# The Design and Implementation of Open vSwitch

Ben Pfaff[*]
Justin Pettit[*]
Teemu Koponen[*]
Ethan J. Jackson[*]
Andy Zhou[*]
Jarno Rajahalme[*]
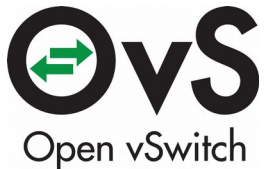Jesse Gross[*]
Alex Wang[*]
Jonathan Stringer[*]
Pravin Shelar[*]
Keith Amidon[†]
Martin Casado[*]
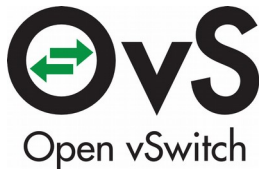
[*]VMware
[†]Awake Networks

Open vSwitch

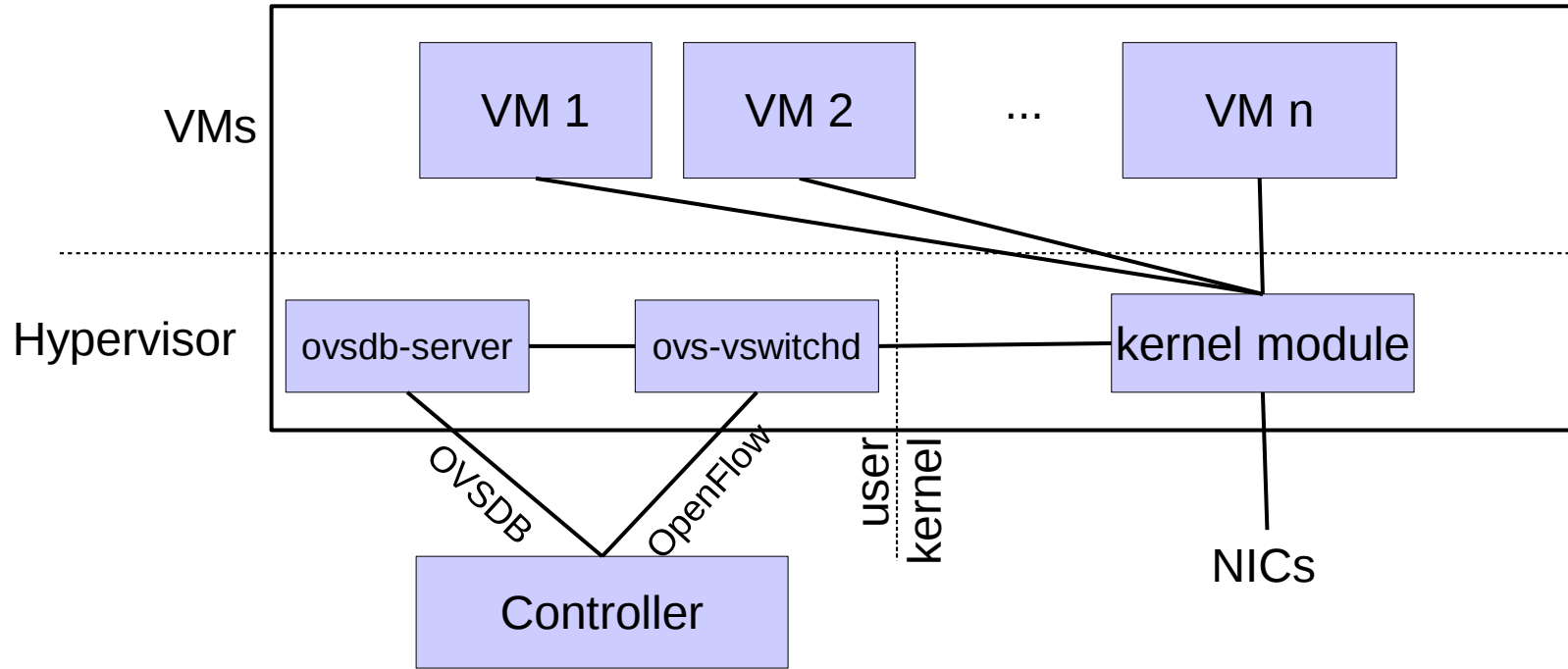# What is Open vSwitch?

From openvswitch.org:

Open vSwitch is a production quality, multilayer virtual switch licensed under the open source Apache 2.0 license.  It is designed to enable massive network automation through programmatic extension, while still supporting standard management interfaces and protocols (e.g. NetFlow, sFlow, SPAN, RSPAN, CLI, LACP, 802.1ag).
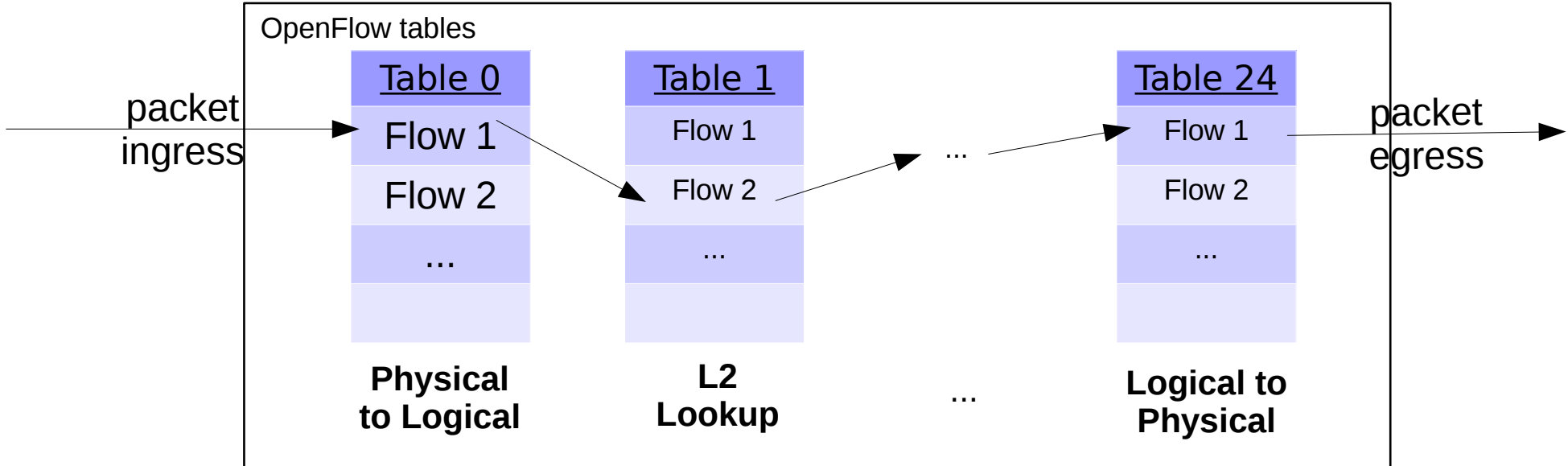
# Where is Open vSwitch Used?

- Broad support:
  - Linux, FreeBSD, NetBSD, Windows, ESX
  - KVM, Xen, Docker, VirtualBox, Hyper-V, …
  - OpenStack, CloudStack, OpenNebula, …
- Widely used:
  - Most popular OpenStack networking backend
  - Default network stack in XenServer
  - 1,440 hits in Google Scholar

# Open vSwitch Architecture

# Network Virtualization Use Case

# Implications for Forwarding Performance

OpenFlow tables

| Table 0 | Table 1 | | Table 24 |
|---------|---------|-----|----------|
| Flow 1 | Flow 1 | | Flow 1 |
| Flow 2 | Flow 2 | ... | Flow 2 |
| ... | ... | | ... |

packet ingress →

packet egress →

Physical to Logical

L2 Lookup

...

Logical to Physical

$k_0$ **hash lookups**

$k_1$ **hash lookups**

...

$k_{24}$ **hash lookups**

**100+ hash lookups per packet for tuple space search?**

OvS
Open vSwitch

# Non-solutions

- All of these helped:
    - Multithreading
    - Userspace RCU
    - Batching packet processing
    - Classifier optimizations
    - Microoptimizations
- None of it helped enough: % versus x.

  **Classification is expensive on general-purpose CPUs!**

# OVS Cache v1: Microflow Cache

Microflow:
- Complete set of packet headers and metadata
- Suitable for hash table
- Shaded data below:

| Eth | IP | TCP | payload |
|-----|----|----|---------|

hit

miss

Microflow Cache

kernel

OpenFlow Tables

userspace

OpenFlow Controller (in theory)

OvS
Open vSwitch

# Speedup with Microflow Cache

Microflow cache
(1 hash lookup)

OpenFlow tables

packet
ingress

| Table 0 | Table 1 | | Table 24 |
|---------|---------|----|----------|
| Flow 1 | Flow 1 | ... | Flow 1 |
| Flow 2 | Flow 2 | | Flow 2 |
| ... | ... | | ... |
| | | | |

packet
egress

Physical to
Logical

L2
Lookup

...

Logical to
Physical

$k_0$ hash
lookups

$k_1$ hash
lookups

...

$k_{24}$ hash
lookups

**From 100+ hash lookups per packet, to just 1!**

⊖vS

Open vSwitch

# Microflow Caching in Practice

- Tremendous speedup for most workloads
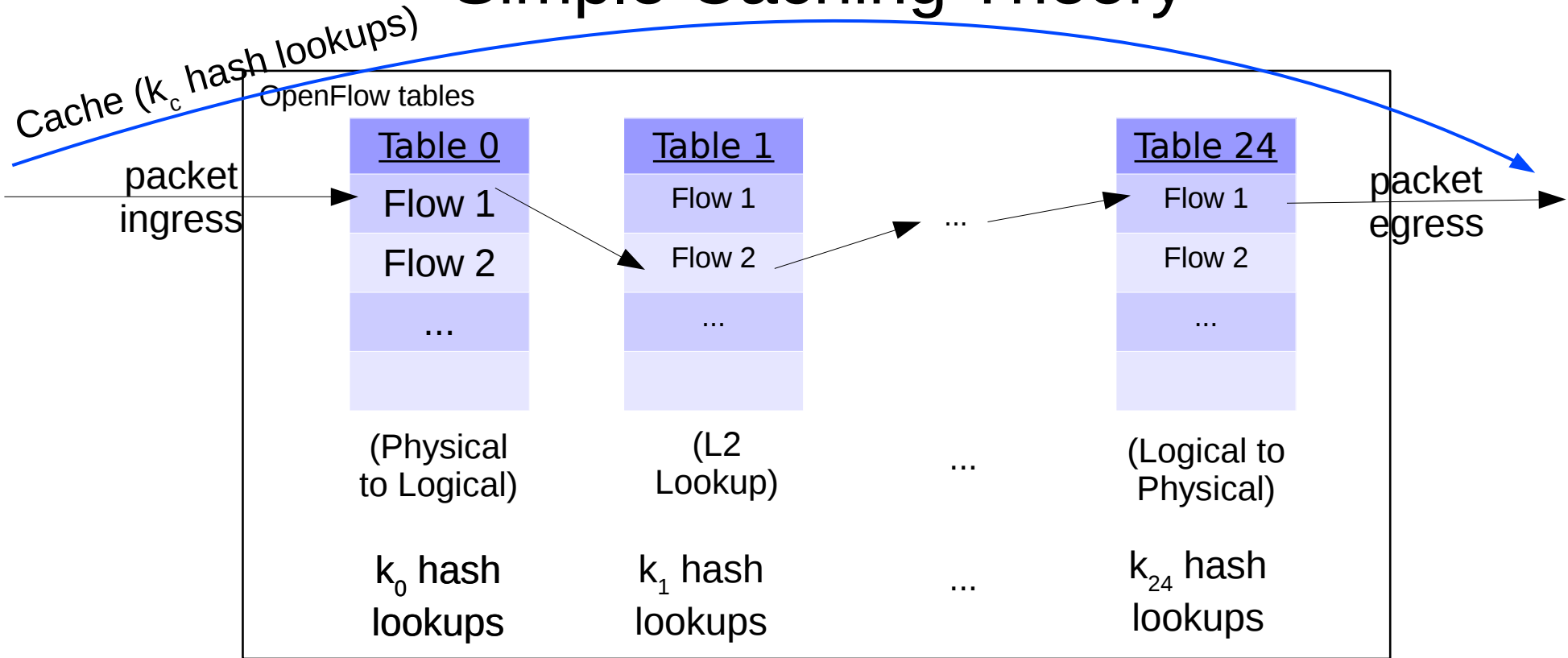
- Problematic traffic patterns:

  - Port scans

    - Malicious

    - Accidental (!)

  - Peer-to-peer rendezvous applications

  - Some kinds of network testing

- All of this traffic has lots of short-lived microflows

  - Fundamental caching problem: low hit rate

# Simple Caching Theory

Cache ($k_c$ hash lookups)

packet ingress

packet egress

OpenFlow tables

| Table 0 | Table 1 | | Table 24 |
|---------|---------|-----|----------|
| Flow 1 | Flow 1 | | Flow 1 |
| Flow 2 | Flow 2 | ... | Flow 2 |
| ... | ... | | ... |
| | | | |

(Physical to Logical)

(L2 Lookup)

...

(Logical to Physical)

$k_0$ hash lookups

$k_1$ hash lookups

...

$k_{24}$ hash lookups

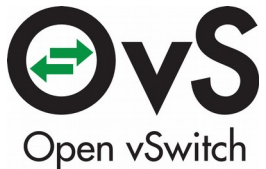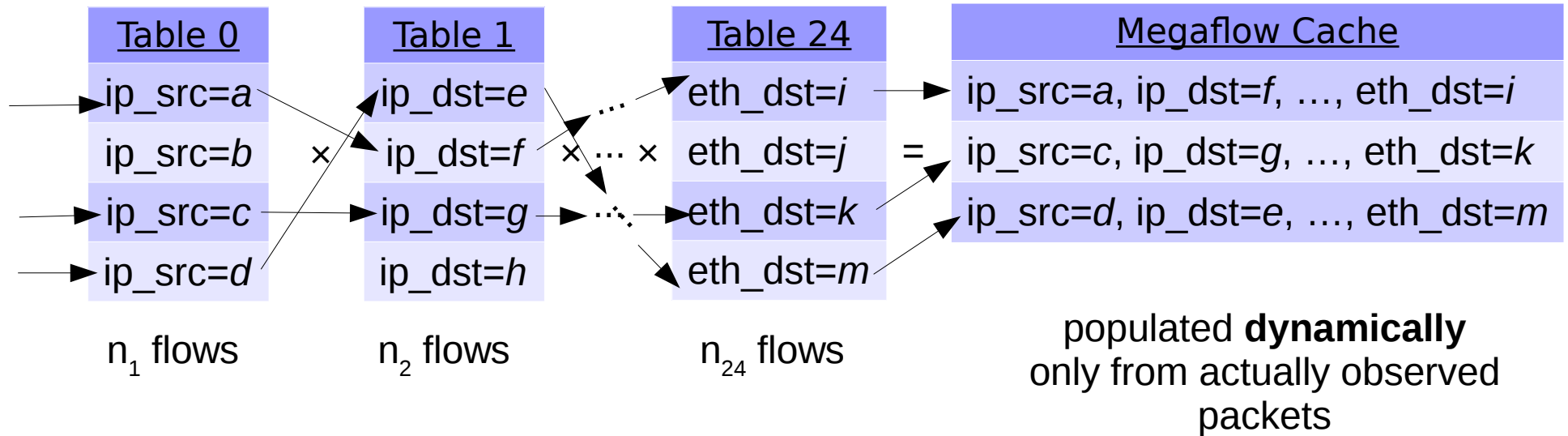**If $k_c$ << $k_0$ + $k_1$ + ... + $k_{24}$: benefit!**

OvS

Open vSwitch

# Naive Approach to Reducing Table Lookups

Combine tables 0...24 into one flow table. Easy! Usually, $k_c << k_0 + k_1 + \ldots + k_{24}$. But:

| Table 0 | | Table 1 | | | | Table 24 | | Table 0+1+...+24 |
|---------|---|---------|---|---|---|----------|---|------------------|
| ip_src=$a$ | | ip_dst=$e$ | | | | eth_dst=$i$ | | ip_src=$a$, ip_dst=$e$, ..., eth_dst=$i$ |
| ip_src=$b$ | × | ip_dst=$f$ | × ⋯ × | | | eth_dst=$j$ | = | ip_src=$a$, ip_dst=$e$, ..., eth_dst=$j$ |
| ip_src=$c$ | | ip_dst=$g$ | | | | eth_dst=$k$ | | ip_src=$a$, ip_dst=$e$, ..., eth_dst=$k$ |
| ip_src=$d$ | | ip_dst=$h$ | | | | eth_dst=$m$ | | ... |
| | | | | | | | | ip_src=$d$, ip_dst=$h$, ..., eth_dst=$k$ |
| | | | | | | | | ip_src=$d$, ip_dst=$h$, ..., eth_dst=$m$ |

$n_1$ flows      $n_2$ flows      $n_{24}$ flows

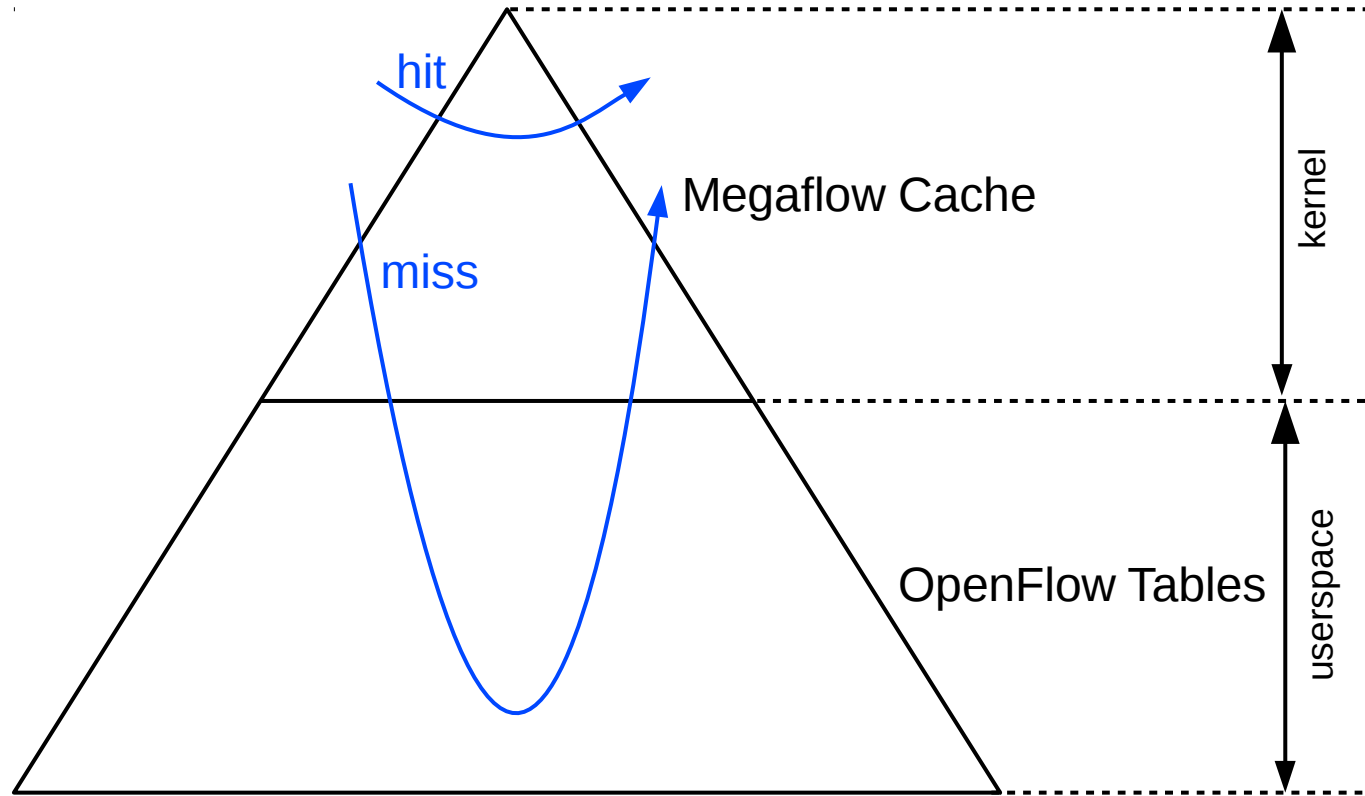up to $n_1 \times n_2 \times \cdots \times n_{24}$ flows

**"Crossproduct Problem"**

# Lazy Approach to Reducing Table Lookups

Solution: Build cache of combined "megaflows" **lazily** as packets arrive.

| Table 0 | Table 1 | Table 24 | Megaflow Cache |
|---------|---------|----------|----------------|
| ip_src=$a$ | ip_dst=$e$ | eth_dst=$i$ | ip_src=$a$, ip_dst=$f$, …, eth_dst=$i$ |
| ip_src=$b$ × | ip_dst=$f$ × … × | eth_dst=$j$ = | ip_src=$c$, ip_dst=$g$, …, eth_dst=$k$ |
| ip_src=$c$ | ip_dst=$g$ | eth_dst=$k$ | ip_src=$d$, ip_dst=$e$, …, eth_dst=$m$ |
| ip_src=$d$ | ip_dst=$h$ | eth_dst=$m$ | |

$n_1$ flows                     $n_2$ flows                     $n_{24}$ flows

populated **dynamically**
only from actually observed
packets

OvS
Open vSwitch

**Same (or better!) table lookups as naive approach.
Traffic locality yields practical cache size.**

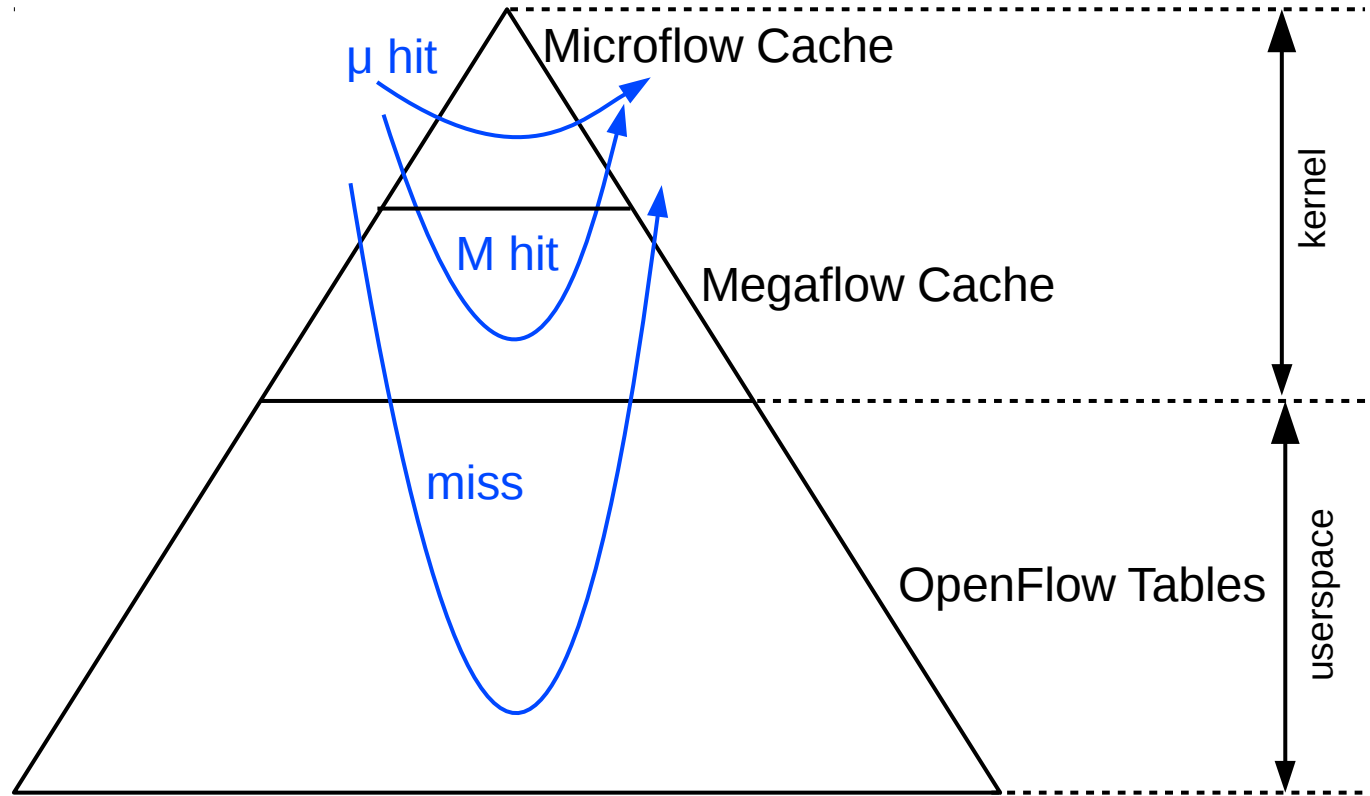# OVS Cache v2: "Megaflow" Cache

# Generating Good Megaflows

- Goal: Less-specific megaflows.

    - Megaflows that match TCP ports are almost like microflows!

    - Naive approach tends to match every field that appears anywhere in flow tables

- Requirements:

    - online

    - fast

- Contribution: Megaflow generation improvements (section 5).

# Megaflow vs. Microflow Cache Performance

- Microflow cache:
  - $k_0 + k_1 + \cdots + k_{24}$ lookups for first packet in microflow
  - 1 lookup for later packets in microflow
- Megaflow cache:
  - $k_c$ lookups for every packet
- $k_c > 1$ is normal, so megaflows perform worse in common case!
- Best of both worlds would be:
  - $k_c$ lookups for first packet in microflow
  - 1 lookup for later packets in microflow

# OVS Cache v3: Dual Caches

# Parting Thoughts

- Common tension: expressibility vs. performance
- OpenFlow is expressive but amenable to high performance
- An application-specific switch would not be built like OVS
  - And would likely be **slower**
- Applications can freely evolve decoupled from performance
- Starting from a more general problem produced better results

Open vSwitch