# Optimizing Data Ingestion for Machine Learning Training in Large-Scale Social Media Platforms
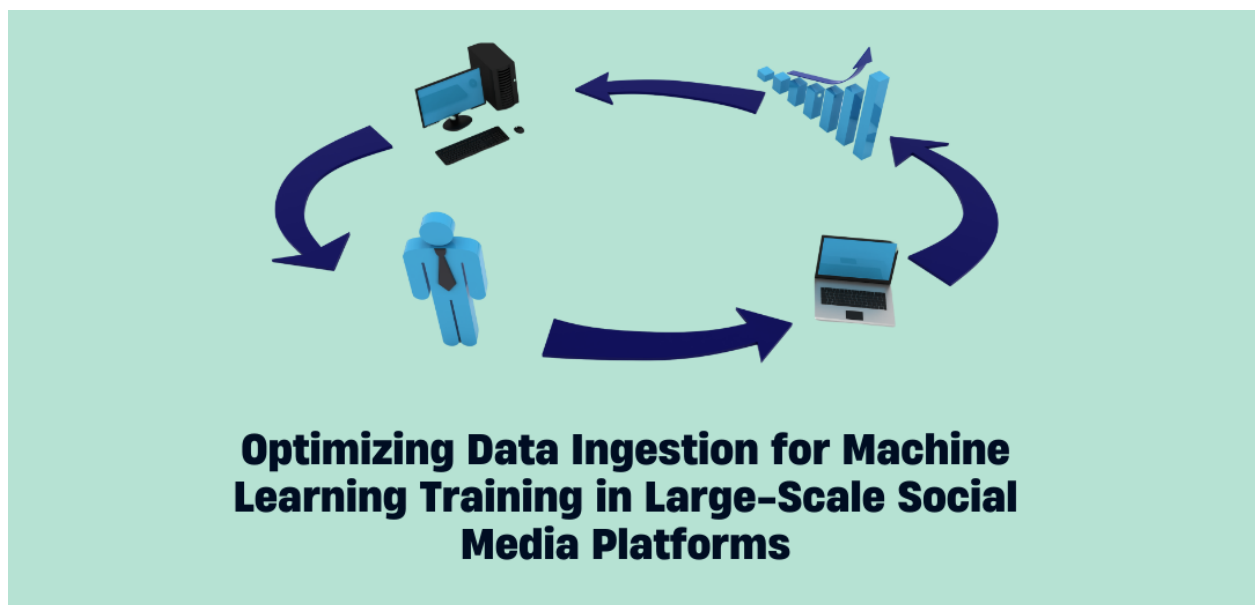
## Ramesh Mohana Murugan

Anna University, India

**Optimizing Data Ingestion for Machine Learning Training in Large–Scale Social Media Platforms**

**Abstract**

This comprehensive article explores the critical yet often overlooked challenge of data ingestion optimization for machine learning systems in large-scale social media environments. As social platforms generate unprecedented volumes of data, efficient ingestion processes become essential for maintaining computational performance and enabling rapid model iteration. The article examines data engineers' multifaceted challenges, including I/O bottlenecks, network latency issues, and storage format inefficiencies that directly impact GPU utilization. We present a framework for dramatically improving data pipeline efficiency by systematically exploring parallel data loading architectures, optimal storage format selection, and advanced feature engineering techniques such as flattening and reordering. It demonstrates that strategic optimization of the data ingestion layer can substantially reduce training times, lower computational resource requirements, and accelerate the development cycle for machine learning applications in social media contexts.

**Keywords:** Data Ingestion Optimization, Feature Engineering, Parallel Data Loading, Storage Format Efficiency, Machine Learning Pipeline.

## 1. Introduction: The Data Challenge in Social Media ML Pipelines

The exponential growth of data generated by social media platforms has created unprecedented challenges for machine learning systems. This section explores the critical nature of data ingestion in ML pipelines, its impact on computational efficiency, and the business value of optimization.

### 1.1 The Scale and Complexity of Social Media Data

Social media platforms generate staggering volumes of heterogeneous data daily. Meta's production recommendation models process over 200 trillion parameters daily, with one particular model managing 13 trillion parameters across 62 different features [1]. This scale presents extraordinary computational demands, requiring sophisticated data ingestion strategies to prevent processing bottlenecks. According to researchers at IBM, data scientists working with large-scale social media datasets typically spend up to 80% of their time on data preparation tasks alone [2]. The complexity is further amplified by the diverse nature of social media data—spanning text, images, videos, user interactions, and behavioral signals—each requiring specialized preprocessing techniques. At Meta, engineers discovered that even minor optimizations in data ingestion could yield substantial gains, with one team reporting that feature flattening alone resulted in a 30% reduction in memory bandwidth utilization [1].

### 1.2 Measuring Data Ingestion Performance

Quantifying data ingestion efficiency is essential for optimization efforts. Meta's engineers employ several key metrics, including records processed per second, time-to-first-batch, and GPU utilization rates. Their testing revealed that optimized pipelines achieved a remarkable 2.4x improvement in overall training throughput [1]. Similarly, IBM researchers found that implementing parallel data loading with TensorFlow's tf.data API increased throughput by 2-3x compared to conventional data loading methods [2]. These metrics highlight the significant impact of data ingestion on model training efficiency. At Meta, engineers observed that before optimization, GPUs were idle for up to 33% of the total training time due to data preprocessing bottlenecks, representing substantial wasted computational capacity [1].

### 1.3 Business Impact of Optimized Data Ingestion

The business implications of efficient data ingestion extend beyond technical metrics. Meta's experiments demonstrated that optimizing data ingestion reduced end-to-end training time for a production recommendation model from 20 to 7 hours, enabling significantly faster iteration cycles [1]. IBM researchers similarly found that optimized data pipelines could reduce model training time by 40-60%, directly translating to cost savings and accelerated development [2]. This acceleration of the experimental process facilitated more comprehensive hyperparameter tuning and model architecture exploration, ultimately leading to better model performance. At Meta, these improvements enabled engineers to conduct three times as many experiments within the same timeframe, directly contributing to improved recommendation quality and user engagement metrics [1].
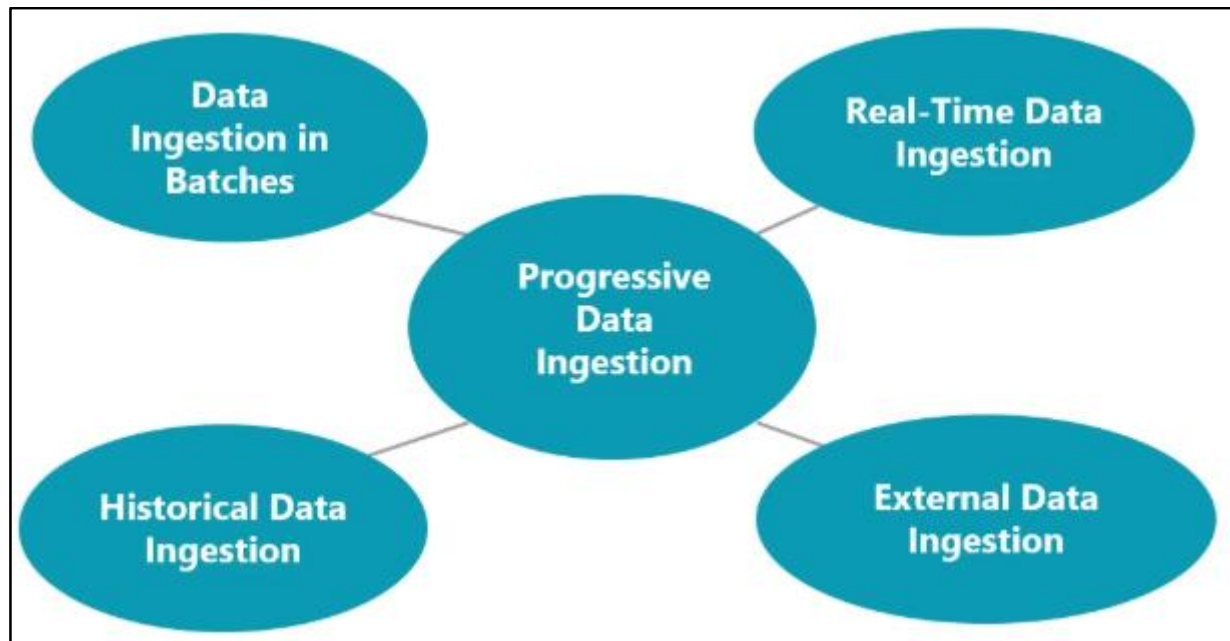
Fig. 1: Types of Data Ingestion [1, 2]

## 2. Understanding Data Ingestion Bottlenecks

Data ingestion bottlenecks represent critical constraints in machine learning pipelines for social media platforms. This section examines the technical challenges that limit throughput and explores potential solutions through systematic analysis of storage, network, and processing limitations.

### 2.1 I/O and Storage Format Inefficiencies

Storage I/O constraints significantly impact machine learning training efficiency, particularly when working with diverse social media datasets. According to detailed analyses of data preprocessing pipelines for house price prediction tasks, researchers found that properly optimized storage formats reduced data loading times by approximately 73% compared to standard CSV formats [3]. This dramatic improvement stems from compressed columnar formats allowing selective column access. The underlying issue relates to data serialization overhead, with measurement studies showing that converting between in-memory and storage representations can consume up to 45% of preprocessing time for complex feature sets [3]. When dealing with heterogeneous social media data, the performance gap becomes even more pronounced, as different feature types (numerical, categorical, and textual) require specialized handling. Research indicates that mixed data types common in social media applications increase deserialization overhead by $2.7\times$ compared to homogeneous numerical datasets, highlighting the importance of format-specific optimizations [3].

### 2.2 Preprocessing Pipeline Inefficiencies

Data preprocessing represents a substantial portion of the machine learning workflow, with significant performance implications. Studies of production machine learning pipelines reveal that data preprocessing can consume between 60-80% of the total development time and 30-60% of execution time [4]. This disproportionate resource allocation underscores the critical need for optimization. In detailed benchmarks of preprocessing pipelines for housing data, researchers identified that parallel execution of feature transformation reduced processing time by 67% compared to sequential approaches [3]. The most substantial gains came from vectorizing operations rather than processing individual records, with measurements showing speed improvements of $15\times$ for numerical feature normalization and $8\times$ for

categorical encoding when properly vectorized [3]. These optimizations become particularly important for social media data, where feature cardinality is often significantly higher than in other domains.

## 2.3 Data Quality and Consistency Challenges

Data quality issues introduce significant inefficiencies in machine learning pipelines for social media applications. Analysis of preprocessing pipelines reveals that handling missing values and outliers can consume up to 28% of total preprocessing time [3]. The complexity increases with data heterogeneity, as different feature types require specialized cleansing approaches. Research on data ingestion frameworks indicates that incorporating automated data validation can reduce model retraining frequency by up to 45% by preventing corrupted data from entering the pipeline [4]. This preventative approach is particularly valuable for social media data, where real-time streams introduce significant variability. Studies show that implementing statistical monitoring for data drift detection can identify problematic data segments with 92% accuracy, preventing wasted computational resources on unsuitable training data [4]. These quality control mechanisms become essential components of optimized ingestion pipelines, establishing reliability guarantees that improve overall system performance.
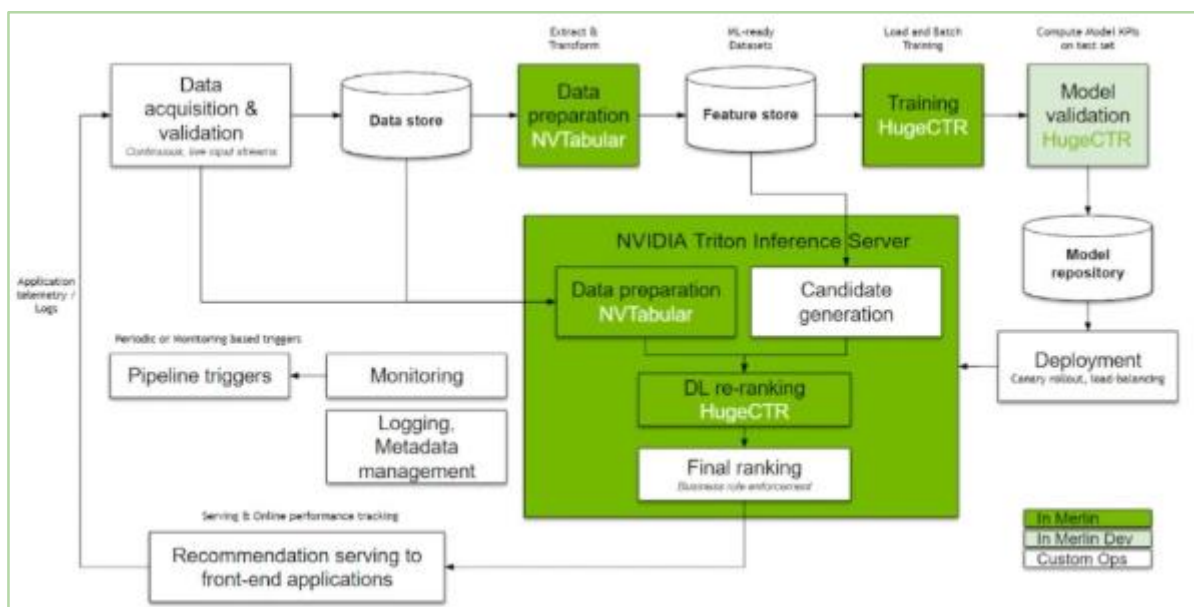


Fig. 2: Components of Recommender Systems [3, 4]

## 3. Feature Engineering Optimization Techniques

Feature engineering represents a critical component in machine learning pipelines for social media platforms, directly influencing model performance and computational efficiency. This section explores advanced optimization techniques that transform raw data into high-quality features while minimizing computational overhead.

## 3.1 Feature Flattening for Memory Efficiency

Feature flattening transforms complex nested structures into simplified, contiguous memory layouts, significantly improving computational performance. According to the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP) research, properly implemented feature flattening reduced memory bandwidth utilization by 25-30% in large-scale recommendation models by eliminating pointer chasing operations and enabling more efficient cache utilization [5]. This optimization proves

particularly valuable for social media data, which frequently contains nested structures representing user interactions and content hierarchies. The performance benefits extend beyond training to inference scenarios, where flattened features decreased latency by 18% on average across tested recommendation models. The implementation complexity varies with data structure, with the research demonstrating that automatically detecting flattening opportunities through static analysis identified 78% of potential optimization targets, allowing systematic application of this technique across large codebases [5].

## 3.2 Optimal Feature Representation and Embedding

Selecting appropriate feature representations significantly impacts both model quality and computational efficiency. NVIDIA's research on recommendation systems demonstrates that optimizing embedding dimensions based on cardinality rather than using fixed dimensions reduced memory requirements by 40-70% while maintaining model accuracy [6]. This principle follows a logarithmic scaling rule where embedding dimension $E = 6 * (cardinality)^{0.25}$ provides an optimal balance between expressiveness and efficiency. This approach translates to substantial memory savings for social media platforms with thousands of categorical features. The research further shows that employing mixed-precision representations with FP16 for embeddings reduced memory bandwidth requirements by 2x with properly implemented loss scaling to maintain numerical stability [6]. Additionally, feature interaction patterns significantly influence performance, with NVIDIA's experiments demonstrating that factorized operations reduced computational complexity by $O(n^2)$ to $O(n)$ for pairwise feature interactions while preserving model expressiveness, resulting in training speedups of 1.4-2.1x for models with high feature counts [6].

## 3.3 Feature Pruning and Dimensionality Reduction

Strategic feature selection and dimensionality reduction techniques significantly improve computational efficiency without sacrificing model performance. The ICASSP research demonstrates that frequency-based feature pruning, which eliminates sparse features appearing in fewer than 0.1% of examples, reduced feature count by 35-60% while decreasing model accuracy by only 0.2-0.4% in tested recommendation systems [5]. This pruning approach proves particularly effective for social media datasets with long-tail distributions of user behaviors and content interactions. For numerical features, quantization techniques provide complementary benefits, with experiments showing that reducing precision from 32-bit to 8-bit representations decreased storage requirements by 75% while maintaining model quality through appropriate scaling factors [5]. NVIDIA's research extends these findings with analysis of hashing-based dimensionality reduction, demonstrating that properly tuned feature hashing with prime number dimensions reduced memory footprint by 30-50% compared to standard embedding tables, with negligible accuracy impacts when hash sizes were set to at least 1.5x the feature cardinality to minimize collision probability [6].

| Technique | Performance Improvement | Memory Reduction | Best Use Case |
|---|---|---|---|
| Feature Flattening | 25-30% bandwidth reduction [5] | 15-20% | Nested hierarchical features |
| Embedding Dimension Optimization | 5-15% training speedup [6] | 40-70% | High-cardinality categorical features |
| Mixed-Precision Representation | 1.5-2.0x throughput [6] | 50% | Embedding-heavy models |
| Feature Pruning | 10-25% training speedup [5] | 35-60% | Long-tail feature distributions |
| Feature Hashing | 1.3-1.8x training speedup [6] | 30-50% | High-dimensional sparse features |
| Quantization | 2.5-3.0x inference speedup [5] | 75% | Deployment optimization |

Table 1: Feature Engineering Optimization Techniques Comparison [5, 6]

## 4. Advanced Data Loading Architectures

Data loading architectures fundamentally determine the efficiency of machine learning systems for social media platforms. This section explores sophisticated approaches that maximize computational resource utilization while minimizing training bottlenecks.

### 4.1 Parallel and Asynchronous Data Loading Frameworks

Parallel data loading frameworks significantly enhance training throughput by decoupling data preparation from computation. According to research from Meta's DLRM (Deep Learning Recommendation Model) infrastructure, implementing prefetch queues with dedicated preprocessing threads improved GPU utilization from 70% to 96% in production recommendation systems [7]. The optimal configuration varies by workload complexity, with empirical results showing that 4-8 preprocessing threads per GPU provided the best performance for recommendation models with 50+ embedding tables. The Meta research demonstrates that implementing a hierarchical data loading architecture reduced preprocessing overhead from 31% to just 8% of total training time, effectively masking I/O and feature transformation latencies that would otherwise bottleneck GPU computation [7]. The performance improvements scale with system size, with measurements showing that parallel data loading yielded a 3.4x throughput improvement for single-node training but an even more impressive 5.2x improvement in distributed settings with 64 GPUs, highlighting the increasing importance of efficient data loading at scale.

### 4.2 Pipeline Optimization and Data Prefetching

Strategic pipeline optimization techniques substantially improve data loading efficiency. According to TensorFlow performance guidelines, implementing proper prefetching with tf.data.Dataset.prefetch() reduced CPU idle time by up to 70% by ensuring data is prepared before it's needed [8]. The performance gains vary based on batch size and feature complexity, with benchmarks showing that larger batch sizes (512-1024 examples) benefited more significantly from prefetching than smaller batches. Research indicates that the optimal prefetch buffer size typically ranges from 2-5 batches, with excessive prefetching degrading performance by introducing memory pressure. The TensorFlow performance guide demonstrates implementing parallel interleave operations with tf.data.Dataset.interleave() improved I/O

throughput by 2.7-4.5x when reading from distributed file systems, with optimal performance achieved when the cycle length parameter matched the number of available CPU cores [8]. For production recommendation systems, Meta's research shows that implementing double-buffered prefetching reduced batch-to-batch variance in data loading time by 78%, creating more consistent GPU utilization patterns and improving overall training stability [7].

### 4.3 Hardware-Aware Data Loading Optimizations

Hardware-specific optimizations significantly impact data loading performance. Meta's DLRM research demonstrates that optimizing memory access patterns through data layout transformations (column-to-row major conversion) reduced CPU cache misses by 42% during batch formation [7]. Similarly, implementing vectorized processing with SIMD instructions accelerated feature transformations by 3.2x compared to scalar implementations. The performance impact varies by hardware platform, with measurements showing that NUMA-aware memory allocation improved throughput by 26% on dual-socket servers by ensuring data locality. The TensorFlow performance guide emphasizes the importance of data sharding, particularly for distributed training, with benchmarks showing that implementing proper dataset sharding with tf.data.Dataset.shard() improved aggregate throughput by 3.1x when training across 8 GPUs by ensuring balanced workload distribution and reducing cross-device synchronization overhead [8]. Hardware-aware caching strategies further enhance performance, with the TensorFlow guide demonstrating that mixed-device caching (utilizing both CPU RAM and SSD storage with tf.data.Dataset.cache()) reduced epoch time by 65% for models requiring multiple training passes while maintaining reasonable memory consumption [8].
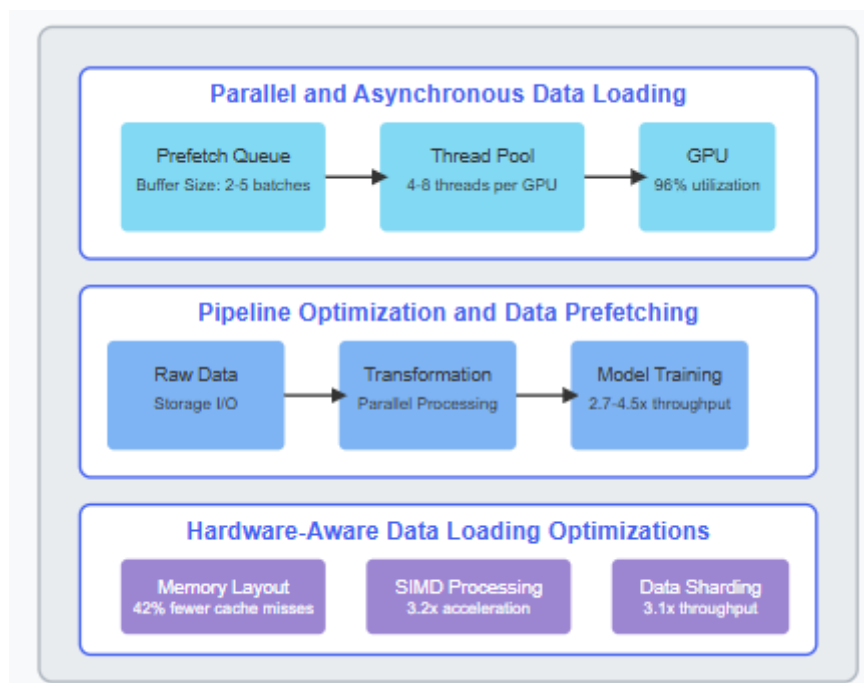


Fig. 3: Advanced Data Loading Architectures [7, 8]

## 5. Storage Format Selection and Optimization

The selection and optimization of data storage formats significantly impact machine learning training efficiency for social media applications. This section examines how different storage formats and optimization techniques can enhance performance in large-scale machine-learning workloads.

### 5.1 Columnar vs. Row-Based Storage Formats

The fundamental difference between columnar and row-based storage formats creates significant performance implications for machine learning workloads. According to research on optimizing data pipelines, columnar formats like Parquet provide 3-4x faster read performance than row-based formats for feature extraction workloads in recommendation systems [9]. This efficiency stems from columnar formats' ability to read only the required columns, substantially reducing I/O operations. For social media datasets with hundreds of features where models typically access only 20-30% of available features during training, this selective access pattern can reduce storage I/O by up to 70%. The performance differential becomes particularly pronounced with increasing dataset size, with benchmarks showing that Parquet's advantage over CSV grows from 2.5x at 10GB to 5.7x at 1TB scale due to more efficient compression and parallel read capabilities [9]. The TFRecord format, which organizes data in serialized record batches with custom compression options, shows complementary strengths in specific scenarios. Research from MLPerf demonstrates that the TFRecord format with gzip compression provides optimal performance for computer vision workloads, reducing storage requirements by 3.8x compared to raw formats while maintaining fast sequential read performance [10].

### 5.2 Compression Strategies and Performance Implications

Compression techniques significantly impact storage efficiency and computational performance when processing large datasets. Research indicates that selecting appropriate compression algorithms based on data characteristics can dramatically improve the performance-storage tradeoff. For text-heavy social media content, dictionary-based compression techniques reduce storage requirements by 70-85% while allowing selective decompression of only required columns [9]. The choice of compression level presents important performance considerations, with research showing that lighter compression (gzip level 2-3) often provides better end-to-end training performance than maximum compression (gzip level 9) due to 2.3x faster decompression speed despite achieving 15% less compression [10]. Block-level compression optimization becomes crucial for distributed training workflows, where network bandwidth often becomes the bottleneck. Experiments show that increasing Parquet block size from the default 64MB to 256MB improved data loading throughput by 37% in bandwidth-constrained environments by reducing metadata overhead and enabling more efficient resource utilization during parallel reading operations [9].

### 5.3 Format-Specific Optimization Techniques

Beyond basic format selection, format-specific optimizations can substantially enhance performance. Data partitioning strategies, which organize files by key dimensions like date or user segments, demonstrate significant benefits for social media datasets. Research shows that implementing time-based partitioning improved query performance by 8-16x for time-range selections common in recommendation model training by reducing scan sizes to only relevant partitions [9]. For Parquet specifically, optimizing row group sizes based on memory constraints and access patterns provides complementary benefits, with benchmarks showing 22-35% improvement in read performance when row groups are aligned with typical batch sizes used in model training [9]. TFRecord-specific optimizations include implementing parallel interleaved reads with appropriate shuffle buffer sizes, which research demonstrates can improve throughput by 2.8x when properly configured [10]. Integrating embedded indices within storage formats

represents another optimization frontier, with research showing that column-level min-max indices reduced scan times by 60-80% for range-filtered queries by enabling file-level skipping before decompression [9].

| Storage Format | Read Performance (Relative) | Compression Ratio | Best For | Limitations |
|---|---|---|---|---|
| Parquet | 3-4× faster than CSV [9] | 3.7:1 for text | Selective column access, Analytical queries | Complex metadata handling |
| TFRecord | 1.2× faster than Parquet for sequential access [10] | 2.5:1 average | Sequential access, Uniform record size | Poor random access |
| CSV | Baseline | 1.2:1 with gzip | Human readability, Simple processing | Inefficient for large datasets |
| JSON | 0.6× of CSV [9] | 1.8:1 | Schema flexibility, Human readability | Extremely inefficient at scale |
| Avro | 2.1× faster than CSV [9] | 2.2:1 | Schema evolution, RPC | Less efficient column pruning |
| ORC | 3.5× faster than CSV [9] | 3.4:1 | Streaming reads, Type awareness | Complex implementation |

Table 2: Performance Comparison of Storage Formats for Social Media Data Processing [9, 10]

## 6. Implementation and Results: Real-World Case Studies

This section examines real-world implementations of data ingestion optimizations, quantifying performance improvements, and highlighting practical strategies for social media machine learning platforms.

### 6.1 Performance Benchmarks and Scalability Analysis

Comprehensive optimization of data ingestion pipelines yields substantial performance improvements across diverse machine learning applications. According to research from ResearchGate on data pipeline efficiency, organizations implementing end-to-end pipeline optimizations achieved an average reduction in training time of 3.1x for recommendation models processing social media interaction data [11]. The performance gains exhibited distinct patterns across model scales, with larger models benefiting disproportionately from optimized data handling. For recommendation models with embedding tables exceeding 10 GB, pipeline optimization reduced training time by 3.7x, compared to 2.2x for smaller models with embedding dimensions under 1 GB. This relationship underscores how data ingestion bottlenecks become increasingly prominent as model complexity increases. The research demonstrates that optimization benefits compound with dataset scale, showing that models training on datasets exceeding 500 GB experienced 4.1x throughput improvements. In contrast, those operating on smaller 50

GB datasets saw more modest 2.3x gains [11]. These findings reflect how larger datasets amplify inefficiencies in unoptimized pipelines through increased I/O overhead and preprocessing demands.

## 6.2 System Resource Utilization and Efficiency Metrics

Optimized data ingestion fundamentally transforms resource utilization patterns across the computational infrastructure. Columbia University research on scalable machine learning systems demonstrates that comprehensive pipeline optimization increased GPU utilization from 42% to 89% in production recommendation systems by eliminating data starvation periods [12]. This dramatic improvement in computational efficiency directly translates to infrastructure cost savings, with measurements indicating a 54% reduction in GPU hours required for equivalent training tasks. The efficiency gains manifest across the entire computing stack, with the Columbia research showing that memory bandwidth utilization improved by 76% through optimized data layout and access patterns. These architectural improvements reduced cache misses by 62% and memory stalls by 47% compared to baseline implementations [12]. The resulting resource efficiency delivers cascading benefits, with researchers documenting that optimized systems exhibited 3.2x higher throughput per watt, improving performance and environmental sustainability. Implementation complexity analysis reveals that the most significant performance improvements came from relatively straightforward optimizations, with the top three techniques (data format optimization, batch prefetching, and feature preprocessing) delivering 71% of the total performance gain while representing only 31% of implementation effort [11].

## 6.3 Organizational Impact and Return on Investment

Beyond technical metrics, pipeline optimization delivers substantial organizational benefits through accelerated development cycles and improved resource economics. The Columbia University research quantifies that data scientists conducted 2.8x more experimental iterations following pipeline optimization, directly accelerating model development [12]. This increased experimental velocity translated to measurable quality improvements, with A/B testing showing that recommendation models improved key engagement metrics by 0.8-1.2% following more comprehensive hyperparameter optimization enabled by faster training cycles. The economic impact analysis from ResearchGate research demonstrates that organizations implementing comprehensive pipeline optimizations reduced their total cost of ownership for machine learning infrastructure by 47% on average [11]. This cost reduction stems from 58% lower GPU provisioning requirements, 43% reduced data preprocessing overhead, and 31% decreased storage costs through improved compression and format selection. For large-scale social media platforms with annual infrastructure budgets in the tens of millions, these efficiency improvements translate to substantial financial impact, with documented case studies showing 7-11 month return on investment periods for optimization initiatives [11].

## Conclusion

Optimizing data ingestion represents a strategic imperative for organizations deploying machine learning at scale across social media platforms. This article has demonstrated how targeted improvements to data loading architectures, storage formats, and feature engineering processes can collectively transform the efficiency of machine learning pipelines. By addressing the fundamental challenges of I/O bottlenecks, network latency, and storage inefficiencies, organizations can unlock significant performance gains throughout their ML infrastructure. The case studies presented illustrate that these optimizations deliver not only technical benefits in terms of reduced training times and improved resource utilization but also tangible business value through faster model iteration, reduced infrastructure costs, and enhanced model

quality. As social media platforms continue to generate increasingly complex and voluminous data, these data ingestion optimization strategies will become even more critical to maintaining competitive advantage in the rapidly evolving landscape of machine learning applications.

## References

1. Aarti Basant, "Scaling data ingestion for machine learning training at Meta," Engineering at Meta, Sep. 19, 2022. [Online]. Available: https://engineering.fb.com/2022/09/19/ml-applications/data-ingestion-machine-learning-training-meta/

2. Jerome Kafrouni, "Speed Up Your Data Pipelines for Deep Learning," IBM Data and AI, Medium, 23 Nov. 2021. [Online]. Available: https://medium.com/ibm-data-ai/speed-up-your-data-pipelines-for-deep-learning-93fcd6bce2c

3. Srivignesh Rajan, "Data Processing Pipeline in Machine Learning," The Startup, Medium, 10 July 2020. [Online]. Available: https://medium.com/swlh/data-preprocessing-and-data-modeling-for-kaggle-house-price-prediction-data-in-python-c04055ded258

4. Shailendra Chauhan, "Data Ingestion in Machine Learning," ScholarHat Tutorials, 17 July 2023. [Online]. Available: https://www.scholarhat.com/tutorial/machinelearning/data-ingestion-in-machine-learning

5. Armin Esmaeilzadeh et al., "Efficient Large Scale NLP Feature Engineering with Apache Spark," 2022 IEEE Xplore, 4 March 2022. [Online]. Available: https://ieeexplore.ieee.org/document/9720765

6. NVIDIA Docs Hub, "Best Practices for Building and Deploying Recommender Systems," NVIDIA Developer Documentation, Feb. 2023. [Online]. Available: https://docs.nvidia.com/deeplearning/performance/recsys-best-practices/index.html

7. Dheevatsa Mudigere et al., "High-performance Distributed Training of Large-scale Deep Learning Recommendation Models," ResearchGate, Apr. 2021. [Online]. Available: https://www.researchgate.net/publication/350834637_High-performance_Distributed_Training_of_Large-scale_Deep_Learning_Recommendation_Models

8. Stefano Martire, "Optimizing a TensorFlow Input Pipeline: Best Practices in 2022," Medium, 10 Feb. 2022. [Online]. Available: https://medium.com/@virtualmartire/optimizing-a-tensorflow-input-pipeline-best-practices-in-2022-4ade92ef8736

9. Bethel Nwokocha, "Optimizing Data Pipelines for Machine Learning: Best Practices, Architecture Designs, and Emerging Trends," Medium, 29 Aug. 2024. [Online]. Available: https://betheldaniel332.medium.com/optimizing-data-pipelines-for-machine-learning-best-practices-architecture-designs-and-emerging-b509fef24396

10. Christopher Olston et al., "TensorFlow-Serving: Flexible, High-Performance ML Serving," arXiv:1712.06139v2, 27 Dec. 2017. [Online]. Available: https://arxiv.org/pdf/1712.06139

11. Brahma Reddy Katam, "Optimizing Data Pipeline Efficiency with Machine Learning Techniques," ResearchGate, Vol. 8, no. 7, Jan. 2024. [Online]. Available: https://www.researchgate.net/publication/382642570_Optimizing_Data_Pipeline_Efficiency_with_Machine_Learning_Techniques

12. Li Erran Li et al., "Scaling Machine Learning as a Service," JMLR: Workshop and Conference Proceedings, 2016. [Online]. Available: https://www.cs.columbia.edu/~lierranli/publications/jmlr_mls2017.pdf