



Ain Shams University - Faculty of Engineering

CESS

CSE351: Distributed Systems –Spring 25 Phase 3

Ali Mohamed Refaat	21P0105
Ezz Eldin Alaa Eldin Omar	21P0100
Karim Sherif Louis	21P0223
Omar Mohamed Korkor	21P0065

Table of Contents

Table of Contents	2
1 Introduction	3
2 Indexer	3
3 Crawler.....	6
3.1 Scalability.....	6
3.2 Crawler Implementation Changes.....	7
4 Master Node	10
5 MongoDB Persistence	14
6 Output Logs	15
6.1 Master Node starting.....	15
6.2 Crawler with no tasks fails & recovery	16
6.3 Dispatching tasks.....	17
6.4 Reassigning tasks.....	17
6.5 Crawler logs	17
6.6 Indexer Start	18
6.7 Indexer search	18

1 Introduction

In this phase of development, our focus was on enhancing the scalability, reliability, and performance of the distributed web crawler system. We significantly improved the indexing pipeline by integrating Elasticsearch, enabling faster and more powerful search capabilities. To validate the system's ability to scale, we deployed multiple crawler instances using Kubernetes, demonstrating a proof of concept for horizontal scalability. Persistent data storage was established through MongoDB, ensuring robustness in task tracking and fault recovery. Furthermore, we enhanced the threading model for improved concurrency and implemented thread safety mechanisms along with fault tolerance features to gracefully handle crawler node failures.

	Name	Instance ID	Instance state	Instance type	Status check	Alarm status	Availability Zone	Public IP
<input type="checkbox"/>	master-node	i-042d59af52dabd3e5	Running	t2.micro	2/2 checks passed	View alarms +	us-east-1a	ec2-54-21
<input type="checkbox"/>	crawler2	i-032ec896e1c2d8fb6	Stopped	t2.micro	-	View alarms +	us-east-1a	-
<input type="checkbox"/>	indexer-node	i-0291dbfdc14a370f5	Running	t2.micro	2/2 checks passed	View alarms +	us-east-1a	ec2-3-80
<input checked="" type="checkbox"/>	elastic-search-...	i-0d4ff1827d82952ea	Running	t2.micro	2/2 checks passed	View alarms +	us-east-1a	ec2-44-21
<input type="checkbox"/>	rabbitmq-server	i-05902baf535c5851a	Running	t2.micro	2/2 checks passed	View alarms +	us-east-1a	ec2-34-21
<input type="checkbox"/>	basic-crawler-...	i-08d2275a512253500	Running	t2.micro	2/2 checks passed	View alarms +	us-east-1a	ec2-44-21
<input type="checkbox"/>	crawler pods	i-0ecbc8a8b81a27f55	Running	t2.micro	2/2 checks passed	View alarms +	us-east-1a	ec2-18-21
<input type="checkbox"/>	mongo-db	i-0ca2f1f74e15f671c	Running	t2.micro	2/2 checks passed	View alarms +	us-east-1a	ec2-54-11

2 Indexer

To improve the capability of indexing, we deployed elastic search on an ec2 as a service.

```
us-east-1.console.aws.amazon.com/ec2-instance-connect/ssh/home?region=us-east-1&connType=standard&instanceId=i-0d4ff1827d82952ea&osUser=ec2-user&sshPort=22&addressF...
google developer inter... | YouTube | Gmail
aws | Search | [Alt+S] | United States (N. Virginia) | ezz21
Last login: Thu May 1 12:00:22 2025 from 197.37.178.13
[ec2-user@ip-172-31-92-167 ~]$ sudo systemctl status elasticsearch.service
● elasticsearch.service - Elasticsearch
   Loaded: loaded (/usr/lib/systemd/system/elasticsearch.service; enabled; preset: disabled)
   Active: active (running) since Fri 2025-05-02 12:32:52 UTC; 54min ago
     Docs: https://www.elastic.co
    Main PID: 2148 (java)
       Tasks: 75 (limit: 1111)
      Memory: 679.6M
         CPU: 55.374s
    CGroup: /system.slice/elasticsearch.service
            └─2148 /usr/share/elasticsearch/jdk/bin/java -Xms4m -Xmx64m -XX:+UseSerialGC -Dcni.name=server -Dcni.script=/usr/share/elasticsearch/bin/elasticsearch -D
              2285 /usr/share/elasticsearch/jdk/bin/java -Des.networkaddress.cache.ttl=60 -Des.networkaddress.cache.negative.ttl=10 -Djava.security.manager=allow -XX
              2304 /usr/share/elasticsearch/modules/x-pack-ml/platform/linux-x86_64/bin/controller

May 02 12:32:03 ip-172-31-92-167.ec2.internal systemd[1]: Starting elasticsearch.service - Elasticsearch...
May 02 12:32:13 ip-172-31-92-167.ec2.internal systemd-entrypoint[2285]: CompileCommand: exclude org/apache/lucene/util/MSBRadixSorter.computeCommonPrefixLengthAndBu
May 02 12:32:13 ip-172-31-92-167.ec2.internal systemd-entrypoint[2285]: CompileCommand: exclude org/apache/lucene/util/RadixSelector.computeCommonPrefixLengthAndBu
May 02 12:32:16 ip-172-31-92-167.ec2.internal systemd-entrypoint[2148]: May 02, 2025 12:32:16 PM sun.util.locale.provider.LocaleProviderAdapter <clinit>
May 02 12:32:16 ip-172-31-92-167.ec2.internal systemd-entrypoint[2148]: WARNING: COMPAT locale provider will be removed in a future release
May 02 12:32:52 ip-172-31-92-167.ec2.internal systemd[1]: Started elasticsearch.service - Elasticsearch.
lines 1-19/19 (END)

i-0d4ff1827d82952ea (elastic-search-instance)
PublicIPs: 44.212.8.36 PrivateIPs: 172.31.92.167
```

Then we integrated that service into our indexer node. Ensuring that elastic search is configured to use stemming and filter stop words for better indexing of the parsed content. Elastic search receives requests in https and needs to be authenticated with a username and password. But for testing we disabled verifying certificates for ease of use.

```
indexer > indexer.py > ElasticIndexer > _init_
13 class ElasticIndexer:
14     def __init__(self, index_name="webpages"):
15         self.index_name = index_name
16         self.es = Elasticsearch(
17             hosts=[{"host": "172.31.92.167", "port": 9200, "scheme": "https"}],
18             basic_auth=("elastic", "sl_ZN64F4UN4uoDw5AIC"),
19             verify_certs=False,
20         )
21         try:
22             print(self.es.info())
23         except Exception as e:
24             print("Failed to connect:", e)
25
26         if not self.es.ping():
27             raise RuntimeError(" Elasticsearch ping failed. Check URL/auth.")
28         logging.info("Elasticsearch cluster is up!")
29
30         settings = {
31             "settings": {
32                 "analysis": {
33                     "analyzer": {
34                         "custom_english_analyzer": {
35                             "type": "custom",
36                             "tokenizer": "standard",
37                             "filter": ["lowercase", "english_stop", "english_stemmer"],
38                         }
39                     },
40                     "filter": {
41                         "english_stop": {"type": "stop", "stopwords": "_english_"},
42                         "english_stemmer": {"type": "stemmer", "language": "english"},
43                     }
44                 }
45             },
46             "index_patterns": ["*"],
47             "type": "webpages"
48         }
49
50         self.es.indices.create(index=self.index_name, body=settings)
51
52         raise
53
54     def add_document(self, url, text):
55         self.es.index(index=self.index_name, document={"url": url, "text": text})
56         logging.info(f"Document indexed: {url}")
57
58     def search(self, keyword):
59         query = {"query": {"match": {"text": keyword}}}
60         response = self.es.search(index=self.index_name, body=query)
61         hits = response.get("hits", {}).get("hits", [])
62         return [hit["_source"]["url"] for hit in hits]
```

It also offers interfaces for adding documents and searching the deployed elastic search.

The rest of our indexer node has the same logic and capabilities as before, receiving messages through the SQS queue from the crawlers by polling the queue.

```
82
83 # IndexerNode class using SQS and ElasticIndexer
84 class IndexerNode:
85     def __init__(self, queue_url):
86         self.queue_url = queue_url
87         self.sqs = boto3.client("sqs", region_name="us-east-1")
88         self.indexer = ElasticIndexer()
89
90     def poll_messages(self):
91         while True:
92             try:
93                 response = self.sqs.receive_message(
94                     QueueUrl=self.queue_url,
95                     MaxNumberOfMessages=5,
96                     WaitTimeSeconds=10,
97                     MessageAttributeNames=["All"],
98                 )
99                 messages = response.get("Messages", [])
100                 if messages:
101                     logging.info(f"Received {len(messages)} messages from SQS")
102                     for msg in messages:
103                         body = json.loads(msg["Body"])
104                         self.indexer.add_document(body["url"], body["text"])
105                         logging.info(f"Indexed document from URL: {body['url']}")
106                         self.sqs.delete_message(
107                             QueueUrl=self.queue_url, ReceiptHandle=msg["ReceiptHandle"]
108                         )
109                         logging.info(f"Deleted message from SQS for URL: {body['url']}")
110             except Exception as e:
111                 logging.error(f"Error while polling messages: {e}")
112                 time.sleep(1)
113
```

```
def start(self):
    logging.info("IndexerNode started. Polling for messages...")
    threading.Thread(target=self.poll_messages, daemon=True).start()

def cli(self):
    while True:
        keyword = input("Enter keyword to search (or 'exit'): ")
        if keyword.lower() == "exit":
            logging.info("Exiting CLI...")
            break
        results = self.indexer.search(keyword)
        if results:
            logging.info(f"Search results for '{keyword}': {results}")
        else:
            logging.info(f"No results found for keyword: {keyword}")

if __name__ == "__main__":
    queue_url = (
        "https://sqs.us-east-1.amazonaws.com/022499012946/crawler-indexer-sqs.fifo"
    )
    node = IndexerNode(queue_url)
    node.start()
    node.cli()
```

3 Crawler

We had a lot of progress both in crawler logic and scalability concerns.

3.1 Scalability

We created another ec2 instance for crawler deployment, but we went on our initial plan of launching a crawler cluster per ec2. The free tier machines that amazon offers aren't computationally powerful so to preserve storage and processing power we deployed k3s, a lightweight distro of Kubernetes. Then we created an image for our crawler .

```
FROM python:3.10-slim
WORKDIR /app
COPY . /app
RUN pip install --no-cache-dir -r requirements.txt
CMD ["python", "crawler.py"]
```

And created a deployment and launched it so crawlers can be easily scaled in this ec2.

```
A newer release of "Amazon Linux" is available.  
Version 2023.7.20250428:  
Run "/usr/bin/dnf check-release-update" for full release and version update info  
  
#  
~\_##### Amazon Linux 2023  
~~\_#####\  
~~\_####|  
~~\_#/ https://aws.amazon.com/linux/amazon-linux-2023  
~~V~'-'->  
~~~  
~~.-.  
~~/_/_/  
~/m/'
```

```
Last login: Fri May 2 14:00:40 2025 from 18.206.107.28  
[ec2-user@ip-172-31-94-112 ~]$ sudo kubectl get deployment  
NAME          READY    UP-TO-DATE    AVAILABLE    AGE  
crawler       2/2      2             2            48m  
[ec2-user@ip-172-31-94-112 ~]$ sudo kubectl get pods  
NAME                                READY   STATUS    RESTARTS   AGE  
crawler-8668d6dc76-qs96f           0/1     Unknown   1          48m  
crawler-8668d6dc76-z6dms           0/1     Unknown   1 (11m ago) 48m  
[ec2-user@ip-172-31-94-112 ~]$ cd crawler  
[ec2-user@ip-172-31-94-112 crawler]$ vim Dockerfile  
[ec2-user@ip-172-31-94-112 crawler]$
```

3.2 Crawler Implementation Changes

Our crawler implementation went through many changes to get better performance and introduce fault tolerance.

The crawler uses threads that work concurrently to provide faster crawling and communication with indexer and master node.

```
191
192     def _handle_message(self, ch, method, properties, body):
193         reporter = Reporter(self.config)
194         try:
195             subtask_id, url, crawler_ip, depth = body.decode().split("|")
196             logging.info(f"Processing {subtask_id} for {url}")
197
198             reporter.report_status(subtask_id, crawler_ip, "PROCESSING")
199
200             html = self.fetcher.fetch_page(url)
201             time.sleep(self.config.crawl_delay)
202
203             if not html:
204                 reporter.report_status(
205                     subtask_id, crawler_ip, "ERROR", {"reason": "fetch failed"}
206                 )
207                 return
208
209             content = self.fetcher.extract_content(html, url)
210             self.save_to_mongo(url, html)
211             logging.info("sent to mongo")
212             self.indexer.send_to_indexer(url, content["text"])
213
214             reporter.report_status(
215                 subtask_id, crawler_ip, "DONE", {"url": url, **content}
216             )
217             reporter.report_urls(
218                 content.get("extracted_urls", [])[:3],
219                 self.queue_name,
220                 subtask_id,
221                 int(depth) + 1,
222             )
223
224
225     except Exception as e:
226         logging.error(f"Message processing failed: {e}")
```

```

def run(self):
    Thread(target=self._heartbeat_loop, daemon=True).start()

    while not self._shutdown_flag:
        try:
            connection, channel = self._setup_rabbitmq()
            channel.basic_consume(
                queue=f"queue_{self.queue_name}",
                on_message_callback=self._handle_message,
                auto_ack=True,
            )
            logging.info(f"Crawler {self.queue_name} started consuming")
            channel.start_consuming()
        except pika.exceptions.AMQPConnectionError:
            logging.error("RabbitMQ connection failed, retrying in 5s...")
            time.sleep(5)
        except Exception as e:
            logging.error(f"Unexpected error: {e}")
            time.sleep(1)
        finally:
            if "connection" in locals() and connection.is_open:
                connection.close()

def stop(self):
    self._shutdown_flag = True
    self.mongo_client.close()

```

The main function of the crawler node is receiving messages from RabbitMQ and extract the needed content. It is handled in the main loop when it consumed a message from the intended queue, it then sends it to the call back function `_handle_message`.

`_handle_message` extracts the content and the URLs with frequent status reports to the master node like processing, done, error if the server refuses to send html content. It also sends the extracted content to the MongoDB for persistence. Status updates and new URLs are sent in different queues so no tight coupling or issues happen in the master node.

The Reporter class is responsible for all reporting related tasks. It handles sending status updates and new URLs to the master node.

```
class Reporter:
    def __init__(self, config: CrawlerConfig):
        self.config = config

    def _get_connection(self):
        return pika.BlockingConnection(
            pika.ConnectionParameters(
                host=self.config.rabbitmq_host,
                port=self.config.rabbitmq_port,
                credentials=pika.PlainCredentials(
                    self.config.rabbitmq_username, self.config.rabbitmq_password
                ),
                heartbeat=300,
                blocked_connection_timeout=300,
            )
        )
```

```
94     def report(self, queue: str, body: Dict[str, Any]):
95         try:
96             connection = self._get_connection()
97             channel = connection.channel()
98             channel.queue_declare(queue=queue, durable=False)
99             channel.basic_publish(exchange="", routing_key=queue, body=json.dumps(body))
100             connection.close()
101         except Exception as e:
102             logging.error(f"Failed to report to {queue}: {e}")
103
104     def report_status(
105         self, subtask_id: str, crawler_ip: str, status: str, data: Optional[Dict] = None
106     ):
107         self.report(
108             "crawler_updates",
109             {
110                 "subtask_id": subtask_id,
111                 "crawler_ip": crawler_ip,
112                 "status": status,
113                 "data": data or {},
114             },
115         )
116
117     def report_urls(
118         self, new_urls: list[str], crawler_ip: str, parent_subtask_id: str, depth: int
119     ):
120         self.report(
121             "crawler_urls",
122             {
123                 "urls": new_urls,
124                 "crawler_ip": crawler_ip,
125                 "parent_subtask_id": parent_subtask_id,
126                 "depth": depth,
127             },
128         )
129         logging.info(
```

Every operation that sends to a queue opens a separate connection and closes it after execution because pika has strict rules with threading. Only one thread can access a connection at a time.

The crawler node also has a thread for heartbeat messages, periodically sending pings to the master so it can be monitored for failures.

```
def _heartbeat_loop(self):
    logging.info("Starting heartbeat loop")
    reporter = Reporter(self.config)
    while not self._shutdown_flag:
        try:
            reporter.report(
                "crawler_registry",
                {"queue_name": self.queue_name, "timestamp": time.time()},
            )
            logging.info(f"Heartbeat sent from {self.queue_name}")
            time.sleep(15)
        except Exception as e:
            logging.error(f"Heartbeat failed: {e}")
            time.sleep(5)
```

4 Master Node

The master node also has a lot of changes to handle failing crawler nodes and redistribution of tasks.

```
49 class RabbitMQManager:
50     _connection_lock = Lock()
51
52     @classmethod
53     def get_channel(cls):
54         """Get a new channel with thread-safe connection"""
55         params = pika.ConnectionParameters(
56             host=RABBITMQ_HOST,
57             port=RABBITMQ_PORT,
58             credentials=RABBITMQ_CREDS,
59             heartbeat=HEARTBEAT_INTERVAL,
60             blocked_connection_timeout=CONNECTION_TIMEOUT,
61             socket_timeout=10,
62         )
63
64         with cls._connection_lock:
65             connection = pika.BlockingConnection(params)
66             channel = connection.channel()
67             channel.exchange_declare(
68                 exchange="crawler_exchange", exchange_type="direct", durable=False
69             )
70             channel.queue_declare(queue="crawler_registry", durable=False)
71             channel.queue_declare(queue="crawler_updates", durable=False)
72             return channel, connection
```

```

73
74 @classmethod
75 def publish_message(cls, exchange, routing_key, body, max_retries=3):
76     """Thread-safe message publishing with retries"""
77     for attempt in range(max_retries):
78         try:
79             channel, connection = cls.get_channel()
80             channel.basic_publish(
81                 exchange=exchange,
82                 routing_key=routing_key,
83                 body=body,
84                 properties=pika.BasicProperties(
85                     delivery_mode=2, # make message persistent
86                 ),
87             )
88             connection.close()
89             return True
90         except Exception as e:
91             logging.error(f"Publish failed (attempt {attempt+1}): {str(e)}")
92             time.sleep(2**attempt)
93     return False
94
95

```

The RabbitMQ class is used to enforce thread safety since the master has a lot of threads running in parallel, so we need to be very careful with open connections.

To allow for scaling, available crawler nodes aren't hardcoded anymore. The master node receives heartbeats on the registration queue. If a heartbeat comes from an unknown crawler the master node registers it as part of its crawler map.

```

master > new-master-node.py > RabbitMQManager > publish_message
162
163
164 def handle_registration(ch, method, properties, body):
165     try:
166         msg = json.loads(body.decode())
167         qn = msg["queue_name"]
168         now = msg.get("timestamp", time.time())
169
170         with threading.Lock():
171             if qn not in crawler_map:
172                 crawler_map[qn] = Crawler(qn)
173                 logging.info(f"Registered crawler {qn}")
174             crawler = crawler_map[qn]
175             crawler.last_ping = now
176             if crawler.status == CrawlerStatus.FAILED:
177                 crawler.status = CrawlerStatus.IDLE
178                 logging.info(f"Crawler {qn} recovered from failed state")
179     except Exception as e:
180         logging.error(f"Error in registration handler: {str(e)}")
181
182
183 def start_registration_listener():
184     while True:
185         try:
186             channel, connection = RabbitMQManager.get_channel()
187             channel.basic_consume(
188                 queue="crawler_registry",
189                 on_message_callback=handle_registration,
190                 auto_ack=True,
191             )
192             logging.info("Registration listener started")
193             channel.start_consuming()
194         except Exception as e:
195             logging.error(f"Registration listener crashed: {str(e)}, restarting in 5s")
196             time.sleep(5)
197             continue

```

```

master > new-master-node.py > RabbitMQManager > publish_message
210 def health_monitor(timeout=60):
211     logging.info("Health monitor started")
212     while True:
213         now = time.time()
214         with threading.Lock():
215
216             for qn, crawler in list(crawler_map.items()):
217                 if (
218                     now - crawler.last_ping > timeout
219                     and crawler.status != CrawlerStatus.FAILED
220                 ):
221                     logging.info(
222                         f"{qn} timed out, reassigning {len(crawler.assigned_subtasks)} subtasks"
223                     )
224                     crawler.status = CrawlerStatus.FAILED
225
226                     if crawler.assigned_subtasks:
227
228                         for st_id in crawler.assigned_subtasks[:]: # Iterate over copy
229                             sub = subtask_map.get(st_id)
230                             if not sub:
231                                 continue
232
233                             sub.status = SubtaskStatus.QUEUED
234                             new_crawler = get_least_loaded_crawler()
235                             if new_crawler:
236                                 sub.assign_to_crawler(new_crawler.ip)
237                                 new_crawler.assign_task(sub.id)
238                                 success = RabbitMQManager.publish_message(
239                                     exchange="crawler_exchange",
240                                     routing_key=new_crawler.ip,
241                                     body=f"{sub.id}|{sub.url}|{new_crawler.ip}|{sub.depth}",
242                                 )
243                                 if success:
244                                     logging.info(
245                                         f"Reassigned {sub.id} to {new_crawler.ip}"

```

There also is health monitor for the crawlers. It checks whether a crawler has been sending heartbeats or not. If a heartbeat has not been sent for 60 seconds, the crawler is considered FAILED and all its tasks are reassigned.

The new master node also handles getting the new URLs sent by the crawlers, checks for depth and assigns them or not based on the level.

```
master > new-master-node.py > RabbitMQManager > publish_message
288
289
290 def handle_crawler_new_urls(ch, method, properties, body):
291     try:
292         message = json.loads(body.decode())
293         parent_subtask_id = message.get("parent_subtask_id")
294         if parent_subtask_id is None:
295             logging.error(f"Malformed message: {message}")
296             return
297         crawler_ip = message["crawler_ip"]
298         urls = message["urls"]
299         depth = int(message["depth"])
300
301         if depth >= 3:
302             logging.info(f"Depth limit reached for {parent_subtask_id}")
303             return
304
305         logging.info(f"New URLs from {crawler_ip}: {len(urls)} URLs")
306
307         with threading.Lock():
308             parent_subtask = subtask_map.get(parent_subtask_id)
309             if not parent_subtask:
310                 logging.warning(f"Unknown parent subtask {parent_subtask_id}")
311                 return
312
313             for url in urls:
314                 new_subtask = Subtask(
315                     parent_task_id=parent_subtask.parent_task_id,
316                     url=url,
317                     depth=depth + 1,
318                 )
319                 subtask_map[new_subtask.id] = new_subtask
320
321                 assigned_crawler = get_least_loaded_crawler()
322                 if assigned_crawler:
323                     new_subtask.assign_to_crawler(assigned_crawler.ip)
324                     assigned_crawler.assign_task(new_subtask.id)
```

```

1         assigned_crawler = get_least_loaded_crawler()
2         if assigned_crawler:
3             new_subtask.assign_to_crawler(assigned_crawler.ip)
4             assigned_crawler.assign_task(new_subtask.id)
5
6         success = RabbitMQManager.publish_message(
7             exchange="crawler_exchange",
8             routing_key=assigned_crawler.ip,
9             body=f"{new_subtask.id}|{new_subtask.url}|{crawler_ip}|{new_subtask.depth}",
10        )
11        if success:
12            logging.info(
13                f"Dispatched {new_subtask.id} to {assigned_crawler.ip}"
14            )
15        else:
16            logging.error(f"Failed to dispatch {new_subtask.id}")
17    except Exception as e:
18        logging.error(f"Error in new URLs handler: {str(e)}")
19
```

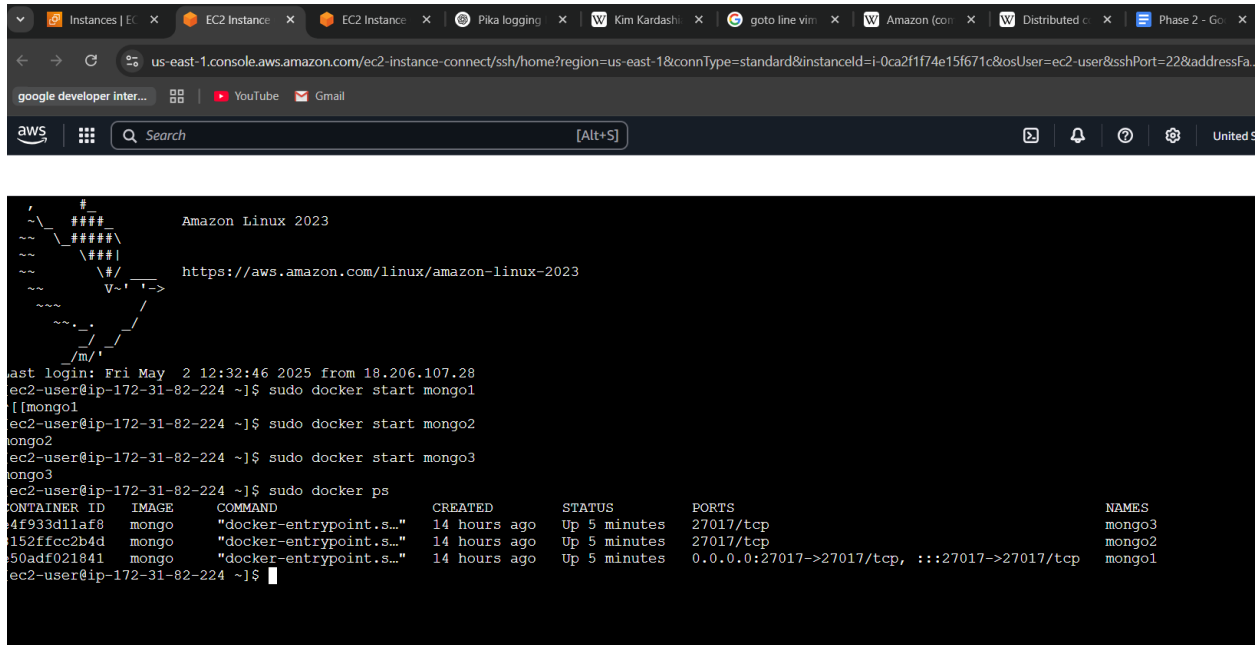
```
6
7 if __name__ == "__main__":
8     Thread(target=start_registration_listener, daemon=True).start()
9     Thread(target=health_monitor, daemon=True).start()
10    Thread(target=start_crawler_update_listener, daemon=True).start()
11    Thread(target=start_crawler_urls_listener, daemon=True).start()
12    master_process()
13
```

All these services are launched as separate threads.

5 MongoDB Persistence

To be able to have a persistent data storage. The number of requests for S3 buckets is low. So we decided to deploy a MongoDB replication set in an ec2 instance.

Using docker containers we ran 3 mongo containers that have a volume attached to them to not lose the data every time the instance doesn't lose data

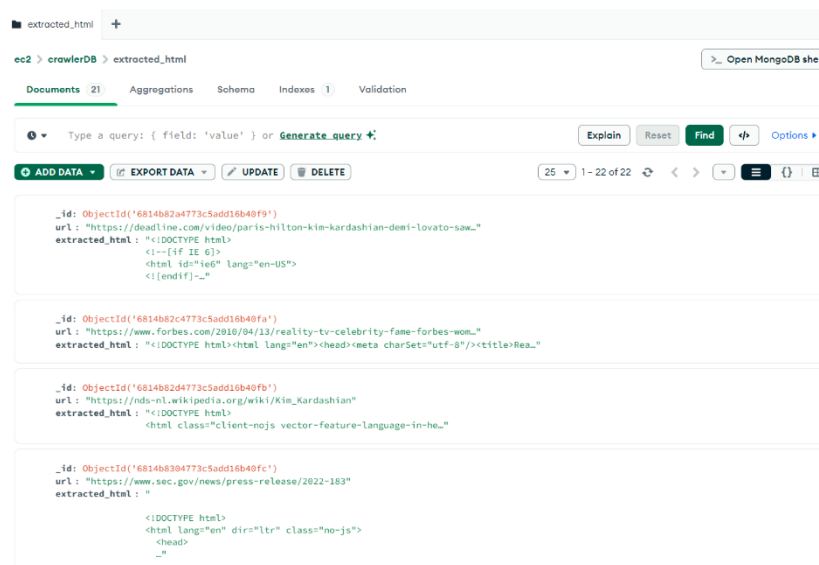


```
us-east-1.console.aws.amazon.com/ec2-instance-connect/ssh/home?region=us-east-1&connType=standard&instanceId=i-0ca2f1f74e15f671c&osUser=ec2-user&sshPort=22&addressFa...
google developer inter...
aws Search [Alt+S]

Amazon Linux 2023
https://aws.amazon.com/linux/amazon-linux-2023

Last login: Fri May 2 12:32:46 2025 from 18.206.107.28
ec2-user@ip-172-31-82-224 ~]$ sudo docker start mongol
[[mongo1
ec2-user@ip-172-31-82-224 ~]$ sudo docker start mongo2
mongo2
ec2-user@ip-172-31-82-224 ~]$ sudo docker start mongo3
mongo3
ec2-user@ip-172-31-82-224 ~]$ sudo docker ps
CONTAINER ID   IMAGE     COMMAND                  CREATED    STATUS    PORTS                               NAMES
4f933d1laf8    mongo    "docker-entrypoint.s..." 14 hours ago Up 5 minutes 27017/tcp    mongo3
152ffcc2b4d    mongo    "docker-entrypoint.s..." 14 hours ago Up 5 minutes 27017/tcp    mongo2
50adf021841    mongo    "docker-entrypoint.s..." 14 hours ago Up 5 minutes 0.0.0.0:27017->27017/tcp, :::27017->27017/tcp    mongo1
ec2-user@ip-172-31-82-224 ~]$
```

MongoDB can be accessed through a GUI MongoDB Compass, official gui for MongoDb.



```
class CrawlerNode:
    def __init__(self, config: CrawlerConfig):
        self.config = config
        self.queue_name = socket.gethostname(socket.gethostname())
        self.fetcher = PageFetcher(config)
        self.indexer = IndexerClient(config.sqs_queue_url)
        self.mongo_client = MongoClient("mongodb://172.31.82.224:27017")
        self.db = self.mongo_client["crawlerDB"]
        self._shutdown_flag = False
```

The crawlers connect to the main mongo container through the pymongo library.

The other containers of the replica set work as backups of the main container, replicating the data and replacing the main container if it falls.

6 Output Logs

6.1 Master Node starting

The master node starts and registers the crawlers from the Kubernetes deployment and a crawler from another ec2

```
2025-05-02 13:22:08,801 - INFO - Found credentials from IAM Role: ezz-kiwi-sqs
2025-05-02 13:22:08,845 - INFO - Health monitor started
2025-05-02 13:22:08,845 - INFO - Master process started
2025-05-02 13:22:08,856 - INFO - Registration listener started
2025-05-02 13:22:08,866 - INFO - Crawler update urls started
2025-05-02 13:22:08,876 - INFO - Crawler update listener started
2025-05-02 13:22:08,885 - INFO - Crawler update urls started
2025-05-02 13:22:20,663 - INFO - Registered crawler 10.42.0.68
2025-05-02 13:22:20,663 - INFO - Registered crawler 10.42.0.67
2025-05-02 13:22:22,251 - INFO - Registered crawler 172.31.84.194
ubuntu@ip-172-31-88-139:~$
```

6.2 Crawler with no tasks fails & recovery

We turned off one of the crawlers and turned it on again.

```
2025-05-02 13:18:35,325 - INFO - Dispatched c16c7863-4bfe-45a1-9f52-e20336fa68dc to 172.31.84.194
2025-05-02 13:18:35,332 - INFO - Update from 172.31.84.194: Subtask c16c7863-4bfe-45a1-9f52-e20336f
2025-05-02 13:18:35,336 - INFO - Dispatched eb758974-f27b-49e7-aacd-4f982e596d60 to 10.42.0.68
2025-05-02 13:18:35,345 - INFO - Dispatched 81625774-acf5-4125-806e-203ace240e10 to 10.42.0.67
2025-05-02 13:18:35,551 - INFO - Update from 172.31.84.194: Subtask eb758974-f27b-49e7-aacd-4f982e5
2025-05-02 13:18:35,552 - INFO - Update from 172.31.84.194: Subtask 81625774-acf5-4125-806e-203ace2
2025-05-02 13:18:36,609 - INFO - Update from 172.31.84.194: Subtask c16c7863-4bfe-45a1-9f52-e20336f
2025-05-02 13:18:36,615 - INFO - Depth limit reached for c16c7863-4bfe-45a1-9f52-e20336fa68dc
2025-05-02 13:18:46,891 - INFO - Update from 172.31.84.194: Subtask eb758974-f27b-49e7-aacd-4f982e5
2025-05-02 13:18:46,892 - INFO - Update from 172.31.84.194: Subtask 81625774-acf5-4125-806e-203ace2
2025-05-02 13:22:08,801 - INFO - Found credentials from IAM Role: ezz-kiwi-sqs
2025-05-02 13:22:08,845 - INFO - Health monitor started
2025-05-02 13:22:08,845 - INFO - Master process started
2025-05-02 13:22:08,856 - INFO - Registration listener started
2025-05-02 13:22:08,866 - INFO - Crawler update urls started
2025-05-02 13:22:08,876 - INFO - Crawler update listener started
2025-05-02 13:22:08,885 - INFO - Crawler update urls started
2025-05-02 13:22:20,663 - INFO - Registered crawler 10.42.0.68
2025-05-02 13:22:20,663 - INFO - Registered crawler 10.42.0.67
2025-05-02 13:22:22,251 - INFO - Registered crawler 172.31.84.194
2025-05-02 13:33:08,850 - INFO - 172.31.84.194 timed out, reassigning 0 subtasks
ubuntu@ip-172-31-88-139:~$
```

```
2025-05-02 13:18:36,615 - INFO - Depth limit reached for c16c7863-4bfe-45a1-9f52-e20336fa68dc
2025-05-02 13:18:46,891 - INFO - Update from 172.31.84.194: Subtask eb758974-f27b-49e7-aacd-4f9
2025-05-02 13:18:46,892 - INFO - Update from 172.31.84.194: Subtask 81625774-acf5-4125-806e-203
2025-05-02 13:22:08,801 - INFO - Found credentials from IAM Role: ezz-kiwi-sqs
2025-05-02 13:22:08,845 - INFO - Health monitor started
2025-05-02 13:22:08,845 - INFO - Master process started
2025-05-02 13:22:08,856 - INFO - Registration listener started
2025-05-02 13:22:08,866 - INFO - Crawler update urls started
2025-05-02 13:22:08,876 - INFO - Crawler update listener started
2025-05-02 13:22:08,885 - INFO - Crawler update urls started
2025-05-02 13:22:20,663 - INFO - Registered crawler 10.42.0.68
2025-05-02 13:22:20,663 - INFO - Registered crawler 10.42.0.67
2025-05-02 13:22:22,251 - INFO - Registered crawler 172.31.84.194
2025-05-02 13:33:08,850 - INFO - 172.31.84.194 timed out, reassigning 0 subtasks
2025-05-02 13:33:46,689 - INFO - Crawler 172.31.84.194 recovered from failed state
```


6.3 Dispatching tasks

```
2025-05-02 13:22:08,801 - INFO - Found credentials from IAM Role: ezz-kiwi-sqs
2025-05-02 13:22:08,845 - INFO - Health monitor started
2025-05-02 13:22:08,845 - INFO - Master process started
2025-05-02 13:22:08,856 - INFO - Registration listener started
2025-05-02 13:22:08,866 - INFO - Crawler update urls started
2025-05-02 13:22:08,876 - INFO - Crawler update listener started
2025-05-02 13:22:08,885 - INFO - Crawler update urls started
2025-05-02 13:22:20,663 - INFO - Registered crawler 10.42.0.68
2025-05-02 13:22:20,663 - INFO - Registered crawler 10.42.0.67
2025-05-02 13:22:22,251 - INFO - Registered crawler 172.31.84.194
2025-05-02 13:33:08,850 - INFO - 172.31.84.194 timed out, reassigning 0 subtasks
2025-05-02 13:33:46,689 - INFO - Crawler 172.31.84.194 recovered from failed state
2025-05-02 13:34:50,581 - INFO - New task: 1 seeds, depth=2
2025-05-02 13:34:50,581 - INFO - Created 1 subtasks for Task dba981cd-6547-48af-9c11-12dd522e470b
2025-05-02 13:34:50,589 - INFO - Dispatched 7cb6f520-202a-48e5-b466-3a659f119bb6 to 10.42.0.68
2025-05-02 13:34:50,981 - INFO - Update from 10.42.0.68: Subtask 7cb6f520-202a-48e5-b466-3a659f119bb6 -> PROCESSING
2025-05-02 13:35:12,698 - INFO - Update from 10.42.0.68: Subtask 7cb6f520-202a-48e5-b466-3a659f119bb6 -> ERROR
2025-05-02 13:37:52,208 - INFO - New task: 1 seeds, depth=2
2025-05-02 13:37:52,208 - INFO - Created 1 subtasks for Task b246077b-6088-40d4-9a97-f241a23101b2
2025-05-02 13:37:52,218 - INFO - Dispatched c97785a9-75c5-4e8b-9b22-c8f9e1432bfe to 10.42.0.68
2025-05-02 13:37:52,363 - INFO - Update from 10.42.0.68: Subtask c97785a9-75c5-4e8b-9b22-c8f9e1432bfe -> PROCESSING
2025-05-02 13:40:19,181 - INFO - Update from 10.42.0.68: Subtask c97785a9-75c5-4e8b-9b22-c8f9e1432bfe -> DONE
2025-05-02 13:40:19,431 - INFO - New URLs from 10.42.0.68: 3 URLs
2025-05-02 13:40:19,442 - INFO - Dispatched dcd448e6-439a-485a-a3a9-31fd1ac6b81b to 10.42.0.68
2025-05-02 13:40:19,449 - INFO - Dispatched 0dcf9c49-23ca-4278-bd52-5c4610e5eb42 to 10.42.0.67
2025-05-02 13:40:19,456 - INFO - Dispatched 40f7e3c1-4407-40b8-940e-a48a990bbf18 to 172.31.84.194
2025-05-02 13:40:19,463 - INFO - Update from 10.42.0.68: Subtask 40f7e3c1-4407-40b8-940e-a48a990bbf18 -> PROCESSING
2025-05-02 13:40:19,848 - INFO - Update from 10.42.0.68: Subtask 0dcf9c49-23ca-4278-bd52-5c4610e5eb42 -> PROCESSING
2025-05-02 13:40:19,849 - INFO - Update from 10.42.0.68: Subtask dcd448e6-439a-485a-a3a9-31fd1ac6b81b -> PROCESSING
2025-05-02 13:40:21,174 - INFO - Update from 10.42.0.68: Subtask 40f7e3c1-4407-40b8-940e-a48a990bbf18 -> DONE
2025-05-02 13:40:21,182 - INFO - Depth limit reached for 40f7e3c1-4407-40b8-940e-a48a990bbf18
ubuntu@ip-172-31-88-139:~$
```

6.4 Reassigning tasks

```
2025-05-02 13:42:08,854 - INFO - 172.31.84.194 timed out, reassigning 1 subtasks
2025-05-02 13:42:08,863 - INFO - Reassigned 40f7e3c1-4407-40b8-940e-a48a990bbf18 to 10.42.0.68
2025-05-02 13:42:09,356 - INFO - Update from 10.42.0.68: Subtask 40f7e3c1-4407-40b8-940e-a48a990bbf18 -> PROCESSING
```

6.5 Crawler logs

```
2025-05-02 13:36:46,809 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:37:01,818 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:37:16,828 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:37:31,837 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:37:46,847 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:38:01,856 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:38:16,867 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:38:31,878 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:38:46,890 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:39:01,898 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:39:16,907 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:39:31,918 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:39:46,929 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:40:01,940 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:40:16,950 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:40:19,456 - INFO - Processing 40f7e3c1-4407-40b8-940e-a48a990bbf18 for https://uk.wikipedia.org/wiki/%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%B5_%D0%BD%D0%B0%D0%B2%D1%8D%D0%BD%D1%8F
2025-05-02 13:40:19,704 - INFO - Fetched page: https://uk.wikipedia.org/wiki/%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%B5_%D0%BD%D0%B0%D0%B2%D1%87%D0%B0%D0%BD%D0%B0%D1%8F
2025-05-02 13:40:21,015 - INFO - Extracted 689 links from https://uk.wikipedia.org/wiki/%D0%9C%D0%B0%D1%88%D0%B8%D0%BD%D0%BD%D0%B5_%D0%BD%D0%B0%D0%B2%D1%87%D0%B0%D0%BD%D0%B0%D1%8F
2025-05-02 13:40:21,097 - INFO - sent to mongo
2025-05-02 13:40:21,183 - INFO - Reported 3 new URLs from 40f7e3c1-4407-40b8-940e-a48a990bbf18 at depth 3
2025-05-02 13:40:31,962 - INFO - Heartbeat sent from 172.31.84.194
2025-05-02 13:40:46,971 - INFO - Heartbeat sent from 172.31.84.194
```

6.6 Indexer Start

```
ec2-user@ip-172-31-87-169:~$ python3 indexer.py
/usr/lib/python3.9/site-packages/requests/_init_.py:87: RequestsDependencyWarning: urllib3 (2.4.0) or chardet (4.0.0) doesn't match a supported version!
  warnings.warn("urllib3 ({}), or chardet ({}), doesn't match a supported version".format(urllib3.__version__, chardet.__version__), RequestsDependencyWarning)
INFO:botocore.credentials:Found credentials from IAM Role: ec2-kivl-sqs
/home/ec2-user/.local/lib/python3.9/site-packages/elasticsearch/_sync/client/_init_.py:399: SecurityWarning: Connecting to 'https://172.31.92.167:9200' using TLS with verify_certs=False is insecure
  transport = transport_class(
/home/ec2-user/.local/lib/python3.9/site-packages/urllib3/connectionpool.py:1097: InsecureRequestWarning: Unverified HTTPS request is being made to host '172.31.92.167'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#tls-warnings
  warnings.warn(
INFO:elastic_transport.transport:GET https://172.31.92.167:9200/ [status:200 duration:0.050s]
('name': 'ip-172-31-92-167-ec2.internal', 'cluster name': 'elasticsearch', 'cluster uuid': 'vxeV3sr7T8u9XuBVCil A', 'version': {'number': '8.12.2', 'build flavor': 'default', 'build type': 'rpm', 'build hash': '48a287ab9497e952de30327444b0809e55d46466', 'build date': '2024-02-19T10:04:32.774273190Z', 'build_snapshot': False, 'lucene_version': '9.9.2', 'minimum_wire_compatibility version': '7.17.0', 'minimum_index_compatibility version': '7.0.0'}, 'tagline': 'You Know, for Search')
/home/ec2-user/.local/lib/python3.9/site-packages/urllib3/connectionpool.py:1097: InsecureRequestWarning: Unverified HTTPS request is being made to host '172.31.92.167'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#tls-warnings
  warnings.warn(
INFO:elastic_transport.transport:HEAD https://172.31.92.167:9200/ [status:200 duration:0.005s]
INFO:root: Elasticsearch cluster is up!
/home/ec2-user/.local/lib/python3.9/site-packages/elasticsearch/_sync/client/_init_.py:54: DeprecationWarning: Passing transport options in the API method is deprecated. Use 'Elasticsearch.options()' instead.
  resp = self.es.indices.create(
/home/ec2-user/.local/lib/python3.9/site-packages/urllib3/connectionpool.py:1097: InsecureRequestWarning: Unverified HTTPS request is being made to host '172.31.92.167'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#tls-warnings
  warnings.warn(
INFO:elastic_transport.transport:PUT https://172.31.92.167:9200/webpages [status:400 duration:0.015s]
INFO:root: Index 'webpages' ready.
INFO:root:IndexingNode started. Polling for messages...
Enter keyword to search (or 'exit'): INFO:root:Received 1 messages from SQS
/home/ec2-user/.local/lib/python3.9/site-packages/urllib3/connectionpool.py:1097: InsecureRequestWarning: Unverified HTTPS request is being made to host '172.31.92.167'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#tls-warnings
  warnings.warn(
INFO:elastic_transport.transport:POST https://172.31.92.167:9200/webpages/_doc [status:201 duration:0.117s]
INFO:root:Document indexed: https://en.wikipedia.org/wiki/Distributed_computing
INFO:root:Indexed document from URL: https://en.wikipedia.org/wiki/Distributed_computing
INFO:root:Deleted message from SQS for URL: https://en.wikipedia.org/wiki/Distributed_computing
INFO:root:Received 1 messages from SQS
/home/ec2-user/.local/lib/python3.9/site-packages/urllib3/connectionpool.py:1097: InsecureRequestWarning: Unverified HTTPS request is being made to host '172.31.92.167'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#tls-warnings
  warnings.warn(
INFO:elastic_transport.transport:POST https://172.31.92.167:9200/webpages/_doc [status:201 duration:0.017s]
INFO:root:Document indexed: https://et.wikipedia.org/wiki/Hajusarvutus
INFO:root:Indexed document from URL: https://et.wikipedia.org/wiki/Hajusarvutus
INFO:root:Deleted message from SQS for URL: https://et.wikipedia.org/wiki/Hajusarvutus
Enter keyword to search (or 'exit'):
```

6.7 Indexer search

```
Enter keyword to search (or 'exit'): distribute
/home/ec2-user/.local/lib/python3.9/site-packages/urllib3/connectionpool.py:1097: InsecureRequestWarning: Unverified HTTPS request is being made to host '172.31.92.167'. Adding certificate verification is strongly advised. See: https://urllib3.readthedocs.io/en/latest/advanced-usage.html#tls-warnings
  warnings.warn(
INFO:elastic_transport.transport:POST https://172.31.92.167:9200/webpages/_search [status:200 duration:0.018s]
INFO:root:Search results for 'distribute': ['https://en.wikipedia.org/wiki/Distributed_computing', 'https://en.wikipedia.org/wiki/Distributed_computing', 'https://ja.wikipedia.org/wiki/AES408066E6495A34E3402B34E34034B34E34034944E34034A54E34034BC4E34034064E3402A34E34034B34E34024B0', 'https://es.wikipedia.org/wiki/Computaci3n distribuida', 'https://en.wikipedia.org/wiki/Amazon_(company)', 'https://web.archive.org/web/20230801073911/https://www.marketscreener.com/quote/stock/AMAZON-COM-1284605/company/', 'https://apps.bostonglobe.com/metro/investigations/spotlight/2019/11/21/seeing-red/convenience-culture-makes-traffic-worse/']
Enter keyword to search (or 'exit'):
```