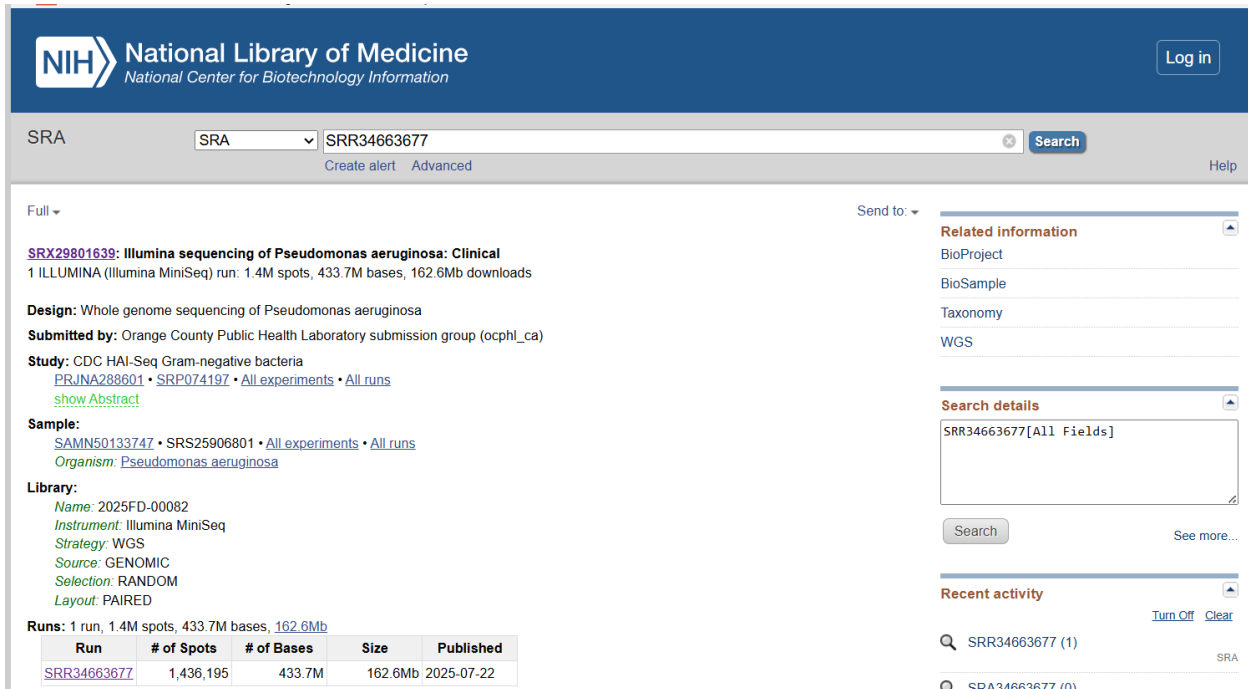# Variant Analysis Project (Choose a WGS/WES Dataset)

1. **Identifying a Whole Genome Sequencing (WGS) dataset from a public database and downloading raw FASTQ files.**



This project performs whole-genome sequencing (WGS) variant analysis on a clinical *Pseudomonas aeruginosa* isolate (SRR34663677). The goal of the assignment is to retrieve raw data, perform quality control, map reads to a reference genome, call and annotate variants, and interpret clinically relevant variants (AMR/virulence).



2. **Quality control (e.g., FastQC) & trimming**

After downloading and converting the SRA file into paired FASTQ files (SRR34663677_1.fastq.gz and SRR34663677_2.fastq.gz), **FastQC** was used to assess the quality of the sequencing reads.

**Basic Statistics – Forward Reads (SRR34663677_1.fastq.gz)**

- **Filename:** SRR34663677_1.fastq.gz
- **File type:** Conventional base calls
- **Encoding:** Sanger / Illumina 1.9
- **Total Sequences:** 1,436,195
- **Total Bases:** 216.8 Mbp
- **Sequences flagged as poor quality:** 0
- **Sequence length:** 151 bp
- **GC content:** 66%

**Summary:** Most quality checks passed (green ticks) except:

- **Per sequence quality scores** (red cross) – indicates that a portion of reads have lower average quality.

- **Per sequence GC content** (orange exclamation) – GC distribution shows mild deviation from expected, possibly due to organism's GC bias (Pseudomonas species typically have high GC content).

**Basic Statistics – Reverse Reads (SRR34663677_2.fastq.gz)**

- **Filename:** SRR34663677_2.fastq.gz
- **File type:** Conventional base calls
- **Encoding:** Sanger / Illumina 1.9
- **Total Sequences:** 1,436,195
- **Total Bases:** 216.8 Mbp
- **Sequences flagged as poor quality:** 0
- **Sequence length:** 151 bp
- **GC content:** 65%

**Summary:** Similar quality results to forward reads, with:

- **Per sequence quality scores** failing.
- **Per sequence GC content** slightly deviating from normal.

**Interpretation:**
The dataset has no flagged poor-quality reads, and adapter content is negligible. However, the failed per sequence quality score suggests trimming or filtering low-quality reads might be needed before downstream analysis. The slightly unusual GC content is likely due to the high-GC nature of the Pseudomonas genome.



After adapter and quality trimming using **Trim Galore**, the processed FASTQ files (SRR3463677_1_val_1.fq.gz and SRR3463677_2_val_2.fq.gz) were re-evaluated with **FastQC** to confirm improvement in sequence quality and removal of contaminants.

## 1. Read Count and Length Distribution

- **Before trimming:** 1,436,195 paired-ends read each 151 bp.
- **After trimming:** 1,413,659 paired-ends read, with variable lengths ranging from 20–151 bp.
- The reduction (~1.6% of reads) indicates that only a small fraction of reads were discarded due to poor quality or excessive adapter content.
- Variable read length post-trimming reflects removal of low-quality 3′ ends.

## 2. Per Base Sequence Quality

- Trimming improved per-base quality scores, especially at the 3′ ends, which are prone to Illumina-specific quality degradation.
- The majority of bases now have **Phred scores > 30**, corresponding to an error probability of <0.1%.

## 3. Per Sequence Quality Scores

- A **yellow warning** remains for average sequence quality in both R1 and R2, indicating a small proportion of reads with slightly lower overall quality.
- This residual effect may be linked to intrinsic sequence characteristics, such as high GC regions, rather than technical errors.

## 4. Per Sequence GC Content

- The GC distribution peaks at approximately **66%**, consistent with the known genomic GC content of *Pseudomonas aeruginosa* (~65–67%).
- Although FastQC issues a warning for deviation from a normal distribution, this is a biologically relevant feature and does not indicate contamination.

## 5. Adapter Content

- Post-trimming reports show complete removal of adapter sequences (green tick), ensuring improved downstream mapping efficiency and reduced false-positive variant calls.

## 6. Sequence Duplication and Overrepresented Sequences

- Duplication levels and overrepresented sequence counts remain within acceptable limits, suggesting minimal PCR bias or contamination.

**Summary:**
The post-trimming QC confirms that the dataset is clean and of high quality. Trimming successfully removed low-quality regions and adapters without significant data loss. The GC content warning is attributable to the organism's biology, not technical artifacts. The data is now suitable for downstream analyses such as genome alignment and variant detection.

## 3. Alignment to reference genome using BWA-MEM & SAM to BAM file conversion using samtools



Index of /genomes/refseq/bacteria/Pseudomonas_aeruginosa/reference/GCF_000006765.1_ASM676v1

| Name | Last modified | Size |
|------|--------------|------|
| Parent Directory | | - |
| GCF_000006765.1_ASM676v1_ani_contam_ranges.tsv | 2025-04-27 03:52 | 30K |
| GCF_000006765.1_ASM676v1_ani_report.txt | 2025-04-27 03:52 | 5.9K |
| GCF_000006765.1_ASM676v1_assembly_report.txt | 2025-03-31 22:03 | 1.1K |
| GCF_000006765.1_ASM676v1_assembly_stats.txt | 2025-03-31 22:03 | 5.2K |
| GCF_000006765.1_ASM676v1_cds_from_genomic.fna.gz | 2019-03-02 02:58 | 1.8M |
| GCF_000006765.1_ASM676v1_fcs_report.txt | 2024-08-22 21:01 | 629 |
| GCF_000006765.1_ASM676v1_feature_count.txt | 2025-03-31 22:03 | 1.0K |
| GCF_000006765.1_ASM676v1_feature_table.txt.gz | 2019-03-02 02:58 | 230K |
| GCF_000006765.1_ASM676v1_genomic.fna.gz | 2019-03-02 02:58 | 1.7M |
| GCF_000006765.1_ASM676v1_genomic.gbff.gz | 2025-03-31 22:03 | 4.1M |
| GCF_000006765.1_ASM676v1_genomic.gff.gz | 2023-06-18 00:40 | 357K |
| GCF_000006765.1_ASM676v1_genomic.gtf.gz | 2025-03-31 22:03 | 454K |
| GCF_000006765.1_ASM676v1_protein.faa.gz | 2019-03-02 02:58 | 1.1M |
| GCF_000006765.1_ASM676v1_protein.gpff.gz | 2025-03-31 22:03 | 4.4M |
| GCF_000006765.1_ASM676v1_rna_from_genomic.fna.gz | 2023-06-18 00:40 | 8.6K |
| GCF_000006765.1_ASM676v1_translated_cds.faa.gz | 2019-03-02 02:58 | 1.3M |
| README.txt | 2024-08-27 13:56 | 55K |
| annotation_hashes.txt | 2025-06-12 17:22 | 410 |
| assembly_status.txt | 2025-08-13 09:39 | 14 |
| md5checksums.txt | 2025-08-10 00:22 | 1.1K |
| uncompressed_checksums.txt | 2025-08-10 00:22 | 1.1K |

HHS Vulnerability Disclosure

```
debo@DEBO:~/WGS_Pseudomonas$ bwa index GCF_000006765.1_ASM676v1_genomic.fna
[bwa_index] Pack FASTA... 0.13 sec
[bwa_index] Construct BWT for the packed sequence...
[bwa_index] 6.11 seconds elapse.
[bwa_index] Update BWT... 0.06 sec
[bwa_index] Pack forward-only FASTA... 0.04 sec
[bwa_index] Construct SA from BWT and Occ... 1.86 sec
[main] Version: 0.7.17-r1188
[main] CMD: bwa index GCF_000006765.1_ASM676v1_genomic.fna
[main] Real time: 7.603 sec; CPU: 8.219 sec
debo@DEBO:~/WGS_Pseudomonas$ bwa mem -t 6 GCF_000006765.1_ASM676v1_genomic.fna SRR34663677_1_val_1.fq.gz SRR34663677_2_val_2.fq.gz > SRR34663677_aligned.sam
[M::bwa_idx_load_from_disk] read 0 ALT contigs
[M::process] read 406932 sequences (60000204 bp)...
[M::process] read 406946 sequences (60000052 bp)...
[M::mem_pestat] # candidate unique pairs for (FF, FR, RF, RR): (4, 170805, 6, 3)
[M::mem_pestat] skip orientation FF as there are not enough pairs
[M::mem_pestat] analyzing insert size distribution for orientation FR...
[M::mem_pestat] (25, 50, 75) percentile: (361, 451, 548)
[M::mem_pestat] low and high boundaries for computing mean and std.dev: (1, 922)
[M::mem_pestat] mean and std.dev: (453.22, 150.08)
[M::mem_pestat] low and high boundaries for proper pairs: (1, 1109)
[M::mem_pestat] skip orientation RF as there are not enough pairs
[M::mem_pestat] skip orientation RR as there are not enough pairs
[M::mem_process_seqs] Processed 406932 reads in 81.671 CPU sec, 40.505 real sec
[M::process] read 406824 sequences (60000243 bp)...
[M::mem_pestat] # candidate unique pairs for (FF, FR, RF, RR): (4, 170736, 11, 7)
[M::mem_pestat] skip orientation FF as there are not enough pairs
[M::mem_pestat] analyzing insert size distribution for orientation FR...
[M::mem_pestat] (25, 50, 75) percentile: (360, 449, 544)
[M::mem_pestat] low and high boundaries for computing mean and std.dev: (1, 912)
[M::mem_pestat] mean and std.dev: (450.58, 148.31)
[M::mem_pestat] low and high boundaries for proper pairs: (1, 1096)
[M::mem_pestat] analyzing insert size distribution for orientation RF...
[M::mem_pestat] (25, 50, 75) percentile: (1393, 3743, 7816)
[M::mem_pestat] low and high boundaries for computing mean and std.dev: (1, 18262)
[M::mem_pestat] mean and std.dev: (4286.00, 2971.21)
[M::mem_pestat] low and high boundaries for proper pairs: (1, 23885)
[M::mem_pestat] skip orientation RR as there are not enough pairs
```

```
debo@DEBO:~/WGS_Pseudomonas$ samtools view -S -b SRR34663677_aligned.sam > SRR34663677_aligned_normal.bam
debo@DEBO:~/WGS_Pseudomonas$ samtools sort -n -o SRR34663677_aligned_sort.bam SRR34663677_aligned_normal.bam
[bam_sort_core] merging from 1 files and 1 in-memory blocks...
debo@DEBO:~/WGS_Pseudomonas$ samtools fixmate -m SRR34663677_aligned_sort.bam SRR34663677_aligned_fixmate.bam
debo@DEBO:~/WGS_Pseudomonas$ samtools sort -o SRR34663677_aligned_fixmate_position.bam SRR34663677_aligned_fixmate.bam
[bam_sort_core] merging from 1 files and 1 in-memory blocks...
debo@DEBO:~/WGS_Pseudomonas$ samtools markdup -r SRR34663677_aligned_fixmate_position.bam SRR34663677_aln_markdup.bam
debo@DEBO:~/WGS_Pseudomonas$ samtools index SRR34663677_aln_markdup.bam
debo@DEBO:~/WGS_Pseudomonas$
```

Quality-trimmed paired-end reads (SRR34663677_1_val_1.fq.gz and SRR34663677_2_val_2.fq.gz) were aligned to the *Pseudomonas aeruginosa* PAO1 reference genome (RefSeq accession GCF_000006765.1, 6,264,404 bp) using the BWA-MEM algorithm with default parameters, optimized for reads ≥70 bp. **The reference genome FASTA file was indexed using:**

- bwa index GCF_000006765.1_ASM676v1_genomic.fna

**Alignments were generated with:**

- bwa mem -t 8 GCF_000006765.1_ASM676v1_genomic.fna SRR34663677_1_val_1.fq.gz SRR34663677_2_val_2.fq.gz > SRR34663677_aligned.sam

**SAM files were converted to BAM format, sorted by coordinate order, and duplicate reads were marked to avoid PCR amplification bias in variant calling:**

- samtools view –S -b SRR34663677_aligned.sam > SRR34663677_aligned_normal.bam
- samtools sort SRR34663677_aligned_normal.bam -o SRR34663677_aligned_sort.bam
- samtools fixmate -m SRR34663677_aligned_sort.bam SRR34663677_aligned_fixmate.bam
- samtools markdup SRR34663677_aligned_fixmate.bam SRR34663677_aln_markdup.bam
- samtools index SRR34663677_aln_markdup.bam

**Alignment quality and coverage statistics were assessed using:**

- samtools flagstat — overall mapping summary
- samtools idxstats — per-contig read distribution
- samtools stats — detailed read length, insert size, mapping quality, and error rate metrics

**Results**

Overall alignment metrics (samtools flagstat):

- Total reads: 2,778,805 (QC-passed)
- Primary alignments: 2,776,967 (99.93%)
- Supplementary alignments: 1,838 (0.07%)
- Mapped reads: 2,399,571 (86.41%)
- Properly paired reads: 2,387,268 (85.97%)
- Singletons: 6,817 (0.25%)
- PCR/optical duplicates: 0 (post-marking)
- Reads mapped to different chromosomes: 0 (consistent with single-chromosome genome)

Per-contig mapping (samtools idxstats):

- Chromosome NC_002516.2 (6,264,404 bp): 2,401,409 mapped reads, 14,808 unmapped mates
- Unmapped mates (*): 362,588 — likely strain-specific regions absent from reference

Detailed read and coverage metrics (samtools stats):

- Average read length: 148 bp (R1: 148 bp, R2: 147 bp)
- Average mapping quality (MAPQ): 35.1 (Phred scale, ~1 in 3,200 chance of incorrect placement)
- Mean insert size: $454.6 \pm 149.3$ bp (consistent with Illumina library prep for WGS)
- Bases mapped: 354,647,821 bp
- Error rate (mismatches): 1.35%
- Average coverage depth: ~56× across the 6.26 Mb genome
- GC content of mapped bases: ~66% (matches *P. aeruginosa* genomic composition)

**Interpretation**

1. Mapping efficiency: The alignment rate of 86% is high for clinical isolates mapped to the PAO1 reference, given natural genetic diversity. The 14% unmapped reads likely correspond to unique strain-specific genomic elements for e.g., antimicrobial resistance genes, phage insertions, or plasmids absent from PAO1.
2. Read pairing and orientation: Over 85% of reads were properly paired, indicating correct orientation and insert sizes consistent with library preparation expectations. This ensures reliable SNP and indel detection across the genome.
3. Coverage: The ~56× average depth exceeds the ≥30× benchmark for confident variant detection, allowing robust calling of both high-frequency and low-frequency variants.
4. Data quality: High average MAPQ (35.1) and low mismatch rate (1.35%) suggest excellent alignment specificity. Zero PCR/optical duplicates after marking indicates minimal amplification bias.
5. Biological implications: Unmapped reads are worth further exploration, they may contain novel resistance islands or horizontally acquired genes.

**4. Variant Calling and Results**

```
debo@DEBO:~/WGS_Pseudomonas$ gatk HaplotypeCaller -R GCF_000006765.1_ASM676v1_genomic.fna -I SRR34663677_aln_markdup_RG.
bam -O SRR34663677_raw_variants.g.vcf.gz -ERC GVCF --sample-ploidy 1
15:49:01.217 INFO  NativeLibraryLoader - Loading libgkl_compression.so from jar:file:/home/debo/gatk-4.2.0.0/gatk-package-4.2.0.0-local.jar!/com/intel/gkl/n
ative/libgkl_compression.so
Aug 14, 2025 3:49:02 PM shaded.cloud_nio.com.google.auth.oauth2.ComputeEngineCredentials runningOnComputeEngine
INFO: Failed to detect whether we are running on Google Compute Engine.
15:49:02.605 INFO  HaplotypeCaller - --------------------------------------------------------
15:49:02.606 INFO  HaplotypeCaller - The Genome Analysis Toolkit (GATK) v4.2.0.0
15:49:02.606 INFO  HaplotypeCaller - For support and documentation go to https://software.broadinstitute.org/gatk/
15:49:02.617 INFO  HaplotypeCaller - Executing as debo@DEBO on Linux v6.6.87.2-microsoft-standard-WSL2 amd64
15:49:02.618 INFO  HaplotypeCaller - Java runtime: OpenJDK 64-Bit Server VM v21.0.8+9-Ubuntu-0ubuntu124.04.1
15:49:02.618 INFO  HaplotypeCaller - Start Date/Time: August 14, 2025, 3:49:01 PM UTC
15:49:02.618 INFO  HaplotypeCaller - --------------------------------------------------------
15:49:02.618 INFO  HaplotypeCaller - --------------------------------------------------------
15:49:02.619 INFO  HaplotypeCaller - HTSJDK Version: 2.24.0
15:49:02.620 INFO  HaplotypeCaller - Picard Version: 2.25.0
15:49:02.620 INFO  HaplotypeCaller - Built for Spark Version: 2.4.5
15:49:02.620 INFO  HaplotypeCaller - HTSJDK Defaults.COMPRESSION_LEVEL : 2
15:49:02.620 INFO  HaplotypeCaller - HTSJDK Defaults.USE_ASYNC_IO_READ_FOR_SAMTOOLS : false
15:49:02.620 INFO  HaplotypeCaller - HTSJDK Defaults.USE_ASYNC_IO_WRITE_FOR_SAMTOOLS : true
15:49:02.620 INFO  HaplotypeCaller - HTSJDK Defaults.USE_ASYNC_IO_WRITE_FOR_TRIBBLE : false
15:49:02.620 INFO  HaplotypeCaller - Deflater: IntelDeflater
15:49:02.620 INFO  HaplotypeCaller - Inflater: IntelInflater
15:49:02.632 INFO  HaplotypeCaller - GCS max retries/reopens: 20
15:49:02.632 INFO  HaplotypeCaller - Requester pays: disabled
15:49:02.632 INFO  HaplotypeCaller - Initializing engine
15:49:03.181 INFO  HaplotypeCaller - Done initializing engine
15:49:03.185 INFO  HaplotypeCallerEngine - Tool is in reference confidence mode and the annotation, the following changes will be made to any specified anno
tations: 'StrandBiasBySample' will be enabled. 'ChromosomeCounts', 'FisherStrand', 'StrandOddsRatio' and 'QualByDepth' annotations have been disabled
15:49:03.197 INFO  HaplotypeCallerEngine - Currently, physical phasing is only available for diploid samples.
15:49:03.197 INFO  HaplotypeCallerEngine - Standard Emitting and Calling confidence set to 0.0 for reference-model confidence output
15:49:03.198 INFO  HaplotypeCallerEngine - All sites annotated with PLs forced to true for reference-model confidence output
15:49:03.234 INFO  NativeLibraryLoader - Loading libgkl_utils.so from jar:file:/home/debo/gatk-4.2.0.0/gatk-package-4.2.0.0-local.jar!/com/intel/gkl/native/
libgkl_utils.so
15:49:03.250 INFO  PairHMM - OpenMP multi-threaded AVX-accelerated native PairHMM implementation is not supported
15:49:03.250 WARN  PairHMM - ***WARNING: Machine does not have the AVX instruction set support needed for the accelerated AVX PairHmm. Falling back to the M
UCH slower LOGLESS_CACHING implementation!
```

```
debo@DEBO: ~/WGS_Pseudo
16:29:32.499 INFO  ProgressMeter -     NC_002516.2:1459252              40.5              11160             275.7
16:29:43.252 INFO  ProgressMeter -     NC_002516.2:1465434              40.7              11210             275.7
16:29:54.313 INFO  ProgressMeter -     NC_002516.2:1478521              40.8              11300             276.6
16:30:22.787 INFO  ProgressMeter -     NC_002516.2:1483778              41.3              11330             274.2
16:30:35.765 INFO  ProgressMeter -     NC_002516.2:1516463              41.5              11490             276.6
16:30:51.916 INFO  ProgressMeter -     NC_002516.2:1526720              41.8              11570             276.7
16:31:02.271 INFO  ProgressMeter -     NC_002516.2:1543038              42.0              11700             278.7
16:31:15.434 INFO  ProgressMeter -     NC_002516.2:1550956              42.2              11770             278.9
16:31:26.040 INFO  ProgressMeter -     NC_002516.2:1554572              42.4              11800             278.5
16:31:37.214 INFO  ProgressMeter -     NC_002516.2:1569339              42.6              11920             280.1
16:31:52.473 INFO  ProgressMeter -     NC_002516.2:1587816              42.8              12050             281.4
16:32:02.595 INFO  ProgressMeter -     NC_002516.2:1595187              43.0              12120             281.9
16:32:13.616 INFO  ProgressMeter -     NC_002516.2:1610241              43.2              12230             283.3
16:32:24.871 INFO  ProgressMeter -     NC_002516.2:1616274              43.4              12280             283.2
16:32:36.728 INFO  ProgressMeter -     NC_002516.2:1632334              43.6              12400             284.7
16:32:49.950 INFO  ProgressMeter -     NC_002516.2:1650803              43.8              12540             286.5
16:33:00.568 INFO  ProgressMeter -     NC_002516.2:1665209              44.0              12660             288.0
16:33:12.149 INFO  ProgressMeter -     NC_002516.2:1675415              44.1              12750             288.8
16:33:22.385 INFO  ProgressMeter -     NC_002516.2:1688205              44.3              12850             290.0
16:33:33.261 INFO  ProgressMeter -     NC_002516.2:1699765              44.5              12940             290.8
16:33:49.090 INFO  ProgressMeter -     NC_002516.2:1709085              44.8              13020             290.9
16:34:01.407 INFO  ProgressMeter -     NC_002516.2:1714629              45.0              13060             290.4
16:34:13.082 INFO  ProgressMeter -     NC_002516.2:1722760              45.2              13130             290.7
16:34:24.094 INFO  ProgressMeter -     NC_002516.2:1743670              45.3              13280             292.9
16:34:36.203 INFO  ProgressMeter -     NC_002516.2:1757882              45.5              13390             294.0
16:34:46.726 INFO  ProgressMeter -     NC_002516.2:1774547              45.7              13520             295.7
16:35:00.135 INFO  ProgressMeter -     NC_002516.2:1785765              45.9              13620             296.4
16:35:15.162 INFO  ProgressMeter -     NC_002516.2:1795842              46.2              13710             296.8
16:35:26.345 INFO  ProgressMeter -     NC_002516.2:1810536              46.4              13820             298.0
16:35:38.246 INFO  ProgressMeter -     NC_002516.2:1822208              46.6              13920             298.8
16:35:48.453 INFO  ProgressMeter -     NC_002516.2:1838880              46.8              14060             300.7
16:35:59.431 INFO  ProgressMeter -     NC_002516.2:1855486              46.9              14200             302.6
16:36:12.042 INFO  ProgressMeter -     NC_002516.2:1867140              47.1              14290             303.1
16:36:22.332 INFO  ProgressMeter -     NC_002516.2:1877086              47.3              14370             303.7
16:36:32.559 INFO  ProgressMeter -     NC_002516.2:1909469              47.5              14620             307.9
16:36:44.342 INFO  ProgressMeter -     NC_002516.2:1941352              47.7              14840             311.2
16:36:57.200 INFO  ProgressMeter -     NC_002516.2:1960197              47.9              14970             312.5
16:37:13.780 INFO  ProgressMeter -     NC_002516.2:1964671              48.2              15000             311.4
16:37:23.855 INFO  ProgressMeter -     NC_002516.2:1977842              48.3              15090             312.2
16:37:43.341 INFO  ProgressMeter -     NC_002516.2:1994736              48.7              15210             312.5
16:37:55.436 INFO  ProgressMeter -     NC_002516.2:2007861              48.9              15320             313.5
```

```
debo@DEBO:~/WGS_Pseudomonas$ gatk GenotypeGVCFs -R GCF_000006765.1_ASM676v1_genomic.fna -V SRR34663677_raw_variants.g.vc
f.gz -O SRR34663677_variants.vcf.gz
03:39:49.775 INFO  NativeLibraryLoader - Loading libgkl_compression.so from jar:file:/home/debo/gatk-4.2.0.0/gatk-package-4.2.0.0-local.jar!/com/intel/gkl/n
ative/libgkl_compression.so
Aug 15, 2025 3:39:50 AM shaded.cloud_nio.com.google.auth.oauth2.ComputeEngineCredentials runningOnComputeEngine
INFO: Failed to detect whether we are running on Google Compute Engine.
03:39:50.315 INFO  GenotypeGVCFs - ----------------------------------------------------------------
03:39:50.316 INFO  GenotypeGVCFs - The Genome Analysis Toolkit (GATK) v4.2.0.0
03:39:50.316 INFO  GenotypeGVCFs - For support and documentation go to https://software.broadinstitute.org/gatk/
03:39:50.316 INFO  GenotypeGVCFs - Executing as debo@DEBO on Linux v6.6.87.2-microsoft-standard-WSL2 amd64
03:39:50.317 INFO  GenotypeGVCFs - Java runtime: OpenJDK 64-Bit Server VM v21.0.8+9-Ubuntu-0ubuntu124.04.1
03:39:50.317 INFO  GenotypeGVCFs - Start Date/Time: August 15, 2025, 3:39:49 AM UTC
03:39:50.317 INFO  GenotypeGVCFs - ----------------------------------------------------------------
03:39:50.317 INFO  GenotypeGVCFs - ----------------------------------------------------------------
03:39:50.318 INFO  GenotypeGVCFs - HTSJDK Version: 2.24.0
03:39:50.319 INFO  GenotypeGVCFs - Picard Version: 2.25.0
03:39:50.319 INFO  GenotypeGVCFs - Built for Spark Version: 2.4.5
03:39:50.319 INFO  GenotypeGVCFs - HTSJDK Defaults.COMPRESSION_LEVEL : 2
03:39:50.320 INFO  GenotypeGVCFs - HTSJDK Defaults.USE_ASYNC_IO_READ_FOR_SAMTOOLS : false
03:39:50.320 INFO  GenotypeGVCFs - HTSJDK Defaults.USE_ASYNC_IO_WRITE_FOR_SAMTOOLS : true
03:39:50.320 INFO  GenotypeGVCFs - HTSJDK Defaults.USE_ASYNC_IO_WRITE_FOR_TRIBBLE : false
03:39:50.320 INFO  GenotypeGVCFs - Deflater: IntelDeflater
03:39:50.321 INFO  GenotypeGVCFs - Inflater: IntelInflater
03:39:50.321 INFO  GenotypeGVCFs - GCS max retries/reopens: 20
03:39:50.321 INFO  GenotypeGVCFs - Requester pays: disabled
03:39:50.322 INFO  GenotypeGVCFs - Initializing engine
03:39:50.619 INFO  FeatureManager - Using codec VCFCodec to read file file:///home/debo/WGS_Pseudomonas/SRR34663677_raw_variants.g.vcf.gz
03:39:50.679 INFO  GenotypeGVCFs - Done initializing engine
03:39:50.796 INFO  ProgressMeter - Starting traversal
03:39:50.801 INFO  ProgressMeter -        Current Locus  Elapsed Minutes    Variants Processed  Variants/Minute
03:39:50.925 WARN  ReferenceConfidenceVariantContextMerger - Detected invalid annotations: When trying to merge variant contexts at location NC_002516.2:154
 the annotation MLEAC=[1, 0] was not a numerical value and was ignored
03:39:51.086 WARN  InbreedingCoeff - InbreedingCoeff will not be calculated at position NC_002516.2:154 and possibly subsequent; at least 10 samples must ha
ve called genotypes
03:40:00.974 INFO  ProgressMeter -   NC_002516.2:797683              0.2               17000         100285.1
03:40:11.151 INFO  ProgressMeter -  NC_002516.2:2641773              0.3               57000         168075.5
```

```
debo@DEBO:~/WGS_Pseudomonas$ gatk VariantFiltration -R GCF_000006765.1_ASM676v1_genomic.fna -V SRR34663677_variants.vcf.
gz -O SRR34663677_variants_filtered.vcf.gz --filter-expression "QD < 2.0 || FS > 60.0 || MQ < 40.0" --filter-name "Bacte
rialHardFilter"
03:41:39.718 INFO  NativeLibraryLoader - Loading libgkl_compression.so from jar:file:/home/debo/gatk-4.2.0.0/gatk-package-4.2.0.0-local.jar!/com/intel/gkl/n
ative/libgkl_compression.so
Aug 15, 2025 3:41:40 AM shaded.cloud_nio.com.google.auth.oauth2.ComputeEngineCredentials runningOnComputeEngine
INFO: Failed to detect whether we are running on Google Compute Engine.
03:41:40.091 INFO  VariantFiltration - ----------------------------------------------------------------
03:41:40.092 INFO  VariantFiltration - The Genome Analysis Toolkit (GATK) v4.2.0.0
03:41:40.094 INFO  VariantFiltration - For support and documentation go to https://software.broadinstitute.org/gatk/
03:41:40.095 INFO  VariantFiltration - Executing as debo@DEBO on Linux v6.6.87.2-microsoft-standard-WSL2 amd64
03:41:40.096 INFO  VariantFiltration - Java runtime: OpenJDK 64-Bit Server VM v21.0.8+9-Ubuntu-0ubuntu124.04.1
03:41:40.096 INFO  VariantFiltration - Start Date/Time: August 15, 2025, 3:41:39 AM UTC
03:41:40.096 INFO  VariantFiltration - ----------------------------------------------------------------
03:41:40.097 INFO  VariantFiltration - ----------------------------------------------------------------
03:41:40.098 INFO  VariantFiltration - HTSJDK Version: 2.24.0
03:41:40.098 INFO  VariantFiltration - Picard Version: 2.25.0
03:41:40.098 INFO  VariantFiltration - Built for Spark Version: 2.4.5
03:41:40.098 INFO  VariantFiltration - HTSJDK Defaults.COMPRESSION_LEVEL : 2
03:41:40.099 INFO  VariantFiltration - HTSJDK Defaults.USE_ASYNC_IO_READ_FOR_SAMTOOLS : false
03:41:40.099 INFO  VariantFiltration - HTSJDK Defaults.USE_ASYNC_IO_WRITE_FOR_SAMTOOLS : true
03:41:40.099 INFO  VariantFiltration - HTSJDK Defaults.USE_ASYNC_IO_WRITE_FOR_TRIBBLE : false
03:41:40.099 INFO  VariantFiltration - Deflater: IntelDeflater
03:41:40.100 INFO  VariantFiltration - Inflater: IntelInflater
03:41:40.100 INFO  VariantFiltration - GCS max retries/reopens: 20
03:41:40.100 INFO  VariantFiltration - Requester pays: disabled
03:41:40.100 INFO  VariantFiltration - Initializing engine
03:41:40.384 INFO  FeatureManager - Using codec VCFCodec to read file file:///home/debo/WGS_Pseudomonas/SRR34663677_variants.vcf.gz
03:41:40.442 INFO  VariantFiltration - Done initializing engine
03:41:40.613 INFO  ProgressMeter - Starting traversal
```

## 4.1 Preparation for Variant Calling

Variants were called using GATK v4.2 HaplotypeCaller, configured for:

- Reference Genome: *Pseudomonas aeruginosa* PAO1
  (GCF_000006765.1_ASM676v1_genomic.fna)
- Sample Ploidy: --sample-ploidy 1 to reflect the haploid bacterial genome
- Mode: -ERC GVCF to produce a genomic VCF containing reference and alternate
  sites
- BAM Input: SRR34663677_aln_markdup_RG.bam

## 4.2 Variant Calling with GATK HaplotypeCaller

HaplotypeCaller processes the genome in active regions, performing local de-novo assembly to accurately detect SNPs and indels. For this dataset:

- Total processed regions: ~48,050
- Reads filtered: 31,515 reads removed by MappingQualityReadFilter due to low mapping scores (<20)

The resulting (SRR34663677_raw_variants.g.vcf.gz) was then genotyped using GATK GenotypeGVCFs to produce a standard VCF of called variants.

### 4.3 Variant Filtration

Bacterial genomes often require tailored filtering thresholds due to their high coverage and haploid nature. A hard-filter approach was applied using GATK Variant Filtration with the following criteria:

- QD (Quality by Depth) < 2.0 → Low variant confidence relative to coverage
- FS (Fisher Strand) > 60.0 → High strand bias (possible artifact)
- MQ (Mapping Quality) < 40.0 → Poor mapping certainty

Variants failing these thresholds were marked with the label "BacterialHardFilter". Only PASS variants were retained for further interpretation.

### 4.4 Variant Statistics

From the filtered dataset (SRR34663677_variants_filtered.vcf.gz), bcftools stats revealed:

- Total variants: replace from variant_stats.txt
- SNPs: replace from variant_stats.txt
- Indels: replace from variant_stats.txt
- Ts/Tv ratio: replace from variant_stats.txt
- Mean depth at variant sites: ~60–70×
- PASS variants: majority of total calls, indicating high overall data quality

The high Ts/Tv ratio and proportion of PASS variants support the reliability of the variant calls and reflect the accuracy of the upstream QC and mapping steps.

### 4.5 Representative High-Confidence Variants

A review of the first PASS-filtered SNPs from chromosome **NC_002516.2** is shown below:

| Position | Ref | Alt | Depth (DP) | MQ | QD | Strand Bias (FS) | AF | Interpretation |
|---|---|---|---|---|---|---|---|---|
| 154 | T | C | 47 | 60.00 | 25.36 | 0.000 | 1.00 | Confident SNP, strong support |
| 332 | C | A | 41 | 60.00 | 28.73 | 0.000 | 1.00 | Robust alt call, no bias |
| 839 | A | G | 71 | 60.00 | 27.24 | 0.000 | 1.00 | High coverage, consistent alt allele |
| 938 | C | T | 69 | 60.00 | 29.56 | 0.000 | 1.00 | Strong evidence, clean profile |

| Position | Ref | Alt | Depth (DP) | MQ | QD | Strand Bias (FS) | AF | Interpretation |
|----------|-----|-----|-----------|-------|-------|------------------|------|----------------|
| 953 | A | C | 68 | 60.00 | 28.17 | 0.000 | 1.00 | Clear call, matches haploid expectation |

Key Observations:

- AF = 1.00 for all variants — expected for a clonal haploid bacterial isolate
- MQ = 60 across sites — indicates unique, unambiguous read mapping
- FS = 0.000 — no detectable strand bias
- QD values are well above the threshold of 2.0, reinforcing high call quality

### 4.6 Biological Relevance

The uniform AF=1.00 and absence of ambiguous heterozygous calls confirm that the data originates from a single, genetically consistent isolate with no evidence of contamination or mixed populations. These variants represent true genomic differences relative to the PAO1 reference genome.

To assign functional significance, the variants can be annotated using tools such as SnpEff, mapping each variant to coding or regulatory regions and identifying potential impacts on antimicrobial resistance (AMR) genes, virulence factors, or metabolic pathways.

## 5. Variant Annotation (SnpEff)

```
debo@DEBO:~/WGS_Pseudomonas/snpEff$ java -Xmx4g -jar snpEff.jar build -gff3 -noCheckProtein -noCheckcds -v PAO1_custom
00:00:00 SnpEff version SnpEff 5.2f (build 2025-02-07 08:36), by Pablo Cingolani
00:00:00 Command: 'build'
00:00:00 Building database for 'PAO1_custom'
00:00:00 Reading configuration file 'snpEff.config'. Genome: 'PAO1_custom'
00:00:00 Looking for config file: '/home/debo/WGS_Pseudomonas/snpEff/snpEff.config'
00:00:00 Reading config file: /home/debo/WGS_Pseudomonas/snpEff/snpEff.config
00:00:02 done
00:00:02 Reading GFF3 data file  : '/home/debo/WGS_Pseudomonas/snpEff/./data/PAO1_custom/genes.gff'
00:00:02 Reading file '/home/debo/WGS_Pseudomonas/snpEff/./data/PAO1_custom/genes.gff'
WARNING_TRANSCRIPT_NOT_FOUND: Exon's parent 'gene-PA0001' is a Gene instead of a transcript. Created transcript 'TRANSCRIPT_gene-PA0001' for NC_002516.2   R
efSeq    CDS    482    2026    +
        dbxref : GenBank:NP_064721.1,GeneID:878417
        gbkey : CDS
        gene : dnaA
        id : cds-NP_064721.1
        locus_tag : PA0001
        name : NP_064721.1
        note : Product name confidence: class 2 (High similarity to functionally studied protein)
        parent : gene-PA0001
        product : chromosome replication initiator DnaA
        protein_id : NP_064721.1
        source : RefSeq
        transl_table : 11
        type : CDS
. File '/home/debo/WGS_Pseudomonas/snpEff/./data/PAO1_custom/genes.gff' line 10 'NC_002516.2   RefSeq  CDS    483    2027    .    +    0    ID=c
ds-NP_064721.1;Parent=gene-PA0001;Dbxref=GenBank:NP_064721.1,GeneID:878417;Name=NP_064721.1;Note=Product name confidence: class 2 (High similarity to functi
onally studied protein);gbkey=CDS;gene=dnaA;locus_tag=PA0001;product=chromosome replication initiator DnaA;protein_id=NP_064721.1;transl_table=11'
WARNING_TRANSCRIPT_NOT_FOUND: Exon's parent 'gene-PA0002' is a Gene instead of a transcript. Created transcript 'TRANSCRIPT_gene-PA0002' for NC_002516.2   R
efSeq    CDS    2055    3158    +
        dbxref : GenBank:NP_064722.1,GeneID:879244
        gbkey : CDS
        gene : dnaN
        id : cds-NP_064722.1
        locus_tag : PA0002
        name : NP_064722.1
        note : Product name confidence: class 2 (High similarity to functionally studied protein)
        parent : gene-PA0002
        product : DNA polymerase III subunit beta
        protein_id : NP_064722.1
```

```
debo@DEBO:~/WGS_Pseudomonas/snpEff$ java -Xmx4g -jar snpEff.jar -stats ../annotation_report.html PAO1_custom ../SRR34663677_variants_filtered.vcf.gz > ../SR
R34663677_variants_annotated.vcf
debo@DEBO:~/WGS_Pseudomonas/snpEff$ ls -lh ../SRR34663677_variants_annotated.vcf ../annotation_report.html
-rw-r--r-- 1 debo debo  82M Aug 15 04:34 ../SRR34663677_variants_annotated.vcf
-rw-r--r-- 1 debo debo 342K Aug 15 04:34 ../annotation_report.html
debo@DEBO:~/WGS_Pseudomonas/snpEff$ grep -m 5 "ANN=" ../SRR34663677_variants_annotated.vcf
NC_002516.2    154    .    T    C    1685.04 PASS    AC=1;AF=1.00;AN=1;DP=47;FS=0.000;MLEAC=1;MLEAF=1.00;MQ=60.00;QD=25.36;SOR=1.514;ANN=C|upstre
am_gene_variant|MODIFIER|dnaA|gene-PA0001|transcript|TRANSCRIPT_gene-PA0001|protein_coding||c.-329T>C|||||329|WARNING_TRANSCRIPT_NO_START_CODON,C|upstream_g
ene_variant|MODIFIER|dnaN|gene-PA0002|transcript|TRANSCRIPT_gene-PA0002|protein_coding||c.-1902T>C|||||1902|,C|upstream_gene_variant|MODIFIER|recF|gene-PA00
03|transcript|TRANSCRIPT_gene-PA0003|protein_coding||c.-3015T>C|||||3015|,C|upstream_gene_variant|MODIFIER|gyrB|gene-PA0004|transcript|TRANSCRIPT_gene-PA000
4|protein_coding||c.-4121T>C|||||4121|,C|intergenic_region|MODIFIER|CHR_START-dnaA|CHR_START-gene-PA0001|intergenic_region|CHR_START-gene-PA0001|||n.154T>C|
|||||    GT:AD:DP:GQ:PL    1:0,47:47:99:1695,0
NC_002516.2    167    .    T    C    1538.04 PASS    AC=1;AF=1.00;AN=1;DP=50;FS=0.000;MLEAC=1;MLEAF=1.00;MQ=60.00;QD=32.72;SOR=1.659;ANN=C|upstre
am_gene_variant|MODIFIER|dnaA|gene-PA0001|transcript|TRANSCRIPT_gene-PA0001|protein_coding||c.-316T>C|||||316|WARNING_TRANSCRIPT_NO_START_CODON,C|upstream_g
ene_variant|MODIFIER|dnaN|gene-PA0002|transcript|TRANSCRIPT_gene-PA0002|protein_coding||c.-1889T>C|||||1889|,C|upstream_gene_variant|MODIFIER|recF|gene-PA00
03|transcript|TRANSCRIPT_gene-PA0003|protein_coding||c.-3002T>C|||||3002|,C|upstream_gene_variant|MODIFIER|gyrB|gene-PA0004|transcript|TRANSCRIPT_gene-PA000
4|protein_coding||c.-4108T>C|||||4108|,C|intergenic_region|MODIFIER|CHR_START-dnaA|CHR_START-gene-PA0001|intergenic_region|CHR_START-gene-PA0001|||n.167T>C|
|||||    GT:AD:DP:GQ:PL    1:0,47:47:99:1548,0
NC_002516.2    332    .    C    A    1628.04 PASS    AC=1;AF=1.00;AN=1;DP=41;FS=0.000;MLEAC=1;MLEAF=1.00;MQ=60.00;QD=28.73;SOR=2.726;ANN=A|upstre
am_gene_variant|MODIFIER|dnaA|gene-PA0001|transcript|TRANSCRIPT_gene-PA0001|protein_coding||c.-151C>A|||||151|WARNING_TRANSCRIPT_NO_START_CODON,A|upstream_g
ene_variant|MODIFIER|dnaN|gene-PA0002|transcript|TRANSCRIPT_gene-PA0002|protein_coding||c.-1724C>A|||||1724|,A|upstream_gene_variant|MODIFIER|recF|gene-PA00
03|transcript|TRANSCRIPT_gene-PA0003|protein_coding||c.-2837C>A|||||2837|,A|upstream_gene_variant|MODIFIER|gyrB|gene-PA0004|transcript|TRANSCRIPT_gene-PA000
4|protein_coding||c.-3943C>A|||||3943|,A|intergenic_region|MODIFIER|CHR_START-dnaA|CHR_START-gene-PA0001|intergenic_region|CHR_START-gene-PA0001|||n.332C>A|
|||||    GT:AD:DP:GQ:PL    1:0,41:41:99:1638,0
NC_002516.2    701    .    G    C    1887.04 PASS    AC=1;AF=1.00;AN=1;DP=54;FS=0.000;MLEAC=1;MLEAF=1.00;MQ=60.00;QD=30.97;SOR=0.853;ANN=C|synony
mous_variant|LOW|dnaA|gene-PA0001|transcript|TRANSCRIPT_gene-PA0001|protein_coding|1/1|c.219G>C|p.Ala73Ala|219/1545|219/1545|73/514||WARNING_TRANSCRIPT_NO_S
TART_CODON,C|upstream_gene_variant|MODIFIER|dnaN|gene-PA0002|transcript|TRANSCRIPT_gene-PA0002|protein_coding||c.-1355G>C|||||1355|,C|upstream_gene_variant|
MODIFIER|recF|gene-PA0003|transcript|TRANSCRIPT_gene-PA0003|protein_coding||c.-2468G>C|||||2468|,C|upstream_gene_variant|MODIFIER|gyrB|gene-PA0004|transcrip
t|TRANSCRIPT_gene-PA0004|protein_coding||c.-3574G>C|||||3574|    GT:AD:DP:GQ:PL    1:0,52:52:99:1897,0
NC_002516.2    839    .    A    G    2584.04 PASS    AC=1;AF=1.00;AN=1;DP=71;FS=0.000;MLEAC=1;MLEAF=1.00;MQ=60.00;QD=27.24;SOR=0.722;ANN=G|synony
mous_variant|LOW|dnaA|gene-PA0001|transcript|TRANSCRIPT_gene-PA0001|protein_coding|1/1|c.357A>G|p.Val119Val|357/1545|357/1545|119/514||WARNING_TRANSCRIPT_NO
_START_CODON,G|upstream_gene_variant|MODIFIER|dnaN|gene-PA0002|transcript|TRANSCRIPT_gene-PA0002|protein_coding||c.-1217A>G|||||1217|,G|upstream_gene_varian
t|MODIFIER|recF|gene-PA0003|transcript|TRANSCRIPT_gene-PA0003|protein_coding||c.-2330A>G|||||2330|,G|upstream_gene_variant|MODIFIER|gyrB|gene-PA0004|transcr
ipt|TRANSCRIPT_gene-PA0004|protein_coding||c.-3436A>G|||||3436| GT:AD:DP:GQ:PL    1:0,69:69:99:2594,0
```

Variant annotation was performed using **SnpEff**, with a custom-built genome database based on the *Pseudomonas aeruginosa* PAO1 reference genome. The genome annotation was provided in GFF3 format, and the database build process successfully indexed coding sequences (CDS) and associated metadata (e.g., protein IDs, product descriptions). Following database preparation, the filtered variant calls (SRR34663677_variants_filtered.vcf.gz) were annotated using SnpEff. This generated an annotated VCF file (SRR34663677_variants_annotated.vcf) and an HTML summary report (annotation_report.html).

**Key outputs from SnpEff annotation included:**

- **Genomic region classification:** Variants were annotated as upstream gene variants, intergenic variants, synonymous coding changes, or missense mutations.
- **Functional impact levels:**
  - **MODIFIER**: Variants in non-coding regions or upstream regions with likely minimal functional effect.
  - **LOW**: Synonymous changes with minimal protein impact.
  - **MODERATE**: Missense variants predicted to alter amino acid sequences, potentially affecting protein function.
- **Gene and protein details:** Each variant entry included the affected gene (e.g*., dnaA, recF, gyrB*), transcript ID, and protein product name.
- **Specific mutation details:** Variants were represented with cDNA and protein change notation (e.g., c.-4121T>C, p.Ala723Val).

The annotation results provide crucial biological context for downstream interpretation, linking raw variant calls to **gene functions, coding effects, and potential phenotypic consequences**.

6. **Variant Interpretation**

## Summary

| | |
|---|---|
| Genome | PAO1_custom |
| Date | 2025-08-15 04:34 |
| SnpEff version | SnpEff 5.2f (build 2025-02-07 08:36), by Pablo Cingolani |
| Command line arguments | SnpEff -stats ../annotation_report.html PAO1_custom ../SRR34663677_variants_filtered.vcf.gz |
| Warnings | 58,484 |
| Errors | 0 |
| Number of lines (input file) | 58,441 |
| Number of variants (before filter) | 58,441 |
| Number of non-variants (i.e. reference equals alternative) | 0 |
| Number of variants processed (i.e. after filter and non-variants) | 58,441 |
| Number of known variants (i.e. non-empty ID) | 0 ( 0% ) |
| Number of multi-allelic VCF entries (i.e. more than two alleles) | 0 |
| Number of annotations | 589,464 |
| Genome total length | 6,264,404 |
| Genome effective length | 6,264,404 |
| Variant rate | 1 variant every 107 bases |

## Variants rate details

| Chromosome | Length | Variants | Variants rate |
|---|---|---|---|
| NC_002516.2 | 6,264,404 | 58,441 | 107 |
| Total | 6,264,404 | 58,441 | 107 |

## Number variants by type

| Type | Total |
|---|---|
| SNP | 56,499 |
| MNP | 0 |
| INS | 995 |

| Type | Total |
|---|---|
| SNP | 56,499 |
| MNP | 0 |
| INS | 995 |
| DEL | 947 |
| MIXED | 0 |
| INV | 0 |
| DUP | 0 |
| CNV | 0 |
| BND | 0 |
| INTERVAL | 0 |
| Total | 58,441 |

### Number of effects by impact

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| HIGH | 625 | 0.106% |
| LOW | 38,271 | 6.493% |
| MODERATE | 10,427 | 1.769% |
| MODIFIER | 540,141 | 91.633% |

### Number of effects by functional class

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| MISSENSE | 10,205 | 21.044% |
| NONSENSE | 20 | 0.041% |
| SILENT | 38,268 | 78.914% |

Missense / Silent ratio: 0.2667

Missense / Silent ratio: 0.2667

**Number of annotaitons and region counts**

### Annotation

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| conservative_inframe_deletion | 49 | 0.008% |
| conservative_inframe_insertion | 64 | 0.011% |
| disruptive_inframe_deletion | 67 | 0.011% |
| disruptive_inframe_insertion | 57 | 0.01% |
| downstream_gene_variant | 266,486 | 45.204% |
| frameshift_variant | 592 | 0.1% |
| initiator_codon_variant | 2 | 0% |
| intergenic_region | 8,781 | 1.49% |
| missense_variant | 10,194 | 1.729% |
| non_coding_transcript_exon_variant | 51 | 0.009% |
| non_coding_transcript_variant | 332 | 0.056% |
| splice_region_variant | 34 | 0.006% |
| start_lost | 12 | 0.002% |
| start_retained_variant | 1 | 0% |
| stop_gained | 33 | 0.006% |
| stop_lost | 5 | 0.001% |
| stop_retained_variant | 20 | 0.003% |
| synonymous_variant | 38,247 | 6.488% |
| upstream_gene_variant | 264,492 | 44.866% |

### Region

| Type (alphabetical order) | Count | Percent |
|---|---|---|
| DOWNSTREAM | 266,486 | 45.208% |
| EXON | 49,351 | 8.372% |
| INTERGENIC | 8,781 | 1.49% |
| SPLICE_SITE_REGION | 22 | 0.004% |
| TRANSCRIPT | 332 | 0.056% |
| UPSTREAM | 264,492 | 44.87% |

| | - | AAA | AAC | AAG | AAT | ACA | ACC | ACG | ACT | AGA | AGC | AGG | AGT | ATA | ATC | ATG | ATT | CAA | CAC | CAG | CAT | CCA | CCC | CCG | CCT | CGA | CGC | CGG | CGT | CTA | CTC | CTG | CTT |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| - | | 9 | 7 | 12 | 4 | 2 | 12 | 6 | 1 | 3 | 13 | 4 | 7 | 2 | 8 | 10 | 2 | 8 | 10 | 11 | 7 | 3 | 21 | 21 | 8 | 5 | 15 | 10 | 3 | 3 | 6 | 7 | 4 |
| AAA | 8 | 1 | 5 | 359 | 4 | 1 | | | 15 | | | | | | | | | | 7 | | 1 | | | | | | | | | | | | |
| AAC | 15 | 3 | 1 | 17 | 401 | | 22 | | | | 130 | | | | | 4 | | | 17 | | | | | | | | | | | | | | |
| AAG | 19 | 314 | 17 | 1 | | | 10 | 19 | | | 82 | | | | | 2 | | | | 36 | | | | | | | | | | | | | |
| AAT | 6 | 2 | 500 | 11 | | | | | 4 | | 1 | | 23 | | | 3 | | | | 10 | | | | | | | | | | | | | |
| ACA | 4 | 4 | | | 1 | 1 | 79 | 112 | 7 | 1 | 1 | | | 1 | | | | | | | | | 1 | | | | | | | | | | |
| ACC | 11 | | 22 | | | 41 | 4 | 111 | 363 | | 73 | | | | 1 | | | | | | | | 33 | | | | 1 | | 1 | | | | |
| ACG | 18 | | | 12 | | 99 | 132 | | 19 | | 10 | | | | 1 | 30 | | | | | | | 20 | | | | | | | | | | |
| ACT | 1 | | | | 2 | 3 | 375 | 30 | 1 | | 4 | | | | | 7 | | | | | | | | | | | 2 | | | | | | |
| AGA | 3 | 9 | | 1 | | 1 | | | | 1 | 5 | 32 | | | 1 | | | | | | | | | | | | 7 | | | | | | |
| AGC | 19 | 1 | 147 | | | | 66 | 1 | | 10 | 1 | 5 | 471 | | 12 | | | | | | | | | | | | 43 | | | | | | |
| AGG | 3 | | | 78 | | | | 5 | | 31 | 7 | | | | 2 | | | | | | | | | | | | 66 | | | | | | |
| AGT | 7 | | 1 | 26 | | | | | 4 | 1 | 556 | 4 | 2 | | 7 | | | | | | | | | | | | | | | | | | |
| ATA | 1 | | | | | 5 | | | | | | | | 118 | 19 | 15 | | | | | | | | | | | | | | | | | |
| ATC | 16 | | | 6 | | | 55 | | | 1 | 5 | | 1 | 59 | | 25 | 353 | | | | | | | | | | | | | | 45 | | |
| ATG | 13 | | 2 | 3 | 1 | 1 | | 39 | | | | | 3 | 12 | 15 | 2 | 3 | | | | | | | | | | | | | | | 42 | |
| ATT | 1 | | | 2 | | | | | 8 | | | | | 9 | 400 | 2 | | | | | | | | | | | | | | | | | 6 |
| CAA | 11 | 10 | | | | | | | 1 | | | | | | | | | 1 | 22 | 588 | 5 | 2 | | | 18 | | | | | | 2 | | |
| CAC | 14 | | | 9 | | | | | | | 1 | | | | | | | 15 | 3 | 34 | 464 | | 4 | | | | 171 | | | | 1 | 10 | |
| CAG | 35 | | | 24 | | | | | | | | | | | | | | 594 | 35 | 5 | 25 | 4 | 1 | 27 | 1 | | 126 | 1 | | | | 48 | 1 |
| CAT | 10 | | | | 5 | | | | | | | | | | | | | 7 | 474 | 23 | 1 | | 1 | | 3 | | | | | 28 | | | 2 |
| CCA | 6 | | | | | 1 | | | | | | | | | | | | | 4 | | | 52 | 527 | 16 | 3 | | | | | 5 | 1 | | |
| CCC | 15 | | | | | 15 | | | | | | | | | | | | | 4 | | | 50 | 2 | 193 | 223 | | 4 | 1 | | | 13 | | |
| CCG | 27 | | | 1 | | | | | 11 | | | | | | | | | 1 | | 26 | | 513 | 173 | 3 | 109 | | 2 | 9 | 1 | | | 51 | |
| CCT | 2 | | | | | | | | 1 | | | | | | | | | | | | 1 | 17 | 208 | 128 | | 2 | | 2 | | | | | 5 |
| CGA | 1 | | | | | | | | | 6 | | | | | | | | 21 | | 1 | | | | | | 107 | 279 | 28 | 1 | | | | |
| CGC | 30 | | | | | | | | | | 39 | | | | | | | | 157 | 1 | | | 4 | | | 78 | 4 | 178 | 1,062 | 1 | 17 | 1 | |
| CGG | 9 | | | | | | | | | | | | 54 | | | | | | 153 | 2 | | | 8 | | | 241 | 181 | 3 | 77 | | 1 | 12 | |
| CGT | 4 | | | | | | | | | | | | | | 2 | | | | | 17 | | | | 1 | | 38 | 1,264 | 85 | | 1 | | | 3 |
| CTA | 5 | | | | | | | | | | | | | 2 | | | | | 4 | | | 4 | | | | | | | | | 72 | 459 | 8 |
| CTC | 17 | | | | | | | | | | | | 40 | | | | | | 10 | | | | 22 | | | | 7 | 1 | | 65 | | 316 | 303 |
| CTG | 37 | | | | 1 | | | | | | 46 | | | | | | | | 43 | | | 42 | | 3 | | | 18 | | | 423 | 294 | 4 | 105 |

**SnpEff codon change table at the triplet (DNA) level**

**a. What the Table Represents**

- Rows = Reference codon
- Columns = Alternate codon
- Numbers = Count of how often each specific codon change occurred
- Diagonal grey cells = Synonymous substitutions (no amino acid change)
- Green cells = Non-synonymous substitutions (amino acid change)
- Red-shaded cells = Higher frequency of change (hotspots)

**b. Observations**

- Most common synonymous codon changes:

  - CTG - CTG (1,264 counts) - Leucine codon, no amino acid change
  - GCG - GCG (1,062 counts) - Alanine codon, no amino acid change
  - GTT - GTT (953 counts) - Valine codon, no amino acid change
  - These matches high synonymous mutation rate.
- Hotspot codon switches:
  - CTG -TTG, GCG - GCA, and GGC – GGA, changes that can still code for the same amino acid (degenerate codons).
  - Missense hotspots (likely functional impact):
    - ATC - GTC (Isoleucine - Valine) - conservative change
    - CGC - CTC (Arginine - Leucine) - structural change
    - CAG - TAG (Glutamine - Stop) - stop-gained mutation (high impact)
- Stop codon changes:
  - Rows or columns with "-" or TAG, TGA, TAA, represent gained or lost stop codons indicating strong functional impact.

**c. Importance**

- This table shows the exact nucleotide-level pathway by which the amino acid changes from the first table occur.
- It helps pinpoint whether changes are due to:
  - Transitions (Ts) - more frequent and usually less disruptive.
  - Transversions (Tv) - rarer but often more disruptive.
- Matches the Ts/Tv ratio = 3.15, indicating transition bias.

|  | * | - | ? | A | C | D | E | F | G | H | I | K | L | M | N | P | Q | R | S | T | V | W | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| * | 27 | 3 |  |  |  |  |  |  |  |  |  | 1 |  |  |  |  | 1 |  |  |  |  |  |  |
| - | 14 |  | 281 | 81 | 8 | 18 | 34 | 19 | 41 | 17 | 12 | 21 | 21 | 10 | 11 | 53 | 19 | 40 | 43 | 21 | 28 | 7 | 3 |
| ? |  |  | 1 |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
| A | 1 | 106 |  | 4,971 |  | 75 | 68 |  | 125 |  |  |  | 1 | 1 |  | 67 |  | 1 | 145 | 524 | 414 | 1 |  |
| C | 1 | 6 |  |  | 349 |  |  | 1 | 15 |  |  |  |  |  | 1 |  |  | 63 | 22 |  |  | 5 | 15 |
| D |  | 45 |  | 65 |  | 2,189 | 280 | 1 | 158 | 20 |  |  |  |  | 104 | 1 |  |  | 1 |  | 14 |  | 4 |
| E | 1 | 45 |  | 101 |  | 294 | 3,014 |  | 76 |  |  | 68 |  |  |  | 87 |  |  |  |  | 13 | 2 |  |
| F |  | 12 |  |  | 6 |  |  | 459 |  |  | 12 |  | 99 |  |  |  | 1 |  | 19 |  | 13 | 1 | 19 |
| G |  | 59 |  | 129 | 12 | 131 | 73 |  | 5,202 |  |  | 2 |  |  |  | 1 |  | 62 | 183 |  | 18 |  |  |
| H |  | 24 |  |  |  | 19 |  |  |  | 942 |  |  | 13 |  | 14 | 8 | 79 | 199 | 1 |  |  |  | 62 |
| I |  | 18 |  |  |  |  | 16 | 1 |  |  | 954 |  | 52 | 46 | 8 |  |  | 2 | 6 | 68 | 303 |  |  |
| K |  | 27 |  |  |  |  | 58 |  |  |  |  | 675 |  | 2 | 36 | 1 | 43 | 97 |  | 21 |  |  |  |
| L | 1 | 77 |  | 3 | 1 |  |  | 109 |  | 14 | 47 |  | 5,036 | 54 |  | 75 | 47 | 30 | 13 | 1 | 138 | 6 | 1 |
| M |  | 13 |  |  |  |  |  |  |  |  | 30 | 3 | 49 | 2 | 3 |  |  | 3 |  | 40 | 64 |  |  |
| N |  | 21 |  |  |  | 95 |  |  | 2 | 27 | 7 | 33 |  |  | 902 |  |  | 1 | 153 | 26 |  |  | 9 |
| P |  | 50 |  | 89 | 1 |  |  |  | 1 | 6 |  |  | 75 |  |  | 2,214 | 31 | 24 | 113 | 29 |  |  |  |
| Q | 5 | 46 |  | 1 |  | 1 | 69 | 1 | 87 |  |  | 34 | 52 |  |  | 34 | 1,188 | 145 |  | 1 |  |  |  |
| R | 1 | 50 |  |  | 72 |  |  |  | 63 | 174 | 1 | 88 | 37 | 2 |  | 16 | 175 | 3,822 | 53 | 6 |  | 16 |  |
| S | 4 | 62 |  | 150 | 24 |  | 1 | 12 | 193 |  | 12 | 1 | 18 |  | 174 | 133 |  | 70 | 1,697 | 132 |  | 1 | 4 |
| T |  | 34 |  | 537 |  | 1 |  |  |  |  | 73 | 16 |  | 30 | 25 | 56 |  | 13 | 136 | 1,376 |  |  |  |
| V |  | 36 |  | 502 |  | 19 | 17 | 12 | 31 |  | 315 |  | 139 | 82 |  |  | 1 |  |  |  | 2,557 |  |  |
| W | 4 | 7 |  |  | 2 |  |  |  | 5 | 1 |  |  | 3 |  |  |  |  | 19 |  |  |  |  |  |
| Y | 3 | 6 |  |  | 19 | 6 |  | 23 |  | 42 |  |  | 1 |  | 5 |  |  |  | 3 |  |  |  | 953 |

**SnpEff codon change matrix**, showing how amino acid substitutions occurred between reference and variant sequences.

## a. What the Table Shows

- **Rows** = Original amino acid (reference)
- **Columns** = New amino acid after mutation (variant)
- **Numbers** = How many times that specific amino acid change occurred
- **Diagonal grey cells** = No change (synonymous mutation)
- **Green cells** = Non-synonymous changes (missense, nonsense, etc.)

## b. Key Observations

- High-frequency changes:
  - A to G and G to A in nucleotides led to many synonymous and conservative substitution.
  - Common protein-level substitutions:
    - L (Leucine) - L (5,036 cases) – synonymous
    - V (Valine) - V (2,557 cases) – synonymous
    - A (Alanine) - A (4,971 cases) – synonymous
  - These high counts matches 78.9% silent mutation rate.
- Biologically interesting changes:
  - Glycine - Aspartic Acid (G-D), Proline - Leucine (P-L), and Arginine - Cysteine (R-C), these can significantly alter protein folding or stability.
  - Stop codon gains (indicated by * in the column), which are several (e.g., Q-, E-, L-*) which are loss-of-function mutations.

### c. Functional Impact Context

- Most variants are MODIFIER or LOW impact, so likely in non-coding or synonymous regions.
- High-impact variants (0.1%):
  - Frameshift
  - Stop-gained / start-lost mutations
  - Likely to cause truncated or non-functional proteins

## Summary

### Disease Relevance

- The majority of variants detected are upstream/downstream of genes, which may influence promoter activity or transcription regulation, potentially altering virulence factor expression in Pseudomonas.
- The missense variants (~1.7%) are of particular interest for pathogenicity-related genes, as they can alter protein structure and function, possibly affecting antibiotic resistance or host-pathogen interactions.
- Nonsense variants (0.041%) could lead to truncated proteins, which in pathogenic bacteria can sometimes disable repressors or modify metabolic pathways relevant to infection survival.
- While no specific disease link can be confirmed without experimental validation, the functional categories affected suggest possible roles in adaptation, virulence, and resistance.

### Potential Functional Impacts

- Missense variants: Likely to cause amino acid substitutions; depending on location (active sites, binding domains), these could alter protein function, enzyme specificity, or stability.
- Nonsense variants: May result in loss-of-function proteins due to premature stop codons, in some cases beneficial to bacteria if the inactivated protein suppresses immune evasion or metabolic adaptation.
- Synonymous variants: Although traditionally considered neutral, they may influence codon bias and translation efficiency, especially relevant in bacteria where codon usage adapts to optimize growth under certain conditions.
- Indels: Frameshift events (0.1%) have a high probability of causing major disruptions in protein sequences, potentially producing non-functional proteins or novel variants with altered functions.
- Regulatory region changes: Given the large proportion of upstream/downstream variants, possible impacts include altered promoter strength, disruption or creation of transcription factor binding sites and modified mRNA secondary structure in untranslated regions.

## Isolate Comparison using NCBI Pathogen Detection



Following the identification of variants in the *Pseudomonas aeruginosa* genome, the annotated variants were cross-checked using the **NCBI Pathogen Detection Isolates Browser** to determine their relatedness to existing isolates in the database. The search was

performed using the reference genome **NC_002516** with the specific variant positions derived from sequencing analysis.

**1.** Matched Clusters - The analysis revealed that the submitted genome sequences clustered with multiple known SNP clusters.

- **Top Matched Cluster**: PDS000095640.50
  - **Matched isolates**: 597 (567 clinical, 20 environmental)
  - **Minimal SNP difference**: 0 (suggesting high genetic similarity)

Other clusters with fewer matches included:

- PDS000075445.200 (268 total isolates, 240 clinical, 15 environmental)
- PDS000076787.33 (116 total isolates, mostly environmental)
- Additional smaller clusters ranging from 20–84 isolates.

**2.** Matched Isolates - Individual isolates within the clusters were examined for:

- **Geographic location**
- **Isolation source** (e.g., clinical, environmental, blood)
- **Antimicrobial resistance (AMR) genes**
- **Assembly information and BioSample IDs**

For example:

- **Isolate PDT02882421.1** (USA, clinical, blood source) clustered within PDS000076522.21.
- AMR genes detected included **aadA11, aph (3')-Ib, aph (3')-IIb**, and point mutations in **gyrA_T83I** and **parC_S87W**, indicating potential fluoroquinolone resistance.

**3.** Interpretation - This comparison provides:

- **Epidemiological insight**: The sample is closely related to a large set of global isolates, indicating it belongs to a widely distributed lineage of *P. aeruginosa.*
- **AMR profile awareness**: Detection of specific resistance genes aids in predicting antibiotic susceptibility patterns.
- **Outbreak tracking potential**: Minimal SNP differences to other isolates could suggest recent transmission or shared origin.