

DataHub (Acryl Data) Backend Engineer Interview Preparation Guide

Comprehensive Solutions & Strategies

November 9, 2025

Contents

1	Introduction	3
1.1	About DataHub	3
1.2	Interview Process	3
1.3	Tech Stack	3
1.4	Key Skills to Demonstrate	3
2	Problem 1: Metadata Lineage Graph	4
2.1	Problem Statement	4
2.2	Solution	4
2.3	Complexity Analysis	5
2.4	Key Insights	5
2.5	Production Enhancements	5
3	Problem 2: Metadata Search Index	6
3.1	Problem Statement	6
3.2	Solution	6
3.3	Complexity Analysis	7
3.4	Key Insights	7
3.5	Production Enhancements	7
4	Problem 3: Schema Evolution Tracking	8
4.1	Problem Statement	8
4.2	Solution	8
4.3	Complexity Analysis	9
4.4	Key Insights	9
4.5	Schema Registry Integration	9
5	Problem 4: Data Quality Monitoring	10
5.1	Problem Statement	10
5.2	Solution	10
5.3	Complexity Analysis	11
5.4	Key Insights	11
5.5	Production Enhancements	11
6	IMPORTANT: Java Concurrency Problems	12
6.1	Why Java for DataHub?	12

7 Problem 5: Thread-Safe Metadata Cache	13
7.1 Problem Statement	13
7.2 Solution	13
7.3 Complexity Analysis	14
7.4 Key Insights	14
7.5 Follow-up Questions	14
8 Problem 6: Producer-Consumer Metadata Ingestion	15
8.1 Problem Statement	15
8.2 Solution	15
8.3 Complexity Analysis	16
8.4 Key Insights	16
8.5 Production Considerations	17
9 Problem 7: Concurrent Batch Processor	18
9.1 Problem Statement	18
9.2 Solution	18
9.3 Complexity Analysis	19
9.4 Key Insights	19
9.5 Optimizations	20
10 Java Concurrency Quick Reference	21
10.1 Essential Concurrency Primitives	21
10.2 Thread-Safe Collections	21
10.3 Common Pitfalls	21
11 System Design Topics	22
11.1 Metadata Catalog Architecture	22
11.2 Scalability Considerations	22
11.3 Data Lineage at Scale	22
12 DataHub-Specific Knowledge	23
12.1 Core Entities	23
12.2 Metadata Aspects	23
12.3 Ingestion Framework	23
13 Behavioral Interview Prep	24
13.1 DataHub/Acryl Data Values	24
13.2 Common Questions	24
13.3 Questions to Ask	24
14 Additional Resources	25
14.1 Technical Learning	25
14.2 Practice Platforms	25
14.3 Data Infrastructure Blogs	25
15 Final Tips	25

1 Introduction

1.1 About DataHub

DataHub is the leading open-source metadata platform that helps organizations discover, understand, and trust their data. Founded by LinkedIn engineers and now commercialized by Acryl Data (founded 2020, ~50 employees), DataHub provides a comprehensive solution for data cataloging, lineage tracking, and data governance.

1.2 Interview Process

Based on typical data infrastructure companies:

1. **Recruiter Screen** - Background and interest (30 min)
2. **Technical Phone Screen** - LeetCode medium + data structures (45-60 min)
3. **Take-Home Assignment** - Metadata API or pipeline (2-4 hours)
4. **Onsite Interviews** (4-5 hours):
 - System Design - Metadata catalog design
 - Coding - Graphs, trees, search algorithms
 - Technical Deep Dive - Past projects
 - Cultural Fit - Team collaboration

1.3 Tech Stack

DataHub's architecture:

- **Backend:** Java (Spring Boot), Python
- **Storage:** MySQL/PostgreSQL, Elasticsearch
- **Graph:** Neo4j for lineage
- **Streaming:** Apache Kafka
- **Frontend:** React, TypeScript
- **Infrastructure:** Kubernetes, Docker

1.4 Key Skills to Demonstrate

- Strong graph algorithm knowledge (critical!)
- Distributed systems and scalability
- Data engineering fundamentals
- Search and indexing expertise
- System design thinking

2 Problem 1: Metadata Lineage Graph

2.1 Problem Statement

DataHub tracks data lineage to understand dependencies between datasets. Implement a system to find all downstream dependencies of a given dataset and detect circular dependencies.

Requirements:

- Build a directed graph of dataset dependencies
- Find all datasets that depend on a given dataset (downstream)
- Detect circular dependencies
- Return dependencies in traversal order

Example:

```

1 lineage = MetadataLineage()
2 lineage.add_dependency("table_c", "table_a") # c depends on a
3 lineage.add_dependency("table_c", "table_b") # c depends on b
4 lineage.add_dependency("table_d", "table_c") # d depends on c
5
6 lineage.get_downstream("table_a") # Returns ["table_c", "table_d"]
7 lineage.has_circular_dependency() # Returns False

```

2.2 Solution

```

1 from typing import List, Dict, Set
2 from collections import defaultdict, deque
3
4 class MetadataLineage:
5     def __init__(self):
6         self.graph = defaultdict(list) # source -> [targets that
depend on it]
7         self.all_nodes = set()
8
9     def add_dependency(self, target: str, source: str) -> None:
10        self.graph[source].append(target)
11        self.all_nodes.add(source)
12        self.all_nodes.add(target)
13
14    def get_downstream(self, dataset: str) -> List[str]:
15        visited = set()
16        queue = deque([dataset])
17
18        while queue:
19            node = queue.popleft()
20            if node in visited:
21                continue
22            visited.add(node)
23
24            for neighbor in self.graph[node]:
25                if neighbor not in visited:
26                    queue.append(neighbor)
27
28        visited.discard(dataset) # Don't include the source dataset
29        return list(visited)

```

```
30
31     def has_circular_dependency(self) -> bool:
32         visited = set()
33         rec_stack = set()
34
35         def dfs(node):
36             visited.add(node)
37             rec_stack.add(node)
38
39             for neighbor in self.graph[node]:
40                 if neighbor not in visited:
41                     if dfs(neighbor):
42                         return True
43                 elif neighbor in rec_stack:
44                     return True
45
46             rec_stack.remove(node)
47             return False
48
49         for node in self.all_nodes:
50             if node not in visited:
51                 if dfs(node):
52                     return True
53
54         return False
```

2.3 Complexity Analysis

- **Time Complexity:** $O(V + E)$ for both operations where V is vertices (datasets), E is edges (dependencies)
- **Space Complexity:** $O(V + E)$ for graph storage plus $O(V)$ for traversal

2.4 Key Insights

1. BFS efficiently finds all reachable nodes (downstream dependencies)
2. DFS with recursion stack detects cycles
3. Directed graph naturally models data lineage (A produces B)
4. Production systems need to handle millions of nodes and edges

2.5 Production Enhancements

- **Graph Database:** Use Neo4j for efficient graph queries at scale
- **Caching:** Cache frequently accessed lineage paths
- **Incremental Updates:** Only recompute affected subgraphs
- **Column-Level Lineage:** Track individual column dependencies
- **Visualization:** Generate visual lineage graphs for users

3 Problem 2: Metadata Search Index

3.1 Problem Statement

DataHub needs efficient search across dataset metadata. Implement a search index that supports prefix matching and ranks results by relevance.

Requirements:

- Index datasets with name, description, and tags
- Support prefix search (e.g., "user" matches "user_events")
- Rank results by relevance (exact match > prefix match)
- Return top K results

Example:

```

1 index = MetadataSearchIndex()
2 index.add_dataset("user_events", "User activity events", ["analytics"])
3 index.add_dataset("user_profile", "User profile data", ["users"])
4 index.add_dataset("product_sales", "Sales transactions", ["sales"])
5
6 index.search("user", k=2) # Returns ["user_events", "user_profile"]
7 index.search("sal", k=1) # Returns ["product_sales"]

```

3.2 Solution

```

1 from typing import List, Dict
2
3 class MetadataSearchIndex:
4     def __init__(self):
5         self.datasets = {} # name -> {description, tags}
6
7     def add_dataset(self, name: str, description: str, tags: List[str]) -> None:
8         self.datasets[name] = {
9             "description": description,
10            "tags": tags
11        }
12
13    def search(self, query: str, k: int) -> List[str]:
14        query_lower = query.lower()
15        scores = []
16
17        for name, metadata in self.datasets.items():
18            score = 0
19            name_lower = name.lower()
20
21            # Exact name match - highest score
22            if name_lower == query_lower:
23                score += 100
24
25            # Name starts with query - high score
26            elif name_lower.startswith(query_lower):
27                score += 50
28
29            # Query in name - medium score

```

```
30         elif query_lower in name_lower:
31             score += 25
32
33     # Check tags
34     for tag in metadata["tags"]:
35         tag_lower = tag.lower()
36         if tag_lower == query_lower:
37             score += 30
38         elif query_lower in tag_lower:
39             score += 10
40
41     # Check description
42     desc_lower = metadata["description"].lower()
43     if query_lower in desc_lower:
44         score += 5
45
46     if score > 0:
47         scores.append((score, name))
48
49     # Sort by score descending, then alphabetically
50     scores.sort(key=lambda x: (-x[0], x[1]))
51
52     # Return top k results
53     return [name for score, name in scores[:k]]
```

3.3 Complexity Analysis

- **Time Complexity:** $O(N \times M)$ where N is datasets, M is avg metadata length; $O(N \log N)$ for sorting
- **Space Complexity:** $O(N \times M)$ for storing all metadata

3.4 Key Insights

1. Simple relevance scoring works for small datasets
2. Exact matches should rank higher than partial matches
3. Production systems use Elasticsearch or similar
4. Consider TF-IDF or BM25 for better ranking

3.5 Production Enhancements

- **Elasticsearch:** Full-text search with analyzers and tokenizers
- **Fuzzy Matching:** Handle typos with edit distance
- **Popularity Boost:** Factor in usage statistics
- **Personalization:** Rank based on user's team or domain
- **Faceted Search:** Filter by type, owner, freshness

4 Problem 3: Schema Evolution Tracking

4.1 Problem Statement

DataHub tracks schema changes over time to ensure compatibility. Implement a system to detect changes and check compatibility between schema versions.

Requirements:

- Store schema versions with metadata
- Detect schema changes (added/removed/modified fields)
- Check backward compatibility (new schema reads old data)
- Check forward compatibility (old schema reads new data)

Example:

```

1 tracker = SchemaEvolution()
2
3 schema_v1 = {"user_id": "int", "name": "string", "email": "string"}
4 schema_v2 = {"user_id": "int", "name": "string", "email": "string",
5               "created_at": "timestamp"}
6
7 tracker.add_schema("users", schema_v1, version=1)
8 tracker.add_schema("users", schema_v2, version=2)
9
10 changes = tracker.get_changes("users", version_from=1, version_to=2)
11 # Returns: {"added": ["created_at"], "removed": [], "modified": []}
12
13 tracker.is_backward_compatible("users", 2, 1) # Returns True

```

4.2 Solution

```

1 from typing import Dict, List
2
3 class SchemaEvolution:
4     def __init__(self):
5         self.schemas = {} # (table_name, version) -> schema dict
6
7     def add_schema(self, table_name: str, schema: Dict[str, str],
8                   version: int) -> None:
9         self.schemas[(table_name, version)] = schema.copy()
10
11     def get_changes(self, table_name: str, version_from: int,
12                     version_to: int) -> Dict[str, List[str]]:
13         schema_old = self.schemas.get((table_name, version_from), {})
14         schema_new = self.schemas.get((table_name, version_to), {})
15
16         added = []
17         removed = []
18         modified = []
19
20         # Find added and modified fields
21         for field, new_type in schema_new.items():
22             if field not in schema_old:
23                 added.append(field)
24             elif schema_old[field] != new_type:

```

```

23             modified.append(field)
24
25     # Find removed fields
26     for field in schema_old:
27         if field not in schema_new:
28             removed.append(field)
29
30     return {
31         "added": added,
32         "removed": removed,
33         "modified": modified
34     }
35
36     def is_backward_compatible(self, table_name: str, new_version: int,
37                               old_version: int) -> bool:
38         changes = self.get_changes(table_name, old_version, new_version
39 )
40
41         # Backward compatible if:
42         # 1. No fields removed (new schema has all old fields)
43         # 2. No field types changed
44         if changes["removed"] or changes["modified"]:
45             return False
46
47     return True

```

4.3 Complexity Analysis

- **Time Complexity:** $O(n)$ where n is number of fields in schemas
- **Space Complexity:** $O(v \times n)$ where v is versions, n is fields per version

4.4 Key Insights

1. Backward compatibility: new code reads old data (adding optional fields OK)
2. Forward compatibility: old code reads new data (more restrictive)
3. Schema evolution is critical for zero-downtime deployments
4. Track not just schema but also statistics and sample values

4.5 Schema Registry Integration

In production DataHub:

- **Avro/Protobuf:** Integration with schema registries
- **Compatibility Modes:** BACKWARD, FORWARD, FULL, NONE
- **Version History:** Complete audit trail of all changes
- **Impact Analysis:** Show which pipelines affected by schema change
- **Alerts:** Notify owners of breaking changes

5 Problem 4: Data Quality Monitoring

5.1 Problem Statement

DataHub monitors data quality metrics to detect anomalies. Implement a system that tracks metrics over time and detects significant deviations using statistical methods.

Requirements:

- Track metrics over time (row count, null percentage, distinct values)
- Calculate rolling statistics (mean, standard deviation)
- Detect anomalies using Z-score (values beyond threshold)
- Return alerts for anomalous metrics

Example:

```
1 monitor = DataQualityMonitor(window_size=5, threshold=2.0)
2
3 # Normal metrics
4 monitor.add_metric("table1", "row_count", 1000)
5 monitor.add_metric("table1", "row_count", 1050)
6 monitor.add_metric("table1", "row_count", 980)
7 monitor.add_metric("table1", "row_count", 1020)
8 monitor.add_metric("table1", "row_count", 1010)
9
10 # Anomalous metric
11 result = monitor.add_metric("table1", "row_count", 5000)
12 # Returns: {"is_anomaly": True, "z_score": 3.5, ...}
```

5.2 Solution

```
1 from typing import Dict
2 from collections import deque
3 import statistics
4
5 class DataQualityMonitor:
6     def __init__(self, window_size: int, threshold: float):
7         self.window_size = window_size
8         self.threshold = threshold
9         self.metrics = {} # (table, metric_name) -> deque of values
10
11     def add_metric(self, table_name: str, metric_name: str, value: float) -> Dict:
12         key = (table_name, metric_name)
13
14         if key not in self.metrics:
15             self.metrics[key] = deque(maxlen=self.window_size)
16
17         history = self.metrics[key]
18
19         # Calculate stats before adding new value
20         is_anomaly = False
21         z_score = 0.0
22         mean = 0.0
23         std_dev = 0.0
24
```

```
25     if len(history) >= 2:
26         mean = statistics.mean(history)
27         std_dev = statistics.stdev(history)
28
29         if std_dev > 0:
30             z_score = abs((value - mean) / std_dev)
31             is_anomaly = z_score > self.threshold
32
33     # Add current value to history
34     history.append(value)
35
36     return {
37         "is_anomaly": is_anomaly,
38         "z_score": z_score,
39         "mean": mean,
40         "std_dev": std_dev
41     }
```

5.3 Complexity Analysis

- **Time Complexity:** $O(w)$ where w is window size (for statistics calculation)
- **Space Complexity:** $O(t \times m \times w)$ where t is tables, m is metrics, w is window size

5.4 Key Insights

1. Z-score detects outliers based on historical pattern
2. Rolling window adapts to changing baselines
3. Threshold of 2.0-3.0 balances sensitivity vs false positives
4. Production systems use more sophisticated methods (ML-based)

5.5 Production Enhancements

- **Multiple Algorithms:** Combine Z-score, IQR, ARIMA, ML models
- **Seasonality:** Account for daily/weekly patterns
- **Multi-metric Correlation:** Detect related anomalies
- **Alert Routing:** Notify data owners via Slack/email
- **Historical Analysis:** Store all anomalies for pattern analysis

6 IMPORTANT: Java Concurrency Problems

6.1 Why Java for DataHub?

DataHub explicitly tests multithreading and concurrency in their first interview round.

Key reasons to use Java (not Python) for concurrency:

- DataHub's backend is Java (Spring Boot)
- Java has robust concurrency primitives
- Python's GIL limits true parallelism
- Interview will expect Java concurrency knowledge

Problems 5-7 are in Java and focus specifically on concurrency patterns used in DataHub's production systems.

7 Problem 5: Thread-Safe Metadata Cache

7.1 Problem Statement

DataHub's metadata service needs a thread-safe cache for frequently accessed metadata. Implement a thread-safe LRU cache that handles concurrent reads and writes.

Requirements:

- Thread-safe get() and put() operations
- LRU eviction policy when capacity is reached
- Support concurrent reads (multiple threads simultaneously)
- Writes must be exclusive (only one writer at a time)
- No data corruption under concurrent access

Example:

```
1 MetadataCache cache = new MetadataCache(2);
2
3 // Thread 1
4 cache.put("dataset1", "metadata1");
5
6 // Thread 2 (concurrent)
7 cache.put("dataset2", "metadata2");
8
9 // Thread 3 (concurrent read)
10 String m = cache.get("dataset1");
```

7.2 Solution

```
1 import java.util.*;
2 import java.util.concurrent.locks.*;
3
4 class MetadataCache {
5     private final int capacity;
6     private final LinkedHashMap<String, String> cache;
7     private final ReadWriteLock lock;
8     private final Lock readLock;
9     private final Lock writeLock;
10
11     public MetadataCache(int capacity) {
12         this.capacity = capacity;
13         this.cache = new LinkedHashMap<String, String>(capacity, 0.75f,
14             true) {
15             @Override
16             protected boolean removeEldestEntry(Map.Entry eldest) {
17                 return size() > MetadataCache.this.capacity;
18             }
19         };
20         this.lock = new ReentrantReadWriteLock();
21         this.readLock = lock.readLock();
22         this.writeLock = lock.writeLock();
23     }
24
25     public String get(String key) {
```

```
25     readLock.lock();
26     try {
27         return cache.get(key);
28     } finally {
29         readLock.unlock();
30     }
31 }
32
33 public void put(String key, String value) {
34     writeLock.lock();
35     try {
36         cache.put(key, value);
37     } finally {
38         writeLock.unlock();
39     }
40 }
41
42 public int size() {
43     readLock.lock();
44     try {
45         return cache.size();
46     } finally {
47         readLock.unlock();
48     }
49 }
50 }
```

7.3 Complexity Analysis

- **Time Complexity:** $O(1)$ for both get and put operations
- **Space Complexity:** $O(\text{capacity})$

7.4 Key Insights

1. **ReadWriteLock:** Allows multiple concurrent readers, but exclusive writer
2. **LinkedHashMap:** Maintains insertion/access order for LRU
3. **accessOrder=true:** Moves accessed entries to end (LRU behavior)
4. **removeEldestEntry:** Automatic eviction when capacity exceeded
5. **Always unlock in finally:** Ensures lock is released even on exception

7.5 Follow-up Questions

1. Why use ReadWriteLock instead of synchronized?
2. What if we needed distributed caching across multiple servers?
3. How would you handle cache invalidation?
4. What are the trade-offs between LRU, LFU, and FIFO?

8 Problem 6: Producer-Consumer Metadata Ingestion

8.1 Problem Statement

DataHub ingests metadata from multiple sources concurrently. Implement a producer-consumer pattern where producers add metadata events to a queue and consumers process them.

Requirements:

- Thread-safe queue for metadata events
- Multiple producers add events concurrently
- Multiple consumers process events concurrently
- Graceful shutdown (process all pending events)
- No events lost or processed twice

Example:

```

1 MetadataIngestionPipeline pipeline = new MetadataIngestionPipeline(10);
2
3 // Start 3 consumer threads
4 pipeline.startConsumers(3);
5
6 // Producers add events
7 pipeline.addEvent(new MetadataEvent("dataset1", "SCHEMA_CHANGE"));
8 pipeline.addEvent(new MetadataEvent("dataset2", "LINEAGE_UPDATE"));
9
10 // Shutdown and wait for completion
11 pipeline.shutdown();

```

8.2 Solution

```

1 import java.util.concurrent.*;
2 import java.util.concurrent.atomic.*;
3
4 class MetadataIngestionPipeline {
5     private final BlockingQueue<MetadataEvent> queue;
6     private ExecutorService executorService;
7     private final AtomicInteger processedCount;
8     private volatile boolean isShutdown;
9     private static final MetadataEvent POISON_PILL =
10         new MetadataEvent("POISON", "PILL");
11
12     public MetadataIngestionPipeline(int queueSize) {
13         this.queue = new ArrayBlockingQueue<>(queueSize);
14         this.processedCount = new AtomicInteger(0);
15         this.isShutdown = false;
16     }
17
18     public void startConsumers(int numConsumers) {
19         executorService = Executors.newFixedThreadPool(numConsumers);
20
21         for (int i = 0; i < numConsumers; i++) {
22             executorService.submit(() -> {
23                 try {
24                     while (true) {

```

```
25         MetadataEvent event = queue.take();
26         if (event == POISON_PILL) {
27             break;
28         }
29         // Process event
30         Thread.sleep(10); // Simulate processing
31         processedCount.incrementAndGet();
32     }
33 } catch (InterruptedException e) {
34     Thread.currentThread().interrupt();
35 }
36 });
37 }
38 }
39
40 public boolean addEvent(MetadataEvent event) throws
41 InterruptedException {
42     if (isShutdown) {
43         return false;
44     }
45     queue.put(event); // Blocks if queue is full
46     return true;
47 }
48
49 public void shutdown() throws InterruptedException {
50     isShutdown = true;
51
52     // Add poison pills for each consumer
53     int consumers = ((ThreadPoolExecutor) executorService).
54     getActiveCount();
55     for (int i = 0; i < consumers; i++) {
56         queue.put(POISON_PILL);
57     }
58
59     executorService.shutdown();
60     executorService.awaitTermination(10, TimeUnit.SECONDS);
61 }
62
63 public int getProcessedCount() {
64     return processedCount.get();
65 }
66 }
```

8.3 Complexity Analysis

- **Time Complexity:** O(1) for add/poll operations
- **Space Complexity:** O(queue_size)

8.4 Key Insights

1. **BlockingQueue:** Thread-safe, handles blocking automatically
2. **Poison Pill:** Sentinel value to signal consumer shutdown
3. **AtomicInteger:** Lock-free counter for thread safety
4. **ExecutorService:** Manages thread lifecycle

5. **Graceful Shutdown:** Process pending events before stopping

8.5 Production Considerations

- **Backpressure:** ArrayBlockingQueue naturally provides backpressure
- **Error Handling:** Wrap event processing in try-catch
- **Monitoring:** Track queue depth, processing rate, error rate
- **Distributed:** Use Kafka for distributed producer-consumer

9 Problem 7: Concurrent Batch Processor

9.1 Problem Statement

DataHub needs to process large batches of metadata updates efficiently. Implement a concurrent batch processor using thread pools.

Requirements:

- Use ExecutorService with thread pool
- Process batches concurrently
- Collect results from all threads safely
- Handle exceptions without crashing
- Proper shutdown and cleanup

Example:

```
1 BatchProcessor processor = new BatchProcessor(4);
2
3 List<String> datasets = Arrays.asList("ds1", "ds2", ..., "ds100");
4 List<String> results = processor.processBatch(datasets);
5
6 processor.shutdown();
```

9.2 Solution

```
1 import java.util.*;
2 import java.util.concurrent.*;
3
4 class BatchProcessor {
5     private final ExecutorService executorService;
6
7     public BatchProcessor(int threadPoolSize) {
8         this.executorService = Executors.newFixedThreadPool(
9             threadPoolSize);
10    }
11
12    public List<String> processBatch(List<String> datasets)
13        throws InterruptedException, ExecutionException {
14        List<Future<String>> futures = new ArrayList<>();
15
16        // Submit all tasks
17        for (String dataset : datasets) {
18            Future<String> future = executorService.submit(() -> {
19                return processDataset(dataset);
20            });
21            futures.add(future);
22        }
23
24        // Collect results
25        List<String> results = new ArrayList<>();
26        for (Future<String> future : futures) {
27            results.add(future.get()); // Blocking
28        }
29    }
30}
```

```
29         return results;
30     }
31
32     public ProcessResult processBatchWithFailures(List<String> datasets
33 ) throws InterruptedException {
34     List<Future<String>> futures = new ArrayList<>();
35
36     for (String dataset : datasets) {
37         futures.add(executorService.submit(() -> processDataset(
38 dataset)));
39     }
40
41     List<String> successes = new ArrayList<>();
42     int failures = 0;
43
44     for (Future<String> future : futures) {
45         try {
46             successes.add(future.get());
47         } catch (ExecutionException e) {
48             failures++;
49         }
50     }
51
52     return new ProcessResult(successes, failures);
53 }
54
55     public void shutdown() throws InterruptedException {
56         executorService.shutdown();
57         executorService.awaitTermination(60, TimeUnit.SECONDS);
58     }
59
60     private String processDataset(String dataset) throws Exception {
61         Thread.sleep(10); // Simulate work
62         if (dataset.contains("fail")) {
63             throw new Exception("Failed: " + dataset);
64         }
65         return "processed_" + dataset;
66     }
}
```

9.3 Complexity Analysis

- **Time Complexity:** $O(n/p)$ where n is items, p is threads (parallelization)
- **Space Complexity:** $O(n)$ for results

9.4 Key Insights

1. **Future:** Represents async computation result
2. **ExecutorService:** Reuses threads, avoids overhead
3. **Thread Pool Sizing:** CPU-bound = cores + 1; I/O-bound = cores \times 2
4. **Exception Handling:** ExecutionException wraps task exceptions
5. **Shutdown:** Always call shutdown() and awaitTermination()

9.5 Optimizations

- **CompletionService:** Process results as they complete (not in order)
- **ForkJoinPool:** For recursive divide-and-conquer tasks
- **Timeout:** Use future.get(timeout, unit) to avoid hanging
- **Batch Size:** Tune for optimal throughput vs latency

10 Java Concurrency Quick Reference

10.1 Essential Concurrency Primitives

- **synchronized**: Basic mutual exclusion (method or block level)
- **ReentrantLock**: Explicit lock with tryLock, fairness options
- **ReadWriteLock**: Many readers OR one writer (cache scenarios)
- **Semaphore**: Limit concurrent access to N threads
- **CountDownLatch**: Wait for N events (one-time barrier)
- **CyclicBarrier**: Synchronize N threads repeatedly
- **Atomic Classes**: Lock-free thread-safe variables

10.2 Thread-Safe Collections

- **ConcurrentHashMap**: Thread-safe map with fine-grained locking
- **CopyOnWriteArrayList**: Thread-safe list (read-heavy workloads)
- **BlockingQueue**: Producer-consumer pattern
- **ConcurrentLinkedQueue**: Non-blocking concurrent queue

10.3 Common Pitfalls

1. **Forgetting unlock**: Always use try-finally with explicit locks
2. **Deadlock**: Acquire locks in consistent order
3. **Race Conditions**: Use atomic operations or synchronization
4. **Ignoring InterruptedException**: Restore interrupt flag
5. **Shared Mutable State**: Minimize or use immutable objects

11 System Design Topics

11.1 Metadata Catalog Architecture

Key components of a metadata platform:

1. **Metadata Store:** Primary database for all metadata
2. **Search Index:** Elasticsearch for fast discovery
3. **Graph Database:** Neo4j for lineage queries
4. **Ingestion:** Kafka for real-time metadata events
5. **API Layer:** RESTful and GraphQL APIs
6. **UI:** React frontend for data discovery

11.2 Scalability Considerations

Critical for enterprise deployments:

- **Sharding:** Partition metadata by domain or data source
- **Caching:** Redis for frequently accessed metadata
- **Read Replicas:** Scale read queries independently
- **Async Processing:** Background jobs for lineage computation
- **Rate Limiting:** Protect from ingestion spikes

11.3 Data Lineage at Scale

Challenges and solutions:

- **Graph Size:** Billions of nodes, trillions of edges
- **Query Performance:** Sub-second lineage traversal
- **Real-time Updates:** Incremental lineage computation
- **Multi-hop Queries:** Efficient path finding algorithms
- **Visualization:** Render large graphs without overwhelming UI

12 DataHub-Specific Knowledge

12.1 Core Entities

DataHub models these entities:

- **Datasets:** Tables, files, topics
- **Charts:** BI visualizations
- **Dashboards:** Collections of charts
- **Data Jobs:** Pipelines, workflows
- **ML Models:** Machine learning models
- **Containers:** Databases, schemas

12.2 Metadata Aspects

Each entity has multiple aspects:

- **Schema:** Column names, types, descriptions
- **Ownership:** Individuals and teams
- **Tags:** Free-form labels
- **Glossary Terms:** Business vocabulary
- **Lineage:** Upstream and downstream dependencies
- **Usage Stats:** Query frequency, users
- **Properties:** Custom key-value pairs

12.3 Ingestion Framework

DataHub supports 100+ data sources:

- **Databases:** MySQL, PostgreSQL, Oracle, SQL Server
- **Warehouses:** Snowflake, BigQuery, Redshift
- **BI Tools:** Looker, Tableau, PowerBI
- **Orchestration:** Airflow, Dagster
- **Streaming:** Kafka, Pulsar

13 Behavioral Interview Prep

13.1 DataHub/Acryl Data Values

- **Data Democratization:** Make data accessible to all
- **Open Source First:** Community-driven development
- **Engineering Excellence:** Scalable, reliable systems
- **Customer Success:** Solve real data discovery problems

13.2 Common Questions

1. Why are you interested in data infrastructure?
2. Tell me about a time you built a scalable distributed system
3. How do you approach system design with conflicting requirements?
4. Describe your experience with graph algorithms
5. What's your approach to contributing to open-source projects?

13.3 Questions to Ask

- What are the biggest technical challenges in scaling DataHub?
- How do you balance open-source community vs commercial features?
- What's the roadmap for column-level lineage?
- How does the team handle on-call and production incidents?
- What opportunities exist for technical leadership growth?

14 Additional Resources

14.1 Technical Learning

- **DataHub GitHub:** Study the codebase and architecture docs
- **Algorithms:** Focus on graph algorithms (BFS, DFS, topological sort)
- **System Design:** "Designing Data-Intensive Applications"
- **Data Engineering:** "Fundamentals of Data Engineering"
- **Similar Projects:** Apache Atlas, Amundsen, LinkedIn WhereHows

14.2 Practice Platforms

- **LeetCode:** Graph problems (medium/hard), trees, search
- **System Design:** Grokking courses, YouTube channels
- **Open Source:** Contribute to DataHub to understand deeply

14.3 Data Infrastructure Blogs

- DataHub Blog (blog.datahubproject.io)
- Acryl Data Engineering Blog
- LinkedIn Engineering Blog
- Netflix Tech Blog (data infrastructure posts)

15 Final Tips

1. **Study DataHub:** Clone the repo, run locally, understand architecture
2. **Graph Mastery:** DataHub is heavily graph-based - know your graph algorithms!
3. **Scalability Mindset:** Think about millions of datasets from day one
4. **Open Source:** Show familiarity with open-source contribution process
5. **Data Passion:** Demonstrate genuine interest in solving data discovery problems
6. **Ask Questions:** Show curiosity about technical decisions and trade-offs

Good luck with your DataHub/Acryl Data interview!