1		ARTICLE				
2						
3	Evolι	tion of a cytoplasmic determinant: evidence for the biochemical				
4	basis of functional evolution of a novel germ line regulator					
5						
6	Leo Blondel <sup>1</sup> , Savandara Besse <sup>1,2</sup> and Cassandra G. Extavour <sup>1,3</sup>					
7						
8	1.	Department of Molecular and Cellular Biology, Harvard University, Cambridge				
9		MA, USA				
10	2.	Current address: Department of Biochemistry and Molecular Medicine, Université				
11		de Montréal, Montréal Québec, Canada				
12	3.	Department of Organismic and Evolutionary Biology, Harvard University,				
13		Cambridge MA, USA				
14						
15	Corresponding author: Cassandra G. Extavour extavour@oeb.harvard.edu					
16						
17						
18						
19						

#### 20 Abstract

21 Germ line specification is essential in sexually reproducing organisms. Despite their 22 critical role, the evolutionary history of the genes that specify animal germ cells is 23 heterogeneous and dynamic. In many insects, the gene oskar is required for the 24 specification of the germ line. However, the germ line role of oskar is thought to be a 25 derived role resulting from co-option from an ancestral somatic role. To address how evolutionary changes in protein sequence could have led to changes in the function of 26 27 Oskar protein that enabled it to regulate germ line specification, we searched for oskar 28 orthologs in 1565 publicly available insect genomic and transcriptomic datasets. The earliest-diverging lineage in which we identified an oskar ortholog was the order 29 30 Zygentoma (silverfish and firebrats), suggesting that oskar originated before the origin of winged insects. We noted some order-specific trends in oskar sequence evolution, 31 32 including whole gene duplications, clade-specific losses, and rapid divergence. An alignment of all known 379 Oskar sequences revealed new highly conserved residues as 33 34 candidates that promote dimerization of the LOTUS domain. Moreover, we identified regions of the OSK domain with conserved predicted RNA binding potential. Furthermore, 35 we show that despite a low overall amino acid conservation, the LOTUS domain shows 36 higher conservation of predicted secondary structure than the OSK domain. Finally, we 37 suggest new key amino acids in the LOTUS domain that may be involved in the previously 38 reported Oskar-Vasa physical interaction that is required for its germ line role. 39

40

Keywords: *oskar*, *vasa*. Drosophila, germ plasm, germ cell, LOTUS domain, RNA
binding, Hidden Markov Models, Hymenoptera, Lepidoptera, Zygentoma

43

## 44 Introduction

With the evolution of obligate multicellularity, many organisms faced a challenge 45 considered a major evolutionary transition: allocating only some cells (germ line) to pass 46 47 on their genetic material to the next generation, relegating the remainder (soma) to death 48 upon death of the organism (reviewed in (Kirk 2005)). This is soma-germ line 49 differentiation, where only cells from the germ line will create the next generation (reviewed in (Kirk 2005)). While there are multiple mechanisms of germ cell specification, 50 51 they can be grouped into two broad categories, induction or inheritance (reviewed in 52 (Extavour and Akam 2003)). Under induction, cells respond to an external signal by adopting germ cell fate. Under the inheritance mechanism, maternally synthesized 53 54 cytoplasmic molecules, collectively called germ plasm, are deposited in the oocyte and "inherited" by a subset of cells during early embryonic divisions. Cells inheriting these 55 56 molecules commit to a germ line fate (reviewed in (Extavour and Akam 2003)).

57

58 The inheritance mechanism in insects that undergo metamorphosis (Holometabola) 59 appears to have evolved by co-option of a key gene, oskar. oskar was first identified in forward genetic screens for axial patterning mutants in Drosophila melanogaster 60 (Lehmann and Nüsslein-Volhard 1986). For the first 20 years following its discovery, 61 62 oskar appeared to be restricted to Drosophilids (Clark, et al. 2007). Its later discovery in the mosquitoes Aedes aegypti, Anopheles gambiae and Culex quinquefasciatus (Juhn 63 and James 2006; Juhn, et al. 2008) and the wasp Nasonia vitripennis (Lynch, et al. 2011) 64 suggested the hypothesis that oskar emerged at the base of the Holometabola, and 65 facilitated the evolution of germ plasm in these insects (Lynch, et al. 2011). However, our 66 subsequent identification of oskar orthologs in the cricket Gryllus bimaculatus (Ewen-67 68 Campen, et al. 2012), and in many additional hemimetabolous insect species (Blondel, et al. 2020), demonstrated that oskar predates the Holometabola, and must be at least 69 70 as old as the major radiation of insects (Misof, et al. 2014). Two secondary losses of oskar 71 from insect genomes have also been reported, in the beetle *Tribolium castaneum* (Lynch, 72 et al. 2011) and the honeybee Apis mellifera (Dearden, et al. 2006), and neither of these 73 insects appear to use germ plasm to establish their germ lines (Nelson 1915; Nagy, et al. 74 1994; Dearden 2006; Schroder 2006). Whether oskar is ubiquitous across all insect

orders, whether it is truly unique to insects, the evidence for or against potential losses or
duplications of the *oskar* locus across insects, and the evolutionary dynamics of the locus,
remain unknown.

78

oskar remains, to our knowledge, the only gene that has been experimentally 79 80 demonstrated to be both necessary and sufficient to induce the formation of functional 81 primordial germ cells (Kim-Ha, et al. 1991; Ephrussi and Lehmann 1992). Thus, in D. 82 melanogaster (Lehmann and Nüsslein-Volhard 1986; Kim-Ha, et al. 1991; Ephrussi and Lehmann 1992) and potentially more broadly in holometabolous insects with germ plasm 83 (Lynch, et al. 2011; Rafiqi, et al. 2020), oskar plays an essential germ line role. However, 84 85 it is clear that oskar's germ line function can evolve rapidly, as even within the genus 86 Drosophila, oskar orthologs from different species cannot always substitute for each other 87 (Webster, et al. 1994; Jones and Macdonald 2007). Moreover, the ancestral function of 88 this gene may have been in the nervous system rather than the germ line (Ewen-Campen, 89 et al. 2012). The current hypothesis is therefore that it was co-opted to play a key role in 90 the acquisition of an inheritance-based germ line specification mechanism approximately 91 300 million years ago (Misof, et al. 2014), in the lineage leading to the Holometabola 92 (Ewen-Campen, et al. 2012). Thus, the case of oskar offers an opportunity to study the 93 evolution of protein function at multiple levels of biological organization, from the genesis 94 of a novel protein, through to potential co-option events and the evolution of functional variation. 95

96

97 Neofunctionalization often correlates with a change in the fitness landscape of the protein 98 sequence caused by novel biochemical constraints imposed by amino acid sequence 99 changes (Sikosek, et al. 2012; Sikosek and Chan 2014). Such potential constraints may 100 be revealed by analyzing the conservation of amino acids, their chemical properties, or 101 structure at the secondary, tertiary or quaternary levels (Sikosek and Chan 2014). Oskar 102 has two well-structured domains conserved across identified orthologs to date (Blondel, 103 et al. 2020): an N-terminal Helix Turn Helix (HTH) domain termed LOTUS with potential 104 RNA binding properties (Anantharaman, et al. 2010; Jeske, et al. 2015; Yang, et al. 2015; 105 Jeske, et al. 2017), and a C-terminal GDSL-lipase-like domain called OSK (Jeske, et al.

106 2015; Yang, et al. 2015) (Figure 1). These two domains are linked by an unstructured 107 highly variable interdomain sequence (Ahuja and Extavour 2014; Jeske, et al. 2015; 108 Yang, et al. 2015). We previously showed that this domain structure is likely the result of 109 a horizontal transfer event of a bacterial GDSL-lipase-like domain, followed by the fusion 110 of this domain with a LOTUS domain in the host genome (Blondel, et al. 2020). 111 Biochemical assays of the properties of the LOTUS and OSK domains provide some 112 clues as to the molecular mechanisms that Oskar uses to assemble germ plasm in D. 113 *melanogaster*. The LOTUS domain is capable of homodimerization (Jeske, et al. 2015; Jeske, et al. 2017), and directly binds and enhances the helicase activity of the ATP-114 dependent DEAD box helicase Vasa, a germ plasm component (Jeske, et al. 2017). The 115 116 OSK domain resembles GDSL lipases in sequence (Jeske, et al. 2015; Yang, et al. 2015; 117 Blondel, et al. 2020), but is predicted to lack enzymatic activity, as the conserved amino 118 acid triad (S200 D202 H205) that defines the active site of these lipases is not conserved 119 in OSK (Anantharaman, et al. 2010; Jeske, et al. 2015; Yang, et al. 2015). Instead, co-120 purification experiments suggest that OSK has RNA binding properties, consistent with 121 its predicted basic surface residues (Jeske, et al. 2015; Yang, et al. 2015). Whether or 122 how changes in the primary sequence of Oskar can explain the evolution of its molecular 123 mechanism or tissue-specific function, remain unknown.

124

125 To date, sequences of approximately 100 oskar orthologs have been reported (Lynch, et 126 al. 2011; Jeske, et al. 2015; Quan and Lynch 2016; Blondel, et al. 2020). However, the 127 vast majority of these are from the Holometabola, and it is thus unclear whether analysis 128 of these sequences alone would have sufficient power to allow extrapolation of 129 conservation and divergence of putative biochemical properties across insects broadly 130 speaking. Multiple hypotheses as to the molecular mechanistic function of particular 131 amino acids in the LOTUS and OSK domains in *D. melanogaster* have been proposed 132 (Jeske, et al. 2015; Yang, et al. 2015; Jeske, et al. 2017), but without sufficient taxon 133 sampling, the potential relevance of these mechanisms to oskar's evolution and function 134 in other insects is unclear.

135

Here we address these outstanding questions by applying a rigorous bioinformatic 136 137 pipeline to generate the most complete collection of oskar sequences to date. By 138 analyzing 1862 Pancrustacean genomes and transcriptomes, we show that oskar likely 139 first arose at least 400 million years ago, before the advent of winged insects (Pterygota). 140 We find that the oskar locus has been lost independently in some insect orders, including 141 near-total absence from the order Hemiptera, and clarify that the absence of oskar from 142 the Bombyx mori and Tribolium castaneum genomes (discussed in Quan and Lynch 143 2016) does not reflect a general absence of oskar from Lepidoptera or Coleoptera. By comparing Oskar sequences in a phylogenetic context, we reveal that distinct biophysical 144 properties of Oskar are associated with Hemimetabola and Holometabola. We use these 145 146 observations to propose testable hypotheses regarding the putative biochemical basis of 147 evolutionary change in Oskar function across insects.

148

## 149 Results

150 HMM-based discovery pipeline yields hundreds of novel oskar orthologs

151 We wished to study the evolution of the *oskar* gene sequence as comprehensively as 152 possible across all insects. To expand our previous collection of nearly 100 orthologous 153 sequences (Blondel, et al. 2020), we designed a new bioinformatics pipeline to scan and 154 search for oskar orthologs across all 1565 NCBI insect transcriptomes and genomes that 155 were publicly available at the time of analysis (Supplementary Table S1; Figure 2; see Methods: Genome and transcriptome pre-processing for NCBI accession numbers and 156 157 additional information). First, we used the HMMER tool suite to build HMM models for 158 each of the LOTUS and OSK domains, using our previously generated multiple sequence 159 alignments (MSA) (Blondel, et al. 2020). We subjected genomes to *in silico* gene model 160 inference using Augustus (Stanke, et al. 2006). We translated the resulting predicted transcripts, as well as the predicted transcripts from RNA-seq datasets, in all six frames. 161 162 We then scanned the resulting protein sequences for the presence of LOTUS and OSK 163 domains using the aforementioned HMM models. Sequences were designated as oskar 164 orthologs based on the same criteria as in our previous study (Blondel, et al. 2020), 165 namely, sequences containing both a LOTUS and an OSK domain (Jeske, et al. 2015), separated by a variable interdomain region. We then aligned all sequences using 166

*hmmalign* and the HMM derived from our previously published full length Oskar alignment
 (Blondel, et al. 2020), and manually curated sequence duplicates and sequences that did
 not align correctly.

170

171 With these methods, we recovered a total of 379 unique oskar sequences from 350 172 unique species. To our knowledge, this comprises the largest collection of oskar orthologs 173 described to date. To determine if oskar orthologs might predate Insecta, we applied the 174 discovery pipeline to all 31 genomes and 266 transcriptomes of non-insect pancrustaceans available at the time of analysis (see Methods: Genomes and 175 176 transcriptomes preprocessing for complete list). However, we did not recover any non-177 insect sequences meeting our criteria for oskar orthologs (Figure 3), strongly suggesting 178 that oskar is restricted to the insect lineage (Lynch, et al. 2011; Ahuja and Extavour 2014).

179

180 We found that 58.65% of RefSeg genomes (78/133), 30.42% of GenBank genomes 181 (94/309), and 21.19% of transcriptomes (238/1123) analyzed contained predicted oskar 182 orthologs (Supplementary Table S1 and Supplementary Figure S1a). Given that detection 183 of putative orthologs is highly dependent on the quality of the genome assembly and 184 annotation, we asked whether there were differences in the assembly statistics of 185 genomes with and without predicted oskar orthologs. We observed a significant difference 186 in N50, L50, number of contigs and number of scaffolds between genomes lacking oskar 187 hits and those where oskar was identified (Mann-Whitney U test p-value < 0.05). 188 Genomes where we did not find oskar showed a significantly higher mean/median contig 189 and scaffold count, smaller contig and scaffold N50 length, larger contig and scaffold L50, 190 and more contigs or scaffolds per genome length, than genomes where we detected an 191 oskar ortholog (Mann Whitney U test p<0.05; Supplementary Figure S2; Supplementary 192 Table S2).

193

194 oskar predates the divergence of Ametabola and other insects

We identified *oskar* orthologs in 15 of the 29 generally recognized (Misof, et al. 2014) insect orders, including eight holometabolous orders, six hemimetabolous orders, and one ametabolous order (Figure 3). This result is consistent with our previous proposals

that *oskar* predates the origins of the Holometabola (Ewen-Campen, et al. 2012; Blondel,
et al. 2020). The novel finding of an *oskar* ortholog from the silverfish *Atelura formicaria*(Zygentoma) allows us to date back the origin of *oskar* further than previous analyses, to
at least 420 million years ago (Misof, et al. 2014), before the divergence of Ametabola
from the remaining insect lineages.

203

204 We then explored the distribution of oskar sequences across insect phylogeny. 205 Interestingly, we identified multiple lineages where oskar appeared to have been lost 206 independently, including confirming the previously reported (Lynch, et al. 2011) losses 207 from the genomes of the red flour beetle Tribolium castaneum, the honeybee Apis 208 mellifera, and the silk moth Bombyx mori (Figure 3). Notably, within Lepidoptera we 209 identified oskar orthologs in only four species, despite the fact that we searched 232 210 available lepidopteran sequence datasets, including 17 well-annotated RefSeq genomes 211 and 135 transcriptomes (Figure 3 and Supplementary Figure S3). In principle, this 212 apparent widespread absence of oskar in Lepidoptera could be due to unusually rapid 213 evolution of the oskar sequence in this lineage, which might render lepidopteran oskar 214 orthologs undetectable by our methods. However, we note that the only four lepidopteran 215 orthologs we detected all belonged to species of the basally branching Adelidae and 216 Palaephatidae families. We therefore favor the interpretation that oskar was lost from a 217 last common ancestor of Meessiidae and Palaphaetidae, approximately 180 million years 218 ago, with the consequence that the majority of extant lepidopteran lineages lack an oskar 219 ortholog (Supplementary Figure S3) (Mitter, et al. 2017; Kawahara, et al. 2019).

220

221 The Hemiptera also appear to have lost oskar, based on our analysis of the 222 datasets 222 available for this clade, including 12 RefSeq genomes and 192 transcriptomes. However, 223 we did identify an oskar ortholog in the Thysanoptera, which is a hemipteran sister group 224 (Misof, et al. 2014). Finally, we identified oskar orthologs in only four of the 11 orders of 225 the Polyneoptera for which data were available. With the exception of Mantodea (13) 226 transcriptomes), the four orders with detectable oskar sequences all had more than ten 227 available sequence datasets (Plecoptera: three genomes and eight transcriptomes; 228 Orthoptera: three genomes and 28 transcriptomes; Phasmatodea: 13 genomes and 31

transcriptomes; Blattodea: five genomes and 51 transcriptomes). The remaining orders
had fewer than eight datasets each available for analysis (Figure 3; Supplementary Table
S1), which could account for the apparent paucity of *oskar* genes in this group. However,
we cannot rule out the possibility that *oskar* in the Polyneoptera may have diverged
beyond our ability to detect it, or that it may have been lost multiple times, as observed
for multiple holometabolous orders.

235

236 As well as multiple convergent losses of oskar, we also uncovered evidence for 237 independent instances of duplication of the oskar locus. We defined a putative duplication 238 instance as two or more oskar sequences (possessing both a LOTUS and OSK domain 239 as per our definition) in the same species that shared less than 80% sequence similarity. 240 All of these events were detected within the Hymenoptera. We therefore performed a 241 phylogenetic analysis of the hymenopteran sequences to test the hypothesis that these 242 were the result of duplication events (Figure 4; Supplementary Figure S4). Our analysis of hymenopteran oskar sequences recovered previously published hymenopteran 243 244 phylogenetic relationships (Peters, et al. 2017). We found that oskar was duplicated in 245 the four Figitidae species studied, a family of parasitoid wasps. Moreover, one out of ten 246 examined Cynipidae species, as well as the only Ceraphronidae species examined, also 247 harbored a duplicated oskar sequence. Multiple oskar duplications were also identified in 248 the Chalcidoid wasps, notably in the Mymaridae (all three species studied), the 249 Eupelmidae (two out of three species), the Aphelinidae (both species) and the 250 Pteromalidae (one out of 17 species). Finally, we identified two additional apparently 251 independent duplication events in the Aculeata, one in the wasp *Polistes fuscatus* (of 29) 252 Vespidae, including three additional Polistes species, two with RefSeq genomes (P. 253 canadensis and P. dominula) in which oskar was identified in single copy), and one in the 254 red imported fire ant Solenopsis invicta (of 41 Formicidae species, including the 255 congeneric S. fugax, with a GenBank genome in which oskar was identified in single 256 copy).

257

258 Evidence for oskar expression in multiple somatic tissues

259 In studied insects to date, oskar is expressed and required in one or both of the germ line 260 (Juhn and James 2006; Juhn, et al. 2008; Lynch, et al. 2011; Lehmann 2016) or the 261 nervous system (Ewen-Campen, et al. 2012; Xu, et al. 2013). We asked whether these 262 expression patterns could be detected in the insects studied here. To this end, we 263 downloaded all available metadata for the transcriptomes analyzed here, to obtain 264 information on the source tissues and developmental stages. We obtained these data for 265 371 out of the 1123 transcriptomes in our analysis, including both holometabolous and 266 hemimetabolous orders (see Methods: TSA metadata parsing and curation). To first 267 explore the distribution of oskar expression in the brain and the germ line, we binned the 268 different tissues reported in the metadata into two categories, brain or germ line. This was 269 done independently of the developmental stage (if that information was included in the 270 metadata) by creating a mapping table and checking the extracted tissues against this 271 table (Supplementary Table S3 at GitHub repository 272 **TableS3 germline brain table.csv**). We then cross referenced our orthology detection 273 with these metadata. We found evidence for oskar expression in the germ line of four 274 orders (Phasmatodea, Hymenoptera, Coleoptera and Diptera), and in the brain of five 275 orders (Orthoptera, Blattodea, Hymenoptera, Coleoptera, Diptera) (see Methods: TSA 276 metadata parsing and curation for details on keyword extractions). In addition, we found 277 evidence of oskar expression in several somatic tissues not previously implicated in 278 studies of oskar expression and function. These tissues included the midgut (Polistes 279 fuscatus, Sitophilus oryzae), fat body (Polistes fuscatus, Arachnocampa luminosa), 280 salivary gland (Culex tarsalis, Anopheles aguasalis, Leptinotarsa decemlineata), venom 281 gland (Culicoides sonorensis, Fopius arisanus), and silk gland (Bactrocera cucurbitae) 282 (Supplementary Figure S5). In terms of developmental stage, only holometabolous 283 insects appeared to express oskar during embryonic, larval or nymphal stages; for all other insects, oskar was detected in transcriptomes derived from adults (Figure 3). 284 285 However, it is important to note that for most species, transcriptomes were available only 286 from adult tissues, rather than from a full range of developmental stages (Supplementary 287 Figure S5). We therefore cannot rule out the possibility that oskar expression at pre-adult stages is also a feature of multiple Hemimetabola. Indeed, we previously reported that 288

oskar is expressed and required in the embryonic nervous system of a cricket, a
 hemimetabolous insect (Ewen-Campen, et al. 2012).

291

## 292 The Long Oskar domain is an evolutionary novelty specific to a subset of Diptera

293 D. melanogaster has two isoforms of Oskar (Markussen, et al. 1995): Short Oskar, 294 containing the LOTUS, OSK and interdomain regions, and Long Oskar, containing all 295 three domains of Short Oskar as well as an additional 5' domain (Supplementary Figure 296 S7). It was previously reported that Long Oskar was absent from *N. vitripennis*, *C. pipiens* 297 and G. bimaculatus (Lynch, et al. 2011; Ewen-Campen, et al. 2012), and within our 298 alignment of Oskar sequences we could only detect the Long Oskar isoform within 299 Diptera. Therefore, using our dataset, we asked when these two isoforms had evolved. 300 We selected the dipteran sequences from our Oskar alignment and then grouped the 301 sequences by family. We plotted the amino acid occupancy at each alignment position 302 (Supplementary Figure S7), and found that Long Oskar predates the Drosophilids, being 303 identified as early as the *Pinpunculidae* (Supplementary Figure S7). Moreover, following 304 the evolution of the Long Oskar isoform, the Long Oskar domain was retained in all 305 families except for the Glossinidae and Scathophagidae. However, given that we 306 identified only eight and two Oskar sequences for these families respectively, we cannot 307 eliminate the possibility that apparent absence of the Long Oskar domain in these groups 308 reflects our small sample size, rather than true evolutionary loss.

309

# 310 The LOTUS and OSK domains evolved differently between hemimetabolous and 311 holometabolous insects

The fact that an *oskar*-dependent germ plasm mode of germ line specification mechanism has been identified only in holometabolous insects suggests that *oskar* may have been co-opted in this clade for this function (Ewen-Campen, et al. 2012). Under this hypothesis, evolution of the *oskar* sequence in the lineage leading to the Holometabola may have changed the physico-chemical properties of Oskar protein, such that it acquired germ plasm nucleation abilities in these insects. To test this hypothesis, we asked whether there were particular sequence features associated with Oskar proteins from 319 holometabolous insects, in which Oskar can assemble germ plasm, and hemimetabolous 320 insects, which lack germ plasm. In particular, we assessed the differential conservation 321 of amino acids at particular positions across Oskar and asked if these might be predicted 322 to change the physico-chemical properties of Oskar in specific ways that could potentially 323 be relevant to germ plasm nucleation. We used the Valdar score (Valdar 2002) as the 324 main conservation indicator for this study (see GitHub file scores.csv), as this metric 325 accounts not only for transition probabilities, stereochemical properties and amino acid 326 frequency gaps, but also for the availability of sequence diversity in the dataset. It 327 computes a weighted score, where sequences from less well-represented clades 328 contribute proportionally more to the score than sequences from overrepresented clades. 329 Due to the highly unbalanced availability of genomic and transcriptomic data between 330 hemimetabolous and holometabolous sequences (Supplementary Table S1; Figure 3) the choice of a weighted score was necessary to avoid biasing the results towards insect 331 332 orders such as Diptera or Hymenoptera. To study the difference between 333 hemimetabolous and holometabolous sequences, we did not use the Valdar score 334 directly, but instead computed the conservation ratio between both groups for each position, which we call the Conservation bias (See Methods: Computation of the 335 336 Conservation Bias). We plotted the conservation bias on the solved three-dimensional 337 crystal structure of the *D. melanogaster* LOTUS and OSK domains (Jeske, et al. 2015; 338 Yang, et al. 2015) to ask whether specific functionally relevant structures showed 339 phylogenetic or other patterns of residue conservation (Figure 5).

340

341 First, we asked if the overall conservation score of the domains was different between 342 holometabolous and hemimetabolous sequences. We observed that the conservation 343 bias for the LOTUS domain was centered around a mean of 1.00, indicating that both 344 Holometabola and Hemimetabola displayed a similar conservation of the LOTUS domain 345 (Figure 5a). For the OSK domain however, the conservation bias was centered around 0.84, indicating that the hemimetabolous sequences displayed a higher level of 346 347 conservation compared to holometabolous sequences (Figure 5a). We then looked at the 348 conservation bias scores in-situ on the LOTUS domain structure. We asked if the amino 349 acids of the  $\beta$  sheets of the LOTUS domain thought to be involved in dimerization of the

350 protein (Jeske, et al. 2015; Yang, et al. 2015) displayed conservation bias. Both  $\beta$  sheets 351 had an overall even bias (mean: 1.03 and 1.05 for  $\beta$ 1 and  $\beta$ 2 respectively) between both 352 groups (Figure 5b). Second, as we had observed that hemimetabolous OSK was more 353 conserved overall than holometabolous OSK, we asked if there were any clear patterns 354 of conservation bias in specific regions of the structure (Figure 5a and b). We found that some of the secondary structures within OSK showed a differential conservation ( $\alpha$ 2: 355 356 0.54,  $\alpha$ 6: 0.42,  $\beta$ 2: 0.52), whereas other structures were within less than 0.1 of the median 357 value for OSK. Moreover, we observed a large pocket of amino acids showing a 358 conservation bias towards hemimetabolous sequences located on the surface of OSK 359 (Figure 5c). This particular area contains the previously reported important amino acids 360 for the RNA binding function of OSK (Jeske, et al. 2015; Yang, et al. 2015) namely, R442, 361 R436 and R576. The electrostatic properties at those positions were conserved in the 362 holometabolous sequences R436: 0.36, R442: 0.29 and R576: 0.81 (Figure 5d), but not 363 in hemimetabolous sequences.

364

365 To gain further insight into the differences in conservation across insects, we reduced the multiple sequence alignment dimensionality using a Multiple Correspondence Analysis 366 367 (MCA), an equivalent of PCA for categorical variables (Lebart, et al. 1984). We performed 368 the dimensionality reduction for the full-length Oskar sequence alignment as well as for 369 the LOTUS and OSK alignments (Supplementary Figure S7). Interestingly, we found that 370 most of the variance in sequence space was due to dipterans and hymenopterans 371 (Supplementary Figure S7). When we considered the OSK domain only, we identified 372 clusters of Drosophilidae, Culicidae and Formicidae sequences (Supplementary Figure S7). This clustering is also reflected for the LOTUS domain, where the Drosophilidae and 373 374 *Culicidae* contribute to a high amount of variance in the first MCA dimension. However, 375 for the LOTUS domain, the *Formicidae* sequences do not cluster away from other Oskar 376 sequences (Supplementary Figure S7). This suggests that the LOTUS domain of Diptera 377 diverged in sequence between Drosophilidae and Culicidae.

378

Evidence for evolution of stronger dimerization potential of the Oskar LOTUS domain inHolometabola

381 The LOTUS domain dimerizes *in vitro* through electrostatic and hydrophobic contacts of 382 Arg215 of the  $\beta$ 2 sheet and Thr195, Asp197 and Leu200 of the  $\alpha$ 2 helix (Jeske, et al. 383 2015; Yang, et al. 2015). To date, however, the biological significance of Oskar dimerization remains unknown. Moreover, the dimerization of the LOTUS domain does 384 385 not appear to be conserved across all Oskar sequences (Jeske, et al. 2015). Specifically, ten LOTUS domains from non-drosophilid species were tested for dimerization, and only 386 387 LOTUS domains from Drosophilidae, Tephritidae and Pteromalidae formed homodimers (Jeske, et al. 2015). The other sequences tested, from Culicidae, Formicidae and 388 389 Gryllidae, remained monomeric under the tested conditions (Jeske, et al. 2015). We 390 selected the LOTUS sequences in our alignment from those six families and placed them 391 into one of two groups, dimeric and monomeric LOTUS, under the assumption that any 392 sequence from that family would conserve the dimerization (or absence thereof) 393 properties previously reported (Jeske, et al. 2015). We asked whether we could detect 394 any evolutionary changes between the two groups in properties of known important 395 dimerization interfaces and residues in our sequence alignment (Jeske, et al. 2015).

396

397 In the *D. melanogaster* structure, two key amino acids, D197 and R215, are predicted to 398 form hydrogen bonds that stabilize the dimer (Jeske, et al. 2015). We found that in the 399 dimer group, the electrostatic properties of these two amino acids are highly conserved 400 (-0.75 for D197 and 0.81 for R215), while in the monomer group the electrostatic 401 interaction is not conserved (0.03 for D197 and -0.11 for R215) (Figure 6e). Given the 402 differential conservation between the two groups, our results support the previous finding 403 that disrupting this interaction prevents dimerization (Jeske, et al. 2015). L200 was 404 previously hypothesized to stabilize the interface via hydrophobic forces (Jeske, et al. 405 2015). We observed that the hydrophobicity of this residue is highly conserved in the 406 dimer group (L200: 0.89), but that in the monomer group this residue is hydrophilic (L200: 2.33) (Figure 6f). In sum, our analyses show that key amino acids in the LOTUS domain 407

408 evolved differently in distinct insect lineages, in a way that may explain why some insect409 LOTUS domains dimerize and some do not.

410

## 411 Conservation of the Oskar-Vasa interaction interface

412 Next, we asked whether we could detect differential conservation of the LOTUS-Vasa 413 interface. It was previously reported that the LOTUS domain of Oskar acts as an interaction domain with Vasa (Jeske, et al. 2017), a key protein with a conserved role in 414 415 the establishment of the animal germ line (Hay, et al. 1990; Lasko 2013). The interaction between Oskar's LOTUS domain and Vasa is through an interaction surface situated in 416 417 the pocket formed by the helices  $\alpha 2$  and  $\alpha 5$  of the LOTUS domain (Figure 6a b and c). 418 Due to the essential role that vasa plays in germ line determination (reviewed in Raz 419 2000; Noce, et al. 2001; Extavour and Akam 2003; Ewen-Campen, et al. 2010; Lasko 420 2013), and the potential co-option of *oskar* to the germ line determination mechanism in 421 Holometabola (Ewen-Campen, et al. 2012), we hypothesized that evolutionary changes 422 in the conservation of the residues of this interface might be detectable between Holometabola and Hemimetabola. First, we observed that the residues of the LOTUS 423 424 domain  $\alpha^2$  and  $\alpha^5$  helices, which directly contact Vasa (Jeske, et al. 2017) were highly 425 conserved overall ( $\alpha$ 2 average Valdar score 0.49;  $\alpha$ 5 Valdar score 0.56) (Figure 6b). 426 Specifically, we observed that the previously reported Vasa interacting amino acids A162 427 and L228 of the LOTUS domain were highly conserved (Valdar score: 0.64 for both 428 residues) (Jeske, et al. 2017). We also noted that Q235 and H227 of the LOTUS domain 429  $\alpha$ 5 helix are likely to be important interaction partners due to their high conservation 430 (Valdar score: 0.90 and 0.90 for both residues) (Figure 6b). Moreover, facing the LOTUS 431 domain H227 is Vasa M540, which may act as a proton donor to form a hydrogen bond 432 between the histidine ring and the sulfur atom of the methionine (Pal and Chakrabarti 2001) (Figure 6b and b'). The LOTUS domain  $\alpha$ 2 helix is overall slightly less conserved 433 434 than the LOTUS domain  $\alpha$ 5 helix (Valdar score: 0.49 vs 0.56) (Figure 6a, b", c"), but 435 hydrophobic properties are conserved on one side of the  $\alpha$ 2 helix (Figure 6c, c') forming 436 a motif of conserved amino acid properties (Figure 6c").

437

438 Previous reports have hypothesized that the *D. melanogaster* LOTUS domain could act 439 as a dsRNA binding domain (Anantharaman, et al. 2010; Callebaut and Mornon 2010). 440 However, in *D. melanogaster*, it was later reported that the LOTUS domain did not bind 441 to nucleotides (Jeske, et al. 2015). Therefore, using our dataset we assessed the potential 442 RNA binding properties of LOTUS domains to test the conservation of this prediction. We 443 used the RNABindR algorithm (Terribilini, et al. 2007) to predict potential RNA binding sites of the LOTUS domain, and computed a conservation score for each position 444 445 (Terribilini, et al. 2007). We found that the  $\alpha$ 5 helix is the location in the LOTUS domain that has the most conserved prediction for RNA binding (Figure 6d). 446

447

Finally, we asked whether the secondary structure of the LOTUS domain might be 448 449 conserved. Secondary structures are often indicative of the tertiary structure of a domain. 450 Therefore, we reasoned that the secondary structure might be conserved even if the 451 sequence varies. We submitted the LOTUS sequences from all identified Oskar orthologs 452 to the Jpred4 servers (Drozdetskiy, et al. 2015) for secondary structure prediction and 453 mapped the results onto the Oskar alignment we obtained. We found that the secondary 454 structure of LOTUS is highly conserved throughout Oskar orthologs, with the exception 455 of the  $\alpha 1$  helix (Supplementary Figure S8) which displays a low conservation score of 0.19 (Figure 6a). 456

457

#### 458 The core of the OSK domain is conserved

459 We asked whether the OSK domain showed any differential conservation across the 460 different parts of the domain. We found that the OSK domain of Oskar showed an overall 461 conservation across all insects, similar to the LOTUS domain (Valdar score: 0.51) (Figure 462 7a). However, the conservation pattern is higher in the core amino acids (Valdar score average of core amino acid: 0.54) when compared to the residues at the surface (Valdar 463 464 score average for surface amino acid: 0.23) (Figure 7a). Despite the overall low 465 conservation of the residues at the surface of the OSK domain, we found that the 466 electrostatic properties are conserved overall (electrostatic conservation score >0; 467 conserved) in the previously reported putative RNA binding pocket (Yang, et al. 2015). 468 However, as previously mentioned, this conservation is stronger in holometabolous

sequences (Figure 5d). These results are in accordance with the potential role of OSK as
an RNA Binding domain in the context of germ plasm assembly (Jeske, et al. 2015; Yang,
et al. 2015). We also submitted the OSK sequences to the same secondary structure
analysis performed on LOTUS. We found that, as for the LOTUS domain, the secondary
structure of OSK is highly conserved throughout all insect sequences analyzed
(Supplementary Figure S8).

475

476 We then asked if the conservation patterns observed at the core of OSK were clustered 477 in sequence motifs. When we looked at the location of the highly conserved amino acids, 478 we found that the conservation was driven by four well-defined sequence motifs (Figure 7c, c', c'', c'''). Given that oskar plays different roles in Holometabola and Hemimetabola, 479 480 we asked whether the conserved OSK motifs showed any difference in conservation 481 between these two groups. Of the four highly conserved OSK core motifs (Figure 7c, c', c", c"), two of them (Figure 7c: Valdar average score: 0.80 and c" Valdar average score: 482 483 0.71) were conserved across all insects, but the other two showed differential 484 conservation between the holometabolous and hemimetabolous sequences (Figure 7c': Valdar score average Holometabola: 0.78; Hemimetabola: 0.58 c":' Valdar score average 485 486 Holometabola: 0.70, Hemimetabola: 0.55). Finally, we noted that only one of the affected 487 OSK domain residues in known loss of function oskar alleles affecting posterior patterning 488 in D. melanogaster, S457, is conserved across all insects (Valdar score: 0.86). This 489 suggests that the role of the other previously reported important amino acids in the 490 function of *D. melanogaster* OSK (Yang, et al. 2015) might not be conserved in other 491 insects (red positions in Figure 7c, c', c''').

492

## 493 Discussion

494 An expanded collection of oskar orthologs

*oskar* provides a powerful case study of functional evolution of a gene with an unusual
genesis (Blondel, et al. 2020). Here, we gathered the most extensive set of orthologous *oskar* sequences to date. However, most insect genomic and transcriptomic data have
been generated from only a few orders, and the vast majority from the Holometabola.
Diptera, Lepidoptera, Coleoptera, Hymenoptera and Hemiptera represent 82% of the

datasets available at the time of this analysis. We emphasize that expanded taxon sampling, particularly for the Hemimetabola, will be critical for further studies of the evolution of protein function across insects. Moreover, only a small proportion (27% for tissue type, 26% for organism stage, and 14% for sex) of the TSA datasets contained usable metadata regarding the stage and tissue type sampled. Future standardization of the nature and format of transcriptomic metadata would also be a worthwhile endeavor that could increase the efficiency and efficacy of future work.

507

### 508 Convergent losses and duplications of oskar in insect evolution

509 A previous report suggested that *oskar* had been lost from the genome of the silk moth 510 B. mori (Lynch, et al. 2011). Our analysis of 232 datasets across 44 of the 126 described 511 lepidopteran families (Kawahara, et al. 2019) strongly suggests that the loss of oskar in 512 the Lepidoptera (butterflies and moths) is not unique to the silk moth, but rather occurred 513 early and repeatedly in lepidopteran evolution. The fact that oskar is a component of the 514 oosome at the posterior of the oocyte (the wasp germ plasm analog (Quan, et al. 2019)) 515 and required for germ cell formation in the wasp Nasonia vitripennis (Lynch, et al. 2011) 516 implies that a common ancestor of Holometabola had already established an oskar-517 dependent inheritance mode of germ line specification. Therefore, the apparent 518 subsequent loss in nearly all Lepidoptera examined of a gene responsible for the 519 establishment of the germ plasm in other Holometabola might seem unexpected. Few 520 studies have directly addressed the molecular mechanisms of germ cell specification in 521 Lepidoptera. In B. mori (Bombicidae), vasa mRNA (Nakao 1999) and protein (Nakao, et 522 al. 2006), and the transcripts of one of four nanos orthologs (nanos-O) (Nakao, et al. 523 2008), have been detected in a region of ventral cortical cytoplasm in pre-blastoderm 524 stage embryos. As putative primordial germ cells form in this location at later stages (Miya 525 1958), some authors have speculated that a germ plasm, located ventrally rather than 526 posteriorly, may specify germ cells in this moth (Toshiki, et al. 2000; Nakao, et al. 2008). 527 However, recent knockdown experiments showed that maternal nanos-O is dispensable 528 for germ cell formation (Nakao and Takasu 2019), consistent with a zygotic, inductive 529 mechanism. In the butterfly Pararge argeria (Nymphalidae), no oskar ortholog has been 530 identified in the genome (Carter, et al. 2013), but the transcripts of one of four identified

531 nanos orthologs (nanos-O) have been detected in a small region of ventral cortical 532 ooplasm, again prompting speculation that this lepidopteran may also deploy a germ 533 plasm (Carter, et al. 2015). We suggest that if these or other Lepidoptera do indeed rely 534 on germ plasm to specify their germ line, they may do so using a germ plasm nucleator other than Oskar. For most studied Lepidoptera, however, classical embryological studies 535 536 report the first appearance of primordial germ cells at post-blastoderm stages, either from 537 the ventral midline of the cellular blastoderm or early germ band (Woodworth 1889; 538 Tomaya 1902; Sehl 1931; Miya 1953, 1958, 1975; Tanaka 1987), from the coelomic sac 539 mesoderm of the abdomen (Johannsen 1929; Eastham 1930; Saito 1937; Presser and 540 Rutschky 1957; Kobayashi and Ando 1984), or from the primary ectoderm of the caudal 541 germ band (Schwangart 1905; Lautenschlager 1932; Ando and Tanaka 1979; Tanaka 542 1987; Guelin 1994) (Figure S6). Taken together, these data suggest that an inductive 543 mechanism may operate to specify germ cells in most moths and butterflies. We 544 speculate that the loss of oskar from most lepidopteran genomes may have facilitated or 545 necessitated secondary reversion to the hypothesized ancestral inductive mechanism for 546 germ line specification.

547

548 Another order with apparent near-total absence of oskar orthologs is the Hemiptera (true 549 bugs), whose sister group Thysanoptera (thrips) nevertheless possesses oskar. This 550 secondary loss of oskar from a last common hemipteran ancestor correlates with the 551 reported post-blastoderm appearance of primordial germ cells in the embryo. Classical 552 studies on most hemipteran species describe germ cell formation as occurring after 553 cellular blastoderm formation, on the inner (yolk-facing) side of the posterior blastoderm 554 surface (Metschnikoff 1866; Witlaczil 1884; Will 1888; Mellanby 1935; Butt 1949; Kelly 555 and Huebner 1989; Heming and Huebner 1994). A notable exception to this is the parthenogenetic pea aphid Acyrthosiphon pisum, for which strong gene expression and 556 557 morphological evidence supports a germ plasm-driven germ cell specification mechanism 558 in both sexual and asexual modes (Miura, et al. 2003; Chang, et al. 2006; Lin, et al. 2014). 559 In contrast, studies of the aphids Aphis plantoides, A. rosea and A. pelargonii describe 560 no germ plasm, and post-blastoderm germ cell formation (Metschnikoff 1866; Witlaczil 561 1884; Will 1888). However, the genomes of all aphids studied here, including A. pisum and three *Aphis* species, appear to lack *oskar*. This suggests that germ plasm assembly
in *A. pisum* either does not require a nucleator molecule or uses a novel non-Oskar
nucleator.

565 In the Hymenoptera (ants, bees, wasps and sawflies), our results strongly suggest that 566 oskar was lost from the genome of the last common ancestor of bees and spheroid wasps (Supplementary Figure S9). Our analysis further suggests multiple additional 567 568 independent losses in as many as 25 other hymenopteran lineages, including some for which good quality RefSeq genomes were available (e.g. the slender twig ant 569 570 Pseudomyrmex gracilis or the wheat stem sawfly Cephus cinctus (Supplementary Figure 571 S9). However, it would be premature to draw strong conclusions about the number of 572 independent losses given the predominance of transcriptome data in the Hymenoptera.

573 In addition to convergent losses of oskar, we also found evidence for clade-specific 574 duplications of oskar in the Hymenoptera. Seven of the nine families containing these 575 putative duplications are families of parasitoid wasps; the remaining two families are ants 576 (Formicidae) and the group of yellowjackets, hornets, and paper wasps (Vespidae) 577 (Figure 4). The phylogenetic relationships of these groups make it highly unlikely that a 578 duplication occurred only once in their last common ancestor, which would be the last 579 common ancestor of all wasps, bees and ants (i.e. Apocrita, all hymenopterans except 580 sawflies) (Supplementary Figure S9). We suggest that the most parsimonious hypothesis 581 is one of three to five independent duplications of oskar, followed by at least nine to 14 582 independent reversions to a single copy, or total loss of the locus (Supplementary Figure 583 9).

584 No notable life history characteristics appear to unite those species with multiple oskar 585 orthologs: they include eusocial and solitary, sting-bearing and stingless, parasitoid and 586 non-parasitic insects. To our knowledge, neither is there anything unique about the germ 587 line specification process in Hymenoptera with one or more than one oskar ortholog. Most 588 Hymenoptera appear to use a germ plasm-driven mechanism to specify germ cells in 589 early blastoderm stage embryos (Supplementary Figure S9 and references therein), and 590 we identified oskar orthologs for all such species described in the embryological literature 591 (Supplementary Figure S9). In the notable example of the honeybee Apis mellifera, in 592 which cytological and molecular evidence suggests germ cell arise from abdominal 593 mesoderm (Bütschli 1870; Nelson 1915; Fleig and Sander 1985, 1986; Zissler 1992; 594 Gutzeit, et al. 1993; Dearden 2006), we identified no oskar ortholog in its well-annotated 595 genome (Supplementary Figure S9), as noted previously by other authors (Lynch, et al. 2011). However, no major differences in germ plasm or pole cell formation have been 596 597 reported in species or families of ants or wasps with duplicated oskar loci, compared with 598 close relatives that possess oskar in single copy (e.g. compare the ants Solenopsis invicta 599 (at least 2 oskars) and Aphaenogaster rudis (1 oskar) (Khila and Abouheif 2008), or the 600 pteromalid wasps Nasonia vitripennis (1 oskar) (Lynch and Desplan 2010; Lynch, et al. 601 2011: Quan, et al. 2019) and Otitesella tsamvi (2 oskars). Thus, future studies that 602 independently abrogate the functions of each paralog individually, will be needed to 603 determine the biological significance, if any, of these oskar duplications.

604

## 605 Functional implications of differential conservation of the LOTUS and OSK domains

606 We have identified novel conserved amino acid positions that we hypothesize are 607 important for the Vasa binding properties of the LOTUS domain and the RNA properties 608 binding of the OSK domain (Figure 6 and 7). Our observation of the conservation of the 609 LOTUS domain  $\alpha^2$  helix is consistent with its previously reported importance LOTUS-610 Vasa binding (Jeske, Müller, and Ephrussi 2017). In the  $\alpha^2$  helix, we also observed high 611 conservation of H227 and Q235. The positions of these residues suggest they may 612 contribute to the interaction between Vasa and LOTUS. We suggest they should therefore 613 be the target of future mutational studies.

614

We also uncovered an interesting new conservation pattern within the OSK domain. The conserved amino acids were more abundant in the core of the domain than on the surface. This differential conservation might be relevant to the acquisition of a germ plasm nucleator role of *oskar* in the Holometabla (Figure 5). We noted that the basic properties of surface residues previously reported for *D. melanogaster* (Yang, et al. 2015) are conserved across insects, which might indicate that the RNA binding properties of OSK observed in *D. melanogaster* (Jeske, et al. 2015; Yang, et al. 2015) are also conserved

throughout holometabolous insects. We speculate that the comparatively low amino acid conservation of the surface residues in Holometabolous OSK domains, which nevertheless display highly conserved basic properties, could have allowed greater flexibility in the co-evolution of specific RNA binding partners for the OSK domains of different lineages.

- 627
- 628 OSK evolved differentially between holometabolous and hemimetabolous insects

629 Finally, we observed a differential conservation of the OSK domain between 630 hemimetabolous and holometabolous insects. Specifically, we found that the OSK 631 sequence was less conserved across the Holometabola than across the Hemimetabola. 632 This observation raises two potential hypotheses regarding the role of the OSK domain 633 in the functional evolution of Oskar. First, perhaps the apparently relaxed purifying 634 selection experienced by OSK in the Holometabola was necessary for the co-option of 635 oskar to a germ plasm nucleation role. Second, Oskar might have a function in the 636 hemimetabolous insects that requires strong conservation of OSK. More studies on the 637 roles and biochemical properties of OSK in hemimetabolous insects will be required to 638 test these hypotheses and further our understanding of the biological relevance of this 639 differential conservation.

640

In conclusion, analysis of the large dataset of novel Oskar sequences presented here provides multiple new testable hypotheses concerning the molecular mechanisms and functional evolution of *oskar*, that will inform future studies on the contribution of this unusual gene to the evolution of animal germ cell specification.

645

#### 646 Materials and Methods

647 Lead contact and materials availability

This study did not generate new unique reagents. This study generated new python3 code and supplementary files referred to below, all of which are available <u>https://github.com/extavourlab/Oskar\_Evolution.</u> Requests for further information and requests for resources and reagents should be directed to and will be fulfilled by Cassandra G. Extavour (<u>extavour@oeb.harvard.edu</u>).

#### 653

### 654 Experimental model and subject details

655 This study used no animal model, nor any cell culture lines. However, it used previously 656 generated genomic and transcriptomic datasets. All the information regarding how those 657 datasets were generated can be found on their respective NCBI pages. The list of all the 658 datasets used study be found in in this can the following files: 659 genome insect database.csv, transcriptome insect database.csv, 660 genome\_crustacean\_database.csv, and transcriptome\_crustacean\_database.csv.

- 661
- 662 Genome and transcriptome preprocessing

We collected all available genome and transcriptome datasets from the NCBI repository 663 664 registered in September 2019 (Figure 2). NCBI maintains two tiers of genomic data: 665 RefSeq, which contains curated and annotated genomes, and GenBank, which contains 666 non-annotated assembled genomic sequences. Transcriptomes are stored in the 667 Transcriptome Shotgun Assembly (TSA) database, with metadata including details on 668 their origin. Among the registered datasets, five genomes were not yet available, and 40 669 transcriptomes were only available in the NCBI Trace repository. As they did not comply 670 with the TSA database standards, they were excluded from the analysis. To search for 671 oskar orthologs in datasets retrieved from GenBank, we needed to generate in silico gene 672 model predictions. We used the genome annotation tool Augustus (Stanke et al. 2006), which requires a Hidden Markov Model (HMM) gene model. To use HMMs producing 673 674 gene models that would be as accurate as possible for non-annotated genomes, we 675 selected the most closely related species (species with the most recent last common 676 ancestor) that possessed an annotated RefSeq genome. We then used the Augustus 677 training tool to build an HMM gene model for each genome.

678

679 We automated this process by creating a series of python scripts that performed the 680 following tasks:

681

1) 1.1\_insect\_database\_builder.py: This script collects the NCBI metadata
 regarding genomes and transcriptomes. Using the NCBI Entrez API, it collects the

684 most up to date information on RefSeq, GenBank, and TSA to generate two CSV 685 files: genome\_insect\_database.csv and transcriptome\_insect\_database.csv.

- 686 2) 1.2\_data\_downloader.py: This is a python wrapper around the *rsync* tool that
   687 downloads the sequence datasets present in the tables created by (1). It
   688 automatically downloads all the available information into a local folder.
- 3) 1.3\_run\_augustus\_training.py: This is a python wrapper around the Augustus
  training tool. It uses the metadata gathered using (1) and the sequence information
  gathered using (2) to build HMM gene models of all RefSeq datasets. It outputs
  sbatch scripts that can be run either locally, or on a SLURM-managed cluster.
  Those scripts will create unique HMM gene models per species.
- 694

695 At the time of this analysis (September 2019), 133 insect genomes were collected from 696 the RefSeq database, 309 genomes from the GenBank database, and 1123 697 transcriptomes from the TSA database. All the accession numbers and metadata are 698 available in the two tables (genome insect database.csv and 699 *transcriptome\_insect\_database.csv*) provided in the supplementary files. This pipeline 700 was repeated for crustaceans and the information can be found in the following two files: 701 genome crustacean database.csv and transcriptome crustacean database.csv.

- 702
- 703 Creation of protein sequence databases

The classical approach for orthology detection compares protein sequences to amino acid
 HMM corresponding to the gene of interest. Since we used three different NCBI
 databases, we performed the following preprocessing actions:

- 707
- RefSeq: well-annotated genomes from NCBI contain gene model translation; no
   extra processing was required.
- 2) GenBank: Using the HMMs created from the RefSeq databases, we created gene
  models for each GenBank genome using Augustus and a custom HMM gene
  model. To choose which HMM gene model to use, we selected the one for each
  insect order that had the highest training accuracy. In the case where an insect
  order did not have any member in the RefSeq database, we used the model of the

most closely related order. We then translated the inferred coding sequences to
create a protein database for each genome. The assignment of the models used
to infer the proteins of each GenBank genome is available in the *Table\_S4\_models.csv* available through the GitHub repository for this study at
<a href="https://github.com/extavourlab/Oskar\_Evolution">https://github.com/extavourlab/Oskar\_Evolution</a>. To automate the process, we
created a custom python script available in the file **1.4\_run\_augustus.py**.

- 3) TSA: Transcriptomes were translated using the emboss tool Transeq (Madeira, et al. 2019). We used this tool with the default parameters, except for the six-frame translation, trim and clean flags. This generated amino acid sequences for each transcript and each potential reading frame.
- 725

## 726 Identification of oskar orthologs

The oskar gene is composed of two conserved domains, LOTUS and OSK, separated by 727 728 a highly variable interdomain linker sequence (Ahuja and Extavour 2014; Jeske, et al. 729 2015; Yang, et al. 2015). To our knowledge, no other gene reported in any domain of life 730 possesses this domain composition (Blondel, et al. 2020). Therefore, here we use the 731 same definition of oskar orthology as in our previous work: a sequence possessing a 732 LOTUS domain followed by an interdomain region, and then an OSK domain (Blondel, et 733 al. 2020). To maximize the number of potential orthologs, we searched each sequence 734 with the previously generated HMM for the LOTUS and OSK domains (Blondel, et al. 735 2020). The presence and order of each domain were then verified for each potential hit 736 and only sequences with the previously defined Oskar structure were kept for further 737 processing. We used the HMMER 3.1 tool suite to build the domain HMM (hmmbuild with 738 default parameters), and then searched the generated protein databases (see Creation 739 of protein sequence databases above) using those models (hmmsearch with default parameters). Hits with an E-value  $\geq 0.05$  were discarded. A summary of all searches 740 741 performed is compiled in **Table S5** searches.csv in the GitHub repository for this study 742 at https://github.com/extavourlab/Oskar\_Evolution.

743

All the hits were then aligned with *hmmalign* with default parameters and the HMM of the full-length Oskar alignment previously generated (Blondel, et al. 2020). The resulting

sequences were automatically processed to remove assembly artifacts, and potential 746 747 isoforms. This filtration step was automated and went as follows: First, the sequences 748 were grouped by taxon. Then each group of sequences was aligned using MUSCLE 749 (Edgar 2004) with default parameters. The Hamming distance (Hamming 1950), a metric that computes the number of different letters between two strings, between each 750 751 sequence in the alignment, was computed. If any group of sequences had a Hamming 752 distance of >80%, then we only kept the sequence with the lowest E-value match. This 753 created a set of sequences containing multiple oskar orthologs per species only if they 754 were the likely product of a gene duplication event. We then used the resulting new alignment to generate a new domain HMM and a new full-length Oskar HMM (using 755 756 *hmmbuild* with default parameters) and ran further iterations of this detection pipeline until 757 we could detect no new oskar orthologs in the available sequence datasets. We called this final set the **filtered set** of sequences and used it in all subsequent orthology 758 759 analyses unless otherwise specified.

760

The Oskar sequences obtained are available in the following supplementary files:
Oskar filtered.aligned.fasta, Oskar filtered.fasta and Oskar consensus.hmm.

The domain definitions for the LOTUS and OSK domains are available in the following
supplementary files: Oskar\_filtered.aligned.LOTUS\_domain.fasta,
LOTUS\_consensus.hmm, Oskar\_filtered.aligned.OSK\_domain.fasta,
OSK consensus.hmm (see 1.5 Oskar tracker.ipynb).

767

768 Correlative analysis of assembly quality and absence of oskar

769 Using the metadata gathered previously from NCBI databases (see Genomes and 770 transcriptomes preprocessing above) we created two pools of source data: genomes 771 where we identified an oskar sequence, and genomes where we failed to find a sequence 772 that met our orthology criteria. We then compared the two distributions for each of the 8 773 available assembly statistics: (1) Contig and (2) Scaffold N50, (3) Contig and (4) Scaffold 774 L50, (5) Contig and (6) Scaffold counts, and (7) Number of Contigs and (8) Scaffolds per genome length. Finally, we performed a Mann-Whitney U statistical analysis to compare 775 776 the means of the two distributions (see 2.1\_Oskar\_discovery\_quality.ipynb).

#### 777

#### 778 TSA metadata parsing and curation

779 Datasets in the TSA database are associated with a biosample object that contains all 780 the metadata surrounding the RNA sequencing acquisitions. These metadata can include 781 information about one or both the tissue of origin and the organism's developmental 782 stage. We first automated the retrieval of these metadata using a custom python script that used the NCBI Entrez API (see 2.3 Oskar tissues stages.ipynb). However, the 783 784 metadata proved to be complex to parse for the following reasons: (1) not all projects had 785 the data entered in the corresponding tag, (2) some data contained typographical errors, 786 and (3) multiple synonyms were used to describe the same thing with different words in 787 different datasets. We therefore created a custom parsing and cleaning pipeline that 788 corrected mistakes and aggregated them into a cohesive set of unique terms that we 789 thought would be most informative to interpret the presence or absence of oskar orthologs 790 (see **2.3** Oskar tissues stages.ipynb to see the mapping table). This strategy 791 sacrificed some of the fine-grained information contained in custom metadata (for 792 example "right leg" became "leg") but allowed us to analyze the expression of oskar using 793 consistent criteria throughout all the datasets. This pipeline generated, for all available 794 datasets, a table of tissues and developmental stages including oskar presence or 795 absence in the dataset (see **Oskar all tissues stages.csv**).

796

797 Dimensionality reduction of Oskar alignment sequence space

798 The Oskar alignment was subjected to a Multiple Correspondence Analysis (MCA). 799 Similar to a PCA, dimension vectors were first computed to maximize the spread of the 800 underlying data in the new dimensions, except that instead of a continuous dataset, each 801 variable (here an amino acid at a given position) contributes to the continuous value on 802 that dimension. Once the projection vectors are computed, each sequence was then 803 mapped onto the dimensions. Each amino acid position (column) in the alignment was 804 considered a dimension with a possible value set of 21 (20 amino acids and gap). We 805 first removed the columns of low information (columns that had less than 30% amino acid 806 occupancy) using trimal (Capella-Gutierrez, et al. 2009) with a cutoff parameter set at 0.3. 807 Then, the alignment was decomposed into its eigenvectors, and projected to the first three

components. To perform this decomposition, we implemented a previously developed
preprocessing method (Rausell, et al. 2010) in a python script (see *MCA.py* and *2.8\_Oskar\_MCA\_Analysis.ipynb*) and performed the eigenvector decomposition with
the previously developed MCA python library (see *Key Resource Table*). We ran the
same algorithm on the LOTUS domain, OSK domain, and full-length Oskar alignments
obtained above (see *Identification of* oskar *orthologs* above).

814

## 815 Phylogenetic inference of Oskar sequences in the Hymenoptera

We aligned all hymenopteran Oskar sequences using PRANK (Loytynoja 2014) with 816 817 default parameters. We then manually annotated duplicated sequences by considering 818 two sequences from the same species that had < 80% amino acid identity, as within-819 species duplications of oskar. We trimmed this alignment to remove all columns with less 820 than 50% occupancy using trimal with the cutoff parameter set at 0.5. To reconstruct the 821 phylogeny of these sequences, we used the maximum likelihood inference software 822 RAxML (Stamatakis 2014) with a gamma-distributed protein model, and activated the flag 823 for auto model selection. We ran 100 bootstraps and then visualized and annotated the 824 obtained tree with Ete3 (Huerta-Cepas, et al. 2016) in a custom ipython notebook (see 825 2.7\_Oskar\_duplication.ipynb).

826

## 827 Calculation of Oskar conservation scores

Using the large set of orthologous Oskar sequences obtained as described above, we computed different conservation scores for each amino acid position. This methodology relies on the hypotheses that if an amino acid, or its associated chemical properties at a particular position in the sequence are important for the structure and/or function of the protein, they will be conserved across evolution. We considered multiple conservation metrics, each highlighting a particular aspect of the protein's properties as described in the followng sections. The scores can be found in the supplementary file **scores.csv**.

835

## 836 Computation of the Valdar score

The Valdar score (Valdar 2002) attempts to account for transition probabilities, stereochemical properties, amino acid frequency gaps, and, particularly essential for this

study, sequence weighting. Due to the heterogeneity of sequence dataset availability, most Oskar sequences occupy only a small portion of insect diversity, primarily Hymenoptera, and Diptera. Sequence weighting allows for the normalization of the influence of each sequence on the score based on how many similar sequences are present in the alignment (Valdar 2002). We implemented the algorithm described in (Valdar 2002) in a python script (see **besse\_blondel\_conservation\_scores.py**), then calculated the conservation scores for the Oskar alignment we generated above.

846

## 847 Computation of the Jensen-Shannon Divergence score

Jensen-Shannon Divergence (JSD) (Lin 1991: Capra and Singh 2007) uses the amino 848 acid and stereochemical properties to infer the "amount" of evolutionary pressure an 849 850 amino acid position may be subject to. This score uses an information theory approach 851 by measuring how much information (in bits) any position in the alignment brings to the 852 overall alignment (Capra and Singh 2007). This score also takes into account neighboring 853 amino acids in calculating the importance of each amino acid. We used the previously 854 published python code to calculate the JSD of our previously generated Oskar alignment 855 (Capra and Singh 2007) (see *score conservation.py*).

856

#### 857 Computation of the Conservation Bias

858 The measure of differences in conservation between the holometabolous and 859 hemimetabolous Oskar sequences presented in the results was done as follows: we first 860 split the alignment into two groups containing the sequences from each clade (see 861 **2.4 Oskar pgc specification.ipynb**). Due to the high heterogeneity in taxon sampling 862 between hemimetabolous and holometabolous insects, we ran a bootstrapped 863 approximation of the conservation scores on holometabolous sequences. We randomly 864 selected N sequences (N = the number of hemimetabolous sequences), computed the 865 Valdar conservation score (see Computation of the Valdar score above), and stored it. 866 After 1000 iterations, we computed the mean conservation score for each position for 867 holometabolous sequences. For hemimetabolous sequences, we directly calculated the 868 Valdar score using the method as described above (see Computation of the Valdar score 869 above). For each position, we then computed what we refer to as the "conservation bias"

870 between Holometabola and Hemimetabola by taking the ratio of the log of the

871 conservation score Holometabola and Hemimetabola. Conservation Bias = log(Valdarholo)

- 872 / log(Valdarhemi) for each position (see **3.4\_LogRatio\_Bootstrap.ipynb**).
- 873

#### 874 Computation of the electrostatic conservation score

875 To study the conservation of electrostatic properties of the Oskar protein we computed 876 implementation of an electrostatic our own conservation score (see 877 besse blondel conservation scores.py). Aspartic acid and Glutamic acid were given a score of -1, Arginine and Lysine a score of 1, and Histidine a score of 0.5. All other 878 879 amino acids were given a score of 0. Then, we summed the electrostatic score for each 880 sequence at each position and divided this raw score by the total number of sequences 881 in the alignment. This computation assigns a score between -1 and 1 at each position, -882 1 being a negative charge conserved across all sequences, and 1 a positive charge.

883

### 884 Computation of the hydrophobic conservation score

885 To study the conservation of hydrophobic properties of the Oskar protein we implemented conservation 886 hydrophobic our own score (see 887 besse\_blondel\_conservation\_scores.py). At each position, each amino acid was given 888 a hydrophobic score taken from a previously published scoring table (Moon and Fleming 889 2011). (This table is implemented in the **besse\_blondel\_conservation\_score.py** file for 890 simplicity.) Scores at each position were then averaged across all sequences. This metric 891 allowed us to measure the hydrophobicity conservation of each position in the alignment 892 and is bounded between 5.39 and -2.20.

893

#### 894 Computation of the RNA binding affinity score

895 RNA binding sites are defined as areas with positively charged residues and hydrophobic 896 residues. To estimate the conservation of RNA binding sites in *oskar* orthologs, we used 897 RNABindR v2.0 (Terribilini, et al. 2007), an algorithm predicting putative RNA binding 898 sites based on sequence information only. We automated the calculation for each 899 sequence by writing a python script that submitted a request to the RNABindR web 800 service (see *RNABindR\_run\_predictions.py*). We then aggregated all results into a 901 scoring matrix, and averaged the score obtained for each position. We call this score the 902 RNABindR score and hypothesize that it reflects the conservation of RNA binding 903 properties of the protein. Importantly, this score was obtained in 2017 for only a subset of 904 219 proteins used in this study (indicated in the supplementary files at: 905 03\_Oskar\_scores\_generation/RNABindR\_raw\_sources). Since then, the RNABindR 906 server has been defunct and we could not repeat those measurements as the source 907 code for this software is unavailable.

908

### 909 Computation of secondary structure conservation

Due to the overall low conservation of the LOTUS domain, we decided to see whether the secondary structure was conserved. To this end, we used the secondary structure prediction algorithm JPred 4 (Drozdetskiy, et al. 2015). Given an amino acid sequence, this tool returns a positional prediction for  $\alpha$ -helix,  $\beta$ -sheet or unstructured. We used the JPred4 web servers to compute the predictions and processed them into a secondary structure alignment (see **2.6\_Oskar\_lotus\_osk\_structures.ipynb**). We then used WebLogo (Crooks, et al. 2004) to visualize the conservation of the secondary structure.

917

#### 918 Visualization of conservation scores

919 We used PyMOL (DeLano 2002) to map the computed conservation scores onto the solved structures of LOTUS and OSK (Jeske, et al. 2015; Jeske, et al. 2017). At the time 920 921 of writing, no full-length Oskar protein structure had been reported. With the caveat that 922 all visualization was done on the structure of the Drosophila melanogaster protein 923 domains, we created a custom python script that augments PyMOL with automatic display 924 and coloring capacities. This script is available as **Oskar pymol visualization.py**, and 925 contains a manual at the beginning of the file. For the OSK domain, we used the structure 926 PDBID: 5A4A, and for the LOTUS domain, PDBID: 5NT7 (Jeske, et al. 2015; Jeske, et 927 al. 2017). The LOTUS structure we used is in complex with Vasa, and in a dimeric form 928 (Jeske, et al. 2017), allowing for easy interpretation of the different conservation scores. 929 For the OSK structure, we removed the residues 399-401 and 604-606 from the PDB file 930 as those amino acids did not align across all sequences and therefore showed highly 931 biased conservation scores.

932

## 933 Statistical analysis

934 All performed module statistical analyses were using the scipy stats 935 (https://www.scipy.org/). Significance thresholds for p-values were set at 0.05. Statistical 936 tests and p-values are reported in the figure legends. All statistical tests can be found in 937 the ipython notebooks mentioned below.

938

## 939 Data and code availability

The study generated a series of python 3 script and python 3 ipython notebook files that perform the entire analysis. All the results presented in this paper can be reproduced by running the aforementioned python 3 code. The primary data, *oskar* orthologs, Oskar alignments, trees, and conservation statistics as well as the code created and used are available as supplementary information. For ease of access, legibility, and reproducibility, the code and datasets have been deposited in a GitHub repository available at <u>https://github.com/extavourlab/Oskar\_Evolution.</u>

947

## 948 Software and libraries

All software and libraries used in this study are published under open source libre licensesand are therefore available to any researcher.

Туре	Name	Version	Source
Software	HMMER	3.1.b2	http://hmmer.org/
Software	PyMOL	1.8.x	https://pymol.org
Software	rsync	3.1.2	http://rsync.samba.org/
Software	Python 3	3.7	https://www.python.org/
Software	Mrbayes	3.2.6	http://nbisweden.github.io/MrBayes/
Software	trimal	1.2rev59	http://trimal.cgenomics.org/
Software	transeq	6.6.0.0	http://emboss.sourceforge.net/apps/cvs/emboss/apps/transeq.html
Software	augustus	2.5.5	http://augustus.gobics.de/
Software	JPred4	4.0	http://www.compbio.dundee.ac.uk/jpred/
Software	RNABindR	2.0	http://ailab1.ist.psu.edu/RNABindR/

Software	Inkscape	0.92.3	https://inkscape.org/
Library	jupyter	4.4.0	https://jupyter.org/
Library	ete3	3.3.1	http://etetoolkit.org
Library	pandas	0.25.1	https://pandas.pydata.org/
Library	mca	1.0.3	https://pypi.org/project/mca/
Library	fuzzywuzzy	0.17.0	https://github.com/seatgeek/fuzzywuzzy
Library	BeautifulSoup4	4.6.3	https://pypi.org/project/beautifulsoup4/
Library	biopython	1.74	https://pypi.org/project/biopython/
Library	numpy	1.16.2	https://www.numpy.org/
Library	seaborn	0.9.0	https://seaborn.pydata.org/
Library	matplotlib	3.0.0	https://matplotlib.org/
Library	scipy	1.1.0	https://www.scipy.org/
Library	progressbar	3.38.0	https://github.com/niltonvolpato/python-progressbar

951

952

## 953 Acknowledgements

954 This work was supported by funds from Harvard University, and support to SB from the

955 Master's in Bioinformatics Program of the University of Bordeaux. We thank members of

956 the Extavour lab for discussion.

957

### 958 Figure Legends

959

960 Figure 1: Overview of Oskar protein structure. The most common isoform of the Oskar 961 protein, Short Oskar, is composed of two well-folded domains, LOTUS and OSK, 962 separated by an interdomain sequence. A second isoform of the protein called Long 963 Oskar is present in some Dipteran insects, which contains a 5' domain as well as the 964 three domains of Short Oskar. Below the schematic representation is a rendering of the 965 previously reported solved structures for the LOTUS (PDBID: 5NT7) and OSK (PDBID: 966 5A4A) domains (Jeske, et al. 2015; Yang, et al. 2015) with a speculative rendering of the 967 unfolded interdomain region shown with a dashed line

968

969 Figure 2: Schematic presentation of the oskar ortholog detection pipeline. 970 Sequences were collected automatically from the three NCBI databases, GenBank 971 (GCA), RefSeq (GCF) and Transcriptome Shotgun Assembly Database (TSA). RefSeq 972 genomes were used to generate Augustus gene model HMMs, which were used to 973 annotate and predict proteins in the non-annotated genomes obtained from GenBank. 974 Transcripts from the TSA database were 6-frame translated using TRANSEQ. Amino acid 975 sequences were consolidated into three protein databases. hmmsearch from the HMMER 976 tool suite was used to search for LOTUS and OSK hits in those sequences. Sequences 977 with hits for both the LOTUS and OSK domains with an E-value <0.05 were annotated as 978 oskar sequences. Sequences were then cleaned to remove duplicates (sequences with 979 <80% sequence similarity coming from the same organism). The resulting sequences 980 were aligned using *hmmalign*, and the process was repeated until no new sequences 981 were identified. Finally, the sequences were consolidated with the dataset metadata into 982 the *oskar* ortholog database that was used for all subsequent analyses.

983

Figure 3: Summary of oskar distribution and expression in insects. Phylogeny from
(Misof, et al. 2014). Symbols in order from left to right: (i) vertical rectangles: grey: no
oskar ortholog was identified in this order. Color (unique for each order): at least one
oskar ortholog was identified in this order. (ii) number of datasets searched. (iii) horizontal
rectangles: proportion of searched datasets in which an oskar ortholog was identified. (iv)

989 pie chart: proportion of oskar sequences identified in RefSeq (GCF) datasets. (v) pie 990 chart: proportion of oskar sequences identified in GenBank (GCA) datasets. (vi) pie chart: 991 proportion of oskar sequences identified in Transcriptome Shotgun Assembly Database 992 (TSA) datasets. (vii) oskar sequences identified in tissue related to germ line 993 (transcriptomes derived from reproductive organs, eggs or embryos). (viii) oskar 994 sequences identified in tissue related to the brain (transcriptomes derived from brain or 995 head). (ix) oskar sequences identified in an egg stage transcriptome. (x) oskar sequences 996 identified in a larval stage transcriptome. (xi) oskar sequences identified in a pupal stage 997 transcriptome. (xii) oskar sequences identified in a nymphal or juvenile stage transcriptome. (xiii) oskar sequences identified in an adult transcriptome. All numbers 998 999 represented graphically here are in Supplementary Table 1. No datasets were available 1000 for Protura, Diplura or Isoptera at the time of analysis.

1001

Figure 4: Phylogenetic reconstruction of hymenopteran Oskar sequences. Phylogenetic tree inferred using RaxML with 100 bootstraps. Each leaf represents an Oskar ortholog. Gray circles: only one Oskar sequence was identified. Red circles: putatively duplicated Oskar sequences identified (sequence similarity <80%). Only families which contained a putative duplication are shown here; see Supplementary Figure S4 for the results of our *oskar* search in the context of a more complete hymenopteran phylogeny.

1009

Figure 5: Differential conservation of amino acids between hemimetabolous and 1010 1011 holometabolous Oskar sequences. (a) Box plot showing the conservation bias for each 1012 of the LOTUS and OSK domains between hemimetabolous and holometabolous Oskar 1013 sequences. Statistical difference was tested using a Mann Whitney U test (p < 0.05). (b) 1014 Ribbon diagram of LOTUS (PDBID: 5NT7) and OSK (PDBID: 5A4A) domain structures, 1015 where each amino acid is colored by conservation bias on the color scale shown in (a). (c, d) Protein surface representation of the OSK domain (PDBID: 5A4A) from two different 1016 1017 angles. Black dashed lines indicate the three amino acids reported previously to be 1018 necessary for OSK binding to RNA in *D. melanogaster* (Jeske, et al. 2015; Yang, et al. 1019 2015). (c) Amino acids colored by conservation bias on the color scale shown in (a).

1020 Cyan: amino acids more highly conserved in hemimetabolous sequences; magenta: 1021 amino acids more highly conserved in holometabolous sequences. **(d)** Amino acids 1022 colored by electrostatic conservation score. Left: hemimetabolous sequences; right: 1023 holometabolous sequences.

1024

1025 Figure 6: Conservation analysis of the LOTUS domain. (a) Ribbon diagram of a 1026 LOTUS domain dimer (cyan/magenta) in complex with two Vasa molecules (yellow) (PDBID: 5NT7) from two different angles. Each LOTUS amino acid is colored based on 1027 1028 its Valdar conservation score. (b, c) Sequence Logo of the  $\alpha$ 5 and  $\alpha$ 2/ $\alpha$ 3 helices respectively generated with WebLogo (Crooks, et al. 2004). Black: hydrophobic residues; 1029 1030 blue: charged residues; green: polar residues. (b', b") Ribbon diagram of the conserved  $\alpha$ 5 helix, with key amino acids displayed as sticks and colored by Valdar conservation 1031 score. Two potential novel Vasa-LOTUS contacts (H227 and Q235) are highlighted with 1032 1033 dashed lines. (c') Ribbon diagram of the conserved  $\alpha 2$  helix, with key amino acids 1034 displayed as sticks and colored by hydrophobicity/hydrophily conservation score. (c") 1035 Ribbon diagram of the conserved  $\alpha 2$  helix, with key amino acids displayed as sticks and 1036 colored by Valdar conservation score. (d) Surface mesh rendering colored with the 1037 RNABindR RNA binding conservation score. (e, f) Ribbon diagram of the LOTUS  $\beta$  sheet dimerization interface. Left: conservation of monomeric LOTUS domains: right: dimeric 1038 1039 LOTUS domains. (e) Amino acids colored by electrostatic conservation score. Dashed lines indicate the key electrostatic interaction thought to stabilize the dimerization. (f) 1040 1041 Amino acids colored by hydrophobicity/hydrophily conservation score. Dashed lines 1042 indicate the key hydrophobic pocket thought to stabilize the dimerization.

1043

**Figure 7: Conservation analysis of the OSK domain. (a)** Ribbon diagram of the OSK domain (PDBID: 5A4A) from two different angles. Each amino acid is colored based on its Valdar conservation score.**(b)** Protein surface representation of the OSK domain colored by Valdar conservation, electrostatic conservation and hydrophobicity/hydrophily conservation score. **(c, c', c'', c''')** Ribbon diagram of newly detected conserved motifs of the OSK domain, showing sequence Logo (bottom row) residues as sticks. Each amino
1050 acid is colored with Valdar conservation scores of holometabolous (top row) and 1051 hemimetabolous (middle row) OSK sequences. Bottom row: sequence Logos of each 1052 conserved motif generated with WebLogo (Crooks, et al. 2004). Black: hydrophobic residues; blue: charged residues; green: polar residues. Red numbers: amino acid 1053 locations of D. melanogaster loss of function oskar alleles leading to the loss of oskar 1054 localization to the posterior pole during embryogenesis (P425S = osk/8) (Kim-Ha, et al. 1055 1991); S452L = osk[255] = osk[7] (Lehmann and Nüsslein-Volhard 1986; Kim-Ha, et al. 1056 1991); S457F = osk[6B10] (Breitwieser, et al. 1996)) or to reduced RNA-binding affinity 1057 of the OSK domain (R436E (Yang, et al. 2015)). 1058

#### 1059 References

Ahuja A, Extavour CG. 2014. Patterns of molecular evolution of the germ line specification gene *oskar* suggest that a novel domain may contribute to functional divergence in *Drosophila*.
Development Genes and Evolution 222:65-77.

- Amy RL. 1961. The embryology of *Habobracon juglandis* (Ashmead). Journal of Morphology1064 109:199-217.
- Anantharaman V, Zhang D, Aravind L. (p15752 co-authors). 2010. OST-HTH: a novel predicted
   RNA-binding domain. Biology direct 5:13.
- Anderson DT, Wood EC. 1968. The morphological basis of embryonic movements in the light
  brown apple moth, *Epiphyas postvittana* (Walk.) (Lepidoptera, Tortricidae). Australian Journal of
  Zoology 16:763-793.
- 1070 Ando H, Tanaka M. 1979. Early embryonic devekopment of the primitive moths, *Enduclyta*
- 1071 *signifer* Walker and *E. excrescens* Butler (Lepidoptera: Hepialidae). International Journal of
- 1072 Insect Morphology and Embryology 9:67-77.
- 1073 Berg GJ, Gassner G. 1978. Fine structure of the blastoderm embryo of the pink bollworm,
- 1074 *Pectinophora gossypiella* (Saunders) (Lepidoptera: gelechiidae). International Journal of Insect 1075 Morphology and Embryology 1:81+105.
- Blondel L, Jones TEM, Extavour CG. 2020. Bacterial contribution to genesis of the novel germline determinant *oskar*. Elife 9:e45539.
- Breitwieser W, Markussen F-H, Horstmann H, Ephrussi A. 1996. Oskar protein interaction with
  Vasa represents an essential step in polar granule assembly. Genes and Development:21792188.
- Bronskill JF. 1959. Embryology of *Pimpla turionellae* (L.) (Hymenoptera: Ichneumonidae).
  Canadian Journal of Zoology 37:655-688.
- Bull AL. 1982. Stages of living embryos in the jewel wasp *Mormoniella (Nasonia) vitripennis*(Walker) (Hymenoptera: Pteromalidae). International Journal of Insect Morphology and
  Embryology 1:1-23.
- Bütschli O. 1870. Zur Entwicklungsgeschichte der Biene. Zeitschrift für WissenschaftlicheZoologie 20:519-564.
- Butt FH. 1949. Embryology of the Milkweed Bug, *Oncopeltus fasciatus* (Hemiptera). Cornell
  Experiment Station Memoir 283:2-43.
- 1090 Callebaut I, Mornon J-P. (p19296 co-authors). 2010. LOTUS, a new domain associated with 1091 small RNA pathways in the germline. Bioinformatics 26:1140-1144.
- 1092 Capella-Gutierrez S, Silla-Martinez JM, Gabaldon T. 2009. trimAl: a tool for automated 1093 alignment trimming in large-scale phylogenetic analyses. Bioinformatics 25:1972-1973.
- Capra JA, Singh M. 2007. Predicting functionally important residues from sequenceconservation. Bioinformatics 23:1875-1882.

- Carter J-M, Baker SC, Pink R, Carter DRF, Collins A, Tomlin J, Gibbs M, Breuker CJ. 2013.
   Unscrambling butterfly oogenesis. BioMedCentral Genomics 14:283-283.
- 1098 Carter JM, Gibbs M, Breuker CJ. 2015. Divergent RNA Localisation Patterns of Maternal Genes 1099 Regulating Embryonic Patterning in the Butterfly Pararge aegeria. PLoS ONE 10:e0144471.
- 1100 Chang CC, Lee WC, Cook CE, Lin GW, Chang T. 2006. Germ-plasm specification and germline 1101 development in the parthenogenetic pea aphid *Acyrthosiphon pisum*: Vasa and Nanos as
- 1102 markers. International Journal of Developmental Biology 50:413-421.
- Chen X-x, Achterberg Cv. 2018. Systematics, Phylogeny, and Evolution of Braconid Wasps: 30
  Years of Progress. Annual Review of Entomology 64:1-24.
- Clark AG, Eisen MB, Smith DR, Bergman CM, Oliver B, Markow TA, Kaufman TC, Kellis M,
  Gelbart W, Iyer VN, et al. (p04187 co-authors). 2007. Evolution of genes and genomes on the *Drosophila* phylogeny. Nature 450:203-218.
- 1108 Crooks GE, Hon G, Chandonia JM, Brenner SE. 2004. WebLogo: a sequence logo generator.1109 Genome Research 14:1188-1190.
- Dearden PK. 2006. Germ cell development in the Honeybee (*Apis mellifera*); vasa and nanos
  expression. BMC Developmental Biology 6:6.
- Dearden PK, Wilson MJ, Sablan L, Osborne PW, Havler M, McNaughton E, Kimura K, Milshina
  NV, Hasselmann M, Gempe T, et al. 2006. Patterns of conservation and change in honey bee
  developmental genes. Genome Research 16:1376-1384.
- 1115 DeLano WL. 2002. Pymol: An Open-Source Molecular Graphics Tool. CCP4 Newsletter on1116 Protein Crystallography 40:82-92.
- Donnell DM, Corley LS, Chen G, Strand MR. 2004. Caste determination in a polyembryonic
  wasp involves inheritance of germ cells. Proceedings of the National Academy of Sciences of
  the United States of America 101:10095-10100.
- 1120 Drozdetskiy A, Cole C, Procter J, Barton GJ. 2015. JPred4: a protein secondary structure 1121 prediction server. Nucleic Acids Research 43:W389-394.
- Eastham LES. 1930. The embryology of *Pieris rapae* Organogeny. Philosophical Transactions
  of the Royal Society of London. Series B: Biological Sciences 219:1-50.
- 1124 Edgar RC. (r20416 co-authors). 2004. MUSCLE: a multiple sequence alignment method with 1125 reduced time and space complexity. BMC Bioinformatics 5:113.
- 1126 Ephrussi A, Lehmann R. 1992. Induction of germ cell formation by *oskar*. Nature 358:387-392.
- Ewen-Campen B, Schwager EE, Extavour CG. 2010. The molecular machinery of germ linespecification. Molecular Reproduction and Development 77:3-18.
- 1129 Ewen-Campen B, Srouji JR, Schwager EE, Extavour CG. 2012. *oskar* Predates the Evolution of 1130 Germ Plasm in Insects. Current Biology 22:2278-2283.

- Extavour CG, Akam ME. 2003. Mechanisms of germ cell specification across the metazoans:
   epigenesis and preformation. Development 130:5869-5884.
- Field J, Ohl M, Kennedy M. 2011. A molecular phylogeny for digger wasps in the tribe Ammophilini (Hymenoptera, Apoidea, Sphecidae). Systematic Entomology 36:732-740.
- Fleig R, Sander K. 1985. Blastoderm development in honey bee embryogenesis as seen in the
   scanning electron microscope. International Journal of Invertebrate Reproduction and
- 1137 Development 8:279-286.
- Fleig R, Sander K. 1986. Embryogenesis of the Honeybee *Apis mellifera* L (Hymenoptera,
  Apidae) an SEM Study. International Journal of Insect Morphology and Embryology 15:449462.
- Fleischmann VG. 1975. Origin and embryonic development of fertile gonads with and without
  pole cells of *Pimpla turionellae* L. (Hymenoptera, Ichneumonidae). Zool. Jb. Anat. Bd. 94:375411.
- 1144 Gatenby JB. 1920. The Cytoplasmic Inclusions of the Germ Cells. Part VI. On the origin and 1145 probable constitution of the germ-cell determinant of *Apanteles glomeratus*, with a note on the 1146 secondary nuclei. Quarterly Journal of Microscopical Science 64:133-153.
- Gatenby JB. 1917a. The embryonic development of *Trichogramma evanescens* Westw.,
  monoembryonic egg parasite of *Donacia simplex*. Quarterly Journal of Microscopical Science
  62:149-187.
- 1150 Gatenby JB. 1918. The segregation of germ cells in *Trichogramma evanescens*. Quarterly 1151 Journal of Microscopical Science 63:161-173.
- 1152 Gatenby JB. 1917b. The segregation of the germ-cells in *Trichogramma evanescens*. Quarterly 1153 Journal of Microscopical Science 62:149-187.
- 1154 Grbic' M. 2000. "Alien" wasps and evolution of development. Bioessays 22:920-932.
- Grbic' M. 2003. Polyembryony in parasitic wasps: evolution of a novel mode of development.
  International Journal of Developmental Biology 47:633-642.
- Grbic' M, Nagy LM, Carroll SB, Strand M. 1996. Polyembryonic development: insect patternformation in a cellularised environment. Development:795-804.
- 1159 Guelin M. 1994. [Activity of W-sex heterochromatin and accumulation of the nuage in nurse 1160 cells of the lepidopteran *Ephestia*]. C. R. Acad. Sci. Paris. Ser. III 317:54-61.
- Gutzeit HO, Zissler D, Fleig R. 1993. Oogenesis in the Honeybee *Apis mellifera* Cytological
   Observations on the Formation and Differentiation of Previtellogenic Ovarian Follicles. Rouxs
   Archives of Developmental Biology 202:181-191.
- Hamming RW. 1950. Error Detecting and Error Correcting Codes. The Bell System TechnicalJournal 29:147-160.

Hay B, Jan LY, Jan YN. 1990. Localization of *vasa*, a component of *Drosophila* polar granules,
in maternal-effect mutants that alter embryonic anteroposterior polarity. Development 109:425433.

- 1169 Hegner RW. 1914. Studies on germ cells. III. The origin of the Keimbahn-determinants in a 1170 parasitic Hymenopteran, *Copidosoma*. Anatomischer Anzeiger 3-4:51-69.
- 1171 Heming BS, Huebner E. 1994. Development of the Germ Cells and Reproductive Primordia in
- 1172 Male and Female Embryos of *Rhodnius prolixus* Stal (Hemiptera, Reduviidae). Canadian 1173 Journal of Zoology 72:1100-1119.
- Huerta-Cepas J, Serra F, Bork P. 2016. ETE 3: Reconstruction, Analysis, and Visualization of
   Phylogenomic Data. Molecular Biology and Evolution 33:1635-1638.
- Jeske M, Bordi M, Glatt S, Muller S, Rybin V, Muller CW, Ephrussi A. 2015. The Crystal
  Structure of the *Drosophila* Germline Inducer Oskar Identifies Two Domains with Distinct Vasa
  Helicase- and RNA-Binding Activities. Cell Reports 12:587-598.
- Jeske M, Muller CW, Ephrussi A. 2017. The LOTUS domain is a conserved DEAD-box RNA
  helicase regulator essential for the recruitment of Vasa to the germ plasm and nuage. Genes
  and Development 31:939-952.
- Johannsen OA. 1929. Some phases in the embryonic development of *Diacrisia virginica* Fabr.
  (Lepidoptera). J. Morphol. Physiol. 2:493-541.
- Jones JR, Macdonald PM. 2007. Oskar controls morphology of polar granules and nuclear
  bodies in *Drosophila*. Development 134:233-236.
- Juhn J, James AA. 2006. *oskar* gene expression in the vector mosquitoes, *Anopheles gambiae*and *Aedes aegypti*. Insect Molecular Biology 15:363-372.
- Juhn J, Marinotti O, Calvo E, James AA. 2008. Gene structure and expression of *nanos* (*nos*)
  and *oskar* (*osk*) orthologues of the vector mosquito, *Culex quinquefasciatus*. Insect Molecular
  Biology 17:545-552.
- 1191 Kawahara AY, Plotkin D, Espeland M, Meusemann K, Toussaint EFA, Donath A, Gimnich F,
- 1192 Frandsen PB, Zwick A, Dos Reis M, et al. 2019. Phylogenomics reveals the evolutionary timing
- and pattern of butterflies and moths. Proceedings of the National Academy of Sciences of the
- 1194 United States of America 116:22657-22663.
- 1195 Kelly GM, Huebner E. 1989. Embryonic development of the hemipteran insect *Rhodnius* 1196 *prolixus*. Journal of Morphology 199:175-196.
- Khila A, Abouheif E. 2008. Reproductive constraint is a developmental mechanism that
  maintains social harmony in advanced ant societies. Proceedings of the National Academy of
  Sciences of the United States of America 105:17884-17889.
- 1200 Kim-Ha J, Smith JL, Macdonald PM. 1991. *oskar* mRNA is localized to the posterior pole of the 1201 *Drosophila* oocyte. Cell 66:23-35.

- Kirk DL. 2005. A twelve-step program for evolving multicellularity and a division of labor.Bioessays 27:299-310.
- Kobayashi Y, Ando H. 1984. Mesodermal Organogenesis in the Embryo of the Primitive Moth, *Neomicropteryx nipponensis* Issiki (Lepidoptera, Micropterygidae). Journal of Morphology
  181:29-47.
- Koscielska MK, Koscielski B. 1987. Early embryonic development of *Tritneptis diprionis*(Chalcidoidea, Hymenoptera). In: Ando H, Jura C, editors. Recent Advances in Insect
  Embryology in Japan and Poland. Tsukuba: Arthropod. Embryol. Soc. Jpn.
- 1210 ISEBU Co. Ltd. p. 207-214.
- Lasko P. 2013. The DEAD-box helicase Vasa: evidence for a multiplicity of functions in RNA
   processes and developmental biology. Biochimica et Biophysica Acta 1829:810-816.
- Lautenschlager F. 1932. Die Embryonalentwicklung der weiblichen Keimdruse bei der Psychide
  Solenobia triquetella. Zool. Jarh. 56:121-162.
- 1215 Lebart L, Morineau A, Warwick KM. 1984. Multivariate Descriptive Statistical Analysis:
- 1216 Correspondence Analysis and Related Techniques for Large Matrices. Chichester, UK: John 1217 Wiley & Sons.
- Lehmann R. 2016. Germ Plasm Biogenesis--An Oskar-Centric Perspective. Current Topics inDevelopmental Biology 116:679-707.
- Lehmann R, Nüsslein-Volhard C. 1986. Abdominal Segmentation, Pole Cell Formation, and
  Embryonic Polarity Require the Localized Activity of *oskar*, a Maternal Gene in *Drosophila*. Cell
  47:144-152.
- Lin GW, Cook CE, Miura T, Chang CC. 2014. Posterior localization of ApVas1 positions the preformed germ plasm in the sexual oviparous pea aphid Acyrthosiphon pisum. EvoDevo 5:18.
- Lin J. 1991. Divergence Measures Based on the Shannon Entropy. IEEE Transactions on Information Theory / Professional Technical Group on Information Theory 37:145-151.
- Loytynoja A. 2014. Phylogeny-aware alignment with PRANK. Methods in Molecular Biology1079:155-170.
- Lynch JA, Desplan C. (p16123 co-authors). 2010. Novel modes of localization and function of *nanos* in the wasp *Nasonia*. Development 137:3813-3821.
- Lynch JA, Özüak O, Khila A, Abouheif E, Desplan C, Roth S. 2011. The Phylogenetic Origin of *oskar* Coincided with the Origin of Maternally Provisioned Germ Plasm and Pole Cells at the
  Base of the Holometabola. PLoS Genetics 7:e1002029.
- Maddison DR, Schultz K-S, Maddison WP. 2007. The Tree of Life Web Project. Zootaxa1668:19-40.
- 1236 Madeira F, Park YM, Lee J, Buso N, Gur T, Madhusoodanan N, Basutkar P, Tivey ARN, Potter
- 1237 SC, Finn RD, et al. 2019. The EMBL-EBI search and sequence analysis tools APIs in 2019.
- 1238 Nucleic Acids Research 47:W636-W641.

- Malm T, Nyman T. 2015. Phylogeny of the symphytan grade of Hymenoptera: new pieces into the old jigsaw(fly) puzzle. Cladistics 31:1-17.
- Markussen FH, Michon AM, Breitwieser W, Ephrussi A. 1995. Translational control of *oskar*generates short OSK, the isoform that induces pole plasm assembly. Development 121:37233732.
- 1244 Mellanby H. 1935. The early embryonic development of *Rhodnius prolixus* (Hemiptera, 1245 Heteroptera). Quarterly Journal of Microscopical Science 78:71-90.
- 1246 Metschnikoff E. 1866. Embryologische Studien an Insekten. Zeit. f. wiss Zool. 16:389-500.
- Misof B, Liu S, Meusemann K, Peters RS, Donath A, Mayer C, Frandsen PB, Ware J, Flouri T,
  Beutel RG, et al. 2014. Phylogenomics resolves the timing and pattern of insect evolution.
  Science 346:763-767.
- Mitter C, Davis DR, Cummings MP. 2017. Phylogeny and Evolution of Lepidoptera. AnnualReview of Entomology 62:265-283.
- Miura T, Braendle C, Shingleton A, Sisk G, Kambhampati S, Stern DL. 2003. A comparison of
   parthenogenetic and sexual embryogenesis of the pea aphid *Acyrthosiphon pisum* (Hemiptera :
   Aphidoidea). Journal of Experimental Zoology Part B-Molecular and Developmental Evolution
   295B:59-81.
- Miya K. 1953. The presumptive genital region at the blastoderm stage of the silkworm egg.Journal of the Faculty of Agriculture of Iwate University:223-227.
- Miya K. 1958. Studies on the embryonic development of the gonad in the silkworm, *Bombyx mori* L. Part I. Differentiation of germ cells. Journal of the Faculty of Agriculture of Iwate
  University 3:436-467.
- Miya K. 1975. Ultrastructural changes of embryonic cells during organogenesis in the silkworm, *Bombyx mori*. I. The Gonad. Journal of the Faculty of Agriculture of Iwate University 12:329338.
- Moon CP, Fleming KG. 2011. Side-chain hydrophobicity scale derived from transmembrane
  protein folding into lipid bilayers. Proceedings of the National Academy of Sciences of the
  United States of America 108:10174-10177.
- Nagy L, Riddiford L, Kiguchi K. 1994. Morphogenesis in the Early Embryo of the Lepidopteran *Bobyx mori*. Developmental Biology:137-151.
- 1269 Nakao H. 1999. Isolation and characterization of a *Bombyx vasa*-like gene. Development Genes1270 and Evolution 209:312-316.
- Nakao H, Hatakeyama M, Lee JM, Shimoda M, Kanda T. 2006. Expression pattern of *Bombyx* vasa-like (BmVLG) protein and its implications in germ cell development. Development Genes
   and Evolution 216:94-99.

- 1274 Nakao H, Matsumoto T, Oba Y, Niimi T, Yaginuma T. 2008. Germ cell specification and early 1275 embryonic patterning in Bombyx mori as revealed by nanos orthologues. Evolution and
- 1276 Development 10:546-554.
- Nakao H, Takasu Y. 2019. Complexities in Bombyx germ cell formation process revealed by
   Bm-nosO (a Bombyx homolog of nanos) knockout. Developmental Biology 445:29-36.
- 1279 Nelson JA. 1915. The embryology of the honey bee. Princeton: Princeton University Press.
- 1280 Noce T, Okamoto-Ito S, Tsunekawa N. 2001. *Vasa* homolog genes in mammalian germ cell 1281 development. Cell Structure and Function 26:131-136.
- Nyman T, Zinovjev AG, Vikberg V, Farrell BD. 2006. Molecular phylogeny of the sawfly
  subfamily Nematinae (Hymenoptera: Tenthredinidae). Systematic Entomology 31:569-583.
- Pal D, Chakrabarti P. 2001. Non-hydrogen bond interactions involving the methionine sulfur
   atom. Journal of Biomolecular Structure and Dynamics 19:115-128.
- 1286 Peters RS, Krogmann L, Mayer C, Donath A, Gunkel S, Meusemann K, Kozlov A,
- Podsiadlowski L, Petersen M, Lanfear R, et al. 2017. Evolutionary History of the Hymenoptera.
  Current Biology 27:1013-1018.
- Peters RS, Niehuis O, Gunkel S, Bläser M, Mayer C, Podsiadlowski L, Kozlov A, Donath A,
  Noort Sv, Liu S, et al. 2018. Transcriptome sequence-based phylogeny of chalcidoid wasps
  (Hymenoptera: Chalcidoidea) reveals a history of rapid radiations, convergence, and
- 1292 evolutionary success. Molecular Phylogenetics and Evolution 120:286-296.
- Presser BD, Rutschky CW. 1957. The embryonic development of the corn earworm, *Heliothis zea* (Boddie) (Lepidoptea, Phalaenidae). Annals of the Entomological Society of America
  50:133-164.
- Prous M, Blank SM, Goulet H, Heibo E, Liston A, Malm T, Nyman T, Schmidt S, Smith DR,
  Vårdal H, et al. 2014. The genera of Nematinae (Hymenoptera, Tenthredinidae). Journal of
  Hymenoptera Research 40:1-69.
- Quan H, Arsala D, Lynch JA. 2019. Transcriptomic and functional analysis of the oosome, a
  unique form of germ plasm in the wasp *Nasonia vitripennis*. BMC Biology 17:78.
- 1301 Quan H, Lynch JA. 2016. The evolution of insect germline specification strategies. Current1302 Opinion in Insect Science 13:99-105.
- 1303 Rafiqi AM, Rajakumar A, Abouheif E. 2020. Origin and elaboration of a major evolutionary
  1304 transition in individuality. Nature 585:239-244.
- Rausell A, Juan D, Pazos F, Valencia A. 2010. Protein interactions and ligand binding: from
  protein subfamilies to functional specificity. Proceedings of the National Academy of Sciences of
  the United States of America 107:1995-2000.
- Raz E. 2000. The function and regulation of *vasa*-like genes in germ-cell development. GenomeBiology 1:1-6.

- 1310 Saito. 1937. On the development of the Tusser, *Antheraea pernyi* Guerin-Meneville, with special
- reference to the comparative embryology of insects. Journal of the Faculty of Agriculture ofHokkaido Imperial University 40:35-109.
- 1313 Schmidt C. 2013. Molecular phylogenetics of ponerine ants (Hymenoptera: Formicidae:
- 1314 Ponerinae). Zootaxa 3647:201-250.
- 1315 Schroder R. 2006. *vasa* mRNA accumulates at the posterior pole during blastoderm formation in 1316 the flour beetle *Tribolium castaneum*. Development, Genes and Evolution 216:277-283.
- 1317 Schwangart F. 1905. Zur Entwickslungsgeschichte der Lepidopteren. Biol. Centralbl. 25:777-1318 789.
- Sehl A. 1931. Furchung und Bildung der Keimanlage bei der Mehlmotte *Ephestia kuehniella*.Zell. Zeit. Morph. U. Okol. 1:429-506.
- Shafiq SA. 1954. A study of the embryonic development of the Gooseberry Sawfly, *Pteronidea ribesii*. Quarterly Journal of Microscopical Science 95:93-114.
- 1323 Sharanowski BJ, Ridenbaugh RD, Piekarski PK, Broad GR, Burke GR, Deans AR, Lemmon AR,

1324 Lemmon ECM, Diehl GJ, Whitfield JB, et al. 2021. Phylogenomics of Ichneumonoidea

- 1325 (Hymenoptera) and implications for evolution of mode of parasitism and viral endogenization.
- 1326 Molecular Phylogenetics and Evolution 156:107023.
- Sikosek T, Chan HS. 2014. Biophysics of protein evolution and evolutionary protein biophysics.
  Journal of the Royal Society, Interface / the Royal Society 11:20140419.
- Sikosek T, Chan HS, Bornberg-Bauer E. 2012. Escape from Adaptive Conflict follows from
   weak functional trade-offs and mutational robustness. Proceedings of the National Academy of
   Sciences of the United States of America 100:14888, 14802
- 1331 Sciences of the United States of America 109:14888-14893.
- Silvestri F. 1906. Contribuzioni alla conoscenza biologica degli Imenotteri parassiti. I. Biologia
  del *Litomastix truncellatus* Dalm. Annali della r. Scuola Superiore di Agricoltura in Portici 6:3-51.
- Silvestri F. 1908. Contribuzioni alla conoscenza degli Imenotteri parassiti. Bollettino del
  Laboratorio di Zoologia Generale e Agraria della r. Scuola Superiore d'Agricoltura (AFTW.
  Facoltà Agraria) in Portici 3:29-84.
- Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis oflarge phylogenies. Bioinformatics 30:1312-1313.
- Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B. 2006. AUGUSTUS: ab initioprediction of alternative transcripts. Nucleic Acids Research 34:W435-439.
- Strand MR, Grbic' M. 1997. The Development and Evolution of Polyembryonic Insects. CurrentTopics in Developmental Biology 35:121-159.
- 1343 Sumitani M, Yamamoto DS, Oishi K, Lee JM, Hatakeyama M. 2003. Germline transformation of 1344 the sawfly, *Athalia rosae* (Hymenoptera: Symphyta), mediated by a piggyBac-derived vector.
- 1345 Insect Biochem Mol Biol 33:449-458.

1346 Tanaka M. 1987. Differentiation and behaviour of Primordial Germ Cells during the Early

1347 Embryonic Development of *Parnassius glacialis* Butler, *Luehdorfia japonica* Leech and *Byasa* 

1348 (Atrophaneura) alcinous alcinous Klug (Lepidoptera: Papilionidae). In: Ando H, Jura C, editors.

- 1349 Recent Advances in Insect Embryology in Japan and Poland. Tsukuba: Arthropod. Embryol.
- 1350 Soc. Jpn.
- 1351 ISEBU Co. Ltd. p. 255-266.

Tawfik MFS. 1957. Alkaline phosphatase in the germ-cell determinant of the egg of *Apanteles*.Journal of Insect Physiology 1:286-291.

- 1354 Terribilini M, Sander JD, Lee JH, Zaback P, Jernigan RL, Honavar V, Dobbs D. 2007.
- 1355 RNABindR: a server for analyzing and predicting RNA-binding sites in proteins. Nucleic Acids1356 Research 35:W578-584.
- 1357 Tomaya K. 1902. On the embryology of the silkworm. Bulletin of the College of Agriculture,1358 Tokyo 5:73-111.
- 1359 Toshiki T, Chantal C, R., Toshio K, Eappen A, Mari K, Natuo K, Jean-Luc T, Bernard M, Gérard
- 1360 C, Paul S, et al. 2000. Germline transformation of the silkworm *Bombyx mori* L. using a piggyBac transposon-derived vector. Nature Biotech.:81-84.
- 1362 Valdar WS. 2002. Scoring residue conservation. Proteins 48:227-241.
- 1363 Vilhelmsen L. 2015. Morphological phylogenetics of the Tenthredinidae (Insecta:Hymenoptera).1364 Invertebrate Systematics 29:164-190.
- Ward PS. 2014. The Phylogeny and Evolution of Ants. Annual Review of Ecology, Evolution,and Systematics 45:23-43.

Ward PS, Blaimer BB, Fisher BL. 2016. A revised phylogenetic classification of the ant
subfamily Formicinae (Hymenoptera: Formicidae), with resurrection of the genera Colobopsis
and Dinomyrmex. Zootaxa 4072:343-357.

- 1370 Webster PJ, Suen J, Macdonald PM. 1994. *Drosophila virilis oskar* transgenes direct body
  1371 patterning but not pole cell formation or maintenance of mRNA localization in *D. melanogaster*.
  1372 Development 120:2027-2037.
- Wiegmann BM, Trautwein MD, Winkler IS, Barr NB, Kim J-W, Lambkin C, Bertone MA, Cassel
  BK, Bayless KM, Heimberg AM, et al. (r40919 co-authors). 2011. Episodic radiations in the fly
  tree of life. Proceedings of the National Academy of Sciences 108:5690-5695.
- 1376 Will L. 1888. Entwicklungsgescshichte der viviparen Aphiden. Zool. Jarh. 3:201-280.
- 1377 Witlaczil E. 1884. Entwicklungsgeschichte der Aphiden. Zeit. f. wiss Zool. 40:559-690.
- Woodworth CW. 1889. Studies on the embryological development of *Euvanessa antiopa*. In:
  Scudder, editor. Butterflies of Eastern United States and Canada. p. 102.

Xu X, Brechbiel JL, Gavis ER. 2013. Dynein-Dependent Transport of *nanos* RNA in *Drosophila* Sensory Neurons Requires Rumpelstiltskin and the Germ Plasm Organizer Oskar. Journal of
 Neuroscience 33:14791-14800.

- 1383 Yang N, Yu Z, Hu M, Wang M, Lehmann R, Xu RM. 2015. Structure of Drosophila Oskar
- 1384 reveals a novel RNA binding protein. Proceedings of the National Academy of Sciences of the 1385 United States of America 112:11541-11546.

1386 Zhurov V, Terzin T, Grbic M. 2004. Early blastomere determines embryo proliferation and caste
 1387 fate in a polyembryonic wasp. Nature 432:764-769.

1388 Zissler D. 1992. From egg to pole cells: ultrastructural aspects of early cleavage and germ cell 1389 determination in insects. Micr. Res. and Tech.:49-74.

1391 Figure 1





1396

# Collect available sequences datasets from NCBI and generate amino acid sequences databases



1398 Figure 3

1399



### 1401 Figure 4



## 1403 Figure 5







Blondel et al. Page 54 of 72

1409	Supplementary Materials							
1410								
1411	Evolution of a cytoplasmic determinant: evidence for the biochemical basis of							
1412	functional evolution of a novel germ line regulator							
1413								
1414	Leo Blondel, Savandara Besse and Cassandra G. Extavour							
1415								
1416	These Supplementary Materials contain the following:							
1417	<ul> <li>Supplementary Tables S1 through S5</li> </ul>							
1418	<ul> <li>Supplementary Figures S1 through S9</li> </ul>							
1419 1420								

### 1421 Supplementary Table Legends

1422

Supplementary Table S1: Number of oskar sequences identified per order and per data source. Each row corresponds to an order and a data source: GCF: RefSeq; GCA: GenBank, TSA: Transcriptome Shotgun Assembly Database. "Filtered hits" column indicates the number of hits after the filtration algorithm described in the Methods is applied. Rightmost column defines the proportion of *oskar* sequences identified, as the number of datasets with a filtered hit divided by the total number of datasets searched.

- 1429
- Supplementary Table S2: Genome quality correlation to oskar identification. Mean and median values for the distributions of each indicated genome quality parameter, in which oskar was (a) or was not (b) identified. The means of both distributions are significantly different for all metrics (Mann Whitney U test, p<0.05). See Supplementary Figure S2 of graphical representation of distributions.
- 1435 1436

1441

1437Supplementary Table S3. Assignment of metadata to germ line or brain categories.1438Thistableisfoundin1439Data>02\_oskar\_analyses/2/3/TableS3\_germline\_brain\_table.csvattheGitHub1440repositoryhttps://github.com/extavourlab/Oskar\_Evolution.theSitHub

Supplementary Table S4. Models used to create protein sequence databases. This 1442 1443 table shows which models were used to run the ab initio gene detection algorithm Augustus as described in Methods and Materials. Column order corresponds to any GCA 1444 dataset of an organism from this order. "Family" column is only used if a member of this 1445 1446 order but of a different family was used. Finally, "augustus model" shows which GCF 1447 dataset or premade augustus model, was used to run the gene prediction. This table is 1448 Data>Tables>TableS4 models.csv the GitHub found in at repository 1449 https://github.com/extavourlab/Oskar Evolution. 1450

Supplementary Table S5. oskar search results master table. This table summarizes 1451 1452 all results of the oskar search performed on each dataset. Each row corresponds to a 1453 dataset. Columns are as follows: Id: the dataset NCBI identifier; Species: the organism's species name; Family name: the organism's family name; Order name: the organism's 1454 1455 order name; Hits: the number of sequences in the dataset found that satisfy our criteria 1456 for oskar orthology: Source: the NCBI database from which this dataset was downloaded: 1457 Filtered\_hits: the number of oskar sequins in remaining the dataset after the filtration process was applied to all identified oskar sequences. For more information on the criteria 1458 1459 used for oskar orthology and the filtration process, please see the Materials and Methods "Identification of This table found 1460 oskar orthologs". is in 1461 Data>Tables>TableS5 models.csv GitHub the repository at https://github.com/extavourlab/Oskar Evolution. 1462

### 1464 Supplementary Figure Legends

1465 1466 Supplementary Figure S1: Summary statistics of the search for oskar orthologs.

1467 (a) Summary of searches and results for each of the three sources of data searched, from left to right: (i) The total number of datasets searched from all three sources (TSA: 1468 Transcriptome Shotgun Assembly Database; GCA: GenBank; GCF: RefSeq); (ii) the 1469 1470 number of filtered oskar sequences identified in each of those datasets; and (iii) the 1471 proportion of filtered oskar sequences identified in each of the three sources. (b) Summary statistics broken down by insect orders. Only orders where an oskar sequence 1472 was identified are shown. From left to right: (iv) The number of oskar sequences identified 1473 1474 in each of the three data sources; (v) the total number of filtered oskar sequences 1475 identified per order; (vi) the proportion of all searched datasets per order where an oskar 1476 sequences was identified. See also Supplementary Table 1

1477

1478 Supplementary Figure S2: Genome and transcriptome quality correlation to oskar

identification. Shown are box plots of the distribution of *oskar* orthologs identified
(ortholog identified or not identified) with respect to multiple genome and transcriptome
quality metrics. For each metric, the means of both distributions were tested for significant
differences using a Mann Whitney U test. A bar with an \* is displayed if the p-value was
less than 0.05. Mean and median values presented in Supplementary Table S2.

1484

Supplementary Figure S3: Evidence for loss of oskar in Lepidoptera. Phylogeny of 1485 1486 the Lepidoptera as per (Kawahara, et al. 2019). Next to each lepidopteran family are shown summary data regarding the status of oskar identification in our searches. Symbols 1487 with column labels in order from left to right: (i) vertical rectangles: grey: no oskar ortholog 1488 1489 was identified in this family; range: at least one oskar ortholog was identified in this order. (ii) number of datasets searched. (iii) horizontal rectangles: proportion of searched 1490 datasets in which an oskar ortholog was identified; colors as in (i); numbers and 1491 1492 proportions at right. (iv) pie chart: proportion of oskar sequences identified in RefSeq 1493 (GCF) datasets: numbers and proportions at right. (v) pie chart: proportion of oskar 1494 sequences identified in GenBank (GCA) datasets; numbers and proportions at right. (vi) 1495 pie chart: proportion of oskar sequences identified in Transcriptome Shotgun Assembly 1496 Database (TSA) datasets; numbers and proportions at right. Circles to the right of some family names indicate that there is literature evidence for involvement of germ plasm 1497 (black) or no germ plasm (white) in germ cell specification. Numbers to the left of the 1498 circles indicate references to the primary literature as follows: [1]: (Kobayashi and Ando 1499 1984); [2] (Ando and Tanaka 1979); [3] (Lautenschlager 1932); [4] (Anderson and Wood 1500 1968); [5] (Tanaka 1987); [6] (Woodworth 1889); [7] (Eastham 1930); [8] (Berg and 1501 1502 Gassner 1978); [9-10] (Sehl 1931; Guelin 1994); [11] (Johannsen 1929); [12] (Presser and Rutschky 1957); [13-23]: (Tomaya 1902; Schwangart 1905; Saito 1937; Miya 1953, 1503 1958, 1975; Nakao 1999; Toshiki, et al. 2000; Nakao, et al. 2006; Nakao, et al. 2008; 1504 1505 Nakao and Takasu 2019). No datasets were available for Urudidea, Sesidea, Alucitidea, 1506 Callidulidea, Mimallonidea, Drepanidea or Lasiocampidea at the time of analysis. 1507

1508 **Supplementary Figure S4: Evidence for duplication of** *oskar* **in Hymenoptera.** 1509 Phylogenetic tree of all hymenopteran Oskar sequences inferred using RaxML with 100 bootstraps. Branch length normalized to show only the topology. Each leaf is an Oskar
ortholog. Gray: only one Oskar sequence was identified in this species. Red: putatively
duplicated Oskar sequences (sequence similarity < 80%; see Methods). Families</li>
containing *oskar* duplications are highlighted as per Figure 4.

1514

Supplementary Figure S5: Tissue and developmental stage metadata analysis of oskar identification in transcriptome datasets. (a) Proportion of analyzed datasets that were sequenced from the developmental stages indicated on the Y axis. (b) Proportion of analyzed datasets per developmental stage in which an oskar ortholog was identified (red). (c) Proportion of analyzed datasets that were sequenced from the tissue type indicated on the Y axis. (d) Proportion of analyzed datasets per tissue type in which an oskar ortholog was identified (red)

1522

Supplementary Figure S6: Multiple Correspondence Analysis (MCA) of full-length
Oskar, the OSK domain and the LOTUS domain. MCA analysis of trimmed (30%
occupancy) alignments for (a) full-length Oskar, (b) the OSK domain and (c) the LOTUS
domain colored by insect order (see legend at right). The alignment was projected onto
the first three main MCA dimensions (1, 2 and 3). Each dot corresponds to one sequence.
Dotted line outlines specific families of interest as discussed in the text

- 1529
- Supplementary Figure S7: Evolution of the structure of Oskar in Diptera. Left: 1530 dipteran phylogeny from (Maddison, et al. 2007; Wiegmann, et al. 2011). Top: schematic 1531 representation of Oskar domain structure. Blue: heatmap showing the overall occupancy 1532 1533 of an amino acid position in the Oskar alignment trimmed for at least 10% overall 1534 occupancy at a given position. For each dipteran family, occupancy at a given position is 1535 defined as (number of non-gap amino acids / number of sequences in that family). If a 3' or 5' extension (defined as a coding sequence unbroken by stop codons, 5' of the first 1536 residue of the LOTUS domain. or 3' of the last residues of the OSK domain but 5' to a 1537 1538 predicted poly-A tail) was detected in a family, a black box outlines the putative domain. Any such identified 5' domains were designated as putative "Long Oskar" domains. 1539 1540

1541 **Supplementary Figure S8: Oskar domains secondary structure conservation.** 1542 Sequence Logo of Jpred4 predictions for LOTUS and OSK domains showing the 1543 conservation of secondary structures, computed with WebLogo (Crooks, et al. 2004). The 1544 height of each letter represents that state's (X, H or B) conservation throughout the 1545 alignment in bits. X (black): unfolded amino acids; H (red): α helices; E (blue): β sheets. 1546 **(a)** Prediction for the LOTUS domain. **(b)** Prediction for the OSK domain. 1547

1548 Supplementary Figure S9. Duplications and losses of oskar in Hymenoptera. Absence (magenta) or presence of oskar orthologs detected in single copy (cyan) or multiple 1549 1550 copies (yellow) in the genomic or transcriptomic datasets examined in this study. Genera 1551 shown in italics indicate individual species searched and are abbreviated simply for space 1552 reasons. Genera shown in regular type (not italics) indicate a summary of the results from multiple congeneric species, which were nearly always consistent within genera; in all 1553 1554 cases where intrageneric results for oskar presence or absence were inconsistent, we 1555 gave precedence for the finding obtained from a genome sequence (GCF or GCA) over

findings obtained from a transcriptome (TSA)), o for Hymenoptera species or genera. For 1556 1557 some species, germ cell specification via germ plasm (black circles) or differentiation from mesoderm (no germ plasm; white circles) has been reported in the literature, with primary 1558 1559 data references indicated by numbers as follows: [1-6]: (Bütschli 1870; Fleig and Sander 1985, 1986; Zissler 1992; Gutzeit, et al. 1993; Dearden 2006); [7]: (Khila and Abouheif 1560 2008); [8-10]: (Bull 1982; Lynch and Desplan 2010; Lynch, et al. 2011); [11]: (Koscielska 1561 and Koscielski 1987); [12-13]: (Silvestri 1906, 1908); [12, 14-21]: (Silvestri 1906; Hegner 1562 1563 1914; Grbic', et al. 1996; Strand and Grbic' 1997; Grbic' 2000, 2003; Donnell, et al. 2004; Zhurov, et al. 2004); [22-25]: (Gatenby 1917a; Gatenby 1917b; Gatenby 1918; Gatenby 1564 1920); [24]: (Amy 1961); [27-28]: (Gatenby 1920; Tawfik 1957); [29-30]: (Bronskill 1959; 1565 Fleischmann 1975); [31]: (Shafiq 1954); [32]: (Sumitani, et al. 2003). Phylogenetic 1566 relationships as per (Nyman, et al. 2006; Field, et al. 2011; Schmidt 2013; Prous, et al. 1567 2014; Ward 2014; Malm and Nyman 2015; Vilhelmsen 2015; Ward, et al. 2016; Peters, 1568 et al. 2017; Chen and Achterberg 2018; Peters, et al. 2018; Sharanowski, et al. 2021). 1569 1570 Evolution of major hymenopteran life history characteristics (eusociality, pollen collecting, stinger, parasitoidism) as per (Peters, et al. 2017). 1571

#### **Supplementary Figure 1** 1573 1574 (ii) Number of oskar sequences identified (i) Number of datasets searched (iii) Proportion of oskar sequences identified а secuences Number of datasets 1123 # of oskar sequence 238 57.89 1000 200 30.04 500 100 % of oskar 309 20 94 21.19 133 78 0 0 0 TSA GCA GCF TSA GCA GCF TSA GCA GCF (v) Number of filtered oskar (iv) Number of oskar sequences (vi) Proportion of datasets b identified per dataset type sequences identified containing oskar 1 Zygentoma 1 25.0 Ephemeroptera 1 14.29 3 Plecoptera 27.27 3 Orthoptera 6.45 2 2 Phasmatodea 2 4.55 1 Blattodea 8 14.29 8 Thysanoptera 10 76.92 5 Psocoptera 5 21.74 128 Hymenoptera 190 47.5 30 32 1 1 Neuroptera 14.29 14 1 Coleoptera 15.89 17 2 Trichoptera 3 30.0 4 Lepidoptera 1.72 4 2 source Mecoptera 2 50.0 TSA GCA Diptera GCF 161 50.31 0 50 100 0 100 150 200 20 40 60 80 100 50 0 1575

# 1577 Supplementary Figure 2

1578



# 1582 Supplementary Figure 3

1583

			i	ii	iii	i	v	`	V	1	vi
	Agathiphagidae			3	0   0%						0   0%
<u>۲</u> ـــــ	Micropterigidae 1	C	)	1	0   0%					ĕ	0   0%
	Heterobathmiidae	<u> </u>		1	0   0%					Ö	0   0%
	Neopseustidae			1	0   0%						0 0%
ЧГ	<ul> <li>Eriocraniidae</li> </ul>	-		3	0   0%						0 0%
	Hepialidae 2	C	)	3	0   0%						0 0%
Ч.Н	Acanthopteroctetidae			3	0   0%						0 0%
	Lophocoronidae			1	0   0%						0 0%
Ч [	Carposinidae			3	0   0%						0   0%
	<ul> <li>Opostegidae</li> <li>Nantiaulidae</li> </ul>			1	0   0%						0   0%
	<ul> <li>Nepliculidae</li> <li>Prodovidao</li> </ul>			1	010%						010%
				4	010%						010%
	Adelidae			1	2   100%						21100%
Ч	Tischeriidae			2	010%						010%
	Palaephatidae			9	2   22%					ě	2   22%
41	Meessiidae			2	0   0%					ŏ	0   0%
	Psychidae 3	C	)	3	0   0%				0 0%	ŏ	0   0%
4	Tineidae	<u> </u>		2	0   0%					Ö	0 0%
٦	Dryadaulidae			2	0   0%						0 0%
	Plutellidae			5	0   0%		0   0%				0 0%
	<ul> <li>Yponomeutidae</li> </ul>			1	0   0%						0 0%
Чг	Urudidea	~						-		-	
	Tortricidae 4	C	)	11	0   0%				0 0%		0 0%
41 7	Casthildae			1	0   0%					•	0   0%
	Sesidea				0.1.00					-	0.1.00
ЧП	Cossidao				010%						010%
	Panilionidae 5	0		7	010%		010%		010%		010%
	- Nymnhalidae 6	X	( –	60	010%		010%		010%		010%
ЧдГ	Pieridae 7	~	5	5	010%		010%		010%		010%
4	Riodinidae	$\cup$	1	2	0   0%			- ŏ	0   0%	_	
	Lycaenidae			2	0   0%			ŏ	0   0%		0 0%
	Hesperiidae			3	0   0%			Ō	0   0%		
ــــــــــــــــــــــــــــــــــــــ	<ul> <li>Cosmopterigidae</li> </ul>			1	0   0%		0   0%				
	- Gelechiidae 8			2	0   0%				0   0%		0 0%
[	- Alucitidea										
	Callidulidea					_		_		_	
	Crambidae			10	0   0%		0   0%		0   0%		0   0%
	Pyralidae 9-1	0 (	)	0	0   0%		0   0%		0   0%	•	0 0%
	- Mimalionidea										
4[	- Drepanidea - Notodontidao			3	01.0%						010%
Ц	Frebidae 11			5	010%				0   0%		0   0%
4	- Lymantriidae			3	0   0%			ă	0   0%	ě	0   0%
Ц Ц	Noctuidae 12	$\sim$		41	0   0%		0   0%	ē	0   0%	ŏ	0   0%
<sub>[</sub>	Geometridae			2	0   0%	_		ē	0   0%	õ	0 0%
Чг	- Lasiocampidea	-									
Ч г—	Bombycidae 13-2	23 ()	)	5	0   0%		0   0%		0   0%		0 0%
	- Saturniidae	_		7	0   0%						0 0%
	- Sphingidae			3	0   0%		0   0%				0 0%

# 1586 Supplementary Figure 4





#### 1590 Supplementary Figure 5







Order

Zygentoma

Plecoptera

Orthoptera

Blattodea

Psocoptera

Neuroptera

Coleoptera

Trichoptera

Lepidoptera

Mecoptera

Diptera

Phasmatodea

Thysanoptera

Hymenoptera

Ephemeroptera

# 1594 Supplementary Figure 6

1595

### a. Full-length Oskar



# 1598 Supplementary Figure 7

1599



# 1602 Supplementary Figure 8





#### 1607 Supplementary Table S1

Insect Order	Source	Number of datasets searched	Total hits	Filtered hits	% of datasets with os <i>kar</i> identified
Archaeognatha	GCA	1	0	0	0
Archaeognatha	TSA	2	0	0	0
Blattodea	GCA	3	1	1	33.33
Blattodea	GCF	2	0	0	0
Blattodea	TSA	51	7	7	13.73
Coleoptera	GCA	12	1	1	8.33
Coleoptera	GCF	9	3	2	22.22
Coleoptera	TSA	86	31	14	16.28
Collembola	TSA	9	0	0	0
Dermaptera	TSA	7	0	0	0
Diptera	GCA	115	63	60	52.17
Diptera	GCF	43	58	43	100
Diptera	TSA	162	72	58	35.8
Embioptera	TSA	5	0	0	0
Ephemeroptera	GCA	2	0	0	0
Ephemeroptera	TSA	5	1	1	20
Grylloblattodea	TSA	2	0	0	0
Hemiptera	GCA	18	0	0	0
Hemiptera	GCF	12	0	0	0
Hemiptera	TSA	192	1	0	0
Hymenoptera	GCA	52	32	30	57.69
Hymenoptera	GCF	47	36	32	68.09
Hymenoptera	TSA	301	157	128	42.52
Lepidoptera	GCA	80	0	0	0
Lepidoptera	GCF	17	0	0	0
Lepidoptera	TSA	135	24	4	2.96

Mantodea	TSA	13	0	0	0
Mantophasmatodea	TSA	2	0	0	0
Mecoptera	TSA	4	2	2	50
Megaloptera	TSA	3	0	0	0
Neuroptera	TSA	7	1	1	14.29
Odonata	GCA	2	0	0	0
Odonata	TSA	7	0	0	0
Orthoptera	GCA	3	0	0	0
Orthoptera	TSA	28	2	2	7.14
Phasmatodea	GCA	13	0	0	0
Phasmatodea	TSA	31	6	2	6.45
Phthiraptera	GCF	1	0	0	0
Phthiraptera	TSA	7	0	0	0
Plecoptera	GCA	3	0	0	0
Plecoptera	TSA	8	3	3	37.5
Psocoptera	TSA	23	5	5	21.74
Raphidioptera	TSA	3	0	0	0
Siphonaptera	GCF	1	0	0	0
Siphonaptera	TSA	4	0	0	0
Strepsiptera	GCA	1	0	0	0
Strepsiptera	TSA	2	0	0	0
Thysanoptera	GCA	1	1	1	100
Thysanoptera	GCF	1	1	1	100
Thysanoptera	TSA	11	10	8	72.73
Trichoptera	GCA	3	1	1	33.33
Trichoptera	TSA	7	2	2	28.57
Zoraptera	TSA	2	0	0	0
Zygentoma	TSA	4	1	1	25
Crustacea	TSA	168	0	0	0

Crustacea	GCF	1	0	0	0
Crustacea	GCA	11	0	0	0

1608 1609 1610

Blondel et al. Page 71 of 72

#### 1611 **Supplementary Table S2**

#### 1612

	Genome parameter	(a) <i>oskar</i> Identified	(b) <i>oskar</i> not identified	ratio (a):(b)
# contins	mean	255,015	43,280	5.89
	median	69,255	20,653	3.35
# apoffoldo	mean	182,706	23,596	7.74
# scanolos	median	40,960	9,398	4.36
contig NEO (bp)	mean	324,036	726,696	0.45
	median	14,052	40,079	0.35
apoffold NEO	mean	2,636,825	5,695,299	0.46
Scallolu NSU	median	96,730	385,460	0.25
contig LEO	mean	40,955	3,701	11.07
	median	6,868	1,300	5.28
apoffold   EQ	mean	27,269	1,500	18.18
Scanolu Lou	median	1,131	191	59.53
# contige per geneme length	mean	0.00060	0.00017	3.53
# contigs per genome length	median	0.00021	0.00009	2.33
# scaffolds per genome	mean	0.00045	0.00009	5.00
length	median	0.00013	0.00005	2.60

1613