
Submitting Next-Generation Sequencing Data to the Division of Antivirals

Guidance for Industry Technical Specifications Document

For questions regarding this technical specifications document, contact
CDER at cdcr-edata@fda.hhs.gov.

**U.S. Department of Health and Human Services
Food and Drug Administration
Center for Drug Evaluation and Research (CDER)**

**July 2026
Technical Specifications Document**

Revision History

Date	Version	Summary of Revisions
July 2019	1.0	Initial Version
July 2026	2.0	Revision/Update

Table of Contents

1.0	INTRODUCTION.....	1
2.0	ACCEPTABLE NGS PLATFORMS.....	2
3.0	SUMMARY OF INFORMATION TO SUBMIT TO THE DIVISION	2
3.1	NGS Protocol.....	2
3.2	Raw NGS Data	2
3.3	Consensus Sequences in FASTA Format.....	3
3.4	Other Sequences.....	3
3.5	Amino Acid Frequency Tables	3
3.6	NGS Report and Summary Tables.....	3
3.7	Conservation Analysis of Reference Sequences	4
4.0	NGS PROTOCOL	4
4.1	General Protocol Design Elements	4
4.2	Sample Preparation and Sequencing	5
4.3	NGS Data Analysis Methods.....	6
4.3.1	<i>Demultiplexing, Quality/Adapter Trimming, and Read Filtering.....</i>	<i>6</i>
4.3.2	<i>Mapping Reads and Assembling Contigs.....</i>	<i>6</i>
4.3.2.1	Mapping sequence reads to a reference sequence and identifying variants.....	7
4.3.2.2	De novo assembly of contigs and identifying variants (when applicable).....	7
5.0	NGS REPORT.....	8
5.1	Reporting Results From Read Mapping and Variant Calling.....	8
5.2	Reporting Results From the De Novo Assembly and Variant Calling.....	8
5.3	Additional Information	9
6.0	NGS FILE TYPES AND SUBMISSION PROCEDURES.....	9
6.1	Recommended File Formats and Submission Methods.....	9
6.2	Naming FASTQ Files.....	10
6.3	Naming Consensus Sequences (FASTA Format).....	11
7.0	AMINO ACID FREQUENCY TABLE EXAMPLE	11

Submitting Next-Generation Sequencing Data to the Division of Antivirals

Guidance for Industry Technical Specifications Document¹

This guidance represents the current thinking of the Food and Drug Administration (FDA or Agency) on this topic. It does not establish any rights for any person and is not binding on FDA or the public. You can use an alternative approach if it satisfies the requirements of the applicable statutes and regulations. To discuss an alternative approach, contact the FDA office responsible for this guidance as listed on the title page.

1.0 INTRODUCTION

The purpose of this technical specifications document is to provide the current thinking of FDA's Division of Antivirals (the Division) regarding the submission of next-generation sequencing (NGS) protocols, data, and analyses in support of resistance assessments for the development of antiviral drugs.

The Division performs independent analyses of all NGS data associated with antiviral drugs, including small molecule and immunoglobulin-based drugs, under development to ensure that drug resistance is carefully characterized and explained in the label of approved antiviral drugs. Provision of comprehensive and accurate resistance information is imperative to inform resistance surveillance efforts, monitor the emergence of novel viral variants that are resistant to approved antiviral drugs, and protect public health. In addition, the resistance information provides important guidance for health care professionals who prescribe antiviral drugs and is included in the drug product information approved by the Division.

NGS is a well-established technology in both research and clinical practice that sponsors frequently employ when performing sequence-based resistance analysis. In contrast to Sanger sequencing, which produces a semiquantitative or consensus output of nucleotide proportions in a viral population, NGS produces nucleotide sequence information for individual viruses within a viral population, frequently resulting in millions of short sequences per sample. Given the inherent complexity of these datasets, this technical specification document is intended to address the absence of standardized bioinformatics protocols for large-scale NGS data analysis.

¹ This technical specifications document has been prepared by the Division of Antivirals in the Center for Drug Evaluation and Research at the Food and Drug Administration. You may submit comments on this guidance at any time. Submit comments to Docket No. FDA-2017-D-6821 (available at <https://www.regulations.gov/docket?D=FDA-2017-D-6821>) (see the instructions for submitting comments in the docket).

Contains Nonbinding Recommendations

In general, FDA's guidance documents do not establish legally enforceable responsibilities. Instead, guidances describe the Agency's current thinking on a topic and should be viewed only as recommendations, unless specific regulatory or statutory requirements are cited. The use of the word *should* in Agency guidances means that something is suggested or recommended, but not required.

2.0 ACCEPTABLE NGS PLATFORMS

The Division will accept data generated from most standard NGS platforms if sponsors adequately describe the sequencing platform and provide the requested data, protocols, and other supporting information described in the following sections. We strongly recommend that sponsors communicate with the Division early during drug development and provide an NGS protocol (including description of data analysis methods) before initiating pivotal clinical trials. In addition, we recommend that sponsors submit a small mock or partial NGS dataset (including raw FASTQ files, consensus sequences in FASTA format, and an amino acid frequency table) before formally submitting a large NGS dataset to ensure that the data are provided in an appropriate format for review. The mock NGS dataset could include results from a small number of clinical trial samples and/or controls/test samples. Sponsors should refer to the following sections for guidance on preparing NGS-related submissions.

3.0 SUMMARY OF INFORMATION TO SUBMIT TO THE DIVISION

3.1 NGS Protocol

Sponsors should submit a detailed NGS protocol that includes wet lab procedures and a description of data analysis methods. The wet lab portion of the NGS protocol should include details on sample collection, nucleic acid extraction, reverse transcription (if applicable), nucleic acid amplification, library preparation, library pooling, and sequencing. Sponsors should describe and provide their rationale for any quality control criteria applied to samples (e.g., minimum viral target copy number to attempt NGS), sample polymerase chain reaction (PCR)/reverse transcription polymerase chain reaction (RT-PCR) products (e.g., minimum concentration and required size), sequencing libraries (e.g., minimum concentration and required size), and/or sequencing results (e.g., numbers of reads, read quality, minimum target gene coverage). Sponsors should include appropriate positive and negative controls during PCR/RT-PCR and sequencing. If available, sponsors should provide validation reports for their NGS methods. See section 4.0, NGS Protocol for additional details.

3.2 Raw NGS Data

Sponsors should provide all raw NGS data in FASTQ format. For de novo assemblies, sponsors should provide all contiguous sequences (contigs) longer than 200 nucleotides in FASTQ format. Sponsors do not need to provide NGS data in other formats (e.g., .ACE, .BAM, or .SAM) unless requested by the Division. FASTQ file names should include participant identification number, visit, and read number (for paired-end sequencing) (see section 6.2). For paired-end sequencing, read 1 and read 2 sequencing data should be provided as separate FASTQ files. Sponsors should provide all FASTQ files from each clinical trial in a single folder. However, in cases in which

Contains Nonbinding Recommendations

multiple reference sequences are used (e.g., for different viral genotypes), we recommend organizing the FASTQ files into subfolders based on viral genotype or reference sequence.

3.3 Consensus Sequences in FASTA Format

Sponsors should provide a single FASTA file containing consensus sequences from all samples in a clinical trial. Sequence identifiers in the FASTA file should include participant identification number and visit (see section 6.3).

3.4 Other Sequences

Sponsors should provide any other sequences (e.g., in FASTA format) needed for data analysis, such as the reference sequences, primer sequences, and/or adapter sequences. In addition, if adapter trimming is performed during data analysis, sponsors should provide the sequences that were trimmed and indicate whether trimming was applied to the 5' and/or 3' end of read 1 and/or read 2 (for paired-end sequencing). In the NGS protocol, sponsors should provide the accession number(s) for the reference sequence(s), if applicable. If a sponsor uses a baseline sequence or consensus sequence from their clinical trial as the reference sequence, the sponsor should describe how the reference sequence was derived.

3.5 Amino Acid Frequency Tables

Sponsors should provide an amino acid frequency table (in XLSX or XPT format) reporting all amino acid substitutions that differ from the reference sequence at frequencies $\geq 1\%$. The frequency table should also include frameshifts, premature stop codons, and in-frame insertions and deletions. It is not necessary to include synonymous (silent) mutations or mutations in noncoding regions unless requested by the Division. The frequency table should include all viral genes that were agreed upon in advance with the Division. See section 7.0 for an example of an amino acid frequency table. Of note, for drugs that target nucleic acids and are impacted by nucleotide mutations (e.g., small interfering RNAs and antisense oligonucleotides), the frequency table should include all nucleotide changes, even if they do not result in an amino acid change.

3.6 NGS Report and Summary Tables

Sponsors should provide an NGS Report (see section 5.0) that includes a comprehensive overview of resistance analysis results and useful summary tables, such as the following:

- a. Tables summarizing the NGS dataset (e.g., number of samples successfully sequenced; number of participants with one or more samples sequenced; number of participants with baseline data, post-baseline data, and paired baseline and post-baseline data; number of samples that met criteria for resistance analysis but failed NGS) both overall and by treatment arm.
- b. Tables listing all baseline amino acid polymorphisms in the viral target or all baseline amino acid polymorphisms at positions of interest in the target, such as amino acid

Contains Nonbinding Recommendations

residues that are in contact/proximity to the drug (based on structural data) or associated with resistance in nonclinical studies. Ideally, the table should include substitutions, premature stop codons, and in-frame insertions and deletions. Sponsors should provide the criteria used to define a baseline polymorphism.

- c. Tables listing all treatment-emergent substitutions (TES) in the viral target or all TES at positions of interest in the target. Ideally, the table should include premature stop codons and in-frame insertions and deletions as well. Sponsors should provide the criteria used to define a TES.
- d. Tables summarizing baseline polymorphisms or TES by key subgroups, such as by viral genotype, drug dose, or participant subgroup (e.g., treatment-naïve versus treatment-experienced).
- e. Tables summarizing the association of baseline amino acid polymorphisms and TES with virologic and clinical outcomes of interest, such as viral nonresponse, viral breakthrough, viral relapse, and/or hospitalization/death.

3.7 Conservation Analysis of Reference Sequences

For drugs targeting a viral protein, sponsors should provide a comprehensive conservation analysis of all amino acid positions in the reference sequence. This analysis can be accomplished by comparing all available sequences represented by the reference sequence from publicly accessible databases and from baseline sequences analyzed during clinical development. Results should include an overall percent identity for each position and list all amino acid polymorphisms in decreasing order of prevalence. This analysis should be provided for the reference sequences for each viral genotype. Please specify version numbers for publicly accessible databases used to obtain reference sequences.

4.0 NGS PROTOCOL

4.1 General Protocol Design Elements

Sponsors should include the following general design elements in NGS protocols:

- a. A description of the participants, time points/visits, and sample types to be analyzed. Include the resistance testing criteria that will be used to determine whether NGS of a sample will be attempted.
- b. A description of the NGS platform to be used.
- c. Target gene or gene region name(s) and size(s) to be analyzed.
- d. General analysis strategy (e.g., identify changes relative to a reference sequence, compare sequences from different time points in the same participant).

Contains Nonbinding Recommendations

- e. The anticipated target gene coverage. We recommend that the target gene has a coverage of greater than or equal to 5,000 reads across the full sequence. However, we recognize that this level of coverage may not be achieved in samples with low viral DNA/RNA levels. Sponsors should identify samples that did not reach this level of coverage. NGS data with minimum coverage below 5,000x may be acceptable when sponsors demonstrate high sequencing quality, provide scientific justification and comprehensive validation data confirming assay performance at the proposed coverage depth, and engage in early consultation with the Division to establish protocol agreement before initiating pivotal trials.
- f. Any negative or positive controls. Negative controls for sample PCR/RT-PCR could include water (no template control) and/or DNA/RNA extracted from a healthy (uninfected) volunteer. Positive controls for sample PCR/RT-PCR, library preparation, and sequencing could include a previously characterized sample from a participant infected with the virus being studied and/or recombinant/synthetic DNA (e.g., a plasmid or DNA fragment).
- g. Any quality control criteria for samples, sample PCR/RT-PCR products, sequencing libraries, sequencing runs, and/or sequencing results for individual samples.

4.2 Sample Preparation and Sequencing

Sponsors should describe the following information concerning sample preparation and sequencing:

- a. Methods for collecting, shipping, and storing samples for sequencing.
- b. Methods for extracting nucleic acids from samples and, in the case of RNA, generating complementary DNA. Methods for denaturing RNA secondary structures and a description of the primers used for reverse transcription should be included.
- c. Methods for quantifying copy numbers of the target gene(s) in each sample (e.g., using quantitative PCR/RT-PCR).
- d. Methods for sample PCR product concentration/size determination and purification.
- e. Negative and positive controls. If negative controls show evidence of contamination, PCR should be repeated for all samples that were processed in parallel with the contaminated negative control.
- f. Methods for NGS library preparation, concentration/size determination, purification, and pooling.
- g. Methods for sequencing, including sequencing instrument, sequencing kit, sequence length, sequencing format (e.g., single-end versus paired-end), et cetera.

Contains Nonbinding Recommendations

NOTE: Many NGS protocols in the published scientific literature include specific (gene-specific primers) or nonspecific (random primers) PCR to increase the concentration of DNA for sequencing. However, although predominant genotypes should be amplified by these approaches, minor variants that are important for resistance may not be amplified to the same extent. Therefore, we recommend that sponsors validate NGS protocols that include a sample PCR step by amplifying and sequencing a pool of known variants with different frequencies (e.g., 0.1%, 1%, 5%, 10%, 50%). If feasible, sponsors can also employ approaches that correct for PCR resampling and PCR/sequencing errors, such as the use of unique identifiers.

4.3 NGS Data Analysis Methods

Sponsors should provide a detailed description of the bioinformatics analysis pipeline(s) used to analyze the sequencing data, generate the amino acid frequency table(s), and obtain consensus sequences. Sponsors may submit this information in the form of a Biocompute Object (<https://www.biocomputeobject.org/>) with prior agreement from the Division. Contract research organizations that use a validated bioinformatics platform across multiple antiviral drug development program areas may consider submitting a drug master file to the FDA. Drug master files are not mandatory but are submitted voluntarily to provide confidential information to the FDA, which can then be referenced by multiple applications without disclosing proprietary information to competitors.

The information in the following sections should be provided in the NGS protocol and summarized in the NGS Report. Currently, it is not necessary for sponsors to provide the actual components of bioinformatics pipelines, specific computational tools, or custom analysis algorithms unless requested by the Division.

4.3.1 Demultiplexing, Quality/Adapter Trimming, and Read Filtering

The sponsor's description of NGS data analysis methods should include the procedures used to add barcodes, separate reads by barcodes, trim barcodes, and any precautions taken to eliminate potential cross-barcode contamination. Sponsors should also describe adapter and quality trimming (e.g., related to base quality and/or degenerate bases). Lastly, sponsors should describe any read filtering (e.g., based on read quality or length after trimming). Sponsors should describe the bioinformatics tools used for these steps and settings for key parameters.

4.3.2 Mapping Reads and Assembling Contigs

Two approaches can be used for short read NGS analyses of data to support the development of antiviral drug products: (1) mapping of short reads to a reference sequence or (2) de novo assembly of short reads to assemble contigs. We recommend that sponsors that plan to use long read technologies to assess antiviral drug resistance reach out to the Division for agreement prior to using the technology.

Contains Nonbinding Recommendations

4.3.2.1 Mapping sequence reads to a reference sequence and identifying variants

Sponsors should consider the following for inclusion in the data analysis methods on mapping sequence reads to a reference sequence and identifying variants:

- a. Indicate the reference sequence(s) used for mapping and provide the nucleotide and amino acid sequences and an accession number, if applicable. Please consult with the Division for prior agreement if standard reference sequences have not been established in the scientific literature.
- b. Indicate the bioinformatics tool(s) used to map reads to the reference sequence and for any post-alignment processing steps, such as local realignment or trimming of PCR primers. Sponsors should describe the bioinformatics tools used for these steps and settings for key parameters. Sponsors should provide specific information about the tool used for mapping, including but not limited to percentage of mismatches tolerated, number of insertions and deletions tolerated, and identity and similarity cutoffs.
- c. Indicate the bioinformatics tool(s) used to identify and filter variants and describe settings for key parameters. Sponsors should indicate the mutation frequency cutoff that was used to distinguish true mutations from background (PCR/sequencing) errors. The optimal mutation frequency cutoff should be determined during NGS assay development and validation and varies by platform but typically ranges from 5 to 15 percent, although other factors (e.g., clinical sample type or input template concentration) may influence this cutoff. Sponsors should describe any other criteria used to identify true mutations (e.g., based on average mutation quality or read direction bias).
- d. Document any PCR/sequencing error correction algorithms applied before variant detection.

4.3.2.2 De novo assembly of contigs and identifying variants (when applicable)

Sponsors should consider the following for inclusion in the data analysis methods on the de novo assembly of contigs and identifying variants:

- a. Indicate the bioinformatics tool(s) used to conduct the de novo assembly of short reads into contigs and describe settings for key parameters. Sponsors should provide specific information about the tool used for contig assembly, including but not limited to percentage of mismatches tolerated, number of insertions and deletions tolerated, and identity and similarity cutoffs.
- b. Indicate the bioinformatics tool(s) used to compare contigs and identify and filter variants and describe settings for key parameters. Describe how variants were distinguished from background (PCR/sequencing) errors.

Contains Nonbinding Recommendations

- c. Indicate the bioinformatics tool(s) and the version number(s) that were used to investigate similarity to known sequences. Identify the database that was queried, provide the search criteria that were used, define the cutoffs used for percent identity, and provide details on how contaminants or other viruses were identified and reported.

5.0 NGS REPORT

Sponsors should include a summary of NGS resistance results in the appropriate Clinical Study Report but refer to a separate NGS Report for detailed NGS results.

5.1 Reporting Results From Read Mapping and Variant Calling

Sponsors should consider the following for inclusion in the NGS Report:

- a. Provide a brief summary of the NGS protocol (including data analysis methods). Include the criteria used for resistance testing, the minimum DNA/RNA copy number for attempting NGS (with brief rationale), and the mutation frequency cutoff used to distinguish true mutations from background (PCR/sequencing) errors (with brief rationale).
- b. Provide an amino acid frequency table (see section 7.0 for an example). These tables will generally be too large for inclusion in the NGS Report and should be provided as a standalone document (in CSV, XLSX, or XPT format).
- c. Provide summary tables of baseline polymorphisms or TES of interest. Sponsors should identify substitutions that they consider to be associated with antiviral resistance and provide their rationale (e.g., based on relative rates of the substitution across study arms, structural information, biochemical/cell culture data, amino acid conservation, published literature).

5.2 Reporting Results From the De Novo Assembly and Variant Calling

Sponsors should consider the following for inclusion in the NGS Report:

- a. Provide contig statistics showing the number of contigs generated from the reads.
- b. Provide a summary of the coverage for each contig.
- c. Provide the number of reads in each contig and the number of reads that did not form contigs.
- d. Describe contaminants or other viruses that were detected during annotation and report the percent representation of each contaminant.

Contains Nonbinding Recommendations

- e. Provide an amino acid frequency table (see section 7.0, Amino Acid Frequency Table Example). These tables will generally be too large for inclusion in the NGS Report and should be provided as a standalone document (in CSV, XLSX, or XPT format).
- f. Provide summary tables of baseline polymorphisms or TES of interest.

5.3 Additional Information

The following additional information is recommended when applicable or requested by the Division:

- a. A summary of missing sequencing data (e.g., numbers/percentages of participants across visits and arms) and the reasons for missing data (e.g., sample not collected, insufficient nucleic acid levels, nucleic acid amplification failed, library preparation failed, sequencing failed, sequencing results did not meet quality control metrics). Every reasonable attempt should be made to minimize the amount of missing data, and wherever possible, failed nucleic acid extraction, reverse transcription (if applicable), amplification, library preparation, and/or sequencing should be repeated at least once.
- b. A summary of basic quality/mapping metrics for all samples that were sequenced. Examples of such metrics include number of reads, number of bases, average read length, average read quality, average base quality, percent of reads with average quality greater than 30, percent of bases with quality scores greater than 30, average number of degenerate bases/read, number of degenerate bases in consensus sequence, percent of mapped reads, minimum/mean/median/maximum target gene sequence coverage, or percent of target sequence with coverage $\geq 100/\geq 1,000/\geq 5,000$. This information can be included in the NGS Report or provided as a standalone table (in CSV, XLSX, or XPT format). Of note, the values above are specific to short read technologies, and quality/coverage cutoffs may vary depending on the sequencing technology used.
- c. A summary of analyses performed to ensure sequencing data integrity (e.g., to rule out potential sample switching or contamination at the study site or central laboratory). These analyses could include comparisons of consensus sequences to each other, to the reference sequence, and to the positive control sequence (e.g., using phylogenetics approaches).

6.0 NGS FILE TYPES AND SUBMISSION PROCEDURES

6.1 Recommended File Formats and Submission Methods

Recommended file formats and submission methods for NGS-related files are shown in Table 1. Sponsors should submit the NGS protocol, NGS Report, reference sequence(s), and primer/adaptor sequences via the FDA Electronic Submissions Gateway (ESG), whereas the raw NGS data (FASTQ format), consensus sequences (FASTA format), and checksum file should be submitted to the Division on an external hard drive (EHD) following the guidelines outlined in the technical specifications document *Transmitting Electronic Submissions Using eCTD*

Contains Nonbinding Recommendations

Specifications (June 2017).² The cover letter and amino acid frequency table should be included in both the ESG and EHD submissions. The cover letter should indicate the total number of samples sequenced, the total number of participants that had samples sequenced, and the total number of FASTQ files submitted. The cover letter should also include a table of contents that describes all submitted files. An md5/checksum file should be provided so that the Division can verify that the files were fully transferred. Hard drive submissions should not include .exe files, which may result in rejection of the submission. If the hard drive is password protected, which is not required or recommended at this time, sponsors should provide the password to the Division ahead of time to allow coordination with the appropriate personnel in the CDER document room.

Table 1. Recommended File Formats and Submission Methods for NGS-Related Files.

NGS-Related File	Recommended Format(s)	Recommended Submission Method(s)
Cover Letter	PDF	EHD and ESG (Module 5.3.5.4)
NGS Protocol	PDF	ESG (Module 5.3.5.4)
NGS Report	PDF	ESG (Module 5.3.5.4)
Amino Acid Frequency Table	CSV, XLSX, XPT	EHD and ESG (Module 5.3.5.4)
Raw NGS Data	FASTQ	EHD*
Consensus Sequences	FASTA	EHD*
Reference Sequence(s)	FASTA	ESG (Module 5.3.5.4)
Primer/Adapter Sequences	CSV, FASTA, TXT, XLSX	ESG (Module 5.3.5.4)
Checksum File	MD5	EHD

*Small NGS datasets can be submitted via EHD or the ESG (Module 5.3.5.4).

6.2 Naming FASTQ Files

FASTQ file names should include participant identification number, visit, and read number (if applicable). For example, FASTQ files from three samples (Baseline, Week 16, and Follow-Up Week 4) from participant 001 in study ABC123, analyzed by paired-end sequencing, could be named as follows:

- a. ABC123-001.Baseline.R1.FASTQ, ABC123-001.Baseline.R2.FASTQ
- b. ABC123-001.Week16.R1.FASTQ, ABC123-001.Week16.R2.FASTQ
- c. ABC123-001.FUWeek4.R1.FASTQ, ABC123-001.FUWeek4.R2.FASTQ

If a sample is re-sequenced due to low sequence quality, we recommend that sponsors submit only the highest quality results to the Division. However, in instances where multiple FASTQ files need to be submitted for a single sample, the FASTQ files could be named as follows:

- a. ABC123-001.Baseline.R1.FASTQ, ABC123-001.Baseline.R2.FASTQ
- b. ABC123-001.Baseline.02.R1.FASTQ, ABC123-001.Baseline.02.R2.FASTQ

² Available at <https://www.fda.gov/media/76812/download?attachment>. See also the Electronic Regulatory Submission and Review web page at <https://www.fda.gov/Drugs/DevelopmentApprovalProcess/FormsSubmissionRequirements/ElectronicSubmissions/default.htm>.

Contains Nonbinding Recommendations

6.3 Naming Consensus Sequences (FASTA Format)

Consensus sequences for each participant and time point should be provided in a single FASTA file for each clinical trial. The study identification number should be included in the FASTA file name. Sequence identifiers in the FASTA file should include participant identification number and visit. If multiple genes are sequenced separately, sponsors should also include the gene name in the sequence identifiers. Lastly, if different reference sequences are used for mapping reads from different virus genotypes/subtypes, sponsors should also include the virus genotype/subtype in the sequence identifiers.

7.0 AMINO ACID FREQUENCY TABLE EXAMPLE

The table below is a hypothetical example of an amino acid frequency table. The table should include all amino acid substitutions that differ from the reference sequence at frequencies $\geq 1\%$. Additional columns that may be helpful include study site identification number, flags for analysis populations and/or stratification factors, flags for virologic and clinical outcomes (e.g., viral breakthrough or hospitalization/death), nucleotide position, viral variant/genotype, flags for baseline polymorphisms and treatment-emergent substitutions, and viral DNA/RNA levels in the sequenced sample. For drugs that target nucleic acids and are impacted by nucleotide mutations (e.g., small interfering RNAs and antisense oligonucleotides), the frequency table should include all nucleotide changes, even if they do not result in an amino acid change.

Table 2. Example of an Amino Acid Frequency Table.

STUDYID	SUBJID	VISIT	ARM	AAPOS	AAREF	AASUB	AACHANGE	TCOV	VCOV	AAFREQ
ABC123	001	D2	Placebo	81	R	K	R81K	4317	156	0.036
ABC123	001	BL	Placebo	98	K	R	K98R	2841	99	0.035
ABC123	001	D2	Placebo	98	K	R	K98R	9487	366	0.039
ABC123	001	D3	Placebo	98	K	R	K98R	9474	378	0.040
ABC123	001	BL	Placebo	120	R	Q	R120Q	4310	200	0.046
ABC123	001	D2	Placebo	120	R	Q	R120Q	12722	470	0.037
ABC123	001	D3	Placebo	120	R	Q	R120Q	12466	489	0.039
ABC123	001	BL	Placebo	147	I	V	I147V	3456	742	0.215
ABC123	001	D2	Placebo	147	I	V	I147V	13456	2709	0.201
ABC123	001	D3	Placebo	147	I	V	I147V	13297	1934	0.145
ABC123	001	BL	Placebo	150	A	V	A150V	3107	43	0.014

STUDYID = study identification number; **SUBJID** = subject identification number; **VISIT** = study visit at which the sample was collected; **ARM** = treatment arm; **AAPOS** = amino acid position in the target protein; **AAREF** = amino acid present at this position in the reference sequence; **AASUB** = amino acid substitution detected by sequencing; **AACHANGE** = amino acid change; **TCOV** = total coverage at the nucleotide site; **VCOV** = variant coverage at the nucleotide site; **AAFREQ** = amino acid frequency of the substitution detected; **BL** = baseline; **D** = day.