u>

6    We used codeml (PAML v4.4, Yang 1997) to fit variants of the M0 model (a single

7    dn/ds ratio, $\omega$) to the 65 drosophilid Ago2 homologues shown in Figure 1. All

8    analyses of sequence evolution excluded the highly-repetitive N-terminal glutamine-

9    rich repeat regions, as these regions are effectively unalignable, and are unlikely to

10   conform to simple models of sequence evolution (Palmer and Obbard 2016). In

11   contrast to the tree topology, which was based on $1^{st}$ and $2^{nd}$ positions only, the

12   alignment for the codeml analysis included all positions. To compare the evolutionary

13   rates of ubiquitously expressed and testis-specific Ago2 paralogues, we fitted a

14   model specifying one $\omega$ for the Ago2 paralogues that were shown to be testis-

15   specific by qPCR (7 homologues), and another $\omega$ for the rest of the tree (58

16   homologues). We also fitted two models to account for rate variation between the

17   *obscura* group Ago2 subclades. The first model specified a separate $\omega$ for the Ago2a

18   subclade (17 homologues), the Ago2e subclade (8 homologues), the Ago2f subclade

19   (5 homologues) and the rest of the tree (35 homologues). The second model

20   additionally incorporated an extra $\omega$ specified for the *D. pseudoobscura-D. persimilis*

21   Ago2a-Ago2b subclade (3 homologues, all of which are testis-specific, in contrast

22   with the rest of the *obscura* group Ago2a subclade). We used Akaike weights to

23   assess which model provided the best fit to the data, given the number of

24   parameters. As mentioned above, the high CUB seen in some Ago2 paralogues may

25   affect PAML analyses by decreasing synonymous site divergence (ds) in those

11

1    lineages, thereby inflating the dn/ds ratio ($\omega$). However, we find no link between

2    levels of codon usage bias and the value of $\omega$, suggesting that codon usage bias is

3    not impacting our PAML analyses.

4

5    <u>Sequencing of Ago2 paralogue haplotypes from *D. subobscura*, *D. obscura* and *D.*</u>

6    <u>*pseudoobscura*</u>

7    To obtain genotype data for the Ago2 paralogues in *D. subobscura*, *D. obscura* and

8    *D. pseudoobscura*, we sequenced the Ago2 paralogues from six males and six

9    females of each species, each from a different wild-collected line (detailed in Table

10    S3, sequence polymorphism data in Appendix S4). We extracted genomic DNA from

11    each individual using the DNeasy Blood and Tissue kit (Qiagen), and amplified and

12    Sanger sequenced each Ago2 paralogue from each individual (reagents and PCR

13    primers as above, sequencing primers detailed in Table S5). We trimmed and

14    assembled Sanger sequence reads using Geneious v.5.6.2

15    (http://www.geneious.com, Kearse et al. 2012), and identified polymorphic sites by

16    eye. After sequencing Ago2a (annotated as a single gene in the *D. pseudoobscura*

17    genome), we discovered two very recent Ago2a paralogues (which we denote

18    Ago2a1 & Ago2a3), which had been cross-amplified. For each *D. pseudoobscura*

19    individual we therefore re-sequenced Ago2a3 using one primer targeted to its

20    neighbouring locus GA22965, and used this sequence to resolve polymorphic sites

21    in the Ago2a1/Ago2a3 composite sequence, thereby gaining both genotypes for

22    each individual. For each Ago2 paralogue, we inferred haplotypes from these

23    sequence data using PHASE (Stephens et al. 2001), apart from the X-linked

24    paralogues (Ago2a1, Ago2a3 & Ago2d) in *D. pseudoobscura* males, for which phase

25    was obtained directly from the sequence data. The hemizygous haploid X-linked

1  sequenced were used in phase inference, and should substantially improve the

2  inferred phasing of female genotypes.

3  To quantify differences between paralogues in their population genetic

4  characteristics, we aligned haplotypes using translational MAFFT (Katoh et al. 2002),

5  and used DnaSP v.5.10.01 (Librado and Rozas 2009) to calculate the following

6  summary statistics for each Ago2 paralogue: $\pi$ (pairwise diversity, with Jukes-Cantor

7  correction as described in Lynch and Crease 1990) at nonsynonymous ($\pi_a$) and

8  synonymous ($\pi_s$) sites, Tajima's *D* (Tajima 1989) and ENC (Wright 1990). To

9  compare the ENC for each gene with the genome as a whole, we used codonW

10  v1.4.2 (Peden 1995) to calculate the ENC for the longest ORF from each gene or

11  transcript in the genomes or transcriptomes of *D. subobscura*, *D. obscura* and *D.*

12  *pseudoobscura* (ORF sets detailed below). In each species, we then compared the

13  ENC values of each Ago2 paralogue with this genome-wide ENC distribution.

14  <u>Testing for positive selection on Ago2 paralogues in the *obscura* group</u>

15  We used McDonald-Kreitman (MK) tests (McDonald and Kreitman 1991) to test for

16  positive selection on each Ago2 paralogue. For each paralogue, we chose an

17  outgroup with divergence at synonymous sites ($K_S$) in the range 0.1-0.2 where

18  possible. However, the prevalence of duplications and losses of Ago2 paralogues in

19  the *obscura* group meant that for some tests no suitably divergent extant outgroup

20  existed. In these cases, we reconstructed hypothetical ancestral sequences using

21  the M0 model provided by codeml from PAML (Yang 1997). To assess the effect of

22  these outgroup choices on our results we repeated each test with another outgroup,

23  and found no effect of outgroup choice on the significance of any tests, and only

24  marginal differences in estimates of $\alpha$ and $\omega_\alpha$ (results of tests using primary and

25  alternative outgroups are detailed in Table S1 & S2).

13

1. A complementary approach to identifying positive selection is to test for reduced

2. diversity at a locus compared with the genome as a whole. To compare the diversity

3. of each *D. pseudoobscura* Ago2 paralogue with the genome-wide distribution of

4. synonymous site diversity, we used genomic data for 12 lines generated by

5. McGaugh et al. 2012. We mapped short reads to the longest ORF for each gene in

6. the R3.2 gene set using Bowtie2 v2.1.0 (Langmead et al. 2009), and estimated

7. synonymous site diversity ($\theta_W$ based on fourfold synonymous sites) at each ORF

8. using PoPoolation (Kofler et al. 2011). We then plotted the distribution of

9. synonymous site diversity, limited to genes in the size range of 0.75kb - 3kb for

10. comparability with the Ago2 paralogues, and compared the fourfold synonymous site

11. diversity levels of each *D. pseudoobscura* Ago2 paralogue with this distribution.

12. Some *D. pseudoobscura* paralogues are located on autosomes (Ago2b, Ago2c &

13. Ago2e) and some on the X chromosome (Ago2a1, Ago2a3 & Ago2d). Therefore,

14. because of the different population genetic expectations for autosomal and X-linked

15. genes (Vicoso and Charlesworth 2006), we examined separate distributions for

16. autosomal and X-linked genes. To provide an additional test for reduced diversity at

17. *D. pseudoobscura* Ago2 paralogues, we performed maximum-likelihood Hudson-

18. Kreitman-Aguadé tests (Wright and Charlesworth 2004), using divergence from *D.*

19. *affinis* and intraspecific polymorphism data for 84 *D. pseudoobscura* loci generated

20. by Haddrill et al. 2010. We performed 63 tests to encompass all one, two, three, four,

21. five and six-way combinations of the paralogues, and calculated Akaike weights from

22. the resulting likelihood estimates to provide an estimate of the level of support for

23. each combination.

24. To infer a genome-wide distribution of synonymous site diversity for *D. obscura* and

25. *D. subobscura*, for which genomic data are unavailable, we used pooled

14

1    transcriptome data from wild-collected adult male flies that had previously been

2    generated for surveys of RNA viruses (van Mierlo et al. 2014; Webster et al. 2016).

3    To generate a *de novo* transcriptome for each species, we assembled short reads

4    with Trinity r20140717 (Grabherr et al. 2011). For each species, we mapped short

5    reads from the pooled sample to the longest ORF for each transcript, estimated

6    synonymous site diversity at each locus using PoPoolation (Kofler et al. 2011), and

7    plotted the distribution of diversity (as described above for *D. pseudoobscura*). The

8    presence of heterozygous sites in males (identified by Sanger sequencing)

9    confirmed that all Ago2 paralogues in *D. subobscura* and *D. obscura* are autosomal:

10   we therefore compared the synonymous site diversity for these paralogues with the

11   autosomal distribution, and do not show the distributions for putatively X-linked

12   genes. Our use of transcriptome data for *D. obscura* and *D. subobscura* will bias the

13   resulting diversity distributions in three ways. First, variation in expression level will

14   cause individuals displaying high levels of expression to be overrepresented among

15   reads, downwardly biasing diversity. Second, highly expressed genes are easier to

16   assemble, and highly expressed genes tend to display lower genetic diversity (Pal et

17   al. 2001; Lemos et al. 2005). Third, high-diversity genes are harder to assemble, *per*

18   *se*. However, as all three biases will tend to artefactually reduce diversity in the

19   genome-wide dataset relative to Ago2, this makes our finding that Ago2 paralogues

20   display unusually low diversity conservative.

21   Identifying selective sweeps in Ago2 paralogues of *D. pseudoobscura*

22   To test whether the unusually low diversity seen in the *D. pseudoobscura* Ago2

23   paralogues is due to recent selection or generally reduced diversity in that region of

24   the genome, we compared diversity at each paralogue to diversity in their

25   neighbouring regions. We obtained sequence data for the 50kb either side of each of

1 these paralogues from the 11 whole genomes detailed in McGaugh et al. 2012

2 (SRA044960.1, SRA044955.2 & SRA044956.1). Note that the very high similarity of

3 these Ago2 paralogues means that they cannot be accurately assembled from short

4 read data, and are not present in the data from McGaugh et al. 2012. For each

5 genome, we therefore replaced the poorly-assembled region corresponding to the

6 paralogue with one of our own Sanger-sequenced haplotypes, making a set of 11 ca.

7 102kb sequences for each paralogue. We aligned these sequences using PRANK

8 (Löytynoja and Goldman 2005) with default settings, and calculated Watterson's θ at

9 all sites in a sliding window across each alignment, with a window size of 5kb and a

10 step of 1kb. For Ago2a1 and Ago2a3, which are located in tandem, we analysed the

11 same genomic region. Since our Ago2 haplotypes were sampled from a different

12 North American population of *D. pseudoobscura* to those of McGaugh et al. 2012, an

13 apparent reduction in local diversity might result from differences in diversity

14 between the two populations. We therefore also repeated these analyses on a

15 dataset in which our Sanger sequenced haplotypes were removed, leaving missing

16 data.

17 To test explicitly for selective sweeps at each region, we used Sweepfinder (Nielsen,

18 Williamson, et al. 2005) to calculate the likelihood and location of a sweep in or near

19 each Ago2 paralogue. We specified a grid size of 20,000, a folded frequency

20 spectrum for all sites, and included invariant sites. To infer the significance of any

21 observed peaks in the composite likelihood ratio, we used ms (Hudson 2002) to

22 generate 1000 samples of 11 sequences under a neutral coalescent model. We

23 generated separate samples for each region surrounding an Ago2 paralogue,

24 conditioning on the number of polymorphic sites observed in that region, the

25 sequence length equal to the alignment length, and an effective population size of

1    $10^6$ (based on a previous estimate for *D. melanogaster* by Li and Stephan 2006). We

2    specified the recombination rate at 5cM/Mb, a conservative value based on previous

3    estimates for *D. pseudoobscura* (McGaugh et al. 2012), which will lead to larger

4    segregating linkage groups and therefore a more stringent significance threshold.

5

6    <u>Results</u>

7    <u>Ago2 has undergone numerous ancient and recent duplications in the *obscura* group</u>

8    Ago2 duplications had previously been noted in *D. pseudoobscura* (Hain et al. 2010),

9    but their age and distribution in other species was unknown. We used BLAST

10    (Altschul et al. 1997) and PCR to identify 65 Ago2 homologues in 39 species

11    sampled across the Drosophilidae, including 30 homologues in 9 *obscura* group

12    species. Using PCR and Sanger sequencing, we verified that the paralogues in *D.*

13    *subobscura*, *D. obscura* and *D. pseudoobscura* are genuine distinct loci, and not

14    artefacts of erroneous assembly. Additionally, we verified that all paralogues

15    possess introns, and so are most likely to be the product of segmental duplication

16    rather than retrotransposition. This is perhaps unsurprising given that segmental

17    duplicates are generally retained at a higher rate than retrotransposed duplicates,

18    despite the rate of retrotransposition being higher than segmental duplication (Hahn,

19    2009).

20    To characterize the relationships between Ago2 homologues in the *obscura* group

21    and the other Drosophilidae, and estimate the date of the duplication events that

22    produced them, we carried out a strict clock Bayesian phylogenetic analysis (Figure

23    1). This showed that there are early diverging Ago2 clades in the *obscura* group: the

24    Ago2e subclade that diverged from other Ago2 paralogues around 21mya (±10 My),

17

1 and the Ago2a and Ago2f subclades that were produced by a gene duplication event

2 around 16mya (±7 My). Subsequently there have been a series of more recent

3 duplications in the *D. pseudoobscura* subgroup Ago2a-d lineage. Using published

4 genomes, transcriptomes and PCR we were unable to identify Ago2e in *D.*

5 *subobscura*, Ago2e or Ago2f in *D. lowei*, or Ago2f in *D. pseudoobscura*, *D. persimilis*

6 and *D. azteca*. While apparent losses may reflect a lack of genomic data (*D.*

7 *subobscura*, *D. lowei* and *D. azteca*), incomplete genome assemblies (*D.*

8 *pseudoobscura* and *D. persimilis*) or unexpressed genes in transcriptome surveys,

9 we attempted to validate the losses observable in *D. pseudoobscura* and *D.*

10 *subobscura* by extensive PCR, and were again unable to recover these genes from

11 those two species.

12 In release 3.03 of the *D. pseudoobscura* genome the paralogues Ago2b-Ago2e have

13 confirmed locations, but Ago2a1 and Ago2a3 (the very recent paralogues newly

14 identified here) lie in tandem on an unplaced contig with a third incomplete copy

15 (Ago2a2) between them. We used PCR to confirm the existence, orientation, and

16 relative positioning of these genes, and to identify the location of this contig, which

17 lies in reverse orientation on chromosome XL-group1a (predicted coordinates

18 3,463,701-3,489,689). We then combined this information with our phylogenetic

19 analysis to reconstruct the positional evolution of *D. pseudoobscura* Ago2

20 paralogues (Figure S1). We found that *D. pseudoobscura* Ago2d is syntenic with *D.*

21 *melanogaster* Ago2, indicating that Ago2d is the ancestral paralogue in this species.

22 We also found that Ago2 paralogues have translocated throughout the *D.*

23 *pseudoobscura* genome (Figure S1), and are situated on autosomes (Ago2b, Ago2c

24 & Ago2e) and both arms of the X chromosome (Ago2a1, Ago2a3 & Ago2d). It should

25 be noted that a lack of genomic data precludes similar synteny analysis for any other

18

1    *obscura* group species; our naming of the Ago2 paralogues in these species as

2    Ago2a (or Ago2a and Ago2b in the case of *D. affinis* and *D. azteca*) reflects their

3    position within the Ago2a subclade, rather than a syntenic relationship or otherwise

4    with *D. pseudoobscura* Ago2a1 and Ago2a3.

5    <u>Ago2 paralogues in *D. subobscura*, *D. obscura* and *D. pseudoobscura* are probably</u>

6    <u>functional</u>

7    Our phylogenetic analysis (Figure 1) revealed that the Ago2 paralogues in the

8    *obscura* group have retained coding sequences for millions of generations, showing

9    that they have remained functional for this period. They have also retained PAZ and

10    PIWI domains and a bilobal structure (characteristic of Argonaute proteins),

11    suggesting that they are part of a functional RNAi pathway. In *D. melanogaster* Ago2

12    plays a key role in antiviral immunity, but is ubiquitously and highly expressed in both

13    males and females, and is not strongly induced by viral challenge (Figure 2a, Aliyari

14    et al. 2008). To test whether this expression pattern has been conserved after Ago2

15    duplication, or whether any Ago2 paralogues have become inducible by viral

16    challenge, we measured the expression of each Ago2 paralogue in female and male

17    *D. subobscura*, *D. obscura* and *D. pseudoobscura* after infection with *Drosophila C*

18    *Virus* (DCV). These species are separated by ~10My of evolution, and represent the

19    three major clades within the *obscura* group. Members of the *obscura* group are

20    highly susceptible to DCV, supporting high viral titres and displaying rapid mortality

21    (Longdon et al. 2015). We found that only one paralogue is expressed in both sexes

22    at a high level in *D. subobscura* (Ago2a), *D. obscura* (Ago2a) and *D. pseudoobscura*

23    (Ago2c). These paralogues show a similar pattern of expression to *D. melanogaster*

24    Ago2, being expressed constitutively throughout the timecourse rather than induced

25    by viral infection (Figure 2). Unexpectedly, and with only one exception, the other

19

1  Ago2 paralogues in all species were expressed exclusively in males (Figure 2b-d),

2  raising the possibility that these duplicates have specialised to a sex-specific role.

3  The one exception was *D. pseudoobscura* Ago2d, which is the ancestral paralogue

4  in this species (inferred by synteny), and for which we could not detect any

5  expression.

6  <u>Ago2 paralogues have repeatedly specialised to the testis</u>

7  To determine whether the strongly male-biased expression pattern is associated with

8  a testis-specific role, we quantified the tissue-specific expression patterns of Ago2

9  paralogues in *D. subobscura*, *D. obscura* and *D. pseudoobscura*. In *D. melanogaster*

10  the single copy of Ago2 was expressed in all adult tissues (Figure 3d), and

11  transcripts were present in the embryo (Figure S2). In *D. subobscura*, *D. obscura*

12  and *D. pseudoobscura*, we found that the Ago2 paralogues exhibited striking

13  differences in their tissue-specific patterns of expression (Figure 3a-c). In each

14  species, one paralogue has retained the ancestral ubiquitous expression pattern in

15  adult tissues. In contrast, every other paralogue was expressed only in the testis,

16  except for the non-expressed *D. pseudoobscura* Ago2d. None of the testis-specific

17  paralogues in *D. pseudoobscura* was detectable in embryos (Figure S2).

18  Interestingly, the ubiquitously expressed paralogue in *D. subobscura* and *D. obscura*

19  is the ancestral gene (Ago2a in both cases, as inferred by synteny with *D.*

20  *melanogaster*), but in *D. pseudoobscura* another paralogue (Ago2c) has evolved the

21  ubiquitous expression pattern, and the ancestral gene (Ago2d) was not expressed at

22  a detectable level in any tissue. When interpreted in the context of the phylogenetic

23  relationships between these paralogues, the most parsimonious explanation is that

24  testis-specificity evolved at least three times: first at the base of the Ago2e clade,

1    second at the base of the Ago2f clade, and third at the base of the *D.*

2    *pseudoobscura-D. persimilis* Ago2a-Ago2b subclade (Figure 1).

3    <u>Testis-specificity is associated with faster protein evolution</u>

4    To test for differences in evolutionary rate between testis-specific and ubiquitously

5    expressed Ago2 paralogues, we fitted sequence evolution models to the set of

6    drosophilid Ago2 sequences depicted in Figure 1 using codeml (PAML, Yang 1997).

7    These tests estimate separate dN/dS ratios ($\omega$) for different subclades in the gene

8    tree, providing a test for differential rates of protein evolution. We found that most

9    support (Akaike weight = 0.99) falls behind a model specifying a different $\omega$ for each

10    obscura group Ago2 subclade, and another separate $\omega$ for the *D. pseudoobscura-D.*

11    *persimilis* Ago2a-Ago2b subclade. Under this model, the testis-specific *D.*

12    *pseudoobscura-D. persimilis* Ago2a-Ago2b subclade has the highest rate of protein

13    evolution ($\omega$=0.32±0.047 SE), followed by the testis-specific Ago2f subclade

14    ($\omega$=0.21±0.014), the ubiquitous Ago2a subclade ($\omega$=0.19±0.012), the testis-specific

15    Ago2e subclade ($\omega$=0.16±0.010), and finally the other Drosophilid Ago2 sequences

16    ($\omega$=0.12±0.002). This shows that the evolution of testis-specificity was accompanied

17    by an increase in the rate of protein evolution following two of the three duplications.

18    We also used the Bayes Empirical Bayes sites test in codeml to identify codons

19    evolving under positive selection across the entire gene tree, and the branch-sites

20    test to identify codons under positive selection in the *obscura* group Ago2 subclade.

21    While we found no positively-selected codons with the sites test, we identified three

22    codons under positive selection (297, 338 & 360) in the *obscura* group Ago2

23    subclade with the branch-sites test (likelihood ratio test M8 vs M8a, p<0.005).

24    <u>McDonald-Kreitman tests identify strong positive selection on *D. pseudoobscura*</u>

25    <u>Ago2e</u>

1 Changes in evolutionary rate after the evolution of testis-specificity may occur as a

2 result of changes in positive selection, or changes in selective constraint. However,

3 unless there are multiple substitutions within single codons, this will be hard to detect

4 using methods such as codeml. Therefore, as a second test for positive selection on

5 Ago2 paralogues in *D. subobscura*, *D. obscura* and *D. pseudoobscura*, we gathered

6 intraspecies polymorphism data for each Ago2 paralogue in these species (Appendix

7 S4), and performed McDonald-Kreitman (MK) tests (Table S1). The MK test uses a

8 comparison of the numbers of fixed differences between species at nonsynonymous

9 (Dn) and synonymous (Ds) sites, and polymorphisms within a species at

10 nonsynonymous (Pn) and synonymous (Ps) sites to infer the action of positive

11 selection. If all mutations are either neutral or strongly deleterious, the Dn/Ds ratio

12 should be approximately equal to the Pn/Ps ratio; however, if there is positive

13 selection, an excess of nonsynonymous differences is expected (McDonald and

14 Kreitman 1991). The majority of MK tests were non-significant (Fisher's exact test,

15 p>0.1), despite often displaying relatively high $K_A/K_S$ ratios e.g. *D. pseudoobscura*

16 Ago2a1 ($K_A/K_S$ =0.34), Ago2b ($K_A/K_S$ =0.43) & Ago2d ($K_A/K_S$ =0.36). However, the

17 low diversity at these loci (<10 polymorphic sites in most cases; see below) means

18 that the MK approach has little power, and that estimates of the proportion of

19 substitutions that are adaptive ($\alpha$) are likely to be poor. In contrast to the other loci,

20 our MK analysis identified strong positive selection acting on *D. pseudoobscura*

21 Ago2e – which has relatively high genetic diversity – with $\alpha$ at 100% ($\alpha$=1.00;

22 Fisher's exact test, p=0.0004). This result is driven by the extreme dearth of

23 nonsynonymous to synonymous polymorphisms (0 Pn to 17 Ps), despite substantial

24 numbers of fixed differences (77 Dn to 120 Ds), and its statistical significance is

25 robust to the choice of outgroup (Table S2).

1 <u>The majority of Ago2 paralogues have extremely low levels of sequence diversity</u>

2 When strong selection acts to reduce genetic diversity at a locus, it can also reduce

3 diversity at linked loci before recombination can break up linkage (Maynard Smith

4 and Haigh 1974). Recent positive selection can therefore be inferred from a

5 reduction in synonymous site diversity compared with other genes. Because MK

6 tests can only detect multiple long-term substitutions, and are hampered by low

7 diversity, diversity-based approaches offer a complementary way to detect very

8 recent strong selection. We therefore compared the synonymous site diversity at

9 each Ago2 paralogue in *D. pseudoobscura* with the distribution of genome-wide

10 synonymous site diversity. We found that all *D. pseudoobscura* paralogues have

11 unusually low diversity relative to other loci: Ago2a1, Ago2b and Ago2c fall into the

12 lowest percentile, Ago2a3 and Ago2d into the 2nd lowest percentile and Ago2e into

13 the 8th lowest percentile (Figure S4). A multi-locus extension of the HKA test (ML-

14 HKA, Wright and Charlesworth 2004) confirmed that the diversity of Ago2a1-Ago2e

15 is significantly lower than the *D. pseudoobscura* genome as a whole (Akaike weight

16 = 0.98).

17 Unfortunately, population-genomic data are not available for *D. subobscura* and *D.*

18 *obscura*, preventing a similar analysis. However, we found similar results for Ago2a

19 and Ago2e when comparing the diversity of *D. subobscura* and *D. obscura* Ago2

20 paralogues to levels of diversity inferred from transcriptome data (data from Webster

21 et al. 2016), suggesting that this effect is not limited to *D. pseudoobscura* and these

22 genes may therefore have been recent targets of selection in multiple species. In *D.*

23 *obscura*, Ago2a and Ago2e fall into the 2nd and 4th lowest diversity percentile

24 respectively, whereas Ago2f falls into the 19th percentile (Figure S4). In *D.*

25 *subobscura*, Ago2a falls into the 7th percentile, whereas Ago2f falls into the 16th

1 percentile (Figure S4). The prevalence of low intraspecific diversity for testis-specific

2 paralogues is consistent with recent selective sweeps, suggesting that positive

3 selection, not merely relaxation of constraint, has contributed to the increased

4 evolutionary rate seen after specialization to the testis.

5 Four out of six *D. pseudoobscura* Ago2 duplicates show a strong signature of recent

6 hard selective sweeps

7 The impact of selection on linked diversity (a selective sweep) is expected to leave a

8 characteristic footprint in local genetic diversity around the site of selection, and this

9 forms the basis of explicit model-based approaches to detect the recent action of

10 positive selection (Nielsen, Bustamante, et al. 2005). For *D. pseudoobscura*,

11 population genomic data for 11 haplotypes is available from McGaugh et al. 2012,

12 permitting an explicit model-based test for recent hard selective sweeps near to

13 Ago2 paralogues. We therefore combined our Ago2 data with 111kb long haplotypes

14 from McGaugh et al. 2012 to analyse the neighbouring region around each

15 paralogue. Ago2a1 and Ago2a3 form a tandem repeat, and were therefore analysed

16 together as a single potential sweep. We found strong evidence for recent selective

17 sweeps at or very close to Ago2a1/3, Ago2b and Ago2c, which display sharp troughs

18 in their diversity levels, and large peaks in the composite likelihood of a sweep,

19 which far exceed a significance threshold derived from coalescent simulation

20 (p<0.01; Figure 4). These localised reductions in diversity remain when our own

21 Ago2 haplotype data are removed, showing the results are robust to the fact that our

22 Ago2 sequence data are derived from a different population to the genome-wide

23 data of McGaugh et al. 2012 (Figure S6; note that sequence data for Ago2

24 paralogues cannot be derived from the data of McGaugh et al. 2012, because of

25 their extreme similarity). In addition, there is ambiguous evidence for a sweep at

24

1    Ago2d, in the form of one significant (p<0.01) likelihood peak just upstream of the

2    paralogue, but two other peaks ~1kb and ~3kb further upstream. There is no

3    evidence for a hard sweep at Ago2e, which has no diversity trough or likelihood

4    peak.

5

6    <u>Discussion</u>

7    <u>Testis-specificity may indicate a loss of antiviral function</u>

8    We have found that Ago2 paralogues in the *obscura* group have repeatedly evolved

9    divergent expression patterns after duplication, with the majority of paralogues

10    specializing to the testis. This is the first report of testis-specificity for any arthropod

11    Ago2, which is ubiquitously expressed in *D. melanogaster* (Celniker et al. 2009), and

12    provides a strong indication that these paralogues have diverged in function. This

13    testis-specificity (Figure 3) suggests that these Argonautes are likely to have lost

14    their ancestral ubiquitous antiviral role. Additionally, the constant level of expression

15    of testis-specific paralogues under DCV infection (Figure 2) suggests that have not

16    evolved an inducible response to viral infection, either restricted to the testis or in

17    other tissues. In contrast, one paralogue in each species has retained the ubiquitous

18    expression pattern seen in *D. melanogaster* (*D. subobscura* Ago2a, *D. obscura*

19    Ago2a & *D. pseudoobscura* Ago2c, Figure 3), suggesting that these paralogues

20    have retained roles in antiviral defence (Li et al. 2002; van Rij et al. 2006), dosage

21    compensation (Menon and Meller 2012) and/or somatic TE suppression (Chung et

22    al. 2008; Czech et al. 2008).

23    <u>Both ubiquitous and testis-specific Ago2 paralogues show evidence of recent</u>

24    <u>positive selection</u>

1   We identified selective sweeps at the ubiquitously expressed Ago2 paralogue in *D.*

2   *pseudoobscura* Ago2c, and very low diversity in the ubiquitously expressed Ago2

3   paralogues of *D. subobscura* and *D. obscura* (Ago2a), suggesting that all of these

4   genes may have recently experienced strong positive selection. Four randomly-

5   chosen testis-specific genes in *D. obscura* and *D. subobscura* do not fall into the

6   low-diversity tails of the genome-wide diversity distributions, suggesting that this is

7   not a general phenomenon of testis-specific expression. This is consistent with

8   previous findings of strong selection and rapid evolution of Ago2 in *D. melanogaster*

9   (Obbard et al. 2006; Obbard, Welch, et al. 2009; Obbard et al. 2011) which has also

10  experienced recent sweeps in *D. melanogaster*, *D. simulans*, and *D. yakuba* (Obbard

11  et al. 2011), and across the *Drosophila* more broadly (Kolaczkowski et al. 2011). It

12  has previously been suggested that this is driven by arms-race coevolution with

13  viruses (Obbard, Gordon, et al. 2009; Kolaczkowski et al. 2011), some of which

14  encode viral suppressors of RNAi (VSRs) that block Ago2 function (Bronkhorst and

15  van Rij 2014). The presence of VSR-encoding viruses, such as Nora virus, in natural

16  *obscura* group populations (Webster et al. 2016), combined with the host-specificity

17  that can be displayed by VSRs (van Mierlo et al. 2014), suggest that arms-race

18  dynamics may also be driving the rapid evolution of ubiquitously expressed Ago2

19  paralogues in the *obscura* group.

20  Potential testis-specific functions

21  In contrast to their ancestral ubiquitous expression pattern, the dominant fate for

22  Ago2 paralogues in the *obscura* group appears to have been specialization to the

23  testis. Paralogues often undergo a brief period of testis-specificity soon after

24  duplication (Assis and Bachtrog 2013; Assis and Bachtrog 2015), and this has given

25  rise to the 'out-of-the-testis' hypothesis, in which new paralogues are initially testis-

26

1   specific before evolving functions in other tissues (Kaessmann 2010). However, two

2   lines of evidence suggest an adaptive basis for the testis-specificity observed for the

3   *obscura* group Ago2 paralogues. First, testis-specificity has been retained for more

4   than 10 million years in Ago2e and Ago2f, in contrast to the broadening of

5   expression over time expected under the out-of-the-testis hypothesis (Kaessmann

6   2010; Assis and Bachtrog 2013). Second, all testis-specific Ago2 paralogues in *D.*

7   *pseudoobscura* show evidence either of long-term positive selection (MK test for the

8   high-diversity Ago2e) or of recent selective sweeps (in low-diversity Ago2a1/3 and

9   Ago2b), and the testis-specific *D. obscura* Ago2e displays a reduction in diversity,

10  potentially driven by selection.

11  Under a subfunctionalization model for Ago2 testis-specialization, five candidate

12  selective pressures seem likely: testis-specific dosage compensation, antiviral

13  defence, gene regulation, TE suppression, and/or the suppression of meiotic drive.

14  Of these, testis-specific dosage compensation seems the least likely to drive testis-

15  specificity because the male-specific lethal (MSL) complex, which Ago2 directs to X-

16  linked genes to carry out dosage compensation in the soma of *D. melanogaster*, is

17  absent from testis (Conrad and Akhtar 2012). Testis-specific antiviral defence seems

18  similarly unlikely, as the only known paternally-transmitted *Drosophila* viruses

19  (Sigmaviruses; Rhabdoviridae) pass through both the male and female gametes

20  (Longdon and Jiggins 2012), and so the potential benefits of testis-specificity seem

21  unclear. Alternatively, testis-specific Ago2 duplicates could be co-evolving with other

22  testis-specific genes through the hairpin RNA pathway, in which siRNAs generated

23  from endogenous hairpin-forming RNAs (hpRNAs) bind Ago2 and regulate the

24  expression of host genes (Okamura et al, 2008). In *D. melanogaster*, hpRNA-derived

25  siRNAs target testis-specific genes involved in male fertility, and coevolve with these

1  targets to maintain base complementarity (Wen et al, 2015). If a similar pathway

2  operates in the *obscura* group, Ago2 paralogues could have specialized to the

3  hpRNA pathway in order to regulate testis-specific genes more effectively.

4  Finally, the suppression of TEs or meiotic drive seem promising candidate selective

5  forces. First, numerous TEs transpose preferentially in the testis, such as *Penelope*

6  in *D. virilis* (Rozhkov et al. 2010) and *copia* in *D. melanogaster* (Pasyukova et al.

7  1997; Morozova et al. 2009), which could impose a selection pressure on Ago2

8  paralogues to provide a testis-specific TE suppression mechanism. It should be

9  noted that all members of the canonical anti-TE Piwi subfamily (Ago3, Aub and Piwi)

10 are also expressed in *obscura* group testis (Figure S3), suggesting that if Ago2

11 paralogues have specialised to suppress TEs, they are doing so alongside the

12 existing TE suppression mechanism. Second, testis-specificity could have evolved to

13 suppress meiotic drive, which is prevalent (in the form of sex-ratio distortion) in the

14 *obscura* group (Gershenson 1928; Sturtevant and Dobzhansky 1936; Wu and

15 Beckenbach 1983; Jaenike 2001; Unckless et al. 2015), and which is suppressed by

16 RNAi-based mechanisms in other species (Tao et al. 2007; Kotelnikov et al. 2009;

17 Gell and Reenan 2013). A high level of meiotic drive in the *obscura* group could

18 therefore impose selection for the evolution of novel suppression mechanisms,

19 leading to the repeated specialization of Ago2 paralogues to the testis.

20 <u>Prospects for novel functions during the evolution of RNAi</u>

21 The functional specialization that we observe for *obscura* group Ago2 paralogues

22 raises the prospect of undiscovered derived functions following Argonaute

23 expansions in other lineages. Ago2 has duplicated frequently across the arthropods,

24 with expansions present in insects (*Drosophila willistoni* (Figure 1) & *Musca*

25 *domestica*, Scott et al. 2014), crustaceans (*Penaeus monodon*, Leebonoi et al. 2015)

1  and chelicerates (*Tetranychus urticae*, *Ixodes scapularis*, *Mesobuthus martensii* &

2  *Parasteatoda tepidariorum*, Palmer and Jiggins 2015). The prevalence of testis-

3  specificity in *obscura* group Ago2 paralogues raises the possibility that specialization

4  to the germline may be more widespread following Argonaute duplication. The

5  expression of Ago2 paralogues has previously been characterized in *P. monodon*,

6  and shows that one paralogue has indeed specialised to the germline of both males

7  and females, but not the testis alone (Leebonoi et al. 2015). Publicly available

8  RNAseq data from the head, gonad and carcass of male and female *Musca*

9  *domestica* (GSE67065, Meisel et al. 2015) suggests that neither *M. domestica* Ago2

10  paralogue has specialised to the testis (Figure S8). However, public data from the

11  head, thorax and abdomen of male and female *D. willistoni* (GSE31723, Meisel et al.

12  2012) shows that one *D. willistoni* Ago2 paralogue (FBgn0212615) is expressed

13  ubiquitously, while the other (FBgn0226485) is expressed only in the male abdomen

14  (Figure S8), consistent with the evolution of testis-specificity after duplication. This

15  raises the possibility that a testis-specific selection pressure may be driving the

16  retention and specialization of Ago2 paralogues across the arthropods.

17  In conclusion, we have identified rapid and repeated evolution of testis-specificity

18  after the duplication of Ago2 in the *obscura* group, associated with low genetic

19  diversity and signatures of strong selection. Ago2 and other RNAi genes have

20  undergone frequent expansions in different eukaryotic lineages (Mukherjee et al.

21  2013; Lewis et al. 2016), and have been shown to switch between ubiquitous and

22  germline- or ovary-specific functions in isolated species. This study provides

23  evidence for the evolution of a new testis-specific RNAi function, and suggests that

24  positive selection may act on young paralogues to drive the rapid evolution of novel

25  RNAi mechanisms across the eukaryotes.

1

## Acknowledgements

## References

Aliyari R, Wu Q, Li H-W, Wang X-H, Li F, Green LD, Han CS, Li W-X, Ding S-W. 2008. Mechanism of induction and suppression of antiviral immunity directed by virus-derived small RNAs in Drosophila. Cell Host Microbe 4:387–397.

Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ. 1997. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. Nucleic Acids Res. 25:3389–3402.

Assis R, Bachtrog D. 2013. Neofunctionalization of young duplicate genes in Drosophila. Proc. Natl. Acad. Sci. 110:17409–17414.

Assis R, Bachtrog D. 2015. Rapid divergence and diversification of mammalian duplicate gene functions. BMC Evol. Biol. 15:138.

Behura SK, Severson DW. 2013. Codon usage bias: causative factors, quantification

methods and genome-wide patterns: with emphasis on insect genomes. Biol. Rev. 88:49–61.

Bronkhorst AW, van Rij RP. 2014. The long and short of antiviral defense: small RNA-based immunity in insects. Curr. Opin. Virol. 7C:19–28.

Brown JB, Boley N, Eisman R, May GE, Stoiber MH, Duff MO, Booth BW, Wen J, Park S, Suzuki AM, et al. 2014. Diversity and dynamics of the Drosophila transcriptome. Nature 512:393–399.

Buck AH, Blaxter M. 2013. Functional diversification of Argonautes in nematodes: an expanding universe. Biochem. Soc. Trans. 41:881–886.

Celniker SE, Dillon LAL, Gerstein MB, Gunsalus KC, Henikoff S, Karpen GH, Kellis M, Lai EC, Lieb JD, Macalpine DM, et al. 2009. Unlocking the secrets of the genome. Nature 459:927–930.

Cerutti H, Casas-Mollano JA. 2006. On the origin and functions of RNA-mediated silencing: from protists to man. Curr. Genet. 50:81–99.

Chung W-J, Okamura K, Martin R, Lai EC. 2008. Endogenous RNA interference provides a somatic defense against Drosophila transposons. Curr. Biol. 18:795–802.

Conrad T, Akhtar A. 2012. Dosage compensation in Drosophila melanogaster: epigenetic fine-tuning of chromosome-wide transcription. Nat. Rev. Genet. 13:123–134.

Czech B, Malone CD, Zhou R, Stark A, Schlingeheyde C, Dus M, Perrimon N, Kellis M, Wohlschlegel JA, Sachidanandam R, et al. 2008. An endogenous small interfering RNA pathway in Drosophila. Nature 453:798–802.

Drummond AJ, Suchard MA, Xie D, Rambaut A. 2012. Bayesian phylogenetics with

BEAUti and the BEAST 1.7. Mol. Biol. Evol. 29:1969–1973.

Eulalio A, Huntzinger E, Izaurralde E. 2008. Getting to the Root of miRNA-Mediated Gene Silencing. Cell 132:9–14.

Finn RD, Mistry J, Tate J, Coggill P, Heger a., Pollington JE, Gavin OL, Gunasekaran P, Ceric G, Forslund K, et al. 2009. The Pfam protein families database. Nucleic Acids Res. 38:D211–D222.

Gell SL, Reenan RA. 2013. Mutations to the piRNA pathway component aubergine enhance meiotic drive of segregation distorter in Drosophila melanogaster. Genetics 193:771–784.

Gershenson S. 1928. A New Sex-Ratio Abnormality in Drosophila obscura. Genetics 13:488–507.

Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q, et al. 2011. Full-length transcriptome assembly from RNA-Seq data without a reference genome. Nat. Biotechnol. 29:644–652.

Haddrill PR, Loewe L, Charlesworth B. 2010. Estimating the parameters of selection on nonsynonymous mutations in Drosophila pseudoobscura and D. miranda. Genetics 185:1381–1396.

Hahn M. 2009. Distinguishing among evolutionary models for the maintenance of gene duplicates. Heredity 100:605-617

Hain D, Bettencourt BR, Okamura K, Csorba T, Meyer W, Jin Z, Biggerstaff J, Siomi H, Hutvagner G, Lai EC, et al. 2010. Natural variation of the amino-terminal glutamine-rich domain in Drosophila argonaute2 is not associated with developmental defects. PLoS One 5:e15264.

Hudson RR. 2002. Generating samples under a Wright-Fisher neutral model of

genetic variation. Bioinformatics 18:337–338.

Jaenike J. 2001. Sex chromosome meiotic drive. Annu. Rev. Ecol. Syst. 32:25–49.

Kaessmann H. 2010. Origins, evolution, and phenotypic impact of new genes. Genome Res. 20:1313–1326.

Katoh K, Misawa K, Kuma K, Miyata T. 2002. MAFFT : a novel method for rapid multiple sequence alignment based on fast Fourier transform. Nucleic Acids Res. 30:3059–3066.

Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S, Cooper A, Markowitz S, Duran C, et al. 2012. Geneious Basic: An integrated and extendable desktop software platform for the organization and analysis of sequence data. Bioinformatics 28:1647–1649.

Kofler R, Orozco-terWengel P, De Maio N, Pandey RV, Nolte V, Futschik A, Kosiol C, Schlötterer C. 2011. PoPoolation: a toolbox for population genetic analysis of next generation sequencing data from pooled individuals. PLoS One 6:e15925.

Kolaczkowski B, Hupalo DN, Kern AD. 2011. Recurrent adaptation in RNA interference genes across the Drosophila phylogeny. Mol. Biol. Evol. 28:1033–1042.

Kotelnikov RN, Klenov MS, Rozovsky YM, Olenina L V., Kibanov M V., Gvozdev V a. 2009. Peculiarities of piRNA-mediated post-transcriptional silencing of Stellate repeats in testes of Drosophila melanogaster. Nucleic Acids Res. 37:3254–3263.

Langmead B, Trapnell C, Pop M, Salzberg SL. 2009. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. Genome Biol. 10:R25.

1   Leebonoi W, Sukthaworn S, Panyim S, Udomkit A. 2015. A novel gonad-specific

2       Argonaute 4 serves as a defense against transposons in the black tiger shrimp

3       Penaeus monodon. Fish Shellfish Immunol. 42:280–288.

4   Lemos B, Bettencourt BR, Meiklejohn CD, Hartl DL. 2005. Evolution of proteins and

5       gene expression levels are coupled in Drosophila and are independently

6       associated with mRNA abundance, protein length, and number of protein-

7       protein interactions. Mol. Biol. Evol. 22:1345–1354.

8   Lewis SH, Salmela H, Obbard DJ. 2016. Duplication and diversification of Dipteran

9       Argonaute genes, and the evolutionary divergence of Piwi and Aubergine.

10      Genome Biol. Evol. 8:507-518.

11  Li H, Li WX, Ding SW. 2002. Induction and suppression of RNA silencing by an

12      animal virus. Science 296:1319–1321.

13  Li H, Stephan W. 2006. Inferring the demographic history and rate of adaptive

14      substitution in Drosophila. PLoS Genet. 2:1580–1589.

15  Librado P, Rozas J. 2009. DnaSP v5: a software for comprehensive analysis of DNA

16      polymorphism data. Bioinformatics 25:1451–1452.

17  Longdon B, Hadfield JD, Day JP, Smith SCL, McGonigle JE, Cogni R, Cao C,

18      Jiggins FM. 2015. The causes and consequences of changes in virulence

19      following pathogen host shifts. PLoS Pathog. 11:e1004728.

20  Longdon B, Jiggins FM. 2012. Vertically transmitted viral endosymbionts of insects:

21      do sigma viruses walk alone? Proc. R. Soc. B 279:3889–3898.

22  Löytynoja A, Goldman N. 2005. An algorithm for progressive multiple alignment of

23      sequences with insertions. Proc. Natl. Acad. Sci. 102:10557–10562.

24  Lu H-L, Tanguy S, Rispe C, Gauthier J-P, Walsh T, Gordon K, Edwards O, Tagu D,

Chang C, Jaubert-Possamai S. 2011. Expansion of genes encoding piRNA-associated argonaute proteins in the pea aphid: diversification of expression profiles in different plastic morphs. PLoS One 6:e28051.

Luteijn MJ, Ketting RF. 2013. PIWI-interacting RNAs: from generation to transgenerational epigenetics. Nat. Rev. Genet. 14:523–534.

Lynch M, Crease TJ. 1990. The analysis of population survey data on DNA sequence variation. Mol. Biol. Evol. 7:377–394.

Marques JT, Carthew RW. 2007. A call to arms: coevolution of animal viruses and host innate immune responses. Trends Genet. 23:359–364.

Maynard Smith J, Haigh J. 1974. The hitch-hiking effect of a favourable gene. Genet. Res. 23:23–35.

McDonald JH, Kreitman M. 1991. Adaptive protein evolution at the Adh locus in Drosophila. Nature 351:652–654.

McGaugh SE, Heil CSS, Manzano-Winkler B, Loewe L, Goldstein S, Himmel TL, Noor MAF. 2012. Recombination modulates how selection affects linked sites in Drosophila. PLoS Biol. 10:e1001422.

Meisel RP, Malone JH, Clark AG. 2012. Disentangling the relationship between sex-biased gene expression and X-linkage. Genome Res. 22:1255–1265.

Meisel RP, Scott JG, Clark AG. 2015. Transcriptome Differences between Alternative Sex Determining Genotypes in the House Fly, Musca domestica. Genome Biol. Evol. 7:2051–2061.

Meister G. 2013. Argonaute proteins: functional insights and emerging roles. Nat. Rev. Genet. 14:447–459.

Menon DU, Meller VH. 2012. A role for siRNA in X-chromosome dosage compensation in Drosophila melanogaster. Genetics 191:1023–1028.

van Mierlo JT, Overheul GJ, Obadia B, van Cleef KWR, Webster CL, Saleh M-C, Obbard DJ, van Rij RP. 2014. Novel Drosophila Viruses Encode Host-Specific Suppressors of RNAi. PLoS Pathog. 10:e1004256.

Miesen P, Girardi E, van Rij RP. 2015. Distinct sets of PIWI proteins produce arbovirus and transposon-derived piRNAs in Aedes aegypti mosquito cells. Nucleic Acids Res. 43:6545–6556.

Morozova T V, Tsybulko EA, Pasyukova EG. 2009. Regularory elements of the copia retrotransposon determine different levels of expression in different organs of males and females of Drosophila melanogaster. Genetika 45:169–177.

Mukherjee K, Campos H, Kolaczkowski B. 2013. Evolution of animal and plant dicers: early parallel duplications and recurrent adaptation of antiviral RNA binding in plants. Mol. Biol. Evol. 30:627–641.

Nielsen R, Bustamante C, Clark AG, Glanowski S, Sackton TB, Hubisz MJ, Fledel-Alon A, Tanenbaum DM, Civello D, White TJ, et al. 2005. A scan for positively selected genes in the genomes of humans and chimpanzees. PLoS Biol. 3:e170.

Nielsen R, Williamson S, Kim Y, Hubisz MJ, Clark AG, Bustamante C. 2005. Genomic scans for selective sweeps using SNP data. Genome Res. 15:1566–1575.

Obbard DJ, Gordon KHJ, Buck AH, Jiggins FM. 2009. The evolution of RNAi as a defence against viruses and transposable elements. Philos. Trans. R. Soc. London Biol. Sci. 364:99–115.

Obbard DJ, Jiggins FM, Bradshaw NJ, Little TJ. 2011. Recent and recurrent selective sweeps of the antiviral RNAi gene Argonaute-2 in three species of Drosophila. Mol. Biol. Evol. 28:1043–1056.

Obbard DJ, Jiggins FM, Halligan DL, Little TJ. 2006. Natural selection drives extremely rapid evolution in antiviral RNAi genes. Curr. Biol. 16:580–585.

Obbard DJ, MacLennan J, Kim KW, Rambaut A, O'Grady PM, Jiggins FM. 2012. Estimating divergence dates and substitution rates in the Drosophila phylogeny. Mol. Biol. Evol. 29:3459–3473.

Obbard DJ, Welch JJ, Kim K-W, Jiggins FM. 2009. Quantifying adaptive evolution in the Drosophila immune system. PLoS Genet. 5:e1000698.

Okamura K, Chung W-J, Ruby JG, Guo H, Bartel D & Lai EC (2008) The Drosophila hairpin RNA pathway generates endogenous short interfering RNAs. Nature 453:803-807.

Pal C, Papp B, Hurst LD. 2001. Highly expressed genes in yeast evolve slowly. Genetics 158:927–931.

Palmer WH, Obbard DJ. 2016. Variation and evolution of the glutamine-rich repeat region of Drosophila Argonaute-2. bioRXiv.

Palmer WJ, Jiggins FM. 2015. Comparative Genomics Reveals the Origins and Diversity of Arthropod Immune Systems. Mol. Biol. Evol. 32:2111–2129.

Pasyukova E, Nuzhdin S, Li W, Flavell AJ. 1997. Germ line transposition of the copia retrotransposon in Drosophila melanogaster is restricted to males by tissue-specific control of copia RNA levels. Mol. Gen. Genet. 255:115–124.

Peden J. 1995. Analysis of codon usage bias. PhD Thesis, University of Nottingham.

van Rij RP, Saleh M-C, Berry B, Foo C, Houk A, Antoniewski C, Andino R. 2006. The RNA silencing endonuclease Argonaute 2 mediates specific antiviral immunity in Drosophila melanogaster. Genes Dev. 20:2985–2995.

Ross RJ, Weiner MM, Lin H. 2014. PIWI proteins and PIWI-interacting RNAs in the soma. Nature 505:353–359.

Rozhkov N V, Aravin AA, Zelentsova ES, Schostak NG, Sachidanandam R, McCombie WR, Hannon GJ, Evgen'ev MB. 2010. Small RNA-based silencing strategies for transposons in the process of invading Drosophila species. RNA 16:1634–1645.

Russo C a, Takezaki N, Nei M. 1995. Molecular phylogeny and divergence times of Drosophilid species. Mol. Biol. Evol. 12:391–404.

Schirle NT, Macrae IJ. 2012. The Crystal Structure of Human Argonaute2. Science 336:1037–1040.

Scott JG, Warren WC, Beukeboom LW, Bopp D, Clark AG, Giers SD, Hediger M, Jones AK, Kasai S, Leichter CA, et al. 2014. Genome of the house fly, Musca domestica L., a global vector of diseases with adaptations to a septic environment. Genome Biol. 15:466–482.

Sienski G, Dönertas D, Brennecke J. 2012. Transcriptional silencing of transposons by Piwi and maelstrom and its impact on chromatin state and gene expression. Cell 151:964–980.

Singh RK, Gase K, Baldwin IT, Pandey SP. 2015. Molecular evolution and diversification of the Argonaute family of proteins in plants. BMC Plant Biol. 15:1–16.

Stephens M, Smith NJ, Donnelly P. 2001. A new statistical method for haplotype

reconstruction from population data. Am. J. Hum. Genet. 68:978–989.

Sturtevant AH, Dobzhansky T. 1936. Geographical Distribution and Cytology of "Sex Ratio" in Drosophila Pseudoobscura and Related Species. Genetics 21:473–490.

Tajima F. 1989. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. Genetics 123:585–595.

Tamura K. 2004. Temporal Patterns of Fruit Fly (Drosophila) Evolution Revealed by Mutation Clocks. Mol. Biol. Evol. 21:36–44.

Tao Y, Araripe L, Kingan SB, Ke Y, Xiao H, Hartl DL. 2007. A sex-ratio meiotic drive system in Drosophila simulans. II: An X-linked distorter. PLoS Biol. 5:2576–2588.

Unckless RL, Larracuente AM, Clark AG. 2015. Sex-ratio meiotic drive and Y-linked resistance in Drosophila affinis. Genetics 199:831–840.

Vicoso B, Charlesworth B. 2006. Evolution on the X chromosome: unusual patterns and processes. Nat. Rev. Genet. 7:645–653.

Webster CL, Longdon B, Lewis SH, Obbard DJ. 2016. Twenty five new viruses associated with the Drosophilidae (Diptera). bioRXiv.

Wen J, Duan H, Bejarano F, Okamura K, Fabian L, Brill JA, Bortolamiol-Becet D, Martin R, Ruby JG, Lai EC. 2015. Adaptive Regulation of Testis Gene Expression and Control of Male Fertility by the Drosophila Harpin RNA Pathway. Mol. Cell 57:165–178.

Wright F. 1990. The "effective number of codons" used in a gene. Gene 87:23–29.

Wright SI, Charlesworth B. 2004. The HKA test revisited: a maximum-likelihood-ratio

1    test of the standard neutral model. Genetics 168:1071–1076.

2  Wu CI, Beckenbach AT. 1983. Evidence for extensive genetic differentiation

3    between the sex-ratio and the standard arrangement of Drosophila

4    pseudobscura and D. persimilis and identification of hybrid sterility factors.

5    Genetics 105:71–86.

6  Yang Z. 1997. PAML: a program package for phylogenetic analysis by maximum

7    likelihood. Comput. Appl. Biosci. 13:555–556.

8

Figure 1: An approximately time-scaled Bayesian gene tree of Ago2 in the Drosophilidae. Duplication events are marked by yellow diamonds, Bayesian posterior support is shown for nodes for which it is less than 100%, and the genes and species that are the focus of the present study are marked in bold. Ago2 has duplicated at least twelve times in the Drosophilidae: seven times in the *obscura* group, twice early in the *melanogaster* group, and once each in the lineages leading to *D. willistoni*, *S. deflexa* and *D. kikkawai*. There has also been a potentially recent duplication of Ago2a on the *D. affinis* / *D. azteca* lineage (~5mya), although the low support for this node may suggest that these paralogues could also nest within the *D. pseudoobscura* / *D. persimilis* expansion, with one paralogue sister to the Ago2a-Ago2b subclade and the other sister to the Ago2c-Ago2d subclade. After duplication, Ago2 paralogues in the *obscura* group have specialised to the testis three times independently (marked with ♂), and have been retained for an extended period of time (>10 My in the case of Ago2e), suggesting an adaptive basis for testis-specificity. The labelling a-e of paralogous clades corresponds to Hain et al. 2010, and is retained for consistency with subsequent publications which also use these labels, while clade f is newly reported here. All genes were identified by BLAST, apart from the following which were found by PCR: *D. teissieri* Ago2; *D. santomea* Ago2; *D. azteca* Ago2a, Ago2b & Ago2e; *D. pseudoobscura* Ago2a1 & Ago2a3.
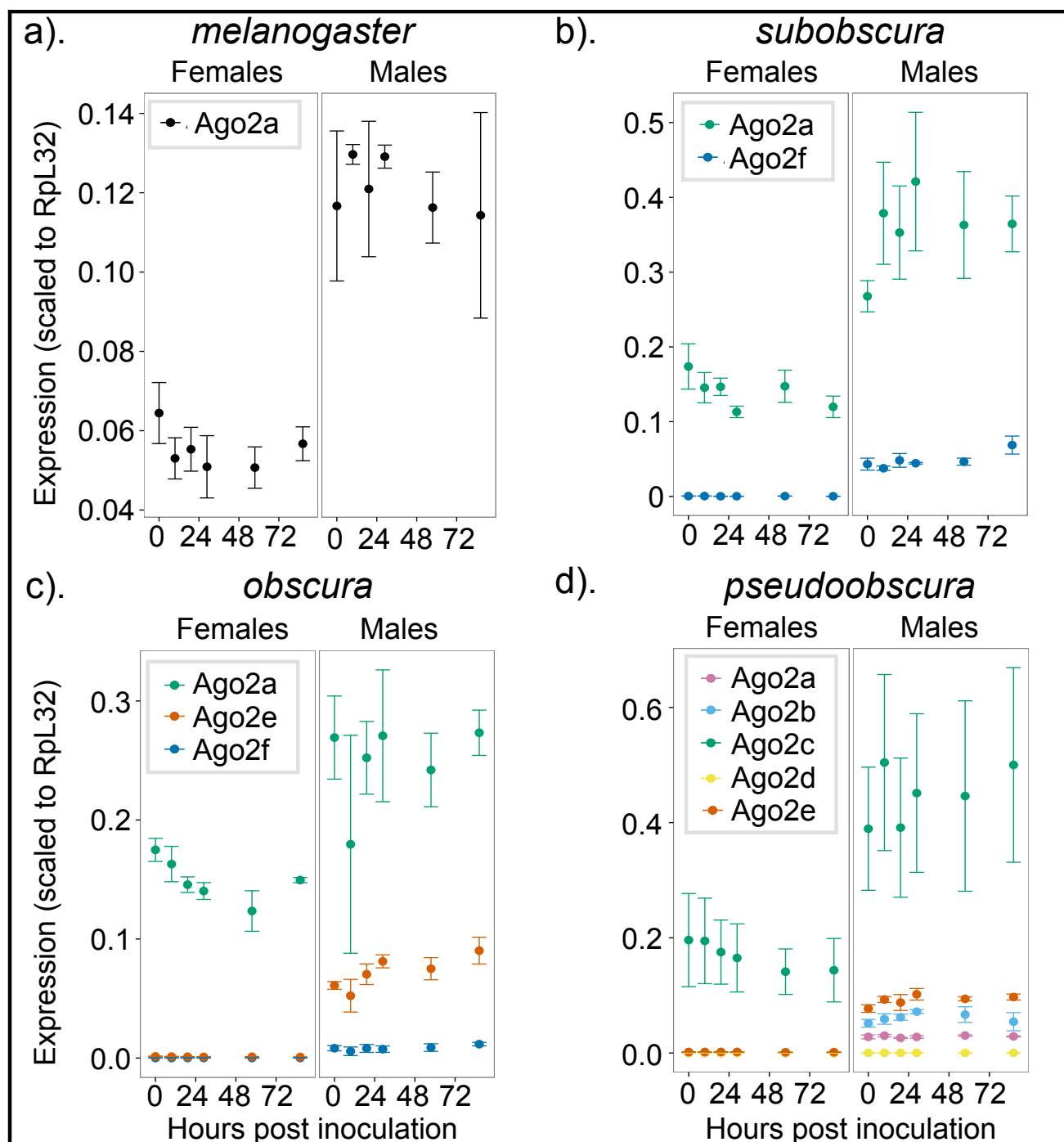
Figure 2: Expression patterns of Ago2 paralogues under challenge with *Drosophila C Virus*. In each *obscura* group species, only one Ago2 paralogue has retained the ancestral pattern of ubiquitous stable expression in each sex (illustrated by *D. melanogaster*). In contrast, all other paralogues are expressed in males only (in *D. pseudoobscura* females, Ago2a, Ago2b, Ago2d & Ago2e are all unexpressed throughout the timecourse). The only exception to this is *D. pseudoobscura* Ago2d, which is unexpressed in either sex. The high degree of sequence similarity between Ago2a1 and Ago2a3 prevented us from amplifying these genes separately in qPCR, and here they are combined as "Ago2a". Error bars indicate 1 standard error estimated from 2 technical replicates in each of three different genetic backgrounds. Apparent differences in expression between sexes and species should be interpreted with caution, as these may be driven by differences in expression levels of the reference gene (RpL32).
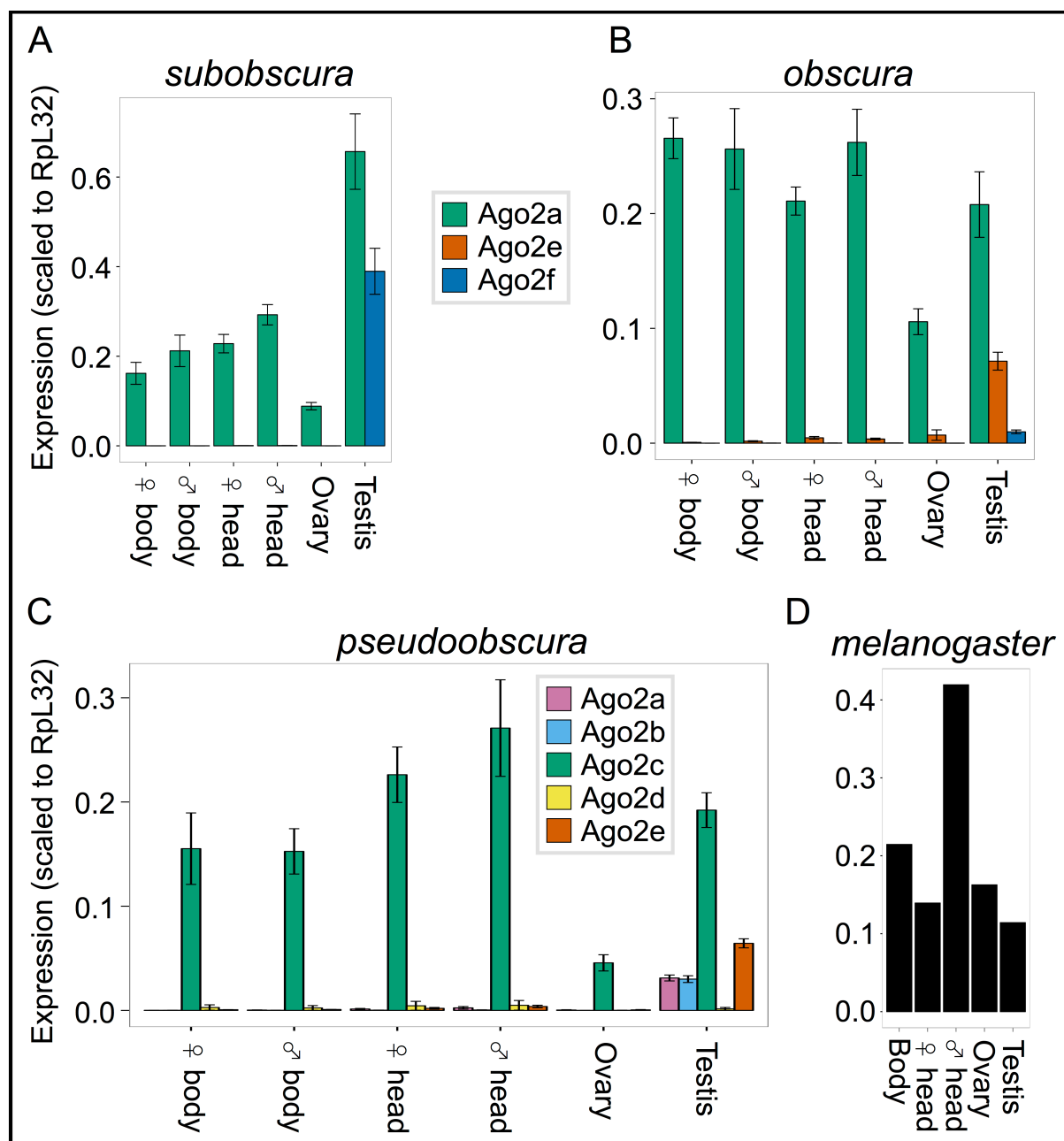
Figure 3: Tissue-specific expression patterns of Ago2 paralogues.
In each of the three *obscura* group species tested, one paralogue has retained the ancestral ubiquitous expression pattern, while the others have specialised to the testis (with the exception of *D. pseudoobscura* Ago2d). The high degree of sequence similarity between Ago2a1 and Ago2a3 prevented us from amplifying these genes separately in qPCR, and here they are combined as "Ago2a". Error bars indicate 1 standard error estimated from 2 technical replicates in each of five different genetic backgrounds. *D. melanogaster* expression levels were taken from a single RNA-seq experiment (Brown et al. 2014).
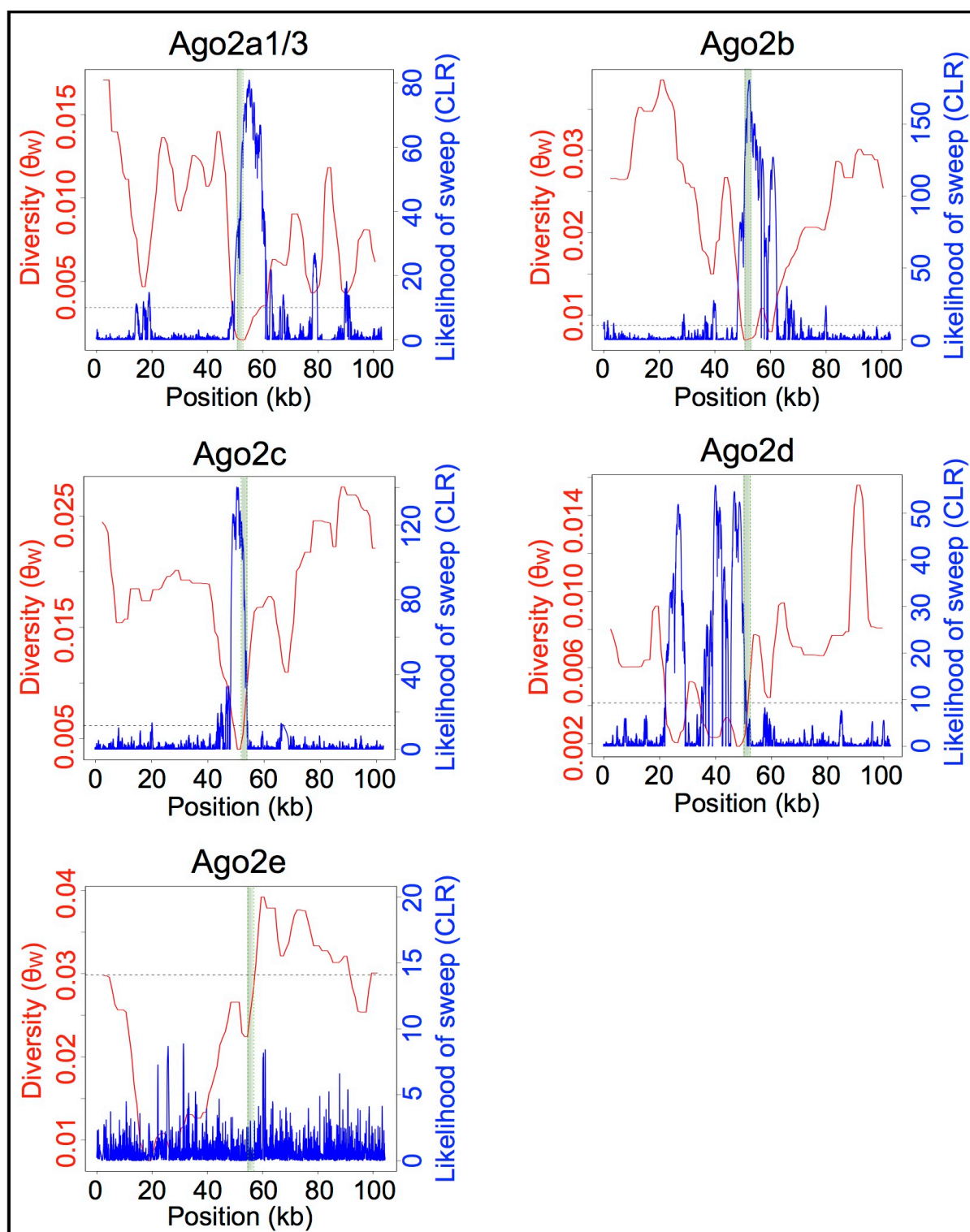
Figure 4: Selective sweeps at *D. pseudoobscura* Ago2 paralogues.
For each paralogue, diversity at all sites (Watterson's θ) is displayed in red, and the likelihood of a sweep centred at that site (composite likelihood ratio, CLR) is displayed in blue. The gene region containing the paralogue is represented by the shaded vertical bar, and the significance threshold for the CLR is displayed by the horizontal dotted line (p<0.01, derived from the 10th-highest CLR out of 1000 coalescent simulations, assuming constant recombination rate and $N_e$). There is strong evidence for sweeps at Ago2a, Ago2b and Ago2c, indicated by troughs in their diversity levels and peaks in the likelihood of a sweep.