

RESEARCH

Open Access



Contrasting patterns of evolutionary constraint and novelty revealed by comparative sperm proteomic analysis in Lepidoptera

Emma Whittington¹, Desiree Forsythe², Kirill Borziak¹, Timothy L. Karr³, James R. Walters⁴ and Steve Dorus^{1*} 

Abstract

Background: Rapid evolution is a hallmark of reproductive genetic systems and arises through the combined processes of sequence divergence, gene gain and loss, and changes in gene and protein expression. While studies aiming to disentangle the molecular ramifications of these processes are progressing, we still know little about the genetic basis of evolutionary transitions in reproductive systems. Here we conduct the first comparative analysis of sperm proteomes in Lepidoptera, a group that exhibits dichotomous spermatogenesis, in which males produce a functional fertilization-competent sperm (eupyrene) and an incompetent sperm morph lacking nuclear DNA (apyrene). Through the integrated application of evolutionary proteomics and genomics, we characterize the genomic patterns potentially associated with the origination and evolution of this unique spermatogenic process and assess the importance of genetic novelty in Lepidopteran sperm biology.

Results: Comparison of the newly characterized Monarch butterfly (*Danaus plexippus*) sperm proteome to those of the Carolina sphinx moth (*Manduca sexta*) and the fruit fly (*Drosophila melanogaster*) demonstrated conservation at the level of protein abundance and post-translational modification within Lepidoptera. In contrast, comparative genomic analyses across insects reveals significant divergence at two levels that differentiate the genetic architecture of sperm in Lepidoptera from other insects. First, a significant reduction in orthology among Monarch sperm genes relative to the remainder of the genome in non-Lepidopteran insect species was observed. Second, a substantial number of sperm proteins were found to be specific to Lepidoptera, in that they lack detectable homology to the genomes of more distantly related insects. Lastly, the functional importance of Lepidoptera specific sperm proteins is broadly supported by their increased abundance relative to proteins conserved across insects.

Conclusions: Our results identify a burst of genetic novelty amongst sperm proteins that may be associated with the origin of heteromorphic spermatogenesis in ancestral Lepidoptera and/or the subsequent evolution of this system. This pattern of genomic diversification is distinct from the remainder of the genome and thus suggests that this transition has had a marked impact on lepidopteran genome evolution. The identification of abundant sperm proteins unique to Lepidoptera, including proteins distinct between specific lineages, will accelerate future functional studies aiming to understand the developmental origin of dichotomous spermatogenesis and the functional diversification of the fertilization incompetent apyrene sperm morph.

Keywords: Spermatogenesis, Lepidoptera, Fertility, Sexual selection, Testis, Mass spectrometry, Parasperms, Apyrene sperm, Positive selection, Genomic

* Correspondence: sdorus@syr.edu

¹Center for Reproductive Evolution, Department of Biology, Syracuse University, Syracuse, NY, USA

Full list of author information is available at the end of the article



Background

Spermatozoa exhibit an exceptional amount of diversity at both the ultrastructure and molecular levels despite their central role in reproduction [1]. One of the least understood peculiarities in sperm variation is the production of heteromorphic sperm via dichotomous spermatogenesis, the developmental process where males produce multiple distinct sperm morphs that differ in their morphology, DNA content and/or other characteristics [2]. Remarkably, one sperm morph is usually fertilization incompetent and often produced in large numbers; such morphs are commonly called “parasperm”, in contrast to fertilizing “eusperm” morphs. Despite the apparent inefficiencies of producing sperm morphs incapable of fertilization, dichotomous spermatogenesis has arisen independently across a broad range of taxa, including insects, brachiopod molluscs and fish. This paradoxical phenomenon, where an investment is made into gametes that will not pass on genetic material to the following generation, has garnered substantial interest, and a variety of hypotheses regarding parasperm function have been postulated [3]. In broad terms, these can be divided into three main functional themes: (1) facilitation, where parasperm aid the capacitation or motility of eusperm in the female reproductive tract, (2) provisioning, where parasperm provide nutrients or other necessary molecules to eusperm, the female or the zygote and (3) mediating postcopulatory sexual selection, where parasperm may serve eusperm either defensively or offensively by delaying female remating, influencing rival sperm, or biasing cryptic female choice. Despite experimental efforts in a number of taxa, a robust determination of parasperm function has yet to be attained.

Dichotomous spermatogenesis was first identified in Lepidoptera [4], the insect order containing butterflies and moths, over a century ago and is intriguing because the parasperm morph (termed apyrene sperm), is anucleate and therefore lacks nuclear DNA. Although it has been suggested that apyrene sperm are the result of a degenerative evolutionary process, several compelling observations suggest that dichotomous spermatogenesis is likely adaptive. First, it has been clearly demonstrated that both sperm morphs are required for successful fertilization in the silkworm moth (*Bombyx mori*) [5]. Second, phylogenetic relationships indicate ancestral origins of dichotomous spermatogenesis and continued maintenance during evolution. For example, dichotomous spermatogenesis is present throughout Lepidoptera, with the sole exception of two species within the most basal suborder of this group. Although multiple independent origins of sperm heteromorphism in Lepidoptera has yet to be formally ruled out, a single ancestral origin is by far the most parsimonious explanation [6].

Third, the ratio of eupyrene to apyrene varies substantively across Lepidoptera but is relatively constant within species, including several cases where apyrene comprise up to 99% of the sperm produced [7]. While variation in the relative production of each sperm morph is not in itself incompatible with stochastic processes, such as drift, it is nearly impossible to reconcile the disproportionate investment in apyrene without acknowledging that they contribute in some fundamental way to reproductive fitness. Although far from definitive, it has also been suggested that this marked variability across species is consistent with ongoing diversifying selection [6]. Arriving at an understanding of apyrene function may be further complicated by the possibility that parasperm are generally more likely to acquire lineage specific functionalities [8].

To better understand the molecular basis of dichotomous spermatogenesis, we recently conducted a proteomic and genomic characterization of sperm in *Manduca sexta* (hereafter *Manduca*) [9]. An important component of our analysis was to determine the taxonomic distribution of sperm proteins, which revealed an unexpectedly high number of proteins that possess little or no homology to proteins outside of Lepidoptera. This pattern is consistent with genetic novelty associated with dichotomous spermatogenesis in Lepidoptera, although we cannot formally rule out relaxation of purifying selection (on apyrene sperm proteins, for example) as an explanation for this marked divergence. Sperm proteins unique to Lepidoptera were also determined to be significantly more abundant than other sperm proteins. Given that apyrene spermatogenesis accounts for 95% of all sperm production in *Manduca* [7], these proteins are likely to be present and function in the more common apyrene sperm morph.

To provide a deeper understanding of the role of genetic novelty and genomic diversification in the evolution of dichotomous spermatogenesis, we have characterized the sperm proteome of the Monarch butterfly (*Danaus plexippus*; hereafter Monarch). In addition to its phylogenetic position and its continued development as a model butterfly species, we have pursued this species because of its distinct mating behavior. Unlike most other Lepidopteran species, male Monarch butterflies employ a strategy of coercive mating, as a consequence female Monarchs remate frequently [10]. In contrast, female remating is rare in *Manduca* and, as in many other Lepidoptera, females attract males via pheromonal calling behavior [11]. Interestingly, cessation of calling appears to be governed by molecular factors present in sperm or seminal fluid [12] and, as a consequence, non-virgin females rarely remate. Despite these behavioral differences, the proportion of eupyrene and apyrene produced

is quite similar between these two species (~95–96%) [7, 13]. Thus, our focus on Monarch is motivated both by their disparate, polyandrous mating system and their utility as a representative butterfly species for comparative analyses with *Manduca*. Therefore, the overarching aims of this study were to (1) characterize the sperm proteome of the Monarch butterfly and compare it with the previously characterized sperm proteome of *Manduca*, (2) contrast patterns of orthology across diverse insect genomes between the sperm proteome and remainder of genes in the genome and (3) analyze genome-wide homology to assess the contribution of evolutionary genetic novelty to Lepidopteran sperm composition.

Methods

Butterfly rearing and sperm purification

Adult male Monarch butterflies, kindly provided by MonarchWatch (Lawrence, Kansas), were dissected between 5 and 10 days post eclosion. The sperm contents of seminal vesicles, including both apyrene and eupyrene sperm, were dissected via a small incision in the mid to distal region of the seminal vesicle. Samples were rinsed in phosphate buffer solution and pelleted via centrifugation (2 min at 15000 rpm) three times to produce a purified sperm sample. Sperm samples from 3 groups of 5 separate males were pooled to form three biological replicates [14].

Protein preparation and 1-dimensional SDS page

Samples were solubilized in 2X LDS sample buffer, as per manufacturers' instructions (Invitrogen, Inc) before quantification via the EZA Protein Quantitation Kit (Invitrogen, Inc). Protein fluorescence was measured using a Typhoon Trio + (Amersham Biosciences/GE Healthcare) with 488 nm excitation and a 610 nm bandpass filter. Fluorescence data was analyzed using the ImageQuant TL software. Three replicates of 25 µg of protein were separated on a 1 mm 10% NuPAGE Novex Bis-Tris Mini Gel set up using the XCell SureLock Mini-Cell system (Invitrogen) as per manufacturer instructions for reduced samples. Following electrophoresis, the gel was stained using SimplyBlue SafeStain (Invitrogen, Inc) and destained as per manufacturer instructions. Each lane on the resulting gel (containing a sample from a single replicate) was sliced into four comparable slices, producing 12 gel fractions for independent tandem mass spectrometry analysis.

Tandem mass spectrometry (MS/MS)

Gel fractions were sliced into 1 mm² pieces for in-gel trypsin digestion. Gel fractions were reduced (DDT) and alkylated (iodoacetamide) before overnight incubation with trypsin at 37 °C. All LC-MS/MS experiments were performed using a Dionex Ultimate 3000 RSLC nanoUPLC

(Thermo Fisher Scientific Inc., Waltham, MA, USA) system and a QExactive Orbitrap mass spectrometer (Thermo Fisher Scientific Inc., Waltham, MA, USA). Separation of peptides was performed by reverse-phase chromatography at a flow rate of 300 nL/min and a Thermo Scientific reverse-phase nano Easy-spray column (Thermo Scientific PepMap C18, 2 µm particle size, 100A pore size, 75 mm i.d. × 50 cm length). Peptides were loaded onto a pre-column (Thermo Scientific PepMap 100 C18, 5 µm particle size, 100A pore size, 300 mm i.d. × 5 mm length) from the Ultimate 3000 autosampler with 0.1% formic acid for 3 min at a flow rate of 10 µL/min. After this period, the column valve was switched to allow elution of peptides from the pre-column onto the analytical column. Solvent A was water plus 0.1% formic acid and solvent B was 80% acetonitrile, 20% water plus 0.1% formic acid. The linear gradient employed was 2–40% B in 30 min. The LC eluant was sprayed into the mass spectrometer by means of an Easy-spray source (Thermo Fisher Scientific Inc.). All m/z values of eluting ions were measured in an Orbitrap mass analyzer, set at a resolution of 70,000. Data dependent scans (Top 20) were employed to automatically isolate and generate fragment ions by higher energy collisional dissociation (HCD) in the quadrupole mass analyzer and measurement of the resulting fragment ions was performed in the Orbitrap analyzer, set at a resolution of 17,500. Peptide ions with charge states of 2+ and above were selected for fragmentation. The mass spectrometry proteomics data have been deposited to the ProteomeXchange Consortium via the PRIDE partner repository with the dataset identifier PXD006454 [15].

MS/MS data analysis

MS/MS data was analyzed using X!Tandem and Comet algorithms within the Trans-Proteomic Pipeline (v 4.8.0) [16]. Spectra were matched against the *D. plexippus* official gene set 2 (OGS2) predicted protein set (downloaded from <http://Monarchbase.umassmed.edu>, last updated in 2012) with a fragment ion mass tolerance of 0.40 Da and a parent monoisotopic mass error of ±10 ppm. For both X!tandem and Comet, iodoacetamide derivative of cysteine was specified as a fixed modification, whereas oxidation of methionine was specified as a variable modification. Two missed cleavages were allowed and non-specific cleavages were excluded from the analysis. False Discovery Rates (FDRs) were estimated using a decoy database of randomized sequence for each protein in the annotated protein database. Peptide identifications were filtered using a greater than 95.0% probability based upon PeptideProphet [17] and the combined probability information from X!Tandem and Comet using Interprophet. Protein assignments were accepted if greater than 99.0%, as specified by the ProteinProphet [18] algorithms respectively. Proteins that contained identical peptides

that could not be differentiated based on MS/MS analysis alone were grouped to satisfy the principles of parsimony. Protein inclusion in the proteome was based on the following stringent criteria: (1) identification in 2 or more biological replicates or (2) identification in a single replicate by 2 or more unique peptides. To identify post-translation modifications (PTMs) of proteins, X!Tandem and Comet were rerun allowing for variable phosphorylation of serine, threonine and tyrosine residues and acetylation of lysine residues. PTM locations were identified using PTMprophet in both the Monarch data presented here and a comparable dataset in *M. sexta* [19].

APEX protein quantitation and analysis

Relative compositional protein abundance was quantified using the APEX Quantitative Proteomics Tool [20]. The training dataset was constructed using fifty proteins with the highest number of uncorrected spectral counts (n_i), and identification probabilities. All 35 physicochemical properties available in the APEX tool were used to predict peptide detection/non-detection. Protein detection probabilities (O_i) were computed using proteins with identification probabilities over 99% and the Random Forest classifier algorithm. APEX protein abundances were calculated using a merged protXML file generated by the ProteinProphet algorithm and highly correlated (all pairwise p values $<9.3 \times 10^{-10}$). The correlation in APEX abundance estimates of orthologous proteins in Monarch and *Manduca* (abundance estimates from Whittington et al. [9]) were normalized, log transformed and assessed using linear regression. Differential protein abundance was analyzed using corrected spectral counts and the R (v 3.0.0) package EdgeR [21]. Results were corrected for multiple testing using the Benjamini-Hochberg method within EdgeR.

Lift-over between *D. plexippus* version 1 and 2 gene sets

Two versions of gene models and corresponding proteins are currently available for *D. plexippus*. Official gene set one (OGS1) was generated using the genome assembly as initially published [22], while the more recent official gene set 2 (OGS2) was generated along with an updated genome assembly [23]. While our proteomic analysis employs the more recent OGS2 gene models, at the time of our analysis only OGS1 gene models were included in publicly available databases for gene function and orthology (e.g. Uniprot and OrthoDB). In order to make use of these public resources, we assigned OGS2 gene models to corresponding OGS1 gene models by sequence alignment. Specifically, OGS2 coding sequences (CDS) were aligned to OGS1 CDS using BLAT [24], requiring 95% identity; the best aligning OGS1 gene model was

assigned as the match for the OGS2 query. In this way, we were able to link predictions of OGS1 gene function and orthology in public databases to OGS2 sequences in our analysis. Of the 584 OGS2 loci identified in the sperm proteome 18 could not be assigned to an OGS1 gene.

Functional annotation and enrichment analysis

Two approaches were employed for functionally annotating *D. plexippus* sperm protein sequences. First, we obtained functional annotations assigned by Uniprot to corresponding *D. plexippus* OGS1 protein sequences (Additional file 1) [25]. Additionally we used the Blast2GO software to assign descriptions of gene function and also gene ontology categories [26]. The entire set of predicted protein sequences from OGS2 were BLASTed against the GenBank non-redundant protein database with results filtered for $E < 10^{-5}$, and also queried against the InterPro functional prediction pipeline [27]. Functional enrichment of Gene Ontology (GO) terms present in the sperm proteome relative to the genomic background was performed using Blast2GO's implementation of a Fisher's exact test with a false discovery rate of 0.01%.

Orthology predictions and analysis

Two approaches were employed for establishing orthology among proteins from different species. First, we used the proteinortho pipeline [28] to assess 3-way orthology between *D. plexippus* OGS2, *M. sexta* OGS1 [29], and *D. melanogaster* (flybase r6.12) gene sets. Proteinortho uses a reciprocal blast approach ($>50\%$ query coverage and $>25\%$ amino acid identity) to group genes with significant sequence similarity into clusters to identify orthologs and paralogs. For each species, genes with multiple protein isoforms were represented by the longest sequence in the proteinortho analysis. *D. melanogaster* and *M. sexta* ortholog predictions were then cross referenced to the published sperm of these two species [9, 30], allowing a three-way assessment of orthology in relation to presence in the sperm proteome. Using proteinortho allowed the direct analysis of the *D. plexippus* OGS2 sequences, which were not analyzed for homology in OrthoDB8 [31]. Potential annotation errors in the Monarch genome were investigated by identifying orthologs between Monarch and *Drosophila* which differed in length by at least 35%. These orthologs were manually curated using BLAST searches against available Lepidoptera and *Drosophila* genes to distinguish putative cases of misannotation from bona fide divergence in length.

A taxonomically broader set of insect ortholog relationships was obtained from OrthoDB8 and used to

assess the proportion of orthologs among sperm proteins relative to the genomic background. A randomized sampling procedure was used to determine the null expectation for the proportion of orthologous proteins found between *D. plexippus* and the queried species. A set of 584 proteins, the number equal to detected *D. plexippus* sperm proteins, was randomly sampled 5000 times from the entire Monarch OGS2 gene set. For each sample, the proportion of genes with an ortholog reported in OrthoDB8 was calculated, yielding a null distribution for the proportion of orthologs expected between *D. plexippus* and the queried species. For each query species, the observed proportion of orthologs in the sperm proteome was compared to this null distribution to determine whether the sperm proteome had a different proportion of orthologs than expected and to assign significance. Comparisons were made to 12 other insect species, reflecting five insect orders: Lepidoptera (*Heliconius melpomene*, *M. sexta*, *Plutella xylostella*, *Bombyx mori*), Diptera (*Drosophila melanogaster*, *Anopheles gambiae*), Hymenoptera (*Apis mellifera*, *Nasonia vitripennis*), Coleoptera (*Tribolium castaneum*, *Dendroctonus ponderosae*), and Hemiptera (*Acyrtosiphon pisum*, *Cimex lectularius*).

Maximum likelihood phylogenetic analysis

The phylogenetic relationships (i.e. topology) among the 13 taxa considered here were taken from [32] (for Lepidoptera) and from [33] (among insect orders). Branch lengths for this topology were determined using maximum likelihood optimization with amino acid sequence data. Thirteen nuclear genes were selected from the set of 1-to-1 orthologous loci provided by the BUSCO Insecta listing from OrthoDB version 9 [34]. Genes were chosen for completeness among the focal species analyzed. The genes used in this analysis correspond to the following OrthoDB9 ortholog groups: EOG090W0153, EOG090W01JK, EOG090W059K, EOG090W05WH, EOG090W06ZM, EOG090W08E4, EOG090W08ZA, EOG090W09XZ, EOG090W0E59, EOG090W0EIQ, EOG090W0F8Q, EOG090W0JMT, EOG090W0JXV. Amino acid sequences were aligned using MUSCLE, with default parameters as implemented in the R package, “msa” [35]. Each alignment was then filtered with Gblocks to remove regions or poor alignment and low representation [36]. After filtering, the alignments yielded a total of 2618 amino acid positions for maximum likelihood analysis. Filtered alignments were concatenated and used as a single dataset for branch length estimation via the R package “phangorn” [37]. Model test comparisons for transition rate matrices were performed, with the optimal model (LG + gamma + invariant class) used for branch length optimization via the “pml.optim” function.

Phylogenetic distribution of sperm proteins

The taxonomic distribution of sperm proteins was determined by BLASTp analyses (statistical cut off of $e < 10^{-5}$ and query coverage of $\geq 50\%$) against the protein data sets of the following taxonomic groupings: butterflies (*Heliconius melpomene*, *Papilio xuthus*, *Lerema accius*), Lepidoptera (Butterflies with *M. sexta*, *Amyleios transittella*, and *Plutella xylostella*), Mecoptera (Lepidoptera with *D. melanogaster*), Mecoptera with *Tribolium castaneum*, and Insecta (all previous taxa as well as: *Apis mellifera*, *Pediculus humanus*, *Acyrtosiphon pisum*, and *Zootermopsis nevadensis*). Lepidopteran species were chosen to maximize species distribution across the full phylogenetic breadth of Lepidoptera, while also utilizing the most comprehensively annotated genomes based on published CEGMA scores (<http://lepbase.org>, [38]). Taxonomically restricted proteins were defined as those identified repeatedly across a given phylogenetic range but without homology in any outgroup species. Proteins exhibiting discontinuous phylogenetic patterns of conservation were considered unresolved.

Maximum likelihood analysis of molecular evolution

Orthology information for the four available *Papilionoidea* was obtained from OrthoDB v9 [39]. Coding sequences corresponding to protein entries for all orthology groups were obtained from Ensembl release 86 for *H. melpomene* and *M. cinxia*, and from lepbase v4 for *D. plexippus* and *P. glaucus*. Translated protein sequences were aligned using the linsi algorithm of MAFFT [40] and reverse translated in frame. Whole phylogeny estimates of dN and dS were obtained using the M1 model as implemented by the PAML software package [41]. Allowing for the absence of no more than one species, evolutionary analyses were conducted for a total of 10,258 orthology groups. Kolmogorov-Smirnov tests were used to compare the distribution of dN between groups of genes; dS was not utilized in these comparisons because synonymous sites were found to be saturated between all of the sequenced *Papilionoidea* genomes. Rapidly evolving sperm proteins were also identified as those in the top 5% of proteins based on dN after the removal of outliers exceeding twice the interquartile range genome-wide.

Results

Monarch sperm proteome

Characterization of the Monarch sperm proteome as part of this study, in conjunction with our previous analysis in *Manduca* [9], allowed us to conduct the first comparative analysis of sperm in Lepidoptera, and in insects more broadly, to begin to assess the origin and evolution of dichotomous spermatogenesis at the genomic level. Tandem mass spectrometry (MS/MS)

analysis of Monarch sperm, purified in triplicate, identified 240 in all three replicates, 140 proteins in two replicates and 553 proteins identified by two or more unique peptides in at least a single replicate. Together this yielded a total of 584 high confidence protein identifications (Additional file 2). Of these, 41% were identified in all three biological replicates. Comparable with our previous analysis of *Manduca* sperm, proteins were identified by an average of 7.9 unique peptides and 21.1 peptide spectral matches. This new dataset thus provides the necessary foundation to refine our understanding of sperm composition at the molecular level in Lepidoptera. (Note: *Drosophila melanogaster* gene names will be used throughout the text where orthologous relationships exist with named genes; otherwise Monarch gene identification numbers will be used.)

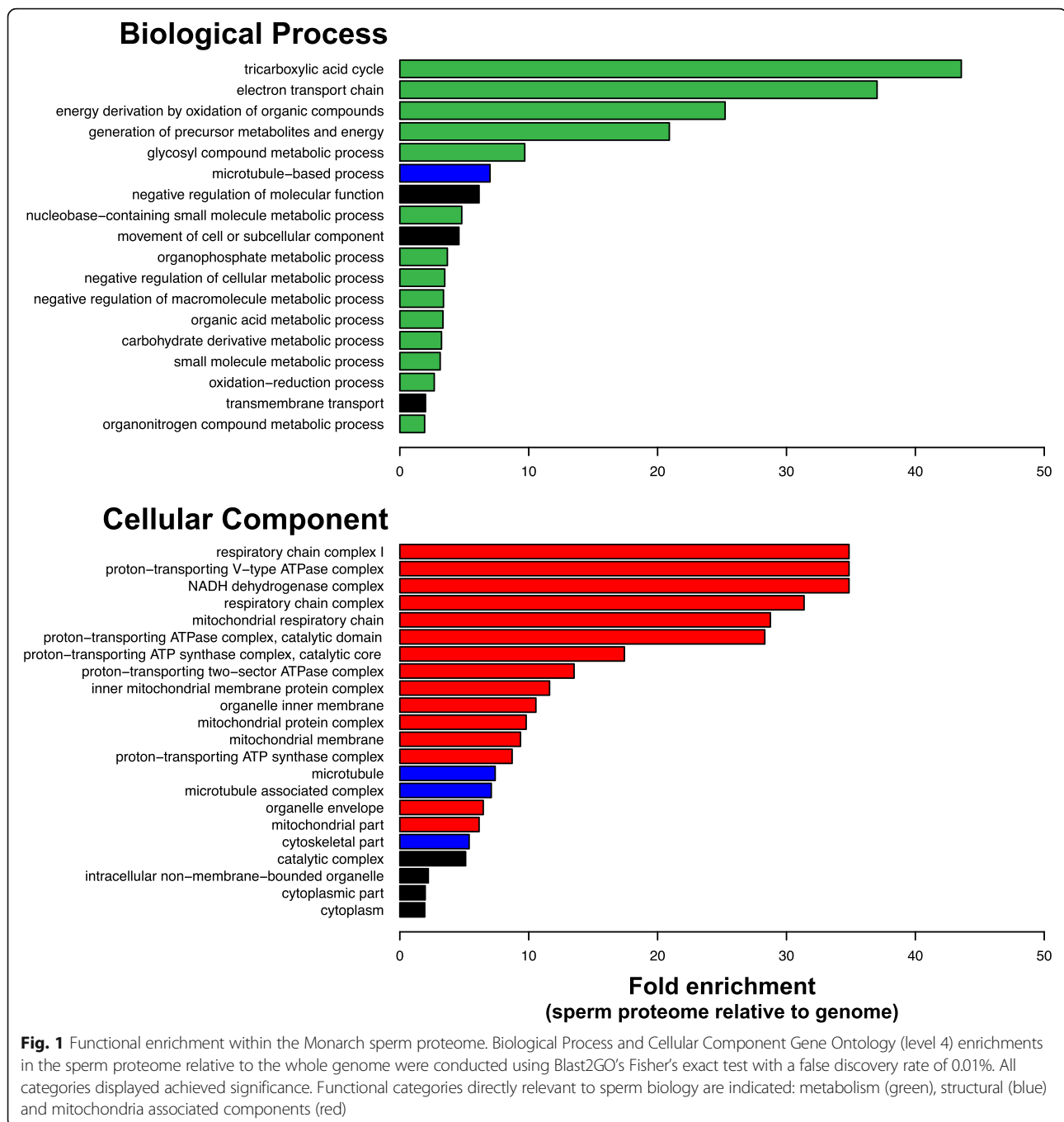
Gene ontology analysis of molecular composition

Gene ontology (GO) analyses were first conducted to confirm the similarity in functional composition between the Monarch and other insect sperm proteomes. Analysis of Biological Process terms revealed a significant enrichment for several metabolic processes, including the tricarboxylic acid (TCA) cycle ($p = 2.22E-16$), electron transport chain ($p = 9.85E-18$), oxidation of organic compounds ($p = 1.33E-25$) and generation of precursor metabolites and energy ($p = 1.09E-30$) (Fig. 1a). GO categories related to the TCA cycle and electron transport have also been identified as enriched in the *Drosophila* and *Manduca* sperm proteomes [9]. Generation of precursor metabolites and energy, and oxidation of organic compounds are also the two most significant enriched GO terms in the *Drosophila* sperm proteome [30]. Thus, broad metabolic functional similarities exist between the well-characterized insect sperm proteomes.

An enrichment of proteins involved in microtubule-based processes was also observed, a finding that is also consistent with previously characterized insect sperm proteomes. Amongst the proteins identified are cut up (ctp), a dynein light chain required for spermatogenesis [42], actin 5 (Act5), which is involved in sperm individualization [43], and DPOGS212342, a member of the recently expanded X-linked *tektin* gene family in *Drosophila* sperm [44]. Although functional annotations are limited amongst the 10% most abundant proteins (see below), several contribute to energetic and metabolic pathways. For example, stress-sensitive B (sesB) and adenine nucleotide translocase 2 (Ant2) are gene duplicates that have been identified in the *Drosophila* sperm proteome and, in the case of Ant2, function specifically in mitochondria during spermatogenesis [45]. Also identified was Bellwether (blw), an ATP synthetase alpha chain which is required for spermatid development [46].

The widespread representation of proteins functioning in mitochondrial energetic pathways is consistent with the contribution of giant, fused mitochondria (i.e. nebenkern) in flagellum development and presence of mitochondrial derivatives in mature spermatozoa (Fig. 1a-b) [47]. In lepidopteran spermatogenesis, the nebenkern divides to form two derivatives, which flank the axoneme during elongation; ultrastructure and size of these derivatives varies greatly between species and between the two sperm morphs [7]. In *Drosophila*, the nebenkern acts as both an organizing center for microtubule polymerization and a source of ATP for axoneme elongation, however it is unclear to what extent these structures contribute to energy required for sperm motility. Of particular note is the identification of porin, a voltage-gated anion channel that localizes to the nebenkern and is critical for sperm mitochondrion organization and individualization [48]. Consistent with these patterns, Cellular Component analysis also revealed a significant enrichment of proteins in a broad set of mitochondrial structures and components, including the respiratory chain complex I ($p = 7.73E-09$), proton-transporting V-type ATPase complex ($p = 9.90E-08$) and the NADH dehydrogenase complex ($p = 7.73E-09$) (Fig. 1b). Aside from those categories relating to mitochondria, a significant enrichment was also observed amongst categories relating to flagellum structure, including microtubule ($p = 5.43E-18$) and cytoskeleton part ($p = 2.54E-12$). These GO categories included the two most abundant proteins in the proteome identified in both Monarch and *Manduca*, beta tubulin 60D (β Tub60D) and alpha tubulin 84B (α Tub84B). α Tub84B is of particular interest as it performs microtubule functions in the post-mitotic spermatocyte, including the formation of the meiotic spindle and sperm tail elongation [49].

Molecular Function GO analysis revealed an enrichment of oxidoreductase proteins acting on NAD(P)H ($p = 7.06E-19$), as well as more moderate enrichments in several categories relating to peptidase activity or regulation of peptidase activity (data not shown). The broad representation of proteins involved in proteolytic activity is worthy of discussion, not solely because these classes of proteins are abundant in other sperm proteomes, but also because proteases are involved in the breakdown of the fibrous sheath surrounding Lepidoptera eupyrene sperm upon transfer to the female [7]. This process has been attributed to a specific ejaculatory duct trypsin-like arginine C-endopeptidase (initiatorin) in the silkworm (*B. mori*) [50] and a similar enzymatic reaction is needed for sperm activation in *Manduca* [51]. Blast2GO analyses identified three serine-type proteases in the top 5% of proteins based on abundance, including a chymotrypsin peptidase (DPOGS213461) and a trypsin



precursor (DPOGS205340). These highly abundant proteases, particularly those that were also identified in *Manduca* (two of the most abundant proteases and 10 in total), are excellent candidates for a sperm activating factor(s) in Lepidoptera.

Conservation of Lepidoptera sperm proteomes

Our previous analysis of *Manduca* was the first foray into the molecular biology of Lepidopteran sperm and was motivated by our interest in the intriguing heteromorphic

sperm system that is found in nearly all species in this order [7]. Here we have aimed to delineate the common molecular components of lepidopteran sperm through comparative analyses. Orthology predictions between the two species identified relationships for 405 (69%) Monarch sperm proteins, of which 369 (91%) were within "one-to-one" orthology groups (Additional file 2). 298 of all orthologs (73.5%) were previously identified by MS/MS in the *Manduca* sperm proteome [9]. An identical analysis in *Drosophila* identified 203 (35%) Monarch proteins with

orthology relationships, of which 166 (82%) were within “one-to-one” orthology groups (Additional file 2). 107 (52.7%) were previously characterized as components of the *Drosophila* sperm proteome [30, 52]. Thus there is a significantly greater overlap in sperm components between the two Lepidopteran species (two tailed Chi-square = 25.55, d.f. = 1, $p < 0.001$), as would be expected given the taxonomic relationship of these species. Additionally, gene duplication does not appear to be a widespread contributor to divergence relating to sperm form or function between Lepidoptera and *Drosophila*. It is also noteworthy that 27 orthologous proteins between Monarch and *Drosophila* were identified that differed substantially in length (>35%). Additional comparative analyses with gene models in other available Lepidoptera and *Drosophila* genomes indicated that 17 of these cases represent bona fide divergence in gene length, while the remainder are likely to represent gene model annotation errors in the Monarch genome. These issues were most commonly the result of inclusion/exclusion of individual exons with adjacent gene models and full gene model fusions (Additional file 2).

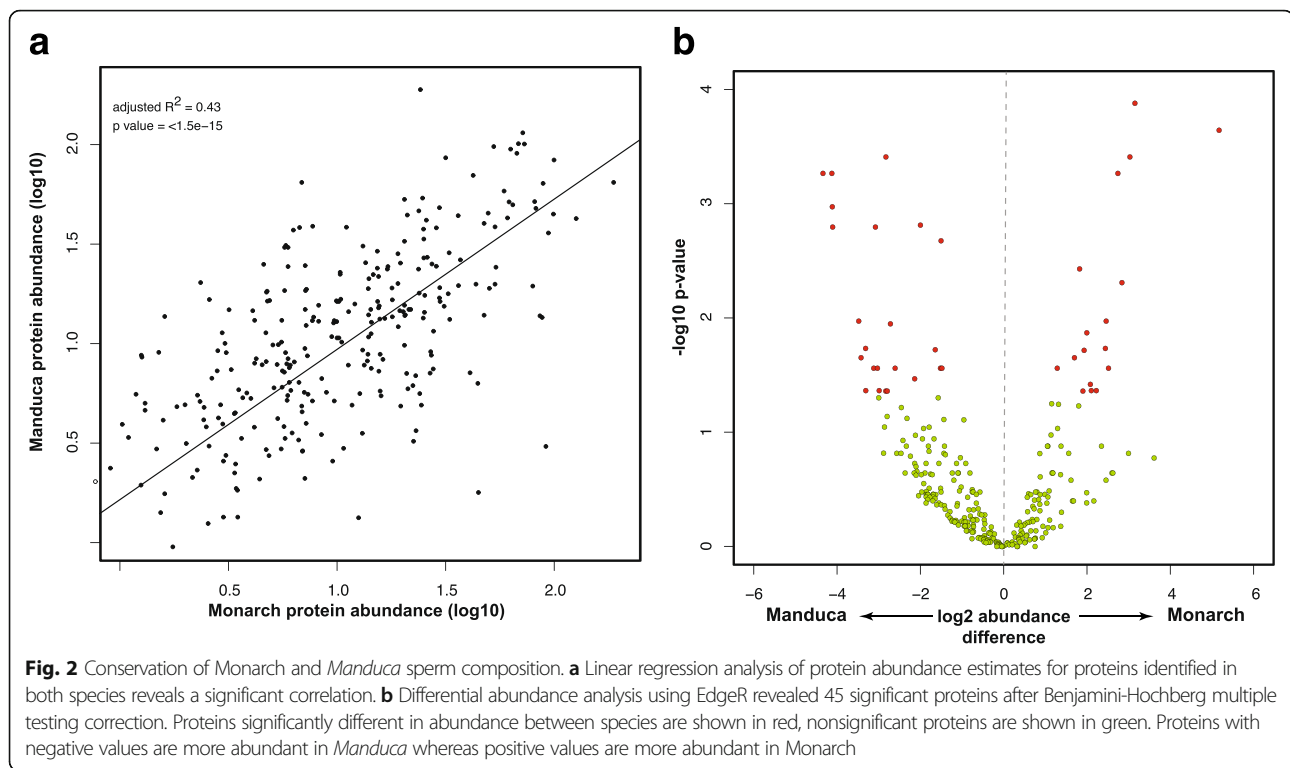
Recent comparative analyses of sperm composition across mammalian orders successfully identified a conserved “core” sperm proteome comprised of more slowly evolving proteins, including a variety of essential structural and metabolic components. To characterize the “core” proteome in insects, we conducted a GO analysis using *Drosophila* orthology, ontology and enrichment data to assess the molecular functionality of the 92 proteins identified in the proteome of all three insect species. This revealed a significant enrichment for proteins involved in cellular respiration ($p = 4.41e-21$), categories associated with energy metabolism, including ATP metabolic process ($p = 1.64e-15$), generation of precursor metabolites and energy ($p = 9.77e-21$), and multiple nucleoside and ribonucleoside metabolic processes. Analysis of cellular component GO terms revealed a significant enrichment for mitochondrion related proteins ($p = 3.72e-22$), respiratory chain complexes ($p = 8.25e-12$), dynein complexes ($p = 1.37e-5$), and axoneme ($p = 3.31e-6$). These GO category enrichments are consistent with a core set of metabolic, energetic, and structural proteins required for general sperm function. Similar sets of core sperm proteins have been identified in previous sperm proteome comparisons [9, 30, 52, 53]. Among this conserved set are several with established reproductive phenotypes in *Drosophila*. This includes proteins associated with sperm individualization, including cullin3 (Cul3) and SKP1-related A (SkpA), which acts in cullin-dependent E3 ubiquitin ligase complex required for caspase activity in sperm individualization [54], gudu, an

Armadillo repeat containing protein [55], and porin (mentioned previously) [48]. Two proteins involved in sperm motility were also identified: dynein axonemal heavy chain 3 (dnah3) [56] and an associated microtubule-binding protein growth arrest specific protein 8 (Gas8) [57].

Comparative analysis of protein abundance

Despite the more proximate link between proteome composition and molecular phenotypes, transcriptomic analyses far outnumber similar research using proteomic approaches. Nonetheless, recent work confirms the utility of comparative evolutionary proteomic studies in identifying both conserved [58] and diversifying proteomic characteristics [59]. We have previously demonstrated a significant correlation in protein abundance between *Manduca* and *Drosophila* sperm, although this analysis was limited by the extent of orthology between these taxa [9]. To further investigate the evolutionary conservation of protein abundance in sperm, a comparison of normalized abundance estimates between Monarch and *Manduca* revealed a significant correlation ($R^2 = 0.43$, $p = < 1 \times 10^{-15}$) (Fig. 2a). We note that this correlation is based on semi-quantitative estimates [20] and would most likely be stronger if more refined absolute quantitative data were available. Several proteins identified as highly abundant in both species are worthy of further mention. Two orthologs of *Sperm leucyl aminopeptidases* (S-LAPs) were identified. S-LAPs are members of a gene family first characterized in *Drosophila* that has recently undergone a dramatic expansion, is testis-specific in expression and encodes the most abundant proteins in the *D. melanogaster* sperm proteome [60]. As would be expected, several microtubule structural components were also amongst the most abundant proteins (top 20), including α Tub84B and tubulin beta 4b chain-like protein, as well as succinate dehydrogenase subunits A and B (SdhA and SdhB), porin, and DPOGS202417, a trypsin precursor that undergoes conserved post translational modification (see below).

We next sought to identify proteins exhibiting differential abundance between the two species. As discussed earlier, Monarch and *Manduca* have distinct mating systems; female Monarch butterflies remate considerably more frequently than *Manduca* females, increasing the potential for sperm competition [10]. These differences may be reflected in molecular diversification in sperm composition between species. An analysis of differential protein abundance identified 45 proteins with significant differences after correction for multiple testing ($P < 0.05$; Fig. 2b), representing 7% of the proteins shared between species (Additional file 3). No directional bias was observed in the number of differentially abundant proteins (one-tail Binomial test; p



value = 0.2757). Several of these proteins are worthy of further discussion given their role in sperm development, function or competitive ability. Proteins identified as more abundant in the Monarch sperm proteome were heavily dominated by mitochondrial NADH dehydrogenase subunits (subunits ND-23, ND-24, ND-39, and ND-51) and other mitochondria-related proteins, including ubiquinol-cytochrome c reductase core protein 2 (UQCR-C2), cytochrome C1 (Cyt-C1), and glutamate oxaloacetate transaminase 2 (Got2). Additionally, two proteins with established sperm phenotypes were identified as more abundant in *Manduca*. These included dynein light chain 90F (Dlc90F), which is required for proper nuclear localization and attachment during sperm differentiation [61], and cut up (ctp), a dynein complex subunit involved in nucleus elongation during spermiogenesis [42]. Serine protease immune response integrator (spirit) is also of interest considering the proposed role of endopeptidases in Lepidoptera sperm activation [50, 51]. Although it would be premature to draw any specific conclusions, some of these proteins play important mechanistic roles in sperm development and function and will be of interest for more targeted functional studies.

Post-translational modification of sperm proteins

During spermatogenesis, the genome is repackaged and condensed on protamines and the cellular machinery required for protein synthesis are expelled.

Consequently, mature sperm cells are considered primarily quiescent [62]. Nonetheless, sperm undergo dynamic molecular transformations after they leave the testis and during their passage through the male and female reproductive tract [63]. One mechanism by which these modifications occur is via post translational modification (PTM), which can play an integral part in the activation of sperm motility and fertilization capacity [64, 65]. Analysis of PTMs in Monarch identified 438 acetylated peptides within 133 proteins. Most notable among these are microtubule proteins, including alpha tubulin 84B (alphaTub84B), beta tubulin 60D (betaTub60D) and dyneins kl-3 and kl-5. Tubulin is a well-known substrate for acetylation, including the highly-conserved acetylation of N-terminus Lysine 40 of alphaTub84B. This modification is essential for normal sperm development, morphology and motility in mice [66]. A similar analysis in *Manduca* identified 111 acetylated peptides within 63 proteins. We found evidence for conserved PTMs within Lepidoptera in 19 proteins (36% of those identified in Monarch), including Lys40 of alphaTub84B.

In contrast to acetylation, only 75 Monarch sperm proteins showed evidence of phosphorylation, 53 of which were also modified in *Manduca* (71%). This included the ortholog of the Y-linked *Drosophila* gene WDY. Although a specific function for WDY in spermatogenesis has yet to be determined, WDY is expressed in a testis-specific

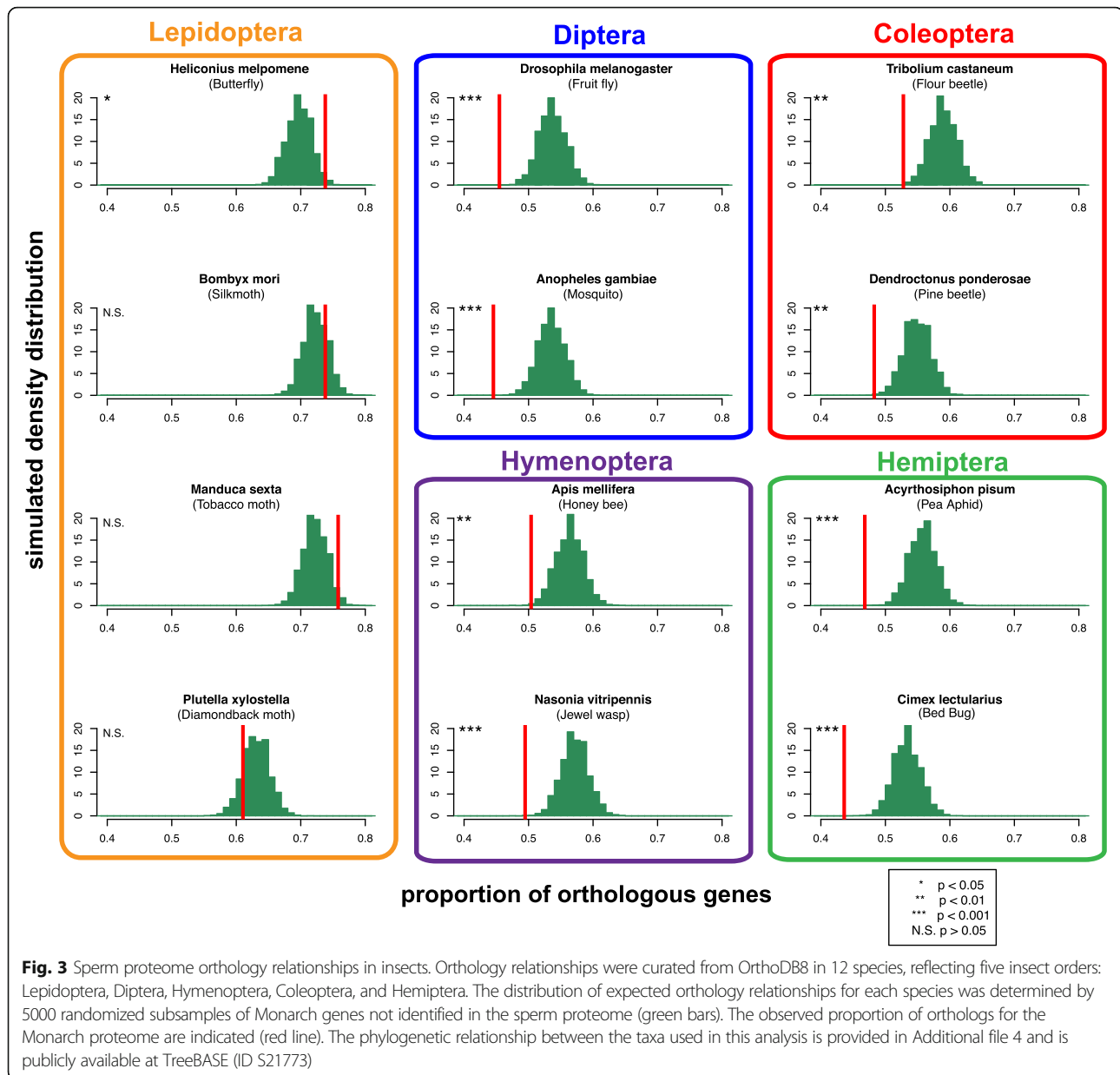
manner and under positive selection in the *D. melanogaster* group [67]. The relative paucity of phosphorylation PTMs may reflect the fact that phosphorylation is one of the more difficult PTMs to identify with certainty via mass spectrometry based proteomics [68]. However, it is also noteworthy that sperm samples in this study were purified from the male seminal vesicle, and thus, before transfer to the female reproductive tract. Although far less is known about the existence of capacitation-like processes in insects, dynamic changes in the mammalian sperm phosphoproteome are associated with sperm capacitation and analogous biochemical alterations might occur within the female reproductive tract of insects [65]. We note that a similar extent of protein phosphorylation has been detected from *Drosophila* sperm samples purified in a similar manner (unpublished data; Whittington and Dorus). Lastly, identical acetylation and phosphorylation PTM patterns were identified for Monarch and *Manduca* HACP012 (DPOGS213379), a putative seminal fluid protein of unknown function previously identified in the Postman butterfly (*Heliconius melpomene*) [69, 70]. The identification of HACP012 in sperm, in the absence of other seminal fluid components, is unexpected but its identification was unambiguous as it was amongst the most abundant 10% of identified Monarch proteins. Seminal protein HACP020 (DPOGS203866), which exhibits signatures of recent adaptive evolution [70], was also identified as highly abundant (5th percentile overall); this suggests that some seminal fluid proteins may also be co-expressed in the testis and establish an association with sperm during spermatogenesis.

Rapid evolution of genetic architecture

Rapid gene evolution [71] and gene gain /loss [72], including de novo gene gain [73], are predominant processes that contribute to the diversification of male reproductive systems. Our previous study identified an enrichment in the number of Lepidoptera specific proteins (i.e. those without homology outside of Lepidoptera) in the sperm proteome relative to other reproductive proteins and non-reproductive tissues. We were unable, however, to determine from a single species whether novel genes contributed to sperm biology more broadly across all Lepidoptera. Here we employed two comparative genomic approaches to confirm and expand upon our original observation. First, we obtained whole-genome orthology relationships between Monarch and nine species, representing five insect orders, and compared the proportion of the sperm proteome with orthologs to the whole genome using a random subsampling approach. No significant differences were observed for three of the four Lepidoptera species analyzed and an excess of orthology amongst sperm proteins was identified in

the Postman butterfly ($p < 0.05$; Fig. 3). In contrast, we identified a significant deficit of sperm orthologs in all comparisons with non-Lepidopteran genomes (all $p < 0.01$). Orthology relationships in OrthoDB are established by a multi-step procedure involving reciprocal best match relationships between species and identity within species to account for gene duplication events since the last common ancestor. As such, the underrepresentation of orthology relationships is unlikely to be accounted for by lineage-specific gene duplication. Therefore, rapid evolution of sperm genes appears to be the most reasonable explanation for the breakdown of reciprocal relationships (see below). This conclusion is consistent with a diverse body of evidence that supports the influence of positive selection on male reproductive genes [71, 74], including those functioning in sperm [52, 75–78]. We note that we cannot rule out the influence of de novo gain but it is currently difficult to assess the contribution of this mechanism to the overall pattern.

The second analysis aimed to characterize the distribution of taxonomically restricted Monarch sperm proteins using BLAST searches across 12 insect species. Based on the analysis above, our a priori expectation was that a substantial number of proteins with identifiable homology amongst Lepidoptera would be absent from more divergent insect species. This analysis identified a total of 45 proteins unique to Monarch, 140 proteins (23.9% of the sperm proteome) with no detectable homology to proteins in non-Lepidopteran insect taxa and 173 proteins conserved across all species surveyed (Fig. 4a). Proteins with discontinuous taxonomic matches ($n = 171$) were considered “unresolved”. Although the number of Monarch-specific proteins is considerably higher than the eight *Manduca*-specific proteins found in our previous study, the number of Lepidoptera specific is comparable to our previous estimate in *Manduca* ($n = 126$). These observations support the hypothesis that a substantial subset of lepidopteran sperm proteins are likely rapidly evolving and thus exhibit little detectable similarity. To pursue this possibility, we calculated nonsynonymous divergence (dN) for 10,212 genes across four species of butterfly and compared dN between Lepidoptera specific sperm proteins, sperm proteins with homology outside of Lepidoptera and the remainder of the genome (Fig. 4b). The average dN of Lepidoptera specific proteins was significantly higher than non-Lepidopteran specific proteins ($D = 0.34$, $p = 5.0 \times 10^{-9}$) and the remainder of the genome ($D = 0.28$, $p = 1.23 \times 10^{-7}$). Interestingly, sperm proteins with homology outside of Lepidoptera also evolve significantly slower than the genome as whole ($D = 0.30$, $p = 3.14 \times 10^{-6}$). Consistent with these trends, 17.7% of Lepidoptera specific sperm proteins were amongst the fastest evolving in the genome (top 5%), compared to only



2.6% of sperm proteins with homology outside of Lepidoptera. In light of the rapid divergence of Lepidoptera specific proteins we next sought to assess their potential contribution to sperm function using protein abundance as a general proxy in the absence of functional annotation for nearly all of these proteins. As was observed in Whittington et al. [9], Lepidopteran specific proteins were found to be significantly more abundant than the remainder of the sperm proteome ($D = 0.20$, $p = 0.0009$, Fig. 4c).

Discussion

Dichotomous spermatogenesis in Lepidoptera, and in particular the production of sperm which do not fertilize

oocytes, has intrigued biologists for over a century. Despite widespread interest, little is known about the functional roles fulfilled by apyrene sperm or why they have been retained in a nearly ubiquitous fashion during the evolution of Lepidoptera. Our comparative proteomic analysis of heteromorphic sperm, a first of its kind, provides important perspective and insights regarding the functional and evolutionary significance of this enigmatic reproductive phenotype. First, our analyses indicate that a substantial number of novel sperm genes are shared amongst Lepidoptera, thus distinguishing them from other insect species without dichotomous spermatogenesis, and suggest they are associated with heteromorphic spermatogenesis and the diversification

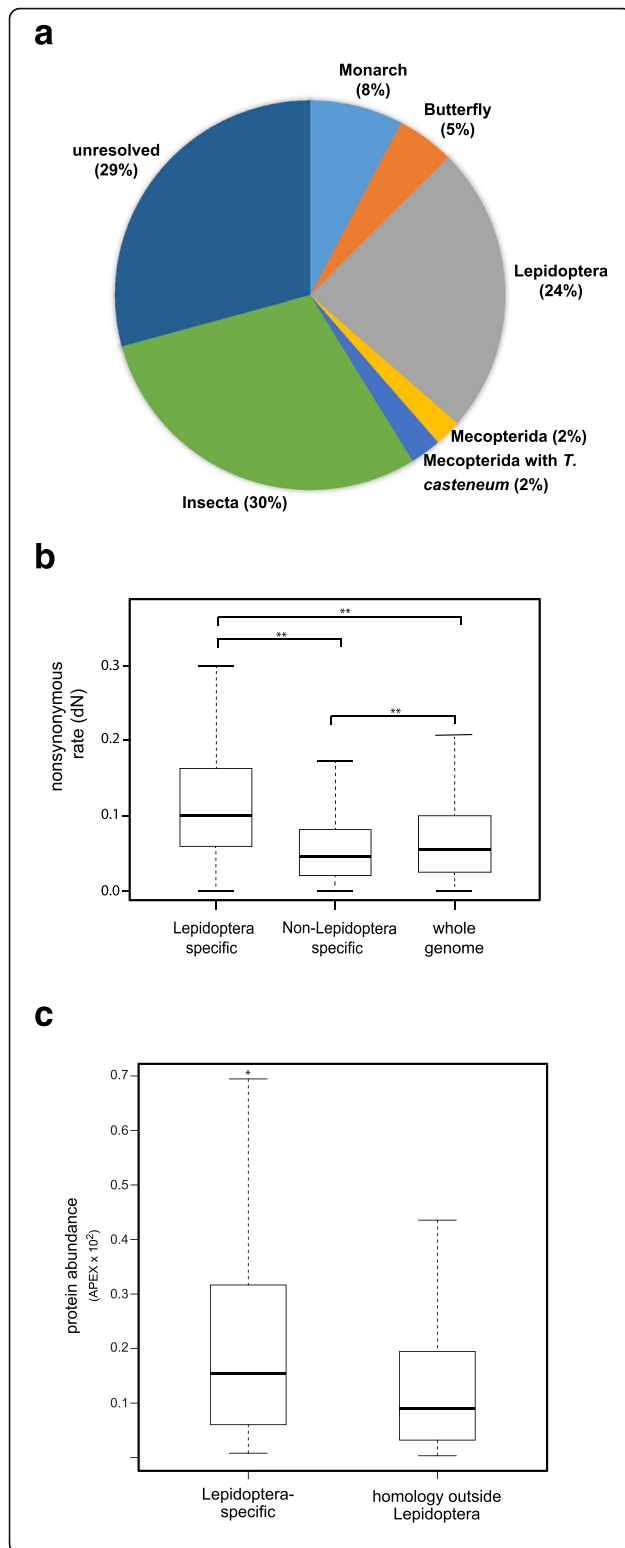


Fig. 4 Taxonomic distribution and evolution of Monarch sperm proteins. **a** Pie chart displaying the taxonomic distribution of proteins homologous to the Monarch sperm proteome and those unique to Monarch. BLAST searches were conducted beginning with closely related butterfly species and sequentially through more divergent species in Mecoptera, Mecoptera plus *Tribolium*, and Insecta. In order to be considered Lepidoptera specific, a protein was required to be present in at least one butterfly other than Monarch and at least one moth species. Proteins with discontinuous taxonomic patterns of homology are included in the category “unresolved”. **b** Box plot showing nonsynonymous divergence (dN) of Monarch proteins across four species of butterfly ($n = 10,212$). Nonsynonymous divergence for sperm proteins identified as specific to Lepidoptera, sperm proteins with homology outside of Lepidoptera and the remainder of the genome are shown. Asterisks (**) indicate p -values less than 1.0×10^{-5} . **c** Box plot displaying the distribution of protein abundance estimates for proteins present only in Lepidoptera and those with homology in other insects. Asterisk (*) indicate p -values less than 0.001

of apyrene and eupyrene sperm. This observation can be attributed, at least in part, to the rapid evolution of Lepidoptera specific sperm genes. It is also possible that de novo gene gain may contribute to this observed genetic novelty, although it is not possible to assess this directly with the genomic and transcriptomic resources currently available in Lepidoptera. Our comparative and quantitative analyses, based on protein abundance measurements in both species, further suggests that some of these proteins contribute to apyrene sperm function and evolution. Given that apyrene sperm constitute the vast majority of cells in our co-mixed samples, it is reasonable to speculate that higher abundance proteins are either present in both sperm morphs or specific to apyrene cells. Confirmation of this will require targeted proteomic analysis of purified apyrene and eupyrene cell populations and will result in a refined set of candidates for further study in relation to apyrene sperm functionality. Ultimately, the comparative analysis of morph-specific sperm proteomes is critical to understanding the functional diversification of the fertilization incompetent apyrene sperm morph and the evolutionary maintenance of dichotomous spermatogenesis.

Conclusion

Our results indicate that the origin of heteromorphic spermatogenesis early in Lepidoptera evolution and/or the subsequent evolution of this system is associated with a burst of genetic novelty that is distinct from patterns of diversification across the remainder of the genome. The evolution of dichotomous spermatogenesis has therefore had a marked impact on Lepidoptera molecular evolution and suggests that focused studies of other reproductive transitions may inform our broader understanding of the evolution of reproductive genetic systems and their contribution to genomic novelty.

Additional files

Additional file 1: Functional Information- Predicted functions of Monarch proteins curated using Uniprot and Blast2GO. (XLSX 29 kb)

Additional file 2: Mass Spectrometry Data- Proteomic data including full Monarch sperm proteome, MS/MS results by replicate, PTM information, orthology relationships and rates of molecular evolution. (XLSX 100 kb)

Additional file 3: Protein Abundance and Quantitative Analyses- APEX protein abundance estimates and differential abundance analyses. (XLSX 54 kb)

Additional file 4: Phylogenetic Results- Phylogeny exhibiting the evolutionary relationship of the thirteen insect species utilized in this study. (PDF 129 kb)

Abbreviations

CDS: Coding Sequence; FDR: False Discovery Rate; GO: Gene Ontology; HCD: Higher energy Collisional Dissociation; LC: Liquid Chromatography; LC-MS/MS: Liquid Chromatography Tandem Mass Spectrometry; MS/MS: Tandem Mass Spectrometry; OGS1: Official Gene Set 1; OGS2: Official Gene Set 2; PTM: Post Translational Modification

Acknowledgements

We thank Monarch Watch and Channing Shives for support in rearing Monarch butterflies and Sheri Skerget for expert technical assistance. Computing for this project was performed on the Syracuse University Crush Virtual Research Cloud and the Community Cluster at the Center for Research Computing at the University of Kansas. We would also like to thank Eric Sedore and Larne Pekowsky of Information Technology Services at Syracuse University and Mike Deery, Renata Feret and Kathryn Lilley at the Cambridge Centre for Proteomics.

Funding

Funding for this study included Syracuse University support to SD, University of Kansas support to JW, and Syracuse University and Marilyn Kerr Fellowships to EW.

Availability of data and materials

Mass spectrometry data is publicly available through the ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) with the dataset identifier PXD006454. Phylogenetic results are publicly available through TreeBASE (<http://treebase.org/treebase-web/home.html>) with the identifier S21773.

Authors' contributions

TLK, JW and SD designed the study; TLK and JW purified samples for MS analysis; EW, DF, KB, JW and SD analyzed the data; EW, TLK, JW and SD wrote the manuscript. All authors have read and approved the final version of this manuscript.

Ethics approval and consent to participate

Not applicable. No field permissions were required.

Consent for publication

Not applicable

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Center for Reproductive Evolution, Department of Biology, Syracuse University, Syracuse, NY, USA. ²Science Education and Society, University of Rhode Island, Kingston, RI, USA. ³Ecology and Evolutionary Biology, Kansas University, Lawrence, KS, USA. ⁴Department of Genomics and Genetic

Resources, Kyoto Institute of Technology, Saga Ippon-cho, Ukyo-ku, Kyoto, Japan.

Received: 30 May 2017 Accepted: 13 November 2017

Published online: 02 December 2017

References

- Pitnick S, Birkhead TR, Hosken DJ. Sperm biology: an evolutionary perspective. 1st ed. Amsterdam: Academic Press/Elsevier; 2009.
- Till-Bottraud I, Joly D, Lachaise D, Snook RR. Pollen and sperm heteromorphism: convergence across kingdoms? *J Evol Biol.* 2005;18:1–18.
- Swallow JG, Wilkinson GS. The long and short of sperm polymorphisms in insects. *Biol Rev Camb Philos Soc.* 2002;77:153–82.
- Meves F. Ueber oligopyrene und apyrene Spermien und über ihre Entstehung, nach Beobachtungen an Paludina und Pygaera. *Arch Für Mikrosk Anat.* 1902;61:1–84.
- Sahara K, Kawamura N. Double copulation of a female with sterile diploid and polyploid males recovers fertility in *Bombyx mori*. *Zygote Camb Engl.* 2002;10:23–9.
- Friedländer M. Control of the eupyrene–apyrene sperm dimorphism in Lepidoptera. *J Insect Physiol.* 1997;43:1085–92.
- Friedländer M, Seth RK, Reynolds SE. Eupyrene and Apyrene sperm: dichotomous spermatogenesis in Lepidoptera. *Adv Insect Physiol.* 2005;32:206–308.
- Snook RR, Hosken DJ, Karr TL. The biology and evolution of polypermery: insights from cellular and functional studies of sperm and centrosomal behavior in the fertilized egg. *Reproduction.* 2011;142:779–92.
- Whittington E, Zhao Q, Borziak K, Walters JR, Dorus S. Characterisation of the *Manduca sexta* sperm proteome: genetic novelty underlying sperm composition in Lepidoptera. *Insect Biochem Mol Biol.* 2015;62:183–93.
- Oberhauser K, Frey D. Coercive mating by overwintering male monarch butterflies. In: Hoth J, Merino L, Oberhauser K, Pisanty I, Price S, Wilkinson T, editors. 1997 North American Conference on the Monarch Butterfly. Canada: Commission for Environmental Cooperation. 1997. pp. 67–78.
- Sasaki M, Riddiford LM. Regulation of reproductive behaviour and egg maturation in the tobacco hawk moth, *Manduca sexta*. *Physiol Entomol.* 1984;9:315–27.
- Stringer IAN, Giebultowicz JM, Riddiford LM. Role of the bursa copulatrix in egg maturation and reproductive behavior of the tobacco hawk moth, *Manduca sexta*. *Int J Invertebr Reprod Dev.* 1985;8:83–91.
- Solensky MJ, Oberhauser KS. Male monarch butterflies, *Danaus plexippus*, adjust ejaculates in response to intensity of sperm competition. *Anim Behav.* 2009;77:465–72.
- Karr TL, Walters JR. Panning for sperm gold: isolation and purification of apyrene and eupyrene sperm from lepidopterans. *Insect Biochem Mol Biol.* 2015;63:152–8.
- Vizcaino JA, Csordas A, del-Toro N, Dianas JA, Griss J, Lavidas I, et al. 2016 update of the PRIDE database and its related tools. *Nucleic Acids Res.* 2016;44:11033.
- Deutsch EW, Mendoza L, Shteynberg D, Farrah T, Lam H, Tasman N, et al. A guided tour of the trans-proteomic pipeline. *Proteomics.* 2010;10:1150–9.
- Keller A, Nesvizhskii AI, Kolker E, Aebersold R. Empirical statistical model to estimate the accuracy of peptide identifications made by MS/MS and database search. *Anal Chem.* 2002;74:5383–92.
- Nesvizhskii AI, Keller A, Kolker E, Aebersold R. A statistical model for identifying proteins by tandem mass spectrometry. *Anal Chem.* 2003;75:4646–58.
- Shteynberg DD, Mendoza L, Slagel J, Lam H, Nesvizhskii AI, Moritz R. PTMProphet: TPP software for validation of modified site locations on post-translationally modified peptides. 60th American Society for Mass Spectrometry (ASMS) Annual Conference, Vancouver, Canada, 2012.
- Braisted JC, Kuntumalla S, Vogel C, Marcotte EM, Rodrigues AR, Wang R, et al. The APEX quantitative proteomics tool: generating protein quantitation estimates from LC-MS/MS proteomics results. *BMC Bioinformatics.* 2008;9:529.
- Robinson MD, McCarthy DJ, Smyth GK. edgeR: a bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics.* 2010;26:139–40.
- Zhan S, Merlin C, Boore JL, Reppert SM. The monarch butterfly genome yields insights into long-distance migration. *Cell.* 2011;147:1171–85.
- Zhan S, Reppert SM. MonarchBase: the monarch butterfly genome database. *Nucleic Acids Res.* 2013;41:D758–63.
- Kent WJ. BLAT—the BLAST-like alignment tool. *Genome Res.* 2002;12:656–64.

25. The UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res.* 2015;43:D204–12.
26. Conesa A, Gotz S, Garcia-Gomez JM, Terol J, Talon M, Robles M. Blast2GO: a universal tool for annotation, visualization and analysis in functional genomics research. *Bioinformatics.* 2005;21:3674–6.
27. Zdobnov EM, Apweiler R. InterProScan—an integration platform for the signature-recognition methods in InterPro. *Bioinforma Oxf Engl.* 2001;17:847–8.
28. Lechner M, Findeiß S, Steiner L, Marz M, Stadler PF, Prohaska SJ. Proteinortho: detection of (co-)orthologs in large-scale analysis. *BMC Bioinformatics.* 2011;12:124.
29. Kanost MR, Arrese EL, Cao X, Chen Y-R, Chellapilla S, Goldsmith MR, et al. Multifaceted biological insights from a draft genome sequence of the tobacco hornworm moth, *Manduca sexta*. *Insect Biochem Mol Biol.* 2016;76:118–47.
30. Wasbrough ER, Dorus S, Hester S, Howard-Murkin J, Lilley K, Wilkin E, et al. The *Drosophila melanogaster* sperm proteome-II (DmSP-II). *J Proteome.* 2010;73:2171–85.
31. Waterhouse RM, Tegenfeldt F, Li J, Zdobnov EM, Kriventseva EV. OrthoDB: a hierarchical catalog of animal, fungal and bacterial orthologs. *Nucleic Acids Res.* 2013;41:D358–65.
32. Kawahara AY, Breinholt JW. Phylogenomics provides strong evidence for relationships of butterflies and moths. *Proc R Soc B Biol Sci.* 2014;281:20140970.
33. Misof B, Liu S, Meusemann K, Peters RS, Donath A, Mayer C, et al. Phylogenomics resolves the timing and pattern of insect evolution. *Science.* 2014;346:763–7.
34. Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics.* 2015;31:3210–2.
35. Bodenhofer U, Bonatesta E, Horejš-Kainrath C, Hochreiter S. msa: an R package for multiple sequence alignment. *Bioinformatics.* 2015;31:3997–9.
36. Castresana J. Selection of conserved blocks from multiple alignments for their use in phylogenetic analysis. *Mol Biol Evol.* 2000;17:540–52.
37. Schliep KP. Phangorn: phylogenetic analysis in R. *Bioinformatics.* 2011;27:592–3.
38. Challis RJ, Kumar S, Dasmahapatra KKK, Jiggins CD, Blaxter M. Lepbase: the Lepidopteran genome database. *bioRxiv 056994*; doi: 10.1101/056994.
39. Zdobnov EM, Tegenfeldt F, Kuznetsov D, Waterhouse RM, Simão FA, Ioannidis P, et al. OrthoDB v9.1: cataloging evolutionary and functional annotations for animal, fungal, plant, archaeal, bacterial and viral orthologs. *Nucleic Acids Res.* 2017;45:D744–9.
40. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol Biol Evol.* 2013;30:772–80.
41. Yang Z. PAML 4: Phylogenetic analysis by maximum likelihood. *Mol Biol Evol.* 2007;24:1586–91.
42. Joti P, Ghosh-Roy A, Ray K. Dynein light chain 1 functions in somatic cyst cells regulate spermatogonial divisions in *Drosophila*. *Sci Rep.* 2011;1:173.
43. Noguchi T. A role for actin dynamics in individualization during spermatogenesis in *Drosophila melanogaster*. *Development.* 2003;130:1805–16.
44. Dorus S, Freeman ZN, Parker ER, Heath BD, Karr TL. Recent origins of sperm genes in *Drosophila*. *Mol Biol Evol.* 2008;25:2157–66.
45. Terhzaz S, Cabrero P, Chintapalli VR, Davies S-A, Dow JAT. Mislocalization of mitochondria and compromised renal function and oxidative stress resistance in *Drosophila* *SexB* mutants. *Physiol Genomics.* 2010;41:33–41.
46. Castrillon DH, Gönczy P, Alexander S, Rawson R, Eberhart CG, Viswanathan S, et al. Toward a molecular genetic analysis of spermatogenesis in *Drosophila melanogaster*: characterization of male-sterile mutants generated by single P element mutagenesis. *Genetics.* 1993;135:489–505.
47. Tokuyasu KT. Dynamics of spermiogenesis in *Drosophila melanogaster*. VI. Significance of “onion” nebenkern formation. *J Ultrastruct Res.* 1975;53:93–112.
48. Park J, Kim Y, Choi S, Koh H, Lee S-H, Kim J-M, et al. *Drosophila* Porin/VDAC affects mitochondrial morphology. *PLoS One.* 2010;5:e13151.
49. Hutchens JA, Hoyle HD, Turner FR, Raff EC. Structurally similar *Drosophila* Alpha-tubulins are functionally distinct *in vivo*. *Mol Biol Cell.* 1997;8:481–500.
50. Osanai M, Kasuga H, Aigaki T. Induction of motility of apyrene spermatozoa and dissociation of Eupyrene sperm bundles of the silkworm, *Bombyx mori*, by initiatorin and trypsin. *Invertebr Reprod Dev.* 1989;15:97–103.
51. Friedländer M, Jeshtadi A, Reynolds SE. The structural mechanism of trypsin-induced intrinsic motility in *Manduca sexta* spermatozoa *in vitro*. *J Insect Physiol.* 2001;47:245–55.
52. Dorus S, Busby SA, Gerike U, Shabanowitz J, Hunt DF, Karr TL. Genomic and functional evolution of the *Drosophila melanogaster* sperm proteome. *Nat Genet.* 2006;38:1440–5.
53. Rettie EC, Dorus S. *Drosophila* sperm proteome evolution: insights from comparative genomic approaches. *Spermatogenesis.* 2012;2:213–23.
54. Arama E, Bader M, Rieckhof GE, Steller H. A ubiquitin ligase complex regulates caspase activation during sperm differentiation in *Drosophila*. *PLoS Biol.* 2007;5:e251. Bach E, editor
55. Cheng W, Ip YT, Xu Z. Gudu, an armadillo repeat-containing protein, is required for spermatogenesis in *Drosophila*. *Gene.* 2013;531:294–300.
56. Karak S, Jacobs JS, Kittelmann M, Spalthoff C, Katana R, Sivan-Loukianova E, et al. Diverse roles of axonemal dyneins in *Drosophila* auditory neuron function and mechanical amplification in hearing. *Sci Rep.* 2015;5:17085.
57. Yeh S-D, Chen Y-J, Chang ACY, Ray R, She B-R, Lee W-S, et al. Isolation and properties of *Gas8*, a growth arrest-specific gene regulated during male gametogenesis to produce a protein associated with the sperm motility apparatus. *J Biol Chem.* 2002;277:6311–7.
58. Bayram HL, Claydon AJ, Brownridge PJ, Hurst JL, Mileham A, Stockley P, et al. Cross-species proteomics in analysis of mammalian sperm proteins. *J Proteome.* 2016;135:38–50.
59. Vicens A, Borziak K, Karr TL, Roldan ERS, Dorus S. Comparative sperm proteomics in mouse species with divergent mating systems. *Mol Biol Evol.* 2017;34:1403–16.
60. Dorus S, Wilkin EC, Karr TL. Expansion and functional diversification of a leucyl aminopeptidase family that encodes the major protein constituents of *Drosophila* sperm. *BMC Genomics.* 2011;12:177.
61. Li MG, Serr M, Newman EA, Hays TS. The *Drosophila* tctex-1 light chain is dispensable for essential cytoplasmic dynein functions but is required during spermatid differentiation. *Mol Biol Cell.* 2004;15:3005–14.
62. Hecht NB. Molecular mechanisms of male germ cell differentiation. *BioEssays.* 1998;20:555–61.
63. McDonough CE, Whittington E, Pitnick S, Dorus S. Proteomics of reproductive systems: towards a molecular understanding of postmating, prezygotic reproductive barriers. *J Proteome.* 2016;135:26–37.
64. Baker MA, Hetherington L, Weinberg A, Naumovski N, Velkov T, Pelzing M, et al. Analysis of phosphopeptide changes as spermatozoa acquire functional competence in the epididymis demonstrates changes in the post-translational modification of Izumo1. *J Proteome Res.* 2012;11:5252–64.
65. Platt MD, Salicioni AM, Hunt DF, Visconti PE. Use of differential isotopic labeling and mass spectrometry to analyze capacitation-associated changes in the phosphorylation status of mouse sperm proteins. *J Proteome Res.* 2009;8:1431–40.
66. Kalebic N, Sorrentino S, Perlas E, Bolasco G, Martinez C, Heppenstall PA. α TAT1 is the major α -tubulin acetyltransferase in mice. *Nat Commun.* 2013;4:1962.
67. Singh ND, Koerich LB, Carvalho AB, Clark AG. Positive and purifying selection on the *Drosophila* Y chromosome. *Mol Biol Evol.* 2012;31:2612–23.
68. Riley NM, Coon JJ. Phosphoproteomics in the age of rapid and deep proteome profiling. *Anal Chem.* 2016;88:74–94.
69. Walters JR, Harrison RG. Combined EST and proteomic analysis identifies rapidly evolving seminal fluid proteins in *Heliconius* butterflies. *Mol Biol Evol.* 2010;27:2000–13.
70. Walters JR, Harrison RG. Decoupling of rapid and adaptive evolution among seminal fluid proteins in *Heliconius* butterflies with divergent mating systems: seminal fluid proteins in *Heliconius* butterflies. *Evolution.* 2011;65:2855–71.
71. Swanson WJ, Vacquier VD. The rapid evolution of reproductive proteins. *Nat Rev Genet.* 2002;3:137–44.
72. Hahn MW, Han MV, Han S-G. Gene family evolution across 12 *Drosophila* genomes. *PLoS Genet.* 2007;3:e197.
73. Zhao L, Saelao P, Jones CD, Begun DJ. Origin and spread of *de novo* genes in *Drosophila melanogaster* populations. *Science.* 2014;343:769–72.
74. Haerty W, Jagadeeshan S, Kulathinal RJ, Wong A, Ravi Ram K, Sirot LK, et al. Evolution in the fast lane: rapidly evolving sex-related genes in *Drosophila*. *Genetics.* 2007;177:1321–35.
75. Vicens A, Lücke L, Roldan ERS. Proteins involved in motility and sperm-egg interaction evolve more rapidly in mouse spermatozoa. *PLoS One.* 2014;9:e91302.
76. Dorus S, Wasbrough ER, Busby J, Wilkin EC, Karr TL. Sperm proteomics reveals intensified selection on mouse sperm membrane and acrosome genes. *Mol Biol Evol.* 2010;27:1235–46.
77. Dean MD, Good JM, Nachman MW. Adaptive evolution of proteins secreted during sperm maturation: an analysis of the mouse epididymal transcriptome. *Mol Biol Evol.* 2008;25:383–92.
78. Vicens A, Gomez Montoto L, Couso-Ferrer F, Sutton KA, Roldan ERS. Sexual selection and the adaptive evolution of PKDREJ protein in primates and rodents. *Mol Hum Reprod.* 2015;21:146–56.