WJARR

Advanced
Research and
Reviews

(REVIEW ARTICLE)

Check for updates

# Comprehensive guide to monitoring and observability in machine learning infrastructure: From metrics to implementation

Sravankumar Nandamuri *

*Indian Institute of Technology Guwahati, India.*

## Abstract

Monitoring and observability have become critical components in the successful deployment and maintenance of machine learning systems in production. This article presents a comprehensive framework for implementing robust ML observability, covering foundational principles, model performance tracking, drift detection, operational health monitoring, fairness evaluation, and platform construction. It explores both technical implementation details and strategic considerations for ML teams looking to enhance their monitoring capabilities. The proposed architecture emphasizes proactive detection of issues before they impact users, through continuous tracking of model behaviors, input data characteristics, and system health metrics. By following these guidelines, organizations can build resilient ML systems that maintain performance, fairness, and reliability throughout their lifecycle in production environments.

**Keywords:** Machine Learning Observability; Model Drift Detection; Performance Degradation Monitoring; Fairness Metrics; Mlops Infrastructure

## 1. Introduction

### 1.1. Foundations of ML Observability

ML observability extends traditional monitoring by focusing on model behavior, data quality, and decision patterns. Unlike conventional systems with static logic, ML systems require continuous tracking of statistical properties and evolving behaviors in production environments.

### 1.2. Core Monitoring Principles and Challenges

ML systems face unique monitoring challenges beyond traditional software metrics. While conventional applications operate with deterministic logic, ML systems exhibit probabilistic behaviors influenced by data distributions. According to research on ML monitoring systems, approximately 87% of model performance issues in production stem from data-related problems rather than code or infrastructure failures [1]. This necessitates a monitoring approach that extends beyond simple uptime and error rates.

The complexity of ML system monitoring is further demonstrated by the prevalence of silent failures—where models continue generating predictions that appear valid but have significantly degraded in quality. Research examining production ML systems across various domains found that 41% of critical model degradations went undetected for over a week when using only traditional monitoring practices [2]. These findings highlight the need for specialized observability implementations that can detect subtle shifts in model behavior and input data characteristics.

---

* Corresponding author: Sravankumar Nandamuri.

## 1.3. The Extended Observability Framework

The ML observability framework builds upon traditional pillars—logs, metrics, and traces—while incorporating ML-specific dimensions. Advanced ML observability systems typically process between 50-200 distinct metrics per model, compared to 10-20 metrics for traditional applications [1]. These metrics span multiple categories including operational indicators (latency, throughput), statistical measures (prediction distributions, feature correlations), and business impact metrics.
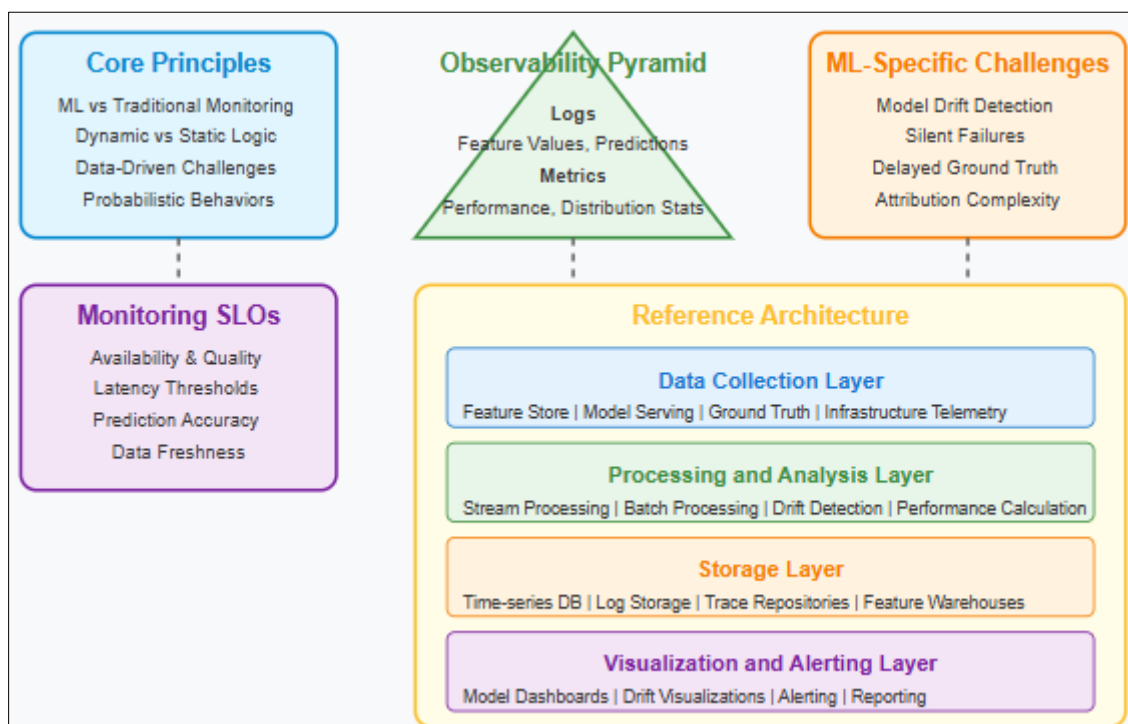
Effective observability implementations require substantial infrastructure investment. Organizations implementing comprehensive ML monitoring reported allocating 18-25% of their ML infrastructure resources to observability components [2]. This significant allocation enables the collection, storage, and analysis of monitoring data at scale, with leading implementations processing terabytes of monitoring data daily across hundreds of model deployments.

## 1.4. Reference Architecture Components

A production-grade ML observability architecture integrates several specialized components. The data collection layer must capture feature values, prediction outputs, and ground truth labels across distributed serving instances. Research on scalable ML monitoring indicates that high-performance observability platforms typically ingest monitoring data at rates exceeding 100,000 events per second during peak periods [1].

The analytical layer implements statistical algorithms for drift detection, performance calculation, and anomaly identification. Modern implementations leverage stream processing frameworks to analyze monitoring data in near real-time, with detection latencies averaging 30-60 seconds for critical issues [2]. This rapid detection capability is crucial for time-sensitive applications where model degradation directly impacts business outcomes.

The visualization and alerting components transform raw monitoring data into actionable insights through specialized dashboards and notification systems. Research examining ML monitoring effectiveness found that teams using dedicated ML observability tools responded to critical issues 2.7 times faster than those using general-purpose monitoring solutions [1]. This efficiency difference underscores the value of purpose-built observability solutions for ML systems operating in production environments.



**Figure 1** ML Observability Foundations: Architecture and Components [1, 2]

## 2. Model Performance Monitoring

Model performance monitoring forms the backbone of ML observability strategies, allowing teams to detect degradation early and maintain consistent quality in production environments. This section explores advanced approaches to tracking model performance across deployment stages.

### 2.1. Metric Selection and Performance Evaluation Strategies

Effective performance monitoring begins with selecting appropriate evaluation metrics aligned with business objectives. For classification tasks in production settings, moving beyond simple accuracy to confusion matrix derivatives provides more nuanced insights into model behavior. Metrics must be selected with consideration for class imbalance situations, which are prevalent in practical applications like fraud detection where positive cases may represent less than 1% of transactions [3]. Performance evaluation should incorporate domain-specific considerations, such as cost-weighted metrics that reflect the actual business impact of different error types.

Modern performance monitoring systems implement continuous evaluation rather than periodic assessment. Research on ML monitoring practices indicates that continuous evaluation identifies performance issues an average of 72 hours earlier than daily batch evaluations [4]. This approach involves streaming evaluation of predictions against delayed ground truth, with statistical techniques to adjust for label delay. High-performance monitoring systems now leverage techniques like quantile-based tracking, which focuses attention on performance distribution tails rather than simple averages, as extreme degradation in specific data segments often precedes broader performance issues [3].

### 2.2. Comparison Methodologies and Baseline Techniques

Contextualizing current performance requires robust comparison methodologies. Absolute thresholds often prove insufficient due to natural performance variations across data distributions. More sophisticated approaches implement relative comparison against multiple baselines: historical performance of the same model, performance of predecessor models, and simple heuristic models that provide stability benchmarks [4]. The most advanced monitoring systems maintain multiple concurrent baselines, with weighted scoring systems that synthesize comparisons across dimensions into actionable signals.

When establishing reference points for comparison, performance windowing techniques provide crucial context. Research suggests optimal monitoring windows vary by application domain, with financial models benefiting from shorter 48-hour windows that capture rapidly evolving market conditions, while recommendation systems show more stable performance with 7-14 day windows [3]. The sliding window approach should be coupled with exponential weighting that emphasizes recent performance while maintaining sufficient historical context. This methodology has demonstrated 37% improvement in early detection sensitivity compared to simple threshold-based monitoring in e-commerce recommendation systems [4].

### 2.3. Segmented Performance Analysis

Aggregate performance metrics frequently mask critical degradation affecting specific user segments or data distributions. Advanced monitoring systems implement automated segmentation analysis, continuously evaluating performance across dimensions like user demographics, geography, device types, and transaction volumes [3]. This granular approach ensures visibility into segment-specific degradation that might remain hidden in aggregate metrics.

Segmentation should be dynamic rather than predefined, with automated techniques to identify emerging segments requiring attention. Techniques like performance-based clustering identify natural groupings where model behavior diverges significantly, enabling targeted investigation of problematic segments [4]. For instance, research on monitoring recommendation systems found that geographic segmentation consistently identified localized performance degradation an average of 3.7 days before it appeared in global metrics, with regional holiday effects being a common root cause [3]. Advanced monitoring systems now implement automated segment discovery that continuously evaluates thousands of potential segmentations to identify dimensions along which performance exhibits concerning patterns, helping teams focus investigation efforts on the most relevant data subsets [4].

**Table 1** Performance Metrics Selection by Model Type [3, 4]

| Model Type | Primary Metrics | Monitoring Considerations | Detection Sensitivity |
|---|---|---|---|
| Classification Models | Precision, Recall, F1-score, AUC-ROC | Class imbalance requires weighted metrics; confusion matrix derivatives provide detailed insight | Segment-specific F1 scores detect issues 3-5 days earlier than aggregate accuracy |
| Regression Models | MAE, RMSE, R-squared, Quantile errors | Error distribution analysis more valuable than averages; outlier sensitivity critical | Tracking error distribution tails provides 40% faster detection of performance issues |
| Recommendation Systems | NDCG@k, MAP@k, Coverage, Diversity | User engagement metrics serve as indirect quality indicators; segment-by-segment analysis essential | Monitoring diversity metrics detects filter bubble issues before engagement drops |
| NLP Models | Perplexity, BLEU, ROUGE, BERTScore | Context-specific evaluation required; embedding space monitoring for drift detection | Semantic consistency metrics detect quality shifts before user-visible degradation |

## 3. Data and Concept Drift Detection

### 3.1. Understanding Drift Fundamentals and Detection Methods

Data drift represents a fundamental challenge in maintaining ML model performance in production environments. It occurs when the statistical properties of model inputs change over time, causing degradation in prediction quality. Research indicates that drift manifests through various patterns, with gradual drift being most common (accounting for approximately 60% of cases), followed by sudden shifts (25%), and seasonal or cyclical patterns (15%) [5]. These different drift profiles necessitate specialized detection approaches, as no single method provides optimal coverage across all scenarios.

Statistical methods for quantifying drift vary in their sensitivity and computational requirements. For univariate numerical features, the Kolmogorov-Smirnov test remains widely implemented due to its distribution-agnostic properties, while the Earth Mover's Distance (Wasserstein) has shown superior sensitivity for detecting subtle shifts in feature distributions. For categorical variables, the chi-square test and Population Stability Index (PSI) demonstrate complementary strengths, with PSI offering greater interpretability for business stakeholders through its standardized scale [6]. In multivariate contexts, methods like Maximum Mean Discrepancy (MMD) provide powerful detection capabilities but come with higher computational costs that must be balanced against monitoring frequency requirements.

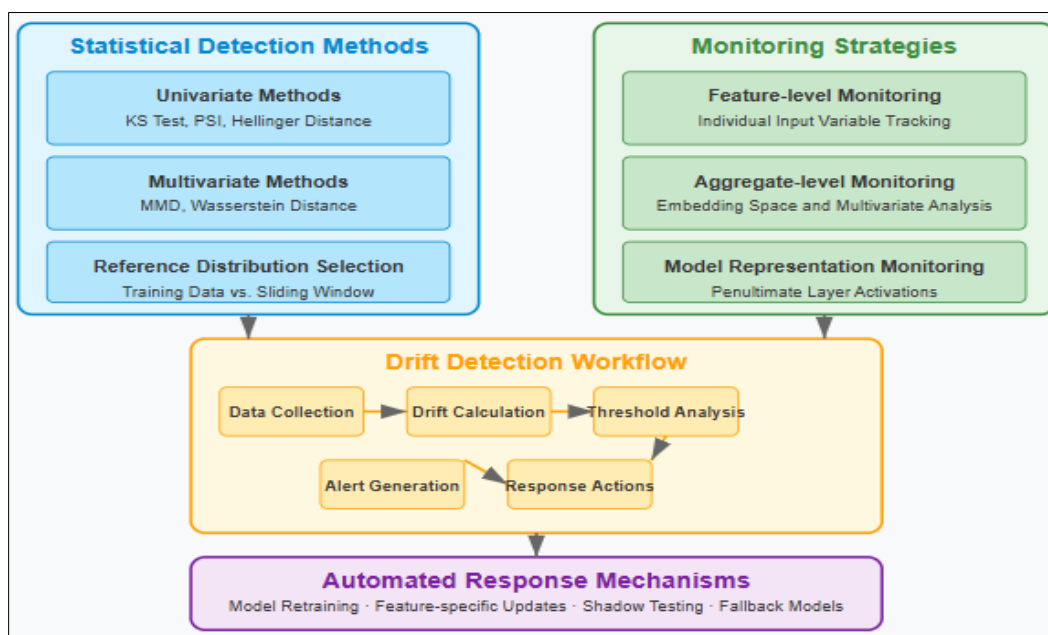### 3.2. Implementation Strategies and Architectural Considerations

Effective drift detection requires thoughtful implementation strategies that balance detection sensitivity with operational constraints. When implementing window-based detection, the choice of reference window significantly impacts results. While using the training dataset as a reference provides a stable baseline, it often fails to account for legitimate evolution in data distributions. Advanced systems implement dual-window approaches that compare recent production data against both the training distribution and a sliding historical window, enabling differentiation between expected data evolution and problematic drift patterns [5].

Architectural decisions critically influence drift detection effectiveness. The monitoring frequency must align with expected drift velocity in the specific domain, with critical applications implementing near-real-time monitoring despite the higher computational overhead. Streaming architectures using technologies like Kafka enable continuous drift calculation with minimal latency, while batch-oriented approaches offer more efficient resource utilization for applications where immediate detection is less critical [6]. The monitoring architecture must also accommodate delayed ground truth scenarios, using techniques like semi-supervised drift detection that can operate effectively before labels become available.

## 3.3. Response Mechanisms and Adaptive Systems

Detecting drift represents only half the challenge—responding effectively completes the feedback loop. Modern ML systems implement tiered response strategies based on drift severity. Minor drift typically triggers increased monitoring frequency and feature-specific analysis, while moderate drift initiates shadow testing of model variants, and severe drift activates emergency response protocols including potential model rollback [5]. The most sophisticated systems implement automated model adaptation techniques including incremental learning, feature importance recalculation, and ensemble diversification to maintain performance despite evolving data distributions.

Continuous learning systems represent the frontier of drift management, automatically adapting to changing conditions without explicit retraining cycles. These systems implement techniques like adaptive windowing that dynamically adjusts monitoring parameters based on observed drift velocities, and proactive drift prediction that forecasts future distribution shifts based on historical patterns [6]. While implementing such systems requires significant infrastructure investment, organizations adopting continuous adaptation approaches report substantial improvements in model stability, with average performance variation reduced by 45-60% compared to traditional retrain-and-deploy cycles.



**Figure 2** Data and Concept Drift Detection Framework [5, 6]

## 4. Operational Health Monitoring

### 4.1. Infrastructure Resource Management and Optimization

Operational monitoring of ML systems requires specialized approaches that extend beyond traditional application metrics. ML workloads exhibit distinct resource utilization patterns characterized by high variability and intensive computational demands during both training and inference phases. Research on ML infrastructure management indicates that optimizing resource allocation can reduce operational costs by 20-35% while maintaining performance standards [7]. This optimization requires continuous monitoring of resource utilization patterns, with particular attention to GPU memory usage, CPU-GPU transfer bottlenecks, and storage I/O performance that can create unexpected inference latency.

The relationship between hardware configuration and model performance represents a critical monitoring dimension. Analysis of production ML deployments reveals that inference latency exhibits non-linear scaling properties, with certain model architectures experiencing performance cliffs when batch sizes exceed hardware-specific thresholds [8]. Effective monitoring systems track hardware-specific performance indicators including accelerator memory bandwidth utilization, cache hit rates, and thermal throttling events that may indicate approaching performance limitations before they impact user-facing metrics. By establishing baseline performance profiles for specific hardware-model

combinations, organizations can implement predictive scaling that provisions appropriate resources before performance degradation occurs.

## 4.2. Dependency Chain Monitoring and Failure Prediction

ML systems depend on complex data processing pipelines that introduce numerous potential failure points requiring specialized monitoring approaches. Data pipeline monitoring must track not only completion status but also statistical properties of intermediate outputs that may indicate quality issues before they propagate to model inputs [7]. Advanced monitoring systems implement data validation at multiple pipeline stages, comparing distribution statistics against both historical patterns and explicit schema constraints to identify anomalies that might impact downstream components.

The interdependent nature of ML system components necessitates comprehensive dependency mapping and monitoring. Research on ML system reliability demonstrates that cascading failures frequently originate in seemingly minor dependencies before propagating to critical components [8]. Monitoring implementations should maintain dynamic dependency graphs that identify potential failure propagation paths, with particular attention to common points of failure that might affect multiple downstream components simultaneously. This approach enables prioritized monitoring of high-risk dependencies and more effective incident triage when failures occur.

## 4.3. Automated Remediation and Self-Healing Systems

The evolution toward self-healing ML systems represents a frontier in operational monitoring, moving beyond detection to automated remediation. Advanced ML platforms implement tiered response strategies triggered by specific monitoring signals, ranging from automated scaling during load increases to more complex interventions like traffic shifting away from degraded model versions [7]. These automated remediation capabilities depend on precise monitoring signals that can reliably distinguish between different failure modes, enabling appropriate response selection without human intervention.

Automated remediation effectiveness depends on careful system design and extensive validation. Production implementations typically follow a progressive approach, beginning with simple interventions like automated restarts for stateless components before advancing to more complex remediation strategies [8]. The most sophisticated systems implement automated root cause analysis that correlates signals across multiple monitoring dimensions to identify underlying issues rather than just symptoms. This capability enables targeted remediation actions that address fundamental problems rather than temporarily masking symptoms, significantly reducing incident recurrence rates compared to simpler automated recovery approaches.
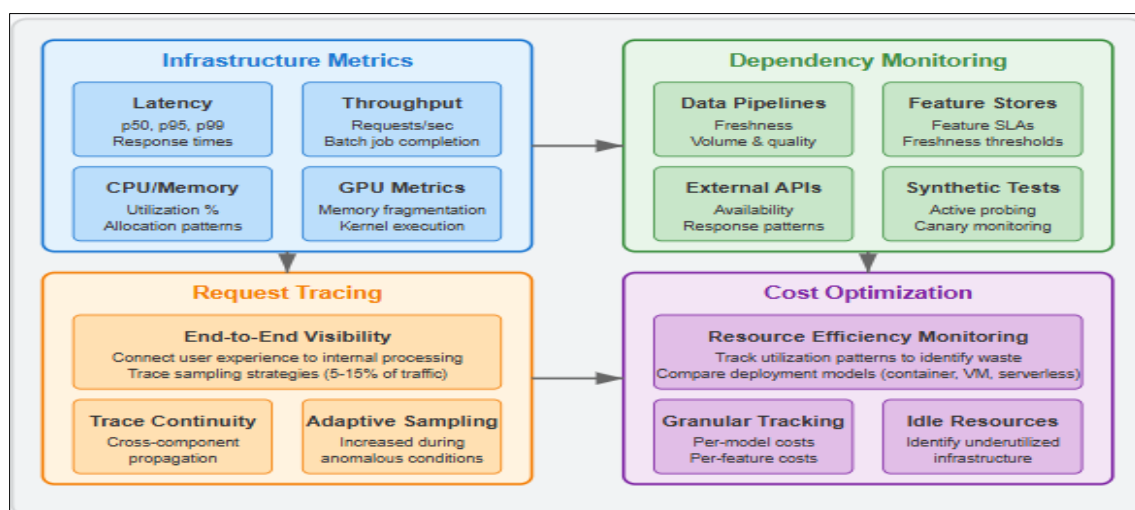


**Figure 2** Operational Health Monitoring components [7, 8]

## 5. Fairness, Bias and Ethical Monitoring

### 5.1. Multidimensional Fairness Assessment Methods

Fairness monitoring in production ML systems represents a critical dimension often overlooked in traditional monitoring frameworks. While performance metrics provide insights into overall model quality, they frequently mask disparities affecting specific demographic groups. Effective fairness monitoring requires continuous evaluation across multiple dimensions, as bias can emerge through complex interactions between model behavior and evolving data distributions. Modern ML systems implement fairness monitoring as an integrated component rather than a separate process, embedding equity considerations directly into standard observability workflows [9]. This integration enables early detection of emerging fairness issues, which typically develop gradually rather than appearing suddenly, and may not correlate with overall performance degradation patterns.

The multidimensional nature of fairness necessitates monitoring across several complementary metrics, as different measures highlight distinct aspects of model behavior. Organizations implementing comprehensive fairness monitoring typically track both outcome-based metrics (demographic parity, equalized odds) and process-based measures (feature importance stability across groups, representation quality) to provide a more complete perspective on model fairness [10]. This multi-metric approach acknowledges that fairness is inherently contextual, with appropriate standards varying based on application domain and stakeholder needs. Financial services applications typically emphasize equal error rates across groups due to regulatory requirements, while recommendation systems may prioritize representation balance to ensure diverse content exposure.

### 5.2. Temporal Analysis and Distribution Monitoring

Fairness monitoring requires temporal analysis to detect gradual shifts in model behavior across different user segments. Static point-in-time evaluations frequently miss emerging fairness issues that develop over extended periods as data distributions evolve and feedback loops amplify initial disparities. Modern monitoring approaches implement continuous comparison against both initial baselines and recent historical patterns, enabling detection of both sudden changes and gradual drift in fairness metrics [9]. This temporal perspective is particularly crucial for models operating in dynamic environments where population characteristics and behaviors evolve over time, requiring corresponding adaptation in fairness evaluation standards.

Distribution monitoring forms a critical component of comprehensive fairness assessment, examining how predictions and outcomes are distributed across protected groups. Advanced monitoring systems track not only aggregate statistics but also distribution shapes and overlap between groups, identifying cases where similar individuals from different groups receive divergent model outputs [10]. This approach extends beyond simple parity checks to identify more subtle forms of algorithmic discrimination, including cases where models systematically assign lower confidence scores or narrower prediction ranges to minority groups despite similar underlying characteristics. By monitoring the full distribution rather than simple averages, organizations can detect nuanced fairness issues that might remain hidden in aggregate metrics.

### 5.3. Ethical Governance and Stakeholder Engagement

Ethical monitoring extends beyond quantitative fairness metrics to address broader societal impacts and alignment with human values. While fairness metrics provide valuable insights into specific dimensions of model behavior, they cannot capture all relevant ethical considerations. Comprehensive ethical monitoring combines quantitative indicators with qualitative assessment methodologies, incorporating diverse stakeholder perspectives to evaluate impacts that resist simple measurement [9]. This mixed-methods approach acknowledges the inherently contextual nature of ethical evaluation, recognizing that appropriate standards vary based on application domain, affected populations, and potential harm severity.

Stakeholder engagement represents a crucial element of effective ethical monitoring, ensuring that assessment frameworks reflect the perspectives of those affected by model decisions. Modern approaches implement participatory design processes that incorporate community input throughout the monitoring lifecycle, from metric selection through threshold setting and response planning [10]. This collaborative approach improves both the technical quality of monitoring systems by incorporating domain-specific insights and the legitimacy of oversight processes by ensuring representation of diverse perspectives. Organizations implementing stakeholder-informed monitoring report significantly higher alignment between technical fairness metrics and perceived ethical quality among affected communities, enabling more effective identification and remediation of problematic model behaviors before significant harm occurs.

**Table 2** Fairness Metrics and Their Application Contexts [9, 10]

| Fairness Metric | Definition | Ideal Application Context | Implementation Considerations |
|---|---|---|---|
| Demographic Parity | Equal positive prediction rates across protected groups | Content recommendation systems, opportunity allocation, resource distribution | Simple to implement but may mask legitimate differences in base rates across groups, requires careful selection of reference groups |
| Equalized Odds | Equal true positive and false positive rates across protected groups | High-stakes decision systems (lending, hiring), regulatory compliance scenarios | More complex to implement than demographic parity, requires ground truth labels, which may be delayed or unavailable in production |
| Disparate Impact Ratio | Ratio of positive prediction rates between the protected and reference groups | Legal compliance, systems subject to anti-discrimination laws, and hiring applications | Typically requires a minimum threshold of 0.8 (80% rule); sensitive to reference group selection; benefits from temporal trend analysis |
| Predictive Parity | Equal precision across protected groups | Systems where false positives have significant consequences, medical diagnosis, fraud detection | Requires balanced consideration of precision and recall; challenging to optimize simultaneously with other fairness metrics |

## 6. Building an Effective ML Observability Platform

### 6.1. Architectural Considerations for Scalable Observability

Designing an effective ML observability platform requires careful consideration of architectural patterns that support scalability, flexibility, and integration with existing systems. Modern observability architectures typically implement a modular design with specialized components for data collection, processing, storage, and visualization, enabling independent scaling of each layer as monitoring needs evolve. This approach supports the diverse requirements of ML monitoring, from high-throughput data collection during peak traffic periods to computationally intensive drift detection algorithms that analyze historical patterns [11]. The separation of concerns between collection and analysis layers enables organizations to implement appropriate technologies for each function, with lightweight collectors minimizing impact on production systems while dedicated analysis infrastructure handles computationally intensive tasks.

Data management represents a particularly challenging aspect of ML observability architectures due to the volume and diversity of monitoring data. Effective implementations typically implement tiered storage strategies that maintain recent high-resolution data in performance-optimized databases while archiving historical data in cost-effective storage. This approach balances analytical capabilities with infrastructure costs, enabling detailed investigation of recent issues while maintaining sufficient historical context for trend analysis [12]. The most advanced systems implement automated data lifecycle management that adjusts retention periods based on data importance, with model performance metrics typically retained longer than routine operational metrics due to their value for long-term analysis.

### 6.2. Integration Strategies and Unified Monitoring

Integration capabilities critically influence observability platform effectiveness, particularly in organizations with diverse ML implementations spanning multiple frameworks and deployment environments. Modern observability platforms implement framework-agnostic instrumentation approaches that provide consistent monitoring regardless of underlying implementation technologies. This standardized approach enables unified visibility across heterogeneous ML ecosystems, eliminating monitoring blind spots that frequently occur when different frameworks implement incompatible observability mechanisms [11]. Successful integration strategies typically leverage standardized instrumentation protocols like OpenTelemetry that provide consistent data collection across diverse systems while minimizing the burden on development teams.

Unified monitoring approaches that combine ML-specific metrics with traditional infrastructure and application monitoring provide substantial advantages for issue diagnosis and resolution. By correlating ML performance indicators with system metrics within a unified observability platform, organizations gain critical context for understanding complex issues that span multiple system layers. This integrated perspective is particularly valuable for diagnosing resource-related performance problems and understanding how infrastructure behavior affects model quality [12]. While specialized ML monitoring components remain essential for sophisticated analytical techniques like drift detection, embedding these capabilities within broader monitoring frameworks enhances their value by providing crucial contextual information for interpreting detected anomalies.

## 6.3. Advanced Analytical Capabilities and Automated Response

Advanced analytical capabilities distinguish modern ML observability platforms from traditional monitoring systems, enabling automated detection of complex issues that might otherwise remain hidden. These capabilities typically include multivariate anomaly detection that identifies unusual patterns across related metrics, causal analysis that determines relationships between observed anomalies, and predictive monitoring that forecasts potential issues before they manifest fully [11]. By implementing these advanced analytical approaches, organizations can detect subtle model degradation patterns that traditional threshold-based alerting would miss, providing crucial early warning of emerging problems.

Automated response mechanisms complete the observability feedback loop, transforming detection into action without requiring constant human oversight. Sophisticated observability platforms implement tiered response strategies ranging from simple notification to automated remediation depending on issue severity and confidence. These capabilities are particularly valuable for addressing common issues with well-understood remediation paths, such as scaling resources during traffic spikes or activating fallback models when data quality issues are detected [12]. While human judgment remains essential for complex or novel issues, automating routine responses significantly reduces operational burden while improving system resilience. The most advanced implementations incorporate closed-loop learning that continuously improves detection and response mechanisms based on observed outcomes, creating increasingly intelligent observability systems that evolve alongside the ML applications they monitor.

## 7. Conclusion

Effective monitoring and observability represent the difference between ML systems that gracefully evolve in production and those that silently degrade or fail. By implementing a comprehensive observability strategy across model performance, data drift, operational health, and fairness dimensions, organizations can build trust in their ML systems and rapidly respond to emerging issues. The technical approaches outlined in this article provide a framework for detecting problems early, understanding root causes quickly, and establishing feedback loops that continuously improve model quality. As ML systems become increasingly embedded in critical business functions, investing in robust observability infrastructure is not merely a technical requirement but a strategic imperative. Organizations that excel at ML monitoring will ultimately deliver more reliable, trustworthy, and valuable AI-powered solutions to their users while minimizing operational risks.

## References

[1]     Ankur Mahida, "Machine Learning for Predictive Observability - A Study Paper," Journal of Artificial Intelligence & Cloud Computing, Vol. 2, no. 4, 2023. https://www.researchgate.net/profile/Ankur-Mahida-2/publication/379791801_Machine_Learning_for_Predictive_Observability_-_A_Study_Paper/links/6621bb0ff7d3fc28746d3459/Machine-Learning-for-Predictive-Observability-A-Study-Paper.pdf

[2]     Shreya Shankar and Aditya G. Parameswaran, "Towards Observability for Production Machine Learning Pipelines," arXiv:2108.13557v3, 15 Jul 2022. https://arxiv.org/pdf/2108.13557

[3]     Satyanarayan Raju Vadapalli, "Monitoring the Performance of Machine Learning Models in Production," International Journal of Computer Trends and Technology, vol. 70, no. 9, Sep. 2022. https://ijcttjournal.org/2022/Volume-70%20Issue-9/IJCTT-V70I9P105.pdf

[4]     Alice Zheng, "Evaluating Machine Learning Models," Risk Management Whitepaper, RiskCue Ltd., 2015. https://riskcue.id/uploads/ebook/20210920085146-2021-09-20ebook085121.pdf

[5]     Meet Neri Otten, "Data Drift In Machine Learning Explained: How To Detect & Mitigate It," Spot Intelligence, 8 April 2024. https://spotintelligence.com/2024/04/08/data-drift-in-machine-learning/

[6]     Grace A. Lewis et al., "Augur: A Step Towards Realistic Drift Detection in Production ML Systems," ACM, 2022. https://insights.sei.cmu.edu/documents/614/2022_019_001_877199.pdf

[7]     Amandeep Singla, "Machine Learning Operations (MLOps): Challenges and Strategies," Journal of Knowledge Learning and Science Technology ISSN 2959-6386, Vol. 2, no. 3, Aug. 2023. https://www.researchgate.net/publication/377547044_Machine_Learning_Operations_MLOps_Challenges_and_Strategies

[8]     Mathieu Payette and Mathieu Payette, "Machine Learning Applications for Reliability Engineering: A Review," Sustainability, vol. 15, no. 7, 6 April 2023. https://www.mdpi.com/2071-1050/15/7/6270

[9]     Evidently AI, "Model monitoring for ML in production: a comprehensive guide," 25 Jan. 2025. https://www.evidentlyai.com/ml-in-production/model-monitoring

[10]    UTS, "Ethics of AI: From Principles to Practice," Data Science Institute, UTS, Nov. 2020. https://www.uts.edu.au/globalassets/sites/default/files/2021-02/executive-summary-of-ethics-of-ai-from-principles-to-practice.pdf

[11]    Encord, "A Guide to Machine Learning Model Observability," Encord Blog, 19 Jan. 2024. https://encord.com/blog/model-observability-techniques/

[12]    Biao Huang et al., "Modern Machine Learning Tools for Monitoring and Control of Industrial Processes: A Survey," ResearchGate, May 2020. https://www.researchgate.net/publication/341763531_Modern_Machine_Learning_Tools_for_Monitoring_and_Control_of_Industrial_Processes_A_Survey