1 **Title: Bacterial contribution to genesis of the novel germ line determinant *oskar***

2

3 **Authors:** *Leo Blondel[1], Tamsin E. M. Jones[2,3] and Cassandra G. Extavour[1,2*]*

4

5 **Affiliations:**

6 1. Department of Molecular and Cellular Biology, Harvard University, 16 Divinity Avenue,

7    Cambridge MA, USA

8 2. Department of Organismic and Evolutionary Biology, Harvard University, 16 Divinity

9    Avenue, Cambridge MA, USA

10 3. Current address: European Bioinformatics Institute, EMBL-EBI, Wellcome Genome

11    Campus, Hinxton, Cambridgeshire, UK

12

13 * Correspondence to extavour@oeb.harvard.edu

14

15 **Abstract:** New cellular functions and developmental processes can evolve by modifying

16 existing genes or creating new genes. New genes can arise not only via duplication or mutation

17 but also by acquiring foreign DNA, also called horizontal gene transfer (HGT). Here we show

18 that HGT likely contributed to the creation of a novel gene indispensable for reproduction in

19 some insects. Long considered a novel gene with unknown origin, *oskar* has evolved to fulfil a

20 crucial role in insect germ cell formation. Our analysis of over 100 Oskar sequences suggests

21 that Oskar arose through a novel gene formation history involving fusion of eukaryotic and

22 prokaryotic sequences. This work shows that highly unusual gene origin processes can birth

23 novel genes that can facilitate evolution of novel developmental mechanisms.

24

25 **Main Text:**

26      **Introduction:** Heritable variation is the raw material of evolutionary change. Genetic

27 variation can arise from mutation and gene duplication of existing genes (*1*), or through *de*

28 *novo* processes (*2*), but the extent to which such novel, or "orphan" genes participate

29 significantly in the evolutionary process is unclear. Mutation of existing cis-regulatory (*3*) or

30 protein coding regions (*4*) can drive evolutionary change in developmental processes.

31 However, recent studies in animals and fungi suggest that new genes can also drive phenotypic

32 change (*5*). Although counterintuitive, novel genes may be integrating continuously into

33 otherwise conserved gene networks, with a higher rate of partner acquisition than subtler

34 variations on preexisting genes (*6*). Moreover, in humans and fruit flies, a large proportion of

35 new genes are expressed in the brain, suggesting their participation in the evolution of major

36 organ systems (*7, 8*). However, while next generation sequencing has improved their

37 discovery, the developmental and evolutionary significance of new genes remains

38 understudied.

39      The mechanism of formation of a new gene may have implications for its function.

40 New genes that arise by duplication, thus possessing the same biophysical properties as their

41 parent genes, have innate potential to participate in preexisting cellular and molecular

42 mechanisms (*1*). However, orphan genes lacking sequence similarity to existing genes must

43 form novel functional molecular relationships with extant genes, in order to persist in the

44 genome. When such genes arise by introduction of foreign DNA into a host genome through

45 horizontal gene transfer (HGT), they may introduce novel, already functional sequence

46 information into a genome. Whether genes created by HGT show a greater propensity to

47 contribute to or enable novel processes is unclear. Endosymbionts in the host germ line

48  cytoplasm (germ line symbionts) could increase the occurrence of evolutionarily relevant HGT

49  events, as foreign DNA integrated into the germ line genome is transferred to the next

50  generation. HGT from bacterial endosymbionts into insect genomes appears widespread,

51  involving transfer of metabolic genes or even larger genomic fragments to the host genome (*9*).

52      Here we examined the evolutionary origins of the *oskar* (*osk*) gene, long considered a

53  novel gene that evolved to be indispensable for insect reproduction (*10*). First discovered in

54  *Drosophila melanogaster* (*11*), *osk* is necessary and sufficient for assembly of germ plasm, a

55  cytoplasmic determinant that specifies the germ line in the embryo. Germ plasm-based germ

56  line specification appears derived within insects, confined to insects that undergo

57  metamorphosis (Holometabola) (*12, 13*). Initially thought exclusive to Diptera (flies and

58  mosquitoes), its discovery in a wasp, another holometabolous insect with germ plasm (*14*), led

59  to the hypothesis that *oskar* originated as a novel gene at the base of the Holometabola

60  approximately 300 Mya, facilitating the evolution of insect germ plasm as a novel

61  developmental mechanism (*14*). However, its subsequent discovery in a cricket (*12*), a basally

62  branching insect without germ plasm (*15*), implied that *osk* was instead at least 50 My older,

63  and that its germ plasm role was derived rather than ancestral (*16*). Despite its orphan gene

64  status, *osk* plays major developmental roles, interacting with the products of many genes highly

65  conserved across animals (*10, 17, 18*). *osk* thus represents an example of a new gene that not

66  only functions within pre-existing gene networks in the nervous system (*12*), but has also

67  evolved into the only animal gene known to be both necessary and sufficient for germ line

68  specification (*19, 20*).

69      The evolutionary origins of this remarkable gene are unknown. Osk contains two

70  biophysically conserved domains, an N-terminal LOTUS domain and a C-terminal hydrolase-

71  like domain called OSK (*17, 21*) (Fig. 1a). A BLASTp search using the full-length *D.*

72  *melanogaster osk* sequence as a query yielded only other holometabolous *osk* genes (E-value <

73  0.01), or hits for the LOTUS or OSK domains (E-value <10) (Supplementary files: BLAST

74  search results). This suggested that full length *osk* was unlikely to be a duplication of any other

75  known gene, prompting us to perform a BLASTp search on each conserved Osk protein

76  domain individually. Strikingly, in our BLASTp search, we recovered no eukaryotic sequences

77  that resembled the OSK domain (E-value < 10) (Supplementary files: BLAST search results).

78          **Results:** To understand this anomaly, we built an alignment of 95 Oskar sequences

79  (Supplementary files: Alignments>OSKAR_FINAL.fasta) and used a custom iterative

80  HMMER sliding window search tool to compare each domain with protein sequences from all

81  domains of life. Sequences most similar to the LOTUS domain were almost exclusively

82  eukaryotic sequences (Supplementary Table 3). In contrast, those most similar to the OSK

83  domain were bacterial, specifically sequences similar to SGNH-like hydrolases (*17, 21*) (Pfam

84  Clan: SGNH_hydrolase - CL0264; Supp. Table 4; Fig. 1b). To visualize their relationships, we

85  graphed the sequence similarity network for the sequences of these domains and their closest

86  hits. We observed that the majority of LOTUS domain sequences clustered within eukaryotic

87  sequences (Fig. 1c). In contrast, OSK domain sequences formed an isolated cluster, a small

88  subset of which formed a connection to bacterial sequences (Fig. 1d). These data are consistent

89  with a previous suggestion, based on BLAST results (*14*), that HGT from a bacterium into an

90  ancestral insect genome may have contributed to the evolution of *osk*. However, this possibility

91  was not adequately addressed by previous analyses, which were based on alignments of full

92  length Osk containing only eukaryotic sequences as outgroups (*12*). To rigorously test this

93  hypothesis, we therefore performed phylogenetic analyses of the two domains independently.

94 A finding that LOTUS sequences branch within eukaryotes, while OSK sequences branch

95 within bacteria, would provide support for the HGT hypothesis.

96       Both Maximum likelihood and Bayesian approaches confirmed this prediction (Fig. 2).

97 As expected, LOTUS sequences from Osk proteins were related to other eukaryotic LOTUS

98 domains, to the exclusion of the only three bacterial sequences with sufficient similarity to

99 include in the analyses (Figs. 2a, S1, S2; see Methods and Supplemental Text). In contrast,

100 OSK domain sequences branched within bacterial sequences (Fig. 2b, S3, S4). Importantly,

101 OSK sequences did not simply form an outgroup to bacterial sequences. Instead, they formed a

102 well-supported clade nested within bacterial GDSL-like lipase sequences. The majority of

103 these bacterial sequences were from the Firmicutes, a bacterial phylum known to include insect

104 germline symbionts (*22, 23*). All other sequences from classified bacterial species, including a

105 clade branching basally to all other sequences, belonged either to the Bacteroidetes or to the

106 Proteobacteria. Members of both of these phyla are also known germline symbionts of insects

107 (*9, 24*) and other arthropods (*25*). In sum, the distinct phylogenetic relationships of the two

108 domains of Oskar are consistent with a bacterial origin for the OSK domain. Further, the

109 specific bacterial clades close to OSK suggest that an ancient arthropod germ line

110 endosymbiont could have been the source of a GDSL-like sequence that was transferred into

111 an ancestral insect genome, and ultimately gave rise to the OSK domain of *oskar*.

112       We then asked if two additional sequence characteristics, GC3 content and codon use,

113 were consistent with distinct domain of life origins for the two Oskar domains (*26*). Under our

114 hypothesis, the HGT event that contributed to *oskar*'s formation would have occurred at least

115 480 Mya, in a common insect ancestor (*27*). We reasoned that if evolutionary time had not

116 completely erased such signatures from the putative bacterially donated sequence (OSK), we

117  might detect differences from the LOTUS domain, and from the host genome. Thus, we

118  performed a parametric analysis of these parameters for 17 well annotated insect genomes

119  (Supplementary Table 5). To quantify the null hypothesis, we calculated an "Intra-Gene

120  distribution" for all genes in the genome, which showed a linear correlation between codon use

121  in the 5' and 3' halves of a given gene. In contrast, the codon use between the LOTUS and

122  OSK domains did not follow this correlation for nearly all measures of codon use (Fig. 3a, 3b,

123  S5). For each genome, we then calculated the residuals of the Intra-Gene distribution and the

124  LOTUS-OSK pair. Pooling the residuals together revealed that the GC3 content was drastically

125  different between the LOTUS and OSK domains, compared to what would be expected within

126  an average gene in that genome (Fig. 3c). Finally, to quantify the codon use difference, we

127  compared the cosine distance in codon use between the LOTUS and OSK domains, with that

128  of the Inter-Gene and Intra-Gene distributions. We found that the LOTUS-OSK distance was

129  closer to that measured between two different, random genes, than between two parts of the

130  same gene (Inter-Gene and Intra-Gene distributions, respectively; Fig. 3d). In sum, whereas

131  most genes have similar codon use across all regions of their coding sequence, the OSK and

132  LOTUS domains of *oskar* use codons in different ways. Together with the phylogenetic and

133  sequence similarity evidence presented above, these analyses are consistent with an HGT

134  origin for the OSK domain (Fig. 4).

135      **Discussion:** While multiple mechanisms can give rise to new genes, HGT is arguably

136  among the least well understood, as it involves multiple genomes and ancient biotic

137  interactions between donor and host organisms that are often difficult to reconstruct. In the

138  case of *oskar*, however, the fact that both germline symbionts (*28*) and HGT events (*9*) are

139  widespread in insects, provides a plausible biological mechanism consistent with our

140  hypothesis that fusion of eukaryotic and bacterial domain sequences led to the birth of this

141  novel gene.

142      Once arisen, novel genes might be expected to disappear rapidly, given that pre-

143  existing gene regulatory networks operated successfully without them (*1*). However, it is clear

144  that new genes can evolve functional connections with existing networks, become essential

145  (*29*), and in some cases lead to new functions (*30*) and contribute to phenotypic diversity (*5*).

146  *oskar* plays multiple critical roles in insect development, from neural patterning (*12, 31*) to

147  oogenesis (*32*). In the Holometabola, a clade of nearly one million extant species (*33*), *oskar*'s

148  co-option to become necessary and sufficient for germ plasm assembly is likely the cell

149  biological mechanism underlying the evolution of this derived mode of insect germ line

150  specification (*12, 14, 16*). Our study thus provides evidence that HGT can not only introduce

151  functional genes into a host genome, but also, by contributing sequences of individual

152  domains, generate genes with entirely novel domain structures that may facilitate the evolution

153  of novel developmental mechanisms.

154 **References**

155

156 1.     J. S. Taylor, J. Raes, Duplication and divergence: the evolution of new genes and old
157      ideas. *Annual review of genetics* **38**, 615-643 (2004).
158 2.     D. Tautz, T. Domazet-Loso, The evolutionary origin of orphan genes. *Nat. Rev. Genet.*
159      **12**, 692-702 (2011).
160 3.     P. J. Wittkopp, G. Kalay, Cis-regulatory elements: molecular mechanisms and
161      evolutionary processes underlying divergence. *Nat. Rev. Genet.* **13**, 59-69 (2011).
162 4.     H. E. Hoekstra, J. A. Coyne, The locus of evolution: evo devo and the genetics of
163      adaptation. *Evolution* **61**, 995-1016 (2007).
164 5.     S. Chen, B. H. Krinsky, M. Long, New genes as drivers of phenotypic evolution. *Nat.*
165      *Rev. Genet.* **14**, 645-660 (2013).
166 6.     W. Zhang, P. Landback, A. R. Gschwend, B. Shen, M. Long, New genes drive the
167      evolution of gene interaction networks in the human and mouse genomes. *Genome*
168      *Biol.* **16**, 202 (2015).
169 7.     Y. E. Zhang, P. Landback, M. Vibranovski, M. Long, New genes expressed in human
170      brains: implications for annotating evolving genomes. *BioEssays* **34**, 982-991 (2012).
171 8.     S. Chen *et al.*, Frequent recent origination of brain genes shaped the evolution of
172      foraging behavior in Drosophila. *Cell Reports* **1**, 118-132 (2012).
173 9.     J. C. Dunning Hotopp *et al.*, Widespread lateral gene transfer from intracellular bacteria
174      to multicellular eukaryotes. *Science (New York, NY)* **317**, 1753-1756 (2007).
175 10.    R. Lehmann, Germ Plasm Biogenesis--An Oskar-Centric Perspective. *Curr. Top. Dev.*
176      *Biol.* **116**, 679-707 (2016).
177 11.    R. Lehmann, C. Nüsslein-Volhard, Abdominal Segmentation, Pole Cell Formation, and
178      Embryonic Polarity Require the Localized Activity of *oskar*, a Maternal Gene in
179      *Drosophila*. *Cell* **47**, 144-152 (1986).
180 12.    B. Ewen-Campen, J. R. Srouji, E. E. Schwager, C. G. Extavour, *oskar* Predates the
181      Evolution of Germ Plasm in Insects. *Curr. Biol.* **22**, 2278-2283 (2012).
182 13.    C. G. Extavour, M. E. Akam, Mechanisms of germ cell specification across the
183      metazoans: epigenesis and preformation. *Development* **130**, 5869-5884 (2003).
184 14.    J. A. Lynch *et al.*, The Phylogenetic Origin of *oskar* Coincided with the Origin of
185      Maternally Provisioned Germ Plasm and Pole Cells at the Base of the Holometabola.
186      *PLoS Genetics* **7**, e1002029 (2011).
187 15.    B. Ewen-Campen, S. Donoughe, D. N. Clarke, C. G. Extavour, Germ cell specification
188      requires zygotic mechanisms rather than germ plasm in a basally branching insect.
189      *Curr. Biol.* **23**, 835-842 (2013).
190 16.    E. Abouheif, Evolution: oskar Reveals Missing Link in Co-optive Evolution. *Curr.*
191      *Biol.* **23**, R24-R25 (2012).
192 17.    M. Jeske *et al.*, The Crystal Structure of the Drosophila Germline Inducer Oskar
193      Identifies Two Domains with Distinct Vasa Helicase- and RNA-Binding Activities.
194      *Cell Reports* **12**, 587-598 (2015).
195 18.    M. Jeske, C. W. Muller, A. Ephrussi, The LOTUS domain is a conserved DEAD-box
196      RNA helicase regulator essential for the recruitment of Vasa to the germ plasm and
197      nuage. *Genes Dev.* **31**, 939-952 (2017).

198 19.    A. Ephrussi, R. Lehmann, Induction of germ cell formation by *oskar*. *Nature* **358**, 387-
199     392 (1992).

200 20.    J. Kim-Ha, J. L. Smith, P. M. Macdonald, *oskar* mRNA is localized to the posterior
201     pole of the *Drosophila* oocyte. *Cell* **66**, 23-35 (1991).

202 21.    N. Yang *et al.*, Structure of Drosophila Oskar reveals a novel RNA binding protein.
203     *Proc. Natl. Acad. Sci. USA* **112**, 11541-11546 (2015).

204 22.    D. Wheeler, A. J. Redding, J. H. Werren, Characterization of an ancient lepidopteran
205     lateral gene transfer. *PLoS ONE* **8**, e59262 (2013).

206 23.    S. T. Chepkemoi *et al.*, Identification of Spiroplasmainsolitum symbionts in Anopheles
207     gambiae. *Wellcome Open Research* **2**, 90 (2017).

208 24.    E. Zchori-Fein, S. J. Perlman, S. E. Kelly, N. Katzir, M. S. Hunter, Characterization of
209     a 'Bacteroidetes' symbiont in Encarsia wasps (Hymenoptera: Aphelinidae): proposal of
210     'Candidatus Cardinium hertigii'. *International Journal of Systematic and Evolutionary*
211     *Microbiology* **54**, 961-968 (2004).

212 25.    E. Zchori-Fein, S. J. Perlman, Distribution of the bacterial symbiont Cardinium in
213     arthropods. *Mol. Ecol.* **13**, 2009-2016 (2004).

214 26.    T. Tuller, Codon bias, tRNA pools and horizontal gene transfer. *Mobile Genetic*
215     *Elements* **1**, 75-77 (2011).

216 27.    B. Misof *et al.*, Phylogenomics resolves the timing and pattern of insect evolution.
217     *Science* **346**, 763-767 (2014).

218 28.    K. Bourtzis, T. A. Miller, Eds., *Insect Symbiosis*, (CRC Press, Boca Raton, FL, 2006),
219     vol. 3, pp. 304.

220 29.    S. Chen, Y. E. Zhang, M. Long, New genes in *Drosophila* quickly become essential.
221     *Science* **330**, 1682-1685 (2010).

222 30.    G. Cornelis *et al.*, Ancestral capture of syncytin-Car1, a fusogenic endogenous
223     retroviral envelope gene involved in placentation and conserved in Carnivora. *Proc.*
224     *Natl. Acad. Sci. USA* **109**, E432-441 (2012).

225 31.    X. Xu, J. L. Brechbiel, E. R. Gavis, Dynein-Dependent Transport of nanos RNA in
226     Drosophila Sensory Neurons Requires Rumpelstiltskin and the Germ Plasm Organizer
227     Oskar. *The Journal of Neuroscience* **33**, 14791-14800 (2013).

228 32.    A. Jenny *et al.*, A translation-independent role of oskar RNA in early Drosophila
229     oogenesis. *Development* **133**, 2827-2833 (2006).

230 33.    J. A. Rees, K. Cranston, Automated assembly of a reference taxonomy for phylogenetic
231     data synthesis. *Biodiversity Data Journal* **5**, e12581 (2017).

232 34.    J. A. Gerlt *et al.*, Enzyme Function Initiative-Enzyme Similarity Tool (EFI-EST): A
233     web tool for generating protein sequence similarity networks. *Biochimica et Biophysica*
234     *Acta (BBA) - Proteins and Proteomics* **1854**, 1019-1037 (2015).

235

241 **Author contributions:** CGME conceived of the project and overall experimental design.

242 TEMJ collected initial transcriptome datasets and identified *oskar* orthologues therein. LB built

243 the HMM model, identified additional orthologues, and performed sequence, phylogenetic,

244 cluster and codon use analyses. LB and CGME interpreted data and wrote the manuscript.

245

246 **Competing interests:** The authors declare no competing interests.

247

248 **Data and materials availability:** All data is available in the main text or the supplementary

249 materials.

250

251

252  **Supplementary Materials:**

253  The Supplementary Information for this paper consists of the following elements:

254  Supplementary figures

255  • Figure S1: LOTUS Domain RaxML Tree.

256  • Figure S2: LOTUS Domain Bayesian Tree.

257  • Figure S3: OSK Domain RaxML Tree.

258  • Figure S4: OSK Domain Bayesian Tree.

259  • Figure S5: AT3/GC3 correlations between the LOTUS and OSK domains.

260  • Figure S6: A3/T3/G3/C3 correlations between the LOTUS and OSK domains.

261  Supplementary tables

262  • Table S1: List of genomes and transcriptomes used for automated *oskar* search.

263  • Table S2: List of *oskar* sequences used in the final alignment.

264  • Table S3: List of sequences used for phylogenetic analysis of the LOTUS domain.

265  • Table S4: List of sequences used for phylogenetic analysis of the OSK domain.

266  • Table S5: List of genomes analyzed for codon use.

267

268  1. Supplementary Discussion

269  (Blondel_Jones_Extavour_HGT_HGT_Paper_SuppInfo_V4_181108.docx)

270  2. Supplementary References

271  (Blondel_Jones_Extavour_HGT_HGT_Paper_SuppInfo_V4_181108.docx)

272  3. Folder titled "Supplementary Information Files" containing the following sub-folders

273     a. Supplementary Information Files>Alignments

274         i. *All sequences identified and analyzed in this study, in FASTA format and*

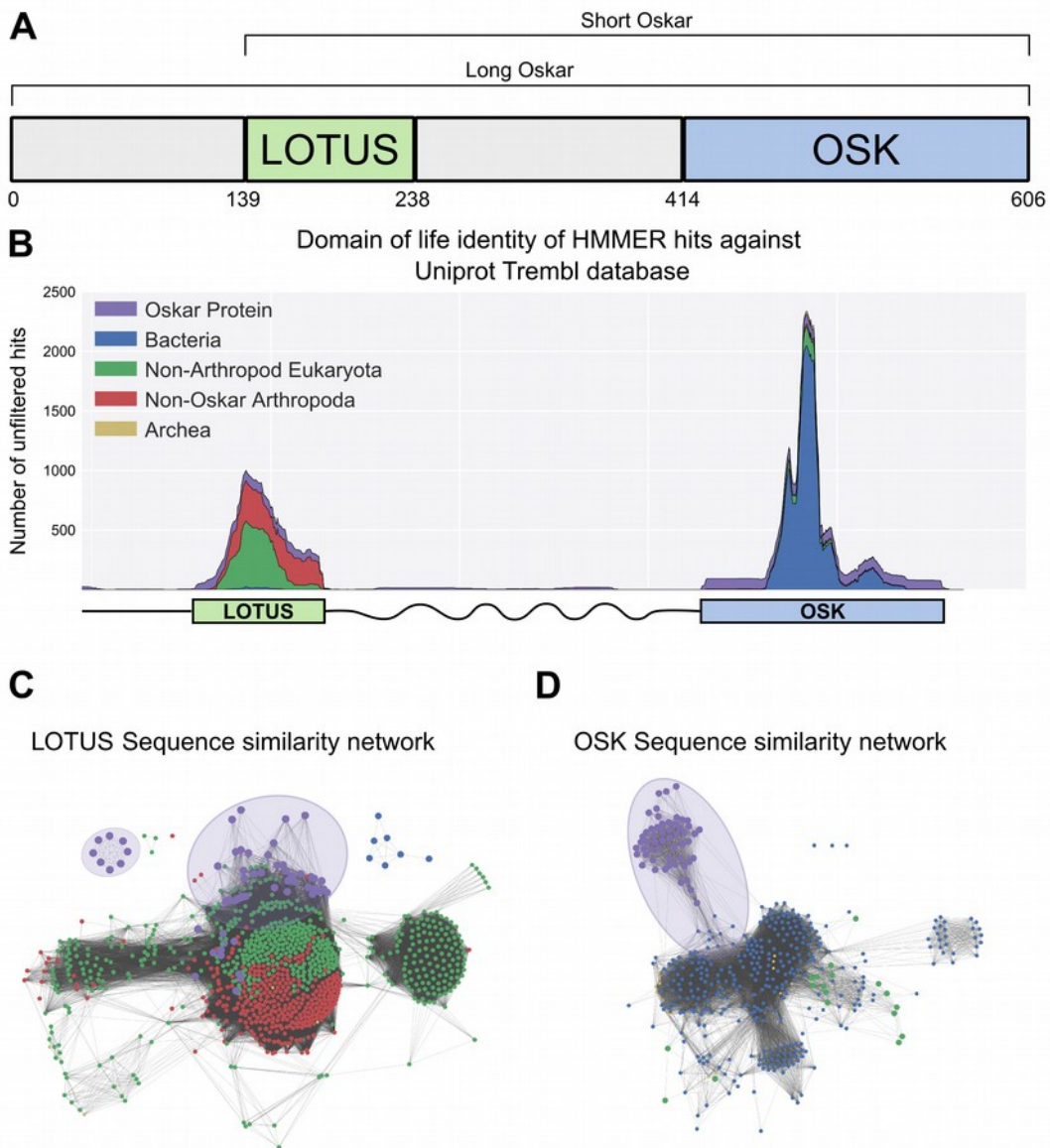275            *with corresponding Alignments*

276     b. Supplementary Information Files>BLAST search results

277         i. *Results of BLASTP searches with full length Oskar, OSK or LOTUS*

278            *domains as queries*

279     c. Supplementary Information Files>Data

280    *i. Necessary files for running the different ipython notebooks:*

281     1. *Taxonomy: Conversion table for UniProt ID to taxon information.*
282      *(uniprot_ID_taxa.tsv )*

283     2. *Codon_Genes: Contains the measured codon frequency for the*
284      *different genomes studied as .csv or .tsv files (organism_name.csv/*
285      *tsv), along with the DNA sequences of LOTUS and OSK domains*
286      *used in the codon use analysis (LOTUS_Seqeuences.gb and*
287      *SGNH_Seqeuences.gb)*

288     3. *Trees: Contains the tree files obtained from RaxML and MrBayes*
289      *phylogenetic analyses of the OSK and LOTUS domains.*

290  d. Supplementary Information Files>HMM

291    *i. HMM models used for iterative searching for sequences similar to full-*
292    *length Oskar, LOTUS and OSK domains*

293  e. Supplementary Information Files>Scripts

294    *i. All custom scripts used to implement the analysis pipelines described.*

295 2. Supplementary Information Files>Tables

296  **a.** Supplementary Tables S1-S5 describing databases searched/analyzed and all
297   search results; Legends in
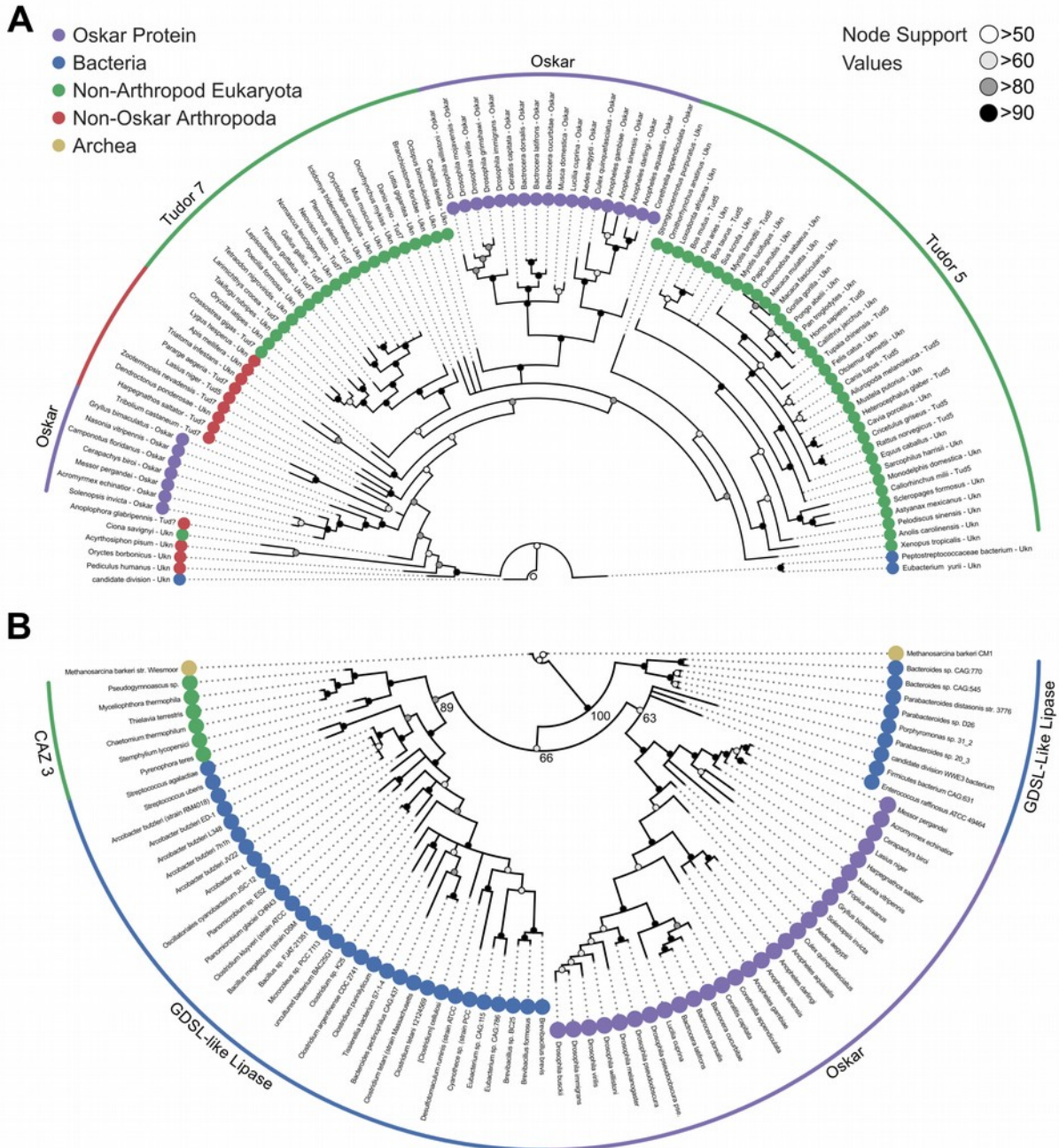298   Blondel_Jones_Extavour_HGT_HGT_Paper_SuppInfo_V4_181108.docx

Figure 1
Blondel, Jones & Extavour



**A**

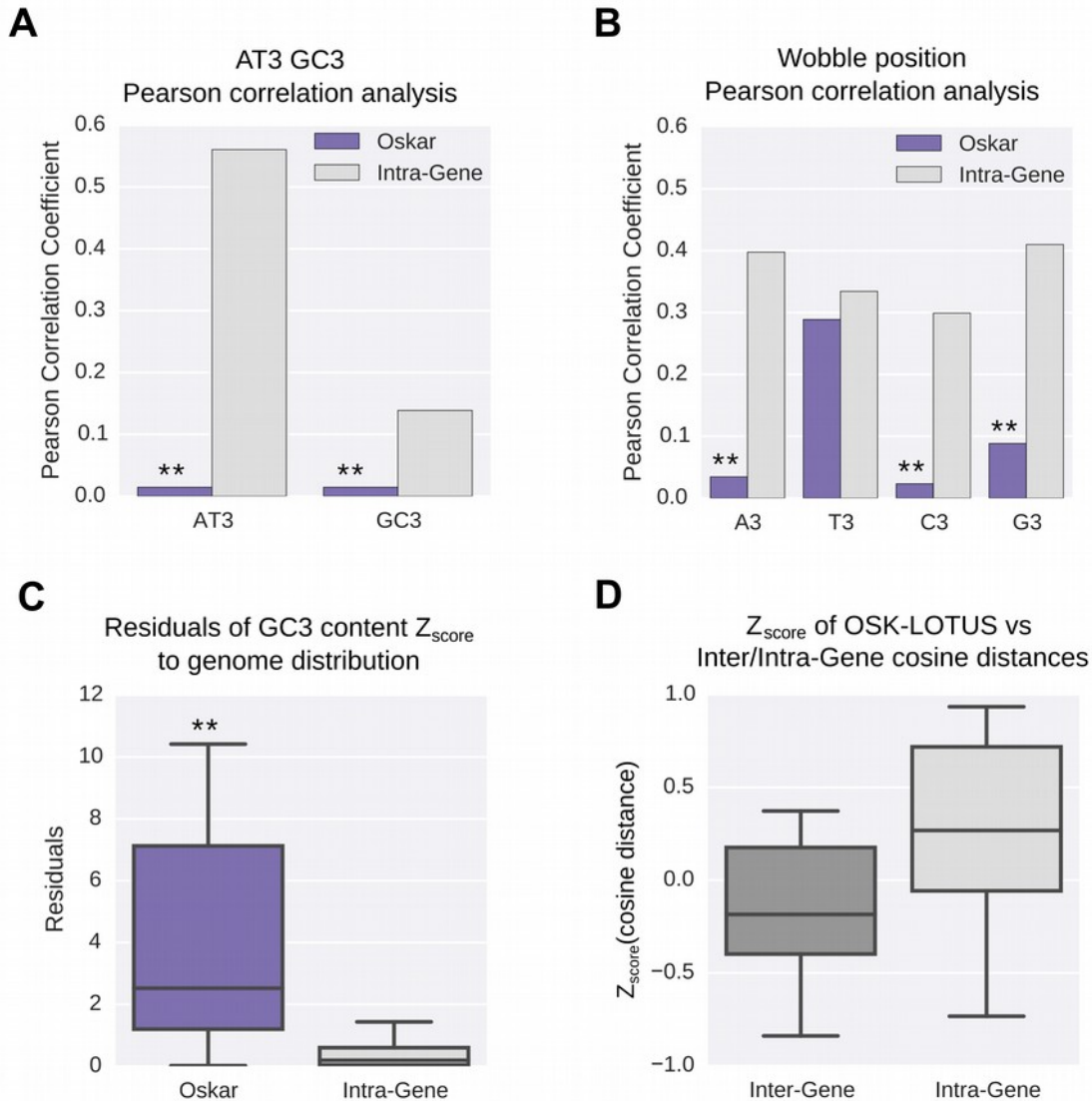Short Oskar

Long Oskar

| | LOTUS | | OSK |

0    139    238    414    606

**B**

Domain of life identity of HMMER hits against
Uniprot Trembl database

Oskar Protein
Bacteria
Non-Arthropod Eukaryota
Non-Oskar Arthropoda
Archea

LOTUS    OSK

**C** LOTUS Sequence similarity network

**D** OSK Sequence similarity network

299

**Figure 1**. **Sequence analysis of the Oskar gene. a**, Schematic representation of the Oskar gene. The LOTUS and OSK hydrolase-like domains are separated by a poorly conserved region of predicted high disorder and variable length between species. In some dipterans, a region 3' to the LOTUS domain is translated to yield a second isoform, called Long Oskar. Residue numbers correspond to the *D. melanogaster* Osk sequence. **b**, Stackplot of domain of life identity of HMMER hits across the protein sequence. For a sliding window of 60 Amino Acids across the protein sequence (X axis), the number of hits in the Trembl (UniProt) database (Y axis) is represented and color coded by domain of life origin (see Methods: Iterative HMMER search of OSK and LOTUS domains), stacked on top of each other. **c, d** EFI-EST[34]-generated graphs of the sequence similarity network of the LOTUS (**c**) and OSK (**d**) domains of Oskar. Sequences were obtained using HMMER against the UniProtKB database. Most Oskar LOTUS sequences cluster within eukaryotes and arthropods. In contrast, Oskar OSK sequences cluster most strongly with a small subset of bacterial sequences.

Figure 2
Blondel, Jones & Extavour

311

**Figure 2**. **Phylogenetic analysis of the LOTUS and OSK domains. a**, Bayesian consensus tree for the LOTUS domain. Three major LOTUS-containing protein families are represented within the tree: Tudor 5, Tudor 7, and Oskar. Oskar LOTUS domains form two clades, one containing only dipterans and one containing all other represented insects (hymenopterans and orthopterans). The tree was rooted to the three bacterial sequences added in the dataset. **b**, Bayesian consensus tree for the OSK domain. The OSK domain is nested within GDSL-like domains of bacterial species from phyla known to contain germ line symbionts in insects. The ten non-Oskar eukaryotic sequences in the analysis form one clade comprising fungal Carbohydrate Active Enzyme 3 (CAZ3) proteins. For Bayesian and RaxML trees with all accession numbers and node support values see Extended Data Figures S1-4.
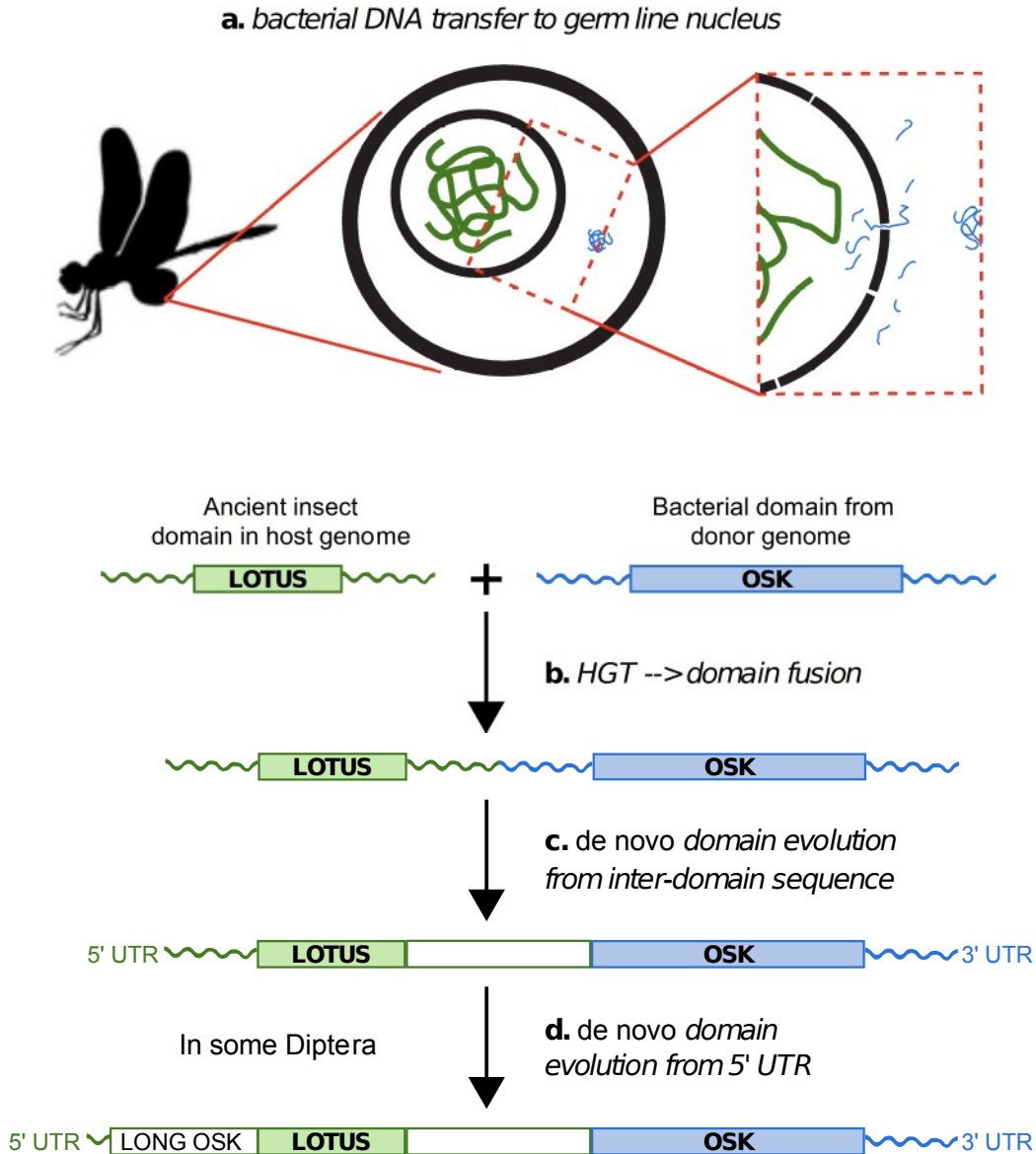
Figure 3
Blondel, Jones & Extavour

**Figure 3. Parametric analysis of codon use for the LOTUS and OSK domains. a**, Pearson correlation analysis
of AT3 and GC3 content for Oskar vs other genes. AT3 and GC3 content are correlated across the sequence of a
gene for all genes in a given genome (grey), but not between the LOTUS and OSK domains of Oskar (purple).
(**: Pearson correlation p-value > 0.1) **b**, Pearson correlation analysis of wobble position identity for the Oskar
gene vs other genes. Wobble position identity content is correlated across the sequence of a gene for all genes in a
given genome (grey) but not between the LOTUS and OSK domains of Oskar (purple), with the exception of T3.
(**: Pearson correlation p-value > 0.1) **c**, Analysis of GC3 content. Measure of the residuals of Z scores for Oskar
gene GC3 content (LOTUS vs OSK) and the Intra-Gene GC3 content. The GC3 content of the LOTUS and OSK
domains does not follow a linear relationship, and the residuals are significantly higher (purple) than those
observed within across the sequences of other genes within a given genome (grey). (** : Mann-Whitney U test p-
value < $10^{-5}$) **d**, Cosine distance analysis of codon frequencies. The distance distribution in codon use between the
LOTUS and OSK domain is less than the measured null distribution distance in codon use between any two
unrelated genes (Inter-Gene; dark grey), but greater than the expected distance within a gene (Intra-Gene; light
grey).

Figure 4
Blondel, Jones & Extavour



**a.** *bacterial DNA transfer to germ line nucleus*

Ancient insect domain in host genome

LOTUS

Bacterial domain from donor genome

OSK

**b.** *HGT --> domain fusion*

LOTUS     OSK

**c.** de novo *domain evolution from inter-domain sequence*

5' UTR     LOTUS     OSK     3' UTR

In some Diptera

**d.** de novo *domain evolution from 5' UTR*

5' UTR     LONG OSK     LOTUS     OSK     3' UTR

334

**Figure 4. Hypothesis for the origin of *oskar*.** Integration of the OSK domain close to a LOTUS domain in an ancestral insect genome. **a**, DNA containing a GDSL-like domain from an endosymbiotic germ line bacterium is transferred to the nucleus of a germ cell in an insect common ancestor. **b**, DNA damage or transposable element activity induces an integration event in the host genome, close to a pre-existing LOTUS-like domain. **c**, The region between the two domains undergoes *de novo* coding evolution, creating an open reading frame with a unique, chimeric domain structure. **d**, In some Diptera, including *D. melanogaster*, part of the 5' UTR of *oskar* undergoes *de novo* coding evolution to form the Long Oskar domain.

342

## Materials and Methods

### BLAST searches of oskar

All BLAST[1] searches were performed using the NCBI BLASTp tool suite on the non-redundant (nr) database. Amino Acid (AA) sequences of *D. melanogaster* full length Oskar (EMBL ID AAF54306.1), as well as the AA sequences for the LOTUS (AA 139-238) and OSK (AA 414-606) domains were used for the BLAST searches, using the default NCBI cut-off parameters. As per NCBI defaults, the E-value cut-off was set at 10. All BLAST searches results are included in the Supplementary files: BLAST search results.

### Hidden Markov Model (HMM) generation and alignments of the OSK and LOTUS domains

101 1KITE transcriptomes[2] (Supplementary Table 1) were downloaded and searched using the local BLAST program (BLAST+) using the tblastn algorithm with default parameters, with Oskar protein sequences of *Drosophila melanogaster, Aedes aegypti, Nasonia vitripennis* and *Gryllus bimaculatus* as queries (EntrezIDs: NP_731295.1, ABC41128.1, NP_001234884.1 and AFV31610.1 respectively). For all of these 1KITE transcriptome searches, predicted protein sequences from transcript data were obtained by in silico translation using the online ExPASy translate tool (https://web.expasy.org/translate/), taking the longest open reading frame. Publicly available sequences in the non-redundant (nr), TSA databases at NCBI, and a then-unpublished transcriptome[3] (kind gift of Matthew Benton and Siegfried Roth, University of Cologne) were subsequently searched using the web-based BLAST tool hosted at NCBI, using the tblastn algorithm with default parameters. Sequences used for queries were the four Oskar proteins described above, and newfound *oskar* sequences from the 1KITE transcriptomes of

366  *Baetis pumilis, Cryptocercus wright,* and *Frankliniella cephalica*. For both searches, *oskar*

367  orthologs were identified by the presence of BLAST hits on the same transcript to both the

368  LOTUS (N-terminal) and OSK (C-terminal) regions of any of the query *oskar* sequences,

369  regardless of E-values. The sequences found were aligned using MUSCLE (8 iterations)[4] into a

370  46-sequence alignment (Supplementary files: Alignments>OSKAR_INITIAL.fasta). From this

371  alignment, the LOTUS and OSK domains were extracted (Supplementary files:

372  Alignments>LOTUS_INITIAL.fasta and Alignments>OSK_INITIAL.fasta) to define the

373  initial Hidden Markov Models (HMM) using the hmmbuild tool from the HMMER tool suite

374  with default parameters[5]. 126 insect genomes and 128 insect transcriptomes (from the

375  Transcriptome Shotgun Assembly TSA database: https://www.ncbi.nlm.nih.gov/Traces/wgs/?

376  view=TSA) were subsequently downloaded from NCBI (download date September 29, 2015 ;

377  Supplementary table 1). Genomes were submitted to Augustus v2.5.5[6] (using the *D.*

378  *melanogaster* exon HMM predictor) and SNAP v2006-07-28[7] (using the default 'fly' HMM)

379  for gene discovery. The resulting nucleotide sequence database comprising all 309 downloaded

380  and annotated genomes and transcriptomes, was then translated in six frames to generate a non-

381  redundant amino acid database (where all sequences with the same amino acid content are

382  merged into one). This process was automated using a series of custom scripts available here:

383  https://github.com/Xqua/Genomes. The non-redundant amino acid database was searched

384  using the HMMER v3.1 tool suite[5] and the HMM for the LOTUS and OSK domains described

385  above. A hit was considered positive if it consisted of a contiguous sequence containing both a

386  LOTUS domain and an OSK domain, with the two domains separated by an inter-domain

387  sequence. We imposed no length, alignment or conservation criteria on the inter-domain

388  sequence, as this is a rapidly-evolving region of Oskar protein with predicted high disorder[8-10].

389  Positive hits were manually curated and added to the main alignment, and the search was

390  performed iteratively until no more new sequences meeting the above criteria were discovered.

391  This resulted in a total of 95 Oskar protein sequences, (see Supplementary Table 2 for the

392  complete list). Using the final resulting alignment (Supplementary Files:

393  Alignments>OSKAR_FINAL.fasta), the LOTUS and OSK domains were extracted from these

394  sequences (Supplementary Files: Alignments>LOTUS_FINAL.fasta and

395  Alignments>OSK_FINAL.fasta), and the final three HMM (for full-length Oskar, OSK, and

396  LOTUS domains) used in subsequent analyses were created using hmmbuild with default

397  parameters (Supplementary files: HMM>OSK.hmm, HMM>LOTUS.hmm and

398  HMM>OSKAR.hmm).

399

400  ***Iterative HMMER search of OSK and LOTUS domains***

401  A reduced version of TrEMBL[11] (v2016-06) was created by concatenating all hits (regardless

402  of E-value) for sequences of the LOTUS domain, the OSK domain and full-length Oskar, using

403  hmmsearch with default parameters and the HMM models created above from the final

404  alignment. This reduced database was created to reduce potential false positive results that

405  might result from the limited size of the sliding window used in the search approach described

406  here. The full-length Oskar alignment of 1133 amino acids (Supplementary files:

407  Alignments>OSKAR_FINAL.fasta) was split into 934 sub-alignments of 60 amino acids each

408  using a sliding window of one amino acid. Each alignment was converted into a HMM using

409  hmmbuild, and searched against the reduced TrEMBL database using hmmsearch using default

410  parameters. Domain of life origin of every hit sequence at each position was recorded.

411  Eukaryotic sequences were further classified as Oskar/Non-Oskar and Arthropod/Non-

412 Arthropod. Finally, for the whole alignment, the counts for each category were saved and

413 plotted in a stack plot representing the proportion of sequences from each category to create

414 Fig. 1b. The python code used for this search is available at https://github.com/Xqua/Iterative-

415 HMMER.

416

417 *Sequence Similarity Networks*

418 LOTUS and OSK domain sequences from the final alignment obtained as described above (see

419 "*Hidden Markov Model (HMM) generation and alignments of the OSK and LOTUS domains*";

420 Supplementary files: Alignments>LOTUS_FINAL.fasta and Alignments>OSK_FINAL.fasta)

421 were searched against TrEMBL[11] (v2016-06) using HMMER. All hits with E-value < 0.01

422 were consolidated into a fasta file that was then entered into the EFI-EST tool[12] using default

423 parameters to generate a sequence similarity network. An alignment score corresponding to

424 30% sequence identity was chosen for the generation of the final sequence similarity network.

425 Finally, the network was graphed using Cytoscape 3[13].

426

427 *Phylogenetic Analysis*

428 For both the LOTUS and OSK domains, in cases where more than one sequence from the same

429 organism was retrieved by the search described above in "*Iterative HMMER Search of OSK*

430 *and LOTUS domains*", only the sequence with the lowest E-value was used for phylogenetic

431 analysis. For the LOTUS domain, the first 97 best hits (lowest E-value) were selected, and the

432 only three bacterial sequences that satisfied an E-value < 0.01 were manually added. For the

433 OSK domain, the first 95 best hits (lowest E-value) were selected, and the only five eukaryotic

434 sequences that satisfied an E-value < 0.01 were manually added. The sequences were filtered

435   to contain only one sequence per species (best E-value kept) generating a set of 100 sequences

436   for the LOTUS domain, and 87 for the OSK domain. Unique identifiers for all sequences used

437   to generate alignments for phylogenetic analysis are available in Supplementary Tables S3, S4.

438   For both datasets, the sequences were then aligned using MUSCLE[4] (8 iterations) and trimmed

439   using trimAl[14] with 70% occupancy. The resulting alignments that were subject to phylogenetic

440   analysis are available in Supplementary Files: Alignments>LOTUS_TREE.fasta and

441   Alignments>OSK_TREE.fasta. For the maximum likelihood tree, we used RaxML v8.2.4[15]

442   with 1000 bootstraps, and the models were selected using the automatic RaxML model

443   selection tool. The substitution model chosen for both domains was LGF. For the Bayesian tree

444   inference, we used MrBayes V3.2.6[16] with a Mixed model (prset aamodel=Mixed) and a

445   gamma distribution (lset rates=Gamma). We ran the MonteCarlo for 4 million generations (std

446   < 0.01) for the OSK domain, and for 3 million generations (std < 0.01) for the LOTUS domain.

447

448   ***Selection of sequences for codon use analysis***

449   To study the codon use of the OSK and LOTUS domains, we chose 17 well-annotated (defined

450   as possessing at least 8,000 annotated genes) insect genomes that included a confidently

451   annotated *oskar* orthologue from the NCBI nucleotide database. The complete list and

452   accession numbers of the sequences used for this analysis is in Supplementary Table 5. This

453   list contains  *oskar* sequences from genomes that were either added to the databases after the

454   first *oskar* sequence search or re-annotated after said search. Therefore the sequences coming

455   from the following organisms are not represented in the final *oskar* alignment: *Harpegnathos*

456   *saltator, Fopius arisanus, Athalia rosae, Orussus abietinus, Stomoxys calcitrans, Bactrocera*

457   *oleae, Neodiprion lecontei.*

458

### *Generation of Intra-Gene distribution of codon use*

460  We wished to determine whether *oskar* differed from the null hypothesis that a given gene

461  would follow similar codon use throughout its sequence. To generate a distribution of codon

462  use similarity across a gene for all genes in the genomes studied, we generated what we named

463  the "Intra-Gene" sequence distribution. Each gene was cut into two fragments at a random

464  position "x" following the rule: $384 < x <$ Length_gene - 384, x modulo 3 = 0 (Corresponding

465  Jupyter notebook file: Scripts>notebook>Codon Analysis AT3 GC3 and A3 T3 G3 C3 Section:

466  4). Thus, we sampled each codon at least twice, preserving the coding frame.

467

### *Fitting a linear model of codon use*

469  Using the Intra-Gene null distribution generated above, we fitted a linear model of codon use

470  frequencies per gene for the wobble position and AT3 GC3 content. To do so, we measured the

471  different frequencies of A3, T3, G3 and C3 (any codon ending in A was counted as A3) and

472  AT3 GC3. Then, we fitted a linear model to the pairs of 5' and 3' regional codon use values for

473  within each gene, obtained from the Intra-Gene distribution described above (conserving the

474  3'/5' position information), and for the OSK and LOTUS domains, for each of the 17 genomes

475  analyzed (Supp Table 3). We then calculated the residuals of the Intra-Gene distribution and

476  the LOTUS-OSK distribution. Finally, we determined the Pearson correlation coefficient for

477  all genomes pooled together, and all *oskar* genes pooled together (Corresponding Jupyter

478  notebook file: Scripts>notebook>Codon Analysis AT3 GC3 and A3 T3 G3 C3 Section: 7 and

479  8).

480

*Calculation of cosine distance*

482 For a given sequence S, we assigned a vector C of dimension 64 (one for each codon). Because

483 the sum of all codon frequencies is 1, C is normalized; we thus used the cosine similarity

484 distance between a given pair of vectors as a metric to quantify the distance in codon use

485 between two sequences. We measured this distance distribution between all the genes in a

486 given genome to create the Inter-Gene distance distribution. Then, we repeated the process but

487 measured the distance between all pairs of genes in the Intra-Gene sequence set per genome.

488 Next, we measured the distance between the LOTUS and OSK domains for each genome.

489 Finally, we determined the Z score of the distance between the LOTUS and OSK domains, and

490 the Inter-Gene and Intra-Gene distance distributions (Corresponding Jupyter notebook file:

491 Scripts>notebook>Cosine Distance Analysis).

492

*Calculation and analysis of the codon use Z_score*

494 For each genome, the codon use frequency for AT3/GC3 and A3/T3/G3/C3 was calculated as

495 described above. Then, Z scores for each sequence from the Intra-Gene, OSK or LOTUS

496 domain sequences were calculated against the corresponding genome frequency distribution.

497 The Z scores were then used to generate the analysis of Pearson correlation coefficients shown

498 in Figures 3, S5 and S6 (Corresponding Jupyter notebook file: Scripts>notebook>Codon

499 Analysis AT3 GC3 and A3 T3 G3 C3 Section: 3, 5 and 6).

500

*Data availability*

502   All sequences discovered using the automatic annotation pipeline described in (M&M HMM

503   and oskar search) are annotated as such in Supplementary Table S2.

504

505   *Code availability*

506   All custom code generated for this study is available in Supplementary Information>Scripts.

507

508

## Methods References

510    1    Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local
511        alignment search tool. *J. Mol. Biol.* **215**, 403-410, doi:10.1016/S0022-2836(05)80360-2
512        (1990).
513    2    Aspöck, H. *et al. 1KITE - 1K Insect Transcriptome Evolution,* <http://www.1kite.org/>
514        (2018).
515    3    Benton, M. A., Kenny, N. J., Conrads, K. H., Roth, S. & Lynch, J. A. Deep, Staged
516        Transcriptomic Resources for the Novel Coleopteran Models Atrachya menetriesi and
517        Callosobruchus maculatus. *PLoS ONE* **11**, e0167431,
518        doi:10.1371/journal.pone.0167431 (2016).
519    4    Edgar, R. C. MUSCLE: a multiple sequence alignment method with reduced time and
520        space complexity. *BMC Bioinformatics* **5**, 113 (2004).
521    5    Eddy, S. R. *HMMER: biosequence analysis using profile hidden Markov models,*
522        <http://hmmer.org/> (2007).
523    6    Stanke, M., Steinkamp, R., Waack, S. & Morgenstern, B. AUGUSTUS: a web server
524        for gene finding in eukaryotes. *Nucleic Acids Res.* **32**, W309-312,
525        doi:10.1093/nar/gkh379 (2004).
526    7    Korf, I. Gene finding in novel genomes. *BMC Bioinformatics* **5**, 59, doi:10.1186/1471-
527        2105-5-59 (2004).
528    8    Jeske, M. *et al.* The Crystal Structure of the Drosophila Germline Inducer Oskar
529        Identifies Two Domains with Distinct Vasa Helicase- and RNA-Binding Activities.
530        *Cell Rep* **12**, 587-598, doi:10.1016/j.celrep.2015.06.055 (2015).
531    9    Ahuja, A. & Extavour, C. G. Patterns of molecular evolution of the germ line
532        specification gene *oskar* suggest that a novel domain may contribute to functional
533        divergence in *Drosophila*. *Dev. Genes Evol.* **222**, 65-77 (2014).
534    10    Yang, N. *et al.* Structure of Drosophila Oskar reveals a novel RNA binding protein.
535        *Proc Natl Acad Sci U S A* **112**, 11541-11546, doi:10.1073/pnas.1515568112 (2015).
536    11    Consortium, U. The Universal Protein Resource (UniProt) 2009. *Nucleic Acids Res.* **37**,
537        D169-174, doi:10.1093/nar/gkn664 (2009).
538    12    Gerlt, J. A. *et al.* Enzyme Function Initiative-Enzyme Similarity Tool (EFI-EST): A
539        web tool for generating protein sequence similarity networks. *Biochim. Biophys. Acta*
540        **1854**, 1019-1037, doi:10.1016/j.bbapap.2015.04.015 (2015).
541    13    Shannon, P. *et al.* Cytoscape: a software environment for integrated models of
542        biomolecular interaction networks. *Genome Res.* **13**, 2498-2504,
543        doi:10.1101/gr.1239303 (2003).
544    14    Capella-Gutierrez, S., Silla-Martinez, J. M. & Gabaldon, T. trimAl: a tool for
545        automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**,
546        1972-1973, doi:10.1093/bioinformatics/btp348 (2009).
547    15    Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
548        large phylogenies. *Bioinformatics* **30**, 1312-1313, doi:10.1093/bioinformatics/btu033
549        (2014).
550    16    Huelsenbeck, J. P. & Ronquist, F. MRBAYES: Bayesian inference of phylogenetic
551        trees. *Bioinformatics* **17**, 754-755 (2001).
552
553

554 **Supplementary Information for**
555
556 Bacterial contribution to genesis of the novel germ line determinant *oskar*
557 *Leo Blondel, Tamsin E. M. Jones and Cassandra G. Extavour*
558
559 The Supplementary Information for this paper consists of the following elements:
560
561 1. Supplementary Discussion (this document)
562 2. Supplementary References (this document)
563 3. Folder titled "Supplementary Information Files" containing the following sub-folders
564     *a.* Supplementary Information Files>Alignments
565         i. *All sequences identified and analyzed in this study, in FASTA format and*
566             *with corresponding Alignments*
567     *b.* Supplementary Information Files>BLAST search results
568         i. *Results of BLASTP searches with full length Oskar, OSK or LOTUS*
569             *domains as queries*
570     *c.* Supplementary Information Files>Data
571         i. *Necessary files for running the different ipython notebooks:*
572             1. *Taxonomy: Conversion table for UniProt ID to taxon information.*
573                *(uniprot_ID_taxa.tsv )*
574             2. *Codon_Genes: Contains the measured codon frequency for the*
575                *different genomes studied as .csv or .tsv files (organism_name.csv/*
576                *tsv), along with the DNA sequences of LOTUS and OSK domains*
577                *used in the codon use analysis (LOTUS_Seqeuences.gb and*
578                *SGNH_Seqeuences.gb)*
579             3. *Trees: Contains the tree files obtained from RaxML and MrBayes*
580                *phylogenetic analyses of the OSK and LOTUS domains.*
581     *d.* Supplementary Information Files>HMM
582         i. *HMM models used for iterative searching for sequences similar to full-*
583             *length Oskar, LOTUS and OSK domains*
584     *e.* Supplementary Information Files>Scripts
585         i. *All custom scripts used to implement the analysis pipelines described.*
586     *f.* Supplementary Information Files>Tables
587         i. *Supplementary Tables S1-S5 describing databases searched/analyzed and*
588             *all search results; Legends in this document*
589
590 Please download Supplementary Information Files here:
591 https://www.dropbox.com/s/q4sd5rty24gxprg/Blondel_Jones_Extavour_HGT_Supplementary
592 %20Information%20Files.zip?dl=0

593 **Supplementary Discussion**

594

595 *Phylogenetic relationships of the Oskar LOTUS domain*

596　　　　LOTUS sequences from non-Oskar proteins that were sufficiently similar to the Osk

597 LOTUS domain to be included in an alignment for phylogenetic analysis, were almost

598 exclusively eukaryotic. (Supplementary Table 3). Only three bacterial sequences matched the

599 LOTUS domain with an E-value < 0.01, and were included in the alignment (Supplementary

600 Table 3). Osk LOTUS domains clustered into two distinct clades, one comprising all Dipteran

601 sequences, and the other comprising all other Osk LOTUS domains examined from both

602 holometabolous and hemimetabolous orders (Fig. 2a). Dipteran Osk LOTUS sequences formed

603 a monophyletic group that branched sister to a clade of LOTUS domains from Tud5 family

604 proteins of non-arthropod animals (NAA). NAA LOTUS domains from Tud7 family members

605 were polyphyletic, but most of them formed a clade branching sister to (Osk LOTUS + NAA

606 Tud5 LOTUS). Non-Dipteran Osk LOTUS domains formed a monophyletic group that was

607 related in a polytomy to the aforementioned (NAA Tud7 LOTUS + (Dipteran Osk LOTUS +

608 NAA Tud5 LOTUS)) clade, and to various arthropod Tud7 family LOTUS domains.

609　　　　The fact that Tud7 LOTUS domains are polyphyletic suggests that arthropod domains

610 in this family may have undergone heterogeneous evolutionary processes relative to their

611 homologues in other animals. The relationships of Dipteran LOTUS sequences were consistent

612 with the current hypothesis for interrelationships between Dipteran species[1] Similarly, among

613 the non-Dipteran Osk LOTUS sequences, the hymenopteran sequences form a clade to the

614 exclusion of the single hemimetabolous sequence (from the cricket *Gryllus bimaculatus*),

615 consistent with the monophyly of Hymenoptera[2]. It is unclear why Dipteran Osk LOTUS

616 domains cluster separately from those of other insect Osk proteins. We speculate that the

617 evolution of the Long Oskar domain[3,4], which appears to be a novelty within Diptera

618 (Supplementary Files: Alignments>OSKAR_FINAL.fasta), may have influenced the evolution

619 of the Osk LOTUS domain in at least some of these insects. Consistent with this hypothesis, of

620 the 17 Dipteran *oskar* genes we examined, the seven *oskar* genes possessing a Long Osk

621 domain clustered into two clades based on the sequences of their LOTUS domain. One of these

622 clades comprised five Drosophila species (*D. willistoni*, *D. mojavensis*, *D. virilis*, *D.*

623 *grimshawi* and *D. immigrans*), and the second was composed of two calyptrate flies from

624 different superfamilies, *Musca domestica* (Muscoidea) and *Lucilia cuprina* (Oestroidea).

625　　　　In summary, the LOTUS domain of Osk proteins is most closely related to a number of

626 other LOTUS domains found in eukaryotic proteins, as would be expected for a gene of animal

627 origin, and the phylogenetic interrelationships of these sequences is largely consistent with the

628 current species or family level trees for the corresponding insects.

629

630 *Phylogenetic relationships of the Oskar OSK domain*

631　　　　The only eukaryotic proteins emerging from the iterative HMMER search for OSK

632 domain sequences that had an E-value < 0.01 were all from fungi. All five of these sequences

633 were annotated as  Carbohydrate Active Enzyme 3 (CAZ3). Most bacterial sequences used in

634 this analysis were annotated as lipases and hydrolases, with a high representation of GDSL-like

635 hydrolases (Supplementary Table S4). OSK sequences formed a monophyletic group but did

636 not branch sister to the other eukaryotic sequences in the analysis. Instead, all CAZ3 sequences

637 formed a clade that was sister to a clade of primarily Firmicutes. We recovered a monophyletic

638 group of Proteobacteria nested within that Firmicutes clade. All Bacteroidetes sequences also

639 formed a monophyletic group, which branched sister to all other sequences except for the two
640 Archaeal sequences in the analysis. Within the OSK clade, the topology of sequence
641 relationships was largely concordant with the species tree for insects [5], as we recovered
642 monophyletic Diptera to the exclusion of other insect species. However, the single orthopteran
643 OSK sequence (from the cricket *Gryllus bimaculatus*) grouped within the Hymenoptera, rather
644 than branching basally to all insect sequences as would be expected for this hemimetabolous
645 sequence.
646

647 **Supplementary Table Legends**
648
649 (see Supplementary Information Files>Tables>Supp TableX)
650
651 **Supplementary Table S1: List of genomes and transcriptomes used for automated *oskar***
652 **search.**
653 List of genomes and transcriptomes that were downloaded, annotated, and searched for *oskar*
654 sequences (*see "Hidden Markov Model (HMM) generation and alignments of the OSK and*
655 *LOTUS domains*" in Methods). The table reports the database provenance (NCBI genome or
656 TSA, or 1KITE database) and the accession number. The TSA accession  ID can be searched
657 using the NCBI TSA browser here: https://www.ncbi.nlm.nih.gov/Traces/wgs/?view=TSA.
658
659 **Supplementary Table S2: List of *oskar* sequences used in the final alignment.**
660 List of accession numbers and database provenance of the sequences used in the final
661 alignments of Oskar analysed herein. The table contains the database provenance (*Type*), the
662 database accession number (*ID*), the species, family and order, and extraction notes.
663
664 **Supplementary Table S3: List of sequences used for phylogenetic analysis of the LOTUS**
665 **domain.**
666 The sequences were obtained by searching the TrEMBL database using hmmsearch and the
667 final HMM generated for LOTUS (Supplementary files: HMM>LOTUS.hmm). Reported are
668 the UniProtID (*Accession Number*), the Domain and Phylum origin of the sequence, the E-
669 value, score and bias given by hmmsearch, and the description of the target from UniProt. To
670 obtain sequences for each entry, either search UniProt directly (https://www.uniprot.org/) or
671 consult the final alignment in Supplementary Files: Alignments>LOTUS_TREE.fasta.
672
673 **Supplementary Table S4: List of sequences used for phylogenetic analysis of the OSK**
674 **domain.**
675 The sequences were obtained by searching the TrEMBL database using hmmsearch and the
676 final HMM generated for OSK (Supplementary files: HMM>OSK.hmm). Reported parameters
677 are as described for Supplementary Table S3. To obtain sequences for each entry, either search
678 UniProt directly (https://www.uniprot.org/) or consult the final alignment in Supplementary
679 Files: Alignments>OSK_TREE.fasta.
680
681 **Supplementary Table S5: List of genomes analyzed for codon use.**
682 This table lists the 17 genomes that were downloaded and analyzed for codon use as described
683 in "*Selection of sequences for codon use analysis*" in Methods. All genomes can be
684 downloaded from https://www.ncbi.nlm.nih.gov/genome/browse#!/overview/. The table lists
685 the species name (*Species*), family (*Family*) and Order (*Order*), NCBI genome accession
686 number (*Genome ID*), and the *oskar* NCBI Nucleotide accession number (*oskar Nucleotide*
687 *ID*).
688

689 **Supplementary References**
690
691
692   1      Kirk-Spriggs, A. H. & Sinclair, B. J.  Vol. 1   (South African National Biodiversity
693           Institute, Pretoria, South Africa, 2017).
694   2      Peters, R. S. *et al.* Evolutionary History of the Hymenoptera. *Curr. Biol.* **27**, 1013-
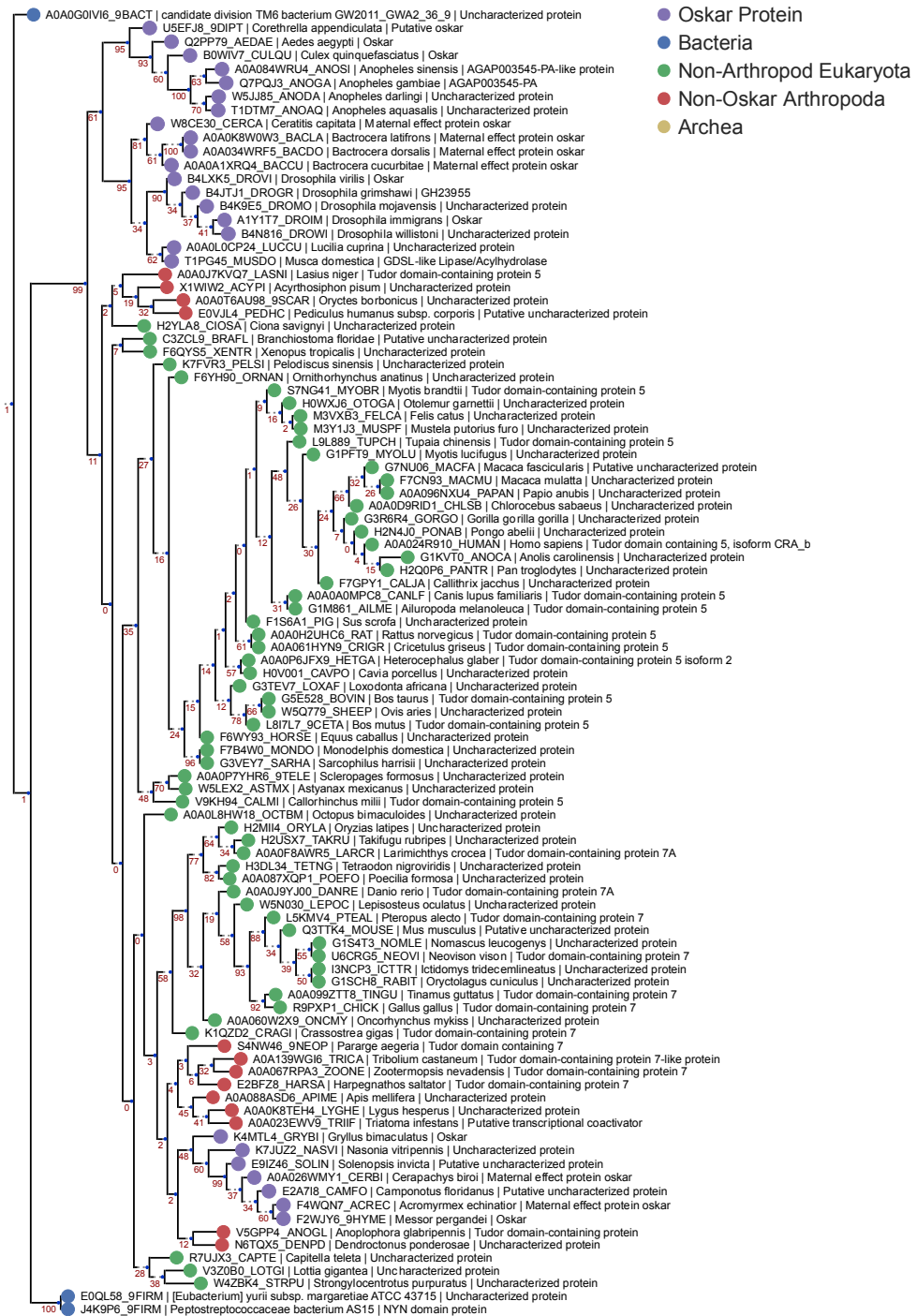695           1018, doi:10.1016/j.cub.2017.01.027 (2017).
696   3      Vanzo, N. F. & Ephrussi, A. Oskar anchoring restricts pole plasm formation to the
697           posterior of the *Drosophila* oocyte. *Development* **129**, 3705-3714 (2002).
698   4      Hurd, T. R. *et al.* Long Oskar Controls Mitochondrial Inheritance in Drosophila
699           melanogaster. *Dev Cell* **39**, 560-571, doi:10.1016/j.devcel.2016.11.004 (2016).
700   5      Misof, B. *et al.* Phylogenomics resolves the timing and pattern of insect evolution.
701           *Science* **346**, 763-767, doi:10.1126/science.1257570 (2014).
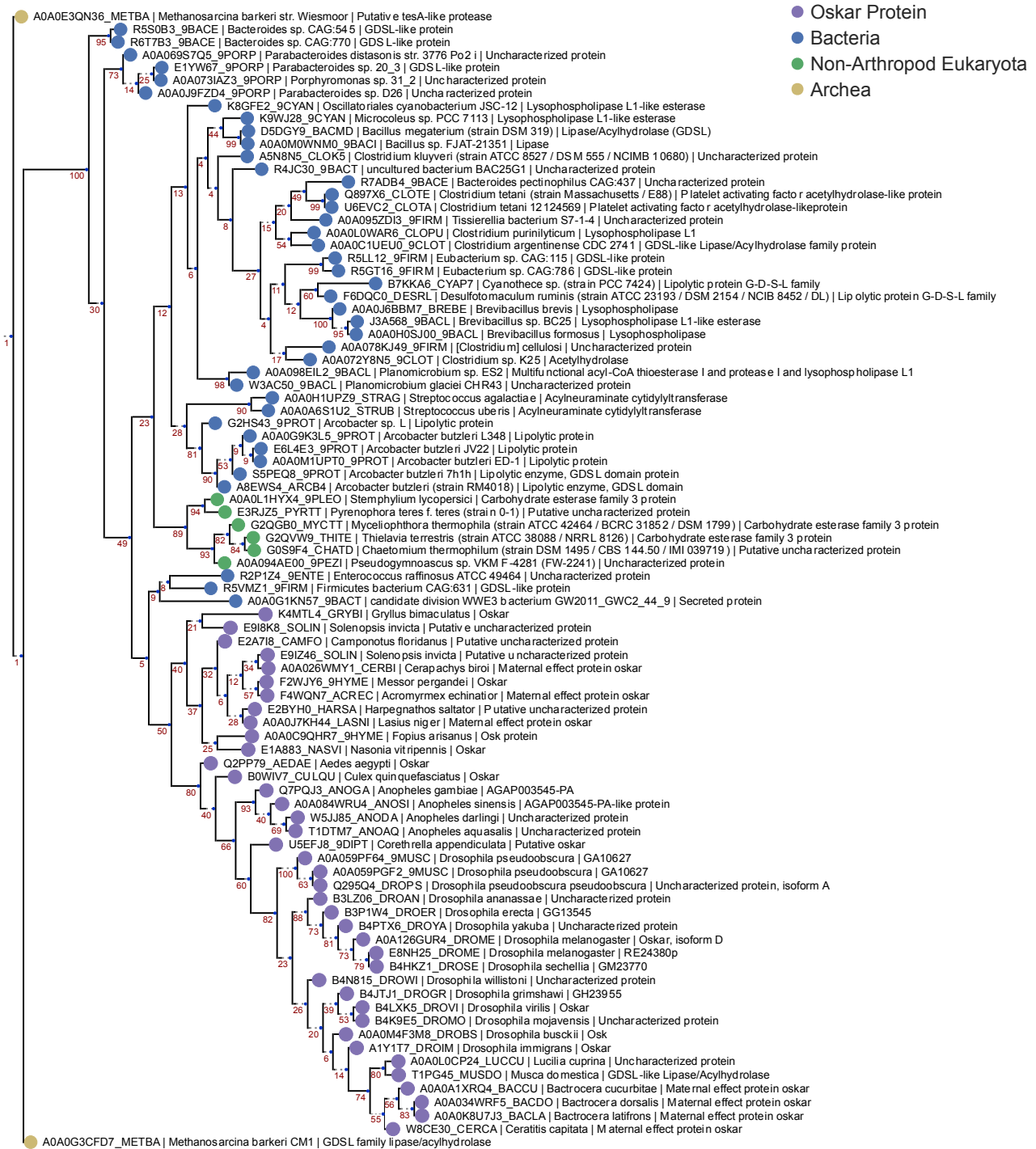702
703

**Oskar Protein**
**Bacteria**
**Non-Arthropod Eukaryota**
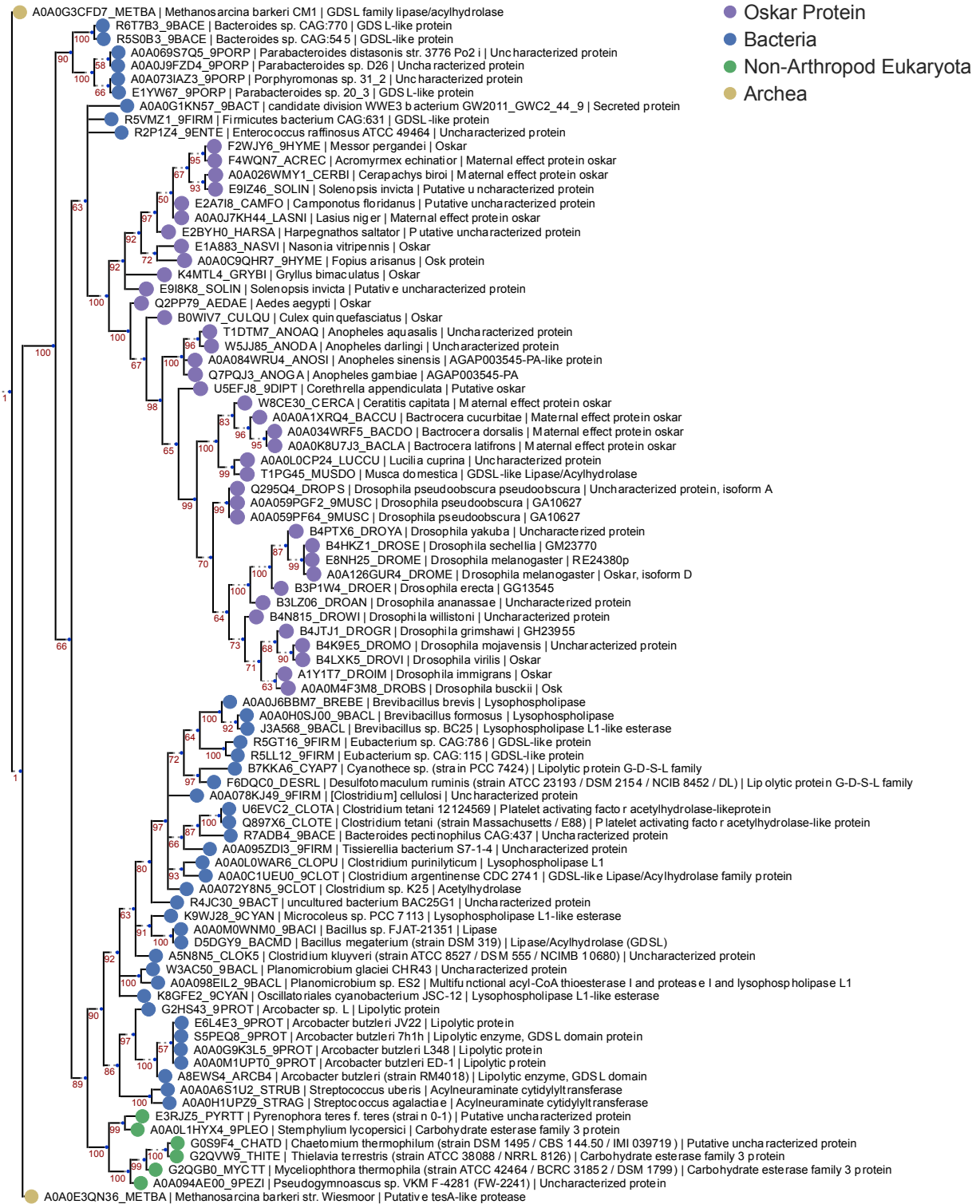**Non-Oskar Arthropoda**
**Archea**

704

**Extended Data Figure S1: LOTUS Domain RaxML Tree.** Phylogenetic tree of the HMMER sequences retrieved
706 from the UniProt database using the LOTUS alignment HMM model. The top 97 hits were selected for phylogenetic
707 analysis, and the only three bacterial sequences found to be a match were added to the alignment manually. The
708 resulting 100 sequences were aligned using MUSCLE with default settings. The sequences were filtered to contain
709 only one sequence per species (best E-value kept) yielding 100 sequences for analysis. Finally, the tree was created
710 using RaxML v8.2.4, using 1000 bootstraps and model selection performed by the RaxML automatic model
711 selection tool. See "Phylogenetic Analysis" in Methods for further detail. Sequences are color-coded as follows:
712 Purple = Oskar; Red = Non-Oskar Arthropod; Green = Non-Arthropod Eukaryote; Blue = Bacteria. Names
713 following leaves display the UniProt accession number followed by the species name and the UniProt protein name.

714

**Extended Data Figure S2: LOTUS Domain Bayesian Tree.** Phylogenetic tree of the HMMER sequences retrieved from the UniProt database using the LOTUS alignment HMM model. 100 sequences were chosen for analysis as described for Supplementary Figure 1. The tree was created using Mr Bayes V3.2.6 using a Mixed model (prset aamodel=Mixed) and a gamma distribution (lset rates=Gamma). The algorithm was allowed to run for 3 million generations to achieve a std < 0.01. See "Phylogenetic Analysis" in Methods for further detail. Sequences are color-coded as follows: Purple = Oskar; Red = Non-Oskar Arthropod; Green = Non-Arthropod Eukaryote; Blue = Bacteria. Names following leaves display the UniProt accession number followed by the species name and the UniProt protein name.
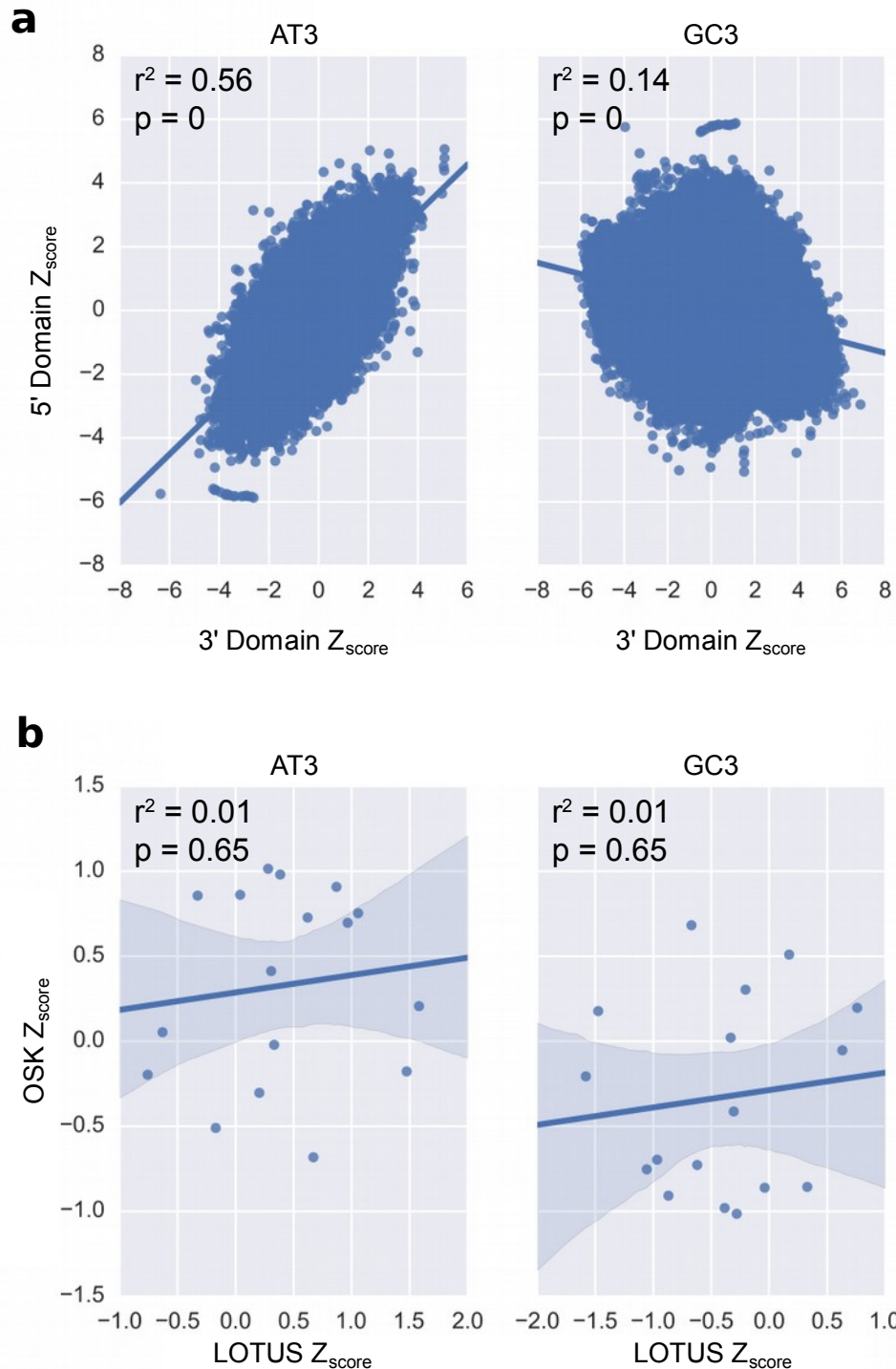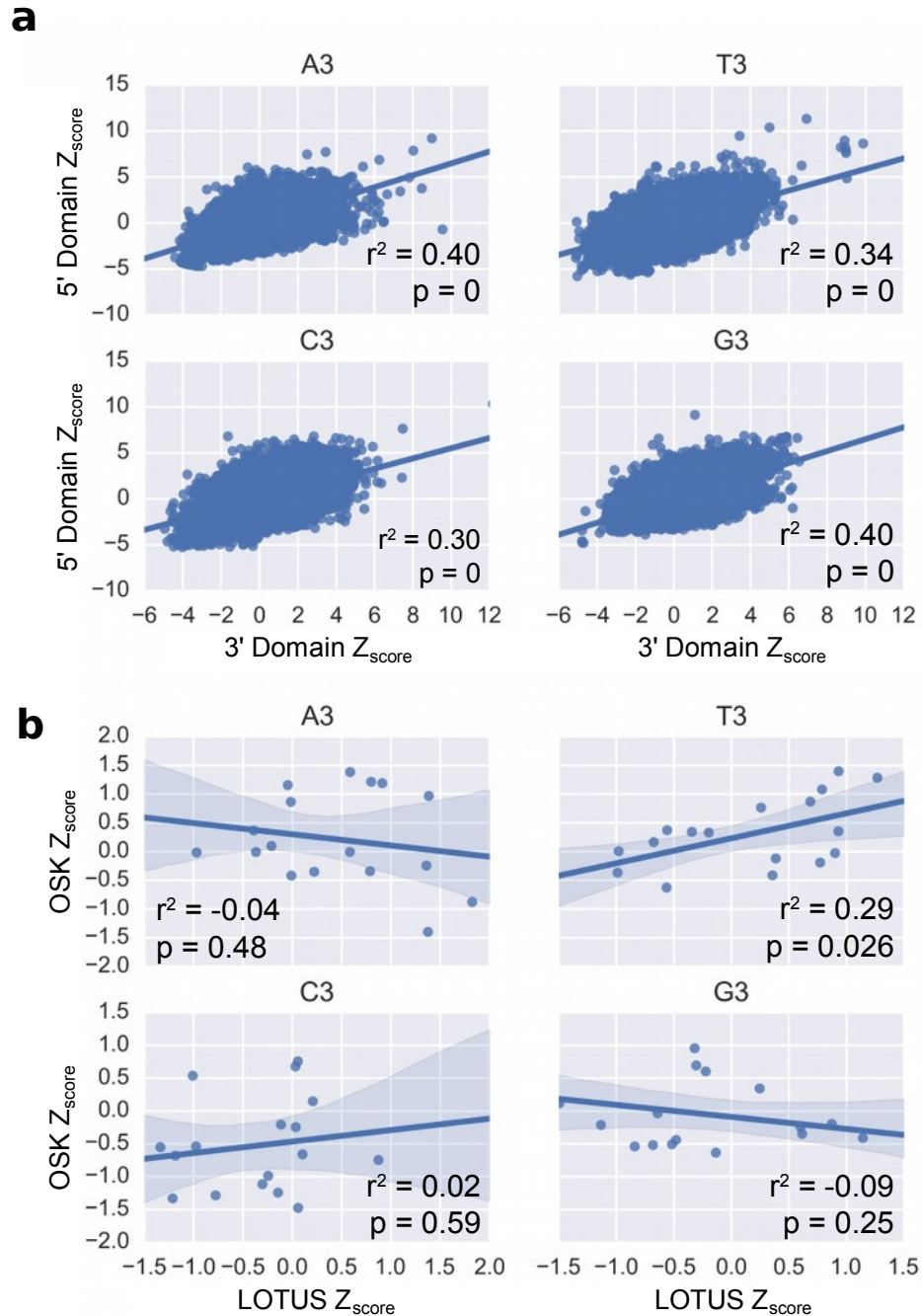
722

**Extended Data Figure S3: OSK Domain RaxML Tree.** Phylogenetic tree of the HMMER sequences retrieved from the UniProt database using the OSK alignment HMM model. The top 95 hits were selected for phylogenetic analysis, and the only five non-Oskar eukaryotic sequences found to be a match were added to the alignment manually. The resulting 100 sequences were aligned using MUSCLE with default settings. The sequences were filtered to contain only one sequence per species (best E-value kept), yielding 87 sequences for analysis. Finally, the tree was created using RaxML v8.2.4, using 1000 bootstraps and model selection performed by the RaxML automatic model selection tool. See "Phylogenetic Analysis" in Methods for further detail. Sequences are color-coded as follows: Purple = Oskar; Red = Non-Oskar Arthropod; Green = Non-Arthropod Eukaryote; Blue = Bacteria. Names following leaves display the UniProt accession number followed by the species name and the UniProt protein name.

**Oskar Protein**
**Bacteria**
**Non-Arthropod Eukaryota**
**Archea**

732

733**Extended Data Figure S4: OSK Domain Bayesian Tree.** Phylogenetic tree of the HMMER sequences hit on the UniProt
734database using the OSK alignment HMM model. 87 sequences were chosen for analysis as described for Supplementary
735Figure 3.The tree was created using Mr Bayes V3.2.6 using a Mixed model (prset aamodel=Mixed) and a gamma
736distribution (lset rates=Gamma). The algorithm was allowed to run for 4 million generations to achieve a std < 0.01. See
737"Phylogenetic Analysis" in Methods for further detail. Sequences are color-coded as follows: Purple = Oskar; Red = Non-
738Oskar Arthropod; Green = Non-Arthropod Eukaryote; Blue = Bacteria. Names following leaves display the UniProt
739accession number followed by the species name and the UniProt protein name.

741 **Extended Data Figure S5: AT3/GC3 correlations between the LOTUS and OSK domains.** (**a**) Intra-Gene distribution
742 scatter plot for the coding sequences of the 17 genomes analyzed. Sequences were cut into two parts as per the description
743 in Methods "Generation of intra-gene distribution of codon use". The AT3 and GC3 codon use was measured and a Z-score
744 was calculated against the genome distribution. Finally, the 5' and 3' "domain" values were plotted against each other and a
745 linear regression was . The AT3 and GC3 content is generally similar in the 5' and 3'regions of all genes across the genome
746 (AT3: $r^2 = 0.56$, p = 0; GC3: $r^2 = 0.14$, p = 0). (**b**) OSK vs LOTUS AT3 and GC3 use across the 17 genomes analyzed. The
747 AT3 and GC3 content Z-scores were calculated against the genome distribution. The AT3 and GC3 content of the two
748 domains of the Oskar gene are not correlated with each other. (AT3: $r^2 = 0.01$, p = 0.65; GC3: $r^2 = 0.01$, p = 0.65).

**Extended Data Figure S6: A3/T3/G3/C3 correlations between the LOTUS and OSK domains. (a)** Intra-Gene distribution scatter plot for the coding sequences of the 17 genomes analyzed. Sequences were cut into two parts as per the description in Methods "Generation of intra-gene distribution of codon use". The A3, T3, G3 and C3 codon use was measured, and Z-score calculations, value plots and linear regression were performed as described for Supplementary Figure 5. The A3, T3 G3 and C3 content is generally similar in the 5' and 3'regions of all genes across the genome (A3: $r^2 = 0.40$, p = 0; T3: $r^2 = 0.34$, p = 0; G3: $r^2 = 0.40$, p = 0; C3: $r^2 = 0.30$, p = 0). **(b)** OSK vs LOTUS A3, T3, G3 and C3 use across the 17 genomes analyzed. The A3, T3, G3 and C3 content Z-score were calculated against the genome distribution. The A3, G3 and C3 content of the two domains of the Oskar gene are not correlated with each other. However, the T3 distribution follows a linear correlation similar to the one found across the Intra-Gene distribution (A3: $r^2 = -0.04$, p = 0.48; T3: $r^2 = 0.29$, p = 0.026; G3: $r^2 = -0.09$, p = 0.25; C3: $r^2 = 0.02$, p = 0.59).