# CLOUD COMPUTING APPLICATIONS

## LOAD BALANCER SCHEMES

Prof. Roy Campbell

# What Does a Server Load Balancer (SLB) Do?

- Gets user to needed resource
  - Server must be available
  - User's "session" must not be broken
    - If user must get to the same resource over and over, the SLB device must ensure that happens (i.e., session persistence)

- In order to do work, SLB must
  - Know servers – IP / port, availability
  - Understand details of some protocols (e.g., FTP, SIP)

- Network Address Translation (NAT)
  - Packets are rewritten as they pass through the SLB device

# Reasons to Load-Balance

- Scale applications / services
- Ease of administration / maintenance
  - Easily and transparently remove physical servers from rotation in order to perform any type of maintenance on that server
- Resource sharing
  - Can run multiple instances of an application / service on a server; could be running on a different port for each instance; can load-balance to different port based on data analyzed

# Load-Balancing Algorithms

- Most predominant
  - **Least connections**: Server with fewest number of flows gets the new flow request
  - **Weighted least connections**: Associate a weight / strength for each server and distribute load across server farm based on the weights of all servers in the farm
  - **Round robin**: Round robin through the servers in server farm
  - **Weighted round robin**: Give each server "weight" number of flows in a row;  weight is set just like it is in weighted least flows
- There are other algorithms that look at or try to predict server load in determining the load of the real server

# How SLB Devices Make Decisions

- The SLB device can make its load-balancing decisions based on several factors
  - Some of these factors can be obtained from the packet headers (i.e., IP address, port numbers)
  - Other factors are obtained by looking at the data beyond the network headers.  Examples:
    - HTTP cookies
    - HTTP URLs
    - SSL client certificates
- The decisions can be based strictly on flow counts, or they can be based on knowledge of application
- For some protocols, like FTP, you must have knowledge of protocol to correctly load-balance (i.e., control and data connection must go to same physical server)
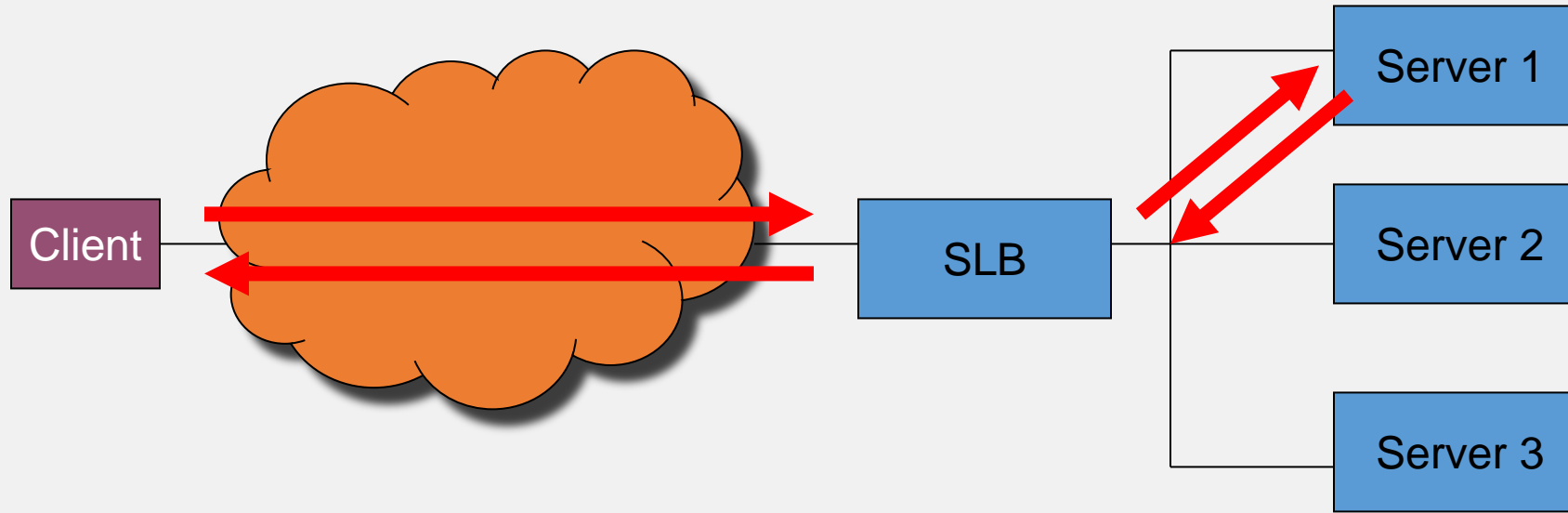
# When a New Flow Arrives

- Determine whether virtual server exists
  - If so, make sure virtual server has available resources
  - If so, then determine level of service needed by that client to that virtual server
    - If virtual machine is configured with particular type of protocol support of session persistence, then do that work
  - Pick a real server for that client
    - The determination of real server is based on flow counts and information about the flow
    - In order to do this, the SLB may need to proxy the flow to get all necessary information for determining the real server; this will be based on the services configured for that virtual server

- If not, the packet is bridged to the correct interface based on Layer 2
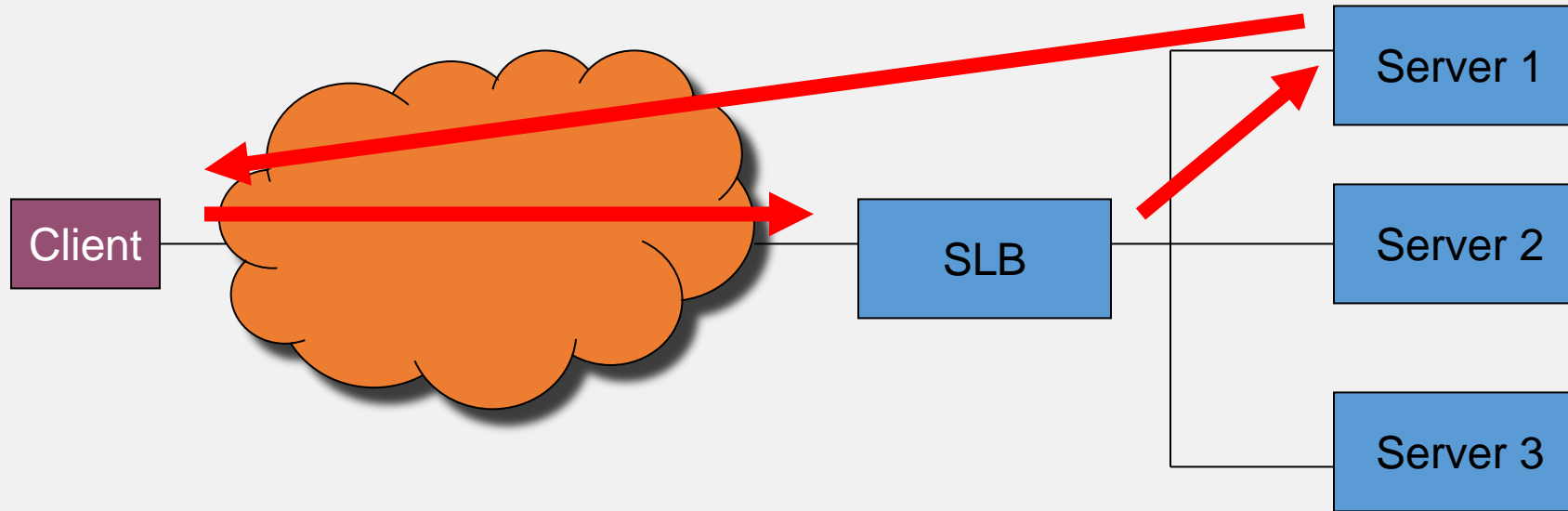
# SLB: Architectures

- Traditional
  - SLB device sits between the Clients and the Servers being load-balanced

- Distributed
  - SLB device sits off to the side and only receives the packets it needs to, based on flow setup and teardown
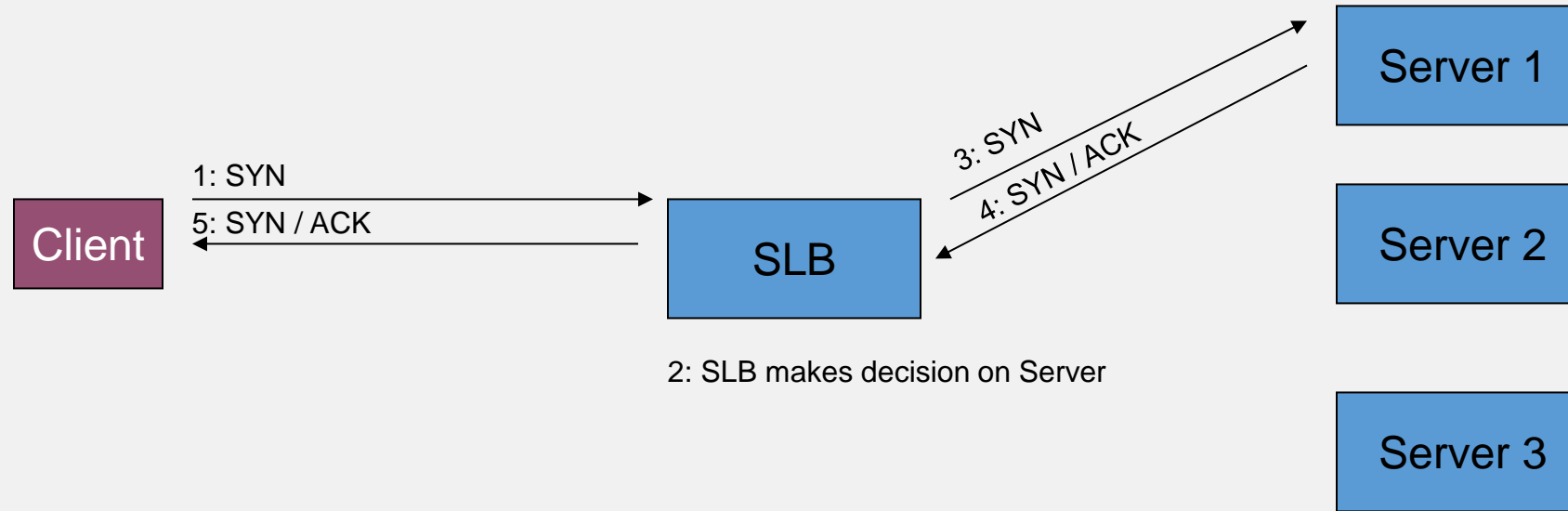
# SLB: Traditional View with NAT

# SLB: Traditional View without NAT
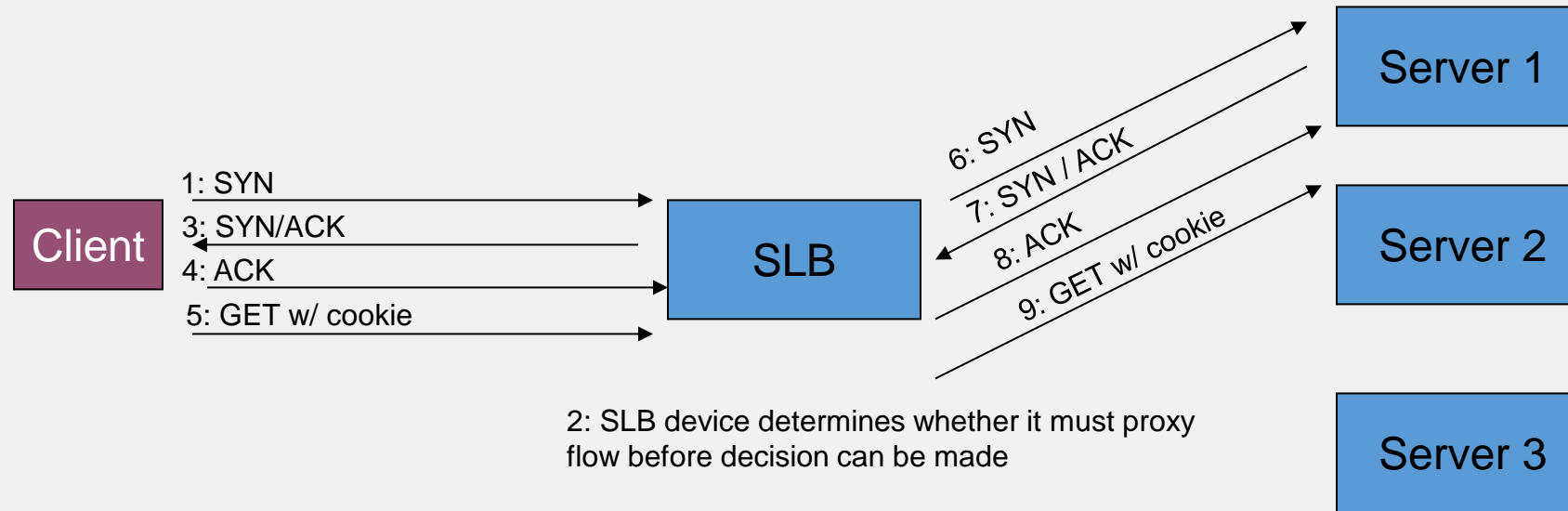
# Load-Balance: Layer 3 / 4

- Look at the destination IP address and port to make a load-balancing decision

- In order to do that, you can determine a real server based on the first packet that arrives

# Layer 3 / 4: Sample Flow

Client

1: SYN

5: SYN / ACK

SLB

3: SYN

4: SYN / ACK

Server 1

Server 2

Server 3
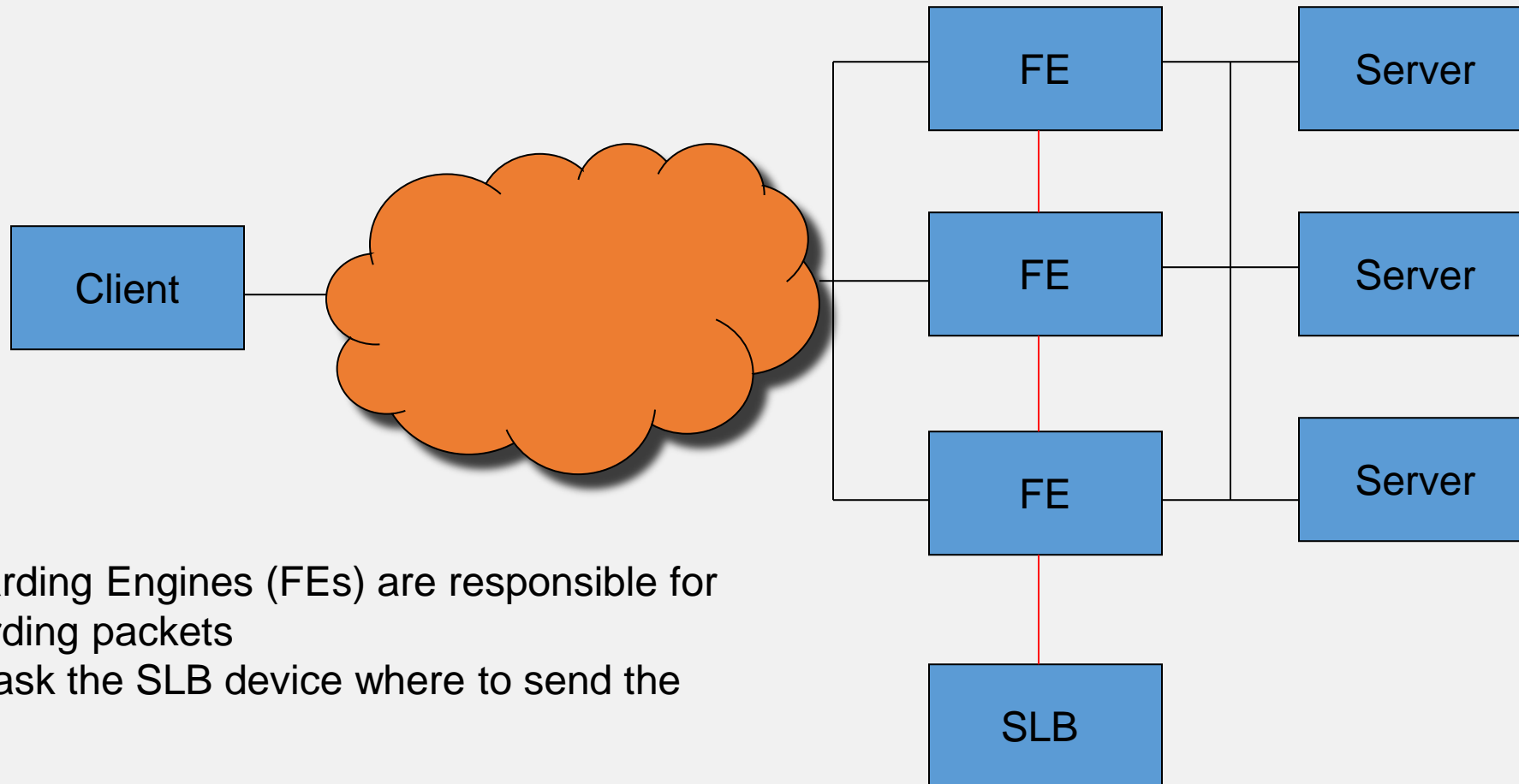
2: SLB makes decision on Server

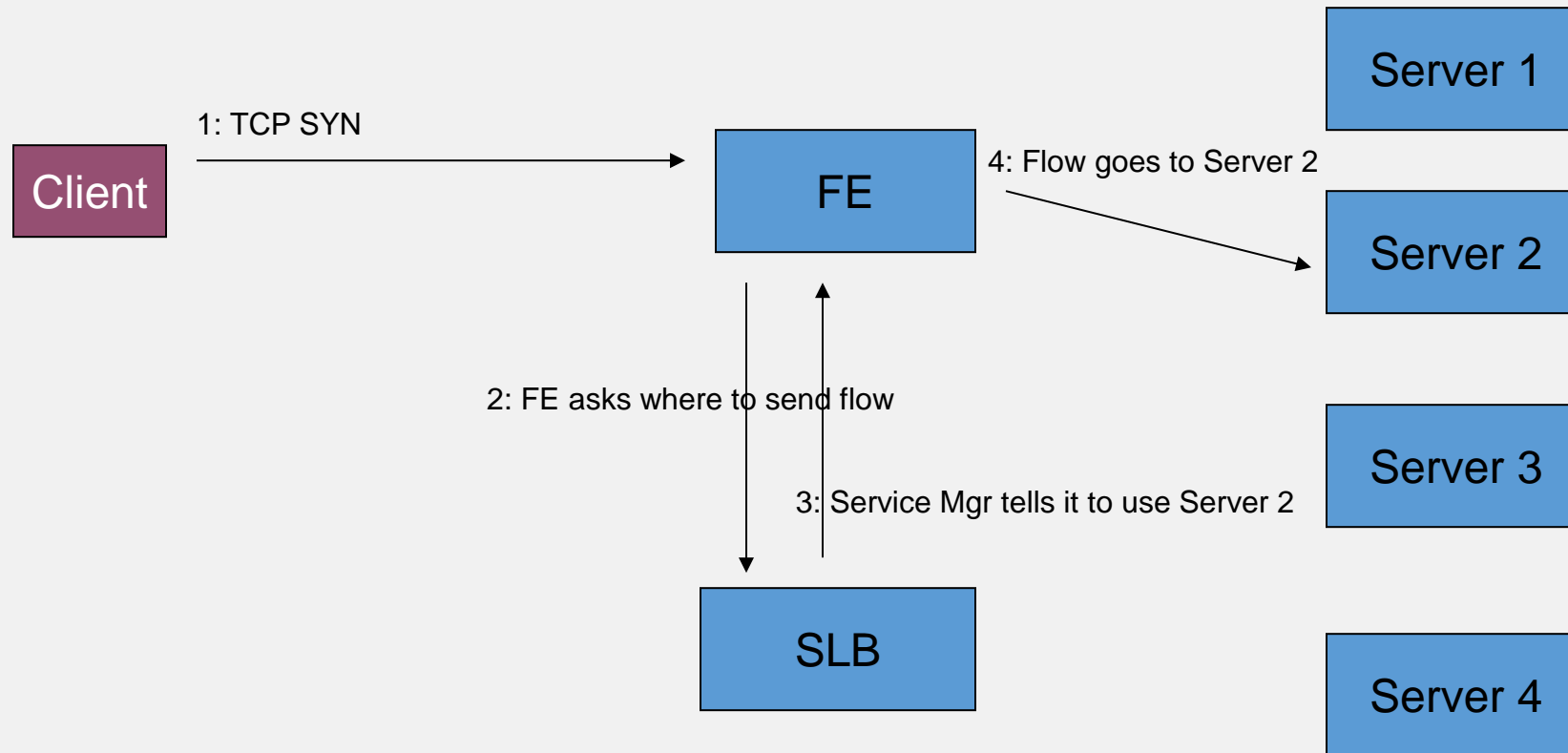Rest of flow continues through HTTP GET and Server response

# Layer 5+: Sample Flow



- Rest of flow continues with Server response
- Note that the flow can be unproxied at this point for efficiency

# SLB: Distributed Architecture



- Forwarding Engines (FEs) are responsible for forwarding packets
- They ask the SLB device where to send the flow

# Distributed Architecture: Sample Flow

**Client**

1: TCP SYN →

**FE**

4: Flow goes to Server 2 →

**Server 1**

**Server 2**

2: FE asks where to send flow

3: Service Mgr tells it to use Server 2

**SLB**

**Server 3**

**Server 4**

- Subsequent packets flow directly from Client to Server 2 through the FE
- The FE must notify the SLB device when the flow ends