# Network Virtualization Overlay Control Protocol Requirements
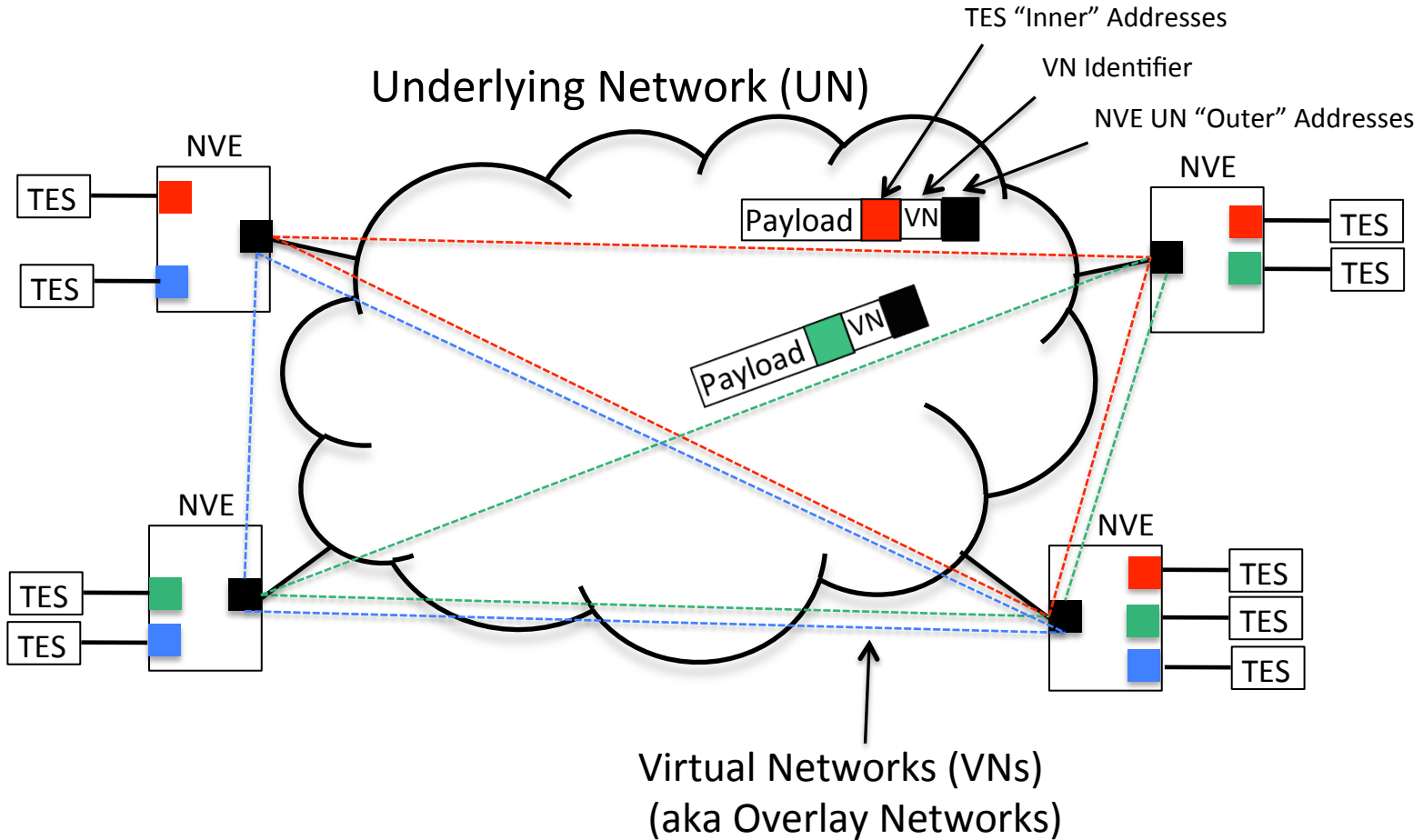
draft-kreeger-nvo3-overlay-cp-00

Lawrence Kreeger, Dinesh Dutt, Thomas Narten, David Black, Murari Sridharan
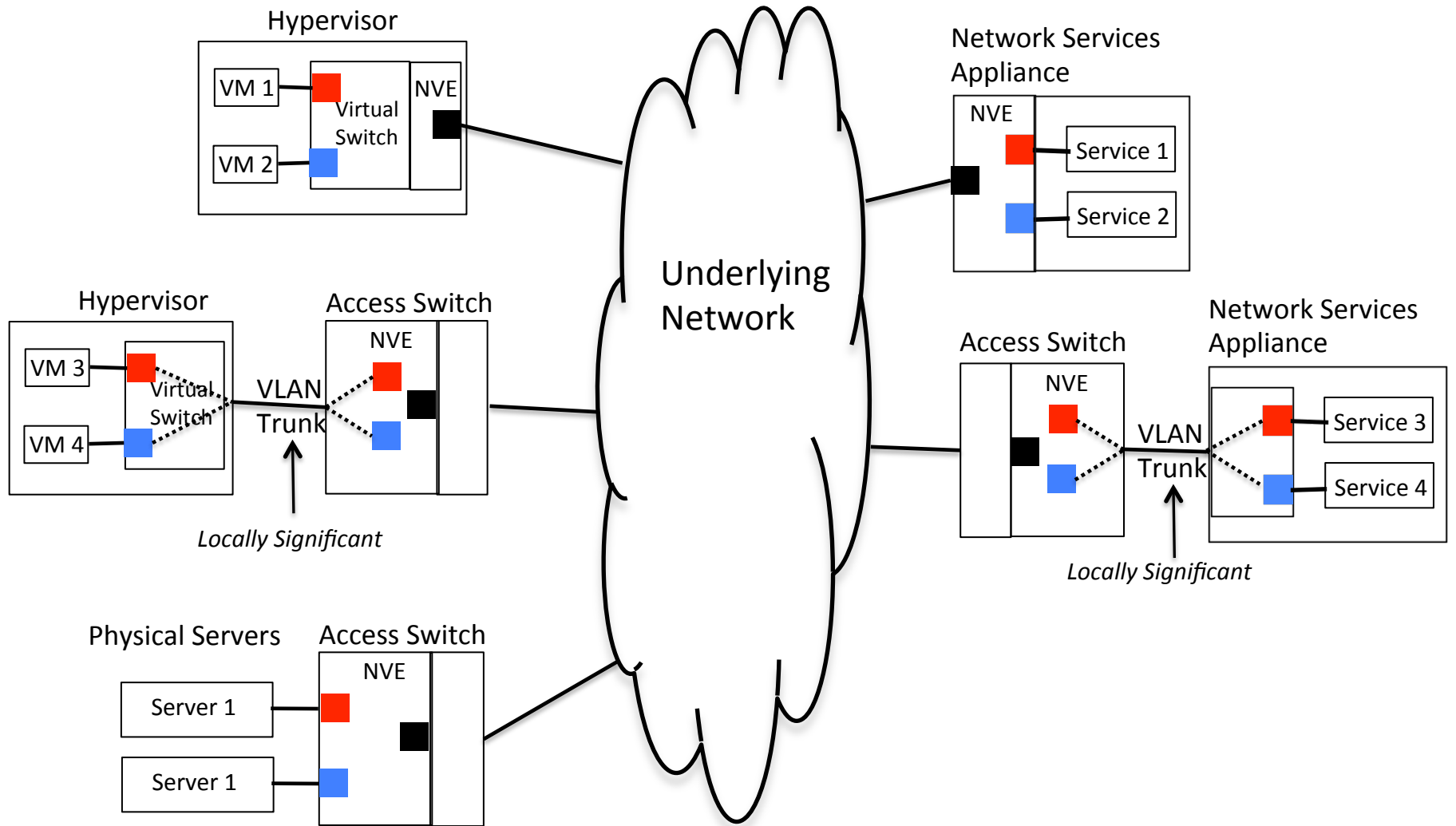
# Purpose

Outline the high level requirements for control protocols needed for overlay virtual networks in highly virtualized data centers.

# Basic Reference Diagram



Network Virtualization Edge (NVE) – (OBP in draft)
Tenant End System (TES) – (End Station in draft)
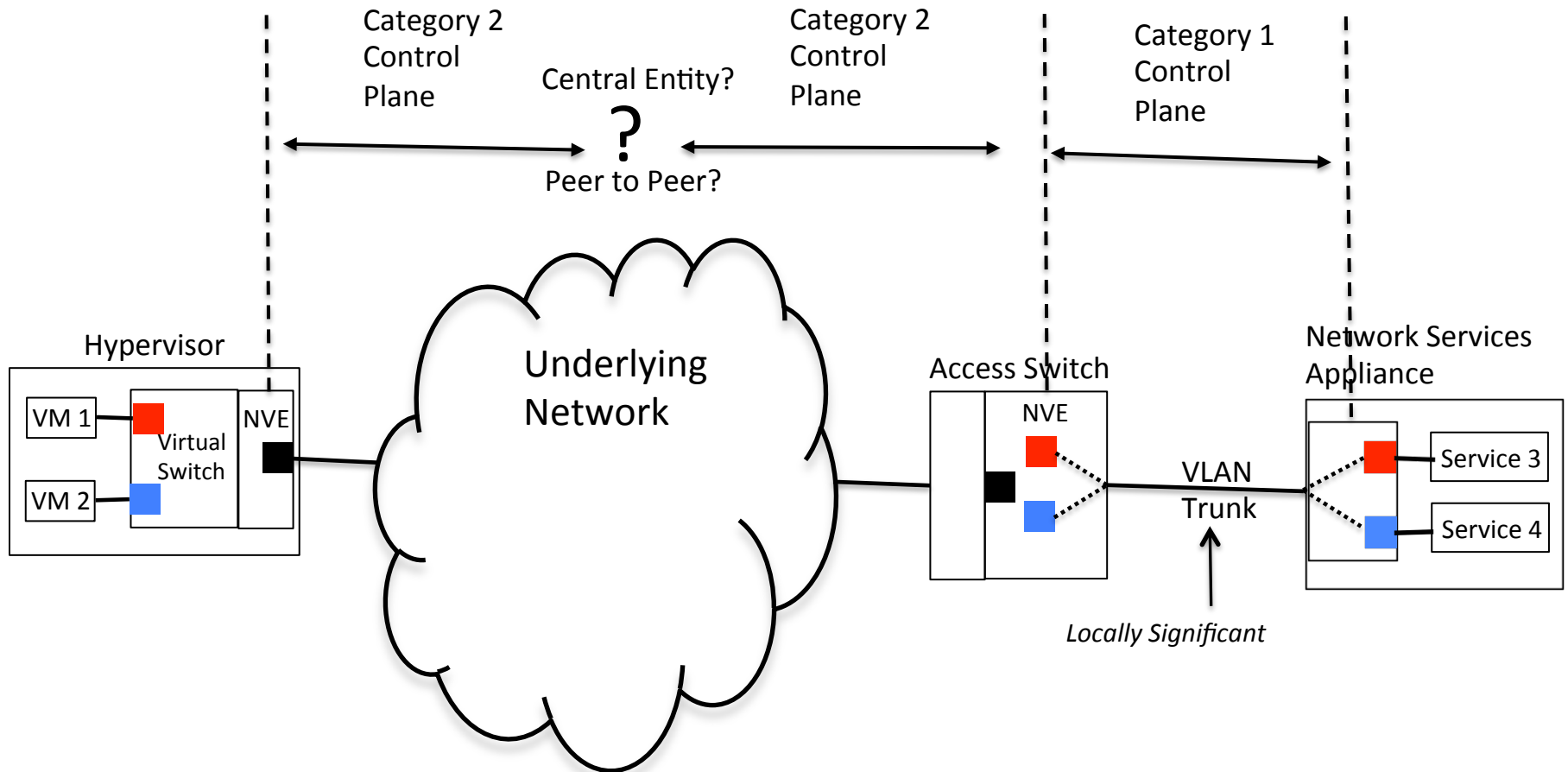
# Possible NVE / TES Scenarios

# Dynamic State Information Needed by an NVE

- Tenant End System (TES) inner address (scoped by Virtual Network (VN)) to outer (Underlying Network (UN)) address of the other Network Virtualization Edge (NVE) used to reach the TES inner address.

- For each VN active on an NVE, a list of UN multicast addresses and/or unicast addresses used to send VN broadcast/multicast packets to other NVEs forwarding to TES for the VN.

- For a given VN, the Virtual Network ID (VN-ID) to use in packets sent across the UN.

- If the TES is not within the same device as the NVE, the NVE needs to know the physical port to reach a given inner address.
  - If multiple VNs are reachable over the same physical port, some kind of tag (e.g. VLAN tag) is needed to keep the VN traffic separated over the wire.

# Two Main Categories of Control Planes

1. For an NVE to obtain dynamic state for communicating with a TES located on a different physical device (e.g. hypervisor or Network Services Appliance).

2. For an NVE to obtain dynamic state for communicating across the Underlying Network to other NVEs.

# Control Plane Category Reference Diagram

Category 2
Control
Plane

Central Entity?

**?**

Category 2
Control
Plane

Category 1
Control
Plane

Peer to Peer?

Hypervisor

VM 1

VM 2

Virtual
Switch

NVE

Underlying
Network

Access Switch

NVE

VLAN
Trunk

Network Services
Appliance
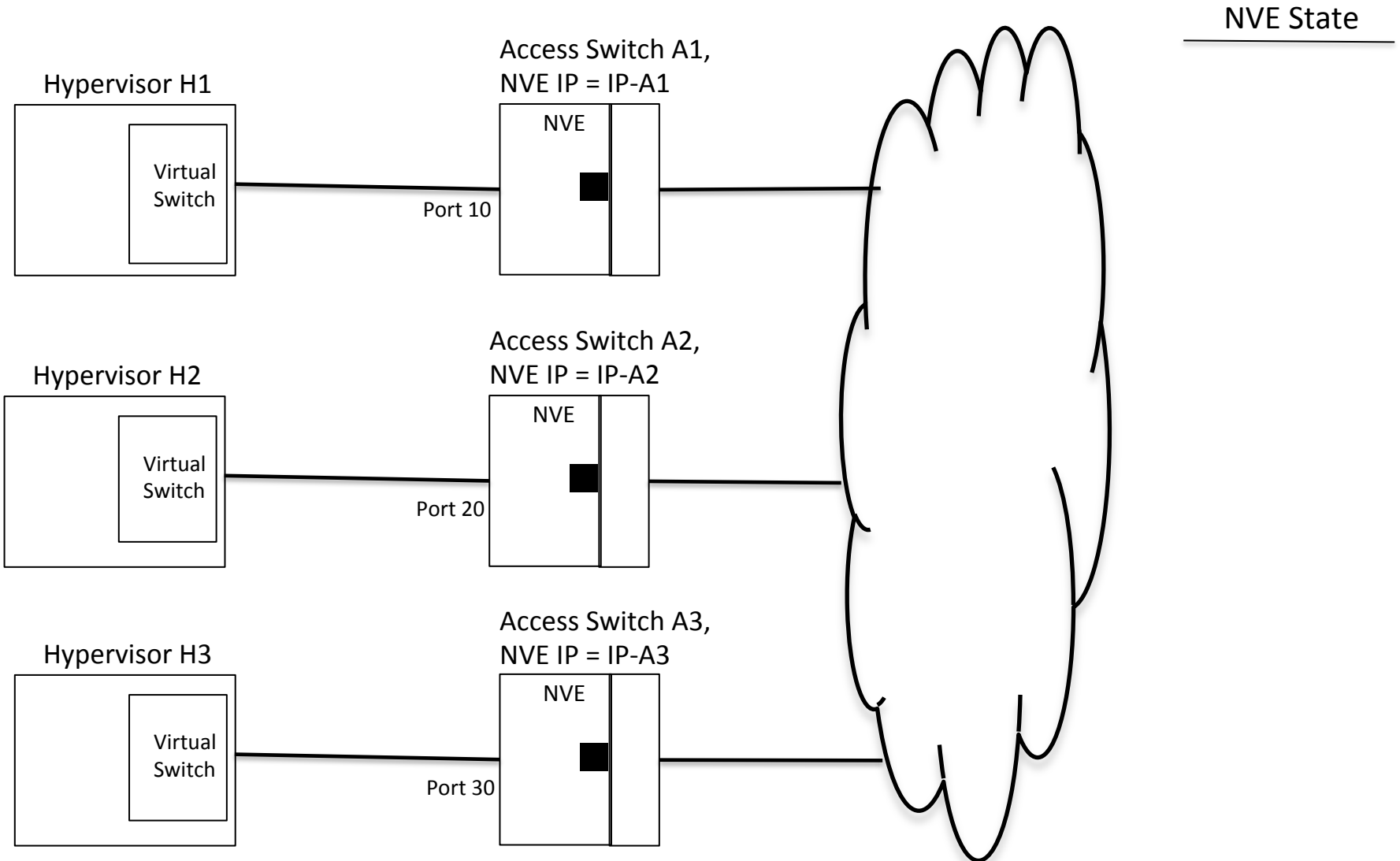
Service 3

Service 4

*Locally Significant*

# Category 2 CP Architecture Possibilities

- Central entity is populated by DC orchestration system
- Central entity is populated by Push from NVE
- Push to NVE from central entity
- Pull from NVE from central entity
- Peer to Peer exchange between NVEs with no central entity
- Central entity could be a monolithic system or a distributed system

# Possible Example CP Scenario

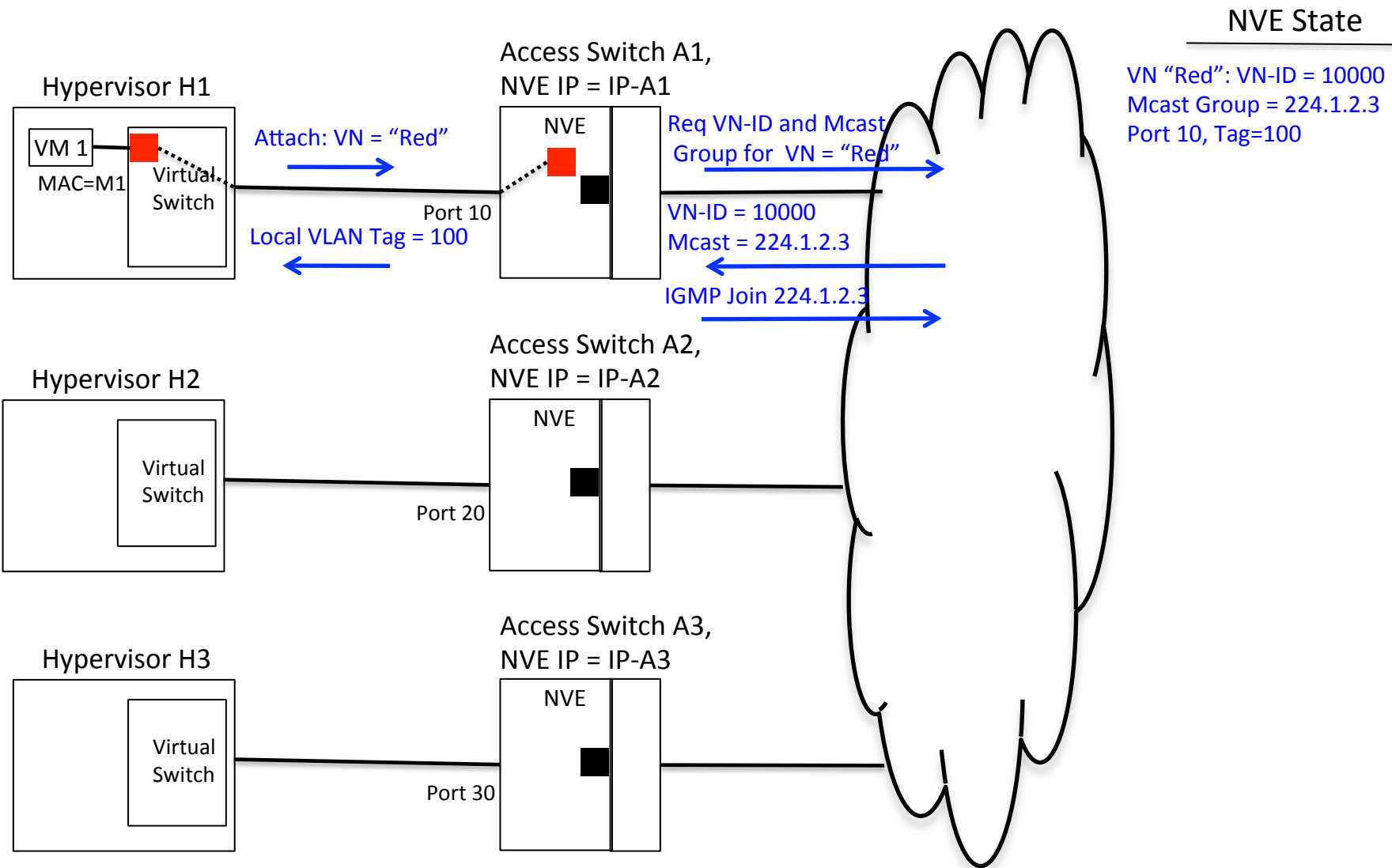This example is not part of the Req draft and is shown for illustrative purposes
Assumes: Central entity with push/pull from NVE, Multicast Enabled IP Underlay

NVE State

Hypervisor H1

Access Switch A1,
NVE IP = IP-A1

Virtual Switch

NVE

Port 10

Hypervisor H2

Access Switch A2,
NVE IP = IP-A2

Virtual Switch

NVE

Port 20

Hypervisor H3

Access Switch A3,
NVE IP = IP-A3

Virtual Switch

NVE

Port 30

# VM 1 comes up on Hypervisor H1, connected the VN "Red"

## H1's Virtual Switch signals to A1 that it needs attachment to VN "Red"

Hypervisor H1

VM 1

MAC=M1

Virtual Switch

Attach: VN = "Red"

Local VLAN Tag = 100

Access Switch A1,
NVE IP = IP-A1

NVE

Port 10

Req VN-ID and Mcast Group for VN = "Red"

VN-ID = 10000
Mcast = 224.1.2.3

IGMP Join 224.1.2.3

NVE State

VN "Red": VN-ID = 10000
Mcast Group = 224.1.2.3
Port 10, Tag=100

Hypervisor H2

Virtual Switch

Access Switch A2,
NVE IP = IP-A2

NVE

Port 20

Hypervisor H3

Virtual Switch

Access Switch A3,
NVE IP = IP-A3

NVE

Port 30

# VM 1 comes up on Hypervisor H1, connected the VN "Red"

## H1's Virtual Switch signals to A1 that MAC M1 is connected to VN "Red"

**NVE State**

VN "Red": VN-ID = 10000
Mcast Group = 224.1.2.3
Port 10, Tag=100
MAC = M1 in "Red" on Port 10

Hypervisor H1

VM 1

MAC=M1

Virtual Switch

Access Switch A1,
NVE IP = IP-A1

NVE

Port 10

Attach: MAC = M1
in VN "Red"

Register MAC = M1 in
VN "Red"  reachable
at IP-A1

Hypervisor H2

Virtual Switch

Access Switch A2,
NVE IP = IP-A2

NVE

Port 20

Hypervisor H3
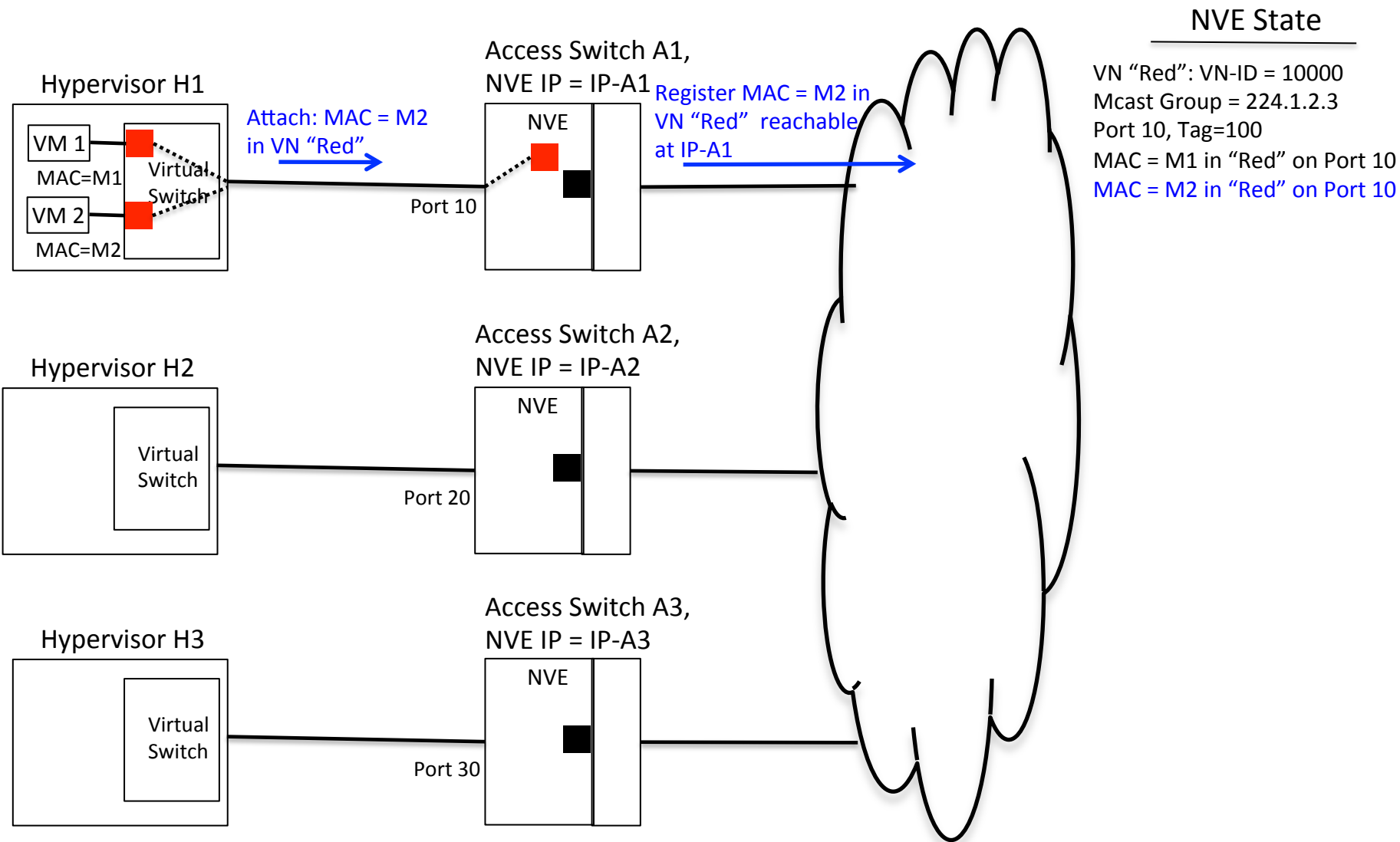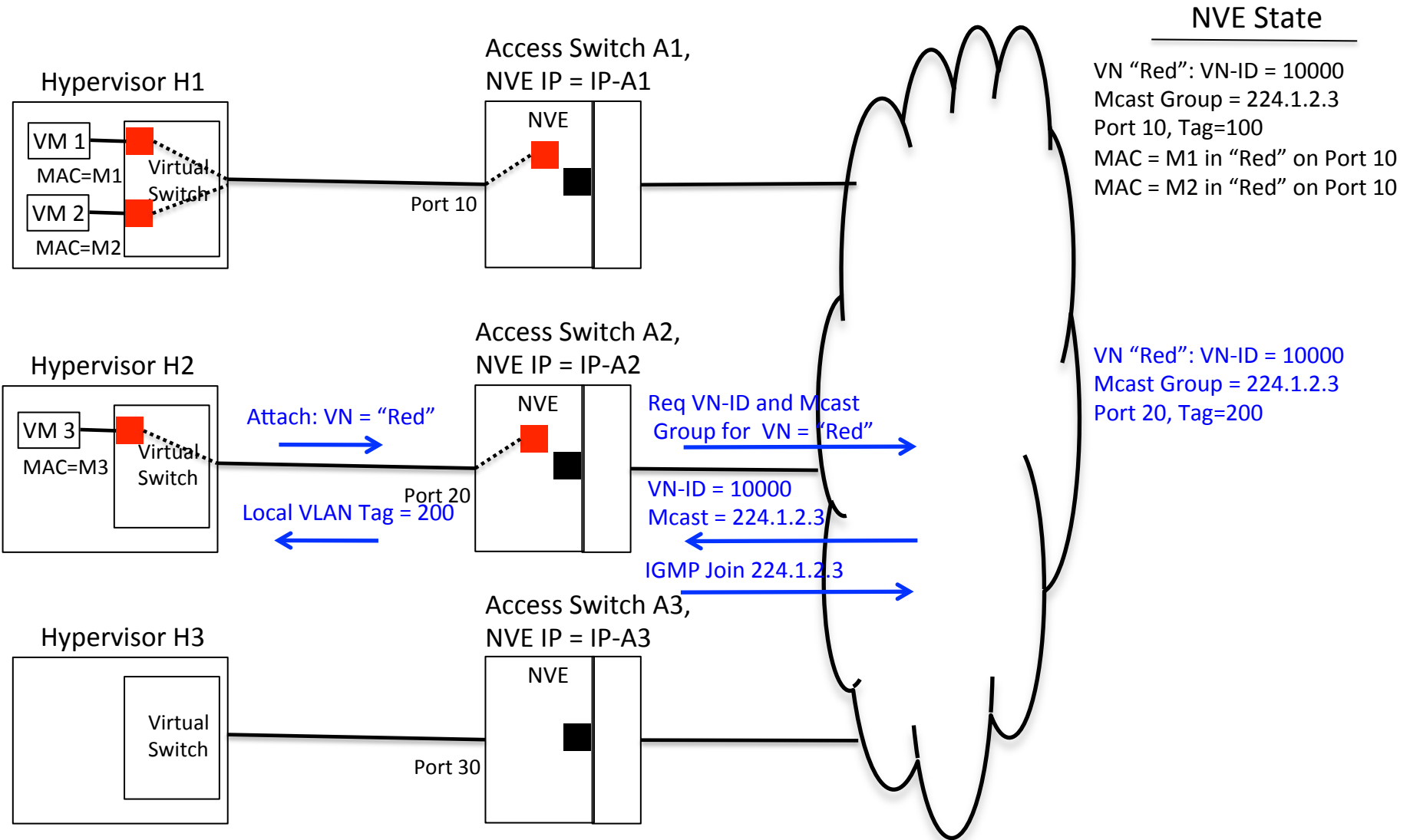
Virtual Switch

Access Switch A3,
NVE IP = IP-A3

NVE

Port 30

# VM 2 comes up on Hypervisor H1, connected the VN "Red"

## H1's Virtual Switch signals to A1 that MAC M2 is connected to VN "Red"



**Hypervisor H1**

VM 1
MAC=M1
Virtual Switch
VM 2
MAC=M2

Attach: MAC = M2 in VN "Red"

**Access Switch A1, NVE IP = IP-A1**

NVE
Port 10

Register MAC = M2 in VN "Red" reachable at IP-A1

**NVE State**

VN "Red": VN-ID = 10000
Mcast Group = 224.1.2.3
Port 10, Tag=100
MAC = M1 in "Red" on Port 10
MAC = M2 in "Red" on Port 10

**Hypervisor H2**

Virtual Switch

**Access Switch A2, NVE IP = IP-A2**

NVE
Port 20

**Hypervisor H3**
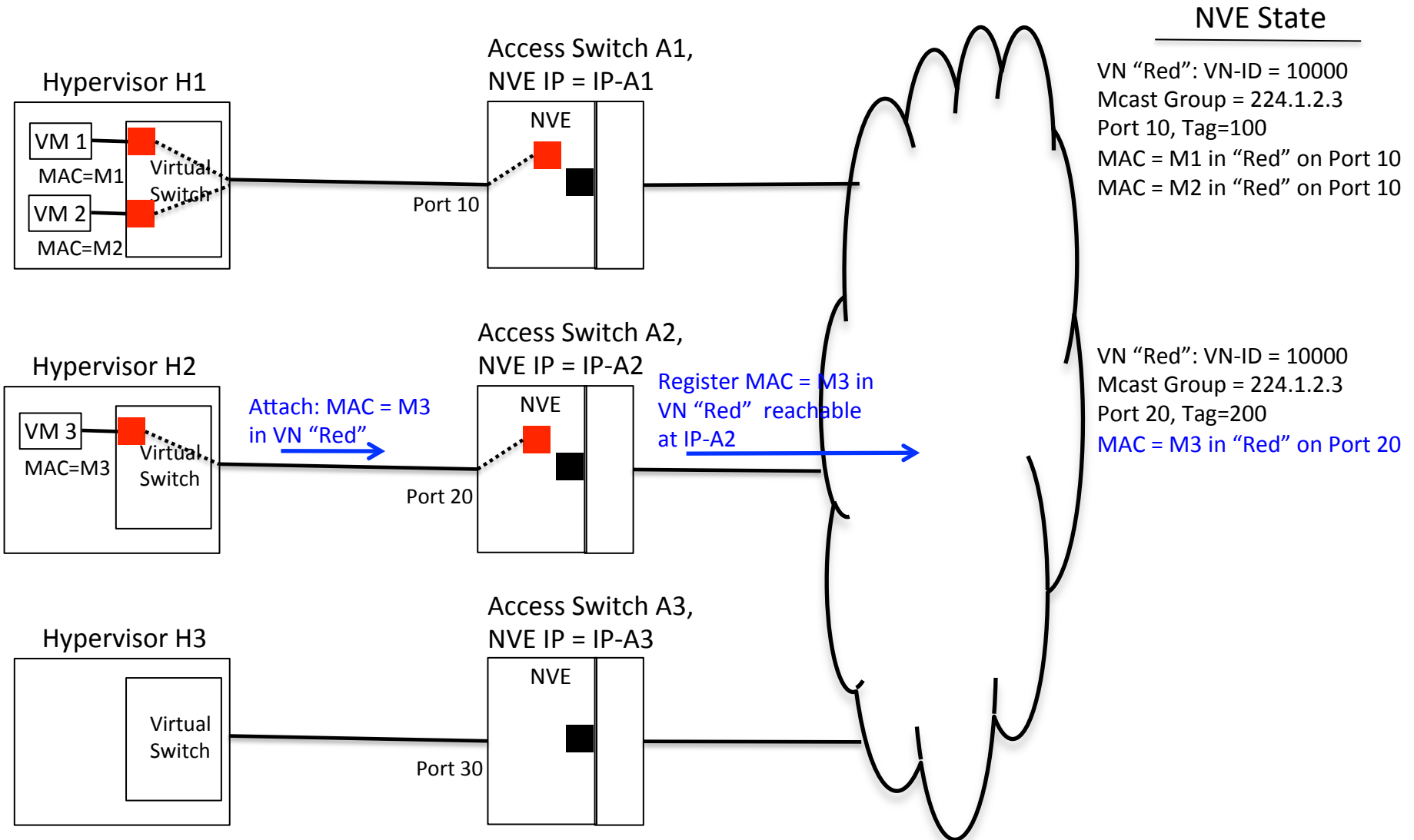
Virtual Switch

**Access Switch A3, NVE IP = IP-A3**

NVE
Port 30

# VM 3 comes up on Hypervisor H2, connected the VN "Red"

## H2's Virtual Switch signals to A2 that it needs attachment to VN "Red"

**NVE State**

VN "Red": VN-ID = 10000
Mcast Group = 224.1.2.3
Port 10, Tag=100
MAC = M1 in "Red" on Port 10
MAC = M2 in "Red" on Port 10

VN "Red": VN-ID = 10000
Mcast Group = 224.1.2.3
Port 20, Tag=200

Hypervisor H1

VM 1
MAC=M1

Virtual Switch

VM 2
MAC=M2

Access Switch A1,
NVE IP = IP-A1

NVE

Port 10

Hypervisor H2

VM 3
MAC=M3

Virtual Switch

Access Switch A2,
NVE IP = IP-A2

NVE

Port 20

Attach: VN = "Red"

Local VLAN Tag = 200

Req VN-ID and Mcast
Group for  VN = "Red"

VN-ID = 10000
Mcast = 224.1.2.3

IGMP Join 224.1.2.3

Hypervisor H3

Virtual Switch

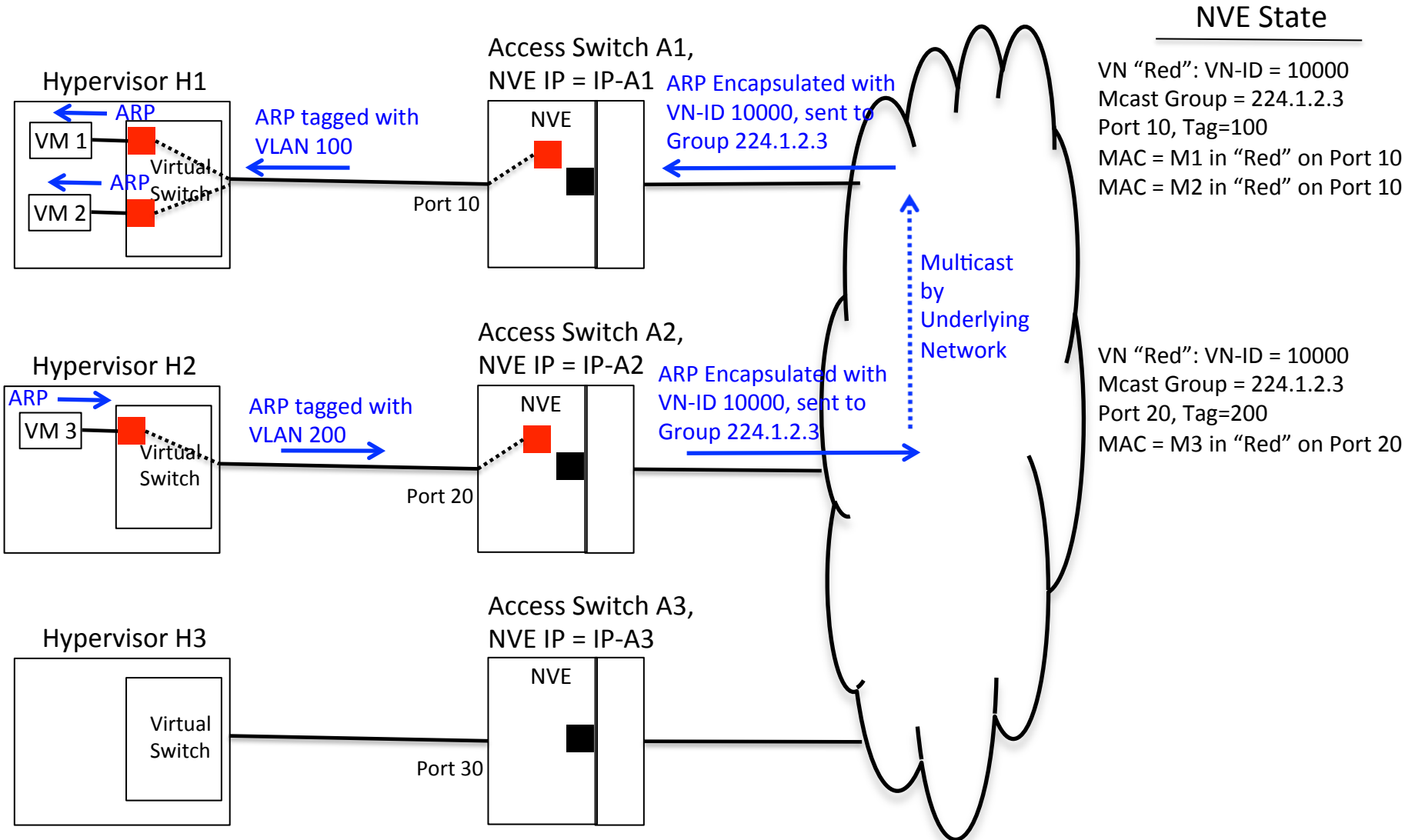Access Switch A3,
NVE IP = IP-A3

NVE

Port 30

# VM 3 comes up on Hypervisor H2, connected the VN "Red"

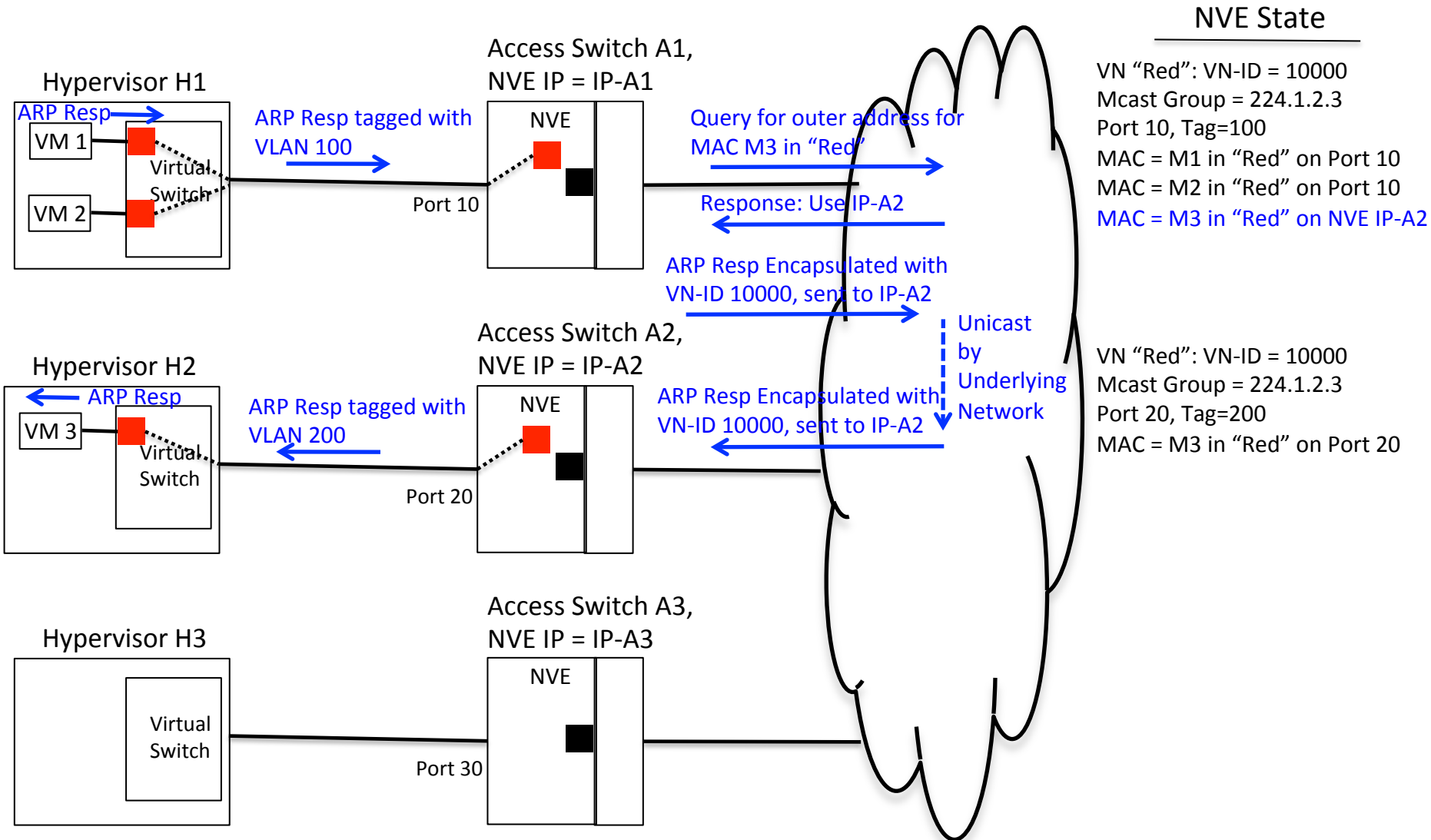## H2's Virtual Switch signals to A2 that MAC M3 is connected to VN "Red"

# VM 3 ARPs for VM1

NVE A2 uses multicast to send the ARP Bcast to all NVEs interested in VN "Red"
NVE A1 Queries to find inner to outer mapping for MAC M3

NVE State

VN "Red": VN-ID = 10000
Mcast Group = 224.1.2.3
Port 10, Tag=100
MAC = M1 in "Red" on Port 10
MAC = M2 in "Red" on Port 10

VN "Red": VN-ID = 10000
Mcast Group = 224.1.2.3
Port 20, Tag=200
MAC = M3 in "Red" on Port 20

Hypervisor H1

Access Switch A1,
NVE IP = IP-A1

ARP Encapsulated with
VN-ID 10000, sent to
Group 224.1.2.3

VM 1

ARP

ARP

VM 2

Virtual Switch

NVE

ARP tagged with
VLAN 100

Port 10

Multicast by Underlying Network

Hypervisor H2

Access Switch A2,
NVE IP = IP-A2

ARP Encapsulated with
VN-ID 10000, sent to
Group 224.1.2.3

ARP

VM 3

Virtual Switch

ARP tagged with
VLAN 200

NVE

Port 20

Hypervisor H3

Access Switch A3,
NVE IP = IP-A3

Virtual Switch

NVE

Port 30

# VM 1 Sends ARP Response to VM3

NVE A1 Queries central entity to find inner to outer mapping for MAC M3
NVE A1 Unicasts ARP Response to A2

# Summary of CP Characteristics

- Lightweight for NVE
  - This means:
    - Low amount of state (only what is needed at the time)
    - Low on complexity (keep it simply)
    - Low on overhead (don't drain resources from NVE)
    - Highly Scalable (don't collapse when scaled)
- Extensible
  - Support multiple address families (e.g. IPv4 and IPv6)
  - Allow addition of new address families
- Quickly reactive to change
  - Support Live Migration of VMs

# Conclusion

- Two Categories of Control Plane protocols are needed to support a dynamic virtualized data center to dynamically build the state needed by an NVE to perform its map+encap and decap +deliver function.

- There are several models of operation possible which the WG will need to decide on.

- To help in deciding, the draft contains important evaluation criteria to use for comparing proposed solutions.