

# Distributed Image Processing System using Cloud Computing

## Contents

Project Scope: .....	2
Objectives: .....	2
Requirements: .....	2
User Stories: .....	3
System Architecture: .....	3
Selected Technologies: .....	4
Why we used these technologies: .....	4
Communication: .....	6
Cost analysis: .....	7
Project plan: .....	10
Sequence diagram: .....	12
Components diagram: .....	13
Infrastructure diagram: .....	13
Network diagram: .....	14

## Project Scope:

The project aims to develop a distributed image processing system utilizing cloud computing technologies. It involves designing and implementing a system capable of distributing image processing tasks across multiple virtual machines in the cloud. The system will support various image processing algorithms such as filtering, edge detection, and color manipulation. It should be scalable to accommodate an increasing workload by adding more virtual machines and should maintain fault tolerance to handle node failures gracefully.

## Objectives:

- Design and implement a distributed image processing system using Python.
- Utilize cloud-based virtual machines for distributed computing.
- Use image processing algorithms including filtering, edge detection, and color manipulation.
- Ensure scalability to accommodate increased workload by adding virtual machines dynamically.
- Ensure fault tolerance by reassigning tasks from failed nodes to operational ones.

## Requirements:

- **Distributed Processing:** The system should distribute image processing tasks across multiple virtual machines in the cloud.
- **Image Processing Algorithms:** Use filtering, edge detection, and colour manipulation algorithms.
- **Scalability:** The system should dynamically scale by adding virtual machines as the workload increases.
- **Fault Tolerance:** The system should handle node failures gracefully by reassigning tasks from failed nodes to operational ones.
- **User Interface:** Develop a user-friendly interface for users to upload images, select processing operations, monitor task progress, and download processed images.
- **Cloud Computing Platform:** Select a suitable cloud computing platform (e.g., AWS, Azure, Google Cloud) for hosting virtual machines.
- **Parallel Processing Framework:** Choose either OpenCL or MPI for parallel processing of image data.

- **Monitoring System:** Implement a monitoring system to track the progress of image processing tasks.
- **Documentation:** Provide comprehensive documentation including system architecture, setup instructions, and user guide.

## User Stories:

- As a user, I want to upload an image to the system for processing.
- As a user, I want to select the type of image processing operation to be performed.
- As a user, I want to download the processed image once the operation is complete.
- As a user, I want to monitor the progress of the image processing task.

## System Architecture:

### 1. **User Interface (UI):**

- Provides an interface for users to interact with the system.
- Allows users to upload images, select processing operations, monitor task progress and view the processed images.

### 2. **Master Node:**

- Virtual machine in the cloud responsible for the application backend before the image processing.
- Handles user requests and generates messages to the worker nodes.
- Divide the image into many segments using image segmentation methods.
- Distributes image processing tasks to worker nodes.
- Manages scalability and fault tolerance.
- Sends back the processed image to the client UI.

### 3. **Worker Nodes:**

- Virtual machines in the cloud responsible for actual image processing using rpc architecture.
- Receive tasks from the backend (master node), perform processing using parallel computing, and return results.

### 4. **Communication Layer:**

- Facilitates communication between different components of the system using TCP and web sockets to RPC communication and we will define it below.
  - Utilizes messaging protocols or frameworks for task distribution and result retrieval.
5. **Monitoring:**
- Monitors system performance, resource utilization, and task progress.

## Selected Technologies:

1. **Cloud Platform:**
  - **Microsoft Azure:** Cloud provider with virtual machines (Azure VMs), storage (Azure Blob Storage), and messaging services (Azure Service Bus).
2. **Programming Language:**
  - **Python:** Selected for its ease of development, rich ecosystem of libraries (e.g., CV for image processing), and suitability for parallel computing.
3. **Parallel Computing Framework:**
  - **MPI (Message Passing Interface):** Enables distributed computing and communication between worker nodes.
4. **Communication:**
  - **RPC with TCP and WebSockets**

## Considerations:

- **Scalability:** Ensure the architecture can scale horizontally by adding more worker nodes dynamically.
- **Fault Tolerance:** Implement mechanisms to handle node failures, such as task reassignment and redundancy.
- **Cost Optimization:** Optimize resource usage to minimize operational costs, especially in cloud environments where costs can scale with usage.

## Why we used these technologies:

1. **Cloud Platform - Microsoft Azure:**

- **Azure VMs:** These provide scalable and customizable virtual machines, allowing the deployment of various applications without worrying about hardware infrastructure.
- **Azure Blob Storage:** Offers scalable object storage for documents, images, videos, and other unstructured data, enabling efficient data management and access.

**Reason for Selection:** Microsoft Azure was chosen for its robust infrastructure, extensive services, and reliable performance. It offers a wide range of scalable solutions that fit the project's requirements, ensuring flexibility and efficiency in deployment and management.

## 2. Programming Language - Python:

- **Ease of Development:** Python's simple and readable syntax makes it easy to write and maintain code, accelerating development cycles.
- **Rich Ecosystem of Libraries:** Python boasts a vast collection of libraries for various purposes, such as computer vision (CV) for image processing, machine learning, data analysis, and more. This wealth of resources enhances productivity and facilitates the implementation of complex functionalities.
- **Suitability for Parallel Computing:** Python supports parallel computing through frameworks like MPI, enabling efficient utilization of resources and faster processing of tasks.

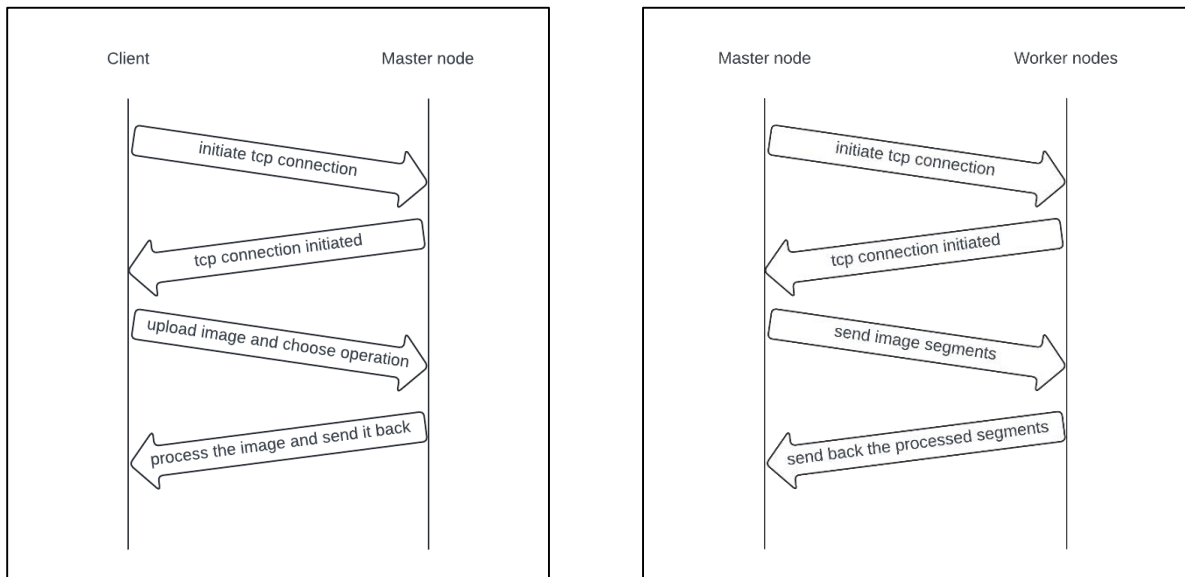
**Reason for Selection:** Python was chosen for its combination of simplicity, versatility, and powerful libraries. Its suitability for parallel computing aligns well with the project's requirements for efficient data processing and analysis.

## 3. Parallel Computing Framework - MPI (Message Passing Interface):

- **Distributed Computing:** MPI enables efficient communication and coordination between multiple processes running on distributed computing systems, allowing for parallel execution of tasks.
- **Scalability:** It supports scaling across multiple nodes, enabling the system to handle large datasets and compute-intensive workloads effectively.
- **Performance:** MPI is known for its high performance and low overhead, making it well-suited for demanding parallel computing tasks.

**Reason for Selection:** MPI was chosen for its proven track record in parallel computing, especially in high-performance computing (HPC) environments. It provides the necessary tools and mechanisms for efficient parallelization of algorithms and processing of large datasets, essential for the project's objectives.

## Communication:



### We will use TCP and WebSockets for RPC Communication

#### 1. Define RPC Protocol:

- Define a custom RPC protocol that specifies how messages are structured and exchanged between the nodes.
- Design the protocol to include message types, method identifiers, parameters, and error handling mechanisms.

#### 2. Choose Python WebSocket Library:

- Select a WebSocket library for Python such as socket or websocket-client.
- These libraries provide easy-to-use APIs for creating WebSocket clients and servers in Python.

#### 3. Establish TCP Connections:

- Use Python's built-in socket module or a third-party library to establish TCP connections between clients and servers.
- TCP connections can be used for long-lived RPC connections where reliability and ordered delivery of data are crucial.

#### **4. WebSocket Server Implementation:**

- Implement a WebSocket server using the chosen WebSocket library to handle incoming WebSocket connections from clients.
- Define WebSocket message handling logic to process RPC requests and send responses back to clients.
- We will use the worker nodes as servers that accept the request from the client.

#### **5. WebSocket Client Implementation:**

- Develop WebSocket client code to connect to the WebSocket server and initiate RPC requests.
- Send RPC request messages over the WebSocket connection and await responses from the server.
- We will use the master node as a client that send requests to the servers that is the worker node.

#### **6. Integrate RPC Logic:**

- Define RPC methods and interfaces that can be invoked over both TCP and WebSocket connections.
- Implement RPC handlers on the server side to execute RPC requests received via TCP or WebSocket and send back responses.

#### **7. Performance Optimization:**

- Optimize TCP settings and configurations for performance, such as adjusting buffer sizes.
- Minimize latency by reducing round-trip times and optimizing data transfer efficiency.

Cost analysis:

Developing a distributed image processing system using cloud computing Application involves various costs, including:

#### **Development Costs:**

- **Software Development:**
  - **Resource Time:**
    - Programming (Python) - Analyzing, designing, coding, and testing the application.
    - Network Engineering - Designing and implementing the network architecture.
  - **Tools and Frameworks:**
    - Python development environment (IDE)
    - Specific libraries for networking, parallel computing, and UI/UX
    - Cloud-based development platform
  - **Testing and Quality Assurance:**
    - Unit testing, integration testing, and user acceptance testing
    - Automated testing tools
- **Documentation:**
  - Creating user manuals, API documentation, and internal technical documentation

#### **Infrastructure Costs:**

- **Deployment:**
  - Cloud-based server
  - Domain name and SSL certificate
  - Load balancer
- **Hosting:**
  - Monthly or annual fees for cloud server or other hosting services
  - Bandwidth costs depending on user activity

#### **Maintenance Costs:**

- **Bug Fixes:** Addressing issues reported by users
- **Feature Enhancements:** Implementing new features and functionality
- **Security Updates:** Maintaining security patches and updates for libraries and frameworks



- Version Control: Managing code changes and releases

**Additional Costs:**

- Project Management: Planning, scheduling, and coordinating development activities
- Legal and Regulatory Compliance: Ensuring compliance with data privacy regulations
- Third-Party Services: APIs, libraries, or other paid services
- Marketing and Promotion: Advertising and promoting the application to attract users

**Cost Estimation:**

Due to the project's scope and varying factors, providing a definitive cost estimate is difficult. However, here's a rough breakdown:

- Development: \$5,000 - \$20,000+
- Infrastructure: \$500 - \$2,000+ per month
- Maintenance: \$1,000 - \$5,000+ per month

**Cost Optimization Strategies:**

- **Open-source libraries and frameworks:**
  - Utilize freely available libraries and frameworks for various functionalities, reducing licensing costs.
- **Cloud-based development and hosting:**
  - Leverage cloud platforms for development and deployment to reduce infrastructure costs and maintenance overhead.
- **Agile development methodology:**
  - Focus on rapid prototyping and iterative development to ensure resource efficiency and early feedback.
- **Community-driven development:**
  - Encourage contributions from open-source communities to leverage shared resources and expertise.

Overall, the cost of developing and maintaining a distributed image processing system using cloud computing Application will depend on various factors like project complexity, team size, and chosen technologies. Implementing cost-optimization strategies can significantly reduce expenses and ensure project viability.

## Project plan:

### Phase 1: Project Planning and Design (2-3 weeks)

#### Tasks:

1. Define project scope, objectives, and requirements.
2. Conduct stakeholder meetings to gather input on project goals.
3. Research cloud computing technologies and select the appropriate platform (AWS, Google Cloud, Azure, etc.).
4. Design system architecture, including components, interactions, and data flows.
5. Determine technologies for parallel processing (MPI, OpenCL).
6. Create a detailed project plan with tasks, responsibilities, and timelines.
7. Draft user stories based on gathered requirements.
8. Document project plan and design decisions.

#### Responsibilities:

- Project Manager: Overall project coordination, stakeholder communication.
- System Architect: System design, technology selection.
- Development Team: Input on technical feasibility, user stories.

#### Timelines:

- Weeks 1-2: Define scope, objectives, and requirements; research and select technologies; design system architecture.
- Week 3: Finalize project plan, document design decisions, and user stories.

### Phase 2: Development of Basic Functionality (2-3 weeks)

#### Tasks:

1. Implement basic image processing operations (filtering, edge detection, color manipulation).
2. Set up cloud environment and provision virtual machines.
3. Develop worker threads for processing tasks.
4. Implement image upload functionality.
5. Develop user interface for basic operations.
6. Conduct unit testing of implemented functionality.

**Responsibilities:**

- Development Team: Implementation of basic functionality, unit testing.
- System Architect: Guidance on system integration.
- Project Manager: Progress tracking, issue resolution.

**Timelines:**

- Weeks 4: Implement basic image processing operations and cloud setup.
- Weeks 5-6: Develop worker threads, image upload functionality, and user interface.

**Phase 3: Development of Advanced Functionality (2-3 weeks)****Tasks:**

1. Implement advanced image processing operations (e.g., feature extraction, object recognition).
2. Develop distributed processing functionality using MPI or OpenCL.
3. Implement scalability features to add more virtual machines dynamically.
4. Incorporate fault tolerance mechanisms to handle node failures.
5. Conduct integration testing of advanced functionality.

**Responsibilities:**

- Development Team: Implementation of advanced functionality, testing.
- System Architect: Guidance on distributed processing and fault tolerance.
- Project Manager: Ensuring adherence to project timeline and goals.

**Timelines:**

- Weeks 7-8: Implement advanced image processing operations and distributed processing.
- Weeks 8-9: Incorporate scalability and fault tolerance features, conduct integration testing.

**Phase 4: Testing, Documentation, and Deployment (2-3 weeks)****Tasks:**

1. Conduct thorough testing of the entire system, including unit, integration, and system testing.
2. Document system design, codebase, and user instructions.
3. Prepare deployment scripts and configurations.
4. Deploy the system to the cloud environment.
5. Perform final system testing and validation.

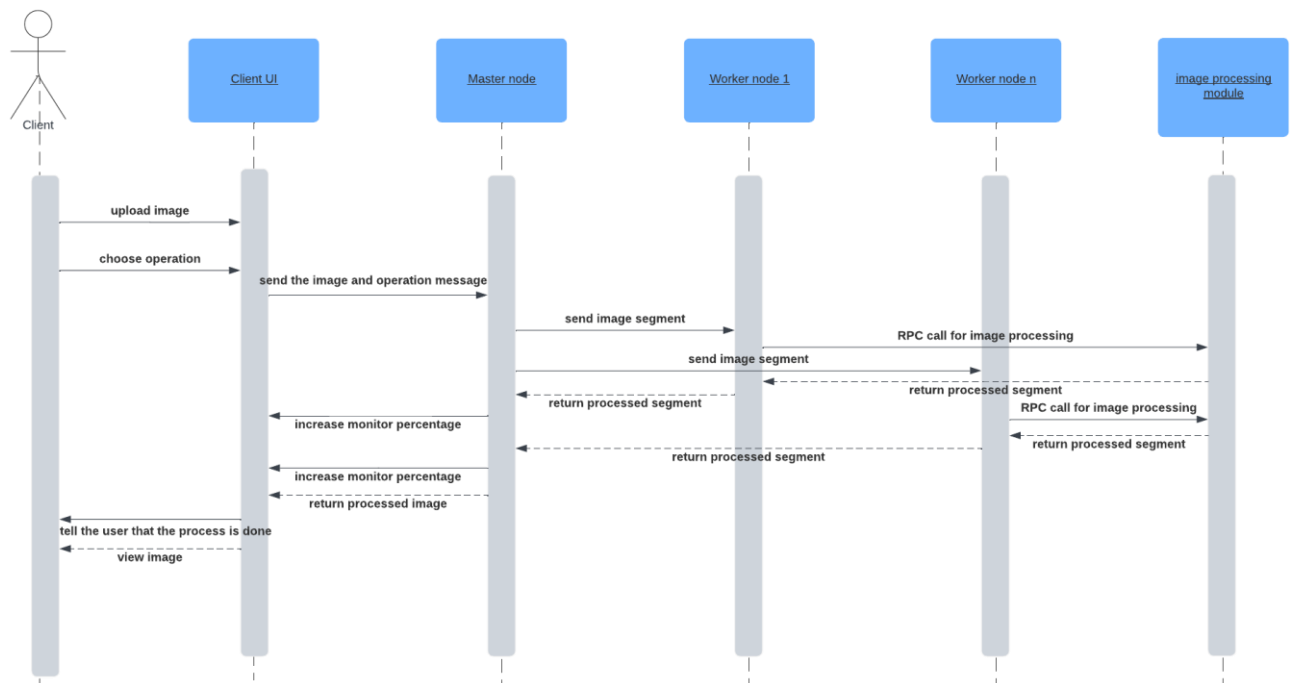
**Responsibilities:**

- Development Team: Testing, documentation, deployment.
- Quality Assurance (QA) Team: Test planning and execution.
- Technical Writer: System documentation.
- Project Manager: Deployment coordination, final validation.

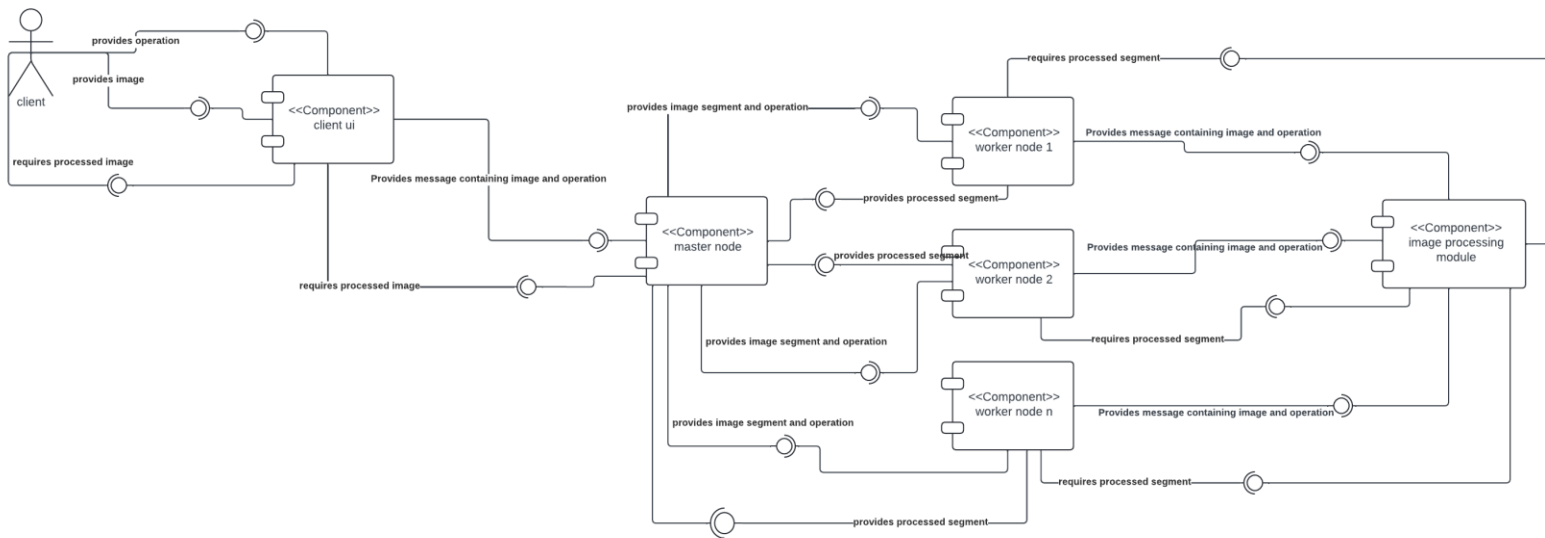
#### Timelines:

- Weeks 10-11: Testing and documentation.
- Weeks 12: Deployment, final testing, and validation.

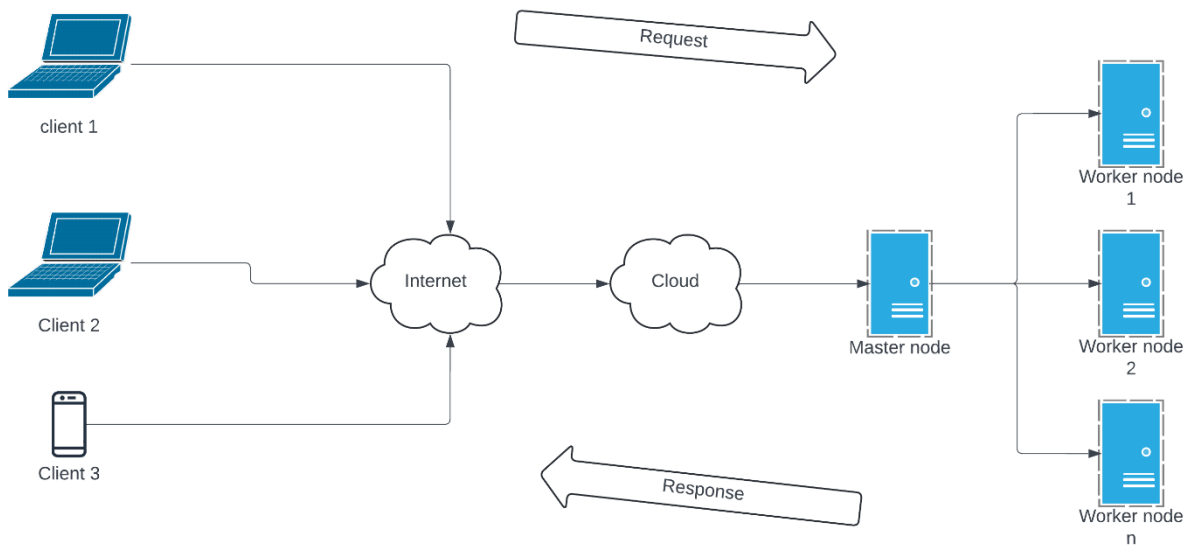
#### Sequence diagram:



## Components diagram:



## Infrastructure diagram:



## Network diagram:

