

## AccuNGS: detecting ultra-rare variants in viruses from clinical samples

Maoz Gelbart<sup>1,#</sup>, Sheri Harari<sup>1,#</sup>, Ya'ara Ben-Ari<sup>1</sup>, Talia Kustin<sup>1</sup>, Moran Meir<sup>1</sup>, Danielle Miller<sup>1</sup>, Orna Mor<sup>2,3</sup> and Adi Stern<sup>1,\*</sup>

<sup>1</sup> School of Molecular Cell Biology and Biotechnology, George S. Wise Faculty of Life Sciences, Tel Aviv University, Tel Aviv, Israel

<sup>2</sup> Central Virology Laboratory, Ministry of Health, Sheba Medical Center, Ramat-Gan, Israel.

<sup>3</sup> Sackler School of Medicine, Tel Aviv University, Tel Aviv, Israel

# Authors equally contributed

\* To whom correspondence should be addressed. Tel: 972-3-6407508; Fax: 972-3-6405984; Email: [sternadi@tau.ac.il](mailto:sternadi@tau.ac.il)

### ABSTRACT

Next generation sequencing is widely used to characterize genetic diversity in a sample, yet is hindered by its relatively low resolution. Particularly, detecting rare genetic variants in clinical samples of viruses is still nearly impossible. Here we describe AccuNGS, an approach that combines error reduction in each sequencing stage with *in silico* error elimination, which enables detection of variants as rare as 1:10,000 or lower. We thoroughly explore AccuNGS background errors and reveal they are mostly generated in the sequencer itself. We demonstrate that as opposed to common assumptions, Illumina paired-end reads are not independent. After applying AccuNGS to an HIV sample taken during acute infection, we reveal that the vast majority of transition variants in the sample segregate at ultra-low frequencies, rendering them undetectable by standard sequencing. These results highlight the early rich accumulation of genetic diversity during viral infection at depths previously unseen.

### INTRODUCTION

Recent advances in high-throughput nucleic acid sequencing have revolutionized our ability to identify the prevalence of minor traits in a heterogeneous sample. Identification of rare single nucleotide variants (SNVs) is important in diverse disciplines spanning post-transcriptional modifications, cancer genetics, non-invasive prenatal diagnoses and microbiology (1-4). SNV identification in virus populations is currently at the heart of many studies monitoring drug resistance, estimating mutation rates, quantifying standing genetic variation and predicting the fitness costs of single mutations (4-10). Accountable quantification of such variants present in clinical specimens requires high template recovery, sufficient sequencing depth, and

discrimination of real minor variants from the background errors of the sequencing process (11,12). However, using the standard next generation sequencing (NGS) protocols may result in significant background error rates. In fact, following typical post-processing of NGS data, mutations that are at frequencies lower than 1-5% are discarded, drastically limiting rare SNV identification (9,13,14). This limitation of NGS has been recently pointed as one of the major gaps in using genotyping to survey resistance mutations and understand HIV treatment failure (15,16).

In the past few years, several innovative experimental approaches were suggested to reduce the background error rates of the NGS process: rolling-circle-based redundant coding of the amplified fragments (1,17-19); consensus sequencing of barcoded genomic fragments (19-24); error reduction by overlapping paired reads in paired-end sequencing (25-27); and usage of improved polymerases (28). Complementary to the library preparation methods, several computational methods were created to facilitate discrimination of true variants from process errors, based on systematic background error modeling (24,29-33). Apart from the usage of overlapping read pairs (ORP), most experimental methods described above are designed for samples with high biomass and are inapplicable for sequencing of viruses from clinical samples, where the biomass of the viruses may be extremely low. Furthermore, these experimental protocols may provide accuracy at the cost of increased technical complexity, and may introduce their own artifacts to the sequencing process (34). On the variant calling side, most variant calling programs model strand-specific sequencing bias but do not properly incorporate the information from the paired region that may have different error characteristics (35). Moreover, it has been suggested that these well-established variant callers do not perform well on clinical virus samples (36).

We therefore sought to create a simple and highly accurate sequencing protocol that will be suitable for sequencing of clinical samples, with a special focus on low-biomass samples of RNA viruses. The motivation was based on the fact that many of these viruses replicate at huge census population sizes with high mutation rates of  $\sim 10^{-4}$ - $10^{-5}$  mutations/base/replication, so viral populations from clinical samples are expected to contain very high levels of heterogeneity (5,37,38). We set out to perform a step-by-step optimization of NGS protocols with an emphasis in each step on high accuracy and high yield. The resulting optimized protocol termed "AccuNGS" is a simple and rapid sequencing method and includes an associated variant caller, which if combined can reliably detect ultra-rare variants at frequencies of 1:10,000 or lower. Using homogenous DNA and RNA samples, we characterized the typical error landscape of the AccuNGS protocol, pinpointed the potential sources that generate errors in our protocol and suggest potential solutions. We applied our method to an RNA sample taken from a patient recently infected with HIV-1 (acute stage, seronegative) to examine the breadth of accumulation of mutations in the viral population in this critical period of the infection.

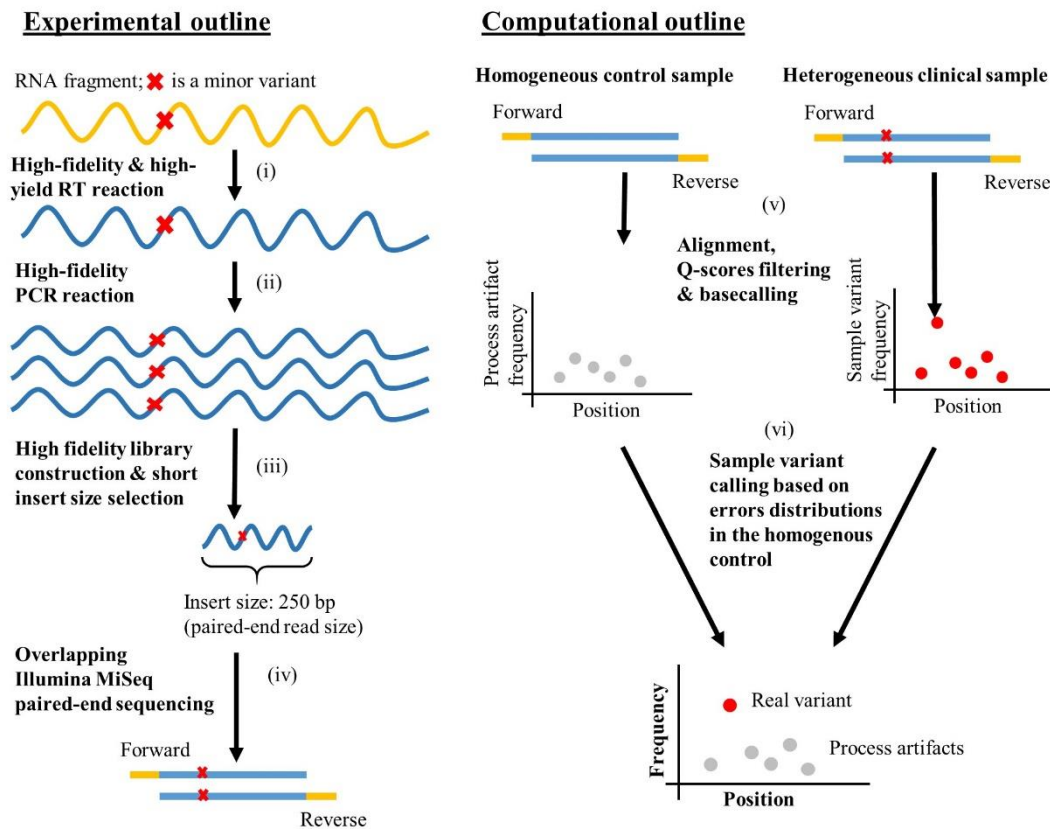
## RESULTS

### Protocol overview

AccuNGS protocol was designed to enable targeted sequencing of a desired genomic region from clinical samples with maximal yield, high accuracy and rapid implementation. It combines several concepts, some of which were previously reported separately. The principles of our protocol are (i) use of high-fidelity polymerases during reverse transcription (RT) and amplification. In particular we chose to use either the SuperScript III or SuperScript IV RT enzymes (reported mean error rate near  $5 \times 10^{-5}$  (29,39)) and the Platinum SuperFi DNA polymerase (error rate of near  $5 \times 10^{-7}$ , (40) and manufacturer at <https://www.thermofisher.com/>); (ii) significant reduction of sequencer errors by overlapping paired-end sequencing; and (iii) a specialized base calling bioinformatics pipeline that incorporates mutation-specific and locus-specific distribution of background errors (see Methods and Fig. 1). Based on the above information we were able to calculate *theoretical* means of  $1.78 \times 10^{-5}$  and  $6.78 \times 10^{-5}$  errors per base introduced in our AccuNGS protocol for DNA and RNA samples, respectively (Table 1). An auxiliary part of our protocol allows for quantification of the actual number of RNA genomes sequenced using uniquely barcoded primers introduced early in the RT step (the "primer-ID" method) (2,20,21,41,42). This is a critical measure when analyzing the levels of diversity present in clinical samples, since low genetic diversity observed in a sample may stem from a small number of sequenced templates rather than from real reduced diversity in the sample. We note that this manuscript focuses on two measures: the mean error rate, reported when we characterize the method, and a cutoff error rate (based on the gamma distribution of errors), which we report as a measure to be used when performing base-calling. Naturally the cutoff value is higher than the mean error.

**Table 1. Mean expected error rates in AccuNGS.**

Step	Error rate	Number of rounds	Expected error
Polymerase Chain Reaction (PCR)	$6.47 \times 10^{-7}$ [ThermoFisher]	40 + 12	$1.68 \times 10^{-5}$
Sequencing	$10^{-6}$ (Q30 filtering & overlapping paired reads)	1	$1 \times 10^{-6}$
<b>Total (DNA starting material)</b>			<b><math>1.78 \times 10^{-5}</math></b>
Reverse Transcriptase (RT)	$3.1 \times 10^{-5} - 6.5 \times 10^{-5}$ (29,39)	1	$4.8 \times 10^{-5}$
<b>Total (RNA starting material)</b>			<b><math>6.58 \times 10^{-5}</math></b>



**Fig. 1. AccuNGS protocol principles.** (i) High fidelity and high-yield RT reaction. The RT primers may be designed with a unique N-base barcode (e.g. "primer-ID") to allow template quantification downstream (ii) High-fidelity PCR reaction (iii) Library construction with size selection for insert size as short as a single paired-end read (iv) Paired-end sequencing where each base in the insert is sequenced twice, once in the forward read and again in the reverse read (v) Alignment of reads, Q-score filtering on both reads and basecalling of both the sample and a homogeneous control (vi) Sample variant calling based on fitted distributions of process errors and position-specific error propensity.

#### AccuNGS error sources analysis at the DNA level

We began by evaluating AccuNGS when a DNA plasmid was used as starting material. Our underlying assumption throughout our working process was that our DNA starting material is homogenous with respect to the theoretical error rate we calculated. This assumption is based on the fact that we used low-copy plasmids that were grown in *E. coli*, and only a single colony was subsequently sequenced. The mutation rate of *E. coli* is in the order of  $1 \times 10^{-10}$  errors/base/replication (23), and sequencing of a single colony ensures only a limited number of replication cycles. Accordingly, error rates in the purified plasmids are expected to be much lower than the expected protocol mean error of  $\sim 1 \times 10^{-5}$ . Thus, errors observed when comparing the

results of the sequencing to our known reference sequence reflect errors created by the library preparation or by the sequencing process itself. All samples underwent basecalling using our specialized bioinformatics pipeline and positions were considered for analysis only if their coverage exceeded 100,000 bases per position (see Methods). Table S1 provides statistics about the number of reads and the distribution of miscalls in each sample.

We thus set out to sequence the HIV-1 pLAI.2 plasmid (43) using AccuNGS at its baseline conditions, including 40 cycles of PCR amplification of a target region with a high-fidelity polymerase, followed by Nextera XT library preparation with a high-fidelity polymerase and size selection of a 250bp insert (see Methods and Table 2). We then compared the baseline AccuNGS error rate to the results of a protocol typically used in clinical (and other) settings, where less focus is put on the fidelity of the process (44). Fig. 2 compares the proportion of errors observed on our plasmid sequence under AccuNGS and under standard sequencing, broken down according to type of mutation (base-to-base). Reassuringly, AccuNGS showed a significant improvement of one to two orders of magnitude over the standard sequencing protocol. For example, the mean A>G error rate went down from  $2.6 \times 10^{-3}$  down to around  $9.2 \times 10^{-5}$ . While this improvement was large, perplexingly it was still an order of magnitude higher than the theoretical error rate we had expected of  $\sim 1 \times 10^{-5}$ . We thus set out to optimize the AccuNGS protocol and try to pinpoint the unexpected source introducing errors into the process.

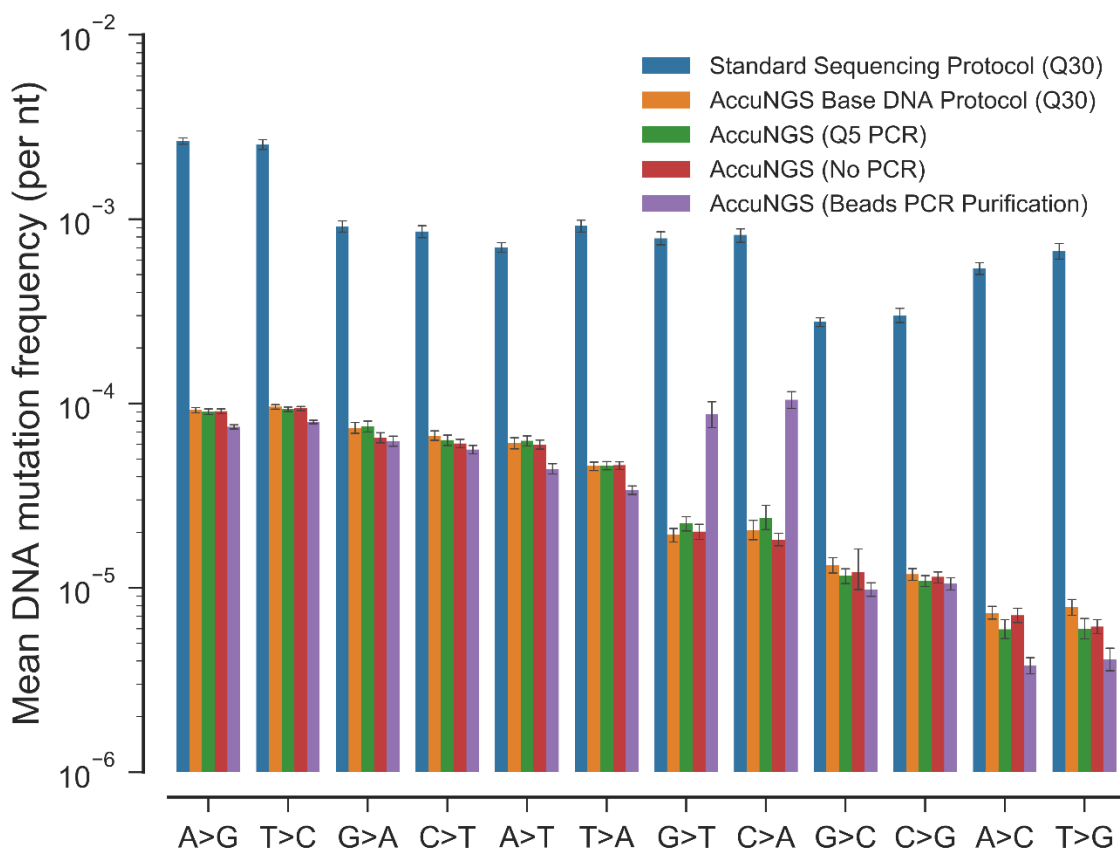
**Sources of process errors based on differential sequencing.** We next performed a set of sequencing trials, whereby at each trial we tested if removing or changing a specific stage of the protocol alleviates some of the observed errors and improves the fidelity of AccuNGS. Our immediate suspect was the PCR amplification step. Due to the exponential nature of PCR, errors introduced at early stages of the amplification will be carried over and have been reported to create a high background error for ultra-deep sequencing (45,46). Accordingly, we were worried that any misspecification of the error rate of the SuperFi DNA Polymerase we use (Table 1) would lead to an inflation in PCR errors. To test this hypothesis we created a sample with no PCR amplification, by harvesting larger quantities of the plasmid from the bacteria. This led to only a slight decrease in the mean error rate of transition mutations (Fig. 2). We hence concluded that the forty PCR cycles that take place in the PCR amplification prior to library preparation do not explain most of the errors of the AccuNGS protocol.

**Table 2. Differential DNA control samples sequenced in this study.**

Sample name	PCR enzyme	PCR cycles	PCR purification method	Library purification method	Cloning bacteria	Tagmentation method	Target region
Baseline	SuperFi (Thermofisher)	40+12	Gel	Beads	<i>E.coli</i> (DH5α)	Nextera XT	Integrase (pLAI)
Q5	Q5 (NEB)	.	.	.	.	.	.
PCR free	.	12	.	.	.	.	Integrase (extended)
Alt. PCR purification (beads)	.	.	Beads	.	.	.	.
Alt. PCR purification (Exosap)	.	.	Exosap	.	.	.	.
Alt. library purification (gel)	.	.	.	Gel	.	.	.
TG1	.	.	.	.	<i>E.coli</i> (TG1)	.	.
Alt. tagmentation	.	.	.	.	.	PCR	.
AmpR	.	.	.	.	.	.	<i>AmpR</i> (pLAI)
RpoB	.	.	.	.	.	.	<i>RpoB</i> (DH5α)
NextSeq (Illumina)	.	.	.	.	.	.	.

"." in a cell indicates a condition that is similar to the baseline; Alt=alternative. All samples besides the NextSeq sample were sequenced on the Illumina MiSeq sequencing platform.

We next focused on various so-called “chemical” processes that take place in AccuNGS library preparation. Mainly we were concerned that using gel extraction for DNA size selection, particularly UV light exposure, may introduce mutations. Indeed, when replacing gel extraction for the PCR products with magnetic beads extraction, we observed a slight reduction of the mean transition error rates (Fig. 2). On the other hand, this sample showed elevated levels of C:G>A:T errors. These errors are often signatures of oxidative stress and are discussed below. Alternative PCR purification using the Exosap cleanup reagent did not show elevated C:G>A:T errors compared to their levels in the baseline protocol, however it showed error levels comparable to the baseline protocol.



**Fig. 2. Mean background error rates of different sequencing protocols at the DNA level.**

AccuNGS dramatically reduces errors present in standard sequencing protocols by almost two orders of magnitude. PCR errors in AccuNGS are negligible based on comparison of PCR and PCR-free sample. Higher rates of G>T and C>A are likely indicative of oxidation (see text for more details). Error bars represent 95% confidence intervals around estimated mean values using 1,000 bootstrap repeats.

We next tested if the source plasmid itself was the major source of observed errors, focusing on the conditions whereby we grew the plasmid. First we tested if the mutations were accumulated naturally due to lack of selection on the HIV genes by sequencing the antibiotic resistance marker *AmpR* on pLAI.2 that is presumably under strong selection against mutations. Next we tested if the plasmid was the source of errors by sequencing the highly conserved *RpoB* gene from *E. coli* itself. Finally we tested if the errors were introduced by the bacteria during plasmid replication by growing the plasmid in an alternative strain of *E. coli* (TG1) with a presumably lower mutation rate (47). However, we observed no change in the error rate distribution in any of these conditions, suggesting that the DNA input was not the major source of minor variants.

We next hypothesized that the tagmentation process in the NexteraXT DNA library preparation kit (Illumina) might be the cause of artifacts. We resorted to a home-made tagmentation protocol

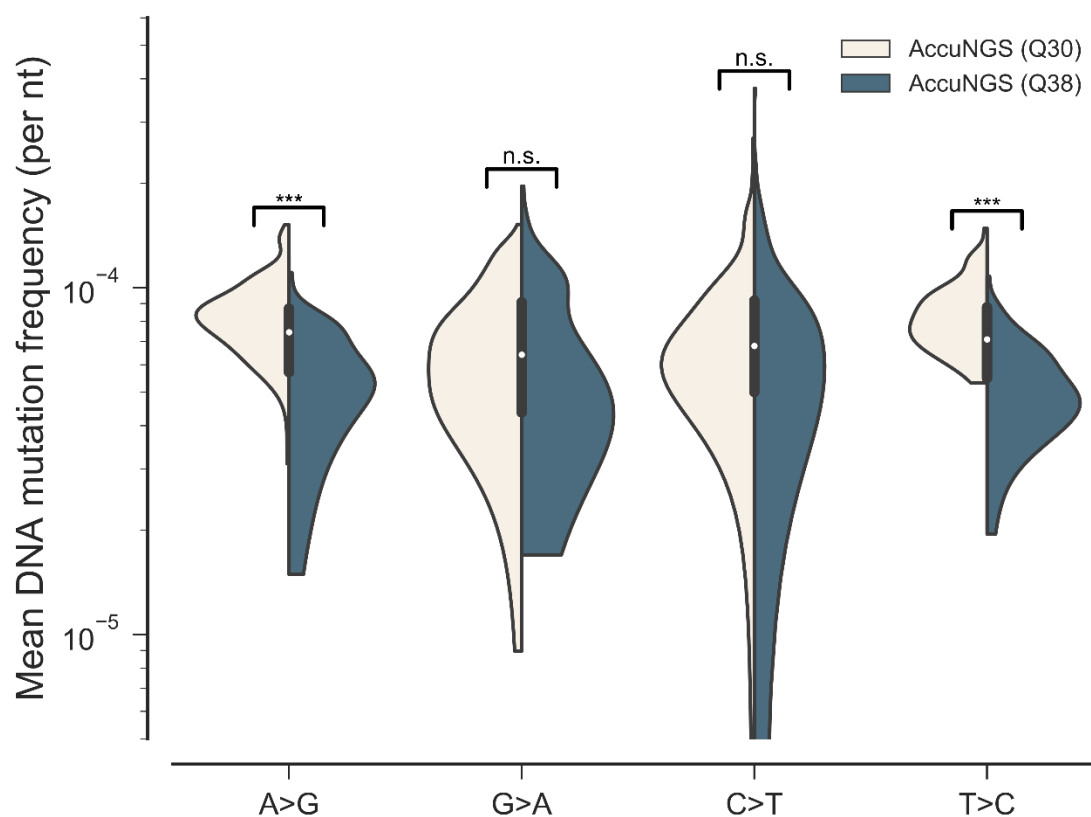
based on introduction of NexteraXT-compatible adapters and indices via PCR amplification of a 250bp fragment of pLAI.2. Yet again, we observed no significant change in the distribution of errors. Replacing the Illumina MiSeq sequencer with the Illumina NextSeq, which is based on a two-channel sequencing process rather than a four-channel process and a different flow cell, resulted in a small increase in error levels. This suggests that the Illumina MiSeq and the higher throughput HiSeq (which employs the same detection method as MiSeq) are slightly more suitable for AccuNGS sequencing.

After having ruled out all sources of error we could conceive of testing, we were still left with the enigma of what causes the errors observed consistently and reproducibly across all the samples we sequenced. We were left with one condition that we could not alter directly: the sequencing step itself, as discussed next.

**Sequencing quality effect.** Each base reported by Illumina sequencers is assigned with a probability of that base being wrong, termed the Q-score. The range of Q-scores reported from the Illumina MiSeq is 0 to 40, which translates to a probability of an erroneous call between 1 and 0.0001, respectively. In the AccuNGS base calling scheme, we consider only sites where the two overlapping reads reported the same base with an average Q-score of 30 or higher, as in (25). Our original interpretation of overlapping reads, in line with previous works (48,49) was that the base called jointly on both reads has a corrected Q-score equals approximately to the sum of the Q-scores from both reads. Accordingly, this means that if we filter for bases with a corrected Q-score of at least 60, this translates to an error probability  $\leq 1 \times 10^{-6}$  per base called, far below our theoretical threshold of  $1.78 \times 10^{-5}$ . We set out to test if this is indeed true. First, we determined whether the independent Q-scores indeed reflect what they are supposed to. When inspecting errors observed on one read only with a Q-score of 30 or higher, we found that their maximal frequency was indeed around  $10^{-3}$  (Fig. S3). Thus, it seems that the individual Q-scores on each read are reflective of the error rates of the process. However, we suspected that the joint Q-scores that we calculate are incorrect – mainly, that each reported Q-score is not independent of the Q-score on the mate read. To test this, we examined whether using a more stringent filter criterion improves AccuNGS results. We applied a very stringent quality filtering of Q38 on the AmpR sample that was sequenced to extreme depth of 1,500,000 bases per position (Table S1). The very high coverage and the good quality of the sequencing allowed us to still retain most sites at a coverage of above 100,000 reads per base (Table S1). We expected that this filtering will improve the results by the difference between twice Q30 (Q60, error probability of  $1 \times 10^{-6}$ ) and twice Q38 (Q76, error probability of  $\sim 2 \times 10^{-8}$ ). This difference translates to an improvement which is far below our observed error rate and hence we did not expect to see any improvement. Surprisingly, we observed a dramatic reduction in the rates of errors for A:T>G:C miscalls, and a modest reduction for C:G>T:A miscalls (Fig. 3, Table 3 and Table S2). We hence concluded that



the assumption that the Q-scores of overlapping reads are independent is an incorrect assumption, and that the sequencer itself is the major source of errors in AccuNGS.



**Fig. 3. The effect of increasing the Q-score filtering threshold on AccuNGS error rates.**

Error distributions for Q30 and Q38 filtering for the AmpR amplicon, presented for each type of transition error. n.s., not significant; \*\*\* $p < 0.0001$

**Effects of surrounding nucleotides on error rates.** Previous studies have indicated that the nucleotides surrounding a called base may influence its propensity to be miscalled (13,28). Indeed, we found that the identity of the surrounding bases sometimes affected the error rate observed with AccuNGS: when focusing on G>A mutations, we observed a higher error when the G was preceded by a C (CpG) and a lower error when the G was preceded by an A (ApG, Fig. S1). The exact same pattern was observed for the reverse complement C>T mutations (data now shown). These phenomena were found in all sequenced samples, suggesting that they aren't affected by any of the differential conditions we tested. When analyzing transversion artifacts, we observed a higher error rate for G:C>T:A mutations compared to all other types of transversions. This enrichment was more prominent in some samples than in the others and is typically associated with oxidative damage (13,50-53). When characterizing the nucleotide context of the C:G>A:T errors, we found that G>T errors occurred more frequently when the mutated G base

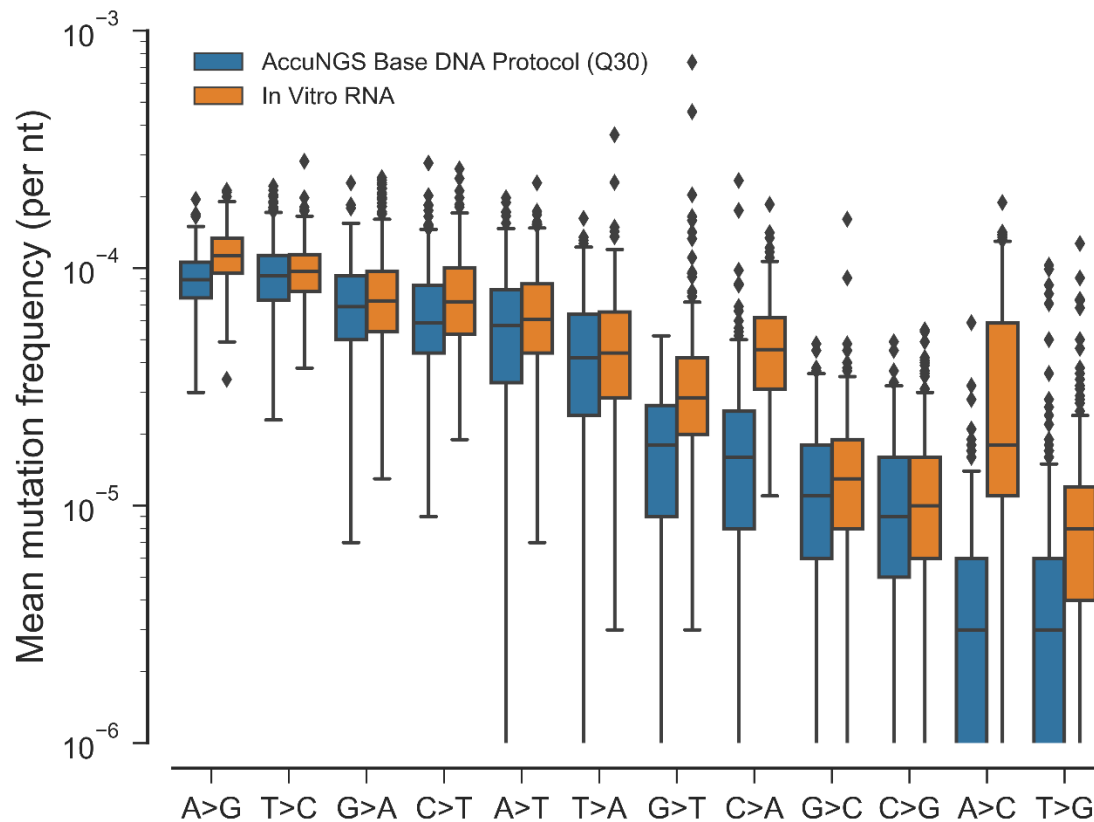
**Table 3. Fitted gamma distributions for different Q-scores cutoffs and substitution types, on the AmpR amplicon.**

Q-score Cutoff	Substitution type	Number of Sites (N)	Shape ( $\kappa$ )	Scale ( $\theta$ )	0.95% of fitted gamma	Mean error rate	Median error rate
Q30	A->G	145	8.855585	9.31E-06	1.32E-04	8.24E-05	8.30E-05
	C->T	172	4.511632	1.69E-05	1.42E-04	7.61E-05	6.49E-05
	G->A	154	5.32	1.33E-05	1.27E-04	7.05E-05	6.69E-05
	T->C	170	12.8402	6.57E-06	1.26E-04	8.43E-05	8.02E-05
Q38	A->G	138	5.863207	8.51E-06	8.79E-05	4.99E-05	4.89E-05
	C->T	164	2.773564	2.49E-05	1.48E-04	6.92E-05	5.90E-05
	G->A	149	2.879898	2.17E-05	1.32E-04	6.25E-05	5.36E-05
	T->C	166	7.09643	7.01E-06	8.39E-05	4.98E-05	4.77E-05

was followed by A or another G (GpA\G). As in the G>A transitions, the reverse complement C>A mutations were more frequent when C was preceded by C or T (C\TpC; Fig. S2).

#### AccuNGS error analysis at the RNA level

Since one of the ultimate goals of the development of AccuNGS was the sequencing of RNA viruses, we set out to characterize how the protocol fares for RNA. In order to obtain a homogenous RNA sample we performed *in-vitro* transcription of a homogeneous plasmid using T7 polymerase, whose error rate has been approximated in the order of  $10^{-6}$  (54), an order of magnitude lower than the error rate observed for DNA with AccuNGS. The RNA was then used as input for reverse transcription reaction with random hexamers using SuperScript III, whose mean error rate has been approximated to be between  $3.1 \times 10^{-5}$  and  $6.5 \times 10^{-5}$  (29,39). We then proceeded with the AccuNGS protocol for DNA as previously described. In this RNA control sample, we expected the observed errors to be the union of those introduced by the DNA part, those introduced during *in-vitro* transcription and those introduced during reverse transcription. With Q30 filtering, the RNA control sample indeed showed a higher mean transition error rate of  $9.52 \times 10^{-5}$  compared to the mean transition error rate of  $8.49 \times 10^{-5}$  in the DNA control sample (Table S1 and Fig. 4). The difference of  $1.03 \times 10^{-5}$  is indeed in line with most additional errors in the RNA control sample stemming from the RT step. By using the difference between the medians of these two control samples, we were able to calculate upper bounds on the base-by-base error rates of the RT used in the process, which were found to be lower for some mutation types than previously reported (Table S3).



**Fig. 4. Error rates of the DNA versus RNA control samples.** Comparison of error distributions reveals that RT most often does not introduce a dramatically high level of error. Boxplots of errors per type of base changes are shown. Raw read bases were filtered when their average Q-score was less than Q30.

#### **Clinical sample sequencing – acute HIV-1 infection**

We next went on to test our method on direct sequencing of a clinical HIV-1 sample. HIV-1 infections typically begin with one to few viruses, indicating that the virus diversity at the population level at the time of infection is very limited (55). Large population sizes coupled with high mutation rates in the order of  $10^{-5}$  mutations/base/replication cycle allow the virus to obtain mutations shortly after infection (38,56). However, HIV-1 populations sequenced shortly after infection, while the patient is at acute infection, have shown very limited diversity (57-59). At this time point, most variants in HIV are expected to be at frequencies below  $10^{-3}$ , thus mostly obscured by the common clinical sequencing protocols' error rates. We obtained a plasma sample from a recently infected HIV-1 patient with laboratory confirmed seroconversion (a negative HIV-1 confirmatory assay followed by a positive test, 2 weeks apart), indicating this patient was likely 15-20 days after infection (so called acute HIV infection (60)). We chose to sequence the *gag* region of the virus (nearly 1800 bases) as it is mostly under purifying selection,

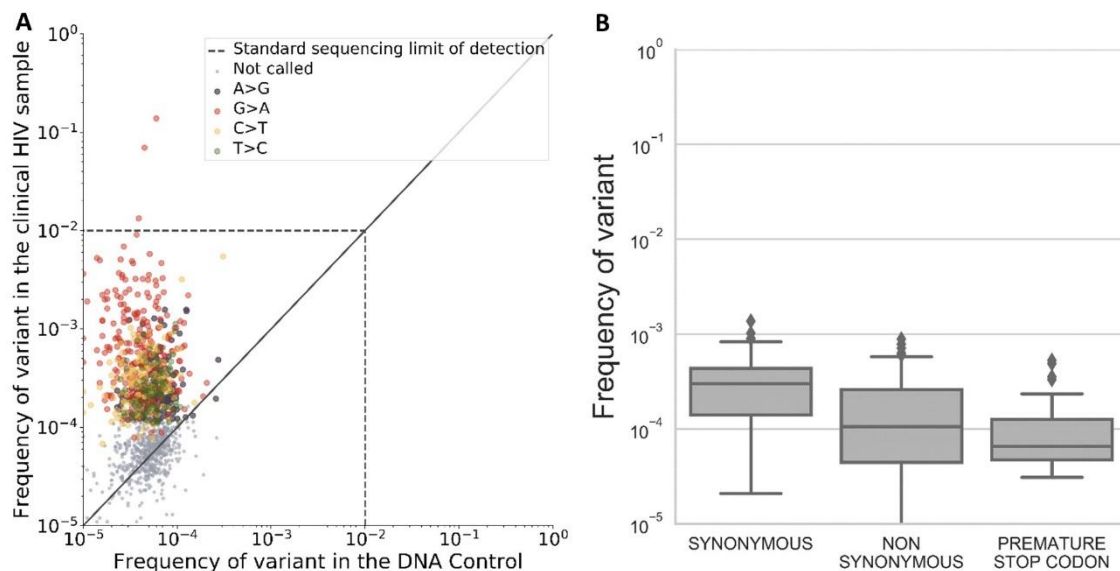
but is also targeted by the HLA component of the adaptive immune system and possibly HLA escape mutations will be seen at this early stage of infection (61). We prepared each RT primer with a unique barcode (“primer ID”), 15 nucleotides long that allowed us to quantify the amount of viruses we have actually sequenced (see Methods and Supplementary Text). We aimed to sequence ~30,000 viruses, which is roughly the inverse of the error rate we obtain with AccuNGS (which is around 1 in 10,000). To this end, we started the protocol with roughly 300,000 viruses, as estimated from the viral load of the sample. Based on RT processivity, we estimated that 10% of the viruses would hence be sequenced (62).

For background errors control, we amplified and sequenced a clonal pLAI.2 DNA in parallel. The sequenced sample and the control had a median coverage of nearly 400,000 bases per position at Q30, allowing us to filter called bases for a minimum of Q38 for each base in the Forward and Reverse reads. As expected, background error rates in the control plasmid were similar to those obtained by our previous controls. Primer-ID analysis revealed that nearly 16,000 viruses were actually sequenced, so variants identified at frequencies of  $10^{-4}$  are likely to represent true diversity (see Supplementary Text). We applied our variant caller on the sample using the pLAI control to serve as the background error distribution. Using a p-value cutoff of 0.01 for each variant called (i.e., this variant is at or above the 99<sup>th</sup> percentile of the gamma distribution for this mutation type), between 40% and 50% of all sequenced positions were identified as containing true transition mutations (Fig. 5A and Table 4). Using the standard sequencing, only *four* transition mutations would have been identified; AccuNGS revealed that several hundreds of transition variants exist at this time of infection at low frequencies of  $10^{-4}$ - $10^{-3}$  (Table 4). When analyzing the type of mutations observed at low frequencies, synonymous mutations were the most prevalent, then non-synonymous mutations and then nonsense mutations (Fig. 5B). This suggests that signals of purifying selection can be already captured at this early time-point of infection and also strongly suggests that the variation observed is true. Interestingly, G>A minor variants were more prevalent than all other mutation classes (Fig. 5). Specifically, G>A variants preceded by A (GpA) were the most prevalent among all G>A variants, followed by G>A variants preceded by G (GpG, Fig. S4). This is possibly evidence for cytosine deamination activity by host APOBEC enzymes on the minus strand during reverse transcription, reflected as excess G>A mutations in the genome of virus (63). This is also in line with the favored editing context of APOBEC3F (63). However, G>A mutations are also the most common replication error of HIV-RT (56), and we cannot rule out that this drives higher frequencies of G>A as observed here.

**Table 4. Variants identified with AccuNGS in the acute HIV sample.** P-values were calculated with AccuNGS variant caller using the specified control sample as control. Sites were considered if their coverage exceeded 25,000x.

Control	#Sites with sufficient coverage	Significance level	Number of variants (Ts+Tv)	Number of transition variants	% Positions with transition variants	Average transition variant frequency
pLAI.2 (DNA)	1,469	5%	1,535	799	54%	8.2E-04
pLAI.2 (DNA)	1,469	1%	1,297	727	49.4%	8.98E-04
In vitro RNA	1,469	5%	1,105	649	44.1%	9.9E-04
In vitro RNA	1,469	1%	984	594	40.4%	1.06E-03

Ts, Transitions; Tv, Transversions.



**Fig. 5. Transition minor variants identified in the acute HIV-1 sample.** (A) Variants plotted against their respective frequency on the DNA control. Variants are colored if they reside in the top 1% of the fitted error distributions based on the DNA control. The vast majority of HIV variants reside below the standard sequencing limit of detection of ~1% (dashed line). (B) Minor C>T transition variants by type. Synonymous mutations are more prevalent than non-synonymous and stop-codon-forming mutations.

## DISCUSSION

Application of next generation sequencing on clinical samples is still limited by the ability to reliably capture minor variants (15,16,36). Here we describe AccuNGS, a simple, rapid and optimized experimental protocol and associated computational pipeline for detecting ultra-rare variants from low-biomass clinical RNA and DNA samples. AccuNGS aims to accurately detect

minor variants present in a population of genomes at frequencies of 1:10,000 or lower, close to the mutation rate of RNA viruses (5). By performing differential sequencing we demonstrate that as opposed to many sequencing protocols, PCR is not a major source of errors in AccuNGS (45,64), and conversely we suggest that the sequencer is a major source for errors, even when correcting the strand-bias using overlapping read pairs. We show that the mean transition error rate of the protocol is around  $7.83 \times 10^{-5}$  when filtering for Q30, and  $5.80 \times 10^{-5}$  when filtering for Q38. These error rates translate to a cutoff of  $1.35 \times 10^{-4}$  and  $1.16 \times 10^{-4}$  based on the 95 percentile of a fitted gamma distribution of all transition errors. Notably when focusing on specific types of transitions (mostly T>C and A>G) the cutoff drops to below  $1 \times 10^{-4}$ .

AccuNGS excels especially when the input is low-biomass heterogeneous RNA. Comparable methods for accurate sequencing such as rolling-circle-based methods typically require extremely high-biomass input, making them irrelevant for clinical virus sequencing. Furthermore we and others have observed that such protocols exhibited a relatively high C>T rate [exceeding  $10^{-4}$  (1,17,18,65), unpublished results]. Such error levels were not recapitulated using AccuNGS, suggesting that these may have been artifacts of the rolling circle approach. A possible alternative for clinical sequencing would be the use of barcoded primers (also known as primer-IDs) during RT reaction, to generate consensus sequences that will correct errors inserted during amplification and sequencing (2). The advantage of the barcoded approaches is that they can also correct for unequal PCR sampling. However, the downside of these approaches is that they require splitting the input sample into numerous reactions, since a barcode has to be attached to each sequencing read (typically spanning 500-600 bases). When the number of initial viruses in the sample is not huge (as typical for limited clinical samples), this is problematic since only a relatively small number of viruses will be sequenced in each reaction. AccuNGS is only limited by the capacity of the RT and PCR reactions (i.e., the length of the targeted sequence that undergoes one RT or PCR reaction), which may span several thousand bases. We also noted that error rates in barcoding-based protocols are comparable to AccuNGS (35,66). We do note that in the AccuNGS approach we recommend primer barcoding only on one end of the amplicon, but this is aimed to understand how many RNA templates were actually sequenced rather than for error correction.

The overlapping read pairs (ORP) concept to reduce sequencing errors was first reported by Chen-Harris et al. (25) and further used by PELE-Seq (26). We note that our approach is novel in that it hinges on the combined use of high fidelity enzymes, ORP and a bioinformatics pipeline that compares variant frequencies to the fitted gamma distribution of process errors (24,27,67,68). Indeed, AccuNGS improves over Chen-Harris et al., and we further show that the use of a high fidelity polymerase is key to bringing down the mean transition error rate by approximately 30% (Fig. S5). We further provide here a step-by-step dissection and optimization

of the sequencing process. This has allowed us to refute the commonly assumed notion that the forward and reverse reads are independent and can be considered technical replicas [e.g. (48,49,69,70)]. The observed improvement in error rates of A:T>G:C transitions when using a more stringent Q-score filtering threshold strongly indicate that the Illumina chemistry, involving low-fidelity polymerases (71), is a major source of sequencer errors that are only partially reflected in the Q-scores, and are usually not modelled. This is supported by the observation that Taq DNA polymerase, which operates during cluster generation in the sequencer (71), tends to make more errors on A:T>G:C (46,64). The lack of improvement in the reciprocal G:C>A:T process errors in spite of more stringent filtering suggests these may stem from Cytosine deamination during thermal cycling, rather than from polymerization during cluster generation (72). Our results also demonstrate that MiSeq to some extent better suits AccuNGS than the newer NextSeq platform, although significant improvement of error rates was also seen on the latter platform. We suggest that the observed error rate of the Q38 samples approaches the limit of detection that can be achieved using any traditional Illumina sequencing protocol. Hopefully in the future sequencing vendors will create a "high-fidelity" sequencing program that incorporates high fidelity enzymes that could minimize process artifacts.

We further note that the modeling process errors has several advantages over position-specific error models (25,30). First, position-specific error models do not perform well when the consensus base in the sample differs from the consensus base in the background homogeneous control. Second, the stochastic nature of process errors at a given position may result in significant differences between technical replicates, highlighting the fragility of an error estimation based on a single base observation. When possible, we strongly recommend using a homogenous control that is as similar as possible to the samples at hand, since this directly allows detecting loci that are highly error prone. In the absence of such a control, any homogenous control (e.g., a plasmid) is useful to control for process errors.

Performing our benchmarking on relatively long genomic regions allowed us to find that some errors tend to occur more in specific contexts than in others. We find that C:G>T:A mutations, that often arise from spontaneous cytosine deamination during thermal cycling (53,64,72), tend to occur when in CpG context (for both C and G), whereas they are less frequent when in CpT (for C>T) and ApG (for G>A) contexts. Similarly, we find that C:G>A:T mutations, which often stem from the formation of 8-oxo-Guanine under oxidative stress (53), are more prevalent in specific contexts. AccuNGS incorporates this observed bias into its variant calling method, based on more-specific background errors distributions.

We used AccuNGS to characterize HIV-1 diversity shortly after infection. To the best of our knowledge, the immediate evolution of HIV-1 following a new infection has been rarely observed (57-59), mainly due to the lack of resolution associated with the common sequencing protocols.

We demonstrate that AccuNGS captures minor transition mutations that mostly segregate at frequencies between 1:100 and 1:10,000, and match the expected properties of real biological variants such as different rates for silent, missense and nonsense mutations. We found that in this sample, G>A mutations were three times more prevalent than C>U mutations (Fig. 5B). In spite of being the most abundant type of mutation induced by the viral reverse transcriptase (6), the high level of observed minor frequencies exceeded our predictions, given the early time point. This may be partially explained by activity of APOBEC3 editing enzymes acting on the antisense of the HIV reverse transcribed genome. The most common APOBEC3 signature observed in our clinical sample (characterized by GpA>ApA mutations) was that of APOBEC3A/D/F/H, unlike a previous study that highlighted the contribution of APOBEC3G to the observed G>A variants in HIV-1 proviral DNA (characterized by GpG>ApG, (73)). It is possible that in our patient there are variants of APOBEC3 that are more active (74), or that in this patient, the viral Vif protein (that encounters host APOBEC3 proteins) lost the activity against some APOBEC3 enzymes (75). Future studies spanning more patients may elucidate this issue.

## Summary

To summarize, we anticipate that using AccuNGS will be highly useful in detecting previously uncharacterized genetic diversity in biological samples. The ease of use of this approach should make it highly amenable for many different studies.

## MATERIALS AND METHODS

### Ethics declaration

The study was approved by the local institutional review board of the Sheba Medical Center (approval number SMC 1765-14) and of Tel-Aviv University. Written informed consent for retention and testing of residual plasma samples was provided by the patients.

### Preparation of plasmids

In order to maintain the plasmid stock as homogenous as possible, plasmids were transformed to a chemically competent bacteria cells [DH5alpha (BioLab, Israel) or TG1 [A kind gift by Itai Benhar (Tel Aviv University, Tel Aviv, Israel)]] by a standard heat-shock protocol. Based on the fact that *Escherichia coli* doubling time is 20 mins in average using rich growing medium (76), a single colony was selected and grown to a maximum of 100 generations. Plasmids were column purified (HiYield™ Plasmid Mini Kit, RBC Bioscience) and stored at -20°C until use.

### Construction of baseline control amplicon

Baseline control amplicon was based on clonal amplification and sequencing of the pLAI.2 plasmid, which contains a full-length HIV-1<sub>LAI</sub> proviral clone (43) (obtained through the NIH AIDS Reagent Program, Division of AIDS, NIAID, NIH: pLAI.2 from Dr. Keith Peden, courtesy of the MRC AIDS Directed Program). The Integrase region of pLAI.2 was amplified using primers:



KLV70 - 5'TTC RGG ATY AGA AGT AAA YAT AGT AAC AG and KLV84 - 5'TCC TGT ATG CAR ACC CCA ATA TG (77). Polymerase chain reaction (PCR) amplification was conducted using the high-fidelity Platinum™ SuperFi™ DNA Polymerase (Invitrogen) in a 50µl reaction using 20-40 ng of the plasmid as input and according to the manufacturer instructions. Amplification in a thermal cycler was performed as follows: initial denaturation for 3min at 98°C, followed by 40 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 60°C and extension for 1min at 72°C, and final extension for 2min at 72°C. An alternative high-fidelity DNA polymerase used was Q5 High-Fidelity DNA Polymerase (New England Biolabs, NEB). PCR cycles were set according to each manufacturer's instructions using the above described PCR program. The Integrase amplicon was gel purified (Wizard® SV Gel and PCR Clean-Up System, Promega) and the concentration determined by Qubit fluorometer (Invitrogen) according to each manufacturer instructions. The purified product was further used for library construction.

For the AmpR sample, the conserved *AmpR* gene was amplified from PLAI.2 plasmid using primers: AmpR FW - 5'AAA GTT CTG CTA TGT GGC GC and AmpR RV - 5'GGT CTG ACA GTT ACC AAT GC. PCR amplification was carried out as described above, except for extension duration of 30sec instead of 1min. Similarly, the conserved *RpoB* gene was amplified from the bacteria genome using the following primers: RpoB FW 5'- ATG GTT TAC TCC TAT ACC GA and RpoB RV 5'- GTG ATC CAG ATC GTT GGT G and the following PCR program: initial denaturation for 3min at 98°C, followed by 40 cycles of denaturation for 10sec at 98°C, annealing for 10sec at 60°C and extension for 4sec at 72°C, and final extension for 2min at 72°C.

#### **Construction of alternative purification amplicons**

The agarose gel purification step of the amplified integrase gene was replaced with other purification methods; (1) For the gel-free sample, the amplified integrase gene was purified using 25µl of AMPure XP beads (0.5X ratio, Beckman Coulter) according to the manufacturer instructions; And (2) For the ExoSap sample, 10µl of the amplified integrase gene were mixed with 4µl of ExoSap (ExoSAP-IT™ PCR Product Cleanup Reagent, Applied Biosystems) and incubated according to the manufacturer instructions. No other changes in the generation of amplicon protocol were made.

#### **Construction of a PCR free control amplicon**

For the PCR-free sample, 10ug of PLAI.2 plasmid was digested using the restriction enzymes: NheI, StuI and XcmI (NEB) according to the manufacturer instructions. A ~1500bp fragment containing the integrase gene was gel purified and concentration was determined by Qubit. The purified product was further used for library construction.

#### **Construction of an RNA control amplicon**

A plasmid containing the full cDNA of Coxsackie virus B3 (CVB3) under a T7 promoter was a kind gift from Marco Vignuzzi (Institut Pasteur, Paris, France). The plasmid was used to generate an

RNA control pool. Ten micrograms of this plasmid were linearized using Sall (NEB), purified by AMPure XP beads (0.5X ratio), and then *in-vitro* transcribed using T7 RNA polymerase (NEB) according to the manufacturer instructions. The transcribed RNA was purified using AMPure XP beads (0.5X ratio) and reverse transcribed with random hexamers using SuperScript III Reverse Transcriptase (Invitrogen) according to the manufacturer instructions. Four microliters of the reverse transcription reaction were used as template for a PCR reaction using primers: CVB FW - 5'GGA GAG AAG GTC AAC TCT ATG GAA GC and CVB RV - 5'TAC CAC CCT GTA GTT CCC CA, which amplify a ~1500bp fragment within the CVB genome. PCR reaction (50µl total) was set and amplified using Platinum™ SuperFi™ as follows: initial denaturation for 3min at 98°C, followed by 40 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 60°C and extension for 15sec at 72°C, and final extension for 2min at 72°C. The CVB amplicon was gel purified and the concentration measured by Qubit. The purified product was further used for library construction.

#### **Construction of clinical HIV-1 amplicon with primer-ID**

Plasma sample from a recently diagnosed HIV-1 patient (clinical sample, ID 83530) with  $>1 \times 10^7$  c/ml HIV-1 viral load was provided by the National HIV Reference Laboratory, Chaim Sheba Medical Center, Ramat-Gan, Israel. The mode of HIV-1 transmission for this patient was MSM, men who have sex with men. HIV-1 viral load was determined with Xpert HIV-1 viral load assay on GeneXpert (Cepheid Inc., Sunnyvale, CA), according to the manufacturer instructions (78). RNA was extracted from 0.5 mL plasma by NucliSENS Easy MAG (Biomerieux, Marcy l'Etoile, France) according to the manufacturer's protocol, eluted in a final volume of 55 µl and stored in -80°C until use. A primer specific to the Gag gene of HIV-1 was designed with a 15 N-bases unique barcode followed by a linker sequence for subsequent PCR, Gag ID RT - 5'TAC CCA TAC GAT GTT CCA GAT TAC GNN NNN NNN NNN NNN NAC TGT ATC ATC TGC TCC TG TRT CT. Based on the measured viral load and sample concentration, 4 µl (containing the genomes of roughly 300,000 viruses) were taken for reverse transcription reaction. RT was performed using SuperScript IV Reverse Transcriptase (Invitrogen) according to the manufacturer instructions with the following adjustments; (1) In order to maximize the primer annealing to the viral RNA, the sample was allowed to cool down gradually from 65°C to room temperature for 10 minutes before it was transferred to ice for 2min; And (2) The reaction was incubated for 30min at 55°C to increase the overall reaction yield. To remove excess primers, the resulting cDNA was purified using AMPure XP beads (0.5X ratio) and eluted with 35µl nuclease-free water. To avoid loss of barcoded primers due to coverage drop at the ends of a read as a result of the tagmentation process, the PCR forward primer was designed with a 60bp overhang so the barcode ("primer ID") is far from read end, Gag ID FW - 5'AAG CGA GGA GCT GTT CAC TGC CAT CCT GGT CGA GCT ACC CAT ACG ATG TTC CAG ATT ACG and Gag ID RV - 5'CTC AAT AAA GCT TGC CTT GAG TGC. PCR amplification was accomplished using Platinum™

SuperFi™ in a 50µl reaction with 33.5µl of the purified cDNA as input using the following conditions: initial denaturation for 3min at 98°C, followed by 40 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 60°C and extension for 1min at 72°C, and final extension for 2min at 72°C. The Gag amplicon was gel purified and the concentration determined by Qubit. The purified product was further used for library construction.

### **Libraries construction**

PCR fragmentation and indexing of samples for sequencing was performed using the Nextera XT DNA Library Prep Kit (Illumina) with the following adjustments to the manufacturer instructions; (1) In order to get a short insert size of ~250bp, 0.85 ng of input DNA was used for tagmentation; (2) No neutralization of the tagmentation buffer was done, as described previously (79); (3) For library amplification of the tagmented DNA, the Nextera XT DNA library prep PCR reagents were replaced with high-fidelity DNA polymerase reagents (the same DNA polymerase that was used for the amplicon generation). The PCR reaction (50µl total) was set as depicted. Directly to the tagmented DNA, index 1 (i5, illumina, 5µl), index 2 (i7, illumina, 5µl), buffer (10µl), high-fidelity DNA polymerase (0.5µl), dNTPs (10mM, 1µl) and nuclease-free water (8.5µl) were added; (4) Amplification was performed with annealing temperature set to 63°C instead of 55°C, as introduced previously (79) and final extension for 2min; (5) Post-amplification clean-up was achieved using AMPure XP beads in a double size-selection manner (80) to remove larger fragments as well as smaller fragments, in order to obtain a narrower size-selection that will maximize the fraction of fully overlapping read pairs. For the first size-selection, 32.5µl of beads (0.65X ratio) were added to bind the large fragments. These beads were separated and discarded. For the second-size selection, 10µl of beads (0.2X ratio) were added to the supernatant to allow binding of intermediate fragments, and the supernatant containing the small fragments was discarded. The intermediate fragments were eluted and their size was determined using a high-sensitivity DNA tape in Tapestation 4200 (Agilent). A mean size of ~370bp, corresponding to the desired insert size of ~250bp, was achieved; And (6) Normalization and pooling was performed manually.

### **Alternative library purification methods**

For the AMPure XP beads-free sample, post-amplification clean-up by double size-selection was replaced with an agarose gel purification of a ~370bp fragment, with no other changes in the library construction protocol.

### **Alternative tagmentation sample**

For the alternative tagmentation sample, a 250bp amplicon within the integrase region was designed, using specific primers with an overhang corresponding to the sequence inserted during the tagmentation step of the NexteraXT DNA library prep kit, NexteraXT free FW - 5'TCG TCG GCA GCG TCA GAT GTG TAT AAG AGA CAG ACT TGT CCA TGC ATG GCT TCT C and NexteraXT free RV - 5'GTC TCG TGG GCT CGG AGA TGT GTA TAA GAG ACA GTC TAT CTG

GCA TGG GTA CCA GCA. PCR reaction was set up using Platinum™ SuperFi™ and carried out as follows: initial denaturation for 3min at 98°C, followed by 40 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 62°C and extension for 15sec at 72°C, and final denaturation for 2min at 72°C. The PCR product was gel purified and the concentration was measured by Qubit. The purified product was indexed by a succeeding PCR amplification using primers corresponding to i5 and i7 NexteraXT primers (IDT) as mentioned previously (79) at a final concentration of 1uM. The PCR reaction was set up using Platinum™ SuperFi™ and amplified as detailed: initial denaturation for 3min at 98°C, followed by 12 cycles of denaturation for 20sec at 98°C, annealing for 30sec at 63°C and extension for 30sec at 72°C, and final extension for 2min at 72°C. Size selection was achieved by gel purification of ~370bp fragments.

### **NextSeq libraries construction**

Illumina NextSeq supports shorter reads than MiSeq. The longest NextSeq read length is 150bp, therefore we selected a shorter insert size of 270bp, compared to the 370bp insert size for the MiSeq platform. The first size selection of the post-NexteraXT amplification cleanup was performed using 42.5µl of AMPure XP beads (0.85X ratio) (80).

### **Sequencing**

Sequencing of all samples (except for the NextSeq sample) was performed on the Illumina MiSeq platform using MiSeq Reagent Kit v2 (500-cycles) [Illumina]. Sequencing of the NextSeq samples was performed on the Illumina NextSeq 500 platform using NextSeq 500/550 High Output Kit (300-cycles) [Illumina].

### **Reads processing and base calling**

The paired-end reads from each control library were aligned against the reference sequence of that control using an in-house script that relies on BLAST command-line tool (81-83). The paired-end reads from the clinical HIV-1 sample were aligned against HIV-1 subtype B HXB2 reference sequence (GenBank accession number K03455.1) and then realigned against the consensus sequence obtained. Bases were called using an in-house script only if the forward and reverse reads agreed and their average Q-score was above an input threshold (30 or 38). At each position, for each alternative base, we calculate mutation frequencies by dividing the number of reads bearing the mutation by sequence coverage. Positions were retained for analysis only if sequenced to a depth of at least 100,000 reads. In order to analyze the errors in the sequencing process we used Python (Anaconda distribution) with the following packages: pandas, matplotlib, seaborn, numpy and stats. Distributions of errors on control plasmids were compared using two-tailed t-test or two-tailed Mann-Whitney U test.

### **Variant calling**

In order to facilitate discrimination of true variants from AccuNGS process artifacts, we created a variant caller based on two principles: (i) positions that exhibit relatively high level of error on a control sample are error-prone for the clinical sample as well; and (ii) process errors on a control

sample follow a gamma distribution. A gamma distribution was fitted for each mutation type in the control sequence. In order to detect and remove outliers from the fitting process we used the “three-sigma-rule”, and positions that showed error higher than three standard deviations from the mean of the fitted distribution were removed. For these rare loci a base was called only if the mutation was more prevalent in the sample by an order of magnitude. For G>A transition mutations, four distinct gamma distributions were fitted, corresponding for all four G>A combinations with preceding nucleotide. Accordingly, for C>T transition mutations four gamma distributions were fitted as well, on the four C>T reverse complement mutations of the G>A mutations. For establishing Figure 5, variants were called on the input RNA sample only if a mutation was in the extreme 1% of the corresponding gamma distribution fitted using the DNA control.

### **Standard sequencing control sample**

Standard control sample of pLAI.2 was taken from a previous study (77). For obtaining mutation frequencies we used the same pipeline as for the AccuNGS samples, but without correcting overlapping paired reads. Positions in this sample were analyzed only if sequenced to a depth of at least 2,000 bases.

### **CODE AVAILABILITY**

We have developed the following computational resources that complement the AccuNGS sequencing protocol:

- (a) Base coverage calculator. AccuNGS relies on overlapping read pairs and high Q-scores for both reads of a pair. The calculator receives as input the length of the target regions and the desired coverage, and outputs the recommended number of reads required for sequencing each sample.
- (b) Computational pipeline for computing the number of unique RNA molecules sequenced, based on primer-ID barcodes (see Supplementary Text).
- (c) Computational pipeline for base-calling and inferring site by site base frequencies.

All resources are freely available at <https://github.com/SternLabTAU/AccuNGS>.

### **ACCESSION NUMBERS**

The datasets generated and reported in this study were deposited in the Sequencing Read Archive (SRA, available at <https://www.ncbi.nlm.nih.gov/sra>), under BioProject PRJNA476431.

### **ACKNOWLEDGEMENTS**

The authors would like to thank the members of the Stern lab for helpful comments and Roy Moscona for helpful discussions and support. This work was supported by the SAIA foundation; by the Israeli Science Foundation [1333/16 to AS]; by the German Israeli Foundation [I-1096-

411.8-2015 to AS]; by the United-States-Israel Binational Science Foundation [2016555 to AS]; by the Edmond J. Safra center for bioinformatics in Tel Aviv University [to MG, TK and DM]; and by the Constantiner Institute for Molecular Genetics in Tel Aviv University [to MG].

## COMPETING INTERESTS

The authors have no competing interests to declare.

## REFERENCES

1. Reid-Bayliss, K.S. and Loeb, L.A. (2017) Accurate RNA consensus sequencing for high-fidelity detection of transcriptional mutagenesis-induced epimutations. *P Natl Acad Sci USA*, **114**, 9415-9420.
2. Salk, J.J., Schmitt, M.W. and Loeb, L.A. (2018) Enhancing the accuracy of next-generation sequencing for detecting rare and subclonal mutations. *Nat Rev Genet*, **19**, 269-285.
3. Achermann, J.C., Domenice, S., Bachega, T.A., Nishi, M.Y. and Mendonca, B.B. (2015) Disorders of sex development: effect of molecular diagnostics. *Nat Rev Endocrinol*, **11**, 478-488.
4. Beerenwinkel, N., Gunthard, H.F., Roth, V. and Metzner, K.J. (2012) Challenges and opportunities in estimating viral genetic diversity from next-generation sequencing data. *Front Microbiol*, **3**, 329.
5. Sanjuan, R., Nebot, M.R., Chirico, N., Mansky, L.M. and Belshaw, R. (2010) Viral mutation rates. *J Virol*, **84**, 9733-9748.
6. Zanini, F., Puller, V., Brodin, J., Albert, J. and Neher, R.A. (2017) In vivo mutation rates and the landscape of fitness costs of HIV-1. *Virus Evol*, **3**, vex003.
7. Brumme, C.J. and Poon, A.F.Y. (2017) Promises and pitfalls of Illumina sequencing for HIV resistance genotyping. *Virus Res*, **239**, 97-105.
8. McElroy, K., Thomas, T. and Luciani, F. (2014) Deep sequencing of evolving pathogen populations: applications, errors, and bioinformatic solutions. *Microb Inform Exp*, **4**, 1.
9. Casadella, M. and Paredes, R. (2017) Deep sequencing for HIV-1 clinical management. *Virus Res*, **239**, 69-81.
10. Theys, K., Feder, A.F., Gelbart, M., Hartl, M., Stern, A. and Pennings, P.S. (2018) Within-patient mutation frequencies reveal fitness costs of CpG dinucleotides and drastic amino acid changes in HIV. *PLoS Genet*, **14**, e1007420.
11. Sims, D., Sudbery, I., Illott, N.E., Heger, A. and Ponting, C.P. (2014) Sequencing depth and coverage: key considerations in genomic analyses. *Nat Rev Genet*, **15**, 121-132.
12. Gallet, R., Fabre, F., Michalakis, Y. and Blanc, S. (2017) The number of target molecules of the amplification step limits accuracy and sensitivity in ultra deep sequencing viral population studies. *J Virol*, **91**, e00561-00517.
13. Meacham, F., Boffelli, D., Dhahbi, J., Martin, D.I., Singer, M. and Pachter, L. (2011) Identification and correction of systematic error in high-throughput sequence data. *BMC Bioinformatics*, **12**, 451.
14. Huber, M., Metzner, K.J., Geissberger, F.D., Shah, C., Leemann, C., Klimkait, T., Boni, J., Trkola, A. and Zagordi, O. (2017) MinVar: A rapid and versatile tool for HIV-1 drug resistance genotyping by deep sequencing. *J Virol Methods*, **240**, 7-13.

15. Boucher, C.A., Bobkova, M.R., Geretti, A.M., Hung, C.C., Kaiser, R., Marcelin, A.G., Streinu-Cercel, A., van Wyk, J., Dorr, P. and Vandamme, A.M. (2018) State of the Art in HIV Drug Resistance: Science and Technology Knowledge Gap. *AIDS Rev*, **20**, 27-42.
16. Clutter, D.S., Jordan, M.R., Bertagnolio, S. and Shafer, R.W. (2016) HIV-1 drug resistance and resistance testing. *Infect Genet Evol*, **46**, 292-307.
17. Acevedo, A., Brodsky, L. and Andino, R. (2014) Mutational and fitness landscapes of an RNA virus revealed through population sequencing. *Nature*, **505**, 686-690.
18. Lou, D.I., Hussmann, J.A., McBee, R.M., Acevedo, A., Andino, R., Press, W.H. and Sawyer, S.L. (2013) High-throughput DNA sequencing errors are reduced by orders of magnitude using circle sequencing. *P Natl Acad Sci USA*, **110**, 19872-19877.
19. Wang, K., Lai, S., Yang, X., Zhu, T., Lu, X., Wu, C.I. and Ruan, J. (2017) Ultrasensitive and high-efficiency screen of de novo low-frequency mutations by o2n-seq. *Nat Commun*, **8**, 15335.
20. Jabara, C.B., Jones, C.D., Roach, J., Anderson, J.A. and Swanstrom, R. (2011) Accurate sampling and deep sequencing of the HIV-1 protease gene using a Primer ID. *P Natl Acad Sci USA*, **108**, 20166-20171.
21. Zhou, S., Jones, C., Mieczkowski, P. and Swanstrom, R. (2015) Primer ID Validates Template Sampling Depth and Greatly Reduces the Error Rate of Next-Generation Sequencing of HIV-1 Genomic RNA Populations. *J Virol*, **89**, 8540-8555.
22. Kennedy, S.R., Schmitt, M.W., Fox, E.J., Kohn, B.F., Salk, J.J., Ahn, E.H., Prindle, M.J., Kuong, K.J., Shen, J.C., Risques, R.A. *et al.* (2014) Detecting ultralow-frequency mutations by Duplex Sequencing. *Nat Protoc*, **9**, 2586-2606.
23. Jee, J., Rasouly, A., Shamovsky, I., Akivis, Y., Steinman, S.R., Mishra, B. and Nudler, E. (2016) Rates and mechanisms of bacterial mutagenesis from maximum-depth sequencing. *Nature*, **534**, 693-696.
24. Newman, A.M., Lovejoy, A.F., Klass, D.M., Kurtz, D.M., Chabon, J.J., Scherer, F., Stehr, H., Liu, C.L., Bratman, S.V., Say, C. *et al.* (2016) Integrated digital error suppression for improved detection of circulating tumor DNA. *Nat Biotechnol*, **34**, 547-555.
25. Chen-Harris, H., Borucki, M.K., Torres, C., Slezak, T.R. and Allen, J.E. (2013) Ultra-deep mutant spectrum profiling: improving sequencing accuracy using overlapping read pairs. *BMC Genomics*, **14**, 96.
26. Preston, J.L., Royall, A.E., Randel, M.A., Sikkink, K.L., Phillips, P.C. and Johnson, E.A. (2016) High-specificity detection of rare alleles with Paired-End Low Error Sequencing (PELE-Seq). *BMC Genomics*, **17**, 464.
27. Schirmer, M., Ijaz, U.Z., D'Amore, R., Hall, N., Sloan, W.T. and Quince, C. (2015) Insight into biases and sequencing errors for amplicon sequencing with the Illumina MiSeq platform. *Nucleic Acids Res*, **43**, e37.
28. Imashimizu, M., Oshima, T., Lubkowska, L. and Kashlev, M. (2013) Direct assessment of transcription fidelity by high-resolution RNA sequencing. *Nucleic Acids Res*, **41**, 9090-9104.
29. Orton, R.J., Wright, C.F., Morelli, M.J., King, D.J., Paton, D.J., King, D.P. and Haydon, D.T. (2015) Distinguishing low frequency mutations from RT-PCR and sequence errors in viral deep sequencing data. *BMC Genomics*, **16**, 229.
30. Gerstung, M., Beisel, C., Rechsteiner, M., Wild, P., Schraml, P., Moch, H. and Beerenwinkel, N. (2012) Reliable detection of subclonal single-nucleotide variants in tumour cell populations. *Nat Commun*, **3**, 811.
31. Wilm, A., Aw, P.P., Bertrand, D., Yeo, G.H., Ong, S.H., Wong, C.H., Khor, C.C., Petric, R., Hibberd, M.L. and Nagarajan, N. (2012) LoFreq: a sequence-quality aware, ultra-

- sensitive variant caller for uncovering cell-population heterogeneity from high-throughput sequencing datasets. *Nucleic Acids Res*, **40**, 11189-11201.
32. Van der Borght, K., Thys, K., Wetzels, Y., Clement, L., Verbist, B., Reumers, J., van Vlijmen, H. and Aerssens, J. (2015) QQ-SNV: single nucleotide variant detection at low frequency by comparing the quality quantiles. *BMC bioinformatics*, **16**, 379.
  33. Yang, X., Charlebois, P., Macalalad, A., Henn, M.R. and Zody, M.C. (2013) V-Phaser 2: variant inference for viral populations. *BMC Genomics*, **14**, 674.
  34. Brodin, J., Hedskog, C., Heddini, A., Benard, E., Neher, R.A., Mild, M. and Albert, J. (2015) Challenges with using primer IDs to improve accuracy of next generation sequencing. *PLoS One*, **10**, e0119123.
  35. Zhang, T.H., Wu, N.C. and Sun, R. (2016) A benchmark study on error-correction by read-pairing and tag-clustering in amplicon-based deep sequencing. *BMC Genomics*, **17**, 108.
  36. McCrone, J.T. and Lauring, A.S. (2016) Measurements of Intrahost Viral Diversity Are Extremely Sensitive to Systematic Errors in Variant Calling. *J Virol*, **90**, 6884-6895.
  37. Roberts, J.D., Bebenek, K. and Kunkel, T.A. (1988) The accuracy of reverse transcriptase from HIV-1. *Science*, **242**, 1171-1173.
  38. Preston, B.D., Poiesz, B.J. and Loeb, L.A. (1988) Fidelity of HIV-1 reverse transcriptase. *Science*, **242**, 1168-1171.
  39. Potter, J., Zheng, W. and Lee, J. (2003) Thermal stability and cDNA synthesis capability of SuperScript III reverse transcriptase. *Focus*, **25**, 19-24.
  40. Barnes, W.M. (1992) The Fidelity of Taq Polymerase Catalyzing Pcr Is Improved by an N-Terminal Deletion. *Gene*, **112**, 29-35.
  41. Seifert, D., Di Giallonardo, F., Topfer, A., Singer, J., Schmutz, S., Gunthard, H.F., Beerenwinkel, N. and Metzner, K.J. (2016) A Comprehensive Analysis of Primer IDs to Study Heterogeneous HIV-1 Populations. *J Mol Biol*, **428**, 238-250.
  42. Kivioja, T., Vaharautio, A., Karlsson, K., Bonke, M., Enge, M., Linnarsson, S. and Taipale, J. (2011) Counting absolute numbers of molecules using unique molecular identifiers. *Nat Methods*, **9**, 72-74.
  43. Peden, K., Emerman, M. and Montagnier, L. (1991) Changes in growth properties on passage in tissue culture of viruses derived from infectious molecular clones of HIV-1LAI, HIV-1MAL, and HIV-1ELI. *Virology*, **185**, 661-672.
  44. Van Laethem, K., Schrooten, Y., Covens, K., Dekeersmaecker, N., De Munter, P., Van Wijngaerden, E., Van Ranst, M. and Vandamme, A.M. (2008) A genotypic assay for the amplification and sequencing of integrase from diverse HIV-1 group M subtypes. *J Virol Methods*, **153**, 176-181.
  45. Brodin, J., Mild, M., Hedskog, C., Sherwood, E., Leitner, T., Andersson, B. and Albert, J. (2013) PCR-induced transitions are the major source of error in cleaned ultra-deep pyrosequencing data. *PLoS One*, **8**, e70388.
  46. Zanini, F., Brodin, J., Albert, J. and Neher, R.A. (2017) Error rates, PCR recombination, and sampling depth in HIV-1 whole genome deep sequencing. *Virus Res*, **239**, 106-114.
  47. Casali, N. (2003), *E. coli Plasmid Vectors*. Springer, pp. 27-48.
  48. Edgar, R.C. and Flyvbjerg, H. (2015) Error filtering, pair assembly and error correction for next-generation sequencing reads. *Bioinformatics*, **31**, 3476-3482.
  49. Zhang, J., Kobert, K., Flouri, T. and Stamatakis, A. (2014) PEAR: a fast and accurate Illumina Paired-End reAd mergeR. *Bioinformatics*, **30**, 614-620.
  50. Costello, M., Pugh, T.J., Fennell, T.J., Stewart, C., Lichtenstein, L., Meldrim, J.C., Fostel, J.L., Friedrich, D.C., Perrin, D., Dionne, D. *et al.* (2013) Discovery and characterization of



- artificial mutations in deep coverage targeted capture sequencing data due to oxidative DNA damage during sample preparation. *Nucleic Acids Res*, **41**, e67.
51. Nishimura, S. (2011) 8-Hydroxyguanine: a base for discovery. *DNA Repair (Amst)*, **10**, 1078-1083.
  52. Park, G., Park, J.K., Shin, S.-H., Jeon, H.-J., Kim, N.K., Kim, Y.J., Shin, H.-T., Lee, E., Lee, K.H. and Son, D.-S. (2017) Characterization of background noise in capture-based targeted sequencing data. *Genome biology*, **18**, 136.
  53. Arbeithuber, B., Makova, K.D. and Tiemann-Boege, I. (2016) Artificial mutations resulting from DNA lesions limit detection levels in ultrasensitive sequencing applications. *DNA Res*, **23**, 547-559.
  54. Sultana, S., Solotchi, M., Ramachandran, A. and Patel, S.S. (2017) Transcriptional fidelities of human mitochondrial POLRMT, yeast mitochondrial Rpo41, and phage T7 single-subunit RNA polymerases. *J Biol Chem*, **292**, 18145-18160.
  55. Keele, B.F., Giorgi, E.E., Salazar-Gonzalez, J.F., Decker, J.M., Pham, K.T., Salazar, M.G., Sun, C., Grayson, T., Wang, S., Li, H. *et al.* (2008) Identification and characterization of transmitted and early founder virus envelopes in primary HIV-1 infection. *P Natl Acad Sci USA*, **105**, 7552-7557.
  56. Abram, M.E., Ferris, A.L., Shao, W., Alvord, W.G. and Hughes, S.H. (2010) Nature, position, and frequency of mutations made in a single cycle of HIV-1 replication. *J Virol*, **84**, 9864-9878.
  57. Zanini, F., Brodin, J., Thebo, L., Lanz, C., Bratt, G., Albert, J. and Neher, R.A. (2015) Population genomics of inpatient HIV-1 evolution. *Elife*, **4**, e11282.
  58. Kijak, G.H., Sanders-Buell, E., Chenine, A.L., Eller, M.A., Goonetilleke, N., Thomas, R., Leviyang, S., Harbolick, E.A., Bose, M., Pham, P. *et al.* (2017) Rare HIV-1 transmitted/founder lineages identified by deep viral sequencing contribute to rapid shifts in dominant quasiespecies during acute and early infection. *Plos Pathog*, **13**, e1006510.
  59. Henn, M.R., Boutwell, C.L., Charlebois, P., Lennon, N.J., Power, K.A., Macalalad, A.R., Berlin, A.M., Malboeuf, C.M., Ryan, E.M., Gnerre, S. *et al.* (2012) Whole genome deep sequencing of HIV-1 reveals the impact of early minor variants upon immune recognition during acute infection. *Plos Pathog*, **8**, e1002529.
  60. Fiebig, E.W., Wright, D.J., Rawal, B.D., Garrett, P.E., Schumacher, R.T., Peddada, L., Heldebrandt, C., Smith, R., Conrad, A., Kleinman, S.H. *et al.* (2003) Dynamics of HIV viremia and antibody seroconversion in plasma donors: implications for diagnosis and staging of primary HIV infection. *AIDS*, **17**, 1871-1879.
  61. Markowitz, M., Louie, M., Hurley, A., Sun, E., Di Mascio, M., Perelson, A.S. and Ho, D.D. (2003) A novel antiviral intervention results in more accurate assessment of human immunodeficiency virus type 1 replication dynamics and T-cell decay in vivo. *J Virol*, **77**, 5037-5038.
  62. Levesque-Sergerie, J.P., Duquette, M., Thibault, C., Delbecchi, L. and Bissonnette, N. (2007) Detection limits of several commercial reverse transcriptase enzymes: impact on the low- and high-abundance transcript levels assessed by quantitative RT-PCR. *Bmc Mol Biol*, **8**, 93.
  63. Bishop, K.N., Holmes, R.K., Sheehy, A.M., Davidson, N.O., Cho, S.J. and Malim, M.H. (2004) Cytidine deamination of retroviral DNA by diverse APOBEC proteins. *Curr Biol*, **14**, 1392-1396.
  64. Potapov, V. and Ong, J.L. (2017) Examining sources of error in PCR by single-molecule sequencing. *PloS one*, **12**, e0169774.

65. Stern, A., Yeh, M.T., Zinger, T., Smith, M., Wright, C., Ling, G., Nielsen, R., Macadam, A. and Andino, R. (2017) The Evolutionary Pathway to Virulence of an RNA Virus. *Cell*, **169**, 35-46 e19.
66. Kou, R., Lam, H., Duan, H., Ye, L., Jongkam, N., Chen, W., Zhang, S. and Li, S. (2016) Benefits and Challenges with Applying Unique Molecular Identifiers in Next Generation Sequencing to Detect Low Frequency Mutations. *PLoS One*, **11**, e0146638.
67. Hao, Y., Xuei, X., Li, L., Nakshatri, H., Edenberg, H.J. and Liu, Y. (2017) RareVar: A Framework for Detecting Low-Frequency Single-Nucleotide Variants. *J Comput Biol*, **24**, 637-646.
68. Xu, C., Gu, X., Padmanabhan, R., Wu, Z., Peng, Q., DiCarlo, J. and Wang, Y. (2018) smCounter2: an accurate low-frequency variant caller for targeted sequencing data with unique molecular identifiers. *bioRxiv*, 281659.
69. Masella, A.P., Bartram, A.K., Truszkowski, J.M., Brown, D.G. and Neufeld, J.D. (2012) PANDAseq: paired-end assembler for illumina sequences. *BMC Bioinformatics*, **13**, 31.
70. Howison, M., Coetzer, M. and Kantor, R. (2018) Measurement error and variant-calling in deep Illumina sequencing of HIV. *bioRxiv*, 276576.
71. Bentley, D.R., Balasubramanian, S., Swerdlow, H.P., Smith, G.P., Milton, J., Brown, C.G., Hall, K.P., Evers, D.J., Barnes, C.L., Bignell, H.R. *et al.* (2008) Accurate whole human genome sequencing using reversible terminator chemistry. *Nature*, **456**, 53-59.
72. Chen, G., Mosier, S., Gocke, C.D., Lin, M.T. and Eshleman, J.R. (2014) Cytosine deamination is a major cause of baseline noise in next-generation sequencing. *Mol Diagn Ther*, **18**, 587-593.
73. Cuevas, J.M., Geller, R., Garijo, R., Lopez-Aldeguer, J. and Sanjuan, R. (2015) Extremely High Mutation Rate of HIV-1 In Vivo. *PLoS Biol*, **13**, e1002251.
74. Harari, A., Ooms, M., Mulder, L.C. and Simon, V. (2009) Polymorphisms and splice variants influence the antiretroviral activity of human APOBEC3H. *J Virol*, **83**, 295-303.
75. Simon, V., Zennou, V., Murray, D., Huang, Y., Ho, D.D. and Bieniasz, P.D. (2005) Natural variation in Vif: differential impact on APOBEC3G/3F and a potential role in HIV-1 diversification. *Plos Pathog*, **1**, e6.
76. Sezonov, G., Joseleau-Petit, D. and D'Ari, R. (2007) Escherichia coli physiology in Luria-Bertani broth. *J Bacteriol*, **189**, 8746-8749.
77. Moscona, R., Ram, D., Wax, M., Bucris, E., Levy, I., Mendelson, E. and Mor, O. (2017) Comparison between next-generation and Sanger-based sequencing for the detection of transmitted drug-resistance mutations among recently infected HIV-1 patients in Israel, 2000–2014. *J INT AIDS SOC*, **20**.
78. Mor, O., Gozlan, Y., Wax, M., Mileguir, F., Rakovsky, A., Noy, B., Mendelson, E. and Levy, I. (2015) Evaluation of the RealTime HIV-1, Xpert HIV-1, and Aptima HIV-1 Quant Dx Assays in Comparison to the NucliSens EasyQ HIV-1 v2.0 Assay for Quantification of HIV-1 Viral Load. *J Clin Microbiol*, **53**, 3458-3465.
79. Baym, M., Kryazhimskiy, S., Lieberman, T.D., Chung, H., Desai, M.M. and Kishony, R. (2015) Inexpensive multiplexed library preparation for megabase-sized genomes. *PLoS One*, **10**, e0128036.
80. Bronner, I.F., Quail, M.A., Turner, D.J. and Swerdlow, H. (2014) Improved Protocols for Illumina Sequencing. *Curr Protoc Hum Genet*, **80**, 18.12.11-42.
81. 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.

82. Altschul, S.F., Madden, T.L., Schaffer, A.A., Zhang, J., Zhang, Z., Miller, W. and Lipman, D.J. (1997) Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Res*, **25**, 3389-3402.
83. Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. and Madden, T.L. (2009) BLAST+: architecture and applications. *BMC Bioinformatics*, **10**, 421.