

1 **A gene expression panel for estimating age in**
2 **males and females of the sleeping sickness**
3 **vector *Glossina morsitans*.**

4 Eric R. Lucas^{1*}, Alistair C. Darby², Stephen J. Torr¹, Martin J. Donnelly^{1,3}

5 ¹Liverpool School of Tropical Medicine, Pembroke Place, Liverpool L3 5QA, UK.

6 ²Institute of Integrative Biology, University of Liverpool, Biosciences Building, Crown Street,
7 Liverpool, L69 7ZB, UK

8 ³Wellcome Sanger Institute, Hinxton, Cambridge CB10 1SA, UK.

9 * Author for correspondence: eric.lucas@lstmed.ac.uk

10 **Abstract**

11 Many vector-borne diseases are controlled by methods that kill the insect vectors responsible
12 for disease transmission. Recording the age structure of vector populations provides
13 information on mortality rates and vectorial capacity, and should form part of the detailed
14 monitoring that occurs in the wake of control programmes, yet tools for obtaining estimates of
15 individual age remain limited. We investigate the potential of using markers of gene
16 expression to predict age in tsetse flies, which are the vectors of deadly and economically
17 damaging African trypanosomiases. We use RNAseq to identify candidate expression
18 markers, and test these markers using qPCR in laboratory-reared *Glossina morsitans*
19 *morsitans* of known age. Measuring the expression of six genes was sufficient to obtain a
20 prediction of age with root mean squared error of less than 8 days, while just two genes were
21 sufficient to classify flies into age categories of ≤ 15 and > 15 days old. Further testing of these
22 markers in field-caught samples and in other species will determine the accuracy of these
23 markers in the field.

24 **Keywords:** tsetse flies, age prediction, transcriptomics, machine learning

25 **1 Introduction**

26 Vector-borne diseases represent major threats to health and livelihood world-wide, being
27 directly responsible for 680,000 deaths annually (Roth et al. 2018), as well as causing huge
28 economic damage to livestock (Eisler et al. 2003, Shaw 2004). Control of the vectors that
29 transmit these diseases is an integral tool for reducing disease burden (Wilson et al. 2020).
30 The metric of success for these control programmes is a reduction in disease burden in the
31 host population. However, when vector control is accompanied by other interventions such as
32 screening and treating the host population for the disease, the contribution of vector control to
33 the subsequent reduction of disease can be hard to determine (World Health Organization
34 2012). Conversely, while the impact on the vector population may not bear a simple
35 relationship to disease burden, it is a direct outcome of vector control. Control efforts should
36 thus be accompanied by detailed monitoring of the targeted vector populations, to estimate
37 impact, to monitor population recovery and to understand the transmission dynamics of the
38 disease. Mostly, monitoring currently relies on counting the number of vectors caught in
39 sentinel traps, which can be greatly affected by trapping method, effort and efficacy, and may
40 only partly reflect the ability of the vector population to transmit disease (Wilson et al. 2015).

41 One aspect of vector monitoring that has been particularly challenging is the quantification of
42 the age-distribution (demographics) of natural populations (Caragata et al. 2011, Cook et al.
43 2006, Sikulu et al. 2010). Estimating vector age is important for two reasons. First, it can
44 provide a measure of the effectiveness of vector control because increased adult mortality
45 should lead to a younger population age structure. Importantly, this measure of control
46 effectiveness is independent of catch size and trapping effort because only the distribution of
47 age needs to be known. Second, in most cases, the probability that an individual vector is
48 infectious for a given disease increases with age (Dye 1992, Woolhouse & Hargrove 1998).
49 Before transmitting the disease, vectors first need to have taken an infected blood meal, and
50 there is then typically a delay between acquisition of infection and onward transmission due
51 to the need for the pathogen to replicate and/or mature. Age grading is therefore useful to
52 determine the proportion of individuals old enough to transmit disease.

53 Tsetse flies (genus *Glossina*) are the vectors of Human African Trypanosomiasis (HAT, or
54 sleeping sickness) and Animal African Trypanosomiasis (AAT, or nagana). HAT is, without
55 treatment, a fatal disease endemic to sub-Saharan Africa (Franco et al. 2014), while AAT
56 presents a major economic burden to rural communities by affecting livestock (Eisler et al.

57 2003). Being a disease primarily of animals and with reservoirs across multiple species, AAT
58 cannot be controlled through treatment alone and is thus highly dependant on vector control
59 (Holmes 2013). *G. morsitans morsitans* is a major vector of AAT in East and Southern Africa
60 and can also transmit HAT (Dale et al. 1995). Catch rates of this species in the wake of vector
61 control can be extremely low (Kgori et al. 2006, Vale et al. 1988, Van den Bossche 1997),
62 making it particularly challenging to conduct ongoing monitoring of important populations. It
63 is therefore all the more important to extract as much information as possible from the limited
64 number of flies obtained.

65 As is the case for all insect vectors, a means to accurately determine the age of tsetse flies is a
66 valuable but elusive goal, and current methods have many shortcomings. Laborious ovary
67 dissections can be used to age females up to their fourth ovarian cycle (Hargrove 2012), but
68 this technique requires specialist dissection skills and cannot be applied to males, despite
69 males being at least as competent at transmission as females, and perhaps more so (Dale et al.
70 1995, Maudlin et al. 1990). Estimates of age based on wing damage (Hargrove 1990) or
71 analysis of pteridines have also been used (Langley et al. 1988, Lehane & Hargrove 1988),
72 but experience in practical applications has shown that measurements in the field vary
73 enormously (for example in mosquitoes: (Lardeux et al. 2000, Penilla et al. 2002)) and cannot
74 be used to reliably estimate age on an individual basis (Hargrove 2020).

75 Here we explore the value of using gene expression to estimate age in tsetse flies. This
76 method has previously been tested in mosquitoes (Caragata et al. 2011, Cook et al. 2006),
77 with encouraging results, but has yet to be applied in tsetse. We use laboratory-reared *G.*
78 *morsitans* as a proof of concept, and show that measuring the expression of just six genes can
79 estimate the age of both male and female tsetse flies with a root mean squared error of less
80 than 8 days. We also trained models to classify tsetse into those younger or older than 15
81 days, since flies younger than 15 days are unlikely to harbour a mature trypanosome infection
82 (Dale et al. 1995), and found that just two genes are sufficient for 95% accurate classification.

83 **2 Methods**

84 **2.1 Sample collection and RNA extraction**

85 *G. morsitans morsitans* individuals were collected from colonies maintained at the Liverpool
86 School of Tropical Medicine. Colonies are kept in meshed boxes (cages) at $26^{\circ}\text{C} \pm 2^{\circ}\text{C}$ and
87 $72 \pm 4\%$ humidity, with a 12hr light-dark photoperiod, and fed three times per week using

88 defibrinated horse blood (TCS Biosciences Ltd., Buckingham, UK) provided through silicon-
89 membrane feeders. Pupae are regularly collected and allowed to emerge to form new cages.
90 Each fly cage contains flies which eclosed over a 2-3 day window, and thus the age of all flies
91 in the cage are known to a precision of either 2 or 3 days. The ages reported here are the
92 middle of the age range (eg: a fly aged 13-15 days or 13-16 days is reported as 14 days old).
93 The age of the samples ranged from 2 to 62 days. While reproductive status of females was
94 not measured precisely, we tried to include a range of physiological states (based on visual
95 inspection of the size of the abdomen) within each age group, so that genes could be identified
96 that are predictive of age in spite of variation caused by the ovarian cycle. Overall, 505 flies
97 were collected (301 female and 204 male, Supplementary Data S1).

98 For sample collection, fly cages were briefly transferred to a cold room (4 °C) where flies to
99 be collected were removed from the cage once quiescent and decapitated. Heads were placed
100 into RNAlater and stored at -20 °C. In case repeated exposure to the cold room created
101 alterations in gene expression, we minimised this exposure by never collecting flies from a
102 given cage more than three times over the course of the experiment. No more than two flies
103 were collected from a cage on a given day, for three reasons. Firstly, we wanted to make sure
104 that flies were obtained from a range of different cages in order to avoid issues of results
105 being confounded by cage of origin (such as an infection specific to one cage of flies). We
106 therefore never obtained more than six flies from a single cage over the course of the
107 experiment. Second, we wanted to minimise the time that samples spent at temperatures
108 above -20 °C after death, limiting the number of samples that could be collected in a single
109 sitting. Third, all flies were collected at the same approximate time of day (morning) to
110 minimise gene expression variation due to circadian cycles (Rund et al. 2011), limiting the
111 number of collections that could be performed on the same day.

112 RNA was extracted from individual fly heads. Single heads contain enough material for RNA
113 sequencing and can easily be removed without the need for precise dissection, providing a
114 quick and convenient tissue for sampling. We avoided the abdomen because of the important
115 effect that sex and the ovarian cycle would have on gene expression in these tissues. RNA
116 extractions were performed using PicoPure kits (Arcturus), increasing the volume of
117 extraction buffer and alcohol to 120µl. cDNA libraries were prepared using SuperScript III
118 Reverse Transcriptase (Invitrogen).

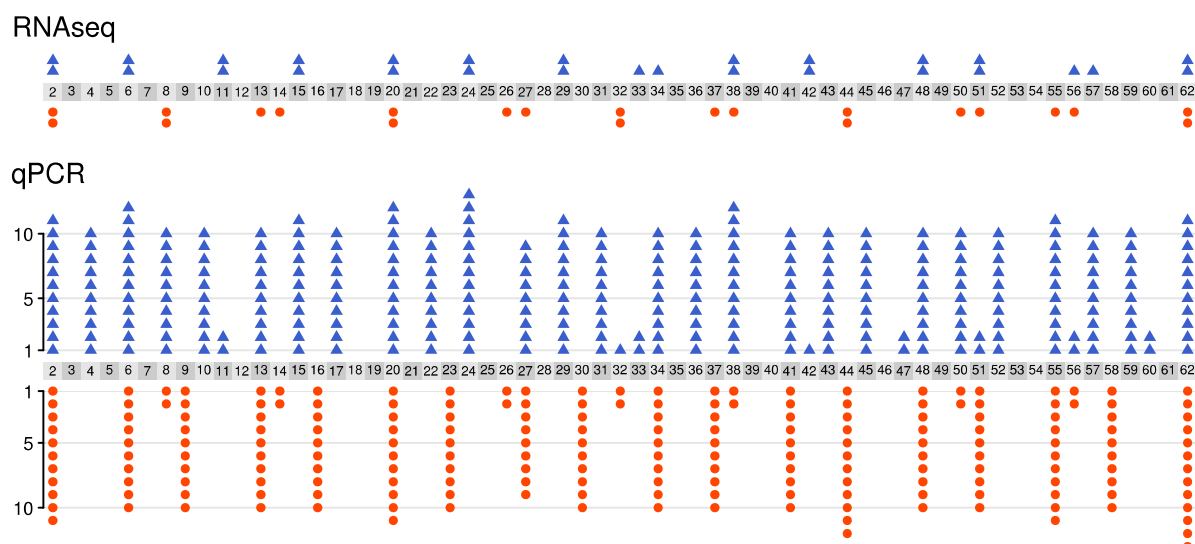


Figure 1: Number of samples used for RNAseq (top; total = 50) and qPCR (bottom; total = 498), split by age category (2 - 62 days old). Individual female and male flies shown as blue triangles and orange circles respectively.

119 2.2 Sequencing

120 cDNA libraries from 22 male and 28 female individual flies ranging in age from 2 to 62 days
121 post-eclosion (Fig. 1, Supplementary Data S1) were sent to the Liverpool Centre for Genomic
122 Research (CGR) for 150bp paired-end sequencing on an Illumina HiSeq 4000 sequencer.
123 Strand-specific library preparation was performed using NEBNext poly A selection and Ultra
124 Directional RNA library preparation kits, producing an average of 23.8 million reads per
125 sample. Reads were then trimmed as part of the CGR's genomic pipeline using Cutadapt
126 version 1.2.1 (Martin 2011) with option -O 3 to remove Illumina adapter sequences, and
127 Sickle version 1.2 (<https://github.com/najoshi/sickle/releases/tag/v1.2>) with a minimum
128 window quality score of 20. Reads shorter than 20 bp after trimming were removed and
129 subsequently unpaired reads were excluded. Data were quality checked using FastQC
130 (Andrews 2010) before analysis.

131 2.3 RNAseq analysis

132 Trimmed reads were aligned to the GmorY1.9 genome using STAR aligner version 2.7.0
133 (Dobin et al. 2013) using the --quantMode GeneCounts option to obtain mapping counts for
134 each gene.

135 Differential expression analysis was performed using the R package *EdgeR* (Robinson et al.
136 2010), with library size normalisation performed using Trimmed Mean of M-values

137 (Robinson & Oshlack 2010) and dispersion calculated with trended and tag-wise estimates.
138 Genes with fewer than 10 reads across all 50 samples were excluded from the analysis. All
139 plotting figures show expression measured as reads per million reads (RPM) from normalised
140 library sizes. Association of gene expression with age and sex was tested using generalised
141 linear modelling (glm) implemented in *edgeR*, with age coded as a continuous variable and
142 sex as a categorical variable. Preliminary analysis found little evidence of an important effect
143 of the number of times a colony was exposed to the cold room on gene expression, but there
144 was a significant effect of the number of days since flies had received a blood meal
145 (Supplementary Data S2). We therefore controlled for days since receiving a blood meal by
146 including it as a fixed continuous factor in the glm. False discovery rate control was set at 1%
147 using the R package *fdrtool* (Klaus & Strimmer 2015).

148 Gene clustering analysis was performed with the *WGCNA* package in R (Langfelder &
149 Horvath 2008), using the normalised read counts generated by *edgeR* and keeping only the
150 5000 genes with the highest variance in expression. We used the hybrid module merging
151 algorithm with a deep split value of 4, a minimum cluster size of 30 and a power parameter of
152 8, followed by module merging using the absolute value of the correlation coefficient between
153 eigengenes as a distance matrix and a merging threshold of 0.2.

154 Prediction of age based on normalised read counts from the RNAseq data was performed
155 using lasso regression implemented with the *glmnet* package in R (Friedman et al. 2010). As
156 the aim was to find genes with consistently high predictive value for age, we explored a range
157 of lasso parameters. This exploratory procedure is recorded in detail in the R script
158 “02_lasso.r” provided on GitHub (<https://github.com/EricRLucas/TsetseAgeMarkers>).

159 **2.4 Primer design and qPCR**

160 Based on the results of the RNAseq analysis, 16 genes were short-listed to be tested as qPCR
161 markers of age in *G. morsitans*, with two further genes being identified as suitable
162 housekeeping genes for our purposes (i.e.: showed minimal variation in expression in the
163 conditions included in our study and no evidence of association with age). Primers were
164 designed for these genes based on the GmorY1.9 genome using NCBI Primer blast (Ye et al.
165 2012). Where possible, amplicons were designed to span exon junctions. Based on testing
166 amplification efficiency using 1:3 serial dilutions, the 10 best primer pairs for age-predictive
167 genes, and the two primer pairs for housekeeping genes, were kept for use in the study and
168 applied to 499 samples (298 females and 201 males), including 44 of the samples used for

169 RNAseq (the remaining 6 samples had too little cDNA left to be included in the qPCR study).
170 One of the samples failed to produce a Ct value for several genes and was therefore excluded
171 from subsequent analysis, leaving 498 samples (Fig. 1). All primers used in this study are
172 listed in Supplementary Data S3.

173 qPCR was run on a AriaMX RealTime PCR instrument in a total volume of 20 μ l, containing
174 10 μ l of SYBR 2x MM, 1.2 μ l of forward primer (5 μ M), 1.2 μ l of reverse primer (5 μ M), 6.6
175 μ l of nuclease-free water and 1 μ l of genomic DNA. Reaction conditions: one cycle of 95°C
176 (3 minutes), 40 cycles of 95°C (10 seconds) and 60°C (10 seconds), one cycle of 95°C (1
177 minute), 55°C (30 seconds) and 95°C (30 seconds, 5 seconds soak time).

178 Missing raw Ct values for age-predictive genes (where the signal never reached the threshold
179 even after 40 cycles) were replaced with the maximum value of 40. Δ Ct values were
180 calculated using the mean Ct of the two housekeeping genes. Where Ct values were missing
181 for either housekeeping gene, normalisation was impossible and the normalised aging gene
182 value was recorded as missing (NA). All samples were run in two technical replicates and the
183 final Δ Ct was taken as the mean of the two replicates. Gene GMOY005321 consistently
184 showed variable Δ Ct values between technical replicates, possibly due to low expression of
185 this gene, and these values were kept unchanged. For all other genes, any gene-sample
186 combinations whose Δ Ct differed by more than 1 between technical replicates were rerun for
187 a third technical replicate, along with both housekeeping genes, providing a third Δ Ct. In most
188 cases, this third Δ Ct was very close to one of the first two and very different from the other,
189 indicating which of the first two technical replicates was wrong. The final Δ Ct was thus taken
190 as the mean of the third replicate and whichever of the first two replicates it was closest to.

191 **2.5 Predicting tsetse age from qPCR data.**

192 Machine learning predictions of tsetse age from qPCR data were performed using the *caret*
193 package in R (<https://cran.r-project.org/package=caret>). The Δ Ct values for each of the 10
194 study genes were used as continuous predictor variables, and sex was included as a
195 categorical predictor variable since some of the genes showed sex-dependent expression.
196 Samples were randomly split into training set (75% of samples) and test set (25% of samples),
197 stratified by sex and age to ensure equal representation of these two variables in the two sets.
198 Due to rounding of sample numbers within each stratification layer, the final numbers in the
199 train and test sets were 380 (76%) and 118 (24%) samples respectively. Model training was
200 performed using three rounds of 10-fold cross-validation. For regression models, whose aim

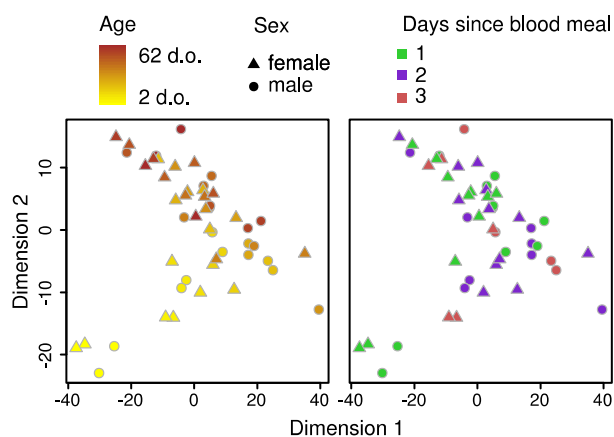


Figure 2: Gene expression clusters primarily by age. Principal component analysis of RNAseq data, coloured by age (left) or days since blood meal (right).

201 is to estimate age as a continuous variable, partial least squares regression (PLS), random
202 forest and extreme gradient boosting (XGB) models were all trained on the data and their
203 predictive accuracies compared. Categorical models were trained to categorise individuals
204 into ≤ 15 and > 15 days old. Simple decision tree, random forest and XGB models were
205 compared for these categorical models.

206 The minimum number of expression markers (genes) required to obtain accurate predictions
207 of age was determined by training the models with different numbers of loci. For each of the
208 random forest and XGB models, the ten genes were ranked according to their variable
209 importance in the full model training described above (sex was found to have a variable
210 importance of 0 in both cases, and was therefore excluded from these models). The models
211 were then trained with all ten genes, the top nine genes, the top eight genes, and so on. For
212 each set of genes, 20 models were trained with a different random split of training and test
213 sets, to account for stochastic variation in model accuracy.

214 All statistical analysis was conducted in R version 3.4.4 (R Core Team 2015). Analysis
215 scripts, qPCR raw data and RNAseq read counts are available on GitHub
216 (<https://github.com/EricRLucas/TsetseAgeMarkers>). Raw sequencing will be submitted to
217 ENA shotgun sequencing archive upon final acceptance of the paper for publication.

218 3 Results

219 We collected 301 female and 204 male *G. morsitans* flies of known age from laboratory
220 colonies, ranging in age from 2 to 62 days old. An initial RNAseq analysis of 28 female and
221 22 male samples showed that gene expression in these samples was primarily affected by age,
222 rather than sex or days since last blood meal (Fig. 2, Supplementary Fig. S1), although this

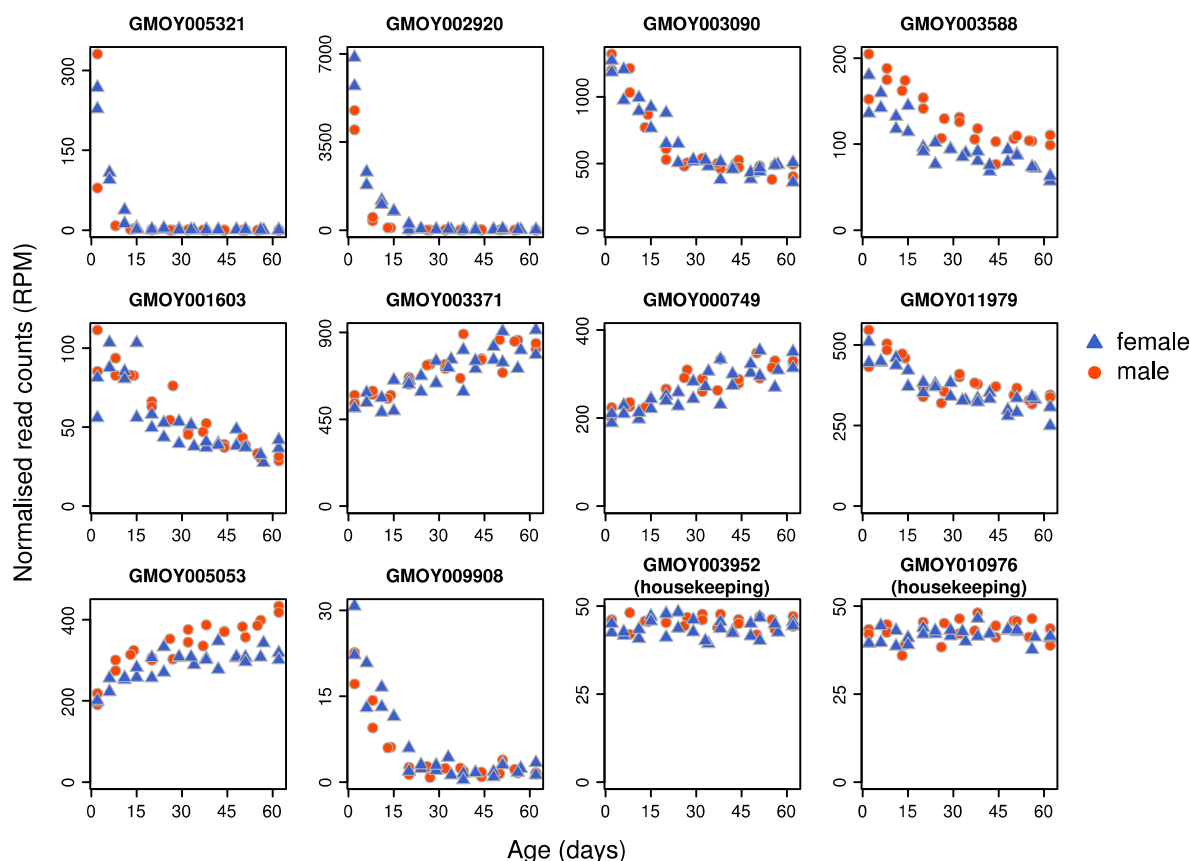


Figure 3: Expression of ten age-related genes and two housekeeping genes from RNAseq data, ordered according to the variable importance in the XGB model (Fig. 4). Very strong early-age expression changes in some genes (eg: *GMOY005321*, *GMOY002920*) allow good discrimination among young individuals, but show little change in later life. Genes with continuous changes (eg: *GMOY003371*, *GMOY000749*) are more gradual and offer more consistent, but less powerful, discrimination at all ages.

223 was primarily due to the strong changes in gene expression found during the first 15 days of
224 life, with older individuals clustering primarily by sex (Supplementary Fig. S2).

225 We identified a set of genes that was likely to provide strong age prediction by looking for
226 genes that: 1. Were strongly correlated with age, or 2. consistently performed well in
227 prediction of age using lasso regression and 3. where possible, belonged to different gene
228 clusters as defined by weighted gene network clustering analysis. We particularly looked for
229 genes showing strong expression changes in older individuals by identifying the genes most
230 differentially expressed when considering only individuals older than 15 days, but even these
231 showed relatively slight changes with age compared to some of the changes seen in the first
232 15 days of life (Fig. 3, Supplementary Fig. S3). Using our criteria, and after testing qPCR
233 primer efficient, we manually picked 10 genes associated with age, and 2 genes with very
234 little variation across samples to serve as housekeeping genes (Figs. 3 and 4).

Gene	Description	Top Drosophila BLAST hit	XGB import.	RF import.	XGB class. import.
GMOY005321	Cuticular protein 49Aa	tr A8DRW0 A8DRW0_DROME Cuticular protein 49Aa	100.0	60.1	3.9
GMOY002920	Cuticular protein 92F	tr Q9VDJ8 Q9VDJ8_DROME Cuticular protein 92F	71.4	72.0	100.0
GMOY003090	Porin	sp Q94920 VDAC_DROME Voltage-dependent anion-selective channel	70.4	68.6	2.6
GMOY003588		tr Q7K188 Q7K188_DROME Protein quiver	67.9	100.0	11.7
GMOY001603	friend of echinoid	tr A0A023GPK8 A0A023GPK8_DROME Friend of echinoid, isoform H	26.5	78.4	2.0
GMOY003371	Elongation factor 1-alpha	sp P05303 EF1A2_DROME Elongation factor 1-alpha 2	18.5	84.3	0.3
GMOY000749		sp P05303 EF1A2_DROME Elongation factor 1-alpha 2	16.7	46.3	0.1
GMOY011979	Vacuolar H ⁺ -ATPase v1 sector subunit E	sp P54611 VATE_DROME V-type proton ATPase subunit E	7.9	38.7	0.1
GMOY005053		tr Q9VG81 Q9VG81_DROME RH49330p	6.8	29.8	6.4
GMOY009908		tr Q9VLZ6 Q9VLZ6_DROME FI24007p1	5.5	38.2	54.8
GMOY003952*	nuclear pore complex component	tr Q7K2X8 Q7K2X8_DROME Nucleoporin at 44A, isoform A	NA	NA	NA
GMOY010976*		sp Q9W123 POF_DROME Protein painting of fourth	NA	NA	NA

Figure 4: Ten age-related genes and two housekeeping genes (denoted with *) were used for qPCR analysis. Gene descriptions are taken from the Contig names in the GmorY1.9 proteome. Top Drosophila BLAST hits obtained by blasting the GmorY1.9 proteome against the *D. melanogaster* swissprot proteome. Variable importance of each gene shown for XGB, random forest (RF) and XGB classifier models trained with all predictor variables.

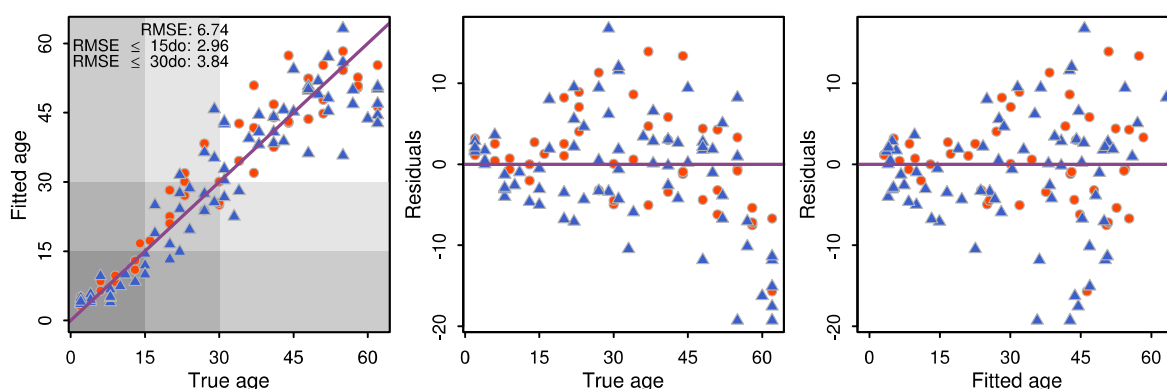


Figure 5: Prediction accuracy of the XGB model was highest (RMSE lowest) for individuals under 15 days old (2.59), and highest when all individuals were considered (6.81). Females are shown as blue triangles and males as orange circles. Purple line shows idealised perfect prediction.

235 We obtained qPCR measurements of expression for these genes from 297 females and 201
236 males (Fig. 1). As expected, expression of all 10 age-related genes was strongly correlated
237 with age (Supplementary Fig. S4) and with the RNAseq data (Supplementary Fig. S5).
238 Principal component analysis of these age-related genes showed that age dominated the first
239 principal component of the data. In particular, samples clustered strongly into those younger
240 and older than 15 days (Supplementary Fig. S6)

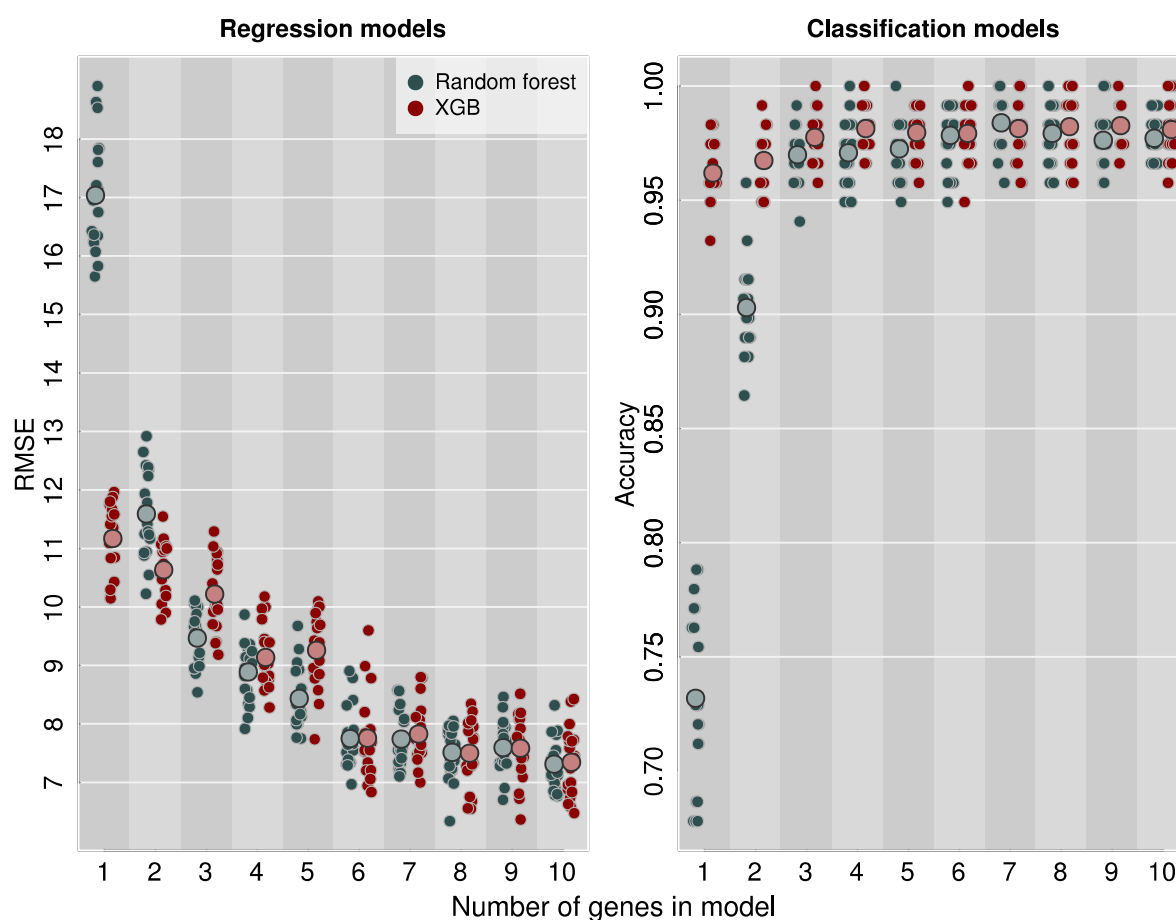
241 The qPCR expression data produced strong overall predictions of age, with predictions being
242 much more accurate in young flies (15 days or younger) compared to older flies. For
243 regression models, PLS provided the poorest predictions of age, while random forest and
244 XGB models performed equally well (Fig. 5, Supplementary Fig. S7). Taking the XGB model
245 as an example, the overall root mean squared error (RMSE) for the final model was 6.74 days,
246 but was 2.96 for individuals ≤ 15 days old. Variable importance for each gene in the random
247 forest and XGB models are shown in Fig. 4. Training the model separately for males and
248 females did not improve prediction accuracy (Supplementary Fig. S8).

249 Models also performed well at classifying samples into age categories of ≤ 15 and > 15 days
250 old (Supplementary Fig. S9). The XGB model performed best in this task, accurately
251 classifying 117 out of 118 samples in the test set.

252 For both the random forest and XGB regression models, prediction accuracy showed little
253 decrease when the variables of least importance were dropped from the models (Fig. 6). In
254 both cases, accuracy remained comparable to that with all 10 genes when only 6 genes were
255 included, with RMSE changing from 7.3 to 7.7 (random forest) or from 7.3 to 7.8 (XGB). In
256 contrast, when moving to 5 genes instead of 6, RMSE changed from 7.7 to 8.4 (random
257 forest) or from 7.8 to 9.3 (XGB). Interestingly, the same 6 genes proved to be sufficient for
258 both model types (GMOY005321, GMOY002920, GMOY003090, GMOY003588,
259 GMOY001603, GMOY003371). For the classification models, even fewer genes were needed
260 (Fig. 6), with just two genes being sufficient for XGB classification accuracy consistently
261 better than 95% (GMOY002920, GMOY009908).

262 **4 Discussion**

263 We have identified a set of gene expression markers that can be used to predict the age of *G.*
264 *morsitans* tsetse flies in the laboratory. Importantly, this method can be applied to both males
265 and females, providing accurate estimates of age in male tsetse. This is particularly important



bennett

Figure 6: Predictive power of XGB and random forest models plateaus after the top 6 genes are included in the models (left). Accuracy of classification models plateaus after top 3 genes are included, with >95% accuracy achievable with only two genes (right). Small points show models run on independent test-train splits of the data (20 replicates per gene number); large points show the mean for each category. Points are jittered on the x axis to show overlapping data.

266 since not only do both male and female tsetse flies transmit trypanosomes, but males appear
267 to be more likely to develop transmissible infections (Dale et al. 1995, Maudlin et al. 1990).
268 Our genetic markers were also unaffected by time since an individual's last blood meal,
269 making them more robust for use on wild-caught individuals, where such factors cannot be
270 controlled. Further work is nevertheless required to test the applicability of these markers in
271 field conditions, since other environmental variables may still affect expression. For example,
272 temperature and humidity were constant in our rearing conditions, and all samples were
273 collected around the same time of day, leaving the possibility that these factors may yet
274 influence the expression of our markers.

275 Like other methods for estimating the age of vectors, prediction accuracy decreases at older
276 ages (Brei et al. 2004, Cook et al. 2006, Cook & Sinkins 2010, Gerade et al. 2004, Liebman et

277 al. 2015, Penilla et al. 2002, Sikulu et al. 2010). In our data, this was because the change in
278 expression with age was much greater in younger compared to older individuals, suggesting
279 that the overall physiology of tsetse changes slowly after a certain life stage, and that there is
280 thus little to detect that can be used for age grading. While we found genes that continued to
281 change in older ages, the rate of change relative to the variance within age groups was not
282 sufficient to achieve the same prediction accuracies as found in younger individuals. While it
283 is likely that more accurate old-age predictions would be achievable using whole-
284 transcriptome methods such as RNAseq, this is too costly to be applied at the scales required
285 for training predictive models. In mosquitoes, spectroscopy-based methods used to estimate
286 age initially suffered from a similar loss of precision at older ages (Liebman et al. 2015,
287 Mayagaya et al. 2009, Sikulu et al. 2010, Sikulu-Lord et al. 2016), but recent studies using
288 machine learning prediction methods have improved prediction accuracies (Lambert et al.
289 2018, Milali et al. 2019). Whether similar performance can be achieved with tsetse should be
290 explored.

291 While we used ten genes in our study, we found that using only the six genes most predictive
292 of age still provided high prediction accuracy, and only two genes were needed for classifying
293 individuals into age groups of ≤ 15 and > 15 days old. By removing four genes from the
294 analysis, qPCR time and costs can be reduced by 1/3 (eight qPCR reactions per sample
295 instead of twelve), while removing eight genes will reduce costs by 2/3. We thus suggest that
296 further studies testing the applicability of these markers in the field restrict themselves to
297 either six or two genes, depending on how precisely age needs to be estimated. Such studies
298 are needed to determine the applicability of these markers in the field, but it would also be
299 interesting to measure the expression of these genes in age-controlled samples of other species
300 of tsetse to determine whether these markers have widespread applicability. Once the field
301 applicability of these markers is confirmed, the technique can be rolled out in the context of
302 monitoring of tsetse control campaigns by comparing the age distribution before and after
303 interventions to confirm that a resulting shift in the population age distribution is observed. In
304 particular, in the wake of a 100% effective campaign, no flies older than the start of the
305 campaign should be found. The resulting data on age structure both before and after control
306 campaigns can then also be used to inform epidemiological models of trypanosomiasis
307 transmission.

308 In conclusion, our study provides a new method for estimating the age of tsetse flies which
309 does not require specialist dissection skills and can be applied to males. The problem remains

310 of finding methods for more accurately estimating age in older individuals. This may involve
311 identifying senescent changes whose rate is steady and consistent enough to be generalisable
312 to any individual in the population.

313 Acknowledgements

314 We are grateful to Rob Leyland for assistance and tuition for the insectary work, and to Tom
315 Churcher and Ben Lambert for valuable discussion on analytical approach. This work was
316 supported by a Liverpool School of Tropical Medicine internal grant (Director's Catalyst
317 Fund) to ERL and a Medical Research Council, UK (MR/T001070/1) grant to MJD and ERL.
318 SJT received support from the UK's Biotechnology and Biological Sciences Research Council
319 (grant numbers: BB/S01375X/1, BB/S00243X/1, BB/P005888/1) and the Bill and Melinda
320 Gates Foundation (INV-001785, OPP1155293).

321 References

- 322 **Andrews S (2010)**. FastQC: A quality control tool for high throughput
323 sequence data, URL <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>
324 **Brei B, Edman JD, Gerade B, Clark JM (2004)**. Relative abundance of two
325 cuticular hydrocarbons indicates whether a mosquito is old enough to
326 transmit malaria parasites, *Journal of medical entomology* 41 : 807-809.
327 **Caragata EP, Poinsignon A, Moreira LA, Johnson PH, Leong YS, Ritchie SA,
328 O'Neill SL, McGraw EA (2011)**. Improved accuracy of the transcriptional
329 profiling method of age grading in *Aedes aegypti* mosquitoes under
330 laboratory and semi-field cage conditions and in the presence of *Wolbachia*
331 infection, *Insect molecular biology* 20 : 215-224.
332 **Cook PE, Hugo LE, Iturbe-Ormaetxe I, Williams CR, Chenoweth SF, Ritchie SA,
333 Ryan PA, Kay BH, Blows MW, O'Neill SL (2006)**. The use of transcriptional
334 profiles to predict adult mosquito age under field conditions, *Proceedings
335 of the National Academy of Sciences of the United States of America* 103 :
336 18060-18065.
337 **Cook PE, Sinkins SP (2010)**. Transcriptional profiling of *Anopheles gambiae*
338 mosquitoes for adult age estimation, *Insect molecular biology* 19 : 745-751.
339 **Dale C, Welburn SC, Maudlin I, Milligan PJM (1995)**. The kinetics of
340 maturation of trypanosome infections in tsetse, *Parasitology* 111 : 187-191.
341 **Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P,
342 Chaisson M, Gingeras TR (2013)**. STAR: ultrafast universal RNA-seq aligner,
343 *Bioinformatics* 29 : 15-21.
344 **Dye C (1992)**. The analysis of parasite transmission by bloodsucking
345 insects, *Annual review of entomology* 37 : 1-19.
346 **Eisler MC, Torr SJ, Coleman PG, Machila N, Morton JF (2003)**. Integrated
347 control of vector-borne diseases of livestock--pyrethroids: panacea or
348 poison?, *Trends in Parasitology* 19 : 341-345.
349 **Franco JR, Simarro PP, Diarra A, Ruiz-Postigo JA, Jannin JG (2014)**. The
350 journey towards elimination of gambiense human African trypanosomiasis: not
351 far, nor easy, *Parasitology* 141 : 748-760.
352 **Friedman J, Hastie T, Tibshirani R (2010)**. Regularization paths for
353 generalized linear models via coordinate descent, *Journal of statistical
354 software* 33 : 1.

355 **Gerade BB, Lee SH, Scott TW, Edman JD, Harrington LC, Kitthawee S, Jones**
356 **JW, Clark JM (2004)**. Field validation of *Aedes aegypti* (Diptera: Culicidae)
357 age estimation by analysis of cuticular hydrocarbons, *Journal of medical*
358 *entomology* 41 : 231-238.
359 **Hargrove JW (1990)**. Age-dependent changes in the probabilities of survival
360 and capture of the tsetse, *Glossina morsitans morsitans* Westwood,
361 *International Journal of Tropical Insect Science* 11 : 323-330.
362 **Hargrove JW (2012)**. Age-specific changes in sperm levels among female
363 tsetse (*Glossina* spp.) with a model for the time course of insemination,
364 *Physiological entomology* 37 : 278-290.
365 **Hargrove JW (2020)**. A model for the relationship between wing fray and
366 chronological and ovarian ages in tsetse (*Glossina* spp), *Medical and*
367 *Veterinary Entomology* 34 : 251-263.
368 **Holmes P (2013)**. Tsetse-transmitted trypanosomes-their biology, disease
369 impact and control, *Journal of invertebrate pathology* 112 : S11-S14.
370 **Kgori P, Modo S, Torr S (2006)**. The use of aerial spraying to eliminate
371 tsetse from the Okavango Delta of Botswana, *Acta Tropica* 99 : 184-199.
372 **Klaus B, Strimmer K (2015)**. fdrtool: Estimation of (local) false discovery
373 rates and higher Criticism, URL <http://CRAN.R-project.org/package=fdrtool>.
374 **Lambert B, Sikulu-Lord MT, Mayagaya VS, Devine G, Dowell F, Churcher TS**
375 **(2018)**. Monitoring the age of mosquito populations using near-infrared
376 spectroscopy, *Scientific reports* 8 : 5274.
377 **Langfelder P, Horvath S (2008)**. WGCNA: an R package for weighted
378 correlation network analysis, *BMC bioinformatics* 9 : 559.
379 **Langley PA, Hall MJR, Felton T, Ceesay M (1988)**. Determining the age of
380 tsetse flies, *Glossina* spp.(Diptera: Glossinidae): an appraisal of the
381 pteridine fluorescence technique, *Bulletin of Entomological Research* 78 :
382 387-395.
383 **Lardeux F, UNG A, Chebret M (2000)**. Spectrofluorometers are not adequate
384 for aging *Aedes* and *Culex* (Diptera: Culicidae) using pteridine
385 fluorescence, *Journal of Medical Entomology* 37 : 769-773.
386 **Lehane MJ, Hargrove J (1988)**. Field experiments on a new method for
387 determining age in tsetse flies (Diptera: Glossinidae), *Ecological*
388 *Entomology* 13 : 319-322.
389 **Liebman K, Swamidoss I, Vizcaino L, Lenhart A, Dowell F, Wirtz R (2015)**.
390 The influence of diet on the use of near-infrared spectroscopy to determine
391 the age of female *Aedes aegypti* mosquitoes, *The American journal of*
392 *tropical medicine and hygiene* 92 : 1070-1075.
393 **Martin M (2011)**. Cutadapt removes adapter sequences from high-throughput
394 sequencing reads, *EMBnet. journal* 17 : 10-12.
395 **Maudlin I, Welburn SC, Milligan P (1990)**. Salivary gland infection: a sex-
396 linked recessive character in tsetse?, *Acta Tropica* 48 : 9-15.
397 **Mayagaya VS, Michel K, Benedict MQ, Killeen GF, Wirtz RA, Ferguson HM,**
398 **Dowell FE (2009)**. Non-destructive determination of age and species of
399 *Anopheles gambiae* sl using near-infrared spectroscopy, *The American journal*
400 *of tropical medicine and hygiene* 81 : 622-630.
401 **Milali MP, Sikulu-Lord MT, Kiware SS, Dowell FE, Corliss GF, Povinelli RJ**
402 **(2019)**. Age grading *An. gambiae* and *An. arabiensis* using near infrared
403 spectra and artificial neural networks, *PLoS one* 14 : e0209451.
404 **Penilla RP, Rodríguez MH, López AD, Viader-Salvadó JM, Sánchez CN (2002)**.
405 Pteridine concentrations differ between insectary-reared and field-
406 collected *Anopheles albimanus* mosquitoes of the same physiological age,
407 *Medical and veterinary entomology* 16 : 225-234.
408 **R Core Team (2015)**. R: A Language and Environment for Statistical
409 Computing, URL <https://www.R-project.org/>.
410 **Robinson MD, McCarthy DJ, Smyth GK (2010)**. edgeR: a Bioconductor package
411 for differential expression analysis of digital gene expression data,
412 *Bioinformatics* 26 : 139-140.
413 **Robinson MD, Oshlack A (2010)**. A scaling normalization method for
414 differential expression analysis of RNA-seq data, *Genome Biology* 11 : R25.

416 **Abd-Allah F, Abdela J, Abdelalim A, et al (2018)**. Global, regional, and
417 national age-sex-specific mortality for 282 causes of death in 195
418 countries and territories, 1980-2017: a systematic analysis for the Global
419 Burden of Disease Study 2017, *The Lancet* 392 : 1736-1788.
420 **Rund SSC, Hou TY, Ward SM, Collins FH, Duffield GE (2011)**. Genome-wide
421 profiling of diel and circadian gene expression in the malaria vector
422 *Anopheles gambiae*, *Proceedings of the National Academy of Sciences* 108 :
423 E421-E430.
424 **Shaw APM (2004)**. *Economics of African Trypanosomiasis*. In: Maudlin, I.;
425 Holmes, P. H. & Miles, M. A. (Ed.), *The Trypanosomiasis*, CABI Publishing.
426 **Sikulu M, Killeen GF, Hugo LE, Ryan PA, Dowell KM, Wirtz RA, Moore SJ,**
427 **Dowell FE (2010)**. Near-infrared spectroscopy as a complementary age grading
428 and species identification tool for African malaria vectors, *Parasites &*
429 *vectors* 3 : 49.
430 **Sikulu-Lord MT, Milali MP, Henry M, Wirtz RA, Hugo LE, Dowell FE, Devine GJ**
431 **(2016)**. Near-Infrared Spectroscopy, a Rapid Method for Predicting the Age
432 of Male and Female Wild-Type and *Wolbachia* Infected *Aedes aegypti*, *PLoS*
433 *Negl Trop Dis* 10 : e0005040.
434 **Vale GA, Lovemore DF, Flint S, Cockbill GF (1988)**. Odour-baited targets to
435 control tsetse flies, *Glossina* spp. (Diptera: Glossinidae), in Zimbabwe,
436 *Bulletin of Entomological Research* 78 : 31-49.
437 **Van den Bossche P (1997)**. The control of *Glossina morsitans morsitans*
438 (Diptera: Glossinidae) in a settled area in Petauke District (Eastern
439 Province, Zambia) using odour-baited targets, *Onderstepoort Journal of*
440 *Veterinary Research* 64 : 251-257.
441 **Wilson AL, Boelaert M, Kleinschmidt I, Pinder M, Scott TW, Tusting LS,**
442 **Lindsay SW (2015)**. Evidence-based vector control? Improving the quality of
443 vector control trials, *Trends in parasitology* 31 : 380-390.
444 **Wilson AL, Courtenay O, Kelly-Hope LA, Scott TW, Takken W, Torr SJ, Lindsay**
445 **SW (2020)**. The importance of vector control for the control and elimination
446 of vector-borne diseases, *PLoS Neglected Tropical Diseases* 14 : e0007831.
447 **Woolhouse MEJ, Hargrove JW (1998)**. On the interpretation of age-prevalence
448 curves for trypanosome infections of tsetse flies, *Parasitology* 116 : 149-
449 156.
450 **World Health Organization (2012)**. Monitoring and evaluation indicators for
451 integrated vector management. Geneva: World Health Organization. .
452 **Ye J, Coulouris G, Zaretskaya I, Cutcutache I, Rozen S, Madden TL (2012)**.
453 Primer-BLAST: a tool to design target-specific primers for polymerase chain
454 reaction, *BMC bioinformatics* 13 : 134.