

2 Mitochondrial DNAs provide insight into trypanosome phylogeny and molecular evolution

3 Authors

4 Kay C*

5 Williams TA

6 Gibson W

7 School of Biological Sciences, University of Bristol, Bristol, UK

8

9 * Corresponding author

10

11 Email addresses:

12 chris.kay@bristol.ac.uk

13 tom.a.williams@bristol.ac.uk

14 w.gibson@bristol.ac.uk

16 Abstract

17 **Background:** Trypanosomes are single-celled eukaryotic parasites characterised by the unique biology of
18 their mitochondrial DNA (mtDNA). African livestock trypanosomes impose a major burden on agriculture
19 across sub-Saharan Africa, but are poorly understood compared to those that cause sleeping sickness and
20 Chagas disease in humans. Here we explore the potential of trypanosome mtDNA to study the evolutionary
21 history of trypanosomes and the molecular evolution of their mtDNAs.

22 **Results:** We used long-read sequencing to completely assemble mtDNAs from four previously
23 uncharacterized African trypanosomes, and leveraged these assemblies to scaffold and assemble a further
24 103 trypanosome mtDNAs from published short-read data. While synteny was largely conserved, there
25 were repeated, independent losses of Complex I genes. Comparison of edited and non-edited genes
26 revealed the impact of RNA editing on nucleotide composition, with non-edited genes approaching the
27 limits of GC loss. African tsetse-transmitted trypanosomes showed high levels of RNA editing compared to
28 other trypanosomes. Whole mtDNA coding regions were used to construct time-resolved phylogenetic
29 trees, revealing deep divergence events among isolates of the pathogens *Trypanosoma brucei* and *T.*
30 *congolense*.

31 **Conclusions:** Our mtDNA data represents a new resource for experimental and evolutionary analyses of
32 trypanosome phylogeny, molecular evolution and function. Molecular clock analyses yielded a timescale for
33 trypanosome evolution congruent with major biogeographical events in Africa and revealed the recent
34 emergence of *Trypanosoma brucei gambiense* and *T. equiperdum*, major human and animal pathogens.

35 Keywords:

36 Trypanosome, kinetoplast, maxicircle, mitochondrial DNA, phylogeny, RNA editing

37 Background

38 Trypanosomes are a group of single-celled eukaryotic flagellates, including important pathogens of humans
39 and their livestock (*Trypanosoma* and *Leishmania*), plants (*Phytomonas*) and insects (*Crithidia*). A
40 distinctive feature of trypanosomes is the compartmentalization of the mitochondrial DNA (mtDNA) into an
41 organelle located at the proximal end of the flagellum, the kinetoplast, which contains a network of
42 interlocked circular DNAs of two types: maxicircles are equivalent to the mitochondrial genome of other
43 eukaryotes and minicircles encode the gRNAs used to edit the maxicircle transcripts(1,2). Thus both mini-
44 and maxicircles are essential for expression of mitochondrial genes. In trypanosomes, mitochondrial
45 transcripts are edited by the insertion or deletion of uridine residues at positions demarcated by short guide
46 RNAs (gRNAs) to yield mRNAs that can be correctly translated(3–5). Why this energetically costly and
47 potentially error prone mRNA processing step evolved, and how, are unanswered questions in
48 trypanosome biology, but RNA editing is found throughout the Kinetoplastea(1,6).

49 MtDNA is widely used in evolutionary, phylogenetic and population genetics analyses and has
50 proved particularly useful as a molecular clock to date speciation events, but extensive RNA editing might
51 potentially undermine the phylogenetic signal. Within the Kinetoplastea, trypanosomes are monophyletic
52 according to phylogenetic trees constructed from nuclear-encoded 18S ribosomal RNA (rRNA) and
53 glycosomal GAPDH genes(7,8), but it has proved difficult to date the emergence of particular lineages, as
54 trypanosomes have no fossil record and are not sufficiently host specific to allow dating by co-speciation
55 with their hosts. Nevertheless, the divergence date of two major groups of pathogenic trypanosomes in
56 Africa (*T. brucei* clade) and South America (*T. cruzi* clade) has been linked to the breakup of Gondwana
57 during the Cretaceous, ~100 Mya(7). The *T. brucei* clade comprises the Salivaria, trypanosomes
58 transmitted via the mouthparts of bloodsucking tsetse flies (*Glossina*) in sub-Saharan Africa, while the *T.*
59 *cruzi* clade contains the agent of Chagas disease, *T. cruzi*, and related New World trypanosomes(9). The
60 100 Mya date has been used to calibrate subsequent trees, e.g. Lewis et al(10) estimated that *T. cruzi*
61 lineages radiated 3.35 Mya and dated the emergence of two hybrid lineages of *T. cruzi* to <60,000 years
62 ago. However, the discovery of trypanosomes from wild animals in Africa that belong to the *T. cruzi* clade
63 suggested the possibility of intercontinental transfer more recently via bats or rodents(11) so dating the

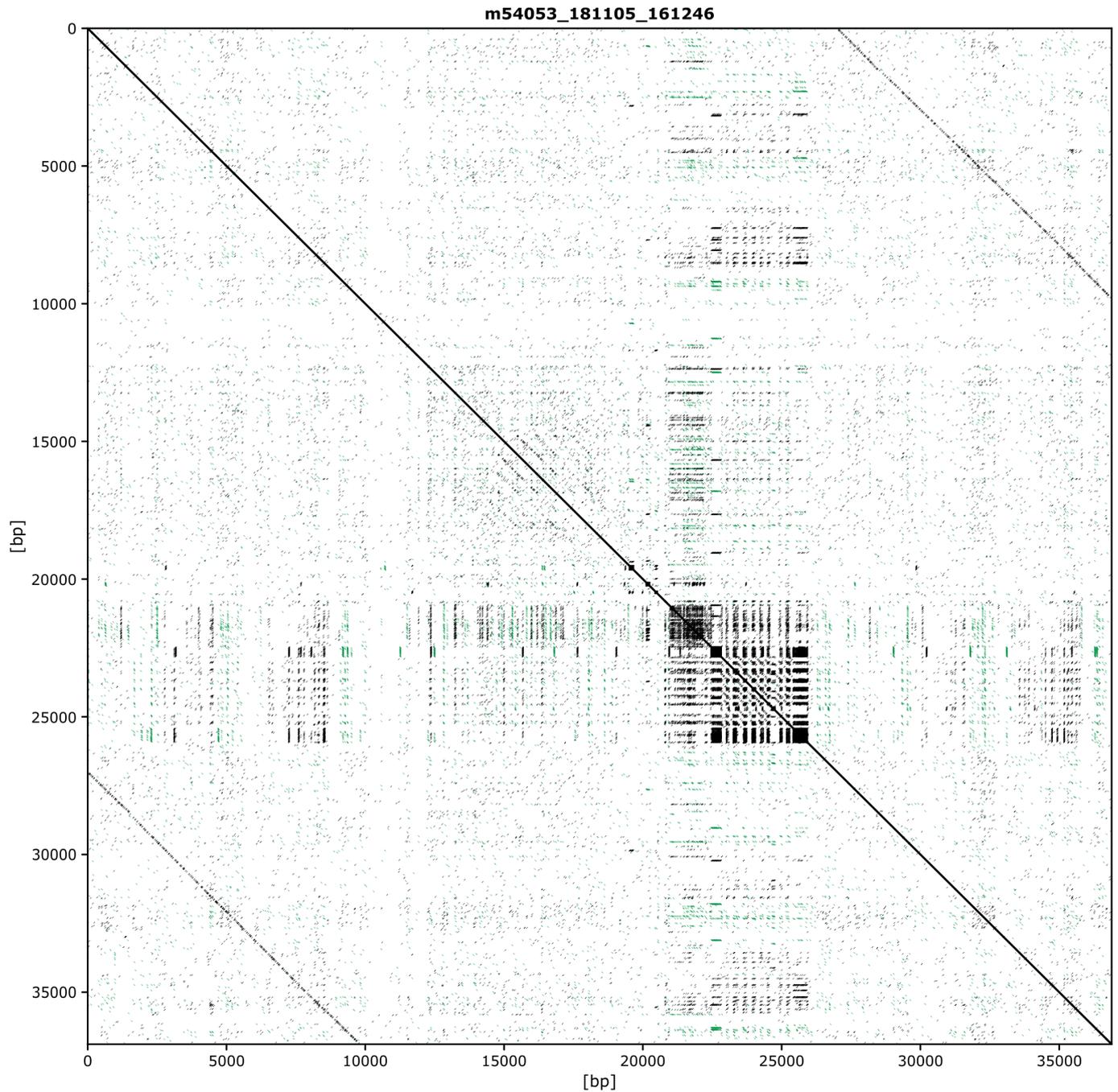
64 emergence of trypanosome lineages remains uncertain. A means to infer origins independent of sparse
65 historical information would give valuable insights into the emergence of different pathogens, as well as
66 provide information on how quickly trypanosomes can switch hosts and vectors, with implications for the
67 emergence of new diseases.

68 Here we have examined the potential of the trypanosome maxicircle for phylogenetic inference and
69 dating. We used long-read sequencing to completely assemble mtDNAs from four previously
70 uncharacterized African trypanosomes and leveraged these assemblies to scaffold and assemble a further
71 103 trypanosome mtDNA coding regions, exploiting the wealth of published short read data. We show that
72 time-resolved phylogenetic trees based on the mtDNA coding region can be used to explore events in the
73 recent history of *Trypanosoma brucei* and *T. congolense*, and infer ages which fit well with historical
74 evidence. Our analyses of edited and non-edited mitochondrial genes indicate very high levels of RNA
75 editing in salivarian trypanosomes, limiting further evolution in this direction without incurring functional
76 costs.

77

78

79



81 **Additional Figure 1. An example of a PacBio read spanning the entire sequence of the trypanosome**
 82 **mitochondrial DNA (maxicircle).** A single read from the *T. congolense* GAM2 readpool is shown dot-
 83 plotted against itself. The highly repetitive short period portion of the non-coding region is visualised as a
 84 densely self-similar region between 20-26 kbp, whilst the longer period portion of the non-coding region
 85 begins at 15 kbp. The remainder of the sequence shown belongs to the coding region. The complete length
 86 of the maxicircle is seen from 0-26.3 kbp, and thereafter begins to repeat. The assembled sequence is
 87 shown in **Additional Figure 2.**

88 Results

89 Whole maxicircle sequences

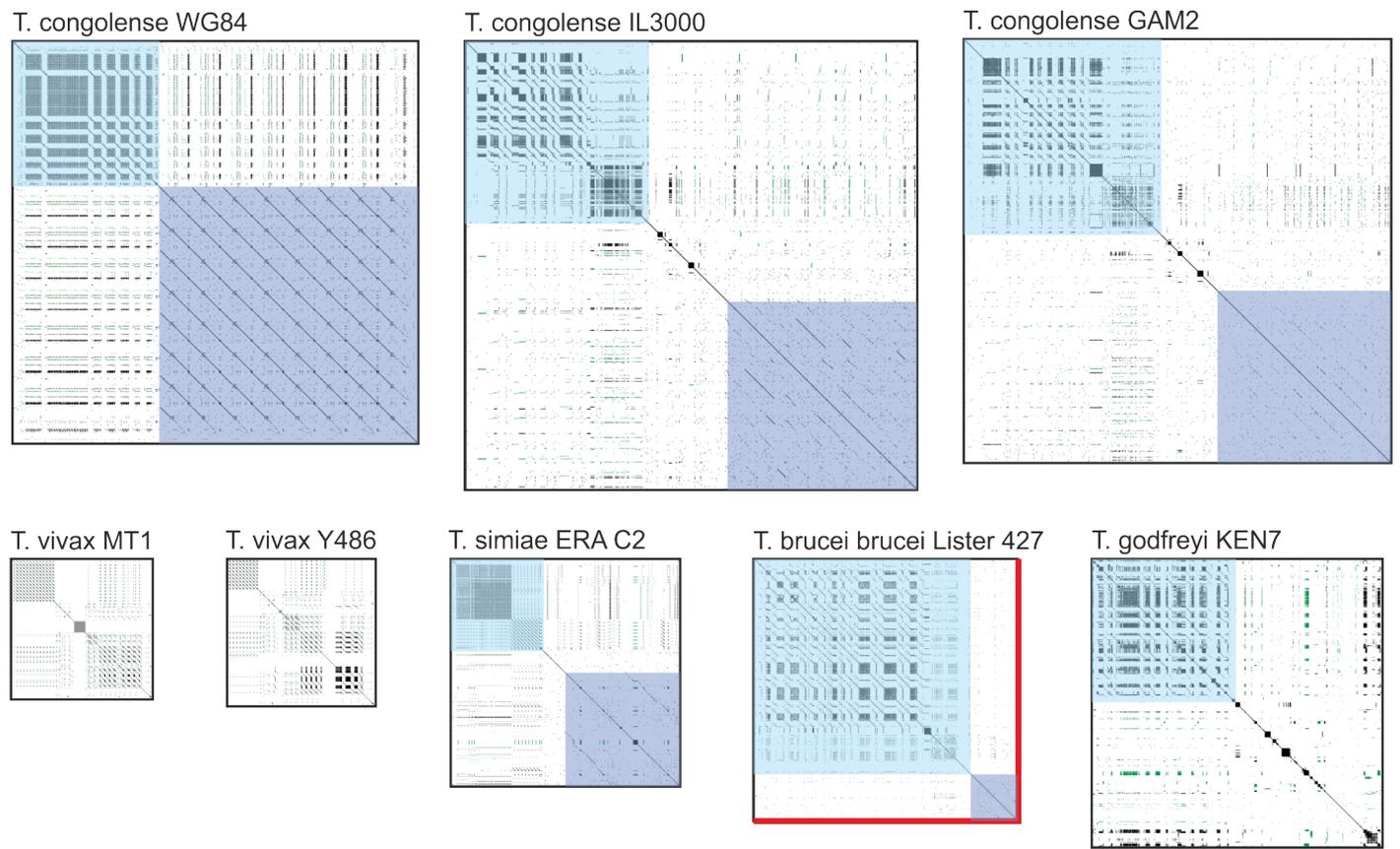
90 We sequenced (PacBio Sequel) mtDNA from four African trypanosomes (*T. congolense* savannah
91 and kilifi, *T. simiae*, and *T. godfreyi*), and assembled complete maxicircles (including the non-coding region,
92 **Additional Table 2**) using the long-read assemblers Canu and Flye(12,13). These novel data include the
93 first complete mtDNAs for *T. simiae*, *T. godfreyi* and the divergent *T. congolense* kilifi subgroup. We
94 assembled two additional complete maxicircles for the genome strains *T. congolense* IL3000 and *T. vivax*
95 Y486 from published data. We then assembled a further 101 maxicircle coding regions from public genome
96 sequence data, using reference sequences or our new assemblies to recover maxicircle reads. In total, we
97 obtained 51 complete maxicircle coding regions for *Trypanosoma brucei*, 34 for *T. congolense*, 3 for *T.*
98 *equiperdum*, 2 for *T. godfreyi*, 1 for *T. grayi*, 2 for *T. simiae*, and 14 for *T. vivax* (**Additional Table 1**). No
99 significant heteroplasmy was detected during sequence assembly.

100 Complete maxicircles ranged between 19.8 kbp (*T. vivax* Y486) and 27.6 kbp (*T. congolense*
101 IL3000), with most of the size variation occurring in the non-coding region (4.6 kbp in *T. vivax* Y486 to 12.6
102 kbp in *T. congolense* IL3000; **Additional Table 2**). The overall GC content was 20.9 – 23.7%, but the GC%
103 of the non-coding region was much lower (14.1% *T. godfreyi* KEN7 to 17.2% in *T. vivax* Y486). No
104 significant correlation was found between coding and non-coding region GC% (n=6, $\rho = -0.20$, $P=0.70$),
105 suggesting that changes to the composition of the non-coding/coding regions are independent. Dot plots of
106 the non-coding regions typically showed two domains: one densely repetitive with short repeats and the
107 other with longer period self-similarity (**Additional Figure 2**). Whilst the organisation of the non-coding
108 region was similar between isolates of the same species (*T. vivax*, *T. congolense*), variation was seen in
109 the fine structure and repeat copy number.

Maxicircle	Sequencing technology	Total size (kbp)	Coding region (kbp)	Non-coding region (kbp)	Whole GC%	Coding region GC%	Non-coding region GC%
<i>T. vivax</i> Y486	Sanger	19.8	15.2	4.6	22.5	24.1	17.2
<i>T. simiae</i> ERA C2	PacBio	22.1	15.0	7.1	22.2	25.3	15.9
<i>T. godfreyi</i> KEN7	PacBio	23.4	14.4	9.0	20.9	25.2	14.1
<i>T. congolense</i> GAM2	PacBio +	26.3	15.0	11.3	23.7	25.9	14.9
<i>T. congolense</i> WG84	Illumina PacBio	26.9	15.0	11.9	22.5	25.9	17.2
<i>T. congolense</i> IL3000	Illumina	27.6	15.0	12.6	21.7	25.8	16.6

113 **Additional Table 2. Properties of the six complete assembled salivarian maxicircles.**

114 Complete sequences were identified by the assembly of circular sequences. For PacBio data, individual
 115 reads spanning the entire maxicircle (**Additional Figure 1**) were used to validate the assembled sequence.



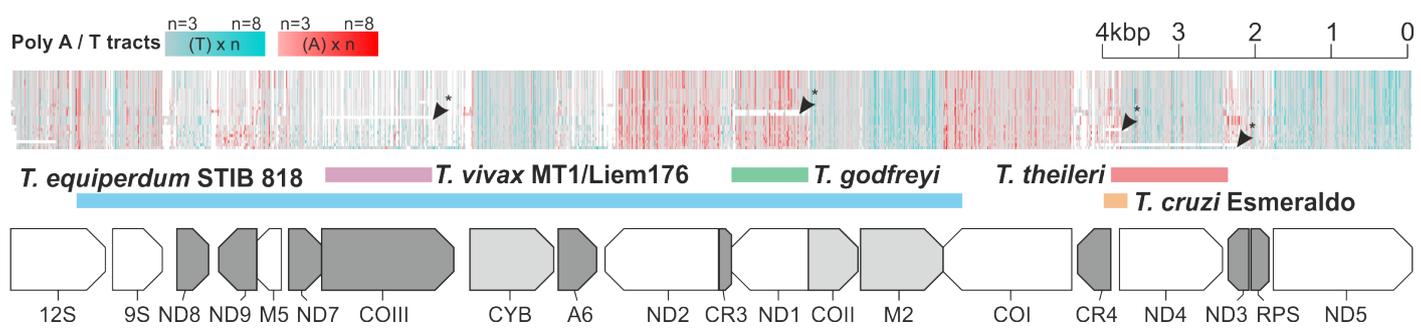
118 **Additional Figure 2. Assembled non-coding regions of salivarian mitochondrial DNAs.**
 119 Reference sequences for the complete non-coding region of *T. vivax* MT1 as well as the truncated (**red**
 120 **line**) non-coding region for *T. b. brucei* Lister 427, are shown to scale against other assembled salivarian
 121 non-coding regions.

122 Independent deletions of Complex I genes

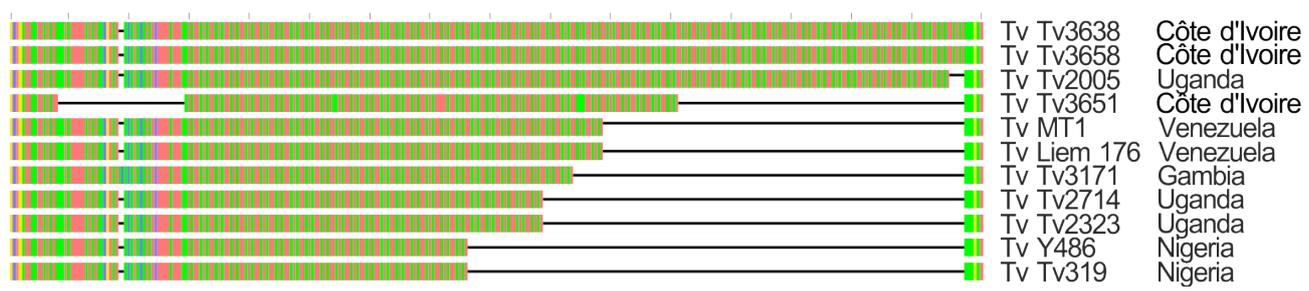
123 Alignment of the coding regions showed overall conservation of synteny (**Figure 1**); however, major
124 gene deletions were evident in *T. godfreyi* (ND1) and *T. theileri* (ND4), as well as previously described
125 deletions in New World *T. vivax*(14), *T. cruzi*(15), and *T. equiperdum*(16). The deletion of ND1 in *T.*
126 *godfreyi* was surprising, as this species undergoes full cyclical development in tsetse flies unlike New World
127 *T. vivax* and *T. equiperdum*, which have both adapted to non-tsetse transmission and evidently do not
128 require a fully functional mitochondrion. The deletion of ND4 in *T. theileri* has also eroded neighbouring
129 genes, CR4 and ND3. Like *T. godfreyi*, *T. theileri* is predicted to require a functional mitochondrion as it
130 completes development to mammal-infective forms in the gut of tabanid flies(17). One possibility is that
131 these deletions represent an early stage of mtDNA reduction in which mitochondrial function is reduced but
132 not abolished.

133 Discounting sequences with segmental gene deletions, the size of the whole coding region (WCR)
134 showed variation across *Trypanosoma* (**Table 1**) and trypanosome WCRs were approximately 1 kbp
135 smaller than the reference sequences for related trypanosomatids *Crithidia* and *Leishmania*. Among the
136 salivarian trypanosomes (lower portion of **Table 1**), coding regions without deletions (n=10) are significantly
137 smaller (Kruskal-Wallis one-way ANOVA on ranks, H=8.0, $P=0.05$) than those of non-salivarian
138 trypanosomes (n=4). These size differences can be traced to changes in the length of edited genes, for if
139 they are summed (**Table 1** 'Σedited') and subtracted from the whole (**Table 1** 'WCR-Σedited'), the
140 remaining coding region is relatively invariant in length (~12.4 kbp).

141 Coding regions frequently contained long homopolymers of either A or T, which appear to relate to
142 reading direction of the gene (**Figure 1**), indicating a strand specific bias. In contrast, the untranslated
143 rRNA genes have low GC content but no directional bias in poly-A/T. Comparatively little expansion and
144 contraction was observed in the intergenic regions, although in *T. vivax* a putative microsatellite was
145 identified between the 9S and ND8 genes (**Additional Figure 3**).



148 **Figure 1. Global overview of maxicircle coding region.** Top: alignment of the mitochondrial DNA coding regions
 149 from 27 isolates (top to bottom, *Trypanosoma brucei* H866, 1829 ALJO, Lister 427, TREU 927, TSW 55, J10, LF1; *T.*
 150 *congolense* IL3000, WG81, GAM2, IL3900, IL3578, ERA D1; *T. simiae* ERA C2; *T. godfreyi* KEN7, ERA F1; *T. vivax*
 151 Liem 176, Y486, Tv2323, *T. cruzi*, CL Brener, Esmeraldo; *T. lewisi*, *T. conorhini*, *T. copemani*, *T. grayi*, *T. theileri*,
 152 *Leishmania tarentolae*). Tracts of poly-T or poly-A are shown coloured turquoise or red respectively. An approximate
 153 scale is shown. Segmental gene deletions in the alignment are indicated by arrows and are also shown below as
 154 coloured bars; the deletion from *T. equiperdum* STIB 818 is also shown for comparison. Bottom: gene order in the
 155 coding region. Non-edited genes in are shown in white, minimally edited genes in light grey and extensively edited
 156 (pan-edited) genes in grey. Editing categories are on the basis of *T. vivax*(14).



160 **Additional Figure 3. Trypanosoma vivax**, an alignment of the intergenic region between 9S and ND8
 161 containing a putative microsatellite. Bases are shown as coloured bands with the top line tick showing 20bp
 162 increments. The sequence ATATA is tandemly repeated between 18 and 51 times in the selected isolates.

Species	WCR (bp)	Edited genes (bp)											Non-edited genes			%AU codons		
		ND8	ND9	ND7	COIII	A6	CR3	CR4	ND3	RPS12	Σ edited	WCR- Σ edited	COI	ND4	ND5*	COI	ND4	ND5*
<i>C. fasciculata</i> Cf-CI	16118	316	293	339	1565	592	101	208	174	131	3618	12500	18		17	35		43
<i>L. tarentolae</i>	16169	276	313	339	1595	556	103	208	188	153	3628	12541	13	9	8	37	59	49
<i>T. theileri</i> Edinburgh TM35	13727	312	370	311	956	331	106			165	2445	11282	19			34		
<i>T. conorhini</i> 025E	15415	303	402	307	1001	337	122	235	219	163	2967	12448		15			51	
<i>T. copemani</i> G1	15215	293	365	290	958	318	114	241	174	157	2796	12419	22	18	22	32	46	41
<i>T. lewisi</i>	15274	295	360	289	989	304	121	235	183	167	2822	12452	16	15	13	36	49	48
<i>T. cruzi</i> CL Brener	15255	289	349	293	959	337	115	238	202	167	2834	12421	19	14	15	35	47	44
<i>T. grayi</i> ANR4	14919	296	369	323	983	340	119	256	203	170	2940	11979	22	18	19	32	48	40
<i>T. vivax</i> Y486	15182	282	295	278	886	321	120	227	224	170	2683	12499	14	9	7	38	58	54
<i>T. godfreyi</i> KEN7	14436	250	288	280	855	311		270	220	160	2634	11802	13	13	11	37	51	49
<i>T. simiae</i> ERA C2	14975	248	275	284	842	298	104	253	218	160	2578	12397	17	20		37	53	
<i>T. simiae</i> tsavo KETRI 3436	14587	247	271		850	297	100	328	213	236	2442	12145	12	9	9	38	55	53
<i>T. brucei</i> brucei Lister 427	14874	276	259	282	868	307	102	222	196	144	2554	12320	18	12	12	36	51	49
<i>T. brucei</i> gambiense 1829 ALJO	14882	276	259	281	868	307	102	224	197	144	2556	12326	19	12	12	35	52	48
<i>T. brucei</i> rhodesiense H866	14880	276	259	281	868	307	102	224	196	144	2555	12325	18	12	12	35	51	48
<i>T. equiperdum</i> BoTat	14883	278	259	285	868	309	103	225	197	144	2565	12318	18	13	12	36	52	48
<i>T. congolense</i> IL3578 (s)	14939	263	286	282	844	295	98	227	232	154	2583	12356	18	11	10	34	54	50
<i>T. congolense</i> IL3900 (f)	15166	262	293	283	1034	299	100	222	226	157	2776	12390	15	12	11	34	54	51
<i>T. congolense</i> ERA D1 (k)	15016	253	299	282	852	374	97	229	224	150	2663	12353	19	14	12	34	51	50
<i>T. congolense</i> IL3000 (s)	14941	263	286	283	846	294	98	227	231	154	2584	12357	17	11	11	34	53	50

163

164

165 **Table 1. Characteristics of edited and non-edited mitochondrial genes in trypanosomatids.**

166 **Left: Variation in coding region size relates to the sequence contribution of edited genes.**

167 Trypanosome whole coding regions (WCR) are shorter than for related trypanosomatids *Leishmania* and
 168 *Crithidia*. These size differences reflect the contraction of edited regions (Σ edited); the remaining coding
 169 region (WCR- Σ edited) is relatively invariant in length. All numbers reflect the length of ungapped
 170 sequences.

171 **Right: Non-edited genes show trends of GC loss which suggest vulnerability to loss of function.**

172 Universally non-edited genes COI, ND4 and the first 500 codons of ND5 (ND5*) were analysed for codon
 173 usage (% AU codons) and possible composition changes (reducible GC%), which is the percentage of
 174 alternate synonymous codons with reduced GC content.

175 Coding regions which are **incomplete** or have **gene deletions** that impact calculations are highlighted.

176 Numbers have been shaded by value order on a column by column basis. Trypanosomes below the solid
 177 line, with non-salivarian above and salivarian below the dashed line.

Species	GC%																T:C																									
	WCR	12S	9S	ND8	ND9	M5	ND7	COIII	CYB	A6	ND2	CR3	ND1	COII	M2	COI	CR4	ND4	ND3	RPS12	ND5	WCR	12S	9S	ND8	ND9	M5	ND7	COIII	CYB	A6	ND2	CR3	ND1	COII	M2	COI	CR4	ND4	ND3	RPS12	ND5
<i>C. fasciculata</i> Cf-CI	26.0	16	19	40	41	15	34	29	25	22	21	38	31	28	21	31	42	23	47	46	25	3	6	5	2	2	7	3	6	6	10	6	5	4	4	8	4	2	5	2	1	6
<i>L. tarentolae</i>	22.9	16	18	50	42	13	32	25	24	21	16	38	26	26	15	29	44	18	45	39	20	4	7	6	2	5	7	3	8	6	11	7	5	4	5	14	4	5	9	2	3	9
<i>T. theileri</i> Edinburgh TM35	29.4	21	19	47	37	22	42	46	28	42	24	41	31	27	22	30					45	27	2	5	5	2	12	4	3	3	5	4	6	8	5	5	5				3	3
<i>T. conorhini</i> 025E	26.8	17	17	48	33	19	42	43	27	41	19	33	29	26	18	31	43	21	36	43	24	3	6	5	2	11	6	2	3	5	4	8	12	5	5	8	5	20	5	8	2	5
<i>T. copemani</i> G1	28.7	17	17	50	39	20	44	47	28	45	22	34	29	27	22	32	44	23	47	47	27	3	6	6	2	9	5	2	3	5	3	6	11	4	5	5	4	29	5	4	2	4
<i>T. lewisi</i>	26.4	18	18	49	39	14	44	44	24	45	19	36	27	28	18	29	43	21	42	43	22	3	5	5	2	9	11	2	3	7	3	7	11	4	5	7	5	27	5	5	3	6
<i>T. cruzi</i> CL Brener	26.1	18	17	49	34	17	43	44	26	37	19	31	26	28	18	30	34	22	36	42	23	3	6	5	2	8	5	2	3	6	4	5	11	4	5	8	4	27	5	6	3	7
<i>T. grayi</i> ANR4	28.0		16	49	38	18	40	46	25	41	21	33	30	25	20	31	44	22	44	41	25	3		6	2	8	4	3	3	5	5	7	12	4	5	6	5	35	4	6	3	5
<i>T. vivax</i> Y486	24.1	14	16	49	47	15	48	52	22	47	15	35	23	22	15	28	44	17	36	49	18	4	10	8	2	5	6	2	3	9	4	9	12	5	6	13	4	25	9	10	2	11
<i>T. godfreyi</i> KEN7	25.2	17	15	55	47	13	48	52	23	46	19			22	16	28	36	20	38	43	22	4	7	9	1	4	5	2	2	8	3	5		8	14	4	19	8	6	3	9	
<i>T. simiae</i> ERA C2	25.3	15	15	55	47	13	46	52	23	49	18	41	26	22	16	28	41	18	40	45	20	4	9	8	1	4	5	2	2	8	3	5	8	4	8	17	4	11	10	6	3	9
<i>T. simiae</i> tsavo KETRI 3436	24.6	15	15	56	48		51			49	17	39	25	22	15	27	35	18	41	48	20	4	9	9	1	4		2		2	6	7	4	7	18	4	11	10	6	2	9	
<i>T. brucei</i> brucei Lister 427	26.1	17	17	53	51	15	46	50	23	45	18	38	26	24	15	29	46	20	46	49	21	4	7	7	2	3	4	2	2	7	3	5	8	4	6	15	4	16	6	5	2	9
<i>T. brucei</i> gambiense 1829 ALJO	26.1	17	17	53	52	15	46	50	23	46	18	38	26	24	15	30	47	20	46	49	21	4	8	7	2	3	4	2	2	7	3	5	8	4	6	14	4	16	6	5	2	10
<i>T. brucei</i> rhodesiense H866	26.2	17	17	53	52	15	46	50	23	45	18	38	26	24	15	30	47	20	46	49	21	4	8	7	2	3	4	2	2	7	3	5	8	4	6	14	4	16	6	5	2	10
<i>T. equiperdum</i> BoTat	26.2	17	17	52	51	16	46	50	24	45	18	38	26	25	15	30	46	20	46	49	21	4	8	8	2	3	4	2	2	7	3	5	8	4	6	14	4	17	7	5	2	9
<i>T. congolense</i> IL3578 (s)	25.8	16	16	52	47	15	46	51	23	49	18	41	27	24	16	29	44	19	41	44	20	4	8	9	1	5	3	2	2	8	3	5	7	3	6	14	4	18	8	7	2	10
<i>T. congolense</i> IL3900 (f)	25.6	16	16	53	45	10	46	48	23	48	19	37	25	23	16	29	44	19	39	43	21	4	8	8	1	5	7	2	2	8	3	5	7	4	6	15	4	17	8	7	3	10
<i>T. congolense</i> ERA D1 (k)	26.6	17	15	55	42	14	47	52	25	45	20	40	28	24	17	30	41	20	42	46	21	4	8	9	1	6	5	2	2	7	3	4	6	3	7	14	3	15	7	5	2	11
<i>T. congolense</i> IL3000 (s)	25.8	17	16	52	47	15	46	52	23	49	18	41	26	24	16	29	44	19	40	44	20	4	8	9	1	5	3	2	2	8	3	5	7	3	6	13	4	18	8	7	2	10

179 **Table 2. Individual mtDNA genes show different trends for GC composition and T:C ratios.**

180 From an alignment of coding regions, aligned sequence regions were extracted and analysed in the
 181 reading direction for GC% (left) and the ratio of T:C (right). Extensively edited genes (dark grey) have
 182 greater GC% than lightly (light grey) or non-edited (white) genes. Likewise, the T:C ratio is very high in
 183 some edited genes e.g. ND9, CR3, CR4. Shading as Table 1.

184 High levels of RNA editing in salivarian trypanosomes

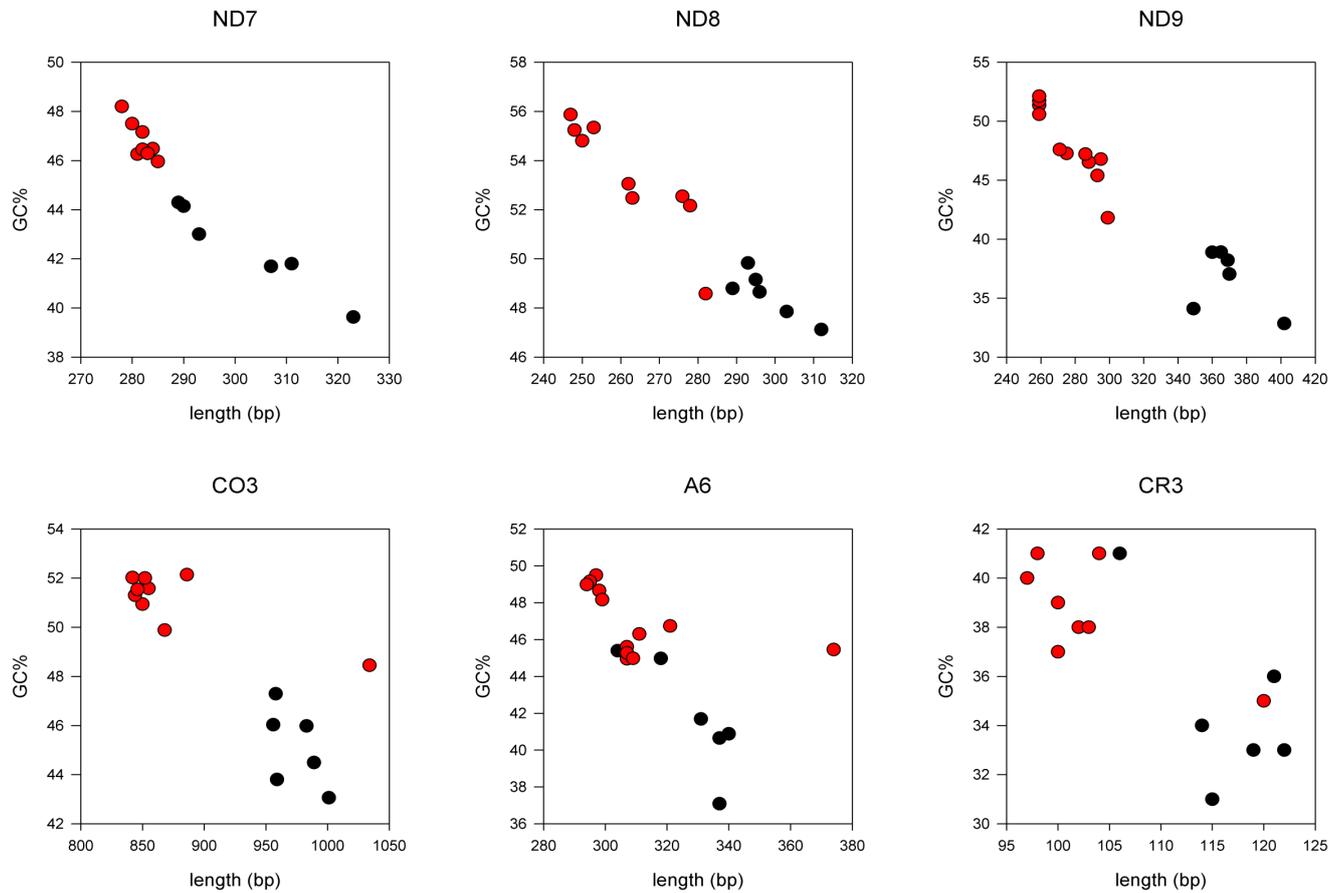
185 The overall GC% of the maxicircle coding region was low (~ 25%) in trypanosomes and other
 186 trypanosomatids, but these mean values conceal the fact that edited genes are far more GC-rich than non-
 187 edited genes (**Table 2**). Comparison of these genes indicates that they are significantly more GC-rich in
 188 **salivarian** compared to non-salivarian trypanosomes (**ND8, ND9, COIII, A6** (n=12,6); **CR3, ND7** (n=11,6);
 189 Kruskal-Wallis one-way ANOVA on ranks all $H > 14$, $P = < 0.001$). These genes also show variation in T:C
 190 ratios, with particularly high T:C ratios in ND9, CR3 and CR4 (**Table 2**).

191 Some edited genes showed a large variation of sequence length (**Table 2**), which was inversely
 192 proportional to GC% ($\rho = \mathbf{ND8}_{(n=20)} -0.91$, $\mathbf{ND9}_{(n=20)} -0.93$, $\mathbf{ND7}_{(n=19)} -0.98$, $\mathbf{COIII}_{(n=20)} -0.96$, $\mathbf{A6}_{(n=20)} -0.95$,
 193 $\mathbf{CR3}_{(n=19)} -0.80$, all $P = < 0.001$, **Figure 2, Additional Figure 4**). Analysis of base composition reveals a
 194 strong proportional correlation of sequence length to T% (**Additional Figure 4**) whilst A% has a weaker
 195 inverse correlation ($\rho = \text{A\%} / \text{T\%}$, $\mathbf{ND8}_{(n=20)} -0.72/0.98$, $\mathbf{ND9}_{(n=20)} -0.65/0.96$, $\mathbf{ND7}_{(n=19)} -0.72/0.99$, $\mathbf{COIII}_{(n=20)} -$
 196 $0.84/0.97$, $\mathbf{A6}_{(n=20)} -0.79/0.96$, $\mathbf{CR3}_{(n=19)} -0.75/0.97$, all $P = < 0.005$). Providing that the translated product
 197 remains similar in size, shorter genes indicate a greater extent of editing in salivarian compared to non-
 198 salivarian trypanosomes; presumably the increase in GC% and decline specifically in T% is offset by U
 199 insertion during RNA editing.

Non-edited genes are approaching limits to GC loss

Unlike pan-edited genes, non-edited and lightly edited genes are characterised by low GC% (**Table 2**). However the T:C ratios for some genes vary significantly between trypanosomes (**Table 2**), especially between **salivarian** and non-salivarian trypanosomes (**12S**_(n=12,5), **9S**_(n=12,6), **CYB**_(n=11,6), **M2**_(n=12,6), **COI**_(n=12,6), **ND4**_(n=12,5), **ND5**_(n=12,6), one-way ANOVA all $F > 70$, $P < 0.001$; **COII**_(n=12,6), $F = 3.6$, $P = 0.06$), whilst other genes do not show significant change (**M5**_(n=11,6) $F = 0.56$, $P = 0.46$; **ND2**_(n=12,6) $F = 2.92$, $P = 0.10$; **ND1**_(n=11,6) $F = 5.48$, $P = 0.03$).

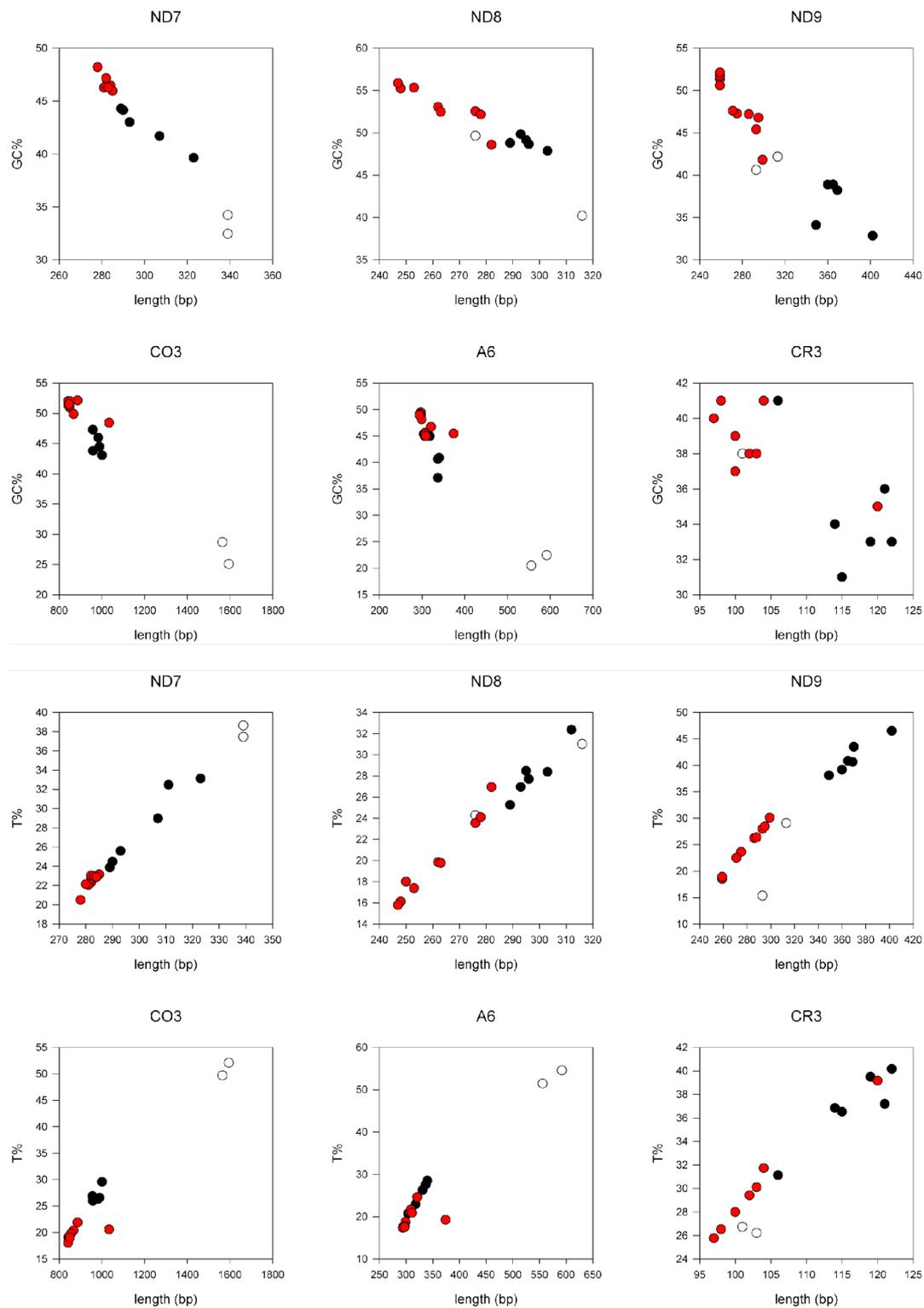
The low GC% of non-edited genes suggests that further reduction might lead to non-synonymous substitutions. Indeed, the total number of AU codons (no G or C) is already high, exceeding 50% in ND4 and ND5, and only a small proportion of remaining codons could be converted to AU codons without incurring translational changes (**Table 1**). Six amino acids (F, I, K, L, N, Y) solely use AU codons anyway, and for amino acids encoded by GC or AU codons, the AU codon was strongly preferred. Thus further GC loss would either result in non-synonymous mutations or introduce a compositional bias in the gene product, suggesting that non-edited genes have reached the limit of GC loss, particularly in salivarians.



226

227 **Figure 2. Correlation between sequence length and GC% for six edited maxicircle genes.**

228 There is an inverse correlation between the length and GC% of edited genes in trypanosome maxicircles,
 229 where salivarian trypanosomes (●) generally have shorter, more GC-rich edited genes than non-salivarian
 230 trypanosomes (●). This relationship also holds true for other trypanosomatid species (**Additional Figure 4**).



232

233

234 **Additional Figure 4. Correlation between sequence length, GC% and T% for six edited maxicircle**
 235 **genes.**

236 Some edited maxicircle genes exhibit transcript length variation with strong correlation between length and

237 T% as well as an inverse correlation for GC%. The weak negative correlation for A% indicates that this is a
 238 strand specific phenomenon consistent with RNA editing, where uridines are inserted back into the
 239 transcript. Key: *Crithidia*, *Leishmania* ○, salivarian ● and non-salivarian trypanosomes ●.

240 Time-resolved phylogeny of African trypanosomes

241 To test whether the discrete mechanisms driving sequence change in the maxicircle coding region would
 242 affect phylogenetic analysis between species of trypanosomes, alignments were prepared of (a) individual
 243 genes, (b) sets of edited and non-edited genes, and (c) the whole coding region (WCR) with and without
 244 edited genes (**Additional Figure 5**). Trees inferred from individual genes (e.g. COI) were congruent,
 245 providing no evidence of recombination between mtDNA loci, and strongly supported the monophyly of
 246 salivarians, although they showed weak resolution in the topology of deeper branches. Sets of non-edited
 247 genes had consistent topology for deeper branches, but were so conserved that the resolving power within
 248 species was limited. Topologies inferred from edited regions alone, which as a whole are faster-evolving,
 249 resolved intraspecific groups confidently but presented conflicting topologies for deeper branches. The
 250 WCR, including edited genes, gave better resolving power (in terms of bootstrap support) than individual
 251 genes or sets of non-edited genes. Therefore, WCRs appear to be useful phylogenetic markers for
 252 trypanosome evolution(18) and were used in subsequent analyses.

253 Aligned WCRs of *T. congolense* (including savannah, forest and kilifi subgroups) and *T. brucei*
 254 (including *T. equiperdum* strains with complete coding regions) were used to construct time-resolved
 255 phylogenies (**Figure 3**). To date species trees, the molecular clock was calibrated using tip isolation dates
 256 (**Additional Data 1**). Best marginal likelihoods were obtained with birth-death models using strict molecular
 257 clocks. Clock rates for *T. brucei* (median 1.81×10^{-7} substitutions/site/year, s/s/y, 95% HPD interval 1.00×10^{-7}
 258 $- 5.56 \times 10^{-7}$ s/s/y) and *T. congolense* (median 7.45×10^{-7} s/s/y, 95% HPD interval 1.43×10^{-8} - 2.89×10^{-6} s/s/y)
 259 were found to be similar in magnitude but significantly different from each other (Kruskal-Wallis one-way
 260 ANOVA on ranks, $H=680$, $P<0.001$). Clock rates calculated alone for the *T. brucei* pan-African clade and *T.*
 261 *congolense* savannah have similar (median *T. brucei* 2.35×10^{-7} s/s/y, *T. congolense* 1.15×10^{-6} s/s/y) but
 262 significantly faster (Kruskal-Wallis one-way ANOVA on ranks, $H>80$, $P<0.001$) rates compared to the
 263 species as a whole (**Additional Data 2**). The rates reported here are similar to those reported for other
 264 mitochondrial clocks(19), and faster than the estimated rate of trypanosome nuclear evolution based on

265 18S data ($\sim 1 \times 10^{-10}$ s/s/y)(20). Rates calculated from regions with non-edited protein coding genes (ND2 <->
266 COI) have lower substitution rates (median, *T. brucei* 1.62×10^{-7} s/s/y, *T. congolense* 4.50×10^{-7} s/s/y)
267 compared to the maxicircle as a whole. However rates for this region are faster than rates predicted for the
268 *T. cruzi* COII–ND1 region ($\sim 2 \times 10^{-8}$ s/s/y)(21).

269 The inferred tree for *T. congolense* shows three clades, deeply separated in time, corresponding to
270 the three known subgroups (**Figure 3**). The kilifi subgroup (Tck) diverged approximately 400 kya (95% HPD
271 interval 29-1200 kya), and the forest (Tcf) and savannah (Tcs) subgroups approximately 115 kya (95%
272 HPD interval 20-680 kya). Most isolates fell in the Tcs clade (**Figure 3**), which was further subdivided into
273 two clades with a divergence date of ~ 4 kya; these clades broadly comprise East (Kenya, Uganda, Ethiopia
274 and Zambia) and West (The Gambia and Burkina Faso) African isolates. Results from hierBAPS and
275 STRUCTURE analyses confirmed these results, and in addition hierBAPS resolved TRT12 from Zambia
276 and IL3578 from Burkina Faso as separate individuals (**Figure 3**); for these analyses, runs with and without
277 admixture models had the same best predicted K, except in the resolution of *T. congolense* subgroups.

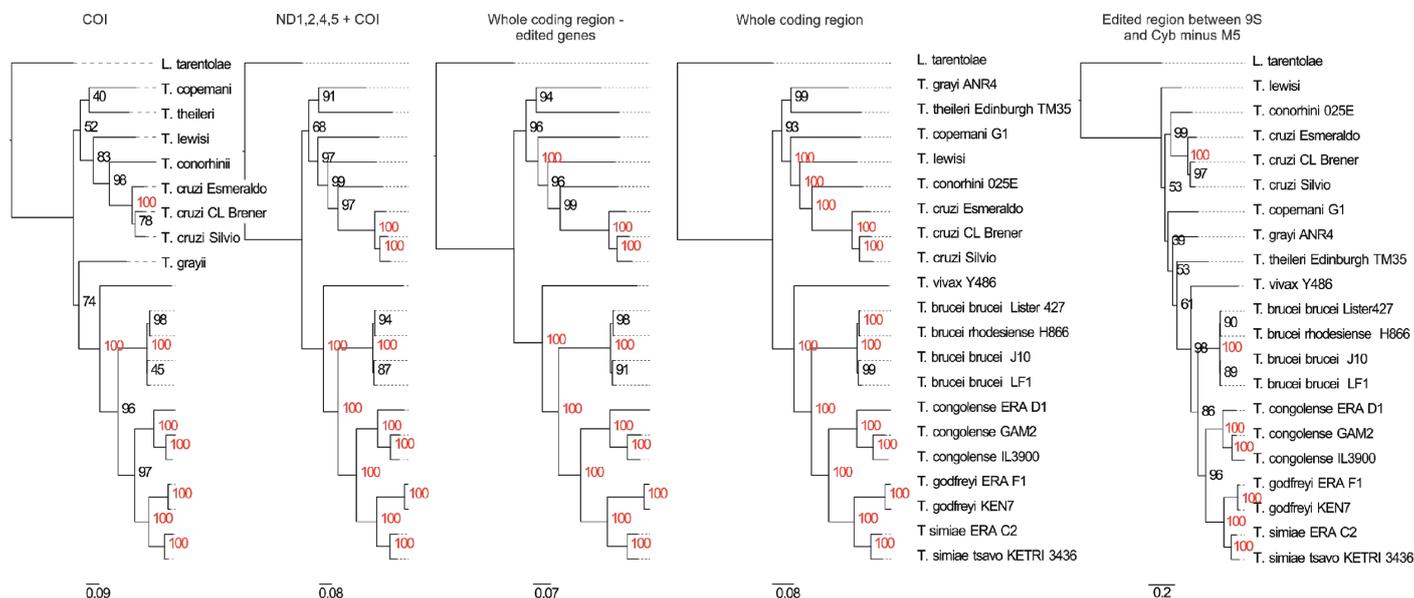
278 The *T. brucei* tree revealed two deeply separated clades, which diverged ~ 108 kya (95% HPD
279 interval 8-325 kya); these two clades also emerged from STRUCTURE analysis (**Figure 3**). Separate from
280 the main clade of *T. brucei* subspecies isolates is a clade containing isolates previously identified as
281 belonging to kiboko (J10, 927) and sindo (LF1) groups on the basis of maxicircle restriction fragment length
282 polymorphisms(22) and COI haplotypes (23); interestingly this clade also contains Old and New World *T.*
283 *equiperdum* isolates (BoTat, Dodola 943, TeAPND1). The main *T. brucei* spp. clade is further subdivided,
284 separating a group of East African isolates containing Lister 427 from a pan-African group ~ 23 kya. A group
285 of largely East African isolates emerge from the pan-African clade ~ 3.5 kya, and this group is also present
286 in STRUCTURE and hierBAPS analyses (**Figure 3**).

287 The inferred trees for *T. brucei* give insights into the evolution of *T. equiperdum* and *T. b. gambiense*.
288 For *T. equiperdum* only three isolates with whole coding regions were included in time-resolved analysis,
289 as sequences of other isolates are either incomplete or have large deletions (STIB841 and STIB842 are
290 truncated in ND5; STIB818 has lost most of the maxicircle (**Figure 1**)). The positions of these other isolates
291 were inferred from maximum likelihood trees on shared sequence (partial 12S, partial COI to partial ND5,
292 ~ 4.7 kbp of sequence), which indicate that STIB841 and STIB842 group with the other *T. equiperdum*
293 isolates with full length sequences, while STIB818 is placed separately (**Additional Figure 6**). The

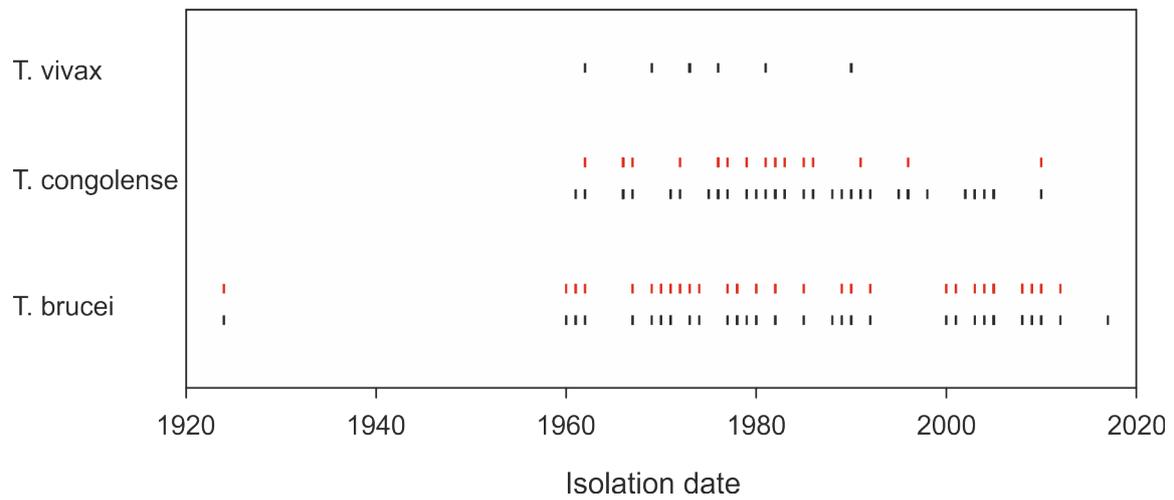
294 divergence date for this main group of *T. equiperdum* from *T. b. brucei* is around 5,000 ya (*T. b. brucei*
295 J10 and *T. equiperdum* BoTat 5,190 ya (95% HPD interval 360-15,800 ya); *T. b. brucei* 927 and *T.*
296 *equiperdum* Dodola 943/TeAPND1 4,310 ya (95% HPD interval 232-9,810 ya)). However, the position of
297 STIB818 suggested by maximum likelihood trees could support a much older origin for this lineage.

298 Origins of both Type 1 and 2 *T. b. gambiense* (*Tbg1* and *Tbg2*) can be estimated from the inferred
299 trees: *Tbg1* 3,240 ya (95% HPD interval 222-8,380 ya); *Tbg2* ~1,000 ya (95% HPD interval 80-3,020 ya).
300 This result is consistent with a published emergence date of 750-9,500 years ago for *Tbg1*, based on
301 estimated mutation rates and the observed number of mutations accumulated per genome in this asexual
302 lineage(24).

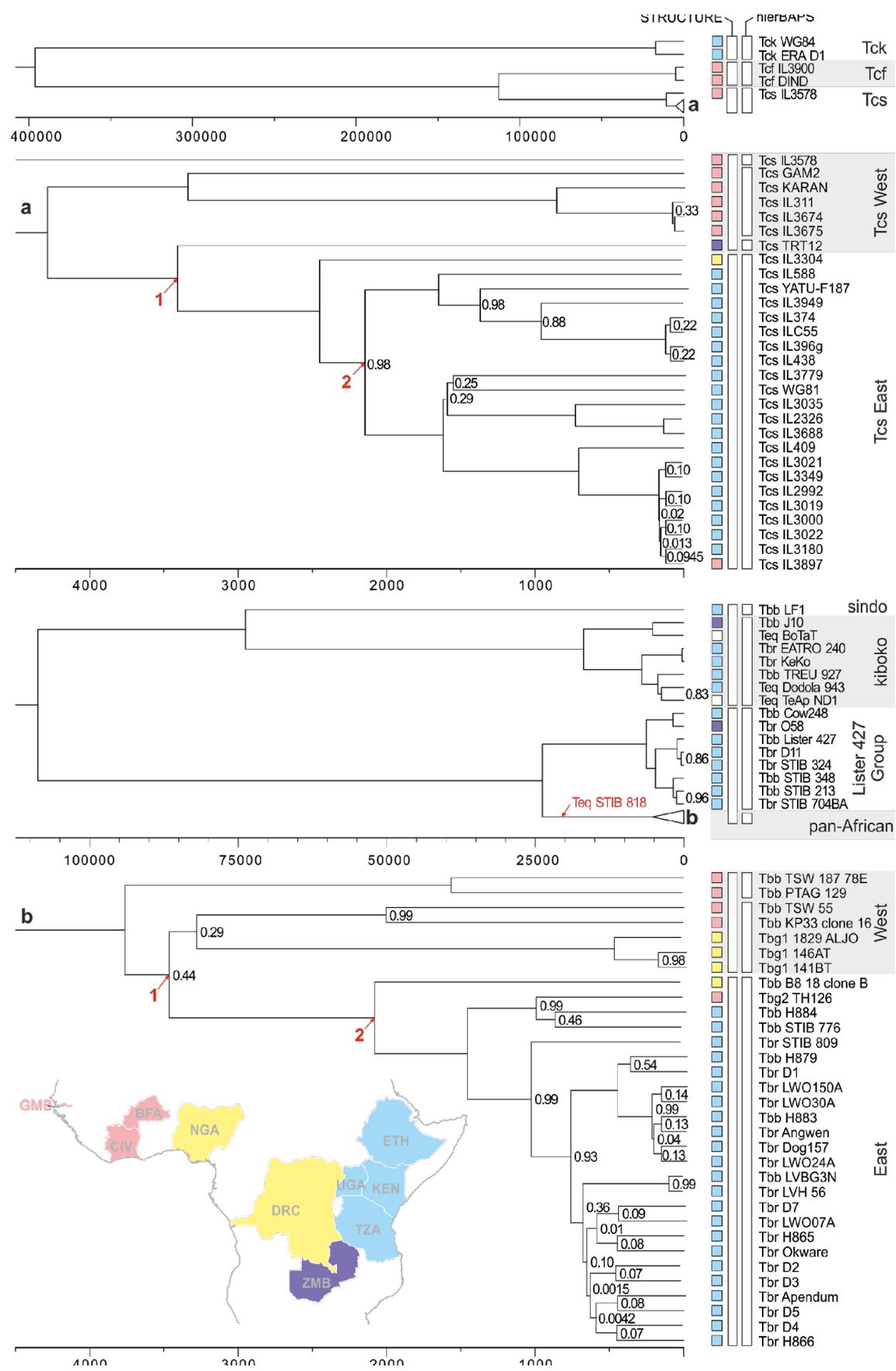
303 Despite the difference in clock rates for each species, the time resolved phylogenies indicate that
304 both *T. congolense* and *T. brucei* underwent a major divergence events simultaneously (**Figure 3: 1**, West
305 to Central Africa, *T. congolense* 3,410 ya, 95% HPD interval 294-13,400 ya; *T. brucei* 3,430 ya, 95% HPD
306 interval 236-9,610 ya; **2**, expansion into East Africa *T. congolense* 2160 ya, 95% HPD interval 160-6,800
307 ya; *T. brucei* 2,070 ya, 95% HPD interval 141-6,210 ya), posing intriguing questions about the causes.
308 Possible reasons include the major climatic changes that have affected the African continent in the past few
309 thousand years, including the gradual desiccation of the Sahara desert (~3.0 kya) and the closure of the
310 Dahomey gap (~4.5 kya)(25), overlaid by movements of wild animals, humans and their livestock in
311 response to ecological changes.



312 **Additional Figure 5. A comparison of maximum likelihood trees inferred from different regions of**
 313 **the trypanosome mitochondrial DNA.** In general using more sequence contributes to higher bootstrap
 314 support for the inferred maximum likelihood topology. If individual genes are used, confidence for the
 315 deepest branches is reduced, and topological variances are observed. Collections of non-edited genes
 316 have a consistent topology but fail to resolve well within species. Use of the WCR, with or without edited
 317 genes, provides better supported trees. If edited genes alone are used, structure within species is well
 318 supported, but multispecies relationships are poorly resolved.

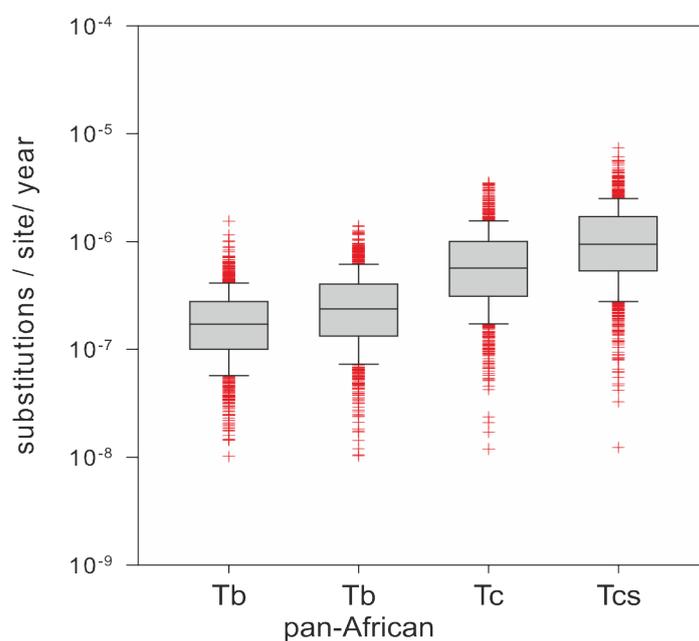


321 **Additional Data 1. Distribution of isolation dates used for inferring time-resolved phylogeny.** A
 322 spread of isolation dates for strains of *T. brucei*, *T. congolense*, *T. equiperdum* and *T. vivax* are shown.
 323 Complete coding regions used for time resolved phylogeny are indicated in red. Multiple complete coding
 324 regions were obtained for *T. vivax* but clocks were not calculated based on the limited range of isolation
 325 dates.

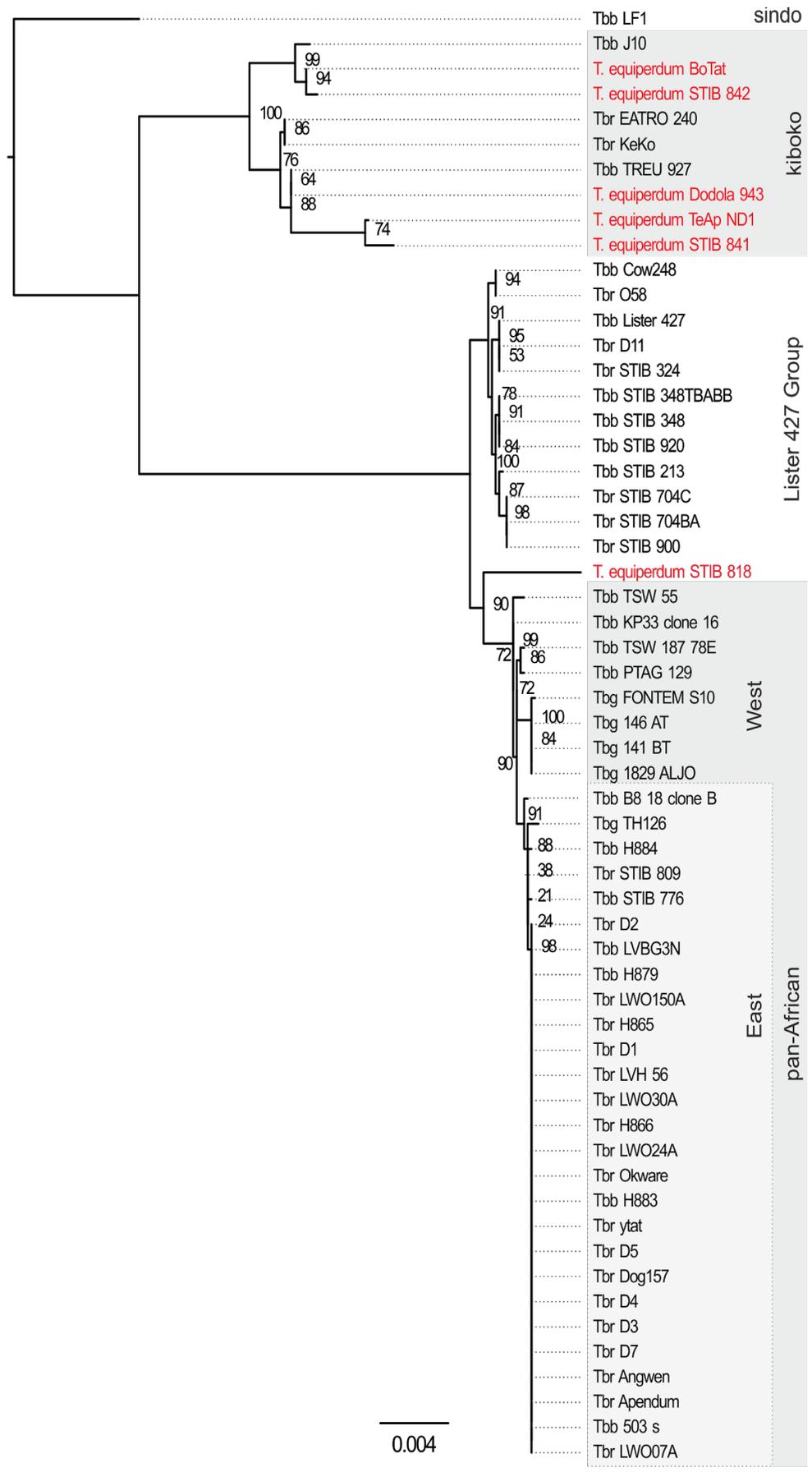


328 **Figure 3. Time resolved phylogenies of *T. congolense* and *T. brucei*.** The savannah (**a**) and pan-
329 African (**b**) clades are expanded below their respective trees. The coloured boxes correspond to countries
330 of origin on the map of Africa (inset). STRUCTURE and hierBAPS groups are indicated by the white boxes.
331 Timelines are in years before present and node values are posterior probabilities <1. Arrowed nodes 1 and
332 2 are discussed in the text. The putative position of *T. equiperdum* STIB818 inferred from an independent
333 ML tree (**Additional Figure 6**) is indicated in b.

Comparison of sampled clock rates



335 **Additional Data 2. Distribution of clock rates sampled from BEAST2 for trypanosome species and**
 336 **subgroups.** Nine hundred evenly sampled clock rates from timed phylogeny runs are shown for *T. brucei*
 337 (Tb) and the pan-African subgroup, as well for *T. congolense* (Tc) and the savannah subgroup (Tcs). Box
 338 and whisker plots show the 10th, 25th, 75th and 90th percentiles with the midline representing the median.



341 **Additional Figure 6. Inferred polyphyly of *T. equiperdum*.** *T. equiperdum* isolates in red font. A
342 maximum likelihood tree inferred from the shared common sequence from the reference sequences of
343 STIB818, STIB841 and STIB842, which have incomplete coding region sequences, and BoTat, Dodola 943
344 and TeAp ND1, which all have complete maxicircle coding regions. Node values represent bootstrap
345 support.

Discussion

The trypanosome maxicircle presents a complex evolutionary system, with several discrete mechanisms bringing about sequence change. Synteny in the coding region is largely conserved, except for segmental gene deletions in some lineages (*T. equiperdum*, New World *T. vivax*, *T. godfreyi*, *T. theileri*), leading to presumed loss of function. Whether complete maxicircle loss, as seen in *T. evansi*, is the inevitable fate for maxicircles with small deletions remains unclear, but the fact that several maxicircles with deletions have been found suggests that maxicircle loss may not happen as a single event.

Assembled maxicircles have low GC content in both coding and non-coding regions. For non-edited genes the remaining permissive mutational space for further GC loss is small, particularly in salivarians, as few mutations would be synonymous and protein composition might already be compromised. The true extent of GC loss in edited genes is cryptic as additional coding information comes from the minicircle gRNAs, however the declines in edited gene length indicate that genes are more extensively edited in salivarian compared to non-salivarian trypanosomes. Given the recent emergence of the salivarian clade this would conflict with the idea that RNA editing is a primitive kinetoplastid feature that is always “on the way out”(1,6).

The observed base composition biases in mtDNA could be driven by the loss of recombination, as GC loss is a feature commonly observed in non-recombining populations(26). Alternatively base composition biases could reflect the metabolic cost and availability of these nucleobases(27). In non-edited genes a strand-specific bias for poly-T as well as selection for AU codons suggests that selection acting at the level of the transcript, such as for translational efficiency or against transcript cost, influences the evolution of these sequences. The rRNA genes have low GC content but as they are not translated are not expected to share the same codon selection pressures. The untranscribed non-coding region has the lowest GC content, but wide variations in the size even within the same species, suggesting that it is not being streamlined for a reduced cost.

The Salivaria appear as a distinct group in the analyses presented, sharing properties of increased T:C ratio in their non-edited genes and shorter edited genes compared to non-salivarian trypanosomes. This disjunction suggests that the salivarians have undergone a period of evolutionary change, perhaps

373 associated with their adaptation to transmission via the salivary route in tsetse. Unfortunately there are no
374 intermediate taxa to sample. Although *T. grayi* is also transmitted by tsetse, this is by the posterior rather
375 than salivarian route, and *T. grayi* is not a close relative of salivarian trypanosomes in phylogenetic
376 analyses(8,9,28).

377 Despite the different evolutionary processes at work in the maxicircle coding region, our analyses
378 demonstrate that it is a useful tool for phylogenetic analysis and a good molecular clock within a species.
379 From population genetics analyses and the consistent phylogeny of isolates using different portions of the
380 maxicircle, recombination appears rare or is restricted to very closely related sequences in the salivarian
381 trypanosomes *T. congolense* and *T. brucei*. This contrasts with *T. cruzi*, where evidence for recombination
382 and heteroplasmy has been presented(29,30). Our analysis of *T. brucei* and *T. congolense* suggests that
383 the maxicircle can be used to probe the recent history and distribution of a species using isolation dates
384 without other assumptions. The dates inferred for the *T. brucei* group fit well with estimations for the date of
385 emergence of the human pathogen *T. b. gambiense* Type 1, previously calculated as 750-9,500 ya, based
386 on estimated mutation rates and the observed number of mutations accumulated per genome in this
387 asexual lineage(24). These dates fit with the development of settled agriculture and burgeoning centres of
388 population in West Africa in the past 10,000 years that favoured the evolution of parasites adapted to
389 human to human transmission. As shown by previous studies(31), *T. equiperdum* is polyphyletic. A new
390 finding here was the emergence of one clade of *T. equiperdum* from the divergent group of *T. b. brucei*
391 associated with wild animal-tsetse transmission cycles in East Africa, referred to as kiboko/sindo group(22).
392 This puts a new perspective on the evolution of *T. equiperdum* from *T. b. brucei*, with an estimated
393 emergence date of *T. equiperdum* of ~5,000 ya. The kiboko/sindo clade itself is estimated to have diverged
394 from the main *T. brucei* clade >108,000 ya.

395 Besides the kiboko/sindo clade, a small group of East African *T. b. brucei* and *T. b. rhodesiense*
396 isolates was clearly separate from the majority of *T. brucei* isolates from sub-Saharan Africa. The human
397 pathogen *T. b. rhodesiense* is characterised by a unique gene, the *SRA* (serum resistance associated)
398 gene, which confers the trait of human infectivity(32). Two major sequence variants of this gene have been
399 identified that distinguish *T. b. rhodesiense* isolates from northern and southern East Africa. Here, *T. b.*
400 *rhodesiense* LVH 56 (northern *SRA* variant) and *T. b. rhodesiense* 058 (southern *SRA* variant) were found
401 in separate clades in the tree (Figure 3), with an estimated divergence time of ~23,000 ya, placing the

emergence of the *SRA* gene, and consequently *T. b. rhodesiense*, earlier than this date.

The dated phylogeny also has ramifications for the evolution of the important livestock pathogen, *T. congolense*. Of the three subgroups of *T. congolense*, kilifi is the earliest diverging, estimated to have split from the forest and savannah subgroups ~400,000 ya. *T. congolense* savannah and forest subgroups diverged more recently around 115,000 ya. The more extensive sampling of the savannah subgroup provides evidence of a split between East and West African isolates about 4,000 ya. The position of TRT12 from Zambia on a long branch at the edge of the East African clade suggests that further subdivisions may emerge with more sampling of the savannah subgroup throughout its geographical range, as already suggested in other studies(33).

From the difference in calculated clock rates for *T. brucei* and *T. congolense*, it is clear that clock rates vary between trypanosome species, which fits with the observation that the rate of nuclear evolution in salivarians is 7-10 fold higher than non-salivarians(20). It is also reasonable to assume that there is rate variation across the coding region. At present, the geological timescale of salivarian divergence is poorly constrained, with most published studies based on a single calibration of divergence between New World and Old World trypanosomes at 100 Mya, coincident with the splitting of Africa and South America. However, it is difficult to exclude the possibility that trypanosome exchange between continents might have occurred much more recently(11), which would have a major impact on inferred rates. The alternative of using isolation dates may provide a useful complementary approach for investigating more recent divergences within trypanosome clades. We show here that isolation dates can be used to explore events in the recent history of a species, and infer ages which fit well with historical evidence (*T. equiperdum*, *T. b. gambiense*). Rate calculations for *T. brucei* and *T. congolense* from different sets of isolation dates are in strong agreement for geographically shared events in recent history, and could be tested further by future analysis of *T. vivax*. Future analyses of deeper trypanosome evolution must address assumptions on how rates are calculated, how rate varies between species, and our confidence in using geological events for speciation barriers. This would put us in a better position to understand the evolution of the salivarian trypanosomes and the genus as a whole, and infer accurate dates for the origins of the group.

428 Conclusions

429 The mtDNA data we present represents a new resource for experimental and evolutionary analyses of
430 trypanosome phylogeny, molecular evolution and function. Despite the different evolutionary processes at
431 work in the maxicircle coding region, our analyses demonstrate that it is a useful tool for phylogenetic
432 analysis and a good molecular clock within a species. Molecular clock analyses yielded a timescale for
433 trypanosome evolution congruent with major biogeographical events in Africa and revealed the recent
434 emergence of *Trypanosoma brucei gambiense* and *T. equiperdum*, major human and animal pathogens.

435 Methods

436 Genomic DNA extraction

437 High molecular weight DNA for genome sequencing was purified from axenically-grown procyclic
438 trypanosomes using a Blood and cell culture kit (Qiagen) and a modification of the manufacturer's yeast
439 cell protocol. Briefly, approximately 5×10^8 trypanosomes were pelleted by centrifugation, washed once
440 with PBS and resuspended in 5 ml lysis buffer containing proteinase and RNAase as per the
441 manufacturer's protocol. Following 1 hour incubation at 50 °C, lysates were centrifuged at 5000 rpm for 5
442 minutes at room temperature in a microfuge to pellet debris before the supernatant was applied to a
443 Genomic-tip 100/G column (Qiagen). Subsequent processing followed the manufacturer's protocol; after
444 isopropanol precipitation, DNA was resuspended in 200 µl 10 mM Tris, 0.1 mM EDTA, pH 8 and stored at 4
445 °C.

446 Sequence data

447 Long read data was obtained on a PacBio Sequel II System, using 1 or 2 cells of a 4 reaction SMRT Cell
448 1M v2 plate per sample, and prepared using the SMRTbell® Template Prep Kit 1.0. Short read sequence
449 data was obtained on an Illumina NovaSeq producing approximately 20 Gbp of 150 bp paired end reads
450 per sample. Reference sequences were obtained from NCBI Refseq database(34). Data in the Sequence
451 Read Archive(35) was recovered using the SRA Toolkit(36). All the data sources used for assembly are

452 listed in **Additional Table 1**.

453 **Assembly**

454 **Illumina assembly:** For species with reference maxicircles, Illumina sequence data was searched using
455 Magic-BLAST v1.4.0(37) for aligning reads. Pooled reads were then assembled using SPAdes v3.13.1(38)
456 and maxicircle contigs identified by BLAST v2.2.31+(39). Where assembly yielded multiple maxicircle
457 contigs, those of >1000 bp were oriented and scaffolded using MeDuSa v1.6(40). For species without close
458 references (e.g. *T. grayi*) NOVOplasty v3.3(41) was used to extend the COI seed region of a related
459 species to yield partial maxicircle sequences.

460 **PacBio assembly:** Long PacBio reads spanning the maxicircle were identified by BLAST; an example
461 maxicircle spanning read is shown in **Additional Figure 1**. These reads were then used to fish additional
462 sequences from the read pool. Reads were then corrected using Canu v1.8(13) and split to less than 12
463 kbp before being assembled with Flye v2.5(12). Illumina read data, where available, were used to polish
464 Flye assembled maxicircles.

465 **Sanger assembly:** Maxicircle reads were identified by BLAST against a reference maxicircle and
466 assembled using CAP3(42).

467 **Assembly assessment:** Reads were aligned to the assembled maxicircle sequences using BWA MEM
468 v0.7.17(43) and visualised in Tablet v1.19.09.03(44). Dot plots were produced in Flexidot v1.06(45).

469 ***T. theileri*:** The *T. theileri* maxicircle sequence was identified from the assembled contig pool by BLAST.

470 **Gene annotation**

471 For partially assembled maxicircles BLAST was used to recover individual non-edited genes. Complete
472 coding regions were prepared using BLAST to crop assembled maxicircles between 12S rRNA and ND5
473 genes. Sequences were aligned using MAFFT v7.427 (46) (coding sequence; G-INS-i, PAM 200, k=2,
474 individual genes); short sequences were discarded. An approximation of gene boundaries for edited genes
475 was made by aligning an annotated coding region of *T. vivax* Y486 to the coding region alignment and
476 cropping sub-alignments on the basis of these annotations. For non-edited protein coding genes, gene
477 boundaries could be determined by predicted open reading frame using translation table 4.

478 Phylogenetics

479 Maximum likelihood trees were inferred using IQ-Tree v1.6.12(47), using ModelFinder(48) to find the best-
480 fitting nucleotide substitution models. Parameters from these runs were used to inform a time-resolved
481 phylogeny using BEAST2 v2.6.1(49). A birth-death model using isolation dates as tip dates and a strict
482 molecular clock was used for both *T. congolense* and *T. brucei*, on the basis of marginal likelihood
483 estimation using the BEAST2 path sampling and ModelTest packages. Each run was sampled every 100
484 iterations over a chain length of 10,000,000 with the first 10% discarded as burn-in; analyses were
485 examined in Tracer v1.7.1(50). Treeannotator v1.10.4 was used to extract a consensus tree from the
486 sampled population, and trees were visualised in FigTree v1.4.2.

487 Population genetics

488 Clustering of isolates into groups was performed by first extracting variable positions from aligned coding
489 regions using SNP-sites(51). Appropriately formatted files were then prepared using PGDSpider v2.1.1.5
490 (52) for later use in hierBAPS(53) and STRUCTURE(54). Job runs in STRUCTURE used 10 iterations of
491 admixture and no admixture models between K 1-8, with 5,000 generations of burn-in and 5,000 sampled,
492 and assumed the maxicircle as a haploid allele. STRUCTURE HARVESTER v0.6.94(54,55) was then used
493 to determine K. Runs in hierBAPS used 4 levels and 20 initial clusters and were run until convergence.

494 Statistics

495 Comparison of sequence properties between species used a representative from each species and species
496 subgroup. For comparing clock rates, 900 evenly sampled clock rates after a 10% burn-in on a MCMC
497 chain length of 10,000,000 were used for the basis of analysis. Where gene or sequence properties are
498 compared, tests for normality (Shapiro-Wilk) and equal variance were first applied to determine an
499 appropriate test of variance (normal, one-way ANOVA; non-normal, Kruskal-Wallis one-way ANOVA on
500 ranks). All correlations used the Pearson correlation coefficient (ρ).

501 Abbreviations

502 The following abbreviations are used for mtDNA genes. **Ribosomal genes**, 12S, 12S rRNA; 9S, 9S rRNA;
503 RPS12, ribosomal protein 12. **Respiratory complex genes**, ND(1-9), NADH dehydrogenase subunits (1-
504 9); CYB, apo-cytochrome b; CO(I-III), cytochrome oxidase subunits (I-III); A6, ATPase subunit 6. **Genes of**
505 **unknown function**, M(2, 5), Maxicircle Unidentified Reading Frame (MURF, 2, 5); CR(3, 4), C-rich
506 Reading frame (3, 4).

507 Declarations

508 Ethics approval

509 Not applicable

510 Consent for publication

511 All authors approved the final version of the MS.

512 Availability of data and materials

513 Public data used to assemble sequences in this study, and the assembled sequences are freely available,
514 and listed by run experiment, Bioproject, and Genbank accession in **Additional Table 1**. Assembled
515 sequences used for our analyses are collected together in **Additional files 1**.

516 **Additional files 1**

517 .zip, compressed archive

518 Assembled maxicircle FASTA formatted sequence data.

519 Competing interests

520 No competing interests to declare.

521 Funding

522 This study was funded by the UK Biotechnology and Biological Sciences Research Council project grant
523 BB/R016437/1. T.A.W. is supported by a Royal Society University Research Fellowship.

524 Authors' contributions

525 C. K. conceived the study, and all authors contributed to the study design, C. K. analysed the data. W.G.
526 provided materials. C.K. wrote the paper with contributions from all authors. All authors agreed on the final
527 manuscript.

528 Acknowledgements

529 This work was funded by the UK Biotechnology and Biological Sciences Research Council

530 References

- 531 1. Lukes J, Hashimi H, Zíková A. Unexplained complexity of the mitochondrial genome and transcriptome
532 in kinetoplastid flagellates. *Curr Genet*. 2005 Nov;48(5):277–99.
- 533 2. Jensen RE, Englund PT. Network news: the replication of kinetoplast DNA. *Annu Rev Microbiol*.
534 2012;66:473–91.
- 535 3. Benne R, Van den Burg J, Brakenhoff JP, Sloof P, Van Boom JH, Tromp MC. Major transcript of the
536 frameshifted *coxII* gene from trypanosome mitochondria contains four nucleotides that are not encoded
537 in the DNA. *Cell*. 1986 Sep 12;46(6):819–26.
- 538 4. Sturm NR, Simpson L. Partially edited mRNAs for cytochrome b and subunit III of cytochrome oxidase
539 from *leishmania tarentolae* mitochondria: RNA editing intermediates. *Cell*. 1990. Vol. 61, p. 871–8.
540 [http://dx.doi.org/10.1016/0092-8674\(90\)90197-m](http://dx.doi.org/10.1016/0092-8674(90)90197-m)
- 541 5. Stuart K, Panigrahi AK. RNA editing: complexity and complications. *Mol Microbiol*. 2002
542 Aug;45(3):591–6.
- 543 6. Simpson L, Maslov DA. Ancient origin of RNA editing in kinetoplastid protozoa. *Curr Opin Genet Dev*.
544 1994 Dec;4(6):887–94.
- 545 7. Stevens JR, Noyes HA, Dover GA, Gibson WC. The ancient and divergent origins of the human

- 546 pathogenic trypanosomes, *Trypanosoma brucei* and *T. cruzi* Parasitology. 1999. 18;107–16
 547 <http://dx.doi.org/10.1017/s0031182098003473>
- 548 8. Hamilton PB, Stevens JR, Gaunt MW, Gidley J, Gibson WC. Trypanosomes are monophyletic:
 549 evidence from genes for glyceraldehyde phosphate dehydrogenase and small subunit ribosomal RNA.
 550 Int J Parasitol. 2004 Nov;34(12):1393–404.
- 551 9. Stevens JR, Gibson WC. The evolution of pathogenic trypanosomes. Cad Saude Publica. 1999
 552 Oct;15(4):673–84.
- 553 10. Lewis MD, Llewellyn MS, Yeo M, Acosta N, Gaunt MW, Miles MA. Recent, independent and
 554 anthropogenic origins of *Trypanosoma cruzi* hybrids. PLoS Negl Trop Dis. 2011 Oct;5(10):e1363.
- 555 11. Hamilton PB, Adams ER, Njiokou F, Gibson WC, Cuny G, Herder S. Phylogenetic analysis reveals the
 556 presence of the *Trypanosoma cruzi* clade in African terrestrial mammals. Infect Genet Evol. 2009
 557 Jan;9(1):81–6.
- 558 12. Kolmogorov M, Yuan J, Lin Y, Pevzner PA. Assembly of long, error-prone reads using repeat graphs.
 559 Nat Biotechnol. 2019 May;37(5):540–6.
- 560 13. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and accurate
 561 long-read assembly via adaptive k-mer weighting and repeat separation. Genome Research. 2017.27;
 562 722–36. <http://dx.doi.org/10.1101/gr.215087.116>
- 563 14. Greif G, Rodriguez M, Reyna-Bello A, Robello C, Alvarez-Valin F. Kinetoplast adaptations in American
 564 strains from *Trypanosoma vivax*. Mutat Res. 2015 Mar;773:69–82.
- 565 15. Westenberger SJ, Cerqueira GC, El-Sayed NM, Zingales B, Campbell DA, Sturm NR. *Trypanosoma*
 566 *cruzi* mitochondrial maxicircles display species- and strain-specific variation and a conserved element
 567 in the non-coding region. BMC Genomics. 2006 Mar 22;7:60.
- 568 16. Lai D-H, Hashimi H, Lun Z-R, Ayala FJ, Lukes J. Adaptations of *Trypanosoma brucei* to gradual loss of
 569 kinetoplast DNA: *Trypanosoma equiperdum* and *Trypanosoma evansi* are petite mutants of *T. brucei*.
 570 Proc Natl Acad Sci U S A. 2008 Feb 12;105(6):1999–2004.
- 571 17. Hoare CA. The trypanosomes of mammals: a zoological monograph. Wiley-Blackwell; 1972. 749 p.
- 572 18. Kaufer A, Barratt J, Stark D, Ellis J. The complete coding region of the maxicircle as a superior
 573 phylogenetic marker for exploring evolutionary relationships between members of the Leishmaniinae.
 574 Infect Genet Evol. 2019 Jun;70:90–100.
- 575 19. Molak M, Ho SYW. Prolonged decay of molecular rate estimates for metazoan mitochondrial DNA.
 576 PeerJ. 2015 Mar 5;3:e821.
- 577 20. Stevens J, Rambaut A. Evolutionary rate differences in trypanosomes. Infect Genet Evol. 2001
 578 Dec;1(2):143–50.
- 579 21. Machado CA, Ayala FJ. Nucleotide sequences provide evidence of genetic exchange among distantly
 580 related lineages of *Trypanosoma cruzi*. Proc Natl Acad Sci U S A. 2001 Jun 19;98(13):7396–401.
- 581 22. Gibson W, Borst P, Fase-Fowler F. Further analysis of intraspecific variation in *Trypanosoma brucei*
 582 using restriction site polymorphisms in the maxi-circle of kinetoplast DNA. Mol Biochem Parasitol. 1985
 583 Apr;15(1):21–36.
- 584 23. Balmer O, Beadell JS, Gibson W, Caccone A. Phylogeography and taxonomy of *Trypanosoma brucei*.
 585 PLoS Negl Trop Dis. 2011 Feb 8;5(2):e961.
- 586 24. Weir W, Capewell P, Foth B, Clucas C, Pountain A, Steketee P, et al. Population genomics reveals the
 587 origin and asexual evolution of human infective trypanosomes. Elife. 2016 Jan 26;5:e11473.

- 588 25. Demenou BB, Doucet J-L, Hardy OJ. History of the fragmentation of the African rain forest in the
589 Dahomey Gap: insight from the demographic history of *Terminalia superba*. *Heredity* . 2018
590 Jun;120(6):547–61.
- 591 26. Lynch M. *The Origins of Genome Architecture*. Sinauer Associates Incorporated; 2007. p. 494
- 592 27. Seward EA, Kelly S. Dietary nitrogen alters codon bias and genome composition in parasitic
593 microorganisms. *Genome Biol.* 2016 Nov 15;17(1):226.
- 594 28. Kelly S, Ivens A, Manna PT, Gibson W, Field MC. A draft genome for the African crocodylian
595 trypanosome *Trypanosoma grayi*. *Scientific Data*. 2014. Vol1 <http://dx.doi.org/10.1038/sdata.2014.24>
- 596 29. Barnabé C, Brenière SF. Scarce events of mitochondrial introgression in *Trypanosoma cruzi*: new case
597 with a Bolivian strain. *Infect Genet Evol.* 2012 Dec;12(8):1879–83.
- 598 30. Messenger LA, Llewellyn MS, Bhattacharyya T, Franzén O, Lewis MD, Ramírez JD, et al. Multiple
599 mitochondrial introgression events and heteroplasmy in *trypanosoma cruzi* revealed by maxicircle
600 MLST and next generation sequencing. *PLoS Negl Trop Dis.* 2012 Apr 10;6(4):e1584.
- 601 31. Cuypers B, Van den Broeck F, Van Reet N, Meehan CJ, Cauchard J, Wilkes JM, et al. Genome-Wide
602 SNP Analysis Reveals Distinct Origins of *Trypanosoma evansi* and *Trypanosoma equiperdum*.
603 *Genome Biol Evol.* 2017 Aug 1;9(8):1990–7.
- 604 32. De Greef C, Imberechts H, Matthyssens G, Van Meirvenne N, Hamers R. A gene expressed only in
605 serum-resistant variants of *Trypanosoma brucei rhodesiense*. *Mol Biochem Parasitol.* 1989
606 Sep;36(2):169–76.
- 607 33. Tihon E, Imamura H, Dujardin J-C, Van Den Abbeele J, Van den Broeck F. Discovery and genomic
608 analyses of hybridization between divergent lineages of *Trypanosoma congolense*, causative agent of
609 Animal African Trypanosomiasis. *Mol Ecol.* 2017 Dec;26(23):6524–38.
- 610 34. O’Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, et al. Reference sequence
611 (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic
612 Acids Res.* 2016 Jan 4;44(D1):D733–45.
- 613 35. Leinonen R, Sugawara H, Shumway M, International Nucleotide Sequence Database Collaboration.
614 The sequence read archive. *Nucleic Acids Res.* 2011 Jan;39(Database issue):D19–21.
- 615 36. SRA Toolkit Development Team. <http://ncbi.github.io/sra-tools/>
- 616 37. Boratyn GM, Thierry-Mieg J, Thierry-Mieg D, Busby B, Madden TL. Magic-BLAST, an accurate RNA-
617 seq aligner for long and short reads. *BMC Bioinformatics.* 2019 Jul 25;20(1):405.
- 618 38. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, et al. SPAdes: a new genome
619 assembly algorithm and its applications to single-cell sequencing. *J Comput Biol.* 2012 May;19(5):455–
620 77.
- 621 39. Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. Basic local alignment search tool. *J Mol Biol.*
622 1990 Oct 5;215(3):403–10.
- 623 40. Bosi E, Donati B, Galardini M, Brunetti S, Sagot M-F, Lió P, et al. MeDuSa: a multi-draft based
624 scaffolder. *Bioinformatics.* 2015 Aug 1;31(15):2443–51.
- 625 41. Dierckxsens N, Mardulyn P, Smits G. NOVOPlasty: de novo assembly of organelle genomes from whole
626 genome data. *Nucleic Acids Research.* 2016. <http://dx.doi.org/10.1093/nar/gkw955>
- 627 42. Huang X. CAP3: A DNA Sequence Assembly Program *Genome Research.* 1999. 9:868–77.
628 <http://dx.doi.org/10.1101/gr.9.9.868>
- 629 43. Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler transform.

- 630 Bioinformatics. 2009 Jul 15;25(14):1754–60.
- 631 44. Milne I, Stephen G, Bayer M, Cock PJA, Pritchard L, Cardle L, et al. Using Tablet for visual exploration
632 of second-generation sequencing data. *Brief Bioinform.* 2013 Mar;14(2):193–202.
- 633 45. Seibt KM, Schmidt T, Heitkam T. FlexiDot: highly customizable, ambiguity-aware dotplots for visual
634 sequence analyses. *Bioinformatics.* 2018 Oct 15;34(20):3575–7.
- 635 46. Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7: improvements in
636 performance and usability. *Mol Biol Evol.* 2013 Apr;30(4):772–80.
- 637 47. Nguyen LT, Schmidt HA, Von Haeseler A, Minh BQ. IQ-TREE: A fast and effective stochastic algorithm
638 for estimating maximum-likelihood phylogenies. *Mol Biol Evol.* 2015;32(1):268–74.
- 639 48. Kalyaanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jeremiin LS. ModelFinder: fast model
640 selection for accurate phylogenetic estimates. *Nat Methods.* 2017 Jun;14(6):587–9.
- 641 49. Bouckaert R, Vaughan TG, Barido-Sottani J, Duchêne S, Fourment M, Gavryushkina A, et al. BEAST
642 2.5: An advanced software platform for Bayesian evolutionary analysis. *PLoS Comput Biol.* 2019
643 Apr;15(4):e1006650.
- 644 50. Rambaut A, Drummond AJ, Xie D, Baele G, Suchard MA. Posterior Summarization in Bayesian
645 Phylogenetics Using Tracer 1.7 *Systematic Biology.* 2018. 67:901–4.
646 <http://dx.doi.org/10.1093/sysbio/syy032>
- 647 51. Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, et al. SNP-sites: rapid efficient
648 extraction of SNPs from multi-FASTA alignments <http://dx.doi.org/10.1101/038190>
- 649 52. Lischer HEL, Excoffier L. PGDSpider: an automated data conversion tool for connecting population
650 genetics and genomics programs. *Bioinformatics.* 2012 Jan 15;28(2):298–9.
- 651 53. Cheng L, Connor TR, Sirén J, Aanensen DM, Corander J. Hierarchical and spatially explicit clustering
652 of DNA sequences with BAPS software. *Mol Biol Evol.* 2013 May;30(5):1224–8.
- 653 54. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus genotype
654 data. *Genetics.* 2000 Jun;155(2):945–59.
- 655 55. Earl DA, vonHoldt BM. STRUCTURE HARVESTER: a website and program for visualizing
656 STRUCTURE output and implementing the Evanno method. *Conservation Genetics Resources.* 2012.
657 4:359–61. <http://dx.doi.org/10.1007/s12686-011-9548-7>
- 658