# Kubernetes Server Client Deployment

This sounds like a comprehensive and well-structured project! To implement this, here's a roadmap you could follow:

## Step 1: Create the Server and Client Scripts

1. **Server.py**:

   - Create two Python scripts (`server1.py` and `server2.py`).
   - Include basic endpoints for handling requests, such as `/status` or `/process`.

2. **Client.py**:

   - Implement a script that sends requests to the servers.
   - Use a round-robin scheduling algorithm to balance the requests.

## Step 2: Dockerize the Applications

1. **Dockerfiles**:

   - Write a `Dockerfile` for each server and the client.
   - Ensure proper installation of dependencies in the Docker images.
   - Example:

   ```dockerfile
   FROM python:3.9
   WORKDIR /app
   COPY . .
   RUN pip install -r requirements.txt
   CMD ["python", "server.py"]
   ```

2. **Build Images**:

   - Build the Docker images using `docker build` and tag them appropriately.

## Step 3: Kubernetes Deployment

1. **Deployments**:

   - Write Kubernetes deployment manifests for the servers and client.
   - Example for a server:

   ```yaml
   apiVersion: apps/v1
   kind: Deployment
   metadata:
     name: server-deployment
   spec:
     replicas: 2
     selector:
       matchLabels:
   ```

```
            app: server
    template:
      metadata:
        labels:
          app: server
      spec:
        containers:
        - name: server
          image: your-docker-image
          ports:
          - containerPort: 5000
```

2. **Services**:

   - Create services for exposing the servers and the client.

3. **Ingress**:

   - Configure an ingress to route requests.

## Step 4: Load Balancing and Scheduling

- Use Kubernetes services with built-in load balancing to distribute requests among the server pods.
- Modify the client to follow a round-robin pattern for direct scheduling (if required).

## Step 5: Monitoring with Prometheus

1. **Prometheus Setup**:

   - Deploy Prometheus on the Kubernetes cluster.
   - Configure it to scrape metrics from the server pods.

2. **Metrics**:

   - Instrument your servers to expose metrics via endpoints like `/metrics` using a library like `prometheus_client`.

3. **Dashboard**:

   - Optionally, integrate Grafana for visualizing metrics.

## Step 6: Test and Refine

- Test the deployment by scaling replicas up and down.
- Monitor the system behavior under different loads using Prometheus.

Would you like help with specific aspects of this, such as Kubernetes manifests, Dockerfiles, or Prometheus configuration?

To implement a **round-robin scheduling** algorithm in your Kubernetes project, you can manage request distribution using:

1. **Client-Side Round Robin (in `client.py`)**:

   - The client script decides which server to send the next request to in a round-robin fashion.

2. **Kubernetes Service Load Balancing**:

   - Kubernetes Services already implement round-robin for load balancing requests across pods. This is achieved by default when using a `ClusterIP` or `LoadBalancer` service type.

Here's how you can implement **round robin** in both approaches:

## 1. Client-Side Round Robin

Modify the `client.py` to distribute requests evenly across the servers.

**Example Implementation:**

```python
import requests
import itertools

# List of server URLs
servers = ["http://server1:5000", "http://server2:5000"]

# Create a round-robin iterator
server_iterator = itertools.cycle(servers)

def send_request(data):
    # Get the next server in the round-robin order
    server = next(server_iterator)
    try:
        # Send a request to the selected server
        response = requests.post(f"{server}/process", json=data)
        print(f"Response from {server}: {response.text}")
    except requests.exceptions.RequestException as e:
        print(f"Error communicating with {server}: {e}")

# Example usage
if __name__ == "__main__":
    for i in range(10):
        send_request({"message": f"Request {i}"})
```

- **Explanation**:
  - The `itertools.cycle()` function cycles through the list of servers endlessly.
  - Each time `next(server_iterator)` is called, it returns the next server in the round-robin order.

## 2. Kubernetes Service Load Balancing

Kubernetes services handle round-robin automatically when you expose multiple pods under the same service.

**Steps:**

1. **Deployment**:
   - Deploy multiple replicas of your server (`server.py`).
   - Kubernetes ensures all replicas are part of a single service.
2. **Service Definition**:
   - Create a Kubernetes service to expose your server pods.
   - Use the `ClusterIP` service type (default) for internal traffic or `LoadBalancer` for external access.
   - Example YAML:

```yaml
apiVersion: v1
kind: Service
metadata:
  name: server-service
spec:
  selector:
    app: server
  ports:
    - protocol: TCP
      port: 5000
      targetPort: 5000
  type: ClusterIP
```

3. **How It Works**:

   - When the client sends requests to the service endpoint (`http://server-service:5000`), Kubernetes routes the requests to the pods using an internal round-robin mechanism.

---

## Combining Both Approaches

You can combine Kubernetes service load balancing with a client-side fallback round-robin mechanism. The client sends requests to multiple service endpoints, which then balance the load across pods.

Would you like a more detailed example or help with setting up these configurations?

Here's a complete setup for your **Kubernetes project with two servers and a client**. This includes the code for the client, servers, Dockerfiles, and Kubernetes manifests.

---

## 1. Server Code (`server.py`)

Each server will serve requests and return its unique identifier.

```python
from flask import Flask, request, jsonify
import socket

app = Flask(__name__)

@app.route('/process', methods=['POST'])
def process_request():
    data = request.get_json()
    hostname = socket.gethostname()  # Get the server's hostname
    response = {
        "server": hostname,
        "message": f"Processed: {data['message']}"
    }
    return jsonify(response)

if __name__ == "__main__":
    app.run(host="0.0.0.0", port=5000)
```

---

## 2. Client Code (`client.py`)

Implements round-robin scheduling for distributing requests to the servers.

```python
import requests
import itertools

# List of server service endpoints
servers = [
    "http://server-service:5000",
    "http://server-service:5000"
]

# Create a round-robin iterator
server_iterator = itertools.cycle(servers)

def send_request(data):
    server = next(server_iterator)  # Get the next server
    try:
        response = requests.post(f"{server}/process", json=data)
        print(f"Response from {server}: {response.json()}")
    except requests.exceptions.RequestException as e:
        print(f"Error communicating with {server}: {e}")

# Example usage
if __name__ == "__main__":
    for i in range(10):
        send_request({"message": f"Request {i}"})
```

## 3. Dockerfile

Use the same `Dockerfile` for both the server and client.

**Dockerfile for Server (`Dockerfile`)**

```dockerfile
# Use Python base image
FROM python:3.9-slim

# Set the working directory
WORKDIR /app

# Copy the current directory into the container
COPY . .

# Install dependencies
RUN pip install flask

# Expose the port for Flask
EXPOSE 5000

# Command to run the server
CMD ["python", "server.py"]
```

**Dockerfile for Client (`Dockerfile`)**

Printed using [ChatGPT to PDF](ChatGPT to PDF), powered by PDFCrowd [HTML to PDF API](HTML to PDF API).

5/13

```dockerfile
# Use Python base image
FROM python:3.9-slim

# Set the working directory
WORKDIR /app

# Copy the current directory into the container
COPY . .

# Install dependencies
RUN pip install requests

# Command to run the client
CMD ["python", "client.py"]
```

## 4. Kubernetes Deployment and Service Manifests

**Server Deployment and Service**

```yaml
# server-deployment.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: server-deployment
spec:
  replicas: 2
  selector:
    matchLabels:
      app: server
  template:
    metadata:
      labels:
        app: server
    spec:
      containers:
      - name: server
        image: <your-dockerhub-username>/server:latest
        ports:
        - containerPort: 5000
---
# server-service.yaml
apiVersion: v1
kind: Service
metadata:
  name: server-service
spec:
  selector:
    app: server
  ports:
    - protocol: TCP
      port: 5000
      targetPort: 5000
  type: ClusterIP
```

**Client Deployment**

```yaml
# client-deployment.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: client-deployment
spec:
  replicas: 1
  selector:
    matchLabels:
      app: client
  template:
    metadata:
      labels:
        app: client
    spec:
      containers:
      - name: client
        image: <your-dockerhub-username>/client:latest
```

## 5. Prometheus Setup

You can use a standard Prometheus Helm chart or create a custom `prometheus.yaml` configuration to monitor the servers.

## Steps to Deploy

1. **Build Docker Images**:

```bash
docker build -t <your-dockerhub-username>/server:latest -f Dockerfile .
docker build -t <your-dockerhub-username>/client:latest -f Dockerfile .
docker push <your-dockerhub-username>/server:latest
docker push <your-dockerhub-username>/client:latest
```

2. **Apply Kubernetes Manifests**:

```bash
kubectl apply -f server-deployment.yaml
kubectl apply -f server-service.yaml
kubectl apply -f client-deployment.yaml
```

3. **Monitor Logs**: Check logs for servers and client to verify the round-robin mechanism:

```bash
kubectl logs deployment/server-deployment
kubectl logs deployment/client-deployment
```

4. **Prometheus Integration**: Expose a `/metrics` endpoint in `server.py` and configure Prometheus to scrape the metrics.

Let me know if you'd like Prometheus configuration details or help with anything else!

To implement **round-robin scheduling** within the **servers**, you can modify the `server.py` code to distribute requests to specific sub-services or processes running within the server.

Here's the **updated full implementation** with round-robin logic inside the server:

## 1. Updated Server Code (`server.py`)

This implementation uses a round-robin approach to distribute tasks to internal workers or processes.

```python
from flask import Flask, request, jsonify
import socket
from itertools import cycle

app = Flask(__name__)

# Simulated workers (could represent sub-services, threads, etc.)
workers = ["Worker1", "Worker2", "Worker3"]

# Round-robin iterator
worker_iterator = cycle(workers)

@app.route('/process', methods=['POST'])
def process_request():
    data = request.get_json()
    hostname = socket.gethostname()  # Get the server's hostname

    # Select the next worker in round-robin order
    worker = next(worker_iterator)

    response = {
        "server": hostname,
        "worker": worker,
        "message": f"Processed by {worker}: {data['message']}"
    }
    return jsonify(response)

if __name__ == "__main__":
    app.run(host="0.0.0.0", port=5000)
```

## 2. Client Code (`client.py`)

The client already implements round-robin scheduling across multiple servers.

```python
import requests
import itertools

# List of server service endpoints
servers = [
    "http://server-service:5000",
```

```python
    "http://server-service:5000"
]

# Create a round-robin iterator
server_iterator = itertools.cycle(servers)

def send_request(data):
    server = next(server_iterator)  # Get the next server
    try:
        response = requests.post(f"{server}/process", json=data)
        print(f"Response from {server}: {response.json()}")
    except requests.exceptions.RequestException as e:
        print(f"Error communicating with {server}: {e}")

# Example usage
if __name__ == "__main__":
    for i in range(10):
        send_request({"message": f"Request {i}"})
```

## 3. Dockerfiles

### Server Dockerfile

```dockerfile
# Use Python base image
FROM python:3.9-slim

# Set the working directory
WORKDIR /app

# Copy the current directory into the container
COPY . .

# Install dependencies
RUN pip install flask

# Expose the port for Flask
EXPOSE 5000

# Command to run the server
CMD ["python", "server.py"]
```

### Client Dockerfile

```dockerfile
# Use Python base image
FROM python:3.9-slim

# Set the working directory
WORKDIR /app

# Copy the current directory into the container
COPY . .

# Install dependencies
RUN pip install requests
```

```
# Command to run the client
CMD ["python", "client.py"]
```

## 4. Kubernetes Manifests

### Server Deployment and Service

```yaml
# server-deployment.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: server-deployment
spec:
  replicas: 2
  selector:
    matchLabels:
      app: server
  template:
    metadata:
      labels:
        app: server
    spec:
      containers:
      - name: server
        image: <your-dockerhub-username>/server:latest
        ports:
        - containerPort: 5000
---
# server-service.yaml
apiVersion: v1
kind: Service
metadata:
  name: server-service
spec:
  selector:
    app: server
  ports:
    - protocol: TCP
      port: 5000
      targetPort: 5000
  type: ClusterIP
```

### Client Deployment

```yaml
# client-deployment.yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: client-deployment
spec:
  replicas: 1
  selector:
    matchLabels:
      app: client
  template:
    metadata:
      labels:
        app: client
```

```yaml
    spec:
      containers:
      - name: client
        image: <your-dockerhub-username>/client:latest
```

## Testing Round-Robin Behavior

1. **Deploy the System**:

```bash
   kubectl apply -f server-deployment.yaml
   kubectl apply -f server-service.yaml
   kubectl apply -f client-deployment.yaml
```

2. **Observe Server Logs**: Check the logs of the server pods to see the round-robin assignment of workers for each request.

```bash
   kubectl logs deployment/server-deployment
```

3. **Check Client Output**: Run the client to confirm that requests are distributed across servers in a round-robin manner, and within each server, requests are further distributed among workers.

## Example Output

### Client Console:

```plaintext
  Response from http://server-service:5000: {'server': 'server-1', 'worker': 'Worker1', 'messa
  ge': 'Processed by Worker1: Request 1'}
  Response from http://server-service:5000: {'server': 'server-2', 'worker': 'Worker2', 'messa
  ge': 'Processed by Worker2: Request 2'}
  Response from http://server-service:5000: {'server': 'server-1', 'worker': 'Worker3', 'messa
  ge': 'Processed by Worker3: Request 3'}
```
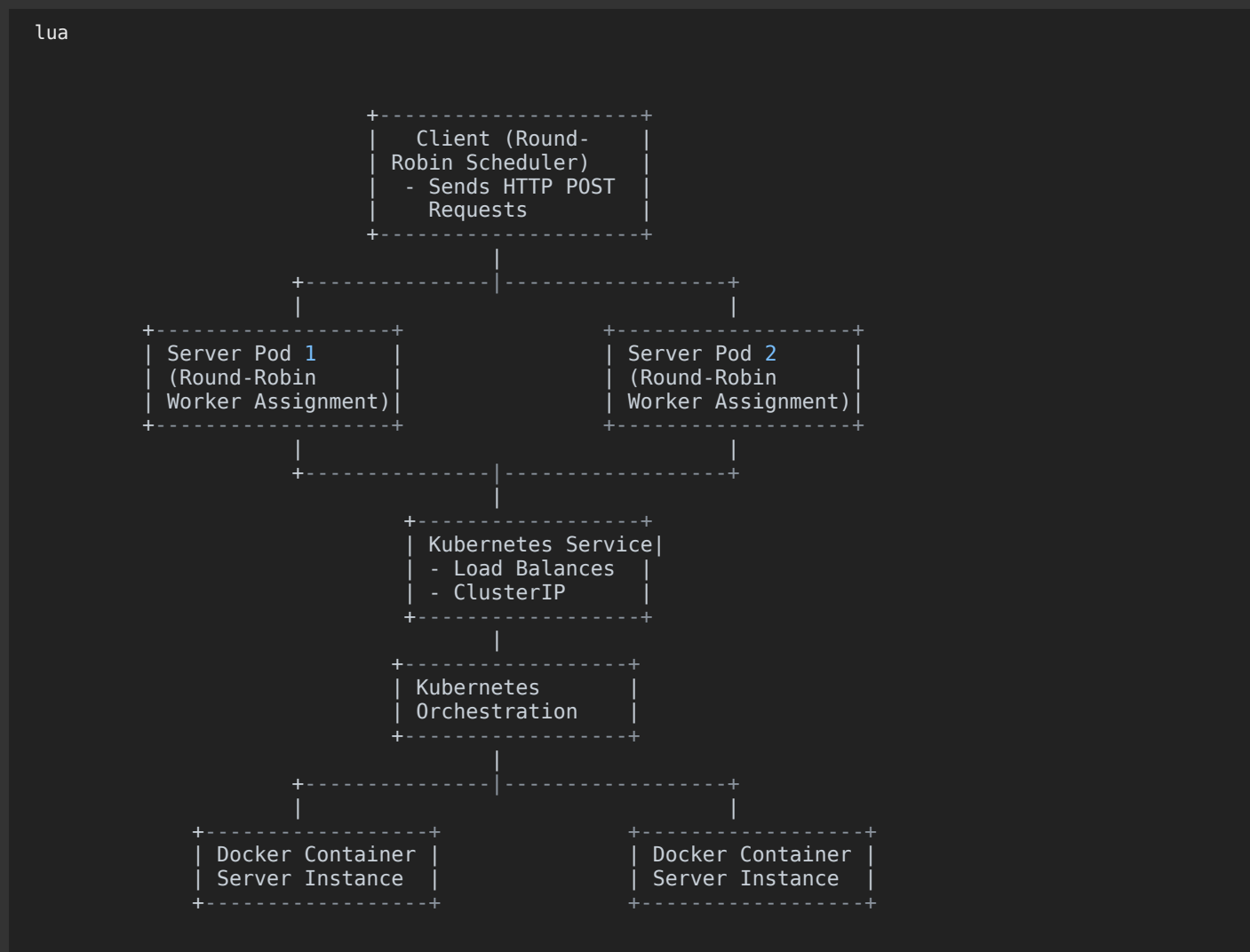
### Server Logs:

```plaintext
  Processed request on server-1 using Worker1
  Processed request on server-2 using Worker2
  Processed request on server-1 using Worker3
```

This implementation handles **round-robin scheduling** at both the client and server levels, ensuring load balancing across servers and internal workers. Let me know if you need additional refinements!

Here's a textual representation of the **data flow diagram flowchart** for the project:

```lua
                    +--------------------+
                    |   Client (Round-   |
                    | Robin Scheduler)   |
                    |  - Sends HTTP POST |
                    |    Requests        |
                    +--------------------+
                              |
                 +------------|------------+
                 |                         |
        +------------------+      +------------------+
        | Server Pod 1     |      | Server Pod 2     |
        | (Round-Robin     |      | (Round-Robin     |
        | Worker Assignment)|      | Worker Assignment)|
        +------------------+      +------------------+
                 |                         |
                 +------------|------------+
                              |
                    +------------------+
                    | Kubernetes Service|
                    | - Load Balances  |
                    | - ClusterIP      |
                    +------------------+
                              |
                    +------------------+
                    | Kubernetes       |
                    | Orchestration    |
                    +------------------+
                              |
                 +------------|------------+
                 |                         |
        +------------------+      +------------------+
        | Docker Container |      | Docker Container |
        | Server Instance  |      | Server Instance  |
        +------------------+      +------------------+
```

## Explanation of Components:

1. **Client (Round-Robin Scheduler)**:
    - Implements round-robin logic to distribute requests across the two servers.
2. **Kubernetes Service**:
    - Acts as a load balancer for requests sent to the server pods.
    - Uses a `ClusterIP` to expose the server pods internally.
3. **Server Pods**:
    - Each pod contains a `server.py` instance that further assigns tasks to workers using round-robin scheduling.
4. **Kubernetes Orchestration**:
    - Manages server pod replication, scaling, and health monitoring.

Printed using [ChatGPT to PDF](#), powered by PDFCrowd [HTML to PDF API](#).

12/13

5. **Docker Containers**:

- Hosts the server instances, which are encapsulated in Docker images.

The **data flow**:

- The client sends requests in a round-robin manner to the Kubernetes service.
- The Kubernetes service distributes requests to server pods.
- Each server pod further processes the requests using its own internal round-robin mechanism among workers.