

Deep kinetoplast genome analyses result in a novel molecular assay for detecting *Trypanosoma brucei gambiense*-specific minicircles

Manon Geerts¹, Zihao Chen², Nicolas Bebronne¹, Nicholas J. Savill², Achim Schnauffer², Philippe Büscher¹, Nick Van Reet^{1,†} and Frederik Van den Broeck^{1,3,*}

¹Department of Biomedical Sciences, Institute of Tropical Medicine, 2000 Antwerp, Belgium, ²Institute of Immunology and Infection Research, University of Edinburgh, Edinburgh EH9 3FL, UK and ³Department of Microbiology, Immunology and Transplantation, Rega Institute for Medical Research, Katholieke Universiteit Leuven, 3000 Leuven, Belgium

Received April 07, 2022; Revised September 28, 2022; Editorial Decision October 05, 2022; Accepted October 06, 2022

ABSTRACT

The World Health Organization targeted *Trypanosoma brucei gambiense* (*Tbg*) human African trypanosomiasis for elimination of transmission by 2030. Sensitive molecular markers that specifically detect *Tbg* type 1 (*Tbg1*) parasites will be important tools to assist in reaching this goal. We aim at improving molecular diagnosis of *Tbg1* infections by targeting the abundant mitochondrial minicircles within the kinetoplast of these parasites. Using Next-Generation Sequencing of total cellular DNA extracts, we assembled and annotated the kinetoplast genome and investigated minicircle sequence diversity in 38 animal- and human-infective trypanosome strains. Computational analyses recognized a total of 241 Minicircle Sequence Classes as *Tbg1*-specific, of which three were shared by the 18 studied *Tbg1* strains. We developed a minicircle-based assay that is applicable on animals and as specific as the *TgsGP*-based assay, the current golden standard for molecular detection of *Tbg1*. The median copy number of the targeted minicircle was equal to eight, suggesting our minicircle-based assay may be used for the sensitive detection of *Tbg1* parasites. Annotation of the targeted minicircle sequence indicated that it encodes genes essential for the survival of the parasite and will thus likely be preserved in natural *Tbg1* populations, the latter ensuring the reliability of our novel diagnostic assay.

INTRODUCTION

Human African trypanosomiasis (HAT), also known as sleeping sickness, is a vector-borne disease caused by two *Trypanosoma brucei* (*Tb*) subspecies and transmitted by tsetse flies. *Trypanosoma brucei rhodesiense* (*Tbr*) causes acute infections in East Africa, whereas *Trypanosoma brucei gambiense* (*Tbg*) causes chronic infections in West and Central Africa (1). *Trypanosoma brucei rhodesiense* and *Tbg* type I (*Tbg1*) are defined by the presence of truncated variant surface glycoprotein (VSG) genes, respectively the serum-resistance-associated (SRA) gene (2) and the *Tbg*-specific glycoprotein (*TgsGP*) gene (3) that play a role in infectivity to humans (4). Unlike *Tbr*, *Tbg1* parasites are genetically homogeneous and form a monophyletic group (5,6). They are responsible for the vast majority of the HAT-cases (85% of the 663 newly reported HAT-cases in 2020) (7). Infections by a third group of human-infective trypanosomes - *Tbg* type II (*Tbg2*) that lack the SRA and *TgsGP* genes - are extremely rare (8). Like the non-human-infective *T. b. brucei* (*Tbb*), *Tbg2* is genetically and phenotypically highly diverse (8).

The WHO targeted *Tbg*-HAT (gHAT) for elimination as a public health problem by 2020 and for global elimination of transmission (EOT) to humans (i.e. zero reported cases) by 2030 (9). Elimination of gHAT as a public health problem has been reached in several countries and HAT foci, and recently Togo and Côte d'Ivoire have been validated as such by WHO (10,11). However, EOT remains challenging due to imperfect diagnostics and the risk of re-emergence from asymptomatic human infections and/or a possible animal reservoir (12–14). Current serological tests based on the *Tbg*-specific LiTat 1.3 and LiTat 1.5 VSG antigens (15–20) still require validation in different animal

*To whom correspondence should be addressed. Tel: +32 32476586; Email: fvandenbroeck@gmail.com

†Co-senior authors: Frederik Van den Broeck and Nick Van Reet.

Present address: Manon Geerts, Fish Eco-Evo-Devo and Conservation, KU Leuven, 3000 Leuven, Belgium; Directorate Taxonomy and Phylogeny, Royal Belgian Institute for Natural Sciences, 1000 Brussels, Belgium.

species. The golden standard for molecular detection of *Tbg1* involves the single-copy *TgsGP* gene (21–23), but the analytical sensitivity of these tests is limited because they target a hemizygous single-copy gene (24). Among the molecular tests, several studies have proposed various genotyping techniques such as isoenzymes, ribosomal genes, VSGs, SNPs and microsatellites (25–32), but these have some disadvantages such as the requirement of multiple PCR reactions or high amounts of input material, without necessarily increasing the sensitivity of *Tbg1* detection. To this end, alternative genetic markers are much needed to reliably and with improved sensitivity demonstrate *Tbg1* infection in humans and in animals, including the tsetse fly vector.

Previous reports indicated that *Tbg1*-specific minicircle sequences exist in the kinetoplast DNA (33,34). The kinetoplast DNA (kDNA), unique to the single mitochondrion of unicellular flagellates of the order Kinetoplastida, is a giant network of dozens of homogenous maxicircles (20–30 kb) interlaced with hundreds to thousands of heterogeneous minicircles (0.5–2.5 kb) (35). Maxicircles are homologous to the mitochondrial genome of other eukaryotes and encode components of the respiratory chain complexes and the mitoribosome. Minicircles generally consist of a ~100 bp conserved sequence region that contains hyper conserved sequences named Conserved Sequence Blocks (CSBs), and a variable region including genes encoding guide RNAs (gRNAs) that are responsible for directing post-transcriptional modification of the maxicircle-encoded messenger RNAs (35). A complete assembly and annotation of the kinetoplast genome of a lab-adapted *Tbb* strain identified 391 minicircle classes, encoding ~1000 gRNA genes (36). Analysis of kDNA minicircles is already used for *Leishmania* detection and differentiation (37,38) and allows subtypes of *Trypanosoma evansi* to be distinguished (39). It has also been proposed for sensitive and specific detection of *Tbg1* infection in humans (34) and animals (33). However, the exact nature of these sequences is unknown and their use for *Tbg1* diagnosis is cumbersome, requiring both nested PCR and DNA hybridization.

Recent advances in the development of bioinformatic tools now facilitate the assembly and annotation of the structurally complex kinetoplast genome (36,40), allowing us to investigate minicircle sequence diversity in tens to hundreds of samples simultaneously (41). Therefore, we used next-generation sequencing of total cellular DNA extracts as a strategy to investigate minicircle sequence diversity in 38 animal- and human-infective trypanosome strains from diverse geographical origins. Following a series of computational analyses, we identify minicircles that were present in all *Tbg1* strains and absent in all *Tbb*, *Tbg2* and *Tbr* strains. Using a newly developed quantitative PCR (qPCR), we demonstrate that our minicircle-based assay reliably identifies *Tbg1* within the Trypanozoon subgenus. Furthermore, we show that the copy number of the targeted minicircle was on average 4-fold higher compared to the hemizygous single-copy *TgsGP* gene, suggesting that our minicircle-based assay may be used for sensitive detection of *Tbg1* infections in humans and animals.

MATERIALS AND METHODS

Ethics statement

Expansion of bloodstream form trypanosome populations in mice received approval from the Animal Ethics Committee of the Institute of Tropical Medicine (DPU2017-1).

DNA extraction and sequencing

To cover a wide geographical range, we included a total of 18 *Tbg1* strains isolated from humans between 1952 and 2008 from Cameroon ($n = 2$), the Democratic Republic of the Congo ($n = 10$), the Congo Brazzaville ($n = 2$), Côte d'Ivoire ($n = 3$) and South Sudan ($n = 1$) (Supplementary Table S1). For comparative purposes, we also included nine *Tbb*, five *Tbg2* and six *Tbr* strains (Supplementary Table S1).

All 38 strains were propagated as bloodstream form populations in OF1 mice (Charles River, Belgium) and purified from the infected mouse blood via DEAE ion exchange chromatography (42). Purified trypanosomes were sedimented by centrifugation ($17\,000 \times g$, 10 min at 4°C). DNA of $50\ \mu\text{l}$ pure trypanosome sediment was extracted using the standard phenol:chloroform method (43), aliquoted at $1\ \text{ng}/\mu\text{l}$ and stored at -20°C . The concentration of extracted DNA was determined using a Qubit 4 Fluorometer (Invitrogen by Thermo Fisher Scientific). Paired-end 150 bp sequences were generated using the DNA nanoball sequencing technology (DNBSEQ™) at the Beijing Genomics Institute (BGI) in Hongkong, China.

Genomic analyses

Paired-end reads were aligned against the *Tbb* TREU927 v4.6 reference genome (available on <https://tritrypdb.org>) using SMALT v0.7.6 (<https://www.sanger.ac.uk/tool/smalt-0/>). A hash index of the reference genome was built using k -mer words of length 13 that were sampled every other position in the genome. Mapping was done using an exhaustive search for alignments with a minimum identity threshold ($-y$) of 80% and a maximum insert size for paired reads of 1,500 bp.

Variant calling and filtering was performed using the Genome Analysis Toolkit (GATK) v4.1.4.1 (44). First, reads were assigned to a single read-group with *AddOrReplaceReadGroups* and duplicated reads were marked with *MarkDuplicates*. Variants were then called for each strain separately with *HaplotypeCaller* using default parameters. The resulting gVCF files for all strains were combined with *CombineGVCFs* to allow joint genotyping with *GenotypeGVCFs*. Single Nucleotide Polymorphisms (SNPs) were extracted from the resulting VCF file with *SelectVariants* and filtered with *VariantFiltration* using the following parameters: $\text{QUAL} < 500$, $\text{DP} < 5$, $\text{QD} < 2.0$, $\text{FS} > 60.0$, $\text{MQ} < 40.0$, $\text{MQRankSum} < -12.5$, $\text{ReadPosRankSum} < -8.0$, $-\text{cluster_window_size } 10$ and $-\text{cluster_size } 3$. Finally, we used BCFtools v1.10.2 (45) to extract bi-allelic SNP sites that were called in all *Tb* strains. Using the resulting set of genome-wide SNPs, we

reconstructed a phylogenetic network with SplitsTree v4 (46) to infer the ancestral relationship among the 38 *Tb* strains.

The species identity of the *Tbg1* and *Tbr* strains was confirmed *in silico* by investigating the presence of the TgsGP and SRA genes, respectively. To this end, MEGAHIT v1.2.9 (47) was used for *de novo* assembly of genomes of all 38 *Tb* strains using default parameters. The presence of TgsGP and SRA genes in the assembled contigs was then confirmed through a local BLAST search (48) using publicly available nucleotide sequences of TgsGP (21) (NCBI accession number: FN555993) and SRA (2) (NCBI accession number: Z37159) with the following parameters: minimum 90% identity, minimum e-value of 0.0001 and minimum alignment length of 500bp.

Assembly of the kinetoplast genome

Reads that did not align to the *Tbb* TREU927 nuclear reference genome were extracted using SAMtools v1.9 (45) and converted to FASTQ format using GATK *SamToFastq*. These unmapped reads were aligned against the 23 kb maxicircle sequence of *Tbb* Lister 427 (GenBank accession id M94286) using SMALT and the same parameters as described above except that the hash index was built with *k*-mer words of length six and the reads were mapped with a minimum identity threshold ($-y$) of 90% and a maximum insert size ($-i$) of 500. SNP calling was done using GATK as described above, and we extracted only those SNPs that passed the quality criteria (see above) and that were present within the maxicircle coding region (1.3–16.3 kb). Similar to the analyses above for the nuclear genome, a phylogenetic network analysis was done using maxicircle coding SNPs with SplitsTree.

Reads that did not align to the maxicircle sequence of *Tbb* Lister 427 were extracted from the alignment file as described above, and used for the assembly of minicircle contigs. Before assembly, sequence reads were trimmed for high quality with fastp v0.20.0 (49) using the following parameters: allow for a maximum of 10% of bases per read that have a phred-scaled base quality below 30, trim bases at either end of the read when their phred-scaled quality is below 30, move a sliding window of 10 bp from front to tail and cut the read once the average phred-scaled base quality drops below 30, and only retain reads with a minimum and maximum length of 100 and 155 bp after trimming, respectively. Using KOMICS v1.8 *assemble* (40), trimmed reads were used for *de novo* assembly of contigs using a *k*-mer list of 99, 109 and 119, and putative minicircle contigs were extracted based on the highly conserved sequence block 3 (CSB3) dodecamer (GGGGTTG[G/A]TGTA) (36,50). Five minicircles had a slight variation of CSB3 (one with GGGGGTGGTGTA found in *Tbg2* strain FEO and four with GGGGTTAGTGTA found in *Tbg1* strains 15BT-relapse, OUSOU, N̄DIMI and ROUPO-VAVOUA-80-MURAZ-14). These minicircles were also retained. KOMICS *circularize* was then used to identify circular minicircle contigs by searching for overlapping fragments at either end of each contig; when a minicircle contig was classified as circular, the overlapping fragment at the start of the contig was removed. Using KOMICS *polish*, all circular

minicircle contigs were oriented by putting the conserved sequence block 1 (CSB1) (GGGCGT[T/G]C) (50) at the start of each contig. One minicircle had a slight variation of CSB1 (GGGCGTGT found in *Tbg2* strain MSUS-CI-78-TSW-157), which was specified to KOMICS to allow proper reorientation for this minicircle. Finally, KOMICS *polish* was also used to remove duplicate sequences, which was achieved by extracting the representative sequences (cluster centroids) of all clusters identified at 97% identity with VSEARCH v2.14.2 (51).

The quality of the minicircle assembly was assessed by re-aligning the unmapped reads to the assembled minicircles using SMALT. Before mapping, we first extended the circularized minicircle sequences by copying the last 150 bp at the start of each sequence to minimize the number of clipped reads at either end of the assembled minicircles, using a custom python script implemented in KOMICS. Following mapping with SMALT with exhaustive search ($-x$) and a percent identity of 97% ($-y$), we have calculated the following metrics using a bash script implemented in KOMICS: number of reads, number of mapped reads, number of properly paired reads, number of reads with mapping quality ≥ 20 , number of CSB3-containing reads, number of mapped CSB3-containing reads and number of perfectly aligned CSB3-containing reads (i.e. alignments without any insertions or deletions). The proportion of perfect alignments of CSB3-containing reads serves as a proxy for the total number of minicircles that were initially present within the DNA sample. All metrics were processed and visualized using the R function *msc.quality* as implemented in the R package rKOMICS (41). In addition, the quality of the assembly was further verified by calling SNPs with BCFtools *mpileup/call*, retaining only SNPs with QUAL ≥ 60 and DP ≥ 30 , and assuming that high-quality assemblies should yield relatively low number of homozygous SNPs.

Finally, using the rKOMICS function *msc.depth*, minicircle copy numbers (MCN) were estimated as the median read depth per minicircle contig divided by the median genome-wide read depth times two (assuming diploidy in all *Tb* strains).

Identification of minicircle sequences unique to *Tbg1*

The diversity and similarity of minicircle sequences within all *Tb* strains were examined with the R package rKOMICS (41). Following visual inspection of length distributions using *msc.length*, we used *preprocess* to retain minicircle sequences that had the expected length (800–1200 bp) and that were successfully circularized. Retained sequences of all samples were concatenated into a single FASTA file and clustered into Minicircle Sequence Classes (MSCs) based on a minimum percent identity (MPI) of 70, 80, 90 and 95–100 with VSEARCH. In order to choose an appropriate MPI for downstream analyses, we processed VSEARCH clustering results with *msc.uc* and inspected - at each MPI - the number of MSCs, the number of perfect alignments and the number of 2-nt and 3-nt gaps. In addition, VSEARCH clustering results were stored into a matrix using the rKOMICS function *msc.matrix*, which records the presence (1) or absence (0) of all MSCs (rows)

for each strain (columns). This matrix was subsequently used to document the number of MSCs per strain with *msc.richness*, to calculate the proportion of MSCs shared between the different *Tb* subspecies with *msc.similarity* and to investigate the ancestry among all *Tb* strains with *msc.pca*. Finally, we used the rKOMICS function *msc.subset* to find MSCs that were present in the *Tbg1* strains and absent in the strains belonging to the other *Tb* subspecies.

Development of minicircle-based quantitative PCR assays

For each of the common *Tbg1*-specific MSC, we extracted the assembled minicircle sequences for all 18 *Tbg1* strains from the alignment using the rKOMICS function *msc.seq*, generated consensus sequences with Jalview v2.11.1.4 (52) and designed primers and probes (Table 1) with the RealTimeDesign qPCR assay design software (LGC, Biosearch Technologies). Probes targeting *Tbg1*-specific MSCs were modified with a FAM dye label at the 5' end and paired with BHQ-1 *plus* at the 3' end.

A 20 μ l reaction mixture contained 1X PerfeCTa qPCR Toughmix (Quantabio), 100 nM of each primer (LGC, Biosearch Technologies), 300 μ M of each probe (LGC, Biosearch Technologies) and 5 μ l of template DNA. The thermal cycling profile consisted of an initial denaturation step at 95°C for 10 min followed by 40 cycles at 95°C for 15 s and 60°C for 1 min. qPCR was conducted on a Q-qPCR Instrument (Quantabio), and detection of the quantification cycle (Cq) was calculated using the Q-qPCR Instrument Software v1.0.2 with the automatic threshold enabled.

Each qPCR targeting *Tbg1*-specific MSCs was multiplexed with the *Trypanozoon*-specific q18S-assay targeting the multi-copy 18S rRNA gene (53), the *Trypanozoon*-specific qGPI-PLC-assay targeting the single-copy glycosylphosphatidylinositol-specific phospholipase C (GPI-PLC) gene (53) and a *Tbg1*-specific qTgsGP-assay, designed to target the single-copy TgsGP gene and avoid amplification of TgsGP-like genes (21). The multi-copy 18S rRNA gene was used as an internal standard for the sensitive detection of *Trypanozoon* DNA. The single-copy GPI-PLC gene was used as an internal standard to determine if sufficient *Trypanozoon* DNA was present to detect a single-copy sequence (54), and for the calculation of relative copy numbers (RCN) of each *Tbg1*-specific MSC (see below). The single-copy TgsGP gene was used as a golden standard for the specific detection of *Tbg1* DNA. The q18S-assay contains a CAL Fluor Orange 56 dye labeled probe paired with BHQ-1 *plus* (53). The qGPI-PLC-assay contains a CAL Fluor Red 610 dye labeled probe paired with BHQ-2 *plus* (53). The qTgsGP-assay contains a Quasar 670 dye labeled probe paired with BHQ-2 *plus*.

The qPCR efficiency and analytical sensitivity were calculated for each qPCR targeting *Tbg1*-specific MSCs in simplex and in quadruplex format. This was done using phenol:chloroform extracted DNA (see above) of two *Tbg1* strains. From these DNA extracts, ten-fold serial dilutions in DEPC-treated water, ranging from 100 pg/ μ l to 1 fg/ μ l, were prepared. Each qPCR was run in quadruplicate for

each DNA dilution. A reaction was considered positive if at least three out of four replicates were positive.

The specificity of the quadruplex qPCR assays was assessed with the phenol:chloroform extracted genomic DNA of 34 *Tbb*, 49 *Tbg1*, 7 *Tbg2*, 15 *Tbr*, 2 *T. equiperdum* and 5 *T. evansi* strains (Supplementary Table S2). The 49 *Tbg1* samples originated from Burkina Faso (2), Cameroon (2), Côte d'Ivoire (6), Congo Brazzaville (3), Democratic Republic of the Congo (35) and South Sudan (1). Note that 26 strains from the Democratic Republic of the Congo were sampled within the context of a treatment outcome study in Mbuji-Mayi (55), with 14 strains sampled from seven patients before (sample name include 'BT') and after ('AT') treatment. Quadruplex qPCR assays were run in duplicate for each DNA extract. The specificity of the assays was further assessed on DNA prepared from man, cattle, dog, goat, horse, sheep and tsetse (*Glossina fuscipes* from Kwamouth, Democratic Republic of the Congo, 2018), all known hosts of *Tbg1* (12).

Relative copy numbers (RCN) of each target were calculated using the Δ Cq-method with qGPI-PLC as reference. This was done by subtracting the Cq-values obtained for each *Tbg1*-specific MSC from the Cq-value obtained for qGPI-PLC. The resulting Δ Cq-value were averaged between replicates and transformed ($2^{\Delta Cq}$) to yield RCNs for each target.

Annotation of minicircle sequences targeted by qPCR assays

Strain *Tbg1* 340AT (MHOM/CD/INRB/2006/21B) isolated in 2006 in Mbuji-Mayi (DRC) (55) (Supplementary Table S1) was selected for representative minicircle annotation because of its comparatively high minicircle complexity (see results). DNA extraction, sequencing and alignment of sequence reads against the *Tbb* TREU927 v4.6 reference genome was done as described above for the other 38 *Tb* isolates initially included in our study. Reads that did not align to the nuclear reference genome were used for the assembly of mitochondrial maxicircles and minicircles with KOMICS, as described above.

Due to the lack of transcriptomic data for *Tbg1* strain 340AT, edited mRNA sequences were predicted following a similar approach as in (40). Edited mRNA sequences for *Tbb* strains Lister 427, EATRO 164 and EATRO 1125 (36) were obtained from GenBank and manually corrected for changes in non-T residues based on alignment of the *Tbg1* 340AT maxicircle with the annotated *Tbb* EATRO 1125 maxicircle (36).

Guide RNA prediction and minicircle annotation were performed with python3 package for *kDNA annotation* (36). The alignments of the gRNAs encoded on the minicircles targeted by the diagnostic assay to their cognate mRNAs were carefully inspected to identify any 'non-redundant' gRNAs, i.e. gRNAs that direct editing events not covered by any other gRNAs. Furthermore, gRNAs 3' of such non-redundant gRNAs were checked for potential premature truncations by the gRNA calling algorithm by carefully examining the editing capacity of their 3' end extensions. Any gRNA genes that were confirmed to be non-redundant were considered essential, as were the minicircles that encoded them.

Table 1. Summary of the various qPCR assays used in this study

Assay		Sequence	Target	Position on target	Expected amplicon length	Reference
qMini1	F	5' TGAGGTCTGAGGTAAGTTCGAAAG 3'	mO-104	52–152	151	this paper
	R	5' TGGATTACTGGTGTCTTCTATTGATAA 3'				
qMini2	P	5' FAM-TTTTCTGGAGAAAACTGTAT-BHQ-1 plus 3'	mO_078	244–431	188	this paper
	F	5' TCTTATGACTGATTTACGAGAATA 3'				
	R	5' GACATAACAGAGGAAAGTGCTC 3'				
qMini3	P	5' FAM-TTGTGGTAAGAGTGATTTAGTAAT-BHQ-1 plus 3'	mO-078	626–810	185	this paper
	F	5' AAACCAACAGAAAAGAGATTGCTTA 3'				
	R	5' ATGGTGATAGAAGTTAGAGATGTGTAG 3'				
qTgsGP	P	5' FAM-TAGATGTAGTATAAGAATTTAAAAT-BHQ-1 plus 3'	<i>TgsGP</i>	753–840	87	this paper
	F	5' GAAGCAGTGGGACCTTAGC 3'				
	R	5' TTTGTGCTCTTGCTTGCTATTAC3'				
q18S	P	5' Quasar 670 -CTCTCCGAACACAGCAGCGACATC-BHQ-21 plus 3'	18S	679–829	150	(53)
	F	5' CGTAGTTGAACTGTGGGCCACGT 3'				
	R	5' ATGCATGACATGCGTAAAGTGAG3'				
qGPI-PLC	P	5' CAL Fluor Orange 560 TCGGACGTGTTTTGACCCACGC-BHQ-1 plus 3'	<i>GPI-PLC</i>	520–626	106	(53)
	F	5' CCCACAACCGTCTCTTTAAAC 3'				
	R	5' GGAGTCGTGCATAAGGGTATTC3'				
	P	5' CAL Fluor Red 610-ACACCACTTTGTAACCTCTGGCAGT-BHQ-1 plus 3'				

F = forward primer, R = reverse primer, P = probe sequence.

RESULTS

Genome analyses confirms the taxon identity of *Tbg1* strains

The genomes of 38 *Tb* strains were sequenced at a median 159× depth (mean = 155, min = 126, max = 178) (Supplementary Table S3). On average 86.3% of the reads (min = 81.25%, max = 92.05%) aligned to the *Tbb* TREU927 nuclear reference genome (Supplementary Table S3). Initial variant discovery with GATK identified a total of 1 558 963 SNPs across the 38 *Tb* strains. Strict quality filtering and the exclusion of multiallelic sites reduced the data set to 316 287 genome-wide bi-allelic SNPs, of which 310 701 SNPs (98.23%) were located within the 11 megabase chromosomes. In addition, joint genotyping identified a total of 150 SNPs within the maxicircle coding region.

Phylogenetic analyses based on genome-wide SNPs and SNPs from the maxicircle coding region confirmed that all 18 *Tbg1* parasites clustered together in a monophyletic group, a prerequisite for downstream analyses that aim at identifying *Tbg1*-specific minicircles (Supplementary Figure S1). Using a local BLAST search of assembled contigs, we also confirmed the presence of the *Tbg1*-specific *TgsGP* gene for the 18 *Tbg1* strains, and the *Tbr*-specific *SRA* gene for the six *Tbr* strains. These two genes were absent for the remaining nine *Tbb* strains and five *Tbg2* strains (Supplementary Table S4).

Assembly and circularization of mitochondrial minicircles

Mitochondrial minicircles were *de novo* assembled, circularized and reoriented for each of the 38 *Tb* strains using the Python package KOMICS. A total of 9076 minicircle contigs were assembled across all 38 *Tb* strains, of which 7156 (78.85%) were successfully circularized (Supplementary Table S5). The length of the majority of

circularized minicircles (7,111 contigs, 99.37%) showed a unimodal distribution around ~1,000 bp (Supplementary Figure S2), which is comparable to the minicircle length found in *Tbb* (36). To verify whether the assembly process was impacted by excluding reads mapped on the nuclear genome or the maxicircle, we performed a BLAST search of the assembled minicircles against the reference genome, retaining only alignments with a minimum 90% identity, an e-value of 0.001 and a minimum length of 150 bp (the read length). This identified a total of 13 contigs that most likely originated from the nuclear genome, as they showed a length between 2608 and 10876 bp (which is much larger than the expected 1000 bp length of minicircles). The length of one contig was 419 bp, which aligned with 96% identity and a length of 151 bp against chromosome 10. It is thus possible that this contig was not fully assembled because of sequence homology with the reference genome. However, this is only 1 contig in 1409 minicircle contigs that were not successfully circularized. Hence, we believe that the impact of our filtering method has a negligible impact on the assembly process.

To validate the quality of the assembly process, sequence reads were aligned to the assembled minicircle contigs and several mapping and genotyping statistics were calculated. First, high-quality assemblies should result in relatively low numbers of homozygous SNPs when reads are aligned against the assembled contigs. Here, we identified a total of 302 homozygous SNPs within 127 minicircle contigs (Supplementary Table S5), which is only 1.4% of all assembled minicircle contigs. Only 48 homozygous SNPs were identified within 17 circularized contigs (0.2% of all circularized minicircles) (Supplementary Table S5). These results show that homozygous SNPs were found for only a fraction of the assembled contigs. Second, on average 96.99% of the sequence reads mapped in proper pairs

and 94.46% aligned with a mapping quality larger than 20 (Supplementary Table S5), indicating that the large majority of mapped reads aligned with a high quality and with the expected orientation to the minicircle assemblies. Third, we calculated the number of aligned reads containing the CSB3 12-mer as a proxy for the total number of minicircles initially present within the DNA sample. This revealed that on average 93.24% of the CSB3-containing reads aligned against the assembled minicircles, and 88.77% aligned perfectly (Supplementary Table S5), suggesting that we were able to retrieve the vast majority of the minicircles. Note that for the assembly and circularization of minicircles, we only used sequencing reads that did not align against the nuclear genome or the coding region of the mitochondrial maxicircle. Therefore, the ~7% of unaligned CSB3-containing reads may have originated from minicircles that have not been (fully) assembled or circularized, the variable region of the maxicircle or the minichromosomes.

Estimation of maxicircle and minicircle copy numbers and network size

To calculate the average number of maxicircles and minicircles per kinetoplast network, we used the coverage (i.e. median read depth) of the diploid nuclear genome. The average genome-wide coverage was 155 per diploid cell (two copies), and the average coverage per haploid sequence was equal to 78 (Supplementary Table S3). The coverage of the mitochondrial maxicircles (coding region only) and minicircles was estimated based on median read depths per 0.1 kb. Average coverage of the maxicircle was 1262 (median = 1325, min = 99, max = 2320) (Supplementary Table S6). This equaled an average copy number of 17 maxicircles (median = 18, min = 1, max = 28) per network (Supplementary Table S6), which is slightly lower compared with previous estimates of 20–50 copies per network (36,56–57). The average copy number of minicircles (MCN) ranged from 2.1 to 37.9 copies per network in *Tbg1* strains, from 0.8 to 32.7 copies per network in *Tbb* strains, from 4 to 11.9 copies per network in *Tbg2* strains and from 3.1 to 11.1 copies per network in *Tbr* strains (Supplementary Table S7). Thus, copy numbers for minicircles within each network varied substantially. The size of each kDNA network was estimated by adding up the estimated copy numbers for all minicircles. These calculations resulted in an estimated average of ~2,100 minicircles per network, ranging between 252 minicircles in the Nabe strain and 4,828 minicircles in the MSUS-CI-78-TSW-157 strain (Supplementary Table S7). These numbers are slightly lower compared to earlier estimates of 5000–10 000 minicircles per network (36,56–58), and future studies should attempt to confirm a lower amount of kDNA in *Tbg1* by other methods, such as quantitation by microscopy.

Sequence diversity and similarity of mitochondrial minicircles

Minicircle sequence diversity was examined using a clustering approach, whereby minicircles were grouped into MSCs based on a minimum percent identity. This was done on the 7111 circularized minicircle contigs of the

expected length, as these would produce the most robust alignments. At 100% identity, a total of 5883 unique MSCs were identified across the 38 *Tb* strains, leaving 1228 MSCs that are shared between two or more isolates. The number of MSCs decreased sharply with decreasing percent identities to a total of 719 MSCs at 70% identity (Supplementary Figure 3A). Regardless of the percent identity used, *Tbg1* parasites contained an average of 103 MSCs per strain (median = 107), which is on average 2.48-fold lower compared to other taxa of the *Trypanozoon* subgenus (Figure 1). Most of the *Tbg1* strains displayed a fairly similar number of MSCs, ranging between 89 and 122 MSCs (Figure 1), with the exception of LiTat 1.5 (50 MSCs) and Bosendja (76 MSCs). The composition of minicircle sequences was further investigated by quantifying the proportion of MSCs unique to *Tbg1*. At 98%–100% identity, the 18 *Tbg1* parasites did not share any MSC with the non-*Tbg1* subspecies (*Tbr*, *Tbg2* and *Tbb*) (Supplementary Figure 3B). Below 98%, there was a steady increase in the proportion of shared MSCs between *Tbg1* and non-*Tbg1* subspecies (Supplementary Figure 3B). At the 98% identity threshold, the *Tbg1* group contained 241 MSCs, none of which were found in the other *Tb* subspecies and three of which were found in all 18 *Tbg1* strains (Figure 2). These three *Tbg1*-specific MSCs were retained as candidate markers for our new molecular test.

Novel multiplex qPCRs including *Tbg1*-specific mitochondrial minicircles as target

Using the RealTimeDesign qPCR assay design software, we successfully designed three simplex qPCR assays (here-after referred to as qMini1, qMini2 and qMini3) targeting two *Tbg1*-specific MSCs (here-after referred to as mO.078 and mO.104) (Table 1). Due to limitations such as ambiguous base count, low GC percentage and low melting temperatures, we were unable to design primers and probes for the third *Tbg1*-specific MSC. The *Tbg1* specificity of the designed primers and probes was confirmed with a BLAST search (48) on TriTrypDB (<https://tritrypdb.org>), with a Primer BLAST search on NCBI (<https://www.ncbi.nlm.nih.gov>) and with the command line search tool grep in sequencing reads generated by this study. The three simplex qPCR assays qMini1, qMini2 and qMini3 were each multiplexed with q18S, qGPI-PLC and qTgsGP (Table 1) to produce three quadruplex reactions (here-after referred to as g-qPCR1, g-qPCR2 and g-qPCR3 when qMini1, qMini2 and qMini3 were included, respectively).

The efficiency and analytical sensitivity of the simplex and quadruplex qPCR assays were investigated using DNA from two *Tbg1* strains (Supplementary Table S8), one with relatively low MCNs (Nabe; average MCN = 1) and one with relatively high MCNs (LOGRA; average MCN = 8). The lower detection limit of 0.05 pg DNA was reached in both strains for qMini1, qMini2 and qMini3 in their respective simplex reactions (Supplementary Figure S4), for q18S in all three g-qPCR assays (Supplementary Figure S4), and for qMini1 and qMini3 in their respective g-qPCR assays, with the exception of qMini3 that achieved the detection limit of 0.5 pg DNA in the Nabe strain (Supplementary Figure S4). This detection limit

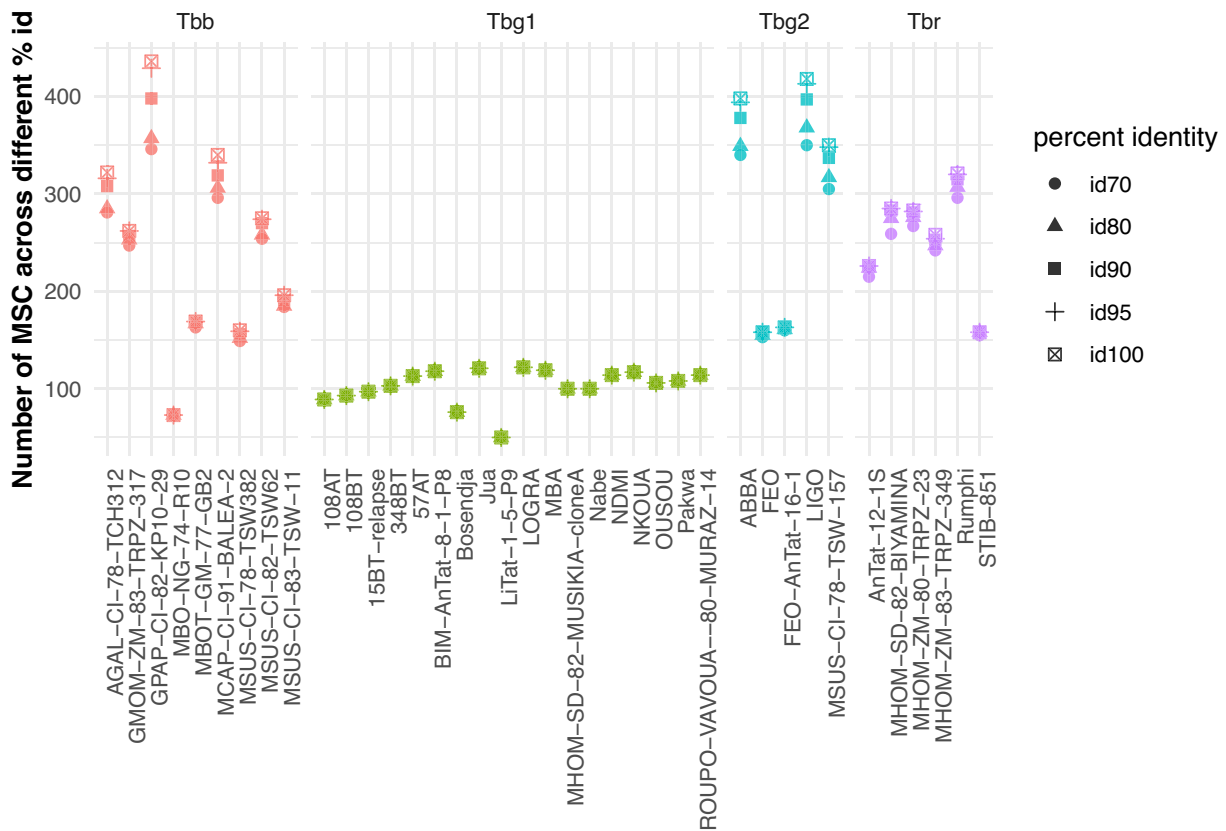


Figure 1. Number of Minicircle Sequence Classes (MSCs) for 38 *Tb* strains. Following *de novo* assembly and circularization, mitochondrial minicircle sequences were clustered into groups of sequences sharing a minimum percent identity (MSCs). Here, we summarized the number of MSCs identified within each *Tb* strain for a range of percent identities (70, 80, 90, 95 and 100). Strains were grouped as *T. b. brucei* (*Tbb*), *T. b. gambiense* type 1 (*Tbg1*), *T. b. gambiense* type 2 (*Tbg2*) and *T. b. rhodesiense* (*Tbr*).

of 0.5 pg DNA was also reached by qTgsGP and/or qGPI-PLC in all three *g*-qPCR assays with the LOGRA strain and in the *g*-qPCR2 and *g*-qPCR3 assays with the Nabe strain (Supplementary Figure S4). The analytical sensitivity of qMini2 was greatly reduced in the *g*-qPCR2 assay to 100 pg DNA with the Nabe strain and 0.5 pg DNA with the LOGRA strain (Supplementary Figure S4). The qPCR efficiency of *g*-qPCR1 and *g*-qPCR3 was estimated between 93% and 106% (Supplementary Figure S4), which is considered acceptable (<https://www.thermofisher.com/content/dam/LifeTech/global/Forms/PDF/real-time-pcr-handbook.pdf>).

The qMini3 assay displays a similar specificity as the qTgsGP assay

To assess the taxon-specificity of qMini1 and qMini3 within the *Trypanozoon* subgenus, a total of 118 DNA extracts were tested with the *g*-qPCR1 and *g*-qPCR3 assays (Supplementary Table S9). Here, qMini2 was excluded because of its low analytical sensitivity in the *g*-qPCR2 assay (see above) and because it targets the same minicircle as qMini3 (Table 1). Six of the 118 DNA extracts were excluded as they didn't react in duplicate with qGPI-PLC, which was used as an internal standard to determine if sufficient *Trypanozoon* DNA was present to detect a single-copy sequence (Supplementary Table S10). The remaining 112 DNA extracts reacted

with the *Trypanozoon*-specific assays q18S and qGPI-PLC (Supplementary Table S9). All 49 *Tbg1* DNA extracts reacted with qTgsGP, qMini1 and qMini3, with the exception of MSUS/CI/82/TSW125_KP1_cloneB (*Tbg1* isolated from a pig in Côte d'Ivoire), ALJO (*Tbg1* isolated from a human patient in DRC) and GUIWI-BOBO80-MURAZ18 (*Tbg1* isolated from a patient in Burkina Faso) that remained negative for qMini1. The qMini1 assay also showed one cross reaction with DNA extracted from the *Tbr* strain Etat 1.2 R, with a Cq-value 5× lower than that of the qGPI-PLC. The qMini3 assay showed no cross reactions or false negative results. In addition, the *g*-qPCR3 assay did not amplify DNA from human and non-human vertebrates that are known to be susceptible for *Tbg1* infection, *in casu* horse, cattle, goat, sheep, pig and dog, and DNA from *Glossina fuscipes*. Also, DNA from other livestock affecting trypanosomes like *T. congolense*, *T. theileri* and *T. vivax* was not amplified (Supplementary Table S9).

The qMini3 assay targets a minicircle with a relatively high, but variable copy number

Relative Copy Numbers (RCN) of the minicircle sequence targeted by qMini3 (RCN_{mO.078}), the target sequence of qTgsGP (RCN_{qTgsGP}) and the target sequence of q18S (RCN_{q18S}) were calculated for each of the 49 *Tbg1* DNA extracts using the ΔCq-method (Supplementary Table S11).

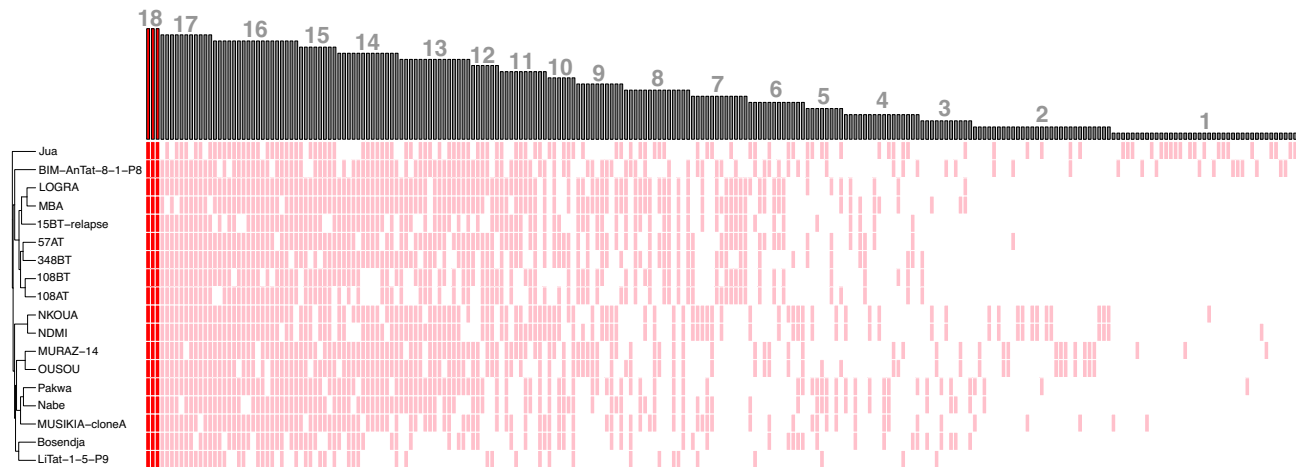


Figure 2. Minicircle sequence diversity in 18 *Tbg1* strains. Clustering at the 98% identity threshold revealed a total of 241 MSCs that were present in at least one of 18 *Tbg1* strains and absent in the *Tbb*, *Tbr* and *Tbg2* strains. Heatmap summarizes the presence/absence of these 241 MSCs (columns) in the 18 *Tbg1* strains (rows). Gray lines indicate the presence of a given MSC in a given *Tbg1* strain. Dark red lines indicate the three MSCs that were found in all 18 *Tbg1* strains. Barplot on top of the plot shows the number of strains harboring a given MSC.

The median $RCN_{mO.078}$ (7.98) was 4.29x higher than the median RCN_{qTgsGP} (1.86) and 2.06x higher than the median RCN_{q18S} (3.88) (Figure 3A). However, $RCN_{mO.078}$ also displayed a larger variation (SD = 9.01, min = 0.77, max = 40.58) compared to RCN_{qTgsGP} (SD = 0.62, min = 0.94, max = 3.47) and RCN_{q18S} (SD = 0.70, min = 1.51, max = 5.47). The $RCN_{mO.078}$ was lower compared to RCN_{qTgsGP} for 4/49 *Tbg1* strains and to RCN_{q18S} for 12/49 *Tbg1* strains. There was no association between $RCN_{mO.078}$ with the year of isolation (Pearson correlation test; $cor = -0.1925748$, $t = -1.331$, $df = 46$, P -value = 0.1897) or the geographical origin of the strain (Kruskal–Wallis test; country: chi-squared = 4.0896, $df = 5$, P -value = 0.5366; region: chi-squared = 3.0453, $df = 2$, P -value = 0.2181). High variation in $RCN_{mO.078}$ was found within one group of *Tbg1* strains isolated from humans between 2005 and 2009 in Mbuji-Mayi (Democratic Republic of the Congo), with a minimum $RCN_{mO.078}$ of 0.8 in the 186BT strain and a maximum $RCN_{mO.078}$ of 36.7 in the 93AT strain. Here, $RCN_{mO.078}$ was significantly higher in strains sampled after treatment (mean $RCN_{mO.078} = 14.9$) compared to the $RCN_{mO.078}$ in strains sampled before treatment (mean $RCN_{mO.078} = 6.9$), although this difference was not significant (Welch two sample t -test on all strains, $t = 2.0305$, $df = 11.944$, $P = 0.06519$ and paired t -test on paired strains, $t = 1.5476$, $df = 5$, $P = 0.183$). The $RCN_{mO.078}$ as calculated using the ΔCq -method was strongly associated with $MCN_{mO.078}$ as calculated using standardized read depths, with a coefficient of determination of 0.91, a slope of 0.32 and a y -intercept of 0.10 (Figure 3B).

The qMini3 assay targets a minicircle containing non-redundant guide RNA genes

Annotation of minicircles was done for strain 340AT. A local BLAST search of assembled contigs revealed the presence of the *Tbg1*-specific *TgsGP* gene, confirming that 340AT is a *Tbg1* strain (Supplementary Table S3). The mitochondrial maxicircle and minicircles were assembled

with KOMICS using sequence reads that did not align to the nuclear reference genome. This resulted in a maxicircle contig of 21 287 bp long (including the entire coding region) and a total of 143 minicircle contigs (including 129 circularized contigs). Hence, 340AT has the highest number of minicircles when compared to the 18 *Tbg1* strains (max. 132 minicircles) initially sequenced in this study (Supplementary Table S7), which was the main motivation for using the 340AT data for representative minicircle annotation.

Annotation of the 143 minicircles revealed that the minicircle targeted by qMini3 encodes four gRNA genes (Figure 4). These gRNAs are involved in editing of the maxicircle genes cytochrome c oxidase subunit 3 (gCOX3(616–656) and gCOX3(341–369)), ATPase subunit 6 (gA6(415–452)) and NADH dehydrogenase subunit 7 (gND7(847–888)) (Figure 4). As expected, the minicircles contain the semi-conserved region characterized by conserved sequence blocks CSB1, CSB2 and CSB3 (36,50), and the gRNA genes are framed by imperfect 18bp inverted repeats (36,59–61).

Next, we investigated if any of the gRNAs encoded by minicircle mO.078 are non-redundant, i.e. whether they direct the editing of sites not covered by any of the other gRNAs encoded in this strain's kDNA. Minicircles encoding only redundant gRNAs might be more prone to loss due to lack of selective pressure, which could make a diagnostic assay based on such minicircles less reliable. Our analyses showed that gRNAs gA6(415–452) and gND7(847–888) direct editing events not covered by any other gRNA (Figure 5), and are thus non-redundant. The two COX3 gRNAs are redundant with gRNAs encoded by other minicircles (results not shown).

DISCUSSION

This study presents a computational investigation of mitochondrial minicircle sequence diversity in trypanosome isolates, resulting in the development of a qPCR assay as a promising new tool for

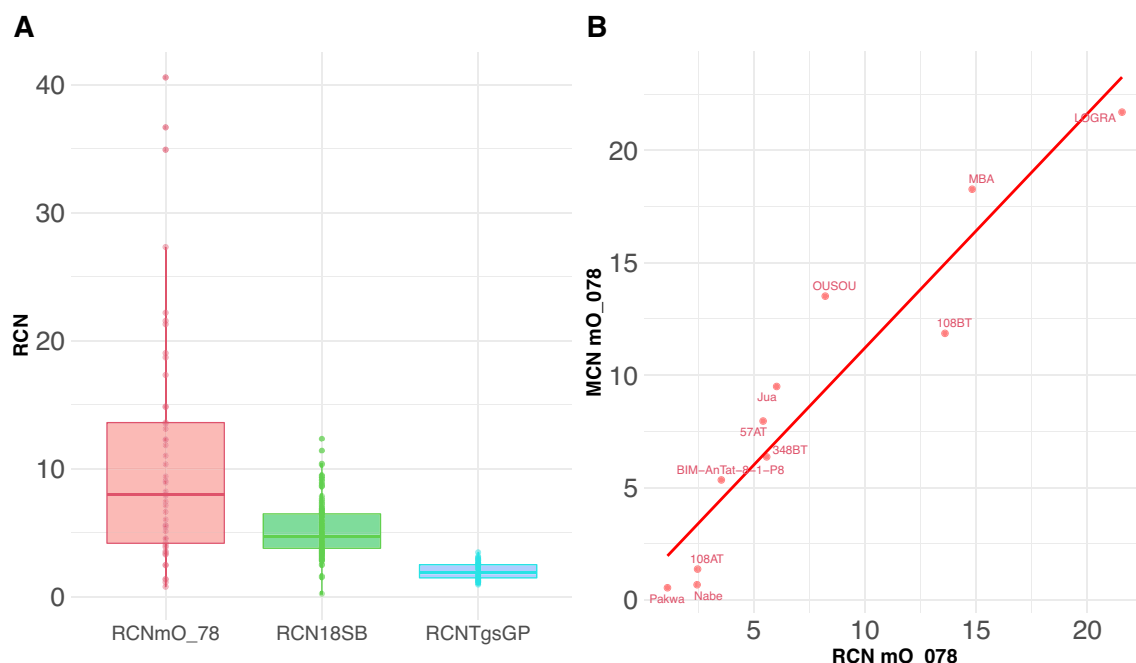


Figure 3. Relative copy numbers of the qMini3 target sequence. (A) Relative copy numbers (RCN) were estimated for the qMini3, q18S and qTgsGP target sequences using the Δ Cq-method with the qGPI-PLC as reference. Boxplots summarize the RCN estimates as calculated for 49 *Tbg1* strains. (B) Scatter plot showing the relationship between RCN (x-axis) and MCN (y-axis). To test whether RCN as calculated using the Δ Cq-method is comparable to the Minicircle Copy Numbers (MCN) as calculated using standardized read depths, RCN and MCN were calculated for the qMini3 target sequence for 11 out of 18 *Tbg1* strains. For these 11 strains, there was sufficient DNA to allow both whole genome sequencing (for MCN calculation) and a qPCR run (for RCN calculation) on the same DNA extract. The remaining seven strains were excluded here as there was not sufficient DNA left for a qPCR run following whole genome sequencing.

sensitive diagnosis of *Tbg1* infections in humans and animals.

Computational and phylogenetic analyses using genome sequencing data revealed that the 18 sequenced *Tbg1* strains are monophyletic—contrary to opinions outlined by (62)—and contain the *TgsGP* gene, an essential precondition for identifying *Tbg1*-specific minicircles. The latter was achieved by grouping *de novo* assembled and circularized minicircles into MSCs according to sequence similarity. This uncovered a variable number of MSCs within the *Trypanozoon* subgenus, with a considerably lower number of MSCs in *Tbg1* compared to the other subspecies. The comparatively lower number of MSC within *Tbg1* may be the result of its asexual evolution (5) that results in the inevitable loss of redundant minicircles due to random genetic drift (63). Hence, the monophyletic origin and the asexual evolution of *Tbg1* may explain the less complex minicircle populations in isolates for this subspecies, with conserved presence of some minicircle classes, which facilitated the identification of *Tbg1*-specific minicircles. Identification of taxon-specific minicircles may prove more challenging for trypanosomatid parasites experiencing occasional recombination (64–66), leading to heterogeneous minicircle populations as a result of biparental inheritance of mitochondrial minicircles (40,67–69).

A total of 241 MSCs were recognized as *Tbg1*-specific, of which three were shared by the 18 studied *Tbg1* strains. For two of the three *Tbg1*-specific minicircles

(mO_078 and mO_104), three molecular assays could be successfully developed and tested. While two qPCR assays were discontinued because of false-negative results or low analytical sensitivities, one qPCR assay targeting minicircle mO_078 was fully specific (i.e. no false positives or false negatives) when tested on DNA of 112 different *Trypanosoma sp.* strains. These results show that the minicircle-based assay is as specific as the assay targeting the *TgsGP* gene, the current golden standard for molecular detection of *Tbg1* parasites (70–76), confirming the taxon-specificity of some kDNA minicircles in *T. b. gambiense* (33,34) and their exploitability in molecular tests as has been described for *Leishmania* (37,38) and *T. evansi* type A and B (39,77–79). Annotation of the mO_078 minicircle demonstrated that it encodes two non-redundant gRNAs that are essential for completing the editing of ATP synthase subunit A6, a gene required for survival in both the bloodstream stage and the insect stage of *Tb* (80,81), and NADH dehydrogenase subunit 7, respectively. Hence, our results indicate that minicircle mO_078 is essential for the survival of the parasite and will most likely be preserved in all natural *Tbg1* strains, ensuring the reliability of a diagnostic assay targeting this minicircle.

Molecular analyses revealed that the *Tbg1*-specific minicircle mO_078 is a multicopy marker for the large majority of strains tested in this study, with a median copy number equal to eight and a maximum copy number of 41. Minicircle mO_078 copy number exceeded that of the multicopy *18S* gene in $\frac{3}{4}$ of the *Tbg1* strains and was as low

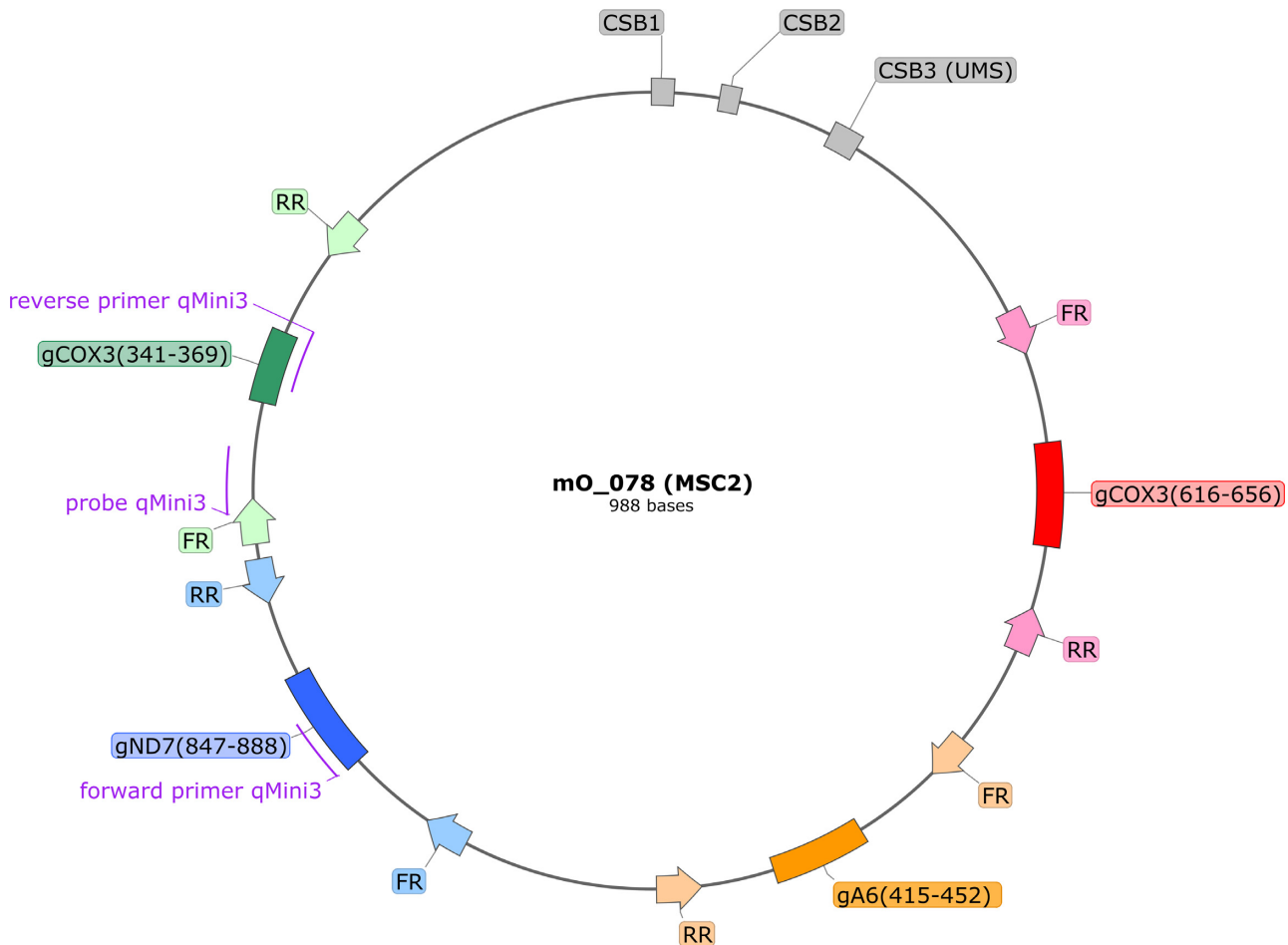


Figure 4. Annotation of minicircle mO.78 targeted by qMini3. Binding sites for the diagnostic PCR primer pair and the corresponding probe for qMini3 are indicated with purple lines. Conserved sequence blocks CSB1, CSB2 and CSB3 (or universal minicircle sequence, UMS) (gray boxes), the four encoded gRNA genes (red, orange, blue and green boxes), and the 18-bp inverted repeats (block arrows) that flank the gRNA genes are also indicated. This figure was generated with SnapGene (<http://www.snapgene.com>).

as the *TgsGP* copy number in only four of the 49 tested *Tbg1* strains. This finding confirmed the multicopy nature of mitochondrial minicircles in *Trypanosoma* and *Leishmania* parasites (36,82). However, our results also indicated that there is a relatively large variation in minicircle copy numbers across strains, suggesting that the detection limit of the minicircle-based assay may depend on the strain being investigated. The detection limit may also depend on the DNA extraction method, as DNA-extraction based on spin columns causes a random loss of small molecules like minicircles (36). Here, we avoided such biases by using the phenol-chloroform extraction method that captures all nucleic acids, although this method may be less amenable for high-throughput processing of human and animal specimens.

The high specificity and generally high copy number of minicircle mO.078 makes this a promising new marker for sensitive diagnosis of *Tbg1* infections. Specifically, within the context of reaching EOT of gHAT by 2030 (9), our minicircle-based assay may prove valuable for studying the role of an animal reservoir in the epidemiology of gHAT (12). Therefore, we propose a new multiplex qPCR assay

(g-qPCR3) that targets the *Tbg1*-specific minicircle mO.078 (serving as sensitive detection of *Tbg1*) in combination with the *Trypanozoon*-specific *18S* (serving as sensitive detection of *Tb s.l.*) and the *Tbg1*-specific *TgsGP* gene (which confirms the trypanosome is human-infective and thus of epidemiological significance). We have shown that the g-qPCR3 assay is applicable on animals, as it does not amplify DNA from other livestock affecting trypanosomes like *T. congolense*, *T. theileri* and *T. vivax*, from the tsetse fly *Glossina fuscipes* and from six livestock species that are known to be susceptible to *Tbg1* infection. Based on these data, we also trust that the specificity will not be compromised when testing blood from wild *Bovidae* and non-human primates (12), although we had no access to specimens from wild fauna to formally test this. The main limitation of our study is that the g-qPCR3 assay was tested on only one *Tbg1* strain isolated from a host other than human, *in casu* a pig from Côte d'Ivoire (83), mainly because of the scarcity of such samples. However, given that minicircle mO.078 is non-redundant and universal (computational analyses confirmed that mO.078 was present in 193 sequenced *Tbg1* strains sampled

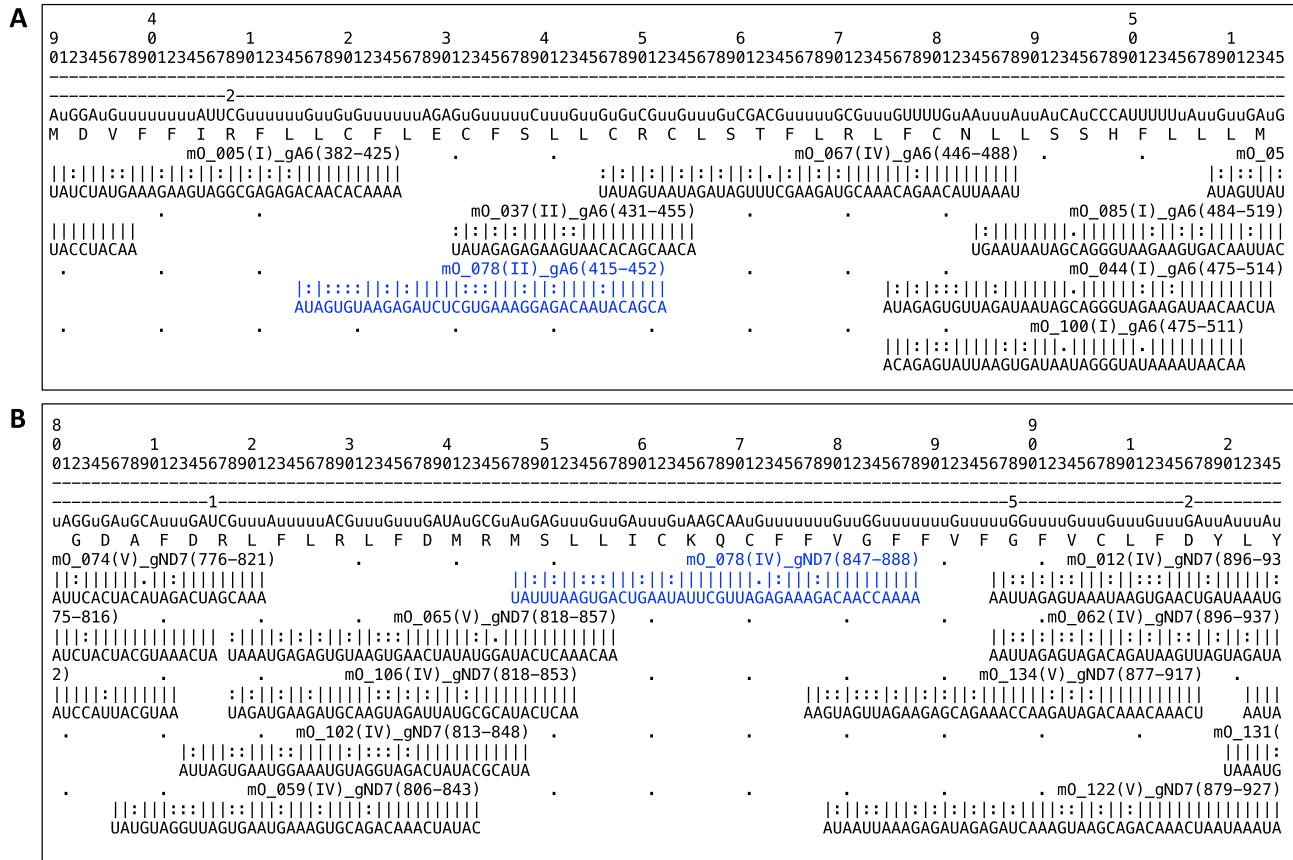


Figure 5. Alignment of the non-redundant gRNAs from the qMini3 targeted minicircle mO_078 to their target mRNAs. The alignment shows the non-redundant gRNAs in blue (A, gA6(415–452); B, gND7(847–888)) as well as the mRNA sequences immediately upstream and downstream, along with the neighbouring gRNAs encoded by other minicircles. Lines 1–3: mRNA position (hundreds, tens, ones); line 4: number of Us that have been deleted from the pre-edited mRNA at this position; line 5: edited mRNA sequence 5' to 3' (lowercase 'u's represent insertions). Note that the 'anchor' sequence at the 5' end of each gRNA cannot direct editing events; line 6: protein sequence. For each gRNA: Line 1: name (mO_name(cassette position)_mRNA(start-end of alignment on mRNA)), underscore characters denote the anchor; line 2: base-pairing: 'l': Watson-Crick basepair, ':': GU basepair, '.': mismatch basepair; line 3: gRNA sequence 3' to 5'.

from humans in 11 different countries (data not shown)), we are confident that the g-qPCR3 assay will be successful at amplifying *Tbg1* DNA from livestock specimens.

In conclusion, this study exemplifies the power of genome assembly and annotation for identifying species-specific multicopy genetic markers. We developed a minicircle-based assay that is as specific as the current golden standard for molecular detection of *Tbg1* infections, and argued that the g-qPCR3 has the diagnostic potential for assessing the importance of an animal reservoir in the epidemiology of gHAT.

DATA AVAILABILITY

Sequence reads generated within the context of this study have been deposited in the European Nucleotide Archive under accession number PRJEB49966 (<https://www.ebi.ac.uk/ena/browser/view/PRJEB49966>). The sequence of the *Tbg1*-specific minicircle mO_078 was deposited to NCBI under accession number OM238297. Scripts used for processing WGS data in this study are available at <https://github.com/GeertsManon/g-qPCR>. All graphical analyses were performed using R v4.0.3 in RStudio v1.4.1103.

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

ACKNOWLEDGEMENTS

We thank Prof. Van Den Abbeele for critical reading of the manuscript, Isabel Saldanha and Steve Torr from the Liverpool School of Tropical Medicine for providing tsetse fly DNA and Dr Vet. Xanthe Helsen and Dr Vet. Wauter Van Deun for providing us with cattle, dog, goat, horse, pig and sheep blood.

FUNDING

P.B. received financial support from the Bill & Melinda Gates Foundation [OPP1174221]; ITM's SOFI program supported by the Flemish Government, Science & Innovation [EWI SOFI-2018 'Cryptic human and animal reservoirs compromise the sustained elimination of gambiense-human African trypanosomiasis in the Democratic Republic of the Congo']; F.V.d.B. was supported by the Department of Economy, Science and Innovation in Flanders and by the Research Foundation

Flanders [1226120N, 1528117N]; A.S. is supported by the UK Medical Research Council Fellowship [MR/L019701/1]; Z.C. is supported by an EASTBIO PhD studentship from the UK Biotechnology and Biological Sciences Research Council.

Conflict of interest statement. None declared.

REFERENCES

- Gibson, W.C. (1986) Will the real *Trypanosoma gambiense* please stand up. *Parasitol. Today*, **2**, 255–257.
- De Greef, C. and Hamers, R. (1994) The serum resistance-associated (SRA) gene of *Trypanosoma brucei rhodesiense* encodes a variant surface glycoprotein-like protein. *Mol. Biochem. Parasitol.*, **68**, 277–284.
- Berberof, M., Pérez-Morga, D. and Pays, E. (2001) A receptor-like flagellar pocket glycoprotein specific to *Trypanosoma brucei gambiense*. *Mol. Biochem. Parasitol.*, **113**, 127–138.
- Pays, E., Vanhollebeke, B., Uzureau, P., Lecordier, L. and Pérez-Morga, D. (2014) The molecular arms race between African trypanosomes and humans. *Nat. Rev. Microbiol.*, **12**, 575–584.
- Weir, W., Capewell, P., Foth, B., Clucas, C., Pountain, A., Steketee, P., Veitch, N., Koffi, M., De Meeüs, T., Kaboré, J. *et al.* (2016) Population genomics reveals the origin and asexual evolution of human infective trypanosomes. *Elife*, **5**, e11473.
- Mathieu-Daudé, F., Stevens, J., Welsh, J., Tibayrenc, M. and McClelland, M. (1995) Genetic diversity and population structure of *Trypanosoma brucei*: clonality versus sexuality. *Mol. Biochem. Parasitol.*, **72**, 89–101.
- WHO (2022) Number of reported cases of human african trypanosomiasis (*T. b. gambiense*). <https://www.who.int/data/gho/data/indicators/indicator-details/GHO/hat-tb-gambiense>, (29 September 2022, date last accessed).
- Jamonneau, V., Truc, P., Grébaud, P., Herder, S., Ravel, S., Solano, P. and De Meeus, T. (2019) *Trypanosoma brucei gambiense* group 2: the unusual suspect. *Trends Parasitol.*, **35**, 983–995.
- Franco, J.-R., Cecchi, G., Priotto, G., Paone, M., Diarra, A., Grout, L., Simarro, P.P., Zhao, W. and Argaw, D. (2020) Monitoring the elimination of human african trypanosomiasis at continental and country level: update to 2018. *PLoS Negl. Trop. Dis.*, **14**, e0008261.
- WHO (2020) Togo is first african country to end sleeping sickness as a public health problem. <https://www.who.int/news/item/27-08-2020-togo-is-first-african-country-to-end-sleeping-sickness-as-a-public-health-problem>, (29 September 2022, date last accessed).
- WHO (2021) WHO validates Côte d'Ivoire for eliminating sleeping sickness as a public health problem. <https://www.who.int/news/item/25-03-2021-who-validates-cote-d-ivoire-for-eliminating-sleeping-sickness-as-a-public-health-problem>, (29 September 2022, date last accessed).
- Büscher, P., Bart, J.-M., Boelaert, M., Bucheton, B., Cecchi, G., Chitnis, N., Courtin, D., Figueiredo, L.M., Franco, J.-R., Grébaud, P. *et al.* (2018) Do cryptic reservoirs threaten gambiense-sleeping sickness elimination? *Trends Parasitol.*, **34**, 197–207.
- Rock, K.S., Ndeffo-Mbah, M.L., Castañón, S., Palmer, C., Pandey, A., Atkins, K.E., Ndung'u, J.M., Hollingsworth, T.D., Galvani, A., Bever, C. *et al.* (2018) Assessing strategies against gambiense sleeping sickness through mathematical modeling. *Clin. Infect. Dis.*, **66**, S286–S292.
- Mehlitz, D. and Molyneux, D.H. (2019) The elimination of *Trypanosoma brucei gambiense*? Challenges of reservoir hosts and transmission cycles: expect the unexpected. *Parasite Epidemiol. Control*, **6**, e00113.
- Bisser, S., Lumbala, C., Nguertoum, E., Kande, V., Flevaud, L., Vatunga, G., Boelaert, M., Büscher, P., Josenando, T., Bessell, P.R. *et al.* (2016) Sensitivity and specificity of a prototype rapid diagnostic test for the detection of *Trypanosoma brucei gambiense* infection: a multi-centric prospective study. *PLoS Negl. Trop. Dis.*, **10**, e0004608.
- Geerts, M., Van Reet, N., Leyten, S., Berghmans, R., Rock, K.S., Coetzer, H.T., E-A Eyssen, L. and Büscher, P. (2020) *Trypanosomabrucei gambiense*-iELISA: a promising new test for the post-elimination monitoring of human African trypanosomiasis. *Clin. Infect. Dis.*, **73**, e2477–e2483.
- Büscher, P., Gilman, Q. and Lejon, V. (2013) Rapid diagnostic test for sleeping sickness. *N. Engl. J. Med.*, **368**, 1069–1070.
- Lejon, V., Büscher, P., Magnus, E., Moons, A., Wouters, I. and Van Meirvenne, N. (1998) A semi-quantitative ELISA for detection of *Trypanosoma brucei gambiense* specific antibodies in serum and cerebrospinal fluid of sleeping sickness patients. *Acta Trop.*, **69**, 151–164.
- Magnus, E., Vervoort, T. and Van Meirvenne, N. (1978) A card-agglutination test with stained trypanosomes (CATT) for the serological diagnosis of *T. b. gambiense* trypanosomiasis. In: *Annales de la Société belge de Médecine Tropicale*. Societe Belge de Medecine Tropicale, Vol. **58**, pp. 169–176.
- Van Meirvenne, N., Magnus, E. and Büscher, P. (1995) Evaluation of variant specific trypanolysis tests for serodiagnosis of human infections with *Trypanosoma brucei gambiense*. *Acta Trop.*, **60**, 189–199.
- Gibson, W., Nemetschke, L. and Ndung'u, J. (2010) Conserved sequence of the TgsGP gene in group 1 *Trypanosoma brucei gambiense*. *Infect. Genet. Evol.*, **10**, 453–458.
- Radwanska, M., Magnus, E., Claes, F., Pérez-Morga, D., Magez, S., Pays, E. and Büscher, P. (2002) Novel primer sequences for polymerase chain reaction-based detection of *Trypanosoma brucei gambiense*. *Am. J. Trop. Med. Hyg.*, **67**, 289–295.
- Compaoré, C.F.A., Ilboudo, H., Kaboré, J., Kaboré, J.W., Camara, O., Bamba, M., Sakande, H., Koné, M., Camara, M., Kaba, D. *et al.* (2020) Analytical sensitivity of loopamp and quantitative real-time PCR on dried blood spots and their potential role in monitoring human African trypanosomiasis elimination. *Exp. Parasitol.*, **219**, 108014.
- Felu, C., Pasture, J., Pays, E. and Pérez-Morga, D. (2007) Diagnostic potential of a conserved genomic rearrangement in the *Trypanosoma brucei gambiense*-specific TGSGP locus. *Am. J. Trop. Med. Hyg.*, **76**, 922–929.
- Agbo, E.E.C., Majiwa, P.A.O., Claassen, H.J.H. and Te Pas, M.F.W. (2002) Molecular variation of *Trypanosoma brucei* subspecies as revealed by AFLP fingerprinting. *Parasitology*, **124**, 349–358.
- Jamonneau, V., Garcia, A., Ravel, S., Cuny, G., Oury, B., Solano, P., N'Guessan, P., N'Dri, L., Sanon, R., Frezil, J.L. *et al.* (2002) Genetic characterization of *Trypanosoma brucei gambiense* and clinical evolution of human african trypanosomiasis in Cote d'Ivoire. *Trop. Med. Int. Health*, **7**, 610–621.
- Gibson, W. (2001) Molecular characterization of field isolates of human pathogenic trypanosomes. *Trop. Med. Int. Health*, **6**, 401–406.
- Auty, H., Anderson, N.E., Picozzi, K., Lembo, T., Mubanga, J., Hoare, R., Fyumagwa, R.D., Mable, B., Hamill, L., Cleaveland, S. *et al.* (2012) Trypanosome diversity in wildlife species from the Serengeti and Luangwa valley ecosystems. *PLoS Negl. Trop. Dis.*, **6**, e1828.
- Biteau, N., Bringaud, F., Gibson, W., Truc, P. and Baltz, T. (2000) Characterization of trypanozoon isolates using a repeated coding sequence and microsatellite markers. *Mol. Biochem. Parasitol.*, **105**, 187–202.
- Agbo, E.C., Majiwa, P.A.O., Claassen, E.J.H.M. and Roos, M.H. (2001) Measure of molecular diversity within the *Trypanosoma brucei* subspecies *Trypanosoma brucei brucei* and *Trypanosoma brucei gambiense* as revealed by genotypic characterization. *Exp. Parasitol.*, **99**, 123–131.
- Bromidge, T., Gibson, W., Hudson, K. and Dukes, P. (1993) Identification of *Trypanosoma brucei gambiense* by PCR amplification of variant surface glycoprotein genes. *Acta Trop.*, **53**, 107–119.
- Richardson, J.B., Lee, K.-Y., Mireji, P., Enyaru, J., Siström, M., Aksoy, S., Zhao, H. and Caccone, A. (2017) Genomic analyses of African trypanozoon strains to assess evolutionary relationships and identify markers for strain identification. *PLoS Negl. Trop. Dis.*, **11**, e0005949.
- Schares, G. and Mehlitz, D. (1996) Sleeping sickness in Zaire: a nested polymerase chain reaction improves the identification of *Trypanosoma* (Trypanozoon) *brucei gambiense* by specific kinetoplast DNA probes. *Trop. Med. Int. Health*, **1**, 59–70.
- Mathieu-Daudé, F., Bicart-See, A., Tibayrenc, M., Brenière, S.-F. and Bosseno, M.-F. (1994) Identification of *Trypanosoma brucei*

- gambiense* group I by a specific kinetoplast DNA probe. *Am. J. Trop. Med. Hyg.*, **50**, 13–19.
35. Lukes, J., Guilbride, D.L., Votýpka, J., Ziková, A., Benne, R. and Englund, P.T. (2002) Kinetoplast DNA network: evolution of an improbable structure. *Eukaryot. Cell*, **1**, 495–502.
 36. Cooper, S., Wadsworth, E.S., Ochsenreiter, T., Ivens, A., Savill, N.J. and Schnauffer, A. (2019) Assembly and annotation of the mitochondrial minicircle genome of a differentiation-competent strain of *Trypanosoma brucei*. *Nucleic Acids Res.*, **47**, 11304–11325.
 37. Singh, N., Curran, M.D., Rastogi, A.K., Middleton, D. and Sundar, S. (1999) Diagnostic PCR with *Leishmania donovani* specificity using sequences from the variable region of kinetoplast minicircle DNA. *Trop. Med. Int. Health*, **4**, 448–453.
 38. Ceccarelli, M., Galluzzi, L., Diotallevi, A., Andreoni, F., Fowler, H., Petersen, C., Vitale, F. and Magnani, M. (2017) The use of kDNA minicircle subclass relative abundance to differentiate between *Leishmania (L.) infantum* and *Leishmania (L.) amazonensis*. *Parasit. Vectors*, **10**, 239.
 39. Borst, P., Fase-Fowler, F. and Gibson, W.C. (1987) Kinetoplast DNA of *Trypanosoma evansi*. *Mol. Biochem. Parasitol.*, **23**, 31–38.
 40. Van den Broeck, F., Savill, N.J., Imamura, H., Sanders, M., Maes, I., Cooper, S., Mateus, D., Jara, M., Adai, V., Arevalo, J. et al. (2020) Ecological divergence and hybridization of Neotropical *Leishmania* parasites. *Proc. Natl. Acad. Sci. U.S.A.*, **117**, 25159–25168.
 41. Geerts, M., Schnauffer, A. and Van den Broeck, F. (2021) rKOMICS: an R package for processing mitochondrial minicircle assemblies in population-scale genome projects. *BMC Bioinf.*, **22**, 468.
 42. Lanham, S.M. and Godfrey, D.G. (1970) Isolation of salivarian trypanosomes from man and other mammals using DEAE-cellulose. *Exp. Parasitol.*, **28**, 521–534.
 43. Sambrook, J. and Russell, D.W. (2006) Purification of nucleic acids by extraction with phenol:chloroform. *CSH Protoc.*, **1**, pdb.prot4455.
 44. McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernysky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M. et al. (2010) The genome analysis toolkit: a mapreduce framework for analyzing next-generation DNA sequencing data. *Genome Res.*, **20**, 1297–1303.
 45. Danecek, P., Bonfield, J.K., Liddle, J., Marshall, J., Ohan, V., Pollard, M.O., Whitwham, A., Keane, T., McCarthy, S.A., Davies, R.M. et al. (2021) Twelve years of SAMtools and BCFtools. *Gigascience*, **10**, giab008.
 46. Huson, D.H. (1998) SplitsTree: analyzing and visualizing evolutionary data. *Bioinformatics*, **14**, 68–73.
 47. Li, D., Liu, C.M., Luo, R., Sadakane, K. and Lam, T.W. (2015) MEGAHIT: an ultra-fast single-node solution for large and complex metagenomics assembly via succinct de Bruijn graph. *Bioinformatics*, **31**, 1674–1676.
 48. Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J. (1990) Basic local alignment search tool. *J. Mol. Biol.*, **215**, 403–410.
 49. Chen, S., Zhou, Y., Chen, Y. and Gu, J. (2018) fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics*, **34**, i884–i890.
 50. Ray, D.S. (1989) Conserved sequence blocks in kinetoplast minicircles from diverse species of trypanosomes. *Mol. Cell. Biol.*, **9**, 1365–1367.
 51. Rognes, T., Flouri, T., Nichols, B., Quince, C. and Mahé, F. (2016) VSEARCH: a versatile open source tool for metagenomics. *PeerJ*, **4**, e2584.
 52. Waterhouse, A.M., Procter, J.B., Martin, D.M.A., Clamp, M. and Barton, G.J. (2009) Jalview version 2—a multiple sequence alignment editor and analysis workbench. *Bioinformatics*, **25**, 1189–1191.
 53. Van Reet, N., Pyana, P.P., Dehou, S., Bebronne, N., Deborggraeve, S. and Büscher, P. (2021) Single nucleotide polymorphisms and copy-number variations in the *Trypanosoma brucei* repeat (TBR) sequence can be used to enhance amplification and genotyping of trypanozoon strains. *PLoS One*, **16**, e0258711.
 54. Picozzi, K., Carrington, M. and Welburn, S.C. (2008) A multiplex PCR that discriminates between *Trypanosoma brucei brucei* and zoonotic *T. b. rhodesiense*. *Exp. Parasitol.*, **118**, 41–46.
 55. Pyana, P.P., Van Reet, N., Ngoyi, D.M., Lukusa, I.N., Shamamba, S.K.B. and Büscher, P. (2014) Melarsoprol sensitivity profile of *Trypanosoma brucei gambiense* isolates from cured and relapsed sleeping sickness patients from the Democratic Republic of the Congo. *PLoS Negl. Trop. Dis.*, **8**, e3212.
 56. Simpson, L. (1987) The mitochondrial genome of kinetoplastid protozoa: genomic organization, transcription, replication, and evolution. *Annu. Rev. Microbiol.*, **41**, 363–382.
 57. Stuart, K. (1979) Kinetoplast DNA of *Trypanosoma brucei*: physical map of the maxicircle. *Plasmid*, **2**, 520–528.
 58. Steinert, M. and Van Assel, S. (1980) Sequence heterogeneity in kinetoplast DNA: reassociation kinetics. *Plasmid*, **3**, 7–17.
 59. Pollard, V.W., Rohrer, S.P., Michelotti, E.F., Hancock, K. and Hajduk, S.L. (1990) Organization of minicircle genes for guide RNAs in *Trypanosoma brucei*. *Cell*, **63**, 783–790.
 60. Hong, M. and Simpson, L. (2003) Genomic organization of *Trypanosoma brucei* kinetoplast DNA minicircles. *Protist*, **154**, 265–279.
 61. Jasmer, D.P. and Stuart, K. (1986) Sequence organization in African trypanosome minicircles is defined by 18 base pair inverted repeats. *Mol. Biochem. Parasitol.*, **18**, 321–331.
 62. Lukeš, J., Kachale, A., Votýpka, J., Butenko, A. and Field, M. (2022) African trypanosome strategies for conquering new hosts and territories: the end of monophyly? *Trends Parasitol.*, **38**, 724–736.
 63. Savill, N.J. and Higgs, P.G. (1999) A theoretical study of random segregation of minicircles in trypanosomatids. *Proc. Roy. Soc. B: Biol. Sci.*, **266**, 611–620.
 64. Goodhead, I., Capewell, P., Wendi Bailey, J., Beament, T., Chance, M., Kay, S., Forrester, S., MacLeod, A., Taylor, M., Noyes, H. et al. (2013) Whole-genome sequencing of *Trypanosoma brucei* reveals introgression between subspecies that is associated with virulence. *MBio*, **4**, e00197.
 65. Tihon, E., Imamura, H., Dujardin, J.-C., Van Den Abbeele, J. and Van den Broeck, F. (2017) Discovery and genomic analyses of hybridization between divergent lineages of *Trypanosoma congolense*, causative agent of animal African trypanosomiasis. *Mol. Ecol.*, **26**, 6524–6538.
 66. Van den Broeck, F., Tavernier, L.J.M., Vermeiren, L., Dujardin, J.C. and Van Den Abbeele, J. (2018) Mitonuclear genomics challenges the theory of clonality in *Trypanosoma congolense*: reply to Tibayrenc and Ayala. *Mol. Ecol.*, **27**, 3425–3431.
 67. Gibson, W. and Garside, L. (1990) Kinetoplast DNA minicircles are inherited from both parents in genetic hybrids of *Trypanosoma brucei*. *Mol. Biochem. Parasitol.*, **42**, 45–53.
 68. Gibson, W., Crow, M. and Kearns, J. (1997) Kinetoplast DNA minicircles are inherited from both parents in genetic crosses of *Trypanosoma brucei*. *Parasitol. Res.*, **83**, 483–488.
 69. Rusman, F., Tomasini, N., Yapur, N.-F., Puebla, A.F., Ragone, P.G. and Diosque, P. (2019) Elucidating diversity in the class composition of the minicircle hypervariable region of *Trypanosoma cruzi*: new perspectives on typing and kDNA inheritance. *PLoS Negl. Trop. Dis.*, **13**, e0007536.
 70. Balyeidhusa, A.S.P., Kironde, F.A.S. and Enyaru, J.C.K. (2012) Apparent lack of a domestic animal reservoir in gambiense sleeping sickness in northwest Uganda. *Vet. Parasitol.*, **187**, 157–167.
 71. Cordon-Obras, C., Berzosa, P., Ndongo-Mabale, N., Bobuakasi, L., Buatiche, J.N., Ndongo-Asumu, P., Benito, A. and Cano, J. (2009) *Trypanosoma brucei gambiense* in domestic livestock of Kogo and Mbini foci (Equatorial Guinea). *Trop. Med. Int. Health*, **14**, 535–541.
 72. Cordon-Obras, C., Garcia-Estébanez, C., Ndongo-Mabale, N., Abaga, S., Ndongo-Asumu, P., Benito, A. and Cano, J. (2010) Screening of *Trypanosoma brucei gambiense* in domestic livestock and tsetse flies from an insular endemic focus (Luba, Equatorial Guinea). *PLoS Negl. Trop. Dis.*, **4**, e704.
 73. Cordon-Obras, C., Rodriguez, Y.F., Fernandez-Martinez, A., Cano, J., Ndongo-Mabale, N., Ncogo-Ada, P., Ndongo-Asumu, P., Aparicio, P., Navarro, M., Benito, A. et al. (2015) Molecular evidence of a *Trypanosoma brucei gambiense* sylvatic cycle in the human african trypanosomiasis foci of equatorial Guinea. *Front. Microbiol.*, **6**, 765.
 74. Simo, G., Fongho, P., Farikou, O., Ndjeto-Tchouli, P.I.N., Tchoumene-Labou, J., Njiokou, F. and Asonganyi, T. (2015) Trypanosome infection rates in tsetse flies in the ‘silent’ sleeping sickness focus of Bafia in the centre region in Cameroon. *Parasit. Vectors*, **8**, 528.
 75. Umeakuana, P.U., Gibson, W., Ezeokonkwo, R.C. and Anene, B.M. (2019) Identification of *Trypanosoma brucei gambiense* in naturally infected dogs in Nigeria. *Parasit. Vectors*, **12**, 420.
 76. Vourchakbé, J., Tioufack, A.A.Z., Mbida, M. and Simo, G. (2020) Trypanosome infections in naturally infected horses and donkeys of

- three active sleeping sickness foci in the south of Chad. *Parasit. Vectors*, **13**, 323.
77. Birhanu, H., Gebrehiwot, T., Goddeeris, B.M., Büscher, P. and Van Reet, N. (2016) New *Trypanosoma evansi* type B isolates from Ethiopian dromedary camels. *PLoS Negl. Trop. Dis.*, **10**, e0004556.
78. Birhanu, H., Fikru, R., Said, M., Kidane, W., Gebrehiwot, T., Hagos, A., Alemu, T., Dawit, T., Berkvens, D., Goddeeris, B.M. *et al.* (2015) Epidemiology of *Trypanosoma evansi* and *Trypanosoma vivax* in domestic animals from selected districts of Tigray and Afar regions, northern Ethiopia. *Parasit. Vectors*, **8**, 212.
79. Boushaki, D., Wallis, J., Van den Broeck, F. and Schnaufer, A. (2022) Molecular Analysis of Trypanosome Infections in Algerian Camels. *Acta Parasitol.*, **67**, 1246–1253.
80. Schnaufer, A., Panigrahi, A.K., Panicucci, B., Igo, R.P. Jr, Salavati, R. and Stuart, K. (2001) An RNA ligase essential for RNA editing and survival of the bloodstream form of *Trypanosoma brucei*. *Science*, **291**, 2159–2162.
81. Ziková, A., Schnaufer, A., Dalley, R.A., Panigrahi, A.K. and Stuart, K.D. (2009) The F(0)F(1)-ATP synthase complex contains novel subunits and is essential for procyclic *Trypanosoma brucei*. *PLoS Pathog.*, **5**, e1000436.
82. Simpson, L., Douglass, S.M., Lake, J.A., Pellegrini, M. and Li, F. (2015) Comparison of the mitochondrial genomes and steady state transcriptomes of two strains of the trypanosomatid parasite, *Leishmania tarentolae*. *PLoS Negl. Trop. Dis.*, **9**, e0003841.
83. Mehlitz, D. (1986) Le réservoir animal de la maladie du sommeil à *Trypanosoma brucei gambiense*. p. 167, ISBN 2-85985-127-5.