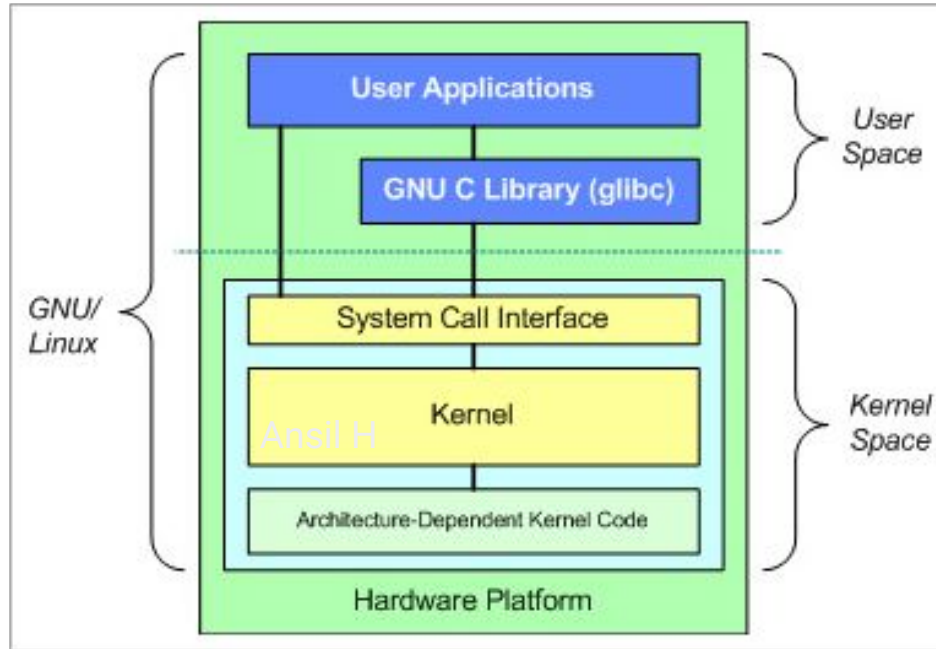


# Container Networking Fundamentals

Ansil H  
Lead SRE@Armorblox



# Linux OS fundamentals



# Why we need containers?

- Process isolation for security
- Resource usage restriction
- Dependency management
- Lifecycle management



# The Building Blocks

- Namespaces
- Cgroups
- Capabilities



# Namespaces

Namespace	Constant	Isolates
Cgroup	CLONE_NEWCGROUP	Cgroup root directory
IPC	CLONE_NEWIPC	System V IPC, POSIX message queues
Network	CLONE_NEWNET	Network devices, stacks, ports, etc.
Mount	CLONE_NEWNS	Mount points
PID	CLONE_NEWPID	Process IDs
User	CLONE_NEWUSER	User and group IDs
UTS	CLONE_NEWUTS	Hostname and NIS domain name



# CGroups or Control Groups

- Resource Limiting
- Prioritization
- Accounting
- Control/Freeze



# Linux Capabilities

Root user can do anything but, can we have an “admin” user with restricted access to some of the privileged operations?

“Split privileged kernel calls and group them into related functionality”

Eg:- The “ping” binary may not work without **CAP\_NET\_RAW** capability.



# Demo - Container from scratch

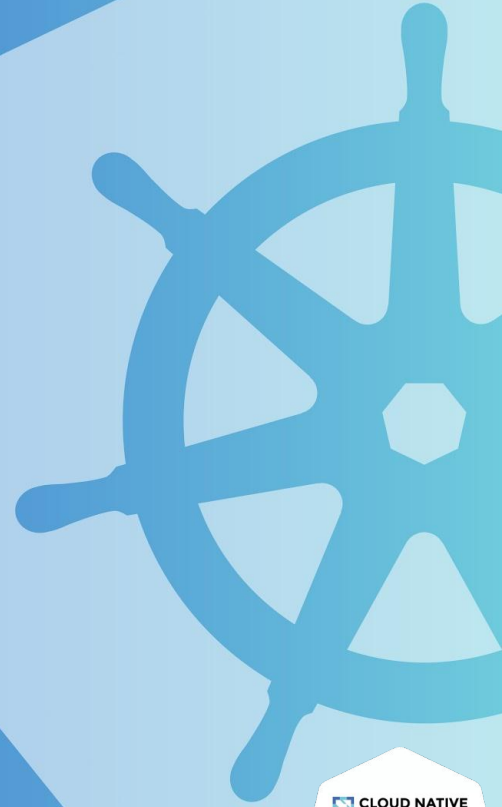
Let's create a container from scratch using the “unshare” command.

In this demo, we will use “busybox” to simulate the needed binaries for the container.





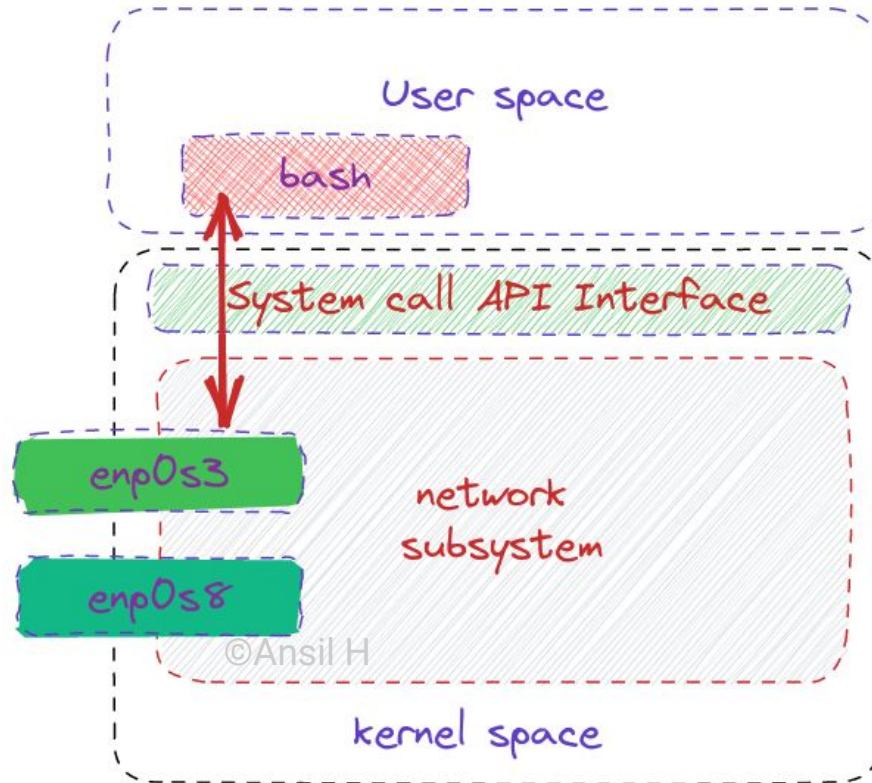
# Container networking



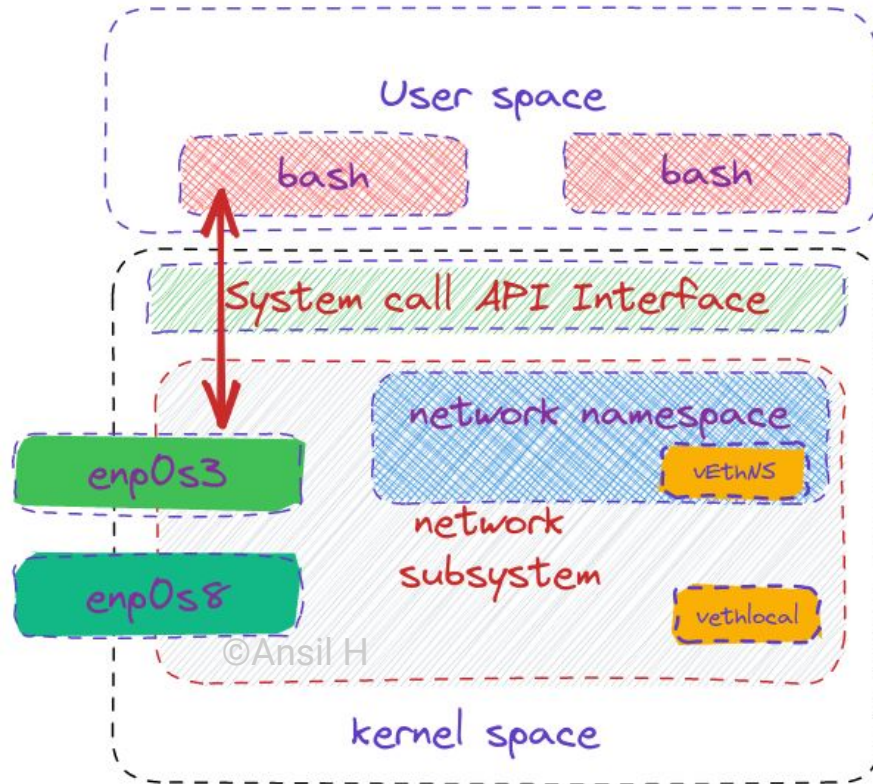
# Bridging



# Networking - Host view



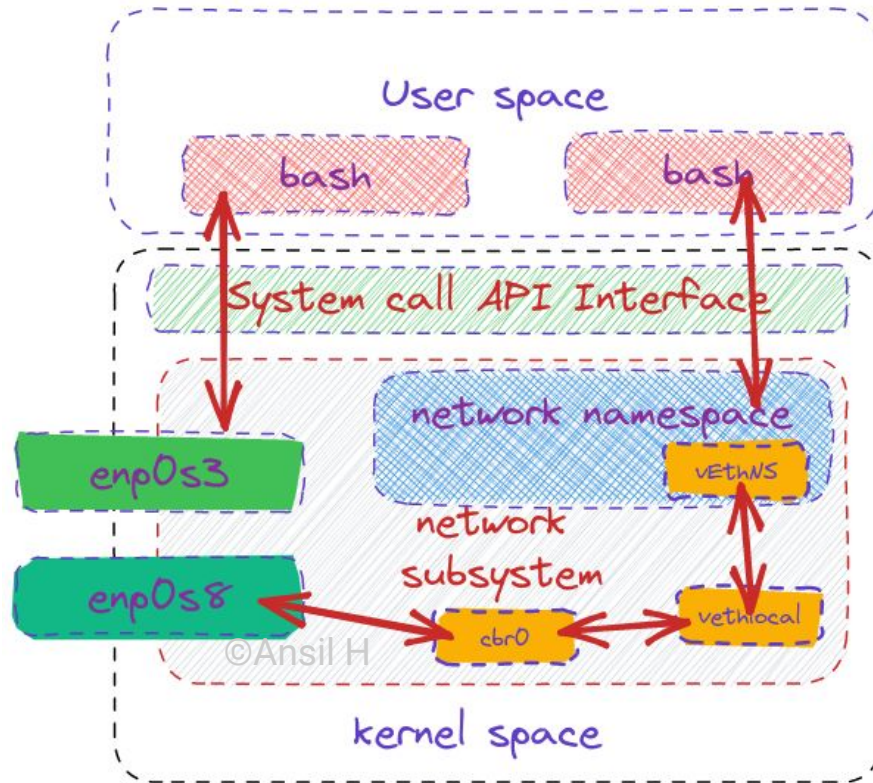
# Network namespace and vETH pair



# vEth Pair

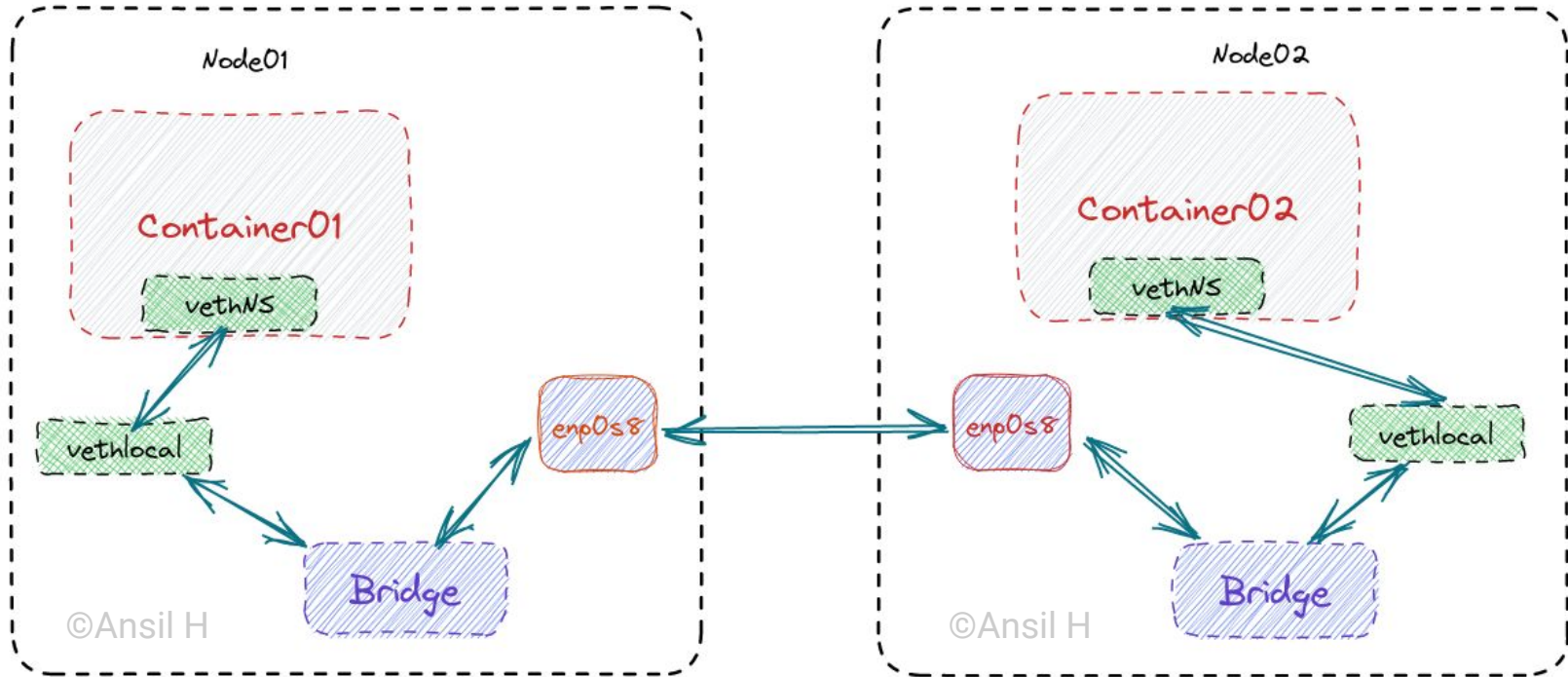


# Bridging





# Bridging - Demo



# Bridging - Demo

In this demo, we will use two hosts called node01 and node02

A container will be started on both nodes using “busybox” and “unshare” commands.

A vEth pair will be created on both nodes

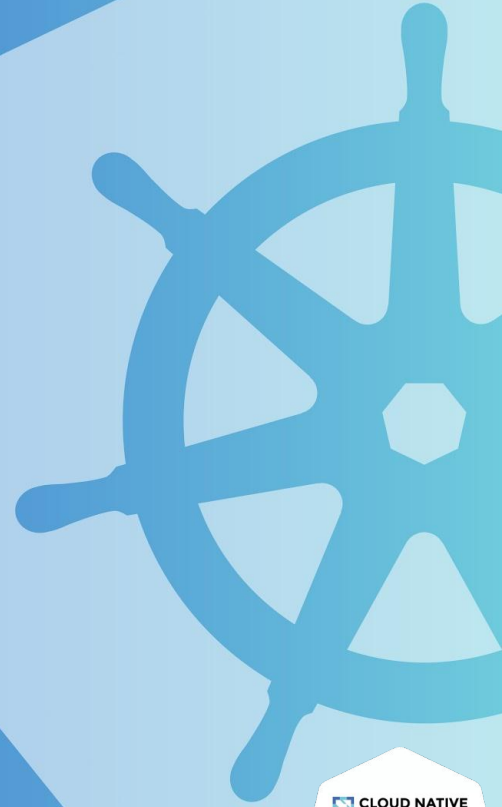
A virtual bridge device will be created on both nodes and will be connected to one of the ethernet device.

Finally one end of the vEth will be connected to the container and other end will be connected to the bridge.

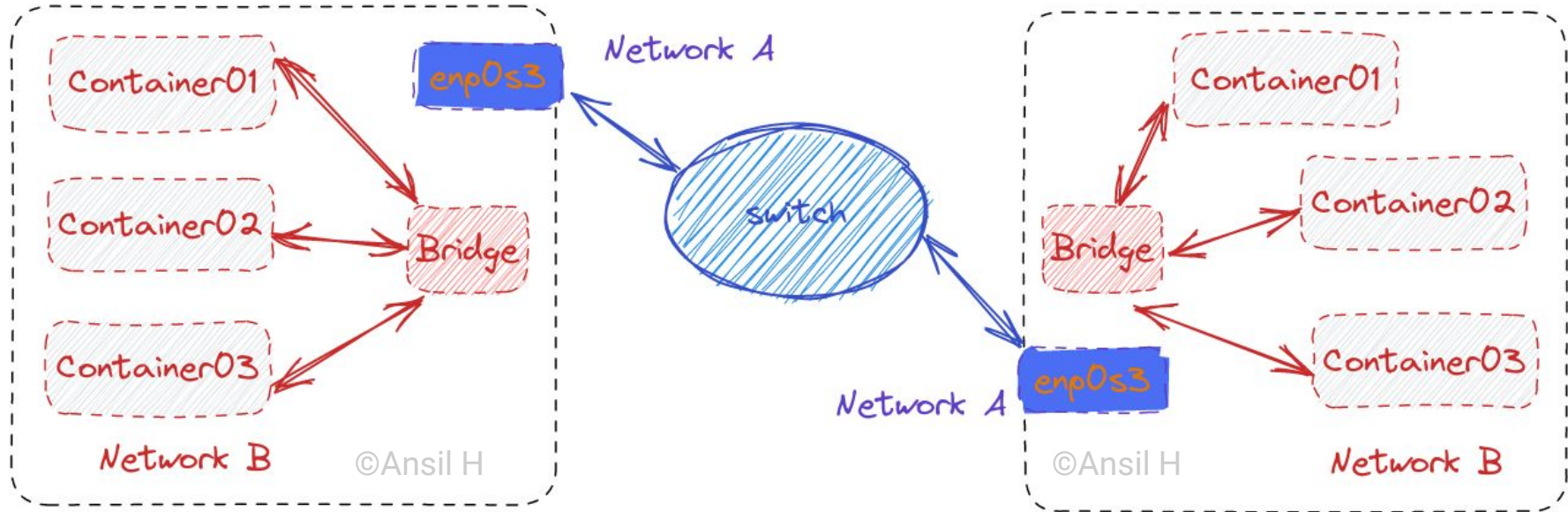




# VXLAN



# Virtual Extensible LAN - VXLAN



# VXLAN

- Packet encapsulation
- TCP packet inside UDP
- Maximum Transfer Unit



# MTU 1450

The default MTU is 1500 in network interfaces.

If we keep the default value inside the container, the MTU may exceed the limit after encapsulation

14 Outer Ethernet header (Mac header)

20 byte Outer IP header

8 byte UDP header

8 byte VXLAN header

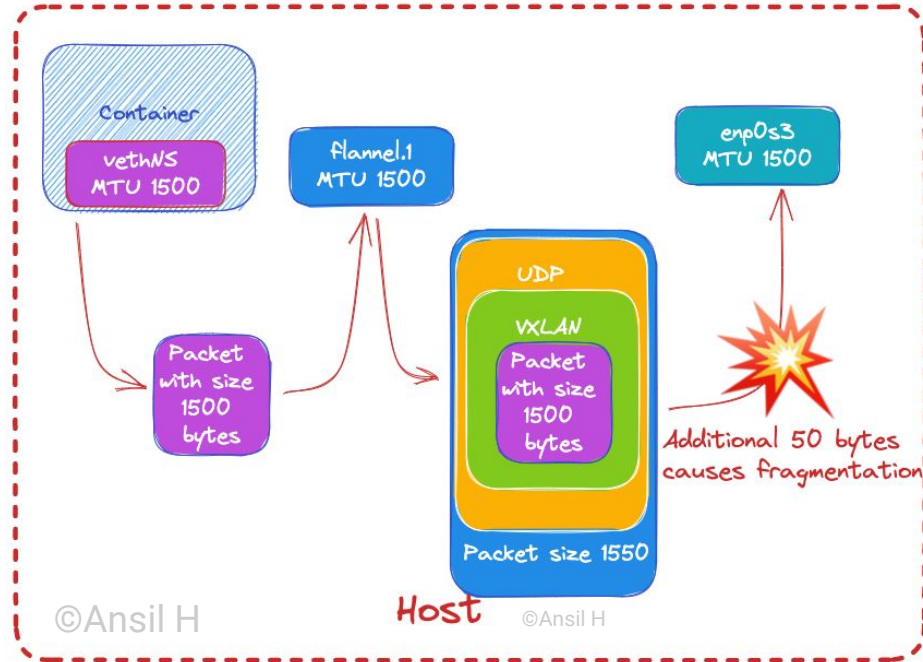
1500 byte payload which includes the original IP header/s.

**Total = 1550 bytes which exceeds the default MTU.**

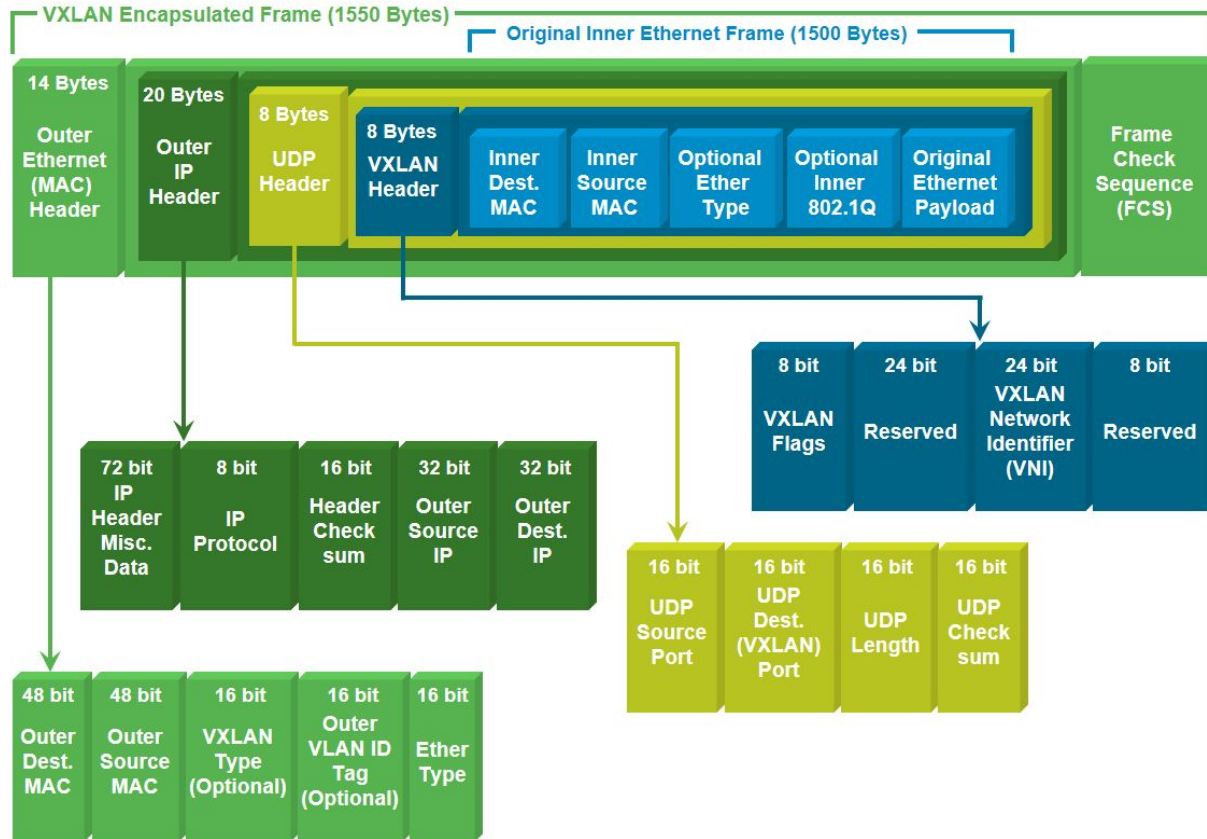
We need to set the MTU to 1450 for the Virtual Eth devices so that the MTU will not exceed after encapsulation,



# VXLAN - MTU



# VXLAN Packet



# Flannel

A network fabric for containers

- Assign subnet to each nodes
- Create flannel VXLAN interface
- Create and sync bridge forwarding database



# VXLAN - Demo

- Create two containers
- Add a veth pair
- Create bridge
- Add veth pair ends to bridge
- Setup etcd
- Start flannel
- Configure IPs and setup packet forwarding





Q&A

