# Clustering, Indexing, and Data Structures for High-Dimensional and Categorical Data: Algorithmic Foundations, Modern Advances, and Scalable Analytic Systems

## Abstract

This survey provides a comprehensive and critical synthesis of contemporary advances in clustering, indexing, and analytic methodologies for high-dimensional and categorical data. Motivated by the widespread emergence of large, complex datasets in domains such as genomics, healthcare, e-commerce, and network analysis, the paper elucidates the fundamental challenges posed by the "curse of dimensionality," data heterogeneity, and the proliferation of categorical and multimodal variables. The scope encompasses key computational paradigms, including nearest neighbor search, clustering, feature selection, and high-dimensional statistical testing, as well as foundational and emerging indexing structures—from traditional spatial trees to compressed, learned, and hybrid neural indexes.

Key contributions include an in-depth analysis of algorithmic strategies tailored to high-dimensional settings, such as ensemble subspace and consensus spectral clustering, robust tensor decompositions, and adaptive index constructions leveraging machine learning. The survey further evaluates space-efficient storage and hardware-accelerated computation, addressing real-time scalability, dynamic adaptation, and resilience to noisy, adversarial, or streaming data. Comprehensive benchmarking, cluster validation, and open-source ecosystem reviews contextualize methodological innovations within system-level performance and reproducibility frameworks.

Conclusions highlight persisting open problems: balancing statistical rigor and computational efficiency, ensuring robustness and interpretability, integrating ethical and privacy considerations, and advancing standardized benchmarking. The survey delineates future research directions—including federated analytics, neural and retrieval-augmented indexing, and unified analytic platforms—emphasizing that adaptive, accountable, and explainable methodologies are essential to harnessing the potential of high-dimensional data across scientific and societal domains.

## 1 Introduction

Recent advancements in artificial intelligence (AI) have led to an exponential growth of research output and practical applications across a broad spectrum of fields. This survey aims to provide a comprehensive and accessible synthesis of current developments in the domain, specifically targeting graduate students, researchers, and professionals in computer science, data science, and allied fields seeking an in-depth understanding of contemporary AI paradigms and their foundational techniques.

A unique aspect of our survey is the emphasis on a unified taxonomy that categorizes approaches according to their core methodologies and cross-domain applicability. We explicitly delineate how this taxonomy advances prior surveys by integrating emerging subfields and highlighting underexplored connections between traditionally disparate research areas. The structured mapping provided herein distinguishes our survey and enables practitioners to more readily identify relevant methods for their specific problem domains.

Moreover, we identify and detail key research gaps uncovered during our systematic review—including areas such as scalability in complex environments, robustness against adversarial conditions, and ethical implications of large-scale deployment. For each, we outline actionable future research opportunities, with the intent to guide ongoing and future work toward addressing these open challenges.

In the following sections, we will introduce a summarized view of the main AI paradigms and their distinguishing features as reflected in Table 1. This table is presented at the start of our technical overview to reinforce key conceptual takeaways. The remainder of the paper systematically maps numbered references to their full bibliographic entries, ensuring accurate traceability and improved reader experience.

### 1.1 Motivation

High-dimensional and categorical data have become pervasive across a broad spectrum of modern analytical domains, driven by rapid advancements in data acquisition, storage, and sensing technologies in fields such as healthcare, genomics, e-commerce, and network analysis [6, 13, 15, 28, 33, 38–40, 49, 52, 53, 60–63, 66, 79, 84, 89, 91–93, 107, 108]. These data types are distinguished not only by an exponential increase in feature dimensionality, but also by the growing prevalence of categorical variables—which are frequently sparse and non-ordinal. This dual trend introduces significant methodological and computational challenges.

One central issue is the so-called "curse of dimensionality," a phenomenon in which distances between data points lose discriminative power as the number of dimensions increases. This undermines the effectiveness of similarity-based techniques, as well as methods for nearest neighbor (NN) search, clustering, and classification [28, 49, 53, 60, 91, 92]. The high-dimensional regime facilitates

noise accumulation, where the signal-to-noise ratio degrades, ultimately diminishing the discriminatory capacity of models even as data and computational resources scale [53, 84, 92, 93]. Beyond these statistical hurdles, high dimensionality incurs significant computational overhead in both data storage and algorithmic execution. Classically efficient indexing and search strategies may deteriorate from logarithmic or sublinear time complexity to linear or superlinear, as they struggle with the combinatorial growth of feature configurations [13, 15, 40, 62, 91].

Categorical variables present further complications. Their sparsity and the absence of inherent distance metrics inhibit the straightforward use of standard statistical and machine learning approaches, often necessitating custom distance metrics, specialized encoding techniques, or novel regularization frameworks [6, 38, 39, 52, 89, 93]. In high-stakes applications—such as medical diagnosis or bioinformatics—interpretability is paramount; however, the opacity of many high-dimensional models further constrains their practical adoption [60, 61, 66, 84, 107]. Therefore, the ongoing methodological imperative is to develop algorithms that scale efficiently while delivering robustness, interpretability, and reliability in both statistical inference and practical decision-making.

Ongoing research has yielded notable progress in addressing these obstacles. Innovations include compressed computation, ensemble learning strategies, high-dimensional data structures for efficient indexing, and methods leveraging spectral, consensus, and regularization principles. Collectively, these developments have extended the boundaries of feasible analysis for large and complex datasets [28, 39, 52, 60, 84, 93]. Nevertheless, as the scale, speed, and heterogeneity of contemporary datasets continue to intensify, foundational challenges remain. This necessitates continuous methodological innovation and rigorous evaluation of advancements within the evolving algorithmic landscape.

## 1.2 Key Concepts and Terminology

Meeting the analytical demands posed by high-dimensional and categorical data necessitates clear definitions of the core computational problems and methodological strategies that underpin modern practice [6, 13, 15, 28, 33, 38–40, 49, 52, 53, 57, 60–63, 66, 79, 84, 89, 91–93, 107, 108]. Among the fundamental primitives are nearest neighbor (NN) and $k$-nearest neighbor (kNN) search, which support similarity-based queries essential for clustering, classification, anomaly detection, and recommender systems. In high-dimensional settings, both exact and approximate NN algorithms are central, with ongoing advances in indexing and pruning techniques, as well as metric learning approaches, to safeguard nearest neighbor structures in the face of sparsity and noise.

Other core analytical tasks include:

**Similarity and Range Search:** These extend NN paradigms to return all objects within a specified distance or similarity threshold from a query. They are pivotal in data mining, information retrieval, and feature-based querying—especially in graph- or spatial-structured data.

**Clustering:** The process of partitioning data into groups that maximize intra-group similarity. Challenges intensify in high-dimensional contexts, where relevant features are often obscured by spurious or noisy information [61, 66, 84, 93].

**Classification:** Assigns category labels to data objects, typically in a supervised framework. The abundance of irrelevant or redundant features in high-dimensional spaces impedes both model accuracy and interpretability.

**Statistical Testing:** In high-dimensional settings, conventional statistical testing must contend with reduced statistical power and inflated type I/II error rates, due to the effects of multiple hypothesis testing and inter-feature dependencies [66].

**Indexing:** Refers to the construction of data structures—such as $k$-d trees, ball trees, cover trees, and emerging learned or adaptive indexes—that expedite various types of queries, even as dimensions proliferate [52, 57, 60, 89, 92].

Frequently, high-dimensional and categorical data analysis requires the interplay among these concepts. For instance, graph-based representations exploit both spatial and relational proximity, while spectral and consensus methods adapt clustering and similarity measures to enhance partition quality and retrieval robustness [6, 38, 40, 84]. Categorical data clustering, in particular, integrates specialized encoding schemes, variable selection, and consensus mechanisms to mitigate the effect of noise from less informative dimensions [6, 38, 52, 93]. Thus, the field employs a multifaceted toolbox, extending foundational concepts to address the distinct analytical challenges posed by complex, high-dimensional datasets.

## 1.3 Scope and Organization

This survey offers a comprehensive synthesis of recent advances in algorithmic, methodological, and system-level approaches for the analysis of high-dimensional and categorical data, with particular emphasis on elucidating the current state of the art and highlighting foundational challenges and opportunities [57, 70, 112]. The review begins with an in-depth analysis of major algorithmic paradigms, including classic and contemporary methods for NN and kNN search, range search, clustering, classification, and statistical testing, each examined through the lens of dimensionality, data heterogeneity, and categorical structure.

Subsequent sections explore indexing methodologies, covering both established data structures and newly emerging approaches such as learned, adaptive, and hybrid indexes, with a focus on computational efficiency, robustness, and adaptability to dynamic data workloads. Special attention is devoted to trends in data compression and representation learning, including advances in compressed computation, symbolic embedding techniques, and spectral models that facilitate scalable and meaningful analytics on massive datasets.

The survey further discusses ensemble and spectral methods, consensus and subspace clustering, and hybrid statistical–machine learning frameworks. Special emphasis is placed on recent techniques such as consensus spectral clustering and self-constrained spectral clustering [57, 112], which have demonstrated strong robustness in the face of high-dimensional noise and uninformative features, and advances in compressed computation that enable direct algorithmic operations on compressed data representations [70]. Each method is critically evaluated for its effectiveness in extracting meaningful structure and mitigating challenges such as dimensionality-induced noise accumulation.

Finally, the survey contextualizes these algorithmic and methodological advances within the broader landscape of practical system integration. It addresses open research questions and emerging trajectories, including dynamic and adaptive computation, interpretable modeling, and resilient, secure indexing strategies for high-dimensional and categorical data analysis. Through a critical engagement with the current literature across these dimensions, this survey aims to provide a foundational orientation for newcomers and a forward-looking roadmap for future research in this rapidly evolving field.

## 2 Clustering High-Dimensional, Categorical, and Mixed Data

Clustering high-dimensional, categorical, and mixed-type data presents unique challenges due to the nature and complexity of the data involved. This section reviews the key algorithmic paradigms, summarizes their main features, highlights existing research gaps, and introduces frameworks that distinguish this survey from previous works.

Despite the progress outlined in Table 1, several actionable research gaps persist within this domain, categorized by data characteristics and application context:

For high-dimensional data, the curse of dimensionality impacts both cluster discernibility and algorithmic efficiency. Existing subspace and projected clustering paradigms remain limited in scalability to truly large feature spaces and often require manual or heuristic selection of relevant dimensions. There is a need for algorithms that can automatically and adaptively identify informative subspaces in ultra-high-dimensional settings, potentially leveraging recent advances in automatic feature selection and representation learning.

When clustering categorical data, current algorithms rely on specialized distance metrics and information-theoretic criteria. Nonetheless, they often struggle with rare categories, missing data, and the integration of domain-specific constraints. Addressing these issues requires the development of robust categorical similarity measures that can naturally handle noise, rare events, and missing values, possibly with the integration of domain knowledge.

Mixed-type data clustering remains an open challenge due to the differing statistical scales of features and the difficulty of balancing continuous and categorical attributes. Most current approaches use simple concatenation techniques or distance-weighting heuristics, which may not perform optimally in heterogeneous settings. Actionable opportunities exist for designing unified, principled objective functions or embeddings that capture the joint structure of mixed features in a more theoretically grounded way.

This survey distinguishes itself from prior reviews by systematically categorizing clustering techniques according to data type compatibility, scalability, and robustness (as summarized in Table 1), and by introducing a unified framework for evaluating algorithm suitability under mixed real-world constraints. This taxonomy and analytical mapping provide researchers and practitioners with a clearer guide for method selection and highlight novel directions for algorithmic innovation.

In summary, while significant advances have been made across different paradigms of clustering high-dimensional, categorical,

and mixed data, substantial gaps remain in terms of scalability, robustness, and holistic support for mixed-type attributes. Continued research—particularly in adaptive dimensionality reduction, domain-aware similarity measures, and unified objective formulations—will be crucial to addressing these open challenges.

### 2.1 Challenges in Clustering High-Dimensional and Categorical Data

Clustering high-dimensional datasets—encompassing continuous, categorical, or mixed types—entails a suite of formidable statistical and computational challenges. Foremost is the phenomenon of noise accumulation: as dimensionality escalates, the distinction between informative and non-informative features blurs, thereby reducing the reliability of traditional similarity measures. This complication is particularly acute in domains like gene expression analysis and text mining, where only a minority of observed variables substantially contribute to cluster separability. Consequently, uninformative features may give rise to diffuse or spurious clusters, especially under conditions of stochastic or adversarial noise [57].

Categorical attributes further amplify these obstacles due to sparsity and high cardinality, making it difficult to define robust distance or similarity metrics. Such issues undermine both distance-based and model-based clustering algorithms [57], weakening their effectiveness and interpretability in real-world applications.

### 2.2 Ensemble Subspace and Consensus Spectral Clustering

To alleviate the curse of dimensionality and limitations of single-view clustering, ensemble subspace approaches and consensus spectral clustering have emerged as prominent strategies. These techniques typically employ feature transformation—such as one-hot encoding for categorical variables—followed by procedures like random projection or subspace sampling to generate diverse, information-rich feature subsets [57, 66]. Through subsampling, clusters may be constructed using only the most relevant dimensions, thereby mitigating the influence of noisy or irrelevant variables.

The ensemble process involves aggregating the results from multiple subspace clusterings, often quantified via co-association matrices and consensus functions (e.g., majority voting), to capitalize on the collective insights of partially independent clusterings [4, 55]. Parallel and distributed computation paradigms are frequently leveraged to ensure scalability.

A notable advancement is the incorporation of feature reweighting, with data-driven measures guiding the assignment of greater importance to features or subspaces associated with high signal-to-noise ratios. This renders ensemble clustering methods not only more robust to noise but also adaptive to heterogeneous feature landscapes [29, 57]. Theoretical analyses demonstrate that these methods achieve statistical consistency and minimax-optimal error rates even as the fraction of truly informative features diminishes—a scenario common in omics and text mining tasks [57, 66]. Empirical results corroborate these theoretical gains, with ensemble and consensus spectral approaches often outperforming baseline methods in genomics and unstructured text clustering tasks [57].

**Table 1: Summary of Major Clustering Paradigms for High-Dimensional, Categorical, and Mixed Data**

| Paradigm | Data Type(s) Supported | Strengths | Limitations |
|---|---|---|---|
| Subspace/Projected Clustering | High-dimensional | Discovers meaningful clusters in relevant feature subsets | Sensitive to subspace selection, scalability concerns |
| Categorical Data Clustering | Categorical | Tailored similarity/distance functions (e.g., k-modes) | May not generalize across domains |
| Mixed Data Clustering | Mixed (numeric + categorical) | Integrates heterogeneous data types (e.g., k-prototypes) | Balancing influence of each type is challenging |
| Spectral Methods | Mostly numeric, extensible | Good for non-convex structures, adaptable to high dimensions | Computational cost, tuning parameters |
| Model-based Clustering | All types (with extensions) | Probabilistic framework, flexible models | Scalability, model selection complexity |

Despite their advantages, consensus-based frameworks show reduced efficacy when data exhibits complex feature dependencies (e.g., spatial, temporal, or network structures) or when dealing with genuinely mixed-type attributes, situations where standard one-hot or projection-based strategies fail to capture generative processes [57]. Furthermore, algorithmic complexity—though mitigated through parallelization—can pose practical limitations in very high-dimensional or resource-constrained environments [57].

## 2.3 Spectral Clustering and Self-Constrained Extensions

Spectral clustering has become a widely adopted method for high-dimensional and categorical datasets, leveraging the global organizational structure encoded within the eigenspaces of similarity or Laplacian matrices [92, 112]. This framework eschews direct modeling of cluster-wise densities, instead utilizing geometric relationships in a transformed, lower-dimensional embedding.

Recent methodological advancements include self-constrained spectral clustering, wherein the canonical objective is augmented with explicit pairwise or label-based constraints. These constraints encode prior knowledge or enforce desired partition properties, implemented through iterative optimization and alternating update rules. This ensures convergence to partitions that honor both intrinsic data similarities and extrinsic supervisory information [112].

Self-constrained extensions are particularly advantageous in semi-supervised contexts and in scenarios requiring alignment with spatial or relational structures—for example, integrating clustering results with spatial databases or graph-indexed data pipelines [112]. Nevertheless, spectral clustering remains sensitive to affinity matrix construction and parameter tuning, necessitating careful preprocessing and validation to ensure reliability [92, 112].

## 2.4 Alternative Clustering Methodologies

The landscape of clustering methods for high-dimensional and mixed-type data encompasses several foundational and emerging paradigms beyond ensemble and spectral approaches. A selection of prominent alternatives includes:

**Hierarchical Clustering**: Agglomerative and divisive variants are valued for their interpretability through dendrograms and the flexibility to resolve clusters at multiple levels of granularity. However, while hierarchical clustering enables insights into relationships between clustered entities, its scalability may be limited in high-dimensional spaces, and robustness can be diminished if parameter choices are not carefully tuned [3, 4, 19, 36, 37, 47, 57, 62, 72, 75, 82, 88, 89, 92, 99, 104, 111, 112]. For instance, hierarchical strategies incorporating adaptive or parameter-free innovations

have been shown to automatically select the number of clusters and improve resilience to noise [19, 99].

**Bayesian and Model-Based Approaches**: These methods, including mixture models, mixed membership models, and tensor-normal mixtures, enable probabilistic cluster allocation with explicit uncertainty quantification. Recent work highlights scalable inference techniques, such as penalized coordinate descent and spectral approximations, which address issues like overparameterization and algorithmic bottlenecks in ultrahigh-dimensional or grouped omics datasets [3, 36, 57, 59, 62, 67–69, 73, 85, 89, 93, 108, 111, 112]. Nevertheless, limitations persist in settings with extreme dimensionality, compositional covariates, or complex dependence structures [36, 62].

**Tensor Clustering**: Tailored for multiway or structured datasets (e.g., omics, neuroimaging), tensor clustering applies models like tensor normal mixtures to improve parsimony, utilize underlying data structure, and yield more interpretable components. Incorporation of penalization and ensemble techniques helps control false positives, enhances variable selection, and accounts for correlation among predictors [36, 57, 59, 67, 93]. Challenges include handling separability assumptions, complex implementations, and tuning in high-dimensional tensors [36, 59].

**Robust and Hybrid Methods**: By integrating distinct clustering criteria (such as density- and partition-based schemes) or leveraging deep learning-based representations, these strategies facilitate adaptation to irregular, non-globular clusters, compositional data, and heterogeneous feature sets [3, 4, 36, 57, 72, 73, 75, 82, 88, 89, 92, 93, 99, 108, 111, 112]. For instance, fully autonomous clustering frameworks [99] and hybrid distance-based probabilistic models [67] afford robust performance and greater flexibility across diverse application domains.

**Deep Clustering Paradigms**: Jointly optimizing feature representations and cluster assignments within neural architectures, deep clustering methods demonstrate notable resilience to noise, initialization sensitivity, and overlapping or complex data structure [3, 57, 73, 82, 108, 111]. This modularity enables development of domain-adaptive, end-to-end frameworks for clustering high-dimensional images, text, time series, or graphs. Nevertheless, issues regarding model interpretability, hyperparameter calibration, and generalization to new domains remain important open challenges [3, 57, 73, 111].

In summary, no single methodology uniformly prevails across all data types or application criteria. The selection of an optimal clustering approach must be carefully matched to the structure, distribution, and analytical objectives pertinent to each dataset [4, 19, 57, 62, 111, 112].

For a concise synopsis, Table 2 provides a structured comparison of principal clustering approaches for high-dimensional, categorical, and mixed data contexts.

## 2.5 Cluster Validation Metrics and Benchmarking

Robust evaluation of clustering results in high-dimensional and mixed-type contexts relies upon comprehensive validation and benchmarking metrics. These include:

**External Indices:** Metrics such as Adjusted Rand Index (ARI), Normalized Mutual Information (NMI), and Cohen's Kappa facilitate quantitative comparison against known ground-truth labels, enhancing comparability across algorithms when gold standards are available [4, 6, 8, 11, 17, 21, 31, 39, 42, 43, 50, 56, 62, 66, 73, 84, 86, 87, 89, 100–102, 104, 110, 111]. For instance, ARI and NMI are frequently used to provide chance-adjusted and information-theoretic perspectives on agreement, respectively, while Cohen's Kappa quantifies agreement beyond random expectation. It is important to note that, as highlighted in recent studies, external metrics can be influenced by class imbalance and label noise, emphasizing the need for careful metric selection and interpretation [8, 17, 89].

**Internal Indices:** Metrics such as the Silhouette coefficient, Davies-Bouldin index, Dunn index, accuracy, AUROC, and F1-score provide model-agnostic assessments of cluster cohesion, separation, and overall quality, independent of external references [1, 4, 6, 11, 17, 19, 21, 36, 39, 42, 43, 47, 50, 56, 62, 66, 73, 84, 87, 89, 95, 100–102, 104, 110]. Recent work has noted limitations of classical indices, such as their sensitivity to noise, cluster size, and shape, and their tendency to suggest a single optimal cluster count, which may not adequately reflect complex underlying structures [1, 95]. Specifically, new correlation-based indices and modifications to existing ones have been proposed to address the identification of multiple local optima in model selection [95].

**Multimodality-Based and Modern Indices:** Recently, measures such as Dip and Silverman's tests have been advocated due to their sensitivity in uncovering true cluster structure, ability to handle irregular cluster shapes, and improved robustness to noise and high-dimensional complexity. These multimodality-based metrics offer a stronger capability to distinguish clusterable data from unclusterable scenarios, and are recommended for general-purpose clusterability assessment [1, 11, 19, 95, 101, 102, 111]. However, it remains an open challenge to develop measures that can universally handle outliers, chaining, and very small clusters efficiently [1].

Despite notable advances, recent studies highlight that many classical metrics are sensitive to noise, class imbalance, and manipulation, underscoring the importance of multimodal validation and robust benchmarking protocols [4, 8, 11, 19, 39, 42, 66, 73, 87, 89, 100, 104]. Challenges such as the reproducibility of performance claims, selection of robust metrics, and transparent, unbiased benchmarking remain at the forefront. This is further complicated in high-dimensional settings by issues such as the curse of dimensionality, dominance of uninformative features, and the instability of benchmarks due to preprocessing bias and dataset selection [17, 57, 73].

Best practice now mandates the joint use of internal and external metrics, meticulous dataset curation with transparency, and open, reproducible benchmarking pipelines to facilitate method development and comparison [4, 8, 36, 43, 57, 62, 66, 84, 111]. Several studies emphasize the necessity for well-documented repositories, comprehensive evaluation protocols, and standardized reporting to ensure fair and meaningful comparative analysis [8, 17, 43, 73]. The development of algorithm-specific and data-type-specific benchmarking frameworks has also been identified as crucial for progress, especially in high-dimensional, categorical, and mixed data contexts [57, 66, 73].

Collectively, these methodological innovations and validation frameworks delineate both the considerable progress and enduring open challenges in the clustering of high-dimensional, categorical, and mixed data. Ongoing advances in interpretability, scalability, and rigorous benchmarking remain essential for the development of effective clustering methodologies and their translation to a broad array of scientific and practical applications.

## 3 Index Structures and Data Representations

This section aims to provide a comprehensive overview of modern index structures and data representations, with a focus on their application to large-scale, AI-driven clustering and retrieval tasks. We address the core question: How do current indexing paradigms and data representation techniques facilitate efficiency, scalability, and adaptability in the context of rapidly growing data and evolving analytics workflows? Our discussion highlights key structures, their challenges, and their downstream relevance to analytics and data science.

### 3.1 Traditional Index Structures

Classic spatial and multidimensional index structures—including R-trees, k-d trees, Quadtrees, Grid indexes, Inverted Indexes, and Column Stores—have long been foundational in database systems for managing multi-attribute and spatial queries. R-trees and their variants are optimal for bounding spatial objects and facilitating efficient range and topological searches, while k-d trees and Quadtrees naturally partition multidimensional or spatial data for point queries and region decompositions. Grid and inverted indexes enable rapid filtering and set operations, with inverted indexes excelling particularly in text and categorical data retrieval. Column stores further separate data by attribute, supporting high compression and swift analytical scans. Despite their versatility and widespread adoption, these structures present significant trade-offs: while highly effective for low to moderate dimensionality, scaling to higher dimensions often incurs substantial costs in storage, maintenance, and query performance, particularly as datasets increase in both volume and complexity [25, 26]. Moreover, in high-throughput or real-time environments, the continual maintenance and updating of indexes can amplify these costs, leading to bottlenecks that undermine their intended efficiency.

As shown in Table 3, the effectiveness and limitations of these index structures are tightly coupled to the underlying data characteristics and query requirements.

**Table 2: Comparison of Principal Clustering Paradigms for High-Dimensional, Categorical, and Mixed Data**

| Methodology | Primary Advantages | Key Limitations |
|---|---|---|
| Ensemble Subspace/Consensus Spectral | Robustness to noise and irrelevant features; scalable via parallelization | Complexity in affinity aggregation; reduced efficacy for data with intricate dependencies or mixed types |
| Spectral (Standard/Self-Constrained) | Captures global structure; accommodates constraints/prior knowledge | Sensitive to affinity matrix and parameter selection; scaling may be nontrivial |
| Hierarchical | Interpretability; flexible resolution | Parameter sensitivity; scalability challenges in high dimensions |
| Bayesian/Model-Based | Probabilistic inference; uncertainty quantification | Overparameterization; bottlenecks in ultrahigh dimensions |
| Tensor Clustering | Exploits multiway data; improved parsimony | Requires structured data; complex implementation |
| Deep Clustering | End-to-end learning; resilience to noise/overlap | Interpretability; hyperparameter tuning; domain transferability |
| Robust/Hybrid | Adaptive to diverse data; handles irregular shapes | Model selection complexity; computational overhead |

**Table 3: Comparison of traditional index structures by usage and limitations**

| Index Type | Primary Use | Dimensionality Support | Key Limitations |
|---|---|---|---|
| R-tree | Spatial objects, range queries | Low/medium | Degrades with high dimensionality |
| k-d tree | Point queries, region search | Low/medium | Poor balance in high-dim spaces |
| Quadtree | 2D/3D spatial partitioning | Low | Scalability issues |
| Grid Index | Numeric filtering | Medium | Inefficient for skewed data |
| Inverted Index | Text/search, categorical | N/A | Poor for numeric/spatial data |
| Column Store | Analytical scans | N/A | Write overhead, schema constraints |

## 3.2 Limitations for High-Dimensional and Categorical Data

Despite their general flexibility, traditional indexes typically underperform as dimensionality grows—a phenomenon often described as the "curse of dimensionality." For instance, R-trees experience increased node overlap and size, resulting in excessive I/O during searches. Similarly, k-d trees become imbalanced with high-dimensional inputs, suffering from sharply reduced partitioning efficiency [25, 26]. Beyond numerical dimensions, most classical indexes struggle to integrate categorical and mixed-type attributes alongside spatial or numerical information; supporting such heterogeneous data often necessitates complex, task-specific adaptations that ultimately compromise generality and performance. This persistent set of limitations has catalyzed the search for unified, extensible indexing frameworks that can accommodate both high-dimensional and heterogeneous data types—achieving flexible indexing and querying without incurring prohibitive design or operational complexity [25, 26].

## 3.3 Modern Memory-Efficient and Compressed Indexes

The explosive growth of data volumes, coupled with physical memory bandwidth limitations, has driven significant advances in compressed and succinct indexing structures. For example, q-gram trees for graph similarity search demonstrate the feasibility of in-memory, space-efficient indexes, achieving storage reductions of 85–95% relative to conventional structures—while maintaining query performance parity. Such indexes synthesize probabilistic and deterministic techniques via hybrid encodings and succinct filters to effectively localize candidate sets for rapid searches, even at scales involving tens of millions of objects [76].

Probabilistic data structures like advanced Bloom filters and Cuckoo filters extend these innovations, achieving near-optimal space usage and constant expected lookup times with explicit trade-offs between false positive rates, insertion, and deletion capabilities [32, 103]. Recent developments in dynamic address management for Cuckoo filters, such as signed-offsets and overlapping windows, have removed classic constraints, establishing new benchmarks for space efficiency and making these structures well-suited for large-scale analytics and scientific workloads [103].

In trie-based indexes, methods such as the Height-Optimized Trie (HOT) and Adaptive Radix Tree (ART) exploit node compression and dynamic fan-out, achieving a notable balance of lookup speed, memory usage, and update efficiency—attributes particularly valuable for real-time, in-memory database systems [23, 83, 96]. Suffix-based and run-length encoded indexes are tailored for repetitive data scenarios, as seen in web archives or genomics, by constructing compact representations that efficiently support factorization and substring queries. These approaches, particularly when leveraging compressed suffix arrays or run-length Burrows–Wheeler Transform (RLBWT), can achieve asymptotically optimal space on repetitive data [14]: High compression ratios, reducing storage requirements by orders of magnitude Efficient substring search and factorization capability Potential for optimal theoretical bounds on repetitive inputs Nonetheless, supporting dynamic updates in compressed indexes can introduce notable overhead in practice [14, 70].

The broader adoption of memory-efficient indexes is not without trade-offs. Striving for optimal compression may inhibit support for dynamic operations or slow update pathways. The intricate engineering required to support range, similarity, and set queries over succinct structures remains a key challenge [70, 96]. As a result, contemporary research continues to seek an optimal balance among compression, adaptability, and query efficiency.

## 3.4 Compressed Computation Paradigm

With uncompressed data volumes increasingly surpassing hardware capabilities, a shift towards the "compressed computation" paradigm has emerged as a necessity. Here, compression is no longer merely a storage or transmission optimization but forms the

basis for direct in-memory computation. The result is not only minimized storage and I/O but also a fundamentally reduced working set during active processing [70].

Critical advances include data structures and algorithms operating natively upon compressed representations—such as run-length compressed suffix arrays, compressed tries, or space-efficient factorization structures—thus circumventing expensive decompression cycles [14, 23, 32, 70, 76, 83, 96, 103]. For example, the direct transformation from RLBWT to LZ77 factorization has enabled self-indexing in extremely limited space, paving the way for on-the-fly analytics in genomics, textual, and scientific archives [14]. Concurrently, compressed suffix trees embedded with witness structures provide unified, online computation of classic text factorizations, achieving linear time and sublinear space while supporting practical query efficiency [96].

Despite these advances, significant challenges persist. Indexes that operate on compressed data must mediate among conflicting objectives: compression ratio, query latency, and support for updates. For instance, partially or fully compressed repositories raise open questions for similarity and range search, as classic indexes typically presuppose uncompressed or partially indexed data [70]. The evolution of adaptive and query-aware compressed indexes—capable of dynamically alternating between compressed and uncompressed representations—constitutes a core frontier for ongoing research.

## 3.5 Learned, Neural, and Adaptive Indexes

Recognizing the limitations of both classical and compressed index approaches for managing high-dimensional, evolving, or heterogeneous datasets, a new generation of learned and adaptive indexes has emerged. Leveraging machine learning, these indexes reinterpret data access as a form of prediction, employing models that estimate the location or probability of a record within the data structure.

Spline-based learned indexes, such as LiLIS, apply error-bounded piecewise linear models to approximate mappings within sorted or spatially partitioned data. These models provide constant-time ($O(1)$) lookup overhead and have been shown, in distributed big data frameworks, to yield dramatic speedups compared to traditional spatial indexes. In particular, LiLIS operates by integrating error-bounded spline-based learned indexes on per-partition data using flexible, spatially-aware partitioning strategies such as R-tree, Quadtree, k-d tree, and grid-based schemes. The mapping of high-dimensional locations into one dimension (e.g., via Z-order curves) enables efficient point, range, $k$-nearest neighbor, and join queries. Experimental evaluation demonstrates that, for distributed spatial workloads (e.g., on Apache Spark), LiLIS significantly outperforms conventional indexes such as Sedona-R-tree, achieving up to orders-of-magnitude faster query speeds and 1.5–2× faster index build times. For example, on real and large synthetic datasets, LiLIS achieves point query times of 82.59 ms and range query times of 468.64 ms, whereas competing methods are much slower, especially for join and $k$NN queries. However, these benefits can be sensitive to partitioning choices (R-tree partitioners excel generally, while k-d/Quadtree best support join operations) and to

query distribution skew. The reliance on model training also introduces additional computational overhead, particularly for massive datasets, and adapting to more complex queries remains an open concern [25, 60].

To summarize key comparative results from recent work on distributed learned indexes:

Model-driven indexing further encompasses approaches in which partitioning mechanisms—such as R-tree, k-d tree, and Z-order curves—are tightly coupled with predictive mappings offered by the learned model. This coupling achieves notable reductions in both index construction and query evaluation times, though it necessitates careful attention to partitioner selection, sensitivity to workload skew, and the costs associated with ongoing model retraining and integration into distributed dataflow platforms [25, 60].

Frameworks for index selection are beginning to adopt online learning paradigms, notably multi-armed bandit and reinforcement learning strategies. For example, recent multi-armed bandit-based approaches remove dependency on DBA intervention and unreliable cost models by directly and consecutively exploring candidate indexes based on observed query performance. These frameworks have demonstrated strong empirical speedups: up to 75% on highly dynamic or shifting workloads and up to 51–59% on hybrid transactional/analytical (HTAP) settings compared to conventional or deep RL-based approaches, while ensuring convergence to (near-)optimal policies [74]. Such methodologies enable continual online adaptation to rapidly fluctuating workload characteristics and transform auto-tuning from a static optimization task into an ongoing learning process.

At a broader system level, frameworks such as annotative indexing aim to subsume both traditional and learned index paradigms within a common, highly modular architecture. Annotative indexes generalize and unify inverted, columnar, and object indexes under a dynamic, transactional model that supports efficient ACID transactions and concurrency for both reads and writes. These designs facilitate expressive query processing over structured, semi-structured (e.g., JSON), heterogeneous, and graph-based data. Annotative indexing is further distinguished by its support for lazy transformation, hybrid neural and sparse retrieval, composability for retrieval augmented generation (RAG), and compatibility with structured and unstructured query modalities. The architecture is capable of scaling to hundreds of concurrent transactional clients and supports advanced operations, including entity lookup, knowledge graph queries, and neural search over both text and other datatypes, all while ensuring modularity and extensibility [26].

Nonetheless, these modern indexing strategies introduce their own unresolved challenges. Theoretical concerns include maintaining bounded errors or guarantees in high-dimensional predictive indexing, and robustness of continuous model retraining. Practical issues persist around ensuring safe concurrent access, effective distributed scaling, adversarial resilience and security, garbage collection, and seamless integration into complex, distributed data systems. Emerging research directions span GPU-accelerated retraining for threshold workloads, exploration of hybrid designs combining learned and classical index primitives, and development of transactional guarantees for broad query types over both structured and unstructured data, with a particular emphasis on scaling and extensibility in heterogeneous deployment environments [25, 26, 60].

**Table 4: Query times (ms) for spatial queries using LiLIS and Sedona-RK on large datasets [25]. "Much slower" indicates high latency or infeasibility for the corresponding method.**

| Method | Point | Range | kNN | Join |
|---|---|---|---|---|
| LiLIS-K | 82.59 | 468.64 | 650.2 | 228581 |
| Sedona-RK | much slower | much slower | 790993 | much slower |

This progression from classical structures through to compressed, adaptive, and learned indexing reflects the dynamic interplay between algorithmic innovation and practical system engineering. The overarching contemporary challenges focus on reconciling scalability, adaptability, and efficiency across diverse and rapidly-growing data scales and modalities—a goal that continues to drive index structure research at all levels.

## 4 Similarity, Range Search, and Graph Querying

### 4.1 Space-Partitioning Indexes for Query Processing

Space-partitioning indexes play a foundational role in efficient distance, similarity, and range query processing over both spatial and non-spatial datasets. These structures—including grid files, k-d trees, R-trees, spatial hashing, and more recently, learned and ensemble-based indexes—enable rapid pruning of the search space by hierarchically or adaptively aggregating data into regions with shared characteristics. This results in significant reductions in computational redundancy during query evaluation. Notably, grid-based methods often surpass tree-based counterparts in performance when appropriate partition strategies are adopted, particularly for point, range, and join queries, due to superior data linearization and lower index traversal overheads [54, 94]. The introduction of machine-learned indexing and hybrid approaches has further advanced performance, with learned indexes employing regression models and space-filling curves to efficiently predict object positions and minimize lookup times. These techniques are particularly effective for high-dimensional or irregularly distributed datasets [18, 25, 71, 94].

Classical methods, however, encounter scalability barriers in contexts characterized by large-scale and highly repetitive datasets. Traditional inverted indexes and spatial structures often suffer from inefficiencies in both indexing and memory footprint [21, 44]. Recent breakthroughs have leveraged repetitiveness in data through compressed suffix arrays, run-length compressed structures, and grammar-compressed partial answers, which have substantially reduced storage requirements while supporting efficient document retrieval and counting [22, 38]. For applications involving online or evolving similarity functions—common in active learning and interactive data analysis—adaptive indexing solutions such as OA-SIS maintain families of locality-sensitive hash (LSH) indexes, dynamically updating them in response to user feedback without costly retraining. This results in heightened responsiveness and improved resource utilization in scenarios where similarity criteria are fluid [25, 48].

Space-partitioning techniques have evolved to address queries over complex multi-attribute datasets, including spatio-textual documents, 3D point clouds with attributes, and genomic sequences. Innovations such as persistent, parallel spatio-textual indexes and compressed attribute-aware spatial indexing facilitate queries across spatial, textual, and temporal dimensions, supporting top-$k$ retrieval and attribute-based filtering with high throughput and efficient updates [18, 21, 44, 47, 70, 97]. At the algorithmic level, secondary partitioning techniques enhance traditional space partitioning by further dividing index cells, enabling duplicate-free and low-latency range and distance queries on spatially extended or non-point objects. For example, recent work [97] partitions each primary cell into secondary partitions defined by the begin and end values of object extents relative to the cell, which reduces both duplication and unnecessary computations in distance-range and join queries and outperforms earlier approaches in empirical comparisons.

The trajectory of research in space-partitioning indexing is determined by the interplay among data distribution, partitioning granularity, compression strategies, and the necessity for adaptation to evolving query patterns and dynamic workloads.

### 4.2 Efficient Index Management and Scaling

With the rise in query volumes and dataset sizes, efficient index management and robust scaling mechanisms are critical for large-scale data retrieval tasks. Avoiding duplicate results is essential; naive methods not only risk multiple reporting but also contribute to redundant computation, especially when dealing with overlapping spatial objects or complex join queries. To mitigate these issues, secondary partitioning has been introduced, notably by dividing each primary partition into secondary partitions determined by object boundaries, such as the begin and end values relative to a partition's spatial extent. This approach, as shown in [97], improves query precision by precisely localizing candidate object sets, reducing unnecessary verifications, and substantially outperforming previous partitioning schemes and leading data-partitioning indexes in empirical comparisons.

Scalability is further advanced through distributed and parallel index architectures that leverage modern cluster-computing environments such as Apache Spark and Flink. Innovations such as lightweight learned index structures, often using spline-based regression and space-filling curve mappings, enable $O(1)$ lookups by predicting object positions within partitions. For instance, each spatial partition can be equipped with a custom-learned index—as in LiLIS [25]—which mitigates the curse of dimensionality and handles data skew through spatially-aware partitioning strategies (e.g., R-tree, Quadtree, KD-tree, or grid-based). Experimental evidence demonstrates that these strategies dramatically reduce both index construction times and query latencies. For example, LiLIS achieves

Clustering, Indexing, and Data Structures for High-Dimensional and Categorical Data: Algorithmic Foundations, Modern Advances, and Scalable Analytics Conference'17, July 2017, Washington, DC, USA

1.5-2 times faster index construction and query speeds that are over an order of magnitude better than traditional methods, while maintaining full compatibility with big data frameworks and diverse spatial query workloads including range, k-nearest-neighbor, and join queries.

Modern dynamic data environments demand indexing mechanisms with strong support for incremental and parallel updates. Persistent spatio-textual indexes, as discussed in [70, 97], support efficient integration of new data, enabling prompt query execution and supporting workflows that require frequent updates—vital for applications such as event recommendation and geo-tagged information retrieval. Notably, the shift towards compressed and rep-index structures allows for direct computation on compressed representations, yielding dramatic space savings for repetitive datasets, although updating these structures remains challenging [70]. Despite significant progress, challenges persist, notably in constructing and maintaining indexes for massive and evolving datasets, and achieving low-overhead rebalancing as data distributions evolve.

The principal approaches are summarized in Table 5, which illustrates the trade-offs between scalability, update efficiency, and distinctive strengths for large-scale spatial query processing.

## 4.3 Graph Analytics and Advanced Query Structures

The increasing prevalence of graph-structured data in sectors such as bioinformatics, social networks, and software engineering necessitates specialized query mechanisms that extend beyond classical spatial or string indexing paradigms. This subsection aims to clarify the current landscape and motivations for advanced indexing and querying in graph analytics, focusing on similarities and trade-offs among leading techniques.

Among the main requirements is the capacity to efficiently resolve similarity queries—such as those based on edit distance or subgraph containment—which remain computationally demanding. Recent advancements in succinct data structures, including q-gram trees with hybrid encoding, have demonstrated substantial reductions in index memory consumption compared to previous filtering methods, while maintaining or improving filtering effectiveness and query speeds [23]. These compact structures blend global and local filtering strategies (such as degree and label-based refinement), enabling efficient navigation of large candidate spaces for graph similarity search. For example, the approach in [23] uses only 5%–15% of the indexing memory of prior methods, introducing enhanced filtering via degree and label filters, and scales to exceptionally large datasets (up to 25 million graphs). However, such in-memory indexes may face practical challenges when applied to even larger or more dynamic databases, as the construction and maintenance costs can grow with increased data volatility.

In the context of directed acyclic graphs (DAGs), efficient querying is enabled by techniques that exploit order, level, and separator-based decompositions. These methods provide strong worst-case performance guarantees, achieving near-optimal query complexities even under adversarial conditions [58]. Specifically, the Partial Order Multiway Search (POMS) algorithm in [58] uses recursive partitioning to achieve a competitive ratio of $O(\log n)$ compared to the optimal, where $n$ is the number of vertices. This generalizes

classical tree search results and shows practical benefits for search models in settings such as debugging and distributed systems. On the other hand, while the theoretical guarantees are robust, practical deployment may be impacted by the computational cost of finding optimal partitions and adapting to evolving DAG structures.

The ongoing convergence of hybrid filtering, succinctness, adaptive partitioning, and strong competitive guarantees reflects broader trends in processing high-dimensional and irregular datasets. Recent work, particularly from 2021 onward [23, 58], highlights both notable progress and emerging trade-offs: highly compact indexes and near-optimal query performance are achievable, but with open challenges regarding adaptability, scalability, and computation overhead for highly dynamic or exceptionally large graphs.

In summary, recent advances in graph analytics offer substantial improvements in query efficiency and memory usage for similarity search and DAG querying. Nonetheless, trade-offs remain between index succinctness, adaptability, and computational overheads. Researchers and practitioners should evaluate these methods based on their specific data scales, update rates, and application requirements.

## 4.4 Unified Perspectives for kNN, Similarity, and Join Operations

A broad analytic perspective reveals that $k$-nearest neighbor (kNN), similarity, range search, and join operations can be interpreted as instances of a unified data retrieval paradigm, especially over large, heterogeneous, or multimodal datasets. Recent empirical syntheses highlight the confluence of several methodological directions:

**Spatial Partitioning:** Organizing the search space hierarchically to prune irrelevant regions. This includes primary and secondary partitioning methods that divide data according to geometric and attribute-based features, shown to accelerate range, kNN, and join queries by reducing unnecessary computations and avoiding duplicate reporting [5, 18, 47, 48, 97].

**Machine Learning-Based Index Construction:** Leveraging regression models, space-filling curves, and other data-driven techniques to predict locations and enhance lookup speeds. Recent learned index structures combine lightweight models, such as splines or neural nets, with classical spatial partitioners, as in LiLIS and LLM-powered index advisors [25, 45, 71]. These data-driven indexes provide superior throughput and index build efficiency for large and dynamic workloads, with the caveat that model training costs and complex query support remain open research challenges.

**Adaptive Query Evaluation:** Dynamically tuning the search process to accommodate data characteristics, distributional shifts, and incoming queries. Approaches such as online metric learning and index parameter adaptation allow for real-time or streaming adjustments, as explored in frameworks like OASIS, which incrementally update similarity functions and reuse index structures to maintain performance and reduce overheads [22, 80, 93].

**Ensemble and Subspace Techniques:** Combining multiple indexing or filtering strategies to mitigate high-dimensional challenges and exploit complementary strengths. Ensemble clustering and consensus methods aggregate results over feature subsets or algorithmic variants, which is especially effective for noisy, high-dimensional, or categorical data [21, 38, 44, 57]. Ensemble schemes,

**Table 5: Comparison of Space-Partitioning Index Strategies for Large-Scale Query Processing**

| Strategy | Scalability | Update Efficiency | Strengths |
|---|---|---|---|
| Tree-Based (e.g., R-tree) | Moderate (suffers in high dimension) | Moderate (requires rebalancing) | General-purpose; established theory |
| Grid-Based | High (especially with proper partitioning) | High (minimal restructuring) | Fast for point/range queries; low traversal overhead |
| Learned/Hybrid | Very High (adapts to data, $O(1)$ lookup) | High (can support incremental updates) | Handles skewed, high-dimensional data; efficient memory use |
| Compressed/Rep-indexes | High (suitable for repetitive data) | Moderate to Low (updates can be complex) | Dramatic space savings for redundant datasets |

parallelization, and local filtering collectively drive state-of-the-art performance for large and challenging similarity, kNN, and join problems [5, 18, 21, 22, 25, 31, 33, 38, 44, 45, 47, 48, 54, 58, 70, 71, 80, 93, 94, 97].

Robustness to noise and high dimensionality is further achieved through parallelism, compression, and ensemble models. Distributed frameworks have unified formerly distinct operations—such as kNN joins, range queries, and similarity joins—into single-session, high-throughput systems, minimizing I/O overhead and enabling resource-efficient knowledge discovery [21, 25, 57, 70]. The empirical record demonstrates that modern systems such as FML-kNN and LiLIS outperform prior MapReduce-style solutions, supporting scalable, robust retrieval across a variety of analytical tasks and data modalities.

For example, grammar-compressed and LCP-based indexes excel on highly repetitive string collections, outperforming naive approaches, but may introduce compromises in index construction time or incremental update capabilities [22, 33, 38, 44]. Meanwhile, machine-learned index structures and consensus-driven, parallelizable clustering approaches have substantially improved scalability and resilience to noise, albeit with increased algorithmic and training complexities [25, 57]. Distributed learned indexes and data-partition specific strategies enable efficient and accurate large-scale spatial and similarity search, particularly as data and workload heterogeneity increase [25, 71, 97]. As such, the methodological integration of space partitioning, local filtering, parallelization, and ensemble learning now underpins the state-of-the-art across similarity, kNN, and join algorithms in massive data environments.

Looking forward, key research challenges involve:

Supporting nonlinear and complex similarity functions, including those accommodating adaptive metrics [93]; Enabling non-parametric and domain-agnostic retrieval that generalizes robustly across workloads and data types [44, 45, 57]; Developing robust, fine-grained incremental index updates to support real-time and streaming scenarios [48, 54, 80]; Standardizing evaluation protocols for multimodal and streaming data, facilitating benchmark-driven progress [5, 38].

Addressing these challenges will catalyze the continued synthesis and advancement of indexing and querying approaches, fully adapted to the evolving demands of dynamic, large-scale, and heterogeneous data landscapes.

## 5 Dimensionality, Data Preprocessing, and Visualization

### 5.1 Data Types and Representational Variety

Modern data science contends with an expanding diversity of data types, including numeric, categorical, temporal, spatial, multimodal, compositional, incomplete, dynamic, and high-variance forms. This variety substantially informs the choice and design of analytical tools by shaping the assumptions underlying algorithmic methods. For example, numeric and continuous variables—ubiquitous across disciplines—facilitate a wide spectrum of quantitative manipulations. In contrast, categorical data, particularly in high-dimensional or sparse contexts as observed in omics or textual datasets, challenge direct statistical analysis and demand well-chosen encoding or embedding methods [72, 75, 99]. Specifically, nominal attributes often require encoding schemes that preserve class informativeness and allow valid correlation or distance-based interpretation [72].

Temporal and sequential datasets further complicate analysis due to the necessity of maintaining order dependencies, affecting similarity computation and clustering methodologies [3, 82]. Spatial data, such as those arising from medical imaging or geographic information systems, impose unique representational requirements that must strike a balance between fidelity, computational efficiency, and the preservation of connectivity or adjacency information [34, 73, 88, 109].

The prevalence of multimodal and compositional data in fields such as systems biology or sensor analytics magnifies these complexities. Compositional data, defined by components representing parts of a whole and summing to a constant, oblige the use of specific transformations—such as log-ratio methods—and purpose-built regression models to ensure inferential validity [62, 102, 104]. Additionally, the challenges posed by incomplete and dynamic datasets—including non-stationarity, time-varying drift, frequent updates, and deletions—necessitate adaptive preprocessing strategies capable of real-time reaction to evolving data [5, 28, 35, 100, 109, 110]. Data representations must also accommodate the practical realities of high variance and high dimensionality, which drive ongoing innovation in domains such as indexing, compression, and scalable embedding frameworks [5, 62, 88, 100, 109].

### 5.2 High-Dimensionality Challenges and Solutions

The widespread occurrence of high-dimensional data exacerbates both statistical and computational hurdles, encapsulated by the "curse of dimensionality." As dimensionality increases, the feature space grows exponentially, rendering conventional notions of distance less meaningful and impairing the performance of algorithms reliant on pairwise proximity [3, 82]. The resulting sparsity and noise accumulation compromise statistical power, heighten overfitting risks, and undermine clustering and learning efficacy. Classic distance metrics such as Euclidean and Manhattan distances, and kernel-based approaches, suffer from degraded discrimination in these settings, raising concerns for both exact and approximate k-nearest neighbor searches, high-dimensional clustering, and analyses of large-scale biological data [6, 38, 39, 57, 60].

To address these phenomena, methodologies that scale and adapt to high-dimensionality have emerged:

**Feature selection and dimensionality reduction:** Linear techniques such as Principal Component Analysis (PCA) and nonlinear methods like t-SNE and UMAP extract salient features and discard noisy or redundant ones [4, 55, 104]. However, recent work highlights that standard dimensionality reduction methods are often vulnerable to scattering noise, which can obscure cluster structures and reduce interpretability. The distance-of-distance (DoD) transformation, for example, has been shown to preprocess neighborhood distances to better separate noise from meaningful clusters in embeddings, significantly improving clustering accuracy, especially in very high-dimensional and low-sample regimes [55]. This approach demonstrates the need for advances specifically targeted at denoising and noise-induced artifact reduction during dimensionality reduction.

**Adaptive metric learning:** Tools including local Mahalanobis transforms and hierarchical subspace models enable more informative similarity calculations under small sample size relative to feature count ($p \gg n$) [39, 82, 104]. Adaptive metrics, such as in double-weighted k-nearest neighbor frameworks, allow the algorithm to individually weight features by informativeness, thereby mitigating the dominance of irrelevant dimensions and improving classification and clustering outcomes in high-dimensional settings [6]. Clustering performance also depends critically on data normalization and scaling; recent studies introduce scaling approaches based on multidimensional shape complexity to enhance the separation of clusters, albeit with additional computational cost [3].

**Ensemble subspace methods:** Aggregation over multiple random or systematically chosen low-dimensional projections mitigates overfitting and stabilizes models [57]. Consensus clustering methods, such as those employing co-association matrices and feature reweighting, have shown substantial improvements in accuracy and robustness against noise, particularly for categorical data with only a minority of informative features. Ensemble approaches also benefit high-dimensional regression and classification, as exemplified by ensemble Lasso or trimmed averaging techniques, which outperform traditional methods with complex, less sparse models and provide more stable results across a range of parameter choices [4, 57].

**Dynamic and streaming data analyses:** Incremental index structures, real-time clustering, and continuous normalization address the demands of evolving datasets [16, 28, 76, 102]. The emergence of learned multi-dimensional indexes introduces machine learning approaches directly into the indexing process, allowing for more adaptive and efficient querying in high-dimensional databases, though challenges in dynamic workloads and precise error bounding persist [60]. Novel data structures for dynamic and streaming updates in geometric and topological spaces have also become central for scalable processing [16, 28, 76].

Despite these advances, many high-dimensionality solutions display sensitivity to specific data distributions and parameterizations. Moreover, practical trade-offs between interpretability, computational cost, and robustness to noise persist as fundamental issues across application domains [22, 50, 102]. Ongoing research also emphasizes the challenge of statistical inference under high-dimensions, such as testing mean vectors or symmetry, particularly with missing data or small sample sizes [22, 50, 102]. The ongoing quest to generalize methods robustly across diverse modalities and to guarantee interpretable, meaningful low-dimensional representations remains central to current research.

## 5.3 Preprocessing and Normalization

Data preprocessing is foundational to robust and reliable analytics, particularly when analyzing high-dimensional, heterogeneous, or noisy datasets. The main objectives are to mitigate noise and outlier effects, normalize feature scales, and ensure compatibility with downstream models.

Standard normalization methods, such as min-max scaling, z-score standardization, and variance-stabilizing transforms, aim to harmonize feature ranges. However, these approaches may be inadequate when faced with outlier-prone or heavy-tailed distributions, or when compositional constraints are present [3, 35, 73, 95]. For compositional data, specialized transformations like log-contrast or isometric log-ratio are necessary to avoid spurious correlations resulting from constant-sum constraints [50, 104]. Notably, [35] addresses preprocessing challenges for nominal and mixed-type attributes by encoding categorical data numerically, allowing for inclusion in quantitative analyses and downstream clustering, even when classical statistical measures may not be well-defined.

The selection of normalization procedure can significantly influence the interpretability and performance of clustering and classification. Recent work [3] proposes determining feature-wise scaling factors by optimizing shape complexity, emphasizing the importance of balancing intra- and inter-cluster distances before clustering. This approach demonstrates improved clustering performance over traditional normalization, particularly in ambiguous scenarios, though it can require expert intervention. In the context of time-series data, careful consideration of normalization choices is likewise critical, as noted in recent comparative taxonomies [73].

Robust outlier detection and adjustment remain crucial in preprocessing, as these steps have a substantial impact on analytical outcomes. Nevertheless, there is a scarcity of standardized benchmarks for evaluating outlier handling and noise mitigation efficacy, underscoring the importance of transparent and reproducible preprocessing pipelines. Domain-specific customization is frequently needed, especially for normalization and duplicate management [8, 28, 38, 50, 100, 104, 110]. The literature demonstrates that modern k-nearest neighbor (kNN) based methods actively incorporate noise-resilient feature selection and adaptive weighting to achieve improved classification on imbalanced or noisy datasets [8, 38]. Furthermore, methods such as consensus spectral clustering on high-dimensional categorical data illustrate the effectiveness of integrating feature reweighting schemes based on informativeness, resulting in robust performance under both stochastic and adversarial noise, even when only a small fraction of features are informative [57].

In streaming and dynamic data environments, preprocessing must address additional challenges: algorithms should efficiently assimilate new information, adapt to evolving data distributions (concept drift), and handle data deletions or reweighting without necessitating a full retraining of models [5, 16, 28, 76]. These constraints are particularly prominent in real-time analytics and online

learning, where the demand for balancing statistical rigor with computational efficiency is ever-present.

## 5.4 Dimensionality Reduction and Visualization Techniques

Dimensionality reduction and visualization remain fundamental tools for extracting informative low-dimensional representations from complex datasets, elucidating latent structures, and supporting both exploratory and inferential analysis. Principal Component Analysis (PCA) and its advanced variants, such as generalized contrastive PCA (gcPCA) and contrastive PCA (cPCA), are central techniques. Unlike conventional cPCA, which requires tuning a sensitive hyperparameter, gcPCA introduces a normalization strategy that penalizes high-variance dimensions, ensuring robust and interpretable separation of patterns across contrasting datasets without hyperparameter dependence [4, 29]. These extensions support accurate distinction of signals across experimental conditions, robust to noise and rank deficiency.

Nonlinear embedding tools, notably t-SNE and UMAP, are widely adopted for visualizing cluster and manifold structures in high-dimensional domains including single-cell genomics, neuroimaging, and image analysis. However, they can be confounded by scattering noise, where random fluctuations obscure cluster boundaries in low-dimensional projections [29]. The distance-of-distance (DoD) transformation addresses this by processing the original distance matrix—by emphasizing differences in local neighborhood distances—to sharpen cluster detection and separate noise as distinct clusters. This method demonstrably improves cluster recovery and classification fidelity (measured via ARI), particularly in high-dimensional or noisy regimes, though it introduces computational overhead and relies on appropriate parameter tuning [55].

Supervised dimensionality reduction and feature selection are advanced not only by penalized regression methods, such as Lasso, Elastic Net, and Adaptive Lasso, but also by ensemble subspace approaches. Ensemble subspace methods, like ensemble Lasso or consensus spectral clustering, aggregate results across diverse randomly-selected feature subsets and employ strategies such as trimmed means or majority voting, achieving higher robustness and accuracy—especially in challenging $p \gg n$ or weak signal scenarios. They are especially effective in the presence of correlated features, noise accumulation, or when only a small subset of features is informative, as confirmed in extensive theoretical and empirical studies [4, 57].

Visualization frameworks have expanded beyond traditional scatterplots to include specialized representations of clusters, graphs, tensors, and multidimensional networks. For example, Flowcube employs interactive matrix-based representations enhanced by direction-based filtering lenses for elucidating complex geographic flows without excessive clutter, enabling richer user-driven pattern discovery in large, spatial datasets [90]. Methods grounded in tensor decomposition and model-based clustering, such as the tensor normal mixture model combined with penalized likelihood and doubly-enhanced EM algorithms, deliver parsimonious and interpretable compression and unsupervised grouping for high-order, structured data with rigorous statistical guarantees [59]. Where data are inherently multiway or graph-structured, compact encodings

and algorithmic advances further support scalable and interactive exploration [5, 16].

Interpretability, transparency, and reproducibility are increasingly emphasized in the field. Contemporary approaches feature explainable feature allocation, use of cluster validity indices, rigorous benchmarking standards, and accessible software tools, ensuring traceability and practical adoption [16, 29, 62, 66, 90, 104]. Nevertheless, persistent challenges include producing consistent embeddings across runs, mitigating batch effects, and scaling methods to ever-larger and more heterogeneous multimodal datasets.

## 6 Feature Selection, Classification, and Vector Modeling

### 6.1 Feature Ranking and Robust Classification

Feature selection and classification in high-dimensional domains—especially when both the number of variables ($p$) and observations ($n$) are large and the data presents high variance or outlier contamination—have undergone substantial advances in recent literature. Key challenges include maintaining robustness to outliers, ensuring interpretability, and achieving computational scalability. Traditional classifiers frequently encounter difficulties in settings where class separation is predominantly due to variance differences or where outliers heavily influence the results.

Recent frameworks incorporate rank-based and subsampling strategies to address these issues. Notably, rank-based classification methods rely on transforming pairwise distances between observations into rank information, enabling classification procedures that are resilient to distributional assumptions and offer enhanced robustness against outlier effects [65]. These frameworks follow a systematic approach composed of (i) computing distance matrices among samples, (ii) applying rank transformations to the distances, and (iii) integrating rank-derived features into classifiers such as quadratic discriminant analysis, facilitating class separation in challenging high-dimensional scenarios. The default use of the $\ell_2$ distance can be adapted to alternate metrics, improving applicability to data with network or non-Euclidean structure, while multi-class problems are handled by direct extension of the methodology.

Complementing rank-based methods, efficient subsampling strategies tailored for high-dimensional settings have also been developed [20]. For example, recent approaches first utilize a random LASSO-based selection to identify a sparse set of active predictors and then employ leverage-score-informed subsampling to select observations that capture the essential structure of the data. This two-stage framework leads to more accurate variable selection and reduction in predictive error (e.g., mean squared prediction error), while substantially reducing computational demands—particularly in large-scale or $p > n$ contexts. Sensitivity analyses reported in these studies confirm that such procedures are robust across a range of algorithmic parameters.

Empirical evidence, including extensive simulation studies and real-world applications such as sentiment classification and blog post comment prediction, shows that both rank-based classifiers and subsampling-based feature selection frameworks can match or surpass state-of-the-art alternatives. This is especially apparent in settings requiring noise resilience or flexible handling of unconventional feature distributions [20, 65]. Nevertheless, several

challenges remain, including algorithmic scalability for extremely large datasets and the principled selection of distance metrics for rank-based methods. These persist as active research directions in high-dimensional robust classification.

## 6.2 Nonparametric and Subdata Selection Methods

This subsection surveys recent advances in nonparametric and subdata selection methods, with emphasis on their objectives and performance in modern high-throughput data analysis settings.

Nonparametric approaches and advanced subdata selection techniques are indispensable for analysis where both the number of observations ($n$) and features ($p$) can be extremely large. The primary objective in this context is to enable robust variable selection and estimation under severe computational and statistical challenges. Traditional LASSO-based selection suffers from diminished efficacy when predictors are highly correlated or the dimensionality is extreme ($p \gg n$), motivating dual-stage frameworks that combine variable screening with informed subdata selection.

A representative dual-stage procedure consists of an initial random LASSO step for variable screening, followed by leverage-score-based sampling to select the most informative data points for estimation. This strategy addresses limitations of classic LASSO, improving both estimation accuracy and computational efficiency, according to recent simulation studies and real applications [20, 36].

To clarify the comparative benefits and limitations of single- versus dual-stage approaches, Table 6 summarizes their core features and trade-offs:

As highlighted above, dual-stage procedures yield systematic improvements in robustness and estimation accuracy compared to conventional single-stage or unstructured methods. However, dual-stage approaches may require additional parameter tuning and careful design, as detailed in recent sensitivity analyses and robustness studies [20]. In scenarios where full data analysis is computationally infeasible, such frameworks are particularly valuable. Remaining challenges for the field include balancing information gain and computational cost, and ensuring scalability as both $n$ and $p$ continue to grow.

Methodological innovation also extends to regression with high-dimensional compositional covariates, leading to hierarchical, mixed, and p-value-free false discovery rate (FDR) control schemes. These newer methods exploit symmetry properties of test statistics under the null, enabling valid inference when conventional significance tests fail due to dimensionality or complex predictor correlation. Theoretical developments secure strong FDR control and optimal asymptotic power, with confirmed practical gains in both synthetic and omics settings [20].

Further, penalized likelihood estimation for high-dimensional mixed-effects models has advanced through coordinate descent algorithms using nonconvex penalties such as SCAD. Results consistently show that SCAD offers superior variable selection accuracy and reduced estimation bias versus LASSO, particularly under high correlation or group structures present in omics and GWAS data [36]. Despite these advances and the availability of open-source implementations, hurdles remain: convergence guarantees for non-Gaussian responses, accelerated parameter tuning, and robust uncertainty quantification are ongoing research areas.

In sum, nonparametric and subdata selection methods have seen significant improvement in both methodology and application, especially with the integration of dual-stage variable and subdata selection. These innovations demonstrate notable trade-offs between statistical efficiency, computational resource requirements, and robustness to complex data structures, and ongoing research continues to address their practical and theoretical challenges.

## 6.3 Statistical Testing in High Dimensions

**Section Objective:** This subsection aims to survey statistical testing methodologies tailored for high-dimensional data, address their practical limitations, and highlight open research challenges in benchmarking and validation for such settings.

Statistical inference in high-dimensional environments requires robust procedures that remain effective when $p$ is large relative to $n$ and dependencies among features are significant. Classical mean vector testing methods—including Hotelling's $T^2$ statistic—deteriorate in reliability as dimensionality increases, resulting in inflated type I error rates and poor statistical power. To overcome these limitations, U-statistic-based techniques have been introduced for one- and two-sample testing paradigms, providing test statistics that converge to $t$-distributions as $p$ becomes large relative to fixed $n$ [22, 50]. These tests obviate the need for resampling or complex adjustments, delivering direct and reliable inference for applications such as neuroimaging and genomics where "large $p$, small $n$" is prevalent. Simulation studies demonstrate strong type I error control and power in both finite-sample and asymptotically high-dimensional scenarios [50].

Recent developments have addressed critical practical aspects, such as missing data and the computational burden of ultra-high dimensions. For high-dimensional mean testing in the presence of missing observations, new test statistics accommodate data missing at random, establishing asymptotic guarantees as both $n$ and $p$ grow [102]. Random projection-based methods exploit the concentration of measure to efficiently estimate null distributions, providing computationally tractable and statistically valid inference in extremely high-dimensional data [22, 39, 104]. However, these approaches often involve trade-offs: while projections can reduce variance and computational cost, they may also obscure weak signals or inflate type II error rates, especially under strong correlations [39].

Multiple comparison corrections and the assessment of clustering validity have also evolved. Contemporary studies stress that proper alignment of statistical assumptions and hypothesis structures is essential for balancing statistical power and controlling error rates [43, 50, 60]. However, benchmarking high-dimensional cluster validity remains challenging, as existing metrics may not generalize well across structured and unstructured data. A prominent open issue is the lack of universally applicable, absolute cluster validity criteria that remain robust in the presence of noise and highly heterogeneous feature sets [43]. In practice, simulation-based benchmarking remains the primary means of interpreting comparative performance, but it suffers from limited reproducibility and sensitivity to generative model assumptions [39, 50].

**Table 6: Comparison of Traditional and Dual-Stage Subdata Selection Methods**

| Method | Variable Selection Stage | Subdata Selection Stage |
| --- | --- | --- |
| Standard LASSO | Single-stage (Lasso only); less robust under high correlation or $p \gg n$ | No explicit subdata selection; potentially inefficient for large $n$ |
| Dual-Stage (Random Lasso + Leverage) | Randomized Lasso for variable screening, enhancing robustness to correlation | Leverage-score sampling finds informative data points, reducing computation and improving accuracy |

Analytic strategies increasingly leverage ensemble and consensus clustering frameworks in high-dimensional contexts. By integrating dimension reduction, feature reweighting, and consensus techniques, these approaches boost robustness to noise and adversarial corruptions. The combination of one-hot encoding, random projection, and spectral consensus mechanisms has been shown to enhance cluster recovery, but at the cost of increased computational demands [50, 104]. Notwithstanding potential for parallelization and empirical success, optimal parameterization and theoretical guarantees in genuinely heterogeneous datasets remain insufficiently explored.

**Example Use Case:** Consider a neuroimaging application with $p \approx 10,000$ voxel-based features and $n$ as few as 20 subjects, where there is missing data in some observations. U-statistic-based $t$-tests [50] can directly test mean differences without requiring resampling, while projection-based inference [22, 104] further accelerates hypothesis testing. For post-hoc cluster validity, practitioners may combine ensemble clustering with absolute validity indices [43], but must be cautious in interpreting results beyond the scope of controlled simulations.

**Open Problems and Research Challenges:** As shown in Table 7, fundamental challenges remain in p-value calibration, testing with missing or incomplete high-dimensional data, robust benchmarking and validation of clustering results, and scalable computation for increasingly complex frameworks. Improving benchmark reproducibility and designing universally applicable validity measures for highly heterogeneous or categorical features are particularly pressing research frontiers.

## 6.4 Vector Space and Distributional Semantic Models

This subsection aims to elucidate the fundamental objectives and current advances in vector space and distributional semantic models, emphasizing the balance among interpretability, predictive accuracy, and scalability in high-dimensional contexts.

Semantic modeling in high-dimensional linguistic or biological contexts demands vector representations that attain an effective balance among interpretability, predictive accuracy, and computational tractability. Distributional semantic models—which encode entities as vectors in high-dimensional spaces—have been instrumental in capturing semantic relationships. Leading approaches such as neural embedding models (e.g., word2vec) and matrix factorization methods (e.g., NMF) provide high predictive accuracy; however, their dense representations typically lack dimension-wise interpretability.

Recent innovations address this shortcoming by proposing dimension selection procedures that directly map naturally occurring attributes (such as specific words) onto dimensions, enabling both interpretability and high accuracy. For instance, Pakzad et al. [64] present a method that selects a subset of the most frequent words as basis dimensions for embedding, achieving interpretable vector spaces without significant sacrifices in accuracy. Applied to the ukWaC corpus, the authors obtained a vector space of $N = 1500$ basis words, and demonstrated on several benchmark datasets (MEN, RG-65, SimLex-999, WordSim353) that reducing basis vectors from 5000 to 1500 incurs only about a 1.5%–2% loss in accuracy, while dramatically improving interpretability. Moreover, interpretability assessments confirm the superiority of these word-based vectors over neural embeddings and NMF baselines.

In the area of database indexing and retrieval, vector model construction leveraging machine learning techniques—including clustering, neural networks, and hybrid systems—underpins efficient and adaptive multi-dimensional index structures. These advances facilitate scalable querying, indexing, and retrieval, accommodating the demands of large, continuously evolving datasets [43].

Despite these advances, several open research challenges persist. Key issues include (i) developing principled metrics for evaluating interpretability in high-dimensional vector spaces, (ii) balancing the trade-off between accuracy and transparency, particularly as the dimensionality is reduced, and (iii) constructing indexing and retrieval systems that adapt efficiently to data scale and distributional shifts. Addressing these challenges will be essential for further progress in scalable, informative, and interpretable analysis within high-dimensional semantic frameworks.

## 7 Benchmarking, Evaluation, and Cluster Validation

This section provides an overview of the main objectives and challenges in the benchmarking, evaluation, and validation of clustering methods. We aim to clarify the key concepts, present standard methodologies, and highlight persistent open problems—especially in high-dimensional settings.

Benchmarking refers to the systematic comparison of algorithms under a variety of controlled scenarios, using standardized datasets and metrics. Evaluation addresses the question of how well a clustering output matches either external ground truth or internal consistency notions, through metrics such as adjusted Rand index, silhouette score, or stability measures. Cluster validation investigates the reliability and statistical significance of discovered clusters, often when ground truth is unavailable.

A central challenge remains the development and selection of meaningful evaluation metrics—particularly as the dimensionality and heterogeneity of data increase. Limitations around metric selection, sensitivity to initialization, and interpretability persist in high-dimensional regimes.

Example use cases include assessing the stability of clustering results on gene expression profiles, or benchmarking clustering algorithms on text corpora with varying topics and vocabulary sizes. In these scenarios, practitioners must carefully select both their evaluation metrics and validation strategies.

**Table 7: Summary of Open Problems in High-Dimensional Statistical Testing and Validation**

| Research Challenge | Limitation | Representative Reference |
|---|---|---|
| Calibration of p-values in ultra-high dimensions | Classical estimators often inflate type I errors for $p \gg n$ | [50] |
| Testing under missing at random (MAR) mechanism | Robustness to missing data requires new test derivations | [102] |
| Cluster validity in heterogeneous, noisy datasets | Few absolute cluster metrics work across data types | [43] |
| Benchmarking and reproducibility of new methods | Evaluation depends on simulation design; lacks consensus benchmarks | [39, 50] |
| Systematic integration of projection-based tests | Risk of masking weak signals or introducing bias in correlated data | [22, 39, 104] |
| Parallel computation for ensemble or consensus testing frameworks | Optimal design and scalability are underexplored in large, real-world settings | [50] |

**Table 8: Overview of Key Open Problems in Clustering Benchmarking and Evaluation**

| Thematic Area | Open Problem | Description |
|---|---|---|
| Benchmarking | Dataset Diversity | Lack of standardized, diverse benchmarks in high-dimensional and heterogeneous settings. |
| Evaluation | Metric Robustness | Selecting or designing evaluation measures resilient to dimensionality and noise. |
| Cluster Validation | Statistical Guarantees | Validating discovered structure without access to ground-truth labels. |
| Cross-cutting | Integration of Methods | Lack of unified frameworks or taxonomies linking benchmarking, evaluation, and validation approaches. |

At the conclusion of each main subsection within this section, we summarize the prominent research challenges in that thematic area. We also propose a stylized taxonomy that integrates validator, benchmarking, and representation-based approaches, to guide future development in this field.

The following subsections elaborate methods, metrics, and systems for benchmarking and validating clustering algorithms, with improved transitions and explicit objectives for each part.

## 7.1 Cluster Validation and Evaluation Metrics

Robust cluster validation underpins the scientific credibility and reproducibility of unsupervised learning methodologies. Two primary paradigms exist for validating clustering results: internal (absolute) and external (relative) measures. Internal validation indices—such as the Silhouette coefficient, Dunn index, and Davies-Bouldin score—evaluate clustering quality without recourse to ground truth labels. These methods efficiently quantify cluster compactness and separation, yet they can be influenced by noise, feature scaling, and data dimensionality. Notably, in high-noise or high-dimensional contexts, these metrics often struggle to distinguish true structure, potentially misrepresenting clusterability, particularly when faced with chaining artifacts, small clusters, or overlapping densities [1, 4, 6, 8, 17, 19, 21, 31, 36, 39, 42, 43, 47, 50, 62, 66, 73, 84, 86, 89, 95, 100–102, 104, 111]. Caution is thus advised against relying solely on internal metrics for conclusive assessment [17, 102].

External indices—including Adjusted Rand Index (ARI), Normalized Mutual Information (NMI), F1-score, and Cohen's Kappa—compare the computed clustering to reference labels, providing more interpretable and objective benchmarking, especially when a gold standard exists [4, 6, 11, 17, 19, 21, 31, 36, 39, 42, 43, 47, 50, 62, 66, 73, 84, 100–102, 110, 111]. Among these, ARI and NMI have consistently demonstrated robustness and discriminative power, outperforming less nuanced measures such as purity [11, 42]. Nevertheless, these metrics possess their own biases, for example towards certain cluster size distributions and cluster counts—a challenge accentuated in multiclass, imbalanced, or high-dimensional data.

Moreover, metrics such as precision at top-$n$ (P@$n$) and area under the ROC curve (AUROC), prevalent in related settings like outlier detection, require meticulous adjustment for dataset imbalance and sampling artifacts to avoid misleading results [17].

Recognizing such limitations, recent advancements have introduced more context-sensitive and adaptive validation strategies. These include indices leveraging correlations between within-cluster and centroid distances, which can highlight multiple plausible clustering solutions, aligning more effectively with real, often hierarchical, data structures [42]. Additionally, multimodality tests—including the Dip test and Silverman's test applied to pairwise distances—offer robust, distribution-agnostic assessments of clusterability, generally outperforming classic indices in differentiating between true signal and noise, though challenges remain for specialized structures such as heavy chaining or tiny clusters [50, 102].

Finally, parameter sensitivity and data preprocessing practices—such as normalization, scaling, and duplicate handling—strongly influence metric reliability. Accordingly, adaptive, data-driven feature scaling procedures and robust software implementations are increasingly integral in contemporary clustering analyses [50, 66].

As summarized in Table 9, selection of appropriate validation metrics must account for data characteristics, application context, and the availability of ground truth labels. No single approach suffices across all scenarios; thus, rigorous studies routinely report multiple metrics and qualify their interpretations.

## 7.2 System-Level and Analytic Metrics

Beyond clustering quality per se, practical deployments—especially at scale—demand careful evaluation of system-level characteristics that impact both efficiency and usability. Key considerations include query efficiency (latency and throughput), memory consumption, index construction/maintenance overhead, scalability with respect to dataset size and dimensionality, and robustness to evolving or noisy data. As datasets grow to millions or billions of points, these aspects become ever more critical for real-world viability [1, 3, 5, 8, 14, 16, 23, 32, 35, 57, 60, 70, 73, 76, 82, 83, 95, 96, 100, 110, 112].

**Table 9: Common Cluster Validation Metrics: Key Properties and Use Cases**

| Metric Type | Example Metrics | Requires Ground Truth | Key Strengths / Limitations |
|---|---|---|---|
| Internal (Absolute) | Silhouette, Dunn, Davies-Bouldin | No | Fast; sensitive to noise/dimensionality; may not detect overlap/chaining |
| External (Relative) | ARI, NMI, F1-score, Cohen's Kappa | Yes | Interpretable w/ labels; can be biased by cluster count/imbalance |
| Multimodality/Novel | Dip test, Silverman's test | No | Robust to noise; less sensitive to structure; challenges in special cases |

**Latency and Throughput:** Real-time analytics and search systems require low query latency and high throughput. Recent advances in approximate nearest neighbor search (ANNS) leverage optimized memory layouts, vector quantization, and adaptive parameter tuning to efficiently balance speed and accuracy [1, 70, 95, 96, 112]. For instance, improved prefetching, cache-aware vector layouts, and selective distance computation [100, 110] can significantly decrease response times for high-dimensional queries. Moreover, space-efficient indexing methods such as succinct structures or compressed representations [14, 16, 23, 32, 76, 83] enable rapid lookups without compromising overall precision, particularly when combined with automation in parameter tuning and hybrid encoding strategies.

**Scalability:** The ability of an algorithm or system to handle increasing amounts of data and higher feature dimensionality is typically measured by its memory footprint and algorithmic time complexity. Scalability is especially relevant for applications involving massive datasets, where hybrid indexing techniques, adaptive partitioning, automaton-based compression, or hierarchical pruning are applied to achieve sublinear or near-linear scaling [1, 14, 16, 23, 32, 57, 76]. Community benchmarks have validated that graph-based and quantization-focused methods can robustly index and search across billions of vectors, with adaptive constructions like learned or tried-based indices further extending to dynamic and high-dimensional data regimes [14, 60].

**Resource Constraints:** For edge devices and embedded environments, rigorous evaluation under memory and computational limitations is essential. To this end, metrics such as "peak $n$ per query" or "memory units per million points" [5, 16, 83, 95] are adopted for fair benchmarking and standardized comparison. Systems based on compressed filters, compact hashing, and specialized data structures (such as windowed Cuckoo filters for online insertion [83]) achieve high memory efficiency and fast operations, supporting scalable and flexible deployment even in constrained hardware settings.

**Robustness:** Adaptive clustering and search frameworks are increasingly assessed for their resilience in the presence of data drift, adversarial noise, or distributional changes. Recent studies highlight the need for algorithms that maintain accuracy and responsiveness as data distributions evolve—either by supporting streaming or online updates [14, 60, 82], or through the design of robust ensemble and consensus techniques [57, 112]. Specialized adaptiveness metrics supplement standard accuracy measures to capture the systems' ability to cope with shifts and perturbations in underlying data.

Systematic benchmarking thus now requires comprehensive reporting of both analytic and engineering metrics, including timing, space usage, accuracy, and robustness across a broad spectrum of operational settings. This holistic evaluation enables practitioners to better assess practical trade-offs and informs the selection or development of clustering and search solutions suitable for large-scale, real-world systems.

## 7.3 Benchmarking Environments and Open-Source Tools

Transparent and reproducible evaluation is anchored in open, standardized benchmarking ecosystems that encompass both implementations and curated datasets. The proliferation of open-source libraries—spanning Python, R, and, more recently, Julia—has increased access to advanced clustering, indexing, and similarity search techniques [5, 11, 17, 21, 22, 29, 32–34, 36, 39, 50, 51, 60, 76, 81–84, 87, 92, 93, 101, 102, 104, 106–108, 111].

**Dataset Repositories:** Resources such as the UCI Machine Learning Repository and OpenML furnish benchmarks across diverse data types, annotated with task-specific and preprocessing information, providing essential transparency and enabling systematic comparative studies across clustering and similarity search tasks [17, 32, 39, 76, 83, 104].

**Simulation Frameworks:** Tools for controlled manipulation of data properties—including class separation, noise, and dimensionality—aid the rigorous comparison of algorithmic performance and highlight sensitivity to dataset idiosyncrasies. These environments support experiments in both synthetic and real-world settings, as demonstrated in studies of high-dimensional mixed-effects models and contrastive data analysis [29, 36, 60, 87].

**Domain-Specific Libraries:** Specialized implementations cater to modalities such as time series, trajectories, and point clouds, offering tailored distance functions and evaluation protocols. Such domain-aware tools are crucial for addressing the unique analytical challenges found in high-dimensional and structured data, as evidenced by scalable and distributed methods for kNN joins and tests for challenging data scenarios [21, 39, 102, 104].

**Reproducibility Artifacts:** The growing adoption of public code repositories, leaderboards, and open challenge datasets has strengthened verifiability of benchmarking efforts. Increasingly, best practices prioritize publishing full experimental configurations, including random seeds, preprocessing scripts, and evaluation parameters, to advance community-wide interpretability and reproducibility [11, 17, 22, 29, 32, 50, 76, 82].

Despite progress, substantial challenges persist. The field continues to lack universally recognized benchmark suites that reflect the complexity of emerging applications, such as those involving graph, tensor, or mixed-type data [60, 81, 84, 92, 106]. Moreover, robust simulation and evaluation environments addressing data corruption, anonymization, and missingness are still needed, with recent research highlighting the importance of modeling such phenomena explicitly [29, 50, 66, 102]. Continued community investment

is essential to curate, annotate, and standardize benchmarks that mirror the intricacies of real-world clustering tasks.

## 7.4 Visualization for Evaluation and Transparency

Visualization is indispensable for both evaluating and communicating the results of clustering and similarity search, bridging the gap between algorithmic output, expert assessment, and end-user trust. The use of 2D and 3D visualization remains fundamental for exploratory analysis, while interactive dashboards now constitute standard practice for both method development and result dissemination [5, 16, 22, 29, 50, 59, 64, 90, 102, 104].

Modern frameworks commonly integrate dimensionality reduction methods—such as t-SNE, UMAP, and contrastive PCA—not merely as visualization tools but also as preprocessing steps that surface latent structure otherwise masked in high-dimensional data. For instance, recent methods like generalized contrastive PCA (gcPCA) enable the robust extraction of low-dimensional patterns that distinguish between complex high-dimensional biological datasets, offering interpretable axes of variation without the ambiguity of parameter tuning or ad hoc decision-making. This approach has proven especially effective in uncovering subtle structure and heterogeneity in data such as neural recordings or gene expression matrices [29]. Dimensionality reduction techniques also serve as foundational tools for validating cluster separation, detecting anomalies, and identifying ambiguous or overlapping subpopulations, where they supplement quantitative metrics with intuitive, expert-driven insight [50, 59, 102, 104].

Emerging solutions employ innovative interaction paradigms. For instance, systems that visualize flows or spatial networks, such as Flowcube, enhance the exploration of origin-destination data by introducing interactive spatial filters that reveal nuanced arrangement patterns across large-scale datasets. The Direction-Based Filtering (DBF) lens in Flowcube, for example, systematically isolates flows along specific spatial directions, unveiling subtle concentration, alignment, or dispersal patterns in geographic movement data that might be missed by conventional automated methods. The user-driven, interactive exploration supported by such paradigms enables the discovery of both prominent and subtle structure in complex datasets, as demonstrated through real-world movement analyses in urban environments [90].

To ensure transparency and reproducibility, the latest best practices emphasize coupling visualization with systematic, script-driven analytics [29, 64, 102]. Open interfaces, reproducible color mappings, and support for exporting or replaying visualization states significantly strengthen interpretability and facilitate peer verification of findings [5, 64]. For example, approaches that select interpretable feature sets or basis vectors, as in interpretable word embedding models, enhance the transparency of downstream visual analyses by ensuring that each dimension corresponds to an easily understood concept, aiding both validation and presentation [64].

However, visual analytics also confront significant and ongoing limitations. These include challenges with ultra-high-dimensional or massively multi-class data, and situations where intrinsic data structure resists intuitive mapping. Such bottlenecks drive research into mixed-modality and interactive visualization methods capable

of scaling alongside increasing analytical and data complexity, as well as the development of advanced frameworks and open-source tools that address the combinatorial and presentational challenges posed by high-resolution or semantically rich datasets [5, 29].

## 8 Data Representation, Storage Optimization, and Hardware Acceleration

This section aims to systematically analyze foundational approaches and current advances in data representation, storage optimization, and hardware acceleration, highlighting their interdependencies and practical impact on AI systems. We begin by outlining the key objectives: (1) to compare techniques for representing high-dimensional AI data; (2) to evaluate storage efficiency strategies; and (3) to examine hardware acceleration mechanisms and their synergy with data structures.

The subsequent discussion transitions from methods for data modeling and compression, through common storage architectures and optimization targets, to hardware designs directly supporting AI workloads. Our goal is to connect these themes, clarifying how representational, software, and hardware elements interact in modern benchmarks and deployments.

Each major subsection will explicitly state its focus and scope, aiding reader orientation. At the conclusion of this section, we synthesize open research questions—particularly concerning evaluation metrics in high-dimensional settings and the integration of representation, validation, and benchmarking methods—within a unified perspective.

Open challenges in this domain include: - Developing scalable metrics for evaluating representation fidelity and downstream task relevance in massive, high-dimensional datasets. - Balancing trade-offs between compression ratio, access latency, and energy efficiency in heterogeneous storage systems. - Enabling adaptive hardware-software co-design that remains robust to shifts in model architecture and data modality.

Finally, Table 10 systematically summarizes principal open problems addressed in this section.

## 8.1 Data Representations for High-Dimensional Analytics

The analytical landscape for high-dimensional and multimodal data demands representations that are both expressive and computationally efficient. Classical strategies have relied on dense, grid-based formats for regular domains; however, in higher dimensions, the storage and computational requirements rapidly become prohibitive, driving the need for more advanced data structures that harness inherent sparsity and structural regularities characteristic of scientific and industrial datasets. Voxel-based encodings, which extend regular grid representations, remain prevalent for 3D spatial data due to their implementation simplicity and direct storage mapping. However, as dimensionality grows or data become increasingly sparse, these encodings exhibit significant memory inefficiency [5].

To mitigate these shortcomings, hierarchical structures have gained prominence. Examples include sparse voxel octrees (SVO), serialized directed acyclic graphs (SVDAG), and a spectrum of dynamic data structures such as OpenVDB, NanoVDB, SPGrid, and

**Table 10: Summary of principal open research challenges in data representation, storage optimization, and hardware acceleration.**

| Area | Open Problem | Key Considerations |
|---|---|---|
| Representation | High-dimensional evaluation | Scalability, metric interpretability, task alignment |
| Storage Optimization | Trade-offs in compression and latency | Efficiency, energy cost, reliability |
| Hardware Acceleration | Co-design for diverse workloads | Flexibility, future-proofing, integration complexity |

DT-Grid. These representations can dramatically reduce memory requirements—often by orders of magnitude—while preserving the capacity for locality-sensitive computations and supporting real-time manipulation [5]. Manifold-based approaches further extend this paradigm by succinctly capturing topological and geometric features, supporting advanced analytical tasks across fields such as computer graphics, computational biology, and scientific simulation.

The effectiveness of these data structures in high-dimensional contexts centers on several critical trade-offs. For instance, static memory layouts such as contiguous arrays deliver high throughput in batch-analytic or streaming scenarios but can be inflexible for adaptive or interactive applications. In contrast, dynamic memory layouts support interactive and adaptive analytics, yet introduce complexity due to sophisticated concurrency controls needed to maintain consistency and performance. Furthermore, hierarchical representations (such as SVO or SVDAG) greatly optimize storage for sparse datasets but at the cost of increased pointer overhead and potentially slower random access, particularly as spatial resolution grows. Dynamic grid solutions like OpenVDB and NanoVDB distinguish themselves by supporting efficient updates and parallelization, but may require substantial engineering effort for integration into existing high-performance workflows [5].

Despite their promise, challenges remain significant. The lack of standardized benchmarks and robust open-source libraries inhibits the objective evaluation of data structures across real-world scenarios. Many current implementations show reduced efficiency when dealing with non-watertight models or datasets that require rich semantic annotations. GPU acceleration, essential for volumetric analytics at scale, is not uniformly supported, which further restricts practical usability. Additionally, a critical obstacle is the absence of mature solutions for efficiently handling extreme resolutions and managing the trade-off between granularity and performance in dynamic or semantically complex contexts [5].

**Open Problems and Research Challenges:**

*Standardized benchmarking and validation present persistent challenges for high-dimensional data analytics. Key unresolved problems include:*
- *Establishing widely-used, standardized datasets and performance metrics for rigorously comparing novel data structures and algorithms.*
- *Developing robust and efficient libraries capable of supporting non-watertight and semantically annotated data in dynamic scenarios.*
- *Achieving scalable GPU-ready implementations that can manage large and heterogeneous volumetric datasets at extreme resolutions.*

- *Designing tools that facilitate the annotation and tracking of semantic information within hierarchical or sparse data representations.*
*Progress in these areas is crucial for enabling reproducible research and for the operational integration of advanced representations into scientific, industrial, and graphics pipelines.*

### 8.2 Space-Efficient Storage Structures

This subsection surveys recent innovations in storage structures that aim to overcome memory and I/O limitations in high-volume, high-dimensional analytics. The objective is to synthesize the main classes of space-efficient structures—probabilistic filters and compressed indexes—highlighting their underlying design choices, trade-offs, and practical relevance for handling massive modern datasets. We focus on guiding questions such as: What are the core algorithmic principles enabling space reduction? How do these structures balance false positives, update capabilities, and adaptability for dynamic or repetitive data? What remain the main open challenges for deployment in evolving and mission-critical environments?

Memory and I/O bottlenecks present fundamental constraints for analytics on high-volume, high-dimensional datasets. Probabilistic summary structures such as Bloom filters, Cuckoo filters, and their many variants provide essential tools by offering efficient, probabilistic set membership queries while greatly reducing memory consumption [14, 23, 32, 70, 76, 83, 96, 103]. Cuckoo filters, in particular, improve upon classical Bloom filters by enabling deletions and supporting tunable false positive rates without sacrificing speed or flexibility [32, 103]. Recent advances have lifted previous architectural constraints—including the necessity for power-of-two bucket counts and rigid layouts—by introducing signed-offset addressing and overlapping window layouts to further reduce overhead and balance load [83]. These improvements make Cuckoo filters especially suitable for domains like genomics and real-time analytics that demand rapid insertions, minimal memory, and high-throughput operations.

Compressed indexes represent another key advance, leveraging succinct data structures, run-length, and grammar-based compression to optimize the space–time tradeoff for a range of workloads [14, 23, 32, 70, 76, 83]. In domains with high data redundancy—such as version-controlled documents, genomic sequences, and large-scale log collections—modern indexes employ innovations like ILCP arrays and grammar-compressed document lists, enabling sublinear or compressed-space retrieval, counting, and ranking. For example, run-length encoded Burrows–Wheeler Transforms can cut working space to a fraction of data size for repetitive texts [76], and compressed suffix trees enable linear-time LZ77/LZ78

factorization in sublinear space [32]. In graph analytics, succinct q-gram tree indexes allow for scalable, in-memory queries on millions of graphs by compactly encoding occurrence patterns—accelerating both similarity search and filtering while requiring as little as 5–15% of the memory used by previous techniques [23]. Height-optimized tries [14] and adaptive radix trees with incremental cracking [96] further advance in-memory indexing for fast analytical workloads.

The interplay of features among widely used space-efficient storage structures is summarized in Table 11, which outlines their distinctive capabilities and recent improvements.

In practice, several persistent challenges temper the theoretical benefits of these advanced data structures: Update costs can be substantial, with many compressed and probabilistic structures struggling to support dynamic workloads without wholesale re-encoding. False positives and query errors inherent in approximate structures limit their use in mission-critical analytics. Dynamic and online compression strategies, essential for real-time or evolving databases, remain an early area of exploration. Effectively handling adversarial or worst-case distributions is another unresolved issue.

For instance, while run-length encoded or grammar-compressed indexes showcase dramatic memory savings and are ideally suited for highly repetitive data, their update and construction speeds often lag behind uncompressed alternatives, especially in dynamic or incremental settings [70, 76]. Similarly, space-optimal Las Vegas dictionaries [103] achieve near-information-theoretic efficiency for static sets but present challenges in adapting to dynamic updates or broader query models. Development of compressed computation algorithms that can operate directly on these representations is ongoing and remains a critical direction for future research [70].

In summary, space-efficient storage structures enable scalable analytics on massive datasets by fundamentally trading accuracy, update flexibility, and ease of implementation for dramatic savings in memory. The key open questions revolve around mitigating practical update costs, lowering query errors, supporting truly dynamic workloads, and extending advances from specific models to unified frameworks for broad real-world deployment.

## 8.3 Hardware and Parallelization for Analytic Scalability

**Objectives and Scope.** This subsection clarifies how modern hardware, parallelization, and distributed systems are leveraged to realize analytic scalability, especially in large-scale, privacy-sensitive, or latency-critical applications. We aim to address the guiding question: *How do recent advances in computational architectures, privacy mechanisms, and parallel methods coalesce to support scalable analytics, and what are the key challenges and tradeoffs encountered?*

Achieving analytic scalability necessitates leveraging modern hardware architectures, distributed systems, and parallelization paradigms. The advent of SIMD-capable CPUs and massively parallel GPUs has fostered a rich ecosystem of algorithms and data structures optimized for hardware acceleration. For example, fine-grained parallelization of index search—applied to both inverted and compressed indexes—uncovers that memory access patterns, cache locality, and SIMD-friendly encoding formats are as pivotal to query performance as the index design itself [13, 19, 41, 66, 78, 102].

It has been empirically established that leaving postings lists uncompressed can maximize traversal speeds; however, compression schemes such as QMX and Simple-8b attain comparable throughput while halving memory requirements, thereby offering a favorable tradeoff for search engine workloads [102].

These scalability concerns extend to distributed and federated environments, where sheer data volumes and stringent privacy constraints preclude centralization. Distributed range query indices [93], privacy-preserving federated learning [24, 56], and hybrid consensus protocols for secure retrieval [42] increasingly depend on sophisticated, decentralized approaches. Within federated analytics, local differential privacy (LDP) and secure, multi-level storage enable privacy-preserving computation, maintaining sub-second latency across thousands of distributed clients [56]. Recent innovations such as federated pseudo-sample clustering [86] illustrate that communication-efficient and privacy-preserving analytics are feasible through the synergy of local summarization, prototype exchange, and robust central aggregation.

Optimization is further enhanced through adaptive load balancing and streaming quantization, especially in domains like high-velocity recommender systems. Here, rapid index updates, cluster balancing, and repair mechanisms empower complex multi-task learning in the presence of continual data drift [21]. The integrated use of real-time streaming index construction with advanced ranking architectures typifies current directions for scalable, high-throughput analytics.

However, extracting optimal performance from hardware and system resources is challenging. Compressed indexes can induce cache bottlenecks; meanwhile, dynamic, parallel query processing—across both document-at-a-time and term-at-a-time paradigms—demands nuanced orchestration for effectiveness and efficiency [13, 102]. A promising avenue is the adoption of learned index structures and adaptive query execution, which dynamically tailor workload strategies to observed hardware characteristics using predictive models [4, 70, 77, 109].

**Key Takeaways.** Synthesizing developments across hardware acceleration, distributed protocols, and privacy-aware federated analytics, recent directions highlight the necessity of co-optimizing data access, encoding, privacy, and adaptability to underlying hardware profiles. These strategies not only sustain sub-second response in distributed and federated scenarios but also provide robustness to data drift, adversarial conditions, and evolving analytic objectives. The interplay between efficient compression, parallelism, decentralized learning, and privacy-by-design principles defines the state of the art in scalable analytic system design.

## 8.4 Adaptive and Online Index Updating

**Objectives and Scope:** This subsection aims to systematically survey and synthesize recent advances in adaptive and online index updating methods, with a focus on their underlying algorithmic frameworks, applicability across data types (relational, high-dimensional, and compressed), and the key open challenges for robustness and efficiency under dynamic workloads. Guiding questions include: How have online learning and feedback-driven mechanisms transformed index adaptation? What are the practical and theoretical

**Table 11: Salient Features of Probabilistic and Compressed Storage Structures**

| Feature | Bloom Filter | Cuckoo Filter | Recent Variants | Use-case Focus |
|---|---|---|---|---|
| Supports Deletions | No | Yes | Variant-specific | |
| Tunable False Positive Rate | By design | Tunable | Adaptive/Variable | |
| Dynamic Resizing | Limited | Possible | Variant-specific | |
| Bucket Structure | Fixed | Power-of-2 (classical) / Relaxed (modern) | Flexible | |
| Use-case Focus | General Set Membership | High-throughput, Frequent Updates | Application-specific | |

constraints when supporting continuous updates, especially in compressed or high-dimensional contexts? What distinguishes the latest approaches, and where do substantial open problems remain?

To maintain agility in ever-changing analytical environments, index structures must accommodate online, dynamic updates and facilitate autonomous tuning. A significant advancement in this direction is the deployment of adaptive and self-tuning indexes, often powered by online machine learning and feedback mechanisms rather than static, manually-tuned configurations. Frameworks based on multi-armed bandits and online learning algorithms allow for the continual exploration and exploitation of possible structural configurations, achieving rapid convergence toward optimal indexing layouts and demonstrating faster adaptation and improved robustness compared to traditional approaches [16, 74, 76]. These developments translate into substantial performance improvements, particularly in hybrid transactional-analytical (HTAP) systems and in responding to dynamic query workloads; for instance, Perera et al. [74] demonstrate that multi-armed bandit learning for online index selection provides up to 75% speed-up in shifting/ad-hoc analytical workloads and 59% in HTAP workloads over static configurations, while offering provable convergence guarantees and outpacing deep reinforcement learning in both speed and stability.

Such adaptability is not exclusive to relational systems. In domains like high-dimensional nearest neighbor search, adaptive algorithms iteratively refine cluster assignments, metric selection, and index organization based on observed data variability and feedback, thus maintaining performance even in adversarial or rapidly evolving settings [60]. Recent surveys, such as Mamun et al. [60], provide a detailed taxonomy that distinguishes multi-dimensional learned indexes by their degree of adaptivity (immutable vs. mutable), data layout dynamics, and model retraining support. Systems like OASIS can maintain families of locality-sensitive hash indices that adapt in real time as underlying similarity measures evolve—capabilities vital for interactive, non-stationary analytic workflows. In this context, supporting dynamic workloads, precise error bounds, and concurrency all remain essential open challenges.

Research into cracking and incremental construction methods (such as those used for Adaptive Radix Trees) reveals that dynamic, workload-driven partial indexing can yield significant construction-time gains without deteriorating query performance [96]. Wu et al. [96] provide practical algorithms for incremental ART construction through data-driven partitioning, minimizing upfront costs and supporting efficient, continuous index evolution. Such approaches highlight the broader trend toward workload-aware and progressive index adaptation.

Nonetheless, efficient online updating remains difficult for compressed or succinct data structures, where balancing, merging, and re-encoding may obscure or counteract performance benefits [14, 32, 76, 96, 103]. For example, Policriti and Prezza [76] show that while run-length encoded BWT-based and LZ77 indexing methods offer substantial reductions in working space for static data, their dynamic variants introduce practical speed bottlenecks. Similarly, Yu [103] provides succinct static set membership structures that approach information-theoretic optimality, but extending these results to efficient dynamic or online variants remains a prominent research challenge. In high-dimensional and topological settings, automaton-based compressions, as reviewed by Boissonnat et al. [16], excel for static and compressible data, but developing dynamic versions poses significant combinatorial obstacles.

Furthermore, ensuring resilience to concept drift, adversarial interactions, and catastrophic forgetting is an unresolved challenge, especially as analytic platforms become more autonomous and must contend with unpredictable, high-throughput streams. It is thus imperative to develop adaptive algorithms that are provably efficient and reliable under continuous workload evolution.

In summary, this section has articulated the major objectives of adaptive and online index updating: achieving autonomous, robust, and workload-aware index evolution across diverse analytic domains. By synthesizing frameworks—ranging from online learning-based adaptation, multi-dimensional and compressible index structures, to practical workload-driven strategies—this survey highlights both the state of the art and persistent research frontiers. The subsequent sections focus on domain-specific applications of these principles and outline open challenges in the pursuit of trustworthy and explainable analytics.

## 9 Multiway Data, Tensor Methods, and Higher-Order Analytics

This section presents the foundational concepts and emergent methodologies for analyzing multiway (i.e., multi-dimensional) data using tensor-based approaches, with the principal objectives of (1) formally characterizing the challenges and opportunities inherent in higher-order data analysis, (2) surveying canonical tensor decompositions and algorithmic structures that underpin modern analytics, and (3) synthesizing recent advances into a unifying framework that distinguishes this survey from prior overviews. We seek to address the following guiding questions: How do tensors expand the representational and computational scope relative to matrix-based methods? Which classes of tensor decomposition offer scalable modeling advantages in high-dimensional data settings?

What are the primary analytic, computational, and practical considerations when deploying such methods across real-world domains? Finally, how does the structure of this survey differ from previous work, and what new taxonomy or organizational principle do we propose for presenting the state of the field?

To facilitate reader orientation, each major subsection will open with a statement of objectives and thematic focus. The narrative is constructed to provide smooth transitions between topics, particularly from data-centric indexing and representation issues to algorithmic and theoretical analyses of higher-order tensor methods. Throughout, we will synthesize technical advances and illustrative, practical examples that clarify key multidimensional concepts and support accessibility for readers from a broad range of technical backgrounds.

At the conclusion of the section, we will summarize the central takeaways and outline open challenges at the intersection of multidimensional data representation, scalable analytic algorithm design, and interdisciplinary applications. This approach aims to provide a comprehensive yet cohesive perspective on higher-order analytics with tensors, establishing a structural and conceptual framework that underscores the novel contributions of this survey.

## 9.1 Prevalence and Application Areas

The rapid expansion of high-dimensional, multi-modal data in scientific and engineering disciplines has driven the extensive adoption of tensor-based methods for advanced data modeling and analysis. In contrast to conventional matrix-based techniques, tensor methodologies are specifically designed to preserve and leverage the intrinsic multiway structure characteristic of contemporary datasets. These datasets, arising from domains such as biomedical imaging (for example, functional MRI or hyperspectral imaging), temporal-spatial time series (including climate models and multi-channel EEG), and complex networked systems (such as multi-relational biological interactions or dynamic social networks) [9], often contain interrelations spanning more than two modes. By exploiting this higher-order structure, tensor models enable richer, more expressive representations of data, thereby uncovering multivariate interactions beyond the scope of pairwise (matrix) approaches.

For instance, in imaging science, tensors can simultaneously encode spatial, temporal, and spectral dimensions. Similarly, in network analysis, hypergraph analogs of tensors facilitate the study of multi-entity relationships, significantly advancing the analytical depth achievable in fields such as genomics and chemometrics [9]. This inherent capacity of tensor methods to capture and model complex relationships underscores the imperative for robust analytical frameworks capable of scaling with and adapting to the escalating complexity of contemporary scientific datasets.

## 9.2 Tensor Decompositions and Higher-Order Methods

At the core of multiway analytics are tensor decomposition techniques, which extend the principles of matrix factorization into higher orders and enable the discovery of latent structures embedded within complex datasets. Among these, the Canonical Polyadic

(CP) and Tucker decompositions are foundational. The CP decomposition represents a tensor as a sum of rank-one components, furnishing interpretable multiway analogs to singular vectors, while the Tucker decomposition generalizes principal component analysis (PCA) to encompass multiple modes by extracting interactions through a core tensor and orthonormal factor matrices [9].

Addressing the inherent nonconvexity and computational complexity of these decompositions, recent algorithmic advances employ strategies such as alternating least squares, gradient-based optimization, and stochastic techniques. These methods capitalize on problem-specific structures and incorporate sophisticated initialization procedures, thereby enhancing convergence properties and robustness to noise.

In addition to classical decompositions, contemporary research has expanded the scope of tensor analytics through higher-order statistical techniques, including tensor singular value decomposition (tensor-SVD), multiway PCA, and independent component analysis (ICA). Each of these frameworks brings distinct advantages for source separation and dimensionality reduction in tensor-formatted data [9]. Furthermore, novel mixture modeling and multi-mode regression approaches have been formulated within the tensor paradigm, empowering researchers to construct expressive models tailored to heterogeneous and structured data streams.

A particularly active research area involves tensor completion and recovery, where the objective is to impute missing entries by leveraging low-rank or structured sparsity assumptions. Such methods are critical for real-world scenarios where datasets are often incomplete or partially observed. While a rich variety of algorithms has emerged, all must contend with the significant challenges imposed by the "curse of dimensionality" and the absence of straightforward low-rank characterizations—factors that make the tensor setting fundamentally more complex than the matrix case.

As shown in Table 12, each decomposition method offers unique trade-offs in terms of modeling capabilities and application suitability within multiway data analysis.

## 9.3 Complexity and Open Challenges

Despite their substantial potential, tensor methods are accompanied by formidable analytical and computational challenges that stem from fundamental aspects of complexity theory and high-dimensional statistics. Notably, unlike matrices, tensors may not possess best low-rank approximations—a phenomenon posing significant obstacles to the design of optimal decomposition algorithms. Central analytical tasks such as low-rank tensor decomposition and rank determination have been shown to be NP-hard in the general case, establishing fundamental barriers for scalable computation [9]. This hardness sharply delineates the limits of what can be achieved algorithmically, especially in large-scale or high-noise data regimes.

A prominent issue is the disparity between what is statistically or information-theoretically achievable, and what current algorithms can compute efficiently. Even when estimators exist with theoretically optimal statistical guarantees, known algorithms may fail to realize these estimates within practical timeframes due to issues such as nonconvexity and local minima.

**Table 12: Comparison of Core Tensor Decomposition Techniques**

| Decomposition | Core Idea | Advantages / Typical Use Cases |
|---|---|---|
| CP Decomposition | Expresses tensor as a sum of rank-one components. | Interpretability, identifies latent factors, applicable in signal processing and topic modeling. |
| Tucker Decomposition | Generalizes PCA to multiway data, yielding a core tensor and factor matrices. | Captures interactions between modes; flexibility in modeling mode-specific variances; used in compression and feature extraction. |
| Tensor-SVD | Generalizes SVD to tensors via multi-linear operations. | Enables robust dimensionality reduction and source separation; effective for multi-modal signal processing. |

Recent research synthesizes methods from optimization, convex geometry, and random matrix theory to navigate these trade-offs. Efforts are ongoing to sharpen our understanding of sample complexity bounds, convergence rates, and the error profiles of different algorithms. Despite these advances, existing approaches often display suboptimal empirical performance, either requiring prohibitive data quantities or exhibiting susceptibility to poor local optima.

The structure of current algorithms for tasks such as clustering or indexing on tensor data tends to be rigid and insufficiently scalable, hindering deployment in practical analytics pipelines [9]. Several key open problems remain at the forefront of the field: designing algorithms that reconcile statistical optimality with computational tractability for high-dimensional and high-order tensors; developing robust initialization and regularization techniques suited to the unique challenges of tensor models; and extending clustering and indexing methodologies that natively operate on, and exploit, multiway tensor structures.

Addressing these open challenges is pivotal to fully realizing the analytical power of tensor and higher-order methods, with substantial implications for their application across varied scientific and engineering domains.

## 10 Applications and Deployment Strategies

This section aims to systematically explore the diverse applications and deployment strategies of the methods surveyed in this work. The primary objective is to distill key deployment patterns and elucidate how emergent techniques are being adapted across different domains. By synthesizing prevailing approaches and discussing practical considerations, we offer a comprehensive reference point for both researchers and practitioners seeking to operationalize these technologies.

We first define the scope of application areas addressed, highlighting the unique characteristics and requirements each domain imposes on deployment design. Subsequently, we examine prevailing deployment strategies, emphasizing practical integration workflows and comparative operational trade-offs. Throughout, we summarize essential takeaways and underscore decision criteria relevant to diverse use cases.

To enhance clarity for readers with differing technical backgrounds, brief illustrative examples are provided alongside complex multidimensional concepts.

Transitioning from foundational domains to more specialized applications, we maintain narrative continuity and highlight thematic connections. This section also seeks to clarify how our analysis synthesizes and extends previous surveys: whereas earlier works often catalog domains independently, here we propose a framework that categorizes deployment strategies by their underlying operational constraints (such as latency tolerance, data privacy, and scalability), offering a novel lens for comparative evaluation.

**Key takeaways:** By systematically categorizing application areas and deployment frameworks, this section aims to provide practical guidance for the real-world adoption of advanced methods in varied operational environments. The synthesized taxonomy establishes a foundation for future comparative studies, emphasizing both commonalities and unique requirements across deployment scenarios.

### 10.1 Application Domains and Case Studies

In recent years, state-of-the-art methods for clustering, indexing, and analytics have been deployed across a wide spectrum of scientific and industrial domains. This proliferation attests not only to the versatility of these techniques but also to the complexity inherent in their large-scale application.

In the fields of genomics and transcriptomics, advanced methodologies such as ensemble subspace regression and penalized mixed models have become instrumental. These tools elucidate molecular subtypes and latent structures within high-dimensional sequencing datasets by effectively balancing interpretability, predictive accuracy, and statistical rigor. Notably, ensemble regression techniques confer robust alternatives to classical penalized models, especially where the dimensionality far exceeds available observations, such as in gene expression biomarker discovery. Here, aggregation across random subspaces mitigates tuning sensitivity and overfitting tendencies [16, 36]. High-dimensional mixed-effects frameworks—augmented with sparsity-inducing penalties such as the smoothly clipped absolute deviation (SCAD)—have further advanced feature selection and inference, particularly within compositional microbiome investigations and genome-wide association study (GWAS) designs. Compared to traditional LASSO approaches, these methods offer superior performance amid clustered or highly correlated predictors [103].

Neuroimaging research, dealing with inherently multiway (tensor) data structures, has seen significant uptake of tensor-based clustering models. By exploiting separable covariance structures, these models enable both computational efficiency and scientific interpretability. The tensor normal mixture model, integrating sparsity-enforcing penalties with customized expectation-maximization procedures, exemplifies this paradigm: it delivers state-of-the-art performance on large neuroimaging datasets while providing principled quantification of cluster uncertainty and sensitivity to initialization [67]. Complementary clusterability diagnostics, grounded in multimodality analyses, serve as robust guides for assessing the intrinsic tendency for cluster formation—thereby cautioning against exclusive reliance on traditional, noise-sensitive internal indices [59].

Text analytics and the digital humanities benefit from innovations in indexing and data compression. Methods that leverage the repetitive structure of large textual corpora—such as run-length Burrows-Wheeler Transform (BWT)-based LZ77 factorization and

succinct membership data structures—substantially reduce memory consumption and computational demands, thus enabling scalable solutions in digital numismatics, linguistics, and large-scale search applications [35, 66, 100]. In chemical informatics, graph-based indexing strategies (e.g., for PubChem-scale datasets) combine hybrid encoding and succinct filtering to achieve notable space reductions and rapid query operations, even in the presence of millions of diverse molecular graphs [39, 104].

In the financial and social sciences, unsupervised learning approaches such as clustering uncover nuanced subpopulations and latent biases that traditional demographic or regression-based analyses may overlook. Notably, large-scale application of K-means clustering to financial wellbeing surveys has revealed patterns—such as explicit mismatches between subjective and objective financial stability—that challenge prevailing assumptions. These findings highlight both methodological opportunities for more informative clustering objectives and the need for mixed-model frameworks to disentangle complex, overlapping constructs [1].

Methodological advances in environmental analytics, EEG/gene clustering, and chemical informatics have been closely tied to the advent of scalable, distributed, and federated analytics platforms. For example, distributed nearest-neighbor systems utilizing Apache Flink, with domain-specific space-filling curve partitioning and granularity-aware load balancing, have enabled efficient analysis of granular smart meter or environmental sensor data—offering superior wall-clock performance relative to traditional central paradigms [76]. Similarly, approximate nearest neighbor search in high-dimensional chemical or image repositories increasingly employs graph-regularized sparse coding and quantization to reconcile recall, speed, and storage footprint [5, 50, 105].

Emerging domains, such as single-cell transcriptomics and clinical subtyping (e.g., diabetes), have driven the adaptation of techniques like generalized contrastive principal component analysis and mixed-membership modeling. Designed to decouple technical artifacts from biological signals, these frameworks produce interpretable axes of variation and robust unsupervised stratification, supporting research in heterogeneous and high-noise environments [32, 110].

The evolution of large-scale algorithms and data structures is intimately linked with augmented capabilities in massive data and graph indexing. As datasets increasingly exceed main memory capacity, techniques including dynamic polygon nearest-neighbor search, adaptive radix trees, voxelized spatial representations, and automaton-based simplex complex compression are indispensable for real-time analytics within both static and dynamic contexts [60, 75, 96, 102]. Current research in multidimensional learned indexes, database cracking, and compressed or low-footprint computation further underscores a dynamic field, where algorithmic, statistical, and hardware constraints motivate the development of novel theoretical models and practical open-source implementations [4, 14, 44, 58].

## 10.2 Large-Scale Deployments and Federated Analytics

As we transition from domain-specific method reviews to crosscutting themes, this section aims to explicitly address the core objectives of the survey: to synthesize emerging technical solutions for scalable, trustworthy analytics and to critically examine the operational and methodological barriers to their deployment. Our focus here is to bridge foundational advances in clustering, indexing, and federated learning with the realities of implementing such tools across diverse institutional settings.

Scaling advanced analytics from domain research to operational deployment introduces both computational and institutional challenges. Federated analytics and privacy-preserving clustering are of growing significance for applications in which data are distributed across independent institutions or geographical zones, subject to legal and governance restrictions on access and sharing. In these contexts, the use of open-source libraries and reproducible workflows is not only best practice, but often essential for enabling trustworthy, cross-institutional scientific collaboration [57].

Deployments at scale typically require algorithms for clustering, indexing, and spatial or graph analysis to function efficiently in distributed or parallelized environments. This demands an intricate balancing act between accuracy, processing speed, and memory resource usage. Empirical benchmarking of open-source range query and graph indexing libraries for high-performance computing has highlighted the importance of context-specific profiling—considering build time, query performance, and memory scaling—as well as the limitations of universal, "one-size-fits-all" strategies. Notably, brute-force or hybrid approaches sometimes demonstrate superior performance over more complex alternatives when operational data fall outside nominal parameter regimes [112]. However, these brute-force methods, while robust, can suffer from prohibitive computational cost and poor scalability in high-dimensional or adversarially noisy settings, whereas some sophisticated algorithms may fail to generalize or degrade sharply when input distributions deviate from assumed models [57, 112]. Failure to account for real-world heterogeneity can thus undermine the reliability of both naive and advanced techniques.

Federated learning introduces additional considerations, including statistical heterogeneity, communication overhead, and privacy preservation. While probabilistic model aggregation and distributed subspace consensus mechanisms have been developed to support inference across disparate data sources, these often face drawbacks: for instance, model drift from non-i.i.d. data, increased synchronization latency, and persistent privacy leakage risks if local updates are insufficiently protected [70].

Crucially, the assurance of reproducibility and the broad dissemination of open-source software, workflow templates, and standardized datasets underpin scientific trust, algorithmic benchmarking, and iterative methodological improvement. Practices such as explicit reporting of statistical validation, computational requirements, and parameter sensitivity facilitate fair comparisons and spur innovation across domains [57]. To ensure transparent evaluation, practitioners should routinely disclose not just strengths but also failure cases and domain-specific limitations of deployed algorithms.

In summary, large-scale deployments require both methodological innovation and transparent, reproducible engineering to overcome the inherent drawbacks of even the most advanced technical approaches. This crosscutting perspective sets the stage for our subsequent detailed survey of federated analytics methods.

## 10.3 Guidelines for Deployment

The objective of this section is to clearly articulate key practical recommendations for reliably deploying analytics solutions on complex, heterogeneous datasets, while highlighting both strengths and notable limitations of current approaches.

Extracting valid scientific and operational insights from complex, heterogeneous datasets requires adhering to principled standards for automation, benchmarking, and statistical validation. Recommendations, along with brief commentary on notable challenges or failure cases, are as follows:

**Automation**: Streamlining feature preprocessing, model selection, and parameter tuning can enable scalable workflows and help maintain interpretability and domain relevance. However, over-automation may obscure crucial data-specific nuances and introduce risks where model outputs become less transparent or less aligned with evolving domain requirements.

**Benchmarking**: Comprehensive benchmarking across diverse datasets and operational conditions, leveraging both internal and external evaluation indices, sensitivity analyses, and simulation-based studies, is essential for assessing clusterability, validity, and model robustness [112]. Yet, benchmarking may overlook certain failure cases, such as poor generalization across highly heterogeneous environments or underperformance on rare subpopulations not reflected in standard benchmarks.

**Statistical Validation**: Rigorous clusterability diagnostics and out-of-sample validation are particularly important in high-noise or high-dimensional environments to avoid spurious discoveries. Nonetheless, in practice, such validation can be challenged by limited access to labeled data, especially in domains with rare ground-truth or highly unbalanced classes, thus impeding the accurate assessment of model performance.

**Transparency and Reproducibility**: Transparent algorithmic reporting and open-source implementations, alongside publishing benchmarks and reproducible code and workflows, are vital for scientific rigor and collaborative development [57]. These practices, however, can be limited by proprietary data, commercial toolchains, or privacy constraints, which may reduce reproducibility in certain real-world scenarios.

**Scalability**: Algorithmic efficiency, memory optimization, and distributed computation must be prioritized. Resource-aware methods—including compressed computation, dynamic data structures, and federated analytics—are increasingly necessary for managing large-scale, heterogeneous datasets [70]. Despite substantial advancements, certain approaches (e.g., ensemble subspace clustering [57]) remain computationally demanding relative to simpler models, and their efficiency may be bounded by hardware or network constraints in very large-scale deployments.

**Interpretability and Ethics**: Ensuring model interpretability, transparent feature selection, and fairness auditing is essential, especially in sensitive biomedical, environmental, or social contexts. Despite ongoing efforts, interpretability and fairness auditing can be hampered by model complexity, latent variable effects, and the absence of reliable fairness metrics tailored to all data domains.

These balanced practices collectively support the development of robust, trustworthy, and adaptive analytics pipelines, while underscoring areas where continued research and practical vigilance

are necessary to ensure effective real-world analytics in the face of evolving data challenges.

## 10.4 Comparison of Representative Large-Scale Deployment Strategies

This section aims to explicitly outline the objectives and provide a structured, critical comparison of deployment strategies employed in high-dimensional and distributed analytics. Our goals are to clarify the trade-offs inherent in each approach, spotlighting not only their advantages but also notable drawbacks and representative scenarios where they may falter or fail.

The successful deployment of large-scale clustering and analytics methods thus hinges on careful alignment between methodological strengths and the practical realities posed by domain data, workflow constraints, and institutional context. While the approaches summarized above offer robust solutions within certain bounds, practitioners must remain vigilant regarding scaling ceilings, performance regressions, or failure cases as systems and data evolve. Continuing advancements in open-source dissemination, standardized evaluation, and adaptive algorithmic design will catalyze further innovation and responsible adoption across scientific and industry domains.

## 11 Crosscutting Themes, Challenges, and Emerging Research Directions

This section aims to explicitly outline the overarching objectives of our analysis: to identify recurring themes, articulate primary challenges, and highlight promising avenues for future research within the surveyed domain. In doing so, we not only underscore the strengths of key approaches but also provide brief, contrasting commentary on their notable drawbacks or documented failure cases, thus ensuring a balanced perspective.

Throughout our synthesis, we observe that while many approaches demonstrate impressive advancements—such as improved scalability and adaptability—certain limitations consistently emerge. For instance, methods that prioritize efficiency may inadvertently compromise robustness or generalizability, often failing in scenarios with high data variability or adversarial conditions. Conversely, models offering high flexibility sometimes incur prohibitive computational costs, impeding their practical deployability.

Additionally, a recurring challenge involves the standardization and interoperability of developed models across diverse application settings. Integration issues and lack of benchmark datasets further compound these obstacles, limiting meaningful performance comparisons and slowing community-wide progress.

By systematically contrasting advantages with corresponding limitations, this section provides deeper clarity on critical research bottlenecks and suggests focal points for forthcoming investigations. All references have been formatted consistently to enhance the professionalism and readability of this synthesis.

### 11.1 Integration and Adaptivity

The escalating volume, complexity, and heterogeneity of contemporary data have accentuated the necessity for adaptive and integrative systems across indexing, clustering, feature selection, similarity search, and statistical modeling. A pronounced trend

**Table 13: Comparison of large-scale deployment strategies for clustering and analytics.**

| Strategy | Advantages | Constraints / Notable Drawbacks | Application Examples |
|---|---|---|---|
| Distributed parallel analytics (e.g., Apache Flink, Spark) | High scalability; fault tolerance; supports massive input volumes | Requires significant infrastructure setup; can introduce resource contention; may need complex partitioning for optimal performance; failure scenarios include network or node bottlenecks impacting throughput | Smart grid analytics, environmental sensor networks |
| Federated clustering/learning | Preserves privacy; data remains local; enables cross-institutional collaboration | High communication costs; statistical heterogeneity can degrade model convergence; aggregation becomes complex with divergent local distributions; failure if local sites cannot synchronize updates | Multi-center biomedical studies, cross-jurisdictional finance |
| Graph-based indexing with hybrid encoding | Space-efficient; supports rapid queries over large, diverse graphs | Computationally expensive index construction; sensitive to parameter tuning; may perform poorly on highly dynamic or evolving graphs; failure can occur when structure does not generalize | Chemical informatics, PubChem-scale search, social network mining |
| Compressed/learned data structures | Drastically reduced memory footprint; competitive accuracy | Implementation complexity; can be highly sensitive to parameter selection; degradation in accuracy if compression loses key features; may struggle with continual updates or streaming data | Text analytics, high-throughput genomics, image retrieval |
| Centralized brute-force/hybrid approaches | Simplicity; robust against certain data irregularities; minimal tuning | Does not scale to massive datasets; high per-node resource demand; often fails on data with high dimensionality or volume (memory/compute limits) | Small- to medium-size or irregular dataset scenarios |

has emerged toward the unification of methodologies traditionally addressed in isolation, including but not limited to the joint handling of clustering and feature selection, spatial and graph indexing, learned and annotative indices, and adaptive tensor models [7, 30, 33, 53, 68, 69, 77, 101, 106, 111]. This movement is primarily motivated by empirical limitations observed in "one-size-fits-all" techniques, which become inadequate as data increases in dimensionality, dynamism, or semantic richness.

For instance, hybrid paradigms combining prototype reduction with learned dimensionality compression empower $k$-NN search to achieve significant gains in speed and accuracy. Recent studies have demonstrated that convex hull selection, stratified sampling, principal component analysis, and auto-encoder-based representations can deliver classification speed-ups of up to 32×, often with minimal or even improved accuracy, but these approaches may be affected by data overlap and class imbalance, requiring flexible representation selection and dynamic parameter tuning [8, 38, 68, 75, 89, 99]. Work on exact $k$NN methods has further highlighted the effectiveness of hybrid and ensemble approaches, showing, for example, that variants leveraging ensemble strategies achieve robust, domain-adaptable performance, especially for complex and high-dimensional data [8, 38, 89]. Addressing class overlap and imbalance simultaneously, recent $k$NN models employ composite weighting schemes, enhancing reliability and robustness across imbalanced and noisy datasets without parameter tuning [8].

Similarly, the integration of feature selection and clustering—especially for mixed-type or high-dimensional datasets—exploits joint optimization and ensemble techniques to reinforce cluster robustness and enhance attribute discrimination, even under adverse conditions such as adversarial noise or low signal-to-noise ratios [73, 93]. Advances in deep clustering link feature learning and cluster assignment in unified or iterative frameworks, improving adaptability to diverse data domains, while partitional and hierarchical methods remain essential for balancing accuracy, efficiency, and interpretability [73, 111]. The importance of evaluating clustering effectiveness via multiple internal and external metrics and accommodating varied evaluation methodologies has also been emphasized in recent surveys [73, 101, 111].

Tensor-based modeling constitutes another pivotal frontier in integrative analytics, offering interpretable and scalable substrates for multiway data prevalent in scientific and engineering applications. Penalized tensor mixture models and scalable decomposition algorithms have been developed to reconcile statistical consistency in clustering with computational scalability, particularly in high-dimensional scenarios [38, 89, 92]. Furthermore, manifold learning perspectives and nonlinear representation approaches offer additional capabilities for capturing intricate high-dimensional structures and heterogeneous experimental conditions; however, these advancements demand stronger theoretical guarantees and enhanced adaptivity at scale [8, 16, 28].

The convergence of indexing paradigms—most notably through annotative and learned indexing—has yielded robust frameworks for unified, scalable data platforms. Annotative indexes generalize over inverted, columnar, and graph-based strategies, supporting transactional, concurrent, and semi-structured workloads, alongside complex knowledge graph scenarios [25, 26, 60]. Recent work provides a detailed taxonomy for classifying learned indexes across dimensions such as design, mutability, data layout, insertion strategy, and supported query types, and outlines open challenges including precise error bounds, efficient (re-)training, concurrency, and GPU acceleration [60]. Annotative indexing further supports ACID-compliant, transactional ingestion and expressive SQL-like querying over heterogeneous JSON data with high concurrency, showing performance advantages in dynamic, large-scale environments and enabling unified management of textual, structured, and vector/graph data [26]. Lightweight distributed learned indices for big spatial data, such as LiLIS, demonstrate dramatic improvements in query speed and index construction, coupled with robust scalability for diverse spatial workloads within existing big data platforms like Spark [25]. Likewise, proximity graph-based frameworks like UNIFY deliver efficient, range-filtered approximate nearest neighbor search with integrated hybrid filtering and automatic strategy selection, optimizing for high scalability and recall in attribute-constrained searches [53]. Such frameworks have been demonstrated to outperform traditional approaches in both efficiency and flexibility, positioning annotative and learned indexes as central solutions for modern adaptive data management systems [25, 26, 43, 60, 97].

## 11.2 Machine Learning for Index and Analytic Optimization

Machine learning has assumed a central role in the optimization of index structures within database management and analytics platforms. Rather than depending exclusively on hand-crafted heuristics or costly offline tuning, modern methodologies treat index management as a learning or decision-making procedure, leveraging workload observation and cost feedback to perpetually adapt [25, 60]. Noteworthy progress has been made through resource-efficient recommendation systems utilizing large language models (LLMs), which synthesize workload characteristics and deduce ideal index schemas with minimal retraining or manual intervention. These systems integrate demonstration pools, scalable inference engines, and domain knowledge injection, attaining recommendation quality and latency that rivals or exceeds traditional index advisors [60, 64]. Key advances include:

- Modeling the index recommendation task for compatibility with few-shot or in-context learning,
- Extracting granular statistics from diverse workloads, and
- Deploying scalable aggregation mechanisms for robustness.

In addition to LLM-driven approaches, online learning frameworks—inspired by bandit algorithms—eliminate dependencies on DBA expertise or traditional query optimizers by utilizing active exploration and exploitation of index alternatives based on real-time performance measurement. These methods guarantee convergence to near-optimal performance, even in rapidly evolving or ad hoc workloads, and frequently surpass the efficiency of both deep reinforcement learning models and static analytics methods [60]. Nonetheless, more sophisticated or hybrid analytic workloads continue to pose challenges related to model expressiveness and rapid adaptation [60, 64].

## 11.3 Transactional and Distributed Perspectives

Recent developments in database management and analytical infrastructures have been characterized by the deepening intertwining of transactional, distributed, and computational paradigms. The capability to robustly execute distributed queries with strong ACID guarantees has become increasingly pivotal given the emergence of vector, graph, and hybrid analytics that necessitate cross-engine and federated access [60, 76, 83]. State-of-the-art systems are required to efficiently manage and query across heterogeneous backends, which often entails seamless integration among graph databases, knowledge graphs, and spatial or textual search engines—all while upholding high standards of performance and correctness [16, 26, 60].

Innovative indexing structures and partitioned system architectures enable distributed query execution and dynamic partitioning, though they introduce new trade-offs involving communication costs, consistency maintenance, and optimization under the constraints of partial or privacy-preserving data access [2, 10, 32, 76, 103]. Federated analytics, in particular, must balance the ideals of openness and collaboration with the complexities of security, transactional integrity, and latency management in geographically dispersive environments [22, 39, 50]. The rising emphasis on ACID properties in open, federated, and multimodal systems reflects a growing awareness of the imperative to integrate transactional guarantees within scalable, adaptive analytic environments [25, 26, 60].

## 11.4 Robustness and Adversarial Resilience

The expansion of indexing and analytics frameworks into sensitive domains—including healthcare, finance, and security—has rendered adversarial and randomized query resilience an essential system property. In high-dimensional, graph-structured, and compressed data environments, deliberate manipulations can degrade system performance and expose critical patterns, particularly when underlying indexes rely on learned or compressed representations [58, 60]. Evaluations have shown that graph-based and tensor analytic models are vulnerable to perturbations, underscoring the need for robust and regularized representations capable of withstanding worst-case and stochastic adversarial behaviors without compromising retrieval or inference quality [9, 58, 70].

For example, privacy-preserving document retrieval has adopted cryptographically fortified indexing, randomized responses, and obfuscation strategies to counteract leakage and inference-based attacks, typically at the cost of increased storage or computational

burden [60, 70]. Likewise, resilient similarity and clustering algorithms for graph and high-dimensional data integrate adversarial feedback, robust scoring metrics, and continual adaptation. However, computational efficiency and universal applicability remain open challenges in scaling robust data analytics [58, 60]. Achieving equilibrium between privacy guarantees, adversarial robustness, and system efficiency is a persisting core problem for the community [70].

## 11.5 Online, Adaptive, and Learned Indexing for Dynamic Workloads

Contemporary workloads, characterized by rapid streaming, immense scale, and frequent evolution, have elevated the importance of online, adaptive indexing systems that can dynamically adjust to shifting data distributions and workload requirements. Emerging learned and hybrid index solutions continuously evolve in response to new data patterns and real-time feedback, outperforming static or heuristically managed systems on throughput and accuracy, especially for streaming and hybrid transactional/analytical processing (HTAP) datasets [25, 60, 70, 74]. Central techniques include:

- Incremental index maintenance,
- Bandit-based adaptation mechanisms, and
- Context-sensitive indexing strategies.

Recent findings indicate that such frameworks can consistently surpass their fixed counterparts, but obstacles persist regarding staleness prevention, resource overhead management, and robust generalization across diverse query and workload types. The integration of multi-armed bandit strategies, adaptive feedback loops, and sophisticated feature extraction are identified as vital constructs for future universally adaptive indexing systems [60, 70].

## 11.6 Societal, Fairness, Privacy, and Ethical Issues

As the survey transitions to crosscutting challenges impacting sensitive, large-scale, or scientific data analytics, it is important to explicitly restate our guiding objective: to critically examine technical advances in adaptive indexing, machine learning, and similarity search in the context of their broader societal, ethical, fairness, and privacy implications. This section synthesizes recent developments and exposes key limitations and risks, facilitating a deeper understanding of how technological design choices reverberate at the societal level.

The integration of advanced analytics, machine learning, and adaptive indexing within domains involving sensitive, personal, or scientific data has brought ethical, fairness, and privacy considerations to the forefront. High-throughput automated decision-making provides scale and efficiency, but also raises significant risks related to bias propagation, privacy violations, and transparency loss if not properly mitigated [11, 13, 19, 31, 40, 42, 45, 46, 51, 54, 63, 78, 86, 87, 92–94, 108].

Recent research emphasizes that to address these risks, modern data systems must satisfy formal privacy criteria such as differential privacy and the use of immutable distributed ledgers; incorporate explainability, auditability, and reproducibility by design; rigorously validate index recommendations and clustering attributions;

and systematically evaluate the societal implications of data accessibility, reusability, and potential algorithmic bias, particularly for high-stakes applications in healthcare, finance, and public policy [19, 35, 40, 45].

Direct technical comparisons reveal several limitations in alternative approaches: For privacy protection, methods employing only traditional anonymization or trusted third parties are increasingly outperformed by decentralized, local differential privacy techniques and distributed ledger architectures, which better mitigate risks of single points of failure and unauthorized data exposure [42, 86]. In clustering and high-dimensional analysis, approaches that directly adapt to noise, feature heterogeneity, and class informativeness (such as consensus subspace clustering and robust ensemble techniques) demonstrate improved fairness and validity over classical methods relying on fixed global parameters [4, 19, 57, 112]. However, enhanced computational cost or the need for extensive feature engineering may limit their deployment. For explainability and auditability, recent index recommendation strategies leveraging large language models and in-context demonstration pools [45] offer significant improvements in transparency and adaptability compared to legacy heuristics or black-box RL methods, yet practical limitations remain for highly complex, dynamic workloads. In addition, reproducibility is best advanced by open benchmarks, transparent data handling, and comprehensive ablation studies; approaches lacking such disclosures face substantial barriers to real-world adoption and validation [45, 54]. Lastly, while techniques such as immutable ledgers and anonymous verification promote compliance and accountability, ongoing research underscores the importance of harmonizing privacy, auditability, and regulatory obligations, especially in data-intensive and dynamic environments.

Transparent methodological disclosure, provenance-aware index construction, and open/reproducible analytics are being advanced to strengthen accountability, while parallel research seeks to harmonize privacy, auditability, and regulatory compliance in data-intensive environments [4, 21, 28, 29, 39, 42, 44, 54, 57, 60, 66, 67, 70, 100, 104, 110, 112].

In summary, ensuring ethical, privacy-preserving, and fair data practices is inseparable from state-of-the-art technical progress. Each new advance should be evaluated not only on performance and scalability, but also on its potential long-term societal consequences, setting a research agenda that foregrounds transparency, accountability, and inclusivity as essential pillars for future work.

## 11.7 Emerging Research Directions

**Objectives:** This section aims to synthesize the key research challenges and future prospects in data analytics and management, providing explicit contrasts between principal approaches and identifying foundational objectives—namely, efficiency, scalability, interpretability, fairness, robustness, and adaptability. Our synthesis anchors on: (1) unifying advances across indexing, retrieval, and analytic modeling; (2) drawing attention to both achievements and unsolved difficulties; and (3) inviting consideration of integrative frameworks bridging current research silos.

*Unified Indexing Architectures: Neural, Hybrid, Annotative, and Compressed.* Neural representations, hybrid model-index combinations, annotative indices, and compressed data structures each enable trade-offs between expressivity, compressibility, and speed. For instance, compressed indexes optimized for repetitive or diverse data types yield significant storage and retrieval efficiencies [25, 60, 64, 70], yet can introduce increased model training complexities and rigidities that impact update flexibilities and query latency [25, 60]. Annotative and interpretable indexing methods [64] offer improved transparency but sometimes incur a modest accuracy penalty or limited adaptability to unseen query types. Open questions include optimal multidimensional trade-offs and guarantees on error bounds, especially for dynamic or adversarial workloads [25, 60].

*Retrieval-Augmented Generation (RAG) and Structured LLM Queries.* Embedding retrieval engines tightly within large language models, as seen in RAG architectures, advances knowledge-grounded query response [57, 60, 64]. However, RAG and structured LLM queries face challenges including seamless integration across disparate vector, relational, and graph indices, and require new unifying interfaces for knowledge graph management and prompt engineering [60, 64]. These frameworks are often susceptible to inconsistency between retrieved evidence and generated content, and to brittleness in the face of rapid schema evolution or incomplete data.

*Unified Statistical–Computational Analytics.* Emerging analytic systems increasingly adjoin scalable tensor and mixture modeling, ensemble clustering (e.g., consensus subspace clustering), and fairness-aware learning, with goals of robust, theoretically-grounded, and practical analytics [57, 60, 70, 92, 112]. Benefits include provable optimality and resilience to noise or partial information [57, 112]. Nonetheless, high computational overhead, increased system complexity, and difficulties in handling nonparametric or mixed-type data present concrete limitations [57, 60]. The pursuit of unifying consensus frameworks and minimax-consistent algorithms remains a primary objective.

*Robust and Scalable Adaptive Systems.* Ensuring robustness to adversarial manipulation, distributional shifts, and environmental changes—especially under federated or streaming architectures—remains challenging [58, 60, 70, 112]. While adaptive learning, privacy-preserving analytics, and composable index synthesis are vital enablers, their deployment is often complicated by trade-offs among privacy, latency, and update adaptability [60, 112]. Constants in algorithmic optimality achieved in some theoretical models do not always translate into simplified practical solutions with minimal computational resource demand [58, 70].

*Crosscutting Challenges and Outlook.* Key cross-domain challenges and promising trajectories are summarized as follows:

(1) **Achieving Efficient, Interpretable, and Adaptive Indexing:** Balancing highly expressive models (such as deep neural indexes) with interpretability and efficient retraining remains a central concern [25, 60, 64].

(2) **Unifying Statistical and Algorithmic Guarantees:** Advancing algorithms that provide rigorous statistical guarantees while being computationally tractable for high-dimensional, large-scale, and heterogeneous data.

(3) **Fairness and Robustness:** Addressing fairness-aware analytics and resilience to noisy, adversarial, or underrepresented data regimes, especially as workloads become increasingly federated and diverse [57, 112].

(4) **General-purpose Integration Frameworks:** Developing common abstractions and taxonomies that span clustering, indexing, and analytic techniques, enabling seamless pipeline construction and hybridization [60].

The ongoing integration of adaptivity, efficiency, fairness, and interpretability across disciplines continues to define the trajectory of fundamental research and is expected to underpin future breakthroughs in data-driven intelligence.

## 12  Synthesis and Conclusion

### 12.1  Restatement of Objectives

This survey set out to achieve the following objectives: (1) to systematically review the state-of-the-art techniques in clustering, indexing, and analytic methods relevant to the targeted application domains, (2) to provide a comparative analysis illuminating the core advantages and limitations of each approach, and (3) to identify open challenges, future directions, and potential synergies across the examined methodologies.

### 12.2  Synthesis of Approaches

The reviewed work demonstrates considerable progress in clustering, indexing, and analytic techniques, each addressing specific aspects of scalability, accuracy, and interpretability. Notably, clustering approaches excel in unsupervised structuring and knowledge extraction, while indexing techniques enhance retrieval efficiency for large-scale datasets. Analytic frameworks further enable actionable insights from complex data streams. Despite distinct emphases, an integrative perspective reveals several common challenges, such as handling high-dimensionality, ensuring scalability, and supporting real-time decision-making.

### 12.3  Proposed Unifying Taxonomy

Based on our analysis, we propose a unifying taxonomy encompassing three interrelated dimensions: 1. Data Characteristics: including dimensionality, heterogeneity, and dynamism. 2. Methodological Axis: spanning clustering granularity, indexing architecture, and analytic depth. 3. Application Targets: ranging from exploration and summarization to predictive analysis.

This taxonomy facilitates systematic comparison and guides both researchers and practitioners in selecting and combining techniques according to specific problem requirements.

### 12.4  Key Challenges and Future Outlook

Key open challenges and promising future directions emerging from the literature are: 1. Developing adaptive algorithms capable of managing evolving and dynamic data distributions. 2. Bridging the gap between unsupervised clustering and supervised analytic frameworks to enhance interpretability and performance. 3. Designing scalable indexing structures that remain efficient with increasing dataset size and complexity. 4. Integrating privacy-preserving and

fairness-aware mechanisms within all methodological stages. 5. Automating parameter selection and hyperparameter tuning to reduce domain expertise barriers. 6. Fostering cross-domain transferability and generalization of developed approaches.

### 12.5  Concluding Remarks

This survey has provided a comprehensive synthesis of clustering, indexing, and analytic methodologies, clarifying their individual and combined roles in addressing contemporary data challenges. By restating objectives, integrating a taxonomy for conceptual cohesion, and outlining explicit challenges and future avenues, we offer a roadmap for continued advancement and cross-fertilization in this rapidly evolving field.

### 12.6  Comparative Review and Synthesis

The contemporary landscape of clustering, indexing, and similarity search for high-dimensional and categorical data is characterized by substantial methodological diversity and paradigm shifts. Traditional hard clustering approaches, such as $k$-means and hierarchical clustering, remain foundational for their simplicity and interpretability. However, these methods encounter significant challenges—including the curse of dimensionality, limited scalability, and sensitivity to noise or parameter selection—when confronted with complex, large-scale, or categorical datasets [76, 77, 86]. In response, modern research has produced a progression of enhanced methodologies: density-based, spectral, consensus, and ensemble clustering techniques, each designed to accommodate heterogeneities in data structure, density, and scale.

Spectral clustering has demonstrated consistently superior performance in high-dimensional contexts, owing to its use of eigenspace transformations that facilitate robust separation and flexible parameterization [82]. Nevertheless, this approach frequently incurs higher computational costs and exhibits increased sensitivity to initialization and hyperparameter configuration [17, 32, 56].

Consensus and ensemble clustering have emerged as pragmatic answers to the instability and ambiguity associated with model selection in high-dimensional or noisy regimes. By aggregating the outputs of multiple clustering executions—employing varying feature projections, subsamples, or foundational algorithms—these strategies capitalize on the "wisdom of the crowd" principle to enhance robustness and accuracy. Theoretical and empirical evidence supports their efficacy in challenging scenarios, such as sub-Gaussian mixtures and mixed-type data [25, 28, 33, 86, 93]. Nonetheless, the computational burden of consensus methods remains a concern, stimulating ongoing research into improving their scalability and refining the minimax optimality of combination rules.

Feature selection and dimensionality reduction are now indispensable for effective clustering and indexing in high-dimensional spaces. Established methods such as Principal Component Analysis (PCA), $t$-SNE, and UMAP remain prevalent for uncovering manifold structures. Yet, these techniques may yield misleading representations under heavy noise or nonlinearity, exemplified by the "scattering noise" phenomenon. Recent advances, such as the distance-of-distance transformation, address these limitations by disentangling structural signals from noise prior to embedding [26].

Moreover, the adoption of ensemble subspace projections, random feature selection, and regularized tensor decompositions—including tensor PCA and tensor-normal mixture models—expands dimensionality reduction techniques to multiway and highly structured data, thereby bolstering both statistical efficiency and scalability [2, 10, 51, 62].

Indexing methodologies are undergoing transformative change with the advent of massive, high-dimensional, and repetitive or categorical datasets. Classical spatial and metric indexes (e.g., $k$d-tree, R-tree) experience sharp performance degradation in very high dimensions or with heterogeneous attribute types. Consequently, contemporary solutions such as graph-based indexes (HNSW, proximity graphs), neural network-based systems, and compressed/text-indexing structures are increasingly adopted [27, 45, 61, 75, 109]. Annotative indexing innovatively integrates paradigms across inverted indexes, graph databases, and knowledge graphs within unified, scalable frameworks, enabling efficient transactional and concurrent querying, as well as supporting dynamic and semi-structured data formats [26]. This multi-paradigm approach supports efficient retrieval for both structured and unstructured data at scale [88]. In parallel, learned indexes and multi-dimensional neural indexing systems exhibit dynamic adaptability, model-driven querying, and robustness to distributional changes and retrieval-augmented generation workflows [3, 25, 35, 82]. The emergence of lightweight distributed learned indexes further underscores advances in minimizing query latencies and index building costs for spatial data, with approaches like LiLIS and ML-based spatial partitioning demonstrating order-of-magnitude improvements in query speeds and throughput [25, 71].

Substantial advances in similarity and range search have followed the evolution from exact $k$-nearest neighbor ($k$NN) algorithms to approximate methods. Notably, the use of product quantization, residual corrections, and graph traversal heuristics has yielded marked improvements in computational scalability. Recent work on graph-based indices has specifically advanced range search performance, with adaptive queries and early-stopping heuristics providing significant speedups, especially in large and dense datasets [61]. Innovations such as minimization residual quantization (MRQ), range-aware filter and hybrid-search algorithms (e.g., UNIFY, HSIG), and specialized index structures for applications like time series and trajectories now support billion-scale, real-time query workloads with reliable recall and efficient resource use [21, 42, 48, 53, 71–73, 78, 98, 109, 111]. Furthermore, robustness to dynamic workloads and adversarial query patterns is increasingly managed through adaptive algorithms and hybrid or ensemble-based indexing strategies [90, 99, 104]. Unified frameworks, like UNIFY, exemplify this by supporting hybrid pre-, post-, and range-filtered search on high-dimensional attribute-rich datasets while maintaining efficient index maintenance and scalability [53].

Tensor analytics, comprising decomposition models and high-order network embedding, represents a frontier in extracting latent structures from multidimensional data arrays as found in omics, neuroscience, and signal processing. Recent algorithms exploit the interplay between statistical and computational constraints to deliver interpretable and consistent factorizations. These methods overcome difficulties such as the lack of best low-rank approximations or the NP-hardness of optimization tasks, effectively balancing parsimony, scalability, and uncertainty quantification [10, 62].

Hardware-aware and compressed computation paradigms further expand the boundaries of feasible analytics by operating on compressed or in-memory representations, vital for petabyte-scale or streaming datasets [12, 83, 108]. These approaches emphasize CPU/GPU affinity, cache locality, and architecture-specific optimizations, as evident in developments related to index compression, efficient filter structures (e.g., windowed cuckoo filters), and compact data structures for document or sequence analysis [30, 83, 92, 107].

Taken together, the field's synthesis highlights a movement towards hybrid, adaptive, and robust systems. Integrating dimensionality reduction, advanced indexing, and multi-perspective clustering is increasingly recognized as essential for the comprehensive analysis of complex, high-dimensional, and categorical data.

## 12.7 Ongoing Challenges and Open Problems

Despite significant advances, several important theoretical and practical challenges continue to impede progress in the field:

Scalability and Expressiveness: Achieving scalable solutions for graph indexing and high-order analytics remains difficult, especially for dynamic, streaming, or extremely large datasets (e.g., containing billions of nodes). Current approaches are often limited by high memory requirements, costly maintenance, and inadequate response times [9, 58]. The challenge is exacerbated by structural complexities inherent in high-dimensional or tensor representations, where NP-hardness and the absence of efficient low-rank approximations often create bottlenecks in both information extraction and computation.

Robustness: Many existing systems fall short in providing robust handling of adversarial input distributions, substantial noise, and distributional shifts. These limitations constrain the deployment of analytics in real-time or adversarially influenced environments [23]. Techniques such as ensemble subspace clustering and consensus spectral methods have made progress toward mitigating noise and adversarial effects, yet robust solutions that scale remain an open area of investigation.

Statistical-Computational Gap: Particularly for high-order and high-dimensional analytics, a persistent gap exists between statistically optimal methods and those attainable by efficient (e.g., polynomial-time) algorithms. Even when theoretically optimal solutions are known, practical implementations may rely on suboptimal initialization and fail to guarantee reliable convergence, especially in complex optimization landscapes [9, 60].

Statistical Rigor vs. Computational Efficiency: Recent advances often prioritize speed or parallelism, at times sacrificing statistical soundness and interpretability. In sensitive domains such as biomedical analytics, failing to maintain statistical consistency or inferential reliability can undermine the trustworthiness of outcomes and restrict real-world applicability [57, 60]. There is a continuing need for methods that balance rapid data processing with transparency and inferential guarantees.

**Table 14: Comparative Overview of Major Methodological Advances in High-Dimensional Analysis**

| Strategy or Method | Domains of Strength | Primary Advantages | Principal Limitations |
|---|---|---|---|
| Traditional Hard Clustering | Numeric, low-dimension | Simplicity, interpretability, fast convergence | Sensitive to noise, non-scalable, poor in high-dimension |
| Spectral Clustering | High-dimensional, networks | Robust separation, adaptable parameterization | High computational cost, initialization sensitivity |
| Consensus/Ensemble Clustering | Heterogeneous, noisy data | Robustness, improved accuracy, model instability handling | Computationally intensive, scaling challenges |
| Dimensionality Reduction | High-dimensional, manifold | Enhanced visualization, subspace recovery | Potential distortion/noise, manifold discontinuity issues |
| Graph-based Indexing | Large-scale, high-dimension | Efficient retrieval, adaptability, multi-paradigm support | Memory overhead, maintenance difficulty |
| Learned/Neural Indexes | Dynamic, large datasets | Model-driven access, adapts to data drift | Training complexity, generalization uncertainty |
| Approximate Similarity Search | Real-time, billion-scale | Fast query, recall-resource trade-offs | Possibly lower accuracy, adversarial vulnerability |
| Tensor Analytics | Multimodal, structured data | Latent pattern discovery, scalability, uncertainty quant. | Stat/comp. gap, convergence obstacles, complexity |
| Hardware-aware Computation | Streaming, petabyte-scale | Efficient memory use, architecture leveraging | Hardware dependency, compression artifacts |

Reproducibility and Benchmarking: The absence of comprehensive and standardized benchmarks, inconsistent data handling practices, and evaluation biases inhibit meaningful comparison of methods. The field critically needs multidimensional benchmarks and open-access repositories to support reproducible research and unbiased progress [57, 60].

Ethical and Societal Considerations: As high-dimensional analysis becomes pervasive in decision-critical domains, issues of fairness, transparency, privacy, and user agency are increasingly urgent. There is intensified demand for algorithmic frameworks that provide guarantees regarding fairness and responsible governance, particularly where analytic outcomes influence individual rights or societal welfare [70, 112].

## 12.8 Future Outlook and Roadmap

Looking ahead, several converging trajectories are anticipated in the evolution of analytic systems and data structures for high-dimensional and categorical data:

Scalability remains a central challenge and priority. Future systems are expected to leverage hybrid architectures that blend compressed, hardware-conscious computation with distributed and cloud-native paradigms, making it feasible to manage both static and streaming massive datasets efficiently [12, 83]. Notably, advances such as highly space- and memory-efficient data structures (e.g., improved Cuckoo filters) and near-optimal dynamic algorithms for problems like vertex cover and matching are setting new practical standards in scalability, flexibility, and update efficiency.

Interpretability will be vital in building trust and driving adoption, especially in critical areas such as medicine, finance, and public policy [57, 60]. Progress here will rely on bridging transparent, explainable outputs from both statistical and machine learning-based models, and on the development and integration of indexes and clustering techniques whose operations are amenable to human scrutiny and reasoning.

Benchmarking and reproducibility form another pillar of future research. There is a growing need to develop comprehensive and standardized benchmarking suites that address clustering, indexing, and similarity search tasks using both synthetic and real-world datasets [57]. Careful benchmarking underpins objective evaluation and fosters reproducible innovation, especially as methods for high-dimensional, categorical, and noisy data become more complex and statistically nuanced.

Ethical and open science integration must be woven into both the algorithmic methodology and practical deployment. Considerations of fairness, privacy, and transparency will ensure that analytic systems serve users equitably and minimize societal risk [70, 112]. As compressed computation and scalable learning become ubiquitous, addressing their impact on society, including open science practices and the mitigation of algorithmic biases, becomes paramount.

Research imperatives for the future include expanded investigation into dynamic graph and tensor analytics, compressed and federated computation, robust scalable clustering for mixed-type and noisy data, and interpretable, learning-augmented indexes [57, 60, 70]. Bridging the divide between statistical rigor, computational efficiency, and societal responsibility stands as a defining challenge for the forthcoming era.

In summation, there is no singularly dominant approach in the high-dimensional analytic landscape. Instead, progress points toward integrated, adaptive, and accountable systems—where advances in dimensionality reduction, indexing, similarity search, and robust clustering are tightly coupled with rigorous benchmarking, interpretability, and societal stewardship. Achieving a cohesive synthesis among statistical excellence, computational scalability, and ethical responsibility will define the future of high-dimensional data analysis systems.

## References

[1] A. Adolfsson, M. Ackerman, and N. C. Brownstein. 2019. To cluster, or not to cluster: An analysis of clusterability methods. *Pattern Recognition* 88 (2019), 13–26. doi:10.1016/j.patcog.2018.10.026

[2] P. Afereidoon. 2025. persiansort: an alternative to mergesort inspired by persian rug. *arXiv preprint arXiv:2505.05775 [cs.DS]* (2025). https://arxiv.org/abs/2505.05775

[3] E. J. Aguilar and V. C. Barbosa. 2023. Shape complexity in cluster analysis. *PLoS ONE* 18, 5 (2023), e0286312. https://journals.plos.org/plosone/article?id=10.1371/journal.pone.0286312

[4] Md Firoz Ahmed, Sujit Kumar Mitra, and Rajdeep Mitra. 2021. Ensemble Linear Subspace Analysis of High-Dimensional Data: Theory and Applications. *Mathematics* 9, 21 (2021), 2669. https://www.mdpi.com/2227-7390/9/21/2669

[5] M. Aleksandrov, P. J. Prentice, and F. Wereszczuk. 2021. Voxelisation Algorithms and Data Structures: A Review. *Sensors* 21, 24 (2021), 8241. https://www.mdpi.com/1424-8220/21/24/8241

[6] Amjad Ali, Zardad Khan, Hailiang Du, and Saeed Aldahmani. 2025. Double weighted k nearest neighbours for binary classification of high dimensional genomic data. *Scientific Reports* 15 (2025), 12681. doi:10.1038/s41598-025-97505-2

[7] Imran Ali, Maria Balta, and Thanos Papadopoulos. 2023. Social media platforms and social enterprise: Bibliometric analysis and systematic review. *International Journal of Information Management* 69 (2023). doi:10.1016/j.ijinfomgt.2022.102510

[8] A. A. Amer, S. D. Ravana, and R. A. A. Habeeb. 2025. Effective k-nearest neighbor models for data classification enhancement. *Journal of Big Data* 12 (2025). https://journalofbigdata.springeropen.com/articles/10.1186/s40537-025-01137-2

[9] Arnab Auddy, Dong Xia, and Ming Yuan. 2024. Tensor Methods in High Dimensional Data Analysis: Opportunities and Challenges. *arXiv preprint arXiv:2405.18412* (2024). https://arxiv.org/abs/2405.18412

[10] L. P. Barnes, S. Cameron, and B. Howard. 2025. On Unbiased Low-Rank Approximation with Minimum Distortion. *arXiv preprint arXiv:2505.09647 [cs.DS]*

(2025). https://arxiv.org/abs/2505.09647

[11] Jean Bertin. 2024. Advancing Similarity Search with GenAI: A Retrieval Augmented Generation Approach. *arXiv preprint arXiv:2501.04006 [cs.IR]* (Dec 2024). https://arxiv.org/abs/2501.04006

[12] Sayan Bhattacharya, Monika Henzinger, and Giuseppe F. Italiano. 2018. Deterministic Fully Dynamic Data Structures for Vertex Cover and Matching. *SIAM J. Comput.* 47, 3 (2018), 859–887. https://dblp.org/rec/journals/siamcomp/BhattacharyaHI18

[13] Xingyan Bin, Jianfei Cui, Wujie Yan, Zhichen Zhao, Xintian Han, Chongyang Yan, Feng Zhang, Xun Zhou, Qi Wu, and Zuotao Liu. 2024. Real-time Indexing for Large-scale Recommendation by Streaming Vector Quantization Retriever. *arXiv preprint arXiv:2501.08695* (2024), 1–20. https://arxiv.org/abs/2501.08695

[14] R. Binna, E. Zangerle, M. Pichl, G. Specht, and V. Leis. 2022. Height Optimized Tries. *ACM Transactions on Database Systems* 47, 1 (2022), 1–46. https://dl.acm.org/doi/10.1145/3506692

[15] Gregory Bint, Anil Maheshwari, Michiel H. M. Smid, and Subhas C. Nandy. 2019. Partial Enclosure Range Searching. *International Journal of Computational Geometry & Applications* 29, 1 (2019), 73–93. https://dblp.org/rec/journals/ijcga/BintMSN19

[16] Jean-Daniel Boissonnat, Karthik C. S., and Sébastien Tavenas. 2017. Building Efficient and Compact Data Structures for Simplicial Complexes. *Algorithmica* 79, 2 (2017), 530–567. doi:10.1007/s00453-017-0373-8

[17] G. O. Campos, A. Zimek, J. Sander, R. J. G. B. Campello, B. Micenková, E. Schubert, I. Assent, and M. E. Houle. 2016. On the evaluation of unsupervised outlier detection: measures, datasets, and an empirical study. *Data Mining and Knowledge Discovery* 30, 4 (2016), 891–927. doi:10.1007/s10618-015-0444-8

[18] A. Chaves Carniel. 2024. Defining and designing spatial queries: the role of spatial relationships. *Geo-spatial Information Science* 27, 6 (2024), 1868–1892. https://www.tandfonline.com/doi/full/10.1080/10095020.2022.2163924

[19] Luyao Chang, Fan Li, Xinzheng Niu, and Jiahui Zhu. 2022. On an improved clustering algorithm based on node density for WSN routing protocol. *Cluster Computing* 25, 4 (2022), 3005–3017. doi:10.1007/s10586-022-03544-z

[20] Vasilis Chasiotis, Lin Wang, and Dimitris Karlis. 2024. Efficient subsampling for high-dimensional data. *arXiv preprint arXiv:2411.06298* (2024). https://arxiv.org/abs/2411.06298

[21] Georgios Chatzigeorgakidis, Sophia Karagiorgou, Spiros Athanasiou, and Spiros Skiadopoulos. 2018. FML-kNN: scalable machine learning on Big Data using k-nearest neighbor joins. *Journal of Big Data* 5 (2018), 4. doi:10.1186/s40537-018-0115-x

[22] B. Chen, F. Chen, J. Wang, and T. Qiu. 2025. An efficient and distribution-free symmetry test for high-dimensional data based on energy statistics and random projections. *Computational Statistics & Data Analysis* 206 (2025), 108123. https://www.sciencedirect.com/science/article/abs/pii/S016794732400207X

[23] X. Chen, H. Huo, J. S. Vitter, Y. Hu, and Q. Zhu. 2021. MSQ-Index: A Succinct Index for Fast Graph Similarity Search. *IEEE Transactions on Knowledge and Data Engineering* 33, 6 (2021), 2654–2668. doi:10.1109/TKDE.2019.2954527

[24] Yaru Chen, Jie Zhou, and Xinglong Luo. 2024. An improved density peaks clustering based on sparrow search algorithm. *Cluster Computing* 27, 8 (2024), 11017–11037. doi:10.1007/s10586-024-04384-9

[25] Z. Chen, W. Hao, Z. Zeng, Y. Wen, L. Shi, Z.-J. Wang, and Y. Zhao. 2025. LiLIS: Enhancing Big Spatial Data Processing with Lightweight Distributed Learned Index. *arXiv preprint arXiv:2504.18883v3* (2025). https://arxiv.org/abs/2504.18883

[26] C. L. A. Clarke. 2024. Annotative Indexing. *arXiv preprint arXiv:2411.06256* (2024). https://arxiv.org/abs/2411.06256

[27] V. Cohen-Addad, B. Guedj, V. Kanade, and G. Rom. 2021. Online k-means Clustering. In *Proceedings of the 24th International Conference on Artificial Intelligence and Statistics (AISTATS) (Proceedings of Machine Learning Research, Vol. 130)*. 1126–1134. https://arxiv.org/abs/1909.06861

[28] Sarita de Berg and Frank Staals. 2025. Nearest Neighbor Searching in a Dynamic Simple Polygon. *arXiv preprint arXiv:2503.03435* (2025), 22. https://arxiv.org/abs/2503.03435

[29] E. F. de Oliveira, P. Garg, J. Hjerling-Leffler, R. Batista-Brito, and L. Sjulson. 2025. Identifying patterns differing between high-dimensional datasets with generalized contrastive PCA. *PLOS Computational Biology* 21, 2 (2025), e1012747. doi:10.1371/journal.pcbi.1012747

[30] Naveen Donthu, Satish Kumar, Nitesh Pandey, and Prashant Gupta. 2021. Forty years of the International Journal of Information Management: A bibliometric analysis. *International Journal of Information Management* 57 (2021), 102307. doi:10.1016/j.ijinfomgt.2020.102307

[31] Simeon Emanuilov and Aleksandar Dimov. 2024. Billion-scale Similarity Search Using a Hybrid Indexing Approach with Advanced Filtering. *Cybernetics and Information Technologies* 24, 4 (2024), 45–58. doi:10.2478/cait-2024-0035

[32] Johannes Fischer, Tomohiro I, Dominik Köppl, and Kunihiko Sadakane. 2018. Lempel-Ziv Factorization Powered by Space Efficient Suffix Trees. *Algorithmica* 80, 7 (2018), 2048–2081. doi:10.1007/s00453-017-0354-y

[33] T. Gagie, A. Hartikainen, K. Karhu, J. Kärkkäinen, G. Navarro, S. J. Puglisi, and J. Sirén. 2017. Document retrieval on repetitive string collections. *Information Retrieval Journal* 20 (2017), 273–303. doi:10.1007/s10791-017-9297-7

[34] A. J. Gallego, J. R. Rico-Juan, and J. J. Valero-Mas. 2022. Efficient k-nearest neighbor search based on clustering and adaptive k values. *Pattern Recognition* 122 (2022), 108356. doi:10.1016/j.patcog.2021.108356

[35] Z. Gniazdowski. 2024. New Approach to Clustering Random Attributes. *Zeszyty Naukowe WWSI* 19, 31 (2024), 41–90. doi:10.48550/arXiv.2412.09748

[36] E. Gorstein, R. Aghdam, and C. Solís-Lemus. 2025. HighDimMixedModels.jl: Robust high-dimensional mixed-effects models across omics data. *PLOS Computational Biology* 21, 1 (2025), e1012143. doi:10.1371/journal.pcbi.1012143

[37] Ralf Hartmut Güting, Suvam Kumar Das, Fabio Valdés, and Suprio Ray. 2025. Exact Trajectory Similarity Search With N-tree: An Efficient Metric Index for kNN and Range Queries. *ACM Transactions on Spatial Algorithms and Systems* 11, 1 (2025), 5:1–5:54. doi:10.1145/3716825

[38] R. K. Halder, M. N. Uddin, M. A. Uddin, S. Aryal, and A. Khraisat. 2024. Enhancing K-nearest neighbor algorithm: a comprehensive review and performance analysis of modifications. *Journal of Big Data* 11, 1, Article 113 (2024). https://journalofbigdata.springeropen.com/articles/10.1186/s40537-024-00973-y

[39] S. W. Harrar and X. Kong. 2022. Recent developments in high-dimensional inference for multivariate data: Parametric, semiparametric and nonparametric approaches. *Journal of Multivariate Analysis* 188 (2022), Article 104855. doi:10.1016/j.jmva.2021.104855

[40] Muhammad Umair Hassan, Xiuyang Zhao, Raheem Sarwar, Naif R. Aljohani, S. M. M. Rahman, K. Muhammad, and M. A. Raza. 2024. SODRet: Instance retrieval using salient object detection for self-service shopping. *Machine Learning with Applications* 15 (2024), 100523. https://www.sciencedirect.com/science/article/pii/S2666827023000762

[41] Majid Hojati, Rob Feick, Steven Roberts, Carson Farmer, and Colin Robertson. 2023. Distributed spatial data sharing: a new model for data ownership and access control. *Journal of Spatial Information Science* 2023, 27 (2023), 1–26. doi:10.5311/JOSIS.2023.27.220

[42] Zainab Iftikhar, Adeel Anjum, Abid Khan, Munam Ali Shah, and Gwanggil Jeon. 2023. Privacy preservation in the internet of vehicles using local differential privacy and IOTA ledger. *Cluster Computing* 26 (2023), 3361–3377. doi:10.1007/s10586-023-04002-0

[43] F. Iglesias, T. Zseby, and A. Zimek. 2020. Absolute Cluster Validity. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42, 9 (2020), 2096–2112. https://ieeexplore.ieee.org/document/8695871

[44] Bharathi B. K. and K. Jaganathan. 2022. The Intrinsic Structure of High-Dimensional Data According to Principal Graphs. *Mathematics* 10, 20 (2022), 3894. https://www.mdpi.com/2227-7390/10/20/3894

[45] Kiyonari Kobayashi, Shusuke Shimbo, and Yuji Matsumoto. 2024. Resource-Efficient Index Advisor Utilizing Large Language Model. *arXiv preprint arXiv:2503.07884* (2024). https://arxiv.org/abs/2503.07884

[46] A. Koudounas, C. Papagiannopoulou, L. Rokach, and S. Papadopoulos. 2020. Gradient-based Learning Methods Extended to Similarity-Based Models for Large-Scale Data. *Journal of Artificial Intelligence Research* 69 (2020), 1209–1247. https://jair.org/index.php/jair/article/view/12192/26600

[47] S. Ladra, M. Rodríguez Luaces, J. R. Parama, and F. Silva-Coira. 2024. Compact and indexed representation for LiDAR point clouds. *International Journal of Geographical Information Science* 27, 4 (2024), 1035–1070. doi:10.1080/10095020.2022.2121664

[48] M. Lawson, W. Gropp, and J. Lofstead. 2021. Exploring Spatial Indexing for Accelerated Feature Retrieval in HPC. *arXiv preprint arXiv:2106.13972* (2021). https://arxiv.org/abs/2106.13972

[49] Kuo-Kai Lee, Wing-Kai Hon, Chung-Shou Liao, Kunihiko Sadakane, and Meng-Tsung Tsai. 2023. Fully Dynamic No-Back-Edge-Traversal Forest via 2D-Range Queries. *International Journal of Computational Geometry & Applications* 33, 1&2 (2023), 43–54. https://dblp.org/rec/journals/ijcga/LeeHLST23

[50] J. Li. 2023. Finite sample t-tests for high-dimensional means. *Journal of Multivariate Analysis* 196 (2023), Article 105183. doi:10.1016/j.jmva.2023.105183

[51] J. Li, B. He, and D. Wang. 2021. A Scalable Random-Walk-Based Network Embedding Algorithm with Local Structural Information. *Journal of Artificial Intelligence Research* 71 (2021), 651–683. https://jair.org/index.php/jair/article/view/12567/26689

[52] Y. Li, R. Zhang, Q. Ma, J. Song, B. Zhang, M. Bai, W. Wang, and Y. Li. 2023. CSD-RkNN: reverse k nearest neighbors queries with category-sensitive distance. *International Journal of Geographical Information Science* 37, 8 (2023), 1709–1730. doi:10.1080/13658816.2023.2249521

[53] Anqi Liang, Pengcheng Zhang, Bin Yao, Zhongpu Chen, Yitong Song, and Guangxu Cheng. 2024. UNIFY: Unified Index for Range Filtered Approximate Nearest Neighbors Search. *arXiv preprint arXiv:2412.02448* (2024). https://arxiv.org/abs/2412.02448

[54] J. Lin and A. Trotman. 2017. The role of index compression in score-at-a-time query evaluation. *Information Retrieval Journal* 20 (2017), 274–314. doi:10.1007/s10791-016-9291-5

[55] J. Liu and M. Vinck. 2022. Improved visualization of high-dimensional data using the distance-of-distance transformation. *PLOS Computational Biology* 18, 12 (2022), e1010764. doi:10.1371/journal.pcbi.1010764

[56] Y. Liu, J. Ding, H. Wang, and Y. Du. 2025. A Clustering Algorithm Based on the Detection of Density Peaks and the Interaction Degree Between Clusters. *Applied Sciences* 15, 7 (2025), 1–19. doi:10.3390/app15073612

[57] S. Loodtoy and V. Yalagandula. 2021. Bibliometric Analysis of International Journal of Information Management. *International Journal of Information Management* (2021). http://repo.lib.jfn.ac.lk/ujrr/bitstream/123456789/4718/2/Bibliometric%20Analysis%20of%20International%20Journal%20of%20Information%20Management.pdf

[58] S. Lu, W. Martens, M. Niewerth, and Y. Tao. 2023. Partial Order Multiway Search. *ACM Transactions on Database Systems* 48, 4 (2023), 1–31. doi:10.1145/3626956

[59] Qing Mai, Xin Zhang, Yuqing Pan, and Kai Deng. 2022. A Doubly Enhanced EM Algorithm for Model-Based Tensor Clustering. *J. Amer. Statist. Assoc.* 117, 540 (2022), 2120–2134. doi:10.1080/01621459.2021.1904959

[60] A.-A. Mamun, H. Wu, Q. He, J. Wang, and W. G. Aref. 2024. A Survey of Learned Indexes for the Multi-dimensional Space. *arXiv preprint arXiv:2403.06456* (2024). https://arxiv.org/abs/2403.06456

[61] Magdalen Dobson Manohar, Taekseung Kim, and Guy E. Blelloch. 2025. Range Retrieval with Graph-Based Indices. *arXiv preprint arXiv:2502.13245* (2025). https://arxiv.org/abs/2502.13245

[62] N. Marco, D. Şentürk, S. Jeste, C. C. DiStefano, A. Dickinson, and D. Telesca. 2024. Flexible regularized estimation in high-dimensional mixed membership models. *Computational Statistics & Data Analysis* 194 (2024), 107931. doi:10.1016/j.csda.2024.107931

[63] Jorge Martinez-Gil. 2022. Evaluation of Code Similarity Search Strategies in Large-Scale Codebases. *Machine Learning with Applications* 10 (2022), 100423. https://www.sciencedirect.com/science/article/pii/S2666827022000868

[64] A. Michalopoulos, D. Tsitsikgos, P. Bouros, N. Mamoulis, and M. Terrovitis. 2025. Efficient Distance Queries on Non-point Data. *ACM Transactions on Spatial Algorithms and Systems* 11, 1 (2025), 1:1–1:37. doi:10.1145/3698194

[65] Xiangbo Mo and Hao Chen. 2024. A new classification framework for high-dimensional data. *arXiv preprint arXiv:2306.15199* (2024). https://arxiv.org/abs/2306.15199

[66] Neda Dousti Mousavi, S. Mostafa Hosseini, and Mahdi Mahmoudi. 2023. Categorical Data Analysis for High-Dimensional Sparse Covariates with Multinomial Responses: An RNA-Seq Cancer Application. *Mathematics* 11, 14 (2023), 3202. https://www.mdpi.com/2227-7390/11/14/3202

[67] Abhinav Natarajan, Maria De Iorio, Andreas Heinecke, Emanuel Mayer, and Simon Glenn. 2024. Cohesion and Repulsion in Bayesian Distance Clustering. *J. Amer. Statist. Assoc.* 119, 546 (2024), 1374–1384. doi:10.1080/01621459.2023.2191821

[68] H. Yepdjio Nkouanga and S. Vajda. 2023. Optimization Strategies for the k-Nearest Neighbor Classifier. *SN Computer Science* 4, 47 (2023). doi:10.1007/s42979-022-01469-3

[69] Daniel Obraczka and Erhard Rahm. 2022. Fast Hubness-Reduced Nearest Neighbor Search for Entity Alignment in Knowledge Graphs. *SN Computer Science* 3, 6 (2022), 501. doi:10.1007/s42979-022-01417-1

[70] A. Pakzad, V. Mehrjou, D. Khosla, and B. Schölkopf. 2021. A Word Selection Method for Producing Interpretable Word Embeddings. *Journal of Artificial Intelligence Research* 71 (2021), 867–900. https://jair.org/index.php/jair/article/download/13353/26748/29105

[71] V. Pandey, A. van Renen, E. T. Zacharatou, A. Kipf, I. Sabek, J. Ding, V. Markl, and A. Kemper. 2023. Enhancing In-Memory Spatial Indexing with Learned Search. *arXiv preprint arXiv:2309.06354* (2023). https://arxiv.org/abs/2309.06354

[72] Y. Pang, X. Zhou, J. Zhang, Q. Sun, and J. Zheng. 2022. Hierarchical electricity time series prediction with cluster analysis and sparse penalty. *Pattern Recognition* 126 (2022), 108599. doi:10.1016/j.patcog.2022.108599

[73] J. Paparrizos, F. Yang, and H. Li. 2024. Bridging the Gap: A Decade Review of Time-Series Clustering Methods. *arXiv preprint arXiv:2412.20582* (2024). https://arxiv.org/abs/2412.20582

[74] R. M. Perera, B. Oetomo, B. I. P. Rubinstein, R. Borovica-Gajic, and M. Roughan. 2023. No DBA? No Regret! Multi-Armed Bandits for Index Tuning of Analytical and HTAP Workloads With Provable Guarantees. *IEEE Transactions on Knowledge and Data Engineering* 35, 12 (2023), 12221–12237. doi:10.1109/TKDE.2023.3271664

[75] Nathan Phelps and Adam Metzler. 2024. An exploratory clustering analysis of the 2016 National Financial Well-Being Survey. *PLOS ONE* 19, 9 (2024), e0309260. doi:10.1371/journal.pone.0309260

[76] Alberto Policriti and Nicola Prezza. 2018. LZ77 Computation Based on the Run-Length Encoded BWT. *Algorithmica* 80, 7 (2018), 1986–2011. doi:10.1007/s00453-017-0379-2

[77] Yifan Qiao, Shiyu Ji, Changhai Wang, Jinjin Shao, and Tao Yang. 2023. Privacy-aware document retrieval with two-level inverted indexing. *Information Retrieval Journal* 26 (2023). doi:10.1007/s10791-023-09428-z

[78] A. Rachwał, A. Popławska, and M. Borys. 2023. Determining the Quality of a Dataset in Clustering Terms. *Applied Sciences* 13, 5 (2023), 1–22. doi:10.3390/app13052942

[79] S. Rahul. 2021. Approximate range counting revisited. *Journal of Computational Geometry* 12, 1 (2021), 183–212. https://jocg.org/index.php/jocg/article/view/3153

[80] S. Ray and B. Nickerson. 2022. Temporally relevant parallel top-k spatial keyword search. *Journal of Spatial Information Science* 24 (2022), 1–36. https://josis.org/index.php/josis/article/view/199

[81] R. Reinanda, E. Meij, and M. de Rijke. 2020. Knowledge Graphs: An Information Retrieval Perspective. *Foundations and Trends® in Information Retrieval* 14, 4 (2020), 289–444. https://www.nowpublishers.com/article/Details/INR-063

[82] M. Z. Rodriguez, C. H. Comin, D. Casanova, O. M. Bruno, D. R. Amancio, L. da F. Costa, and F. A. Rodrigues. 2019. Clustering algorithms: A comparative approach. *PLoS ONE* 14, 1 (2019), e0210236. doi:10.1371/journal.pone.0210236

[83] J. E. Schmitz, J. Zentgraf, and S. Rahmann. 2025. Smaller and More Flexible Cuckoo Filters. *arXiv preprint arXiv:2505.05847* (2025). https://arxiv.org/abs/2505.05847

[84] Patrick Schäfer, Jakob Brand, Ulf Leser, Botao Peng, and Themis Palpanas. 2024. Fast and Exact Similarity Search in less than a Blink of an Eye. *arXiv preprint arXiv:2411.17483* (Dec. 2024). https://arxiv.org/abs/2411.17483

[85] Nijaguna Gollara Siddappa and Thippeswamy Kampalappa. 2020. Imbalance Data Classification Using Local Mahalanobis Distance Learning Based on Nearest Neighbor. *SN Computer Science* 1, 76 (2020), 1–15. doi:10.1007/s42979-020-0085-x

[86] S. Song and X. Liang. 2024. Federated Pseudo-Sample Clustering Algorithm: A Label-Personalized Federated Learning Scheme Based on Image Clustering. *Applied Sciences* 14, 6 (2024), 1–18. doi:10.3390/app14062345

[87] Liyang Sun, Yujing Wang, Zejian Wang, Xinyi Wu, Xiangming Dou, Jinji Li, Yicheng Bai, Xuerui Wang, Weinan Zhang, Yong Yu, and Zhenguo Li. 2024. The Disruption Index Measures Displacement Between a Paper and Its Citations. *arXiv preprint arXiv:2504.04677* (2024). https://arxiv.org/abs/2504.04677

[88] B. Tang, H. He, and S. Zhang. 2020. MCENN: A variant of extended nearest neighbor method for pattern recognition. *Pattern Recognition Letters* 133 (2020), 116–122. doi:10.1016/j.patrec.2020.01.015

[89] S. Uddin, I. Haque, H. Lu, M. A. Moni, and E. Gide. 2022. Comparative performance analysis of K-nearest neighbour (KNN) algorithm and its different variants for disease prediction. *Scientific Reports* 12 (2022). https://www.nature.com/articles/s41598-022-10358-x

[90] Katerina Vrotsou, Georg Fuchs, Natalia Andrienko, and Gennady Andrienko. 2017. An Interactive Approach for Exploration of Flows Through Direction-Based Filtering. *Journal of Geovisualization and Spatial Analysis* 1, 1 (2017), 1–21. doi:10.1007/s41651-017-0001-7

[91] H. Wang and Q. Zeng. 2021. Unit-disk range searching and applications. *Journal of Computational Geometry* 12, 1 (2021), 381–417. https://jocg.org/index.php/jocg/article/view/4015

[92] H. Wang, J. Zhang, Y. Wei, Y. Wang, X. Zhang, and J. Pei. 2023. Neural Similarity Search on Supergraph Containment. *IEEE Transactions on Knowledge and Data Engineering* 35, 11 (2023), 11200–11214. doi:10.1109/TKDE.2023.3279920

[93] S. Wang, L. Qin, J. X. Yu, R. Jin, and L. Chang. 2020. Continuously Adaptive Similarity Search. *ACM Transactions on Information Systems* 38, 3 (2020), 28:1–28:28. https://dl.acm.org/doi/10.1145/3318464.3380601

[94] H. Wei, P. Li, H. Gao, and C. Wang. 2017. String Similarity Search: A Hash-Based Approach. *IEEE Transactions on Knowledge and Data Engineering* 29, 7 (2017), 1371–1385. doi:10.1109/TKDE.2017.2692024

[95] N. Wiroonsri. 2024. Clustering performance analysis using a new correlation-based cluster validity index. *Pattern Recognition* 145 (2024), 109910. doi:10.1016/j.patcog.2023.109910

[96] G. Wu, J. Zhang, J. Fu, and J. Wang. 2022. A case study for Adaptive Radix Tree index. *Information Systems* 106 (2022), 101920. https://www.sciencedirect.com/science/article/abs/pii/S0306437921001228

[97] Y. Wu, X. Zhou, Y. Zhang, L. Ma, and J. Fan. 2024. Automatic Index Tuning: A Survey. *IEEE Transactions on Knowledge and Data Engineering* 36, 12 (2024), 7657–7676. https://ieeexplore.ieee.org/document/10582533

[98] Jie Xue, Yuan Li, Saladi Rahul, and Ravi Janardan. 2020. Searching for the closest-pair in a query translate. *Journal of Computational Geometry* 11, 2 (2020), 1–33. doi:10.20382/jocg.v11i2a3

[99] J. Yang and C.-T. Lin. 2025. Autonomous clustering by fast find of mass and distance peaks. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 47, 1 (2025), 1–14. doi:10.1109/TPAMI.2025.40031325

[100] Mingyu Yang, Wentao Li, and Wei Wang. 2025. Fast High-dimensional Approximate Nearest Neighbor Search with Efficient Index Time and Space. *arXiv preprint arXiv:2411.06158* (2025), 8. https://arxiv.org/abs/2411.06158

[101] H. Yin, A. Aryani, S. Petrie, A. Nambissan, A. Astudillo, and S. Cao. 2024. A Rapid Review of Clustering Algorithms. *arXiv preprint arXiv:2401.07389* (2024). https://arxiv.org/abs/2401.07389

[102] Y. Yin. 2021. Test for high-dimensional mean vector under missing observations. *Journal of Multivariate Analysis* 186 (2021), Article 104797. doi:10.1016/j.jmva.2021.104797

[103] Huacheng Yu. 2022. Nearly Optimal Static Las Vegas Succinct Dictionary. *SIAM J. Comput.* 51, 3 (2022), 174–249. doi:10.1137/20M1363649

[104] P. Yuan, C. Jin, and G. Li. 2024. FDR control for linear log-contrast models with high-dimensional compositional covariates. *Computational Statistics Data Analysis* 197 (2024), 107973. https://www.sciencedirect.com/science/article/abs/

pii/S0167947324000574

[105] Z. Yuan and C. L. Philip Chen. 2023. Forgetful Forests: Data Structures for Machine Learning on Data Streams with Incremental Computation and Filtering. *Algorithms* 16, 6 (2023), 278. doi:10.3390/algorithms16060278

[106] R. Zanibbi, B. Mansouri, and A. Agarwal. 2025. Mathematical Information Retrieval: Search and Question Answering. *Foundations and Trends® in Information Retrieval* 19, 1–2 (2025), 1–190. https://www.nowpublishers.com/article/Details/INR-095

[107] D. Zhang, Y. Huang, H. Wang, D. Yang, Z. He, and J. Xu. 2021. Continuous Trajectory Similarity Search for Online Outlier Detection. *IEEE Transactions on Knowledge and Data Engineering* 33, 10 (2021), 3405–3419. doi:10.1109/TKDE.2020.3046670

[108] J. Zhang, J. Tang, C. Ma, X. Chen, Y. Liu, and J. Li. 2018. Fast and Flexible Top-k Similarity Search on Large Networks. *ACM Transactions on Information Systems* 36, 2 (2018), 14:1–14:34. doi:10.1145/3086695

[109] Y. Zhang, M. Xiang, and B. Yang. 2017. Graph regularized nonnegative sparse coding using incoherent dictionary for approximate nearest neighbor search. *Pattern Recognition* 70 (Oct. 2017), 75–88. doi:10.1016/j.patcog.2017.05.004

[110] Xiaoyao Zhong, Haotian Li, Jiabao Jin, Mingyu Yang, Deming Chu, Xiangyu Wang, Zhitao Shen, Wei Jia, George Gu, Yi Xie, Xuemin Lin, Heng Tao Shen, Jingkuan Song, and Peng Cheng. 2025. VSAG: An Optimized Search Framework for Graph-based Approximate Nearest Neighbor Search. *arXiv preprint arXiv:2503.17911* (2025), 16. https://arxiv.org/abs/2503.17911

[111] S. Zhou, H. Xu, Z. Zheng, J. Chen, Z. Li, J. Bu, J. Wu, X. Wang, W. Zhu, and M. Ester. 2022. A Comprehensive Survey on Deep Clustering: Taxonomy, Challenges, and Future Directions. *arXiv preprint arXiv:2206.07579* (2022). https://arxiv.org/abs/2206.07579

[112] F. Zhu, Y. Kou, X. Jia, and Y. Zhu. 2023. An Efficient and Robust Semantic Hashing Framework for Similarity Search. *ACM Transactions on Information Systems* 41, 2 (2023), 1–30. doi:10.1145/3570725