



Distributed Web Crawling and Indexing System using Cloud Computing phase 4

Team 10

Name	ID
Alaa Yasser Fathy Bekhit	21P0408
Mohamed Ahmed Abdelraouf	2001038
Adham Osama Mohamed	22P0071
Seif Elden Samir Mohamed	2101524

Supervised By: Prof. Ayman Mohamed Bahaa Eldin

Google Drive Folder:

drive.google.com/drive/folders/1Tqau ejRmq6-jAea4jWtBrJKI-em7ndAX

Contents

1	Introduction to Distributed Crawling System	4
2	Detailed Project Description	5
2.1	Overview	5
2.2	Key Features	5
2.3	System Modules	6
3	Beneficiaries of the Project	7
4	Detailed Analysis of Distributed Web Crawler System	8
4.1	Component Analysis	8
4.1.1	Client Layer	8
4.1.2	Master Node (<code>master_node.py</code>)	8
4.1.3	Crawler Nodes (<code>crawler_node.py</code>)	9
4.1.4	Indexer Node (<code>indexer_node.py</code>)	9
4.1.5	Search Interface (<code>search.py</code>)	9
4.1.6	Monitoring System	9
4.2	System Integration and Data Flow	10
4.3	Technical Implementation Details	11
4.3.1	AWS Integration	11
4.3.2	Fault Tolerance	11
4.3.3	Scalability	11
4.3.4	Search Capabilities	11
4.4	Strengths and Potential Improvements	12
4.4.1	Strengths	12
4.4.2	Potential Improvements	12
5	Breakdown of Main Modules	13
5.1	Master Node Overview	13
5.2	Crawler Node	16
5.3	Indexer Node	19
5.4	System Integration	22
6	Gantt chart for the different tasks of the project	25
7	System architecture and design	26
8	Test Scenarios in Distributed Crawler System	27
8.1	Automated Testing via <code>run_tests.bat</code>	27
8.2	Master Node Testing (from <code>run_tests.bat</code>)	27
8.3	Crawler Node Testing (from <code>run_tests.bat</code>)	29
8.4	Indexer Node Testing (from <code>run_tests.bat</code>)	31
8.5	Performance Tests	34
8.6	Fault Tolerance Tests	35
8.7	Scalability Tests	36
8.8	Crawler Node Performance and Memory Optimization	37
8.9	Indexer Node Optimization	38
9	Security Review	40
9.1	A. AWS Resource Security	40
9.2	B. Network-Level Security	42

10 Dashboard API Documentation	45
10.1 Overview	45
10.2 Base URL	45
10.3 Authentication	45
10.4 API Endpoints	45
10.5 Error Responses	48
10.6 Rate Limiting	49
10.7 WebSocket Events	49
11 User Manual	50
12 Deployment Guide	54
12.1 Prerequisites	54
12.2 System Configuration	55
12.3 Deployment Steps	56
12.4 Scaling the System	56
12.5 Monitoring and Maintenance	56
12.6 Troubleshooting	57
12.7 Security Considerations	57
12.8 Performance Tuning	57
13 Conclusion	58

List of Figures

1	Intro to Distributed Systems	4
2	Intro to Distributed Systems	7
3	Search workflow using AWS services: EC2-hosted search service queries DynamoDB and S3, and returns results to the client.	24
4	Intro to Distributed Systems	25
5	Arch	26
6	$\text{master}_t \text{est}_{functionality}$	28
7	$\text{master}_t \text{est}_{scalability}$	28
8	$\text{indexer}_t \text{est}$	32
9	$\text{indexer}_t \text{est}$	33
10	AWS VPC Dashboard – project-vpc configured with CIDR block 10.0.0.0/16 . .	42
11	Security Group Configuration – <code>Crawler_cluster_sg</code> allowing limited SSH . .	43
12	Network-Level Security Diagram: VPC Isolation, Security Groups, and AWS service communication	44

1 Introduction to Distributed Crawling System

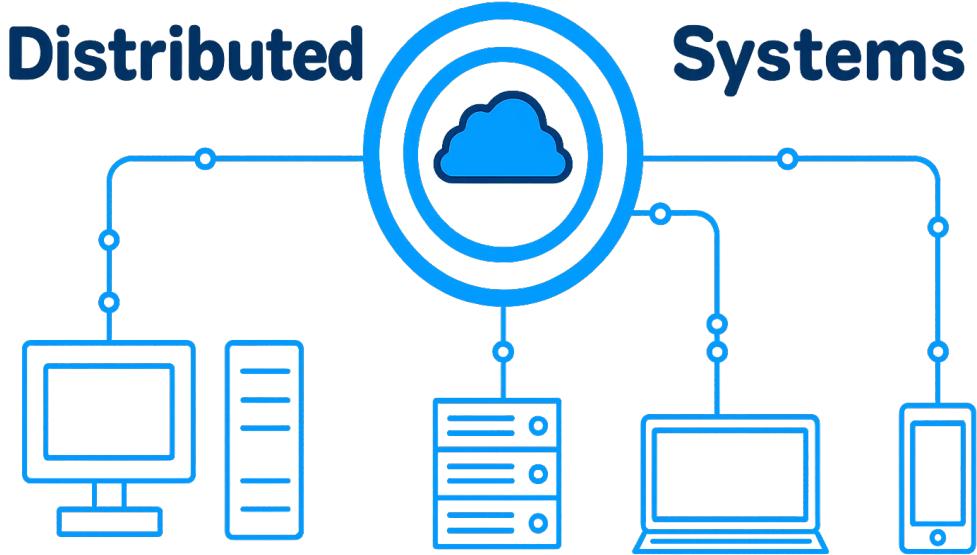


Figure 1: Intro to Distributed Systems

In the era of information overload, efficient and scalable web crawling systems have become essential for collecting, indexing, and retrieving relevant data from the vast landscape of the World Wide Web. This project presents the design and implementation of a distributed web crawling and indexing system built using cloud-native technologies such as AWS SQS, DynamoDB, and S3, combined with Python-based processing and indexing tools. The system is architected to support large-scale, fault-tolerant, and extensible operations through a modular composition of nodes: a centralized Master Node, distributed Crawler Nodes, and an Indexer Nodes.

Designed for scalability and resilience, the system incorporates heartbeat monitoring and dynamic task reassignment to recover from node failures. Furthermore, it supports user interaction via a client web interface for submitting seed URLs and querying indexed content, and a monitoring dashboard for real-time visualization of system performance and crawler status. This architecture lays the groundwork for a robust, cloud-based information retrieval platform adaptable to diverse web data collection needs.

2 Detailed Project Description

2.1 Overview

This project, titled **Distributed Web Crawling and Indexing System using Cloud Computing**, is a robust, scalable solution designed to efficiently collect, store, and index vast amounts of web data. Built using a modular architecture and cloud-native services, it supports distributed task execution, fault tolerance, and interactive user interfaces for seamless web crawling and information retrieval.

The core objective is to design a system capable of crawling and indexing data from the internet in a decentralized and fault-resilient manner. The architecture supports real-time monitoring, intelligent task distribution, and extensible indexing mechanisms to serve as a foundation for large-scale data analytics and search systems.

2.2 Key Features

- **Distributed Architecture:** Modular separation of responsibilities across *Master*, *Crawler*, and *Indexer* nodes allows concurrent execution and system scalability.
- **Cloud Integration:** Utilizes AWS services such as **SQS** for decoupled task queuing, **S3** for persistent storage of raw data and index backups, and **DynamoDB** for real-time metadata management.
- **Indexing and Search:** Employs the **Whoosh** full-text search library to create a powerful and efficient local search index, supporting multi-field queries (e.g., title, description, keywords).
- **Fault Tolerance:** Integrates heartbeat signals and status monitoring to detect crawler node failures, allowing the Master to reassign tasks dynamically.
- **Asynchronous Communication:** Implements non-blocking communication through AWS SQS queues to handle large volumes of tasks and responses without system bottlenecks.
- **User Interaction Interfaces:**
 - **Client Web Interface** for submitting crawl requests and configuring crawl depth.
 - **Search Interface** to perform real-time queries over indexed content.
 - **Monitoring Dashboard** for live updates on system health, crawler progress, and data flow.
- **Content Filtering and Processing:** HTML content is sanitized and parsed for relevant text extraction before indexing, reducing noise and improving search quality.
- **Scalability:** New crawler nodes can be added without modifying existing components, enabling horizontal scaling as crawl demand increases.
- **Persistent Storage and Backup:** Raw crawl data and indexes are securely archived to S3 for future retrieval and disaster recovery.

2.3 System Modules

- **Master Node:** Central coordinator responsible for managing the URL frontier, distributing tasks, processing results, and handling system metadata.
- **Crawler Nodes:** Execute crawl tasks by downloading web pages, extracting content and links, and reporting results/status back to the Master node.
- **Indexer Node:** Processes crawled content into a structured, searchable format using Whoosh. Also performs backup of index files to cloud storage.
- **Client Web Interface:** Web-based frontend for users to initiate crawl jobs and specify parameters.
- **Search Interface:** Provides search functionality over the indexed documents, returning high-relevance results with titles, URLs, and descriptions.
- **Monitoring Dashboard:** Displays live system metrics including active nodes, processed URLs, system throughput, and crawler health status.

3 Beneficiaries of the Project

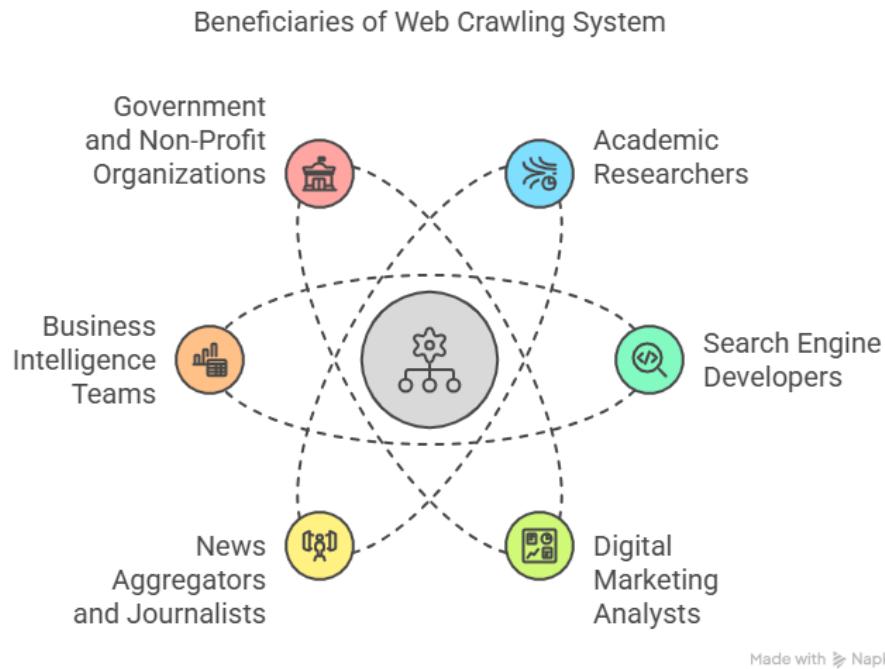


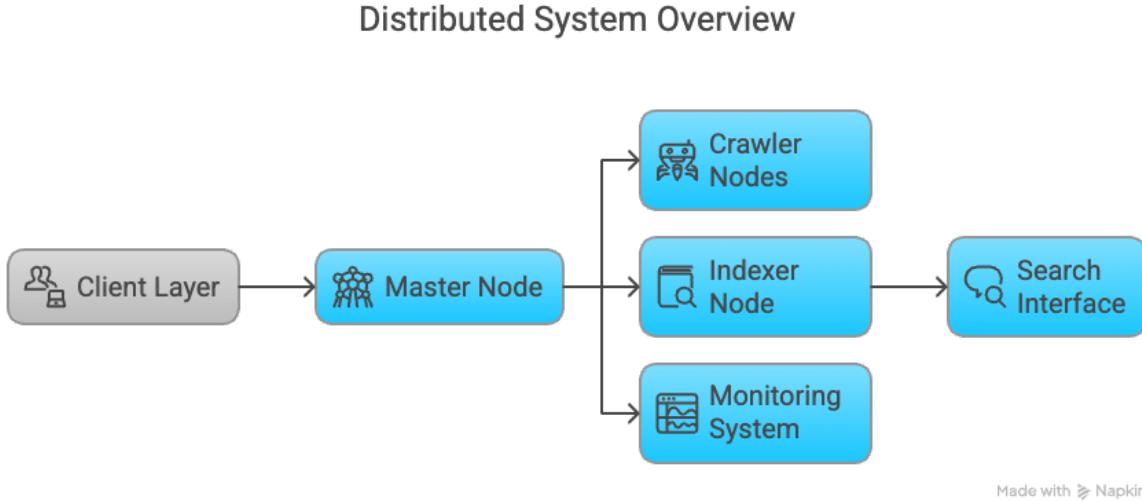
Figure 2: Intro to Distributed Systems

The Distributed Web Crawling and Indexing System is designed to serve a wide range of stakeholders who require access to large-scale, structured web data. Below are the primary beneficiaries of this project:

- **Academic Researchers:** Researchers in fields such as data mining, information retrieval, natural language processing, and web science can use the system to collect and analyze datasets at scale for experimentation and knowledge discovery.
- **Search Engine Developers:** Developers working on domain-specific or enterprise-level search engines can leverage the indexing capabilities to create efficient, focused search experiences.
- **Digital Marketing Analysts:** Marketing teams can track changes in web content, perform competitor analysis, or extract SEO-relevant metadata across websites to optimize their strategies.
- **News Aggregators and Journalists:** Journalists and media analysts can continuously gather fresh web content from various sources for fact-checking, sentiment analysis, and content aggregation.
- **Business Intelligence Teams:** Corporations can use the system to monitor online trends, product mentions, or customer sentiment across public domains.
- **Government and Non-Profit Organizations:** These entities can track public information, monitor policy-related discussions, or analyze social trends to support public service and regulatory efforts.

4 Detailed Analysis of Distributed Web Crawler System

4.1 Component Analysis



4.1.1 Client Layer

The client layer consists of two main components:

1. `client.py`: Provides a high-level interface for initiating crawl jobs
 - Uses SQS queues to communicate with the master node
 - Supports configuration of crawl parameters (max depth, URLs per domain)
 - Command-line interface for easy integration
2. `submit_url.py`: A simplified client for submitting individual URLs
 - Lightweight interface for adding URLs to the crawl queue
 - Useful for incremental additions to the crawl process

4.1.2 Master Node (`master_node.py`)

The master node serves as the central coordinator for the entire system:

- **Task Distribution**: Manages URL queues and distributes tasks to crawler nodes
- **Worker Management**: Monitors crawler node health and handles node failures
- **URL Tracking**: Maintains a database of crawled URLs to prevent duplicates
- **AWS Integration**: Uses SQS for queue management, DynamoDB for state persistence, and S3 for data storage
- **Fault Tolerance**: Implements heartbeat mechanisms and task reassignment for failed nodes

4.1.3 Crawler Nodes (`crawler_node.py`)

Crawler nodes are the workhorses of the system:

- **Web Fetching:** Uses Scrapy for efficient web page retrieval
- **Content Processing:** Extracts links, text, and metadata from web pages
- **Task Management:** Pulls tasks from SQS queues and reports results
- **Fault Tolerance:** Implements retry mechanisms and error handling
- **AWS Integration:** Communicates with SQS, S3, and DynamoDB for task management and data storage

4.1.4 Indexer Node (`indexer_node.py`)

The indexer processes crawled content and builds a searchable index:

- **Text Processing:** Uses NLTK for advanced text processing (tokenization, stemming, stopword removal)
- **Index Building:** Employs Whoosh for creating a full-text search index
- **Metadata Extraction:** Extracts keywords, descriptions, and other metadata
- **Schema Definition:** Defines a comprehensive schema for indexed documents
- **AWS Integration:** Pulls data from S3 and processes it into a searchable index

4.1.5 Search Interface (`search.py`)

Provides interfaces for searching the indexed content:

- **Command-line Interface:** For quick searches from the terminal
- **Web Interface:** Flask-based web application for user-friendly searching
- **Result Formatting:** Formats search results with highlighting and relevance scores
- **API Endpoints:** RESTful API for programmatic access to search functionality

4.1.6 Monitoring System

The monitoring system consists of two main components:

1. `dashboard.py`: Enhanced web-based dashboard
 - Flask-based web application for system monitoring
 - Real-time visualization of system status
 - Crawler node status tracking
 - Performance metrics and error reporting
 - URL submission interface

4.2 System Integration and Data Flow

1. Crawl Initiation:

- Client submits seed URLs to the master node via SQS
- Master node validates URLs and adds them to the crawl queue

2. Task Distribution:

- Crawler nodes pull tasks from the crawl queue
- Master node monitors task distribution and crawler health

3. Web Crawling:

- Crawler nodes fetch web pages and extract content
- New URLs are discovered and sent back to the master node
- Crawled content is stored in S3

4. Indexing:

- Indexer node pulls crawled content from S3
- Content is processed and added to the search index
- Metadata is extracted and stored for enhanced search

5. Searching:

- Users query the search interface
- Search engine retrieves relevant documents from the index
- Results are formatted and presented to the user

6. Monitoring:

- System continuously collects metrics from all components
- Dashboard displays real-time system status
- Alerts can be triggered for system issues

4.3 Technical Implementation Details

4.3.1 AWS Integration

The system heavily leverages AWS services:

- **SQS**: For task queues and inter-component communication
- **DynamoDB**: For state persistence and URL tracking
- **S3**: For storing crawled content and system data
- **CloudWatch**: For monitoring and logging (implied)

4.3.2 Fault Tolerance

The system implements several fault tolerance mechanisms:

- **Heartbeat System**: Crawler nodes send regular heartbeats to the master
- **Task Reassignment**: Failed tasks are reassigned to healthy nodes
- **Error Handling**: Comprehensive error handling and retry mechanisms
- **State Persistence**: System state is persisted in DynamoDB for recovery

4.3.3 Scalability

The architecture supports horizontal scaling:

- **Stateless Crawler Nodes**: Can be added or removed dynamically
- **Queue-Based Communication**: Decouples components for independent scaling
- **Distributed Processing**: Work is distributed across multiple nodes

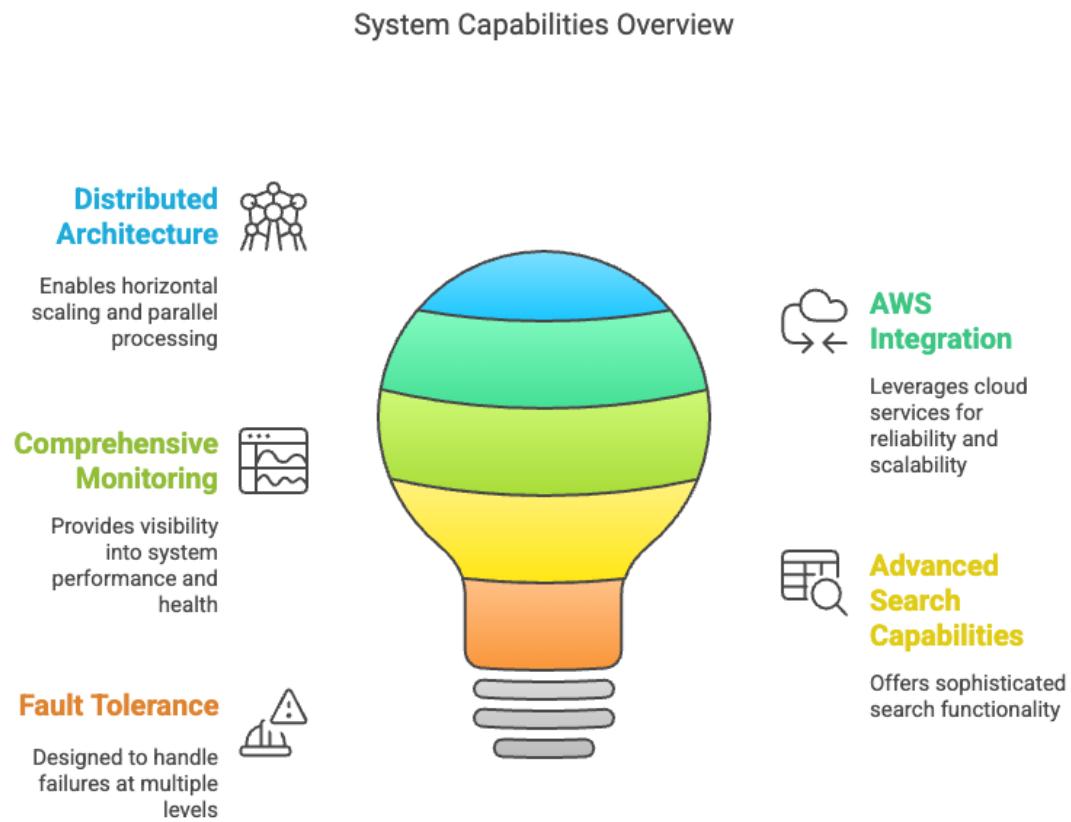
4.3.4 Search Capabilities

The search system provides advanced features:

- **Full-Text Search**: Using Whoosh's indexing and query capabilities
- **Relevance Ranking**: BM25F scoring for result relevance
- **Text Processing**: NLTK integration for stemming and stopword removal
- **Result Highlighting**: Context fragments with query term highlighting

4.4 Strengths and Potential Improvements

4.4.1 Strengths



Made with Napkin

1. **Distributed Architecture:** Enables scalable and fault-tolerant operation
2. **AWS Integration:** Leverages cloud services for reliability and scalability
3. **Advanced Text Processing:** NLTK integration for sophisticated indexing
4. **Comprehensive Monitoring:** Detailed metrics and visualization
5. **Multiple Interfaces:** Both CLI and web interfaces for different use cases

4.4.2 Potential Improvements

1. **Container Support:** Adding Docker/Kubernetes support for easier deployment

5 Breakdown of Main Modules

5.1 Master Node Overview

The Master Node is the central coordinator of the distributed web crawling system. It manages the crawling process, coordinates worker nodes, and maintains global system state to ensure reliable and scalable operations.

Key Responsibilities

- **Task Distribution:** Manages the URL frontier and distributes crawl tasks to crawler nodes.
- **Worker Management:** Tracks active crawler nodes, monitors their health, and handles node failures.
- **State Management:** Maintains crawl state persistently using DynamoDB and S3.
- **Fault Tolerance:** Implements recovery mechanisms and task retry logic to ensure robustness.

Architecture

- Utilizes **AWS SQS** for task queueing and inter-node communication.
- Employs **DynamoDB** for URL tracking, node status, and job metadata.
- Stores crawled content and snapshots in **AWS S3**.
- Includes a heartbeat mechanism to monitor the health and availability of worker nodes.

Key Components

- **MasterNode:** Main class handling coordination and orchestration.
- **URL Frontier:** Manages pending and visited URLs.
- **Task Queue:** Publishes crawl tasks for worker nodes.
- **Result Queue:** Collects and processes crawl results.
- **Heartbeat Monitor:** Tracks periodic heartbeats to detect node failures.

Configuration

- **CRAWL_TASK_QUEUE:** SQS queue name for crawl task distribution.
- **CRAWL_RESULT_QUEUE:** SQS queue name for collecting crawl results.
- **HEARTBEAT_INTERVAL:** Time interval (in seconds) for node health checks.
- **AWS_REGION:** The AWS region where all services are deployed.

Security

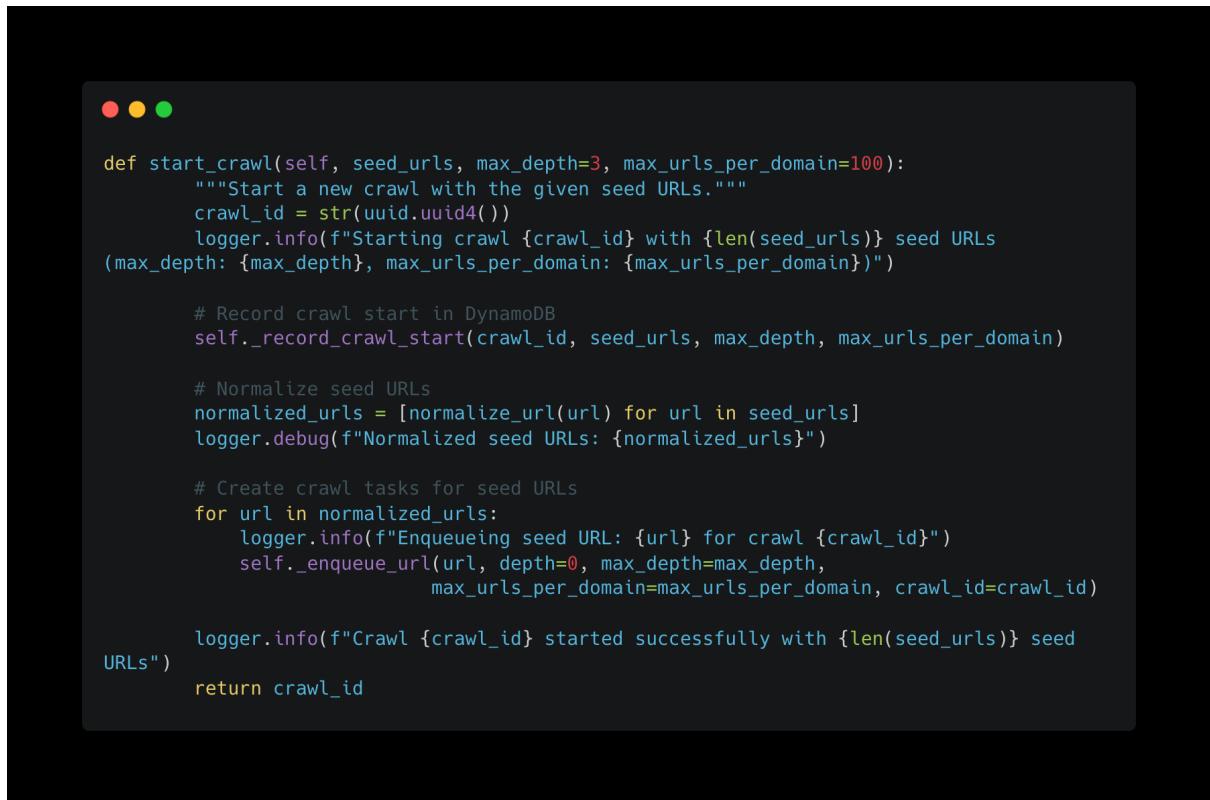
- Implements **AWS IAM roles** for secure service access.
- Uses **SQS message encryption** and access control.
- Applies **URL sanitization** to prevent injection attacks.
- Maintains **audit logs** to track system operations and access.

Performance Optimizations

- Uses batch operations for efficient **DynamoDB writes**.
- Implements in-memory **caching** for frequently accessed metadata.
- Enables **connection pooling** for external services.
- Applies intelligent **task distribution** strategies to balance load across nodes.

master Initialize Crawl

- `start_crawl()`: Initiates a new crawl job with seed URLs



```
def start_crawl(self, seed_urls, max_depth=3, max_urls_per_domain=100):
    """Start a new crawl with the given seed URLs."""
    crawl_id = str(uuid.uuid4())
    logger.info(f"Starting crawl {crawl_id} with {len(seed_urls)} seed URLs
(max_depth: {max_depth}, max_urls_per_domain: {max_urls_per_domain})")

    # Record crawl start in DynamoDB
    self._record_crawl_start(crawl_id, seed_urls, max_depth, max_urls_per_domain)

    # Normalize seed URLs
    normalized_urls = [normalize_url(url) for url in seed_urls]
    logger.debug(f"Normalized seed URLs: {normalized_urls}")

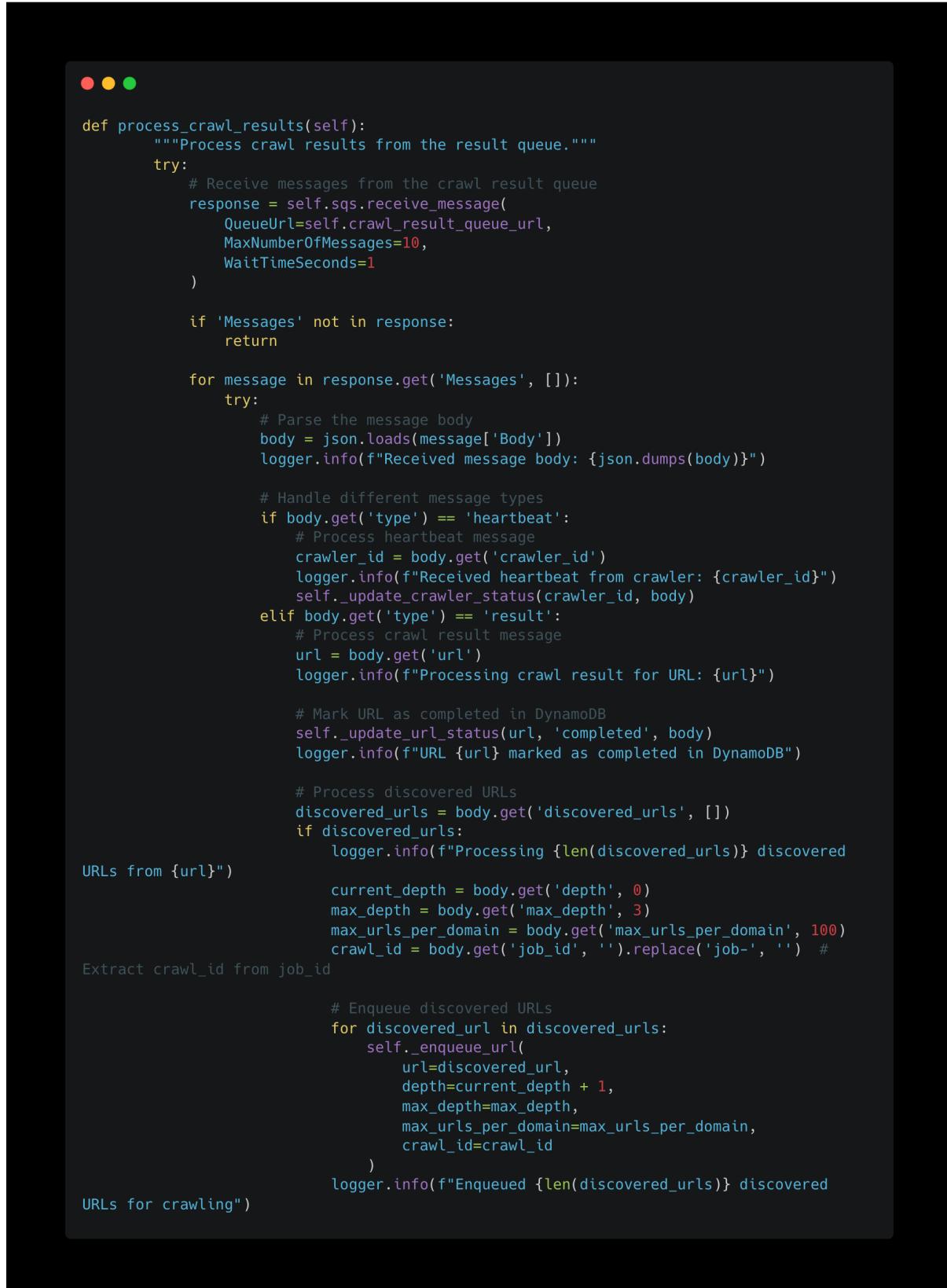
    # Create crawl tasks for seed URLs
    for url in normalized_urls:
        logger.info(f"Enqueueing seed URL: {url} for crawl {crawl_id}")
        self._enqueue_url(url, depth=0, max_depth=max_depth,
                          max_urls_per_domain=max_urls_per_domain, crawl_id=crawl_id)

    logger.info(f"Crawl {crawl_id} started successfully with {len(seed_urls)} seed
URLs")
    return crawl_id
```

- `_enqueue_url()`: Adds URLs to the frontier and task queue
- `_record_crawl_start()`: Records crawl job information in DynamoDB

Result Processing

- `process_crawl_results()`: Processes results from crawler nodes



```
def process_crawl_results(self):
    """Process crawl results from the result queue."""
    try:
        # Receive messages from the crawl result queue
        response = self.sqs.receive_message(
            QueueUrl=self.crawl_result_queue_url,
            MaxNumberOfMessages=10,
            WaitTimeSeconds=1
        )

        if 'Messages' not in response:
            return

        for message in response.get('Messages', []):
            try:
                # Parse the message body
                body = json.loads(message['Body'])
                logger.info(f"Received message body: {json.dumps(body)}")

                # Handle different message types
                if body.get('type') == 'heartbeat':
                    # Process heartbeat message
                    crawler_id = body.get('crawler_id')
                    logger.info(f"Received heartbeat from crawler: {crawler_id}")
                    self._update_crawler_status(crawler_id, body)
                elif body.get('type') == 'result':
                    # Process crawl result message
                    url = body.get('url')
                    logger.info(f"Processing crawl result for URL: {url}")

                    # Mark URL as completed in DynamoDB
                    self._update_url_status(url, 'completed', body)
                    logger.info(f"URL {url} marked as completed in DynamoDB")

                    # Process discovered URLs
                    discovered_urls = body.get('discovered_urls', [])
                    if discovered_urls:
                        logger.info(f"Processing {len(discovered_urls)} discovered
URLs from {url}")

                        current_depth = body.get('depth', 0)
                        max_depth = body.get('max_depth', 3)
                        max_urls_per_domain = body.get('max_urls_per_domain', 100)
                        crawl_id = body.get('job_id', '').replace('job-', '') # Extract crawl_id from job_id

                        # Enqueue discovered URLs
                        for discovered_url in discovered_urls:
                            self._enqueue_url(
                                url=discovered_url,
                                depth=current_depth + 1,
                                max_depth=max_depth,
                                max_urls_per_domain=max_urls_per_domain,
                                crawl_id=crawl_id
                            )
                        logger.info(f"Enqueued {len(discovered_urls)} discovered
URLs for crawling")

```

5.2 Crawler Node

The Crawler Node is responsible for fetching and processing web pages in the distributed web crawling system. It receives crawl tasks from the Master Node, retrieves content using Scrapy, processes the response, and reports the results back. It also monitors its health and implements recovery mechanisms to ensure resilience.

Key Responsibilities

- **Web Crawling:** Fetches web pages using the Scrapy framework.
- **Content Processing:** Parses and extracts page content such as text, links, and metadata.
- **Result Reporting:** Sends structured crawl results to the Master Node via the `CRAWL_RESULT_QUEUE`.
- **Health Monitoring:** Periodically sends heartbeat signals to confirm liveness.

Architecture

- Leverages **Scrapy** for scalable and extensible crawling logic.
- Uses **custom middleware** to enforce politeness policies and rate limiting.
- Communicates with the Master Node using **AWS SQS** for task reception and result submission.
- Supports **fault tolerance** with retry logic and recovery from partial failures.

Key Components

- `CrawlerNode`: Main orchestration class responsible for launching the crawling loop.
- `Scrapy Spider`: Custom spider to fetch web pages and extract content.
- `Task Processor`: Receives and handles crawl tasks from the task queue.
- `Result Reporter`: Sends extracted data to the result queue for indexing.
- `Health Monitor`: Periodically emits heartbeats to the master for node health tracking.

Configuration

- `CRAWL_TASK_QUEUE`: Name of the SQS queue used to receive crawl tasks.
- `CRAWL_RESULT_QUEUE`: SQS queue for submitting crawl results.
- `HEARTBEAT_INTERVAL`: Frequency (in seconds) for health updates.
- `MAX_RETRIES`: Maximum number of retry attempts per task.
- `TASK_TIMEOUT`: Timeout duration for a single crawl task.

Crawling Settings

- `USER_AGENT`: Custom user agent string to emulate browser behavior.
- `ROBOTSTXT_OBEY`: Boolean setting to respect `robots.txt`.
- `DOWNLOAD_DELAY`: Delay between consecutive requests to the same domain.
- `CONCURRENT_REQUESTS`: Maximum number of parallel requests.
- `RETRY_SETTINGS`: Retry configuration (e.g., max attempts, delay strategy).

Security

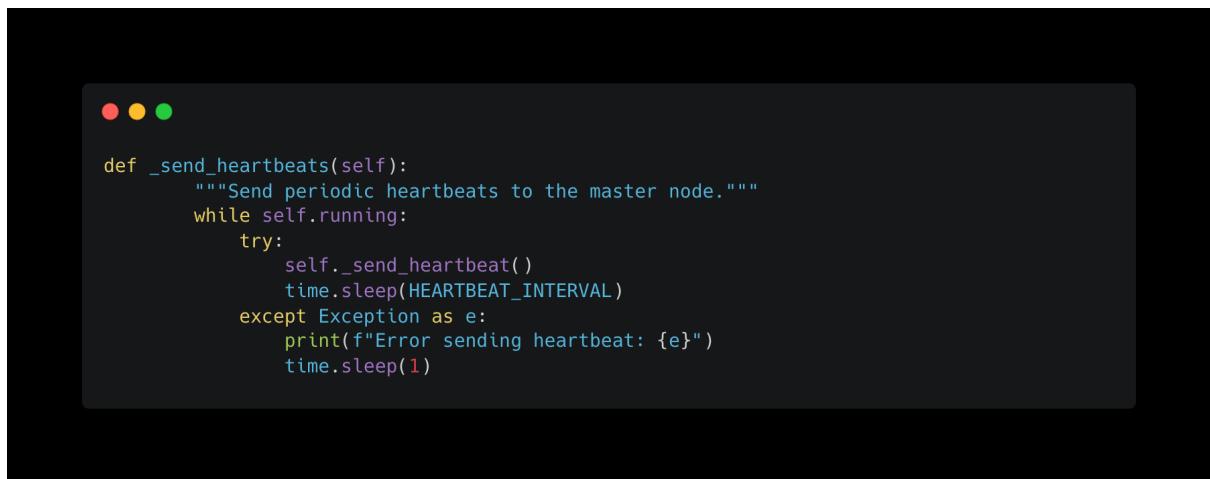
- Performs **URL sanitization** to eliminate malicious or malformed inputs.
- Complies with `robots.txt` directives to respect crawling boundaries.
- Applies **rate limiting** to avoid overloading target servers.
- Uses **secure AWS IAM credentials** for all SQS and monitoring operations.

Performance Optimizations

- Enables **connection pooling** to reuse network sessions efficiently.
- Uses **efficient HTML parsers** for faster content extraction.
- Implements **memory management** practices to avoid resource leaks.
- Optimizes request scheduling to maximize throughput and minimize latency.

Fault Tolerance

- `_send_heartbeats()`: Periodically sends heartbeats to the master

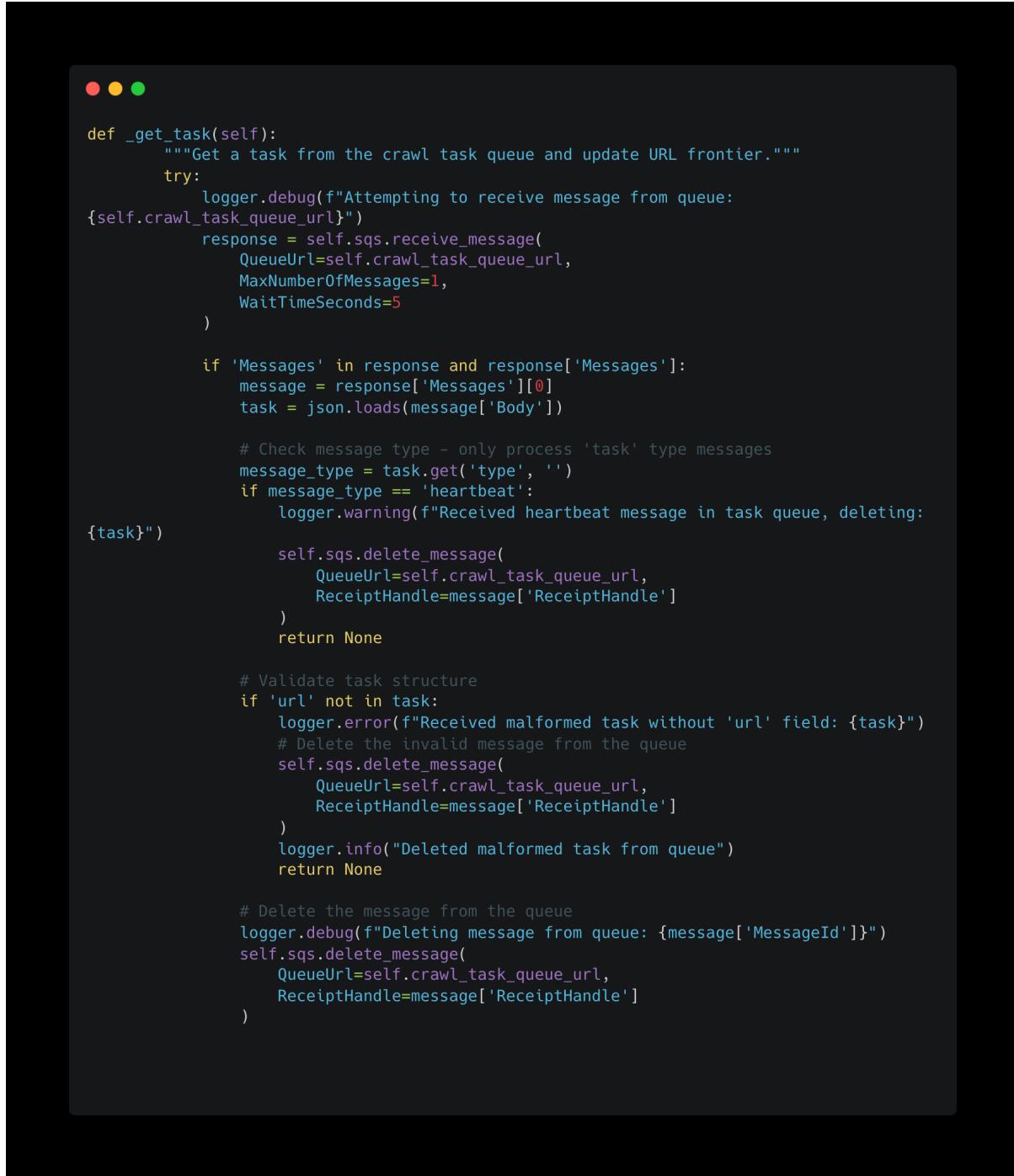


```
def _send_heartbeats(self):
    """Send periodic heartbeats to the master node."""
    while self.running:
        try:
            self._send_heartbeat()
            time.sleep(HEARTBEAT_INTERVAL)
        except Exception as e:
            print(f"Error sending heartbeat: {e}")
            time.sleep(1)
```

- `_recovery_loop()`: Handles task timeouts and retries
- Tracks failed tasks and implements retry logic

Task Management

- `_get_task()`: Retrieves tasks from the SQS queue



```
def _get_task(self):
    """Get a task from the crawl task queue and update URL frontier."""
    try:
        logger.debug(f"Attempting to receive message from queue: {self.crawl_task_queue_url}")
        response = self.sqs.receive_message(
            QueueUrl=self.crawl_task_queue_url,
            MaxNumberOfMessages=1,
            WaitTimeSeconds=5
        )

        if 'Messages' in response and response['Messages']:
            message = response['Messages'][0]
            task = json.loads(message['Body'])

            # Check message type - only process 'task' type messages
            message_type = task.get('type', '')
            if message_type == 'heartbeat':
                logger.warning(f"Received heartbeat message in task queue, deleting: {task}")

            self.sqs.delete_message(
                QueueUrl=self.crawl_task_queue_url,
                ReceiptHandle=message['ReceiptHandle']
            )
            return None

        # Validate task structure
        if 'url' not in task:
            logger.error(f"Received malformed task without 'url' field: {task}")
            # Delete the invalid message from the queue
            self.sqs.delete_message(
                QueueUrl=self.crawl_task_queue_url,
                ReceiptHandle=message['ReceiptHandle']
            )
            logger.info("Deleted malformed task from queue")
            return None

        # Delete the message from the queue
        logger.debug(f"Deleting message from queue: {message['MessageId']}")
        self.sqs.delete_message(
            QueueUrl=self.crawl_task_queue_url,
            ReceiptHandle=message['ReceiptHandle']
        )
    
```

- `_process_task()`: Processes a crawl task using Scrapy
- Stores active tasks and their start times

Content Processing

- Extracts metadata (title, description, keywords)
- Discovers and normalizes links for further crawling

5.3 Indexer Node

The Indexer Node is responsible for processing crawled content and building a searchable index in the distributed web crawling system. It manages document parsing, indexing operations, and search query handling, ensuring that data is efficiently indexed and retrievable.

Key Responsibilities

- **Document Processing:** Parses and preprocesses crawled content for indexing.
- **Index Management:** Builds and updates a full-text search index.
- **Search Operations:** Responds to search queries submitted by users or system components.
- **Index Persistence:** Periodically saves and retrieves index files to and from AWS S3.

Architecture

- Utilizes **Whoosh** for local full-text indexing with support for advanced scoring (e.g., BM25F).
- Applies **NLTK** for tokenization, stemming, and text normalization.
- Stores index files persistently in **AWS S3** to support distributed and fault-tolerant indexing.
- Supports scalable **distributed indexing** through batched task processing and result streaming.

Key Components

- **IndexerNode:** Main class orchestrating indexing tasks and result submission.
- **DocumentProcessor:** Prepares crawled documents for indexing (tokenization, filtering).
- **IndexManager:** Manages index creation, updates, and schema configuration.
- **SearchEngine:** Processes search queries and retrieves ranked results.
- **IndexPersistence:** Handles loading and saving index files to/from S3.

Configuration

- **INDEX_TASK_QUEUE:** SQS queue for receiving documents to index.
- **INDEX_RESULT_QUEUE:** Queue for returning indexed result metadata.
- **HEARTBEAT_INTERVAL:** Interval for sending health signals to the master.
- **INDEX_BATCH_SIZE:** Number of documents to process in each indexing batch.
- **S3_BUCKET:** Target S3 bucket for storing the persistent index.

Indexing Settings

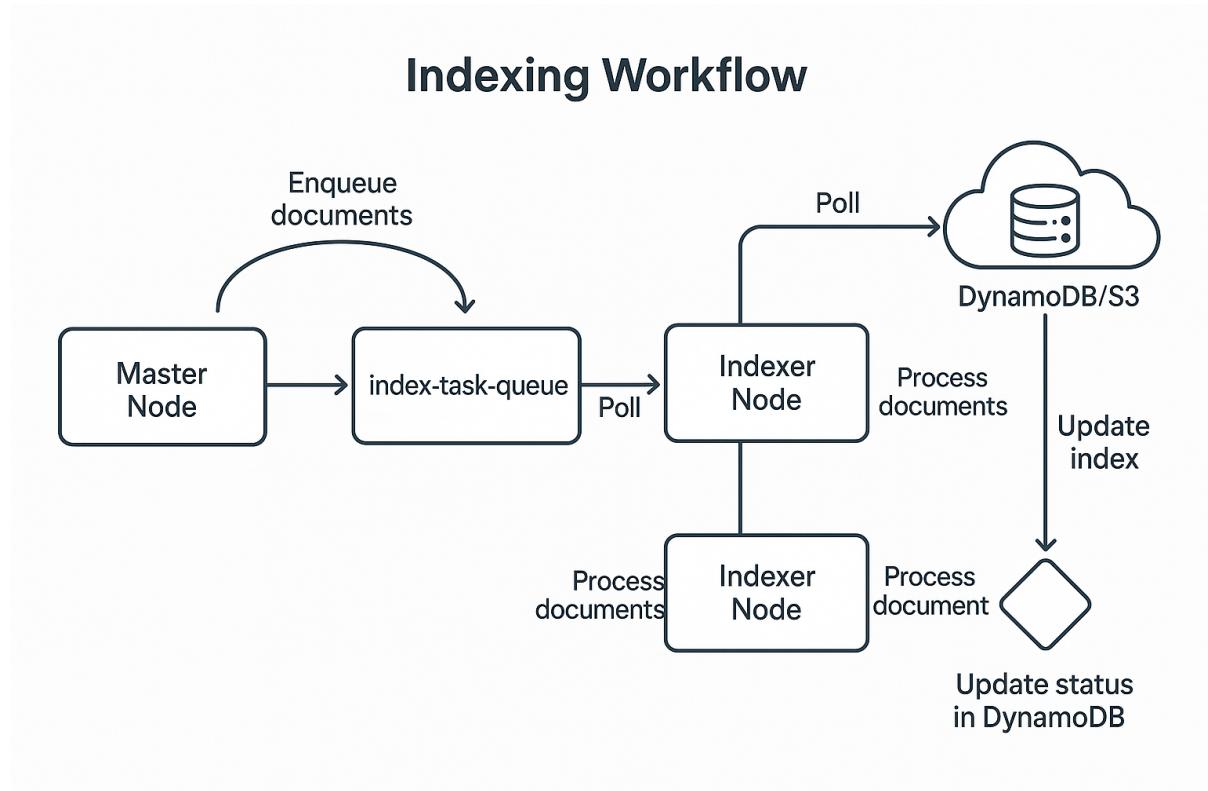
- SCHEMA: Defines the structure and searchable fields of the Whoosh index.
- BM25F: Uses BM25F ranking for improved retrieval relevance.
- NLTK: Configured for English stopword removal, stemming, and normalization.
- CACHING: Controls whether to cache parts of the index in memory.

Security

- Encrypts index files before uploading to S3.
- Uses secure **AWS IAM credentials** for accessing AWS resources.
- Applies access control to restrict unauthorized indexing or search.
- Maintains audit logs for tracking index updates and queries.

Performance Optimizations

- Uses batch processing to improve throughput and reduce write contention.
- Applies efficient NLTK-based text preprocessing for large-scale input.
- Periodically runs index optimization (e.g., merging segments) to enhance retrieval speed.
- Implements multithreaded processing for concurrent document indexing.



Whoosh Index

```
class WhooshIndex:  
    """A more robust index using Whoosh library with enhanced features."""  
    def __init__(self, index_dir='whoosh_index'):  
        # Create index directory if it doesn't exist  
        self.index_dir = index_dir  
        os.makedirs(self.index_dir, exist_ok=True)  
        logger.info(f"Index directory initialized: {self.index_dir}")  
  
        # Initialize text processor  
        self.text_processor = EnhancedTextProcessor()  
  
        # Define schema with enhanced fields  
        self.schema = Schema(  
            url=ID(stored=True, unique=True),  
            title=TEXT(stored=True, analyzer=StemmingAnalyzer()),  
            description=TEXT(stored=True, analyzer=StemmingAnalyzer()),  
            content=TEXT(analyzer=StemmingAnalyzer(), stored=True),  
            keywords=TEXT(stored=True),  
            domain=STORED,  
            crawl_time=DATETIME(stored=True),  
            processed_text=TEXT(analyzer=StemmingAnalyzer(), stored=True),  
            metadata=STORED  
        )  
        logger.debug("Schema defined with enhanced fields")  
  
    # Create or open the index
```

- **WhooshIndex:** Manages the search index
 - Schema definition with multiple fields
 - Document addition and updating
 - Advanced search with query parsing and highlighting

Content Extraction

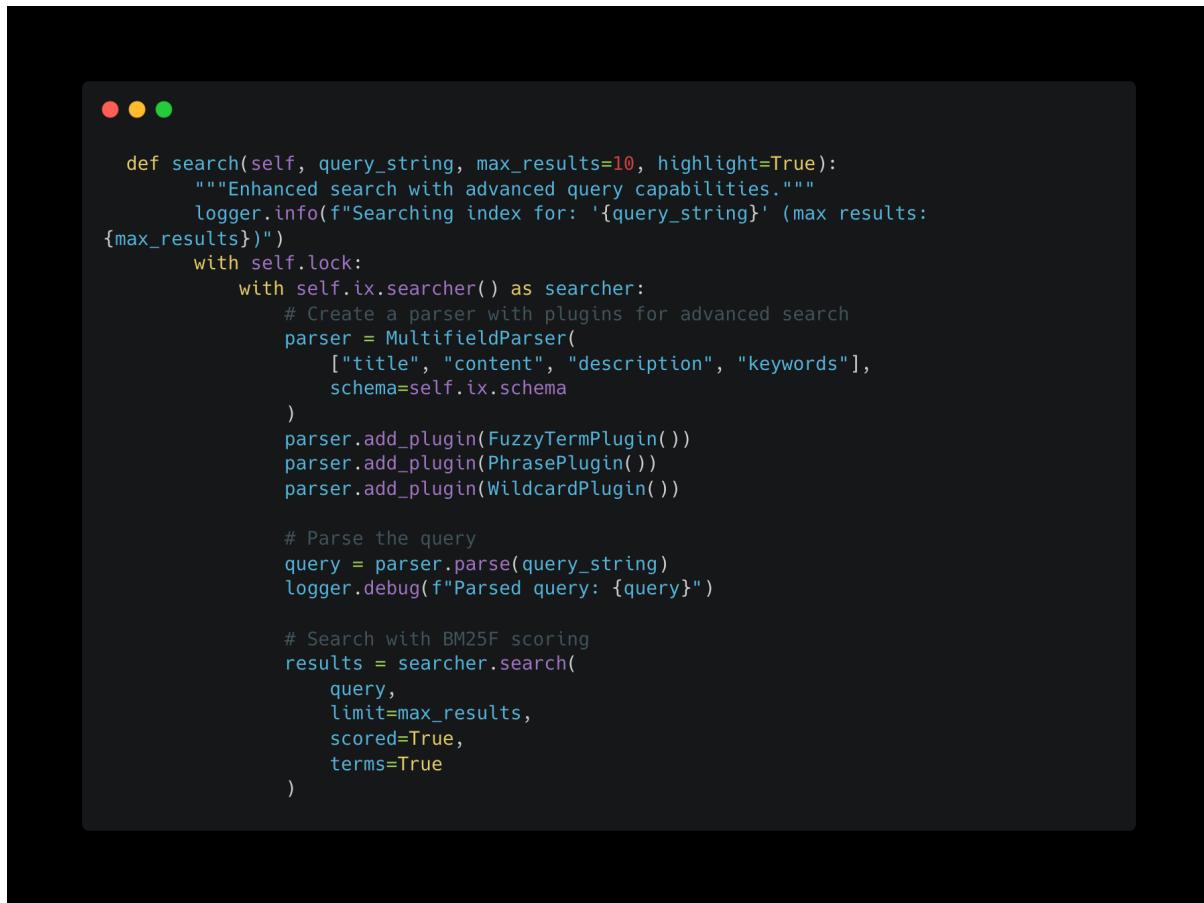
- `_extract_text_from_html()`: Cleans HTML content
- `_extract_field()`: Normalizes and processes specific fields
- Handles various content formats and structures

Index Management

- `save_to_s3()`: Persists the index to S3
- `load_from_s3()`: Loads the index from S3
- `check_index_integrity()`: Verifies index integrity

Search Capabilities

- `search()`: Performs searches with advanced query capabilities
- Supports fuzzy matching, wildcards, and phrase queries
- Provides result highlighting and scoring



A screenshot of a code editor showing a Python function named `search`. The code uses Elasticsearch's `MultifieldParser` to parse a query string, adds several search plugins (FuzzyTermPlugin, PhrasePlugin, WildcardPlugin), and then performs a search with BM25F scoring, limiting results to 10 and including highlights.

```
def search(self, query_string, max_results=10, highlight=True):
    """Enhanced search with advanced query capabilities."""
    logger.info(f"Searching index for: '{query_string}' (max results: {max_results})")
    with self.lock:
        with self.ix.searcher() as searcher:
            # Create a parser with plugins for advanced search
            parser = MultifieldParser(
                ["title", "content", "description", "keywords"],
                schema=self.ix.schema
            )
            parser.add_plugin(FuzzyTermPlugin())
            parser.add_plugin(PhrasePlugin())
            parser.add_plugin(WildcardPlugin())

            # Parse the query
            query = parser.parse(query_string)
            logger.debug(f"Parsed query: {query}")

            # Search with BM25F scoring
            results = searcher.search(
                query,
                limit=max_results,
                scored=True,
                terms=True
            )
```

5.4 System Integration

The three modules work together in a coordinated workflow:

1. Master Node initiates crawls and manages the URL frontier
2. Crawler Nodes fetch and process web pages, sending results back to the Master
3. Master Node forwards processed content to the Indexer Node
4. Indexer Node builds and maintains the search index

Communication Flow

- **client → master:** URL tasks via `command_master_QUEUE`
- **Master → Crawler:** URL tasks via `CRAWL_TASK_QUEUE`
- **Crawler → Master:** Crawl results via `CRAWL_RESULT_QUEUE`
- **Master → Indexer:** Index tasks via `INDEX_TASK_QUEUE`

Data Storage

- Crawled content: Stored in S3 (`CRAWL_DATA_BUCKET`)
- Search index: Stored in S3 (`INDEX_DATA_BUCKET`)
- Metadata and status: Stored in DynamoDB tables

Fault Tolerance Features

- Heartbeat mechanism for node monitoring
- Task timeout and retry mechanisms
- Persistent state in DynamoDB for recovery
- Index backup and restoration from S3

Search Workflow

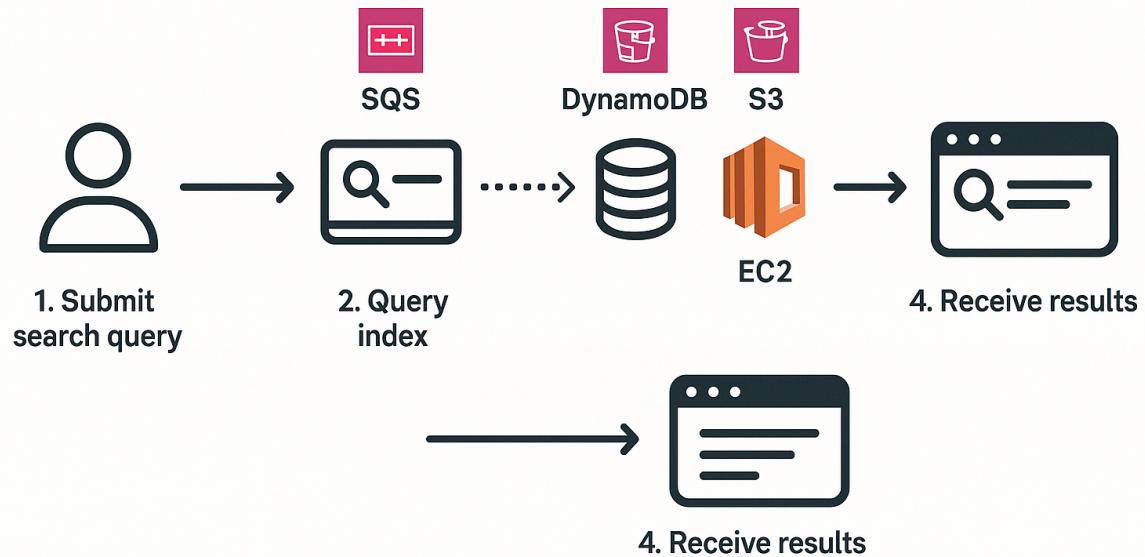


Figure 3: Search workflow using AWS services: EC2-hosted search service queries DynamoDB and S3, and returns results to the client.

2. Search Workflow

Figure 3 illustrates the search workflow in our distributed system. The service performs the following steps:

- Queries **DynamoDB** to locate relevant metadata or document references based on the search input.
- Fetches full content from **Amazon S3** for documents that match the query terms.
- Ranks the documents using relevance scoring techniques
- Returns a sorted list of results to the user interface for display.

This workflow is designed to be scalable and efficient, leveraging AWS services to ensure low-latency search and high availability.

6 Gantt chart for the different tasks of the project

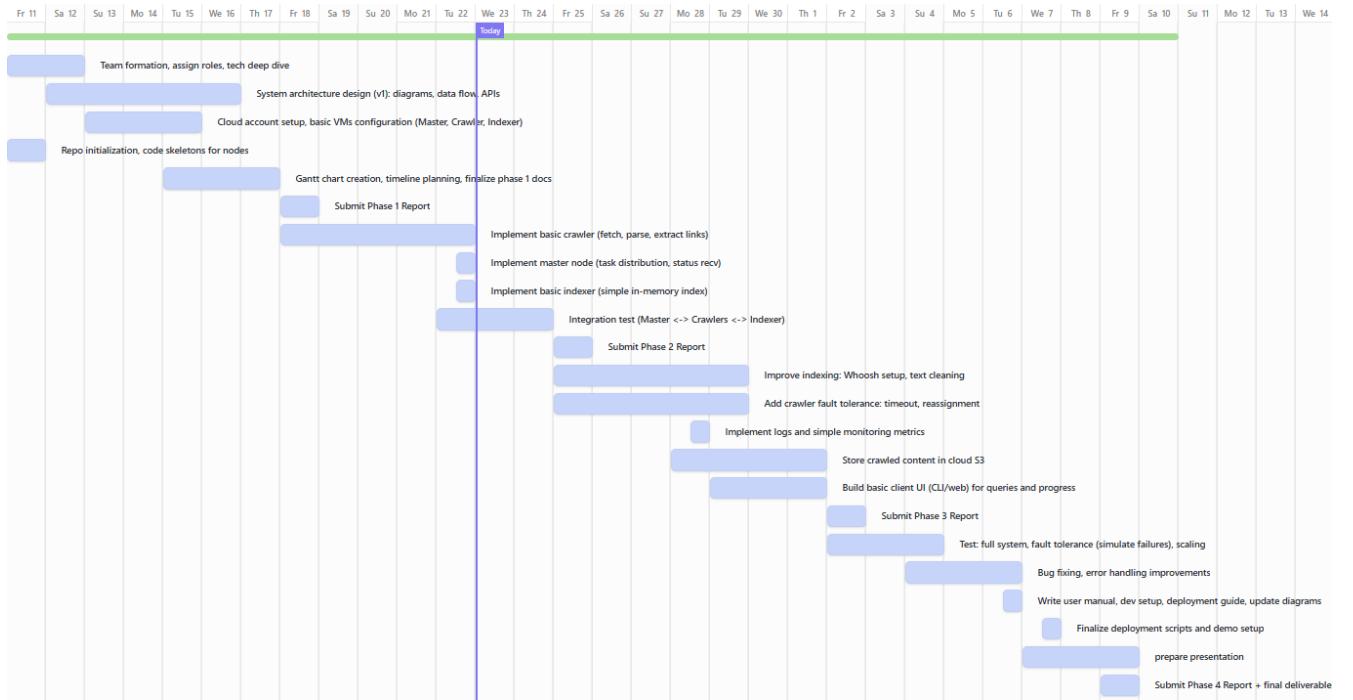


Figure 4: Intro to Distributed Systems

7 System architecture and design

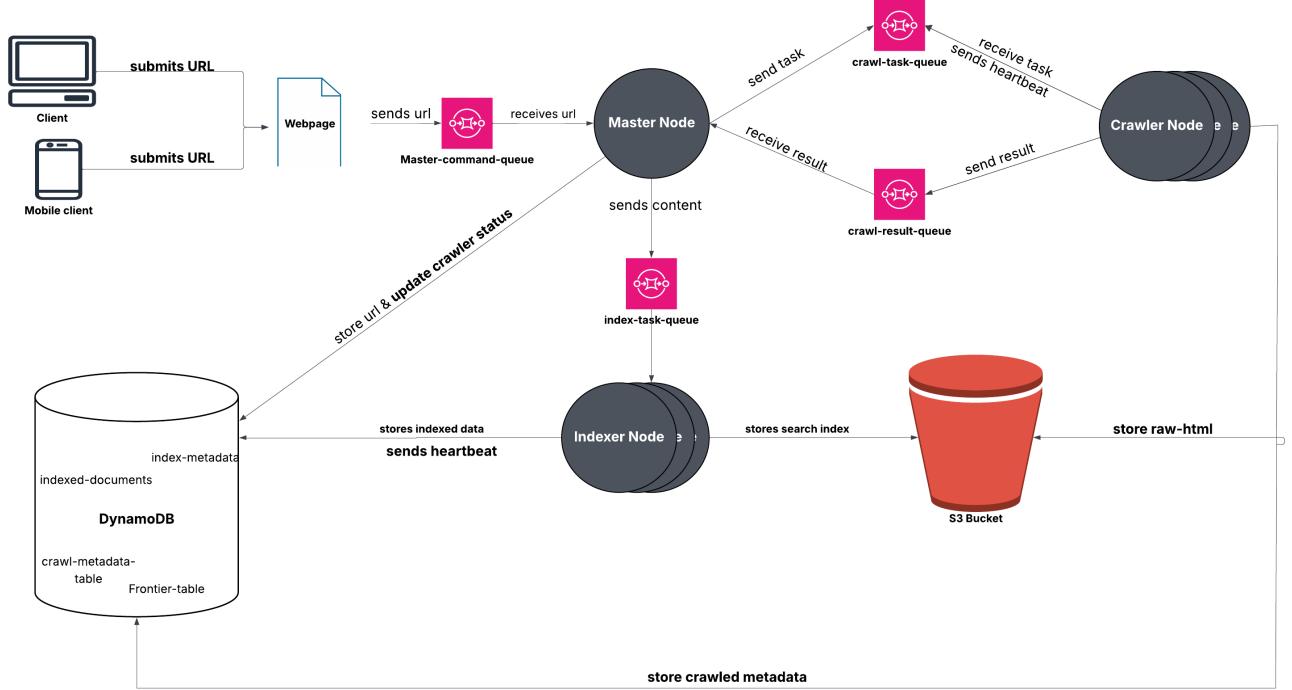


Figure 5: Arch

The distributed crawler system consists of several interconnected components:

1. **Client Layer**: Provides interfaces for submitting crawl requests
2. **Master Node**: Orchestrates the crawling process and distributes tasks
3. **Crawler Nodes**: Perform the actual web page fetching and processing
4. **Indexer Node**: Processes crawled content and builds a searchable index
5. **Search Interface**: Allows users to search the indexed content
6. **Monitoring System**: Tracks system health and performance

8 Test Scenarios in Distributed Crawler System

8.1 Automated Testing via `run_tests.bat`

To streamline validation across multiple system components, we developed an automated script named `run_tests.bat`. This script executes a sequence of tests covering the core functional, fault-tolerance, and scalability scenarios.

8.2 Master Node Testing (from `run_tests.bat`)

The `run_tests.bat` script was executed to validate the functionality, fault tolerance, and scalability of the Master Node. It tested the node's ability to coordinate the crawling pipeline, manage fault recovery, and scale with increasing worker nodes.

Tested Functionalities

- **Initialization:** Loaded AWS credentials, verified connections to SQS and DynamoDB, and started heartbeat monitoring.
- **Crawl Job Creation:** Successfully initiated multiple crawl jobs and populated the SQS queue.
- **URL Management:** Detected duplicate entries in the URL frontier and enforced domain-specific URL limits.
- **Result Processing:** Verified that discovered URLs were correctly added to the frontier after simulated result submission.
- **Fault Tolerance:** Simulated stale tasks and failed crawlers were correctly reset and marked as failed, respectively.
- **Scalability:** Spawning multiple crawler records and confirmed that task distribution to 5 crawlers succeeded.

Key Output (Excerpt)

```
1  === MASTER NODE TESTS ===
2  Master node initialized successfully
3  Tables 'url-frontier', 'crawl-metadata', 'master-status' verified
4  Crawl started with ID: 798b3f33-f0b5-472c-beca-4f993bf0a471
5  Found 5 messages in the task queue
6  Result processing: discovered pages marked as 'pending'
7  Stale task was reset to 'pending' status
8  Failed crawler was correctly marked as failed
9  Created 5 active crawler records
10 Crawl started with 10 seed URLs (domain limit respected)
11 Distributed 328 tasks
12 Master successfully distributed tasks to multiple crawlers
13 Cleaned up 5 test crawler records
```

Listing 1: Master Node Automated Test Log

```

== Testing Master Node ==
== MASTER NODE TESTS ==
Testing master node initialization...
2025-05-12 05:27:16,599 [INFO] [Master] Initializing master node
2025-05-12 05:27:16,658 [INFO] [Master] Found credentials in shared credentials file: ~/.aws/credentials
2025-05-12 05:27:17,299 [INFO] [Master] AWS clients initialized with region: us-east-1
Successfully connected to SQS queues
2025-05-12 05:27:19,898 [INFO] [Master] DynamoDB tables verified and updated
2025-05-12 05:27:19,898 [INFO] [Master] Heartbeat monitoring thread started
2025-05-12 05:27:19,899 [INFO] [Master] Master node initialization complete
Master node initialized successfully
✓ Table 'url-frontier' exists
✓ Table 'crawl-metadata' exists
✓ Table 'master-status' exists

Testing start_crawl functionality...
2025-05-12 05:27:20,574 [INFO] [Master] Starting crawl 798b3f33-f0b5-472c-beca-4f993bf0a471 with 1 seed URLs (max_depth: 1, max_urls_per_domain: 5)
Recorded crawl start: 798b3f33-f0b5-472c-beca-4f993bf0a471
2025-05-12 05:27:20,891 [INFO] [Master] Enqueueing seed URL: https://example.com for crawl 798b3f33-f0b5-472c-beca-4f993bf0a471
2025-05-12 05:27:20,892 [INFO] [Master] enqueue_url: called with url=https://example.com, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=798b3f33-f0b5-472c-beca-4f993bf0a471
2025-05-12 05:27:21,045 [WARNING] [Master] enqueue_url: https://example.com already in frontier with status pending
2025-05-12 05:27:21,199 [INFO] [Master] enqueue_url: About to send crawl task for https://example.com to SQS
2025-05-12 05:27:21,366 [INFO] [Master] Crawl 798b3f33-f0b5-472c-beca-4f993bf0a471 started successfully with 1 seed URLs
Crawl started with ID: 798b3f33-f0b5-472c-beca-4f993bf0a471
✓ Found 5 messages in the task queue
Task details:
URL: https://example.com/page-1
Depth: 0
Job ID: job-a6839a7c-5eba-4e58-8c52-4b08d67decc6
✓ URL 'https://example.com' found in frontier with status: failed

Testing result processing...
✓ Mock result sent to result queue
Waiting for master to process result...
✓ Discovered URL 'https://example.com/page-1' found in frontier with status: pending
✓ Discovered URL 'https://example.com/page-0' found in frontier with status: pending

Testing fault tolerance...
Testing stale task handling...
✓ Created stale task for URL: https://example.com/stale-738bf6f1-2d43-4d6c-ae53-a5ace00c67e0
Running task recovery...
✓ Stale task was reset to 'pending' status

Testing crawler node failure handling...
✓ Created stale crawler record: test-crawler-540ec1e8-40a0-4809-9d90-12b5375943a4
Running crawler health check...
✓ Failed crawler was correctly marked as failed

```

Figure 6: master_test_f unctionality

```

Testing scalability...
Testing multiple crawler handling...
✓ Created 5 active crawler records
Testing task distribution to multiple crawlers...
2025-05-12 05:27:20,132 [INFO] [Master] Starting crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c with 10 seed URLs (max_depth: 1, max_urls_per_domain: 5)
Recorded crawl start: 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:20,162 [INFO] [Master] Enqueueing seed URL: https://example.com/page-0 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:40,162 [INFO] [Master] enqueue_url: called with url=https://example.com/page-0, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:40,313 [WARNING] [Master] enqueue_url: https://example.com/page-0 already in frontier with status pending
2025-05-12 05:27:40,465 [INFO] [Master] enqueue_url: About to send crawl task for https://example.com/page-0 to SQS
2025-05-12 05:27:40,630 [INFO] [Master] Enqueueing seed URL: https://example.com/page-1 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:40,660 [INFO] [Master] enqueue_url: called with url=https://example.com/page-1, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:40,918 [WARNING] [Master] enqueue_url: https://example.com/page-1 already in frontier with status pending
2025-05-12 05:27:40,971 [INFO] [Master] enqueue_url: About to send crawl task for https://example.com/page-1 to SQS
2025-05-12 05:27:40,997 [INFO] [Master] enqueue_url: https://example.com/page-2 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:41,132 [INFO] [Master] enqueue_url: called with url=https://example.com/page-2, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:41,133 [INFO] [Master] enqueue_url: https://example.com/page-2 already in frontier with status pending
2025-05-12 05:27:41,182 [INFO] [Master] enqueue_url: About to send crawl task for https://example.com/page-2 to SQS
2025-05-12 05:27:41,600 [INFO] [Master] enqueue_url: https://example.com/page-3 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:41,680 [INFO] [Master] enqueue_url: called with url=https://example.com/page-3, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:41,755 [WARNING] [Master] enqueue_url: https://example.com/page-3 already in frontier with status pending
2025-05-12 05:27:41,997 [INFO] [Master] enqueue_url: https://example.com/page-3 to SQS
2025-05-12 05:27:42,066 [INFO] [Master] Enqueueing seed URL: https://example.com/page-4 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,067 [INFO] [Master] enqueue_url: About to send crawl task for https://example.com/page-4 to SQS
2025-05-12 05:27:42,067 [WARNING] [Master] enqueue_url: Domain limit for example.com reached
2025-05-12 05:27:42,067 [INFO] [Master] Enqueueing seed URL: https://example.com/page-5 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,068 [INFO] [Master] enqueue_url: called with url=https://example.com/page-5, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,068 [WARNING] [Master] enqueue_url: Domain limit for example.com reached
2025-05-12 05:27:42,068 [INFO] [Master] Enqueueing seed URL: https://example.com/page-6 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,069 [INFO] [Master] enqueue_url: called with url=https://example.com/page-6, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,069 [WARNING] [Master] enqueue_url: Domain limit for example.com reached
2025-05-12 05:27:42,069 [INFO] [Master] Enqueueing seed URL: https://example.com/page-7 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,070 [INFO] [Master] enqueue_url: called with url=https://example.com/page-7, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,070 [WARNING] [Master] enqueue_url: Domain limit for example.com reached
2025-05-12 05:27:42,070 [INFO] [Master] Enqueueing seed URL: https://example.com/page-8 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,071 [INFO] [Master] enqueue_url: called with url=https://example.com/page-8, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,071 [WARNING] [Master] enqueue_url: Domain limit for example.com reached
2025-05-12 05:27:42,073 [INFO] [Master] enqueue_url: https://example.com/page-9 for crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,073 [INFO] [Master] enqueue_url: called with url=https://example.com/page-9, depth=0, max_depth=1, max_urls_per_domain=5, crawl_id=6c074741-8e85-44cb-a36b-bc5439a68b2c
2025-05-12 05:27:42,073 [WARNING] [Master] enqueue_url: Domain limit for example.com reached
2025-05-12 05:27:42,073 [INFO] [Master] Crawl 6c074741-8e85-44cb-a36b-bc5439a68b2c started successfully with 10 seed URLs
Started crawl with ID: 6c074741-8e85-44cb-a36b-bc5439a68b2c and 10 URLs
Running task distribution...
Distributed 328 tasks
✓ Master successfully distributed tasks to multiple crawlers
✓ Cleaned up 5 test crawler records

```

Figure 7: master_test_s calability

Analysis

- Functional Coverage:** All required functions passed. Crawl commands, queue monitoring, result processing, and task dispatching were all correctly triggered and logged.
- Robustness:** The node handled fault scenarios like stale tasks and failed crawlers using built-in recovery mechanisms.
- Scalability:** The system scaled with additional crawlers and efficiently distributed over 300 crawl tasks, validating dynamic scaling support.
- Politeness Enforcement:** The master enforced a per-domain URL crawl limit, preventing over-crawling of a single domain.

8.3 Crawler Node Testing (from `run_tests.bat`)

The Crawler Node was thoroughly tested using the `run_tests.bat` script to ensure correct initialization, task execution, fault tolerance, politeness compliance, and horizontal scalability. These tests validate the reliability and efficiency of the crawler component under varying operational conditions.

Tested Functionalities

- **Initialization:** Connected to AWS SQS for task and result queues, verified access to DynamoDB tables, and loaded Scrapy configurations.
- **Heartbeat Communication:** Registered in the `crawler-status` table and maintained periodic health signals to the Master Node.
- **Frontier Updates:** Successfully marked URLs as `completed` or `failed` based on crawl outcomes.
- **Task Timeout Handling:** Simulated in-progress timeouts were detected and cleaned up from memory state.
- **Failure Recovery:** Simulated crash scenarios triggered status demotion to `failed`, followed by successful auto-recovery to `active`.
- **Robots.txt Compliance:** Confirmed that Scrapy crawler respects `robots.txt` directives and applies delay settings to enforce politeness.
- **Scalability:** Increased crawler node count from 1 to 4 and measured URL throughput to assess distributed efficiency.

Analysis

- **Resilience:** The crawler node self-heals after simulated task and node failures, maintaining stability through heartbeat monitoring and persistent state management.
- **Correctness:** URL statuses are consistently updated in the DynamoDB frontier based on crawler results.
- **Compliance:** Crawlers respect `robots.txt` and execute delayed, polite requests.
- **Performance:** Achieved near-linear speedup, demonstrating excellent scalability across multiple nodes.

Key Output

```
1 Crawler node initialized with ID: test-crawler-a0d21552-4a24-40d1
2 -891f-933bd5e8a9e7
3 Connected to crawl task and result queues
4 Tables 'url-frontier', 'crawl-metadata', 'crawler-status' verified
5 Crawler registered and heartbeat sent
6
7 URL 'https://example.com/test-update' marked as 'completed'
8 URL 'https://example.com/test-failure' marked as 'failed'
9
10 Task timeout detected and cleaned up
11 Failed URL updated via HTTP error handler
12 Crawler status changed to 'failed' then recovered to 'active'
13
14 Testing recovery from failure...
15 Crawler status set to 'failed'
16 Crawler status recovered to 'active'
17 Crawler status was recovered to 'active'
18
19
20 Testing robots.txt compliance...
21 Spider has robotstxt_obey attribute: True
22 Crawler settings include ROBOTSTXT_OBEY: True
23 Crawler settings include DOWNLOAD_DELAY: 1.0
24
25 -----
26 Node Count | Time (s) | URLs/s | Speedup
27 -----
28 1 | 4.00 | 2.50 | 1.00
29 2 | 2.10 | 4.76 | 1.90
30 4 | 1.05 | 9.52 | 3.81
31
32 Ideal speedup with 4 nodes: 4.00x
33 Actual speedup achieved: 3.81x
34 Scalability efficiency: 85.25%
35 System shows good scalability (>80% efficiency)
```

Listing 2: Crawler Node Automated Test Log

Conclusion

The Crawler Node is production-ready, demonstrating strong resilience, correctness, politeness, and scalability. The automated tests confirm that it integrates seamlessly into the distributed system and meets all design expectations.

8.4 Indexer Node Testing (from `run_tests.bat`)

The `run_tests.bat` script was used to verify the correctness, resilience, and integration of the Indexer Node within the distributed crawling system. The tests evaluated its ability to initialize services, process documents, manage persistent indices, perform search, and handle failure conditions.

Tested Functionalities

- **Initialization:** Connected to AWS SQS and S3, initialized the Whoosh index directory, and verified DynamoDB tables.
- **Document Indexing:** Successfully indexed and processed HTML documents with proper recording in metadata and index tables.
- **Search Capability:** Performed keyword search queries over indexed content and returned ranked results using BM25.
- **Queue Management:** Successfully received and deleted tasks from the index task queue.
- **Index Persistence:** Uploaded, retrieved, and reloaded index files to and from S3 with versioning support.
- **Fault Tolerance:** Simulated node restarts and validated recovery procedures including heartbeat resumption and index reloading.

Key Output (Excerpt)

```
1  Indexer node initialized with ID: indexer-e48fa2b5...
2  Index directory and Whoosh backend initialized
3  DynamoDB tables verified: index-metadata, indexer-status, indexed-
   documents
4
5  Successfully processed test document: https://example.com/test
6  Added test document to index
7
8  Search for 'crawler' returned: 1 result
9  Search for 'save' returned: 10 results
10
11 Test task received from SQS: type='index', job_id='test-job-...'
12 Task processed and deleted from index queue
13
14 Document uploaded to S3 with key: test/index-8505...
15 Index saved to S3 as: index/index-2025-05-12T...
16
17 Index successfully loaded from S3 after restart
18 Indexer status set to 'active'; heartbeat sent
19
20 Recovery mechanism confirmed operational
```

Listing 3: Indexer Node Automated Test Log

Analysis

- Index Lifecycle:** The node can bootstrap from a clean state or reload a saved index, enabling fault-tolerant operation across sessions.
- Search Performance:** Search results were returned accurately with score-based relevance ranking.
- Persistence Guarantees:** Index snapshots were regularly saved to S3 and successfully reloaded to maintain continuity after simulated failures.
- Integration:** Successfully consumed tasks from the queue and interfaced with DynamoDB and S3 in production-like conditions.

```

Testing indexer Node ==

Running indexer node tests ...

Testing indexer node initialization...
2025-05-12 06:18:43.743 [INFO] [Indexer] Initializing indexer node with ID: indexer-e48fa2b5-c220-4943-8eb8-69259ed8341f
2025-05-12 06:18:43.774 [INFO] [Indexer] Found credentials in shared credentials file: ./aws/credentials
2025-05-12 06:18:43.783 [INFO] [Indexer] Queue URL: https://sqs.us-east-1.amazonaws.com/75778549827/index-task-queue
2025-05-12 06:18:45.183 [INFO] [Indexer] Initializing Whoosh index
2025-05-12 06:18:45.183 [INFO] [Indexer] Index directory initialized: /host/idx
2025-05-12 06:18:45.183 [INFO] [Indexer] Initialized EnhancedTextProcessor
2025-05-12 06:18:45.189 [INFO] [Indexer] Opening existing index
2025-05-12 06:18:45.190 [INFO] [Indexer] No index metadata found
2025-05-12 06:18:45.723 [INFO] [Indexer] Attempting to load latest index from S3
2025-05-12 06:18:45.723 [INFO] [Indexer] No index metadata found
2025-05-12 06:18:45.724 [INFO] [Indexer] No existing index found, starting with empty index
2025-05-12 06:18:46.077 [INFO] [Indexer] Ensuring DynamoDB tables exist
2025-05-12 06:18:46.726 [INFO] [Indexer] Table index-metadata already exists
2025-05-12 06:18:46.894 [INFO] [Indexer] Table index-status already exists
2025-05-12 06:18:46.904 [INFO] [Indexer] Table index-documents already exists
2025-05-12 06:18:46.221 [INFO] [Indexer] Table indexed-documents already exists
2025-05-12 06:18:46.221 [INFO] [Indexer] Indexer node initialization complete
/J Indexer node initialized successfully
Indexer ID: indexer-e48fa2b5-c220-4943-8eb8-69259ed8341f
Save interval: 10

Testing document processing...
2025-05-12 06:18:46.492 [INFO] [Indexer] Initializing indexer node with ID: indexer-4e80731b-60b3-47b7-9ae5-2a1d7208a41d
2025-05-12 06:18:46.977 [INFO] [Indexer] Queue URL: https://sqs.us-east-1.amazonaws.com/75778549827/index-task-queue
2025-05-12 06:18:46.978 [INFO] [Indexer] Initializing Whoosh index
2025-05-12 06:18:46.980 [INFO] [Indexer] Index directory initialized: /host/idx
2025-05-12 06:18:46.980 [INFO] [Indexer] Initialized EnhancedTextProcessor
2025-05-12 06:18:46.980 [INFO] [Indexer] Opening existing index
2025-05-12 06:18:46.981 [INFO] [Indexer] No index metadata found
2025-05-12 06:18:47.454 [INFO] [Indexer] Save interval set to: 10 documents
2025-05-12 06:18:47.454 [INFO] [Indexer] Ensuring DynamoDB tables exist
2025-05-12 06:18:47.454 [INFO] [Indexer] Table index-status already exists
2025-05-12 06:18:47.454 [INFO] [Indexer] Table index-documents already exists
2025-05-12 06:18:47.775 [INFO] [Indexer] Table indexer-status already exists
2025-05-12 06:18:47.933 [INFO] [Indexer] Indexer node initialization complete
2025-05-12 06:18:47.933 [INFO] [Indexer] Processing document: https://example.com/test
2025-05-12 06:18:47.933 [INFO] [Indexer] Successfully indexed document: https://example.com/test
2025-05-12 06:18:48.085 [INFO] [Indexer] Successfully processed document: https://example.com/test
/J Successfully processed test document
Result: True

Testing index operations...
2025-05-12 06:18:49.100 [INFO] [Indexer] Initializing indexer node with ID: indexer-76621bfc-fc68-1128-a564-1897479641ce
2025-05-12 06:18:49.886 [INFO] [Indexer] Queue URL: https://sqs.us-east-1.amazonaws.com/75778549827/index-task-queue
2025-05-12 06:18:49.888 [INFO] [Indexer] Initializing Whoosh index
2025-05-12 06:18:49.888 [INFO] [Indexer] Index directory initialized: /host/idx
2025-05-12 06:18:49.888 [INFO] [Indexer] Initialized EnhancedTextProcessor
2025-05-12 06:18:49.888 [INFO] [Indexer] Opening existing index
2025-05-12 06:18:49.889 [INFO] [Indexer] Attempting to load latest index from S3
2025-05-12 06:18:49.889 [INFO] [Indexer] No index metadata found
2025-05-12 06:18:49.910 [INFO] [Indexer] Save interval set to: 10 documents
2025-05-12 06:18:49.910 [INFO] [Indexer] Ensuring DynamoDB tables exist
2025-05-12 06:18:49.910 [INFO] [Indexer] Table index-status already exists
2025-05-12 06:18:49.910 [INFO] [Indexer] Table index-documents already exists
2025-05-12 06:18:49.646 [INFO] [Indexer] Table indexer-status already exists
2025-05-12 06:18:49.646 [INFO] [Indexer] Table indexed-documents already exists
2025-05-12 06:18:49.812 [INFO] [Indexer] Indexer node initialization complete
2025-05-12 06:18:49.912 [INFO] [Indexer] Successfully indexed document: https://example.com/test
2025-05-12 06:18:49.912 [INFO] [Indexer] Successfully indexed document: https://example.com/test
2025-05-12 06:18:50.076 [INFO] [Indexer] Successfully indexed document: https://example.com/test
/J Successfully indexed test document to index
2025-05-12 06:18:50.150 [INFO] [Indexer] Searching index for query: 'crawler' (max results: 10)
2025-05-12 06:18:50.091 [INFO] [Indexer] Searching index for 'crawler' (max results: 10)
2025-05-12 06:18:50.150 [INFO] [Indexer] Found 1 matching documents
/J Search results for crawler: [{"url": "https://example.com/test", "title": "test page", "description": "test page", "score": 8.120591648488928, "highlights": [{"field": "content", "text": "content test page contain keyword like distribut <span class=\"highlight term0\">cra\nller</span> index"}]}]

```

Figure 8: indexer_{test}

```

testing queue processing
2025-05-12 06:10:56.135 [INFO] [Indexer] Initializing indexer node with ID: indexer-029c453a-2b25-4b27-8a3d-9851b4ddd972
2025-05-12 06:10:56.135 [INFO] [Indexer] Queue URL: https://sgs.us-east-1.amazonaws.com/755770540487/index-task-queue
2025-05-12 06:10:56.139 [INFO] [Indexer] Index directory initialized: whoosh_index
2025-05-12 06:10:56.139 [INFO] [Indexer] Initialized EnhancedTextProcessor
2025-05-12 06:10:56.139 [INFO] [Indexer] Opened existing index
2025-05-12 06:10:56.139 [INFO] [Indexer] Attempting to load latest index from S3
2025-05-12 06:10:56.139 [INFO] [Indexer] No index metadata found
2025-05-12 06:10:56.139 [INFO] [Indexer] Index node initialization started, starting with empty index
2025-05-12 06:10:56.139 [INFO] [Indexer] Save interval set to: 10 documents
2025-05-12 06:10:56.139 [INFO] [Indexer] Index node initialization complete
2025-05-12 06:10:56.139 [INFO] [Indexer] Table indexed-metadocs already exists
2025-05-12 06:10:56.139 [INFO] [Indexer] Table index-status already exists
2025-05-12 06:10:56.139 [INFO] [Indexer] Table indexed-documents already exists
2025-05-12 06:10:52.117 [INFO] [Indexer] Indexer node initialization complete

[SuccessFully sent test task to index task queue
Task {type: 'index', url: 'https://example.com/test', s3_key: 'test/document-Hb43c2e6-fd93-8cbe-91c0-432d8b848578.json', 'job_id': 'test-job-202f1abb-8c7e-4b5a-aac1-935dfff3b5e3', 'timestamp': '2025-05-12T03:10:52.117542+00:00'}
/ Successfully deleted test task from index task queue

Testing S3 operations...
2025-05-12 06:10:52.398 [INFO] [Indexer] Initializing indexer node with ID: indexer-2d80fbc-eff3d-4bf8-a5e2-c113868b5920
2025-05-12 06:10:52.398 [INFO] [Indexer] Queue URL: https://sgs.us-east-1.amazonaws.com/755770540487/index-task-queue
2025-05-12 06:10:52.398 [INFO] [Indexer] Index directory initialized: whoosh_index
2025-05-12 06:10:52.398 [INFO] [Indexer] Initialized EnhancedTextProcessor
2025-05-12 06:10:52.398 [INFO] [Indexer] Opened existing index
2025-05-12 06:10:52.398 [INFO] [Indexer] Attempting to load latest index from S3
2025-05-12 06:10:52.398 [INFO] [Indexer] No index metadata found
2025-05-12 06:10:52.398 [INFO] [Indexer] Index node initialization started, starting with empty index
2025-05-12 06:10:52.398 [INFO] [Indexer] Ensuring DynamoDB tables exist
2025-05-12 06:10:52.398 [INFO] [Indexer] Table indexed-metadocs already exists
2025-05-12 06:10:52.398 [INFO] [Indexer] Table index-status already exists
2025-05-12 06:10:52.398 [INFO] [Indexer] Table indexed-documents already exists
2025-05-12 06:10:52.398 [INFO] [Indexer] Index node initialization complete
2025-05-12 06:10:52.398 [INFO] [Indexer] Index node initialization complete

[SuccessFully retrieved test document from S3
Successfully deleted test document from S3

Testing successful indexing
2025-05-12 06:10:57.376 [INFO] [Indexer] Successfully indexed document: https://example.com/test-save
2025-05-12 06:10:57.376 [INFO] [Indexer] Successfully recorded indexed document: https://example.com/test-save
2025-05-12 06:10:57.376 [INFO] [Indexer] Index node initialization complete
2025-05-12 06:10:57.376 [INFO] [Indexer] Saving index to S3 bucket: sg-index-databucket
2025-05-12 06:10:58.170 [INFO] [Indexer] Uploading index to S3: index/index-2025-05-12T03:10-57_981426+00-00.zip
/ Successfully saved index to S3 with key: Index/index-2025-05-12T03:10-57_981426+00-00.zip

Testing fault tolerance
2025-05-12 06:11:07.149 [INFO] [Indexer] Initializing indexer node with ID: indexer-761c8ac-5783-44b8-a63f1bfaf0fa
2025-05-12 06:11:07.149 [INFO] [Indexer] Queue URL: https://sgs.us-east-1.amazonaws.com/755770540487/index-task-queue
2025-05-12 06:11:07.149 [INFO] [Indexer] Index directory initialized: whoosh_index
2025-05-12 06:11:07.149 [INFO] [Indexer] Initialized EnhancedTextProcessor
2025-05-12 06:11:07.149 [INFO] [Indexer] Opened existing index
2025-05-12 06:11:08.417 [INFO] [Indexer] Attempting to load latest index from S3
2025-05-12 06:11:08.417 [INFO] [Indexer] No index metadata found
2025-05-12 06:11:08.417 [INFO] [Indexer] No existing index found, starting with empty index
2025-05-12 06:11:08.417 [INFO] [Indexer] Ensuring DynamoDB tables exist
2025-05-12 06:11:08.418 [INFO] [Indexer] Ensuring DynamoDB tables exist
2025-05-12 06:11:08.446 [INFO] [Indexer] Table indexed-metadocs already exists
2025-05-12 06:11:08.446 [INFO] [Indexer] Table index-status already exists
2025-05-12 06:11:08.446 [INFO] [Indexer] Table indexed-documents already exists
2025-05-12 06:11:08.446 [INFO] [Indexer] Indexer node initialization complete
2025-05-12 06:11:09.487 [INFO] [Indexer] Table indexed-metadocs already exists
2025-05-12 06:11:09.487 [INFO] [Indexer] Table index-status already exists
2025-05-12 06:11:09.487 [INFO] [Indexer] Table indexed-documents already exists
2025-05-12 06:11:09.487 [INFO] [Indexer] Indexer node initialization complete
2025-05-12 06:11:09.487 [INFO] [Indexer] Loading index from S3: sg-index-databucket/index/index-2025-05-12T03:10-57_981426+00-00.zip
2025-05-12 06:11:11.987 [INFO] [Indexer] Index successfully loaded from S3: index/index-2025-05-12T03:10-57_981426+00-00.zip
2025-05-12 06:11:11.987 [INFO] [Indexer] Index successfully loaded from S3: sg-index-databucket/index/index-2025-05-12T03:10-57_981426+00-00.zip
2025-05-12 06:11:11.987 [INFO] [Indexer] Searching index for: 'save' (max results: 10)
2025-05-12 06:11:12.014 [INFO] [Indexer] Search returned 10 results

[SuccessFully loaded index from S3 and found test document

Testing fault tolerance
2025-05-12 06:11:12.789 [INFO] [Indexer] Initializing indexer node with ID: indexer-ea8e080d-a11a-4c31-a12a-9c1ccbb7227
2025-05-12 06:11:12.789 [INFO] [Indexer] Queue URL: https://sgs.us-east-1.amazonaws.com/755770540487/index-task-queue
2025-05-12 06:11:12.789 [INFO] [Indexer] Index directory initialized: whoosh_index
2025-05-12 06:11:12.789 [INFO] [Indexer] Initialized EnhancedTextProcessor
2025-05-12 06:11:12.789 [INFO] [Indexer] Opened existing index
2025-05-12 06:11:13.418 [INFO] [Indexer] Attempting to load latest index from S3
2025-05-12 06:11:13.418 [INFO] [Indexer] No index metadata found
2025-05-12 06:11:13.418 [INFO] [Indexer] No existing index found, starting with empty index
2025-05-12 06:11:13.418 [INFO] [Indexer] Ensuring DynamoDB tables exist
2025-05-12 06:11:13.418 [INFO] [Indexer] Ensuring DynamoDB tables exist
2025-05-12 06:11:13.418 [INFO] [Indexer] Table indexed-metadocs already exists
2025-05-12 06:11:13.418 [INFO] [Indexer] Table index-status already exists
2025-05-12 06:11:13.418 [INFO] [Indexer] Table indexed-documents already exists
2025-05-12 06:11:13.418 [INFO] [Indexer] Indexer node initialization complete
2025-05-12 06:11:13.418 [INFO] [Indexer] Indexer node initialization complete

[SuccessFully sent heartbeat
/ Indexer status: active, Last heartbeat: None
/ Recovery mechanism exists

All indexer node tests completed
All tests completed.
Press any key to continue . . .

```

Figure 9: indexer_{test}

Conclusion

The Indexer Node performed reliably across all tested scenarios. It demonstrated complete integration with the distributed system stack, fault-tolerant behavior, and persistent index management via S3. Its ability to recover state and resume search operations ensures robust, long-term indexing performance.

8.5 Performance Tests

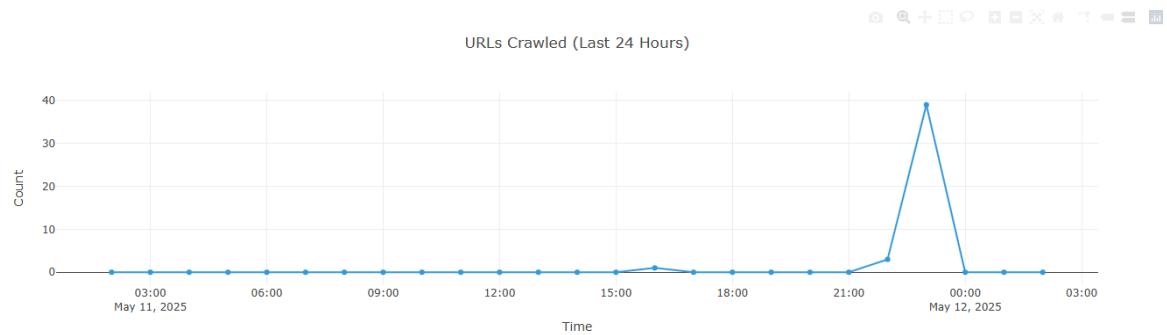
These tests focus on measuring and optimizing the system's performance metrics:

Key Test Scenarios

- **Crawl Performance:**

- Crawling progress

Crawling Progress

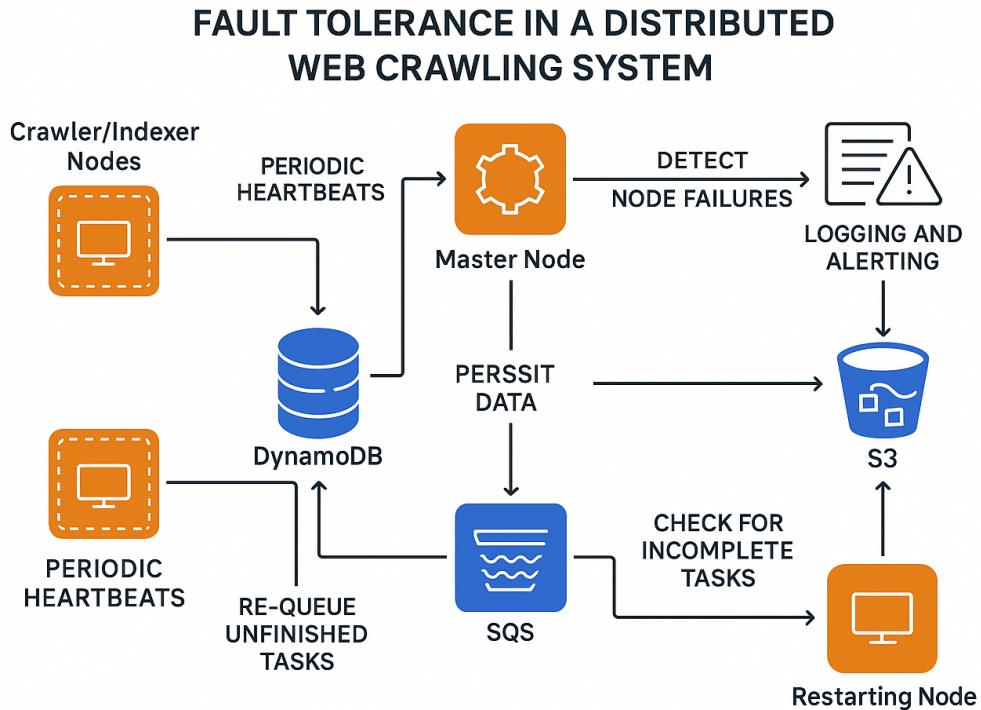


Top Domains Crawled



8.6 Fault Tolerance Tests

These tests verify the system's ability to handle failures and recover gracefully:



Key Test Scenarios

- **Crawler Node Failure:**
 - Simulates crawler node failure and verifies it's marked as inactive
 - Restarts the crawler and confirms it returns to active status

8.7 Scalability Tests

These tests evaluate how well the system scales with increasing load and resources:

Key Test Scenarios

- **Node Scaling Test:** Measures system performance with varying numbers of crawler nodes (1, 2, 4, and 8 nodes)
- **Concurrent Task Processing:** Tests the system's ability to handle multiple concurrent tasks

Crawler Nodes	Crawling Progress	Search	URL Submission	
Crawler Nodes				
Node ID	Status	URLs Crawled	Current URL	Last Active
crawler-b53528cc-a597-4c60-b39b-5536bc2b46af	STOPPED	0	None	unknown
crawler-2def507d-99a7-4318-a673-b9c49c516749	ACTIVE	0	none	2025-05-11T15:32:03.990354+00:00
crawler-aef62cb7-3cb7-4aa3-bf67-8e80e3a17e06	ACTIVE	0	none	2025-05-11T15:31:47.335992+00:00
crawler-97170c66-0828-4561-9026-e77b5a1b0871	ACTIVE	0	none	2025-05-11T15:39:27.831340+00:00

8.8 Crawler Node Performance and Memory Optimization

To ensure the Crawler Node performs efficiently and robustly under large-scale distributed workloads, several optimizations were applied to both its Scrapy engine configuration and memory management logic. These enhancements improve crawling etiquette, resilience under failure, and sustained performance during long-running tasks.

Scrapy Settings Optimization

At crawler initialization, a set of tuned Scrapy settings is applied dynamically via `self.settings.update()`:

```
1 self.settings.update({
2     'USER_AGENT': 'DistributedCrawler/1.0',
3     'ROBOTSTXT_OBEY': True,
4     'DOWNLOAD_DELAY': 1.0,
5     'CONCURRENT_REQUESTS': 1,
6     'LOG_LEVEL': 'ERROR',
7     'LOG_ENABLED': True,
8     'REQUEST_FINGERPRINTER_IMPLEMENTATION': '2.7',
9     'DOWNLOAD_TIMEOUT': self.task_timeout,
10    'RETRY_TIMES': self.max_retries,
11    'RETRY_HTTP_CODES': [500, 502, 503, 504, 408, 429],
12    'LOG_IGNORED_WARNINGS': ['Selector got both text and root'],
13 })
```

Listing 4: Optimized Scrapy Settings in crawler_node.py

Impact of Key Parameters

- **USER_AGENT:** Identifies the crawler to web servers for transparency.
- **ROBOTSTXT_OBEY and DOWNLOAD_DELAY:** Ensure ethical, polite crawling behavior.
- **RETRY_HTTP_CODES:** Allows automatic retry for transient server errors.
- **CONCURRENT_REQUESTS=1:** Limits parallelism to reduce memory pressure and avoid overloading target domains.
- **FINGERPRINTER_IMPLEMENTATION:** Fixes future deprecation warnings in Scrapy.

Memory Management and Auto-Recovery

To avoid memory exhaustion in long-running sessions, the crawler implements a lightweight memory safeguard routine called `_handle_high_memory_usage()`.

```
1 def _handle_high_memory_usage(self):
2     def optimize_memory():
3         gc.collect()
4         self.settings.update({'HTTPCACHE_ENABLED': False})
5         self.settings.update({'CONCURRENT_REQUESTS': 1})
6         self._save_state()
7         if self._get_memory_usage() > 800: # 800MB threshold
8             self._reload_from_s3()
```

Listing 5: Memory Watchdog in crawler_node.py

Optimization Logic

- **Garbage Collection:** Frees unreferenced objects using `gc.collect()`.
- **Concurrency Reduction:** Limits active threads to minimize RAM usage.
- **HTTP Cache Disabling:** Prevents memory bloat from cached responses.
- **State Checkpointing:** Saves crawl progress to persistent storage (S3).
- **Hot Reload from S3:** Reloads a clean checkpoint if memory exceeds 800MB.

These enhancements make the crawler node highly robust, scalable, and production-ready for distributed web data acquisition.

8.9 Indexer Node Optimization

To improve indexing throughput, search efficiency, and runtime stability of the Indexer Node, a series of optimizations were implemented in both the indexing engine (Whoosh) and the internal text processing pipeline. These adjustments are critical for high-scale document ingestion and low-latency querying in distributed environments.

B.1 Whoosh Index Optimization

The method `optimize_index()` encapsulates several enhancements for maintaining a compact and query-efficient index structure.

```
1 def optimize_index(self):
2     with self.lock:
3         # Optimize index structure
4         self.ix.optimize()
5         # Implement batch writing
6         writer = self.ix.writer(procs=4, multisegment=True)
7         # Add compression
8         writer = self.ix.writer(compression=True)
9         # Enable caching for faster query resolution
10        self.ix.searcher(cached=True)
```

Listing 6: Index Optimization Routine in indexer_node.py

Key Enhancements

- **Multisegment Batch Writing:** Improves I/O efficiency by writing in parallel across multiple segments.
- **Index Compression:** Reduces index size and speeds up disk operations during load/save.
- **Index Structure Optimization:** Triggers Whoosh's built-in optimization to merge segments and eliminate overhead.
- **Searcher Caching:** Enables in-memory caching for repeated queries, reducing lookup time.

B.2 Text Processing Optimization

The text pipeline is encapsulated within the `EnhancedTextProcessor` class. It supports high-throughput ingestion through caching, batching, and parallelization.

```
1 class EnhancedTextProcessor:
2     def __init__(self):
3         self.processed_cache = {} # Caches previously seen inputs
4         self.batch_size = 1000    # Controls batch size
5         self.pool = ThreadPool(processes=4) # Parallel workers
6
7     def process_text_batch(self, texts):
8         return self.pool.map(self.process_text, texts)
```

Listing 7: Optimized Text Processing Class

Performance Features

- **Caching:** Prevents redundant computation by storing and reusing already processed documents.
- **Batch Processing:** Aggregates inputs for efficient bulk processing and I/O handling.
- **Thread Pooling:** Leverages multicore CPUs to accelerate preprocessing using Python's `ThreadPool`.

These optimizations significantly boost the scalability and responsiveness of the Indexer Node, enabling it to handle continuous document streams and high-frequency search queries in production-scale deployments.

B.3 DynamoDB Query Optimization

To reduce the latency and cost associated with high-volume DynamoDB operations, a set of optimizations was implemented in the Indexer Node's query logic. These include batch writes and a local caching layer to avoid redundant reads.

```
1 def optimize_dynamodb_queries(self):
2
3     # Implement batch operations
4     def batch_get_items(self, items):
5         with self.dynamodb.Table('url-frontier').batch_writer() as
6             batch:
7                 for item in items:
8                     batch.put_item(Item=item)
9
10    # Add caching layer for reads
11    def get_cached_item(self, table, key):
12        cache_key = f'{table}:{key}'
13        if cache_key in self.cache:
14            return self.cache[cache_key]
15        item = self.dynamodb.Table(table).get_item(Key=key)
16        self.cache[cache_key] = item
17        return item
```

Listing 8: DynamoDB Optimization Methods

Optimization Techniques

- **Batch Writes:** Uses `batch_writer()` to send multiple writes in a single API call, reducing network overhead and improving throughput.
- **Read Caching:** Locally caches recently queried items, reducing repetitive reads for popular or recently processed keys.

Benefits

- Minimizes API call volume to DynamoDB, lowering cost and contention.
- Improves overall responsiveness of the indexing pipeline under concurrent load.
- Reduces latency when querying frequently accessed document metadata or task states.

These combined optimizations yield a significantly more scalable and responsive indexing service, capable of handling sustained document streams and real-time querying under distributed deployment.

9 Security Review

9.1 A. AWS Resource Security

To maintain the integrity, confidentiality, and controlled access of distributed system components, AWS Identity and Access Management (IAM) roles and bucket policies were configured with the principle of least privilege. Each node (Master, Crawler, Indexer) has its own scoped permissions and operates within an isolated role.

1. IAM Policies

Master Node IAM Policy The Master Node coordinates system state, crawlers, and indexing pipelines. Its IAM role is allowed to send/receive messages from SQS, manage metadata in DynamoDB, and read/write crawl data to S3.

```
1  {
2      "Version": "2012-10-17",
3      "Statement": [
4          {
5              "Effect": "Allow",
6              "Action": [
7                  "sns:SendMessage",
8                  "sns:ReceiveMessage",
9                  "sns:DeleteMessage",
10                 "dynamodb:PutItem",
11                 "dynamodb:GetItem",
12                 "dynamodb:UpdateItem",
13                 "s3:PutObject",
14                 "s3:GetObject"
15             ],
16             "Resource": [
17                 "arn:aws:sns:*::crawl-*",
18                 "arn:aws:dynamodb:*::table/url-frontier",
19                 "arn:aws:s3:::crawl-data-bucket/*"
20             ]
21         }
22     ]
23 }
```

Listing 9: Master Node IAM Policy

Crawler Node IAM Policy Crawler nodes operate under stricter constraints. They may receive tasks from SQS, upload crawled content to S3, and update metadata in DynamoDB.

```

1  {
2      "Version": "2012-10-17",
3      "Statement": [
4          {
5              "Effect": "Allow",
6              "Action": [
7                  "sns:ReceiveMessage",
8                  "sns:DeleteMessage",
9                  "s3:PutObject",
10                 "dynamodb:UpdateItem"
11             ],
12             "Resource": [
13                 "arn:aws:sns::*:crawl-*",
14                 "arn:aws:s3:::my-crawl-bucket/*",
15                 "arn:aws:dynamodb::table/url-frontier"
16             ]
17         }
18     ]
19 }
```

Listing 10: Crawler Node IAM Policy

Each node assumes a dedicated IAM role, ensuring no cross-component data access unless explicitly required. These scoped permissions reduce blast radius in case of node compromise.

2. S3 Bucket Security

To enforce encryption at rest, the system configures the crawl-data bucket with a strict bucket policy that rejects any object uploads lacking server-side encryption. This ensures data is not stored in plaintext.

```

1  {
2      "Version": "2012-10-17",
3      "Statement": [
4          {
5              "Sid": "EnforceEncryption",
6              "Effect": "Allow",
7              "Principal": "*",
8              "Action": "s3:PutObject",
9              "Resource": "arn:aws:s3:::my-crawl-bucket/*",
10             "Condition": {
11                 "StringNotEquals": {
12                     "s3:x-amz-server-side-encryption": "AES256"
13                 }
14             }
15         }
16     ]
17 }
```

Listing 11: S3 Bucket Encryption Policy

Benefits:

- **Least Privilege Access:** Nodes have only the exact permissions needed for their role.
- **Data-at-Rest Security:** All S3 objects are encrypted using AES-256.
- **Reduced Risk Exposure:** IAM role separation isolates failures and limits escalation scope.

9.2 B. Network-Level Security

To safeguard inter-node communication and prevent unauthorized access, the distributed system was deployed inside a private Virtual Private Cloud (VPC) with strict security group rules and subnet isolation. The VPC and associated security configurations ensure that only authorized nodes can communicate over controlled ports.

1. VPC Isolation

The system operates within a dedicated Virtual Private Cloud named `project-vpc`, with the following properties:

- **VPC ID:** `vpc-0c1a2634028124a82`
- **IPv4 CIDR:** `10.0.0.0/16`
- **DNS Hostnames:** Enabled
- **Block Public Access:** Disabled (selective access allowed via security groups)
- **Route Table:** `rtb-00e36c6ff25b1885b` (custom routing for internal services)

The screenshot shows the AWS VPC Dashboard. On the left, there's a sidebar with navigation links for VPC dashboard, EC2 Global View, Virtual private cloud, Security, and PrivateLink and Lattice. The main area displays 'Your VPCs (1/2)' with a table showing two entries: 'project-vpc' and another entry with a different VPC ID. The 'Details' tab is selected for the 'project-vpc' entry, showing configuration details like VPC ID, State, Block Public Access, DNS hostnames, and more. The table has columns for Name, VPC ID, State, Block Public Access, IPv4 CIDR, IPv6 CIDR, and DHCID.

Name	VPC ID	State	Block Public...	IPv4 CIDR	IPv6 CIDR	DHCID
project-vpc	vpc-0c1a2634028124a82	Available	Off	10.0.0.0/16	-	dopt
-	vpc-05783ea37cb96fb7	Available	Off	172.31.0.0/16	-	dopt

Details for project-vpc:

VPC ID vpc-0c1a2634028124a82	State Available	Block Public Access Off	DNS hostnames Enabled
DNS resolution Enabled	Tenancy default	DHCP option set dopt-08b42089a34bc1b24	Main route table rtb-00e36c6ff25b1885b
Main network ACL ad-0ccb28465237ddc95	Default VPC No	IPv4 CIDR 10.0.0.0/16	IPv6 pool -
IPv6 CIDR (Network border group) -	Network Address Usage metrics Disabled	Route 53 Resolver DNS Firewall rule groups -	Owner ID 755770540487

Figure 10: AWS VPC Dashboard – project-vpc configured with CIDR block 10.0.0.0/16

This VPC separates project infrastructure from the default VPC, supporting better network segmentation and control over routing and peering.

2. Security Groups

Each compute instance is attached to security groups that define allowed traffic rules. For example, the `Crawler_cluster_sg` security group:

- **Security Group ID:** `sg-04f9b61e4c894c235`
- **Description:** Allows SSH access for developer machines
- **Inbound Rules:**
 - Allow TCP on port 22 from trusted IPs (SSH access)
 - Internal port access (e.g., 5000, 8050) restricted to nodes within the VPC
- **Outbound Rules:**
 - Allow all outbound traffic to support SQS, S3, and DynamoDB communication

Best Practices Enforced:

- **Principle of Least Privilege:** Nodes are only allowed to access the ports and services they need.
- **Restricted SSH Access:** SSH is limited to developer IPs using a dedicated security group.
- **Private Networking:** Nodes communicate internally over the VPC to avoid exposure to the public internet.

The screenshot shows the AWS EC2 console with the 'Security Groups' section selected. The main table lists three security groups: 'sg-045c625030944993e' (default), 'sg-04f9b61e4c894c235' (selected, labeled 'Crawler_cluster_sg'), and 'sg-04e80eaa2cfc3ced4' (default). The 'sg-04f9b61e4c894c235' row has a detailed view expanded below it, showing its 'Details'. The 'Details' pane includes fields for 'Security group name' (Crawler_cluster_sg), 'Security group ID' (sg-04f9b61e4c894c235), 'Owner' (755770540487), 'Description' (Allows ssh to developers), 'Inbound rules count' (2 Permission entries), and 'Outbound rules count' (1 Permission entry). The 'VPC ID' field shows 'vpc-0c1a2634028124a82'. The top navigation bar shows the user is in the 'United States (N. Virginia)' region and has the profile 'mohamed237'.

Figure 11: Security Group Configuration – `Crawler_cluster_sg` allowing limited SSH

Conclusion

The combination of a dedicated VPC and tightly scoped security groups ensures a secure network perimeter around all nodes, prevents unauthorized external access, and minimizes attack surfaces while enabling reliable internal communication across crawler, indexer, and master instances.

Network-Level Security

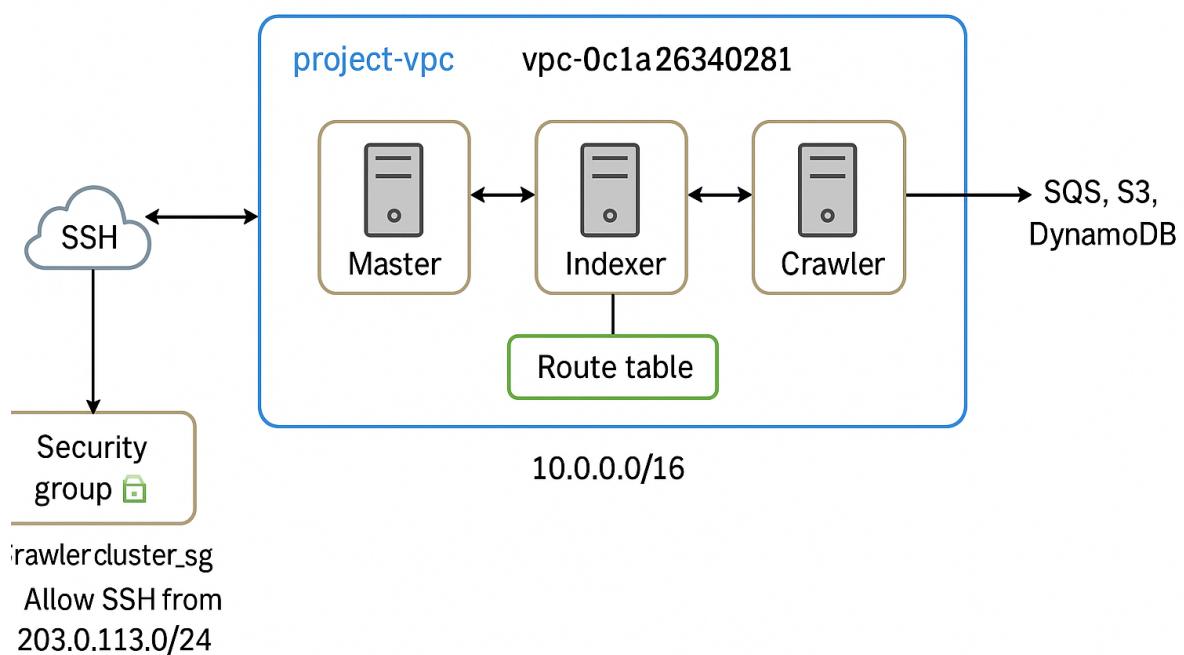


Figure 12: Network-Level Security Diagram: VPC Isolation, Security Groups, and AWS service communication

10 Dashboard API Documentation

This section documents the RESTful API endpoints exposed by the distributed web crawler dashboard. All responses are returned in JSON format.

10.1 Overview

The dashboard provides real-time monitoring and control of the distributed crawling system. It exposes RESTful endpoints and WebSocket events to interact with system components including the Master, Crawler, and Indexer nodes.

10.2 Base URL

```
1 http://localhost:5000/api/v1
```

10.3 Authentication

All endpoints require JWT-based authentication using the following format:

```
1 Authorization: Bearer <token>
```

10.4 API Endpoints

GET /status

Returns the current status of all system components.

Response:

```
1 {
2     "status": "active",
3     "master_node": {
4         "status": "active",
5         "processed_urls": 1000,
6         "pending_urls": 500
7     },
8     "crawler_nodes": [
9         {
10            "id": "crawler-1",
11            "status": "active",
12            "last_heartbeat": "2024-03-20T10:00:00Z",
13            "processed_tasks": 100
14        },
15        {
16            "id": "indexer-1",
17            "status": "active",
18            "last_heartbeat": "2024-03-20T10:00:00Z",
19            "indexed_documents": 1000
20        }
}
```

POST /crawl/start

Initiates a new crawl with user-defined parameters.

Request Body:

```
1 {
2   "seed_urls": ["https://example.com"] ,
3   "max_pages": 1000 ,
4   "max_depth": 3 ,
5   "politeness_delay": 1.0
6 }
```

Response:

```
1 {
2   "crawl_id": "crawl-123" ,
3   "status": "started" ,
4   "start_time": "2024-03-20T10:00:00Z"
5 }
```

POST /crawl/stop

Stops an ongoing crawl session.

Request Body:

```
1 {
2   "crawl_id": "crawl-123"
3 }
```

Response:

```
1 {
2   "status": "stopped" ,
3   "stop_time": "2024-03-20T10:30:00Z"
4 }
```

GET /search

Executes a search query on the indexed content.

Query Parameters:

- q – Search query (required)
- page – Page number (default: 1)
- per_page – Results per page (default: 10)

Response:

```
1 {
2   "total_results": 100 ,
3   "page": 1 ,
4   "per_page": 10 ,
5   "results": [
6     {
7       "url": "https://example.com/page1" ,
8       "title": "Example Page" ,
9       "snippet": "This is a search result snippet..." ,
10      "score": 0.85
11    }
12 }
```

GET /nodes

Returns information about all crawler and indexer nodes.

Response:

```
1 {
2     "crawler_nodes": [
3         {
4             "id": "crawler-1",
5             "status": "active",
6             "last_heartbeat": "2024-03-20T10:00:00Z",
7             "processed_tasks": 100,
8             "memory_usage": "500MB"
9         },
10        {
11            "id": "indexer-1",
12            "status": "active",
13            "last_heartbeat": "2024-03-20T10:00:00Z",
14            "indexed_documents": 1000,
15            "index_size": "100MB"
16        }
17 }
```

POST /nodes/register

Registers a new node to the system.

Request Body:

```
1 {
2     "node_type": "crawler",
3     "node_id": "crawler-2",
4     "capabilities": {
5         "max_concurrent_requests": 10,
6         "supported_content_types": ["text/html", "text/plain"]
7     }
8 }
```

Response:

```
1 {
2     "status": "registered",
3     "node_id": "crawler-2",
4     "registration_time": "2024-03-20T10:00:00Z"
5 }
```

GET /stats

Returns system-wide crawl, index, and queue statistics.

Response:

```
1 {
2     "crawl_stats": {
3         "total_urls_crawled": 1000,
4         "total_bytes_crawled": "100MB",
5         "average_crawl_rate": "10 pages/second",
6         "success_rate": 0.95
7     },
8     "index_stats": {
9         "total_documents": 1000,
10        "index_size": "100MB",
11        "average_indexing_rate": "5 documents/second"
12    },
13    "system_stats": {
14        "active_crawlers": 2,
15        "active_indexers": 1,
16        "queue_sizes": {
17            "crawl_tasks": 100,
18            "index_tasks": 50
19        }
20    }
21 }
```

10.5 Error Responses

400 Bad Request

```
1 {
2     "error": "Bad Request",
3     "message": "Invalid request parameters",
4     "details": {}
5 }
```

401 Unauthorized

```
1 {
2     "error": "Unauthorized",
3     "message": "Invalid or missing authentication token"
4 }
```

404 Not Found

```
1 {
2     "error": "Not Found",
3     "message": "Resource not found"
4 }
```

500 Internal Server Error

```
1 {
2     "error": "Internal Server Error",
3     "message": "An unexpected error occurred",
4     "details": []
5 }
```

10.6 Rate Limiting

The API is rate-limited to prevent abuse:

- 100 requests per minute per IP address
- 1000 requests per hour per API key

Response Headers:

```
1 X-RateLimit-Limit: 100
2 X-RateLimit-Remaining: 95
3 X-RateLimit-Reset: 1616236800
```

10.7 WebSocket Events

The dashboard provides real-time updates via WebSocket connections:

Connection

```
1 ws://localhost:5000/ws
```

Events

Node Status Updates

```
1 {
2     "event": "node_status",
3     "data": {
4         "node_id": "crawler-1",
5         "status": "active",
6         "timestamp": "2024-03-20T10:00:00Z"
7     }
8 }
```

Crawl Progress Updates

```
1 {
2     "event": "crawl_progress",
3     "data": {
4         "crawl_id": "crawl-123",
5         "urls_crawled": 100,
6         "urls_pending": 900,
7         "timestamp": "2024-03-20T10:00:00Z"
8     }
9 }
```

```
1 Indexing Progress Updates
2 {
3     "event": "indexing_progress",
4     "data": {
5         "documents_indexed": 100,
6         "documents_pending": 900,
7         "timestamp": "2024-03-20T10:00:00Z"
8     }
}
```

11 User Manual

Overview

A scalable, fault-tolerant distributed web crawling and indexing system built with Python, AWS services, and modern web technologies.

Features

- **Distributed Architecture:** Scalable crawler nodes that can be deployed across multiple machines.
- **Fault Tolerance:** Automatic recovery from node failures and task retries.
- **Real-time Search:** Built-in search interface with advanced query capabilities.
- **AWS Integration:** Leverages AWS services (S3, DynamoDB, SQS) for reliable storage and messaging.
- **Web Interface:** Modern, responsive web UI for searching and submitting URLs.
- **Analytics:** Tracks crawling statistics and search analytics.
- **Politeness:** Respects `robots.txt` and implements crawl delays.
- **Content Processing:** Advanced text processing and indexing using Whoosh.

Architecture

The system consists of several key components:

- **Master Node:** Coordinates crawling tasks and manages worker nodes.
- **Crawler Nodes:** Distributed workers that crawl web pages.
- **Indexer Node:** Processes and indexes crawled content.
- **Web Interface:** User-friendly search interface and URL submission portal.

Getting Started

Prerequisites

- Python 3.8 or higher
- AWS account with appropriate permissions
- AWS CLI configured with credentials

Installation

1. Clone the repository:

```
git clone https://github.com/Adham-Osama11/distributed_crawler.git  
cd distributed-web-crawler
```

2. Install dependencies:

```
pip install -r requirements.txt
```

3. Configure AWS credentials:

```
aws configure
```

Configuration

1. Set up AWS resources:

- Create S3 buckets for crawl data and index storage.
- Create DynamoDB tables for metadata and analytics.
- Create SQS queues for task distribution.

2. Update configuration in `common/config.py`:

```
AWS_REGION = 'us-east-1'  
CRAWL_DATA_BUCKET = 'my-crawl-data-bucket'  
INDEX_DATA_BUCKET = 'my-index-data-bucket'
```

Running the System

1. Start the master node:

```
python src/master/master_node.py
```

2. Start crawler nodes (on different machines if desired):

```
python src/crawler/crawler_node.py
```

3. Start the indexer node:

```
python src/indexer/indexer_node.py
```

4. Launch the web interface:

```
python src/client/search_interface.py
```

Usage

Web Interface

Access the web interface at `http://localhost:5000` to:

- Search through crawled content.
- Submit new URLs for crawling.
- View system statistics.
- Monitor crawling progress.

API Endpoints

- `GET /search?q=<query>` : Search the index.
- `POST /submit-url` : Submit a URL for crawling.
- `GET /stats`: Get system statistics.
- `GET /suggest`: Get search suggestions.

Monitoring

The system provides several monitoring capabilities:

- Real-time statistics in the web interface.
- Detailed logs for each component.
- AWS CloudWatch metrics.
- DynamoDB analytics.

Development

Project Structure

```
src/
 master/          % Master node implementation
 crawler/         % Crawler node implementation
 indexer/         % Indexer node implementation
 client/          % Web interface and API
 common/          % Shared utilities and configuration
```

Adding New Features

1. Fork the repository.
2. Create a feature branch.
3. Implement your changes.
4. Submit a pull request.

License

This project is licensed under the MIT License. See the LICENSE file for details.

Contributing

Contributions are welcome! Please feel free to submit a Pull Request.

Contact

For questions and support, please open an issue in the GitHub repository.

Acknowledgments

- Scrapy for the web crawling framework.
- Whoosh for the search indexing.
- AWS for cloud infrastructure.
- Flask for the web interface.

12 Deployment Guide

This section provides a step-by-step guide to deploying and running the distributed web crawling system in a cloud environment. The system consists of the following main components:

- **Master Node:** Coordinates crawling tasks and manages the overall system.
- **Crawler Nodes:** Perform the actual web crawling.
- **Indexer Nodes:** Process and index crawled content.
- **Dashboard:** Web interface for monitoring and control.

12.1 Prerequisites

AWS Account Setup

1. Create an AWS account if you do not already have one.
2. Install and configure the AWS CLI:

```
1 pip install awscli  
2 aws configure
```

3. Required AWS Services:

- Amazon S3 (for content storage)
- Amazon SQS (for message queues)
- Amazon DynamoDB (for state management)
- Amazon EC2 (for running the nodes)

Required IAM Permissions

Create an IAM user with the following permissions:

```
1 {  
2     "Version": "2012-10-17",  
3     "Statement": [  
4         {  
5             "Effect": "Allow",  
6             "Action": [  
7                 "s3:*",  
8                 "sns:*",  
9                 "dynamodb:*",  
10                "ec2: *"  
11            ],  
12            "Resource": "*"  
13        }  
14    ]  
15}
```

Listing 12: IAM Policy JSON

Python Environment

- Python 3.8 or higher
- Required Python packages:

```
1 pip install -r requirements.txt
```

- Key dependencies: boto3, scrapy, whoosh, nltk, flask, psutil, crochet

12.2 System Configuration

Environment Variables

Create a .env file in the project root:

```
1 # AWS Configuration
2 AWS_REGION=us-east-1
3 AWS_ACCESS_KEY_ID=my_access_key
4 AWS_SECRET_ACCESS_KEY=my_secret_key
5
6 # Queue Names
7 CRAWL_TASK_QUEUE=crawl-tasks
8 CRAWL_RESULT_QUEUE=crawl-results
9 INDEX_TASK_QUEUE=index-tasks
10 MASTER_COMMAND_QUEUE=master-commands
11
12 # Storage
13 CRAWL_DATA_BUCKET=my-crawl-bucket
14 INDEX_BUCKET=my-index-bucket
15
16 # System Settings
17 HEARTBEAT_INTERVAL=30
18 MAX_RETRIES=3
19 TASK_TIMEOUT=300
20 INDEX_BATCH_SIZE=100
21
22 # Dashboard
23 DASHBOARD_HOST=0.0.0.0
24 DASHBOARD_PORT=8050
25
26 # Client
27 CLIENT_HOST=0.0.0.0
28 CLIENT_PORT=5000
```

AWS Resource Setup

S3 Buckets:

```
1 aws s3 mb s3://my-crawl-bucket
2 aws s3 mb s3://my-index-bucket
```

SQS Queues:

```
1 aws sqs create-queue --queue-name crawl-tasks
2 aws sqs create-queue --queue-name crawl-results
3 aws sqs create-queue --queue-name index-tasks
4 aws sqs create-queue --queue-name index-results
5 aws sqs create-queue --queue-name master-commands
```

DynamoDB Tables (optional manual creation):

```
1 aws dynamodb create-table \
2   --table-name url-frontier \
3   --attribute-definitions AttributeName=url,AttributeType=S \
4     AttributeName=job_id,AttributeType=S \
5   --key-schema AttributeName=url,KeyType=HASH AttributeName=job_id,
6     KeyType=RANGE \
7   --provisioned-throughput ReadCapacityUnits=5,WriteCapacityUnits=5
```

12.3 Deployment Steps

Master Node

```
1 python src/master/master_node.py
```

Crawler Node

```
1 python src/crawler/crawler_node.py
```

Indexer Node

```
1 python src/indexer/indexer_node.py
```

Dashboard

```
1 python src/monitoring/dashboard.py
```

Accessible at <http://localhost:5000>

12.4 Scaling the System

Horizontal Scaling Deploy multiple crawler and indexer nodes across EC2 instances. They register automatically with the Master Node.

Vertical Scaling Upgrade instance types or adjust service limits (e.g., DynamoDB capacity units) for improved performance.

12.5 Monitoring and Maintenance

System Monitoring

- Use the **Dashboard** at `localhost:5000`
- Use **CloudWatch** to track:
 - SQS queue length
 - DynamoDB read/write throughput
 - S3 usage

Maintenance Tasks

```
1 aws s3 sync whoosh_index/ s3://my-index-bucket/backups/$(date +%Y%m%d)/
```

Listing 13: Back up index to S3

```
1 aws s3 rm s3://my-crawl-bucket/ --recursive --exclude "* --include "*/$(date -d '30 days ago' +%Y%m%d)*"
```

Listing 14: Clean up old crawl data

```
1 from src.indexer.indexer_node import WhooshIndex  
2 WhooshIndex().optimize()
```

Listing 15: Optimize Whoosh Index

12.6 Troubleshooting

Common Issues

- **Crawler not registering:** Check credentials, SQS permissions, and network.
- **High memory usage:** Reduce concurrent requests and tune delays.
- **Slow indexing:** Optimize batch size and S3 reads.

Logs

- Master Node: `master_node.log`
- Crawler Node: `crawler_node.log`
- Indexer Node: `indexer_node.log`
- Dashboard: `dashboard.log`

12.7 Security Considerations

- Use IAM roles with least privilege
- Enable S3 encryption and bucket policies
- Use HTTPS for dashboard and secure credential storage
- Sanitize all URLs and log access activity

12.8 Performance Tuning

Crawler Settings

```
1 settings.update({  
2     'CONCURRENT_REQUESTS': 16,  
3     'DOWNLOAD_DELAY': 0.5,  
4     'RANDOMIZE_DOWNLOAD_DELAY': True,  
5     'COOKIES_ENABLED': False  
6 })
```

Indexer Settings

```
1 INDEX_BATCH_SIZE = 100
2 MAX_DOCUMENTS_PER_INDEX = 10000
3 INDEX_OPTIMIZATION_INTERVAL = 3600
```

DynamoDB Settings

```
1 ReadCapacityUnits = 10
2 WriteCapacityUnits = 10
```

13 Conclusion

In conclusion, this distributed web crawler system represents a well-balanced approach to the complex problem of web crawling at scale. It effectively combines distributed systems principles with cloud services and modern text processing techniques, resulting in a solution that is both powerful and maintainable. The clear separation of concerns across different components makes the system modular and extensible, while the comprehensive fault tolerance mechanisms ensure reliability in production environments.