HIGHCLASS: EFFICIENT METAGENOMIC CLASSIFICATION VIA QUALITY-AWARE TOKEN MAPPING AND SPARSIFIED INDEXING

Anonymous authors

Paper under double-blind review

ABSTRACT

Metagenomic classification requires both high accuracy and computational efficiency to process the exponentially growing volume of sequencing data. We present *HighClass*, a novel classification framework that fundamentally transforms the computational paradigm through variable-length token indexing, quality-aware scoring, and learned sparsification.

Our key innovation replaces alignment operations with hash-based token mapping, achieving $O(|\mathcal{T}|)$ complexity while maintaining competitive accuracy. We establish rigorous theoretical foundations: (1) generalization bounds proving $O(\sqrt{V|\mathcal{Y}|/n})$ convergence for vocabulary size V and $|\mathcal{Y}|$ taxa; (2) concentration inequalities under exponential α -mixing with explicit dependency factors; (3) consistency guarantees for maximum likelihood classification under identifiability conditions. HighClass achieves 85.1% F1 on CAMI II—within 1.5% of state-of-the-art—while delivering 4.2× speedup and 68% memory reduction. Variable-length tokens provide 6.8 percentage points improvement over fixed k-mers through superior pattern capture. Quality-aware scoring with learned sensitivity $\eta=1.8$ optimally weights sequencing evidence. Gradient-based sparsification retains 32% of genomic regions while preserving 94% accuracy.

Beyond empirical gains, our work establishes the first comprehensive theory of token-based genomic classification, providing uniform convergence guarantees and explicit characterization of dependency effects through α -mixing analysis. These results transform sequence classification from heuristic approaches to principled methods with provable guarantees.

1 Introduction

Metagenomic sequencing generates unprecedented volumes of data requiring classification at rates exceeding 10^{10} reads per day in clinical and environmental applications (Lloyd-Price et al., 2019; Thompson et al., 2017; Gardy & Loman, 2018). The fundamental challenge is determining the taxonomic origin of each read $X \in \mathcal{X}$ with respect to a reference database \mathcal{D} while maintaining both accuracy and computational tractability.

1.1 THE FUNDAMENTAL TRADE-OFF

Current metagenomic classifiers fall into two paradigmatic categories, each with inherent limitations:

Alignment-based methods solve classification through explicit sequence alignment, typically employing seed-and-extend strategies. For a read $X \in \Sigma^m$ and database \mathcal{D} , practical implementations achieve high accuracy but spend most time in seed-and-extend steps with effective per-read cost $\mathcal{O}(m\log n + k\log k)$ where n is index size and k the number of k-mer matches. The computational burden becomes prohibitive for modern datasets exceeding 10^{10} reads.

Alignment-free methods bypass explicit alignment using k-mer indexing or minimizer schemes. While achieving $\mathcal{O}(m)$ query complexity, they sacrifice accuracy through information loss during the fixed-length decomposition, particularly problematic for: (i) reads with heterogeneous quality

profiles, (ii) closely related taxa differing by subtle variations, and (iii) novel organisms absent from the reference.

This dichotomy has persisted for over a decade, with incremental improvements failing to bridge the accuracy-speed gap. While recent hardware-based approaches like in-storage computing (Ghiasi et al., 2022; 2023; 2024) have shown promise in accelerating genomic analysis by moving computation closer to data, they require specialized hardware infrastructure. We argue that the fundamental limitation can be addressed algorithmically by treating genomic sequences not as monolithic strings but as compositions of quality-aware tokens that exploit both structural patterns and uncertainty information.

1.2 Our Approach: Theoretical Framework for Token Dependencies

Token-based classification presents three fundamental dependency structures requiring rigorous theoretical treatment: (i) overlapping extraction windows where adjacent tokens share subsequences, (ii) sequential correlations arising from genomic structure, and (iii) spatially clustered quality patterns from sequencing technology.

We develop a comprehensive theoretical framework addressing these dependencies through exponential α -mixing analysis. Our approach establishes concentration inequalities that explicitly quantify dependency effects, proving that token scores concentrate around their expectations with variance factor $(1+2C/\gamma)$ where C and γ characterize the mixing rate. This theoretical foundation, combined with empirical validation on CAMI II data showing $\gamma\approx 0.15$, demonstrates that our multi-stage architecture with 32,000-token vocabulary provides sufficient statistical redundancy for robust classification.

Our theoretical contributions (Section 4) include: (i) generalization bounds via Rademacher complexity establishing $O(\sqrt{V|\mathcal{Y}|/n})$ convergence rate, (ii) concentration inequalities under α -mixing quantifying token dependency impact, and (iii) consistency results proving asymptotic optimality of maximum likelihood classification.

1.3 OUR CONTRIBUTIONS

We establish three fundamental advances in metagenomic classification:

- (1) Theoretical Foundations. We develop the first rigorous theoretical framework for token-based genomic classification. Our generalization bound (Theorem 6) establishes uniform convergence at rate $O(\sqrt{V|\mathcal{Y}|/n})$ for vocabulary size V, taxonomic classes $|\mathcal{Y}|$, and sample size n. Concentration inequalities (Lemma 7) quantify the impact of token dependencies through exponential α -mixing coefficients. Consistency results (Theorem 8) prove asymptotic optimality of maximum likelihood classification under identifiability conditions.
- (2) Algorithmic Innovation. HighClass transforms the computational paradigm by replacing alignment operations with hash-based token mapping. Variable-length tokens capture discriminative genomic patterns with superior statistical power compared to fixed k-mers, achieving $O(|\mathcal{T}|)$ query complexity. Quality-aware scoring with learned sensitivity $\eta \approx 1.8$ optimally weights evidence based on sequencing confidence. Gradient-based sparsification reduces the index to 32% of original size while preserving 94% accuracy through principled feature selection.
- (3) Empirical Excellence. Comprehensive evaluation on CAMI II demonstrates 85.1% F1 with $4.2\times$ speedup and 68% memory reduction compared to state-of-the-art methods. Rigorous ablation studies isolate contributions: variable-length tokens provide 6.8 percentage points over k-mers, quality weighting adds 1.9 points, and sparsification enables deployment with 6.8GB memory. Statistical validation includes 95% confidence intervals, Wilcoxon signed-rank tests with Holm-Bonferroni correction, and Cohen's d effect sizes quantifying practical significance.
- Our work synthesizes QA-Token (Gollwitzer et al., 2025) vocabularies, MetaTrinity's (Gollwitzer et al., 2023) multi-stage architecture, and gradient-based sparsification inspired by genome sparsification techniques (Alser et al., 2024) into a unified theoretical and practical framework that advances the state of metagenomic classification.

2 RELATED WORK

2.1 FOUNDATION TECHNOLOGIES

Our work builds upon three key recent advances:

QA-Token (Gollwitzer et al., 2025) provides quality-aware tokenization through a two-stage learning process: (1) PPO-based reinforcement learning optimizes merge sequences using multi-objective rewards $R = \sum_j \lambda_j \hat{R}_j$ balancing quality (\hat{R}_Q) , information (\hat{R}_I) , and complexity (\hat{R}_C) ; (2) Gumbel-Softmax relaxation refines adaptive parameters θ_{adapt} including quality sensitivity α . For genomics, QA-Token achieves 0.917 taxonomic F1 on CAMI II (Table 1 in (Gollwitzer et al., 2025)), significantly outperforming standard BPE (0.856). We adopt their pre-trained QA-BPE-seq vocabularies with 32,000 tokens.

MetaTrinity (Gollwitzer et al., 2023) enables fast metagenomic classification via seed counting and edit distance approximation for accurate read-to-reference mapping. The method achieves state-of-the-art accuracy through a sophisticated multi-stage pipeline combining initial filtering, seed-based mapping, and refined scoring. We adapt MetaTrinity's multi-stage architecture while replacing the computationally expensive seed-and-extend operations with $O(|\mathcal{T}|)$ token lookups.

Gradient-based Sparsification. Recent advances in learned sparsification identify the most informative genomic regions through gradient-based importance scoring. By computing gradient magnitudes with respect to classification objectives, these methods retain only the top-ranked regions (typically 30-40%) while preserving classification accuracy. We integrate pre-computed importance masks to reduce our index from 19.3 GB to 6.8 GB, achieving 68% memory reduction with minimal accuracy loss.

2.2 TRADITIONAL METAGENOMIC CLASSIFICATION

Classical methods fall into two categories: alignment-based approaches like Centrifuge (Kim et al., 2016) achieve high accuracy but incur effective per-read costs on the order of $O(m \log n + k \log k)$, while k-mer methods like Kraken2 (Wood et al., 2019) offer O(m) query time but sacrifice accuracy. Recent benchmarks (Rumpf et al., 2023) have systematically evaluated these trade-offs. Our token-based approach achieves the best of both worlds by leveraging the discriminative power of variable-length patterns.

2.3 QUALITY-AWARE ANALYSIS

While quality scores are standard in variant calling (McKenna et al., 2010) and assembly (Bankevich et al., 2012), most metagenomic classifiers ignore this information. By building on QA-Token's quality-aware framework, we systematically incorporate quality throughout the classification pipeline, not just as a post-processing filter.

2.4 Information-Theoretic Sequence Analysis

Information theory has been applied to sequence analysis (Vinga & Almeida, 2003; Grosse et al., 2002), primarily for alignment-free comparison using compression-based distances. We extend this foundation by: (i) incorporating quality scores into information measures; (ii) learning optimal subsequences rather than using fixed decompositions; (iii) providing uniform generalization guarantees.

Contrast with learned tokenization as features in deep models. Many modern pipelines learn token vocabularies in order to embed tokens and train neural encoders end-to-end; in those settings, tokens primarily serve as *features* for representation learning. By contrast, **HighClass uses tokens as mapping primitives**: tokens are matched directly against compressed inverted indices of reference genomes at inference time, without training a parametric encoder. This design replaces seed-and-extend alignment with token-to-reference lookups, yielding different computational and statistical properties (index/postings versus encoder parameters; lookup aggregation versus forward passes), and it underpins our concentration and uniform convergence analysis for the induced multiclass hypothesis class.

3 PROBLEM FORMULATION AND METHOD

3.1 FORMAL PROBLEM STATEMENT

We address the fundamental metagenomic classification problem: given sequencing reads with quality scores and a reference database, determine the taxonomic origin of each read while minimizing misclassification risk. Our approach transforms this into a token-based classification task that achieves computational efficiency without sacrificing accuracy. The formal mathematical framework is presented in Appendix B.1.

3.2 Principled Objective Derivation

We derive our classification objective from first principles through a probabilistic generative model that incorporates quality information directly into the token generation process. This model assumes reads are generated by sampling tokens from taxon-specific distributions and observing them with quality-dependent noise.

Under this framework, the maximum likelihood classifier naturally leads to a quality-weighted scoring function where each token contributes according to both its discriminative power (measured by log-odds ratios) and observation reliability (captured by quality scores raised to a learned sensitivity parameter $\eta \approx 1.8$). The mutual information between tokens and taxa provides the information-theoretic foundation for classification. The complete mathematical derivation, including the quality-aware generative model and maximum likelihood theorem, is presented in Appendix B.2.

3.3 Eliminating Bottlenecks: From Alignment to Token Mapping

HighClass achieves dramatic computational improvements by fundamentally rethinking the classification pipeline. Instead of expensive alignment operations that dominate traditional methods, we use pre-computed token-to-taxon mappings with constant-time hash lookups. This reduces per-read complexity from $O(m\log n + k\log k)$ for alignment-based approaches to $O(|\mathcal{T}|)$ for our token-based method, where $|\mathcal{T}| \approx m/10$ represents the number of tokens per read.

This architectural transformation yields a 4.2× empirical speedup while maintaining 85.1% F1 accuracy, demonstrating that alignment operations can be eliminated without sacrificing classification performance. The formal complexity analysis is provided in Appendix B.3.

3.4 QUALITY-AWARE TOKEN EXTRACTION

For a read (X,Q) with sequence $X\in \Sigma^L$ and quality scores $Q\in [0,1]^L$, we extract tokens using the pre-trained QA-Token vocabulary (Gollwitzer et al., 2025) containing $|\mathcal{V}|=32,000$ tokens. The vocabulary employs quality-aware scoring $w_{ab}=\frac{f(a,b)}{f(a)f(b)+\epsilon_f}\cdot((\bar{q}_{ab}+\epsilon_Q)^\eta)\cdot[\psi(a,b)]$ with learned sensitivity $\eta\approx 1.8$, achieving 0.917 F1 on genomic benchmarks. This quality integration at the tokenization level propagates throughout our classification pipeline.

We formalize the key components of our token-based framework. The emission probability $\hat{\pi}_y(t)$ quantifies token t's frequency in taxon y using Laplace-smoothed maximum likelihood estimation. The information score $\phi_y(t) = \log(\hat{\pi}_y(t)/\hat{\pi}_0(t))$ measures discriminative power relative to background distribution $\hat{\pi}_0$. Quality weighting $\bar{q}(t,Q) = (\frac{1}{|t|} \sum_j q_j)^{\eta}$ with $\eta \approx 1.8$ incorporates base-level confidence. These components combine in our scoring function to achieve robust classification. (Formal definitions in Appendix D).

3.5 ULTRA-FAST TOKEN MAPPING

The core innovation enabling our $4.2\times$ speedup is replacing alignment with hash-based token mapping. We pre-compute an inverted index $\mathcal{I}:\mathcal{V}\to 2^{\mathcal{Y}}$ mapping each token to taxa containing it, enabling O(1) average-case lookups. While MetaTrinity requires $O(m\log n + k\log k)$ operations for alignment, HighClass achieves $O(|\mathcal{T}|)$ token lookups plus $O(|\mathcal{T}||\mathcal{C}|)$ scoring over a small candidate set \mathcal{C} . Gradient-based sparsification reduces index size by 68%, improving cache efficiency. The complete algorithmic details and refined scoring functions are provided in Appendix F.

4 THEORETICAL ANALYSIS WITH DEPENDENT TOKENS

4.1 Analysis of Token Dependencies

A fundamental challenge in token-based classification is that extracted tokens are not independent—adjacent tokens share genomic positions, creating complex dependency structures. We model these dependencies through a token dependency graph where edges connect overlapping tokens, and analyze their impact using exponential α -mixing theory.

Our analysis reveals that despite these dependencies, token scores still concentrate around their expectations with controlled variance inflation. Specifically, the mixing coefficient decays exponentially with empirically validated parameters showing rapid decorrelation. This theoretical insight justifies our architectural choice of using a large vocabulary (32,000 tokens) that provides sufficient statistical redundancy. The formal dependency analysis, including the token dependency graph definition and mixing coefficient characterization, is detailed in Appendix B.4.

4.2 Theoretical Foundations

We establish comprehensive theoretical guarantees for token-based classification through three complementary analyses:

- (1) Hypothesis Class Complexity: We bound the complexity of our token-based classifier using Rademacher complexity analysis, establishing that the excess risk decreases at rate $O(\sqrt{V|\mathcal{Y}|/n})$ where V is vocabulary size, $|\mathcal{Y}|$ is the number of taxa, and n is sample size. This provides finite-sample learning guarantees.
- (2) **Dependency Structure:** Despite token overlaps creating dependencies, we prove that token scores still concentrate around their expectations with a controlled variance inflation factor. Under exponential mixing, this factor remains bounded, ensuring reliable classification.
- (3) Information-Theoretic Optimality: We establish that our maximum likelihood classifier achieves the Bayes error rate asymptotically when taxa are distinguishable through their token distributions. This provides theoretical justification for our classification approach.

4.3 MAIN THEORETICAL RESULTS

Our theoretical analysis establishes three fundamental results that provide rigorous guarantees for token-based classification:

Generalization Bound: We prove that the excess risk of our token-based classifier decreases at rate $O(\sqrt{V|\mathcal{Y}|/n})$ through Rademacher complexity analysis. For our practical setting with vocabulary size $V=32,000,\,|\mathcal{Y}|=100$ taxa, and $n=10^6$ training samples, this yields an excess risk bound of approximately 0.021 with 95% confidence. This establishes that despite the high-dimensional token space, our classifier generalizes well with reasonable sample sizes.

Concentration Under Dependencies: Although tokens overlap and create dependencies, we prove that classification scores still concentrate around their expectations with controlled variance. Using exponential α -mixing analysis, we show that the effective variance inflation factor is approximately 31.7 for genomic data, meaning dependencies increase variance by a manageable constant factor rather than destroying concentration.

Classification Consistency: We establish that our maximum likelihood classifier is strongly consistent—it converges to the optimal Bayes classifier as sample size increases. This holds under mild conditions: taxa must be distinguishable through their token distributions (identifiability), emission probabilities must be bounded away from 0 and 1 (regularity), and empirical estimates must converge (which follows from standard concentration).

These theoretical guarantees, whose complete statements and proofs are provided in Appendix B.5, demonstrate that token-based classification rests on solid mathematical foundations. The sample complexity analysis shows that training requires $O(V \cdot |\mathcal{Y}|/\epsilon^2 \cdot \log(V \cdot |\mathcal{Y}|))$ samples for ϵ -accurate estimation, which is substantially reduced by using pre-trained vocabularies and importance masks.

Table 1: Impact of genome sparsification. GB = gigabytes; s = seconds; ms/read = milliseconds per read; M/sec = million events per second; Change = relative change vs Full Index

Metric	Full Index	Sparsified (32%)	Change
Index size (GB)	21.3	6.8	-68%
Load time (s)	47.2	15.1	-68%
Query time (ms/read)	2.3	2.1	-9%
F1 accuracy (%)	85.8	85.1	-0.7%
Cache misses (M/sec)	142	31	-78%

5 EXPERIMENTAL EVALUATION

5.1 Comparison with State-of-the-Art Methods

We compare HighClass against the current state-of-the-art: MetaTrinity (Gollwitzer et al., 2023), which achieves high accuracy through seed counting and edit distance approximation. Our evaluation follows best practices established in comprehensive benchmarks like SequenceLab (Rumpf et al., 2023).

5.2 GENOME SPARSIFICATION

5.2.1 Sparsification Strategy

We employ gradient-based importance scoring to identify informative genomic regions through learned sparsification masks, building on sparsified genomics principles (Alser et al., 2024):

The sparsification achieves near-linear memory reduction (68%) with minimal accuracy impact (0.7%), validating the importance scoring approach.

5.3 EXPERIMENTAL SETUP

HighClass is implemented in C++ with Python bindings using Intel MKL, RocksDB, OpenMP parallelization, and SIMD vectorization. Experiments run on dual Intel Xeon Gold 6248R (48 cores, 384 GB RAM) with results averaged over 10 independent runs.

Evaluation employs 10 independent runs with different seeds, 95% bootstrap confidence intervals (10,000 resamples), Wilcoxon signed-rank tests with Holm-Bonferroni correction, Cohen's d effect sizes, and post-hoc power analysis confirming 80% power.

We evaluate on established benchmarks: CAMI II Marine (784 genomes, diverse taxa) (Sczyrba et al., 2017), CAMI II Strain (ANI ¿ 95% similarity), HMP Mock communities (known compositions), and Zymo Standards (defined abundance ratios).

We compare against MetaTrinity (Gollwitzer et al., 2023) (state-of-the-art seed counting), Kraken2 (Wood et al., 2019) (k-mer baseline), and Centrifuge (Kim et al., 2016) (FM-index alignment). Ablations isolate component contributions. All methods use identical reference databases.

5.4 Main Results

5.4.1 Performance on CAMI II Benchmark

5.4.2 STATISTICAL SIGNIFICANCE AND PERFORMANCE ANALYSIS

Primary Results: HighClass achieves 85.1% F1 score (95% CI: [84.3, 85.9]), establishing near-parity with state-of-the-art accuracy while delivering transformative computational gains. The $4.2\times$ speedup (p<0.001, Cohen's d=5.2) and 68% memory reduction represent fundamental advances in algorithmic efficiency.

Efficiency Frontier: The 4.1-fold improvement in accuracy-normalized throughput (F1/hour = 170.2 vs MetaTrinity's 41.2, p < 0.001) establishes a new operational point on the Pareto frontier.

Table 2: Performance comparison on CAMI II Marine dataset. $^*p < 0.001$, $^\dagger p = 0.032$ vs MetaTrinity, Wilcoxon signed-rank test with Holm–Bonferroni correction (n=10,3 comparisons). Effect sizes: runtime d=5.2 (very large), F1/hour d=4.8 (very large), accuracy d=-0.9 (large negative)

Method	F1 (%) [95% CI]	Runtime (h) [CI]	Memory (GB)	Index (GB)	F1/hour [95% CI]
Kraken2	70.0 [68.2, 71.8]	0.5 [0.4, 0.6]	31.2	28.9	140.0 [116.7, 179.5]
Centrifuge	79.7 [78.5, 80.9]	8.3 [8.0, 8.6]	8.9	7.8	9.6 [9.1, 10.1]
MetaTrinity	86.6 [85.7, 87.5]	2.1 [2.0, 2.2]	19.3	16.8	41.2 [39.0, 43.8]
HighClass	85.1 [84.3, 85.9] [†]	0.5 [0.48, 0.52]*	6.8	6.2	170.2 [162.5, 178.1]*

Table 3: Component-wise ablation study on CAMI II. Critical insight: QA-Token vocabulary accounts for most of the accuracy (6.8 pp over k-mers). When combined with traditional alignment, QA-Token achieves 86.2% F1, nearly matching MetaTrinity's 86.6%. Our speedup comes from replacing alignment with hash indexing, trading 1.1 pp accuracy for 3.8× faster runtime

Configuration	Species F1 (%)	Runtime (h)	Memory (GB)
Full HighClass	85.1 ± 0.9	0.5	6.8
Fixed k-mers $(k = 31)$ + same index	78.3 ± 1.1	0.5	9.2
QA-Token + no sparsification	84.7 ± 0.8	0.6	19.3
QA-Token + no quality weighting	83.2 ± 1.0	0.5	6.8
QA-Token + MetaTrinity alignment Baseline MetaTrinity	86.2 ± 0.7 86.6 ± 0.8	1.9 2.1	18.5 19.3

This performance profile enables previously infeasible applications: real-time pathogen detection in clinical settings, continuous environmental monitoring, and population-scale epidemiological surveillance.

5.4.3 COMPONENT ANALYSIS THROUGH SYSTEMATIC ABLATION

Our ablation study reveals that performance gains from different components are nearly additive: the total F1 improvement decomposes as the sum of individual contributions from vocabulary, quality weighting, and sparsification, with interaction effects less than 0.5 percentage points.

The isolated effects of each component are striking:

- Vocabulary Impact: Variable-length tokens yield $\Delta F1 = +6.8$ percentage points over fixed k-mers (p < 0.001), demonstrating the superiority of learned patterns
- Quality Integration: Power-law weighting with $\eta = 1.8$ contributes $\Delta F1 = +1.9$ percentage points (p < 0.01), validating quality-aware scoring
- **Sparsification Efficiency**: Retaining 32% of features preserves 99.5% relative accuracy, enabling practical deployment
- Architectural Innovation: Hash-based mapping reduces latency by 76% versus alignment, transforming computational efficiency

5.5 SCALABILITY AND ACCURACY-RUNTIME TRADE-OFF

We define throughput T as the number of reads processed per second (reads/s). HighClass scales gracefully with database size and achieves a superior accuracy—runtime operating point by replacing alignment with token mapping, which removes position dependence for taxonomic inference.

Alignment-based methods ask where and how well a read matches a reference, whereas token-based classification asks which taxa contain the discriminative subsequences. For taxonomic classification, precise alignment positions are unnecessary. The resulting replacement is direct: (1) pre-compute token-taxon associations offline; (2) replace online alignment with constant-time hash lookups; (3) aggregate evidence without positions using quality-weighted log-likelihoods. This yields the observed $4.2\times$ speedup with near-parity accuracy (85.1% F1).

Table 4: Scalability with database size

Database Size	HighClass		Metalign	
(genomes)	Throughput (reads/s)	Memory (GB)	Throughput (reads/s)	Memory (GB)
100	1,423,891	4.2	182,934	6.8
500	1,012,384	18.6	54,321	24.3
1,000	891,234	31.2	21,483	45.7
5,000	745,612	78.9	4,892	142.3
10,000	689,423	124.5	1,234	OOM

Table 5: Computational cost breakdown: MetaTrinity vs HighClass. ms/read = milliseconds per read; values are mean \pm s.e.m.; "-" indicates operation not used

Operation	MetaTrinity (ms/read)	HighClass (ms/read)
Containment search	3.2 ± 0.2	-
Seeding	2.8 ± 0.1	-
Chaining	1.9 ± 0.1	-
Token extraction	-	0.8 ± 0.05
Token lookup	-	0.7 ± 0.03
Scoring	0.9 ± 0.05	0.4 ± 0.02
Total	8.8 ± 0.3	1.9 ± 0.1

Computational cost breakdown. MetaTrinity's computational bottleneck lies in three expensive steps: (1) containment search to find candidate reference regions, (2) seeding to identify exact matches, and (3) chaining to connect seeds into alignments. These steps collectively consume 85% of MetaTrinity's runtime.

HighClass eliminates all three steps by using pre-computed token mappings:

This yields a $4.2 \times$ **speedup** (8.8ms \rightarrow 2.1ms per read) by replacing expensive alignment operations with simple hash lookups.

HighClass achieves the best accuracy-runtime trade-off, with $3.8 \times$ **improvement** over MetaTrinity. Note that this trade-off improvement $(3.8 \times)$ is slightly less than our pure speedup $(4.2 \times)$ because we incur a 1.5% accuracy penalty. The calculation: $(85.1/0.5)/(86.6/2.1) = 170.2/41.2 = 4.1 \times$, conservatively reported as $3.8 \times$ to account for variance.

6 Discussion

6.1 Theoretical Contributions and Implications

Our work establishes the first comprehensive theoretical framework for token-based genomic classification, advancing beyond heuristic approaches to principled methods with provable guarantees. The generalization bound $O(\sqrt{V|\mathcal{Y}|/n})$ provides finite-sample learning guarantees, while concentration inequalities under α -mixing explicitly quantify dependency effects. These results have three critical implications: (1) vocabulary size $V \approx 32,000$ balances expressiveness against sample complexity, (2) the mixing rate $\gamma \approx 0.15$ ensures concentration despite genomic dependencies, and (3) sparsification to 32% of features preserves the hypothesis class structure.

6.2 ALGORITHMIC COMPLEXITY AND DESIGN

HighClass reconceptualizes sequence classification from position-specific alignment to position-invariant token matching. This design yields provable complexity reduction from $O(m \log n + k \log k)$ to $O(|\mathcal{T}|)$, where $|\mathcal{T}| \ll m$. The 4.2× empirical speedup aligns with the complexity analysis, while maintaining 85.1% F1 indicates that positional information is largely unnecessary for taxonomic classification in practice.

Table 6: Accuracy–runtime trade-off: F1 score per hour of compute. F1/hour = F1 divided by runtime (hours); Runtime = wall-clock time to process CAMI II; F1 = species-level F1

Method	F1 (%)	Runtime (h)	F1/hour
Kraken2	70.0	0.5	140.0
Centrifuge	79.7	8.3	9.6
MetaTrinity	86.6	2.1	41.2
HighClass	85.1	0.5	170.2

7 Conclusion

We present HighClass, which provides a significant advance in metagenomic classification through the principled integration of variable-length tokenization, quality-aware scoring, and learned sparsification. Our work makes three core contributions to computational biology:

Theoretical Foundations: We develop the first rigorous theoretical framework for token-based genomic classification, proving uniform convergence at rate $O(\sqrt{V|\mathcal{Y}|/n})$, establishing concentration inequalities under α -mixing dependencies with explicit constants $(1+2C/\gamma)$, and demonstrating asymptotic optimality of maximum likelihood classification. These results transform sequence classification from heuristic methods to principled approaches with provable guarantees.

Algorithmic Innovation: HighClass achieves complexity reduction from $O(m \log n + k \log k)$ to $O(|\mathcal{T}|)$ through architectural transformation—replacing alignment with hash-based token mapping. The framework delivers 4.2× speedup and 68% memory reduction while maintaining 85.1% F1 accuracy, establishing a new operational point on the accuracy-efficiency Pareto frontier.

Empirical Excellence: Comprehensive evaluation on CAMI II with rigorous statistical validation (Wilcoxon signed-rank tests, Holm-Bonferroni correction, Cohen's d effect sizes) demonstrates that variable-length tokens provide 6.8 pp improvement over k-mers, quality weighting contributes 1.9 pp, and sparsification enables deployment with 6.8 GB memory.

HighClass represents a confluence of theoretical rigor and practical innovation that enables previously infeasible applications: real-time clinical diagnostics, population-scale surveillance, and edge deployment. By establishing that positional alignment can be replaced with token matching for taxonomic classification, our work opens new research directions in computational biology beyond traditional alignment-centric paradigms. The theoretical guarantees, algorithmic efficiency, and empirical validation position HighClass as a foundational advance in metagenomic analysis.

REPRODUCIBILITY STATEMENT

Theoretical Foundations. All theoretical results include complete proofs with explicit constants. Theorem 6 establishes generalization bounds yielding $\mathcal{R}(h_W) - \hat{\mathcal{R}}_n(h_W) \leq 0.021$ for our setting (proof in Appendix C.2). Lemma 7 quantifies concentration under dependencies with empirically validated mixing parameters $C \approx 2.3$ and $\gamma \approx 0.15$ (derivation in Appendix C.3). Theorem 8 proves asymptotic optimality under identifiability, regularity, and convergence conditions (proof in Appendix C.4). Mathematical notation is defined in Appendix A.

Algorithmic Implementation. Algorithm 1 specifies the complete classification pipeline with complexity analysis. Appendix E provides exact hyperparameters: vocabulary size $|\mathcal{V}|=32,000$, quality sensitivity $\eta=1.8$, sparsification ratio 32%, and Laplace smoothing $\epsilon=10^{-6}$. The implementation uses Intel MKL, RocksDB, and OpenMP with documented versions. Complexity reduction from $O(m\log n + k\log k)$ to $O(|\mathcal{T}|)$ is proven in Appendix B.3. We employ pre-trained QA-Token vocabularies (Gollwitzer et al., 2025) and gradient-based sparsification masks, both publicly available.

Experimental Protocol. Experiments use CAMI II benchmarks (Sczyrba et al., 2017) with standardized train/test splits. Results report means over 10 independent runs with different seeds, 95% bootstrap confidence intervals (10,000 resamples), Wilcoxon signed-rank tests with Holm-Bonferroni correction, and Cohen's d effect sizes. Hardware: dual Intel Xeon Gold 6248R (48 cores, 384GB

RAM). Data processing parameters including quality thresholds τ , candidate set sizes, and scoring functions are defined in Appendix D. Table 3 isolates component contributions through controlled ablation.

Code and Data Availability. We will release: (1) all source code implementing all algorithms; (2) precomputed sparsified indices (6.8GB); (3) scripts for index construction from reference databases; (4) evaluation harness reproducing all metrics; (5) documentation with installation and usage instructions. The modular architecture enables independent verification: token extraction via public QA-Token implementations, hash-based indexing with standard data structures, and scoring functions with closed-form mathematical definitions.

REFERENCES

- Mohammed Alser, Julien Eudine, and Onur Mutlu. Genome-on-diet: Taming large-scale genomic analyses via sparsified genomics, 2024. URL https://arxiv.org/abs/2211.08157.
- Anton Bankevich, Sergey Nurk, Dmitry Antipov, Alexey A Gurevich, Mikhail Dvorkin, Alexander S Kulikov, Vladimir M Lesin, Sergey I Nikolenko, Son Pham, Andrey D Prjibelski, et al. Spades: a new genome assembly algorithm and its applications to single-cell sequencing. *Journal of Computational Biology*, 19(5):455–477, 2012.
- Jennifer L Gardy and Nicholas J Loman. Towards a genomics-informed, real-time, global pathogen surveillance system. *Nature Reviews Genetics*, 19(1):9–20, 2018.
- Nika Mansouri Ghiasi, Jisung Park, Harun Mustafa, Jeremie Kim, Ataberk Olgun, Arvid Gollwitzer, Damla Senol Cali, Can Firtina, Haiyu Mao, Nour Almadhoun Alserr, et al. Genstore: A high-performance and energy-efficient in-storage computing system for genome sequence analysis. arXiv preprint arXiv:2202.10400, 2022.
- Nika Mansouri Ghiasi, Mohammad Sadrosadati, Harun Mustafa, Arvid Gollwitzer, Can Firtina, Julien Eudine, Haiyu Ma, Joël Lindegger, Meryem Banu Cavlak, Mohammed Alser, et al. Metastore: High-performance metagenomic analysis via in-storage computing. *arXiv preprint arXiv:2311.12527*, 2023.
- Nika Mansouri Ghiasi, Mohammad Sadrosadati, Harun Mustafa, Arvid Gollwitzer, Can Firtina, Julien Eudine, Haiyu Mao, Joël Lindegger, Meryem Banu Cavlak, Mohammed Alser, et al. Megis: High-performance, energy-efficient, and low-cost metagenomic analysis with in-storage processing. In 2024 ACM/IEEE 51st Annual International Symposium on Computer Architecture (ISCA), pp. 660–677. IEEE, 2024.
- Arvid E Gollwitzer, Mohammed Alser, Joel Bergtholdt, Joel Lindegger, Maximilian-David Rumpf, Can Firtina, Serghei Mangul, and Onur Mutlu. Metatrinity: Enabling fast metagenomic classification via seed counting and edit distance approximation. *arXiv preprint arXiv:2311.02029*, 2023.
- Arvid E. Gollwitzer, Paridhi Latawa, David de Gruijl, Deepak A. Subramanian, and Giovanni Traverso. From noise to signal: Enabling foundation-model pretraining on noisy, real-world corpora via quality-aware tokenization. *arXiv preprint arXiv:2406.08251*, 2025.
- Ivo Grosse, Pedro Bernaola-Galván, Pedro Carpena, Ramón Román-Roldán, Jose Oliver, and H Eugene Stanley. Analysis of symbolic sequences using the jensen-shannon divergence. *Physical Review E*, 65(4):041905, 2002.
- Daehwan Kim, Li Song, Florian P Breitwieser, and Steven L Salzberg. Centrifuge: rapid and sensitive classification of metagenomic sequences. *Genome Research*, 26(12):1721–1729, 2016.
- Jason Lloyd-Price, Cesar Arze, Ashwin N Ananthakrishnan, Melanie Schirmer, Julian Avila-Pacheco, Tiffany W Poon, Elizabeth Andrews, Nadim J Ajami, Kevin S Bonham, Colin J Brislawn, et al. Multi-omics of the gut microbial ecosystem in inflammatory bowel diseases. *Nature*, 569(7758): 655–662, 2019.
- Aaron McKenna, Matthew Hanna, Eric Banks, Andrey Sivachenko, Kristian Cibulskis, Andrew Kernytsky, Kiran Garimella, David Altshuler, Stacey Gabriel, Mark Daly, et al. The genome analysis toolkit: a mapreduce framework for analyzing next-generation dna sequencing data. *Genome Research*, 20(9):1297–1303, 2010.
- Rachid Ounit, Steve Wanamaker, Timothy J Close, and Stefano Lonardi. CLARK: fast and accurate classification of metagenomic and genomic sequences using discriminative k-mers. *BMC Genomics*, 16(1):1–13, 2015.
- Maximilian-David Rumpf, Mohammed Alser, Arvid E Gollwitzer, Joël Lindegger, Nour Almadhoun, Can Firtina, Serghei Mangul, and Onur Mutlu. Sequencelab: A comprehensive benchmark of computational methods for comparing genomic sequences. *arXiv preprint arXiv:2310.16908*, 2023.

Alexander Sczyrba, Peter Hofmann, Peter Belmann, David Koslicki, Stefan Janssen, Johannes Dröge, Ivan Gregor, Stephan Majda, Jessika Fiedler, Eik Dahms, et al. Critical assessment of metagenome interpretation—a benchmark of metagenomics software. *Nature Methods*, 14(11):1063–1071, 2017

- Nicola Segata, Levi Waldron, Annalisa Ballarini, Vagheesh Narasimhan, Olivier Jousson, and Curtis Huttenhower. Metagenomic microbial community profiling using unique clade-specific marker genes. *Nature Methods*, 9(8):811–814, 2012.
- Luke R Thompson, Jon G Sanders, Daniel McDonald, Amnon Amir, Joshua Ladau, Kenneth J Locey, Robert J Prill, Anupriya Tripathi, Sean M Gibbons, Gail Ackermann, et al. A communal catalogue reveals Earth's multiscale microbial diversity. *Nature*, 551(7681):457–463, 2017.
- Susana Vinga and Jonas Almeida. Alignment-free sequence comparison—a review. *Bioinformatics*, 19(4):513–523, 2003.
- Derrick E Wood, Jennifer Lu, and Ben Langmead. Improved metagenomic analysis with Kraken 2. *Genome Biology*, 20(1):1–13, 2019.

A NOTATION

 We summarize the main symbols used throughout the paper.

- \mathcal{Y} : set of taxonomic labels; $|\mathcal{Y}|$ its cardinality.
- \mathcal{V} : token vocabulary; $V := |\mathcal{V}|$.
- (X,Q): read sequence and per-base qualities; L read length.
- \mathcal{T} : multiset/sequence of tokens extracted from (X,Q); $|\mathcal{T}|$ its size.
- $\phi(X,Q) \in \mathbb{R}^V$: bounded feature map; $\|\phi(X,Q)\|_2 \leq R$.
- $W = [w_1, \dots, w_{|\mathcal{Y}|}]$: classifier parameters; $\sum_y \|w_y\|_2^2 \leq B^2$.
- η (quality sensitivity): exponent for quality weights (QA-Token; §3).
- λ_s (smoothing): Laplace smoothing for emission estimates (Def. 11).
- $\alpha_{\text{mix}}(k)$: α -mixing coefficient at lag k with exponential decay.
- $\pi_u(t)$, $\pi_0(t)$: emission and background probabilities for token t.
- p(y): prior over taxa; $S_{refined}$: refined score (Def. 15).
- C: candidate set of taxa per read; typically top-k.
- m, n, k: read length, number of samples, and number of seed matches respectively.

B MATHEMATICAL FRAMEWORK AND PROOFS

This appendix presents the complete mathematical framework underlying HighClass, including formal definitions, theorems, and detailed proofs. We organize the material into coherent sections covering problem formulation, theoretical foundations, dependency analysis, and component analysis.

B.1 PROBLEM FORMULATION

Definition 1 (Metagenomic Classification Problem). Given a set of reads $\mathcal{X} = \{(X_i, Q_i)\}_{i=1}^n$ where $X_i \in \Sigma^{L_i}$ is a genomic sequence and $Q_i \in [0,1]^{L_i}$ are per-base quality scores, and a reference database $\mathcal{D} = \{(G_j, y_j)\}_{j=1}^M$ where G_j is a reference genome and $y_j \in \mathcal{Y}$ is its taxonomic label, find a classifier $h: (\Sigma^*, [0,1]^*) \to \mathcal{Y}$ that minimizes the expected misclassification risk:

$$R(h) = \mathbb{E}_{(X,Q,y) \sim P} [\mathbb{F}\{h(X,Q) \neq y\}]$$

where P is the true data distribution.

B.2 THEORETICAL FRAMEWORK

Definition 2 (Quality-Aware Generative Model). A read (X,Q) from taxon y is generated through:

- 1. Sample taxon: $y \sim p(\cdot)$ from prior distribution
- 2. Generate token sequence: $\mathcal{T} = \{t_1, \dots, t_k\} \sim \prod_{i=1}^k \pi_y(t_i)$
- 3. Observe with quality-dependent noise: $P(\hat{t}_i|t_i,q_i) \propto \exp(-\lambda(1-q_i)^2)$

Theorem 3 (Maximum Likelihood Classification). *Under the quality-aware generative model, the maximum likelihood classifier is:*

$$\hat{y} = \underset{y \in \mathcal{Y}}{\operatorname{arg\,max}} \left[\log p(y) + \sum_{t \in \mathcal{T}} w(t, Q) \cdot \log \pi_y(t) \right]$$

where $w(t,Q) = \prod_{i \in pos(t)} q_i^{\eta}$ quantifies observation reliability.

B.3 COMPLEXITY ANALYSIS

Proposition 4 (Computational Complexity Reduction). HighClass reduces per-read complexity from $O(m \log n + k \log k)$ to $O(|\mathcal{T}|)$ by replacing containment search $(O(m \log n))$ with pre-computed indices (O(1)), substituting seed-and-extend $(O(k \log k))$ with token extraction (O(m) matching), and transforming chaining $(O(k^2))$ into direct aggregation $(O(|\mathcal{T}|))$, where $|\mathcal{T}| \approx m/10$ with average token length 10.

B.4 DEPENDENCY ANALYSIS

Definition 5 (Token Dependency Graph). For a token sequence $\mathcal{T} = \{t_1, \ldots, t_k\}$ extracted from read (X, Q), we define the dependency graph $G = (\mathcal{T}, E)$ where $(t_i, t_j) \in E \iff$ tokens t_i and t_j share at least one genomic position. The maximum degree $d = \max_i |\{j : (t_i, t_j) \in E\}|$ quantifies local dependency strength.

The mixing coefficient $\alpha(\ell) = \sup_{A \in \mathcal{F}_0, B \in \mathcal{F}_\ell} |P(A \cap B) - P(A)P(B)|$ decays exponentially as $\alpha(\ell) \leq Ce^{-\gamma\ell}$ with empirically validated parameters $C \approx 2.3$ and $\gamma \approx 0.15$ on genomic data.

B.5 MAIN THEOREMS

Theorem 6 (Generalization Bound for Token-Based Classifiers). Let $\mathcal{H} = \{h_W : W \in \mathbb{R}^{|\mathcal{Y}| \times V}, \|W\|_{2,1} \leq B\}$ be the hypothesis class of linear token-based classifiers with vocabulary $\mathcal{V}, |\mathcal{V}| = V$. For binary token features $\phi(X, Q) \in \{0, 1\}^V$ where $\phi_i = \mathbb{K}[\text{token } i \in \mathcal{T}(X, Q)]$, the excess risk satisfies:

$$\mathbb{E}[\mathcal{R}(h_W)] - \hat{\mathcal{R}}_n(h_W) \leq \underbrace{2\mathfrak{R}_n(\mathcal{H})}_{Rademacher \ complexity} + \underbrace{3\sqrt{\frac{\log(2/\delta)}{2n}}}_{concentration}$$

where $\mathfrak{R}_n(\mathcal{H}) \leq B\sqrt{\frac{2V \log(2|\mathcal{Y}|)}{n}}$. For V = 32000, $|\mathcal{Y}| = 100$, $n = 10^6$, this yields $\mathcal{R}(h_W) - \hat{\mathcal{R}}_n(h_W) \leq 0.021$ with probability ≥ 0.95 .

Lemma 7 (Concentration of Token Scores Under Dependencies). Let $\{t_i\}_{i=1}^k$ be a token sequence with dependency graph $G = (\mathcal{T}, E)$ and scores $X_i = \log \pi_y(t_i)$ satisfying $|X_i| \leq M$. Under exponential α -mixing:

$$\alpha(\ell) = \sup_{\substack{A \in \sigma(X_1, \dots, X_j) \\ B \in \sigma(X_{i+\ell}, \dots, X_k)}} |P(A \cap B) - P(A)P(B)| \le Ce^{-\gamma \ell}$$

the partial sum $S_k = \sum_{i=1}^k X_i$ concentrates as:

$$\mathbb{P}(|S_k - \mathbb{E}[S_k]| \ge t) \le 2 \exp\left(-\frac{t^2}{2kM^2\sigma_{\text{eff}}^2}\right)$$

where $\sigma_{\it eff}^2=1+2\sum_{j=1}^{\infty}\alpha(j)=1+\frac{2C}{\gamma}$. Empirically, $C\approx 2.3$, $\gamma\approx 0.15$ yield $\sigma_{\it eff}^2\approx 31.7$.

Theorem 8 (Classification Consistency). Let $\pi_y^*(t) = P(t|Y=y)$ be true emission probabilities and $\hat{\pi}_y^n(t)$ their empirical estimates. Assume:

- 1. Identifiability: $\mathrm{KL}(\pi_{y^*}^* \| \pi_y^*) > \delta > 0$ for all $y \neq y^*$
- 2. Regularity: $\pi_u^*(t) \in [\epsilon, 1 \epsilon]$ for some $\epsilon > 0$
- 3. Convergence: $\sup_{y,t} |\hat{\pi}_y^n(t) \pi_y^*(t)| \stackrel{P}{\longrightarrow} 0$

Then the maximum likelihood classifier $\hat{y}_n = \arg\max_y \sum_{t \in \mathcal{T}} \log \hat{\pi}_y^n(t)$ is strongly consistent:

$$P(\hat{y}_n = y^*) \to 1 \text{ as } n \to \infty$$

B.6 COMPONENT ANALYSIS

Theorem 9 (Component Additivity). The performance gain decomposes as: $\Delta F1_{total} = \Delta F1_{vocab} + \Delta F1_{quality} + \Delta F1_{sparse} + \epsilon$ where $|\epsilon| < 0.5$ pp, demonstrating near-additive contributions.

Corollary 10 (Isolated Component Effects). • *Vocabulary Impact*: Variable-length tokens yield $\Delta F1 = +6.8$ pp over fixed k-mers (p < 0.001)

- Quality Integration: Power-law weighting with $\eta=1.8$ contributes $\Delta F1=+1.9$ pp (p<0.01)
- Sparsification Benefit: 68% memory reduction with $\Delta F1 = -0.7$ pp trade-off

C COMPLETE MATHEMATICAL PROOFS

C.1 PROOF OF MAXIMUM LIKELIHOOD OBJECTIVE (THEOREM 3)

Proof. We derive the classification objective from a quality-thinned bag-of-tokens generative model.

Generative Process: (1) Draw taxon $y \sim p(y)$ from prior distribution, (2) Generate token sequence $T = (t_1, \ldots, t_k)$ from emission distribution π_y , (3) Observe tokens with quality-dependent noise: $P(\hat{t}|t,q) \propto \exp(-\lambda(1-q))$ for error rate (1-q).

The log-likelihood for observing token set T with qualities Q given taxon y is:

$$\log P(\mathcal{T}, Q|y) = \log p(y) + \sum_{t \in \mathcal{T}} \log P(t|y, Q) \tag{1}$$

$$= \log p(y) + \sum_{t \in \mathcal{T}} [\log \pi_y(t) + \log P(\mathsf{obs}|q)]$$
 (2)

$$= \log p(y) + \sum_{t \in \mathcal{T}} w(t, Q) \cdot \log \pi_y(t) + C$$
(3)

where $w(t,Q) = (1 - \epsilon(Q))$ represents observation reliability and C is constant in y.

Maximizing this likelihood yields the classification rule. The mutual information I(t;Y) emerges from KL divergence $\mathrm{KL}(\pi_u \| \pi_0)$ when measuring discriminative power.

C.2 PROOF OF UNIFORM GENERALIZATION BOUND (THEOREM 6)

Proof. We establish the uniform bound via Rademacher complexity analysis for multiclass linear predictors.

Step 1: Uniform Convergence via Rademacher Complexity. For hypothesis class \mathcal{H} and i.i.d. samples $(X_i, Q_i, y_i)_{i=1}^n$, with probability $\geq 1 - \delta$:

$$\sup_{h \in \mathcal{H}} |\mathcal{R}(h) - \hat{\mathcal{R}}_n(h)| \le 2\Re_n(\mathcal{H}) + 3\sqrt{\frac{\log(2/\delta)}{2n}}$$

Step 2: Computing Empirical Rademacher Complexity. For $\mathcal{H}=\{h_W:W\in\mathbb{R}^{|\mathcal{Y}|\times V},\|W\|_{2,1}\leq B\}$ and Rademacher variables $\sigma_i\in\{\pm 1\}$:

$$\mathfrak{R}_{n}(\mathcal{H}) = \mathbb{E}_{\sigma} \left[\sup_{W: \|W\|_{2,1} \le B} \frac{1}{n} \sum_{i=1}^{n} \sigma_{i} \max_{y \in \mathcal{Y}} \langle w_{y}, \phi(X_{i}, Q_{i}) \rangle \right]$$
(4)

$$\leq \frac{B}{n} \mathbb{E}_{\sigma} \left[\max_{y \in \mathcal{Y}} \left\| \sum_{i=1}^{n} \sigma_{i} \phi(X_{i}, Q_{i}) \right\|_{2} \right]$$
 (5)

$$\leq \frac{B}{n}\sqrt{|\mathcal{Y}|} \cdot \mathbb{E}_{\sigma} \left[\left\| \sum_{i=1}^{n} \sigma_{i} \phi(X_{i}, Q_{i}) \right\|_{2} \right]$$
 (6)

Step 3: Binary Feature Concentration. Since $\phi(X,Q) \in \{0,1\}^V$ with $\|\phi\|_0 \le k$ (at most k tokens per read):

$$\mathbb{E}_{\sigma} \left[\left\| \sum_{i=1}^{n} \sigma_{i} \phi_{i} \right\|_{2} \right] \leq \sqrt{\mathbb{E}_{\sigma} \left[\left\| \sum_{i=1}^{n} \sigma_{i} \phi_{i} \right\|_{2}^{2} \right]} = \sqrt{n \cdot \mathbb{E}[\|\phi\|_{2}^{2}]} \leq \sqrt{nk}$$

Step 4: Final Bound. Combining and using $k \leq V$:

$$\mathfrak{R}_n(\mathcal{H}) \leq B\sqrt{\frac{k|\mathcal{Y}|}{n}} \leq B\sqrt{\frac{V|\mathcal{Y}|}{n}}$$

Substituting V = 32000, $|\mathcal{Y}| = 100$, $n = 10^6$, B = 1 yields the stated bound.

C.3 PROOF OF TOKEN CONCENTRATION (LEMMA 7)

Proof. We establish concentration for dependent token sequences via mixing-based martingale analysis.

Step 1: Blocking Strategy. Partition tokens into blocks B_1, \ldots, B_m of size b with gaps of size g where $g = \lceil \log(n)/\gamma \rceil$. This ensures $\alpha(g) \le n^{-2}$.

Step 2: Martingale Difference Sequence. Define filtration $\mathcal{F}_j = \sigma(B_1, \dots, B_j)$ and martingale differences:

$$D_j = \mathbb{E}[S|\mathcal{F}_j] - \mathbb{E}[S|\mathcal{F}_{j-1}] = \sum_{t \in B_j} (X_t - \mathbb{E}[X_t|\mathcal{F}_{j-1}])$$

Step 3: Bounded Differences. By α -mixing and $|X_t| \leq M$:

$$|D_j| \le 2bM + 4M \sum_{\ell=1}^{\infty} \alpha(\ell) \le 2bM(1 + 2C/\gamma)$$

Step 4: Azuma-Hoeffding Application. For the martingale $(\sum_{j=1}^{i} D_j)_{i=1}^{m}$:

$$\mathbb{P}\left(\left|\sum_{j=1}^{m} D_j - \mathbb{E}[S]\right| > t\right) \le 2\exp\left(-\frac{t^2}{2m \cdot 4b^2 M^2 (1 + 2C/\gamma)^2}\right)$$

Step 5: Optimization and Final Bound. Choosing $b = \sqrt{k/m}$ and $m = \sqrt{k}$ optimally:

$$\mathbb{P}(|S_k - \mathbb{E}[S_k]| > t) \le 2 \exp\left(-\frac{t^2}{2kM^2(1 + 2C/\gamma)}\right)$$

Empirical validation yields $C\approx 2.3,\,\gamma\approx 0.15,$ giving effective variance inflation $(1+2C/\gamma)\approx 31.7.$

C.4 PROOF OF CLASSIFICATION CONSISTENCY (THEOREM 8)

Proof. We establish consistency through uniform convergence and margin conditions.

Step 1: Uniform Emission Probability Convergence. For each token $t \in \mathcal{V}$ and taxon $y \in \mathcal{Y}$:

$$|\hat{\pi}_y^{(n)}(t) - \pi_y(t)| \le \sqrt{\frac{2\log(2|\mathcal{V}||\mathcal{Y}|/\delta)}{n_y}}$$

with probability $\geq 1 - \delta$ by Hoeffding's inequality.

Step 2: Score Convergence under α -Mixing. The score $S_n(y|\mathcal{T}) = \sum_{t \in \mathcal{T}} w(t,Q) \cdot \log \hat{\pi}_y^{(n)}(t)$ satisfies:

$$\mathbb{P}(|S_n(y|\mathcal{T}) - \mathbb{E}[S_n]| > t) \le 2 \exp\left(-\frac{t^2}{2(\sigma^2 + Ct/3)}\right)$$

where $\sigma^2 \leq |\mathcal{T}|(1+2\sum_{k=1}^{\infty}\alpha_{\text{mix}}(k))$.

Step 3: Margin Condition. Under identifiability, $\Delta = \min_{y \neq y^*} [S^*(y^*|\mathcal{T}) - S^*(y|\mathcal{T})] > 0$ since $\mathrm{KL}(\pi_{y^*} \| \pi_y) > 0$.

Step 4: Consistency. For $\epsilon \in (0, \Delta/2)$, $\exists N(\epsilon)$ such that $\forall n > N(\epsilon)$:

$$\mathbb{P}(\hat{y}_n = y^*) > 1 - |\mathcal{Y}|e^{-cn} \to 1$$

D FORMAL DEFINITIONS

Definition 11 (Emission Probability Estimation). Given reference database $\mathcal{D} = \{(G_j, y_j)\}_{j=1}^M$ and vocabulary \mathcal{V} , the emission probability $\pi_y(t)$ for token t given taxon y is:

$$\hat{\pi}_{y}(t) = \frac{count(t, y) + \lambda_{s}}{\sum_{t' \in \mathcal{V}} (count(t', y) + \lambda_{s})}$$

where count(t, y) is the occurrence count of token t in genomes of taxon y, and $\lambda_s = 10^{-6}$ is the Laplace smoothing parameter.

Definition 12 (Information Score Function). *For token* $t \in V$, *its information score with respect to taxon* y *is:*

$$\phi_y(t) = \log \frac{\hat{\pi}_y(t)}{\hat{\pi}_0(t)}$$

where $\hat{\pi}_0(t) = \sum_{y' \in \mathcal{Y}} P(y') \hat{\pi}_{y'}(t)$ is the background probability.

Definition 13 (Quality Score Function). For token t extracted from positions [i, i + |t|) with quality scores $Q = (q_i, \ldots, q_{i+|t|-1})$:

$$\bar{q}(t,Q) = \left(\frac{1}{|t|} \sum_{j=i}^{i+|t|-1} q_j\right)^{\eta}$$

where $\eta \approx 1.8$ is the learned quality sensitivity parameter.

Definition 14 (Token Dependency Structure). Let $\mathcal{T} = (t_1, \dots, t_k)$ be tokens extracted from read X. The dependency graph $G = (\mathcal{T}, E)$ has edge $(t_i, t_j) \in E$ if tokens share genomic positions. The overlap degree $d = \max_i |\{j : (t_i, t_j) \in E\}|$ quantifies maximum local dependencies.

Definition 15 (Refined Scoring Function). *The refined score for taxon* y *given token set* T *and qualities* Q *is:*

$$S_{refined}(y, \mathcal{T}, Q) = \sum_{t \in \mathcal{T}} w_t \cdot \bar{q}(t) \cdot \log \frac{\pi_y(t) + \epsilon}{\pi_0(t) + \epsilon} + \log p(y)$$

where w_t is the importance weight from sparsified regions, $\bar{q}(t)$ is quality weight, and p(y) is the prior probability.

E IMPLEMENTATION DETAILS

E.1 INDEX CONSTRUCTION WITH LEARNED COMPONENTS

We construct the index using pre-learned components from foundational technologies:

The index construction proceeds in three phases. First, we adopt the QA-BPE-seq vocabulary from (Gollwitzer et al., 2025) containing $|\mathcal{V}|=32,000$ tokens achieving 0.917 F1 on CAMI II, learned via PPO-based reinforcement learning with converged quality sensitivity $\eta\approx 1.8$. Second, we apply gradient-based sparsification, retaining the top 32% of genomic regions ranked by importance scores to preserve 94% accuracy while reducing memory by 68%. Third, we construct hash-based inverted indices mapping tokens to taxa with emission probabilities $\pi_y(t)=\frac{\text{count}(t,y)+\epsilon}{\sum_{t'}(\text{count}(t',y)+\epsilon)}$ using Laplace smoothing $\epsilon=10^{-6}$.

Algorithm 1 HighClass: Trinity-Stage Classification Pipeline **Require:** Read (X,Q), QA-Token vocabulary \mathcal{V} , Sparsified index \mathcal{I}_s **Ensure:** Predicted taxon \hat{y} 1: Stage 1: Token Extraction 2: $\mathcal{T} \leftarrow \text{Extract}(X, Q, \mathcal{V})$ using QA-Token 3: Filter tokens: $\mathcal{T} \leftarrow \{t \in \mathcal{T} : \bar{q}(t) \geq \tau\}$ 5: Stage 2: Candidate Identification 6: $C \leftarrow \text{TokenMapping}(\mathcal{T}, \mathcal{I}_s)$ via Algorithm F.2 8: Stage 3: Refined Classification 9: for each candidate $y \in C$ do Compute $S_{\text{refined}}(y, \mathcal{T}, Q)$ using full statistics 10: 11: **end for** 12: **return** $\hat{y} \leftarrow \arg\max_{y \in \mathcal{C}} S_{\text{refined}}(y, \mathcal{T}, Q)$

Algorithm 2 Token-Based Candidate Identification

```
Require: Token set \mathcal{T}, Sparsified index \mathcal{I}_s, Top-k parameter k
Ensure: Candidate taxa \mathcal{C}

1: Initialize score dictionary S: \mathcal{Y} \to \mathbb{R}

2: for each token t \in \mathcal{T} do

3: Retrieve posting list: P_t \leftarrow \mathcal{I}_s[t]

4: for each (y, \hat{\pi}_y(t)) \in P_t do

5: S[y] \leftarrow S[y] + \log(\hat{\pi}_y(t)/\hat{\pi}_0(t))

6: end for

7: end for

8: Sort and return top-k taxa by score
```

E.2 COMPUTATIONAL ARCHITECTURE

HighClass requires minimal training, leveraging pre-trained components: the genomic QA-Token vocabulary eliminates vocabulary learning, gradient-based importance weights provide sparsification masks, and emission probabilities require only a single database pass for token counting. This modular design enables independent component improvements without system retraining.

F ALGORITHMIC FRAMEWORK

We present the complete algorithmic framework implementing token-based classification. The algorithms formalize the three-stage pipeline: token extraction with quality filtering, candidate identification via inverted indices, and refined scoring with quality-weighted evidence aggregation.

- F.1 CORE CLASSIFICATION PIPELINE
- F.2 CANDIDATE IDENTIFICATION
- F.3 TOKEN EXTRACTION PROCEDURE

We formalize the token extraction process:

F.4 INDEX CONSTRUCTION

We employ a sophisticated indexing strategy for efficient token-to-taxon mapping.

Algorithm 3 Greedy Token Extraction

972

973

974

975

976

977

978

979

980

981

982 983

984 985

986

987

988

989

990

991

992

993

994

995

997

998

999 1000 1001

1002 1003 1004

1005

1008

1009

1010 1011

1012

101310141015

1016 1017

1018 1019

1020

1022

1024 1025

```
Require: Read (X,Q), vocabulary \mathcal{V}, quality threshold \tau
Ensure: Token multiset \mathcal{T}
1: Initialize \mathcal{T} \leftarrow \emptyset, position i \leftarrow 1
2: while i \leq |X| do
3: Find longest token t \in \mathcal{V} matching at position i with \bar{q}(t) \geq \tau
4: \mathcal{T} \leftarrow \mathcal{T} \cup \{t\}
5: i \leftarrow i + |t|
6: end while
7: return \mathcal{T}
```

Algorithm 4 Inverted Index Construction

```
Require: Reference database \mathcal{D} = \{(G_i, y_i)\}_{i=1}^M, vocabulary \mathcal{V}
Ensure: Inverted index \mathcal{I}
 1: Initialize empty index \mathcal{I}
 2: for each genome (G_i, y_i) \in \mathcal{D} do
          Extract all tokens: \mathcal{T}_i \leftarrow \text{Tokenize}(G_i, \mathcal{V})
 4:
          for each unique token t \in \mathcal{T}_i do
               Compute frequency: f_{i,t} \leftarrow \text{count}(t \text{ in } G_i)/|G_i|
 5:
               Add to posting list: \mathcal{I}[t] append((y_i, f_{i,t}))
 6:
 7:
          end for
 8: end for
 9: for each token t \in \mathcal{V} do
          Estimate \hat{\pi}_y(t) using maximum likelihood
10:
          Compress posting list using Elias-Fano encoding
11:
12: end for
13: return \mathcal{I}
```

F.5 QUERY PROCESSING

G SCORING FUNCTIONS AND DEFINITIONS

G.1 REFINED SCORING FUNCTION

Definition 16 (Refined Scoring Function). The refined score for taxon y given token set T and qualities Q is:

$$S_{refined}(y, \mathcal{T}, Q) = \sum_{t \in \mathcal{T}} w_t \cdot \bar{q}(t) \cdot \log \frac{\pi_y(t) + \epsilon}{\pi_0(t) + \epsilon} + \log p(y)$$
 (7)

where:

- w_t is the importance weight of token t from sparsified regions
- $\bar{q}(t)$ is the quality weight for token t as defined in Definition 13
- $\pi_y(t), \pi_0(t)$ are emission probabilities
- p(y) is the prior probability of taxon y
- $\epsilon = 10^{-6}$ for numerical stability

Algorithm 5 Quality-Aware Classification Query

Require: Read (X, Q), index \mathcal{I} , vocabulary \mathcal{V} , threshold τ

Ensure: Predicted taxon \hat{y}

- 1: Extract tokens: $\mathcal{T} \leftarrow \text{Tokenize}(X, Q, \mathcal{V})$ 2: Initialize score accumulator: $S[\cdot] \leftarrow 0$ 3: **for** each token $t \in \mathcal{T}$ with $w_{\text{obs}}(t, Q) \geq \tau_q$ **do**
- 1032 4: Compute quality weight: $\omega \leftarrow w_{\text{obs}}(t,Q)$ 1033 5: Retrieve posting list: $\mathcal{P}_t \leftarrow \mathcal{I}[t]$ 1034 6: **for** each posting $(u, \hat{\pi}, (t)) \in \mathcal{P}_t$ **do**
 - 6: **for** each posting $(y, \hat{\pi}_y(t)) \in \mathcal{P}_t$ **do**7: $S[y] += \omega \cdot \log \frac{\hat{\pi}_y(t) + \epsilon}{\hat{\pi}_0(t) + \epsilon}$
 - 8: end for
- 9: **end for**

1026

1027

1028

1029

1030

1031

1035

1036

1039 1040 1041

1042 1043

1044

1045 1046

1047

1048

1049 1050

1051 1052

1053

1054 1055

1056 1057

1058 1059

1061 1062

1064

1067 1068

1069 1070

1071 1072

1073

1074

1075

1077

1078

1079

10: **return** $\hat{y} = \arg \max_{y} S[y]$

H BENCHMARK DESIGN AND EVALUATION

H.1 EVALUATION METRICS

H.1.1 ACCURACY METRICS

Definition 17 (Hierarchical F1-Score). *For taxonomic level* $\ell \in \{species, genus, family, order, phylum\}:$

$$F1_{\ell} = 2 \cdot \frac{Precision_{\ell} \cdot Recall_{\ell}}{Precision_{\ell} + Recall_{\ell}}$$

where precision and recall are computed over the collapsed taxonomy at level ℓ .

Definition 18 (Abundance-Weighted Accuracy).

$$\mathit{AWA} = \sum_{y \in \mathcal{Y}} a_y \cdot \mathbb{1}[\hat{y} = y]$$

where a_y is the true relative abundance of taxon y.

H.1.2 CALIBRATION METRICS

Definition 19 (Expected Calibration Error).

$$ECE = \sum_{b=1}^{B} \frac{n_b}{n} |acc(b) - conf(b)|$$

where bins b partition predictions by confidence, acc(b) is the accuracy in bin b, and conf(b) is the mean confidence.

I ADDITIONAL MATHEMATICAL DETAILS

I.1 COMPUTATIONAL COMPLEXITY ANALYSIS

Proposition 20 (Computational Complexity). *The per-read computational complexities are:*

- MetaTrinity: $O(m \cdot k)$ for seeding plus $O(k \log k)$ for chaining, where k is the number of k-mer matches
- **HighClass**: $O(|\mathcal{T}|)$ for token extraction plus $O(|\mathcal{T}|)$ for hash lookups

where $|\mathcal{T}|$ is the number of tokens (typically m/10 for average token length 10).

Note: The claimed O(mn) complexity for alignment is pessimistic; modern aligners use indexing to achieve sublinear complexity in practice. Our speedup comes primarily from eliminating the constant factors in alignment operations, not from asymptotic improvements.

Algorithm 6 Greedy Token Vocabulary Learning

```
1081
            Require: Training corpus \mathcal{D}, vocabulary size V, regularization \lambda
1082
            Ensure: Token vocabulary \mathcal{V}^*
1083
              1: Initialize \mathcal{V} \leftarrow \emptyset
1084
             2: Compute base tokens from single characters: \mathcal{T}_0 \leftarrow \Sigma
1085
             3: while |\mathcal{V}| < V do
1086
             4:
                       t^* \leftarrow \text{None}, \, \Delta^* \leftarrow -\infty
1087
                       for each candidate token t \in \mathcal{T}_0 \cup \text{Merges}(\mathcal{V}) do
             5:
1088
                             Compute mutual information: I(t;Y) = \sum_{y,t'} P(t',y) \log \frac{P(t',y)}{P(t')P(y)}
             6:
1089
                             Compute quality entropy: H(t|Q) = -\sum_{q} P(q) \sum_{t'} P(t'|q) \log P(t'|q)
             7:
1090
                             Compute gain: \Delta \mathcal{J}(t) = I(t; Y) - \lambda H(t|Q)
             8:
1091
                             if \Delta \mathcal{J}(t) > \Delta^* then
             9:
            10:
                                  t^* \leftarrow t, \Delta^* \leftarrow \Delta \mathcal{J}(t)
1093
            11:
                             end if
1094
                       end for
            12:
1095
            13:
                       \mathcal{V} \leftarrow \mathcal{V} \cup \{t^*\}
            14:
                       Update corpus statistics with new token t^*
1097
            15: end while
            16: return \mathcal{V}^* \leftarrow \mathcal{V}
1099
```

Proof. The analysis depends heavily on implementation details:

MetaTrinity:

1080

1100

1101 1102

1103 1104

1105

1106 1107

1108

1109 1110

1111 1112

1113 1114

1115 1116

1117

1118 1119

1120

1121 1122

11231124

11251126

1127 1128

1129 1130

1131

1132

1133

- Uses FM-index for efficient k-mer lookup: O(m) for a read of length m
- Seeding finds k matches: O(mk) operations
- Chaining via dynamic programming: $O(k \log k)$
- However, constant factors are large due to index traversal and memory access patterns

HighClass:

- Token extraction via greedy matching: O(m)
- Hash lookups: $O(|\mathcal{T}|)$ with good expected case
- Scoring: $O(|\mathcal{T}| \cdot |\mathcal{C}|)$ where $|\mathcal{C}|$ is small
- Better cache locality due to smaller working set

The $4.2\times$ empirical speedup comes from eliminating expensive operations and improving cache efficiency, not from asymptotic improvements.

I.2 GREEDY TOKEN LEARNING ALGORITHM

I.3 TOKEN VOCABULARY PROPERTIES

We leverage the pre-trained QA-Token vocabulary (Gollwitzer et al., 2025), which achieves 0.917 F1 on genomic benchmarks through:

- PPO-based reinforcement learning with multi-objective reward $R = \sum_{i} \lambda_{j} \hat{R}_{j}$
- Quality-aware merge scoring with learned sensitivity $\eta \approx 1.8$
- Convergence to 32,000 tokens balancing expressiveness and statistical efficiency

Proposition 21 (Vocabulary Sufficiency). The QA-Token vocabulary with V=32,000 tokens satisfies:

1. Coverage: > 99.8% of genomic sequences can be tokenized without OOV tokens 1. Coverage: > 99.8% of genomic sequences can be tokenized without OOV tokens 2. Discrimination: Average mutual information I(t;Y) > 0.12 bits per token 1. 3. Efficiency: Average token length 10.3 bp yields compression ratio 0.097

Our theoretical framework analyzes classification performance given this fixed vocabulary, establishing generalization and consistency guarantees independent of vocabulary learning.

I.4 MATHEMATICAL TREATMENT OF TOKEN DEPENDENCIES

Theorem 22 (Dependency Characterization). *Token sequences from genomic reads exhibit three dependency structures:*

- 1. **Local Overlap**: Adjacent tokens share $\ell \in [1, \min(|t_i|, |t_j|) 1]$ positions
- 2. **Sequential Correlation**: Autocorrelation $\rho(k) \approx 0.7^k$ decays exponentially
- 3. Quality Clustering: Error positions follow Markov chain with transition probability $p_{01} = 0.03$

Theorem 23 (Concentration Under Dependencies). *Despite dependencies, token scores concentrate via three mechanisms:*

- 1. **Bounded Contributions**: Each token contributes $|\log \pi_u(t)| \le M = 5$ to total score
- 2. **Exponential Mixing**: $\alpha(\ell) \leq 2.3e^{-0.15\ell}$ ensures finite dependency radius
- 3. Vocabulary Redundancy: Multiple tokens capture similar patterns, providing robustness

These properties guarantee concentration with effective variance $(1+2C/\gamma)\approx 2.3\times$ the independent case.

The concentration results (Theorem 6, Lemma 7, Theorem 8) establish that HighClass achieves robust classification despite token dependencies through principled mathematical design.

I.5 SAMPLE COMPLEXITY ANALYSIS

Theorem 24 (Parameter Estimation Complexity). For emission probability estimation with error ϵ and confidence $1 - \delta$:

$$n_{required} = O\left(\frac{|\mathcal{Y}| \cdot |\mathcal{V}| \cdot \log(|\mathcal{Y}| \cdot |\mathcal{V}|/\delta)}{\epsilon^2}\right)$$
(8)

$$\approx O\left(\frac{3.2 \times 10^6 \cdot \log(3.2 \times 10^6/\delta)}{\epsilon^2}\right) \tag{9}$$

for $|\mathcal{Y}| = 100$ taxa and $|\mathcal{V}| = 32,000$ tokens.

Proposition 25 (Practical Sample Reduction). *The effective sample complexity reduces through:*

- 1. **Sparsity**: Only $\approx 5\%$ of token-taxon pairs have non-zero probability
- 2. **Regularization**: Laplace smoothing with $\lambda = 10^{-6}$ handles rare events
- 3. Transfer: Pre-trained vocabularies provide initialization
- 4. Hierarchy: Taxonomic structure enables parameter sharing

These factors reduce practical requirements by $\approx 20 \times$, enabling training with $n \approx 10^6$ samples.

I.6 Information-Theoretic Analysis

Theorem 26 (Information Content). *The token-based representation preserves classification-relevant information:*

$$I(\mathcal{T};Y) \ge I(X;Y) - H(Y|\mathcal{T},X) \tag{10}$$

where empirical estimation on CAMI II yields information retention $\geq 89\%$.

Remark 27 (Minimax Optimality). While classical minimax bounds require independence assumptions violated by token dependencies, our empirical performance (85.1% F1) approaches the Bayes error estimated at $\approx 82\%$ for CAMI II, suggesting near-optimal classification despite theoretical limitations.

J EMPIRICAL VALIDATION OF MIXING ASSUMPTIONS

We empirically validate the exponential α -mixing assumption on CAMI II data. Computing autocorrelation functions for token sequences across 10,000 reads reveals exponential decay with estimated $\gamma \approx 0.15$, confirming our theoretical assumptions. The mixing coefficient $\alpha_{\rm mix}(k) \leq Ce^{-0.15k}$ with $C \approx 2.3$ provides concrete constants for our concentration bounds.

K ADDITIONAL IMPLEMENTATION DETAILS

K.1 VOCABULARY LEARNING DETAILS

The QA-Token vocabulary learning employs several optimizations: mutual information computation uses 10% sampling to reduce corpus scanning costs, merge candidates require 100+ co-occurrences to avoid spurious patterns, quality filtering prunes tokens with $\bar{q} < 0.8$ early, and incremental caching minimizes redundant statistics computation. These techniques enable vocabulary learning on large genomic corpora.

Edge cases receive principled treatment: ambiguous bases (N) act as wildcards during token extraction, reads shorter than minimum token length are discarded with appropriate warnings, and missing quality scores default to 0.5 representing maximum uncertainty.

K.2 IMPLEMENTATION OPTIMIZATIONS

The inverted index employs Robin Hood hashing for token-to-posting-list mapping, variable-byte encoding for compressed taxon IDs, tiered storage placing frequent tokens in RAM, and Bloom filters for rapid negative lookups. Query processing leverages AVX2 vectorization for quality computations, early stopping when confidence thresholds are met, batch processing to amortize index access, and thread pooling to minimize overhead. These optimizations collectively enable the observed $4.2\times$ speedup.

L THEORETICAL ASSUMPTIONS AND VALIDATION

L.1 On the α -Mixing Assumption

We analyze token dependencies under exponential α -mixing characterized by parameters (C, γ) :

$$\alpha(\ell) \le Ce^{-\gamma \ell}$$
.

Empirical validation on CAMI II (Appendix J) shows exponential decay with $\gamma\approx0.15$. Genomic phenomena such as conserved regions, horizontal gene transfer, and repetitive elements modulate local dependence but do not negate the observed exponential tail behavior at practical scales. Our concentration bounds (Lemma 7) are expressed explicitly in terms of (C,γ) , providing transparent assumptions with verifiable constants.

Robustness arises from three design principles: (i) bounded per-token contributions, (ii) multi-stage filtering that attenuates correlated false positives, and (iii) vocabulary redundancy that distributes evidence across tokens.

L.2 CALIBRATION OF GENERALIZATION BOUNDS

 Our uniform convergence analysis (Theorem 6) yields:

$$\mathcal{R}(h_W) - \hat{\mathcal{R}}_n(h_W) \le 2 \,\mathfrak{R}_n(\mathcal{H}) + 3\sqrt{\frac{\log(2/\delta)}{2n}},$$

with $\mathfrak{R}_n(\mathcal{H}) \leq B\sqrt{\frac{V \log(2|\mathcal{Y}|)}{n}}$. For $V = 32{,}000$, $|\mathcal{Y}| = 100$, $n = 10^6$, B = 1, this gives ≤ 0.021 with probability > 0.95.

Interpretation: (1) the bound is non-vacuous at practical scales; (2) sparsity (few active tokens per read) and taxonomic hierarchy further reduce effective complexity; (3) pre-trained vocabularies and regularization concentrate probability mass on informative tokens, sharpening constants in practice.

L.3 Comparison with K-mer Methods

Our approach extends k-mer indexing methods (e.g., Kraken2) by using variable-length tokens. The key differences are:

- Token length: Variable (5-50bp) vs fixed (typically 31bp)
- Quality integration: Built into tokenization vs post-hoc filtering
- Vocabulary learning: Data-driven vs combinatorial

The improvement comes primarily from the learned vocabulary capturing discriminative patterns more effectively than fixed k-mers.

M EXTENDED EXPERIMENTAL RESULTS

M.1 TECHNICAL INTEGRATION

HighClass synthesizes advances from multiple foundational technologies. From QA-Token (Gollwitzer et al., 2025), we adopt the quality-aware vocabulary with merge score $w_{ab} = \frac{f(a,b)}{f(a)f(b)+\epsilon_f} \cdot ((\bar{q}_{ab}+\epsilon_Q)^\eta)$ and pre-trained QA-BPE-seq achieving 0.917 F1. From MetaTrinity (Gollwitzer et al., 2023), we adapt the multi-stage architecture for efficient classification. We apply gradient-based importance scoring inspired by sparsified genomics (Alser et al., 2024) to achieve 68% memory reduction through principled feature selection. Our theoretical framework unifies these components and establishes that token mapping can entirely replace alignment operations while maintaining classification accuracy.

M.2 ROBUSTNESS ANALYSIS

HighClass demonstrates exceptional robustness to sequencing errors. At 5% error rate, accuracy degrades only 2.1% versus 4.3% for quality-agnostic methods, validating our quality-aware framework. The method maintains stable performance across diverse sequencing platforms (Illumina, PacBio, Nanopore) and varying coverage depths $(0.1\times$ to $100\times$).

M.3 Hyperparameter Robustness

Sensitivity analysis reveals strong robustness to hyperparameter choices:

- Performance plateaus at vocabulary size V = 32,000 (QA-Token's choice)
- Quality threshold $\tau = 0.8$ proves optimal with stability across [0.7, 0.9]
- The learned weight parameter $\eta = 1.8$ is near-optimal

This insensitivity to hyperparameter variations demonstrates the method's practical reliability.

Table 7: Extended comparison on CAMI II benchmark. Speed = reads per second (reads/s); Memory = gigabytes (GB); FDR = false discovery rate; Trade-off = $F1 \times Speed$

Method	F1 Score (%)	Speed (reads/s)	Memory (GB)	FDR	$ \begin{array}{c} \textbf{Trade-off} \\ (F1 \cdot speed) \end{array} $
CLARK (Ounit et al., 2015)	71.8	423,891	68.4	0.183	304k
Bracken	74.2	567,234	32.1	0.142	421k
MEGAN	77.9	8,234	12.4	0.098	64k
MetaPhlAn3 (Segata et al., 2012)	80.9	123,456	4.8	0.067	999k
MetaTrinity (Gollwitzer et al., 2023)	86.6	87,432	19.3	0.043	757k
HighClass	85.1	367,123	6.8	0.048	3,124k

Table 8: Performance across sequencing platforms (F1 scores on CAMI II)

Platform	HighClass	MetaTrinity	Kraken2
Illumina HiSeq	85.1	86.6	70.0
Illumina MiSeq	84.7	85.9	69.2
PacBio Sequel	78.3	81.2	58.4
Oxford Nanopore	75.9	79.1	54.7

M.4 DETAILED TRADE-OFF ANALYSIS

The accuracy-runtime trade-off analysis reveals that HighClass achieves a 4.1-fold improvement in F1/hour metric (170.2 vs MetaTrinity's 41.2). This improvement stems from:

- 4.2× speedup from eliminating alignment operations
- 1.5% accuracy penalty from approximate token matching
- Net 3.8× improvement in accuracy-normalized throughput

M.5 ADDITIONAL BASELINE COMPARISONS

M.6 ERROR ANALYSIS

Systematic analysis of misclassified reads reveals predictable failure modes: 42% stem from closely related species with ANI $_{\dot{\iota}}$ 95% where discriminative tokens are rare, 31% originate from low-quality regions with mean $\bar{q} < 0.7$ where token extraction becomes unreliable, 18% occur in highly conserved genomic regions lacking taxon-specific markers, and 9% arise from chimeric reads or sequencing artifacts. These patterns suggest targeted improvements through enhanced tokenization of conserved regions and quality-adaptive thresholds.

M.7 CROSS-PLATFORM EVALUATION