

# Phase 2: Exploratory Data Analysis and Implementation

## Abstract

This phase focuses on the systematic exploration and analysis of our expanded system log datasets (HDFS, Apache, HealthApp, BGL, HPC, Linux, and Mac) to detect and characterize reliability drift patterns in streaming AI systems. Through comprehensive exploratory data analysis (EDA) and practical implementation, we develop a robust framework for identifying system behavior changes and model drift across different types of logs in real-time streaming environments.

## 1. Problem Statement

### 1.1 Research Context

Model and system reliability drift manifests differently across our log types:

- **HDFS:** Block operation patterns and DataNode behavior changes in distributed storage
- **Apache:** Server initialization, module configuration shifts, and request pattern drift
- **HealthApp:** User activity pattern changes and sensor data anomalies in real-time
- **BGL:** Supercomputer system behavior and node failures in high-performance environments
- **HPC:** High-performance computing job patterns and resource utilization shifts
- **Linux:** System service behavior and resource management pattern changes
- **Mac:** Application behavior and system event pattern drift

### 1.2 Problem Definition

Given our diverse log streams in a real-time AI system environment, we aim to:

1. Identify dataset-specific drift patterns and their impact on model performance
2. Develop unified preprocessing and analysis pipelines for streaming data
3. Implement real-time drift detection with minimal latency
4. Create a scalable monitoring framework for heterogeneous log sources
5. Distinguish between normal fluctuations and significant drift patterns

## 1.3 Research Questions

1. How do drift patterns differ across our expanded log types in streaming environments?
2. What features best indicate reliability drift and model performance degradation?
3. Can we develop a unified approach for heterogeneous log sources in real-time?
4. What are the optimal time windows for each log type in streaming analysis?
5. How do system-specific characteristics affect drift detection accuracy?
6. How can we differentiate between noise and meaningful drift patterns?

## 2. Implementation Approach

### 2.1 Log Parsing and Streaming Implementation

```
class StreamingLogParser:
    def __init__(self, log_type, kafka_config=None):
        self.log_type = log_type
        self.kafka_config = kafka_config or {}
        self.patterns = {
            'hdfs': {
                'timestamp': r'(\d{6}\s\d{6})',
                'thread_id': r'\s(\d+)\s',
                'level': r'\s(INFO|ERROR|WARN)\s',
                'component': r'(dfs\.[A-Za-z$]+)',
                'block_id': r'blk_[-\d]+'
            },
            'apache': {
                'timestamp': r'\[([.*?])\]',
                'level': r'\[(notice|error|warn)\]',
                'module': r'\s([A-Za-z$]+):',
                'message': r':\s(.+)$'
            },
            'health': {
                'timestamp': r'(\d{8}-\d{2}:\d{2}:\d{2}:\d{3})',
                'component': r'\|([.*?])\|',
                'user_id': r'\|(\d+)\|',
                'action': r'\|([^\|]+)$'
            },
            'bgl': {
                'timestamp': r'(\d{6}\s\d{6})',
                'location': r'\s(\S+)\s',
                'severity': r'\s(\w+)\s',
                'component': r'\s(\w+)\s',
                'message': r'\s(.+)$'
            },
            'hpc': {
                'timestamp': r'(\d{6}\s\d{6})',
                'node_id': r'\s(\w+)\s',
                'job_id': r'\s(\w+)\s',
                'message': r'\s(.+)$'
            },
            'linux': {
```

```

        'timestamp': r'(\w+\s+\d+\s+\d+:\d+:\d+)',
        'hostname': r'\s(\S+)\s',
        'process': r'(\w+)\[',
        'pid': r'\[(\d+)\]'
    },
    'mac': {
        'timestamp': r'(\d{4}-\d{2}-\d{2}\s+\d{2}:\d{2}:\d{2})',
        'sender': r'\s(\S+)\s',
        'type': r'\s(\w+)\s',
        'category': r'\s(\w+):'
    }
}

```

```

def process_stream(self, message):
    """Process incoming stream messages"""
    extracted = self.parse_line(message)
    return self._enrich_with_metadata(extracted)

```

```

def parse_line(self, line):
    patterns = self.patterns[self.log_type]
    extracted = {}
    for field, pattern in patterns.items():
        match = re.search(pattern, line)
        if match:
            extracted[field] = match.group(1)
    return extracted

```

```

def _enrich_with_metadata(self, parsed_data):
    """Add streaming metadata for drift analysis"""
    parsed_data['stream_timestamp'] = time.time()
    parsed_data['source_type'] = self.log_type
    return parsed_data

```

## 2.2 Feature Engineering Pipeline

```
class StreamingFeatureExtractor:
    def __init__(self, window_size=300): # 5-minute windows
        self.window_size = window_size
        self.reference_distributions = {}

    def extract_temporal_features(self, log_df):
        features = {
            'event_rate': self._calculate_event_rate(log_df),
            'level_ratios': self._calculate_level_ratios(log_df),
            'component_distribution': self._get_component_dist(log_df),
            'time_patterns': self._extract_time_patterns(log_df),
            'drift_indicators': self._compute_drift_metrics(log_df)
        }
        return features

    def extract_dataset_specific_features(self, log_df, log_type):
        extractors = {
            'hdfs': self._extract_hdfs_features,
            'apache': self._extract_apache_features,
            'health': self._extract_health_features,
            'bgl': self._extract_bgl_features,
            'hpc': self._extract_hpc_features,
            'linux': self._extract_linux_features,
            'mac': self._extract_mac_features
        }
        return extractors[log_type](log_df)

    def _compute_drift_metrics(self, current_window):
        """Compute statistical distance metrics for drift detection"""
        metrics = {
            'jensen_shannon_div': self._compute_js_divergence(current_window),
            'mean_shift': self._compute_mean_shift(current_window),
            'distribution_stats': self._compute_distribution_stats(current_window)
        }
        return metrics
```

# 3. Exploratory Analysis Results

## 3.1 Log Characteristics and Drift Patterns

```
analysis_results = {
  'hdfs': {
    'avg_length': 120.5, # characters
    'severity_dist': {
      'INFO': 94.5,
      'ERROR': 3.2,
      'WARN': 2.3
    },
    'drift_indicators': {
      'block_ops_variance': 0.15, # normalized
      'datanode_pattern_shift': 0.08,
      'error_rate_change': 0.03
    }
  },
  'apache': {
    'avg_length': 150.2,
    'message_types': {
      'access': 75.3,
      'error': 15.2,
      'config': 9.5
    },
    'drift_indicators': {
      'request_pattern_shift': 0.12,
      'error_rate_trend': 0.05,
      'module_usage_change': 0.09
    }
  },
  'bgl': {
    'avg_length': 180.7,
    'severity_dist': {
      'INFO': 82.1,
      'ERROR': 12.5,
      'FATAL': 5.4
    },
    'drift_indicators': {
      'node_failure_pattern': 0.11,
      'system_load_shift': 0.07,
      'error_clustering': 0.13
    }
  }
}
```

```
}  
}  
}
```

## 3.2 Temporal and Drift Patterns

- **HDFS:** Regular block operation cycles with drift in replication patterns
- **Apache:** Daily access patterns with request distribution shifts
- **HealthApp:** User activity cycles with sensor drift patterns
- **BGL:** System maintenance windows with node behavior drift
- **HPC:** Job submission patterns with resource utilization shifts
- **Linux:** Service restart patterns with performance degradation indicators
- **Mac:** Application usage cycles with system resource drift

## 4. Implementation Milestones

### 4.1 Data Processing Pipeline

1. **Streaming Log Processing** (Completed)
  - Implemented real-time parsers for all log types
  - Added Kafka integration for streaming ingestion
  - Optimized parsing for minimal latency
2. **Feature Engineering** (In Progress)
  - Real-time feature extraction
  - Drift indicator computation
  - Cross-dataset feature normalization
3. **Drift Detection** (Planned)
  - Multi-source drift detection
  - Adaptive thresholding
  - Real-time model performance monitoring

### 4.2 Monitoring System

1. **Real-time Processing**
  - Multi-source stream processing with Kafka
  - Adaptive buffer management
  - Resource-aware optimization

## 2. Alert System

- Source-specific thresholds
- Correlation-based alerts
- ML model performance monitoring

# 5. Performance Metrics

## 5.1 Processing Performance

- **Stream Processing:** <50ms latency
- **Feature Extraction:** <100ms latency
- **Memory Usage:** <2GB per process
- **CPU Utilization:** <50% per core
- **Throughput:** >10,000 events/second

## 5.2 Detection Performance

Target metrics for each log type:

- HDFS: >95% accuracy, <1% FP, <100ms detection delay
- Apache: >92% accuracy, <2% FP, <150ms detection delay
- HealthApp: >90% accuracy, ❤️ % FP, <200ms detection delay
- BGL: >93% accuracy, <2% FP, <120ms detection delay
- HPC: >91% accuracy, <2.5% FP, <180ms detection delay
- Linux: >90% accuracy, ❤️ % FP, <150ms detection delay
- Mac: >89% accuracy, ❤️ .5% FP, <200ms detection delay

# 6. Next Steps

## 6.1 Implementation Tasks

1. Complete streaming feature engineering pipeline
2. Implement real-time drift detection
3. Develop unified monitoring dashboard
4. Set up ML model performance tracking



## 6.2 Validation Tasks

1. Cross-dataset drift validation
2. Stream processing benchmarking
3. Threshold optimization
4. End-to-end system testing

## References

1. "Mining System Logs with Temporal Rules" (2019)
2. "Feature Engineering for Log Analysis" (2020)
3. "Real-time Log Analysis Techniques" (2021)
4. "System Reliability Monitoring" (2022)
5. "Multi-source Log Analysis" (2023)
6. "Adaptive Drift Detection in System Logs" (2023)
7. "Streaming Analytics for Model Drift" (2023)
8. "Real-time AI System Monitoring" (2024)