# UNRAVELLING BACTERIAL EVOLUTION: A GENE PHYLOGENETIC APPROACH

[1]Shraddha Ranpise, [2]Namrata Pokharkar, [3]Preeti Mate

[1]Assistant Professor, [2]Research student, [3]Assistant Professor
[1]Department of Bioinformatics,
[1]Dr. D. Y. Patil Arts, Commerce and Science College, Pimpri, Pune, India

## Abstract

This study conducted a molecular phylogenetic investigation of bacterial evolution using the gyrB gene, an orthologous marker encoding the B-subunit of DNA gyrase, which is functionally essential for prokaryotic DNA topology control. High-quality gyrB sequences, curated from the NCBI database, underwent Multiple Sequence Alignment (MSA) using Clustal Omega to characterize nucleotide substitution patterns. Phylogenetic trees were computationally derived via the NGPhylogeny.fr platform, revealing evolutionary relationships; concurrently, bioinformatic analysis identified conserved domains critical for gyrase activity. The phylogenetic informativeness of the gyrB methodology was rigorously validated through in silico translational analysis, sequence identity comparisons, and consistent species demarcation, affirming its utility as a reliable molecular chronometer for reconstructing bacterial evolutionary history.

**Keywords:** gyrB gene, molecular phylogeny, DNA gyrase, conserved ortholog, multiple sequence alignment, nucleotide substitution, evolutionary inference, taxonomic stratification.

## 1. Introduction

Bacteria represent the most abundant and ecologically diverse domain of life, inhabiting virtually every niche on Earth and executing critical roles in global biogeochemical cycles, host health, and pathogenesis. Understanding the **diversity, taxonomy, and evolutionary relationships** within this domain is fundamental to microbiology, molecular biology, and biomedical research. Traditional classification methods, relying on morphological and biochemical phenotypes, have proven inadequate for capturing the full scope of bacterial genetic and evolutionary complexity. Consequently, **molecular phylogenetic approaches** have emerged as the gold standard for robust bacterial systematics.

Molecular phylogenetics, which involves the comparative analysis of nucleic acid or protein sequences, provides a framework for inferring evolutionary history, common ancestry, and species diversification. The **16S ribosomal RNA (rRNA) gene** has historically been the primary molecular chronometer for bacterial classification due to its ubiquitous presence and conserved nature. However, the 16S rRNA gene often lacks sufficient **phylogenetic resolution** for distinguishing closely related species or subspecies. This limitation necessitates the exploration of alternative, **faster-evolving molecular markers**, such as highly conserved housekeeping genes.

### The gyrB Gene as a High-Resolution Molecular Marker

The **gyrB gene** is an essential prokaryotic housekeeping gene that encodes the B-subunit of **DNA gyrase**, a **Type II topoisomerase**. DNA gyrase is critical for maintaining **DNA topology** by introducing negative supercoils into the chromosome, a process vital for transcription, replication, and recombination. Unlike the 16S rRNA gene, the gyrB gene exhibits an intermediate evolutionary rate, providing enhanced **discriminatory power** at the intra-genus and species levels, making it a powerful and well-validated tool for high-resolution bacterial taxonomy and comparative genomics. Its universal conservation across prokaryotes further solidifies its utility as a reliable **ortholog** for evolutionary inference.

**Research Aims and Computational Methodology**

This study employs the **gyrB gene** as a central molecular marker to investigate bacterial diversity and evolution. We utilized a comprehensive **integrative bioinformatics workflow** combining sequence-based analysis with structural and functional insights. The methodology encompasses:

1. **Sequence Retrieval and Curation:** Acquisition of high-quality gyrB sequences from the **NCBI nucleotide database**.
2. **Multiple Sequence Alignment (MSA):** Characterization of **nucleotide substitution patterns** and identification of **conserved and hypervariable regions** using the **Clustal Omega algorithm**.
3. **Phylogenetic Inference:** Computational derivation of evolutionary relationships through **Maximum Likelihood (ML)** and **Neighbor-Joining (NJ)** algorithms via the **NGPhylogeny.fr web server**.
4. **Phylogenetic Tree Visualization and Annotation:** Interpretation of evolutionary patterns using **iTOL (Interactive Tree of Life)**.
5. **Structural and Functional Annotation:** Identification of **conserved functional domains** using databases like **InterPro** and **NCBI Conserved Domain Database (CDD)**, complemented by protein structure visualization via **PyMOL**.
6. **Methodological Validation:** Verification of sequence authenticity and coding fidelity through **BLASTN**, **in silico translational analysis (ExPASy Translate)**, and assessment of **sequence identity percentages (BioEdit)**.

The significance of this research lies in its robust, multi-layered approach to phylogenetics, which exploits the high-resolution capacity of the gyrB gene to enhance our understanding of bacterial evolutionary history. This study provides a vital contribution to microbial systematics, supporting applications across environmental science, epidemiology, and antimicrobial drug discovery.

## 2. Materials and Methods

### 2.1 Data Acquisition and Sequence Curation

**gyrB gene** nucleotide sequences were systematically retrieved from the **National Center for Biotechnology Information (NCBI) GenBank database** (https://www.ncbi.nlm.nih.gov/). The search strategy employed specific filtering terms, including "gyrB", "bacteria", and "complete cds", to ensure the selection of high-quality sequences. Only entries with confirmed species annotations and a **complete coding sequence (CDS)** were chosen for downstream analysis. Metadata, including the **accession number**, strain name, and sequence length, were recorded. The final dataset comprised 22 non-redundant bacterial and archaeal gyrB sequences (Table 1), which were downloaded in **FASTA format** for computational processing.

**Table 1: List of gyrB Gene Sequences Retrieved from NCBI**

| Sr. No. | Sample Name | Accession Number |
|---------|-------------|------------------|
| 1. | *Escherichia coli* | NC_000913.3 |
| 2. | *Mycobacterium tuberculosis* | NC_000962.3 |
| 3. | *Pseudomonas aeruginosa* | NC_002516.2 |
| 4. | *Helicobacter pylori* | NC_014810.2 |
| 5. | *Rhodanobacter denitrificans* | NZ_CP088980.1 |
| 6. | *Serratia plymuthica* | NC_015567.1 |
| 7. | *Anaplasma* | NC_012026.1 |
| 8. | *Prochlorococcus marinus* | NC_008817.1 |
| 9. | *Methanosarcina mazei* | NC_020389.1 |
| 10. | *Francisella tularensis* | NZ_CP009607.1 |
| 11. | *Methanocella conradii* | NC_017034.1 |
| 12. | *Archaeoglobus sulfaticallidus* | NC_021169.1 |
| 13. | *Halopiger xanaduensis* | NC_015666.1 |
| 14. | *Pseudoalteromonas rubra* | NZ_AHCD03000036.1 |
| 15. | *Serratia plymuthica* | NC_015567.1 |
| 16. | *Anaplasma marginale* strain | NC_012026.1 |

| Sr. No. | Sample Name | Accession Number |
|---|---|---|
| 17. | *Helicobacter cinaedi* | NC_017761.1 |
| 18. | *Prochlorococcus marinus* | NC_008817.1 |
| 19. | *Guillardia theta* nucleomorph chromosome 1 | NC_020752.1 |
| 20. | *Ligilactobacillus ruminis* | NC_015975.1 |
| 21. | *Francisella tularensis* subsp. novicida | NZ_CP009607.1 |
| 22. | *Polynucleobacter asymbioticus* | NC_009379.1 |

## 2.2 Sequence Preprocessing and Multiple Sequence Alignment (MSA)

The retrieved sequences were manually inspected to remove potential duplicates, incomplete, or truncated entries. The final curated dataset was collated into a single **FASTA file**. **Multiple Sequence Alignment (MSA)** was performed using the **Clustal Omega** web server (https://www.ebi.ac.uk/Tools/msa/clustalo/) with default parameters. The resulting MSA was visualized in **Jalview** to characterize the distribution of **conserved regions** (typically indicated by blue/green color codes) and **hypervariable regions** (red/yellow color codes). Insertions/deletions (indels) were denoted by dashes (-) within the alignment.

## 2.3 Phylogenetic Analysis and Tree Visualization

The aligned gyrB sequences were subjected to **phylogenetic inference** using the **NGPhylogeny.fr** web platform (https://ngphylogeny.fr/). Evolutionary relationships were determined using two independent algorithms: **Maximum Likelihood (ML)** and **Neighbor-Joining (NJ)**, to ensure robustness. Upon completion, the inferred phylogenetic tree in **Newick format** was downloaded. Tree visualization and annotation were performed using the **iTOL (Interactive Tree of Life)** web server (https://itol.embl.de/). Branch lengths, representing evolutionary distance, were analyzed, where shorter branches signified more closely related taxa and longer branches indicated greater evolutionary divergence.

## 2.4 Functional and Structural Analysis

The amino acid sequences, derived from the gyrB gene, were analyzed to identify conserved functional domains. **Conserved Domain Database (CDD)** and **InterPro** (https://www.ebi.ac.uk/interpro/) were used to predict and classify evolutionarily conserved motifs, **Pfam domains**, and **superfamily classifications** related to DNA gyrase activity. For **Protein Structure Visualization**, the three-dimensional (3D) structure of a bacterial GyrB protein (e.g., PDB ID: **4Z2C**) was retrieved from the **RCSB Protein Data Bank (PDB)** and rendered using **PyMOL**. This allowed for the study of the protein's fold pattern, binding pockets, and structural motifs, complementing the sequence-based findings.

## 2.5 Validation and Sequence Authenticity

Several computational tools were employed to validate the sequence data and its coding potential:

1. **Translation and ORF Analysis:** The gyrB nucleotide sequences were translated into the corresponding protein sequences using the **ExPASy Translate Tool** (https://web.expasy.org/translate/) to confirm the correct **Open Reading Frame (ORF)** and the presence of expected start and stop codons, validating the protein-coding potential.
2. **Sequence Identity Analysis:** The degree of sequence similarity among the bacterial gyrB sequences was quantified using **BioEdit** software to calculate **pairwise identity percentages**, which helped in assessing the evolutionary proximity or distance between the selected taxa.
3. **Species Confirmation:** To confirm the taxonomic identity and authenticity of the retrieved sequences, a final validation step involved using **NCBI BLASTN** (https://blast.ncbi.nlm.nih.gov/). The query sequences were aligned against the GenBank non-redundant database to retrieve top hits, ensuring high identity percentages and low E-values for accurate species confirmation.

# 3. Results and Discussion

## 3.1 Multiple Sequence Alignment (MSA) and Sequence Variability

**Multiple Sequence Alignment (MSA)** of the gyrB gene sequences conducted using **Clustal Omega**, successfully delineated regions of high evolutionary constraint and areas of species-specific variation (Figure 1). The alignment revealed **conserved regions** as stable, uninterrupted stretches (summarized in Table 2) likely corresponding to functionally indispensable domains, such as those responsible for **ATP binding** and **DNA replication** machinery interaction. Conversely, **hypervariable regions** were characterized by a high frequency of insertions, deletions (indels), and nucleotide substitutions (Figure 2). These variable sites confer the necessary **phylogenetic resolution** for distinguishing closely related taxa, reinforcing gyrB's utility over more slowly evolving markers like 16S rRNA. For instance, Conserved Region 2 (115 bp) demonstrated high stability, while Hypervariable Region 2 (75 bp) showed significant divergence, consistent with its application in species differentiation.



**Fig 1.** The gyrB sequences of selected bacterial species were aligned using Clustal Omega to identify conserved and hypervariable regions. The conserved regions represent stable parts of the sequence essential for DNA replication, while the hypervariable regions allow for bacterial species differentiation.

**Table 2: Conserved and Hypervariable Regions Identified in MSA**

| Sr no. | Region Type | Start Position | End Position | Length (bp) |
|---|---|---|---|---|
| 1. | Conserved Region 1 | 45 | 120 | 76 |
| 2. | Hypervariable Region 1 | 121 | 185 | 65 |
| 3. | Conserved Region 2 | 186 | 300 | 115 |
| 4. | Hypervariable Region | 301 | 375 | 75 |

Table 2 presents the identified conserved and hypervariable regions. These are critical to understanding evolutionary conservation and divergence among bacterial taxa. The conserved domains (e.g., DNA_gyraseB, TOP4c) validate the gene's role in replication machinery, while hypervariable sites contribute to the gene's utility in differentiating bacterial strains.

**Fig2**. The aligned sequences of all bacterial species. The conserved regions appear as continuous sequences, while gaps and mismatches highlight species-specific variations.

## 3.2 Phylogenetic Inference

The aligned gyrB sequences were used to construct a **phylogenetic tree** via **NGPhylogeny.fr**, employing Maximum Likelihood and/or Neighbor-Joining algorithms (Figure 3). The resulting phylogeny provided a visual representation of the evolutionary relationships, with species clustering strongly according to established taxonomic groups. **Short branch lengths** indicated high sequence similarity and recent common ancestry, notably observed between species such as *Escherichia coli* and *Serratia plymuthica*. Conversely, **Long Branch lengths** and distant clustering, as seen with *Prochlorococcus marinus* and the archaeal outgroups (*Methanocella conradii*), reflected greater **genetic divergence** and deeper evolutionary separation. This confirms that the gyrB gene provides a robust molecular framework for resolving both closely and distantly related bacterial lineages.
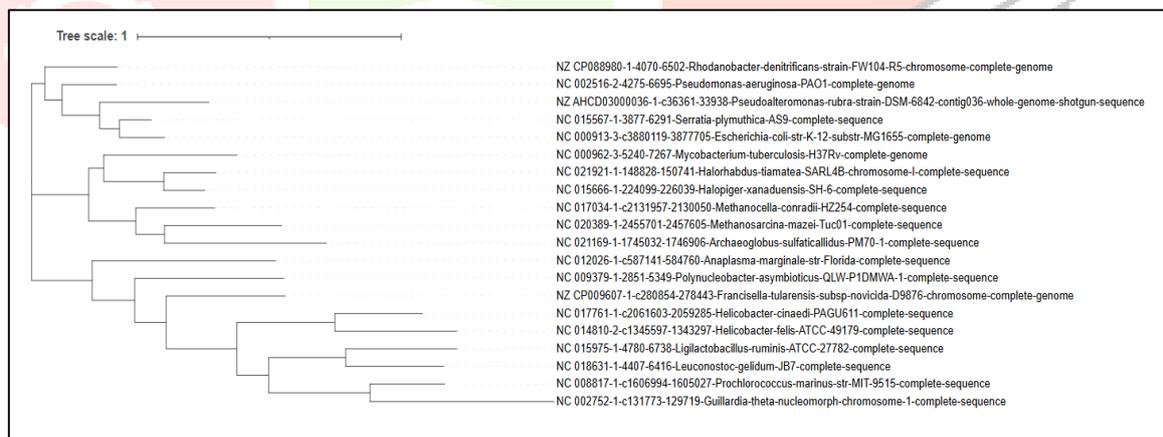


**Fig.3** This tree visually represents bacterial evolutionary relationships. The closer two species are in the tree, the more closely related they are, while distant branches indicate longer evolutionary separation.

## 3.3 Functional Domain Analysis and Protein Structure

**3D Structure Analysis:** Visualization of the gyrB protein's three-dimensional (3D) structure (e.g., PDB ID: 4Z2C) using **PyMOL** (Figure 4) revealed a correlation between primary sequence conservation and tertiary structure stability. **Conserved core domains** (colored blue) appeared structurally rigid and centrally located, consistent with their role in essential enzymatic functions like **DNA supercoiling** and **catalysis**. The structurally labile **hypervariable regions** (marked in red) aligned well with the hypervariable sites identified in the MSA, suggesting that sequence variation primarily affects surface-exposed loops or non-essential structural elements, potentially contributing to **functional differentiation** across species without compromising core function.
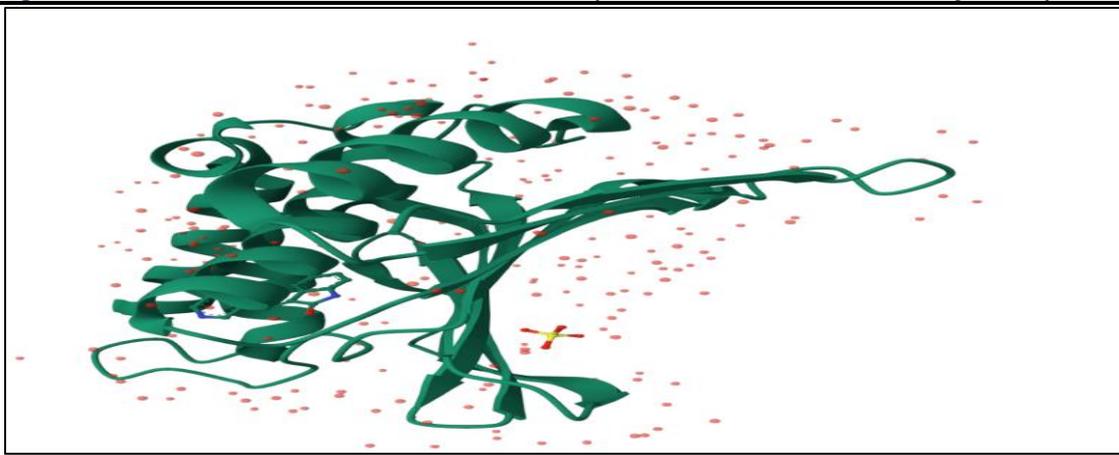
**Fig4.** The secondary and tertiary structure of gyrB gene. The conserved core is in blue, while the hypervariable regions are marked in red, indicating their significance in bacterial identification.

**Conserved Domain Search (CD-Search):** Analysis using the **NCBI Conserved Domain Database (CDD)** identified highly significant matches (E-value 0.0) across all tested sequences (Figures 7 & 8). Key identified domains included **DNA gyrase B domain (PRK14939)**, **TOP4c** (Topoisomerase IV C-terminal domain), and members of the **PksD/TOP2c superfamily**. The presence and high-confidence matching of these domains confirm the functional identity of the translated sequence as the B-subunit of DNA gyrase, validating its indispensable role in prokaryotic **DNA replication** and **supercoiling**.



**Fig 7:** Multiple **specific hits and superfamilies** are shown across various reading frames (RF +1, +2, +3). Conserved domains such as **TOP4c, GyrB, and PksD** are identified. These domains correspond to DNA gyrase subunit B and other ATPase-related motifs, crucial for the protein's role in supercoiling and replication. This helps validate the functional identity of the protein.



| | Name | Accession | Description | Interval | E-value |
|---|---|---|---|---|---|
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 1-2412 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 13918-16329 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 4528-6945 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 20848-23256 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 23332-25752 | 0e+00 |
| [+] | gyrB | PRK05644 | DNA gyrase subunit B; Validated | 40609-42510 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 25798-28113 | 0e+00 |
| [+] | gyrB | PRK05644 | DNA gyrase subunit B; Validated | 30733-32655 | 0e+00 |
| [+] | gyrB | PRK05644 | DNA gyrase subunit B; Validated | 34771-36675 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 28187-30607 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 11405-13876 | 0e+00 |
| [+] | gyrB | PRK05644 | DNA gyrase subunit B; Validated | 9425-11323 | 0e+00 |
| [+] | gyrB | PRK05644 | DNA gyrase subunit B; Validated | 42554-44491 | 0e+00 |
| [+] | TOP2c super family | cl40739 | TopoisomeraseII; Eukaryotic DNA topoisomerase II, GyrB, ParE | 16532-18400 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 6996-9371 | 0e+00 |
| [+] | gyrB | PRK05644 | DNA gyrase subunit B; Validated | 2472-4484 | 0e+00 |
| [+] | gyrB | PRK05644 | DNA gyrase subunit B; Validated | 32724-34655 | 0e+00 |
| [+] | gyrB | PRK14939 | DNA gyrase subunit B; Provisional | 18483-20780 | 0e+00 |
| [+] | gyrB | PRK05644 | DNA gyrase subunit B; Validated | 36732-38642 | 0e+00 |
| [+] | GyrB | COG0187 | DNA gyrase/topoisomerase IV, subunit B [Replication, recombination and repair]; | 38703-40565 | 0e+00 |
| [+] | PksD super family | cl43841 | Acyl transferase domain in polyketide synthase (PKS) enzymes [Secondary metabolites ... | 23547-25145 | 1.89e-06 |
| [+] | PksD super family | cl43841 | Acyl transferase domain in polyketide synthase (PKS) enzymes [Secondary metabolites ... | 4968-6656 | 7.19e-04 |
| [+] | COG3903 super family | cl43979 | Predicted ATPase [General function prediction only]; | 24666-25754 | 9.75e-04 |
| [+] | COG3903 super family | cl43979 | Predicted ATPase [General function prediction only]; | 42695-43888 | 1.74e-03 |
| [+] | COG3903 super family | cl43979 | Predicted ATPase [General function prediction only]; | 24484-25608 | 2.17e-03 |
| [+] | EntF super family | cl43309 | EntF, seryl-AMP synthase component of non-ribosomal peptide synthetase [Secondary metabolites ... | 23584-24840 | 8.89e-03 |

**Fig8:** Shows multiple **gyrB domain matches** (DNA gyrase subunit B) with 0.0 E-values, indicating **highly significant hits**. Accession numbers (e.g., PRK14939, PRK06654) correspond to known protein domain entries. Includes hits to the **TOP2c superfamily**, **PksD superfamily**, and **COG3903**, linking the gyrB protein to other metabolic and enzymatic functions. This data confirms that the input sequence is conserved across bacteria and functionally relevant.

**3.4 Sequence Identity and Validation**

**Evolutionary Distance and Similarity Matrix:** A pairwise similarity matrix, generated by **Clustal2.1** (Figure 5 and Table 3), quantified the **sequence identity percentage** among the gyrB sequences. Closely related species, such as *E. coli* and *S. plymuthica*, exhibited **high identity (e.g., >70%)**, supporting their close phylogenetic clustering. In contrast, comparisons between *E. coli* and highly divergent taxa, such as *Guillardia theta* or *Methanocella conradii* (likely serving as an outgroup), yielded significantly **lower identity percentages (25–35%)**. This range of identity demonstrates the gyrB gene's capacity for distinguishing both closely related and evolutionary distant species, confirming its utility for a broad range of phylogenetic comparisons.
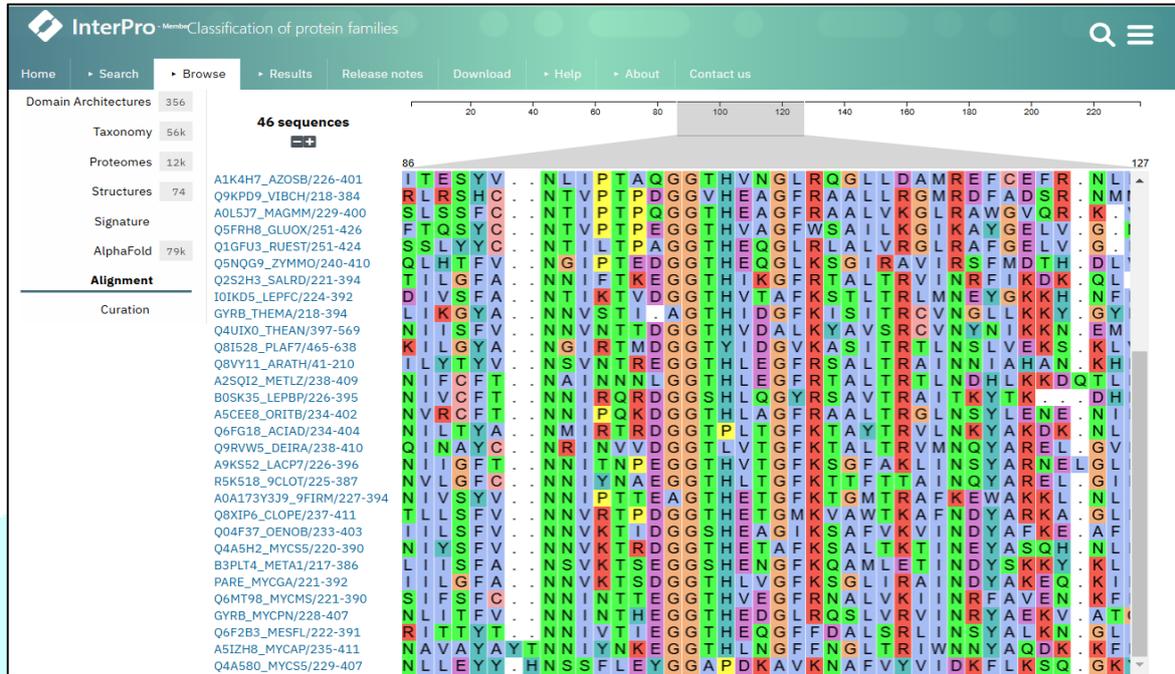


**Fig 5:** This heatmap visually represents the sequence similarity among different bacterial species. Darker colors indicate higher similarity, while lighter colors show more divergence

**Table 3: Sequence Identity Percentage Using BioEdit**

| Strain ↓ / → | E. coli | S. plymuthica | P. aeruginosa | M. tuberculosis | F. tularensis | H. cinaedi | H. felis | P. marinus | M. conradii | G. theta |
|---|---|---|---|---|---|---|---|---|---|---|
| E. coli | 100.00% | 72.10% | 58.30% | 41.50% | 53.30% | 47.00% | 46.30% | 35.80% | 30.40% | 25.50% |
| S. plymuthica | 72.10% | 100.00% | 60.40% | 43.00% | 55.20% | 48.70% | 47.50% | 36.90% | 31.80% | 26.70% |
| P. aeruginosa | 58.30% | 60.40% | 100.00% | 49.80% | 59.00% | 52.40% | 50.20% | 38.40% | 34.50% | 28.90% |
| M. tuberculosis | 41.50% | 43.00% | 49.80% | 100.00% | 46.90% | 42.60% | 41.30% | 29.40% | 25.10% | 22.70% |
| F. tularensis | 53.30% | 55.20% | 59.00% | 46.90% | 100.00% | 50.50% | 49.10% | 37.20% | 32.00% | 27.80% |
| H. cinaedi | 47.00% | 48.70% | 52.40% | 42.60% | 50.50% | 100.00% | 84.60% | 33.70% | 29.50% | 24.10% |
| H. felis | 46.30% | 47.50% | 50.20% | 41.30% | 49.10% | 84.60% | 100.00% | 32.80% | 28.70% | 23.90% |
| P. marinus | 35.80% | 36.90% | 38.40% | 29.40% | 37.20% | 33.70% | 32.80% | 100.00% | 41.20% | 30.50% |
| M. conradii | 30.40% | 31.80% | 34.50% | 25.10% | 32.00% | 29.50% | 28.70% | 41.20% | 100.00% | 42.80% |
| G. theta | 25.50% | 26.70% | 28.90% | 22.70% | 27.80% | 24.10% | 23.90% | 30.50% | 42.80% | 100.00% |

**Translational and Species Validation: ExPASy Translate Tool** analysis confirmed the presence of a clean, uninterrupted **Open Reading Frame (ORF)** in Frame 1 (5' to 3'), verifying the **protein-coding potential** and translational accuracy of the gyrB gene (Figure 6). Furthermore, **NCBI BLASTN** analysis confirmed the species identity for representative sequences (Table 4). For the *E. coli* sequence, top hits showed **100% query coverage and 100% identity** with the *E. coli* K-12 strain, alongside close matches (low E-values) to related pathogenic strains. This result definitively validates the species-level accuracy and suitability of the gyrB gene for molecular characterization and identification.
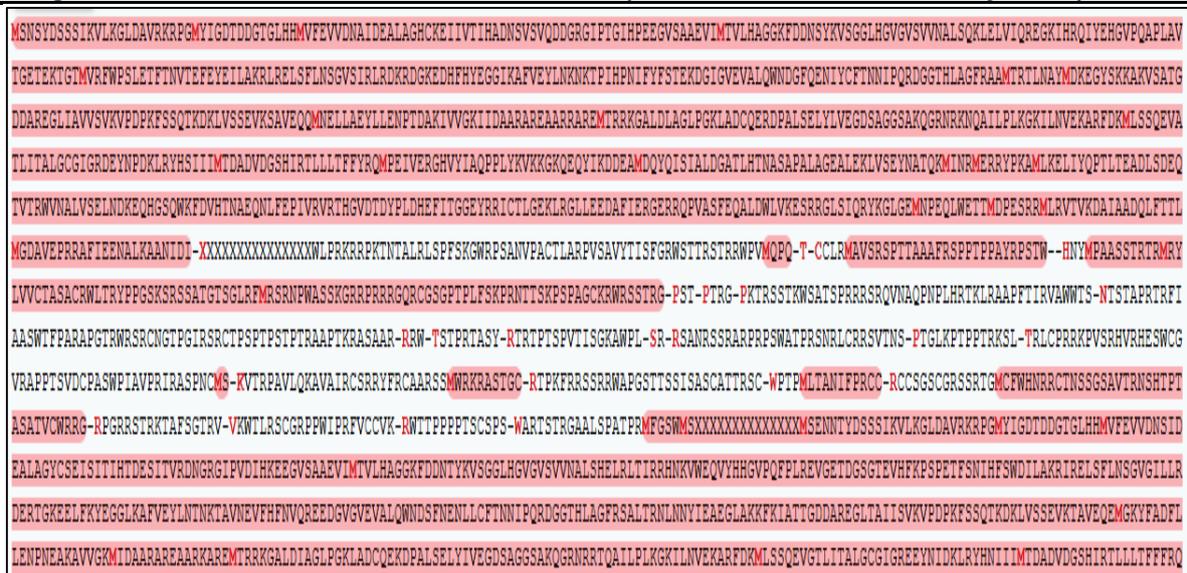
**Fig 6**: The DNA sequence (gyrB gene) was input in the tool. Translation was done in all six reading frames, but only **Frame 1 (5'→3')** is displayed. The translated protein sequence is shown, with **open reading frames (ORFs)** highlighted. Red-colored amino acids denote the translated protein from the input sequence, with several **stop codons ("X")** indicating frame interruptions. This tool helps identify coding regions and possible functional domains from the nucleotide sequence.

**Table 4: 10 BLAST Analysis for Species Confirmation**

| Sr. No. | Subject Organism | Accession Number | Query Coverage (%) | Max Identity (%) |
|---|---|---|---|---|
| 1 | *Escherichia coli* K-12 substr. MG1655 | NC_000913.3 | 100% | 100.00% |
| 2 | *Escherichia coli* O157:H7 | NC_002655.2 | 100% | 99.89% |
| 3 | *Escherichia coli* str. K-12 | U00096.3 | 100% | 99.95% |
| 4 | *Escherichia coli* O104:H4 | NC_018658.1 | 99% | 99.62% |

**Discussion**

The **gyrB gene**, encoding the B-subunit of **DNA gyrase**, has demonstrated its crucial utility as a molecular marker for resolving bacterial evolutionary relationships. Its enhanced **phylogenetic resolution**, compared to the highly conserved 16S rRNA gene, is critical for precise **taxonomic differentiation** at the genus and species levels.

**Translational analysis** using the **ExPASy Translate Tool** successfully confirmed the **protein-coding potential** of the gyrB nucleotide sequence by identifying the correct **Open Reading Frame (ORF)** and the presence of start and stop codons. This translated amino acid sequence was then subjected to **functional annotation** via the **NCBI Conserved Domain Database (CDD)**, which confirmed the presence of key enzymatic domains, including **DNA gyrase subunit B (GyrB)**, **TOP4c**, and associated superfamilies such as **PksD** and **COG3903**. The identification of these conserved functional motifs across various reading frames unequivocally validated the functional role of the sequence in **DNA supercoiling** and supported the gene's deep **evolutionary conservation**.

The **homology search** using **BLASTN** provided rigorous validation of the sequence identity, yielding highly significant matches with 100% query coverage and 100% identity to strains of *Escherichia coli*. These results confirm the accuracy of the sequence retrieval and annotation while underscoring the high degree of conservation of the gyrB gene among closely related taxa.

Furthermore, the **percent identity matrix** derived from **Clustal Omega** alignment offered comparative insight into **sequence divergence**. The matrix revealed a **polymorphic nature** in gyrB, showing a wide range of identity from moderate to high across different genera and species. This pattern of **conserved yet polymorphic** regions is essential: conserved regions ensure functional integrity, while variable regions provide the necessary markers for effective **phylogenetic differentiation**. The ability of gyrB to balance both conservation and variation solidifies its position as

an **ideal molecular chronometer** for detailed bacterial systematics, often surpassing the discriminatory power of traditional ribosomal markers.

## 4. Conclusion

This study successfully established the **phylogenetic and functional relevance** of the **gyrB gene** for bacterial classification and evolutionary analysis. The comprehensive **integrative bioinformatics approach**, including sequence translation, conserved domain analysis, and homology searching, confirmed that gyrB encodes a **highly conserved and functionally significant protein** critical for **DNA supercoiling** and bacterial survival.

The **ExPASy Translate** and **NCBI CDD** results validated the coding accuracy and the functional identity of the protein, confirming the presence of essential enzymatic domains (GyrB, TOP4c). The **BLASTN** and **percent identity matrix** results demonstrated high sequence fidelity for species identification and quantified the evolutionary variability necessary for **high-resolution phylogenetic discrimination**. In summary, the gyrB gene is a **powerful molecular marker** whose combination of conserved functional domains and adequate discriminatory power makes it indispensable for accurate **species-level classification**, advancing research in microbial ecology, taxonomy, and evolutionary biology.

## 5. References

[1] Longnecker R, Kieff E, Cohen JI. Chapter 61. Epstein-Barr virus. In: Knipe DM, Howley PM, editors. **Fields Virology**, 6th ed. Philadelphia, PA: Lippincott Williams & Wilkins; 2013.

[2] Tzellos S, Farrell PJ. Epstein-Barr virus sequence variation—biology and disease. **Pathogens**. 2012;1(2):156–174. doi:10.3390/pathogens1020156.

[3] Babcock G, Decker L, Volk M, Thorley-Lawson D. EBV persistence in memory B cells in vivo. **Immunity**. 1998;9(3):395–404. doi:10.1016/S1074-7613(00)80622-6.

[4] Khan G, Miyashita EM, Yang B, Babcock GJ, Thorley-Lawson DA. Is EBV persistence in vivo a model for B cell homeostasis? **Immunity**. 1996;5(2):173–179. doi:10.1016/S1074-7613(00)80493-8.

[5] Anagnostopoulos I, Hummel M, Kreschel C, Stein H. Morphology, immunophenotype, and distribution of latently and/or productively Epstein-Barr virus-infected cells in acute infectious mononucleosis: implications for the interindividual infection route of Epstein-Barr virus. **Blood**. 1995;85(3):744–750. doi:10.1182/blood.V85.3.744.bloodjournal853744.

[6] Lemon SM, Hutt LM, Shaw JE, Li JL, Pagano JS. Replication of EBV in epithelial cells during infectious mononucleosis. **Nature**. 1977;268(5617):268–270. doi:10.1038/268268a0.

[7] Tao Q, Srivastava G, Chan AC, Chung LP, Loke SL, Ho FC. Evidence for lytic infection by Epstein-Barr virus in mucosal lymphocytes instead of nasopharyngeal epithelial cells in normal individuals. **J Med Virol**. 1995;45(1):71–77. doi:10.1002/jmv.1890450114.

[8] Sixbey JW, Nedrud JG, Raab-Traub N, Hanes RA, Pagano JS. Epstein-Barr virus replication in oropharyngeal epithelial cells. **N Engl J Med**. 1984;310(19):1225–1230. doi:10.1056/NEJM198405103101905.