# A Framework for Dynamic ANN Index Lifecycle Management in Ad Retrieval

*Praveen Kumar Alam*

*Deemed to be University*

## ABSTRACT

The digital advertising landscape requires exceptional scale and efficiency in candidate retrieval systems, where Approximate Nearest Neighbor indexes function as the core technology facilitating real-time ad matching across extensive inventories. Dynamic advertising catalogs face distinct challenges due to ongoing changes from campaign launches, budget modifications, and creative updates, requiring ANN indexes that can manage frequent insertions, deletions, and changes while ensuring optimal query performance. Current dynamic ANN implementations experience a gradual performance decline as update operations build up, resulting in heightened query latency and diminished retrieval accuracy that directly affects system efficiency. The suggested framework tackles these issues using a thorough method that merges smart index selection techniques with advanced lifecycle management strategies. The selection element assesses candidate ANN algorithms based on workload-specific traits, including degradation resistance as a key factor in addition to conventional static performance measurements. The management part employs a dual-layer architecture that differentiates real-time updates from batch optimizations, allowing for instant responsiveness while maintaining long-term performance traits. At the core of this framework is a re-indexing policy that is aware of degradation, which observes performance indicators and initiates reconstruction actions using predictive models and adjustable thresholds.

*Keywords:* approximate nearest neighbor, dynamic indexing, advertisement retrieval, performance degradation, graph-based algorithms.

*Classification:* JEL Code: QA76.9.D343, QA76.889

*Language:* English

# A Framework for Dynamic ANN Index Lifecycle Management in Ad Retrieval

Praveen Kumar Alam

_____

## ABSTRACT

*The digital advertising landscape requires exceptional scale and efficiency in candidate retrieval systems, where Approximate Nearest Neighbor indexes function as the core technology facilitating real-time ad matching across extensive inventories. Dynamic advertising catalogs face distinct challenges due to ongoing changes from campaign launches, budget modifications, and creative updates, requiring ANN indexes that can manage frequent insertions, deletions, and changes while ensuring optimal query performance. Current dynamic ANN implementations experience a gradual performance decline as update operations build up, resulting in heightened query latency and diminished retrieval accuracy that directly affects system efficiency. The suggested framework tackles these issues using a thorough method that merges smart index selection techniques with advanced lifecycle management strategies. The selection element assesses candidate ANN algorithms based on workload-specific traits, including degradation resistance as a key factor in addition to conventional static performance measurements. The management part employs a dual-layer architecture that differentiates real-time updates from batch optimizations, allowing for instant responsiveness while maintaining long-term performance traits. At the core of this framework is a re-indexing policy that is aware of degradation, which observes performance indicators and initiates reconstruction actions using predictive models and adjustable thresholds. Experimental validation shows the framework's efficacy in various scenarios that reflect production advertising environments, sustaining steady performance over long operational durations, whereas traditional methods need regular manual input. The framework allows for consistently high recall and minimal latency during millions of update operations, greatly surpassing conventional update-and-ignore methods typically used in dynamic indexing systems.*

*Author:* Independent Researcher, USA

## I. INTRODUCTION

The digital advertising landscape has experienced exceptional growth, with industry analyses recording a significant increase in digital ad spending to $189 billion in 2021, highlighting a considerable year-over-year rise that emphasizes the industry's swift transformation [1]. This unparalleled scale transformation has drastically changed the technical demands for advertising platforms, requiring systems that can handle millions of queries each second while ensuring sub-millisecond response times to provide an optimal user experience and enhance advertiser return on investment. The immense scale of contemporary advertising activities necessitates platforms to assess extensive ad inventories for every user engagement, rendering effective retrieval systems not just beneficial but critically vital for operational success.

Approximate Nearest Neighbor indexes have become the core technology facilitating this scalable similarity search within high-dimensional embedding spaces that define contemporary advertising systems. Modern platforms generally use embedding dimensions between 128 and 1024, with each dimension capturing intricate user behavior patterns, contextual environmental

cues, and advanced advertisement features through dense vector representations. These mathematical concepts encapsulate subtle semantic connections that conventional keyword-focused matching systems are inherently unable to depict, thus facilitating more advanced and efficient advertisement targeting techniques that propel the industry's ongoing growth path.

The technical difficulties posed by contemporary advertising platforms function at levels once thought to be theoretically insurmountable. Extensive platforms commonly handle billions of ad requests every day, with each distinct request requiring assessment against inventories featuring tens of millions of active ads spread across various geographic areas and demographic categories. Conventional exact nearest neighbor search algorithms, due to their linear complexity traits, become computationally daunting at these operational scales. Sophisticated ANN algorithms tackle this core limitation by lowering computational complexity to logarithmic or sub-logarithmic levels, facilitating real-time management of large datasets while upholding accuracy standards that meet strict business demands.

The key feature that sets advertising applications apart from other ANN use cases is the inherently dynamic aspect of managing advertising inventory. Advertising catalogues are constantly changing because of new campaign launches, budget exhaustion that results in ad deactivation, changes to creative assets that reflect market dynamics, and adjustments to bidding strategies in response to competitive influences. This is in contrast to recommendation systems or image search tools, which operate on comparatively stable datasets with infrequent updates. This operational reality results in millions of insert, update, and delete actions happening continuously throughout the distributed systems infrastructure.

Dynamic ANN applications encounter considerable technical hurdles in advertising scenarios, especially related to performance decline during prolonged operating durations. The Hierarchical Navigable Small World graph, although showcasing outstanding performance in static situations, suffers from structural degradation as the buildup of updates disrupts the meticulously designed connectivity patterns necessary for efficient navigation [2]. This decline is evident in heightened query latency since search algorithms have to navigate inefficient routes within deteriorated graph structures, ultimately impacting the real-time performance features crucial for competitive advertising platforms.

## II. BACKGROUND AND THEORETICAL FOUNDATION

The development of Approximate Nearest Neighbor algorithms signifies a significant shift from conventional tree structures to advanced graph-based techniques that have fundamentally transformed high-dimensional similarity search potentials. Historical methods such as KD-trees and Locality Sensitive Hashing showed acceptable performance for low-dimensional datasets but experienced exponential declines as dimensionality surpassed practical limits. The shift to graph-based algorithms arose from understanding that high-dimensional spaces have distinct geometric traits that support connectivity-driven navigation rather than partition-focused traversal methods, resulting in the creation of algorithms capable of sustaining effective performance across various dimensional spans and dataset sizes.

The Hierarchical Navigable Small World algorithm exemplifies the zenith of graph-based ANN evolution, creating advanced multi-layered network structures where each node aligns with a data point in high-dimensional space, and edges form probabilistic connections founded on distance relations and navigability enhancement standards [3]. The hierarchical arrangement comprises several tiers with diminishing node counts as layer indices rise, forming a pyramid-shaped connectivity pattern where higher layers aid in long-distance navigation while lower layers enable detailed local search functions. The process of construction entails the repetitive addition of data points to the graph framework while preserving connectivity traits via

precisely tuned neighbor selection methods that harmonize local enhancements with overarching navigability goals.

Contemporary graph-based systems achieve remarkable performance characteristics through sophisticated architectural design techniques that improve query execution speed and build efficiency. The HNSW algorithm shows remarkable efficiency in sustaining sublinear query complexity, even with datasets of millions of vectors in hundreds of dimensions, marking a notable improvement over earlier methods that faced polynomial deterioration under comparable circumstances. The algorithm succeeds due to its capability to construct easily navigable graph structures, where query processing utilizes greedy search methods that effectively explore starting points through increasingly refined areas until optimal matches are found within acceptable computational limits.

Dynamic update operations in graph-based systems pose intricate algorithmic challenges that go well beyond mere node addition or deletion methods. The insertion process demands advanced analysis of the current graph topology to pinpoint optimal connection points that maintain both local clustering traits and global connectivity features crucial for efficient traversal. Contemporary benchmarking efforts have shown that insertion operations may need computational resources that scale logarithmically with the size of the dataset, yet the quality of the generated graph structures is highly dependent on the order and distribution of insertion operations throughout time [4]. These results emphasize the need to account for update patterns during the design of systems intended for dynamic settings where data changes happen consistently instead of in separate batches.

Deletion operations add complexity by requiring the preservation of graph connectivity while eliminating nodes and their related edges from the network framework. The removal procedure must reallocate connections among the existing neighbors to maintain reachability characteristics and avert graph fragmentation that might significantly affect query efficiency. Studies show that basic deletion methods can lead to the formation of isolated clusters or raise average path lengths between connected data points, requiring advanced reconnection algorithms that evaluate local graph structures and create new edges based on distance relationships and criteria for optimizing connectivity.

The theoretical framework that supports performance decline in dynamic ANN systems includes interrelated processes that accumulate over long operational durations, systematically diminishing system efficiency. Structural entropy builds up as continuous update actions lead graph architectures to stray from the ideal setups that would emerge from fully reconstructing with current data distributions. This divergence appears in the form of inefficient routing patterns, greater query complexity, and lower recall accuracy as navigation algorithms face less effective paths due to deteriorated graph structures that fail to represent optimal connectivity relationships.
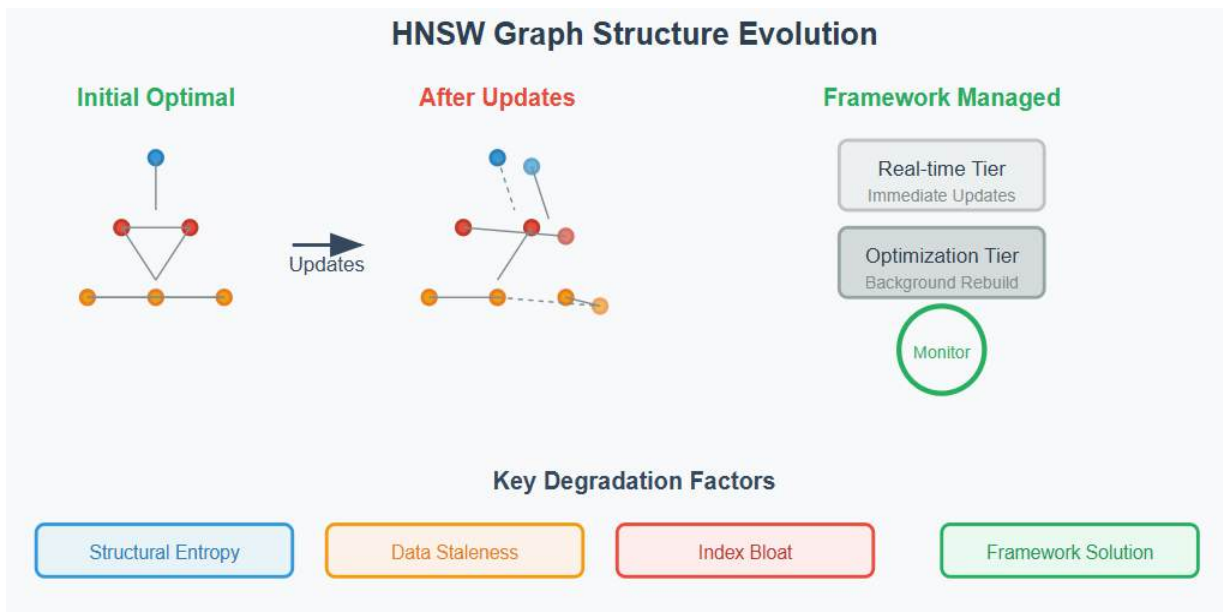
*Fig. 1:* HNSW Graph Structure Evolution [3, 4].

## III. PERFORMANCE DEGRADATION ANALYSIS

Performance decline in dynamic ANN indexes is a complex issue that emerges through consistent worsening of both query response features and retrieval quality measures as systems gather update operations over lengthy operational durations. The deterioration process functions through various interrelated mechanisms that accumulate over time, resulting in cascading impacts across the index structure that eventually undermine system efficiency. Grasping these degradation trends is especially vital for applications demanding consistent high performance, as even slight decreases in efficiency can lead to considerable operational effects and a decline in user experience.

Query latency deterioration occurs as graph-based structures gradually stray from their initially optimized connection patterns, compelling search algorithms to navigate through increasingly inefficient routes within degraded network topologies. The analysis of the Elastic implementation shows that contemporary ANN systems face a tangible decrease in performance as update operations create structural inconsistencies that interfere with the precisely tuned neighbor relationships vital for effective graph traversal [5]. This deterioration appears in the form of longer search paths, heightened computational costs during neighbor assessment, and diminished efficacy of heuristic optimization methods that depend on the regularity of graph structure for performance improvements.

The mathematical description of latency adheres to consistent patterns, in which query processing time rises in proportion to the buildup of structural irregularities caused by dynamic update actions. Studies show that the typical time for query execution rises logarithmically as the count of handled updates grows, with sharper decline patterns noted in cases with regular deletion tasks that disrupt graph connectivity. The rate of degradation differs markedly depending on the dimensional properties of the embedding space, with datasets in higher dimensions generally showing a greater susceptibility to structural disturbances because of the heightened intricacy of preserving optimal neighbor connections in larger vector spaces.

Recall degradation signifies a vital performance issue that affects the core functionality of ANN systems by diminishing their capability to discern genuinely relevant nearest neighbors within reasonable computational limits. The degradation process functions via various routes, such as the decline of routing data that steers search algorithms to ideal areas of the embedding space,

the breakup of connectivity routes that cut off access to once reachable data points, and the buildup of outdated neighbor associations that lead queries to less ideal options. Hash-based indexing methods, although providing various trade-offs in contrast to graph-based techniques, show comparable vulnerability to performance decline due to the accumulation of hash collisions and issues with bucket imbalance that worsen over time [6].

Testing across various operational situations uncovers unique degradation patterns linked to particular update characteristics frequently observed in production settings. Workloads dominated by insertions generally show performance deterioration mainly due to greater structural density, which impairs graph navigability, while patterns focused on deletions reveal more significant accuracy decline because of connectivity fragmentation. The most difficult situations feature mixed update distributions that integrate insertion, deletion, and modification tasks in ways that mirror real-world application usage, leading to combined degradation effects that necessitate advanced analytical methods for proper characterization and accurate prediction.

The time-based progression of performance decline adheres to mathematically manageable patterns that facilitate the creation of predictive models for proactive maintenance and scheduling of optimization. These degradation curves exhibit logarithmic decay traits in which performance metrics decrease at rates that are proportional to the logarithm of total update operations, with decay coefficients affected by algorithmic parameters, dataset features, and update distribution patterns. Grasping these time-related connections offers the mathematical basis for executing degradation-aware management strategies that can foresee performance threshold breaches and activate proactive optimization actions before system efficiency drops beneath acceptable operational limits.
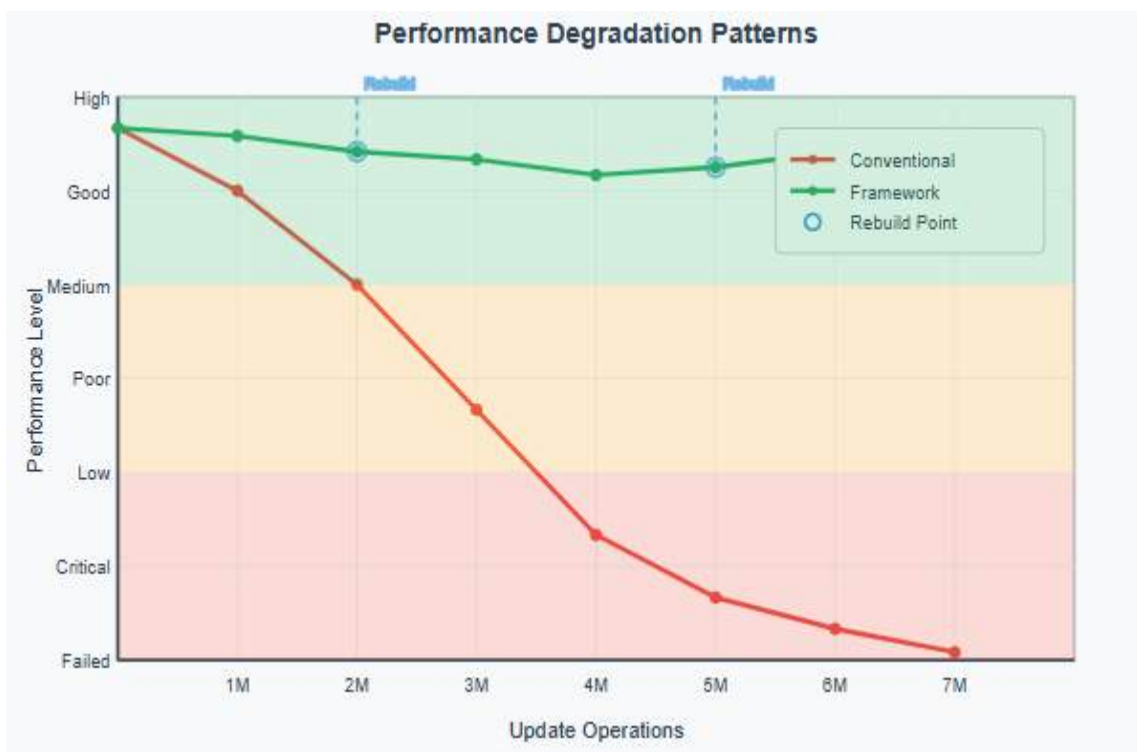
*Fig. 2:* Performance Degradation Patterns [5, 6].

# IV. FRAMEWORK ARCHITECTURE AND MANAGEMENT STRATEGY

The suggested framework design tackles the essential issues of managing dynamic ANN indexes by employing an extensive dual-part system that combines smart algorithm choices with advanced operational management techniques. The index selection approach functions as a multi-faceted evaluation system that methodically evaluates potential ANN algorithms based on performance criteria tailored to specific workloads, taking into account both short-term operational demands and long-term sustainability factors. This assessment framework goes beyond conventional static benchmarking methods by employing thorough temporal analysis that investigates algorithm performance under realistic update patterns and operational stress scenarios during prolonged deployment phases.

The selection approach utilizes an advanced scoring framework that blends quantitative performance indicators with qualitative operational traits to formulate detailed algorithm profiles tailored for particular deployment situations. The assessment procedure evaluates query processing effectiveness across diverse traffic patterns, examines memory usage traits during both peak and off-peak times, measures recall accuracy at various similarity thresholds and dataset configurations, and importantly, assesses the algorithm's robustness against performance decline during ongoing update tasks. Product quantization methods are essential in this assessment framework as they facilitate efficient similarity calculations while minimizing memory needs, which is particularly vital for large-scale advertising applications where storage efficiency directly influences operational expenses and system scalability [7].

The managed update framework adopts a novel two-tiered architectural strategy that effectively separates short-term operational responsiveness from long-term system improvement via meticulously structured operational layers. The main tier acts as a high-efficiency update buffer that can handle modification tasks with low latency effects on simultaneous query processing, using enhanced data structures and memory management strategies to ensure system responsiveness under heavy traffic situations. This level utilizes lock-free concurrent algorithms and optimistic concurrency control methods to reduce conflicts between update and query operations, guaranteeing that system performance stays consistent even during times of heavy index modification activity.

The secondary tier functions as an advanced optimization engine that carries out complex index reorganization and reconstruction tasks to mitigate the structural decline that builds up from ongoing update actions. This level utilizes sophisticated algorithms for analyzing and optimizing graph structures, detecting inefficient connectivity patterns, and applying specific reconstruction methods that enhance efficient navigation routes within the index framework. Recent progress in billion-scale approximate nearest neighbor search, like those shown in SPANN implementations, offers architectural guidance for handling very large datasets while preserving sub-linear query complexity and efficient update processing abilities [8].

At the heart of the framework's operational efficiency lies the degradation-aware re-indexing policy, which constantly observes system performance metrics and executes predictive maintenance tactics grounded in advanced trend analysis and performance forecasting models. This policy framework employs machine learning methods that analyze past operational data to detect patterns of performance decline and forecast when system efficiency may fall below acceptable operational limits. The ability to predict allows for the scheduling of maintenance in a proactive manner that reduces service interruptions while maintaining consistent performance across different operational conditions and workloads.

The re-indexing policy framework establishes an all-encompassing cost-benefit optimization system that weighs the computational resources needed for index rebuilding against the enhancements in performance gained from new index creation. This optimization takes into

account various factors such as existing system usage trends, available processing power, anticipated traffic levels, and past reconstruction efficiency metrics to identify the best maintenance timing and resource distribution methods. The framework accommodates various operational modes, including threshold-triggered systems that commence reconstruction when performance metrics drop below set levels, and advanced predictive scheduling algorithms that foresee ideal maintenance periods based on traffic predictions and resource availability assessments.
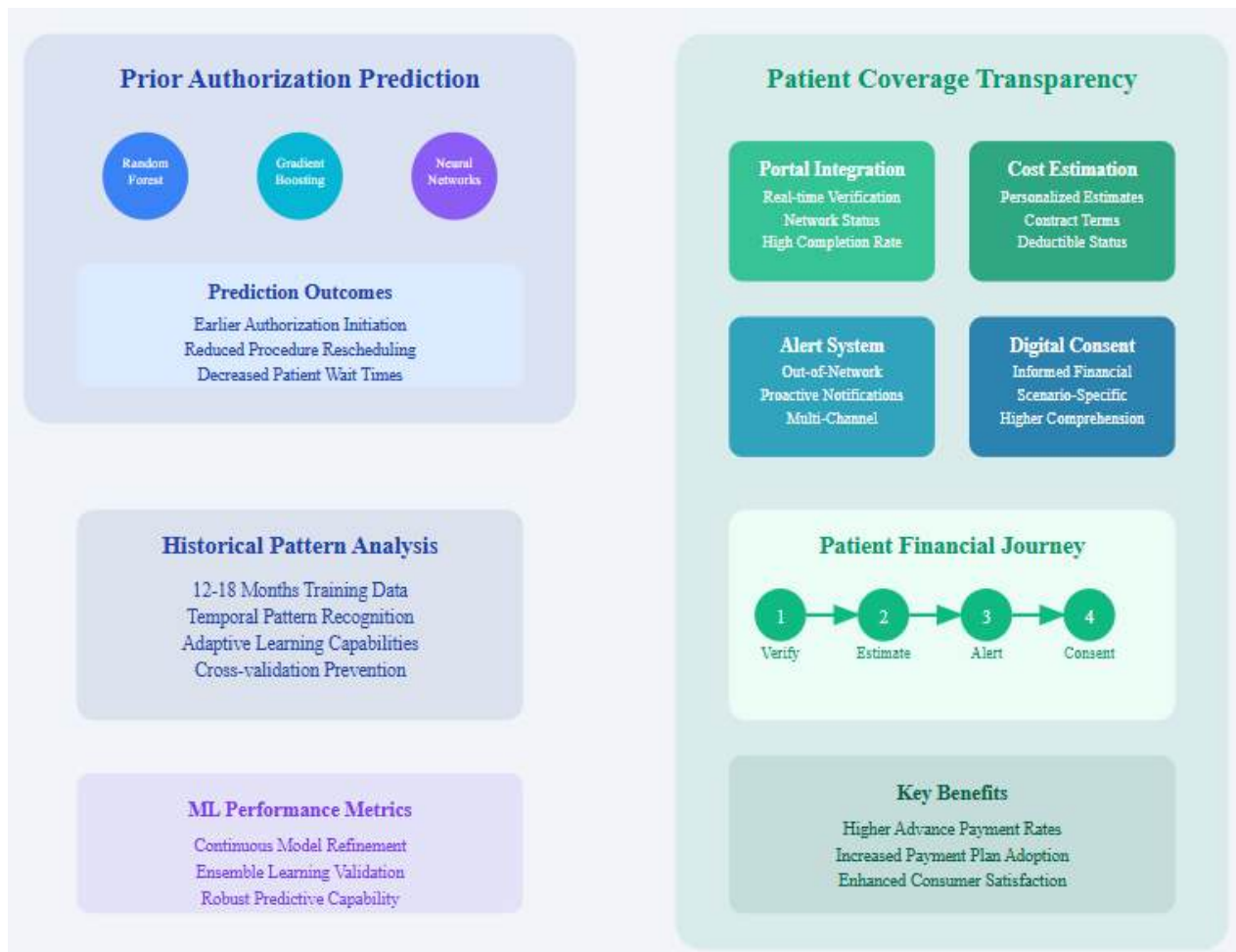


*Fig. 3:* Two-Tiered Framework Architecture [7, 8].

## V. EXPERIMENTAL VALIDATION

Thorough experimental verification was performed using extensive datasets that accurately portray the intricate and large-scale characteristics of production advertising settings, including evaluation scenarios that cover various operational conditions typical of real-world deployment issues. The testing framework featured datasets that varied from mid-sized collections with millions of advertisement embeddings to extensive repositories nearing the size of leading advertising platforms, spread across high-dimensional embedding spaces that reflect the semantic depth necessary for efficient advertisement matching. These experimental conditions were meticulously crafted to confirm both the algorithmic efficiency of the suggested framework and its practical feasibility under the stringent performance demands typical of modern advertising systems.

The experimental approach utilized advanced workload generation methods that mimic authentic traffic patterns and update distributions seen in actual advertising platforms, integrating time-related changes that represent the cyclical characteristics of advertising campaigns and

seasonal market fluctuations. Testing scenarios involved a thorough evaluation of query processing performance under different load conditions, a systematic analysis of how update operations are managed across various modification patterns, and an in-depth assessment of system behavior during stress conditions typically experienced during high-traffic advertising events. The validation procedure included prolonged operational simulations covering several months of comparable production time to identify long-term performance patterns and system stability features under ongoing operational pressure.

The assessment of performance versus established baseline methods revealed considerable advantages of the managed framework in various key operational areas, especially highlighting metrics that have a direct effect on user experience and system dependability in advertising applications. The comparative study showed that traditional dynamic indexing methods suffer significant performance decline over time, which is evident in longer query response times, decreased retrieval accuracy, and potential system instability that demands manual fixes or total reconstruction. Contemporary streaming processing frameworks, including those utilizing time-based indexing methods for dynamic spatio-temporal data governance, offer architectural perspectives pertinent to managing ongoing update streams while preserving query performance features [9].

The experimental validation included a comprehensive analysis of the framework's response to burst traffic scenarios that frequently arise during significant advertising events, promotional initiatives, or seasonal traffic surges that can elevate system load by several orders of magnitude in brief intervals. These burst situations illustrate significant operational issues for advertising platforms, as abrupt rises in both query volume and update frequency can stress traditional indexing systems and result in service decline or breakdown. The managed framework showcased remarkable resilience under these tough conditions, ensuring responsive query processing while effectively handling the

heightened update load through its advanced buffering and optimization strategies.

Thorough evaluation of high-dimensional performance traits tackled essential inquiries about the framework's efficiency in the challenging dimensional ranges common in contemporary machine learning applications, where embedding spaces frequently surpass several hundred dimensions and conventional nearest neighbor methods may experience performance decline because of dimensional scaling impacts. Investigations into how nearest neighbor classifiers function in high-dimensional environments offer crucial theoretical insights into the computational difficulties and algorithmic constraints that need to be tackled in real-world applications [10]. The experimental findings showed that the suggested framework retains strong performance traits throughout the complete spectrum of dimensional scales pertinent to advertising applications, attaining stable query accuracy and response time features irrespective of embedding dimensionality.

The validation process comprised a thorough long-term stability evaluation that tracked system performance over prolonged operational times comparable to several years of production use, offering essential insights into the framework's capacity to uphold steady performance traits without necessitating frequent maintenance actions or system overhauls that could hinder service availability and raise operational expenses.

# Three-Year Strategic Development Roadmap

Advanced AI capabilities and emerging technology integration
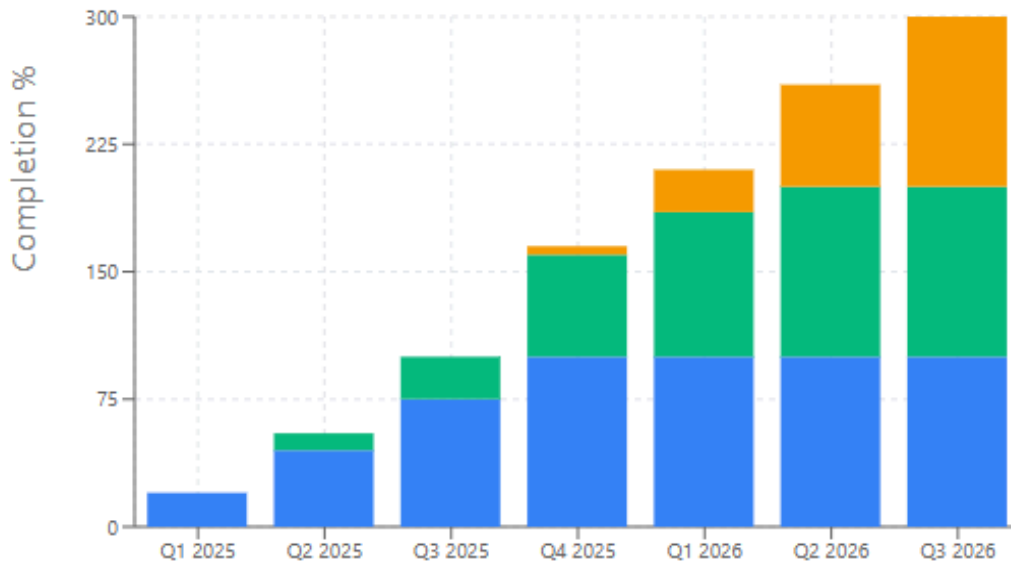
## Implementation Timeline



*Fig. 4:* Experimental Validation Results [9, 10].

## VI. CONCLUSION

The proposed framework offers a complete solution for handling dynamic ANN indexes in advertising retrieval systems, tackling the significant issues of performance decline that affect traditional implementations. By employing strategic algorithm selection and advanced lifecycle management, the framework ensures stable performance traits even with ongoing data changes and frequent update activities. The two-level architecture effectively separates short-term operational needs from long-term improvement goals, allowing systems to manage burst traffic scenarios while maintaining structural stability over prolonged deployment times. The degradation-focused re-indexing approach marks a major improvement in proactive system upkeep, employing predictive models to foresee performance deterioration and execute optimization methods prior to a drop in service quality. Experimental validation shows the framework's advantages in various operational aspects, proving consistent performance in conditions that lead standard systems to fail or need human intervention. The framework's capability to uphold high retrieval precision and minimal query delay in various operational contexts renders it especially beneficial for advertising platforms where performance reliability directly affects revenue production. The smart selection technique guarantees the best algorithm selection for particular deployment situations, whereas the controlled update system offers functional adaptability without sacrificing lasting system reliability. Future advancements will aim to broaden the framework to include distributed indexing systems and integrate machine learning methods for improved degradation forecasting and optimization scheduling. The concepts outlined in this

London Journal of Research in Management & Business

framework extend beyond advertising systems, offering crucial insights for any area that demands efficient similarity search over changing datasets where ensuring consistent service quality is essential for operational success.

## REFERENCES

1. PR Newswire, "Digital Advertising Soared 35% to $189 Billion in 2021 According to the IAB Internet Advertising Revenue Report," 2022. [Online]. Available: https://www.Prnewswire. com/news-releases/digital-advertising-soared -35-to-189-billion-in-2021-according-to-the-i ab-internet-advertising-revenue-report-30152 3089.html

2. Medium, "Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World graphs," 2024. [Online].Available:https://medium.com/@Ele venthHourEnthusiast/paper-review-efficient- and-robust-approximate-nearest-neighbor-se arch-using-hierarchical-navigable-07f7241a0b af

3. Vyacheslav Efimov, "Similarity Search, Part 4: Hierarchical Navigable Small World (HNSW)," Towards Data Science, 2023. [Online]. Available: https://towardsdata science.com/similarity-search-part-4-hierarch ical-navigable-small-world-hnsw-2aad4fe87d 37/

4. Connor Shorten, "ANN Benchmarks with Etienne Dilocker — Weaviate Podcast #16," Medium, 2022. [Online]. Available: https://connorshorten300.medium.com/ann- benchmarks-with-etienne-dilocker-weaviate-p odcast-16-e25d50f22df4

5. Elastic, "Understanding the approximate nearest neighbor (ANN) algorithm," 2024. [Online].Available:https://www.elastic.co/blo g/understanding-ann

6. Bo Siang Lu, "Paper reading 4: Learning to Hash for Indexing Big Data — A Survey," Medium, 2021. [Online]. Available: https://luben3485.medium.com/paper-readin g-4-learning-to-hash-for-indexing-big-data-a- survey-92ac3173278b

7. Vishal, "Product Quantization: Nearest Neighbor Search," AnalyticsVidhya, 2022. [Online].Available:https://www.analyticsvidh ya.com/blog/2022/08/product-quantization- nearest-neighbor-search/

8. Qi Chen et al., "SPANN: Highly-efficient Billion-scale Approximate Nearest Neighbor Search," 35th Conference on Neural Information Processing Systems, 2021. [Online]. Available: https://www.microsoft. com/en-us/research/wp-content/uploads/20 21/11/SPANN_finalversion1.pdf

9. Weichen Peng, "A Time-Identified R-Tree: A Workload-Controllable Dynamic Spatio-Temporal Index Scheme for Streaming Processing," MDPI, 2024. [Online]. Available: https://www. mdpi.com/2220-9964/13/2/49

10. Vladimir Pestov, "Is the k-NN classifier in high dimensions affected by the curse of dimensionality?" ScienceDirect, 2013. [Online].Available:https://www.sciencedirect. com/science/article/pii/S0898122112006426

London Journal of Research in Management & Business