

Cloud Computing (CC) Presentation

Load Balancing & Its Strategies



Introduction

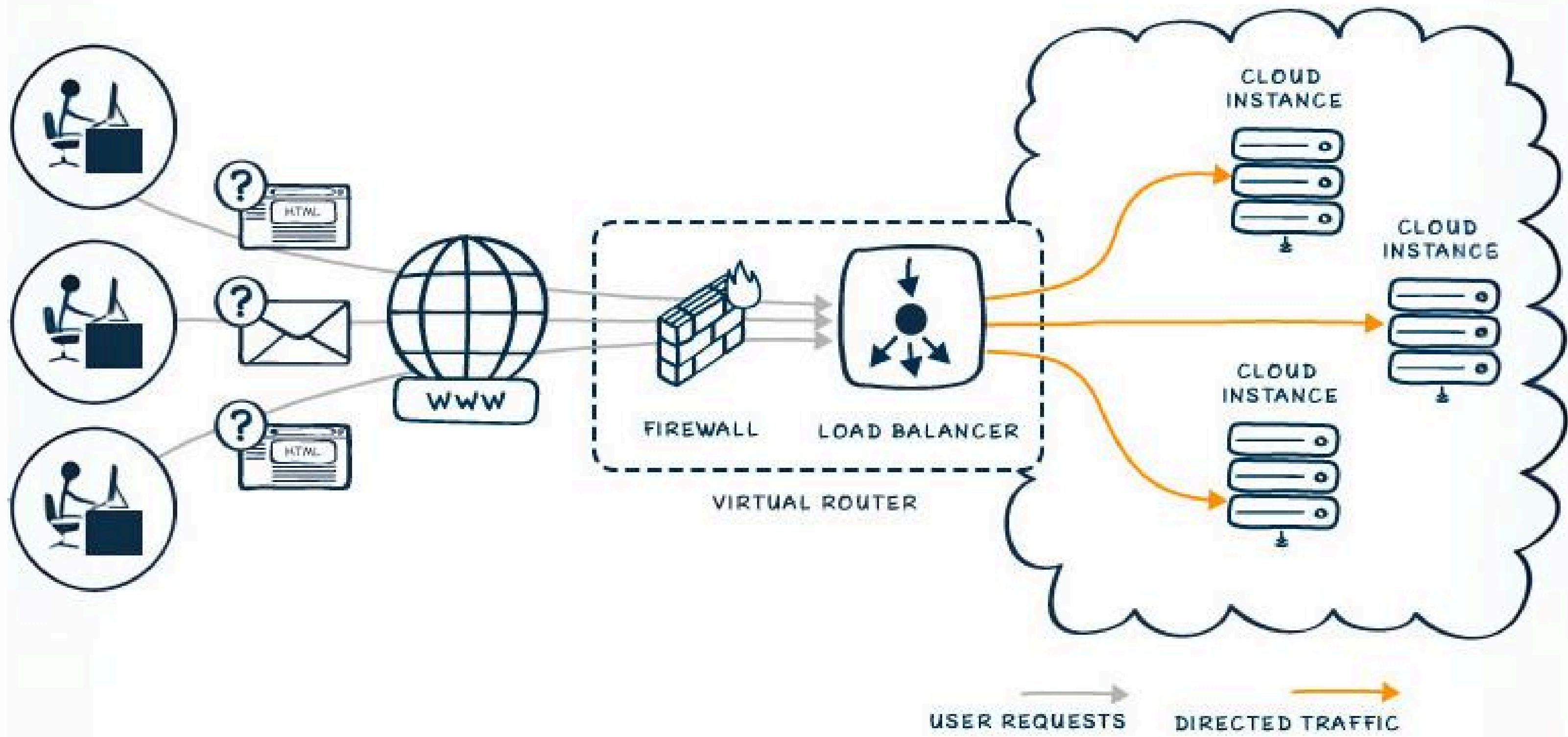
What is a Load Balancer?

- A networking device or software application
- Distributes and balances incoming traffic among the horizontally scaled homogeneous servers

Advantages

- Higher availability
- Higher scalability
- More efficient utilization of servers
- Better application security
- Better application performance





Availability

- Server failure can increase application downtime
- Load balancers increase the **fault tolerance** of systems
- Automatically detects server problems and redirects client traffic to available servers
- Make these tasks easier:
 - Run application server maintenance or upgrades **without application downtime**
 - Provide automatic **disaster recovery** to backup sites
 - Perform **health checks** and prevent issues that can cause downtime



Scalability

- **Direct network traffic** intelligently among multiple servers
- Handle millions of client requests
- Easier horizontal scaling
- Functionalities:
 - Prevents **traffic bottlenecks** at any one server
 - Predicts application traffic so that you can **add or remove** different **servers**, if needed
 - Adds **redundancy** to your system so that you can **scale with confidence**



Security

- **Built-in security features** to add another layer of security
- Useful tool to deal with **DDoS attacks**
- Concurrent requests which cause server failure are redirected.
- Handles the following:
 - **Monitor traffic and block malicious content**
 - Automatically **redirect attack traffic** to multiple backend servers to minimize impact
 - Route traffic through a **group of network firewalls** for additional security



Performance

- Improve application performance by:
 - Increasing response time
 - Reducing network latency.
- Critical tasks such as:
 - **Distribute the load evenly** between servers to improve application performance
 - Redirect client requests to a geographically closer server to **reduce latency**
 - Ensure the **reliability** and **performance** of physical and virtual computing resources



Load Balancing Strategies

What are LB Strategies?

- A load balancing strategy/algorithm is the set of rules that a load balancer follows
- Goal is to determine the best server for each of the different client requests
- For a request “X” which server “Y” is the most optimal to be routed to

Types of LB Strategies

- Load balancing algorithms fall into two main categories:
 - Static LB Strategies - Static load balancing strategies/algorithms follow fixed rules and are independent of the current server state
 - Dynamic LB Strategies - Dynamic load balancing strategies/algorithms examine the current state of the servers before distributing traffic, which includes no. of connections, response time of the server, etc.

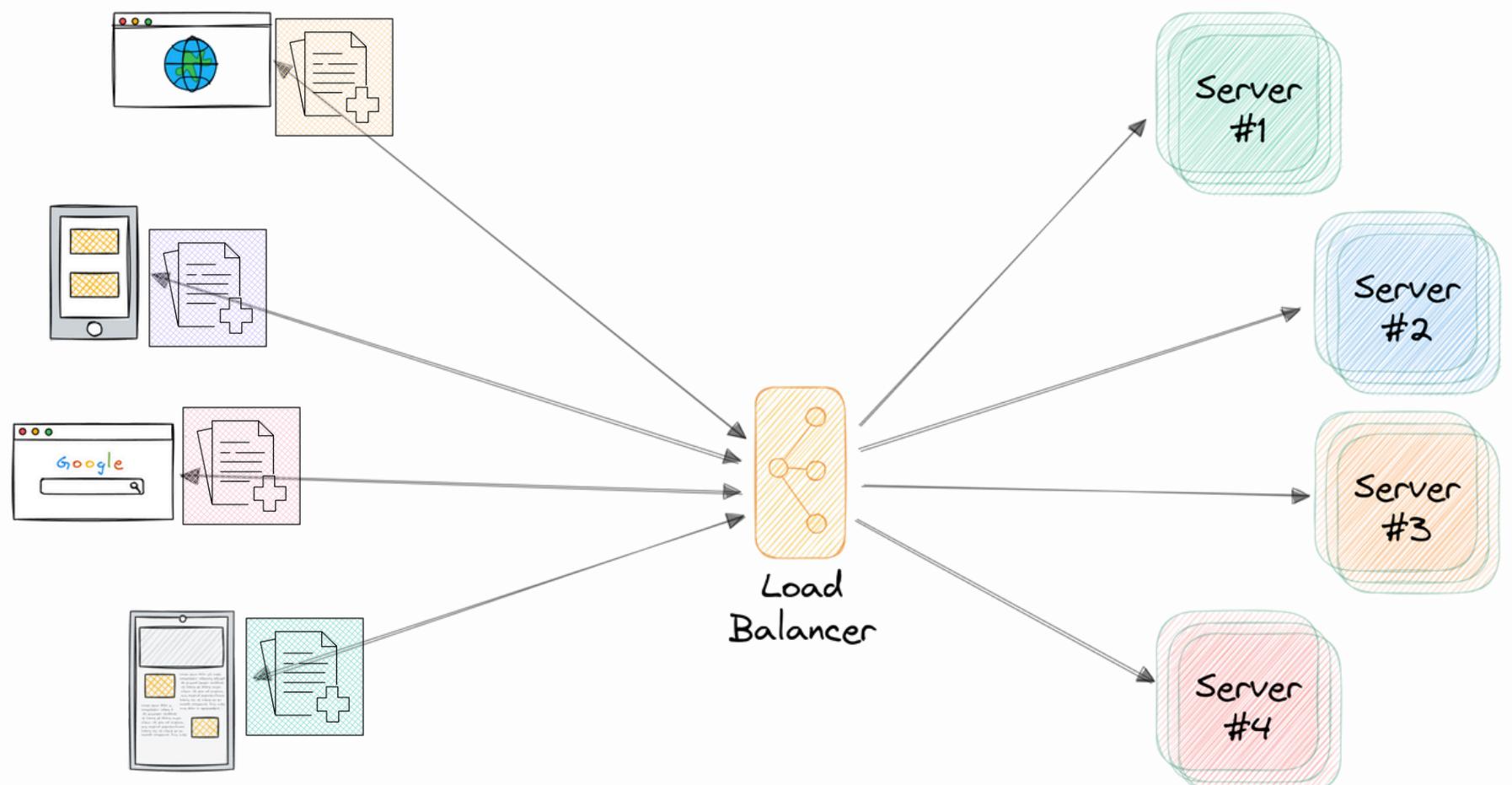
Implementation Details

- Protocol (HTTP, WebSocket, etc) agnostic load balancing
- Forwarding TCP segments instead of HTTP, WS, etc requests and responses
- Client or server is oblivious to whether the requests and responses are going through a load-balancer or to the server or client respectively directly
- Load balancer handles two TCP connections for each client-to-server connection, one between client and LB and the other between LB and one of the servers
- IO between both the connections can be handled in user-space (infinite loop reading and writing) or kernel-space (sendfile(2) and copy_file_range(2) syscalls in Linux)



(S) Random

- Distributes incoming requests randomly across servers
- Simple and easy to implement
- Pros:
 - No need for tracking server load
 - Even distribution over time if request volume is high
- Cons:
 - Can lead to uneven distribution for smaller volumes
 - Doesn't account for server capacity or load



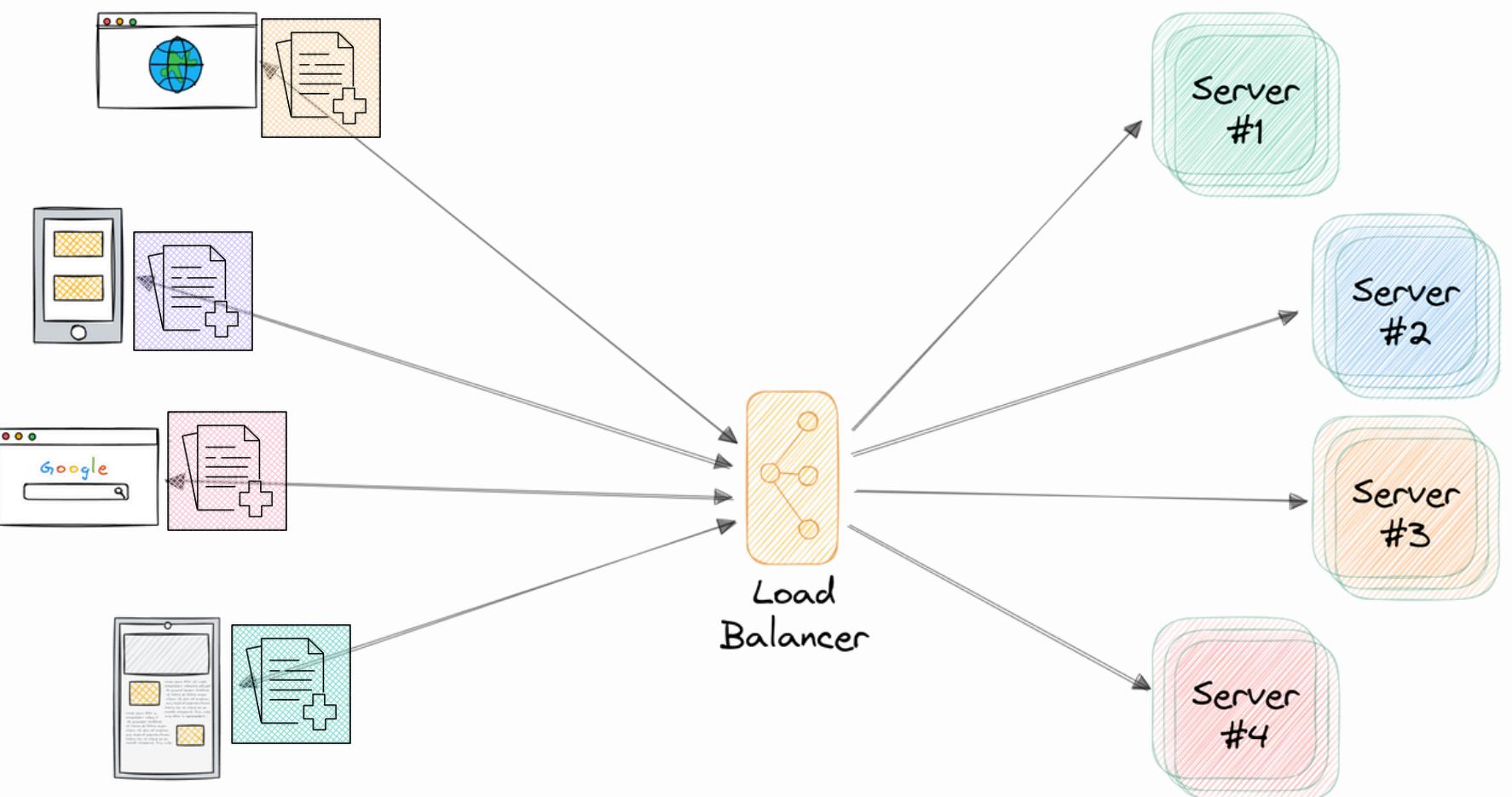


```
● ● ●

1 type RandomStrategy struct {
2     servers []ServerAddr
3 }
4
5 func NewRandomStrategy(servers []ServerAddr) *RandomStrategy {
6     return &RandomStrategy{servers}
7 }
8
9 func (strategy *RandomStrategy) ServerAddr(ClientAddr string) (ServerAddr, error) {
10    if len(strategy.servers) == 0 {
11        return "", errors.New("no servers available")
12    }
13
14    server := strategy.servers[rand.IntN(len(strategy.servers))]
15
16    return server, nil
17 }
18
19 func (strategy *RandomStrategy) Connected(addr ServerAddr) {}
20 func (strategy *RandomStrategy) Disconnected(addr ServerAddr) {}
```

(S) Round Robin

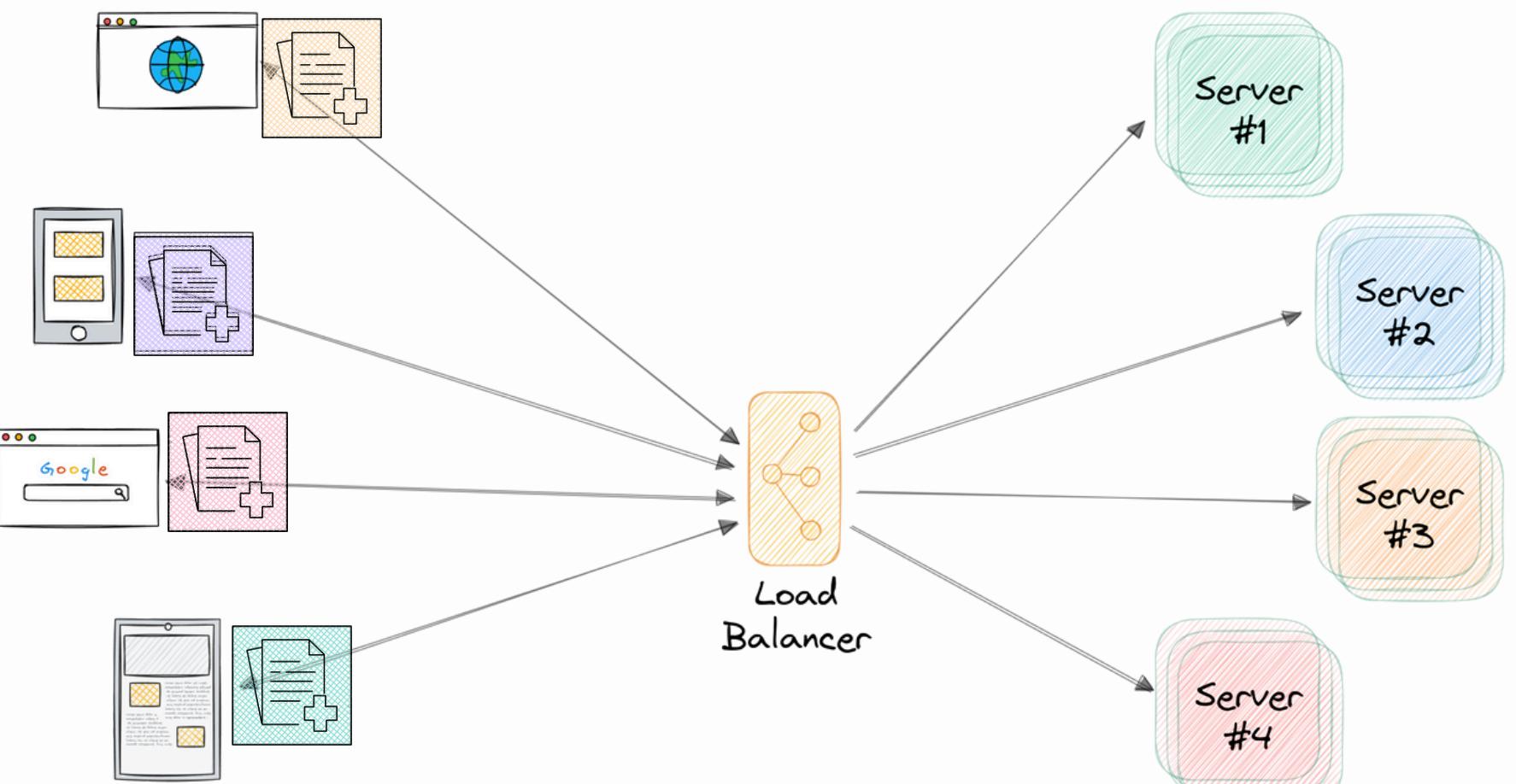
- Sequentially distributes requests across servers in a circular order
- Equal distribution of traffic
- Pros:
 - Simple to implement
 - Ensures that no server is overloaded early on
- Cons:
 - Doesn't account for server load or capacity differences
 - Can cause delays if a server becomes overloaded



```
1 type RoundRobinStrategy struct {
2     servers    []ServerAddr
3     nextIndex int
4 }
5
6 func NewRoundRobinStrategy(servers []ServerAddr) *RoundRobinStrategy {
7     return &RoundRobinStrategy{servers, 0}
8 }
9
10 func (strategy *RoundRobinStrategy) ServerAddr(ClientAddr string) (ServerAddr, error) {
11     if len(strategy.servers) == 0 {
12         return "", errors.New("no servers available")
13     }
14
15     server := strategy.servers[strategy.nextIndex]
16     strategy.nextIndex = (strategy.nextIndex + 1) % len(strategy.servers)
17
18     return server, nil
19 }
20
21 func (strategy *RoundRobinStrategy) Connected(addr ServerAddr) {}
22 func (strategy *RoundRobinStrategy) Disconnected(addr ServerAddr) {}
23
```

(S) Weighted Round Robin

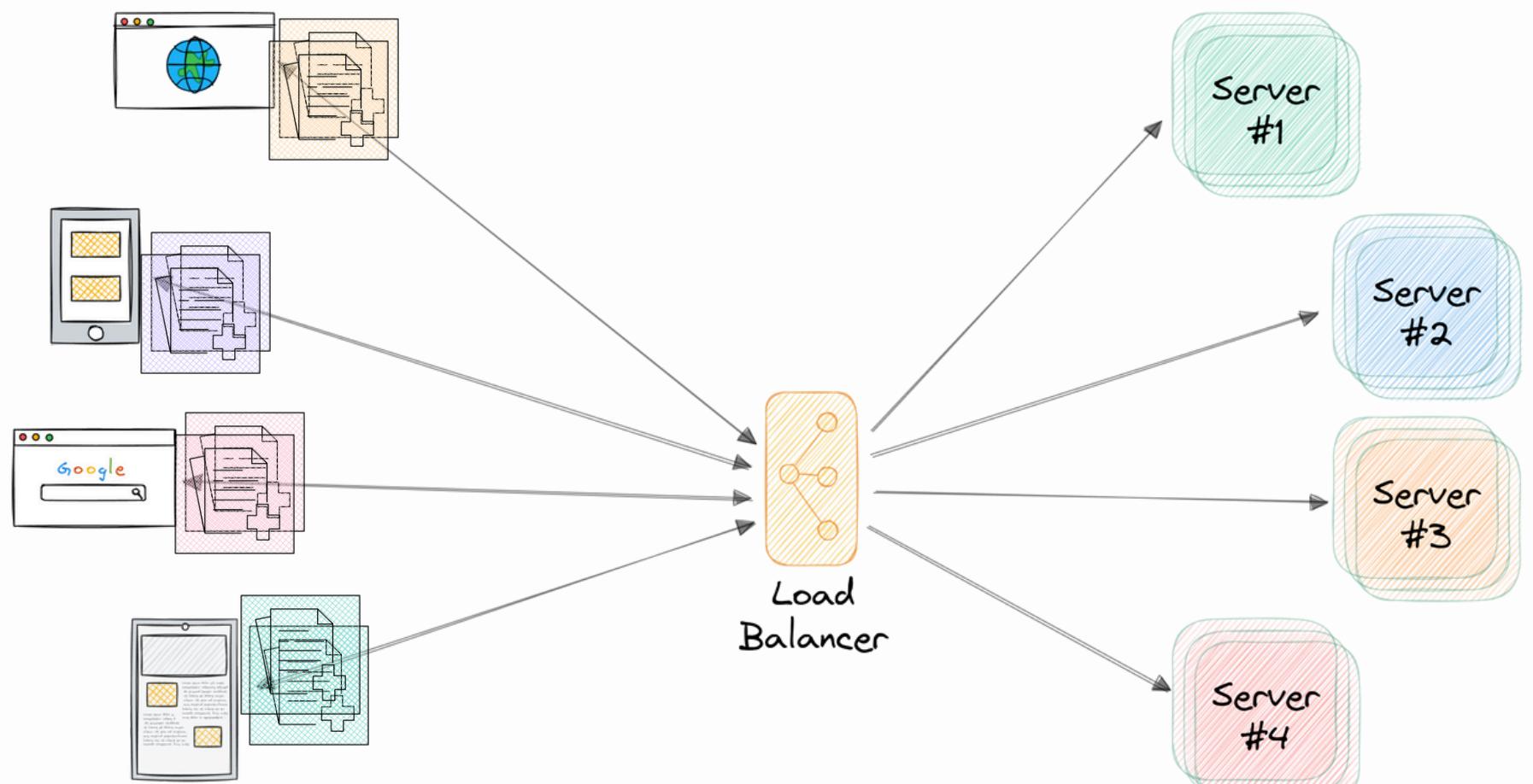
- Distributes requests based on preassigned weights to each server
- Servers with higher capacity receive more requests
- Pros:
 - Accounts for server performance differences
 - Customizable for different server configurations
 - Leads to better resource utilization in a heterogeneous server environment
- Cons:
 - Weights need to be manually calculated and assigned
 - Not ideal for dynamic load changes



```
1 type WeightedRoundRobinStrategy struct {
2     servers      []ServerAddr
3     weights      []float32
4     nextIndex    int
5     selectedTimes int
6 }
7
8 func NewWeightedRoundRobinStrategy(servers []ServerAddr, weights []float32) *WeightedRoundRobinStrategy {
9     // Normalize the weights to 1
10    sum := float32(0)
11    for _, weight := range weights {
12        sum += weight
13    }
14
15    for i, weight := range weights {
16        weights[i] = weight / sum
17    }
18
19    return &WeightedRoundRobinStrategy{servers, weights, 0, 0}
20 }
21
22 func (strategy *WeightedRoundRobinStrategy) ServerAddr(ClientAddr string) (ServerAddr, error) {
23     if len(strategy.servers) == 0 {
24         return "", errors.New("no servers available")
25     }
26
27     server := strategy.servers[strategy.nextIndex]
28     strategy.selectedTimes++
29
30     if strategy.selectedTimes >= int(strategy.weights[strategy.nextIndex]*WeightFactor) {
31         strategy.nextIndex = (strategy.nextIndex + 1) % len(strategy.servers)
32         strategy.selectedTimes = 0
33     }
34
35     return server, nil
36 }
37
38 func (strategy *WeightedRoundRobinStrategy) Connected(addr ServerAddr) {}
39 func (strategy *WeightedRoundRobinStrategy) Disconnected(addr ServerAddr) {}
```

(S) IP Hash

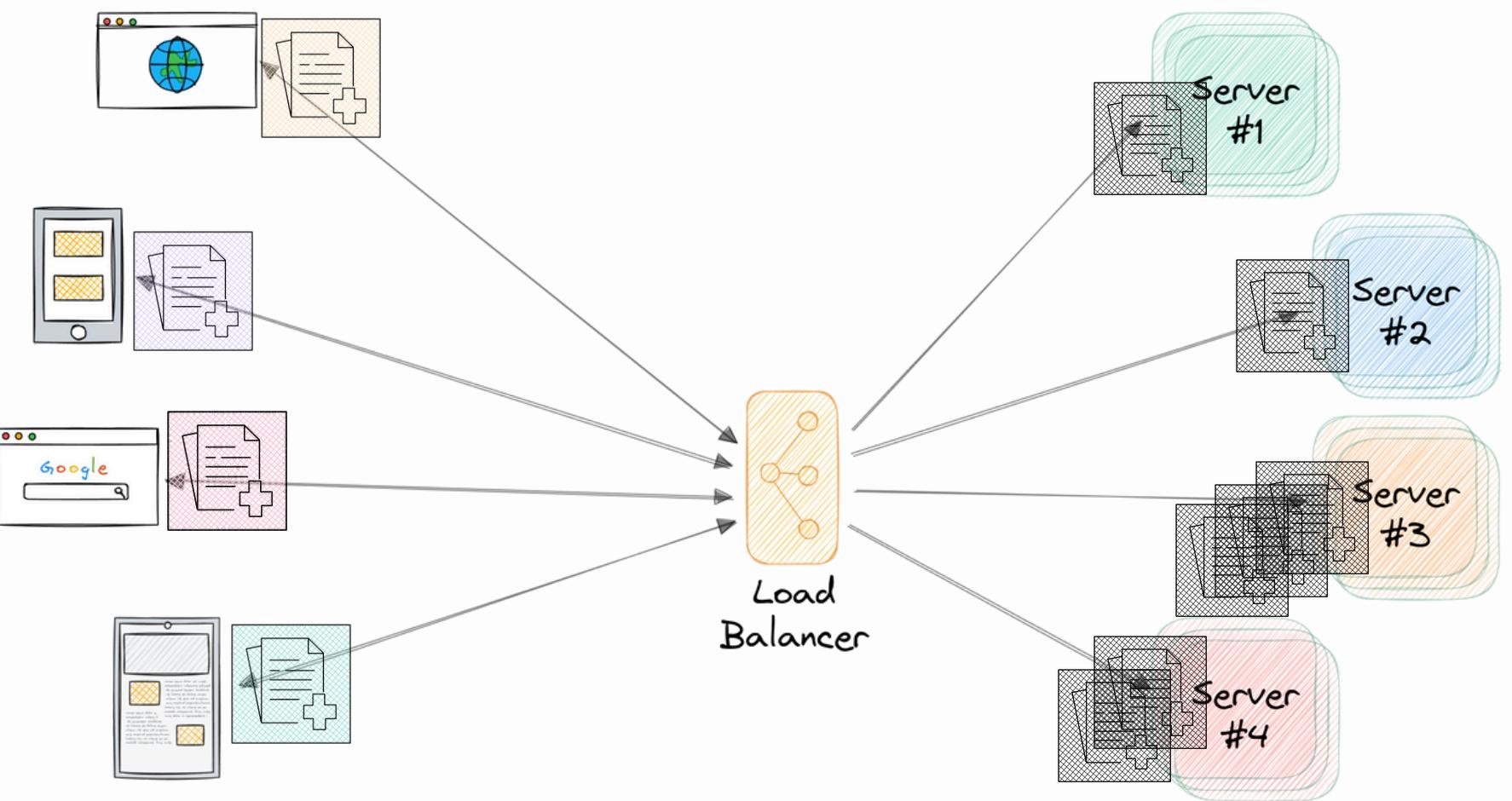
- Uses client IP addresses to consistently direct requests to the same server
- Ensures session persistence or sticky sessions
- Pros:
 - Useful for session-based applications
 - Simple hash-based mapping
- Cons:
 - Load distribution might become uneven if client IPs are clustered and the hash function is weak
 - Not suitable for environments with changing client IPs (e.g., mobile clients) as the session persistence has no benefits in this case



```
1 type IPHashStrategy struct {
2     servers []ServerAddr
3 }
4
5 func NewIPHashStrategy(servers []ServerAddr) *IPHashStrategy {
6     return &IPHashStrategy{servers}
7 }
8
9 func (strategy *IPHashStrategy) ServerAddr(ClientAddr string) (ServerAddr, error) {
10    if len(strategy.servers) == 0 {
11        return "", errors.New("no servers available")
12    }
13
14    fragments := strings.Split(ClientAddr, ":")
15    addr := strings.Join(fragments[:len(fragments)-1], ":")
16
17    // Hash the IP address (Based on the Sum of Bytes of the Client Address) to select a server.
18    // TODO: Implement a better hashing algorithm.
19    hash := 0
20    for i := 0; i < len(addr); i++ {
21        hash += int(addr[i])
22    }
23
24    server := strategy.servers[hash%len(strategy.servers)]
25
26    return server, nil
27 }
28
29 func (strategy *IPHashStrategy) Connected(addr ServerAddr) {}
30 func (strategy *IPHashStrategy) Disconnected(addr ServerAddr) {}
```

(D) Least Connection

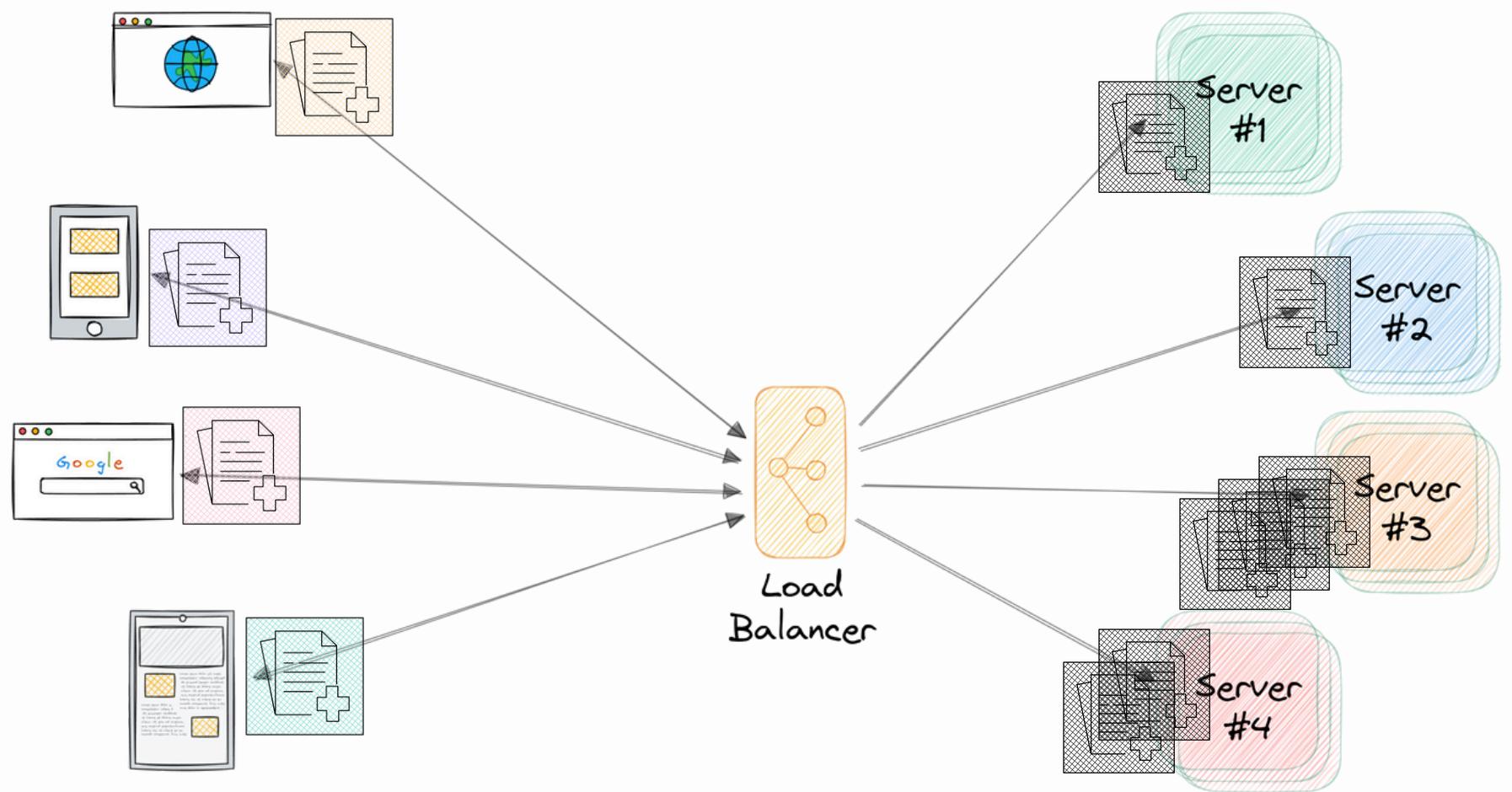
- Directs requests to the server with the fewest active connections
- Active connections provides a good heuristic for server load. Getting server load information like CPU, memory, etc usage repeatedly can induce latency in the system which LBs should aim to minimize.
- Pros:
 - Ensures that servers with less load handle more requests
 - Better performance during traffic spikes
- Cons:
 - More computationally intensive as connections need to be constantly tracked
 - Not always effective with varying request sizes. Active connections is only a heuristic.



```
 1  type LeastConnectionsStrategy struct {
 2      servers      []ServerAddr
 3      connections []int
 4  }
 5
 6  func NewLeastConnectionsStrategy(servers []ServerAddr) *LeastConnectionsStrategy {
 7      connections := make([]int, len(servers))
 8      for i := range connections {
 9          connections[i] = 0
10      }
11
12      return &LeastConnectionsStrategy{servers, connections}
13  }
14
15  func (strategy *LeastConnectionsStrategy) ServerAddr(ClientAddr string) (ServerAddr, error) {
16      if len(strategy.servers) == 0 {
17          return "", errors.New("no servers available")
18      }
19
20      minConnections := strategy.connections[0]
21      minIndex := 0
22
23      for i, connections := range strategy.connections {
24          if connections < minConnections {
25              minConnections = connections
26              minIndex = i
27          }
28      }
29
30      server := strategy.servers[minIndex]
31      strategy.connections[minIndex]++
32
33      return server, nil
34  }
35
36  func (strategy *LeastConnectionsStrategy) Connected(addr ServerAddr) {}
37
38  func (strategy *LeastConnectionsStrategy) Disconnected(addr ServerAddr) {
39      log.Println(strategy.connections)
40      for i, serverAddr := range strategy.servers {
41          if serverAddr == addr && strategy.connections[i] > 0 {
42              strategy.connections[i]--
43              break
44          }
45      }
46      log.Println(strategy.connections)
47  }
```

(D) Weighted Least Connection

- Similar to Least Connections but considers server weights, giving more connections to higher-capacity servers
- Combines load-based balancing with server capacity
- Performs better in both a heterogeneous server environment and one with varying response times
- Pros:
 - More precise load distribution for different server capabilities
 - Suitable for mixed environments (servers with different specs)
- Cons:
 - Higher complexity in terms of configuration and tracking



```
● ● ●
1 type WeightedLeastConnectionsStrategy struct {
2     servers      []ServerAddr
3     connections  []int
4     weights      []float32
5 }
6
7 func WeightedNewLeastConnectionsStrategy(servers []ServerAddr, weights []float32) *WeightedLeastConnectionsStrategy {
8     connections := make([]int, len(servers))
9     for i := range connections {
10         connections[i] = 0
11     }
12
13     // Normalize the weights to 1
14     sum := float32(0)
15     for _, weight := range weights {
16         sum += weight
17     }
18
19     for i, weight := range weights {
20         weights[i] = weight / sum
21     }
22
23     return &WeightedLeastConnectionsStrategy{servers, connections, weights}
24 }
25
26 func (strategy *WeightedLeastConnectionsStrategy) ServerAddr(ClientAddr string) (ServerAddr, error) {
27     if len(strategy.servers) == 0 {
28         return "", errors.New("no servers available")
29     }
30
31     maxRatio := float32(0.0)
32     maxIndex := 0
33
34     totalConnections := 0
35     for _, connections := range strategy.connections {
36         totalConnections += connections
37     }
38
39     // Below code works for Least Connection Strategy
40     for i, connections := range strategy.connections {
41         normalizedConnections := float32(connections) / (float32(totalConnections) + 1e-6)
42
43         ratio := float32(strategy.weights[i]) / (float32(normalizedConnections) + 1e-6)
44
45         if ratio > maxRatio {
46             maxRatio = ratio
47             maxIndex = i
48         }
49     }
50
51     server := strategy.servers[maxIndex]
52     strategy.connections[maxIndex]++
53
54     return server, nil
55 }
56
57 func (strategy *WeightedLeastConnectionsStrategy) Connected(addr ServerAddr) {}
58
59 func (strategy *WeightedLeastConnectionsStrategy) Disconnected(addr ServerAddr) {
60     log.Println(strategy.connections)
61     for i, serverAddr := range strategy.servers {
62         if serverAddr == addr && strategy.connections[i] > 0 {
63             strategy.connections[i]--
64             break
65         }
66     }
67     log.Println(strategy.connections)
68 }
69
```

Pitfalls

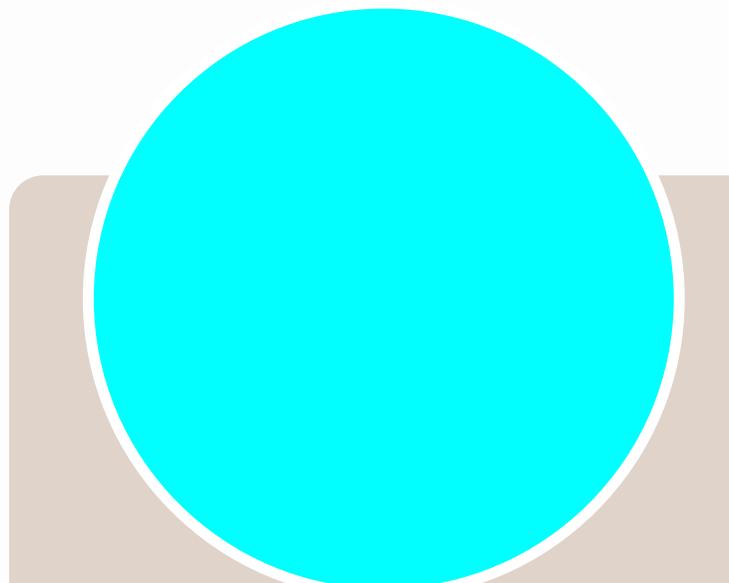
1. Connections handles by Load Balancer
 - a. LB handles 2x TCP connections
 - b. So the connection limit on the LB server can be a bottleneck
2. LB servers need to be scaled horizontally too
3. An orchestrator service is needed for monitoring server health in a more granular way
4. Dynamic handling of servers (on-the-go addition & deletion of servers)

Team Members



Vinayak Anand

S20220010243



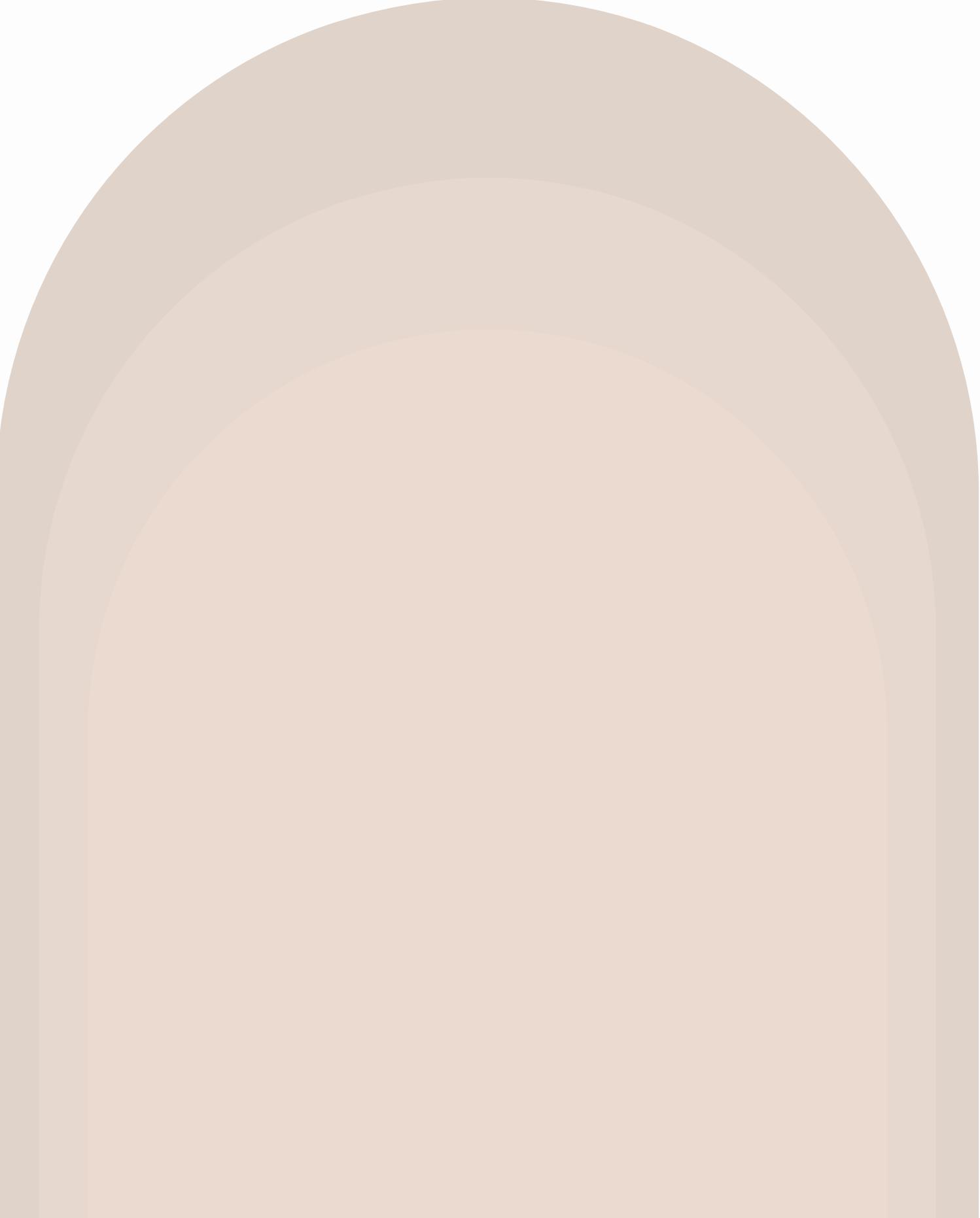
Saurabh Pal

S20220010196



Bishwajeet Sahoo

S20220010039



Thank You...

