# VIRTUAL NETWORK OVER TRILL
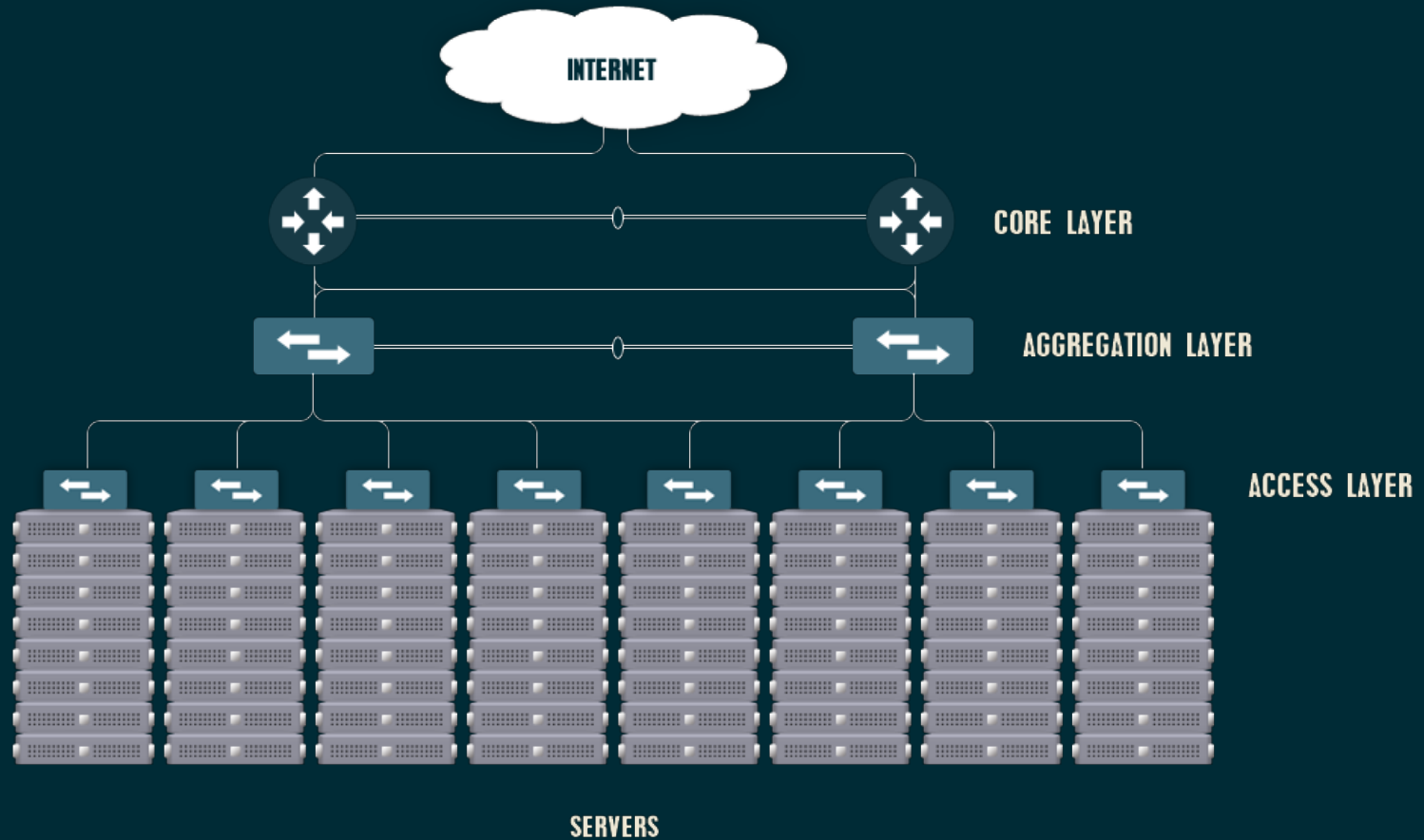
## DESIGN, IMPLEMENTATION AND DEMONSTRATION

**William Dauchy** - Gandi.net

**Kernel Recipes 2013**

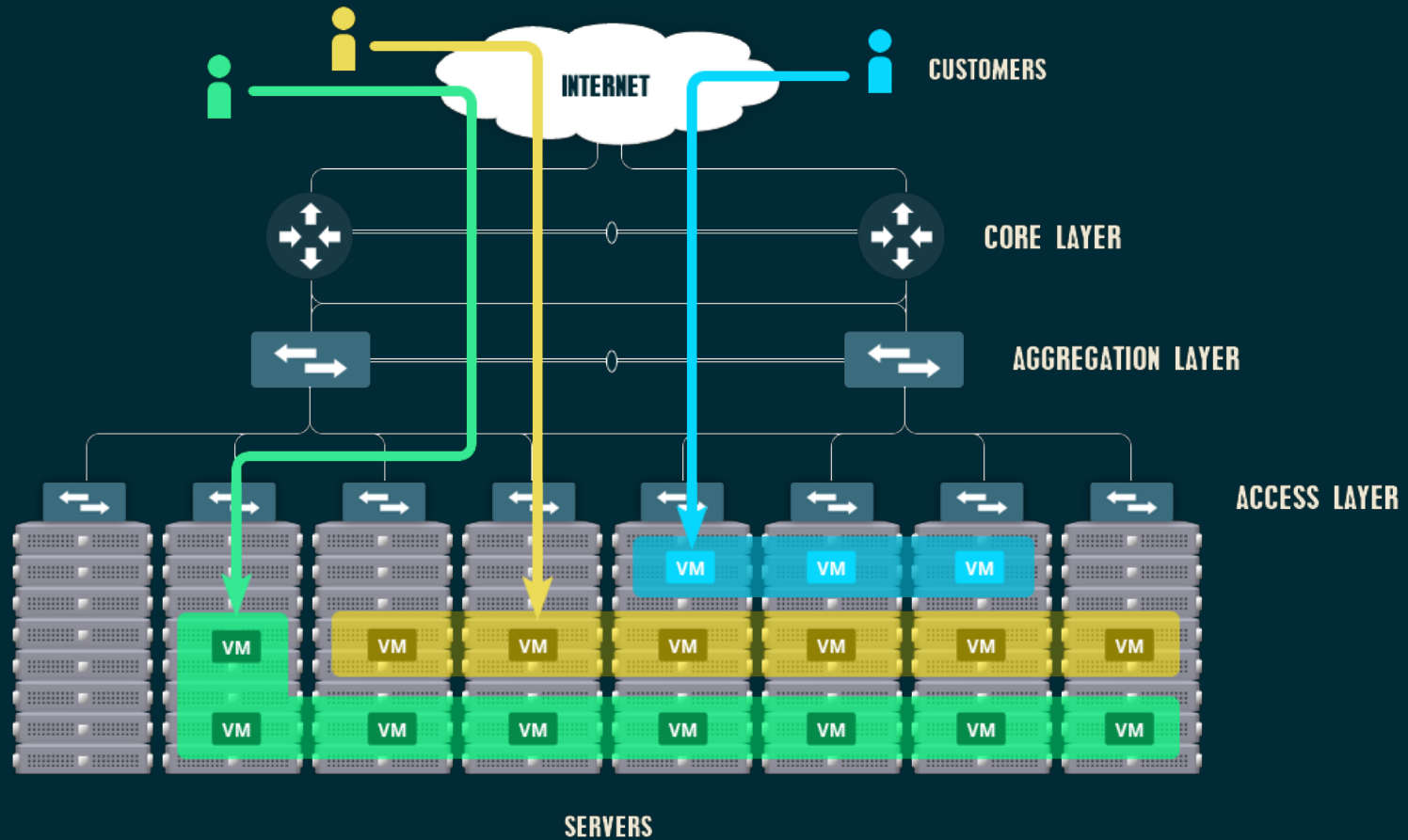gandi.net     Kernel Recipes

# CONVENTIONAL DATA CENTER



INTERNET

CORE LAYER

AGGREGATION LAYER

ACCESS LAYER

SERVERS

# MAIN GOAL

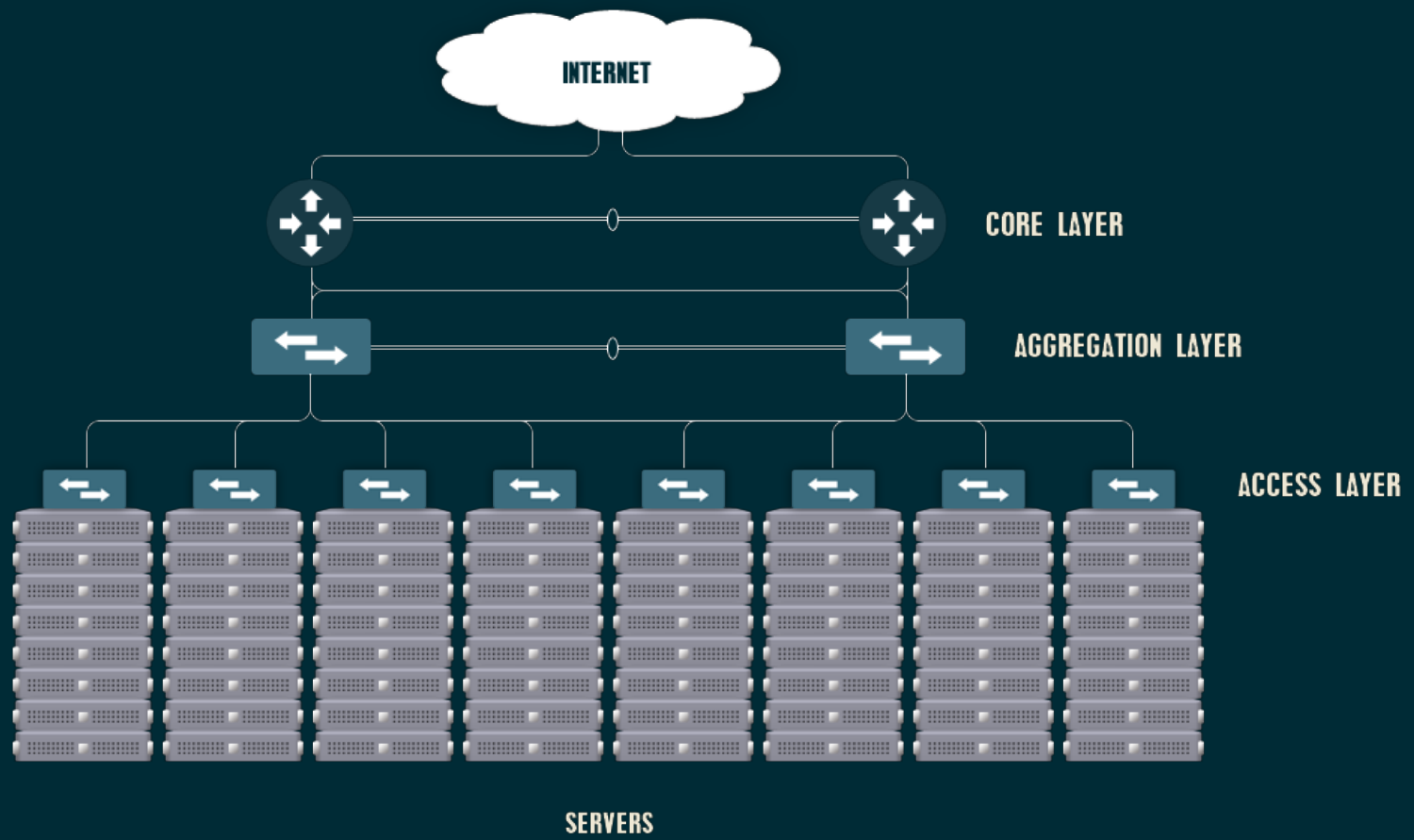- provide large scale multi-tenancy

# REQUIREMENTS

- Seamless VM mobility
- Easy management
- Layer 2 core scaling
- Fault resiliance
- VLAN scalability

# LAYER 2 - SWITCHING BENEFITS

- Management simplified + Plug & play
- Seamless Virtual Machine mobility
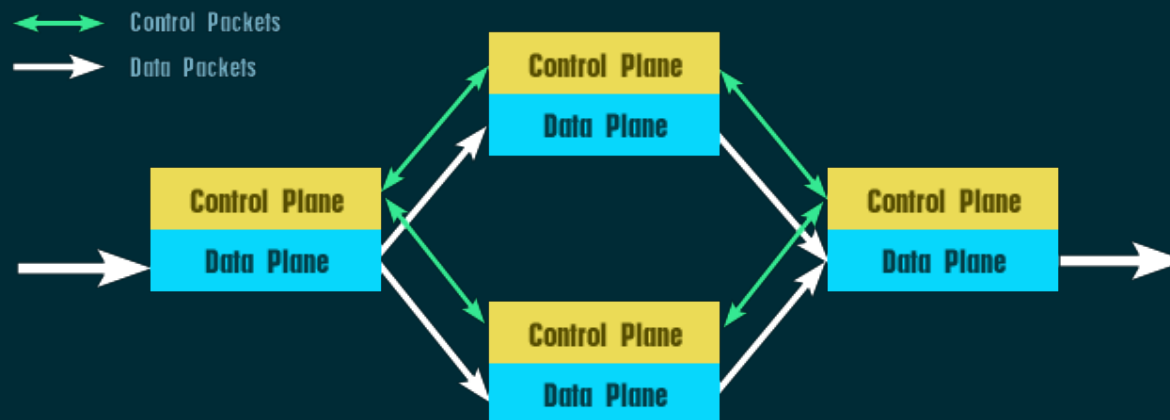- Auto learning + determistic failover

# LAYER 2 - SWITCHING LIMITATION

- A large number of tenants implies
  - a huge number of MAC address in switch table (TCAM overflow)
  - ARP storm at nodes
- STP to ensure a loop free topology
  - blocking redundant paths
  - Core-computes required, recomputes when topology changes
- Number of VLANs is limited to 4096

INTERNET

CORE LAYER

AGGREGATION LAYER

ACCESS LAYER

SERVERS

# WHAT IS TRILL

- New device: RBridge
  - Control plane
  - Data plane



- Encapsulate native frames in a transport header
- Providing a hop count and nickname
- Route the encapsulated frames using IS-IS
- Decapsulate native frames before delivery

# IETF STANDARD

- RFC 5556 Transparent Interconnection of Lots of Links (TRILL): Problem and Applicability Statement
- RFC 6325 Routing Bridges (RBridges): Base Protocol Specification
- RFC 6326 Transparent Interconnection of Lots of Links (TRILL) Use of IS-IS
- RFC 6327 Routing Bridges (RBridges): Adjacency
- RFC 6439 Routing Bridges (RBridges): Appointed Forwarders
- RFC 6361 PPP Transparent Interconnection of Lots of Links (TRILL) Protocol Control Protocol
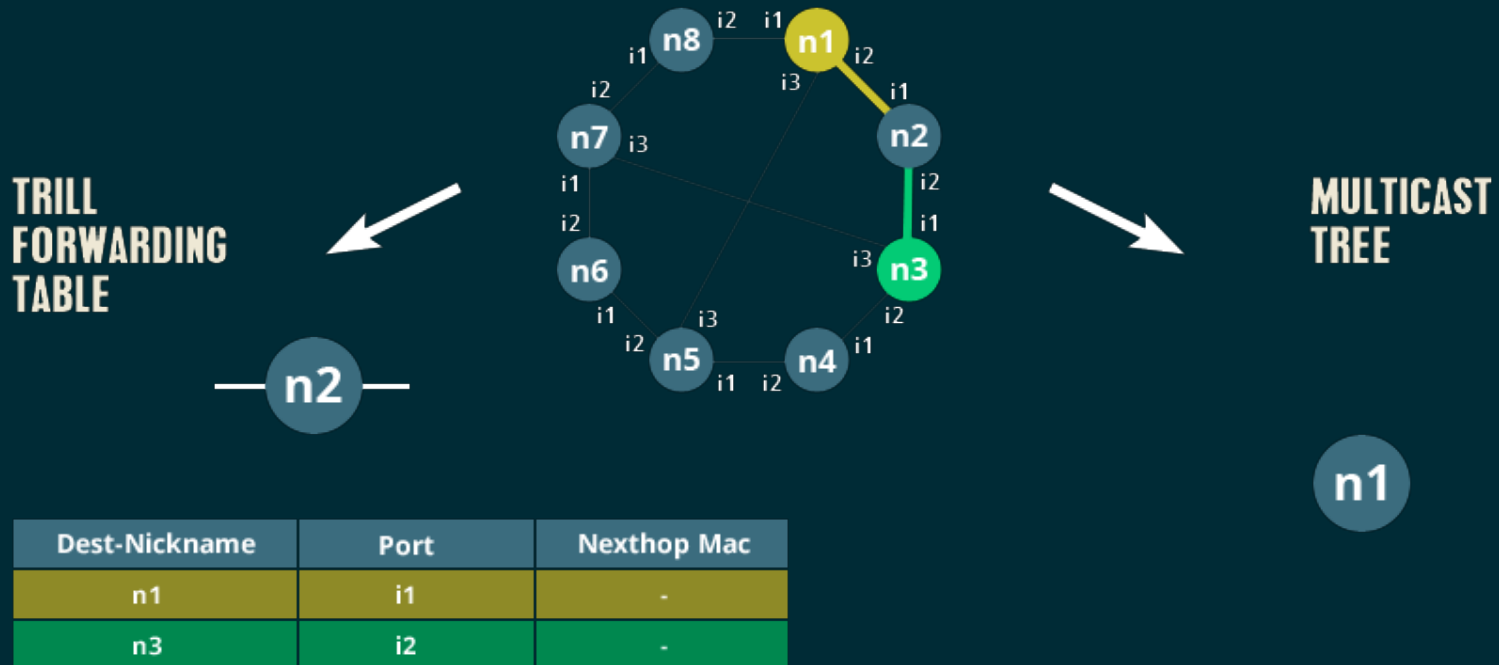
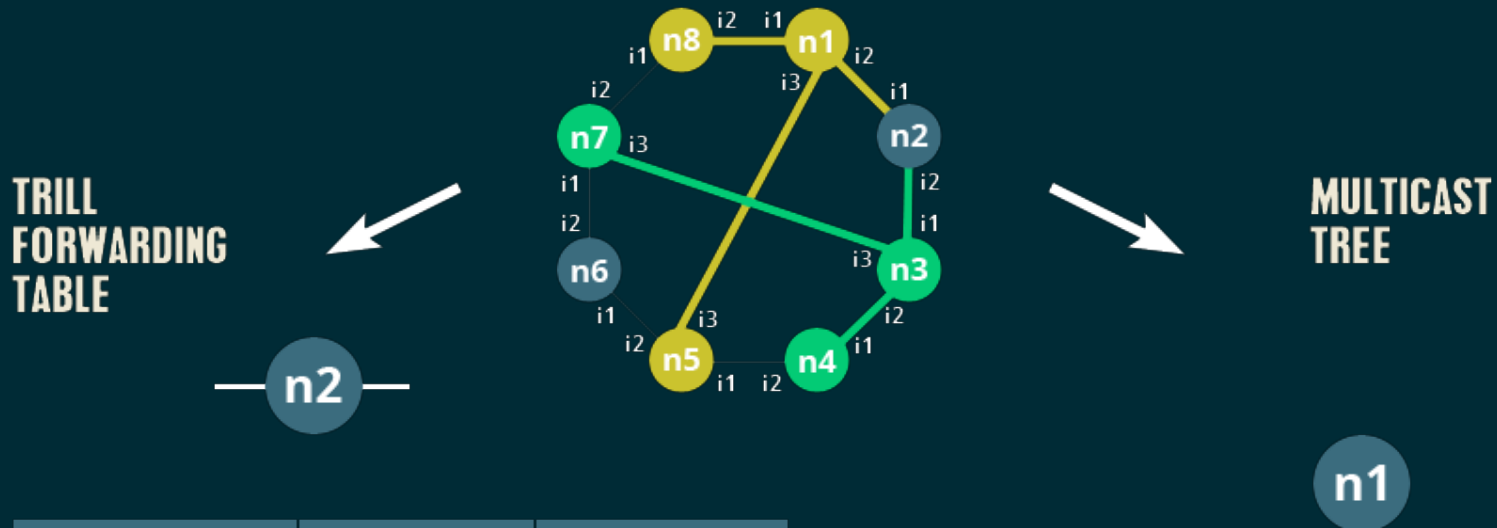# DESIGN AND IMPLEMENTATION

# CONTROL PLANE

## unicast building



**TRILL FORWARDING TABLE**

**MULTICAST TREE**

| Dest-Nickname | Port | Nexthop Mac |
| --- | --- | --- |

# CONTROL PLANE

## unicast building - first iteration



**TRILL FORWARDING TABLE**

**MULTICAST TREE**

| Dest-Nickname | Port | Nexthop Mac |
|---|---|---|
| n1 | i1 | - |
| n3 | i2 | - |

# CONTROL PLANE

unicast building - second iteration



**TRILL FORWARDING TABLE**

**MULTICAST TREE**

| Dest-Nickname | Port | Nexthop Mac |
|---|---|---|
| n1 | i1 | - |
| n3 | i2 | - |
| n8 | i1 | MAC - n1 |
| n5 | i1 | MAC - n1 |
| n4 | i2 | MAC - n3 |
| n7 | i2 | MAC - n3 |

# CONTROL PLANE

unicast building - third iteration



**TRILL FORWARDING TABLE**

**MULTICAST TREE**

| Dest-Nickname | Port | Nexthop Mac |
|---|---|---|
| n1 | i1 | - |
| n3 | i2 | - |
| n8 | i1 | MAC - n1 |
| n5 | i1 | MAC - n1 |
| n4 | i2 | MAC - n3 |
| n7 | i2 | MAC - n3 |
| n6 | i1 | MAC - n1 |
| n6 | i2 | MAC - n3 |

# CONTROL PLANE

## unicast building - final result

TRILL FORWARDING TABLE

MULTICAST TREE

| Dest-Nickname | Port | Nexthop Mac |
|---|---|---|
| n1 | i1 | - |
| n3 | i2 | - |
| n8 | i1 | MAC - n1 |
| n5 | i1 | MAC - n1 |
| n4 | i2 | MAC - n3 |
| n7 | i2 | MAC - n3 |
| n6 | i1 | MAC - n1 |
| n6 | i2 | MAC - n3 |

ECMP:
Equal Cost
Multipath

# CONTROL PLANE



TRILL
FORWARDING
TABLE

MULTICAST
TREE

| Dest-Nickname | Port | Nexthop Mac |
|---|---|---|
| n1 | i1 | - |
| n3 | i2 | - |
| n8 | i1 | MAC - n1 |
| n5 | i1 | MAC - n1 |
| n4 | i2 | MAC - n3 |
| n7 | i2 | MAC - n3 |
| n6 | i1 | MAC - n1 |
| n6 | i2 | MAC - n3 |

# CONTROL PLANE

## multicast building - first iteration

**TRILL FORWARDING TABLE**

**MULTICAST TREE**



| Dest-Nickname | Port | Nexthop Mac |
|---|---|---|
| n1 | i1 | - |
| n3 | i2 | - |
| n8 | i1 | MAC - n1 |
| n5 | i1 | MAC - n1 |
| n4 | i2 | MAC - n3 |
| n7 | i2 | MAC - n3 |
| n6 | i1 | MAC - n1 |
| n6 | i2 | MAC - n3 |

# CONTROL PLANE

multicast building - final iteration



## TRILL FORWARDING TABLE

| Dest-Nickname | Port | Nexthop Mac |
|---|---|---|
| n1 | i1 | - |
| n3 | i2 | - |
| n8 | i1 | MAC - n1 |
| n5 | i1 | MAC - n1 |
| n4 | i2 | MAC - n3 |
| n7 | i2 | MAC - n3 |
| n6 | i1 | MAC - n1 |
| n6 | i2 | MAC - n3 |

## MULTICAST TREE

# DATA PLANE

# DATA PLANE



VM 1  VM 2  VM 3  VM 4  VM 5  VM 6

VM1 sends the native frame

RB1  RB2

RB3  RB4

| | MAC - RB3 | MAC - RB2 | |
|---|---|---|---|
| | MAC - RB1 | MAC - RB3 | |
| | NickName RB2 | NickName RB2 | |
| | Nickname RB1 | Nickname RB1 | |
| MAC - VM6 | MAC - VM6 | MAC - VM6 | MAC - VM6 |
| MAC - VM1 | MAC - VM1 | MAC - VM1 | MAC - VM1 |
| IP - VM6 | IP - VM6 | IP - VM6 | IP - VM6 |
| IP - VM1 | IP - VM1 | IP - VM1 | IP - VM1 |
| DATA | DATA | DATA | DATA |

# DATA PLANE

VM 1 VM 2 VM 3    VM 4 VM 5 VM 6

RB1    RB2

RB1 adds the TRILL header
and the external mac header

RB3    RB4

| MAC - VM6 | MAC - RB3 | MAC - RB2 | MAC - VM6 |
|---|---|---|---|
| MAC - VM1 | MAC - RB1 | MAC - RB3 | MAC - VM1 |
| IP - VM6 | NickName RB2 | NickName RB2 | IP - VM6 |
| IP - VM1 | Nickname RB1 | Nickname RB1 | IP - VM1 |
| DATA | MAC - VM6 | MAC - VM6 | DATA |
| | MAC - VM1 | MAC - VM1 | |
| | IP - VM6 | IP - VM6 | |
| | IP - VM1 | IP - VM1 | |
| | DATA | DATA | |

# DATA PLANE

# IMPLEMENTATION - SENDING

# IMPLEMENTATION - RECEIVING

# LAYER 2 - SWITCHING LIMITATION

- A large number of tenants implies
  - a huge number of MAC address in switch table
  - ARP storm at nodes
- STP to ensure a loop free topology
  - blocking redundant paths
  - Core-computes required, recomputes when topology changes
- Number of VLANs is limited to 4096
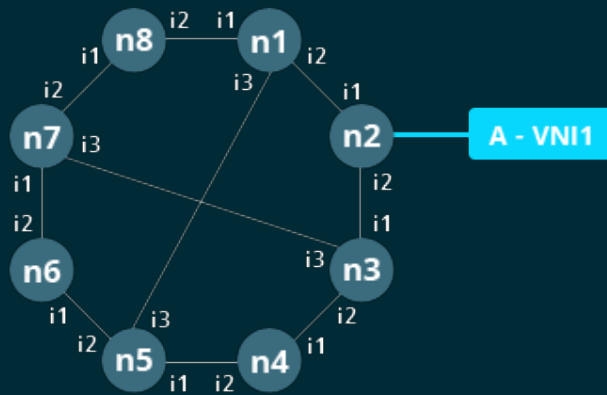
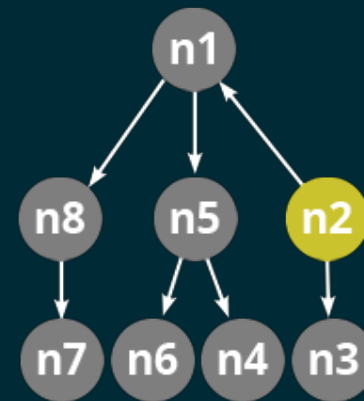# TRILL + VNI = VNT

Virtual Network over TRILL

# VNT FRAME FORMAT

# VNI LIFE



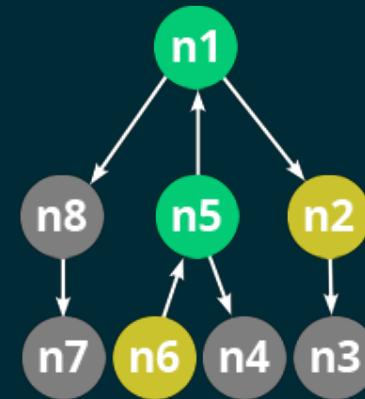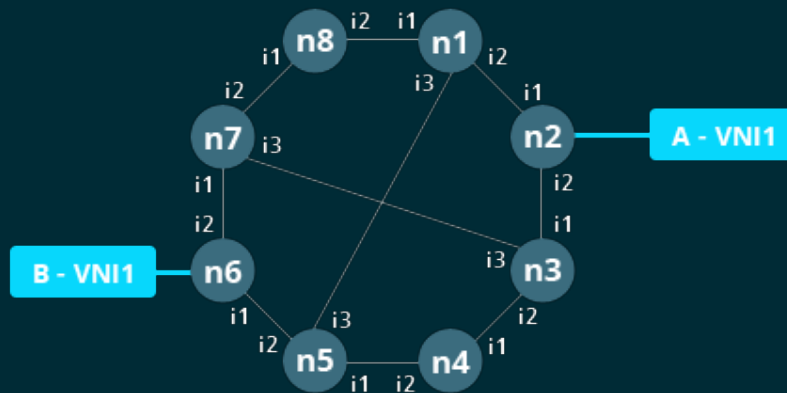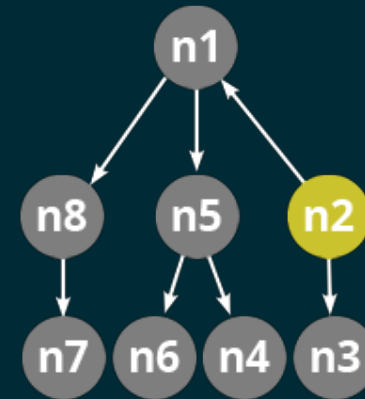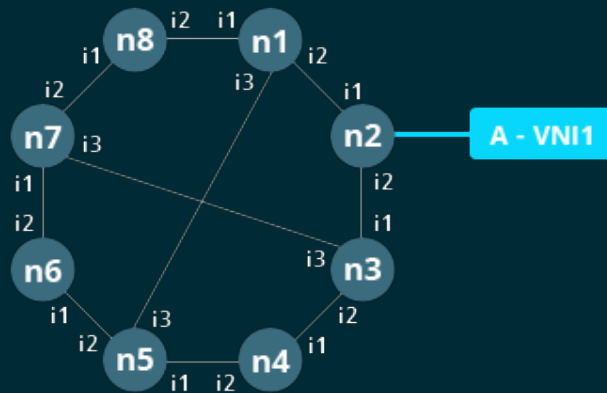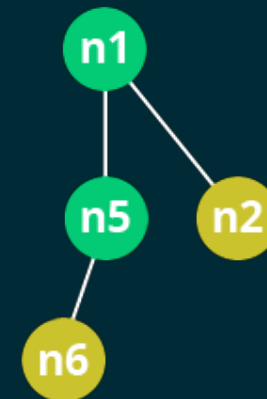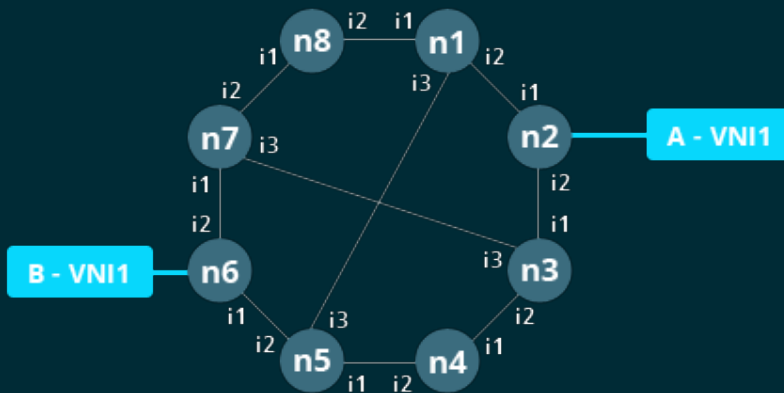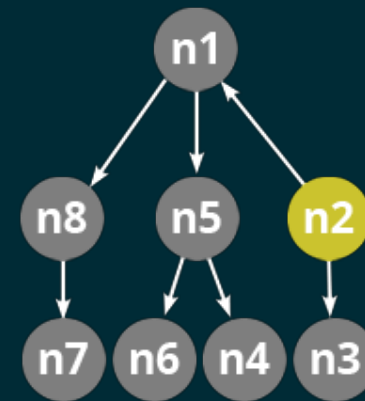| | |
|---|---|
| →(green) | VNI added to locally supported VNI |
| →(orange) | VNI deleted from locally supported VNI |
| →(red) | VNI received for the first time on interface i |
| →(cyan) | VNI revoked from neighbor |

# VNI TOPOLOGY BUILDING

## TOPOLOGY



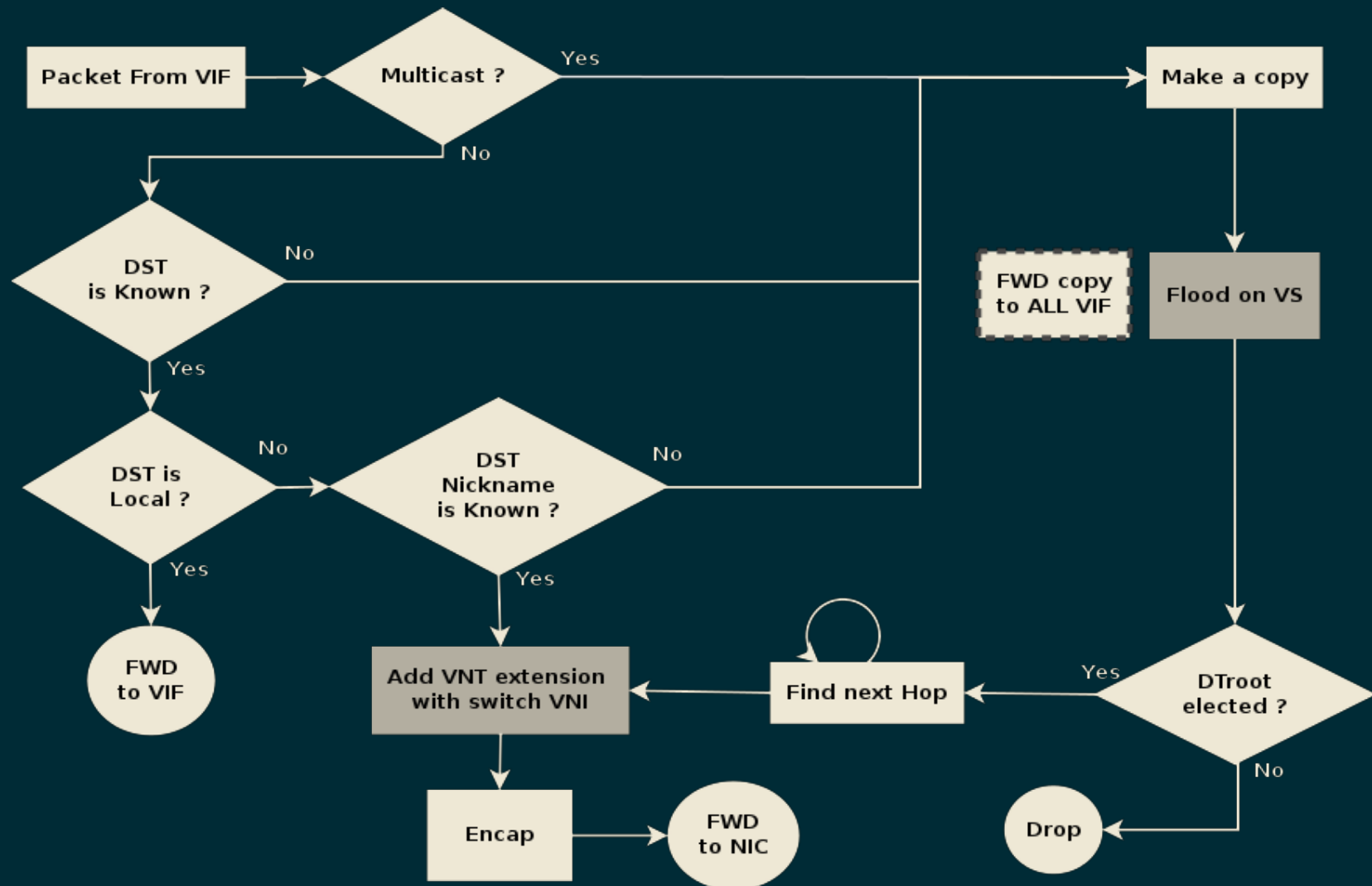## MULTICAST TREE

# VNI TOPOLOGY BUILDING
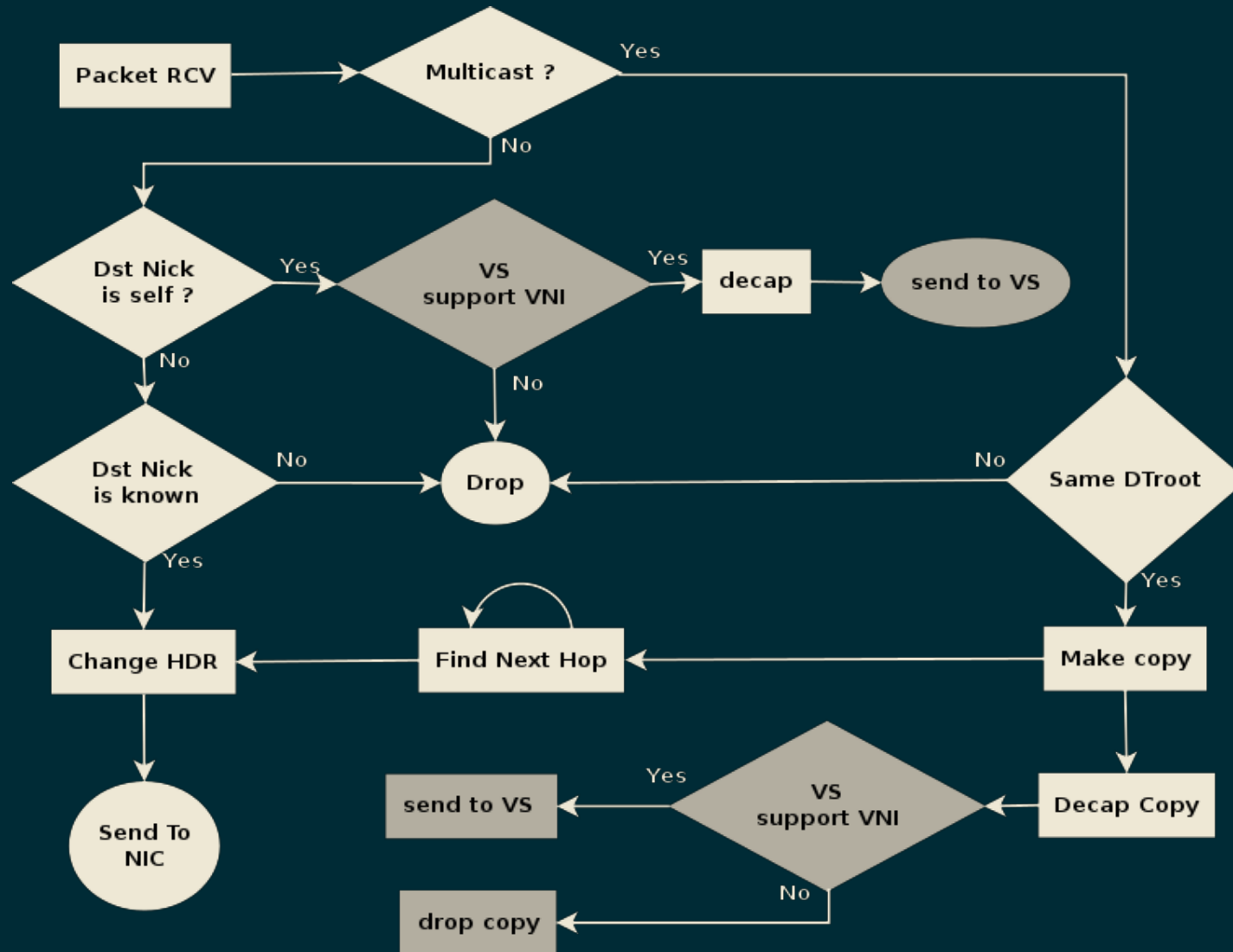
# VNI TOPOLOGY BUILDING
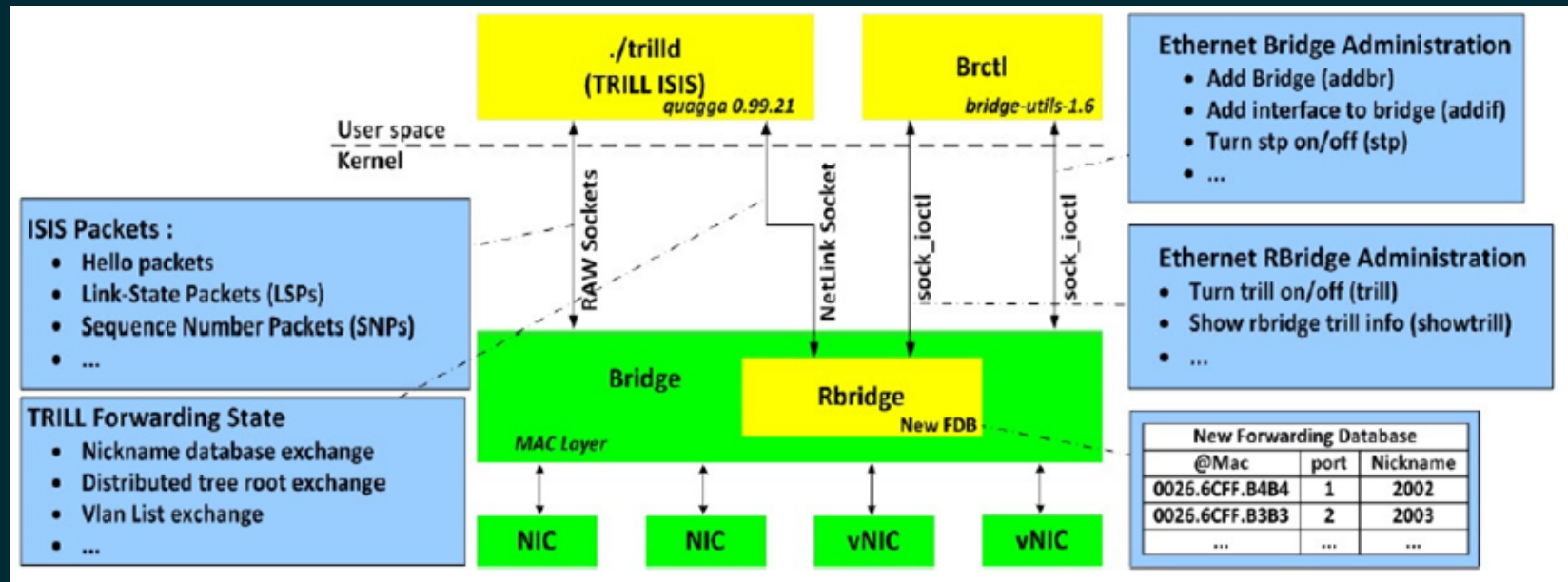
## TOPOLOGY



## MULTICAST TREE

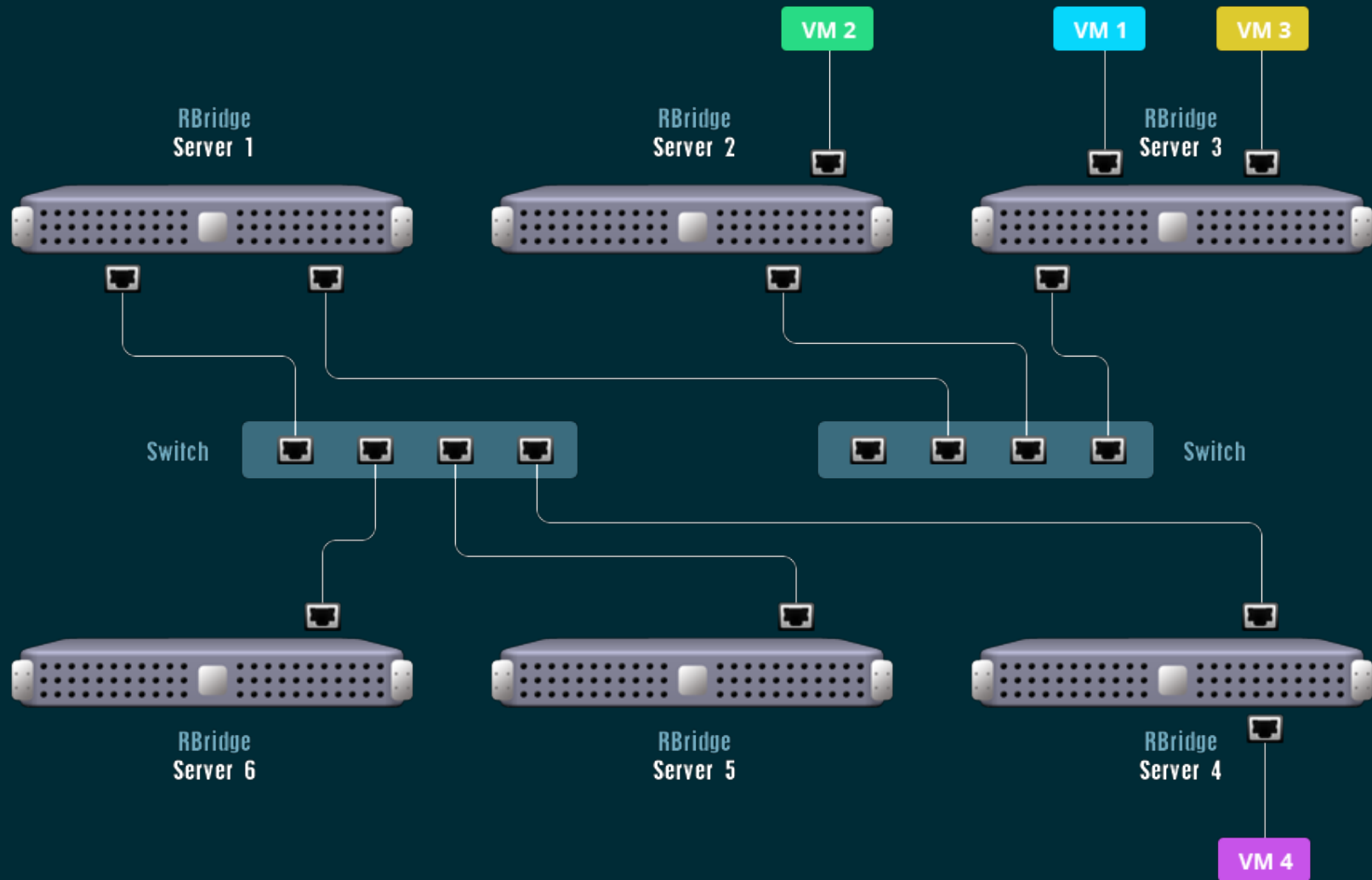# IMPLEMENTATION WITH VNI - SENDING

# IMPLEMENTATION WITH VNI - RECEIVING

# LINUX BIG PICTURE

DEMONSTRATION

# SCREENCAST

**screencast**

(live explanation to understand what's going on)

# PH.D. STUDY

Ahmed Amamou - **ahmed@gandi.net**

"Network isolation for Virtualized Datacenters"

University Pierre & Marie Curie - GANDI SAS

project still in development and cleaning

TRILL sources: **github.com/Gandi/ktrill**

VNT: still two research projects working on it - drafts

# GANDI.NET

Gandi Hosting - **gandi.net/hosting**

William Dauchy - **william@gandi.net**

slides **pres.gandi.net/kr2013**