



# How We Built Magic Transit

# Agenda

- Who are we
- Who is Cloudflare
- What is Magic Transit
- Designing The Product
- How Magic Transit Works
- Questions

# Who Are We?

Erich Heine

Systems Engineer

*erich@cloudflare.com*

Connor Jones

Systems Engineer

*conjones@cloudflare.com*



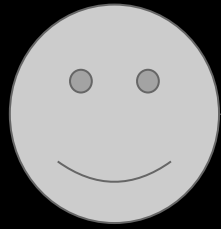
{erich, conjones}@cloudflare.com

# Who is Cloudflare?

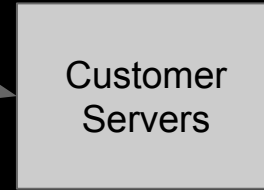
“Helping build a  
better internet”

# But How???

- GLOBAL Anycast network
- Provide security & performance services:
  - DDoS
  - DNS
  - Spectrum
  - CDN, WAF, Workers
  - ...



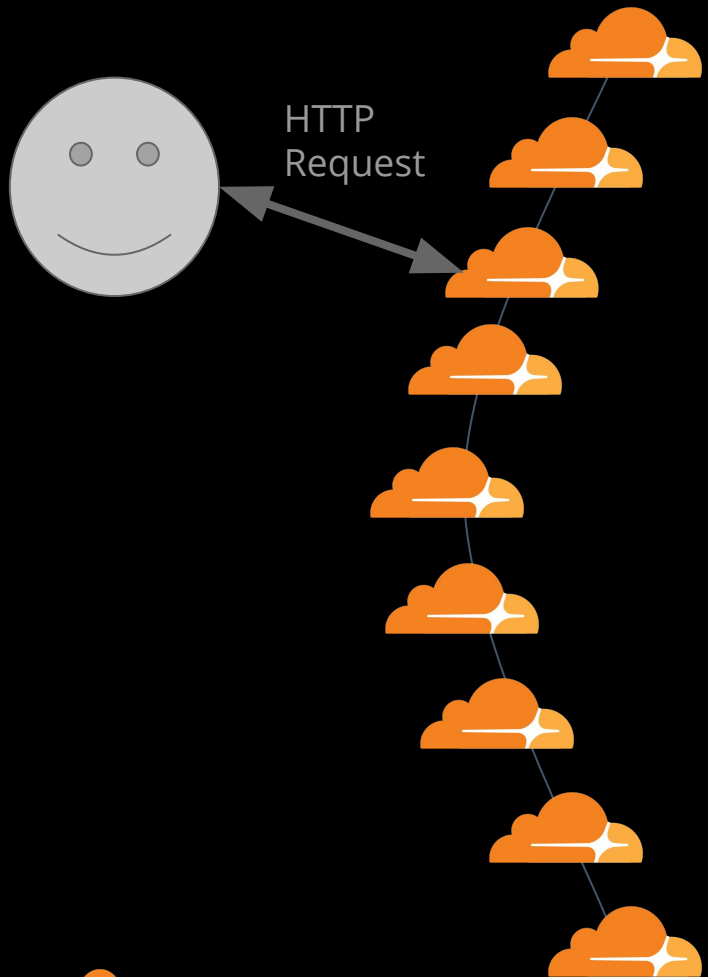
HTTP  
Request



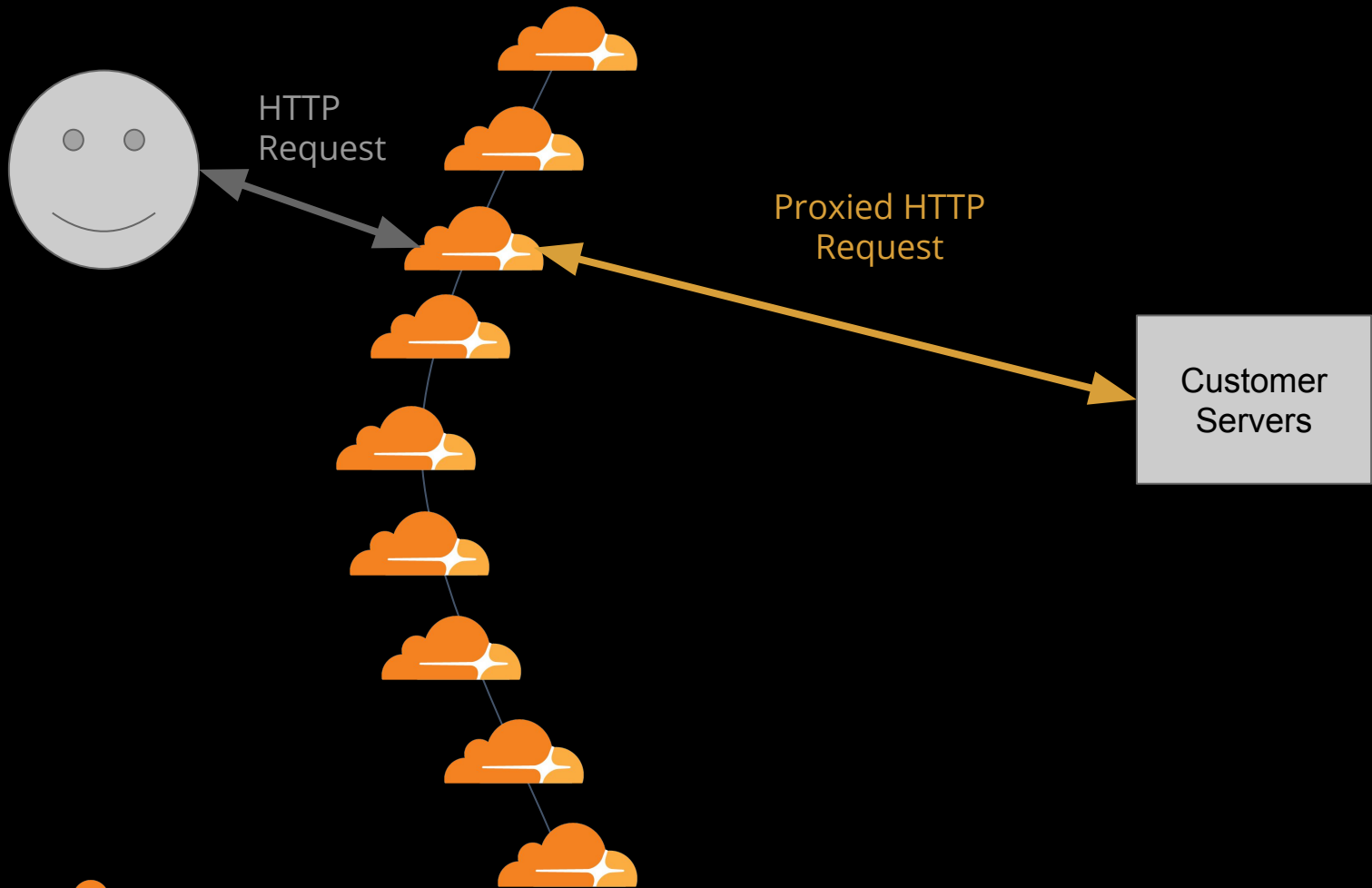


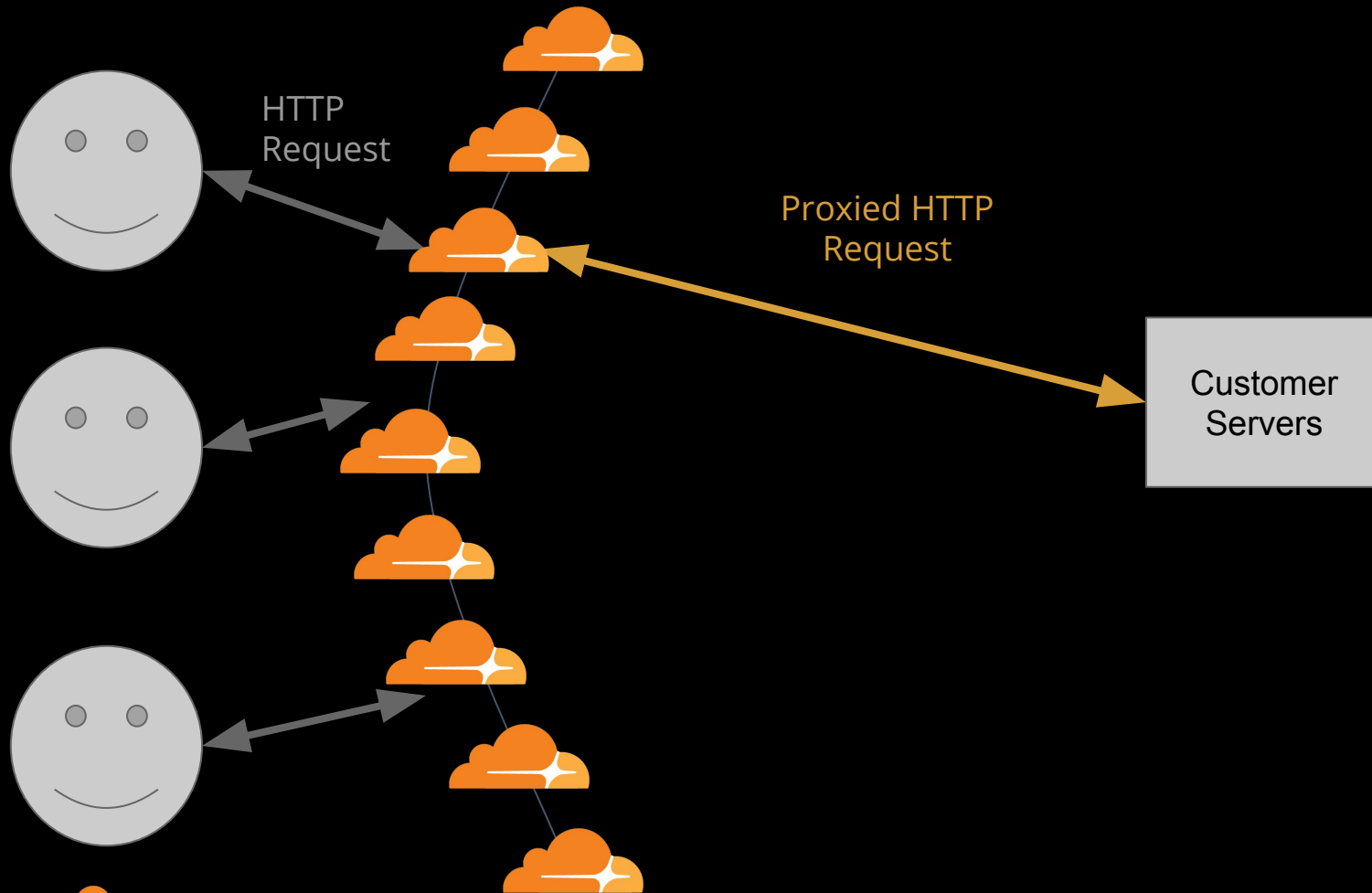


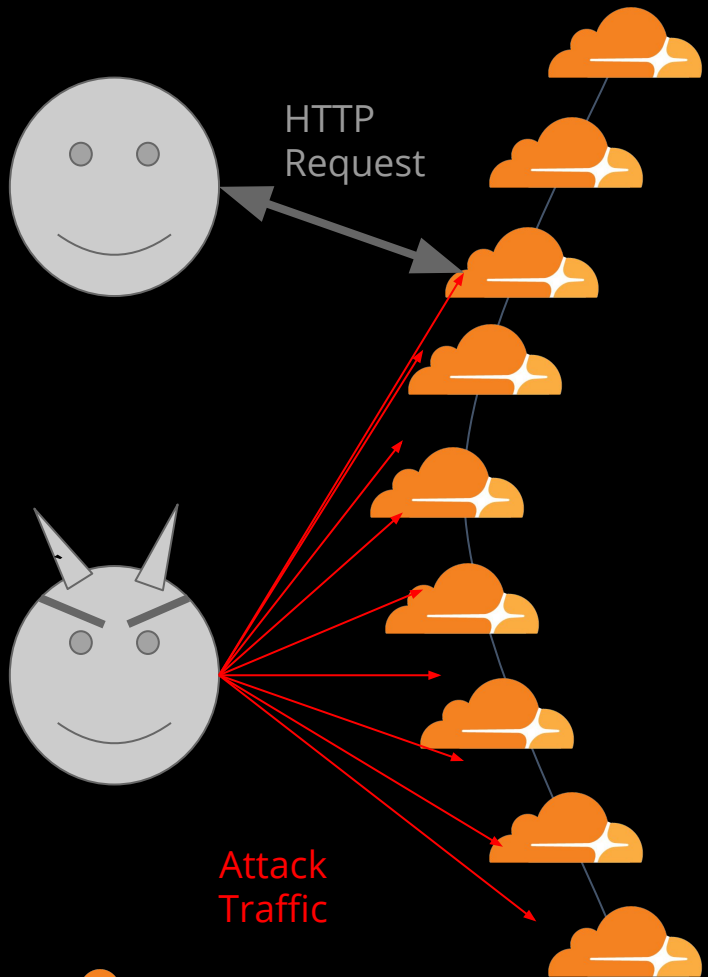
Customer  
Servers



Customer  
Servers







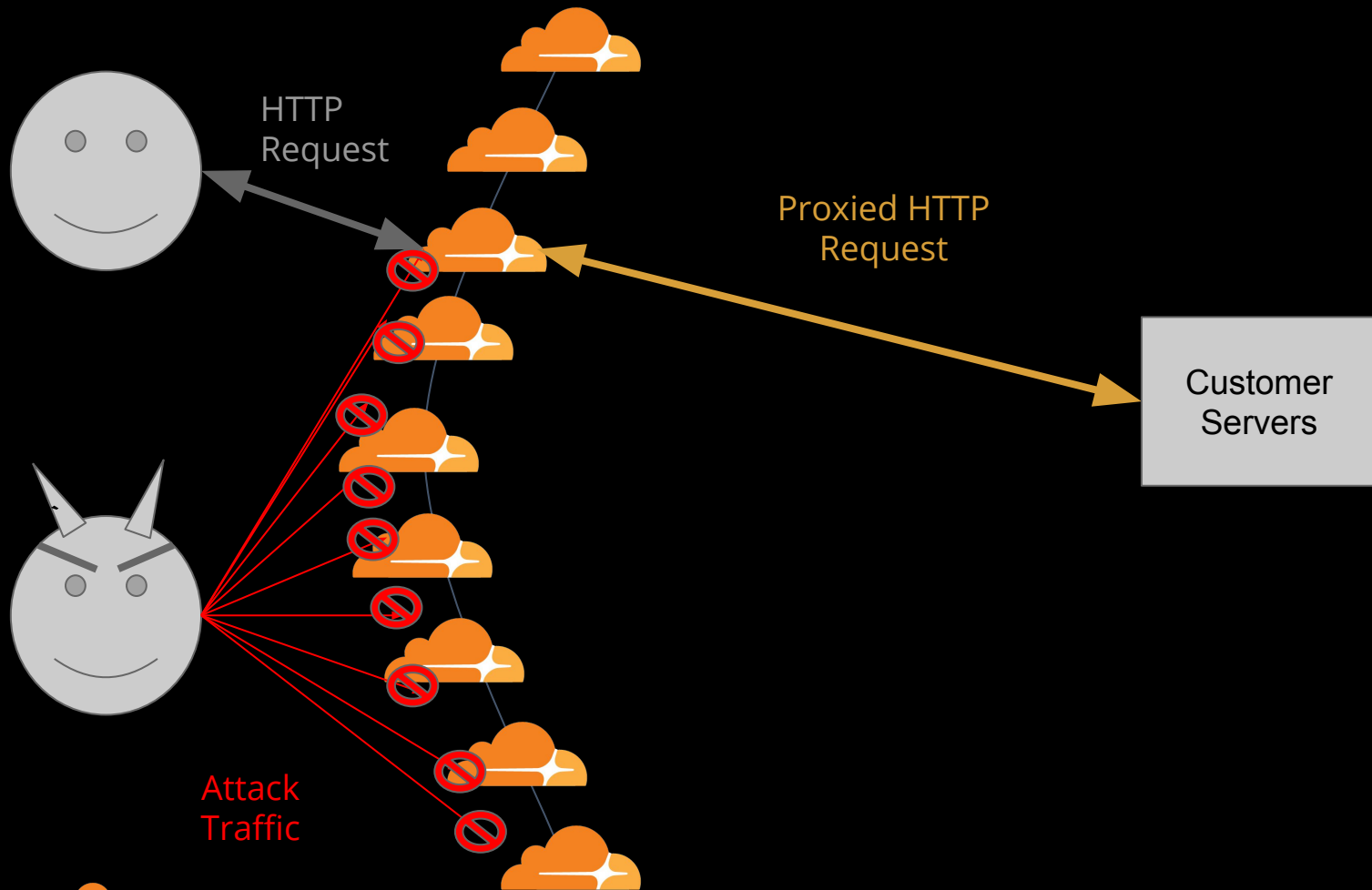
Proxied HTTP Request

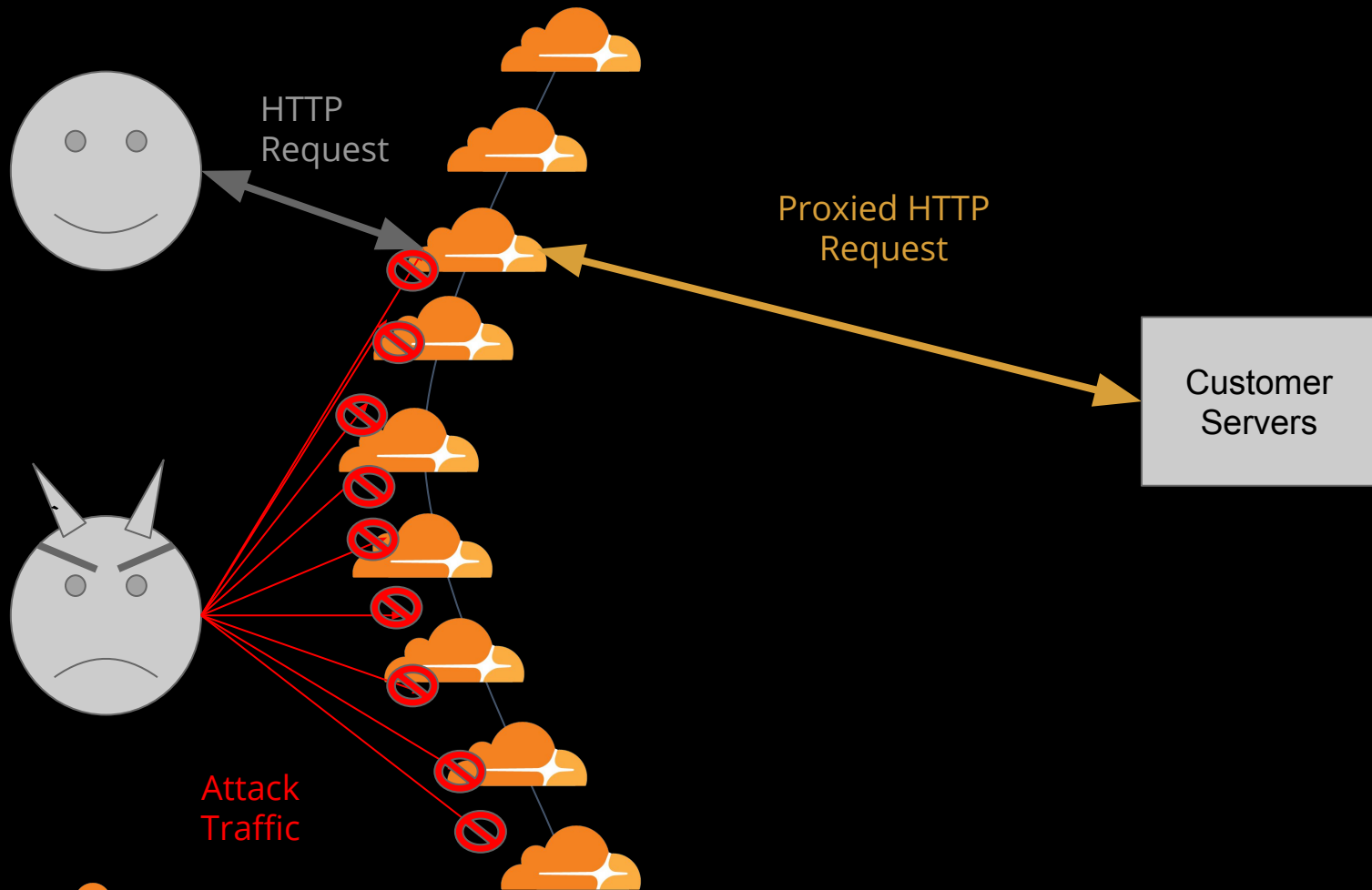
Customer Servers

Attack Traffic



{erich, conjones}@cloudflare.com





# Customers to Cloudflare:

Can you do this with all my traffic?



## But why?

- Consolidation of vendors
- Network functions in the network (not the boxes)
- Proxying is not always the answer

## Cloudflare to customers:

- Much of this operates at L3 already
- We have a big network
- We'll build you a Magic Transit

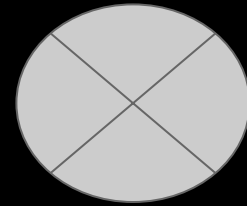
# What is Magic Transit?

# Magic Transit

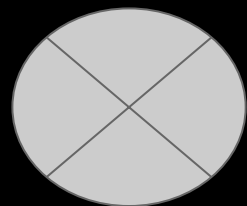
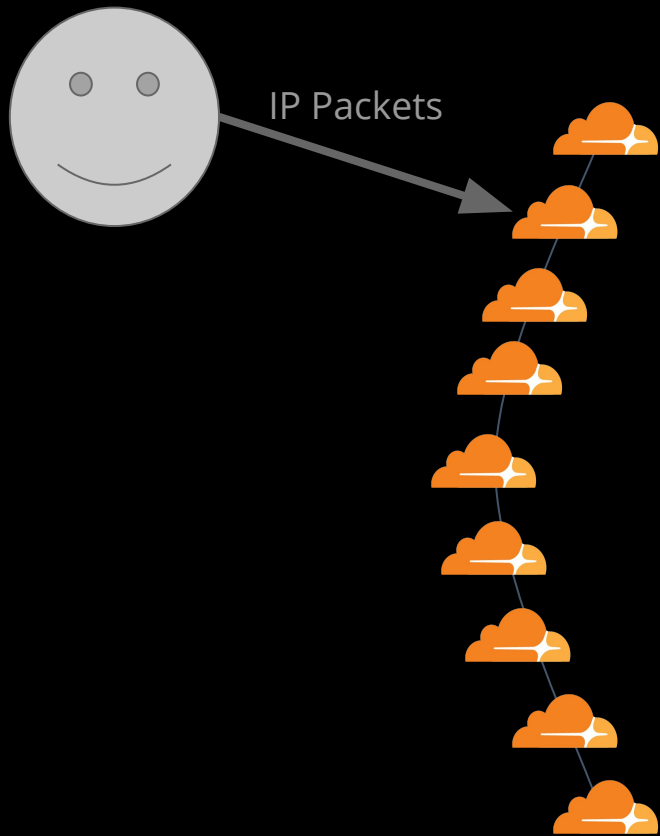
- Network functions as a service
- L3 DDoS scrubbing
- Firewall

By:

- Existing DDoS solution
- Advertises customer prefixes on BGP

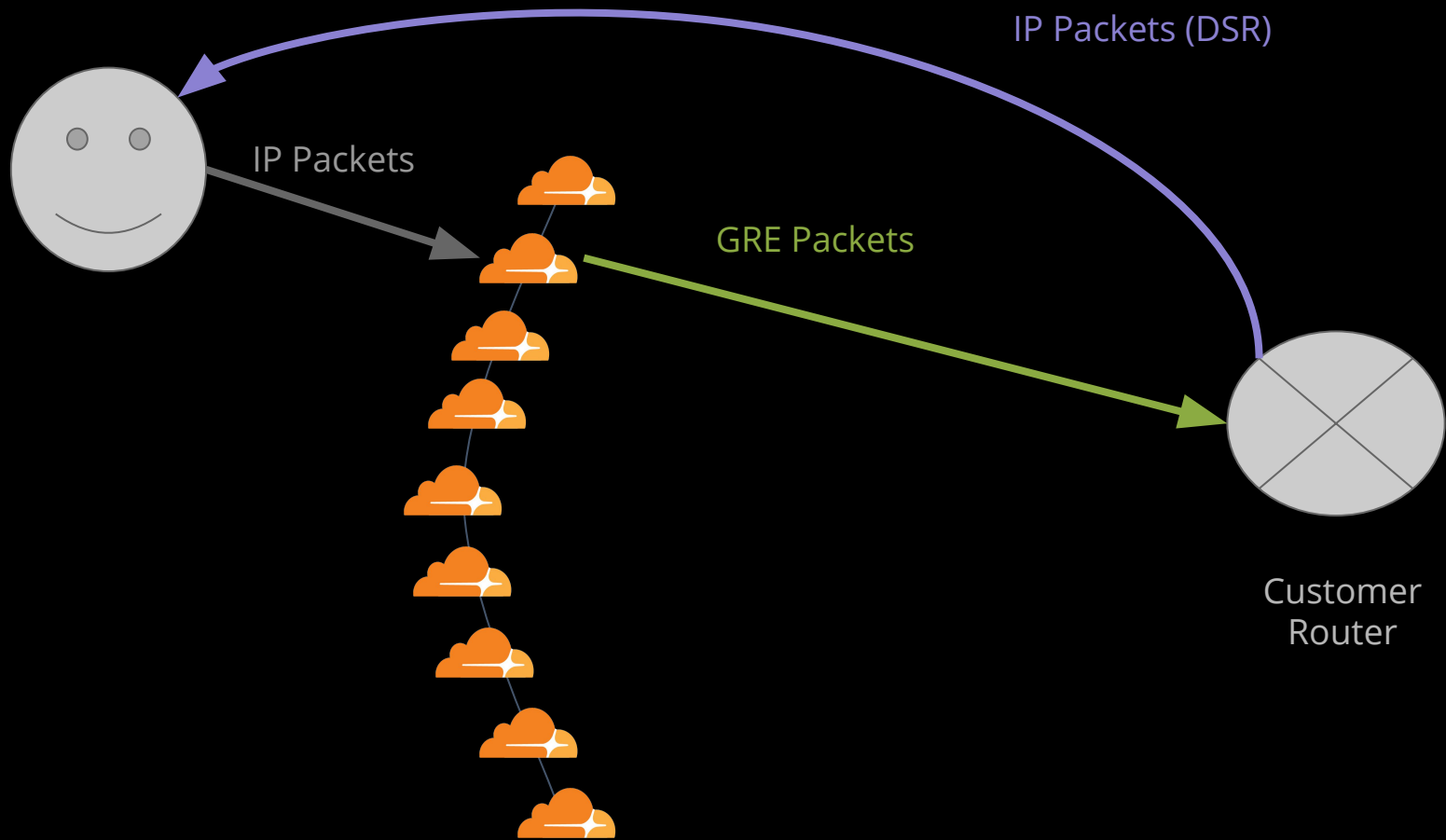


Customer  
Router

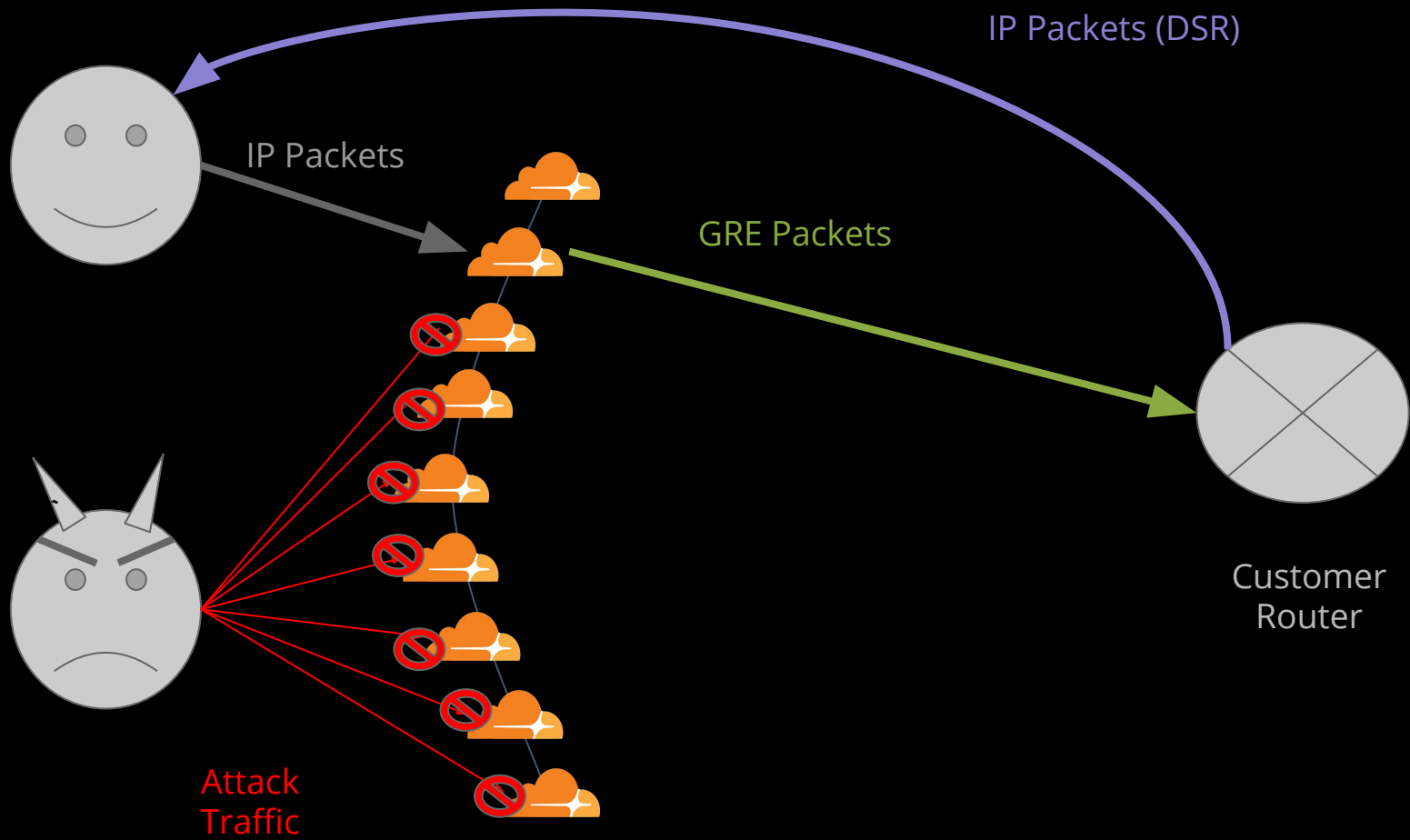


Customer  
Router









# Designing the product

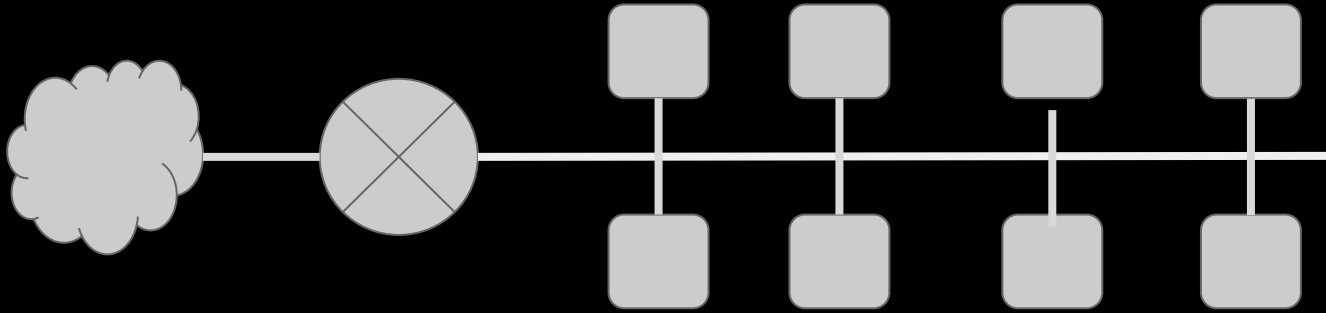
Proof of concept

{erich, conjones}@cloudflare.com

# Will it even work?

- Does GRE actually work that way?
- Does it scale?
- How do we fit it in our platform?

# A Cloudflare POP



# Challenges

- The metal is a router
- Uniform servers
- Customer isolation
- Security
- Timeline + skillset

# Technology options

- VRFs
- XDP
- Namespaces

# VRF

## Pros:

- Routing isolation
- Stays at L3

## Cons:

- Lacking expertise
- Unclear iptables interaction

# XDP

## Pros:

- Bypass a lot of processing
- Lots of expertise
- “Just Code”

## Cons:

- It's a lot of code
- Time investment



# Namespaces

## Pros:

- Isolation
- Expertise
- Bash!

## Cons:

- Lots of veths
- Complex setup

And the winner is...

# Namespaces

- Proof of concept in a few weeks
- Surprisingly small shell script
- Dataplane performance is pretty good

# Namespaces

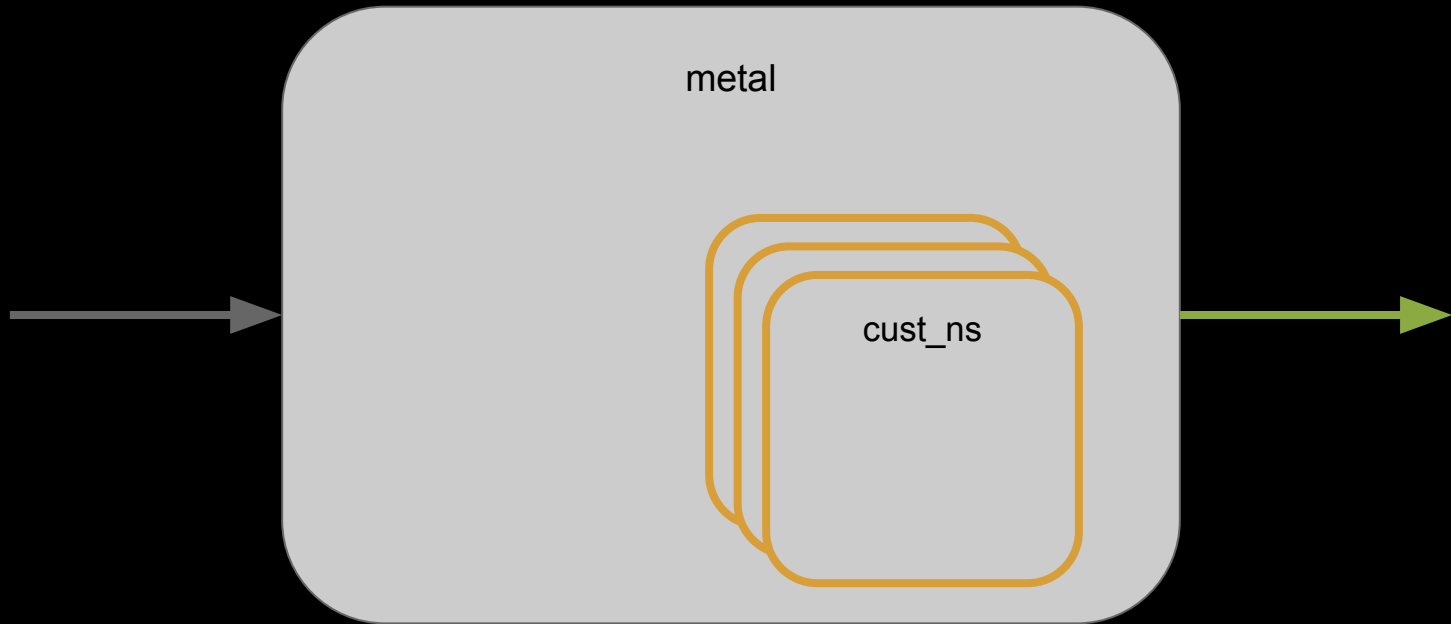
- Proof of concept in a few weeks
- Surprisingly small shell script
- **Dataplane performance is pretty good**

We can use this for the product!

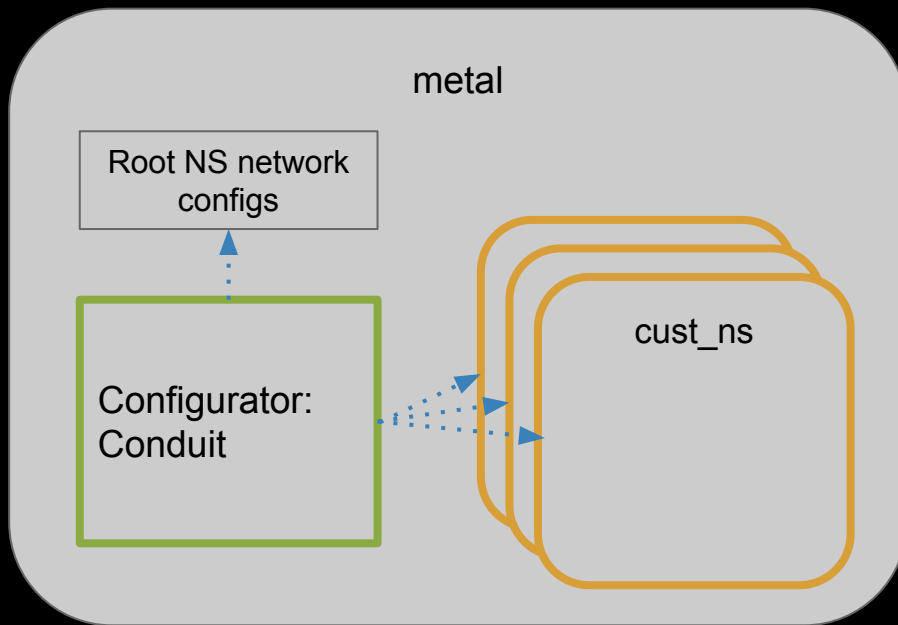
# Designing the product

The actual product

# One namespace per customer



# Configuration daemon



# Configurator: Conduit

- Daemon written in Go
- Uses QS as config distribution
- Can crash without (big) consequence



# Behind The Curtain Of Magic Transit

Shuffling packets from eyeballs to origins

{erich, conjones}@cloudflare.com

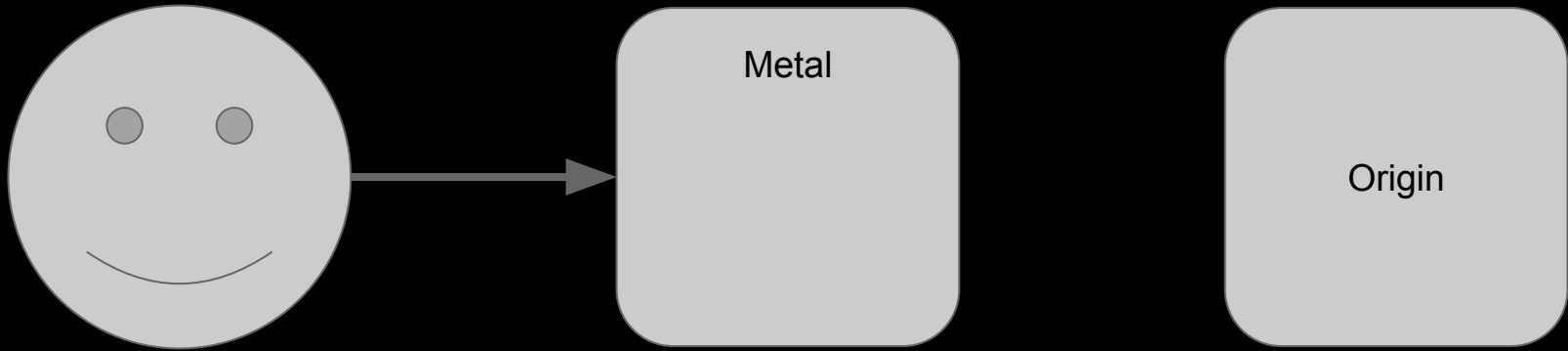
# How Packets Get To The Origin



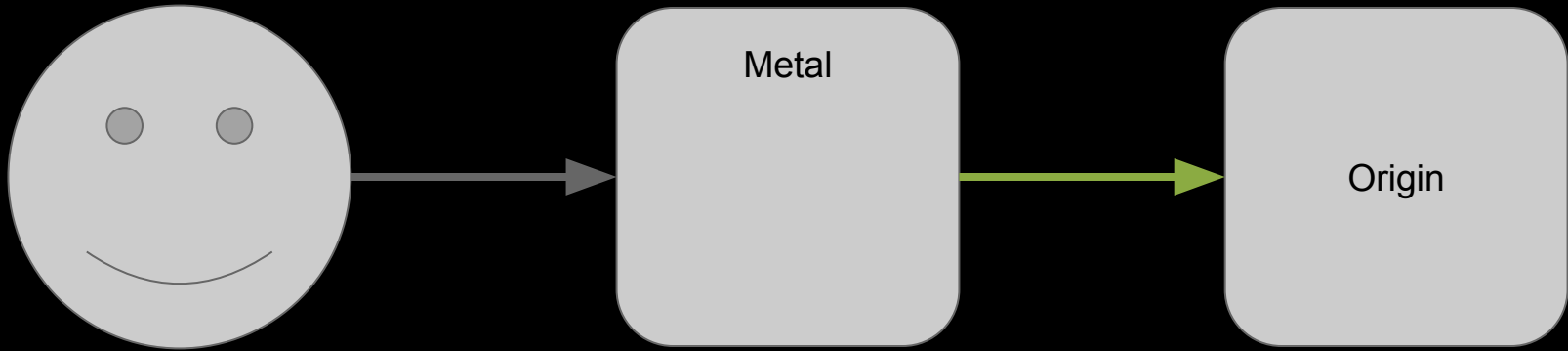
Metal

Origin

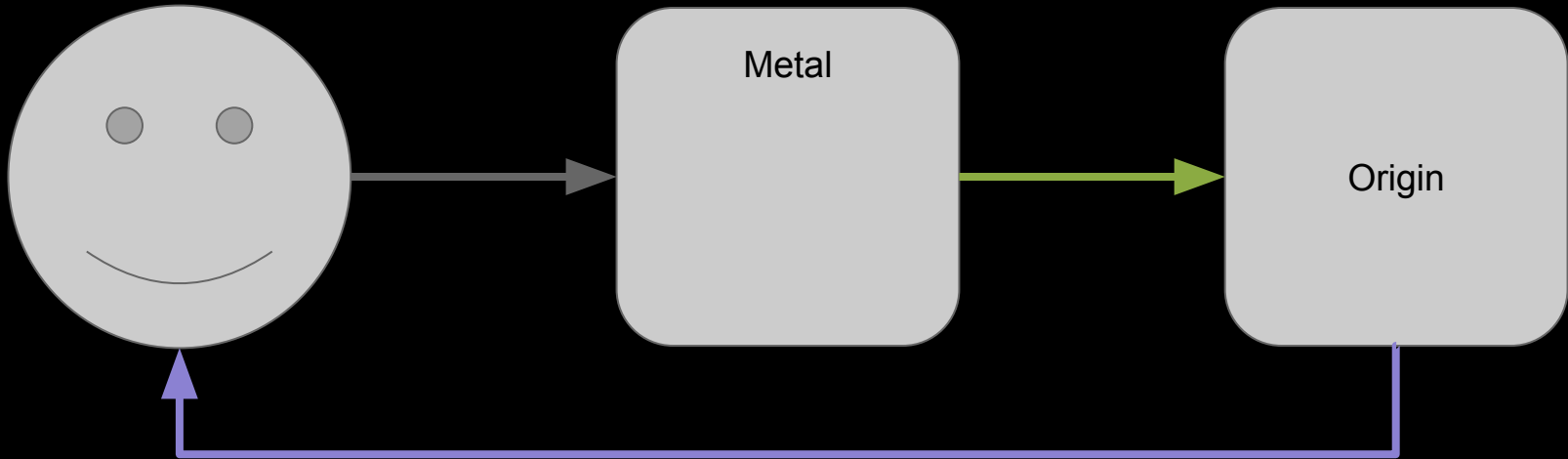
# Packets Arrive At The Edge



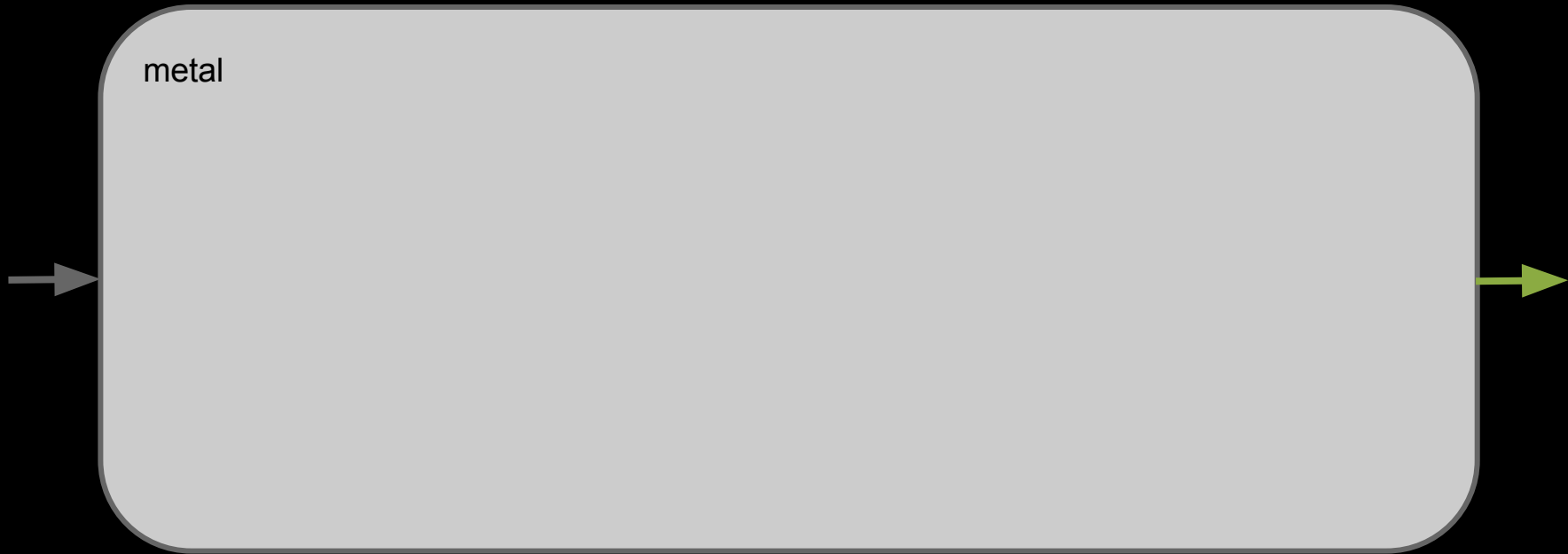
# Magic Happens



# Origin Responds to Eyeball



# An Edge Metal



# One Namespace Per Customer

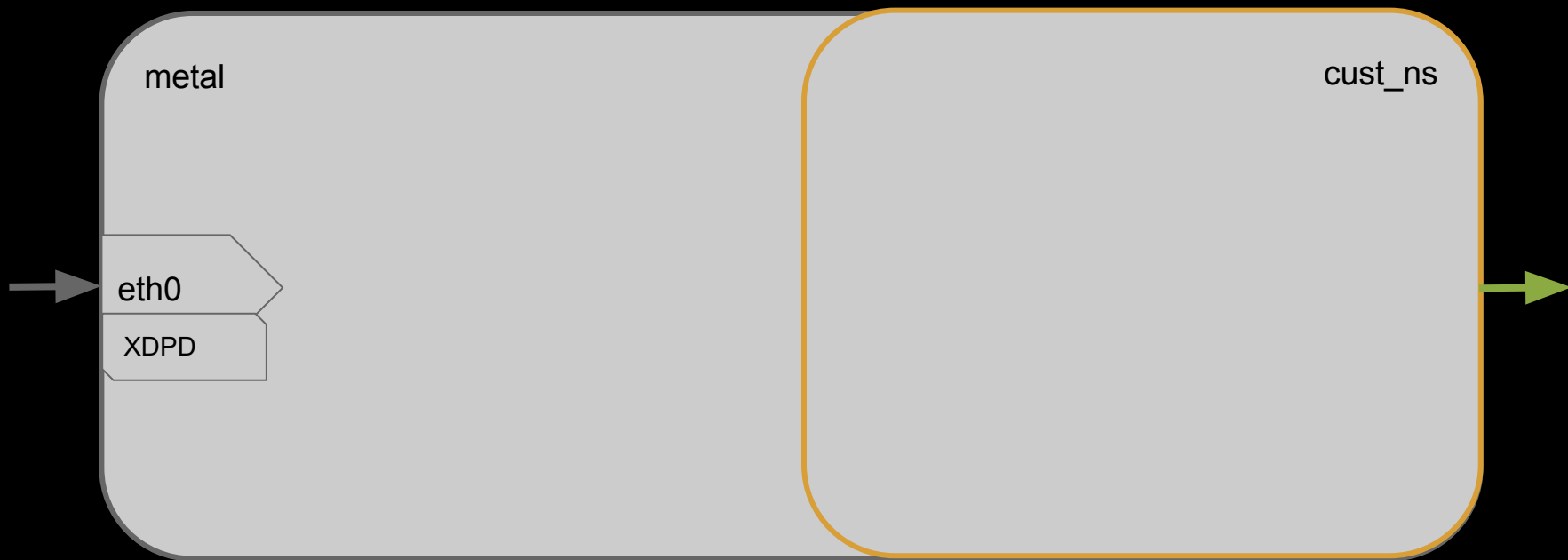


# Packets Arrive

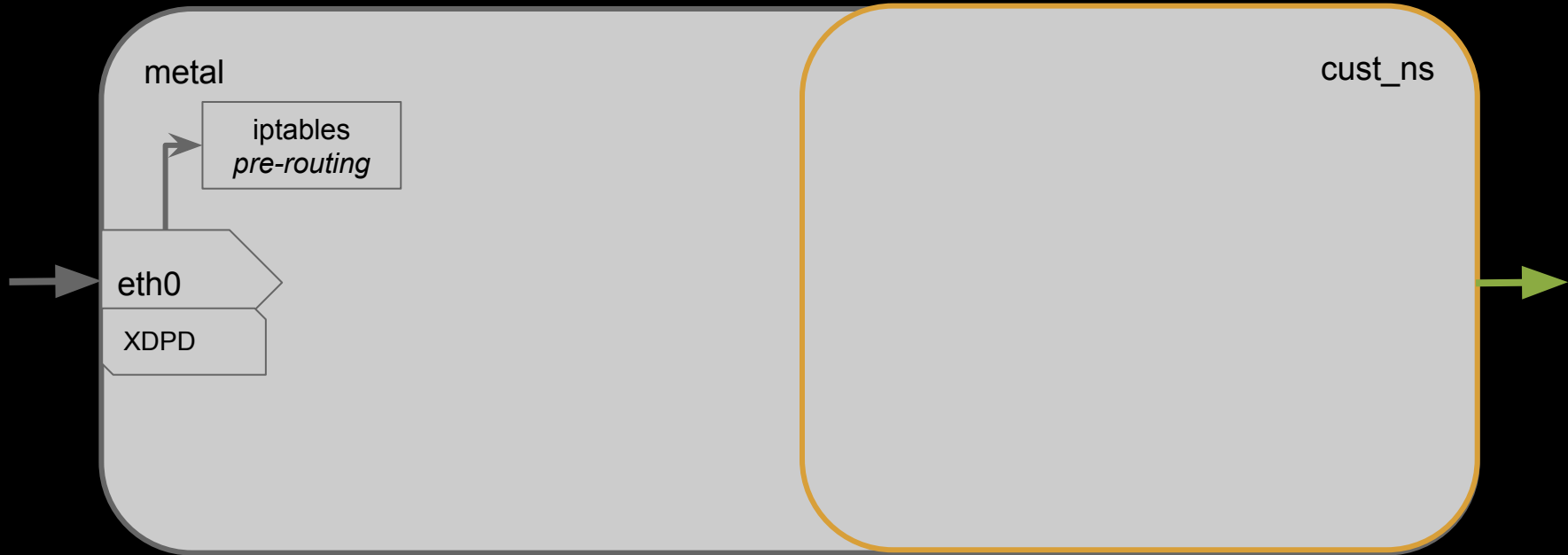




# A Lil XDP



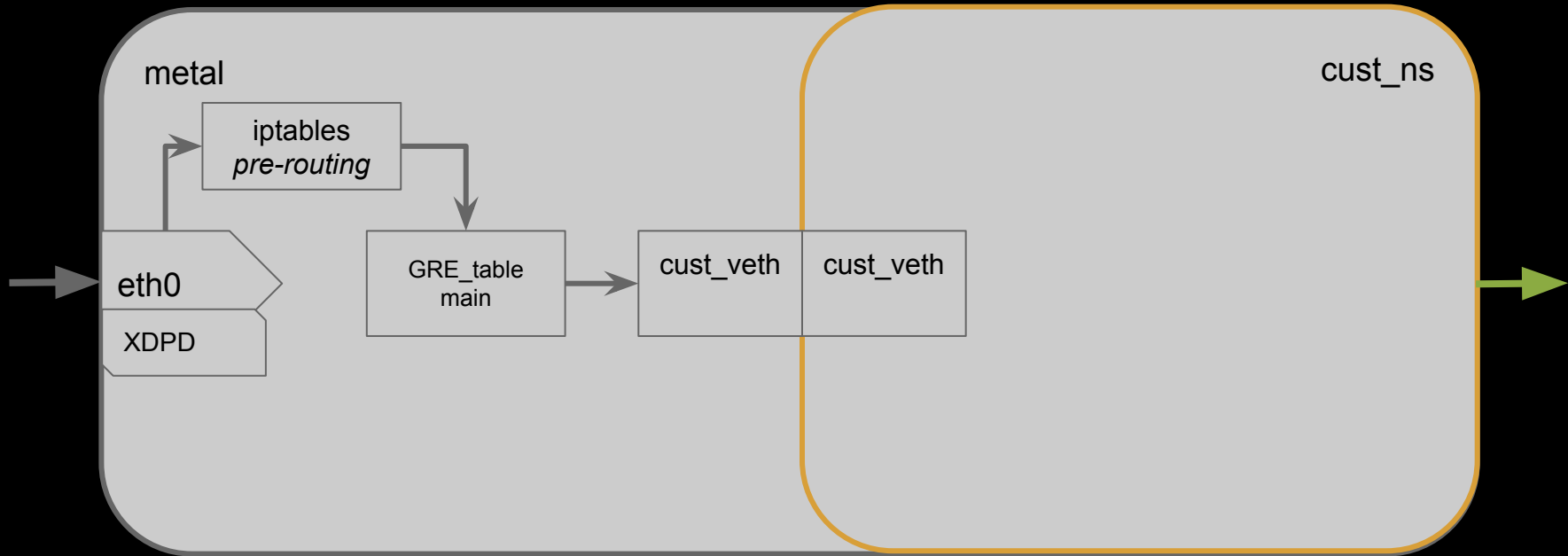
# Marking Packets



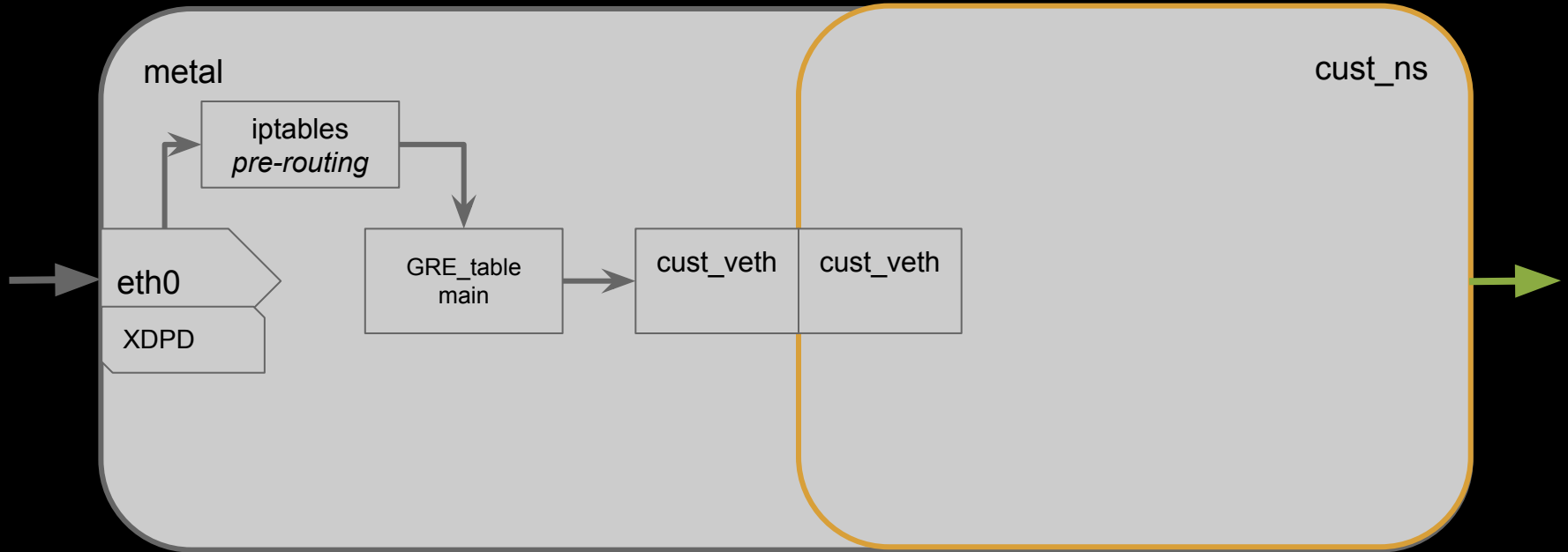
# Route Tables with an S



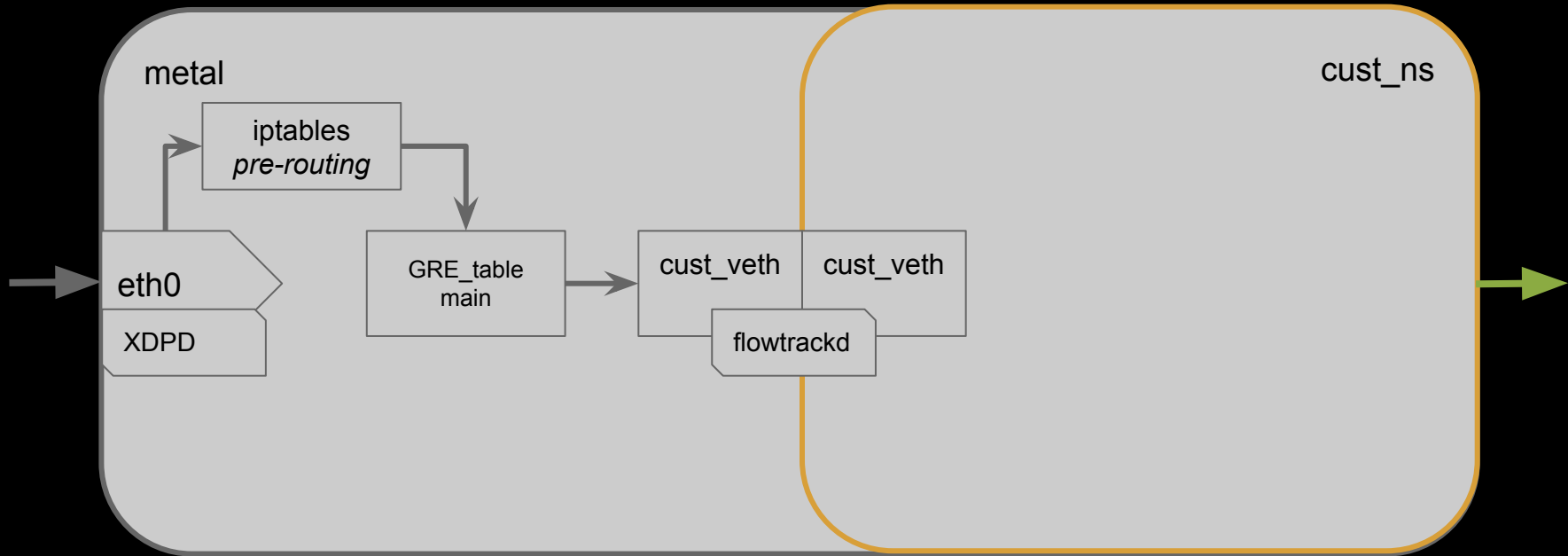
# Getting Traffic Into The Namespace



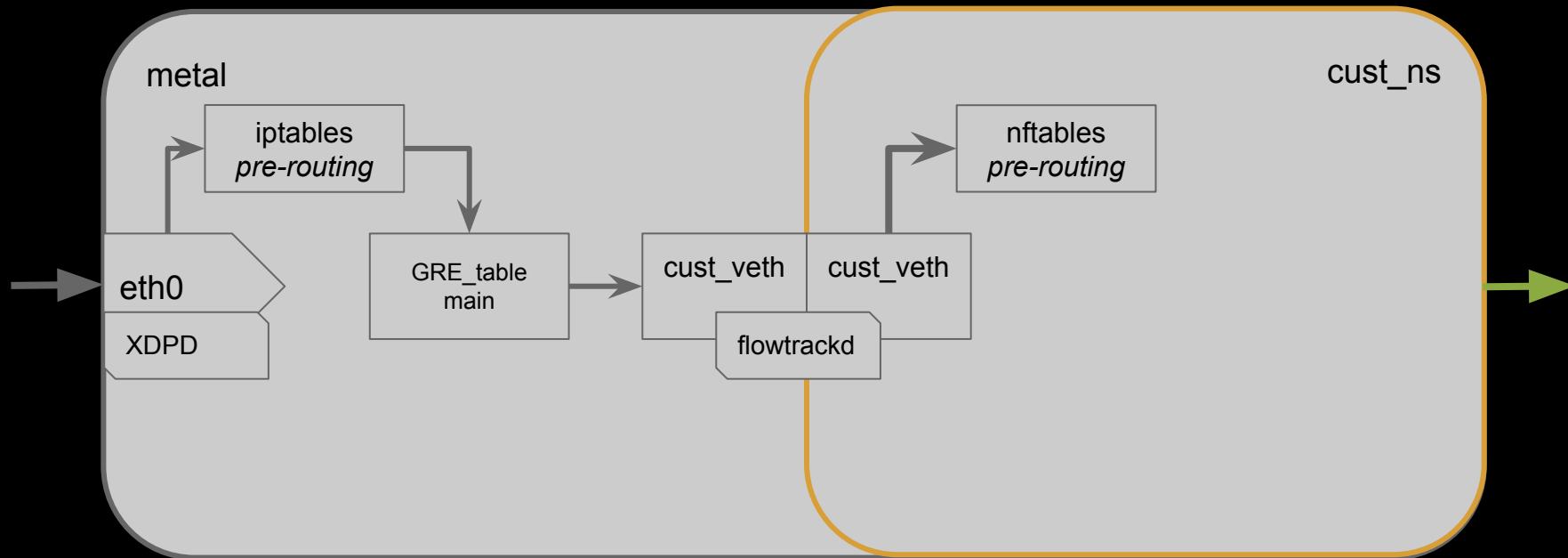
# veth IPs...aaand its gone



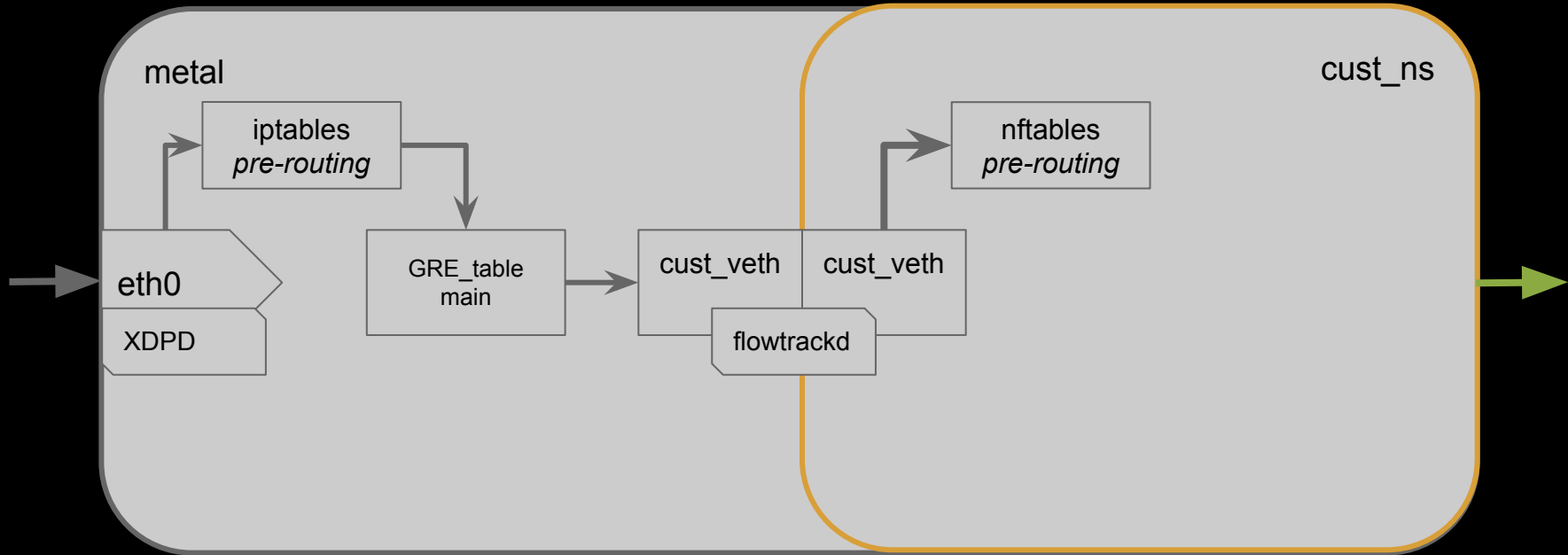
# Connection tracking with half the connection, with AF\_XDP



# Your Traffic, Your FW Rules

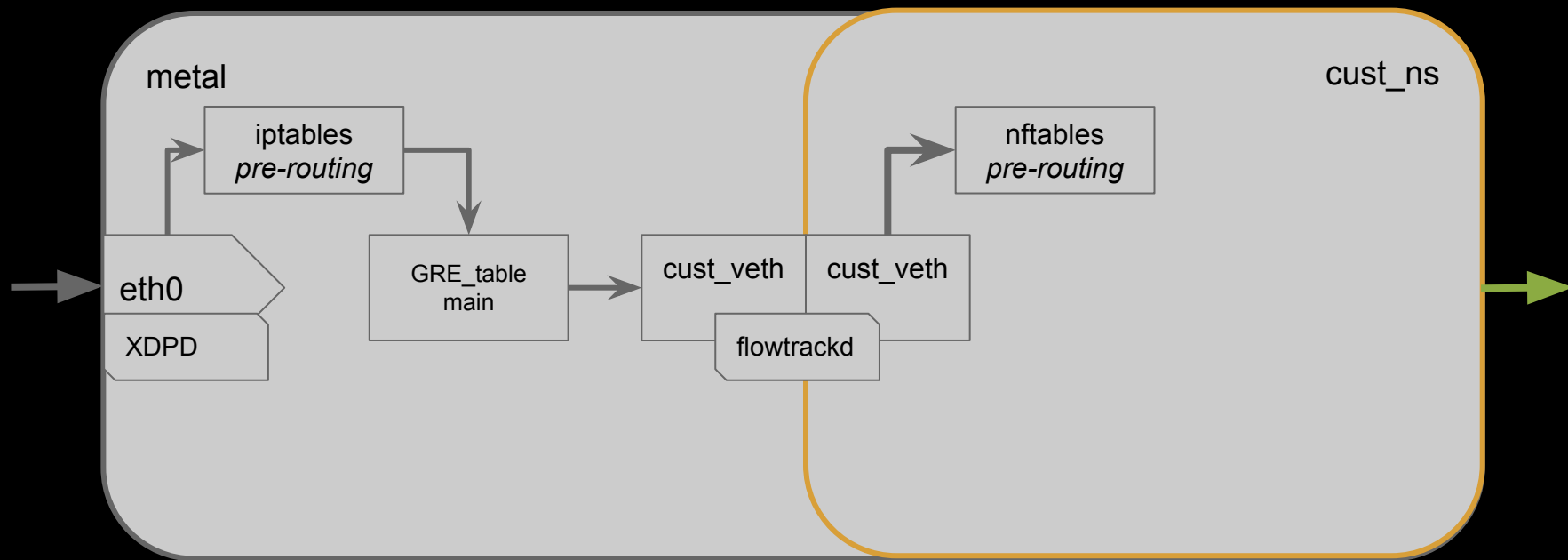


# Dummy nft rule + conntrack

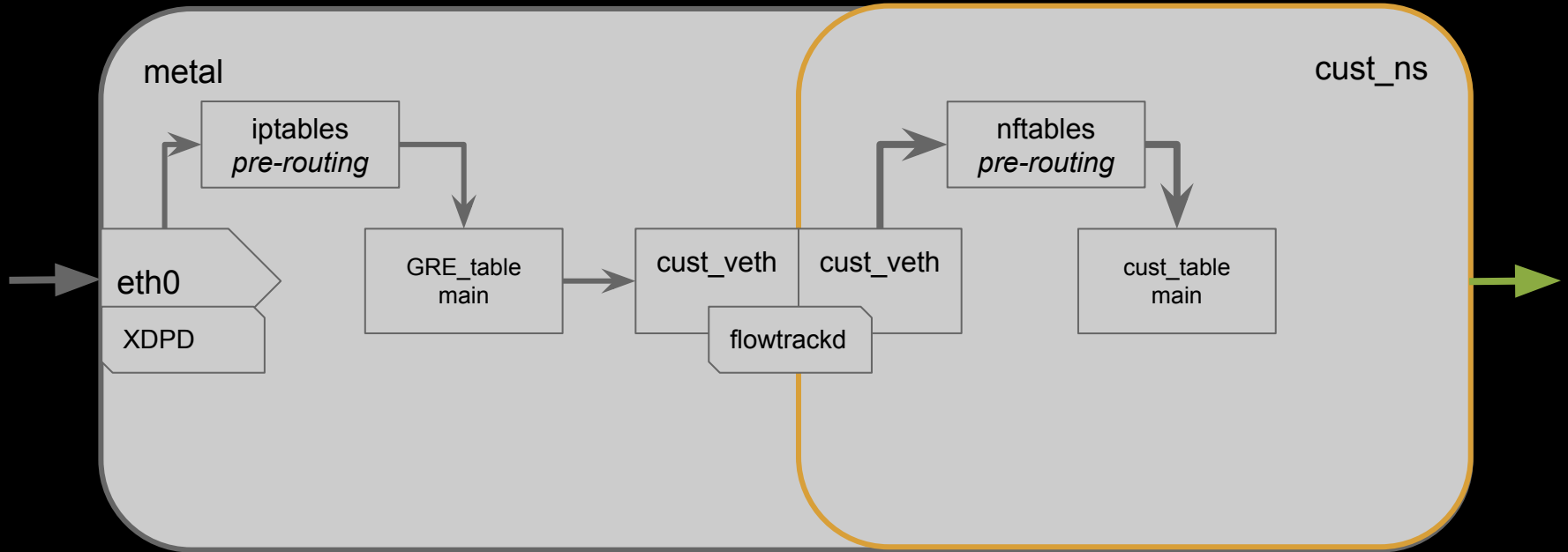




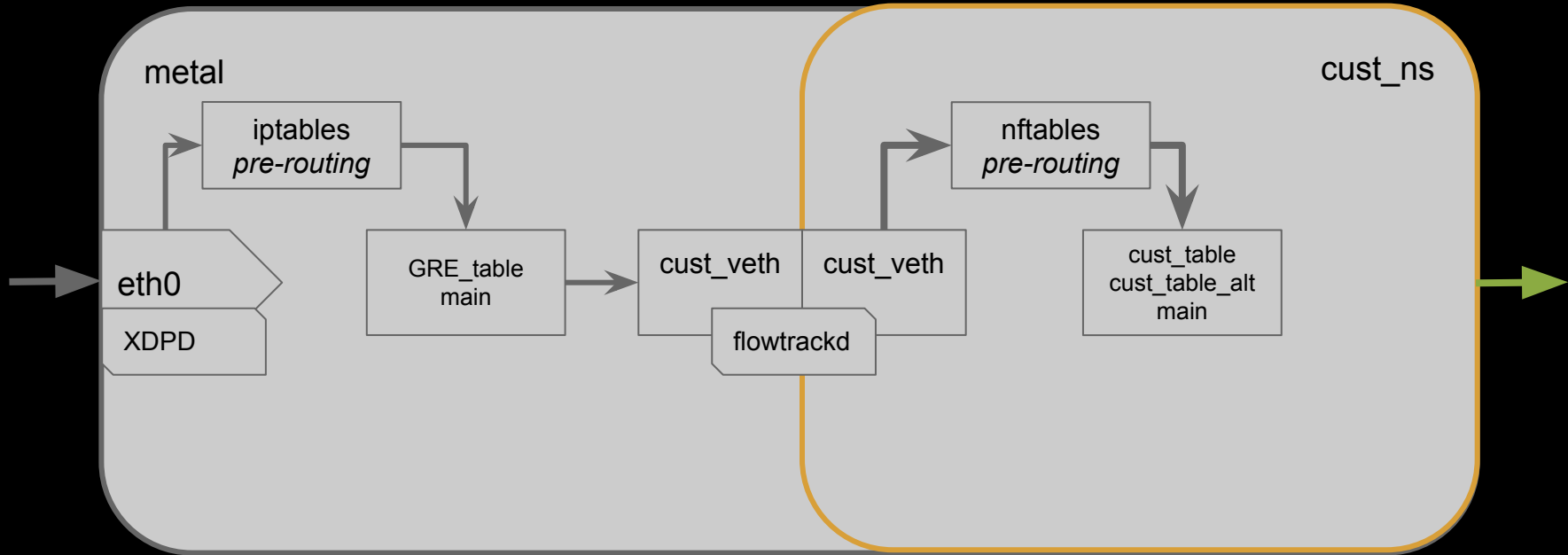
# Dummy nft rule + notrack



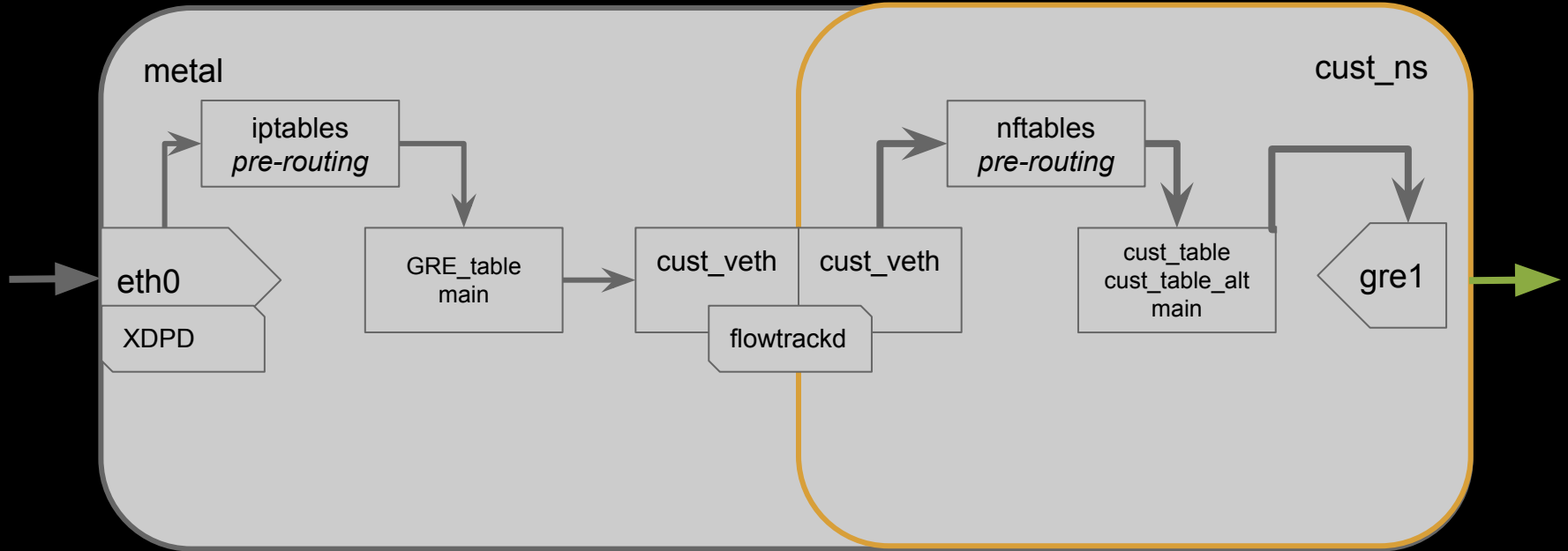
# Route Tables With an S, The Second, pt 1



# Route Tables With an S, The Second, pt 2

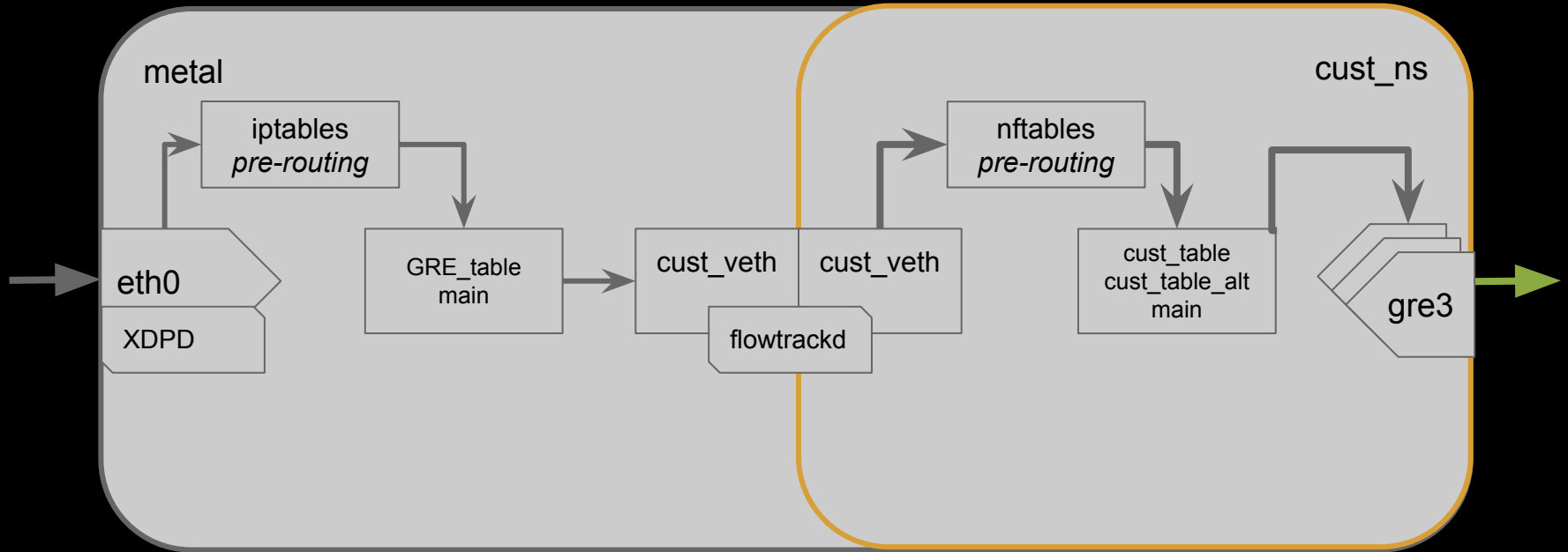


# Wrapped With A Bow



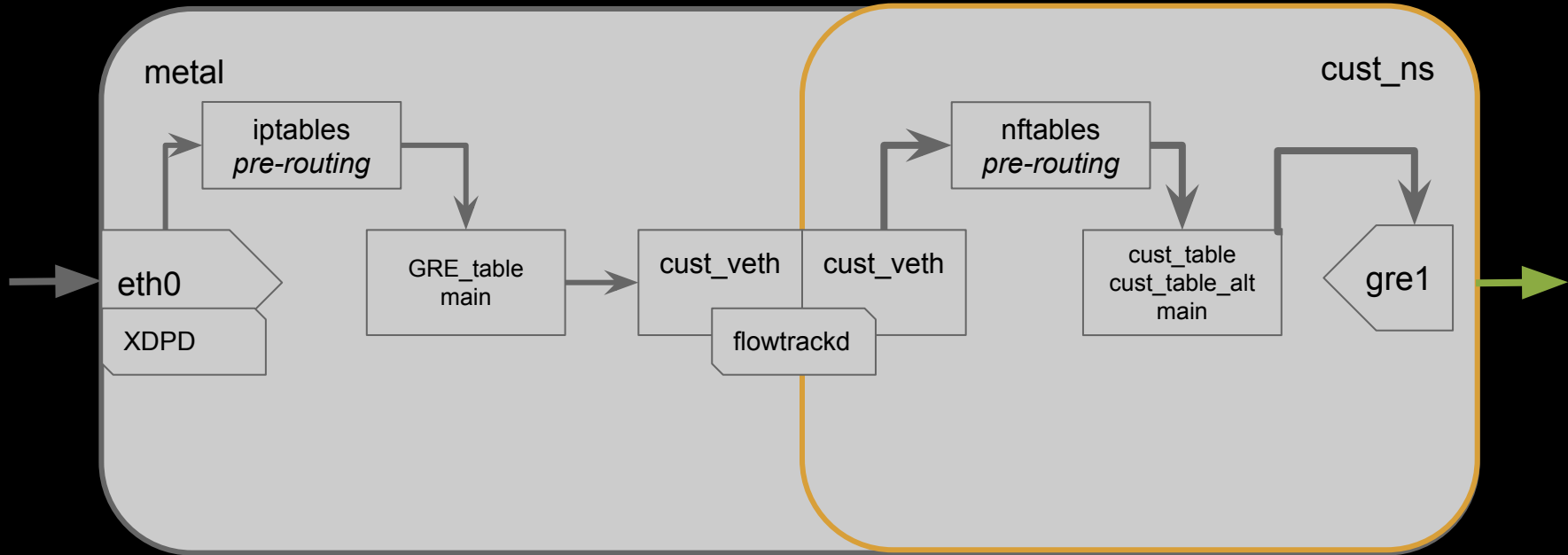
# Devices, Devices, Devices

## *ECMP and Route Priorities*

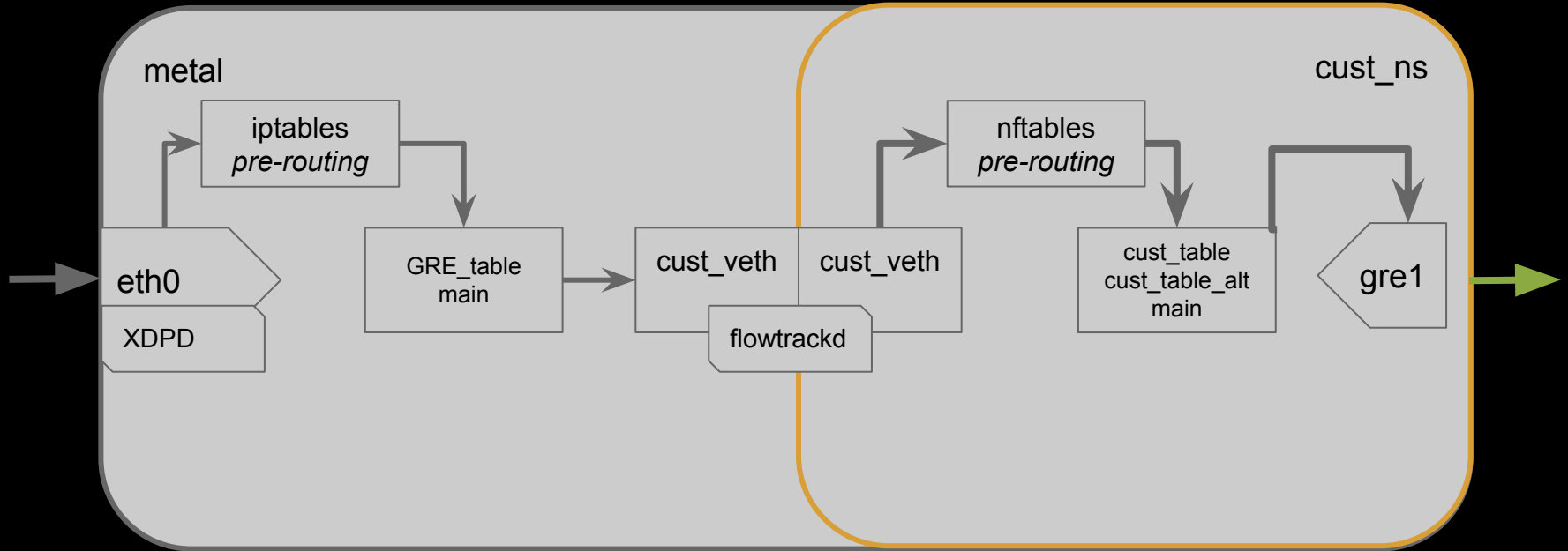


# Devices, Devices, Devices

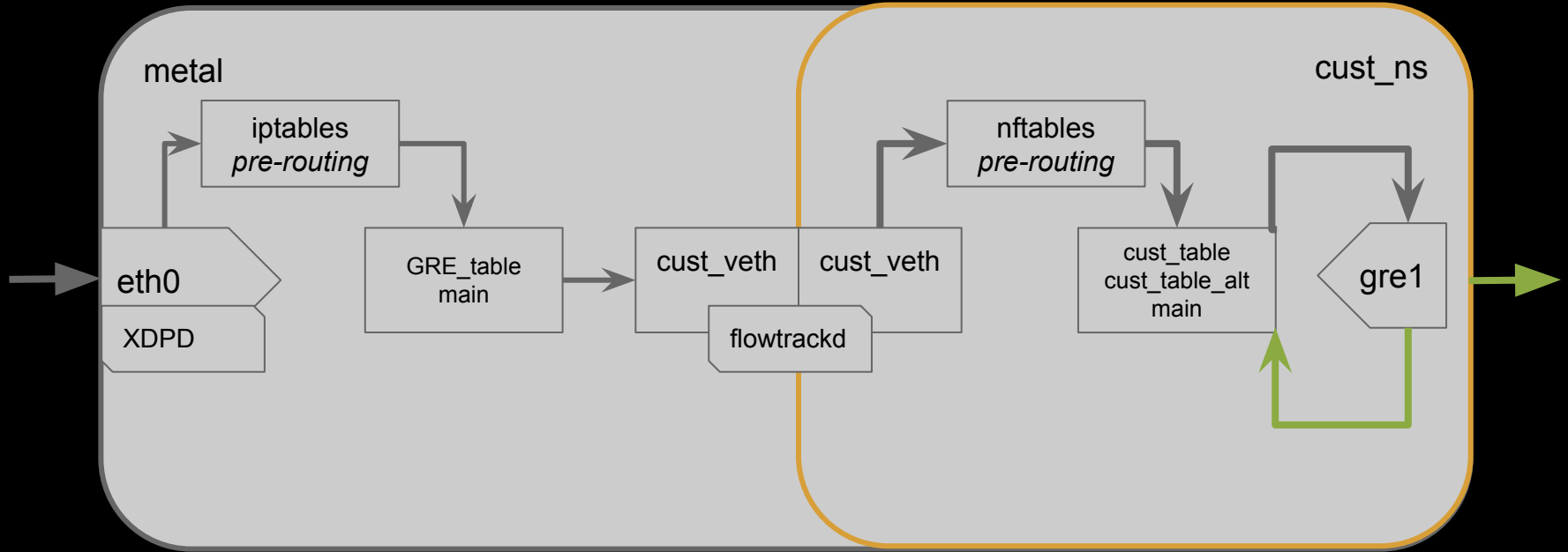
## MAGIC-3



# Focusing On One Device

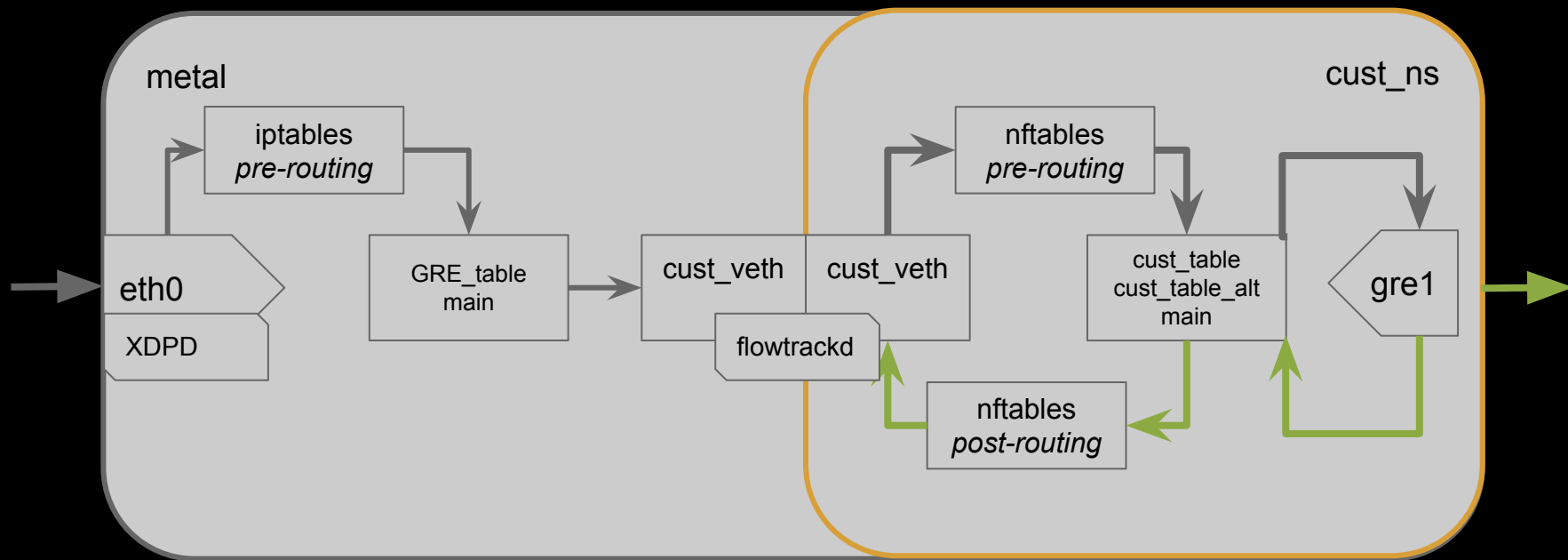


# Default Route

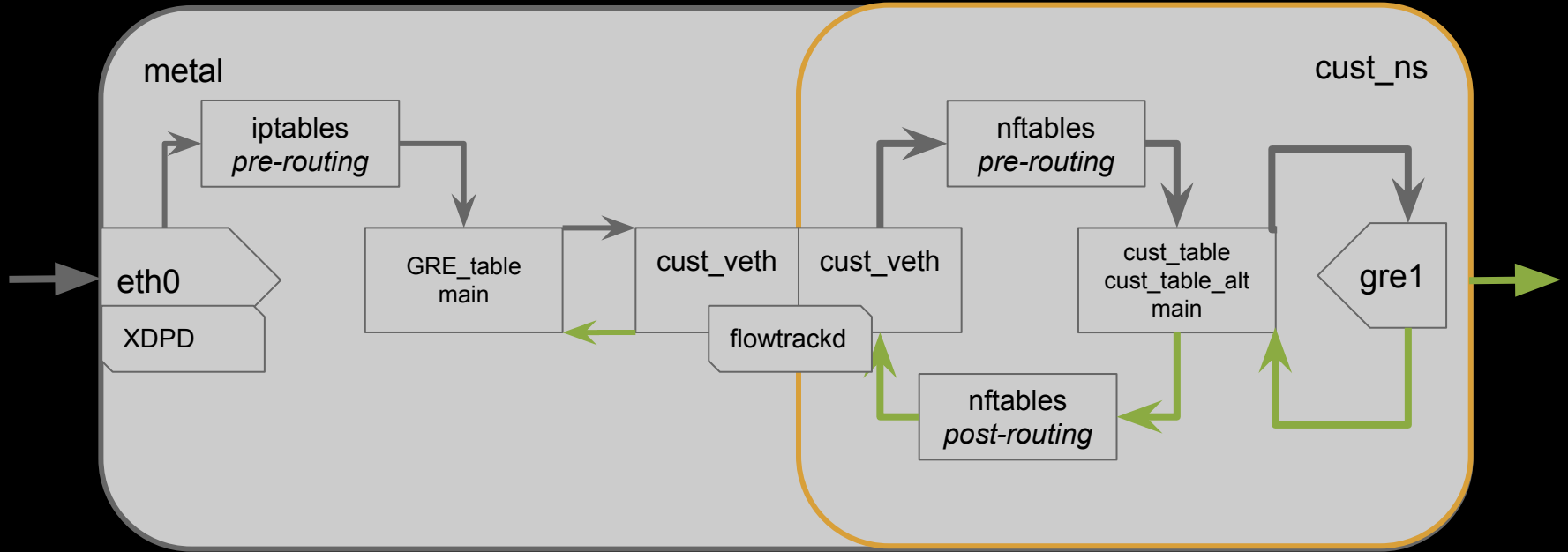




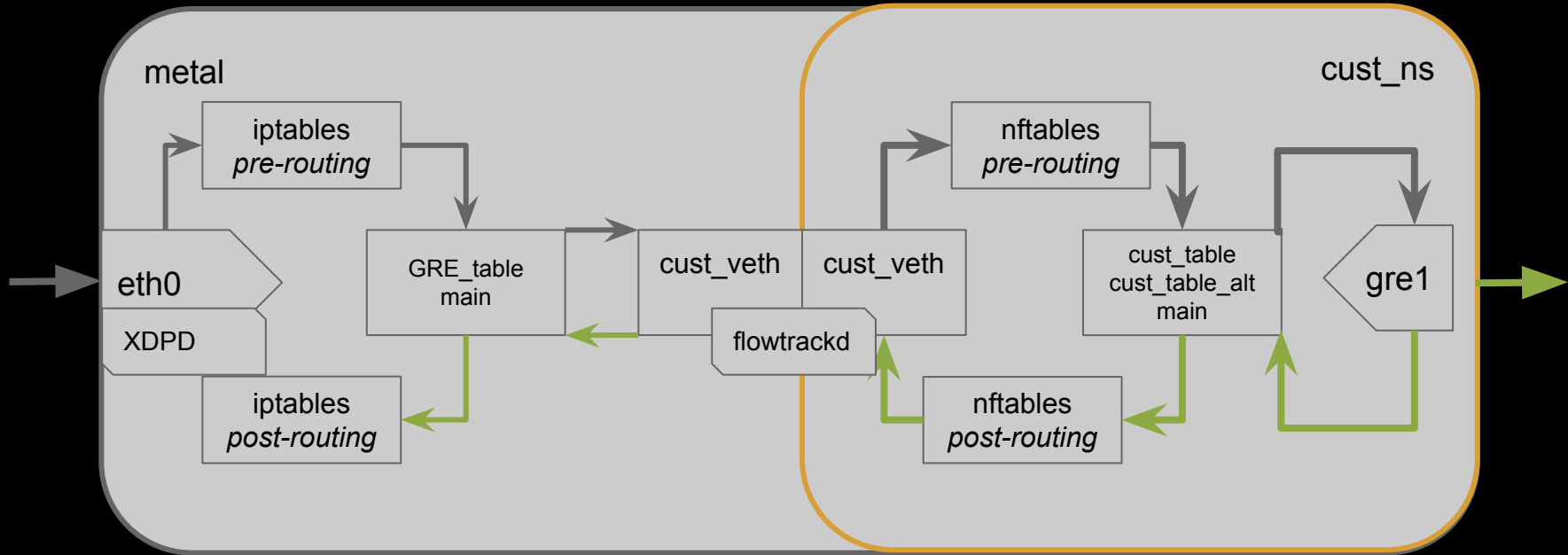
# Outgoing Rules



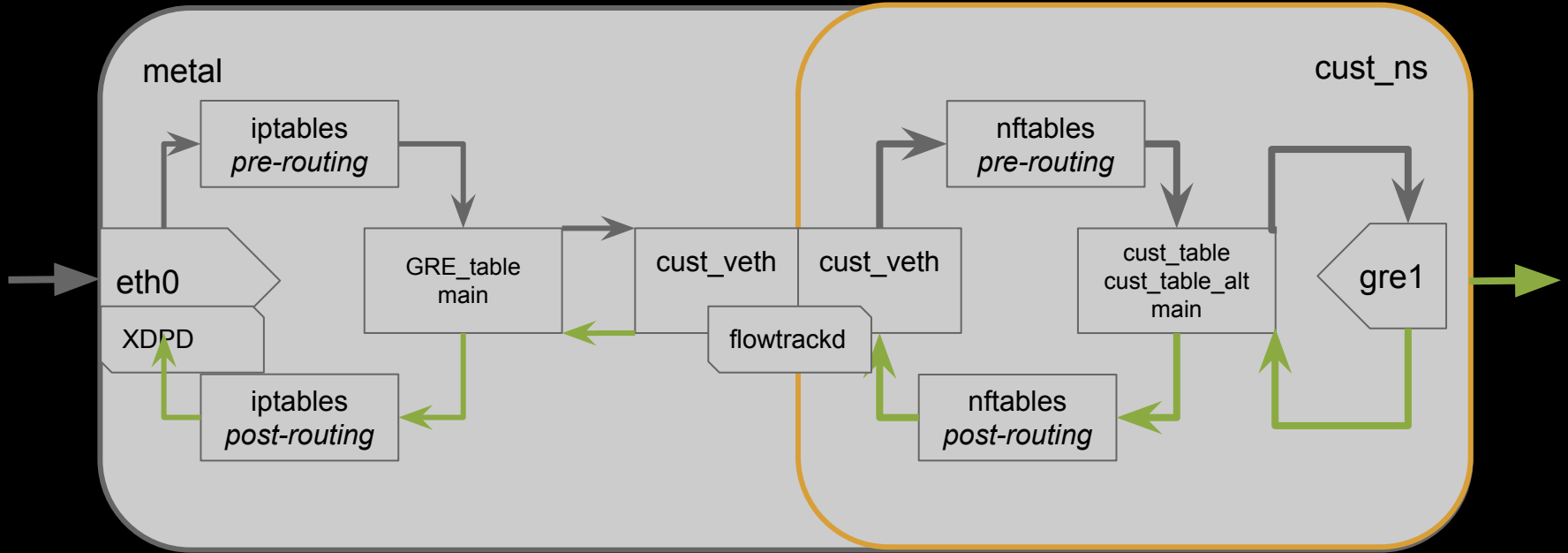
# Traverse The veth, again



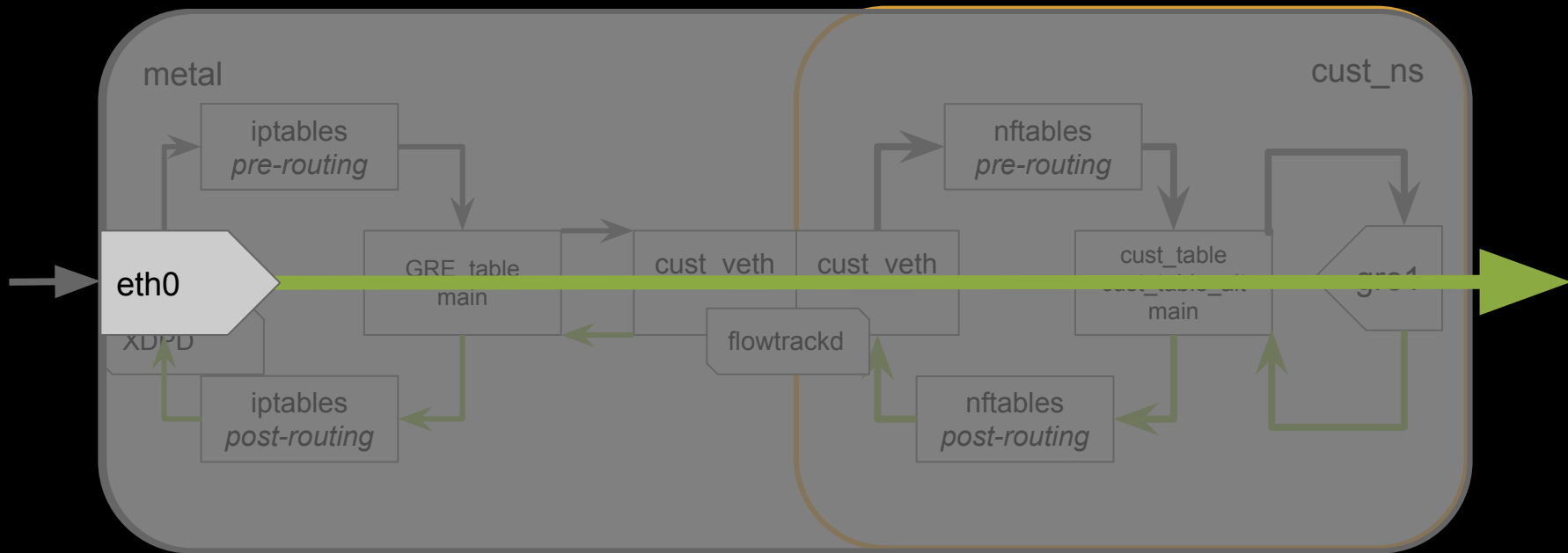
# Outgoing rules



# NIC Again

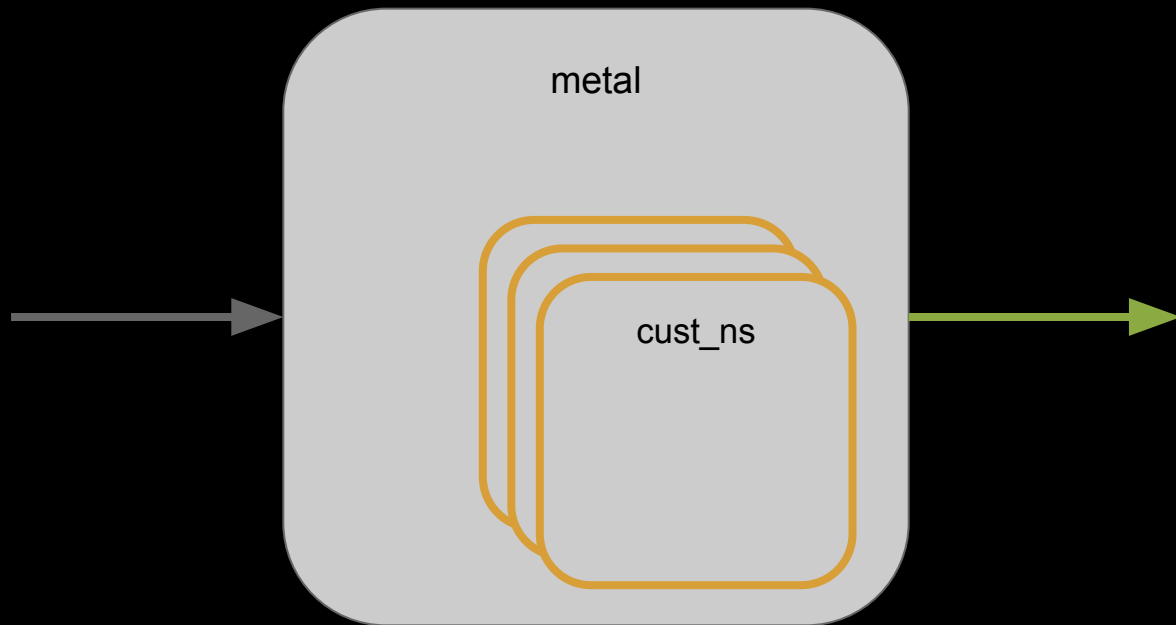


# shipit!

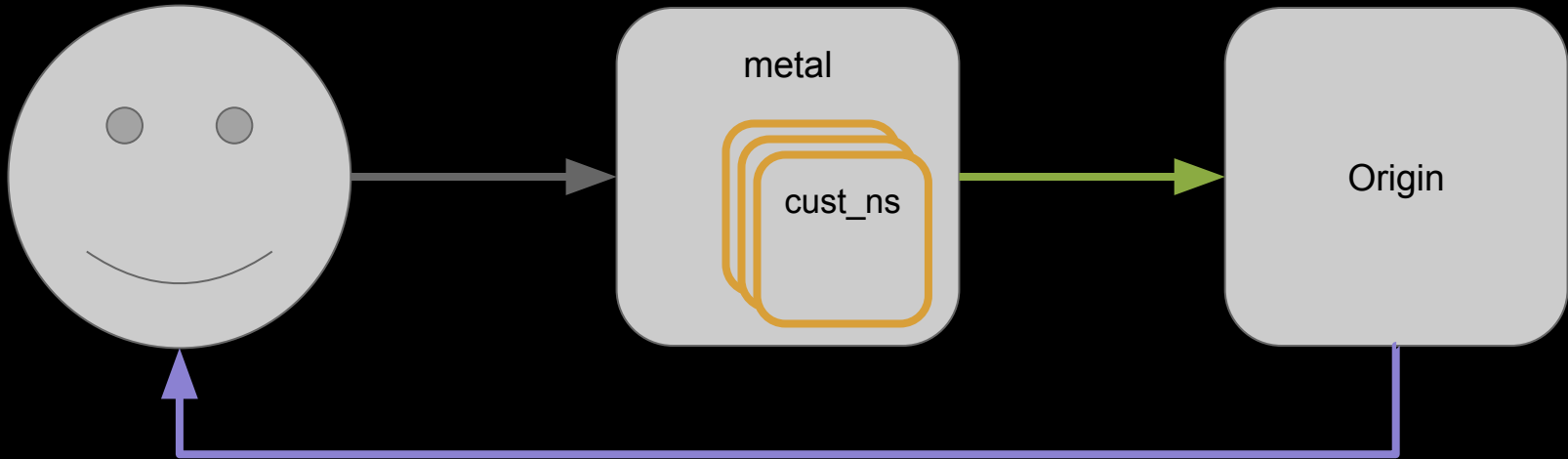


# Whoops! We bricked a metal!

## *Slow-loop/Fast-loop Config Updates*



<todo>quippy gre-ping title</todo>



# Open Questions

- tx/rx checksum offloading
- Disappearing IPs
- ARP?



# Questions?

{erich, conjones}@cloudflare.com