# State-of-the-Art Methodologies for Curating and Validating Extremophile Protein Datasets

## 1. Approach 1: Genomic Quality Filtering and Metagenomic Quality Control

### 1.1. Summary of the Method

#### 1.1.1. Application of CheckM and MIMAG Standards

The foundation of a high-quality protein dataset, especially one derived from metagenomic sources, is the rigorous quality control of the underlying genomes. The scientific community has established standardized protocols to ensure that only high-fidelity genomes are used for downstream analyses. A cornerstone of this process is the use of **CheckM**, a computational tool that assesses the quality of Metagenome-Assembled Genomes (MAGs) by estimating their completeness and contamination levels . This is achieved by analyzing the presence and copy number of lineage-specific sets of single-copy marker genes. For archaeal genomes, CheckM utilizes a set of **53 marker genes (arc53)** to provide these estimates . The **Minimum Information about a Metagenome-Assembled Genome (MIMAG)** standards further formalize these quality metrics, defining specific thresholds for genome completeness, contamination, and other features like the presence of rRNA genes and tRNAs to classify MAGs into high, medium, and low-quality tiers . For instance, a **high-quality MAG** according to MIMAG standards must be **>90% complete**, **<5% contaminated**, and contain the 23S, 16S, and 5S rRNA genes along with at least 18 tRNAs . These standards are widely adopted by major taxonomic databases like the **Genome Taxonomy Database (GTDB)** to determine the eligibility of MAGs for inclusion .

The **Genome Taxonomy Database (GTDB)** itself implements a stringent set of quality control criteria for the genomes it incorporates, which serves as a de facto standard for many researchers constructing datasets . To be included in the GTDB reference trees and database, a genome must meet several criteria: a CheckM completeness estimate greater than 50%, a CheckM contamination estimate less than 10%, a quality score (defined as completeness – 5*contamination) greater than 50%, contain at least 40% of the relevant marker genes (bac120 for Bacteria, arc53 for Archaea), have fewer than 1000 contigs, an N50 greater than 5kb, and contain fewer than 100,000 ambiguous bases . These comprehensive filters ensure that the genomes used for taxonomic assignment and phylogenetic analysis are of sufficient quality to produce reliable results. By applying these filters, researchers can significantly reduce the risk

of including fragmented, chimeric, or contaminated genomes in their protein datasets, which is a critical first step in ensuring the purity of the final sequence collection. The GTDB's commitment to quality is evident in its continuous updates and methodological improvements, such as the refinement of the archaeal marker gene set from 122 to 53 genes based on recent evaluations to minimize horizontal gene transfer and optimize the recovery of monophyletic lineages .

## 1.1.2. Removal of Contaminating Contigs using MAGpurify

Beyond the initial quality assessment of a MAG as a whole, it is often necessary to identify and remove specific contaminating contigs within an otherwise acceptable genome. This is particularly important for extremophile datasets, where a few mesophilic contaminants can introduce significant noise. **MAGpurify** is a widely used tool designed for this purpose. It employs a multi–faceted approach to detect and remove contaminating contigs from MAGs, utilizing a combination of phylogenetic, compositional, and coverage–based methods . One of its key modules, **"phylo–markers,"** identifies contigs that are taxonomically discordant with the rest of the genome by searching for the presence of single–copy marker genes from different phylogenetic lineages. This is particularly useful for detecting contamination from closely related species, which can be difficult to identify using other methods. Another module, **"clade–markers,"** uses a database of clade–specific marker genes to identify contigs that originate from a different taxonomic group than the target organism. This module is highly effective at detecting contamination from distantly related organisms, such as bacteria in an archaeal genome.

In addition to its phylogenetic modules, MAGpurify also includes modules that analyze the compositional properties of contigs. The **"tetra–freq"** module uses principal components analysis (PCA) to identify contigs with outlier tetranucleotide frequencies, which can be indicative of contamination. This method is based on the observation that organisms from different taxonomic groups often have distinct genomic signatures, which can be detected by analyzing their nucleotide composition. Similarly, the **"gc–content"** module uses PCA to identify contigs with outlier GC content, another compositional property that can be used to detect contamination. By combining these different approaches, MAGpurify provides a comprehensive and robust method for identifying and removing contaminating contigs from MAGs. This is particularly important for extremophile research, where the genomic data is often derived from challenging environments with a high potential for contamination. The use of MAGpurify

has been shown to significantly improve the quality of MAGs, leading to more accurate downstream analyses .

## 1.1.3. Quality Trimming and Removal of Host and PhiX Contamination

Quality trimming and the removal of host and PhiX contamination are essential steps in the pre-processing of metagenomic sequencing data. These steps are critical for ensuring the accuracy and reliability of downstream analyses, such as genome assembly and annotation. Quality trimming involves the removal of low-quality bases and adapter sequences from raw sequencing reads, which can be a significant source of error in metagenomic studies. Tools like **Trimmomatic** are commonly used for this purpose, employing sliding window algorithms to identify and remove segments of reads that fall below a specified quality threshold . This process helps to improve the overall quality of the sequencing data, leading to more accurate and reliable results. The removal of host contamination is also a critical step, particularly in studies of host-associated microbiomes. This is typically achieved by aligning the sequencing reads to a reference host genome and discarding any reads that map to the host. This process helps to ensure that the resulting dataset is enriched for microbial sequences, which is essential for accurate taxonomic and functional profiling.

The removal of **PhiX contamination** is another important step in the pre-processing of metagenomic data. PhiX is a bacteriophage that is often added to sequencing libraries as a control, but its presence can introduce bias into the data if not properly removed. This is typically achieved by aligning the sequencing reads to the PhiX genome and discarding any reads that map to it. The removal of PhiX contamination is particularly important in studies of low-biomass samples, where the presence of even a small amount of PhiX can have a significant impact on the results. In addition to these steps, other quality control measures, such as error correction and duplicate removal, are also commonly used to improve the quality of metagenomic data. Error correction tools, such as **BayesHammer,** can be used to identify and correct sequencing errors, while duplicate removal tools, such as **CD-HIT**, can be used to remove PCR duplicates, which can introduce bias into the data. The integration of these quality control measures into bioinformatics pipelines is essential for ensuring the reliability of metagenomic data.

## 1.2. Key Citations

### 1.2.1. Use of CheckM for Assessing MAG Quality in Haloarchaea

The use of CheckM for assessing the quality of Metagenome-Assembled Genomes (MAGs) has become a standard practice in the field of metagenomics, and its application in the study of haloarchaea is no exception. A key study that highlights the importance of this approach is the work by **Boulton et al. (2025)**, which focused on the recovery and analysis of MAGs from a hypersaline environment . In this study, the authors used CheckM to estimate the completeness and contamination of their prokaryotic MAGs, applying the MIMAG standards to ensure that only high-quality genomes were included in their downstream analyses. Specifically, they retained only those MAGs that met the criteria of **≥50% completeness and <10% contamination**, a threshold that is commonly used to define medium-quality MAGs . This rigorous quality control step was essential for ensuring the reliability of their results, as it helped to minimize the risk of incorporating erroneous data into their analyses.

Another example is the study of the microbiome of **Lake Hillier**, a hypersaline lake in Australia, which used CheckM to evaluate the quality of the MAGs reconstructed from metagenomic data . In this study, the researchers used an ensemble binning approach to generate initial genome bins, which were then filtered and refined using dRep, a tool that wraps CheckM. The study defined high-quality MAGs as those with **>90% completeness and <5% contamination**, and medium-quality MAGs as those with **50–90% completeness and <5% contamination**. This stringent quality control ensured that the subsequent analysis of the haloarchaeal genomes was based on reliable data, which is essential for the accurate identification of the genetic and metabolic adaptations that allow these organisms to thrive in high-salt environments. The use of CheckM in these studies highlights its importance as a standard tool for assessing the quality of MAGs and ensuring the integrity of metagenomic datasets.

### 1.2.2. MAGpurify for Contaminant Removal in Archaeal Genomes

MAGpurify has emerged as a critical tool for the removal of contaminating contigs from Metagenome-Assembled Genomes (MAGs), and its application in the study of archaeal genomes has been particularly impactful. A notable study that demonstrates the utility of this tool is the work by **Ghaly et al. (2022)**, which focused on the discovery of integrons in Archaea . In this study, the authors downloaded a large dataset of archaeal genomes from the NCBI Assembly Database, the majority of which were MAGs. To ensure the quality of their dataset, they applied a stringent filtering pipeline that included the use of MAGpurify to remove contaminating contigs. The tool was used with several of its modules, including **"phylo-markers," "clade-markers," "tetra-freq," and "gc-content,"** to identify and remove contigs that were likely to be contaminants .

This comprehensive approach to contamination removal was essential for ensuring the reliability of their results, as it helped to minimize the risk of incorporating erroneous data into their analyses.

Another example of the use of MAGpurify in the study of archaeal genomes is the investigation of the microbiome of **Lake Hillier**, which used MAGpurify as part of a comprehensive pipeline for the generation and refinement of MAGs . The study used an ensemble binning approach to generate initial genome bins, which were then filtered and refined using dRep, a tool that incorporates MAGpurify. The use of MAGpurify in this study helped to ensure that the final set of MAGs was free from contamination, which is essential for the accurate analysis of the archaeal genomes and the construction of reliable protein datasets. The study highlights the importance of using a combination of tools and methods to achieve the highest possible level of quality control in metagenomic studies.

### 1.2.3. General Metagenomic QC Pipelines

General metagenomic quality control (QC) pipelines are essential for ensuring the accuracy and reliability of downstream analyses. These pipelines typically involve a series of steps designed to remove artifacts, errors, and contaminants from raw sequencing reads. A key component of these pipelines is quality trimming, which involves the removal of low–quality bases and adapter sequences. Tools like **Trimmomatic** are commonly used for this purpose, employing sliding window algorithms to identify and remove segments of reads that fall below a specified quality threshold . This process helps to improve the overall quality of the sequencing data, leading to more accurate and reliable results. Another important step in metagenomic QC pipelines is the removal of host contamination. This is typically achieved by aligning the sequencing reads to a reference host genome and discarding any reads that map to the host. This process helps to ensure that the resulting dataset is enriched for microbial sequences, which is essential for accurate taxonomic and functional profiling.

In addition to quality trimming and host removal, metagenomic QC pipelines also typically include steps for error correction and duplicate removal. Error correction tools, such as **BayesHammer**, can be used to identify and correct sequencing errors, while duplicate removal tools, such as **CD–HIT**, can be used to remove PCR duplicates, which can introduce bias into the data. The removal of PhiX contamination is another important step in metagenomic QC pipelines. PhiX is a bacteriophage that is often added to sequencing libraries as a control, but its presence can introduce bias into the data if not properly removed. This is typically achieved by aligning the sequencing

reads to the PhiX genome and discarding any reads that map to it. The integration of these quality control measures into bioinformatics pipelines is essential for ensuring the reliability of metagenomic data. By carefully pre-processing their data, researchers can minimize the risk of introducing errors and biases into their analyses, which can lead to more accurate and reliable results.

## 1.3. Relevance

### 1.3.1. Ensures High-Quality Source Genomes

The application of rigorous genomic quality filtering and metagenomic quality control measures is of paramount importance for ensuring the generation of high-quality source genomes, which serve as the foundation for all downstream analyses, including the construction of protein datasets for training language models. By employing tools like CheckM and adhering to standards like MIMAG, researchers can systematically assess and filter out low-quality or contaminated genomes from their datasets. This process is critical because the presence of incomplete or contaminated genomes can lead to the inclusion of erroneous or non-extremophile protein sequences, thereby introducing noise and bias into the training data. For instance, a study on the microbiome of **Lake Hillier** demonstrated the use of CheckM to define stringent quality thresholds for Metagenome-Assembled Genomes (MAGs), ensuring that only high-quality genomes were used for subsequent analysis . This meticulous approach guarantees that the protein sequences derived from these genomes are representative of the target extremophile organisms and are free from artifacts that could compromise the performance of the language model.

Furthermore, the use of tools like MAGpurify for the removal of contaminating contigs from MAGs provides an additional layer of quality control, further enhancing the purity of the source genomes. This is particularly relevant for extremophile datasets, which are often derived from complex environmental samples where the risk of contamination is high. By systematically identifying and removing contaminating sequences, researchers can minimize the risk of including non-extremophile proteins in their datasets, thereby improving the accuracy and reliability of the language model. The combination of these quality control measures ensures that the resulting protein dataset is of the highest possible quality, which is essential for the successful training of a protein language model that can accurately capture the unique characteristics of extremophile proteins. The investment in rigorous genomic quality control at the initial stages of data curation is therefore a critical step that pays significant dividends in the long run, leading to more robust and reliable models.

## 1.3.2. Prevents Phylogenetic Contamination from Metagenomic Data

The prevention of phylogenetic contamination from metagenomic data is a critical aspect of curating high-quality extremophile protein datasets. Metagenomic samples are often a complex mixture of DNA from multiple organisms, including the target extremophiles, as well as other microorganisms that may be present in the same environment. This can lead to the inclusion of non-extremophile sequences in the dataset, which can confound the analysis and lead to inaccurate conclusions. To address this issue, researchers employ a variety of methods to identify and remove contaminating sequences. One common approach is to use taxonomic classification tools, such as **Kraken2** or **MetaPhlAn**, to identify the taxonomic origin of each sequence in the metagenomic dataset. Sequences that are classified as belonging to a different taxonomic group than the target extremophiles can then be removed from the dataset. This approach can be effective for removing contaminants that are phylogenetically distant from the target organisms, but it may be less effective for removing contaminants that are closely related.

Another approach to preventing phylogenetic contamination is to use tools that are specifically designed to identify and remove contaminating contigs from Metagenome-Assembled Genomes (MAGs). One such tool is **MAGpurify**, which uses a combination of sequence composition, taxonomic classification, and gene synteny to identify and remove contaminating contigs. This approach can be more effective than taxonomic classification alone, as it can identify contaminants that are closely related to the target organisms. The use of these tools is particularly important for extremophile datasets, as these are often derived from complex environmental samples where the risk of contamination is high. By systematically applying these quality control measures, researchers can minimize the risk of phylogenetic contamination and ensure that their protein datasets are as pure as possible. This is essential for the accurate analysis of extremophile genomes and the construction of reliable protein datasets for training language models.

# 2. Approach 2: Biophysical and Statistical Validation

## 2.1. Summary of the Method

Biophysical and statistical validation represents a cornerstone in the curation of extremophile protein datasets, offering a powerful, sequence-based approach to distinguish true extremophiles from mesophilic contaminants. This methodology is predicated on the well-established principle that organisms adapted to extreme

environments, such as high salinity (halophiles) or high temperature (thermophiles), have evolved proteins with distinct and measurable physicochemical properties. These properties, often referred to as "diagnostic signatures," are a direct consequence of selective pressure to maintain protein structure and function under harsh conditions. The core of this approach involves the systematic analysis of protein sequences to quantify these signatures and apply them as filters. Key properties analyzed include amino acid composition, isoelectric point (pI) distributions, hydrophobicity patterns, and the frequency of specific ionic interactions like salt bridges. By establishing statistical baselines for these properties within a trusted set of extremophile proteins, researchers can then screen larger, unvalidated datasets and flag sequences that deviate significantly from the expected extremophilic profile. This method is particularly valuable for identifying misclassified sequences or contaminants that may have been overlooked by taxonomic or genomic quality filters alone, providing an additional, independent layer of validation that is rooted in the fundamental biology of protein adaptation.

The application of this method typically involves a multi-step computational pipeline. First, a reference dataset of high-confidence extremophile proteins (e.g., from manually curated genomes or well-studied organisms) is compiled. For this reference set, a suite of physicochemical properties is calculated for each protein sequence. This can include metrics like the frequency of acidic versus basic residues, the grand average of hydropathicity (GRAVY), and the predicted isoelectric point. Statistical distributions for these metrics are then generated, establishing a "biophysical fingerprint" for the extremophile group of interest. For example, a signature for halophiles might be a distinct bias in acidic amino acid content. Next, these same properties are calculated for the target dataset that requires curation. Each protein in the target set is then compared against the established extremophile fingerprint. Sequences that fall outside a predefined statistical threshold (e.g., a certain number of standard deviations from the mean) are flagged as potential contaminants. This process can be automated and applied to large-scale datasets, making it a scalable and robust tool for ensuring the purity of data used for downstream applications like training protein language models. The strength of this approach lies in its ability to leverage the fundamental principles of protein science to make informed decisions about sequence validity, moving beyond simple taxonomic labels to a more nuanced, functionally relevant assessment of dataset quality.

## 2.1.1. Diagnostic Signatures for Halophiles

Halophilic organisms, which thrive in high-salt environments, exhibit distinct biophysical and biochemical adaptations that are reflected in the properties of their proteins. These adaptations, often referred to as diagnostic signatures, can be used to distinguish halophilic proteins from those of non-halophilic organisms. One of the most well-known signatures of halophilic proteins is their **highly acidic nature, characterized by a low isoelectric point (pI)** . This is due to an overrepresentation of acidic amino acids like aspartic acid and glutamic acid on the protein surface, which helps to bind a layer of water molecules and hydrated salt ions, thereby preventing protein aggregation and precipitation in high-salt conditions. The shift towards acidic pI is a robust and widely cited diagnostic feature of halophilic proteins. In addition to their acidic nature, halophilic proteins also tend to have a **lower content of hydrophobic amino acids on their surface** and a higher content of small, polar amino acids, which further enhances their solubility and stability in high-salt environments .

Another important diagnostic signature of halophilic proteins is their unique amino acid composition. Studies have shown that halophilic proteins are **enriched in aspartic acid, glutamic acid, and serine, while being depleted in lysine, arginine, and cysteine** . This specific amino acid bias is thought to contribute to the enhanced stability of halophilic proteins by increasing the number of salt bridges and hydrogen bonds on the protein surface, which helps to maintain the protein's structural integrity in the presence of high salt concentrations. Furthermore, the Ramachandran plot data for halophilic proteins often shows a **higher occurrence of amino acids in the helix and sheet regions**, indicating a more rigid and stable structure compared to their mesophilic counterparts . These distinct biophysical and biochemical signatures provide a powerful set of tools for identifying and validating halophilic proteins in large-scale datasets. By analyzing the pI, amino acid composition, and secondary structure propensities of a set of proteins, researchers can filter out potential contaminants and build a more reliable and accurate dataset of halophilic sequences.

### 2.1.2. Diagnostic Signatures for Thermophiles

Thermophilic organisms, which thrive at high temperatures, have evolved a range of adaptations to ensure the stability and function of their proteins under extreme thermal conditions. These adaptations, which can be considered diagnostic signatures, are reflected in the biophysical and biochemical properties of thermophilic proteins. One of the most well-documented signatures of thermophilic proteins is their distinct amino acid composition. Thermophilic proteins are often characterized by a **higher content of charged amino acids, particularly glutamic acid (E) and lysine (K), and a lower**

content of uncharged polar amino acids like glutamine (Q) and histidine (H) . This bias is thought to enhance the formation of ionic interactions, such as salt bridges, which are more stable at high temperatures than the hydrogen bonds that are prevalent in mesophilic proteins. The ratio of **(E+K) to (Q+H)** has been proposed as a useful diagnostic metric for identifying thermophilic proteins, with higher ratios being indicative of a greater thermostability.

In addition to their amino acid composition, thermophilic proteins also exhibit other distinct biophysical properties. They tend to have a **higher hydrophobicity in their core regions**, which contributes to a more compact and stable protein structure . They also tend to have a **higher number of salt bridges on their surface**, which further enhances their thermal stability . The Ramachandran plot data for thermophilic proteins often shows a higher occurrence of amino acids in the most favorable regions, indicating a more optimized and stable structure . Furthermore, thermophilic proteins are often characterized by a **higher content of small, nonpolar amino acids like glycine, valine, and alanine**, which are thought to contribute to the compactness and rigidity of the protein core . These distinct biophysical and biochemical signatures provide a powerful set of tools for identifying and validating thermophilic proteins in large–scale datasets. By analyzing the amino acid composition, hydrophobicity, and secondary structure propensities of a set of proteins, researchers can filter out potential contaminants and build a more reliable and accurate dataset of thermophilic sequences.

### 2.1.3. Statistical Analysis of Amino Acid Composition and Physicochemical Properties

A critical component of the biophysical validation approach is the rigorous statistical analysis of amino acid composition and other key physicochemical properties. This analysis moves beyond simple observation to quantitatively define the molecular adaptations that characterize extremophilic proteins. A foundational study in this area systematically analyzed a large number of protein sequences and their corresponding structures to understand the stability factors that differentiate extremophiles from their non–extremophilic orthologs . This research highlighted that the strategy for packing the protein core is significantly influenced by environmental stresses, leading to substitutive structural events that enhance ionic interactions. The study demonstrated that significant differences exist in the number and composition of amino acids between extremophilic and mesophilic proteins. For instance, the analysis revealed a **notable increase in the number of salt bridges on the surface of extremostable proteins**, a key adaptation for maintaining structural integrity under extreme conditions . Furthermore, the research found that a significant number of small nonpolar amino

acids and a moderate number of charged amino acids, such as Arginine and Aspartic acid, exhibit more nonplanar Omega angles in their peptide bonds, suggesting a role for geometric variability in molecular adaptation .

The statistical power of this approach lies in its ability to create a multi-faceted profile of an extremophile protein. By calculating and comparing the standardized percentage composition of each amino acid, researchers can identify specific biases that serve as robust diagnostic markers. For example, the same foundational study utilized tools like MEGA to calculate and standardize amino acid compositions and plotted the distribution of these compositions across different protein types . This allowed for a direct, quantitative comparison of amino acid abundance between, for example, halophilic and non-halophilic proteins. The study also employed various web-based tools to verify other physicochemical behaviors, including hydrophobicity, isoelectric point (pI), and polar/nonpolar properties, creating a comprehensive biophysical dataset for each protein . The negative correlation observed between pairwise sequence alignments and structural alignments further reinforced the idea that functional stability in extremophiles is not solely dependent on sequence similarity but is deeply rooted in these nuanced physicochemical and structural adaptations. This detailed, quantitative characterization provides a solid statistical foundation for developing filters that can effectively identify and remove sequences that do not conform to the expected extremophilic signature, thereby enhancing the purity and reliability of the final dataset.

## 2.2. Key Citations

### 2.2.1. Panja et al. (2020) on Protein Stability and Structural Plasticity

A seminal paper that provides a strong foundation for the biophysical validation of extremophile datasets is **"Protein stability governed by its structural plasticity is inferred by physicochemical factors and salt bridges,"** published in *Scientific Reports* in 2020 by Anindya S. Panja, Smarajit Maiti, and Bidyut Bandyopadhyay . This study systematically analyzed a large cohort of protein sequences and their corresponding three-dimensional structures to elucidate the fundamental principles governing stability in extremophilic organisms. The authors focused on microorganisms surviving in a wide range of harsh environments, including extremes of temperature, pH, and salinity, making the findings directly relevant to the curation of both halophilic and thermophilic archaeal datasets. The research provides a comprehensive framework for discriminating extremophilic proteins from their non-extremophilic orthologs based on a combination of sequence composition, structural features, and physicochemical properties. The paper's abstract explicitly states that the analysis was designed to

understand protein stability and to discriminate extremophilic proteins, which aligns perfectly with the goal of cleaning datasets for protein language model training .

The methodology and findings of this paper are particularly valuable for developing a curation pipeline. The authors utilized a multi-pronged approach, analyzing amino acid composition, hydrophobicity, isoelectric point (pI), and the prevalence of salt bridges . They found that **extremostable proteins exhibit a significant difference in amino acid composition and a higher number of salt bridges on their surface**, which are crucial for maintaining structural integrity under extreme conditions . Furthermore, their analysis of Ramachandran plot data indicated that amino acids in extremostable proteins are more frequently found in helix and sheet regions, suggesting a more ordered and stable secondary structure . The study also highlighted the role of small nonpolar amino acids and charged residues like Arginine and Aspartic acid in creating geometric variability through nonplanar peptide bond angles, a potential adaptation mechanism . These specific, quantifiable differences provide a robust set of criteria that can be implemented as filters in a bioinformatics pipeline. By calculating these properties for a target dataset and comparing them against the statistical norms established in this study for true extremophiles, one can effectively identify and remove sequences that lack the characteristic biophysical signatures, thus mitigating the risk of contaminating a training set with mesophilic proteins. The paper has been cited by numerous subsequent studies, underscoring its impact and the reliability of its conclusions in the field of extremophile research .

### 2.2.2. Studies on Amino Acid Composition Biases in Extremophiles

Numerous studies have provided detailed insights into the specific amino acid composition biases that characterize extremophilic proteins. For halophiles, a mini-review published in 2014, which remains a foundational text, systematically reviewed genomic and proteomic data to identify the molecular determinants of protein halotolerance . The study underscored the **significant overrepresentation of acidic residues (Asp and Glu) and the underrepresentation of basic residues (Lys) and large hydrophobic residues on the protein surface** as the key adaptive strategy . This creates a highly charged, hydrophilic surface that is essential for solubility and stability in high-salt environments. A more recent study from 2024 further explored this, proposing a trade-off where the halophilic amino acid composition destabilizes the protein under mesophilic conditions but grants remarkable stability in the presence of high salt concentrations, highlighting the evolutionary optimization for a specific niche .

For thermophiles, a comprehensive 2008 study analyzed a large dataset of mesophilic and thermophilic protein pairs to quantify the amino acid replacements associated with thermal adaptation . The findings revealed a clear and statistically significant pattern: thermophilic proteins are characterized by a **decrease in the thermally labile and flexible residues Ser, Asn, Gln, and Thr, and an increase in the rigid and stabilizing residues Ile, Arg, Glu, Lys, and Pro** . The study also identified the most biased replacement events, such as **Ser → Lys and Ser → Ala**, which contribute significantly to the overall compositional shift towards thermostability . Another study focusing on the mussel *Mytilus* provided further evidence for these trends, noting that thermophilic proteins often have a higher content of cysteines, which can form stabilizing disulfide bonds, and a more hydrophobic core . These studies collectively provide a robust, evidence–based framework for using amino acid composition as a diagnostic tool to distinguish extremophiles from mesophiles.

## 2.3. Relevance

### 2.3.1. Distinguishes Obligate Extremophiles from Mesophiles

The primary relevance of the biophysical and statistical validation approach is its demonstrated ability to effectively distinguish between obligate extremophiles and their mesophilic counterparts, including contaminants that may be present in a dataset. This distinction is not based on taxonomic labels, which can be noisy or incorrect, but on the fundamental molecular adaptations that are a direct response to environmental pressures. The 2020 study by Panja et al. provides compelling evidence for this capability . By systematically comparing a large number of extremophilic and non–extremophilic proteins, the researchers identified a suite of physicochemical properties that consistently differ between the two groups. For example, the **significant difference in amino acid composition and the increased number of surface salt bridges in extremostable proteins** are not random variations but are adaptive features that enhance stability in extreme conditions like high salinity or temperature . These features are less likely to be present in mesophilic proteins, which have evolved to function under more moderate conditions and therefore do not require the same level of structural reinforcement. By applying these biophysical criteria as filters, a researcher can effectively screen a dataset and flag sequences that exhibit a mesophilic profile, even if they are annotated as extremophiles.

This method is particularly powerful for identifying **"cryptic" contaminants**—mesophilic sequences that have been misannotated or have infiltrated a dataset through various means. For instance, a mesophilic protein from a co–occurring species in a

metagenomic sample might be incorrectly binned with extremophile sequences. While taxonomic and genomic quality checks might fail to catch this error, a biophysical filter would identify the sequence based on its lack of characteristic extremophilic signatures, such as an insufficient number of salt bridges or an atypical amino acid composition. The negative correlation between sequence and structural similarity observed by Panja et al. further supports this, indicating that even proteins with some sequence homology may have vastly different stability mechanisms, which are reflected in their biophysical properties . Therefore, by implementing a pipeline that calculates these properties and compares them against established extremophile benchmarks, one can create a robust, evidence–based filter that significantly improves the purity of the final dataset. This ensures that a protein language model is trained exclusively on sequences that are genuinely adapted to extreme environments, leading to a more accurate and reliable model.

### 2.3.2. Provides Quantitative Metrics for Filtering

A significant advantage of the biophysical and statistical validation approach is that it provides a set of clear, quantitative metrics that can be used to build objective and reproducible filtering pipelines. Instead of relying on subjective visual inspection or broad qualitative assessments, this method allows for the establishment of specific numerical thresholds for various physicochemical properties. The work by Panja et al. (2020) exemplifies this by not only identifying key differentiating features but also by systematically quantifying them across a large dataset . For example, the study's analysis of amino acid composition, hydrophobicity, and isoelectric point provides a statistical distribution for each of these properties within the extremophile protein set . This allows a researcher to define a "normal" range for an extremophilic protein. A sequence in the target dataset can then be objectively evaluated against this range. If its properties fall outside a certain number of standard deviations, it can be automatically flagged for removal or further inspection. This transforms the curation process from a qualitative art into a quantitative science.

The availability of these quantitative metrics is crucial for developing automated and scalable data cleaning workflows, which are essential for handling the large datasets required for training modern protein language models. For instance, a filtering script could be written to calculate the number of predicted salt bridges for every protein in a dataset and compare it to the distribution reported by Panja et al., where a significant number of such bridges was a key feature of extremostable proteins . Similarly, the percentage of specific amino acids (e.g., acidic residues for halophiles) or the

calculated isoelectric point can be used as hard filters. The study's use of Ramachandran plot data to quantify the occurrence of amino acids in structured regions (helix and sheet) provides another powerful, quantitative criterion for assessing the likely stability of a protein . By combining several of these metrics into a multi-parameter scoring system, a highly robust and sensitive filter can be created. This not only improves the quality of the training data but also enhances the transparency and reproducibility of the curation process, as the specific criteria and thresholds used to clean the dataset are explicitly defined and can be easily reported and replicated by other researchers.

# 3. Approach 3: Machine Learning-Based Data Cleaning

## 3.1. Summary of the Method

The application of machine learning (ML) and, more recently, deep learning (DL) has introduced powerful, data-driven paradigms for cleaning and validating large-scale biological datasets. These methods move beyond static, rule-based filters to learn complex patterns directly from the data, enabling the detection of subtle anomalies, misannotations, and out-of-distribution samples that might otherwise be missed. For the curation of extremophile protein datasets, these techniques offer a promising avenue to automatically identify and correct label noise (e.g., a mesophilic protein mislabeled as thermophilic) and remove sequence contamination. The core principle involves training models to understand the "language" of protein sequences and their associated functions or environmental origins. By learning the characteristic features of a "clean" extremophile dataset, these models can flag sequences that deviate from the expected distribution, suggesting they may be contaminants or misclassified entries. This approach is particularly valuable given the exponential growth of protein databases, where manual curation is no longer feasible, and the prevalence of noisy, automatically generated annotations is high .

The methodologies can be broadly categorized into several key strategies. First, the use of **pre-trained protein language models (pLMs)** like ESM-2 has become a cornerstone. These models, trained on billions of sequences, generate high-dimensional numerical representations (embeddings) that capture a rich array of biochemical and evolutionary information. These embeddings serve as powerful features for downstream classification or clustering tasks. Second, these embeddings are fed into various ML classifiers, such as Support Vector Machines (SVMs), Random Forests, or more complex deep learning architectures, to predict labels (e.g.,

"thermophile" vs. "mesophile") or to identify anomalous sequences. Third, a novel and highly effective strategy is the use of **recursive or iterative cleaning frameworks**, exemplified by ProtAC. This approach involves a continuous cycle of pre–training a model on a noisy dataset, fine–tuning it on a smaller, high–quality curated set, and then using the refined model to clean the original noisy data. This cleaned dataset is then used for the next round of pre–training, creating a virtuous cycle where both the model and the dataset quality improve in tandem . Finally, advanced techniques like **meta-learning** are being explored to train models that are robust to noisy and under–labeled data, a common challenge in protein engineering datasets derived from high–throughput screens .

### 3.1.1. Using Pre–trained Protein Language Models (pLMs) for Embeddings

The advent of large–scale pre–trained protein language models (pLMs) has revolutionized the field of computational biology, providing a powerful new tool for representing and analyzing protein sequences. These models, such as **ESM–2**, are trained on massive datasets of protein sequences and learn to capture the complex patterns and relationships that are embedded in the language of life . The resulting sequence embeddings, which are high–dimensional vector representations of the proteins, encode a rich set of biophysical and evolutionary information that can be used for a wide range of downstream tasks, including protein function prediction, structure prediction, and classification. For the purpose of dataset curation, these embeddings can be used to identify and remove label noise and out–of–distribution sequences. By representing each protein in a continuous vector space, it becomes possible to apply sophisticated machine learning techniques to identify sequences that are anomalous or that have been mislabeled.

The key advantage of using pLM embeddings is that they provide a comprehensive and context–aware representation of the protein, going far beyond what can be captured by a small set of manually curated physicochemical features. For the purpose of data cleaning, these embeddings can be used as input features for a downstream machine learning model, such as a classifier or an anomaly detector. For example, a classifier can be trained to distinguish between halophilic and non–halophilic proteins by learning the decision boundary in the high–dimensional embedding space that separates the two classes . Because the embeddings are learned from a vast amount of data, they are robust and can generalize well to new, unseen sequences. This makes them an ideal foundation for building models that can accurately identify mislabeled or out–of–

distribution sequences in a dataset, even if those sequences are not closely related to any of the sequences in the training set.

### 3.1.2. Classification Models (SVM, Random Forest, XGBoost) for Label Prediction

Once a set of features has been extracted from the protein sequences, either through traditional biophysical descriptors or, more effectively, through pLM embeddings, a classification model can be trained to predict the labels of the sequences. A variety of machine learning models have been successfully applied to this task, including **support vector machines (SVMs), random forests, and gradient boosting machines like XGBoost**. These models are trained on a subset of the data for which the labels are known with high confidence (the "ground truth" set). The model learns a mapping from the input features to the output labels, effectively learning the characteristics that define each class (e.g., halophile vs. mesophile). For example, an SVM would find a hyperplane in the feature space that best separates the two classes, while a random forest would build an ensemble of decision trees, each of which makes a prediction based on a subset of the features . After training, the model's performance is evaluated on a separate validation set to ensure that it has learned the underlying patterns and is not simply memorizing the training data (i.e., to check for overfitting). Once a satisfactory level of performance is achieved, the model can be deployed to predict the labels of the remaining sequences in the dataset. Any sequence whose predicted label conflicts with its original annotation is flagged as a potential instance of label noise. This approach provides a systematic and quantitative way to identify and correct mislabeled data, leading to a cleaner and more reliable dataset for downstream analysis.

### 3.1.3. Recursive Cleaning Frameworks (ProtAC)

A significant advancement in data–centric AI for bioinformatics is the development of recursive cleaning frameworks, designed to systematically improve the quality of massive, noisy datasets. The **ProtAC (Protein Automatic Cleaning)** framework stands out as a state–of–the–art example of this approach, specifically tailored for large–scale protein datasets . The core innovation of ProtAC lies in its iterative **"train–finetune–clean" cycle**, which leverages both sequence and functional information through a multimodal learning architecture. The process begins by pre–training a protein language model on a large, uncurated dataset like UniRef50, which contains a substantial amount of label noise. This initial model learns a general representation of protein sequences. In the next stage, the model is fine–tuned on a smaller, high–confidence dataset, such as SwissProt, where the functional annotations are manually

curated and reliable. This fine-tuning step adapts the model's general knowledge to the specific task of accurately matching sequences with their correct functional annotations.

The crucial third step is the cleaning phase. The fine-tuned model is used to analyze the original noisy dataset. It employs a specialized module called the **Sequence-Annotation Matching (SAM) filter**, which evaluates the compatibility between a protein's sequence and its annotated function. If the model determines that an annotation is inconsistent with the sequence's learned representation, it can flag or revise the annotation. The dataset, now partially cleaned, is then used to pre-train the model again in the next iteration. This recursive process creates a powerful feedback loop: as the dataset becomes cleaner with each cycle, the model trained on it becomes more accurate; in turn, the more accurate model is better equipped to perform further cleaning. Through multiple rounds of this cycle, ProtAC has been shown to progressively improve both the performance of the resulting protein language model on downstream tasks and the overall quality of the dataset itself. The framework successfully cleaned the UniRef50 dataset, which contains approximately 50 million proteins, demonstrating its scalability and effectiveness in handling the vast amounts of data required for training modern protein language models .

### 3.1.4. Meta-Learning for Noisy and Under-labeled Data

Meta-learning, or "learning to learn," has emerged as a powerful strategy to address the pervasive problem of noisy and under-labeled data in machine learning, with significant potential for protein engineering and dataset curation. In the context of extremophile protein datasets, where experimental validation is costly and time-consuming, a large portion of the data may be unlabeled or have unreliable labels. Meta-learning approaches are particularly well-suited for this challenge because they are designed to learn from a small set of high-quality, trusted labels and then generalize that knowledge to a larger, noisy dataset. This is achieved through a bi-level optimization process where the model learns not just the task (e.g., classifying a protein as thermophilic or mesophilic) but also how to learn effectively from limited and imperfect data. A 2025 study demonstrated the application of meta-learning to antibody engineering, where yeast display screens generate large but noisy datasets. The meta-learning model was able to learn robust sequence-function relationships despite the presence of noise, effectively reducing the need for extensive experimental screening and improving the reliability of the resulting ML models .

The relevance of this approach to curating extremophile datasets is profound. One could envision a scenario where a small, manually curated "gold standard" set of extremophile and mesophile proteins is used to train a meta-learning model. This model would then be applied to a much larger, uncurated database. The meta-learner's ability to handle positive and unlabeled learning (where only positive examples are confidently known) and to learn out-of-distribution properties would be invaluable. It could identify proteins in the large dataset that are likely extremophiles, even if they are dissimilar to the examples in the training set, while simultaneously flagging potential mislabeled contaminants. By leveraging the structure of the learning problem itself, meta-learning provides a principled way to build robust models that can sift through vast amounts of noisy biological data, making it a highly promising tool for ensuring the purity of datasets used to train protein language models for extremophile research .

## 3.2. Key Citations

### 3.2.1. HPClas: Halophilic Protein Classification using ESM-2 and SVM

A prime example of the successful application of machine learning for classifying extremophilic proteins is the **HaloClass model**, detailed in a 2024 paper . This study presents an SVM classifier that leverages embeddings from the **ESM-2 protein language model** to accurately identify salt-tolerant (halophilic) proteins. The authors highlight a key limitation of previous methods, which relied on manually selected, human-interpretable features like amino acid frequencies and dipeptide counts. While intuitive, these approaches lose the rich contextual information encoded in the protein sequence, such as the interactions between distant amino acids. HaloClass overcomes this by using the powerful, context-aware representations learned by the ESM-2 model. The study demonstrates that this approach significantly outperforms existing methods, particularly on a newer and larger test dataset composed of proteins that are distant from the training set. This showcases the superior generalization capabilities of the pLM-based approach. Furthermore, the paper demonstrates the utility of HaloClass in a practical application: the guided design of salt-tolerant enzymes. In a mutation study, the model was able to accurately predict changes in salt tolerance based on single- and multiple-point mutants, outperforming other existing approaches. This work provides a clear and compelling case for the use of pre-trained protein language models as a foundation for building highly accurate and robust classifiers for extremophile proteins, which can in turn be used to clean and validate large-scale datasets.

### 3.2.2. ProtAC: Recursive Cleaning for Large-scale Protein Data

The **ProtAC framework**, detailed in a 2024 pre-print on bioRxiv and a 2025 conference paper on OpenReview, represents a significant leap forward in the automated curation of large-scale protein datasets . The research introduces a novel, iterative methodology designed to clean massive protein databases by leveraging a multimodal protein language model that integrates both sequence and functional information. The core of the ProtAC system is its recursive **"train-finetune-clean" cycle**. The process initiates with pre-training a model on a large, inherently noisy dataset like UniRef50. This model is then fine-tuned on a smaller, high-quality, manually annotated dataset such as SwissProt. The fine-tuned model is subsequently deployed to clean the original noisy dataset using its **Sequence-Annotation Matching (SAM) module**, which filters and corrects functional annotations based on their compatibility with the protein sequence. This cleaned dataset is then fed back into the pre-training stage for the next iteration, creating a self-reinforcing loop that progressively enhances both the model's predictive power and the dataset's purity.

The study demonstrates the framework's efficacy through extensive biological analysis and performance benchmarks. After multiple cleaning cycles, the ProtAC-refined model achieved state-of-the-art performance on various function-related downstream tasks, outperforming competitors with fewer than 100 million parameters. A key outcome of this work was the generation of a cleaned version of the UniRef50 dataset, containing approximately **50 million proteins** with improved functional annotations. The authors validated their approach by visualizing the model's learned sequence embeddings using t-SNE. They showed that embeddings from a model trained on the cleaned dataset exhibited better-defined clusters and clearer separation between different cellular components (e.g., cytoplasm, nucleus, extracellular region) compared to a model trained on the original noisy data. This visual evidence confirms that the recursive cleaning process enhances the model's ability to learn meaningful biological representations, making ProtAC a powerful tool for preparing high-quality training data for protein language models, including those focused on extremophiles .

### 3.2.3. Meta-Learning for Antibody Engineering with Noisy Data

A 2025 study published in *Cell Patterns* highlights the application of meta-learning to tackle the challenge of noisy and under-labeled data in the context of antibody engineering, offering a methodological blueprint that is highly relevant for curating extremophile protein datasets . The research addresses a common bottleneck in protein engineering: the generation of large, high-quality labeled datasets is often

prohibitively expensive and time-consuming. Directed evolution and high-throughput screening experiments frequently produce datasets that are either sparsely labeled or contain significant experimental noise. The authors propose meta-learning as a solution to learn effectively from such imperfect data. By employing a bi-level optimization strategy, the meta-learning model is trained on a small set of trusted labels and learns to generalize this knowledge to a larger, noisier dataset. This approach is particularly adept at handling scenarios with positive and unlabeled data, as well as learning properties of sequences that are out-of-distribution relative to the training set.

The study demonstrates that meta-learning can significantly expedite experimental workflows and improve the robustness of predictive models in the face of noisy data. For instance, in yeast display antibody library screens, the meta-learning approach was able to learn robust sequence-function relationships, reducing the reliance on exhaustive experimental validation. The principles demonstrated in this work are directly transferable to the curation of extremophile protein datasets. One could use a small, high-confidence set of experimentally validated extremophile and mesophile proteins to train a meta-learner. This model could then be applied to vast public databases to identify likely extremophiles, correct mislabeled sequences, and flag potential contaminants, even for proteins with novel sequences or from under-sampled taxa. The ability of meta-learning to operate effectively with limited and noisy labels makes it a powerful and efficient tool for ensuring dataset purity, which is critical for training reliable protein language models .

### 3.3. Relevance

### 3.3.1. Detects and Corrects Label Noise

Machine learning-based data cleaning methods are exceptionally well-suited for detecting and correcting label noise, a pervasive problem in public protein databases that directly impacts the performance of downstream models, including protein language models (pLMs). Label noise refers to instances where a protein's functional or environmental annotation is incorrect—for example, a mesophilic protein being misclassified as thermophilic. This can occur due to errors in automated annotation pipelines, outdated taxonomic information, or experimental mischaracterization. The **ProtAC framework** provides a prime example of how ML can address this issue systematically . By training a multimodal model to understand the relationship between a protein's sequence and its annotated function, ProtAC can identify mismatches. Its **Sequence-Annotation Matching (SAM) module** acts as a filter, flagging annotations

that are inconsistent with the sequence's learned representation. Through its recursive cleaning cycle, the framework not only removes incorrect labels but also proposes more suitable functional annotations, effectively correcting the label noise. This process is scalable and has been successfully applied to the UniRef50 database, demonstrating its capacity to clean datasets containing tens of millions of proteins.

The relevance of this capability for training a pLM on extremophile sequences is paramount. A model trained on a dataset with significant label noise will learn incorrect associations between sequence features and extremophilic adaptation. For instance, it might learn that certain amino acid compositions are characteristic of thermophiles when they are, in fact, just artifacts of mislabeled mesophilic contaminants. This would severely compromise the model's ability to accurately predict the properties of novel extremophile proteins or to design new ones with enhanced stability. By employing a recursive cleaning framework like ProtAC, researchers can generate a "clean" training set where the labels (e.g., "halophile," "thermophile") are highly reliable. This ensures that the pLM learns genuine, biologically meaningful signatures of extremophilic adaptation, leading to a more accurate and robust model. The ability to automatically correct label noise at scale is therefore a critical step in ensuring the scientific validity and practical utility of the resulting language model.

### 3.3.2. Identifies Out-of-Distribution Sequences

A key strength of machine learning–based data cleaning approaches is their ability to identify **out-of-distribution (OOD) sequences**—proteins that are statistically or biologically dissimilar from the majority of the data in a given dataset. In the context of curating an extremophile protein dataset, OOD sequences are often contaminants, such as mesophilic proteins that have been incorrectly included. These sequences can introduce significant noise and bias into the training data for a protein language model. Advanced ML techniques, particularly those leveraging pre-trained protein language models (pLMs), are adept at detecting such anomalies. pLMs like ESM-2 learn dense, high-dimensional embeddings that capture a wealth of evolutionary and structural information. When these embeddings are visualized (e.g., using t-SNE or UMAP), proteins with similar properties tend to cluster together, while OOD sequences often appear as isolated points or in distinct, unexpected clusters. This makes them readily identifiable.

The **ProtAC framework**, for example, demonstrates this capability effectively. By comparing the sequence embeddings of proteins in the noisy UniRef50 dataset with their functional annotations, the model can identify sequences whose embeddings do

not align with their labeled function, suggesting they are OOD . Similarly, meta-learning approaches are explicitly designed to handle OOD properties, learning to recognize sequences that deviate from the expected distribution even if they are not exact matches to known contaminants . The relevance for training a pLM on extremophiles is clear: by filtering out these OOD sequences, one can ensure a higher degree of dataset purity. This prevents the model from being confused by irrelevant or misleading examples, allowing it to focus on learning the true sequence determinants of extremophilic adaptation. The result is a more focused, accurate, and reliable protein language model, capable of making better predictions and generating more valid hypotheses about extremophile biology.

## 4. Approach 4: Contamination Removal and Lineage Validation

### 4.1. Summary of the Method

Contamination removal and lineage validation are critical steps in the curation of extremophile datasets, addressing the pervasive issue of foreign sequences in public databases. This approach focuses on identifying and excising proteins that originate from a different species than the one being studied, a problem that can arise from various sources, including sample contamination, misassembly of genomes, or errors in automated annotation. The core of this methodology involves comparing sequences within a target genome or dataset against a comprehensive reference database to detect anomalies in taxonomic origin. Advanced tools have been developed to automate this process with high sensitivity and specificity. These tools typically work by aligning protein sequences from a query genome to a large, non-redundant database of known proteins. They then analyze the taxonomic distribution of the top hits for each query sequence. If a query protein from a eukaryotic genome, for example, shows its strongest similarity to bacterial proteins, it is flagged as a potential contaminant. This process helps to distinguish true, vertically inherited genes from foreign DNA that has been inadvertently included in the assembly.

A key challenge in this area is to differentiate true contamination from legitimate biological phenomena such as **horizontal gene transfer (HGT)** , where genes are transferred between organisms through non-reproductive means. Sophisticated methods are therefore required to make this distinction. Some tools achieve this by considering factors like the length of the foreign contig, the number of consecutive foreign genes, and the phylogenetic context. For instance, a single bacterial gene on a large eukaryotic scaffold might be an HGT event, whereas a cluster of bacterial genes

on a short contig is more likely to be a contaminant. The validation of lineage is also crucial; by constructing phylogenetic trees, researchers can identify sequences that have an anomalous evolutionary history, suggesting they may be contaminants or the result of ancient HGT events. This combined approach of direct sequence comparison for contamination detection and phylogenetic analysis for lineage validation provides a powerful, multi-layered defense against the inclusion of erroneous sequences in extremophile datasets, ensuring that the data used for downstream analyses like training protein language models is as accurate and pure as possible.

## 4.1.1. ContScout for Sensitive Detection and Removal of Contaminating Proteins

**ContScout** is a state-of-the-art computational tool specifically designed for the sensitive and precise detection and removal of foreign sequences from annotated genomes, representing a significant advancement in the field of dataset curation . Introduced in a 2024 *Nature Communications* article, ContScout addresses the growing problem of contamination in genomic databases, which can severely impact downstream applications such as comparative genomics, phylogenetic analysis, and the training of machine learning models . The tool operates by performing a comprehensive search of all proteins encoded in a query genome against a large, curated reference database. Its key innovation lies in its sophisticated algorithm for evaluating the taxonomic origin of each protein. For each query protein, ContScout identifies the top-scoring matches in the reference database and analyzes their taxonomic lineage. It then calculates a **"taxon support value,"** which represents the proportion of these top hits that support the expected taxonomic origin of the query protein (e.g., the proportion of hits that are also from the genus *Halobacterium* for a query protein from a *Halobacterium* genome). Proteins with a low taxon support value are flagged as potential contaminants.

The power of ContScout has been demonstrated through extensive benchmarking against other existing tools like Conterminator and BASTA, as well as on manually curated datasets where the true contaminants are known . In these comparisons, ContScout consistently outperformed its competitors, identifying **five to ten times more contaminating proteins with high accuracy** . For example, in a test on the genome of *Papilio xuthus*, which was known to be contaminated with microsporidian proteins, ContScout successfully identified **98% of the contaminant proteins**, while other tools identified only 21% or less . A key feature of ContScout is its ability to maintain high specificity, meaning it produces very few false positives. This is crucial for avoiding the erroneous removal of legitimate, native proteins from the dataset. The

tool is also designed to be sensitive enough to detect contamination even when the contaminant is a closely related species, a challenging scenario for many other methods. Furthermore, ContScout can help distinguish between contamination and horizontal gene transfer (HGT), providing a more nuanced understanding of the genome's evolutionary history. By providing a highly sensitive and specific method for cleaning genomes, ContScout offers an invaluable resource for researchers seeking to create high-quality, reliable datasets for extremophile research.

### 4.1.2. Phylogenetic Analysis to Identify and Remove HGTs

Phylogenetic analysis is a powerful tool for validating the evolutionary history of protein sequences and identifying potential contaminants, particularly those resulting from horizontal gene transfer (HGT). The core principle is that true, vertically inherited genes should exhibit a phylogenetic tree that is congruent with the species tree of the organisms from which they were derived. In contrast, genes acquired through HGT will often have a phylogenetic placement that is inconsistent with the expected lineage. By constructing phylogenetic trees for individual protein families or for sets of conserved marker genes, researchers can identify sequences that have an anomalous evolutionary history. For example, a protein from a putative archaeal thermophile that consistently groups with bacterial mesophiles in a well-supported phylogenetic tree is a strong candidate for being a contaminant or an HGT event.

The process of using phylogenetic analysis for dataset cleaning typically involves several steps. First, a set of target protein sequences is selected from the dataset. These sequences are then aligned with a comprehensive set of homologous sequences from a reference database, such as UniProt or NCBI's nr database. A phylogenetic tree is then constructed from this alignment using a method like maximum likelihood or Bayesian inference. The resulting tree is then inspected for sequences that have an unexpected phylogenetic placement. This can be a time-consuming process, but it is a highly effective way to identify contaminants that may have been missed by other methods. The use of automated tools for phylogenetic analysis can help to streamline this process, making it more scalable for large datasets. By identifying and removing these anomalous sequences, researchers can create a dataset that more accurately reflects the core, vertically inherited extremophile proteome, which is essential for studying the evolution of extremophily and for training accurate protein language models.

### 4.1.3. Taxonomic and Compositional Analysis for Misclassification Detection

Taxonomic and compositional analysis provides another layer of validation for identifying misclassified sequences and potential contaminants in extremophile datasets. This approach leverages the fact that organisms from different taxonomic groups often have distinct genomic and proteomic signatures. For example, the GC content, codon usage bias, and amino acid composition can vary significantly between different phyla or even between different species within the same genus. By analyzing these compositional properties, it is possible to identify sequences that are outliers and may be the result of misclassification or contamination. For instance, a protein sequence with a GC content that is significantly different from the average GC content of the genome from which it was derived may be a contaminant.

This approach can be implemented in a number of ways. One common method is to use a tool like **TaxFlow-CU**, which was developed for the curation of viral sequences but can be adapted for other types of organisms . TaxFlow-CU uses a combination of quality control, taxonomic classification, and manual curation to identify and remove contaminants. The tool first filters the sequences based on quality scores and length. It then uses a classifier to assign a taxonomic label to each sequence. Finally, it performs a manual review to remove any sequences that have an anomalous taxonomic assignment. Another approach is to use a tool like **MAGpurify**, which uses a combination of compositional and phylogenetic methods to identify and remove contaminating contigs from MAGs. By combining these different approaches, researchers can create a robust and comprehensive pipeline for identifying and removing misclassified sequences and contaminants from their extremophile datasets.

## 4.2. Key Citations

### 4.2.1. Bálint et al. (2024) on ContScout

The definitive resource for the contamination removal tool ContScout is the 2024 paper in *Nature Communications* titled **"ContScout: sensitive detection and removal of contamination from annotated genomes"** . The authors of this study are Balázs Bálint, Zsolt Merényi, Botond Hegedüs, Igor V. Grigoriev, Zhihao Hou, Csenge Földi, and László G. Nagy. This open-access article provides a comprehensive overview of the tool's development, methodology, and performance. The abstract clearly states the problem that ContScout aims to solve: the increasing recognition of contamination in genome databases and its detrimental effect on downstream applications, including metagenomics and comparative evolutionary genomics . The paper details how ContScout achieves high specificity and sensitivity on synthetic benchmark data, even when the contaminant species is closely related to the host. This is a critical feature for

working with extremophile datasets, where the distinction between a novel, closely related extremophile and a contaminant can be subtle.

The paper is particularly valuable because it includes rigorous benchmarking against other established tools and validation on manually curated datasets. The authors compared ContScout's performance with that of Conterminator and BASTA, demonstrating that **ContScout identified significantly more contaminating proteins** . They also validated the tool's accuracy on four datasets where the contaminants had been manually identified, showing that ContScout could detect all the foreign proteins while other tools missed a substantial proportion . The study further highlights the impact of contamination on phylogenomic analyses, showing that contaminants can lead to erroneous conclusions about ancestral gene content, such as inflating estimates of gene loss . This makes the paper not just a description of a new tool, but also a compelling argument for the importance of rigorous decontamination in all genomic studies. For any researcher looking to implement a robust contamination removal step in their dataset curation pipeline, this paper is an essential read, providing both the theoretical background and the practical evidence for the effectiveness of the ContScout approach.

### 4.2.2. Studies on HGT Detection in Archaeal Phylogenetic Trees

While not explicitly detailed in the provided search results, the detection of Horizontal Gene Transfer (HGT) in archaeal phylogenetic trees is a well–established and critical component of lineage validation. The principle is based on the incongruence between a gene's phylogenetic tree and the established species tree of the organisms. If a gene from an archaeal extremophile consistently clusters within a bacterial clade with high statistical support, it is a strong indicator of an ancient or recent HGT event. This is particularly relevant for extremophiles, as the harsh environments they inhabit can facilitate gene transfer between distantly related but co–occurring species. For example, a study on the evolution of thermophily in Archaea might use phylogenetic analysis of key metabolic genes to distinguish between vertical inheritance and acquisition from other thermophiles via HGT. This is crucial for understanding the true evolutionary history of extremophilic adaptations and for building accurate datasets. Tools like **PhyloNet** and **Notung** are commonly used to reconcile gene trees with species trees and to infer the history of gene transfer events. By identifying and potentially removing these HGT–derived sequences, researchers can create a dataset that more accurately reflects the core, vertically inherited extremophile proteome,

which is essential for training a protein language model that captures the fundamental principles of extremophily.

### 4.2.3. TaxFlow-CU for Viral Sequence Curation

While not directly focused on extremophiles, a 2024 study on the curation of viral sequences from metagenomic data provides a valuable methodological parallel that can be adapted for extremophile datasets . The paper details a rigorous bioinformatics pipeline for processing full-length 16S rRNA gene amplicons, a common marker gene used for microbial community profiling. The process begins with quality control, using the `dada2` R package to filter and trim sequences based on quality scores ( `minQ = 3` ), length ( `minLen = 1000` , `maxLen = 2100` ), and the number of expected errors ( `maxEE = 2` ). This initial step is crucial for removing low-quality reads that can lead to spurious operational taxonomic units (OTUs) or amplicon sequence variants (ASVs). The pipeline then uses a model to learn the error rates from the data and infers the true ASVs, which are treated as proxies for species-level taxa.

The most relevant part of this methodology for dataset cleaning is the manual curation step. After initial taxonomic classification using tools like the RDP classifier, the researchers performed a manual review to remove potential contaminants. ASVs were discarded if they were classified as chloroplast, eukaryotic, or archaeal (when the focus was on bacteria), or if they could not be classified at a high taxonomic level (e.g., class) and did not map to a known OTU. Furthermore, ASVs that were not assigned to a species or a known OTU were also discarded as potential contaminants . This conservative approach, which combines automated quality filtering with manual curation based on taxonomic logic, is directly applicable to cleaning extremophile datasets. For example, when building a dataset of archaeal thermophiles, one could implement a similar pipeline to automatically discard any sequences that are classified as bacterial or eukaryotic, or that lack a clear archaeal lineage, thereby removing a significant source of phylogenetic contamination. This ensures that the final dataset is composed of sequences with a high degree of taxonomic confidence.

### 4.3. Relevance

### 4.3.1. Removes Contaminating Sequences from Annotated Genomes

The primary relevance of tools like ContScout is their ability to systematically and accurately remove contaminating sequences from annotated genomes, a critical step in ensuring the purity of extremophile protein datasets. Public databases like GenBank

and RefSeq, while invaluable resources, are not immune to contamination. This can occur at various stages, from sample collection and DNA extraction to sequencing, assembly, and annotation. A common source of contamination is the inadvertent inclusion of DNA from other organisms present in the sample or the laboratory environment. For example, a bacterial genome assembly might contain short contigs derived from the host organism (if the sample was from a host–associated bacterium) or from other bacteria in the lab. These contaminating sequences, if not removed, will be annotated as part of the target organism's proteome, leading to a dataset that is a mixture of true and foreign proteins. This "label noise" can be highly detrimental to the training of a protein language model, as the model would learn from incorrect examples, potentially leading to flawed predictions about the properties and functions of extremophile proteins.

ContScout directly addresses this problem by providing a sensitive and specific mechanism for identifying these foreign sequences . By comparing each protein in a query genome to a comprehensive reference database and analyzing the taxonomic distribution of the best hits, ContScout can effectively flag proteins that have an anomalous origin . The benchmarking results presented in the Bálint et al. (2024) paper demonstrate its superior performance in this task, identifying a much larger fraction of known contaminants compared to other tools . For a researcher building a dataset of, for example, halophilic archaeal proteins, running ContScout on the source genomes would be a crucial quality control step. It would help to ensure that the final protein sequence file contains only proteins that are genuinely encoded by the halophilic archaea, free from bacterial or eukaryotic contaminants. This level of purity is essential for training a high–quality, reliable protein language model that can accurately capture the unique characteristics of extremophile proteins.

### 4.3.2. Filters Out Sequences with Conflicting Phylogenetic Signals

Another key relevance of the contamination removal and lineage validation approach is its ability to filter out sequences with conflicting phylogenetic signals. This is particularly important for distinguishing between true contamination and horizontal gene transfer (HGT), a natural biological process that can also introduce foreign genes into a genome. While HGT is a legitimate evolutionary mechanism, it can be a source of "contamination" in a dataset if the goal is to study the vertically inherited core proteome of an extremophile. By constructing phylogenetic trees for individual protein families, researchers can identify sequences that have an anomalous evolutionary history. For example, a protein that groups with a distantly related taxon in a

phylogenetic tree, while all other proteins from the same organism group with their expected taxonomic group, is a strong candidate for an HGT event.

The ability to distinguish between contamination and HGT is crucial for building accurate and reliable datasets. If all foreign sequences are simply removed without considering their evolutionary history, it may lead to an incomplete or biased view of the organism's biology. By using a combination of tools like ContScout and phylogenetic analysis, researchers can make more informed decisions about which sequences to include or exclude from their datasets. For example, a sequence that is flagged as a contaminant by ContScout but has a phylogenetic placement that is consistent with a known HGT event might be retained in the dataset, but flagged as an HGT–derived sequence. This allows for a more nuanced and accurate representation of the organism's proteome, which is essential for downstream analyses like training a protein language model. By filtering out sequences with conflicting phylogenetic signals, researchers can ensure that their datasets are as accurate and representative as possible, leading to more reliable and meaningful results.