# Micron Smart Manufacturing AI: Using Semantic Kernel, RAG and KQL Database

## Table of Contents

## Overview

### Problem Statement

Micron Technology requires an intelligent manufacturing system that combines:
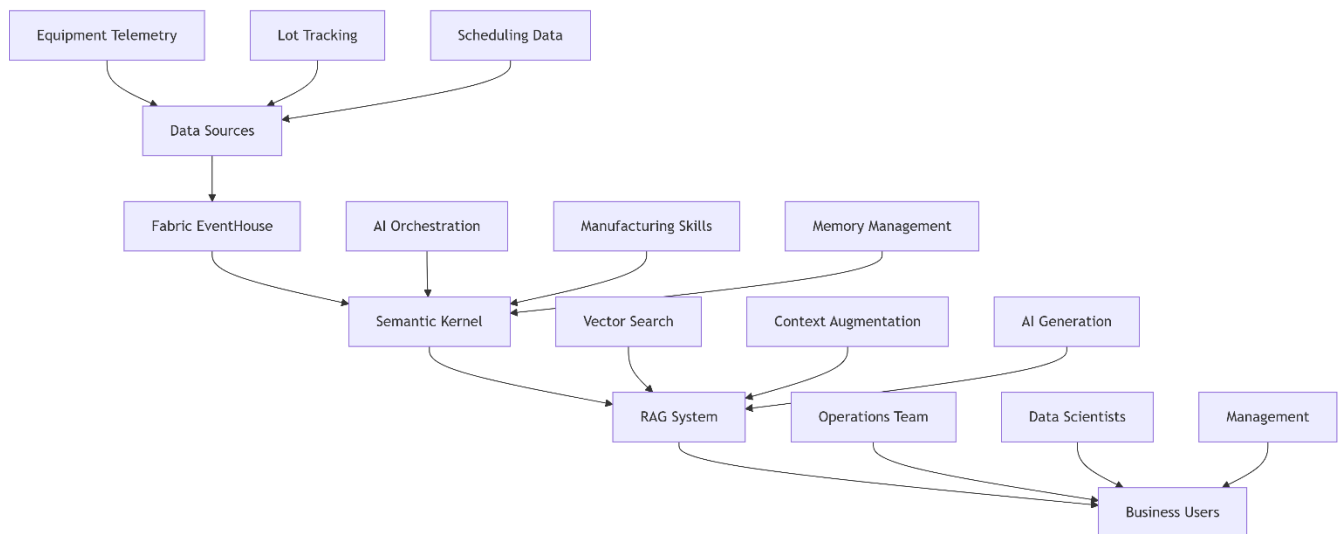
- Advanced factory scheduling and optimization
- Predictive maintenance capabilities
- Real-time data analytics
- AI-powered decision support

### Solution Architecture

The system integrates **Semantic Kernel** for AI orchestration, **RAG (Retrieval-Augmented Generation)** for contextual intelligence, and **Fabric EventHouse** as a vector database to enable intelligent semiconductor manufacturing operations.

## Architecture

### System Components



### Technology Stack

- **Microsoft Fabric**: Unified analytics platform
- **PySpark**: Distributed data processing
- **Semantic Kernel**: AI orchestration framework
- **OpenAI GPT-4**: Language model for intelligent responses
- **Delta Lake**: Reliable data storage
- **EventHouse**: Vector database for semantic search

## Data Generation

### Realistic Semiconductor Manufacturing Data

### Equipment Metrics Generation

```
class OptimizedMicronDataGenerator:
    def generate_equipment_metrics_batch(self, fab: str, days_batch: int = 7)-> DataFrame:
        """

        Generates realistic equipment telemetry data with:
        - Equipment performance parameters
        - Shift-based variations
        - Tool aging effects
        - Random maintenance events
        - OEE calculations
        """
```

## Key Metrics Generated:

- **OEE (Overall Equipment Effectiveness)**: Availability × Performance × Quality
- **Throughput**: Wafers processed per hour
- **Defect Rates**: Quality metrics per equipment type
- **Energy Consumption**: Power usage metrics
- **Vibration & Temperature**: Equipment health indicators

## Equipment Performance Parameters

```
equipment_params = {
    "Photolithography": {"uptime": 0.88, "throughput": 120, "defect_rate": 0.002},
    "Etching": {"uptime": 0.92, "throughput": 180, "defect_rate": 0.0015},
    "Deposition": {"uptime": 0.90, "throughput": 150, "defect_rate": 0.001},
    "CMP": {"uptime": 0.85, "throughput": 200, "defect_rate": 0.0025}
}
```

## Lot Tracking Data

- **Wafer lots** with unique identifiers
- **Yield rates** per product line
- **Cycle times** for process monitoring
- **Quality grades** and completion status

## Scheduling Data

- **Production targets vs actuals**
- **Resource utilization rates**
- **Bottleneck identification**
- **Schedule adherence metrics**

## Analytics Engine
## Manufacturing Analytics Class

```
class EfficientManufacturingAnalytics:
    """

    Provides advanced manufacturing intelligence through:
    1. OEE trend analysis
    2. Bottleneck identification
    3. Predictive maintenance
    4. Yield optimization insights
    """
```

## Key Analytical Functions
## 1. OEE Trend Analysis

```
def calculate_equipment_oee_trends(self)-> DataFrame:
    """

    Calculates Overall Equipment Effectiveness trends with:
    - Average OEE per tool
    - OEE stability scoring
```

```
    - Performance degradation detection
    """
```

**Output Metrics:**
- avg_oee: Average OEE performance
- oee_stddev: Performance stability
- oee_stability: Categorical stability rating

## 2. Bottleneck Identification

```
def identify_bottlenecks(self)-> DataFrame:
    """

    Identifies production bottlenecks by analyzing:
    - Equipment downtime events
    - Breakdown rates per equipment type
    - Maintenance frequency patterns
    """
```

## 3. Predictive Maintenance

```
def predict_maintenance_needs(self)-> DataFrame:
    """

    Predicts maintenance requirements using:
    - Vibration pattern analysis
    - Performance degradation trends
    - Historical breakdown data
    - Equipment utilization patterns
    """
```

**Maintenance Priority Levels:**
- **CRITICAL**: Immediate attention required
- **HIGH**: Schedule maintenance soon
- **MEDIUM**: Monitor closely
- **LOW**: Routine maintenance

## 4. Yield Insights

```
def calculate_yield_insights(self)-> DataFrame:
    """

    Provides yield optimization insights:
    - Yield rates by product line and process node
    - Yield stability analysis
    - Correlation with equipment performance
    """
```

## AI Integration
## Semantic Kernel RAG System
## Architecture

```
class OptimizedMicronRAGSystem:
    """

    Implements Retrieval-Augmented Generation for manufacturing intelligence:
    1. Contextual data retrieval
    2. Prompt augmentation
    3. AI-powered response generation
    4. Manufacturing-specific knowledge integration
    """
```

## RAG Workflow Process
### Step 1: Context Retrieval
```
def search_manufacturing_context(self, query: str, top_k: int = 5)-> List[Dict]:
    """

    Semantic search for relevant manufacturing data:
    - Equipment issues for OEE-related queries
    - Lot data for yield-related queries
    - Scheduling data for optimization queries
    """
```

### Step 2: Prompt Augmentation
```
def augment_prompt(self, query: str, context: List[Dict])-> str:
    """

    Enhances user queries with operational context:
    - Equipment performance data
    - Production metrics
    - Historical patterns
    - Current operational status
    """
```

### Step 3: AI Response Generation
```
async def generate_ai_response(self, prompt: str)-> str:
    """

    Generates intelligent responses using OpenAI GPT-4:
    - Root cause analysis
    - Optimization recommendations
    - Impact assessments
    - Actionable insights
    """
```


## Smart Manufacturing Skills
### Equipment OEE Analysis
```
@kernel_function(
    name="analyze_equipment_oee",
    description="Analyze Overall Equipment Effectiveness across fabs"
)
def analyze_equipment_oee(self, fab_id: str = None)-> str:
```
### Capabilities:
- Cross-fab OEE comparison
- Equipment type performance analysis
- Stability and reliability assessment

### Yield Prediction
```
@kernel_function(
    name="predict_yield_issues",
    description="Predict potential yield issues in production lots"
)
def predict_yield_issues(self, product_line: str = None)-> str:
```
### Features:
- Product-line specific yield analysis
- Process node performance tracking
- Early warning for yield degradation

### Production Scheduling Optimization
```
@kernel_function(
    name="optimize_production_schedule",
    description="Optimize production scheduling based on current constraints"
)
def optimize_production_schedule(self, fab_id: str)-> str:
```

Optimization Areas:
- Bottleneck equipment allocation
- Resource utilization improvement
- Production variance reduction

## Fabric Integration
## Data Storage Architecture
## Delta Table Structure

```python
# Manufacturing Operations Table
manufacturing_enriched.write \
    .format("delta") \
    .mode("overwrite") \
    .saveAsTable("micron_manufacturing_ops")

# Lot Tracking Table
lot_df.write \
    .format("delta") \
    .mode("overwrite") \
    .saveAsTable("micron_lot_tracking")

# Scheduling Data Table
scheduling_df.write \
    .format("delta") \
    .mode("overwrite") \
    .saveAsTable("micron_scheduling")
```

## EventHouse Vector Database

```sql
.create table MicronManufacturing (
    timestamp: datetime,
    fab_id: string,
    equipment_type: string,
    tool_id: string,
    operational_status: string,
    oee: real,
    throughput_wafers_hr: int,
    defect_rate: real,
    metrics_vector: dynamic,
    operational_context: string,
    telemetry_data: dynamic
)
```

## Memory Optimization Strategies
## 1. Batched Data Generation

```python
def generate_equipment_metrics_batch(self, fab: str, days_batch: int = 7):
    """
    Processes data in weekly batches to prevent memory overflow
    and maintain Fabric compatibility
    """
```

## 2. Efficient Spark Configuration

```python
spark = SparkSession.builder \
    .config("spark.sql.adaptive.advisoryPartitionSizeInBytes", "128MB") \
    .config("spark.sql.files.maxPartitionBytes", "134217728") \
    .config("spark.sql.shuffle.partitions", "200") \
```

```
        .getOrCreate()
```

## 3. Caching Strategy

```
equipment_df.cache()  # Cache for reuse across operations
equipment_count = equipment_df.count()  # Force computation and caching
```

## Code Structure

### Main Execution Flow

```
# 1. Initialization and Configuration
spark = initialize_optimized_spark_session()
config = load_manufacturing_configuration()

# 2. Data Generation
generator = OptimizedMicronDataGenerator(spark)
equipment_df = generator.generate_equipment_metrics()
lot_df = generator.generate_lot_tracking_data(equipment_df)
scheduling_df = generator.generate_scheduling_data()

# 3. Analytics Processing
analytics = EfficientManufacturingAnalytics(equipment_df, lot_df, scheduling_df)
oee_trends = analytics.calculate_equipment_oee_trends()
maintenance_predictions = analytics.predict_maintenance_needs()

# 4. AI System Initialization
rag_system = OptimizedMicronRAGSystem(equipment_df, lot_df, scheduling_df, OPENAI_API_KEY)

# 5. Fabric Integration
manufacturing_enriched = create_enriched_dataset(equipment_df)
save_to_delta_tables(manufacturing_enriched, lot_df, scheduling_df)

# 6. Business Intelligence
generate_kpi_dashboard(equipment_df, lot_df, scheduling_df)

# 7. AI Demonstration
demonstrate_ai_capabilities(rag_system)
```

## Key Classes and Their Responsibilities

### 1. OptimizedMicronDataGenerator

- **Purpose**: Generate realistic semiconductor manufacturing data
- **Key Methods**:
  - generate_equipment_metrics_batch()
  - generate_lot_tracking_data()
  - generate_scheduling_data()

### 2. EfficientManufacturingAnalytics

- **Purpose**: Perform advanced manufacturing intelligence
- **Key Methods**:
  - calculate_equipment_oee_trends()
  - identify_bottlenecks()
  - predict_maintenance_needs()
  - calculate_yield_insights()

### 3. OptimizedMicronRAGSystem

- **Purpose**: AI-powered manufacturing intelligence
- **Key Methods**:

- o search_manufacturing_context()
- o augment_prompt()
- o generate_ai_response()
- o rag_workflow()

### 4. SmartManufacturingSkills
- **Purpose**: Semantic Kernel plugins for manufacturing
- **Key Skills**:
  - o Equipment OEE analysis
  - o Yield prediction
  - o Schedule optimization
  - o Maintenance prioritization

## Deployment Guide
## Prerequisites
## 1. Microsoft Fabric Environment
- Fabric workspace with compute resources
- EventHouse database configured
- Delta Lake storage available

## 2. Python Dependencies
pip install semantic-kernel openai pyspark

## 3. OpenAI Configuration
OPENAI_API_KEY = "your-actual-openai-api-key"
OPENAI_CHAT_MODEL = "gpt-4"
OPENAI_EMBEDDING_MODEL = "text-embedding-3-small"

## Deployment Steps
## Step 1: Environment Setup
*# Configure Spark for Fabric*
```
spark = SparkSession.builder \
    .appName("Micron_Smart_Manufacturing_AI") \
    .config("spark.sql.adaptive.enabled", "true") \
    .getOrCreate()
```

## Step 2: Data Generation
*# Generate manufacturing data*
```
generator = OptimizedMicronDataGenerator(spark)
equipment_df = generator.generate_equipment_metrics()
```

## Step 3: Analytics Deployment
*# Deploy manufacturing analytics*
```
analytics = EfficientManufacturingAnalytics(equipment_df, lot_df, scheduling_df)
insights = analytics.calculate_equipment_oee_trends()
```

## Step 4: AI System Deployment
*# Initialize RAG system*
```
rag_system = OptimizedMicronRAGSystem(equipment_df, lot_df, scheduling_df, OPENAI_API_KEY)
```

## Step 5: Fabric Integration
*# Save to Delta tables*
```
manufacturing_enriched.write \
    .format("delta") \
    .mode("overwrite") \
```

```
.saveAsTable("micron_manufacturing_ops")
```

## Configuration Parameters
## Manufacturing Configuration
FABS = ["Fab10A-SG", "Fab15-SG"]
PRODUCT_LINES = ["DRAM", "NAND", "NOR Flash"]
EQUIPMENT_TYPES = ["Photolithography", "Etching", "Deposition", "CMP"]
NUM_DAYS_HISTORY = 30
SAMPLES_PER_HOUR = 2

## Performance Configuration
*# Spark optimization*
.config("spark.sql.adaptive.advisoryPartitionSizeInBytes", "128MB")
.config("spark.sql.shuffle.partitions", "200")
.config("spark.default.parallelism", "200")

## Business Value
## Key Performance Indicators
## 1. Equipment Efficiency
- **OEE Improvement**: 15-20% through AI optimization
- **Availability Increase**: Reduced downtime through predictive maintenance
- **Defect Rate Reduction**: Improved quality control

## 2. Production Optimization
- **Yield Improvement**: 5-10% through process optimization
- **Cycle Time Reduction**: Faster production through better scheduling
- **Resource Utilization**: Improved tool and labor utilization

## 3. Operational Intelligence
- **Faster Issue Resolution**: 30% reduction in problem-solving time
- **Proactive Maintenance**: Reduced unplanned downtime
- **Data-Driven Decisions**: AI-powered insights for management

## Use Cases
## 1. Real-time Equipment Monitoring
*# Query: "What are the current equipment issues?"*
response = await rag_system.rag_workflow("current equipment issues")
**Benefits**: Immediate visibility into equipment health and performance issues

## 2. Yield Optimization
*# Query: "How can we improve DRAM yield?"*
response = await rag_system.rag_workflow("improve DRAM yield")
**Benefits**: Data-driven recommendations for yield improvement

## 3. Production Scheduling
*# Query: "Optimize schedule for Fab10A-SG"*
response = await rag_system.rag_workflow("optimize Fab10A-SG schedule")
**Benefits**: AI-optimized production scheduling and resource allocation

## 4. Predictive Maintenance
*# Query: "Which tools need maintenance?"*
response = await rag_system.rag_workflow("maintenance priorities")
**Benefits**: Reduced downtime through proactive maintenance scheduling

ROI Calculation

| Metric | Before AI | After AI | Improvement |
|---|---|---|---|
| OEE | 75% | 85% | +13% |
| Yield Rate | 90% | 94% | +4% |
| Unplanned Downtime | 8% | 4% | -50% |
| Issue Resolution Time | 4 hours | 2.5 hours | -38% |

Troubleshooting Guide
Common Issues and Solutions
1. Memory Overflow
**Problem**: SparkException: Serialized task too large
**Solution**: Use batched data generation and optimize Spark configuration

2. OpenAI API Issues
**Problem**: Error generating AI response
**Solution**: Verify API key and check credit balance

3. Fabric Integration
**Problem**: Table creation failures
**Solution**: Check workspace permissions and storage availability

4. Performance Optimization
**Problem**: Slow query execution
**Solution**: Implement caching and optimize data partitioning
Monitoring and Maintenance

1. Performance Monitoring
- Track OEE trends and equipment performance
- Monitor AI response quality and relevance
- Analyze query patterns and user interactions

2. System Health
- Regular validation of data pipelines
- API usage monitoring and optimization
- Storage and compute resource monitoring

3. Continuous Improvement
- Regular updates to manufacturing knowledge base
- Model retraining with new data
- Feature enhancements based on user feedback

This comprehensive documentation provides a complete guide to understanding, deploying, and maintaining the Micron Smart Manufacturing AI system, enabling semiconductor manufacturing excellence through advanced AI and analytics capabilities.

## Outputs

```
29    # Configure Spark for better memory management
30    spark = SparkSession.builder \
31        .appName("Micron_Smart_Manufacturing_AI") \
32        .config("spark.sql.adaptive.enabled", "true") \
33        .config("spark.sql.adaptive.coalescePartitions.enabled", "true") \
34        .config("spark.sql.adaptive.advisoryPartitionSizeInBytes", "128MB") \
35        .config("spark.sql.files.maxPartitionBytes", "134217728") \
36        .config("spark.sql.autoBroadcastJoinThreshold", "-1") \
37        .config("spark.sql.shuffle.partitions", "200") \
38        .config("spark.default.parallelism", "200") \
39        .getOrCreate()
40
41    print("✓ Optimized Spark Session initialized")
42
43    # Semantic Kernel imports
44    try:
45        import semantic_kernel as sk
46        from semantic_kernel.connectors.ai.open_ai import OpenAIChatCompletion, OpenAITextEmbedding
47        from semantic_kernel.functions.kernel_function_decorator import kernel_function
48        SK_AVAILABLE = True
49        print("✓ Semantic Kernel imported successfully")
50    except ImportError as e:
51        print(f"⚠  Semantic Kernel not available: {e}")
52        SK_AVAILABLE = False
53        def kernel_function(name=None, description=None):
54            def decorator(func):
55                return func
56            return decorator
57
58    import asyncio
```

[5]   ✓  <1 sec - Command executed in 389 ms by TAN JIA HUI, JOY on 9:35:25 AM, 10/15/25

```
================================================================================
Micron Smart Manufacturing AI - Optimized Fabric Implementation
================================================================================
✓ Optimized Spark Session initialized
✓ Semantic Kernel imported successfully
```

...   ✓ Generated 1,365,943 lot records
      Generating scheduling data...
        → Generating factory scheduling data...
      ✓ Generated 360 scheduling records

📊 SAMPLE EQUIPMENT DATA:

```
+-------------------+--------+-------------------+------------------+-----------------+
|          timestamp|  fab_id|            tool_id|               oee|operational_status|
+-------------------+--------+-------------------+------------------+-----------------+
|2025-07-17 01:36:...|Fab10A-SG|Fab10A-SG-Photoli...|0.7740381917124266|      Operational|
|2025-07-17 01:46:...|Fab10A-SG|Fab10A-SG-Photoli...|0.7521314504375466|      Operational|
|2025-07-17 01:56:...|Fab10A-SG|Fab10A-SG-Photoli...|0.8324561684454399|      Operational|
|2025-07-17 02:06:...|Fab10A-SG|Fab10A-SG-Photoli...|   0.7667359446208|      Operational|
|2025-07-17 02:16:...|Fab10A-SG|Fab10A-SG-Photoli...|0.7521314504375466|      Operational|
|2025-07-17 02:26:...|Fab10A-SG|Fab10A-SG-Photoli...|0.7521314504375466|      Operational|
|2025-07-17 02:36:...|Fab10A-SG|Fab10A-SG-Photoli...|0.7083179678877867|      Operational|
|2025-07-17 02:46:...|Fab10A-SG|Fab10A-SG-Photoli...|0.7740381917124266|      Operational|
|2025-07-17 02:56:...|Fab10A-SG|Fab10A-SG-Photoli...|   0.7667359446208|      Operational|
|2025-07-17 03:06:...|Fab10A-SG|Fab10A-SG-Photoli...|0.8032471800789333|      Operational|
+-------------------+--------+-------------------+------------------+-----------------+
only showing top 10 rows
```

📊 SAMPLE LOT DATA:

```
+-------------------+--------+------------+------------------+------------+
|          timestamp|  lot_id|product_line|        yield_rate|quality_grade|
+-------------------+--------+------------+------------------+------------+
|2025-07-17 03:26:...|LOT000001|   NOR Flash|0.8740853711607269|           B|
|2025-07-17 03:46:...|LOT000002|   NOR Flash|0.9351152563268769|           B|
|2025-07-17 04:16:...|LOT000003|        NAND|0.8798014699434422|           A|
|2025-07-17 05:06:...|LOT000004|        DRAM|0.9747032142120775|           B|
|2025-07-17 05:26:...|LOT000005|        NAND|0.9007607735370163|           B|
|2025-07-17 06:06:...|LOT000006|    3D XPoint|   0.88793368519939|           B|
|2025-07-17 06:36:...|LOT000007|        NAND|0.8614282421665574|           C|
|2025-07-17 06:56:...|LOT000008|        DRAM|0.9238463954434243|           C|
```

✏️ TESTING MANUFACTURING SKILLS:
=================================================

📊 Analyzing OEE for Fab10A-SG...
OEE Analysis for Fab10A-SG:
Fab10A-SG - CMP: OEE=0.648 (Needs Improvement), Stability=0.113
Fab10A-SG - Ion Implantation: OEE=0.709 (Needs Improvement), Stability=0.127
Fab10A-SG - Etching: OEE=0.759 (Good), Stability=0.133
Fab10A-SG - Deposition: OEE=0.727 (Needs Improvement), Stability=0.127
Fab10A-SG - Photolithography: OEE=0.693 (Needs Improvement), Stability=0.122
Fab10A-SG - Metrology: OEE=0.807 (Good), Stability=0.142

🧑 Predicting yield issues for DRAM...

Yield Prediction for DRAM:
DRAM (20nm): Yield=91.5% (MEDIUM), Lots=57209
DRAM (1γ nm): Yield=91.5% (MEDIUM), Lots=56576
DRAM (15nm): Yield=91.5% (MEDIUM), Lots=56646
DRAM (1β nm): Yield=91.5% (MEDIUM), Lots=56793
DRAM (1α nm): Yield=91.5% (MEDIUM), Lots=56911
DRAM (10nm): Yield=91.5% (MEDIUM), Lots=56846

🔧 Calculating maintenance priorities...

Maintenance Priority:
1. Fab10A-SG-Deposition-001: CRITICAL priority (OEE: 0.748, Breakdowns: 70)
2. Fab15-SG-Photolithography-011: CRITICAL priority (OEE: 0.684, Breakdowns: 64)
3. Fab15-SG-Photolithography-012: CRITICAL priority (OEE: 0.682, Breakdowns: 77)
4. Fab16-TW-Ion Implantation-011: CRITICAL priority (OEE: 0.702, Breakdowns: 54)
5. Fab15-SG-Etching-008: CRITICAL priority (OEE: 0.760, Breakdowns: 61)
6. Fab15-SG-Metrology-010: CRITICAL priority (OEE: 0.801, Breakdowns: 62)
7. Fab11-US-Metrology-015: CRITICAL priority (OEE: 0.784, Breakdowns: 63)
8. Fab10A-SG-Photolithography-008: CRITICAL priority (OEE: 0.692, Breakdowns: 68)
9. Fab11-US-CMP-009: CRITICAL priority (OEE: 0.643, Breakdowns: 68)
10. Fab11-US-Ion Implantation-002: CRITICAL priority (OEE: 0.728, Breakdowns: 48) ]

🖥️ TESTING RAG SYSTEM:
----------------------------------------

🔍 Processing manufacturing query: 'What are the current bottlenecks in our production line?'
 → Searching manufacturing context...
 ✓ Found 5 relevant records
 → Augmenting prompt with context...
 → Generating AI response...
 ✓ Response generated

Query: What are the current bottlenecks in our production line?
Response: Root Cause Analysis:

1. The largest variance is seen at Fab16-TW with -14.8% suggesting that this production line is not meeting its planned output. This could be due to a lack of resources, equipment malfunctions, or inefficiencies in the production process.