1

2

3

4

5

# Advancing age grading techniques for *Glossina morsitans morsitans*, vectors of African trypanosomiasis, through mid-infrared spectroscopy and machine learning

9

Mauro Pazmiño-Betancourth[1], Ivan Casas Gómez-Uribarri[1], Karina Mondragon-Shem[2], Simon A Babayan[1], Francesco Baldini[1,3*, $], Lee Rafuse Haines[2,4,*,$]

[1]School of Biodiversity, One Health and Veterinary Medicine, University of Glasgow, UK, G12 8QQ

[2]Department of Vector Biology, Liverpool School of Tropical Medicine, UK, L3 5QA

[3] Ifakara Health Institute, Environmental Health, and Ecological Sciences Department, Morogoro, United Republic of Tanzania

[4] Department of Biological Sciences, University of Notre Dame, USA, 46556

*Francesco.Baldini@glasgow.ac.uk, lhaines@nd.edu

$These authors equally supervised the work

19

**Corresponding authors:**

Francesco.Baldini@glasgow.ac.uk, lhaines@nd.edu

**Short title (70 characters):**

Mid-infrared spectroscopy and machine learning for tsetse surveillance

# Abstract

Tsetse are the insects responsible for transmitting African trypanosomes, which cause sleeping sickness in humans and animal trypanosomiasis in wildlife and livestock. Knowing the age of these flies is important when assessing the effectiveness of vector control programs and modelling disease risk. However, current methods to assess fly age are labour-intensive, slow, and often inaccurate as skilled personnel are in short supply. Mid-infrared spectroscopy (MIRS), a fast and cost-effective tool to accurately estimate several biological traits of insects, offers a promising alternative. This is achieved by characterising the biochemical composition of the insect cuticle using infrared light coupled with machine learning algorithms to estimate the traits of interest.

We tested the performance of MIRS in estimating tsetse sex and age for the first time using spectra obtained from their cuticle. We used 541 insectary-reared *Glossina m. morsitans* of two different age groups for males (5 and 7 weeks) and three age groups for females (3 days, 5 weeks, and 7 weeks). Spectra were collected from the head, thorax, and abdomen of each sample. Machine learning models differentiated between male and female flies with a 96% accuracy and predicted the age group with 94% and 87% accuracy for males and females, respectively. The key infrared regions important for discriminating sex and age classification were characteristic of lipid and protein content. Our results support the use of MIRS as a fast and accurate way to identify tsetse sex and age with minimal pre-processing. Further validation using wild-caught tsetse can pave the way for this technique to be implemented as a routine surveillance tool in vector control programmes.

47    **Author summary (150-200)**

48    Male and female tsetse transmit the parasites that cause sleeping sickness in humans and

49    nagana in livestock. To control these diseases, knowing the age of these flies is important, as

50    it helps evaluate the efficacy of control measures and assess disease risk. However, current

51    age-grading methods are laborious, often unreliable, and in the case of male tsetse, highly

52    inaccurate. This study explores a novel approach that uses mid-infrared spectroscopy (MIRS)

53    to estimate the age of individual tsetse. Machine learning can detect signatures in MIRS that

54    help identify the composition of a fly's cuticle, which differs between sexes and changes as

55    they age.

56    We trained machine learning models that distinguished male from female flies with 96%

57    accuracy and predicted the correct age group with 94% accuracy for males and 87% accuracy

58    for females. MIRS offers a fast and reliable way to identify tsetse sex and age with minimal

59    preparation. If this method is successfully validated with wild flies, it holds the potential to

60    vastly increase the accuracy of the way we monitor and combat these disease-carrying

61    insects, thus offering significant advantages in our efforts to control them.

62

63

64

65

66

## Introduction

Tsetse are blood-feeding flies that can transmit trypanosome parasites of human and animal concern. There are two parasite species that cause Human African Trypanosomiasis (HAT), or sleeping sickness, and infected patients can die if they do not receive treatment. The promising decline of cases in endemic areas[1] in recent years is due to ongoing disease and vector control efforts, but continued support is critical to ensure the success of disease elimination programmes. However, Rhodesiense HAT (the more severe form) is still a concern due to livestock and wildlife forming part of its transmission cycle. Animal African trypanosomiasis (AAT) affects wildlife and domestic animals, causing three million cattle deaths/year with agricultural losses nearing US$ 5 billion/year[2]. Both female and male tsetse can transmit trypanosomes, but only adult flies older than 20 days post-emergence that have ingested blood from a parasite-infected host can be infectious. Tsetse age is therefore crucial for estimating transmission risk and the efficacy of vector control programmes. Accurate age grading in the field is crucial for disease monitoring and evaluation operations. An effective vector control intervention, which does not discriminate against age, overall will reduce the average age of tsetse populations. For example, if in an area of ongoing vector control only young flies are caught, this suggests newly emerged flies in the area, whereas capturing older flies either indicates fly reinvasion from outside the intervention zone or intervention failure.

Tsetse age grading for female flies currently relies on performing a labour-intensive ovarian dissection, which requires the use of a microscope and an experienced dissector. Female tsetse give birth to a larva every 9 days[3] throughout life, and the four ovarioles develop in a specific, predictable sequence; as each egg descends into the uterus, it leaves behind a scar

90    (named 'relic') that can be microscopically identified[4]. No new relics are created after the

91    4th ovarian cycle, thus limiting the confidence of this method in flies older than seven

92    weeks[4]. Furthermore, factors such as nutritional stress[5] , tsetse strain[6] and

93    temperature[7] can affect the length of this 9-day process, and even with adjustments, the

94    method can be imprecise. Ovarian dissections are time consuming and need to be performed

95    while the tsetse is still 'fresh', and tissues maintain their form. After death, flies quickly

96    become dehydrated and age grading is no longer possible by this method. This makes it

97    difficult to process large numbers of flies when monitoring control interventions.

98    The current situation is worse for male tsetse, as there are no dependable methods for age-

99    grading them. Wing fray analysis in either wild male[8] or female flies is unreliable as artifacts

100   can be introduced through trapping protocols. Other approaches like tsetse eye pigment

101   (pteridine) analysis[9] and gene expression [10] are too complex or costly for routine use in

102   field settings. Thus, all current age-grading methods are either too imprecise, laborious, or

103   expensive.

104   Mid-infrared spectroscopy (MIRS) has proven to be a versatile technique for determining

105   mosquito age and species in both insectary-reared and field-collected mosquitoes[11–14].

106   MIRS quantifies the energy a molecule absorbs based on its molecular vibrations[15,16]. As

107   the insect surface is covered with a complex mixture of cuticular proteins, polysaccharides,

108   wax and other lipids, this tool provides a way to detect the differences between different

109   samples. The chemical composition of male and female cuticles, as well as different species-

110   specific signatures, can be resolved alongside more transient aspects such as cuticular

111   changes over time[17]. Scanning a dried insect sample with MIRS is fast (1-2 minutes), and

112   when combined with the use of machine learning (ML) algorithms, it provides a powerful

113   toolbox for researchers to rapidly assess vector populations with minimum sample processing

114   and high accuracy.

115   In this study, we use ML to estimate the age and sex from MIRS of different fly tissues

116   collected from insectary-reared tsetse (*Glossina morsitans morsitans)* of known age and sex.

117   We also identified the regions of the tsetse mid-infrared spectrum associated with age and

118   sex, to elucidate the biological basis of our model predictions.

119

# Methods

## Tsetse rearing

122   An age-stratified colony of *Glossina morsitans morsitans* Westwood, established in 2004 at

123   the Liverpool School of Tropical Medicine (LSTM), UK, was daily maintained under the

124   following conditions: 26 – 28 °C, 68 – 78 % humidity and a 12 h/12 h light/dark cycle. Tsetse

125   were fed three times a week on sterile defibrinated horse blood (TCS Biosciences Ltd,

126   Buckingham, UK) using a silicon membrane feeding system.

## Tsetse sampling strategy and desiccation

128   Young, unmated female flies were first collected from emerging pupal pots as male

129   emergence is delayed, and male collection was timed after the females had emerged. Both

130   teneral (unfed, newly emerged) female and male collections were  isolated from each other

131   to prevent potential cuticular contamination with contact sex pheromones (cuticular

132   hydrocarbon) during mating[18].

133   We collected 354 female and 187 male teneral tsetse from the LSTM colony in total for

134   analysis. At specific ages, tsetse were killed with chloroform-soaked cotton, placed on a thin

135   layer of cotton wool inside a 15 ml falcon tube half-filled with silica gel beads, sealed and then

136   stored at 4°C until required.  Desiccated tsetse were transferred to 96-well plates in

137   preparation for shipping to the University of Glasgow. Upon analysis, dried flies were

138   dissected into three sections: head, thorax and abdomen using dissection tweezers.

## Infrared Spectroscopy

140   Spectra from individual heads, thoraces and abdomens were taken by Attenuated Total

141   Reflection (ATR) FT-IR spectroscopy using a Bruker ALPHA II spectrometer equipped with a

142   Globar lamp, a deuterated L-alanine doped triglycene sulphate (DLaTGS) detector, a

143   Potassium Bromide (KBr) beam splitter, and a diamond ATR accessory (Bruker Platinum ATR

144   Unit A225). Twenty-four scans were collected at room temperature between 4000 and 400

145   $cm^{-1}$ with 4 $cm^{-1}$ resolution per sample. When measuring the tsetse samples, we made efforts

146   to avoid practices that introduce sources of bias such as: not always measuring first young

147   and then old samples, or first females and then males. Low-quality spectra were discarded

148   using a custom script designed for mosquito spectra [11,19].

## Machine learning analysis

150   Spectra were centred around a mean of 0 and scaled to a standard deviation of 1 prior to any

151   analysis. Uniform Manifold Approximation and Projection (UMAP) was applied for clustering

152   analysis. Sex and age groups were binarized using one hot encoding[20]. First, we shuffled

153   and split the dataset into the training (80%) and test sets (20%), stratified by sex and age

154   groups (Supplementary Material Table S1). The training set was used to compute baseline

155   performance of four machine learning algorithms: Logistic Regression (LR), Random Forest

156  (RF), Support vector machine (SVC) and Classification and Regression Tree (CART) using 10-

157  fold cross validation and default parameter settings on the training set. Additionally, a

158  permutation score test was performed to evaluate if there was a dependency between the

159  features (absorbance of each wavenumber) and classes (sex and age groups) (Supplementary

160  Material Fig S1). The best model was then optimized using hyperparameter tuning, which

161  consists in choosing a set of optimal values for the model hyperparameters to maximize its

162  performance. The remaining 20% of the data (the test set) was used for the final evaluation

163  of the optimized models. The individual metrics used to evaluate the models were accuracy,

164  sensitivity, and specificity. Machine learning was performed using Python 3.10 and scikit-learn

165  1.2.2.

## Data and code availability

167  The infrared spectral data generated for this study have been deposited in the Enlighten

168  database and are available at http://dx.doi.org/10.5525/gla.researchdata.1564.

169  All code to reproduce the machine learning analysis and figures is available at

170  https://github.com/maurocolapso/Pazmino_TsetseMIRS_2023.git

## Results

## Optimization of tsetse desiccation

173  To understand how long it took for tsetse in different nutritional states to fully dehydrate

174  (which is key to avoid the noise in the spectra caused by the water signal), we placed

175  individual tsetse into 15ml tubes containing a deep layer of silica gel under a thin cap of cotton

176  wool. Fly weight loss was daily recorded until it stabilised. Unfed flies rapidly desiccated within

177  24h, while fully engorged, bloodfed male and female flies took over three days to dehydrate
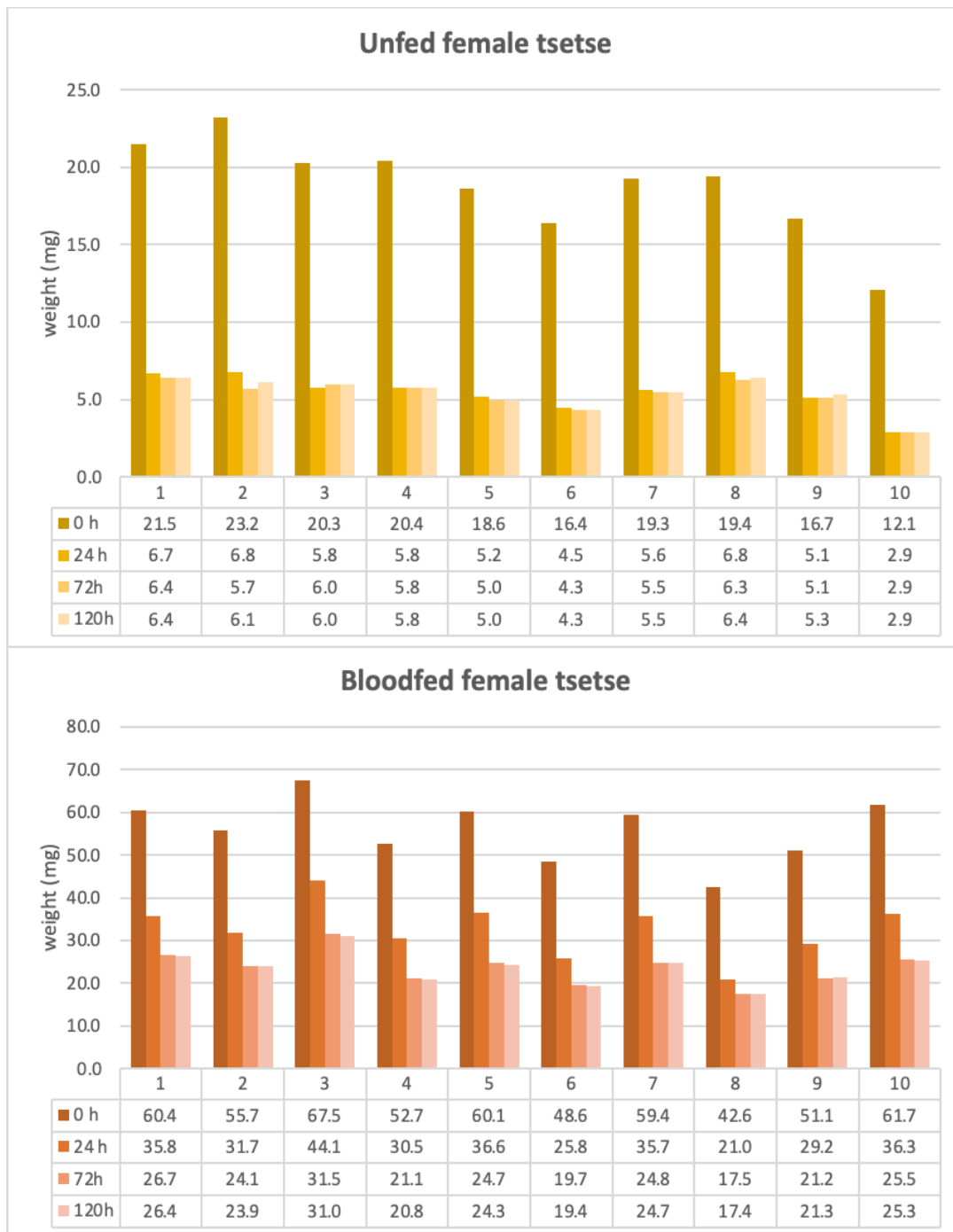
178     the water-rich meal. Based on this data we adopted a standardized ~72h of desiccation on

179     silica for all flies subjected to MIRS analysis (**Error! Reference source not found.**)

180

181 **Table 1.** Desiccation time test for unfed and bloodfed female and male tsetse



**Unfed female tsetse**

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| ■ 0 h | 21.5 | 23.2 | 20.3 | 20.4 | 18.6 | 16.4 | 19.3 | 19.4 | 16.7 | 12.1 |
| ■ 24 h | 6.7 | 6.8 | 5.8 | 5.8 | 5.2 | 4.5 | 5.6 | 6.8 | 5.1 | 2.9 |
| ■ 72h | 6.4 | 5.7 | 6.0 | 5.8 | 5.0 | 4.3 | 5.5 | 6.3 | 5.1 | 2.9 |
| 120h | 6.4 | 6.1 | 6.0 | 5.8 | 5.0 | 4.3 | 5.5 | 6.4 | 5.3 | 2.9 |

**Bloodfed female tsetse**

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| ■ 0 h | 60.4 | 55.7 | 67.5 | 52.7 | 60.1 | 48.6 | 59.4 | 42.6 | 51.1 | 61.7 |
| ■ 24 h | 35.8 | 31.7 | 44.1 | 30.5 | 36.6 | 25.8 | 35.7 | 21.0 | 29.2 | 36.3 |
| ■ 72h | 26.7 | 24.1 | 31.5 | 21.1 | 24.7 | 19.7 | 24.8 | 17.5 | 21.2 | 25.5 |
| 120h | 26.4 | 23.9 | 31.0 | 20.8 | 24.3 | 19.4 | 24.7 | 17.4 | 21.3 | 25.3 |

182

183

## Differences between tsetse tissues

185 Initial tests focused on finding the best body regions or tissues to give a high signal clarity

186 when doing spectrometric readings, as the large size tsetse presented novel logistical

187 challenges. Because wild-caught flies are likely to acquire foreign hydrocarbons from mating,

188    blood feeding, or the resting environment, we sampled zones of the cuticle expected to show

189    the least contamination (**Error! Reference source not found.**).
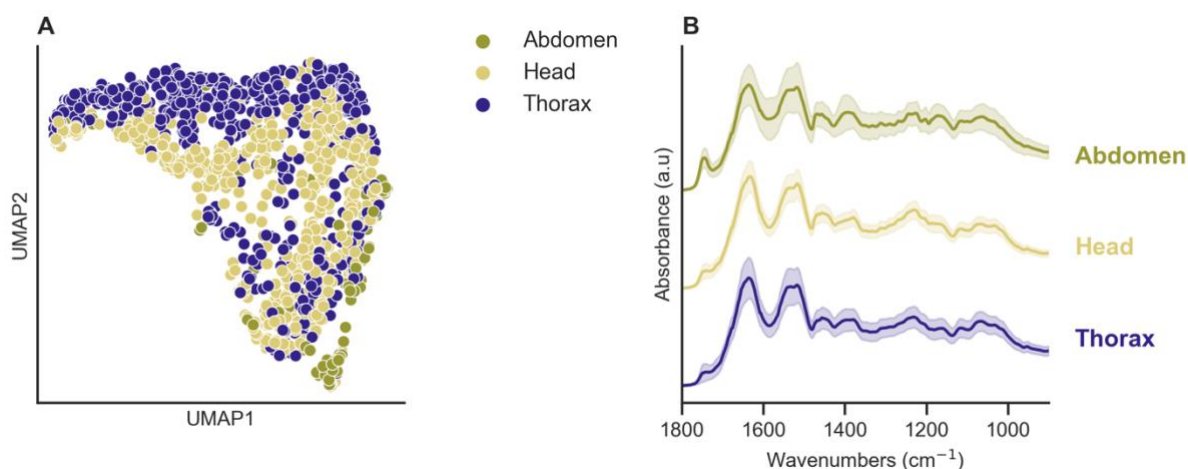
190



191

192    **Fig 1.** Tsetse biology and ecology suggest the heads and dorsal side of male tsetse or the
193    lateral side of the thorax in both sexes would be the best areas (blue circles) to detect
194    individual cuticular hydrocarbons.

195

196    We further investigated the variation between spectra of different tissues. Spectra from fly

197    abdomens differed substantially those from heads and thoraces (Fig 2A), showing lower

198    intensity and a higher variability, especially in the 1800 to 900 cm$^{-1}$ region (Fig 2B). Moreover,

199    visual inspection of the abdomens indicated that despite ~60 days in a sealed anhydrous

200    environment, complete desiccation was not achieved, particularly if the fly had ingested a

201    large blood volume prior to collection. This residual horse blood and water could be driving

202    the greater variability of the abdominal spectra compared to the other tissues. On top of that,

203    previous work in other insects showed the thorax as a target tissue for MIRS. Consequently,

204    we decided to focus our analysis on the spectra obtained from heads and thoraces only. A

205    total of 1071 spectra were obtained by scanning the heads and lateral part of the thoraces of

206    541 flies of different ages (Fig 1, Table *2*).

207

208



**Fig 2**. **Spectra comparison from the abdomen, head, and thorax. A)** Uniform Manifold Approximation and Projection (UMAP) of the abdomens, heads, and thoraces showed that the spectra collected from abdomens formed a separate cluster. Abdomens (olive green), head (yellow), thorax (purple) **B)** High variability of the spectra from abdomens (olive green line) showed that the sources of those inconsistencies were of low intensity and great variability at some wavelengths (primarily in the 1800 to 900 $cm^{-1}$ region) compared to spectra from head (yellow line) and thorax (purple line). Spectra shown in panel B have been manually shifted across the Y-axis for ease of comparison.

218

219

<div align="center">Table 2. Summary of aggregated samples sizes.</div>

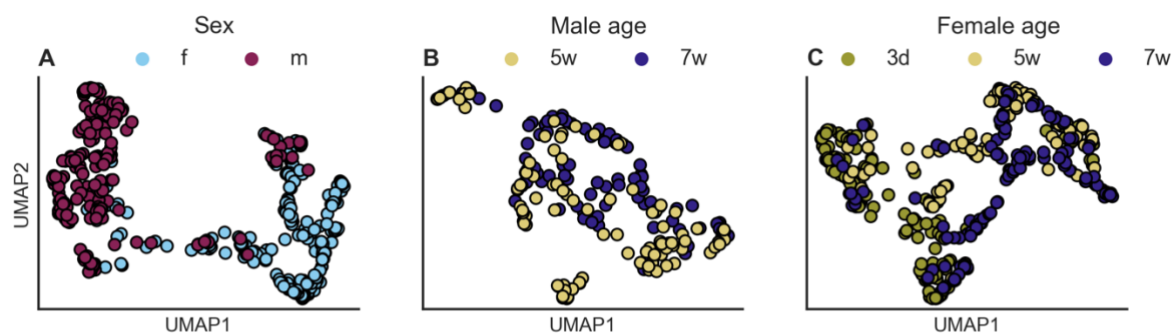| Sex | Age | Tissue | # spectra | # samples |
|---|---|---|---|---|
| **Female** | 3 days | Head | 133 | 136 |
| | | Thorax | 136 | |
| | 5 weeks | Head | 92 | 96 |
| | | Thorax | 96 | |
| | 7 weeks | Head | 120 | 122 |
| | | Thorax | 122 | |
| Male | 5 weeks | Head | 94 | 94 |
| | | Thorax | 93 | |
| | 7 weeks | Head | 93 | 93 |
| | | Thorax | 92 | |
| **Total number of samples** | | | **1071** | **541** |

## Differences between sexes and age groups

We used the unsupervised machine learning algorithm Uniform Manifold Approximation and

Projection (UMAP) to investigate whether the spectra from fly heads (

Fig **3** A-C) and thoraces (

Fig **3** D-F) differed between flies of different sex and age. Most of the male flies produced

different spectra than females, with the thorax showing clearer clusters with fewer samples

overlapping between them (**Fig 3** A and D). For age groups, there were not clusters in males

regardless the tissue (**Fig 3** B and E). In females, there were a distinct cluster composed of old

flies (5 and 7 weeks) when using the thorax, however, there was a high overlap between

samples from different age groups (**Fig 3** C and F). These results show that MIRS contains

biochemical information associated with sex and age as expected from relative changes in the
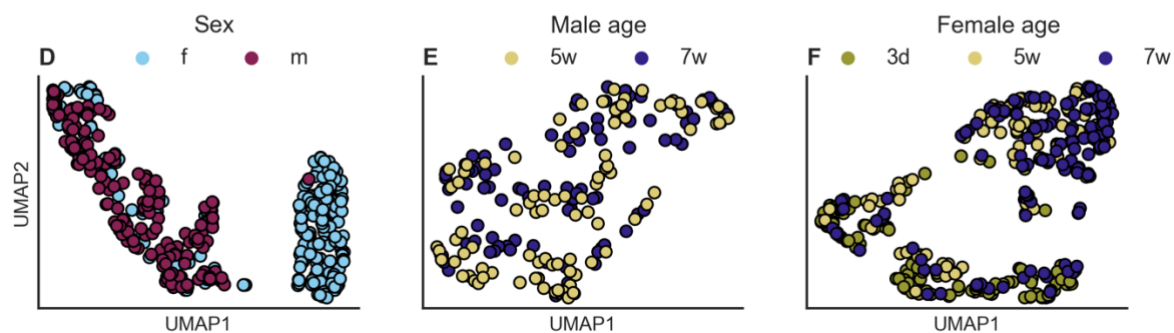
cuticular composition of tsetse.

**Fig 3. MIRS spectra according to tsetse sex and age from specific tissues.** Unsupervised clustering of MIRS measurements using Uniform Manifold Approximation and Projection of MIRS in two-dimensional space using the heads and thorax. Samples are coloured by: **A, D)** sex (females: blue, males: purple). **B, E)** Males coloured by age (5 weeks: yellow, 7 weeks: dark blue). **C, F)** Females coloured by age (3 days: olive green, 5 weeks: yellow, 7 weeks: dark blue)
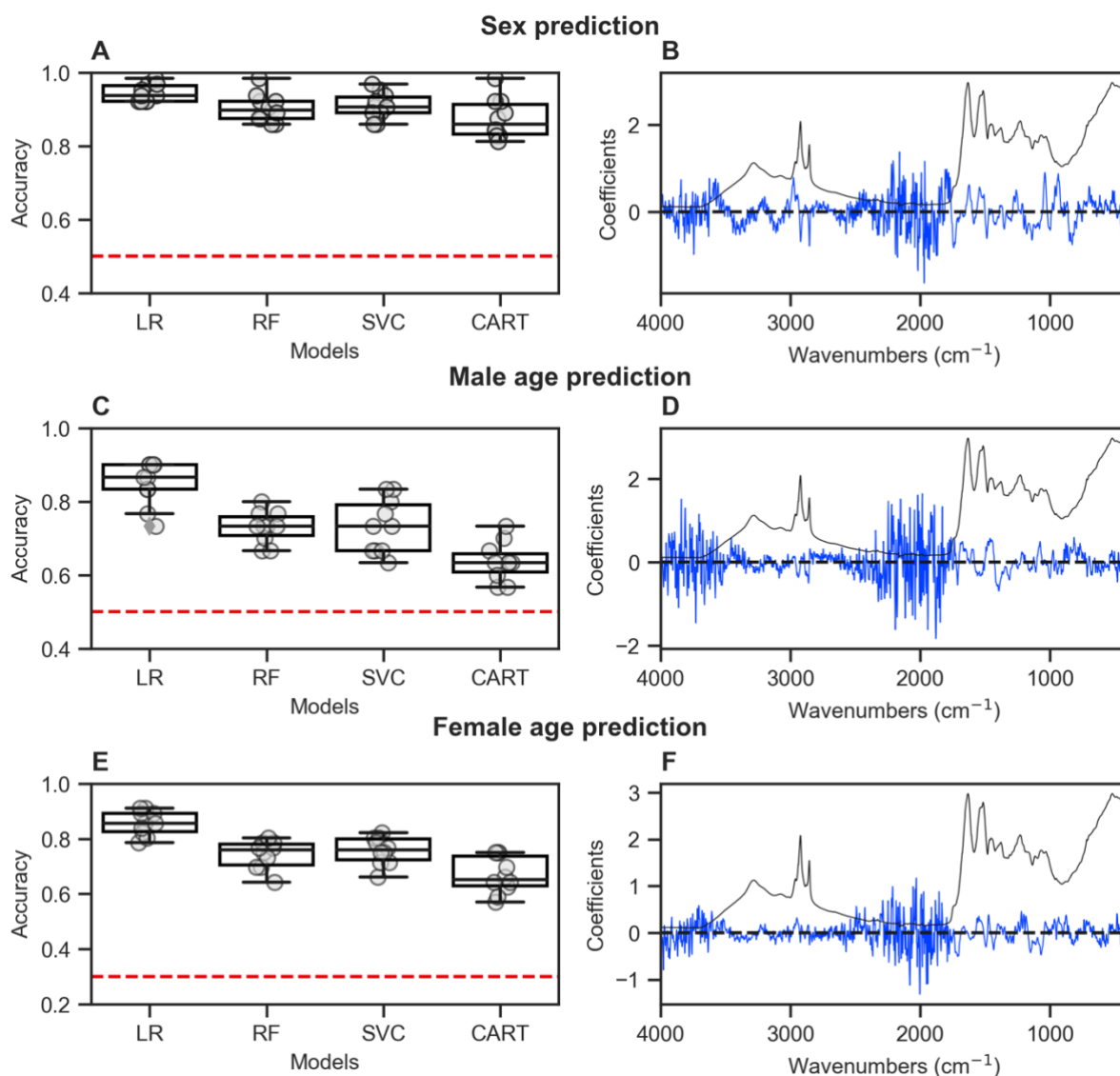
## Sex and age prediction using the complete spectral data

To identify tsetse sex and age-specific patterns within our MIRS dataset, we compared logistic

regression (LR), Random Forest (RF), support vector machine (SVC) and the Classification and

Regression Tree (CART) algorithms. Among these, LR had the highest accuracy (Fig 4A) when

estimating the sex of five- and seven-week-old flies. Training accuracy was 94% when using

both head (Fig 4A) and thorax (Supplementary Material Table S2). Similar accuracies were

obtained in the test set (head = 99%, thorax = 94%, Supplementary Material Table S2). Logistic

Regression was also the most accurate algorithm for identifying age groups among flies of the

250  same sex (Fig 4C). For males, the thorax was marginally better at age prediction with an

251  accuracy of 88% compared to 85% for the head (Supplementary Material Table S2). Similar

252  performance was found in the test set with 92% and 89% for thorax and head, respectively

253  (Supplementary Material Table S2). In females, even though there was some difference in

254  accuracy between the head and thorax on the training set (head = 86%, thorax = 92%),

255  accuracy on the test set was similar for both tissues at 93% (Supplementary Material Table

256  S2). While these initial results suggest that infrared spectra could be used to predict key

257  biological traits of tsetse flies, further analysis of model coefficients suggested that the

258  predictions were being based mostly on flat regions of the spectra, between $4000 - 3750$ $cm^{-1}$

259  $^{1}$ and $2250 - 1800$ $cm^{-1}$ (Fig 4 B, D and E), which are unlikely to contain biochemical

260  information associated with insect cuticle[11] and are primarily used to monitor the presence

261  of $CO_2$ in the environment[15]. This phenomenon was observed with all predictive algorithms

262  regardless of what tissue was used. To further investigate this, we applied the framework by

263  Eid et al. [21]. Briefly, we divided the spectrum into three parts: two regions known to contain

264  vibrations from key chemical bonds ($3500 - 2500$ $cm^{-1}$ and $1800 - 600$ $cm^{-1}$) and one region

265  where no chemical information associated with insect cuticle is expected ($2500 - 1800$ $cm^{-1}$).

266  We then compared the accuracy of four algorithms: Logistic regression, SVM with two kernels

267  (RBF and linear) and Random Forest on each region. While the biochemical fingerprint regions

268  ($3500 - 2500$ $cm^{-1}$, $1800 - 600$ $cm^{-1}$) gave variable prediction accuracies ($60 - 96\%$), when

269  using the region with no chemical information associated with insect cuticle ($2500 - 1800$ $cm^{-1}$

270  $^{1}$), two algorithms (Logistic Regression and SVM with a linear kernel) could still predict

271  different traits with high accuracy ($83 - 94\%$), indicating possible overfitting (Supplementary

272  Table S3). To produce more generalisable models, we therefore chose to base our predictions

273    on the spectral region of 1800 − 600 cm$^{-1}$, which is known to contain the most relevant
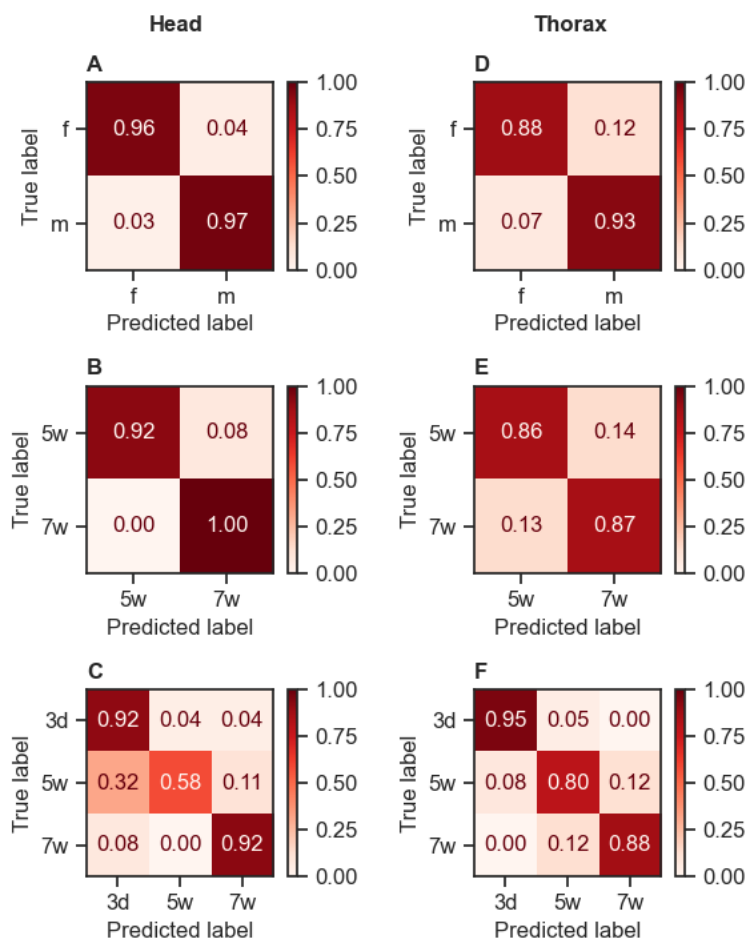
274    biochemical information in insects[12–14].



**Fig 4. Prediction of tsetse sex and age using MIRS.** Model performance on the training set of various ML models (LR: Logistic regression, RF: random forest, SVC: support vector machine and CART: decision tree classifier) for sex and age prediction using the heads of tsetse (**A, C, E**). Boxplots show the distribution of accuracies using 10-fold cross-validation. The horizontal dashed red line indicates a 0.5 accuracy for binary predictions (**A, C**) and 0.3 for a three-class prediction (**E**). Coefficients of the best model (blue line) plotted against the mean spectra of tsetse (**B, D, F**) show how the model relies on the 4000 − 3500 and 2500 −1800 cm$^{-1}$ regions for prediction, which are lacking key biological information.
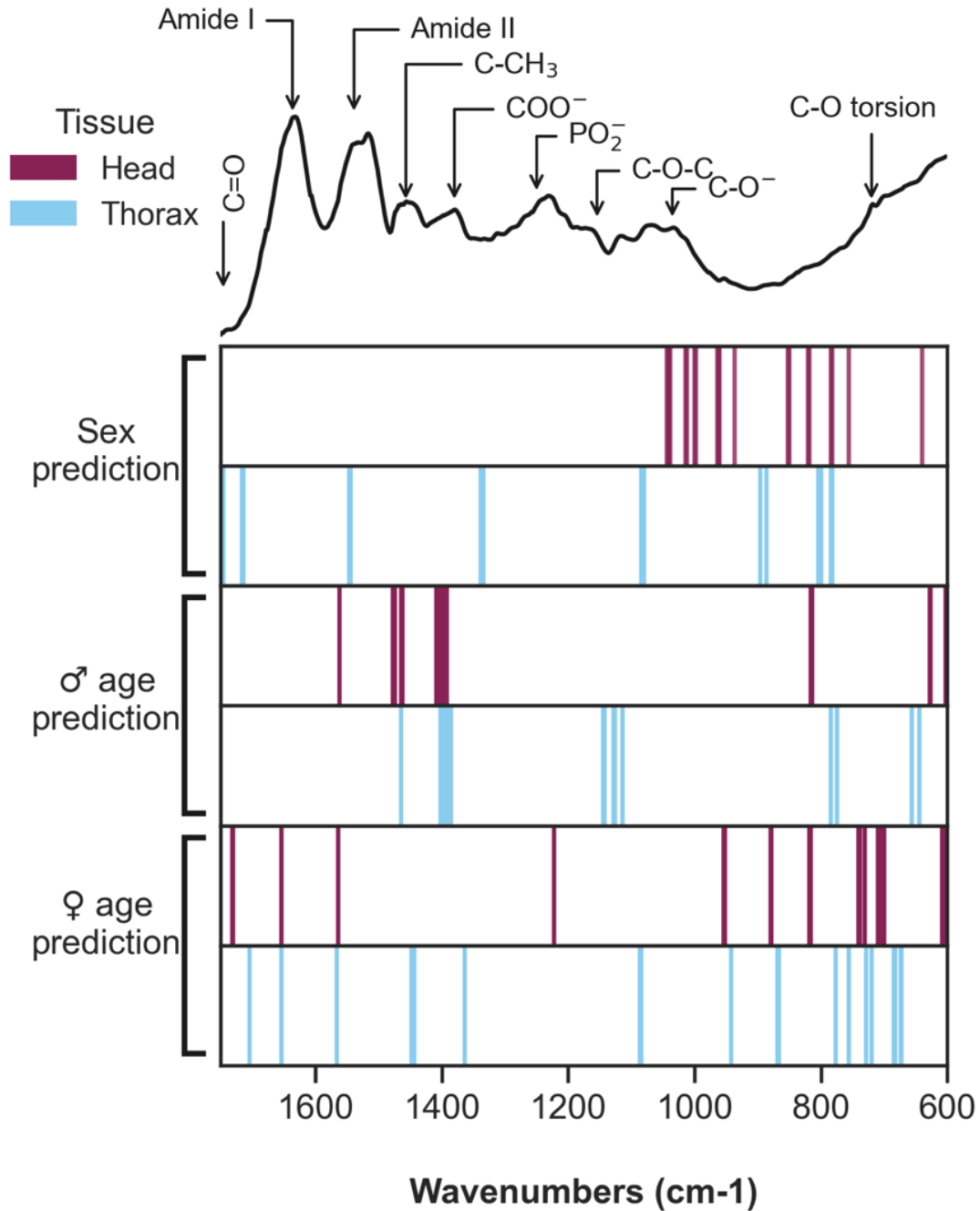
## Sex and age prediction using the biochemical fingerprint region of the spectra

When considering only the spectral region from 1750 - 600 cm$^{-1}$, the accuracy of predicting fly sex and age marginally declined regardless of the algorithm used for analysis (Supplementary Material Fig S2). Logistic regression was able to clearly predict sex using the spectra from the heads with an accuracy of 96% (Fig 5A). The most informative wavenumbers appeared around the 1000 and 800 cm$^{-1}$ areas of the spectra (Fig 6). Logistic regression was able to predict 5-week vs. 7-week-old males using spectral data from the head with an accuracy of 95% (Fig 5B). Most of the coefficients used by this model were in the 1636 and 1400 cm$^{-1}$ region (Fig 6**Fig 6**). Finally, age prediction in females was better when using the thorax. Young teneral flies were identified by the model with over 90% accuracy and older flies with 80% accuracy (Fig 5F). Like males, the important wavenumbers were in the same range, the 1750 to 1450 cm$^{-1}$ 800 to 600 cm$^{-1}$region (Fig 6). However, when using spectral data from the head, the model struggled to identify the 5-week-old age group; an accuracy of 58% (Fig 5C) was obtained, which was likely influenced by several 5-week-old samples being misclassified as 3-day-olds. A summary of the performance of Logistic Regression is shown in Table 3 and wavenumber importance and their assignments are presented in Supplementary Table S4. These results suggest that MIRS-ML is a promising approach when using the tsetse head or thorax to reliably produce quality spectra for sex and age prediction of laboratory-reared flies.

**Fig 5. Confusion matrix for predicting tsetse sex and age using reduced number of wavenumbers.** Accurate identification of females (f) and males (m) (**A, D**) and two-week age difference (5 weeks (5w) vs 7 weeks (7w) old) in male flies (**B, E**). Spectra from the thoraces of young female flies (3d post emergence) compared to older female flies (5 weeks (5w) and 7 weeks (7w) old (**C, F**)

**Fig 6. Important wavenumbers for predicting tsetse sex and age change depending on the trait predicted.** Coloured lines represent the position of the most informative wavenumbers used by the models to predict sex, male age, and female age. Lines are coloured depending on the tissue used for MIRS: head (purple), thorax (light blue). Example spectra with band assignments is added on the top for reference.

318  **Table 3.** Accuracy, sensitivity and specificity of tsetse sex and age prediction on males and
319  females in the training set and test set

|  | Tissue | Accuracy (train set) | Accuracy (test set) |
|---|---|---|---|
| **Sex prediction** | Head | 0.95 ± 0.04 | 0.96 |
|  | Thorax | 0.93 ± 0.03 | 0.90 |
| **Males age prediction** | Head | 0.85 ± 0.06 | 0.94 |
|  | Thorax | 0.82 ± 0.06 | 0.86 |
| **Females age prediction** | Head | 0.84 ± 0.05 | 0.83 |
|  | Thorax | 0.85 ± 0.04 | 0.87 |

320

# Discussion

322  Here, we showed for the first time how a MIRS + ML toolbox can be applied to predict the sex

323  and age of desiccated insectary-reared tsetse. The spectra collected from the head and

324  thorax, but not the abdomen, allow accurate sex prediction. Age grading was successful in

325  both sexes, even when flies were only two weeks apart in age.  When using exclusively the

326  thorax, this toolbox can easily differentiate between females and males using the infrared

327  region related to lipids and carbohydrates. Interestingly, predictions using the head identified

328  a narrow spectral region related to lipids as the most informative. It has been previously

329  reported that *G. pallidipes* females possess a higher amount of cuticular lipids than males[22],

330  which is likely linked to the female sex pheromone that constitutes the main cuticular

331  hydrocarbon. Considering this potential bias, we did not mix the two sexes for analysis since

332  the signal difference can mask the differences between the ages of each sex.

333  When analysing the most important regions for age grading in both males and females, some

334  clear patterns emerged depending on the tissue and biological trait. In male flies, the $C-CH_3$

335 and $COO^-$ bands were consistently important in age grading for all tissues. However, the bands

336 related to proteins and lipids and the $-(CH_2)$-rock functional group related to wax was

337 important across female tissues. Characterizing the informative and predominant

338 wavenumbers is an important for understanding the association between age and absorption

339 bands, which can be used to optimize data collection or model generalisation. An early

340 staining method showed a relationship between cuticular layers in the thorax from laboratory

341 and field caught flies[23]. Other methods using gene expression panels have also found that

342 genes related to cuticular proteins were important for age grading. One study used RNAseq

343 to analyse gene expression associated with age and sex in *G. m. morsitans* that were sourced

344 from the same colony at LSTM [7]. Out of the ten genes shortlisted in the study, two proved

345 to be enough for accurate age classification, one of these being cuticular protein 92F

346 (GMOY002920). A second cuticular protein, 49Aa (GMOY005321), was also part of the list

347 [10]. Previous work using MIRS with other insect vectors also reported differences in female

348 cuticles between very young and old individuals, and the model predicted 3-day old females

349 with minimal misclassification. However, when differentiating between 5- and 7-week-olds,

350 the misclassification between both classes increased.

351 When we used the complete spectra for training, we found that LR and SVM with a linear

352 kernel used the region from 2500 – 1800 $cm^{-1}$ to predict sex and age, which does not contain

353 any biochemical information related to insect cuticle. To ensure the algorithms learn from the

354 biochemical differences between sexes and age groups, we restricted the inputs to specific

355 spectral regions and limit the features the model uses. The strength of machine learning lies

356 in finding patterns to separate classes; however, patterns can arise from confounding effects

357 of contamination by water and $CO_2$ rather than from the structural constituents of the

358 specimen. It is important to diagnose and assess what the model is learning to rule out any

359    bias and avoid overfitting. In spectroscopy data, variation between samples (i.e., baseline

360    offset, variation on $CO_2$ levels during different days when measuring) was robust enough for

361    the model to accurately classify age and sex.

362    When determining the feasibility of using different tsetse tissues for analysis, the abdomen

363    showed inconsistent spectra compared to the head and thorax, which might be caused by the

364    presence of blood from previous meals and incomplete desiccation. However, the

365    information from tsetse abdomens could still be used to identify blood meal sources, as

366    demonstrated by the application of MIRS with *Anopheles* mosquitoes [24]

367    In summary, our results provide proof-of-principle for how MIRS can detect cuticular signals

368    linked to ageing in tsetse. Future validation of this technique using field samples is needed,

369    where environmental cues (naturally minimised in housed insect colonies) impact ageing

370    rates. The next step will be to test the MIRS toolbox against wild tsetse collected from

371    endemic areas, and preferably a region currently implementing vector control strategies. The

372    machine learning models we describe here need to be further refined using more insectary-

373    reared flies alongside a small complementary set of field samples (age-graded when trapped)

374    to be able to confirm the efficacy and accuracy of this technology in the field [12].

## Conclusions

376    Our data strongly support the use of MIRS for high-accuracy age grading of both male and

377    female *Glossina spp.* reared under insectary conditions. The protocol's robustness, minimal

378    maintenance, cost-effectiveness, and speed make it an ideal technique for vector surveillance

379    programmes in resource-limited settings, and implementation will strengthen ongoing

380    control efforts to control transmission of African trypanosomiasis.

381

## Acknowledgements

383     We are grateful to Jonathan Thornton, the dedicated technician overseeing the tsetse

384     colony at LSTM, for his invaluable assistance in procuring tsetse flies for this work.

385

## Author contributions

387     Conceptualization:  F.B., L.R.H.

388     Data curation: M.P., K.M.S.

389     Formal analysis: M.P

390     Funding acquisition: F.B, L.R.H.

391     Investigation: M.P, I.C, K.M.S.

392     Methodology: F.B., L.R.H.

393     Project administration: F.B., L.R.H.

394     Resources: F.B., L.R.H.

395     Supervision: F.B., L.R.H

396     Visualization: M.P, K.M.S.

397     Writing – original draft: M.P.

398     Writing – review & editing: M.P, K.M.S., I.C, S.A.B., F.B, L.R.H.

399

# Funding

# Competing interests

408  The authors declare that they have not competing interests.

# References

410  1.  Franco JR, Cecchi G, Paone M, Diarra A, Grout L, Kadima Ebeja A, et al. The elimination of human
411      African trypanosomiasis: Achievements in relation to WHO road map targets for 2020. PLoS
412      Negl Trop Dis. 2022;16: e0010047. doi:10.1371/journal.pntd.0010047

413  2.  The disease | Programme Against African Trypanosomosis (PAAT) | Food and Agriculture
414      Organization of the United Nations. [cited 20 Oct 2023]. Available:
415      https://www.fao.org/paat/the-programme/the-disease/en/

416  3.  Tobe SS, Langley PA. Reproductive physiology of glossina. Annual Review of Entomology.
417      1978;23: 283–307. doi:10.1146/annurev.en.23.010178.001435

418  4.  Hargrove JW. A model for the relationship between wing fray and chronological and ovarian
419      ages in tsetse (Glossina spp). Medical and Veterinary Entomology. 2020;34: 251–263.
420      doi:10.1111/mve.12439

421  5.  English S, Barreaux AMG, Leyland R, Lord JS, Hargrove JW, Vale GA, et al. Investigating the
422      unaccounted ones: insights on age-dependent reproductive loss in a viviparous fly. Frontiers in
423      Ecology and Evolution. 2023;11. Available:
424      https://www.frontiersin.org/articles/10.3389/fevo.2023.1057474

425  6.  Pagabeleguem S, Ravel S, Dicko AH, Vreysen MJB, Parker A, Takac P, et al. Influence of
426      temperature and relative humidity on survival and fecundity of three tsetse strains. Parasites &
427      Vectors. 2016;9: 520. doi:10.1186/s13071-016-1805-x

428  7.  Jackson CHN. The biology of tsetse flies. Biol Rev Camb Philos Soc. 1949;24: 174–199.
429      doi:10.1111/j.1469-185x.1949.tb00574.x

8.  Jackson CHN. An artificially isolated generation of tsetse flies (diptera). Bulletin of Entomological Research. 1946;37: 291–299. doi:10.1017/S0007485300022203

9.  Lehane MJ, Mail TS. Determining the age of adult male and female Glossina morsitans morsitans using a new technique. Ecological Entomology. 1985;10: 219–224. doi:10.1111/j.1365-2311.1985.tb00551.x

10. Lucas ER, Darby AC, Torr SJ, Donnelly MJ. A gene expression panel for estimating age in males and females of the sleeping sickness vector Glossina morsitans. PLOS Neglected Tropical Diseases. 2021;15: 1–15. doi:10.1371/journal.pntd.0009797

11. González Jiménez M, Babayan SA, Khazaeli P, Doyle M, Walton F, Reedy E, et al. Prediction of mosquito species and population age structure using mid-infrared spectroscopy and supervised machine learning [version 3; peer review: 2 approved]. Wellcome Open Research. 2019. doi:10.12688/wellcomeopenres.15201.3

12. Siria DJ, Sanou R, Mitton J, Mwanga EP, Niang A, Sare I, et al. Rapid age-grading and species identification of natural mosquitoes for malaria surveillance. Nature Communications. 2022;13: 1–9. doi:10.1038/s41467-022-28980-8

13. Sroute L, Byrd BD, Huffman SW. Classification of mosquitoes with infrared spectroscopy and partial least squares-discriminant analysis. Appl Spectrosc. 2020;74: 900–912. doi:10.1177/0003702820915729

14. Khoshmanesh A, Christensen D, Perez-Guaita D, Iturbe-Ormaetxe I, O'Neill SL, McNaughton D, et al. Screening of wolbachia endosymbiont infection in aedes aegypti mosquitoes using attenuated total reflection mid-infrared spectroscopy. Anal Chem. 2017;89: 5285–5293. doi:10.1021/acs.analchem.6b04827

15. Stuart BH. Infrared spectroscopy: fundamentals and applications. Methods. 2004. doi:10.1002/0470011149

16. Baker MJ, Trevisan J, Bassan P, Bhargava R, Butler HJ, Dorling KM, et al. Using Fourier transform IR spectroscopy to analyze biological materials. Nature Protocols. 2014;9: 1771–1791. doi:10.1038/nprot.2014.110

17. Johnson BJ, Hugo LE, Churcher TS, Ong OTW, Devine GJ. Mosquito age grading and vector-control programmes. Trends in Parasitology. 2020;36: 39–51. doi:10.1016/j.pt.2019.10.011

18. Langley PA, Coates TW, Carlson DA. Sex recognition pheromone in the tsetse fly Glossina pallidipes Austen. Experientia. 1982;38: 473–475. doi:10.1007/BF01952645

19. Babayan S, Gonzalez M. SimonAB/Gonzalez-Jimenez_MIRS: First public release. Zenodo; 2019. doi:10.5281/ZENODO.2609356

20. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine learning in python. the Journal of machine Learning research. 2011;12: 2825–2830.

21. Eid F-E, Elmarakeby HA, Chan YA, Fornelos N, ElHefnawi M, Van Allen EM, et al. Systematic auditing is essential to debiasing machine learning in biology. Commun Biol. 2021;4: 1–9. doi:10.1038/s42003-021-01674-5

468    22.  Jurenka R, Terblanche JS, Jaco Klok C, Chown SL, Krafsur ES. Cuticular lipid mass and desiccation
469           rates in Glossina pallidipes: interpopulation variation. Physiological Entomology. 2007;32: 287–
470           293. doi:10.1111/j.1365-3032.2007.00571.x

471    23.  Schlein Y. Age grading of tsetse flies by the cuticular growth layers in the thoracic phragma.
472           Annals of Tropical Medicine & Parasitology. 1979;73: 297–298.
473           doi:10.1080/00034983.1979.11687262

474    24.  Mwanga EP, Mapua SA, Siria DJ, Ngowo HS, Nangacha F, Mgando J, et al. Using mid-infrared
475           spectroscopy and supervised machine-learning to identify vertebrate blood meals in the
476           malaria vector, Anopheles arabiensis. Malaria Journal. 2019;18: 187. doi:10.1186/s12936-019-
477           2822-y

478