

Pervasive sequence-level variation in the transcriptome of *Plasmodium falciparum*

Bruhad Dave, Abhishek Kanyal, D.V. Mamatharani and Krishanpal Karmodiya¹*

Department of Biology, Indian Institute of Science Education and Research, Dr. Homi Bhabha Road, Pashan, Pune 411008, Maharashtra, India

Received July 01, 2021; Revised March 09, 2022; Editorial Decision April 28, 2022; Accepted May 14, 2022

ABSTRACT

Single-nucleotide variations (SNVs) in RNA, arising from co- and post-transcriptional phenomena including transcription errors and RNA-editing, are well studied in a range of organisms. In the malaria parasite *Plasmodium falciparum*, stage-specific and non-specific gene-expression variations accompany the parasite's array of developmental and morphological phenotypes over the course of its complex life cycle. However, the extent, rate and effect of sequence-level variation in the parasite's transcriptome are unknown. Here, we report the presence of pervasive, non-specific SNVs in the *P. falciparum* transcriptome. SNV rates for a gene were correlated to gene length ($r \approx 0.65-0.7$) but not to the AT-content of that gene. Global SNV rates for the *P. falciparum* lines we used, and for publicly available *P. vivax* and *P. falciparum* clinical isolate datasets, were of the order of 10^{-3} per base, $\sim 10\times$ higher than rates we calculated for bacterial datasets. These variations may reflect an intrinsic transcriptional error rate in the parasite, and RNA editing may be responsible for a subset of them. This seemingly characteristic property of the parasite may have implications for clinical outcomes and the basic biology and evolution of *P. falciparum* and parasite biology more broadly. We anticipate that our study will prompt further investigations into the exact sources, consequences and possible adaptive roles of these SNVs.

INTRODUCTION

Fidelity in the transcription of DNA into RNA and the correct translation of mRNAs into proteins is crucial. Accurately made proteins produce 'correct' phenotypes and ensure that the cell survives. DNA mutations represent a significant source of variation in the flow of genetic information, and many affect phenotypes and become established as single-nucleotide polymorphisms (SNPs) in a given pop-

ulation of cells. For unicellular parasites, these SNPs are a conceivably essential way to adapt to life in their respective hosts over many generations. Proteomes are comparatively more robust to variations, but the transcriptional landscape provides much scope for diversification.

A previous body of work has shown an inherent heterogeneity in the gene-expression levels as well as copy-number variation (1) in *Plasmodium falciparum*, a parasite that still affects over 200 million people worldwide (2). These observations have been reported across multiple conditions and levels – in untreated parasite cultures (3) and as a response to physiological-like stressors (4); at the population (5) and the single-cell (6) level, and between clinical isolates and lab-adapted cultures (1): Antigenic variation in *P. falciparum* is well-documented (7), and the parasite exhibits alternative splicing (8–11); such transcriptional variation in essentially clonal populations represents another potential layer of complexity that is likely to affect clinical outcomes (12,13), and studies have shown that heterogeneity through gene expression variation serves as a population-level survival strategy in unicellular organisms (14,15).

Sequence-level variation—single-nucleotide variations (SNVs) and insertion-deletion events (indels)—is another potential source of population-level transcriptomic diversity. In the form of transcriptional error rates and RNA editing, it has been extensively described in systems including bacteria (16,17), yeast (18), *Caenorhabditis elegans* (19), cephalopods (20), rodents (21,22) and humans (23,24). However, studies on such heterogeneity in the *P. falciparum* transcriptome are largely missing and such variations have the potential to impact downstream sources of transcriptional heterogeneity, which may facilitate stress adaptation – these might in turn allow for persister populations to survive under drug regimes and eventually evolve drug-resistance. To investigate the extent and rate of SNVs in *P. falciparum*, we analysed transcriptomic data from a lab-adapted parasite culture (strain 3D7), to obtain three datasets – an untreated control, a temperature-stressed culture and a drug-stressed culture. We also analysed RNAseq data from three drug-resistant *P. falciparum* lines sourced from the MR4 repository, namely MRA.1236, MRA.1240 and MRA.1241. We also performed whole-genome se-

*To whom correspondence should be addressed. Tel: +91 20 25908195; Email: krish@iiserpune.ac.in

quencing corresponding to each of the six resulting datasets, which allowed us to accurately discard transcriptomic sequence variations arising from genomic SNPs, without the need to use predictive or consensus-based methods. We then used REDIttools 2.0 (25) to perform empirical variant-calling against the *P. falciparum* reference genome (v.41), and we found SNV rates on the order of 10^{-3} per base, a metric that was consistent across all six *P. falciparum* samples. In this work, we describe the spectrum of base substitutions and their predicted functional effects.

MATERIALS AND METHODS

Parasite cultures

P. falciparum strain 3D7 was cultured as previously described (26). Briefly, parasites were cultured in RPMI1640 medium supplemented with 25 mM HEPES, 0.5% Albumin I, 1.77 mM sodium bicarbonate, 100 μ M hypoxanthine and 12.5 μ g ml⁻¹ gentamicin sulfate at 37°C. Parasites were sub-cultured after every 2 days. Subculturing was done by splitting the flask into multiple flasks in order to maintain parasitemia around 5%. Hematocrit was maintained to 1–1.5% by adding freshly washed O⁺ve human RBC isolated from healthy human donors. Synchronization was done with the help of 5% sorbitol in the ring stage. Late-stage synchronization was performed using the Percoll density gradient method (63%). Parasitemia was monitored using Giemsa staining of thin blood smear.

Stress induction

Parasites were subjected to heat and therapeutic (artemisinin treatment) stresses for 6 hours from late ring (~17 h) to early trophozoite (~23 h) stage as described earlier (27). Briefly, double synchronization was carried out to achieve tight synchronization of parasite stages. Parasites were exposed to heat stress (40°C for 6 h) and artemisinin stress (30 nM for 6 h).

RNA sequencing

Parasites were harvested for RNA isolation after 6 h of stress induction. Total RNA was isolated using TRIzol reagent according to the protocol. DNase treated RNA was used for cDNA synthesis. Quality of the RNA was verified using Agilent Bioanalyzer 2100. The cDNA libraries were prepared for samples using Illumina TruSeq RNA library preparation kit. Transcriptome sequencing was performed using Illumina NextSeq 550 system in house at IISER Pune with a standard flow cell.

Whole genome sequencing

Plasmodium genome DNA was isolated using the genome DNA isolation kit. DNA concentrations were measured on the Qubit double-stranded DNA (dsDNA) HS assay kit (Invitrogen). Libraries for paired-end sequencing were constructed from DNA extracts ranging from <50 ng/ml to 0.2 ng/ μ l, using the Illumina NexteraXT kit (FC-131-1024, Illumina, CA, USA). The pooled NexteraXT libraries were loaded onto an Illumina NextSeq 550 system in house at IISER, Pune with a standard flow cell.

DATA ANALYSIS, SNV CALLING AND DOWNSTREAM ANALYSIS

Quality control and alignment

We first checked the sequencing quality for each RNAseq and WGS sample by running each fastq file through FastQC (28), and we trimmed the dataset to exclude positions that had a phred score of <25 using Trim Galore (29) v0.6.6 (with cutadapt (30) v3.2). We then aligned the trimmed RNAseq fastq files to the *Plasmodium falciparum* 3D7 genome (v. 41 from Ensembl) using STAR (31) v2.7.6a in 2-pass mode, and then indexed the resulting BAM files using samtools (32,33) v1.10. To align the WGS data, we used BWA (bwa mem) (34) v0.7.17-r1198, following which we converted the resulting SAM file to a BAM file, which we sorted by coordinates and indexed using samtools. The command-line options used here are provided in the Supplementary Information.

SNV calling and filtering

In order to use REDIttools 2.0 with python 3.6, we modified the python scripts we used to update the syntax to python 3 wherever appropriate. We then ran REDIttools 2.0 using minimal command line options (Supplementary information) on the aligned RNAseq data and WGS data, using the RNA-table obtained from the former run as an input to the latter to specify the positions to be assayed in the WGS. The output generated by REDIttools is a tab-delimited table containing data about each position read by the program in the NGS data, including the chromosome, position, reference nucleotide, coverage at that position, average read quality at that position, alternate nucleotides found at that position (if any) for the transcriptomic data (the first nine columns) and similar information for WGS data (the remaining columns). In order to annotate the RNA-table with the DNA-table, we used Annotate_with_DNA.py, which is provided as part of the REDIttools suite of scripts. This script writes the WGS data for each locus in the RNA-table filling in columns 10 and further. We filtered this combined output table using *awk* command-line to exclude those positions where: (i) the WGS data showed an SNP, (ii) the frequency of base changes in RNAseq data was <0.1, (iii) the RNAseq data was invariant, (iv) the RNAseq coverage was <5 reads and (v) the WGS coverage was <10 reads. We then used samtools to remove duplicate reads in the aligned RNAseq data as described in Supplementary Information. Then, we reapplied REDIttools to this deduplicated RNAseq file, using the previously generated, DNA-annotated and filtered RNA-table as a region file along with the previously noted REDIttools options. We filtered and used the output tables from this REDIttools run for further analysis.

Downstream analysis

We used another script called AnnotateTable.py (also provided with REDIttools 2.0) to annotate each SNV in the final REDIttools table with gene IDs, using a sorted and trimmed gtf file (*P. falciparum* v41, from Ensembl) using the command lines (Supplementary information). Using custom python scripts, we calculated the number of SNVs oc-

curing in each gene and the relative frequency of SNV occurrence in each gene and extracted the nucleotide sequences on either side of each changed position in order to analyse patterns in the flanking nucleotides. We also arrayed each SNV on the gene in which it occurred, dividing each gene into 100 equal-width bins, and constructed a histogram based on binning each SNV in order to analyse any positional bias of SNV occurrence on the gene body

We converted the gene-annotated RNA_DNA_deduplicated file into a format resembling variant call format to use as input for snpEff (35), which we used according to the documentation to perform functional annotation of the SNVs in our output file. We processed the output VCF files generated by snpEff using snpSift (36) to filter the file and retain data related to predicted amino-acid changes. We analysed this filtered VCF file using a custom script to visualise the spectrum of amino-acid changes in each sample. We did not retain any annotations of the types ‘upstream variant’ and ‘downstream variant’, since snpEff defines upstream and downstream regions as 5 kilobases in length, which is too large for the relatively compact genome of *P. falciparum*. SNV rates per base were calculated as the number of SNVs in the filtered output divided by the genome length in nucleotides. We used Salmon (37) to obtain gene abundance estimates in transcripts per million.

Analysis of samples sourced from SRA

NGS (RNAseq) datasets for *P. falciparum* patient-isolates from Mali, *Plasmodium vivax* liver-stages (mixed cultures and hypnozoites), *Plasmodium vivax* blood-stages, *Escherichia coli*, and *Bacillus subtilis* were downloaded using SRA toolkit (38) as fastq files. We used genome versions PvP01 for *P. vivax*, ASM584v2 for *Escherichia coli* and ASM608879v1 for *B. subtilis*. The bacteria datasets were aligned using BWA, while the *Plasmodium* spp. datasets were aligned using STAR. We analysed each dataset with REDIttools 2.0, filtering the output files to remove positions where; (i) the frequency of base changes in RNAseq data was <0.1, (ii) the RNAseq data was invariant, (iii) the RNAseq coverage was <5 reads and (iv) the average phred-score was <25. SNV rates per base were calculated as the number of SNVs in the filtered output divided by the genome length in nucleotides.

Circos plot generation

We used the R package RCircos (39) to generate circos representations for the SNV distribution on the whole-genome scale.

Statistical analysis and error bars

A Pearson’s chi-squared test was performed to assess whether the distribution of frequency of SNV types (by base-substitution, as % of total) was significantly different from a uniform distribution (with the assumption that all base-substitutions are equally likely to occur). The test was performed in R using the inbuilt chisq.test() function. A Dunnett’s test was used to calculate statistical significance for the differences in average SNV rates between the 3D7 parasite cultures and the three MRA parasite lines. This test

was performed in R using the DunnettTest() function from the R library DescTools (40), using PF3D7_Ctrl as the control set. For Figures 2–4, error bars indicate 95% confidence interval for replicates and the number of RNA-sequencing replicates for each sample is as follows:

- *P. falciparum* MRA_1236: 3 replicates
- *P. falciparum* MRA_1240: 3 replicates
- *P. falciparum* MRA_1241: 2 replicates
- *P. falciparum* 3D7 Control: 2 replicates
- *P. falciparum* 3D7 Drug-stressed: 2 replicates
- *P. falciparum* 3D7 Temperature-stressed: 2 replicates
- *P. falciparum* Mali Isolates [PRJNA498885]: 3 replicates
- *P. vivax* [PRJNA422240] (mixed culture): 2 replicates
- *P. vivax* [PRJNA422240] (hypnozoites only): 2 replicates
- *P. vivax* [PRJNA515743]: 4 replicates
- *E. coli* [PRJNA592142]: 3 replicates
- *B. subtilis* [PRJNA592142]: 3 replicates

RESULTS

A global view of transcriptional sequence-variation

We applied REDIttools 2.0 (a tool originally designed to detect RNA-editing) to both the RNAseq data and the WGS data for each sample (Methods). We removed SNVs arising from genomic single nucleotide polymorphisms (SNPs) as well as discarded positions where RNAseq data showed no variation. Additionally, we ascertained optimum cutoffs for WGS coverage (Supplementary Figure S1, Supplementary Table S3), RNAseq coverage (Supplementary Figure S2, Supplementary Table S4) and frequency of variation (the number of reads supporting a variant nucleotide divided by the total number of reads covering that position) (Supplementary Figure S3, Supplementary Table S5). SNVs not supported by at least 10 WGS reads and 5 RNAseq reads were removed. Finally, we retained only those SNVs whose frequency of variation was greater than or equal to 0.1. In our estimation, this combination of cutoffs ensures that sequencing errors and other technical errors are largely filtered out, and mitigates any technical variability. We noted that the rates of occurrence of SNVs—which were spread out all over the genome (Figure 1A)—were not dependent on sequencing and alignment statistics of the sample, indicating that they are consequences of biological characteristics of the parasite, rather than technical properties of the sequencing runs and alignment methods (Supplementary Table S1).

Following filtering, each sample showed $\sim 3 \times 10^4$ variant positions (Supplementary Table S1). Annotating SNVs with the genes in which they occurred showed that they were located all across the transcriptome, affecting ~ 3660 of 5700 genes in the reference annotation (Supplementary Table S2) in all samples. To check whether the SNVs were biased toward either end of the transcript, we constructed an average histogram of SNV frequency, dividing each affected gene into bins of equal width, and then assigning bins to each SNV. We observed that SNVs occupied the whole averaged gene length and did not show a specific positional bias along the gene body (Figure 1B). Notably, replicate datasets for a given strain showed only a small amount of overlap, further indicating the pervasive, non-specific nature of the SNVs (Supplementary Figure S4). We computed the Sum of SNV Frequencies (Σ SNV Freq) for a given gene as the

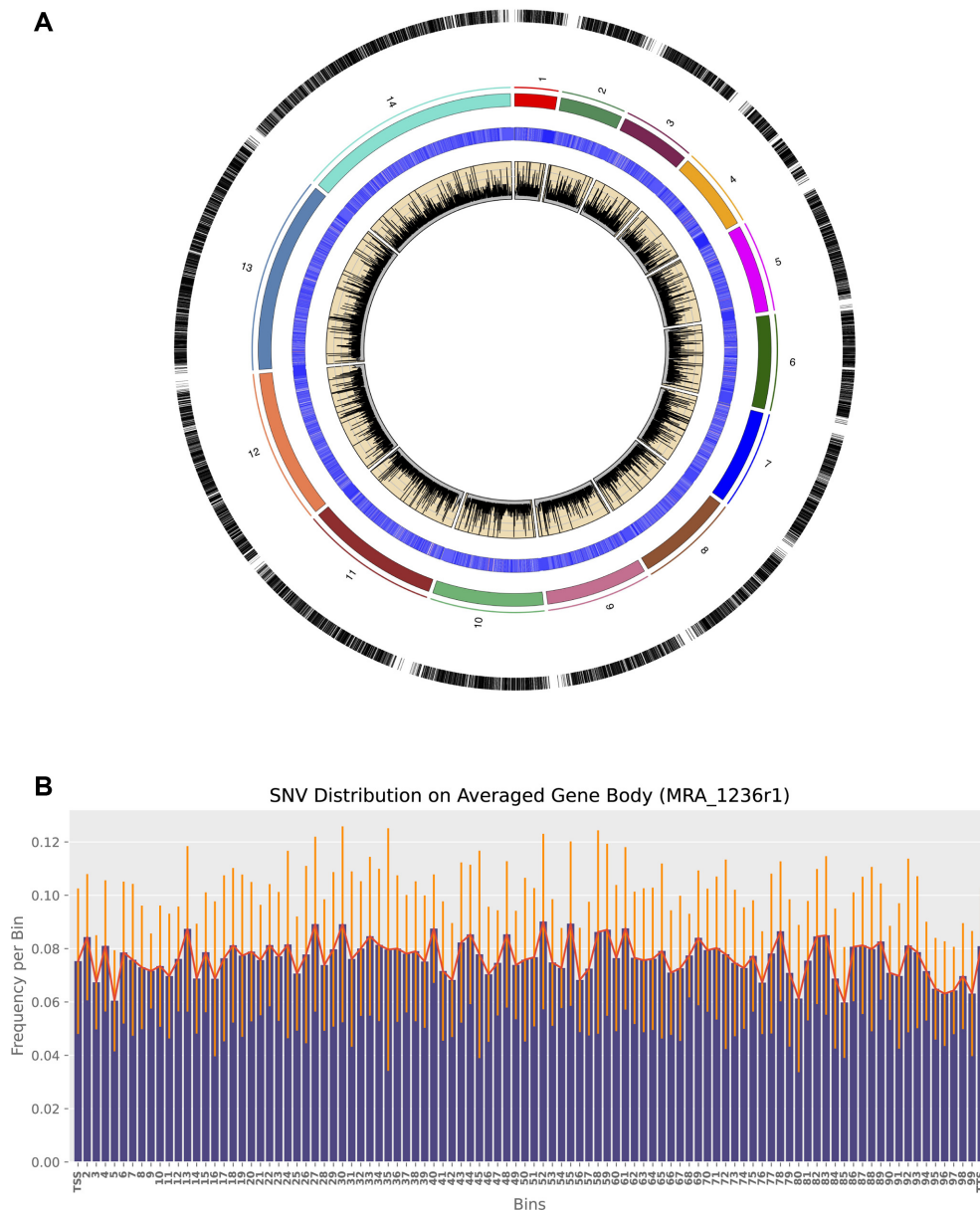


Figure 1. A global view of single-nucleotide variations in *Plasmodium falciparum* (data shown for a representative sample, MRA_1236 replicate 1). Plot showing the pervasive nature of SNVs. From outward: Track 1: SNV positions represented as vertical lines; Track 2: Ideogram of chromosomes, proportional to chromosome lengths; Track 3: Heatmap of RNaseq coverage at recorded positions, converted to \log_{10} scale for visualization; Track 4: Area plot of frequency of variation at recorded positions. Histogram showing distribution of SNVs on an averaged gene body.

sum of variation frequencies at each variant locus on a gene. Variation frequency is the number of variant reads divided by the total reads covering that locus. Σ SNV Freq showed no linear correlation to the gene abundance estimates (in transcripts per million) of that gene (Supplementary Figure S5). We similarly observed that Σ SNV Freq was not significantly correlated (Pearson's coefficient of correlation $r \sim 0.2$) to the AT-content of that gene (the number of A- or T-nucleotides divided by the length of that gene in base pairs) (Supplementary Figure S6). We did, however, observe that both the Σ SNV Freq and the number of variant loci on a gene had a positive (Pearson's $r = 0.7$) linear correlation with gene length (Supplementary Figure S9).

The spectrum of SNV types and effects

In order to understand the variations in greater detail, we characterised the range of nucleotide substitutions and their probable functional effects. We observed that A-to-G changes and T-to-C changes predominated, each representing on average 28.6% and 20.3% of the total number of SNVs found (Figure 2A). A-to-T and T-to-A substitutions ($\sim 14\%$ and $\sim 13\%$ respectively, of the total number of SNVs) were the second most abundant on average. G-to-C and C-to-G substitutions were the least abundant SNV type (Figure 2A, Supplementary Table S6). These proportions were relatively well-conserved between the drug-resistant MRA lines as well as the drug-sensitive 3D7 cul-

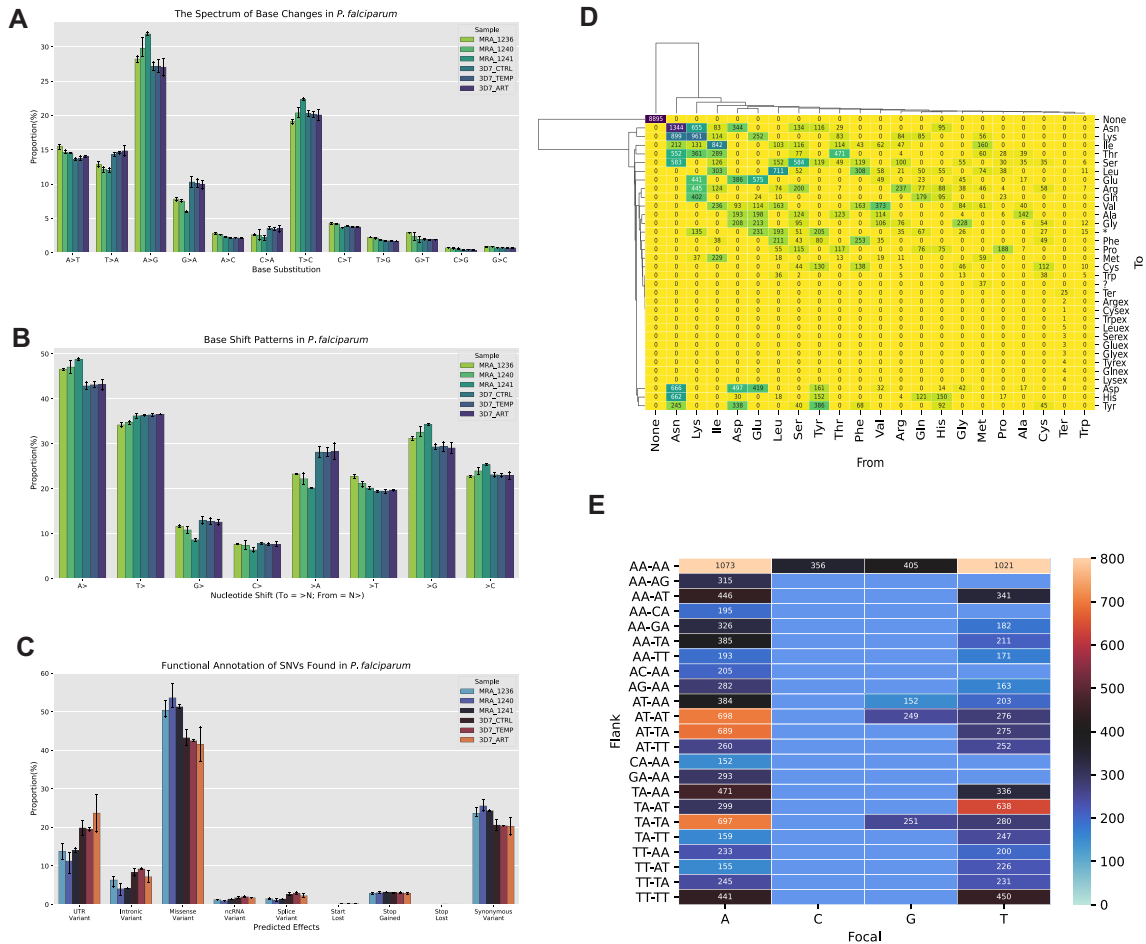


Figure 2. SNV types and effects. (A) Base changes (percentage of total). (B) Base shift patterns (%>X denotes the proportion of SNVs where base X changed to another base; %>X denotes the proportion of SNVs where the reference base was X). The distribution of predicted functional effects: (D) representative heatmap of the spectrum of amino acid changes. (E) Representative heatmap of the abundance of the most common dinucleotide patterns flanking each focal (original/reference) base. (A–C) Error bars represent 95% confidence intervals; number of replicates = 2 (3D7_CTRL, 3D7_TEMP, 3D7_ART, MRA_1241), 3 (MRA_1236, MRA_1240). (D) and (E) show data from sample MRA_1236 replicate 1.

ture. To check whether these base-change characteristics depended on environmental stresses, we subjected the 3D7 *P. falciparum* 3D7 line to temperature stress (at 40°C) and separately, to mild drug stress (dihydroartemisinin at 30 nM), each for six hours at the early trophozoite stage. However, these stresses had a minimal effect on the base-change profile, and the proportions of nucleotide substitutions were remarkably similar between the control *P. falciparum* 3D7 culture and the two stressed cultures (Figure 2A). With the assumption that all base substitutions are equally probable, we performed a chi-squared test and found that the frequency distribution of base substitutions we observed differed significantly from the expected uniform distribution ($P = 5.8868 \times 10^{217} x^2 = 1521.3$, $df = 167$). Interestingly, As and Ts changed most frequently, as might be expected due to the AT-richness of the *P. falciparum* genome, but Cs and especially Gs were misincorporated most often by percentage (Figure 2B, Supplementary Table S7). Given the bias of the spectrum of base-substitutions toward the A-to-G and T-to-C substitution types, we speculated that RNA editing may be responsible for a subset of SNVs. We tested this hypothesis by searching for sequence-motifs centered

on or around SNVs as well as performing a BLAST-based search for potential RNA-editing enzymes in the proteome of the parasite. We used human and *Trypanosoma brucei* enzymes known to be RNA-editors as query sequences for the latter analysis, but we did not find any high-confidence candidate RNA-editors in *P. falciparum*, nor any conserved motifs (Supplementary Table S8).

Further, we annotated effects to the SNVs using snpEff (35) to investigate the range of predicted functional consequences, and filtered the snpEff outputs using snpSift (36) to analyse amino acid change patterns. A small percentage were found in intronic regions, possibly reflecting alternative splicing in some transcripts (Figure 2C, Supplementary Table S9). The majority of SNVs in the coding region were missense (mean value 47.8%), with about 22.7% of them being synonymous, and about 3% being nonsense changes (Figure 2C). These proportions were also well conserved between all three MRA lines, the 3D7 control and two 3D7 stressed cultures, with missense variants being more common across MRA lines (51.8%) than 3D7 lines (42.4%). Asparagine was consistently the most changed amino acid, likely due to its abundance in the parasite's proteome, with

lysine, isoleucine, aspartate, and glutamate rounding out the most changed residues, although a proportion of SNVs resulted in synonymous codon changes (Figure 2D, Supplementary Figure S7). We also tested whether focal nucleotides had characteristic flanking sequences that might increase their propensity to be changed. To this end, we extracted a pentanucleotide sequence for each variant position, taking two nucleotides on either side of each focal nucleotide, and quantified the most abundant flanking sequences (Figure 2E, Supplementary Figure S8). We found that an ‘AA_AA’ or a ‘TT_TT’ pattern accounted for most of the SNVs we observed. This pattern was most abundant around all focal nucleotides, but patterns other than these were very rare when the focal base was a C or a G. We further observed that an A or a T seemed a necessary part of both 5'- and 3'-flanking dinucleotides in all of the most abundant flanking sequences. For focal As and Ts, ‘AT_TA’ (for focal As) and ‘TA_AT’ (for focal Ts), i.e. patterns forming pentanucleotide sequences of alternating As and Ts, were also well represented. We note that these patterns are likely to be further reflections of the parasite’s AT-rich genome.

Transcriptional sequence-variation in *Plasmodium falciparum*

To investigate the extent of the SNVs in each *P. falciparum* line, we calculated a per-kilobase (/kb) rate of change by dividing the total number of SNVs found by the length of the *P. falciparum* genome. We obtained an average rate of ~ 1.48 variations/kb, i.e. $\sim 1.48 \times 10^{-3}$ variations per base. This unexpectedly high variation rate was also well-conserved between various *P. falciparum* lines, with strain- and stress-specific rates ranging from 1.29 to 1.78 variations/kb (Figure 3).

Transcriptional sequence-variation is higher in *Plasmodium* than bacteria

Transcriptional sequence-level variation is often attributable to error rates associated with transcription. These rates are noted to range from $\sim 10^{-5}$ – 10^{-6} per base in yeast (18), similar rates in *C. elegans* (19), and comparable or higher rates in bacteria ranging from 10^{-4} (17) to 10^{-5} – 10^{-6} (16) errors per base. To investigate the relative differences between the SNV rate in *Plasmodium falciparum* and other organisms, we retrieved transcriptomic data describing *E. coli*, *B. subtilis* [PRJNA592142], *P. vivax* (liver stages [PRJNA422240] and blood stages [PRJNA515743]), and *P. falciparum* patient isolates [PRJNA498885] from SRA and performed identical calculations to arrive at SNV rates per nucleotide and per kb using REDIttools. We used filtering parameters similar to the ones described above, except for the WGS-coverage filter. We observed an average SNV rate per base of 7.23×10^{-4} for *E. coli*, 3.49×10^{-4} for *B. subtilis*, 4.35×10^{-3} for *P. vivax* schizonts + hypnozoite mixed sample, 1.80×10^{-3} for *P. vivax* hypnozoites, $\sim 10^{-3}$ for *P. vivax* blood stages cultured *in vivo* in simians, and 7.42×10^{-3} in *P. falciparum* patient isolates (Figure 4).

Since these values came from transcriptome-only analyses and SNP exclusion was not possible (this was reflected in

the much higher rates we observed for the *Plasmodium spp.* samples from SRA), we also calculated the SNV rates of in-house samples without SNP exclusion ($\sim 2.1 \times 10^{-3}$ across all sample replicates) (Figure 4). Interestingly, we observed that the SNV rates in *Plasmodium* species were consistently an order of magnitude higher than in bacteria, for which the reported error rates had been the highest to date (to the best of our knowledge).

DISCUSSION

In this work, we show that SNVs arising from base substitutions occur pervasively and non-specifically at a rate of the order of one every kilobase in the transcriptomes *P. falciparum* across treatment conditions and strains. We also show that a majority of these SNVs are predicted to have a functional impact. Given their non-specific occurrence, we speculate that these SNVs are likely reflections of RNA Pol II errors.

It is also possible that *P. falciparum* has mechanisms facilitating RNA editing, another potential source of a subset of the SNVs we report. While our search did not yield any distinct sequence motifs, nor any high-confidence RNA-editing enzymes candidates (Supplementary Table S8), RNA editing may still be occurring in *P. falciparum* – the phenomenon in the parasite might be more akin to what is termed promiscuous editing (41) in humans, wherein repetitive elements in the human transcriptome, such as Alu elements, are widely edited. Given that the *P. falciparum* genome and transcriptome are AT(/AU)-rich, it is conceivable that such a form of RNA editing may be occurring in relatively low-complexity regions of the parasite transcriptome. However, the fact that replicate datasets of any given sample showed an overlap in SNVs of no greater than $\sim 10\%$ indicates that RNA-editing is unlikely to be (solely) responsible for this pervasive variation.

This aforementioned AT-richness is also a characteristic of the *Plasmodium* genus, and specifically *P. falciparum* that research is still unable to fully explain. Our data shows that A and T are the most likely to change (as would be expected from an AT-rich starting point), Gs and Cs are more likely to be misincorporated (Figure 2C). As a result, it seems that the net effect of the SNVs in *P. falciparum* is of a compensating nature, in the context of nucleotide bias. We also observe such a pattern when comparing %GC values for WGS data and the corresponding RNAseq datasets, with the %GC rising a little in the transcriptome (Supplementary Table S10). Previous work (20) suggests that a large number of SNVs (as highly specific, recoding RNA editing events) in cephalopods represents a paradigm of low levels of genetic mutation, and correspondingly high levels of transcriptomic mutations (which provide the requisite protein diversity). In *Plasmodium*, if the numerous SNVs we report actually do lead to a nucleotide-bias compensation, then an analogous paradigm may be in play, and this may, in part, explain why the *Plasmodium* genome itself retains its AT-richness.

Indeed, this AT-richness was a preeminent concern to us since it opened up the possibility that the abundance of transcriptional SNVs we observed could be down to a handful of confounding factors. The first of these was a bio-

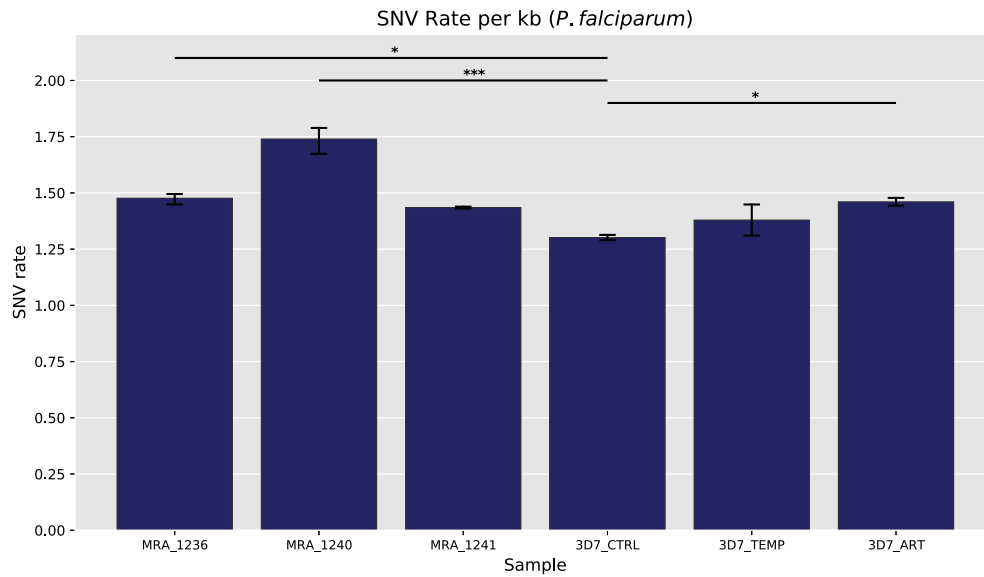


Figure 3. SNV rates in MRA and 3D7 *Plasmodium falciparum* lines. *** denotes significance at $P < 0.001$ (two-tailed) and *denotes $P < 0.05$ as calculated using a Dunett's test with PF3D7_Ctrl as the control group. Error bars represent 95% confidence intervals; number of replicates = 2 (3D7_CTRL, 3D7_ART, 3D7_TEMP, MRA_1241), 3 (MRA_1236, MRA_1240).

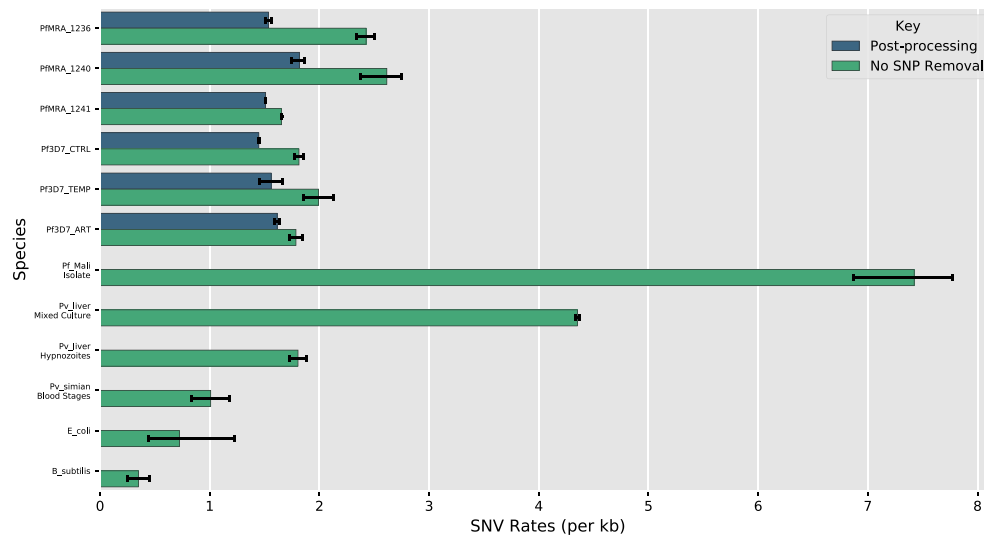


Figure 4. SNV rates of MR4 and 3D7 *Plasmodium falciparum* lines as compared with those for other samples and species. Error bars represent 95% confidence intervals; number of replicates = 2 (3D7_CTRL, 3D7_ART, 3D7_TEMP, MRA_1241, *P. vivax* Mixed Culture, *P. vivax* Hypnozoites), 3 (MRA_1236, MRA_1240, PF_Mali Isolate, *E. coli*, *B. subtilis*), 4 (*P. vivax* Blood Stages).

logical issue—the possibility of the parasite’s RNA pol II exhibiting slippage or becoming more error-prone in AT-rich swathes of the genome. Recent work showed that the translational machinery of *P. falciparum* can handle the AT-richness of its genome effectively (42). In one of their experiments, the authors noted that the abundance of polyA-containing reporter mRNAs, as measured using Real-Time Quantitative Reverse Transcription PCR, was comparable to that of control (non-polyA-containing) reporter mRNAs in *P. falciparum*. This implies that the transcriptional machinery in the parasite is likely also adept at handling the AT-richness of the parasite’s genome. This fact, taken together with our data would seem to indicate that most of

the SNVs we observed, while they might be reflections of an inherent transcriptional (i.e. RNA Pol II mediated) error rate, are not entirely explicable by simply the relatively lower complexity of the parasite’s transcriptome. It also leaves open the question of whether this points to the possibility that the parasite RNA pol II has evolved to this state of stability on an AT-rich landscape.

Another, larger source of concern was the possibility of the parasite’s AT-richness causing technical errors during the library preparation phase of NGS. However, our results (especially our observation that the AT-richness of a gene is not significantly correlated to the variation in that gene) suggest that the analysis strategies and filtering we

employed have mitigated such false-positives to a significant degree. We also sought to reduce technical errors in general, namely errors in sequencing, errors associated with PCR, and errors made by the reverse transcriptase enzyme during cDNA preparation. Quality trimming of the raw data (Methods), as well as the quality cutoff built into REDItools2.0, mitigated sequencing errors by omitting low quality data, i.e. ambiguous sequence calls, from further analysis. We used read depth cutoffs at two filtering stages: first, when filtering the initial variant calls, and second, when filtering variant calls made on deduplicated data (Methods). These steps mitigated errors introduced into reads during PCR, since such errors would not pass the read depth cutoff following data deduplication. Lastly, as described above in the second instance of data filtration, we retained only those SNVs where the RNA coverage was at least 5 reads: the aforementioned deduplication means that the SNVs thus retained were covered by at least 5 *unique* reads. We expect that this significantly lowers the probability of reverse transcriptase errors being retained in the final set of variants, since it is extremely unlikely that the enzyme will make an error corresponding to the same gene locus on several unique reads, which by definition do not share a common start and end position in alignment. We thus expect that quality-control of the raw data and stringent filtering of the results of our analysis have together significantly mitigated technical artefacts, and the SNVs we report are biological in origin.

An intriguing observation we made in this study was the apparent lack of a positive association and indeed a negative, non-linear relationship between Σ SNV Freq and gene expression (Supplementary Figure S5). Unfortunately, we do not have an exact explanation for this observation. One possibility is the presence of an as-yet unknown RNA-surveillance mechanism that controls the number of errors occurring in transcripts associated with a gene. An analogous mechanism has been predicted in *E. coli*, where transcriptional error rates were found to be lowest in highly abundant proteins, on which selection is expected to act more strongly and in which the consequences of high error rates would likely be most significant (43). Our observations indicate the possible presence of such a mechanism in *P. falciparum*, although we are unable to elucidate its precise nature in the present work, and we look forward to interpretations from the field.

In summary, just as clonal *P. falciparum* cultures exhibit an inherent variation in gene expression levels (1,3–6), our results suggest that heterogeneity at the sequence level could add a layer of complexity to the overall diversity of the parasite's eventual phenotype. We speculate that it could be another source of variation characteristic of the parasite, conceivably arising from transcription errors, allowing a population of genetically identical cells to be phenotypically plastic to stresses or challenges, and facilitate a bet-hedging strategy in the face of various stressors. Recent work showed that the transcriptional and translational machinery of *P. falciparum* could handle the AT-richness of its genome effectively (42). This fact, taken together with our data would seem to indicate that most of the SNVs we observed, while they might be reflections of an inherent transcriptional (i.e. RNA Pol II mediated) error rate, are likely

not entirely explicable by simply the relatively lower complexity of the parasite's transcriptome. Therefore, we anticipate that our observations, which may be generalisable for other pathogenic parasite, will lead to further investigations into the exact source(s) and consequences of the pervasive SNVs that seem characteristic of *P. falciparum*, with regard to its basic biology, possible clinical implications of such variation, and its potential interplay with the previously reported phenomena of gene-expression level variation as well as structural variations in the *Plasmodium* transcriptome.

DATA AVAILABILITY

All custom code written for the analysis described herein as well as for generating figures is deposited at <https://github.com/bruhad-dave/Contextualize-SNVs>. REDItools 2.0 (25) is available at <https://github.com/tizianoflati/reditools2.0>. Other code used in this work is included in Supplementary Information. Raw RNAseq data, along with the corresponding final list (filtered, deduplicated) of variants as REDItools 2.0 output tables as well as raw gene abundance estimates as calculated by Salmon are deposited in GEO with accession GSE179055. Raw WGS data is deposited in SRA with BioProject ID PRJNA741726.

SUPPLEMENTARY DATA

Supplementary Data are available at NARGAB Online.

ACKNOWLEDGEMENTS

We thank Prof. Sutirth Dey for his valuable inputs and suggestions on this work. We are thankful to PARAM Brahma High Performance Computing facility at IISER, Pune for their support. The following reagents were obtained through BEI Resources (www.mr4.org), NIAID, NIH: *Plasmodium falciparum*, Strain IPC 3445 (MRA-1236), Strain IPC 5202 (MRA-1240), Strain IPC 4912 (MRA-1241), contributed by Didier Menard.

Author contributions: B.D. designed, performed experiments, and analyzed data. A.K. and D.V.M. cultured *P. falciparum* and generated NGS data. B.D. and K.K. wrote the manuscript. K.K. planned, coordinated, and supervised the project. All authors read and approved the final manuscript

FUNDING

DBT-Genome Engineering Technologies program [BT/PR25858/GET/119/169/2017 to KK] from the Government of India. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Conflict of interest statement. None declared.

REFERENCES

- Mackinnon, M.J., Li, J., Mok, S., Kortok, M.M., Marsh, K., Preiser, P.R. and Bozdech, Z. (2009) Comparative transcriptional and genomic analysis of *Plasmodium falciparum* field isolates. *PLoS Pathog.*, **5**, e1000644.
- World Malaria Report (2020) <https://www.who.int/teams/global-malaria-programme/reports/world-malaria-report-2020>.

3. Rovira-Graells, N., Gupta, A.P., Planet, E., Crowley, V.M., Mok, S., Ribas de Pouplana, L., Preiser, P.R., Bozdech, Z. and Cortés, A. (2012) Transcriptional variation in the malaria parasite *Plasmodium falciparum*. *Genome Res.*, **22**, 925–938.
4. Rawat, M., Srivastava, A., Johri, S., Gupta, I. and Karmodiya, K. (2021) Single-cell RNA sequencing reveals cellular heterogeneity and stage transition under temperature stress in synchronized *Plasmodium falciparum* cells. *Microbiol. Spectr.*, **9**, e00008–e00021.
5. Tarr, S.J., Díaz-Ingelmo, O., Stewart, L.B., Hocking, S.E., Murray, L., Duffy, C.W., Otto, T.D., Chappell, L., Rayner, J.C., Awandare, G.A. *et al.* (2018) Schizont transcriptome variation among clinical isolates and laboratory-adapted clones of the malaria parasite *Plasmodium falciparum*. *BMC Genomics*, **19**, 894.
6. Reid, A.J., Talman, A.M., Bennett, H.M., Gomes, A.R., Sanders, M.J., Illingworth, C.J.R., Billker, O., Berriman, M. and Lawniczak, M.K. (2018) Single-cell RNA-seq reveals hidden transcriptional variation in malaria parasites. *Elife*, **7**, e33105.
7. Scherf, A., Lopez-Rubio, J.J. and Riviere, L. (2008) Antigenic variation in *Plasmodium falciparum*. *Annu. Rev. Microbiol.*, **62**, 445–470.
8. Yeoh, L.M., Goodman, C.D., Mollard, V., McHugh, E., Lee, V.V., Sturm, A., Cozijnsen, A., McFadden, G.I. and Ralph, S.A. (2019) Alternative splicing is required for stage differentiation in malaria parasites. *Genome Biol.*, **20**, 151.
9. Sorber, K., Dimon, M.T. and DeRisi, J.L. (2011) RNA-Seq analysis of splicing in *Plasmodium falciparum* uncovers new splice junctions, alternative splicing and splicing of antisense transcripts. *Nucleic Acids Res.*, **39**, 3820–3835.
10. Gabriel, H.B., de Azevedo, M.F., Palmisano, G., Wunderlich, G., Kimura, E.A., Katzin, A.M. and Alves, J.M.P. (2015) Single-target high-throughput transcription analyses reveal high levels of alternative splicing present in the FPPS/GGPPS from *Plasmodium falciparum*. *Sci. Rep.*, **5**, 18429.
11. Iriko, H., Jin, L., Kaneko, O., Takeo, S., Han, E.-T., Tachibana, M., Otsuki, H., Torii, M. and Tsuboi, T. (2009) A small-scale systematic analysis of alternative splicing in *Plasmodium falciparum*. *Parasitol. Int.*, **58**, 196–199.
12. Hoo, R., Bruske, E., Dimonte, S., Zhu, L., Mordmüller, B., Sim, B.K.L., Krensner, P.G., Hoffman, S.L., Bozdech, Z., Frank, M. *et al.* (2019) Transcriptome profiling reveals functional variation in *Plasmodium falciparum* parasites from controlled human malaria infection studies. *EBioMedicine*, **48**, 442–452.
13. Jr, D.A.M., Pochet, N., Krupka, M., Williams, C., Seydel, K., Taylor, T.E., Peer, Y.V., de Regev, A., Wirth, D., Daily, J.P. *et al.* (2012) Transcriptional profiling of *Plasmodium falciparum* parasites from patients with severe malaria identifies distinct low vs. high parasitemic clusters. *PLoS One*, **7**, e40739.
14. Martins, B.M.C. and Locke, J.C.W. (2015) Microbial individuality: how single-cell heterogeneity enables population level strategies. *Curr. Opin. Microbiol.*, **24**, 104–112.
15. Goldman, S.L., MacKay, M., Afshinnekoo, E., Melnick, A.M., Wu, S. and Mason, C.E. (2019) The impact of heterogeneity on single-cell sequencing. *Front. Genet.*, **10**, 8.
16. Li, W. and Lynch, M. (2020) Universally high transcript error rates in bacteria. *Elife*, **9**, e54898.
17. Traverse, C.C. and Ochman, H. (2016) Conserved rates and patterns of transcription errors across bacterial growth states and lifestyles. *Proc. Natl Acad. Sci. U.S.A.*, **113**, 3311–3316.
18. Gout, J.-F., Li, W., Fritsch, C., Li, A., Haroon, S., Singh, L., Hua, D., Fazelinia, H., Smith, Z., Seeholzer, S. *et al.* (2017) The landscape of transcription errors in eukaryotic cells. *Sci. Adv.*, **3**, e1701484.
19. Gout, J.-F., Thomas, W.K., Smith, Z., Okamoto, K. and Lynch, M. (2013) Large-scale detection of in vivo transcription errors. *Proc. Natl Acad. Sci. U.S.A.*, **110**, 18584–18589.
20. Liscovitch-Brauer, N., Alon, S., Porath, H.T., Elstein, B., Unger, R., Ziv, T., Admon, A., Levanon, E.Y., Rosenthal, J.J.C. and Eisenberg, E. (2017) Trade-off between transcriptome plasticity and genome evolution in cephalopods. *Cell*, **169**, 191–202.
21. Licht, K., Kapoor, U., Amman, F., Picardi, E., Martin, D., Bajad, P. and Jantsch, M.F. (2019) A high resolution A-to-I editing map in the mouse identifies editing events controlled by pre-mRNA splicing. *Genome Res.*, **29**, 1453–1463.
22. Levitsky, L.I., Kliuchnikova, A.A., Kuznetsova, K.G., Karpov, D.S., Ivanov, M.V., Pyatnitskiy, M.A., Kalinina, O.V., Gorshkov, M.V. and Moshkovskii, S.A. (2019) Adenosine-to-Inosine RNA editing in mouse and human brain proteomes. *Proteomics*, **19**, 1900195.
23. Gott, J.M. and Emeson, R.B. (2000) Functions and mechanisms of rna editing. *Annu. Rev. Genet.*, **34**, 499–531.
24. Blow, M., Futreal, P.A., Wooster, R. and Stratton, M.R. (2004) A survey of RNA editing in human brain. *Genome Res.*, **14**, 2379–2387.
25. Flati, T., Gioiosa, S., Spallanzani, N., Tagliaferri, I., Diroma, M.A., Pesole, G., Chillemi, G., Picardi, E. and Castrignanò, T. (2020) HPC-REDIttools: a novel HPC-aware tool for improved large scale RNA-editing analysis. *BMC Bioinf.*, **21**, 353.
26. Radfar, A., Méndez, D., Moneriz, C., Linares, M., Marín-García, P., Puyet, A., Díez, A. and Bautista, J.M. (2009) Synchronous culture of *Plasmodium falciparum* at high parasitemia levels. *Nat. Protoc.*, **4**, 1899–1915.
27. Rawat, M., Kanyal, A., Sahasrabudhe, A., Vembar, S.S., Lopez-Rubio, J.-J. and Karmodiya, K. (2021) Histone acetyltransferase pf_{gcn5} regulates stress responsive and artemisinin resistance related genes in *Plasmodium falciparum*. *Sci. Rep.*, **11**, 852.
28. Andrews, S. (2010) FastQC: a quality control tool for high throughput sequence data [online]. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>, (last accessed: March 14, 2021).
29. Krueger, F. (2021) Trim Galore: a wrapper tool around Cutadapt and Fastqc to consistently apply quality and adapter trimming to Fastq files, with some extra functionality for MspI-digested Rbbs-type (reduced representation Bisulfite-seq) libraries. https://www.bioinformatics.babraham.ac.uk/projects/trim_galore/, (last accessed: March 14, 2021).
30. Martin, M. (2011) Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal*, **17**, 10–12.
31. Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M. and Gingeras, T.R. (2013) STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, **29**, 15–21.
32. Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R. and 1000 Genome Project Data Processing Subgroup (2009) The sequence alignment/map format and SAMtools. *Bioinformatics*, **25**, 2078–2079.
33. Li, H. (2011) A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*, **27**, 2987–2993.
34. Li, H. (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. bioRxiv: <https://arxiv.org/abs/1303.3997>, 18 May 2021, preprint: not peer reviewed.
35. Cingolani, P., Platts, A., Wang, L.L., Coon, M., Nguyen, T., Wang, L., Land, S.J., Lu, X. and Ruden, D.M. (2012) A program for annotating and predicting the effects of single nucleotide polymorphisms, snpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)*, **6**, 80–92.
36. Cingolani, P., Patel, V.M., Coon, M., Nguyen, T., Land, S.J., Ruden, D.M. and Lu, X. (2012) Using *Drosophila melanogaster* as a model for genotoxic chemical mutational studies with a new program, *SnpSift*. *Front. Genet.*, **3**, 35.
37. Patro, R., Duggal, G., Love, M.I., Irizarry, R.A. and Kingsford, C. (2017) Salmon provides fast and bias-aware quantification of transcript expression. *Nat. Methods*, **14**, 417–419.
38. SRA Toolkit Development Team (2021) <https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software>, (last accessed: March 14, 2021).
39. Zhang, H., Meltzer, P. and Davis, S. (2013) RCircos: an R package for circos 2D track plots. *BMC Bioinf.*, **14**, 244.
40. Signorell, A., Ken Aho, K., Alfons, A., Anderegg, N., Aragon, T., Arachchige, C., Antti Arppe, A., Baddeley, A., Barton, K. *et al.* (2021) DescTools: tools for descriptive statistics. R package version 0.99.44, <https://cran.r-project.org/package=DescTools>, (last accessed: December 02, 2021).
41. Wahlstedt, H. and Ohman, M. (2011) Site-selective versus promiscuous A-to-I editing. *Wiley Interdiscip. Rev. RNA*, **2**, 761–771.
42. Pavlovic Djuranovic, S., Erath, J., Andrews, R.J., Bayguinov, P.O., Chung, J.J., Chalker, D.L., Fitzpatrick, J.A., Moss, W.N., Szczesny, P. and Djuranovic, S. (2020) *Plasmodium falciparum* translational machinery condones polyadenosine repeats. *Elife*, **9**, e57799.
43. Meer, K.M., Nelson, P.G., Xiong, K. and Masel, J. (2020) High transcriptional error rates vary as a function of gene expression level. *Genome Biol. Evol.*, **12**, 3754–3761.