



UNIVERSITY OF LEEDS

This is a repository copy of *Mass spectrometric characterisation of the major peptides of the male ejaculatory duct, including a glycopeptide with an unusual zwitterionic glycosylation*.

White Rose Research Online URL for this paper:
<https://eprints.whiterose.ac.uk/175554/>

Version: Accepted Version

Article:

Sturm, S, Dowle, A, Audsley, N et al. (1 more author) (2021) Mass spectrometric characterisation of the major peptides of the male ejaculatory duct, including a glycopeptide with an unusual zwitterionic glycosylation. *Journal of Proteomics*, 246. 104307. ISSN 1874-3919

<https://doi.org/10.1016/j.jprot.2021.104307>

© 2021, Elsevier. This manuscript version is made available under the CC-BY-NC-ND 4.0 license <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

Reuse

This article is distributed under the terms of the Creative Commons Attribution-NonCommercial-NoDerivs (CC BY-NC-ND) licence. This licence only allows you to download this work and share it with others as long as you credit the authors, but you can't change the article in any way or use it commercially. More information and the full terms of the licence here: <https://creativecommons.org/licenses/>

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.



eprints@whiterose.ac.uk
<https://eprints.whiterose.ac.uk/>

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19

Mass spectrometric characterization of the major peptides of the male ejaculatory duct, including a glycopeptide with an unusual zwitterionic glycosylation

Sebastian Sturm^a, Adam Dowle^b, Neil Audsley^c, R. Elwyn Isaac^{a*}

^a School of Biology, University of Leeds, Leeds LS2 9JT, UK.

^b Bioscience Technology Facility, Department of Biology, University of York, Wentworth Way, York YO10 5DD, UK. Email:adam.dowle@york.ac.uk

^c Institute for Agri-Food Research and Innovation, Newcastle University, Newcastle Upon-Tyne, NE1 7RU, UK. Email:neil.audsley@newcastle.ac.uk

*Corresponding author.

Email address: r.e.isaac@leeds.ac.uk

Declarations of interest: none.

20 ABSTRACT

21 Peptides present in the seminal fluid of *Drosophila melanogaster* can function as
22 antimicrobial agents, enzyme inhibitors and as pheromones that elicit physiological and
23 behavioural responses in the post-mated female. Understanding the molecular interactions by
24 which these peptides influence reproduction requires detailed knowledge of their molecular
25 structures. However, this information is often lacking and cannot be gleaned from just gene
26 sequences and standard proteomic data. We now report the native structures of four seminal
27 fluid peptides (andropin, CG42782, Met75C and Acp54A1) from the ejaculatory duct of male
28 *D. melanogaster*. The mature CG42782, Met75C and Acp54A1 peptides each have a cyclic
29 structure formed by a disulfide bond, which will reduce conformational freedom and enhance
30 metabolic stability. In addition, the presence of a penultimate Pro in CG42782 and Met75C
31 will help prevent degradation by carboxypeptidases. Met75C has undergone more extensive
32 post-translational modifications with the formation of an N-terminal pyroglutamyl residue
33 and the attachment of a mucin-like O-glycan to the side chain of Thr₄. Both of these
34 modifications are expected to further enhance the stability of the secreted peptide. The glycan
35 has a rare zwitterionic structure comprising an O-linked N-acetyl hexosamine, a hexose and,
36 unusually, phosphoethanolamine. A survey of various genomes showed that *andropin*,
37 *CG42782*, and *Acp54A1* are relatively recent genes and are restricted to the *melanogaster*
38 subgroup. *Met75C*, however, was also found in members of the *obscura* species groups and
39 in *Scaptodrosophila lebanonensis*. *Andropin* is related to the cecropin gene family and
40 probably arose by tandem gene duplication, whereas *CG42782*, *Met75C* and *Acp54A1*
41 possibly emerged *de novo*. We speculate that the post-translational modifications that we
42 report for these gene products will be important not only for a biological function, but also
43 for metabolic stability and might also facilitate transport across tissue barriers, such as the
44 blood-brain barrier of the female insect.

45

46 *Keywords:* Seminal fluid peptides, novel genes, andropin, CG42782, Met75C, Acp54A1,
47 male accessory organs, ejaculatory duct, *Drosophila melanogaster*, glycopeptide,
48 phosphoethanolamine

49 1. Introduction

50 For several decades, the polyandrous *Drosophila melanogaster* with its amenable genetics
51 and readily accessible bioinformation, has proven to be an excellent animal model to study
52 the physiological role and evolutionary significance of seminal fluid (SF) proteins/peptides
53 (SFPs)[1, 2]. Such studies have provided genetic and molecular insights into the role of SFPs
54 in areas such as securing paternity, manipulation of female behaviour, sexual conflict, sperm
55 storage, sperm competition and postcopulatory sexual selection [2-4]. The male *Drosophila*
56 accessory sex glands, responsible for the synthesis and secretion of most of the SFPs, include
57 the paired male accessory glands (MAGs) or paragonia, the ejaculatory duct (ED) and the
58 ejaculatory bulb (EB) [5, 6]. The MAGs are relatively large elongated sacs that open below
59 the vasa deferentia into the distended section of the ED that connects to the seminal vesicles.
60 The wall of the MAG includes two types of secretory cells: a single layer of cuboidal
61 binucleated cells with larger vacuolated secondary cells interspersed towards the apical
62 region. The ED, consisting of a single layer of large glandular polygonal cells, tapers as it
63 approaches and traverses the EB and discharges in the aedeagus (equivalent to penis). The
64 EB is a thick-walled muscular organ that contains fluid of a composition different to the
65 ED[5]. The ED is capable of strong peristaltic contractions, which presumably forces the
66 MAG and ED gland secretions together with sperm into the EB.

67 Several studies have used mass spectrometry-based proteomic analysis and transcriptomic
68 data of the *Drosophila* MAGs, ED and EB to identify candidate SFPs [7-10]. Confirmation
69 that many of these are actually transferred to the female on mating was obtained by using
70 isotope labelling of proteins to distinguish between male and female derived molecules in the
71 female reproductive tract after mating [8]. In addition, indirect evidence of transfer has been
72 obtained from quantitative proteomics of male reproductive glands before and after mating
73 [10]. We can conclude from the proteomic studies and also from high-throughput expression
74 data that many of the prominent SFPs are expressed either exclusively in the male
75 reproductive tract or are at least highly enriched in these tissues. Proteomics of soluble
76 proteins from separated male tissues, together with cell localisation data from *in situ*
77 hybridisation, tagged protein expression in transgenic flies and immuno-histochemical
78 studies, indicate that the MAGs, ED and EB all contribute to the synthesis of SFPs [9, 10].

79 Only a few of the SFPs of *D. melanogaster* have been biochemically characterised (e.g.
80 [7]). For other SFPs, a biochemical function can sometimes be inferred from sequence
81 homology to well characterised proteins. [1]. A number of SFPs can be classified as enzymes,

82 (e.g. proteases and lipases), binding proteins (e.g. lectins) or mating plug components [2, 10].
83 Another significant subset of SFPs are smaller proteins or peptides (molecular mass of <15
84 kDa) that include pheromones, antimicrobials, protease inhibitors as well as odorant binding
85 proteins [1, 2]. The best known of these molecules is the sex peptide (SP), a 36-mer peptide
86 pheromone that is responsible for eliciting multiple physiological and behavioural responses
87 in the post-mated female, including a reluctance to re-mate for several days and increased
88 oviposition [11]. In contrast, other *Drosophila* SF peptides are poorly characterised, both
89 structurally and physiologically.

90 Proteomics, involving analysing the masses of tryptic peptide fragments, has
91 revolutionised the way we identify and, in some instances, quantify protein products of the
92 male reproductive glands [12]. These approaches, however, are primarily used to identify
93 proteins and are not suited for the structural characterisation of native SFPs. The reliance on
94 convenient tryptic cleavage sites and the fact that many SFPs are post-translationally
95 modified, conspire to hinder the acquisition of structural information [13]. Furthermore, some
96 proteomic methods involve the fractionation of proteins prior to proteolysis, which runs the
97 risk of omitting small proteins and peptides from the analysis.

98 In this paper we report the application of peptidomics, a widely used method for the
99 analysis of regulatory peptides in insect nervous and endocrine tissues, to investigate the
100 structure of SF peptides originating from the ED of male *D. melanogaster*. This approach
101 preferentially extracts peptides in their native form, whilst inhibiting endogenous protease
102 activity during the extraction procedure and uses trypsin only selectively after peptide
103 identification [14-16]. We have structurally characterised four major ED peptides with a
104 molecular mass of < 4 kDa that are destined for the SF. All four are taxon-restricted and
105 products of novel genes. One prominent ED peptide is a highly post-translationally modified
106 product of two genes *Met75Ca* and *Met75Cb*. The mature form of *Met75C* is a 30-mer
107 peptide with a rare zwitterionic O-glycan. We speculate that the post-translational
108 modifications we describe are likely to be important not only for a biological function, but
109 also for extracellular stability and might also facilitate transport across tissue barriers, such as
110 the blood-brain barrier of the female insect.

111 **2. Methods**

112 *2.1. Insects*

113 The Dahomey strain of *D. melanogaster* was provided by Professor T. Chapman
114 (University of East Anglia, U.K.). The wild-type Canton-S and *w¹¹¹⁸* strains were from
115 Professor Y-J Kim (Gwangju Institute of Science and Technology, Gwangju 500-712,
116 Republic of Korea). Insects were maintained on oatmeal/molasses/agar medium at 25°C in a
117 12:12 light-dark cycle.

118 2.2. *Preparation of tissue samples for mass spectrometry*

119 Tissues were dissected from 3-8-day-old unmated males in cold insect saline
120 containing 150mM NaCl, 10mM KCl, 3.9mM NaHCO₃, 3.5mM MgCl₂, 1.3mM CaCl₂
121 adjusted to pH 7.2. Individual tissue samples were rinsed in purified water and subsequently
122 placed on a MALDI sample plate for direct tissue profiling. Alternatively, 10 organs were
123 transferred into a vial containing 10 µl of extraction solution (50% aqueous methanol
124 containing 0.1% formic acid (FA)). Tissue collections were sonicated in a water bath and
125 subsequently centrifuged for 15 min at 13,000 rpm. Supernatants were transferred into clean
126 vials and stored at -20°C until required for analysis.

127 2.3. *Reduction, alkylation, tryptic digestion, sulfonation and purification of peptides*

128 Peptides were subjected to cystine reduction by dithiothreitol (Sigma-Aldrich
129 Company Ltd. Gillingham, Dorset, U.K.) and alkylation by iodoacetic acid (Sigma-Aldrich
130 Company Ltd.) followed by enzymatic digestion using trypsin (Sequencing Grade Modified
131 Trypsin, Promega U.K. Ltd., Southampton, U.K.) as described previously [14]. Briefly, tissue
132 supernatants were concentrated by using a SpeedVac, and adjusted to a volume of 25 µl with
133 50 mM ammonium bicarbonate (ABC) buffer, pH 8.2. Disulfide reduction was performed by
134 adding 200 mM dithiothreitol (DTT) in ABC buffer to a final concentration of 10 mM DTT
135 at 37 °C for 1 h. Carbamidomethylation of reduced cysteines was performed by adding
136 200 mM iodoacetamide (IAA) in ABC buffer to a final concentration of 40 mM IAA at room
137 temperature for 1 h. Unreacted IAA was inactivated/precipitated by adding DTT at room
138 temperature to a final concentration of 40 mM DTT for 15 min. Extracts with reduced and
139 carbamidomethylated peptides were adjusted to 100 µl of 0.5 % aqueous FA prior to
140 desalting and purification. For digestion of proteins, 0.1 µl of a 1 µg/µl solution of trypsin in
141 50 mM acetic acid (Sequencing Grade Modified Trypsin, Promega) was added to the reduced
142 and alkylated extract and incubated at 37 °C for 16 h.

143 .

144 For *de novo* sequencing, peptides were sulfonated using 4-sulfophenyl isothiocyanate
145 (SPITC; Sigma-Aldrich Company Ltd.) following the protocol of Sturm *et al.* [15]. Briefly,
146 25 µl of a 10 mg/ml SPITC solution in 20 mM NaHCO₃ (pH 9.0) was added to 25 µl of a
147 digested extract. The reaction mixture was incubated in a thermo-mixer (Eppendorf,
148 Hamburg, Germany) at 55 °C and 300 rpm for 1 h. After incubation the reaction was
149 terminated by adding 10 µl of 1% aqueous FA.

150 The reduced and alkylated, or additionally digested and SPITC derivatized peptides
151 were purified and desalted using custom made C18 spin columns using cut-outs of C18
152 Empore 3M extraction discs (3M, USA; gifted by IVA Analysentechnik e.K., Meerbusch,
153 Germany) packed into 200 µl pipette tips (Eppendorf, Germany) as described in [17]. The
154 sorbent was activated with 100 µl of 80% aqueous acetonitrile (ACN) containing 0.1% FA
155 and equilibrated with 100 µl 0.1% aqueous FA by centrifuging the spin columns at 2000 rpm.
156 Samples were diluted with 100 µl of 0.1% aqueous FA, loaded on the spin column and
157 washed with 200 µl of 0.1% aqueous FA. Peptides were eluted with 5 µl of 80% aqueous
158 ACN containing 0.1% FA step-by-step directly onto the MALDI sample plate using a
159 modified 20 ml plastic syringe.

160 2.5. Matrix-assisted laser desorption ionisation mass spectrometry (MALDI-TOF 161 MS)

162 Mass spectra were acquired in positive reflector mode using Bruker ultrafleX III or
163 ultrafleXtreme mass spectrometers (Bruker Daltonics, Bremen, Germany). Samples were
164 mixed with 0.5 µl 2,5-Dihydroxybenzoic acid (DHB, 10 mg/ml in 20% acetonitrile, 1% FA)
165 on the MALDI sample plate and dried using a gentle stream of hot air. External calibration
166 was conducted using calibration mixtures containing bradykinin¹⁻⁷, angiotensin I, angiotensin
167 II, substance P, bombesin, ACTH clip¹⁻¹⁷, ACTH clip¹⁸⁻³⁹, somatostatin 28, insulin and
168 ubiquitin I (Bruker Daltonics).

169 2.6. Acetylation of amines

170 The reaction mixture was freshly prepared by mixing acetic anhydride (Sigma-
171 Aldrich Company Ltd.) and methanol (Fisher Chemicals U.K., Loughborough, U.K.) in a
172 ratio of 1:3. A small droplet of the solution was added to cover the surface of the dried
173 disrupted tissue on the MALDI sample plate and rapidly evaporated with a hair dryer after 30
174 seconds. Subsequently 0.5 µl of DHB matrix solution was added on the heated sample and
175 left to dry.

176 2.7. *Fourier Transform Ion Cyclotron Resonance mass spectrometry (FT-ICR MS)*

177 A FT-ICR MS analysis was performed using a solariX XR FT MS (Bruker Daltonics)
178 with a 9.4 T superconducting magnet. Peptide solutions were diluted 1:20 into 50% aqueous
179 acetonitrile containing 1% FA before introduction by TriVersa NanoMate (Advion
180 BioSciences, Ithaca, NY) in positive-ion mode. The applied voltage was adjusted between
181 1.4-1.7 kV to achieve a stable ion current. A 120°C nitrogen dry gas was supplied at 1.3
182 L/min to aid desolvation. Instrument control and data acquisition used Compass 1.4 (Bruker
183 Daltonics). Spectra were generated by the accumulation of 20 scans with 0.2 s ion cooling
184 time and 0.5 s scan time with 400K data points recorded. Peptide precursors were manually
185 selected for isolation and subsequent fragmentation by collision induced dissociation in the
186 hexapole (Q-CID) with argon as the collision gas. Collision energies were optimized for each
187 peptide. Spectra were processed using DataAnalysis version 4.0 (Bruker Daltonics). Mass
188 deconvolution was performed using version 2.0 of the SNAP averaging algorithm (C 4.9384
189 %, N 1.3577 %, O 1.4773 %, S 0.0417 %, H 7.7583 %).

190 2.8 *Bioinformatics*

191 Sequence databases were searched using BLASTP and TBLASTN algorithms [18]
192 (<http://blast.ncbi.nlm.nih.gov/Blast.cgi>). Nucleotide sequences were translated using ExPASy
193 Translate (Swiss Institute of Bioinformatics; <http://web.expasy.org/translate/>). Retrieved
194 protein sequences were submitted to the SignalP-5.0 server
195 (<http://www.cbs.dtu.dk/services/SignalP/>) to establish the occurrence of a signal peptide for
196 secretion and the probable cleavage site for signal peptidase [19]. Peptide sequences were
197 aligned using ClustalOmega [20] (<https://www.ebi.ac.uk/Tools/msa/clustalo/>). The predicted
198 cleavage sites of the signal peptidase differed between the Met75C orthologs resulting in
199 varying lengths of the signal peptide and the propeptide. As a result of this, homologous
200 sequence position may be part of the signal peptide in one species, but of the propeptide in
201 another. Therefore, we manually adjusted and optimized the alignment for Met75C sequence
202 logos to provide a consensus signal and propeptide. Sequence logos were created with
203 Weblogo Version 3.7.4 (<http://weblogo.threeplusone.com>) using equiprobable composition
204 [21].

205 A maximum likelihood phylogeny was calculated with MEGA X[22] using the LG model
206 with gamma distribution and five categories [23]. Percentage similarities and identities of

207 sequences were calculated using MatGAT 2.01 [24] applying the in-build alignment
208 algorithm using default settings and the BLOSUM50 substitution matrix.

209 **3. Results**

210 *3.1 Peptidomics of the ejaculatory duct*

211 We initially used MALDI-TOF MS to investigate the peptide inventory of tissues associated
212 with the male reproductive tract in three different strains of *D. melanogaster*. The male
213 accessory glands (MAG), dilated and narrow regions of the anterior ejaculatory duct (ED),
214 ejaculatory bulb (EB), posterior ED, testis, and seminal vesicle (SV) of individual flies were
215 analysed. The spectra for the different regions of the ED (Fig. 1) feature a common pattern of
216 ion signals that are absent in tissue samples from the MAG, the EB, SV and testis (data not
217 shown). We found highly similar and reproducible patterns, in the laboratory strains *w¹¹¹⁸*,
218 Dahomey and Canton-S, only varying in the relative abundance of the ion signals between
219 different individuals. To identify the structures of the detected ion signals, individual tissue
220 sections and a peptide extract of 10 EDs that had been treated with DDT and IAA to reduce
221 and alkylate cysteines, were subjected to MALDI-TOF MSMS. The same peptide extract was
222 also analysed by FT-ICR MS. Ion signals in the MALDI-TOF MS spectra of the ED were
223 identified as peptide products of the following genes: *CG34098* (*Acp4A1*; Supplementary Fig.
224 1), *CG42782* (Supplementary Fig. 2), *CG1361* (*andropin*; Supplementary Fig. 3), *CG32197*
225 and *CG18064* (*Met75Ca* and *Met75Cb*; Fig. 2). All five genes code for proteins that are
226 predicted to have a secretory signal peptide and a signal peptidase cleavage site. The ion signal
227 of andropin was found to be consistently lower than the other three gene products. Surprisingly,
228 three distinct ion signals at m/z 3451.3, m/z 3816.4 and m/z 3939.7 (Fig. 1b) were all identified
229 as peptide forms encoded by *Met75Ca* and *Met75Cb* (Table 1). The sequences of the predicted
230 prepropeptide products of *Met75Ca* and *Met75Cb* are identical and, hereafter, are referred to
231 as *Met75C*. A literature search indicated that the mass differences between *Met75C*-derived
232 ion signals and the fragmentation pattern of the peptide ion at m/z 3939.7 (Fig. 2) is likely due
233 to the presence of a post-translational modification comprising an N-acetyl hexosamine
234 (HexNAc), a hexose (Hex) and a phosphoethanolamine (PEA) [25, 26]. The ED peptide
235 Dup99B reported by Saudan *et al.*[27] was not detected in any of the MALDI-TOF MS
236 analyses. However, a mass match corresponding to the glycosylated form of Dup99B was
237 found in the FT-ICR MS analysis. Table 1 summarises the detected ion signals and the deduced
238 sequences for the native and alkylated peptides found in the ED by MALDI-TOF MSMS and
239 FT-ICR MS.

240 3.2. *Glycosylation of Met75C*

241 In order to determine the glycosylated residue, the reduced and alkylated extract of 10 EDs
242 was digested with trypsin and subsequently derivatized with 4-sulfophenyl isothiocyanate
243 (SPITC), a procedure that sulfonates N-termini and facilitates sequencing due to the
244 generation of abundant fragments of the γ -ion series in MALDI-TOF MSMS. Figure 3 shows
245 a MALDI-TOF mass spectrum of the derivatised ED tryptic peptides and the sequences of
246 selected peptides confirmed by MSMS. Two ion signals at m/z 774.3 and m/z 936.4 were
247 identified as glycosylated tryptic peptides derived from the N-terminus of Met75C (Fig. 3).
248 Fragmentation of the ion signal at m/z 774.3 (data not shown) resulted in a complex pattern
249 due to the co-fragmentation of the HexNAc-modified N-terminal tryptic peptide (Fig. 3 a₁)
250 with the C-terminal tryptic peptide of Met75C (Fig. 3 a₂). Fragmentation of the ion signal at
251 m/z 936.4 (Fig. 3), however, confirmed the glycosylation comprising a distal hexose and a
252 proximal N-acetyl hexose (Fig. 4) attached to a residue in the N-terminal pentapeptide
253 sequence (pQIATR). Thr₄ and Arg₅ are possible sites of attachment for an O-linked and N-
254 linked glycan, respectively. However, glycosylation of an N of Arg is unlikely since this
255 post-translational modification has been confirmed only for bacteria and the presence of the
256 glycan structure at this position would likely interfere with tryptic cleavage of the Arg-Gln
257 peptide bond [28]. We therefore conclude that Thr₄ is the glycosylation site.

258 3.3. *Glycan substitution with phosphoethanolamine (PEA)*

259 Phosphoethanolamine (PEA) is a zwitterionic structure exhibiting a negative charge at the
260 glycosidic-attached phosphate group and a distal positive charge at the amine. In order to
261 confirm this structure, we performed on-plate derivatisation of single ED tissues with a
262 solution of acetic anhydride. Acetic anhydride acetylates primary amines resulting in a mass
263 shift of +42 Da for free N-termini, each lysine and PEA. Acetylation resulted in the expected
264 number of mass shifts, namely, three times +42 Da for the Acp54A1 and two times +42 Da
265 for CG42782 (Fig. 5). The glycosylated Met75C, which has one Lys at position 12, but no
266 free amino group at the N-terminus because of the cyclic N-terminal pyroglutamic acid,
267 featured two mass shifts of +42 Da after derivatisation, which indicates the presence of an
268 additional primary amine within the glycan structure.

269 3.4. *Proposed glycan structure*

270 We sought further clarification of the glycan structure using high-resolution FT-ICR
271 MS, but the fragmentation of the Met75C parent ion using electron-capture dissociation or
272 collision induced dissociation did not provide evidence of the location of the glycan.
273 However, low energy qCID fragmentation (12 V) resulted in the multi-charged fragments m/z
274 779.54434 ($z=5$) and m/z 974.17824 ($z=4$). These ions correspond to Met75C with the
275 remaining HexNAc-PEA and a loss of a hexose (deconvoluted $z=1$; 3893.69259 and
276 3893.69113) respectively. This suggests that both the hexose (Hex) and the PEA are attached
277 to the N-acetylhexosamine (HexNAc) rather than in a linear manner (Fig. 6).

278 3.5. Evolutionary history of the ejaculatory duct peptides

279 To investigate the evolutionary history of the ED peptides, we screened publicly
280 available genome and transcriptome databases for sequences related to *Acp51A1* (CG34098),
281 CG42782, *andropin* and *Met75C*. Our analysis shows that *andropin* orthologs are present in
282 the *melanogaster* subgroup as well as the *eugracilis*, *takahashii* and *suzukii* subgroups (Fig. 7)
283 and, as noted previously, the predicted mature peptides have diverged with non-identical
284 residues replaced by conservative substitutions [29] (Supplementary material:
285 *Andropin_alignment.fas*).

286 CG42782 orthologues were only found in members of the *melanogaster* subgroup (Fig.
287 7). These orthologues have the two absolutely conserved Cys residues that form the disulfide
288 bridge, as well as four other conserved residues within the ring structure (Supplementary
289 material: CG42782_alignment.fas). In total, seven out of ten residues of the peptide ring are
290 either conserved or are substituted by amino acids with similar side chain properties. The
291 predicted C-termini of the peptides, however, show structural divergence with four (*D. erecta*,
292 *D. simulans*, *D. sechellia* and *D. mauritania*) out of six predicted to have a C-terminal Gly,
293 which we expect to be utilised to convert the carboxyl group of the penultimate residue to an
294 amide.

295 Orthologs of *Acp51A1* (CG34098) have been found in *D. simulans*, *D. sechellia* and *D.*
296 *yakuba* of the *melanogaster* subgroup, but not in *D. erecta*. Additional *Acp51A1* orthologs
297 could also be assigned to two species outside of the *melanogaster* subgroup, namely *D.*
298 *biarmipes* and *D. takahashi* (Fig. 7). The ring forming Cys residues are absolutely conserved
299 as are two Gly residues within the cyclic structure (Supplementary material:
300 *Acp54A1_alignment.fas*).

301 In our database screening for *Met75C*-related genes, we found numerous orthologs in
 302 members of the subgenus *Sophophora*, as well as multiple paralogs in some of these species.
 303 However, we could not find *Met75C*-related genes in *D. willistoni* and the subgenus *Drosophila*
 304 (e.g., *D. albomicans*, *D. grimshawi*, *D. mojavenensis* and *D. virilis*; Fig. 7). Nevertheless, we
 305 detected an ortholog in *Scaptodrosophila lebanonensis* that had 51.0% amino acid similarity
 306 and 42.6% identity with *D. melanogaster Met75C*. *D. melanogaster* has two paralogous genes
 307 (*CG18064, Met75Cb* and *CG32197, Met75Ca*), encoding identical *Met75C* prepropeptides,
 308 that are adjacent to each other and separated by 2611 bp on the left arm of chromosome 3.
 309 Further species in which we found paralogs are *D. suzukii*, *D. subpulchrella*, *D. takahashii*, *D.*
 310 *rhopaloea*, *D. leontia*, *D. bocki*, *D. kikkawai*, *D. nikananu*, *D. bipectinata*, *D. azteca* and *D.*
 311 *athabasca*, each with two paralogs as well as *D. sechellia* and *D. subobscura* containing three.
 312 Phylogenetic analysis of the sequences suggests that most of these paralogs arose from
 313 independent gene duplication events with the exception of a common duplication in the
 314 ancestral line of *D. leontia*, *D. bocki*, and *D. kikkawai*, as well as a common evolutionary event
 315 shared by *D. azteca* and *D. athabasca* (Supplementary Fig.4). Interestingly, duplication in these
 316 clades leads to a diversification of one daughter gene while the other one remains more
 317 ancestral (see corresponding percent similarity/identity Supplementary Fig.4.). The average
 318 distance between paralogs is around 1800 bp with the exception of *D. azteca* and *D. athabasca*,
 319 where the paralogs are located on different chromosomes, namely chromosome 2 and XR as a
 320 result of chromosomal rearrangement [30]. There is a high degree of sequence conservation
 321 amongst the *Met75C* orthologues at the N-terminal glutamine, the two cysteines, a Gly and
 322 several Arg and Tyr residues within the ring structure (Fig. 8). Notably, there is a common C-
 323 terminal sequence comprising conservative hydrophobic residues on either side of a Pro. The
 324 Thr in *D. melanogaster Met75C* to which a glycan is attached, is also conserved in the majority
 325 of the predicted peptides.

326 4. Discussion

327 Proteomic studies that exploit the prowess of *Drosophila melanogaster* as a model
 328 organism have, for some time, provided a wealth of information on the protein composition
 329 and chemical complexity of SF. Clearly, structural characterisation of individual SFPs is
 330 critical for advancing our understanding of the molecular mechanisms by which these
 331 molecules influence reproduction [7-9]. This aspect has often been neglected in proteomic
 332 studies of the reproductive tissues of *Drosophila*, which had the primary objective of
 333 identifying the presence or absence of a soluble protein in a tissue extract. By applying a

334 peptidomic strategy that has been used widely to characterise peptides of the insect nervous
335 and endocrine systems, we have determined the native structures of four SF peptides
336 extracted from the *D. melanogaster* ED. Of particular interest is the revelation that Met75C is
337 O-glycosylated, a peptide modification that is not uncommon in *Drosophila*. It is found in
338 several proline-rich antibacterial peptides (drosocin and matured Pro-domain of attacin C)
339 and the ED peptide, DUP99B [56, 27, 61, 62]. Drosocin occurs as both single and doubly
340 glycosylated forms and glycosylation is essential for full antibacterial activity [56].

341 *CG42782* encodes a 46-mer prepropeptide that is predicted to be processed by
342 cleavage of the signal peptide at the Ser-Tyr bond to generate a 23-mer mature secreted
343 peptide. The molecular ion $[M+H]^+$ at m/z 3034.0 and its fragmentation pattern confirms this
344 prediction and establishes that the peptide forms a ring structure through the formation of a
345 disulphide bond between Cys⁷ and Cys¹⁸. The presence of a penultimate Pro residue is
346 expected to protect the peptide from general carboxypeptidase attack and the cyclic peptide is
347 also expected to hinder degradation by proteolytic enzymes [31, 32]. We found no evidence
348 that *CG42782* is present in the genomes outside of the *melanogaster* subgroup (Fig.7).

349 The mature Acp54A1 (CG34098) peptide found in the ED starts with a methionine,
350 which is consistent with the predicted signal cleavage site at the Ala-Met bond. The peptide
351 comprises 17 amino acids, 12 of which form a ring by virtue of a disulphide between Cys₃
352 and Cys₁₂, which will reduce vulnerability to peptidase attack [32]. Orthologs of *Acp54A1*
353 were found in members of the *melanogaster* subgroup and also in *D. biarmipes* and *D.*
354 *takahashii* (Fig. 7). The presence of *Acp54A1* in *D. yakuba* and the absence of the gene from
355 *D. erecta* is discordant with a previous report that the gene could not be found in *D. yakuba*
356 [33]. The failure to identify an Acp54A1 orthologue in *D. erecta* might be explained by gene
357 loss in this species.

358 The andropin molecular ion $[M+H]^+$ at m/z 3749.0 (Supplementary Fig. 3) matched
359 the predicted mass of the mature 34 residue peptide, generated after cleavage of the signal
360 peptide at the Ala-Val bond and established the absence of any further post-translational
361 modification. Prior studies have shown that *Andropin* expression is male specific and is
362 restricted to the ED. *Andropin* is closely linked to the *Cecropin* gene family on the right arm
363 of chromosome 3 and likely to have arisen from duplication of a common ancestral *Cecropin*
364 gene with subsequent divergence in both sequence and tissue expression [34]. The only
365 sequence similarity between the *Andropin* and *Cecropin* protein products lies in the signal

366 peptide of the preproprotein [29]. The synthetic andropin peptide has moderate anti-bacterial
367 activity towards Gram-positive bacteria, but, unlike the cecropin peptides, is not active
368 against Gram-negative bacteria [34]. It has been reported that the ED is the source of a
369 second antibacterial compound in addition to andropin. It is conceivable that one or more of
370 the other ED [35].

371 Met75C is a 30-mer peptide encoded by two intron-less *D. melanogaster* genes
372 (*CG18064*, *Met75Cb* and *CG32197*, *Met75Ca*) that are tandem duplicates. Met75C was
373 identified previously as an ED peptide from tryptic fragment mass analysis and the peptide
374 has been confirmed as a *bona fide* SF component that is passed to the female reproductive
375 tract during copulation [8, 10]. We confirm the prediction that the signal peptide is cleaved at
376 the Ala-Gln bond and reveal that the glutamine is cyclised to pyroglutamate in the mature
377 peptide. The C-terminal half of the peptide, like both CG42782 and Acp51A1, is a ring
378 structure formed by a disulphide bond between Cys₁₄ and Cys₂₇. Of special interest, is the
379 discovery that Met75C is decorated with an O-glycan attached to Thr₄ in the N-terminal
380 sequence pQIATR. Our MS analysis predicts that the glycan comprises an N-
381 acetylhexosamine and that this sugar is extended by a single hexose and the zwitterionic
382 phosphoethanolamine. This type of modification, generating a positively and negatively
383 charged glycan, is rare, but has been described previously for O-linked glycosylations in the
384 nest material of the wasp, *Vespula germanica* and the glycopeptide, noctilisin, in the venom
385 gland of *Sirex noctilio* [36, 37] (Fig.6). Phosphoethanolamine has also been reported for N-
386 linked glycans in extracts of the royal jelly, the venom gland and homogenates of larvae of
387 the honey bee, *Apis mellifera* [38] as well as for glycosphingolipids of the pupae of the blow
388 fly *Calliphora vicina* [39, 40].

389 The most common O-glycosylation of secreted proteins in *D. melanogaster* is of the
390 mucin-like O-glycan type that is initiated by the attachment of N-acetylgalactosamine
391 (GalNAc) to the side chain of either Thr or Ser [41, 42, 62]. This reaction is catalysed by
392 polypeptide N-acetylgalactosaminyltransferases (PGANTs), of which there are 14 in *D.*
393 *melanogaster* [41, 42]. This family of type II membrane enzymes are resident in the Golgi
394 and can display different substrate specificities and expression patterns. Transcriptomic data
395 has shown that the MAG/ED tissues strongly express the genes encoding PGANT4 and
396 PGANT9, suggesting that one or both enzymes are involved in initiating O-glycan synthesis
397 in these tissues [43].

398 A common extension to the GalNAc is the addition of galactose in a β -1,3-linkage
399 carried out by core β -1,3-galactosyltransferases (C1GalTs)[44]. In *Drosophila*, there are 11
400 genes for this family of transferases and of these, two are strongly and exclusively expressed
401 in MAG/ED tissues [43]. Since it is estimated that around 90% of *D. melanogaster* O-glycans
402 are of the mucin-type and given the strong expression of the aforementioned transferases in
403 the male reproductive tract, it is very likely that this biosynthetic pathway is responsible for
404 the O-glycosylation of Met75C.

405 A comparison of the predicted primary structures of the Met75C peptides from
406 various *Drosophila* species in the subgenus *Sophophora* reveals that about half of the amino
407 acids are highly conserved, whereas significant amino acid divergence has occurred in other
408 positions both within and outside the peptide ring. Two plausible, and not necessarily
409 mutually exclusive explanations for the conservation, are constraints imposed by structural
410 and conformational requirements for a presumed physiological role in reproduction and/or
411 the need to protect the integrity of the peptide from peptidases, not only in the seminal fluid
412 but also in the female [45]. The reduced conformational freedom resulting from peptide
413 cyclisation can confer specificity to molecular interactions, as well as enhanced stability [32].
414 An N-terminal glutamine that cyclises to pyroglutamate, the peptide ring and Pro residues are
415 typical features that can provide protection from peptidases [31, 46]. The glycan with its
416 zwitterionic phosphoethanolamine provides a charged carbohydrate shield, which is expected
417 to provide protection of the N-terminal section of Met75C from proteolysis.

418 Met75C is not the only glycosylated peptide secreted by the ED of *D. melanogaster*.
419 DUP99B is a 31-mer peptide that has a branched N-linked glycan, comprising two molecules
420 of N-acetylglucosamine, three of mannose and two of fucose, attached at Asn₄. DUP99B is a
421 paralogue of the MAG sex peptide (SP), and like SP, can elevate oviposition and inhibit male
422 receptivity when injected into virgin females [27, 47]. A synthetic non-glycosylated form of
423 DUP99B is also able to activate the SP receptor expressed in mammalian cells, albeit with a
424 10-fold reduction in potency [48]. Interestingly, the presence of the glycan in the natural
425 DUP99B increases the potency of injected peptide in eliciting post-mating female
426 responses[27], suggesting that the glycan influences receptor binding and/or has an *in vivo*
427 role facilitating ligand access to the female SP receptor. Although there is no sequence
428 homology between Met75C and DUP99B, there are compelling structural analogies in that
429 both peptides have a pyroglutamyl N-terminus followed by a short peptide carrying a glycan,
430 which is then separated by a linker sequence from a peptide ring of similar size that

431 dominates the C-terminal half of both peptides. It has been reported that DUP99B can attach
432 to sperm, perhaps by virtue of the N-linked glycan [27]. Whether the zwitterionic glycan
433 moiety of Met75C facilitates a similar adhesion to *D. melanogaster* sperm proteins or binding
434 to other components of the ejaculate, such as the lectin-like SF proteins, is not known, but
435 worthy of investigation. There is no functional information available for Met75C, but the
436 aforementioned structural comparison with a biologically active DUP99B invites speculation
437 that Met75C might serve as a pheromonic peptide that signals via the female central nervous
438 system to influence the physiology or behaviour of the post-mated female. It is worth noting
439 that glycosylation of mammalian neuroactive peptides is known to not only confer metabolic
440 stability, but also increase penetration of the blood-brain barrier, which has encouraged the
441 development of glycopeptide drugs that target the central nervous system [49, 50].

442 All five of the genes discussed here are taxon-restricted, without detectable homologs
443 in non-Drosophilidae species and belong to a class commonly referred to as ‘new’ or ‘young’
444 protein-coding genes [51, 52]. Met75C is the most ancient of these five genes and was found
445 in *Scaptodrosophila lebanonensis*, which diverged from *Drosophila melanogaster* about 70-
446 74 Mya [53, 54]. The remaining genes appear to have originated in the *Sophophora* subgenus
447 after split of the *willistoni* species group from the *obscura* and *melanogaster* lineages. The
448 most recent genes are CG42782 and Acp54A1, estimated to have emerged about 8.8-12.8
449 Mya and 22.3 - 35.6 Mya, respectively [53, 55]. There are several known mechanisms for
450 new-gene origination with gene duplication recognised as the dominant source [52].

451 *Andropin* is most probably a product of a duplication of an ancestral cecropin gene, but the
452 origins of CG42782, Acp54A1 and Met75Ca/b are unclear. It has been suggested that such SF
453 peptide genes might have arisen *de novo* from non-coding DNA, since the intergenic and
454 intronic DNA of *D. melanogaster* harbors thousands of small short open reading frames
455 (ORF) that have the potential to generate secreted peptides [51]. If, during the course of
456 evolution, these ORFs recruit appropriate regulatory elements to drive expression in the male
457 reproductive tissues, then these *de novo* genes could rapidly evolve lineage-specific roles in
458 reproductive physiology. In addition to the structural requirement of a signal peptide for
459 translocation to the secretory pathway, the new peptide will also benefit functionally from
460 features, like those described for Acp54A1, CG42782 and Met75C, that will reduce
461 susceptibility to degradation by extracellular peptidases that the peptides will encounter in the
462 ejaculate and in female reproductive tissues.

463

464

465 **Credit authorship contribution statement**

466 **S. Sturm:** Methodology, Investigation, Formal analysis, Writing- Original draft preparation,
 467 Writing- Reviewing and Editing. **A. Dowle:** Data collection, Formal analysis, Writing-
 468 Reviewing and Editing. **R.E. Isaac:** Conceptualization, Writing- Original draft preparation,
 469 Writing- Reviewing and Editing. **N. Audsley:** Conceptualization, Writing- Reviewing and
 470 Editing.

471 **Acknowledgements**

472 The York Centre of Excellence in Mass Spectrometry was created thanks to a major capital
 473 investment through Science City York, supported by Yorkshire Forward with funds from the
 474 Northern Way Initiative, and subsequent support from EPSRC (EP/K039660/1;
 475 EP/M028127/1). We thank Reinhard Predel (University of Cologne) for providing access to
 476 the MALDI-TOF MS. REI and NA thank the Leverhulme Trust for support (RPG-2020-368).
 477 SS and REI thank Carole Sowden for technical assistance and Divya Ramesh for critical
 478 reading of the manuscript.

479

480 **References**

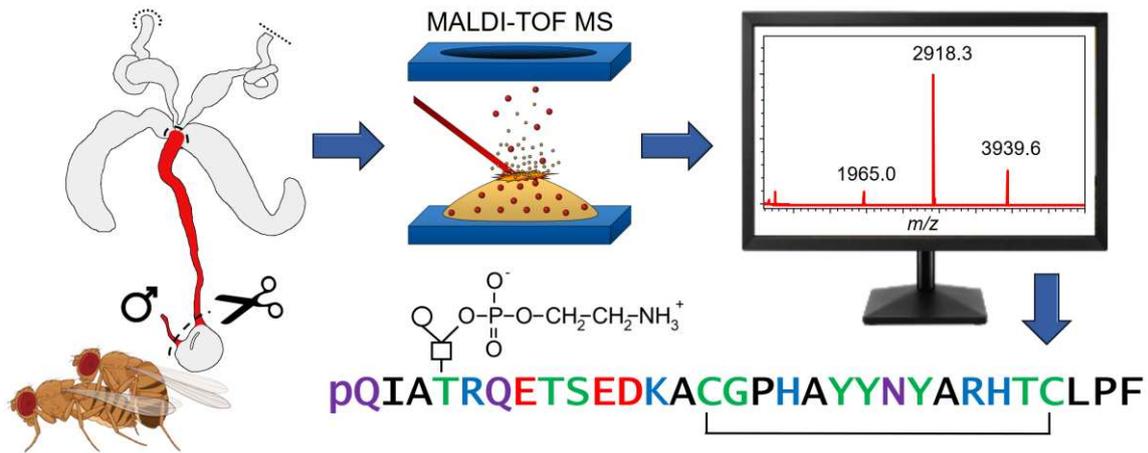
- 481 1. Avila, F.W., et al., *Insect seminal fluid proteins: identification and function*. Annu
 482 Rev Entomol, 2011. **56**: p. 21-40.
- 483 2. Wigby, S., et al., *The Drosophila seminal proteome and its role in postcopulatory*
 484 *sexual selection*. Philos Trans R Soc Lond B Biol Sci, 2020. **375**(1813): p. 20200072.
- 485 3. Sirot, L.K., et al., *Sexual Conflict and Seminal Fluid Proteins: A Dynamic Landscape*
 486 *of Sexual Interactions*. Cold Spring Harbor Perspectives in Biology, 2015. **7**(2).
- 487 4. Ramm, S.A., *Seminal fluid and accessory male investment in sperm competition*.
 488 Philos Trans R Soc Lond B Biol Sci, 2020. **375**(1813): p. 20200068.
- 489 5. Miller, A., *The internal anatomy and histology of the imago of Drosophila*
 490 *melanogaster*, in *Biology of Drosophila*, M. Demerec, Editor. 1950, John Wiley &
 491 Sons: London. p. 420-534.
- 492 6. Chen, P.S., *The Functional-Morphology and Biochemistry of Insect Male Accessory-*
 493 *Glands and Their Secretions*. Annual Review of Entomology, 1984. **29**: p. 233-255.
- 494 7. Walker, M.J., et al., *Proteomic identification of Drosophila melanogaster male*
 495 *accessory gland proteins, including a pro-cathepsin and a soluble gamma-glutamyl*
 496 *transpeptidase*. Proteome Sci, 2006. **4**: p. 9.
- 497 8. Findlay, G.D., et al., *Proteomics reveals novel Drosophila seminal fluid proteins*
 498 *transferred at mating*. PLoS Biol, 2008. **6**(7): p. e178.
- 499 9. Takemori, N. and M.T. Yamamoto, *Proteome mapping of the Drosophila*
 500 *melanogaster male reproductive system*. Proteomics, 2009. **9**(9): p. 2484-93.
- 501 10. Sepil, I., et al., *Quantitative Proteomics Identification of Seminal Fluid Proteins in*
 502 *Male Drosophila melanogaster*. Mol Cell Proteomics, 2019. **18**(Suppl 1): p. S46-S58.

- 503 11. Chen, P.S., et al., *A male accessory gland peptide that regulates reproductive*
504 *behavior of female D. melanogaster*. Cell, 1988. **54**(3): p. 291-8.
- 505 12. Findlay, G.D. and W.J. Swanson, *Proteomics enhances evolutionary and functional*
506 *analysis of reproductive proteins*. Bioessays, 2010. **32**(1): p. 26-36.
- 507 13. Fricker, L.D., et al., *Peptidomics: identification and quantification of endogenous*
508 *peptides in neuroendocrine tissues*. Mass Spectrom Rev, 2006. **25**(2): p. 327-44.
- 509 14. Sturm, S. and R. Predel, *Mass spectrometric identification, sequence evolution, and*
510 *intraspecific variability of dimeric peptides encoded by cockroach akh genes*. Anal
511 Bioanal Chem, 2015. **407**(6): p. 1685-93.
- 512 15. Sturm, S., et al., *Agatoxin-like peptides in the neuroendocrine system of the honey bee*
513 *and other insects*. J Proteomics, 2016. **132**: p. 77-84.
- 514 16. Sturm, S., et al., *The structure of the Drosophila melanogaster sex peptide:*
515 *Identification of hydroxylated isoleucine and a strain variation in the pattern of amino*
516 *acid hydroxylation*. Insect Biochem Mol Biol, 2020. **124**: p. 103414.
- 517 17. Rappsilber, J., M. Mann, and Y. Ishihama, *Protocol for micro-purification,*
518 *enrichment, pre-fractionation and storage of peptides for proteomics using StageTips*.
519 Nat Protoc, 2007. **2**(8): p. 1896-906.
- 520 18. Altschul, S.F., et al., *Gapped BLAST and PSI-BLAST: a new generation of protein*
521 *database search programs*. Nucleic Acids Res, 1997. **25**(17): p. 3389-402.
- 522 19. Almagro Armenteros, J.J., et al., *SignalP 5.0 improves signal peptide predictions*
523 *using deep neural networks*. Nat Biotechnol, 2019. **37**(4): p. 420-423.
- 524 20. Sievers, F. and D.G. Higgins, *Clustal Omega for making accurate alignments of many*
525 *protein sequences*. Protein Sci, 2018. **27**(1): p. 135-145.
- 526 21. Crooks, G.E., et al., *WebLogo: a sequence logo generator*. Genome Res, 2004. **14**(6):
527 p. 1188-90.
- 528 22. Kumar, S., et al., *MEGA X: Molecular Evolutionary Genetics Analysis across*
529 *Computing Platforms*. Mol Biol Evol, 2018. **35**(6): p. 1547-1549.
- 530 23. Le, S.Q. and O. Gascuel, *An improved general amino acid replacement matrix*. Mol
531 Biol Evol, 2008. **25**(7): p. 1307-20.
- 532 24. Campanella, J.J., L. Bitincka, and J. Smalley, *MatGAT: an application that generates*
533 *similarity/identity matrices using protein or DNA sequences*. BMC Bioinformatics,
534 2003. **4**: p. 29.
- 535 25. Conboy, J.J. and J.D. Henion, *The determination of glycopeptides by liquid*
536 *chromatography/mass spectrometry with collision-induced dissociation*. J Am Soc
537 Mass Spectrom, 1992. **3**(8): p. 804-14.
- 538 26. Huddleston, M.J., M.F. Bean, and S.A. Carr, *Collisional fragmentation of*
539 *glycopeptides by electrospray ionization LC/MS and LC/MS/MS: methods for*
540 *selective detection of glycopeptides in protein digests*. Anal Chem, 1993. **65**(7): p.
541 877-84.
- 542 27. Saudan, P., et al., *Ductus ejaculatorius peptide 99B (DUP99B), a novel Drosophila*
543 *melanogaster sex-peptide pheromone*. Eur J Biochem, 2002. **269**(3): p. 989-97.
- 544 28. Pan, X., J. Luo, and S. Li, *Bacteria-Catalyzed Arginine Glycosylation in Pathogens*
545 *and Host*. Front Cell Infect Microbiol, 2020. **10**: p. 185.
- 546 29. Date-Ito, A., et al., *Rapid evolution of the male-specific antibacterial protein*
547 *andropin gene in Drosophila*. J Mol Evol, 2002. **54**(5): p. 665-70.
- 548 30. Schaeffer, S.W., et al., *Polytene chromosomal maps of 11 Drosophila species: the*
549 *order of genomic scaffolds inferred from genetic and physical maps*. Genetics, 2008.
550 **179**(3): p. 1601-55.
- 551 31. Walter, R., W.H. Simmons, and T. Yoshimoto, *Proline specific endo- and*
552 *exopeptidases*. Mol Cell Biochem, 1980. **30**(2): p. 111-27.

- 553 32. Dougherty, P.G., A. Sahni, and D. Pei, *Understanding Cell Penetration of Cyclic*
554 *Peptides*. Chem Rev, 2019. **119**(17): p. 10241-10287.
- 555 33. Begun, D.J. and H.A. Lindfors, *Rapid evolution of genomic Acp complement in the*
556 *melanogaster subgroup of Drosophila*. Mol Biol Evol, 2005. **22**(10): p. 2010-21.
- 557 34. Samakovlis, C., et al., *The andropin gene and its product, a male-specific*
558 *antibacterial peptide in Drosophila melanogaster*. EMBO J, 1991. **10**(1): p. 163-9.
- 559 35. Lung, O., L. Kuo, and M.F. Wolfner, *Drosophila males transfer antibacterial*
560 *proteins from their accessory gland and ejaculatory duct to their mates*. J Insect
561 Physiol, 2001. **47**(6): p. 617-622.
- 562 36. Maes, E., et al., *Major O-glycans from the nest of Vespula germanica contain*
563 *phospho-ethanolamine*. Carbohydr Res, 2005. **340**(11): p. 1852-8.
- 564 37. Bordeaux, J.M., et al., *Noctilisin, a Venom Glycopeptide of Sirex noctilio*
565 *(Hymenoptera: Siricidae), Causes Needle Wilt and Defense Gene Responses in Pines*.
566 J Econ Entomol, 2014. **107**(5): p. 1931-45.
- 567 38. Hykollari, A., et al., *Tissue-specific glycosylation in the honeybee: Analysis of the N-*
568 *glycomes of Apis mellifera larvae and venom*. Biochim Biophys Acta Gen Subj, 2019.
569 **1863**(11): p. 129409.
- 570 39. Helling, F., et al., *Glycosphingolipids in insects. The amphoteric moiety, N-*
571 *acetylglucosamine-linked phosphoethanolamine, distinguishes a group of ceramide*
572 *oligosaccharides from the pupae of Calliphora vicina (Insecta: Diptera)*. Eur J
573 Biochem, 1991. **200**(2): p. 409-21.
- 574 40. Weske, B., et al., *Glycosphingolipids in insects. Chemical structures of two variants*
575 *of a glucuronic-acid-containing ceramide hexasaccharide from a pupae of Calliphora*
576 *vicina (Insecta: Diptera), distinguished by a N-acetylglucosamine-bound*
577 *phosphoethanolamine sidechain*. Eur J Biochem, 1990. **191**(2): p. 379-88.
- 578 41. Zhu, F., D. Li, and K. Chen, *Structures and functions of invertebrate glycosylation*.
579 Open Biol, 2019. **9**(1): p. 180232.
- 580 42. Zhang, L. and K.G. Ten Hagen, *O-Linked glycosylation in Drosophila melanogaster*.
581 Curr Opin Struct Biol, 2019. **56**: p. 139-145.
- 582 43. Leader, D.P., et al., *FlyAtlas 2: a new version of the Drosophila melanogaster*
583 *expression atlas with RNA-Seq, miRNA-Seq and sex-specific data*. Nucleic Acids Res,
584 2018. **46**(D1): p. D809-D815.
- 585 44. Muller, R., et al., *Characterization of mucin-type core-1 beta1-3*
586 *galactosyltransferase homologous enzymes in Drosophila melanogaster*. FEBS J,
587 2005. **272**(17): p. 4295-305.
- 588 45. Laflamme, B.A. and M.F. Wolfner, *Identification and function of proteolysis*
589 *regulators in seminal fluid*. Mol Reprod Dev, 2013. **80**(2): p. 80-101.
- 590 46. Isaac, R.E. and N. Audsley, *Insect Peptide Hormones*, in *Amino Acids, Peptides and*
591 *Proteins in Organic Chemistry*, A.B. Hughes, Editor. 2010, Wiley-VCH Verlag
592 GmbH & Co. p. 575-595.
- 593 47. Rexhepaj, A., et al., *The sex-peptide DUP99B is expressed in the male ejaculatory*
594 *duct and in the cardia of both sexes*. Eur J Biochem, 2003. **270**(21): p. 4306-14.
- 595 48. Kim, Y.J., et al., *MIPs are ancestral ligands for the sex peptide receptor*. Proc Natl
596 Acad Sci U S A, 2010. **107**(14): p. 6520-5.
- 597 49. Apostol, C.R., M. Hay, and R. Polt, *Glycopeptide drugs: A pharmacological*
598 *dimension between "Small Molecules" and "Biologics"*. Peptides, 2020. **131**: p.
599 170369.
- 600 50. Lefever, M., et al., *Structural Requirements for CNS Active Opioid Glycopeptides*. J
601 Med Chem, 2015. **58**(15): p. 5728-41.

- 602 51. Begun, D.J., et al., *Recently evolved genes identified from Drosophila yakuba and D.*
603 *erecta accessory gland expressed sequence tags.* Genetics, 2006. **172**(3): p. 1675-81.
- 604 52. Zhou, Q., et al., *On the origin of new genes in Drosophila.* Genome Res, 2008. **18**(9):
605 p. 1446-55.
- 606 53. Russo, C.A., et al., *Phylogenetic analysis and a time tree for a large drosophilid data*
607 *set (Diptera: Drosophilidae).* Zoological Journal of the Linnean Society, 2013.
608 **169**(4): p. 765-775.
- 609 54. Gao, J.-j., et al., *Phylogenetic relationships between Sophophora and Lordiphosa,*
610 *with proposition of a hypothesis on the vicariant divergences of tropical lineages*
611 *between the Old and New Worlds in the family Drosophilidae.* Molecular
612 phylogenetics and evolution, 2011. **60**(1): p. 98-107.
- 613 55. Tamura, K., S. Subramanian, and S. Kumar, *Temporal patterns of fruit fly*
614 *(Drosophila) evolution revealed by mutation clocks.* Molecular biology and evolution,
615 2004. **21**(1): p. 36-44.
- 616 56. Bulet, P., et al., *A novel inducible antibacterial peptide of Drosophila carries an O-*
617 *glycosylated substitution.* J Biol Chem, 1993. **268**(20): p. 14893-7.
- 618 57. van der Linde, K., et al., *A supermatrix-based molecular phylogeny of the family*
619 *Drosophilidae.* Genet Res (Camb), 2010. **92**(1): p. 25-38.
- 620 58. Seetharam, A.S. and G.W. Stuart, *Whole genome phylogeny for 21 Drosophila*
621 *species using predicted 2b-RAD fragments.* PeerJ, 2013. **1**: p. e226.
- 622 59. Suvorov, A., et al., *Widespread introgression across a phylogeny of 155 Drosophila*
623 *genomes.* bioRxiv, 2020.
- 624 60. Kim, B.Y., et al., *Highly contiguous assemblies of 101 drosophilid genomes.* bioRxiv,
625 2020.
- 626 61. Nishihara, S. *Functional analysis of glycosylation using Drosophila melanogaster.*
627 Glycoconjugate J, 2020. **37**: p.1–14.
- 628 62. Rabel, D., et al., *Primary Structure and in Vitro Antibacterial Properties of the*
629 *Drosophila melanogaster Attacin C Pro-domain.* J Biol Chem, 2004. **279**(15): p.
630 14853-9.

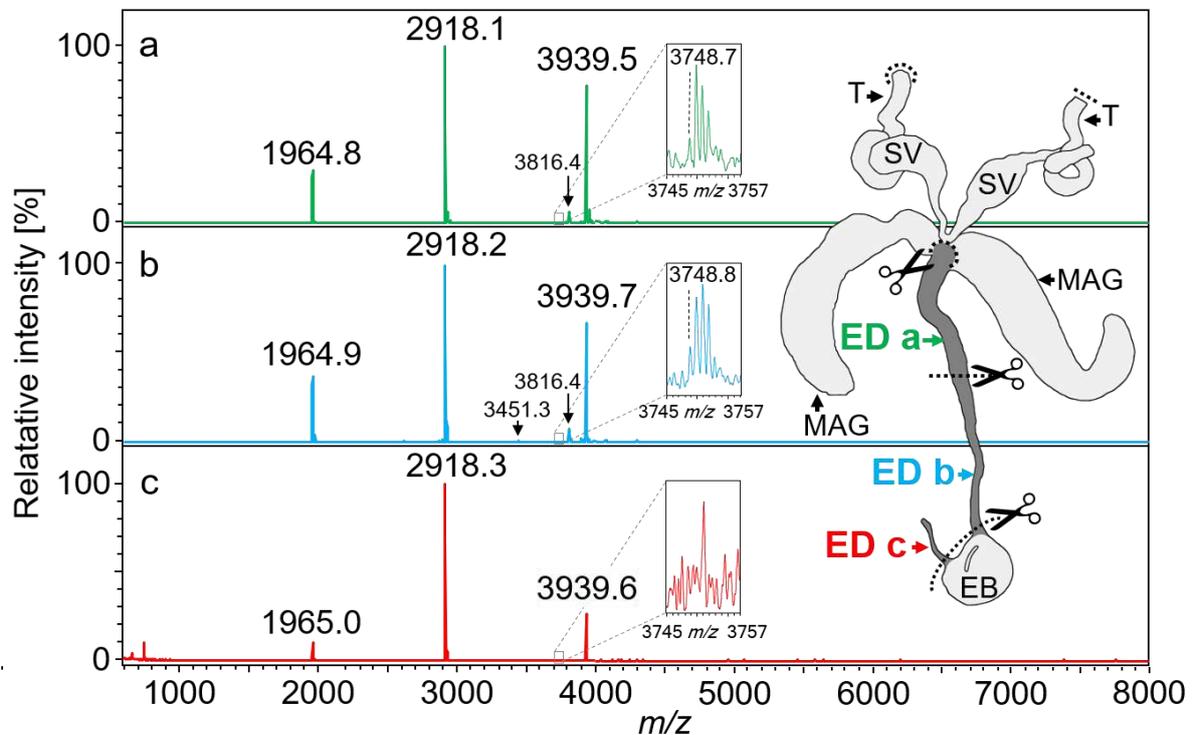
631



632

633 **Graphic abstract**

634

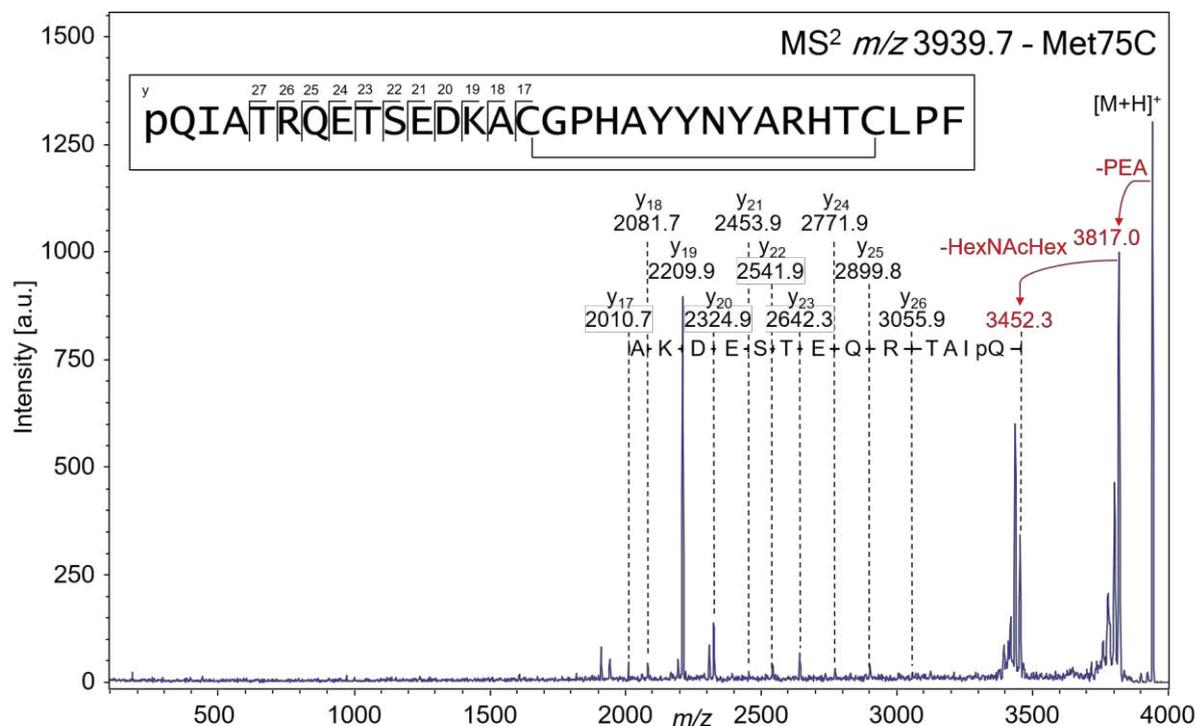


635

636 **Fig. 1. MALDI-TOF mass spectra of a single male ejaculatory duct (ED).** Spectra from
 637 (a) the dilated and (b) narrow regions of the anterior ejaculatory duct (ED), and (c) the
 638 posterior ED of a w^{1118} male. These regions are defined, relative to the other male
 639 reproductive tissues (T, testis; SV, seminal vesicle; MAG, male accessory gland) in the
 640 overlaid cartoon. Spectra were recorded in the mass ranges m/z 600-4000 and 3000-8000 and
 641 merged. All tissue sections gave prominent ion signals at m/z 1964.9, 2918.2 and 3939.6,
 642 corresponding to the molecular ions of Acp54A1 (CG34098), CG42782 and Met75Ca/b
 643 (CG18064 and CG32197), respectively. The minor ion signals at m/z 3451.3 and 3816.4 as

644 well as the abundant ion at m/z 3939.7 produced highly similar fragment patterns (see Fig. 2),
 645 hence considered to be peptide forms. A low abundant ion signal at m/z 3748.7 was noted in
 646 some spectra (see inset). This ion signal was more abundant in spectra of the ED extract and
 647 was assigned to andropin by MALDI-TOF MSMS.

648



649

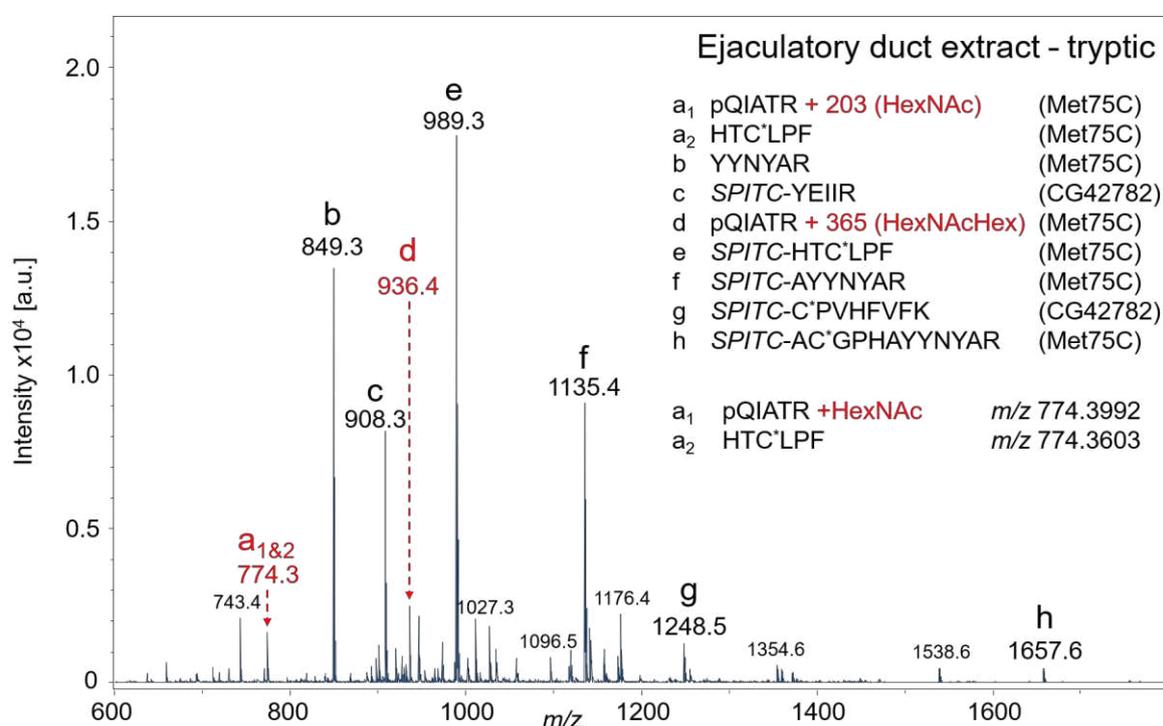
650 Fig. 2. MALDI-TOF MS fragment spectrum of the Met75C precursor ion at m/z 3939.7
 651 obtained from a single ED. The spectrum of the non-reduced and non-derivatized peptide
 652 reveals the loss of phosphoethanolamine (PEA) and a disaccharide (HexNAcHex) as
 653 indicated. The fragment spectrum also comprises the y-ion series of the linear section of
 654 Met75C.

655

656 Table 1. Gene symbols, names, theoretical masses ($[M+H]^+$), deconvoluted experimental
 657 masses obtained from FT-ICR MS as well as the corresponding mass difference (Δ Mass) and
 658 sequences of ED peptides. Abbreviations: Carbamidomethylated cysteine, C*; pyroglutamic
 659 acid, pQ; N-acetylhexosamine, HexNAc; hexose, Hex; phosphoethanolamine, PEA; fucose,
 660 Fuc.

| Gene symbol | Name | [M+H] ⁺ | FT-ICR MS | Δ Mass | Sequence |
|-----------------|-----------|--------------------|-----------|--------|---|
| CG34098 | Acp54A1 | 1964.9768 | | | MKCRLGFVKRGGQCTWP |
| | | 2081.0354 | 2081.0351 | 0.0003 | MKC*RLGFVKRGGQC*TWP |
| CG42782 | | 2918.3575 | | | YEIIRQCPVHFVFKNNYCQYQPM |
| | | 3034.4161 | 3034.4153 | 0.0008 | YEIIRQC*PVHFVFKNNYC*QYQPM |
| CG32197/CG18064 | Met75Ca/b | 3451.5583 | | | pQIATRQETSEDKACGPHAYNYARHTCLPF |
| | | 3567.6168 | 3567.6134 | 0.0034 | pQIATRQETSEDKAC*GPHAYNYARHTC*LPF |
| CG32197/CG18064 | Met75Ca/b | 3654.6376 | | | pQIAT(HexNAc)RQETSEDKACGPHAYNYARHTCLPF |
| | | 3770.6962 | 3770.6917 | 0.0045 | pQIAT(HexNAc)RQETSEDKAC*GPHAYNYARHTC*LPF |
| CG32197/CG18064 | Met75Ca/b | 3816.6905 | | | pQIAT(HexNAcHex)RQETSEDKACGPHAYNYARHTCLPF |
| | | 3932.7490 | 3932.7483 | 0.0007 | pQIAT(HexNAcHex)RQETSEDKAC*GPHAYNYARHTC*LPF |
| CG32197/CG18064 | Met75Ca/b | 3939.6990 | | | pQIAT(HexNAcHexPEA)RQETSEDKACGPHAYNYARHTCLPF |
| | | 4055.7576 | 4055.7571 | 0.0005 | pQIAT(HexNAcHexPEA)RQETSEDKAC*GPHAYNYARHTC*LPF |
| CG1361 | Andropin | 3749.0789 | 3749.0764 | 0.0025 | VFIDILDKVENAIHNAQVIGIFAKPFPEKLINPK |
| CG33495 | DUP99b | 4930.2041 | | | pQDRN(2HexNAc3Hex2Fuc)DTEWIQSQKDREKWCRLNLGPYLGGRG |
| | | 5046.2627 | 5046.2587 | 0.0040 | pQDRN(2HexNAc3Hex2Fuc)DTEWIQSQKDREKWC*RLNLGPYLGGRG* |

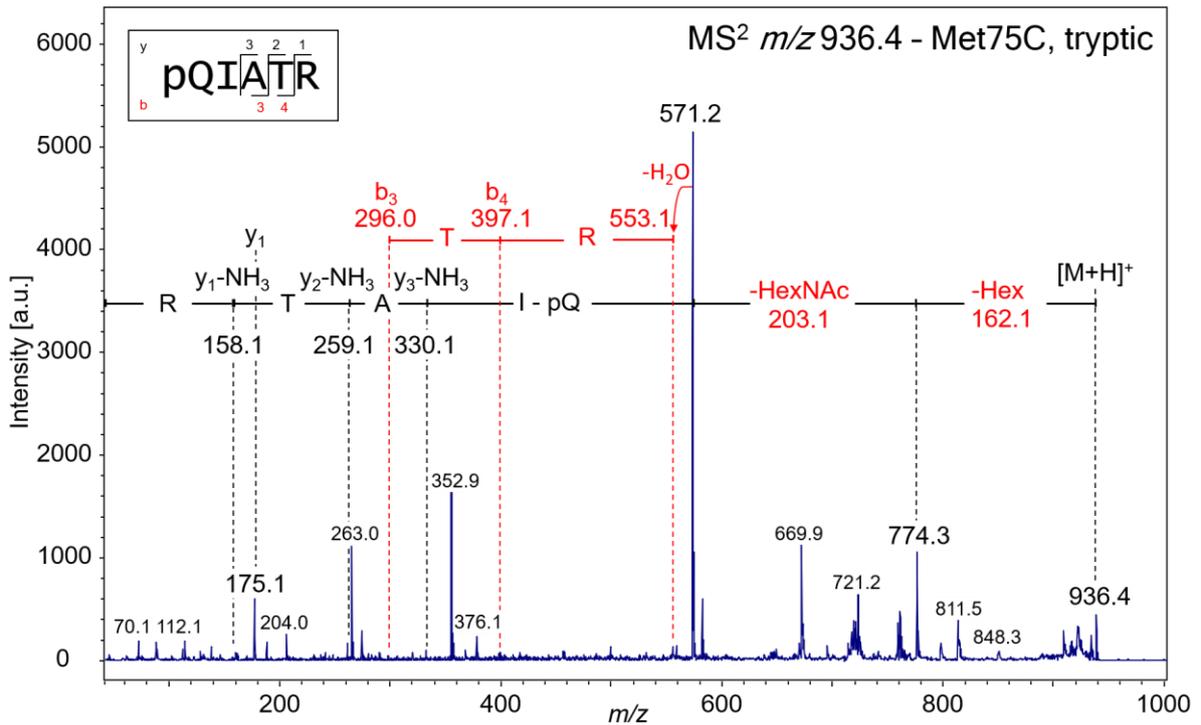
661



662

663 Fig. 3. MALDI-TOF mass spectrum of the extracted ED peptides after reduction,
 664 alkylation, tryptic digestion and sulfonation . The identity of the tryptic fragments (a-h)
 665 was confirmed by MSMS analysis.

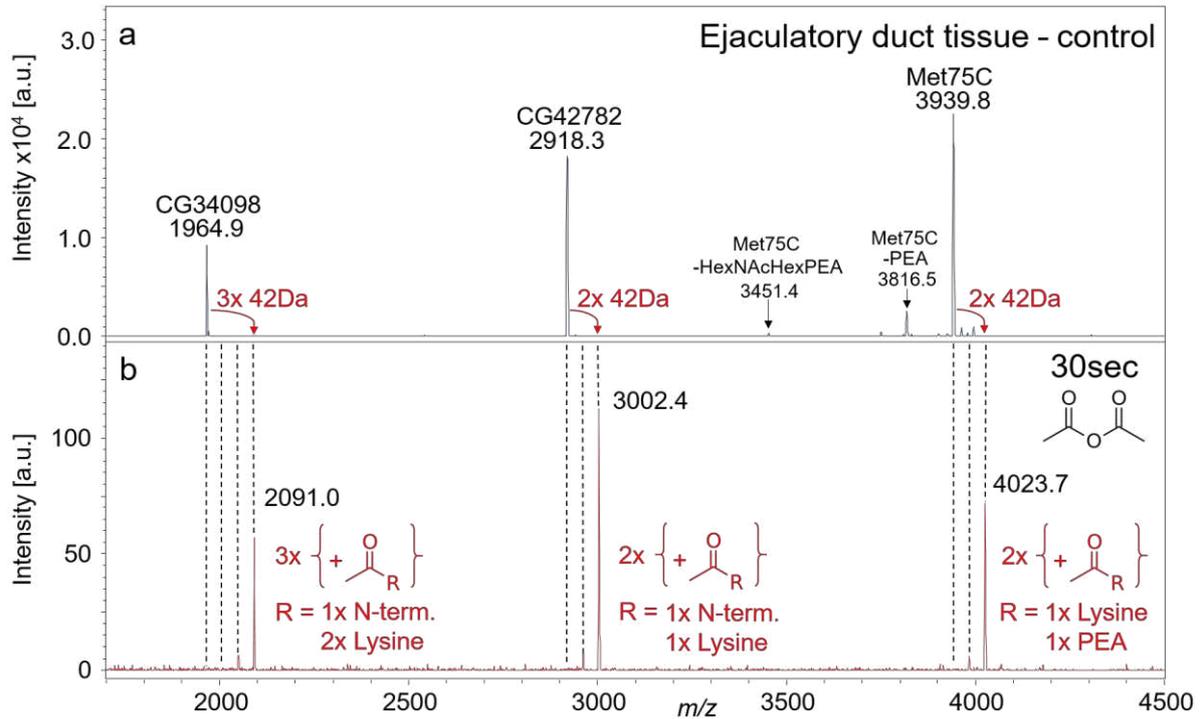
666



667

668 Fig. 4. MALDI-TOF MS fragment spectrum of the Met75C tryptic peptide ion [M+H]⁺,
 669 m/z 936.4. The loss of a hexose and N-acetylhexosamine from the peptide and the y-ion
 670 series for the pentapeptide, pQIATR, are indicated.

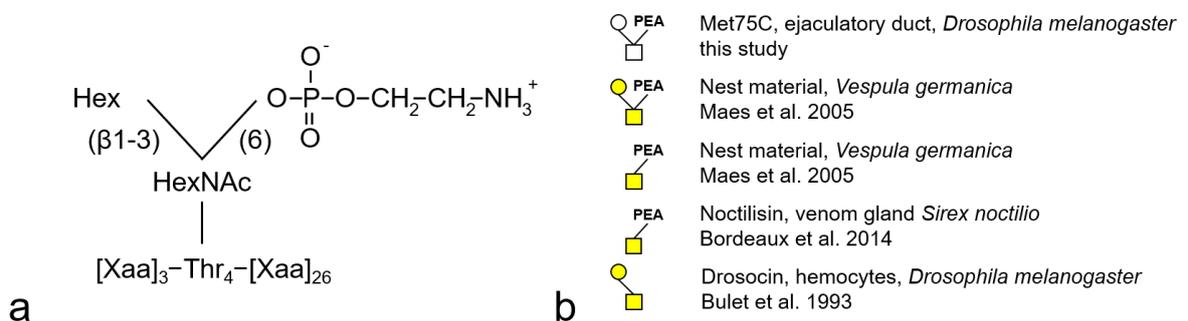
671



672

673 Fig. 5 MALDI-TOF mass spectra by direct tissue analysis of a single male ejaculatory
 674 duct (ED) either untreated (a) or treated (b) for 30 seconds with acetic anhydride to
 675 acetylate primary amines. The increase in mass of the molecular ions for Acp54A1,
 676 CG42782 and Met75C by the addition of two or three acetyl groups (42 Da) is labelled with
 677 red arrows. Dashed lines indicate the mass shifts between precursors, partially acetylated
 678 peptides and final products.

679

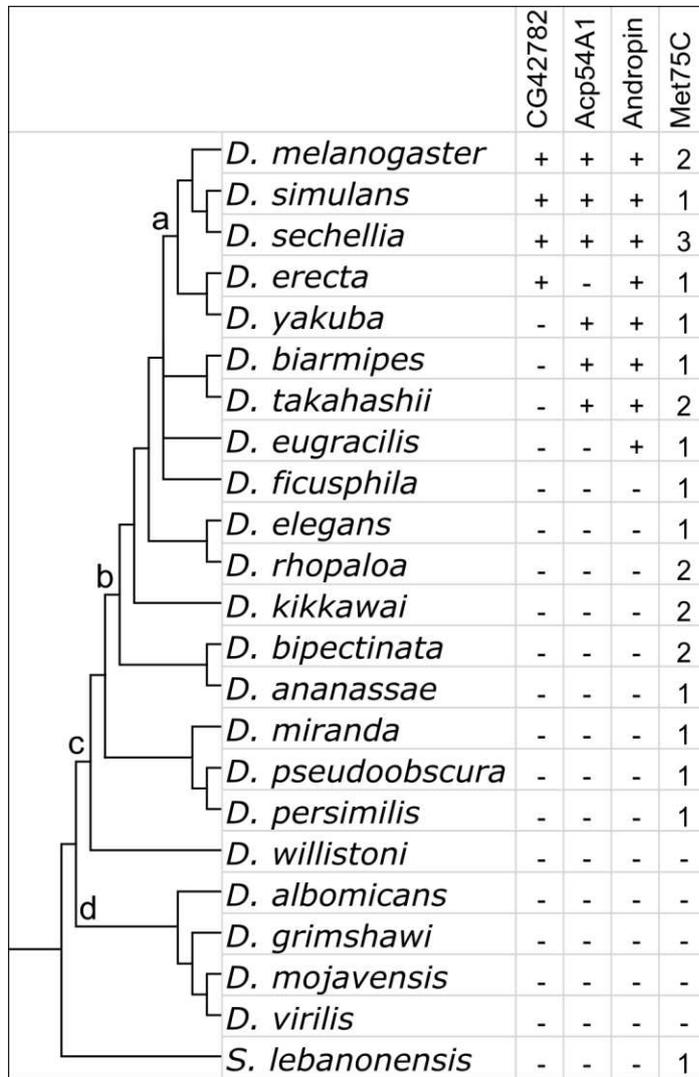


680

681 Fig. 6. The structure of the glycan of Met75C and similar insect mucin-like O-glycans.
 682 (a) The proposed structure of the zwitterionic Met75C glycan attached to Thr and (b)
 683 structures of selected insect O-glycans that have N-acetylgalactosamine (GalNAc, yellow

684 square) directly linked to Thr or Ser. GalNAc are extended with galactose (yellow circle) and
 685 phosphoethanolamine (PEA). Maes *et al.* 2005[36]; Bulet *et al.* 1993[56]; Bordeaux *et al.*
 686 1965[37].

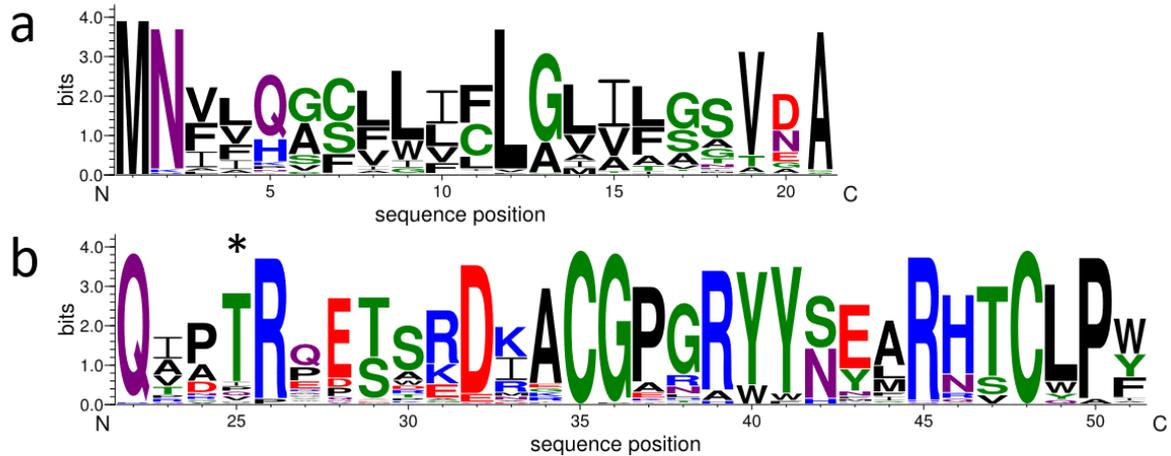
687



688

689 **Fig. 7. Phylogram showing the distribution of orthologs in the genomes of traditional**
 690 **representatives of various *Drosophila* groups.** Topology adapted from [57-60] with
 691 disputed clades illustrated as polytomy. Names of clades referred to in the text are labelled at
 692 the nodes (a) *melanogaster* subgroup (b) *melanogaster* group (c) subgenus *Sophophora* (d)
 693 subgenus *Drosophila*. Corresponding accession numbers and sequences are given in the
 694 supplementary material.

695



696

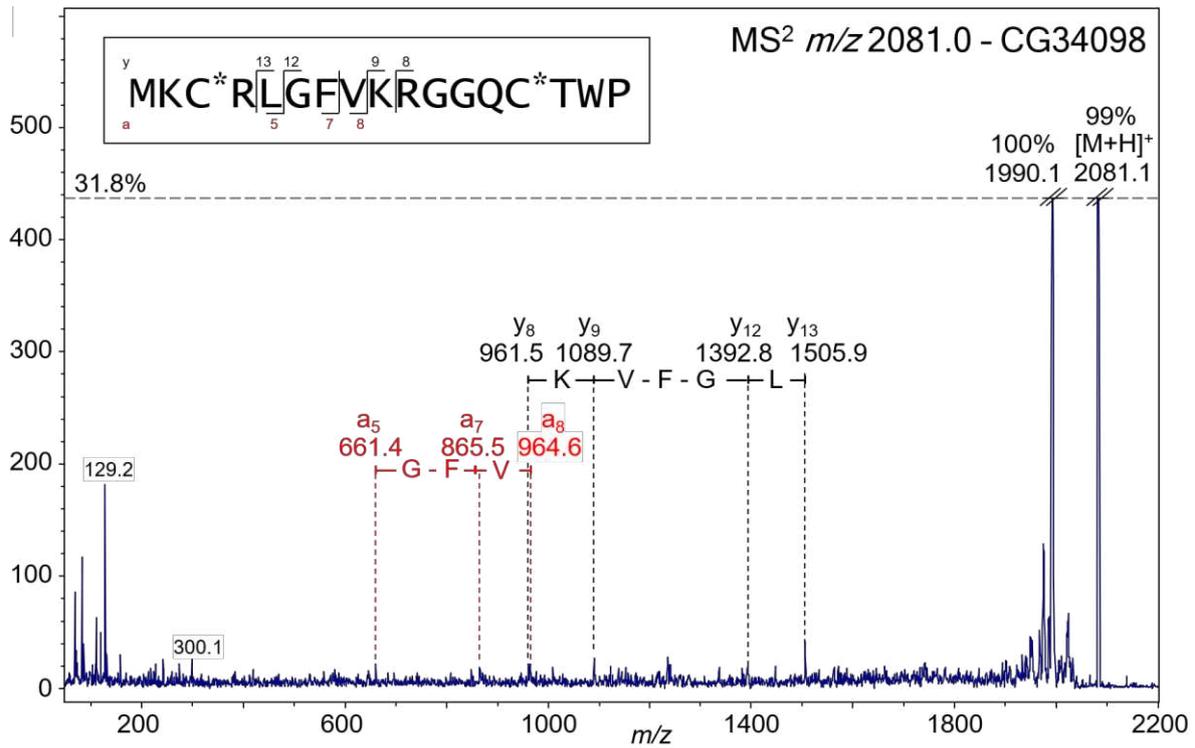
697 **Fig. 8. A graphical representation of the conservation of amino acids of 79 aligned**
 698 **sequences of Met75 ortho- and paralogs predicted from genome data of 51 *Drosophila***
 699 ***species and *Scaptodrosophila lebanonensis*.*** The overall height of a stack indicates the
 700 sequence conservation at that position, while the relative height of letters indicates the
 701 frequency of occurrence of amino acids in one particular position. Consensus range of the
 702 predicted signal peptides (a) and the predicted mature peptides after signal peptide cleavage
 703 (b). Amino acid residues are colored according to their biochemical properties; polar (green),
 704 neutral (purple), basic (blue), acidic (red), hydrophobic (black). The proposed glycosylation
 705 site in *D. melanogaster* is indicated by an asterisk. The width of each stack is proportional to
 706 the fraction of valid symbols in that position. Rare sequence positions due to insertions or
 707 extension were excluded applying a 30%-cut-off. A complete sequence logo and the
 708 corresponding alignment are given in the supplementary material (Supplementary Fig.4;
 709 Met75C_alignment.fas).

710

711

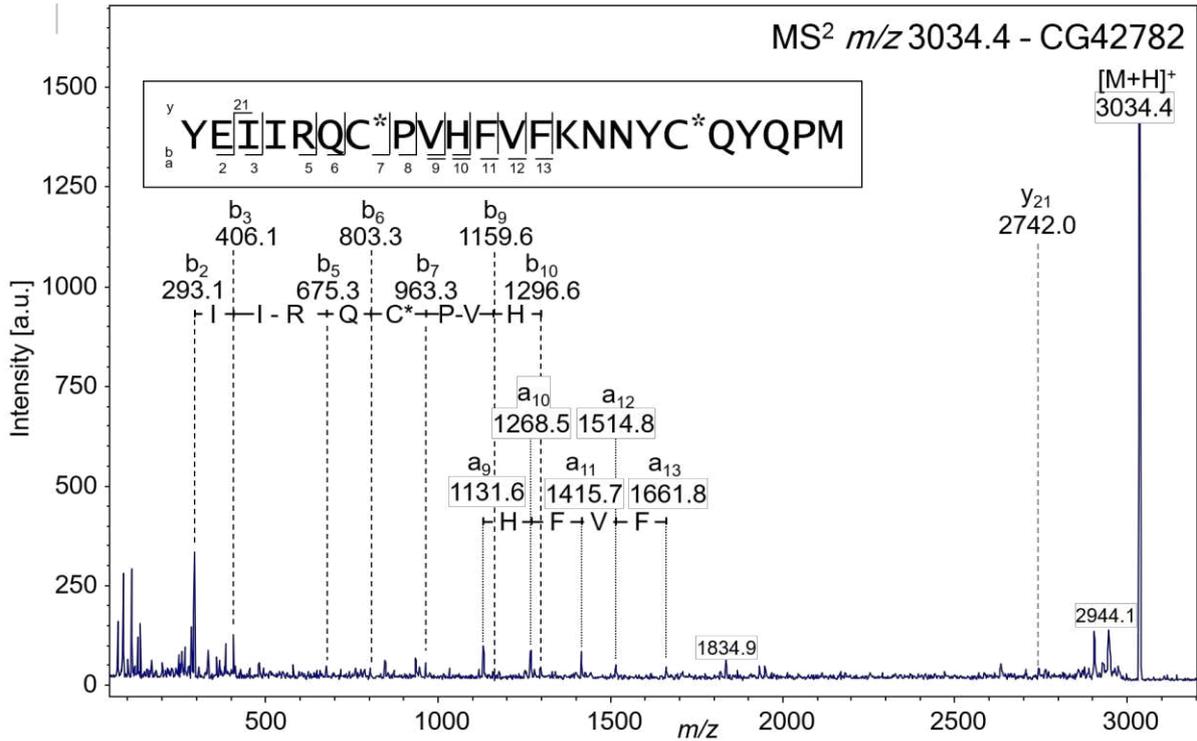
712

713 Supplementary material



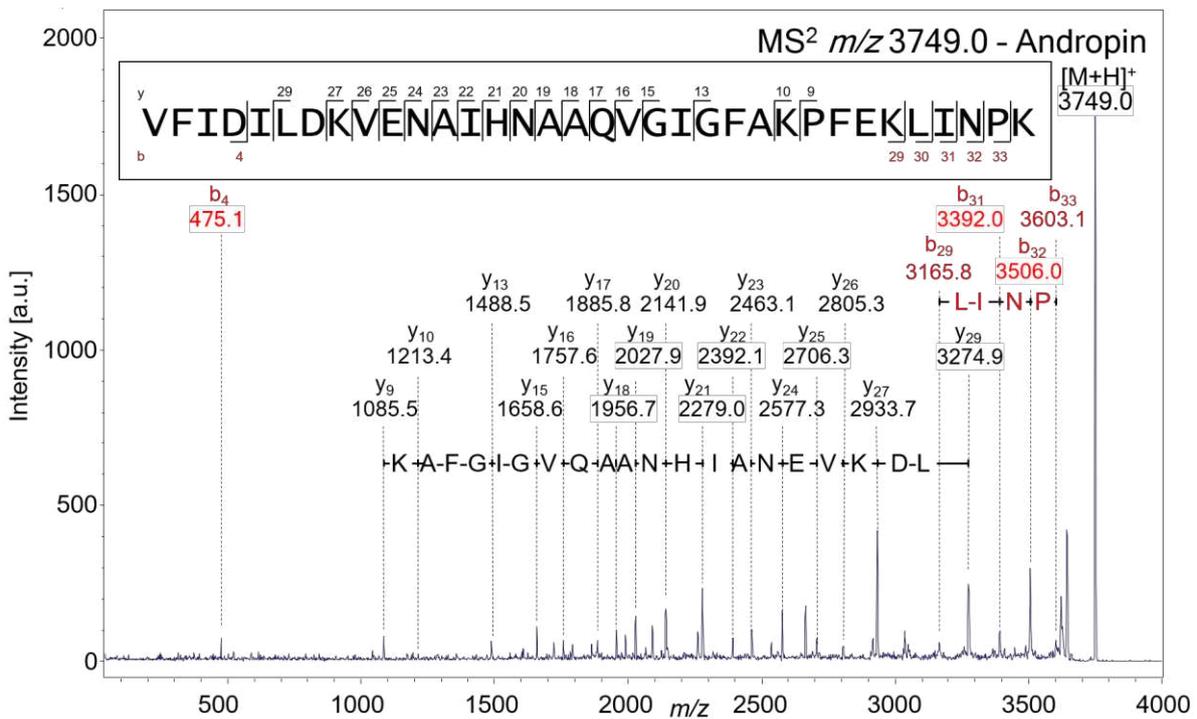
714

715 **Supplementary Fig. 1. MALDI-TOF MSMS fragment spectrum of CG34098 from a**
 716 **reduced and alkylated ED extract.** Product-ion spectrum resulting from fragmentation of
 717 the isolated molecular ion at m/z 2081.1, assigned to the CG34098 peptide. A selection of
 718 diagnostic fragments is labelled in the spectrum and the sequence coverage by all observed *a*-
 719 and *y*-series fragment ions is given in the schematic representation in the inset.



720

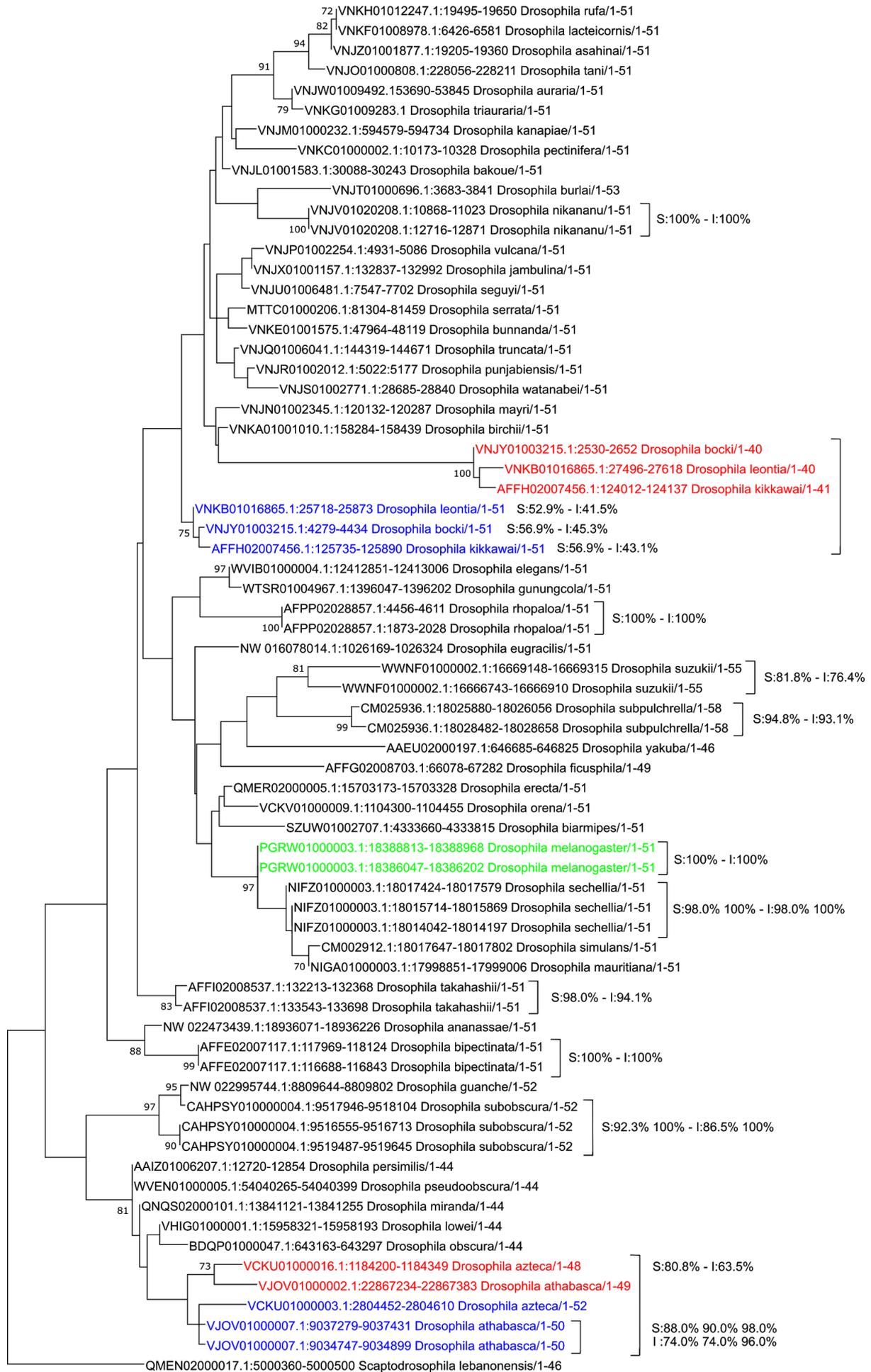
721 **Supplementary Fig. 2. MALDI-TOF MSMS fragment spectrum of CG42782 from a**
 722 **reduced and alkylated ED extract.** Product-ion spectrum resulting from fragmentation of
 723 the isolated molecular ion at m/z 3034.4, assigned to the CG42782 peptide. A selection of
 724 diagnostic fragments is labelled in the spectrum and the sequence coverage by all observed a -
 725 b - and y -series fragment ions is given in the schematic representation in the inset.



726

727 **Supplementary Fig. 3. MALDI-TOF MSMS fragment spectrum of andropin from a**
728 **reduced and alkylated ED extract.** Product-ion spectrum resulting from fragmentation of
729 the isolated molecular ion at m/z 3749.0 assigned to the andropin peptide. A selection of
730 diagnostic fragments is labelled in the spectrum and the sequence coverage by all observed *b*-
731 and *y*-type fragment ions is given in the schematic representation in the inset.

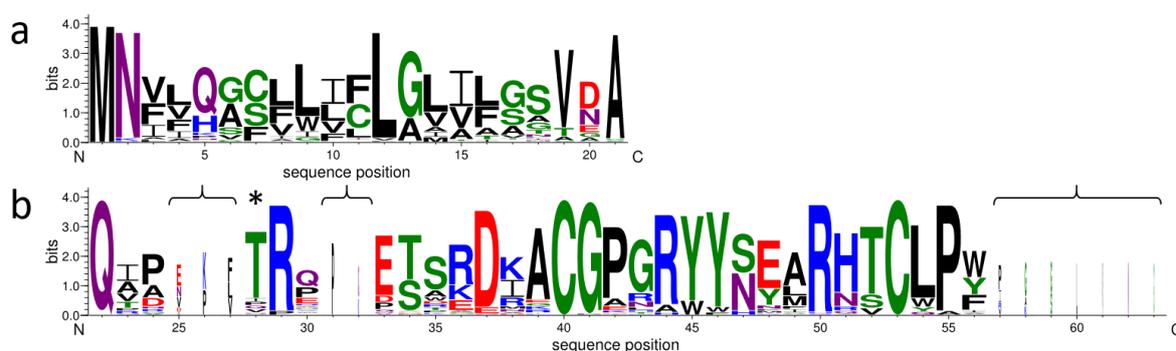
732



734 **Supplementary Fig. 4. Phylogenetic tree of Met75C orthologs and paralogs.** The
 735 evolutionary history of 69 Met75C-related genes with 63 positions from 52 species was
 736 inferred using the Maximum Likelihood method and the Le Gascuel (2008) amino acid
 737 substitution model with five discrete gamma categories. The tree with the highest log
 738 likelihood (-2333.99) is shown. Sequences of a common species have been included if they
 739 were located in the same contig to exclude sequences which are no paralogs with the
 740 exception of *D. azteca* and *D. athabasca*, where paralogs are located at different
 741 chromosomes. Gene duplication events are highlighted with square brackets labeled with the
 742 pairwise percent similarity (S) / identity (I). Note the presence of multiple independent gene
 743 duplication events with the exception of a common duplication in the ancestral line of *D.*
 744 *leontia*, *D. bocki*, and *D. kikkawai*, as well as a common evolutionary event shared by *D.*
 745 *azteca* and *D. athabasca* (highlighted in red and blue). Nodes which occurred in at least 70%
 746 out of 500 bootstrap replicates are labeled at the branches.

747

748



749

750 **Supplementary Fig. 4. . The complete graphical representation of the conservation of**
 751 **amino acids of 79 aligned sequences of Met75C ortho- and paralogs predicted from**
 752 **genome data of 51 *Drosophila* species and *Scaptodrosophila lebanonensis*.** The overall
 753 height of a stack indicates the sequence conservation at that position, while the relative height
 754 of letters indicates the frequency of occurrence of amino acids in one particular position. (a)
 755 the predicted signal peptide for the majority of sequences and (b) the predicted mature
 756 peptides after signal peptide cleavage. Amino acid residues are colored according to their
 757 biochemical properties; polar (green), neutral (purple), basic (blue), acidic (red), hydrophobic
 758 (black). The proposed glycosylation site in *D. melanogaster* is indicated by an asterisk. The

759 width of each stack is proportional to the fraction of valid symbols in that position. Rare
760 sequence positions due to insertions or extensions are marked with braces.

761

762

763

764