



Phase 1

CSE354, Distributed Computing

Name:	Ali Tarek Abdelmonim	ID:	21P0123
Name:	Mohamed Mostafa Mamdouh	ID:	21P0244
Name:	Mohamed Walid Hemly	ID:	21P0266
Name:	Tsneam Ahmed Eliwa Zky	ID:	21P0284

Presented To: Dr Ayman Bahaa Dr Hossam Mohamed Eng. Alaa Hamdy Eng. Mostafa Ashraf

Phase No: (1)

Date: 19 / 4 /2025

Table of Contents

Introduction	5
Objective and Overview.....	5
Project's purpose.....	5
The problem that needs to be solved	5
The desired functionality	5
Team Roles and Responsibilities.....	6
Roles	6
Requirement Refinement	6
1. Distributed Crawling	6
2. Web Page Indexing	6
3. Fault Tolerance.....	6
4. Scalability.....	7
5. Politeness and robots.txt Compliance	7
6. Basic Search Functionality	7
7. User Stories	7
Technology Selection	8
List of Chosen Technologies	8
Justification	8
Programming Language: Python	8
Web Crawling Library: Beautiful Soup.....	8
Indexing/Search Library: Elasticsearch.....	8
Cloud Computing Platform: Google Cloud Platform (GCP)	8
Task Queue Service: GCP Pub/Sub	8
Persistent Storage: Google Cloud Storage.....	8
System Architecture Design.....	9
1. High-Level Architecture Diagrams.....	9
• Purpose:	9
Content:	9
2. Component Interaction	9
Purpose:	9
Content:	9
3. Data Flow.....	9

Purpose:	9
Content:	9
4. Fault Tolerance Mechanism	10
Purpose:	10
Content:	10
5. Deployment	10
Purpose:	10
Content:	10
6. Sequence for a Crawl Job	10
Purpose:	10
Content:	10
More Details on System Architecture Design	11
1. Client Layer (User Interface).	11
2. Master Node (Controller/Scheduler)	11
3. Crawler Nodes.....	11
4. Indexer Nodes	11
5. Task Queue (Google Cloud Pub/Sub)	12
6. Persistent Storage (Google Cloud Storage)	12
Data Flow Overview	12
Project Planning	13
Cloud Environment Setup	14
Set up accounts with the chosen cloud provider.....	14
Configure basic cloud resources (initial VMs for Master, Crawler, Indexer for testing, Task Queue service, Storage service).....	14
1. Service Account for Instances	15
2. Master VM.....	15
3. Crawler VMs.....	16
4. Indexer VMs.....	17
Initiate a New Web Crawl	18
Pub/Sub	18
ElasticSearch.....	19
VM instance templates	19
VM instance groups.....	20

Artifact registry repository and images	20
Cloud Storage	21
Initial Code Repository and Skeleton Code.....	21
Conclusion	23
References.....	23

Introduction

The advancement of technology has led to an exponential growth in the volume of online information, creating both opportunities and challenges in accessing and organizing data efficiently. To address this, our project focuses on designing and developing a Distributed Web Crawling and Indexing System using Cloud Computing.

This project aims to leverage distributed computing to perform large-scale web crawling and indexing tasks. By utilizing cloud-based virtual machines and distributed task queues, the system will efficiently crawl websites, process their contents, and create a searchable index of web pages. Emphasis is placed on ensuring fault tolerance and scalability, allowing the system to maintain functionality despite node failures and to adapt to increasing workloads.

In Phase 1, we establish the foundation of this project by refining the requirements, finalizing the technology stack, designing the system architecture, and setting up the cloud environment. This phase ensures a clear direction and robust framework for implementing the system in subsequent phases.

The following report outlines the objectives, selected technologies, detailed system architecture, and project plan for Phase 1, laying the groundwork for a scalable and resilient distributed system.

Objective and Overview

Project's purpose

The purpose of this project is to design and implement a Distributed Web Crawling and Indexing System that utilizes cloud computing technologies. By leveraging the power of distributed systems and virtualized resources, the aim is to create a scalable and efficient solution for processing and organizing vast amounts of online information.

The problem that needs to be solved

The problem being addressed stems from the rapid growth of data available on the web, which presents challenges in accessing, indexing, and managing such large volumes of information. The system is intended to be both scalable and fault-tolerant, ensuring reliability even when parts of the system experience failure. This will allow for uninterrupted operation and an adaptable infrastructure that can handle increasing workloads.

The desired functionality

The system's functionality revolves around distributed web crawling, content processing, and indexing. Crawling involves fetching and analysing web pages across multiple virtual machines, while processing transforms the extracted content into meaningful data structures. Finally, the indexing module organizes this information into a searchable format, enabling users to efficiently query and retrieve relevant web pages.

Team Roles and Responsibilities

Roles

- **Architect:** Designs the system architecture, focusing on interactions and fault tolerance strategies, this role is assigned to **Ali Tarek**
- **Crawler Lead:** Implements the distributed web crawling functionality and ensures politeness measures, this role is assigned to **Tsneam Ahmed**
- **Indexer Lead:** Develops the indexing mechanism and search functionality, this role is assigned to **Mohamed Mostafa**
- **Cloud Infrastructure Lead:** Configures and manages cloud computing resources, this role is assigned to **Mohamed Mostafa**
- **Tester/Documentation Lead:** Tests system components and ensures thorough documentation, this role is assigned to **Mohamed Walid**

Requirement Refinement

Requirement Refinement is a crucial step to ensure the project's objectives, scope, and functionalities are clearly defined and well understood. This stage involves interpreting the project specifications, identifying key requirements, and refining them for implementation. Below are the refined requirements for the Distributed Web Crawling and Indexing System using Cloud Computing.

1. Distributed Crawling

The system must efficiently distribute crawling tasks across multiple virtual machines hosted in a cloud environment. This requirement ensures scalability, as the system can process large volumes of URLs simultaneously, with the ability to add more virtual machines as needed.

Key points include:

- Assigning crawling tasks dynamically to available crawler nodes.
- Implementing mechanisms for fault tolerance to reassign tasks if a crawler node fails.
- Respecting website policies through adherence to robots.txt directives.

2. Web Page Indexing

The system must implement an indexing mechanism to process crawled web pages and create a searchable index. This includes:

- Extracting meaningful content such as text, keywords, and metadata from crawled pages.
- Structuring the data in an organized, efficient manner for fast querying.
- Ensuring scalability and reliability by utilizing cloud-based storage solutions.

3. Fault Tolerance

The system must be resilient to failures, with robust mechanisms to maintain continuous operation despite node failures. Specific strategies include:

- **Crawler Node Failure:** Monitoring crawler nodes and reassigning tasks from failed nodes to operational ones.
- **Indexer Node Failure:** Implementing strategies to ensure index integrity and recovering from node failures without data loss.
- **Data Persistence:** Storing crawled data and indexes in durable cloud storage to avoid data loss.

4. Scalability

The system must be designed to scale effortlessly as the workload increases. This involves:

- Adding more crawlers or indexer nodes to the system.
- Ensuring efficient distribution of tasks and resources.
- Using cloud-based services for dynamic scalability.

5. Politeness and robots.txt Compliance

The crawler nodes must adhere to robots.txt directives specified by websites, ensuring ethical and polite crawling practices. This includes implementing crawl delays to avoid overloading servers.

6. Basic Search Functionality

The system must provide users with a simple search interface to query the indexed content. Features include:

- Keyword-based search with relevance ranking.
- Returning URLs and metadata of matching web pages.

7. User Stories

User stories help define the functional requirements from a user's perspective. include:

- "As a user, I want to initiate a web crawl by providing a list of seed URLs."
- "As a user, I want to specify crawling parameters such as depth limit or domain restrictions."
- "As a user, I want the system to continue crawling even if some crawler nodes fail."
- "As a user, I want to search for keywords and retrieve relevant web pages from the indexed content."
- "As a user, I want to monitor the progress of the web crawling and indexing process."

Technology Selection

List of Chosen Technologies

- **Programming language:** Python (rich libraries for web crawling, indexing, and distributed computing).
- Web crawling libraries: Beautiful Soup
- Indexing/search libraries: Elasticsearch
- Cloud computing platform: GCP
- Task queue service: GCP Pub/Sub
- **Persistent storage:** Google Cloud Storage

Justification

Programming Language: Python

Python is an excellent choice for this project due to its versatility and a rich ecosystem of libraries that simplify tasks like web crawling, indexing, and distributed computing. Its clean syntax and extensive documentation allow for faster development and easier debugging. Python's compatibility with numerous frameworks and cloud-based APIs makes it ideal for implementing distributed systems efficiently.

Web Crawling Library: Beautiful Soup

Beautiful Soup is a powerful library for parsing HTML and XML documents. It allows easy navigation and extraction of data from web pages. It's lightweight, making it suitable for handling the parsing needs of our system, such as extracting text and URLs from web pages, and works seamlessly with other Python libraries for additional functionalities.

Indexing/Search Library: Elasticsearch

Elasticsearch is highly scalable and provides fast full-text search capabilities, which makes it a perfect fit for indexing and searching large datasets. Its distributed nature aligns with our system's design, ensuring fault tolerance and high availability. Elasticsearch also supports advanced querying, such as relevance ranking and aggregations, which enhances the search functionality of your system.

Cloud Computing Platform: Google Cloud Platform (GCP)

GCP provides a wide range of cloud services and is well-suited for distributed applications. It offers cost-effective virtual machines, robust networking capabilities, and scalable storage solutions, making it an optimal choice for deploying your system. GCP's integration with Pub/Sub for message queuing and its compatibility with Terraform further streamline the development process.

Task Queue Service: GCP Pub/Sub

GCP Pub/Sub: A cloud-based messaging service designed for scalability and reliability. It ensures efficient task distribution and provides strong fault tolerance guarantees. Using Pub/Sub for inter-component communication simplifies the orchestration of crawling and indexing tasks.

Persistent Storage: Google Cloud Storage

Google Cloud Storage offers a scalable and durable solution for storing crawled data and indexed content. Its high availability and multi-region storage options ensure data persistence and reliability. It integrates seamlessly with other GCP services, streamlining the data pipeline from crawling to storage and indexing.

System Architecture Design

1. High-Level Architecture Diagrams

Purpose: Provide a clear overview of the entire system and its major components.

Content:

- Include components: Client Layer (UI), Master Node, Crawler Nodes, Indexer Nodes, Task Queue (Pub/Sub), and Persistent Storage (GCS).
- Show data flow:
 - User submits crawl jobs to Master Node via UI.
 - Master Node assigns tasks to Crawler Nodes via Pub/Sub.
 - Crawler Nodes send processed content to Pub/Sub for Indexer Nodes.
 - Indexer Nodes store indexed data and serve queries from UI.

We Used arrows to indicate communication and data exchange.

2. Component Interaction

Purpose: Illustrate how system components interact with each other during the crawling and indexing process.

Content:

- Depict interactions:
 - UI submits crawl jobs to Master Node.
 - Master Node publishes tasks to Pub/Sub, which Crawler Nodes subscribe to.
 - Crawler Nodes fetch URLs, process data, and publish results to Pub/Sub for indexing.
 - Indexer Nodes subscribe to processed data, build indexes, and store them in GCS.
 - UI queries indexed data via Indexer Nodes.
 - Include key APIs and communication protocols.

3. Data Flow

Purpose: Show the flow of data between components throughout the system.

Content:

- Seed URLs and parameters flow from UI to Master Node.
- Master Node breaks tasks and publishes to Pub/Sub for Crawler Nodes.
- Crawler Nodes process data and publish it back to Pub/Sub for Indexer Nodes.
- Indexer Nodes build indexes and store them in GCS.
- Indexed data serves search queries from the UI.

Highlight different types of data (raw HTML, processed text, indexed content).

4. Fault Tolerance Mechanism

Purpose: Highlight strategies for handling component failures.

Content:

- Show:
 - Task reassignment from failed Crawler Nodes by the Master Node.
 - Data replication and recovery for Indexer Nodes.
 - Reliable messaging in Pub/Sub.
 - Durable storage in GCS for persistence.
 - Illustrate heartbeat monitoring and task timeout mechanisms.

5. Deployment

Purpose: Visualize the system's deployment in the cloud.

Content:

- Include cloud resources:
 - VM instances for Master Node, Crawler Nodes, and Indexer Nodes.
 - GCP Pub/Sub for messaging.
 - GCS for storage.
- Specify configuration details (e.g., VM types, network setup).

6. Sequence for a Crawl Job

Purpose: Detail the steps involved in initiating and executing a crawl job.

Content:

- Show:
 - UI sending job request to Master Node.
 - Master Node assigning tasks to Crawler Nodes.
 - Crawler Nodes processing URLs and sending results to Indexer Nodes.
 - Indexer Nodes building the index and storing it in GCS.
 - UI querying indexed data from Indexer Nodes.

More Details on System Architecture Design

1. Client Layer (User Interface)

- **Purpose:** Allows users to initiate web crawling jobs, set parameters, view progress, and search indexed data.
- Features:
 - Input seed URLs, crawl depth limit, and domain restrictions.
 - Submit queries to search indexed content.
 - Monitor progress of crawling and indexing.
 - Technology:
 - Web application built with Flask for backend and HTML/TailwindCSS for frontend.
 - Publishes crawl jobs to the UI-Master topic via GCP Pub/Sub, which the master node subscribes to.

2. Master Node (Controller/Scheduler)

- **Purpose:** Manages task distribution and system orchestration.
- Responsibilities:
 - Divide seed URLs into tasks and assign them to Crawler Nodes.
 - Monitor worker nodes and reassign tasks if failures occur.
- Publish crawl tasks to the **CRAWL_TASKS_TOPIC_ID** and subscribe to the **NEW_CRAWL_JOB_SUBSCRIPTION_ID**.
- Technology:
 - Python with Google Cloud Pub/Sub APIs for task distribution.
 - Reads seed URLs from Google Cloud Storage and handles dynamic job submissions.

3. Crawler Nodes

- Purpose: Perform web crawling tasks.
- Responsibilities:
 - Fetch web pages using **requests**.
 - Parse HTML with Beautiful Soup and extract text content and new URLs.
 - Save raw HTML and processed text to Google Cloud Storage.
 - Publish new URLs to UI-Master topic and processed content to Crawler-Indexer topic for indexing.
- Technology:
 - Python with libraries like Beautiful Soup and requests.
 - Utilize GCP Pub/Sub for task communication and GCS for storage.

4. Indexer Nodes

- **Purpose:** Build and manage the searchable index.
- Responsibilities:
 - Subscribe to Pub/Sub for processed content messages from Crawler Nodes.
 - Download processed text from Google Cloud Storage.
 - Create a robust index using Elasticsearch for search functionality.
 - Handle queries from the Client Layer.
- Technology:

- Python and Elasticsearch for indexing and search.
- Persistent index storage in Google Cloud Storage.

5. Task Queue (Google Cloud Pub/Sub)

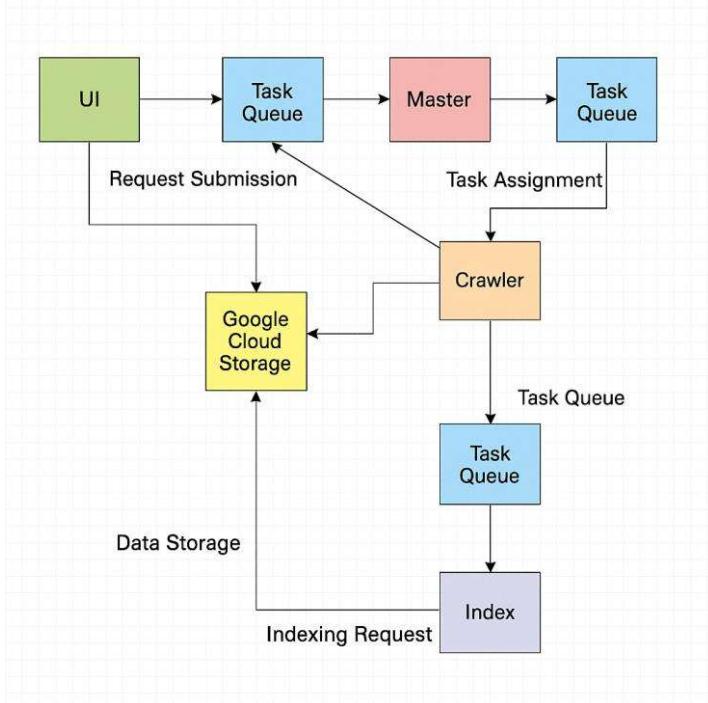
- **Purpose:** Facilitates communication between system components.
- Responsibilities:
 - Master Node publishes crawl tasks to **Master-Crawler topic**.
 - Crawler Nodes publish processed data to **Crawler-Indexer topic**.
 - Reliable delivery ensures fault tolerance.

6. Persistent Storage (Google Cloud Storage)

- **Purpose:** Durable and scalable storage for data.
- Responsibilities:
 - Store raw HTML, extracted text, and indexed data.
 - Ensure persistence for system fault tolerance and recovery.

Data Flow Overview

1. Job Submission:
 - Client sends crawl jobs with parameters via Pub/Sub.
 - Master Node processes and assigns tasks to Crawler Nodes.
 -
2. Crawling:
 - Crawler Nodes fetch pages, extract data, and store results in GCS.
 - Publish new URLs and processed content to Pub/Sub for further processing.
3. Indexing:
 - Indexer Nodes ingest processed content from GCS and build the index.
4. Querying:
 - Client queries indexed data via the Indexer Nodes.



Project Planning

Detailed Project Plan

Task	1	2	3	4
Team Roles & Tech Stack Finalized				
System Architecture + Diagrams				
Cloud VM Setup				
Git Repo & Skeleton Code				
Master Node Implementation				
Crawler Node Implementation				
Indexer & Search Dev				
Fault Tolerance (Crawler + Indexer)				
Logging & Monitoring				

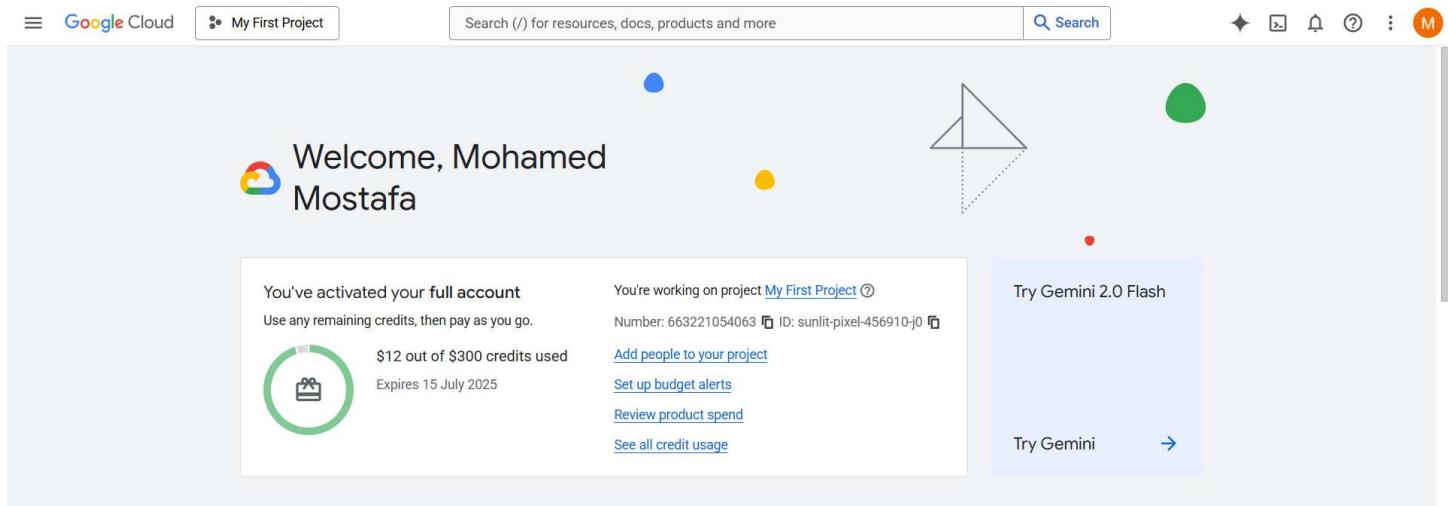
Cloud Environment Setup

Set up accounts with the chosen cloud provider.

Google Cloud Platform (GCP) under the project: Distributed Project.

Project ID: sunlit-pixel-456910-j0

Project Number: 663221054063



Configure basic cloud resources (initial VMs for Master, Crawler, Indexer for testing, Task Queue service, Storage service).

1. Service Account for Instances

The screenshot shows the Google Cloud IAM and admin interface. On the left, a sidebar lists various IAM components: IAM, PAM, Principal access bound..., Organisations, Identity and organisation, Policy troubleshooter, Policy analyser, Organisation policies, Service accounts (which is selected), Workload Identity Fede..., Workforce identity fede..., Labels, Tags, Settings, Privacy and security, and Identity-Aware Proxy. The main panel displays the 'ui-client' service account details. It includes fields for Name (ui-client) and Description, both with 'Save' buttons. Below these are Email (ui-client@sunlit-pixel-456910-j0.iam.gserviceaccount.com) and Unique ID (113019955715541366290). Under 'Service account status', it shows that the account is Enabled and has a 'Disable service account' button. A 'Advanced settings' section is partially visible at the bottom.

2. Master VM

The screenshot shows the Google Cloud Compute Engine interface. On the left, a sidebar lists categories: Overview, Virtual machines (selected), VM instances (selected), Instance templates, Sole-tenant nodes, Machine images, TPUs, Committed-use discou..., Reservations, Migrate to Virtual Mach..., Storage, Instance groups (selected), Instance groups, Health checks, VM Manager (selected), Patch, Marketplace, and Release notes. The main panel shows the 'master-node' VM instance details. It includes tabs for Details, Observability, OS info, and Screenshot. Under Details, there are SSH and Connect to serial console buttons. A note says 'Connecting to serial ports is disabled'. The Logs section shows Logging and Serial port 1 (console). The Basic information table provides the following details:

Name	master-node
Instance ID	3748868954772195615
Description	None
Type	Instance
Status	Stopped
Creation time	Apr 18, 2025, 6:04:02 pm UTC+02:00
Location	us-central1-c
Instance template	None
In use by	None
Physical host	None
Maintenance status	—
Reservations	Automatically chosen

At the bottom, there is a 'Equivalent code' button.

3. Crawler VMs

The screenshot shows the 'crawler-new' instance details page in the Google Cloud Compute Engine interface. The left sidebar contains navigation links for Compute Engine, Overview, Virtual machines, VM instances, Instance templates (which is selected), Sole-tenant nodes, Machine images, TPUs, Committed-use discounts, Reservations, Migrate to Virtual Machine, Storage, Instance groups (selected), Instance groups, Health checks, VM Manager (selected), Patch, Marketplace, and Release notes. The main content area has tabs for Basic information, Machine configuration, and Networking. The 'Basic information' tab displays the following details:

Name	crawler-new
Type	Instance template
Creation time	Apr 18, 2025, 7:05:49 pm UTC+02:00
In use by	None
Location	us-central1
Reservations	Automatically choose
Labels	container... : cos-stable...
Tags	—
Placement policy	No policy
Confidential VM service	Disabled

The 'Machine configuration' tab displays the following details:

Machine type	e2-medium
Minimum CPU platform	None
Architecture	—
vCPUs to core ratio	—
Custom visible cores	—
All-core turbo-only mode	—
Display device	Disabled
GPUs	None

The 'Networking' tab displays the following details:

Public DNS PTR record	None
Total egress bandwidth tier	—

4. Indexer VMs

Compute Engine

Overview

Virtual machines

- VM instances
- Instance templates**
- Sole-tenant nodes
- Machine images
- TPUs
- Committed-use discou...
- Reservations
- Migrate to Virtual Mach...

Storage

Instance groups

- Instance groups
- Health checks

VM Manager

- Patch
- Marketplace
- Release notes

← indexer-new + Create VM ⚡ Create similar ⚡ Create instance group 🗑 Delete

Basic information

Name	indexer-new
Type	Instance template
Creation time	Apr 18, 2025, 7:06:40 pm UTC+02:00
In use by	None
Location	us-central1
Reservations	Automatically choose
Labels	container-... : cos-stable...
Tags ⓘ	—
Placement policy	No policy ⓘ
Confidential VM service ⓘ	Disabled

Machine configuration

Machine type	e2-medium
Minimum CPU platform	None
Architecture	—
vCPUs to core ratio	—
Custom visible cores ⓘ	—
All-core turbo-only mode ⓘ	—
Display device	Disabled
GPUs	None

Networking

Public DNS PTR record	None
Total egress bandwidth tier	—

Initiate a New Web Crawl

Success: Crawl job submitted. Task ID: dc82472c-a183-4393-b208-32e549ecd979

Initiate a New Web Crawl

Seed URLs (at least one required)

https://example.com



Crawl Depth Limit (Optional)

e.g., 3 (0 or empty means crawl only seed URLs)

Maximum link depth from seed URLs. Leave empty or 0 to only process seeds.

Restrict Crawl to Domain (Optional)

e.g., example.com

If set, the crawler will only follow links within this domain.

Start Crawl Job

Crawl Configuration Submitted:

Seed URLs:

- https://www.youtube.com/

Depth Limit:

Domain Restriction: None

Pub/Sub

Pub/Sub / Topics

Pub/Sub	Topics	+ CREATE TOPIC	DELETE	SHOW INFO PANEL	LEARN
Topics	LIST METRICS				
Subscriptions					
Snapshots					
Schemas					
Pub/Sub Lite					
Lite reservations					
Lite topics					
Lite Subscriptions					

Filter Filter topics

Topic ID ↑	Encryption key	Topic name	Retention	Ingestion source	⋮
Crawler-Indexer	Google-managed	projects/sunlit-pixel-456910-j0/topics/Crawler-Indexer	—	—	⋮
Master-Crawler	Google-managed	projects/sunlit-pixel-456910-j0/topics/Master-Crawler	—	—	⋮
UI-Master	Google-managed	projects/sunlit-pixel-456910-j0/topics/UI-Master	—	—	⋮

ElasticSearch

The screenshot shows the Elasticsearch interface with the 'Elasticsearch' tab selected. The left sidebar includes links for Overview, Dev Tools, Kibana, Discover, Dashboards, Content, Index Management (which is selected), Connectors, Web Crawlers, Build, Playground, Search applications, Relevance, Inference Endpoints, Synonyms, and Other tools. The main area displays search results for three documents:

- Document id: https://app.flowcv.com/**

Field	Contents
* content	"Google Maps When you have eliminated the JavaScript , whatever remains must be an empty page. Enable JavaScr:"
* url	"https://app.flowcv.com/"
- Document id: https://www.google.com:443/maps**

Field	Contents
* content	"Google Maps When you have eliminated the JavaScript , whatever remains must be an empty page. Enable JavaScr:"
* url	"https://www.google.com:443/maps"
- Document id: https://www.youtube.com/**

Field	Contents
* content	"YouTube About Press Copyright Contact us Creators Advertise Developers Terms Privacy Policy & Safety How You'...
* url	"https://www.youtube.com/"
- Document id: https://www.plex.tv/**

Field	Contents
* content	"Stream Movies & Find Shows On The Best Streaming Services Changing the site language for you. Please wait .."
* url	"https://www.plex.tv/"

VM instance templates

The screenshot shows the Google Cloud Compute Engine interface with the 'My First Project' selected. The left sidebar includes links for Compute Engine (selected), Overview, Virtual machines, VM instances, Instance templates (which is selected), Sole-tenant nodes, Machine images, TPUs, Committed-use discounts, Reservations, Migrate to Virtual Machines, Storage, and Instance groups. The main area displays the 'Instance templates' page with the following information:

Instance templates are saved VM configurations used to create identical VMs, either individually or as part of managed instance groups. [Learn more](#)

Filter instance templates:

Name	Machine type	Image	Disk type	Location	Placement policy	In use by	Creation time	Actions
crawler-new	e2-medium	cos-stable-117-18613-164109	Balanced persistent disk	us-central1	No policy	crawler-mig	Apr 18, 2025, 7:05:49 pm UTC+02:00	⋮
indexer-new	e2-medium	cos-stable-117-18613-164109	Balanced persistent disk	us-central1	No policy	indexer-mig	Apr 18, 2025, 7:06:40 pm UTC+02:00	⋮

VM instance groups

The screenshot shows the Google Cloud Compute Engine Instance Groups page. On the left, there's a sidebar with navigation links like Overview, VM instances, Instance templates, Sole-tenant nodes, Machine images, TPUs, Committed-use discounts, Reservations, Migrate to Virtual Machines, Storage, and Instance groups. Under Instance groups, 'Instance groups' is selected. The main area has tabs for 'Instance groups' (selected), 'Create instance group', 'Refresh', and 'Delete'. A search bar at the top right shows 'topics' and a 'Search' button. Below the tabs, a brief description of instance groups is provided, followed by a table listing two entries:

Status	Name	Instances	Template	Group type	Creation time	Recommendation	Auto-scaling	Zone	In Use By
<input checked="" type="checkbox"/>	crawler-mig	1	crawler-new (Regional)	Managed	Apr 18, 2025, 7:07:42 pm UTC+02:00	On: Target CPU utilisation 80%		us-central1-c	
<input checked="" type="checkbox"/>	indexer-mig	1	indexer-new (Regional)	Managed	Apr 18, 2025, 7:08:20 pm UTC+02:00	On: Target CPU utilisation 80%		us-central1-c	

Artifact registry repository and images

The screenshot shows the Google Cloud Artifact Registry Repositories page. The sidebar includes 'Artifact Registry' (selected), 'Repositories' (selected), and 'Settings'. The main area shows a breadcrumb path: 'me-west1-docker.pkg.dev > sunlit-pixel-456910-j0 > master-node'. The title is 'Images for master-node'. There are buttons for 'DELETE', 'EDIT REPOSITORY', and 'SETUP INSTRUCTIONS'. Below the title, 'Repository details' are listed: Format Docker and Type Standard. A 'SHOW MORE' link is visible. A 'Filter' input field is present. A table lists three images:

<input checked="" type="checkbox"/>	Name	Connection	Created	Updated
<input checked="" type="checkbox"/>	crawler-node	—	2 days ago	1 hour ago
<input checked="" type="checkbox"/>	indexer-node	—	2 days ago	1 hour ago
<input checked="" type="checkbox"/>	master-node	—	2 days ago	1 hour ago

Cloud Storage

The screenshot shows the Google Cloud Storage interface. On the left, a sidebar navigation includes 'Overview', 'Buckets' (selected), 'Monitoring', 'Settings', 'Storage Intelligence', 'Insights datasets', and 'Configuration'. The main area displays 'Bucket details' for 'webcrawlingproject'. Key information shown includes 'Location: me-west1 (Tel Aviv)', 'Storage class: Standard', 'Public access: Not public', and 'Protection: Soft delete'. Below this, tabs for 'Objects', 'Configuration', 'Permissions', 'Protection', 'Lifecycle', 'Observability', 'New', 'Inventory Reports', and 'Operations' are visible. A 'Folder browser' section shows the directory structure: 'webcrawlingproject' containing 'crawl_tasks/', 'processed_text/', 'raw_html/', and 'seeds/'. A table lists these items with columns for Name, Size, Type, Created, Storage class, Last modified, and Public ac.

Initial Code Repository and Skeleton Code

The screenshot shows a GitHub repository page for 'Distributed-Web-Crawling-and-Indexing-System-GCP'. The repository is private. It has 1 branch and 0 tags. The main branch is 'main'. The repository was last updated 5 days ago by user 'mohamedmostafam0'. There are 10 commits. The repository contains files like 'src', 'terraform', 'CSE 354 Project.pdf', 'CSE354 Distributed Computing Project Deliver...', and 'README.md'. The 'README' file is currently selected. The content of the README file is:

```
Distributed-Web-Crawling-and-Indexing-System-GCP
```

```
PROJECT/
├── README.md
├── CSE 354 Project.pdf
└── CSE354 Distributed Computing Project Deliverables.pdf

├── Documentation/

└── src/
    ├── scripts/
    │   ├── .env
    │   ├── crawler_node.py
    │   ├── indexer_node.py
    │   ├── master_node.py
    │   └── requirements.txt

    └── UI/
        ├── main.py
        ├── templates/
        └── index.html

└── terraform/
    ├── Scripts/
    │   ├── crawler_bootstrap.sh
    │   ├── indexer_bootstrap.sh
    │   └── master_bootstrap.sh

    ├── compute.tf
    ├── main.tf
    ├── monitoring.tf
    ├── network.tf
    ├── outputs.tf
    ├── provider.tf
    ├── PubSub.tf
    ├── storage.tf
    ├── terraform.tfvars
    └── variables.tf
```

Conclusion

Conclusion In Phase 1 of this project, significant progress has been made in defining the objectives, requirements, and foundational elements essential for successful implementation. We have outlined the problem to be solved, the project's purpose, and the desired functionalities, ensuring clarity in our goals. Detailed roles and responsibilities have been assigned to the team to ensure efficient collaboration. By selecting a robust technology stack, including tools like Python, BeautifulSoup, Elasticsearch, and Google Cloud Platform, the project is positioned for high performance, scalability, and fault tolerance. Furthermore, the system architecture has been meticulously designed to support distributed crawling, indexing, and deployment. These efforts lay a strong foundation for the subsequent phases, ensuring that we remain on track to achieve the project's objectives effectively and innovatively.

References

https://cloud-based_web_crawler_architecture_cloud_lab_ucm.pdf

<https://cloud.google.com/architecture/framework/reliability/build-highly-available-systems>

<https://www.robotstxt.org/>

<https://www.crummy.com/software/BeautifulSoup/>

<https://cloud.google.com/?hl=en>

<https://developer.hashicorp.com/terraform>

<https://cloud.google.com/pubsub?hl=en>

<https://cloud.google.com/storage?hl=en>

<https://cloud.google.com/appengine?hl=en>