

1 **The *G.m. morsitans* (Diptera: Glossinidae) genome as a source of**
2 **microsatellite markers for other tsetse fly (*Glossina*) species**

3

4 **Chaz Hyseni^{1*}, Jon S. Beadell¹, Zaneli Gomez Ocampo¹, Johnson O. Ouma², Loyce M. Okedi³, Michael**

5 **W. Gaunt⁴, Adalgisa Caccone¹**

6

7

8 ¹Department of Ecology and Evolutionary Biology, Yale University, New Haven, Connecticut, USA

9 ²Trypanosomiasis Research Centre, Kenya Agricultural Research Institute, Kikuyu, Kenya

10 ³ National Livestock Resources Research Institute, Tororo, Uganda

11 ⁴ London School of Tropical Medicine, London, United Kingdom

12

13

14

15

16

17

18 *E-mail: chaz.hyseni@aya.yale.edu

19 **Abstract**

20 We searched the *Glossina morsitans morsitans* genome for short sequence repeats (SSR) in order to
21 adapt polymorphic microsatellite markers to other species of *Glossina*, *G. fuscipes fuscipes* and *G.*
22 *pallidipes*, two major vectors of African trypanosomiasis. We tested 30 loci containing perfect di-, tri-, or
23 tetranucleotide repeats. We identified seven polymorphic loci that amplified across both *G.f. fuscipes*
24 and *G. pallidipes* samples, as well as seven additional loci that were variable in just one species. Five of
25 these fourteen loci were homozygous in males of one or both species and are likely to be X-linked.
26 Although the success rate of adapting SSR markers from the *G.m. morsitans* genome for use in other
27 species was not very high, this process yielded several polymorphic markers that should be useful in
28 future studies of tsetse ecology and evolution.

29

30

31

32

33

34

35

36

37

38

39 Trypanosomiasis is responsible for thousands of infected people and millions of cattle deaths with
40 serious repercussions for public health and economic development in 37 countries in sub-Saharan Africa
41 (<http://www.fao.org/Ag/againfo/programmes/en/paat/disease.html>). The tsetse species *Glossina*
42 *fuscipes* and *Glossina pallidipes* are primary vectors of both human and animal forms of African
43 trypanosomiasis (sleeping sickness and nagana, respectively). Molecular tools are increasingly being
44 used to garner information about the population structure, reproductive behavior, and other aspects of
45 tsetse ecology and life history (reviewed in Krafur 2009) pertinent to the design of more effective
46 vector control strategies. While microsatellite markers have been characterized for both *G. fuscipes*
47 *fuscipes* (Abila et al. 2008, Brown et al. 2008) and *G. pallidipes* (Ouma et al. 2003, Ouma et al. 2006),
48 many of these loci are difficult to score and exhibit high null allele frequencies. For *G.f. fuscipes*, a few
49 additional markers have been isolated as a result of recent efforts to adapt loci from other tsetse species
50 (Beadell et al. 2010, unpubl. data). Approximately 13 autosomal loci are now available for *G.f. fuscipes*
51 and 8 for *G. pallidipes*. Additional autosomal markers are required for exploring fine-scale structuring of
52 populations, for increasing the certainty of estimates of effective population size, as well as quantifying
53 dispersal and gene flow. Vector control would also benefit from further insight into the reproductive
54 behavior of tsetse. Development of X-linked loci would help provide such insight.

55 Recently, the Wellcome Trust Sanger Institute completed whole genome shotgun sequencing of *G.*
56 *morsitans morsitans* (http://www.sanger.ac.uk/Projects/G_morsitans/). Here, we explore the utility of
57 using the *G.m. morsitans* genome as a source of short sequence repeats (SSR) for use in other tsetse
58 species by testing a panel of SSRs for cross-amplification and variability in *G.f. fuscipes* and *G. pallidipes*.
59 We used both *MISA* (MicroSatellite identification tool, Thiel et al. 2003) and *Msatfinder* (Thurston &
60 Field 2005) to search the *G.m. morsitans* genome. An initial screening for 1-, 2-, 3-, 4-, 5- and 6-
61 nucleotide SSRs greater than 30bp in length produced over 6,000 results. From these we selected di-,
62 tri- and tetranucleotide SSRs 50-200bp in size with at least 100bp of sequence flanking either side and

63 used a script written in the stream editor program, *sed*, to extract the up- and downstream flanking
64 regions. These sequences were run through *Msatcommander* (Faircloth 2008), which includes a Primer3
65 (Rozen & Skaletsky 2000) extension. Primers were designed in *Msatcommander* by modifying the
66 default Primer3 settings. Thirty loci were identified for cross-amplification testing.

67 Loci were tested on samples collected from sites in western Uganda: Kabunkanga (*G.f. fuscipes*,
68 n=22) and Murchison Falls (*G. pallidipes*, n=23). Six of the twenty-three *G.p.* samples and eight of the
69 twenty-two *G.f.f.* samples were male. Individual flies were sexed by inspection of genitalia in order to
70 permit identification of X-linked loci. Since tsetse males are heterogametic, male individuals are
71 expected to be homozygous in X-linked loci.

72 DNA was extracted from fly legs using the Qiagen DNEasy Tissue Kit (Qiagen, Inc.) as per the
73 manufacturer's protocol. We used the M13-tailed primer method (Boutin-Ganache et al. 2001) to obtain
74 individual genotypes. Forward primers were 5'-tailed with a 15-mer M13 sequence (5'-
75 **TCCCAGTCACGACGT**-3') and in addition to the unmodified reverse primers, reactions included a third
76 primer with a fluorescent dye attached to the same M13 sequence. We prepared 12.5 μ L reactions using
77 1x PCR Buffer (50mM KCl, 10mM Tris-HCl, pH 8.3) and 2 mM $MgCl_2$ (Applied Biosystems), 0.22 mM each
78 dNTP and 3 μ g BSA (New England Biolabs), 5 pmol fluorescently-labeled M13 primer, 5 pmol reverse
79 primer, 0.4 pmol M13-tailed forward primer and 0.5 units AmpliTaq Gold DNA Polymerase (Applied
80 Biosystems). PCR amplification was done via touchdown thermal cycling: after the initial denaturation
81 (95°C for 5 min), reactions cycled through 95°C for 30s, 60-49°C (1°C decrement/cycle) for 25s and 72°C
82 for 30s in the first step, going through an additional 40 cycles of 95°C, 48°C and 72°C, and finally a 15-
83 minute extension at 72°C. Annealing temperatures of 53°C to 48°C in the second step and 60-54°C to 60-
84 49°C in the first step were tested to ensure that lower annealing temperatures did not result in stronger
85 PCR product and lower null allele frequencies due to non-specific amplification. We tested several loci

86 using 5 pmol each of fluorescently-labeled forward primer and unlabeled reverse primer, which along
87 with lower annealing temperatures enhanced PCR amplification compared to the above three-primer
88 method. PCR amplicons were run through an Applied Biosystems 3730xl DNA Analyzer and the resulting
89 data were analyzed using *Genemarker 1.91* (SoftGenetics).

90 Observed and expected heterozygosity values were calculated using *Genalex 6.2* (Peakall & Smouse
91 2006). We used *Genepop 4.0.10* (Rousset 2008) to test for departures from Hardy-Weinberg equilibrium
92 (HWE) and to determine whether any linkage disequilibrium (LD) was present among loci. Loci
93 developed in this study were also tested for linkage with autosomal loci developed for prior studies
94 (*G.f.f.*: Beadell et al. 2010; *G.p.*: Ouma et al. 2003, Ouma et al. 2006). We estimated null allele
95 frequencies in *Microchecker* (Oosterhout et al. 2004) using the method of Oosterhout et al. (2006) and
96 in *FreeNA* (Chapuis & Estoup 2007) using the method of Dempster et al. (1977). Calculations in *Genepop*
97 were performed using a burn-in of 10,000 and 1,000 batches with 10,000 iterations per batch. We
98 performed 10,000 randomizations in *Microchecker* and *FreeNA*.

99 Twenty-four of the thirty loci selected from the *G.m. morsitans* genome amplified in either *G.f.f.*,
100 *G.p.* or both species. Of these 24, only 15 loci were easily scored and polymorphic in at least one species.
101 We only reported 14 loci here (see Table 1), because locus GmmC22, polymorphic in *G.p.*, appeared to
102 be physically linked to a gene (GenBank accession no. AM940018.1). Although *G.p.* and *G.m.m.* share
103 more recent ancestry than *G.f.f.* and *G.m.m.* (Dyer et al. 2008), transfer of loci from *G.m.m.* to *G.p.* was
104 no more successful than transfer to *G.f.f.* We recovered 10 polymorphic markers for *G.p.* and 11
105 polymorphic markers for *G.f.f.*; only 7 loci were polymorphic in both species. In *G.f.f.*, the success rate of
106 adapting loci from the *G.m.morsitans* genome (11 polymorphic loci out of 30 tested) was lower
107 compared with the success rate of isolating polymorphic loci from a fuscipes-specific library (17
108 polymorphic loci out of 20 tested; unpubl. data).

109 Allele numbers varied from 2 to 8 for *G.f.f.* and 2 to 16 for *G.p.* Observed heterozygosities ranged
110 from 0.09 to 0.86 for *G.f.f.* and 0.13 to 0.87 for *G.p.* After applying sequential Bonferroni correction
111 (Holm 1979), the only significant LD among loci characterized in this study was observed between
112 GmmL17 and GmmP07 in the *G.p.* samples. When taking into account previously developed loci,
113 significant linkage was, again, only observed in *G.p.*, between GmmL11 and locus B115 of Ouma et al
114 (2006). The only locus to exhibit significant deviation from HWE after sequential Bonferroni correction
115 was locus GmmL17 in *G.f.f.*

116 Based on the observation that all male samples are homozygous, locus GmmL17 is likely X-linked.
117 Using the criterion of male homozygosity, we identified five loci with potential X-linkage: GmmC15,
118 GmmD03, GmmF10, GmmL17 and GmmP07. While *G.p.* males were homozygous for locus GmmC15,
119 however, there was one heterozygous male in *G.f.f.* X-linkage of loci GmmF10, GmmL17 and GmmP07 is
120 corroborated by low p-values for HWE tests and high (above 0.1) null allele frequencies, and loci
121 GmmL17 and GmmP07 were, at least in *G.p.*, also significantly linked to each other. Fluorescence in situ
122 hybridization (FISH) has been used to map microsatellite markers (Stratikopoulos et al. 2008). FISH will
123 be employed in future studies to confirm location in the X chromosome of the loci we suspect to be X-
124 linked.

125 The putative X-linked loci we identified will be valuable for remating studies, where SSRs are used to
126 estimate the number of unique male contributors to female sperm pools (Bonizzoni et al. 2002). The
127 other markers reported here should provide additional power for resolving fine-scale spatial structuring
128 of tsetse populations and determining the temporal stability of these populations. These autosomal
129 markers, in addition to the 13 that are already developed for *G.f.f.* and 8 for *G.p.*, will also allow more
130 accurate estimation of effective population size, of gene flow and migration rates, and will aid in

131 characterizing other facets of tsetse ecology and life history relevant to the design and improvement of
132 vector control strategies.

133 ***Acknowledgement***

134 Funding was provided by the NIH (R01AI068932) and WHO-TDR (80132) grants.

135 ***References***

- 136 Abila PP, Slotman MA, Parmakelis A, *et al.* (2008) High levels of genetic differentiation between Ugandan
137 *Glossina fuscipes fuscipes* populations separated by Lake Kyoga. *PLoS Neglected Tropical*
138 *Diseases*, **2**, e242.
- 139 Beadell JS, Hyseni C, Abila PP, *et al.* (2010) Phylogeography and population structure of *Glossina fuscipes*
140 *fuscipes* in Uganda: implications for control of tsetse. *PLoS Neglected Tropical Diseases*, **4**, e636.
- 141 Bonizzoni M, Katsoyannos BI, Marguerie R, *et al.* (2002) Microsatellite analysis reveals remating by wild
142 Mediterranean fruit fly females, *Ceratitis capitata*. *Molecular Ecology*, **11**, 1915-1921.
- 143 Boutin-Ganache I, Raposo M, Raymond M, Deschepper CF (2001) M13-tailed primers improve the
144 readability and usability of microsatellite analyses performed with two different allele-sizing
145 methods. *Biotechniques*, **31**, 24-27.
- 146 Brown JE, Komatsu KJ, Abila PP, *et al.* (2008) Polymorphic microsatellite markers for the tsetse fly
147 *Glossina fuscipes fuscipes* (Diptera: Glossinidae), a vector of human African trypanosomiasis.
148 *Molecular Ecology Resources*, **8**, 1506-1508.
- 149 Chapuis MP, Estoup A (2007) Microsatellite null alleles and estimation of population differentiation.
150 *Molecular Biology and Evolution*, **24**, 621-631.
- 151 Dempster AP, Laird NM, Rubin DB (1977) Maximum likelihood from incomplete data via the EM
152 algorithm. *Journal of the Royal Statistical Society Series B-Methodological*, **39**, 1-38.

- 153 Dyer NA, Lawton SP, Ravel S, *et al.* (2008) Molecular phylogenetics of tsetse flies (Diptera: Glossinidae)
154 based on mitochondrial (COI, 16S, ND2) and nuclear ribosomal DNA sequences, with an
155 emphasis on the palpalis group. *Molecular Phylogenetics and Evolution*, **49**, 227-239.
- 156 Faircloth BC (2008) MSATCOMMANDER: detection of microsatellite repeat arrays and automated, locus-
157 specific primer design. *Molecular Ecology Resources*, **8**, 92-94.
- 158 Holm S (1979) A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics*,
159 **6**, 65-70.
- 160 Krafur ES (2009) Tsetse flies: genetics, evolution, and role as vectors. *Infection, Genetics and Evolution*,
161 **9**, 124-141.
- 162 Ouma JO, Cummings MA, Jones KC, Krafur ES (2003) Characterization of microsatellite markers in the
163 tsetse fly, *Glossina pallidipes* (Diptera: Glossinidae). *Molecular Ecology Notes*, **3**, 450-453.
- 164 Ouma JO, Marquez JG, Krafur ES (2006) New polymorphic microsatellites in *Glossina pallidipes* (Diptera:
165 Glossinidae) and their cross-amplification in other tsetse fly taxa. *Biochemical Genetics*, **44**, 471-
166 477.
- 167 Peakall R, Smouse PE (2006) GENALEX 6: genetic analysis in Excel. Population genetic software for
168 teaching and research. *Molecular Ecology Notes*, **6**, 288-295.
- 169 Rousset F (2008) GENEPOP '007: a complete re-implementation of the GENEPOP software for Windows
170 and Linux. *Molecular Ecology Resources*, **8**, 103-106.
- 171 Rozen S, Skaletsky H (2000) Primer3 on the www for general users and for biologist programmers. In:
172 *Bioinformatics Methods and Protocols: Methods in Molecular Biology* (eds Krawetz S, Misener S),
173 pp. 365-386. Humana Press, Totowa, New Jersey.
- 174 Stratikopoulos EE, Augustinos AA, Petalas YG, *et al.* (2008) An integrated genetic and cytogenetic map
175 for the Mediterranean fruit fly, *Ceratitidis capitata*, based on microsatellite and morphological
176 markers. *Genetica*, **133**, 147-157.

- 177 Thiel T, Michalek W, Varshney RK, Graner A (2003) Exploiting EST databases for the development and
178 characterization of gene-derived SSR-markers in barley (*Hordeum vulgare* L.). *Theoretical and*
179 *Applied Genetics*, **106**, 411-422.
- 180 Thurston M, Field D (2005) MSATFINDER: Detection and characterization of microsatellites.
- 181 Van Oosterhout C, Hutchinson WF, Wills DPM, Shipley P (2004) MICRO-CHECKER: software for
182 identifying and correcting genotyping errors in microsatellite data. *Molecular Ecology Notes*, **4**,
183 535-538.
- 184 van Oosterhout C, Weetman D, Hutchinson WF (2006) Estimation and adjustment of microsatellite null
185 alleles in nonequilibrium populations. *Molecular Ecology Notes*, **6**, 255-256.

Table 1.

Locus	Primer Sequence (5'-3')	G.m.m genome: Sequence ID	G.m.m genome: (Motif)Units [Fragment Size]	Species	Allele Size Range	No. of Flies (% Amplified)	No. of Alleles	H _o	H _E	H-W(P)	X-linked	Null (O et al)	Null (D et al)
<i>GmmA06</i>	F: ACTCCATGTTATGTTCTGTC R: TGCCTTAGTTGAGAACTCTGC	Tfly_23-t487a06.q1k	(AC)31 [209]	<i>G.f.f</i>	146-174	22 (100%)	7	0.727	0.699	0.728		-0.035	0.000
				<i>G.p</i>	154-166	23 (100%)	3	0.130	0.127	1.000		-0.066	0.000
<i>GmmB20</i>	F: AAATGCATGTCTAACTGTCCG R: AGCAAAAGGCCAACTAAAGTGATG	Tfly_23-t572b20.q1k	(GT)33 [222]	<i>G.f.f</i>	171-181	22 (100%)	3	0.500	0.635	0.419		0.100	0.068
				<i>G.p</i>	221-281	23 (100%)	16	0.870	0.895	0.344		0.004	0.000
<i>GmmC15</i>	F: ACTGCATCTGCCTCTGTCG R: TGAACGAGAAAAATGTGAATGGTAAG	Tfly_23-t601c15.p1k	(ATT)20 [208]	<i>G.f.f</i>	193-211	22 (100%)	4	0.409	0.506	0.178		0.080	0.068
				<i>G.p</i>	190-211	23 (100%)	5	0.609	0.723	0.086	-	0.063	0.066
<i>GmmC17</i>	F: TGCGCTTTGAACGGAACG R: CTATGCCGCTGGCTTATC	Tfly_23-t506c17.p1k	(ATGT)14 [224]	<i>G.f.f</i>	184-188	22 (100%)	2	0.091	0.089	1.000		-0.047	0.000
				<i>G.p</i>	190-202	23 (100%)	2	0.304	0.322	1.000		0.016	0.011
<i>GmmD03</i>	F: TGCACTTACCGATTGCAC R: GTTGAAAGTTGGTTGTGACGC	Tfly_23-t605d03.q1k	(AAT)20 [189]	<i>G.f.f</i>	151-160	22 (100%)	3	0.136	0.132	1.000	-	-0.070	0.000
				<i>G.p</i>						No amplification			
<i>GmmD15</i>	F: GCATCACACTTTGCTTGCG R: CGTTGGAAACTAGACATCTCACG	Tfly_23-t535d15.q1k	(GT)59 [227]	<i>G.f.f</i>	149-159	22 (100%)	5	0.545	0.606	0.409		0.017	0.056
				<i>G.p</i>						No amplification			
<i>GmmF10</i>	F: TGCCTTTGATAGAGAAACCATC R: ACCTGGACACTTATACCGCTC	Tfly_23-t606f10.q1k	(AC)42 [248]	<i>G.f.f</i>								No amplification	
				<i>G.p</i>	186-200	23 (100%)	6	0.565	0.773	0.049	-	0.120	0.117
<i>GmmH09</i>	F: ACCCAGATAACCTATATTGCTCG R: CGTTCAGGCAGATACGAAAATTG	Tfly_23-t508h09.p1k	(AT)27 [191]	<i>G.f.f</i>	148-182	22 (100%)	7	0.864	0.829	0.121		-0.042	0.000
				<i>G.p</i>						Monomorphic			
<i>GmmK06</i>	F: TAACGTGCATGTGCGTGTG R: CCATCAATACGAGCAGACCG	Tfly_23-t504k06.q1k	(ATGT)12 [154]	<i>G.f.f</i>								Monomorphic	
				<i>G.p</i>	123-131	23 (100%)	3	0.435	0.414	1.000		-0.056	0.000
<i>GmmK22</i>	F: ACGCTTACGTTTCCGTTACAC R: AAGCTAACCGAACCAGCAC	Tfly_23-t513k22.p1k	(GTT)19 [228]	<i>G.f.f</i>								Monomorphic	
				<i>G.p</i>	192-198	23 (100%)	3	0.478	0.569	0.288		0.069	0.038
<i>GmmL03</i>	F: ACAGTCCAATTTTCGCCCG R: GGCCAACAATGTCATAAACCG	Tfly_23-t545l03.q1k	(AT)28 [210]	<i>G.f.f</i>	181-187	21 (95.5%)	3	0.381	0.528	0.290		0.129	0.079
				<i>G.p</i>						Monomorphic			
<i>GmmL11</i>	F: CCACCACTAACAACGACAGC R: TGGCTGGTTACAAGATTGCAC	Tfly_23-t516l11.p1k	(AT)27 [249]	<i>G.f.f</i>	216-246	22 (100%)	8	0.682	0.736	0.194		0.020	0.032
				<i>G.p</i>	250-252	23 (100%)	2	0.696	0.464	0.020		-0.448	0.000
<i>GmmL17</i>	F: CGTACATGCAAGGCAGAGC R: TCAACTGAAACCGAAAGAGC	Tfly_23-t608l17.q1k	(GGTT)10 [273]	<i>G.f.f</i>	283-299	22 (100%)	4	0.227	0.496	0.001*	+	0.233	0.191
				<i>G.p</i>	277-285	23 (100%)	2	0.217	0.322	0.165	+	0.129	0.092
<i>GmmP07</i>	F: ACTGACATATTGAGTTGAAAGGGG R: TCTTCCGTTAAATACAGAGTGACAG	Tfly_23-t516p07.q1k	(ATT)18 [234]	<i>G.f.f</i>	213-234	22 (100%)	4	0.227	0.357	0.038	+	0.137	0.121
				<i>G.p</i>	231-246	22 (100%)	4	0.348	0.563	0.010	+	0.168	0.142

Allele size ranges represent size ranges of amplified DNA minus the 15-mer M13 sequence. The asterisk (*) indicates significance after sequential Bonferroni correction. The plus sign (+) indicates a putative X-linked locus, based on the absence of heterozygous males from both species. The minus sign (-) indicates a potentially X-linked locus based on only one species. Heterozygosity calculated using *Genalex* 6.2: H_o - observed heterozygosity and H_E - expected heterozygosity. P-values for the Hardy-Weinberg (H-W(P)) equilibrium test computed with *GenePop* 4.0.10. Null allele frequencies estimated in *Microchecker* based on Oosterhout et al. (2006) (Null (O et al)) and *FreeNA* based on Dempster et al. (1977) (Null (D et al)).