

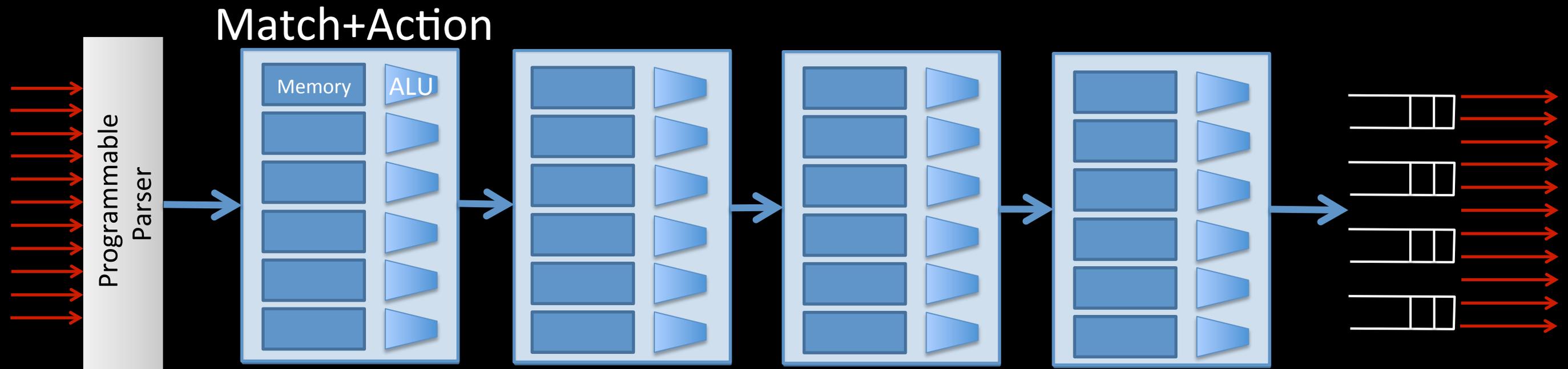


Programming the Forwarding Plane

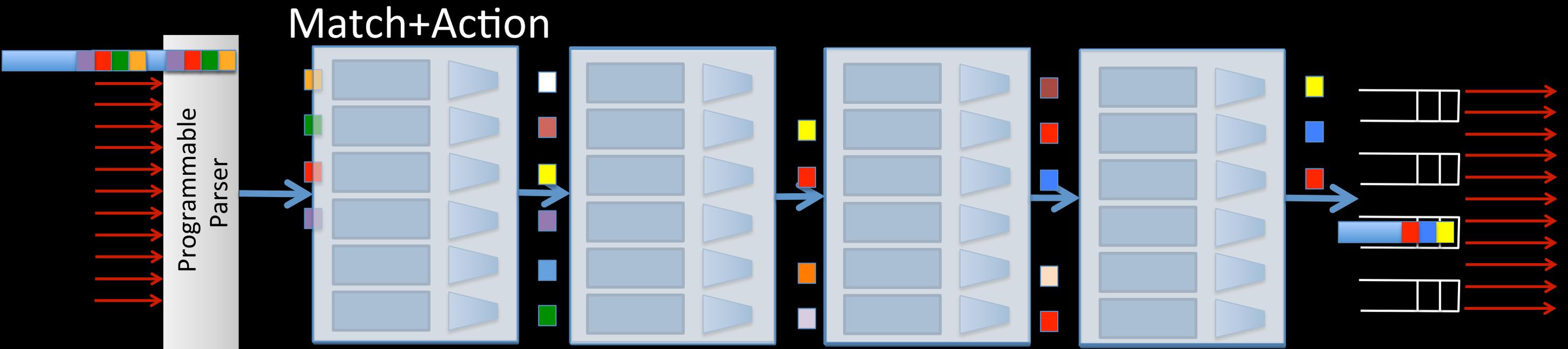
Nick McKeown

PISA: Protocol Independent Switch Architecture

[Sigcomm 2013]



PISA: Protocol Independent Switch Architecture



P4 and PISA

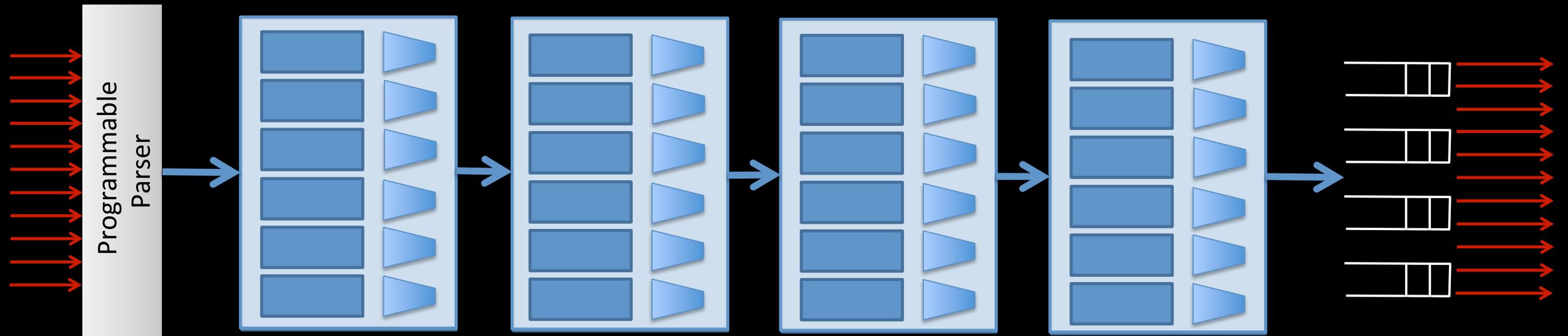
P4 code



Compiler 



Compiler Target



[ACM CCR 2014]
 "Best Paper 2014"

P4: Programming Protocol-Independent Packet Processors

Pat Bosshart¹, Dan Daly^{*}, Glen Gibb¹, Martin Izzard¹, Nick McKeown[‡], Jennifer Rexford^{**}, Cole Schlesinger^{**}, Dan Talayco[†], Amin Vahdat^{*}, George Varghese[§], David Walker^{**},
[†]Barefoot Networks ^{*}Intel [‡]Stanford University ^{**}Princeton University [§]Google [§]Microsoft Research

ABSTRACT

P4 is a high-level language for programming protocol-independent packet processors. P4 works in conjunction with SDN control protocols like OpenFlow. In its current form, OpenFlow explicitly specifies protocol headers on which it operates. This set has grown from 12 to 41 fields in a few years, increasing the complexity of the specification while still not providing the flexibility to add new headers. In this paper we propose P4 as a strawman proposal for how OpenFlow should evolve in the future. We have three goals: (1) Reconfigurability in the field: Programmers should be able to change the way switches process packets once they are deployed. (2) Protocol independence: Switches should not be tied to any specific network protocols. (3) Target independence: Programmers should be able to describe packet-processing functionality independently of the specifics of the underlying hardware. As an example, we describe how to use P4 to configure a switch to add a new hierarchical label.

multiple stages of rule tables, to allow switches to expose more of their capabilities to the controller. The proliferation of new header fields shows no signs of stopping. For example, data-center network operators increasingly want to apply new forms of packet encapsulation (e.g., NVGRE, VXLAN, and STT), for which they resort to deploying software switches that are easier to extend with new functionality. Rather than repeatedly extending the OpenFlow specification, we argue that future switches should support flexible mechanisms for parsing packets and matching header fields, allowing controller applications to leverage these capabilities through a common, open interface (i.e., a new "OpenFlow 2.0" API). Such a general, extensible approach would be simpler, more elegant, and more future-proof than today's OpenFlow 1.x standard.

1. INTRODUCTION

Software-Defined Networking (SDN) gives operators programmatic control over their networks. In SDN, the control plane is physically separate from the forwarding plane, and one control plane controls multiple forwarding devices. While forwarding devices could be programmed in many ways, having a common, open, vendor-agnostic interface (like OpenFlow) enables a control plane to control forwarding devices from different hardware and software vendors.

Version	Date	Header Fields
OF 1.0		12 fields (Ethernet, TCP/IPv4)
OF 1.1	Dec 2009	15 fields (MPLS, inter-table metadata)
OF 1.2	Feb 2011	36 fields (ARP, ICMP, IPv6, etc.)
OF 1.3	Dec 2011	40 fields
OF 1.4	Jun 2012	41 fields

Table 1: Fields recognized by the OpenFlow standard

The OpenFlow interface started simple, with the abstraction of a single table of rules that could match packets on a dozen header fields (e.g., MAC addresses, IP addresses, protocol, TCP/UDP port numbers, etc.). Over the past five years, the specification has grown increasingly complicated (see Table 1), with many

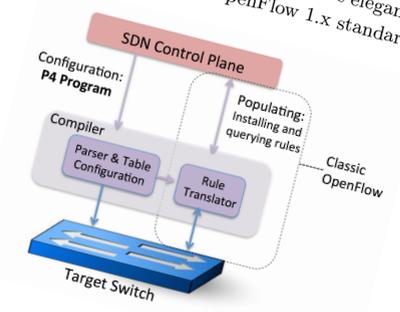


Figure 1: P4 is a language to configure switches.

Recent chip designs demonstrate that such flexibility can be achieved in custom ASICs at terabit speeds [1, 2, 3]. Programming this new generation of switch chips is far from easy. Each chip has its own low-level interface, akin to microcode programming. In this paper, we sketch the design of a higher-level language for programming P4, an independent Packet Processors (P4) Programming Language. We describe the relationship between P4—its syntax, semantics, and how packets are processed—and the underlying hardware, as OpenFlow.



Update on P4 Language Ecosystem

P4.org – P4 Language Consortium

The screenshot shows a web browser window with the URL p4.org. The page features a navigation bar with links for SPEC, CODE, NEWS, JOIN US, and BLOG. The main content area is titled 'BOARD MEMBERS' and lists two members: Nick McKeown from Stanford University and Jennifer Rexford from Princeton University. Below this, a section titled 'Field Reconfigurable' explains that P4 allows network engineers to change switch configurations after deployment. To the right, there is a code snippet for ingress control and a 'TRY IT' button with a GitHub link.

BOARD MEMBERS
Two Board members oversee the consortium:


Nick McKeown
Stanford University


Jennifer Rexford
Princeton University

Field Reconfigurable
P4 allows network engineers to change the way their switches process packets after they are deployed.

```
control ingress {  
  apply (routing);  
}
```

 **TRY IT** Get the code from GitHub

P4.org – P4 Language Consortium

The screenshot shows the P4.org website with a navigation bar containing links for SPEC, CODE, NEWS, JOIN US, and BLOG. The main content area features a list of events and a section titled 'Field Reconfigurable' which includes a description, a code snippet, and a 'TRY IT' button with a GitHub link.

- Regular P4 meetings
- Full-day tutorial at Sig
- 2nd P4 Workshop at S
- 1st P4 Boot camp for PhD students November 19-20
- 1st P4 Developers Day November 19

Field Reconfigurable
P4 allows network engineers to change the way their switches process packets after they are deployed.

```
control ingress {  
    apply (routing) ;  
}
```

[TRY IT](#) Get the code from GitHub



P4 Consortium – P4.org



Operators



Systems



Targets



Academia

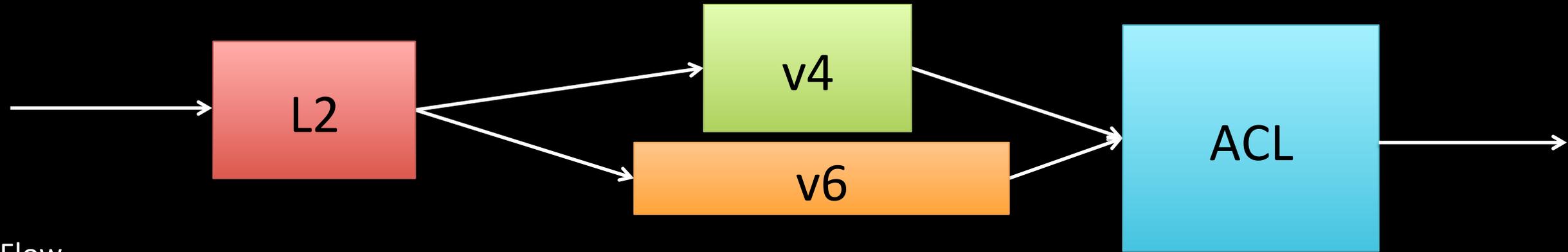


Mapping P4 programs to compiler target

Lavanya Jose, Lisa Yan, George Varghese, NM

[NSDI 2015]

Naïve Mapping: Control Flow Graph



Control Flow

Switch Pipeline

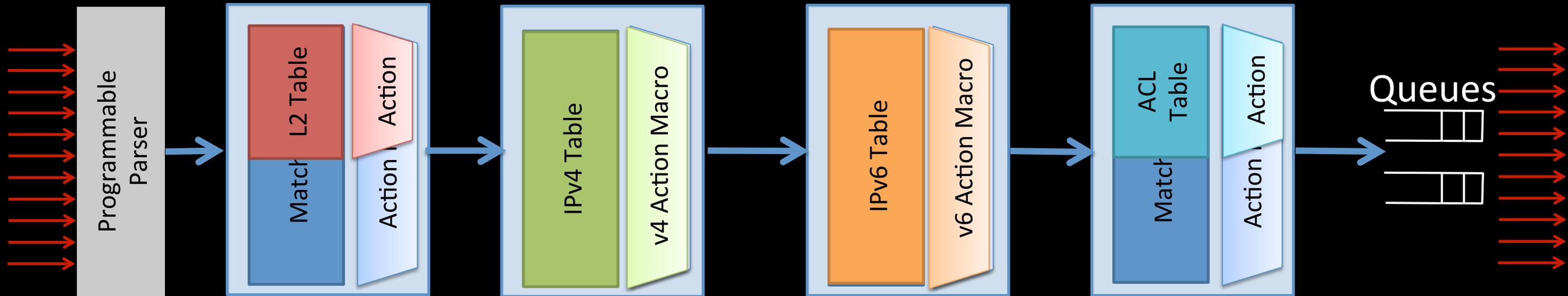
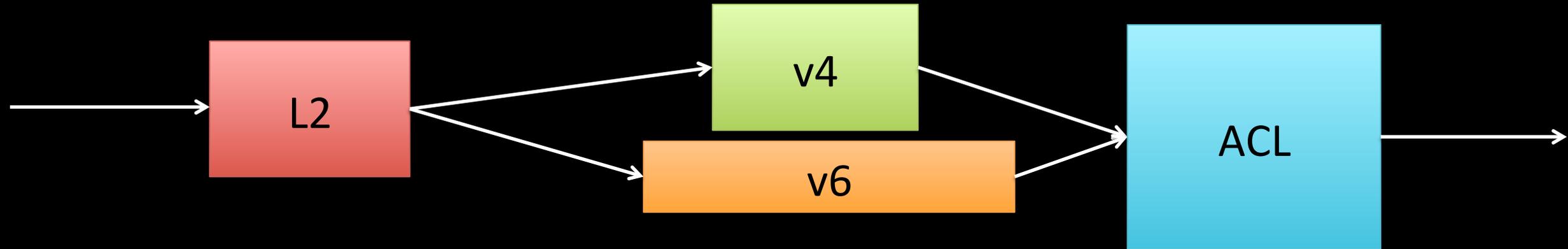
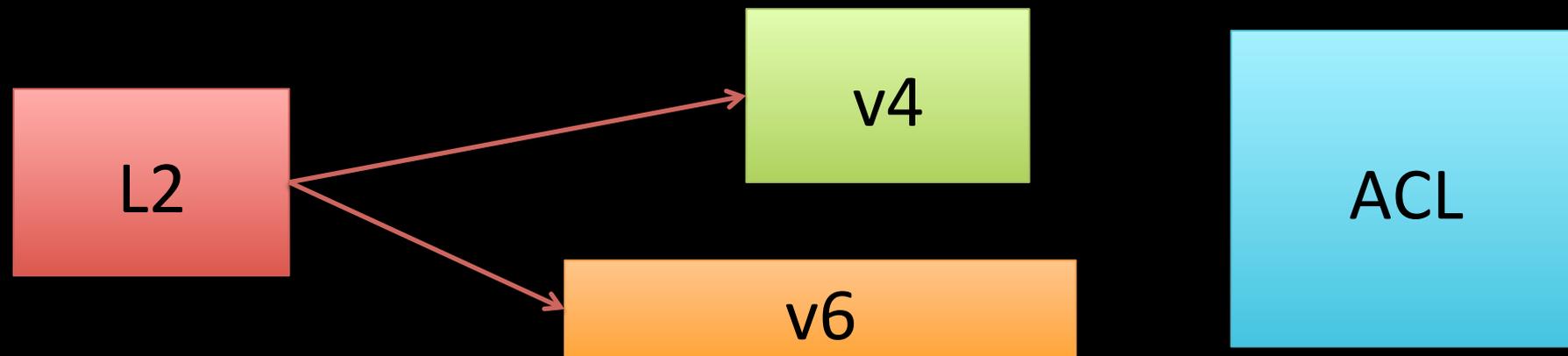


Table Dependency Graph (TDG)

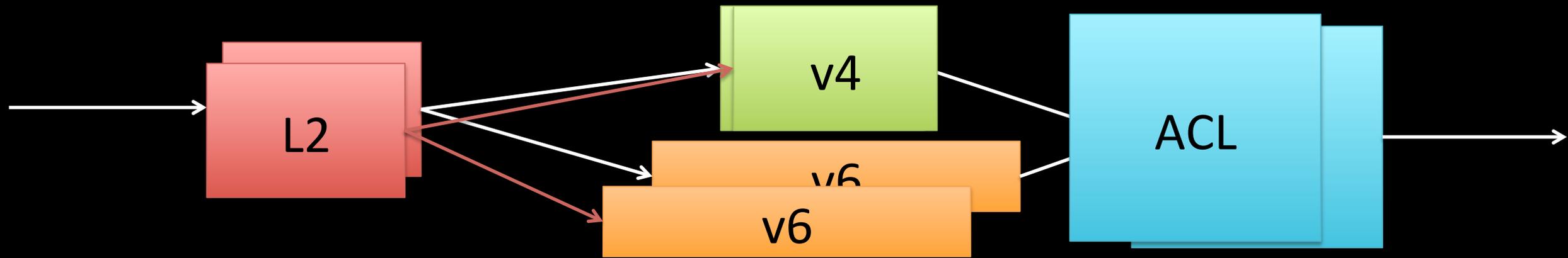


Control Flow Graph

Table Dependency Graph

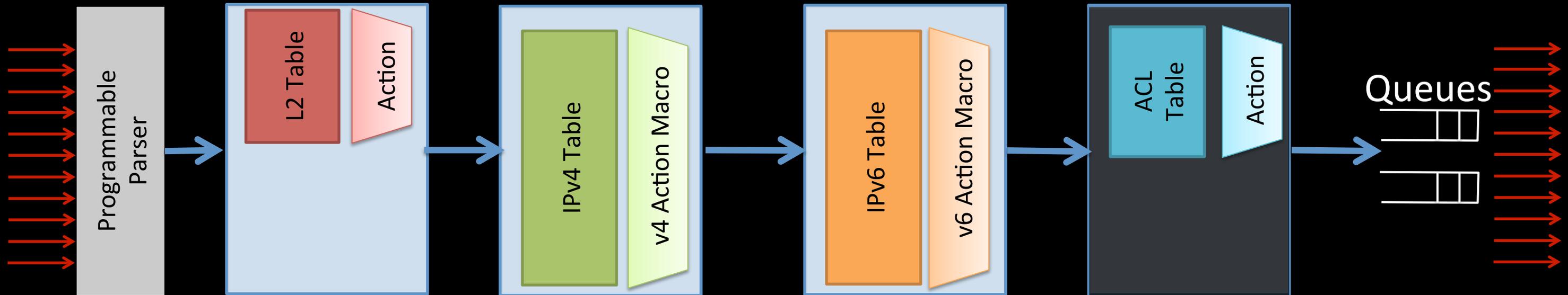


Efficient Mapping: TDG

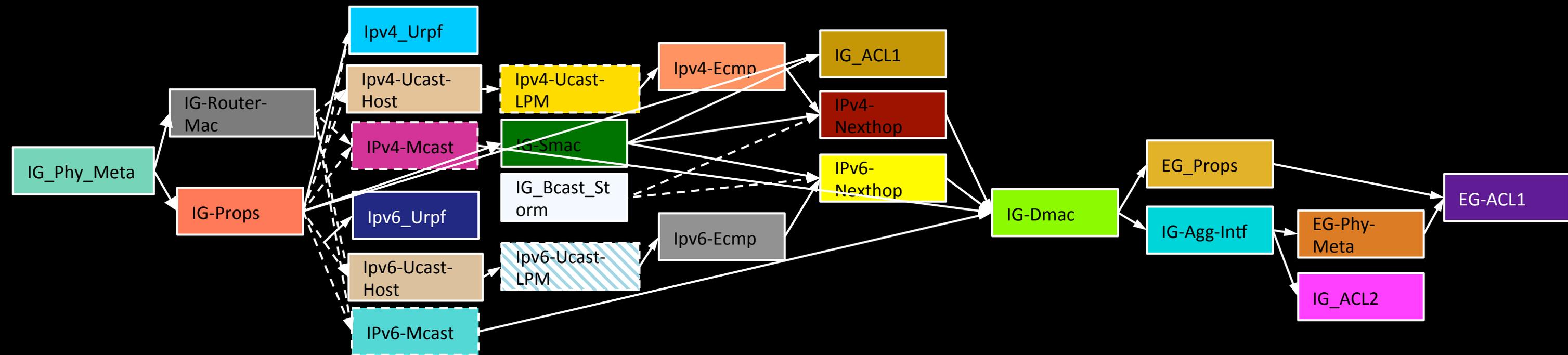


Control Flow Graph

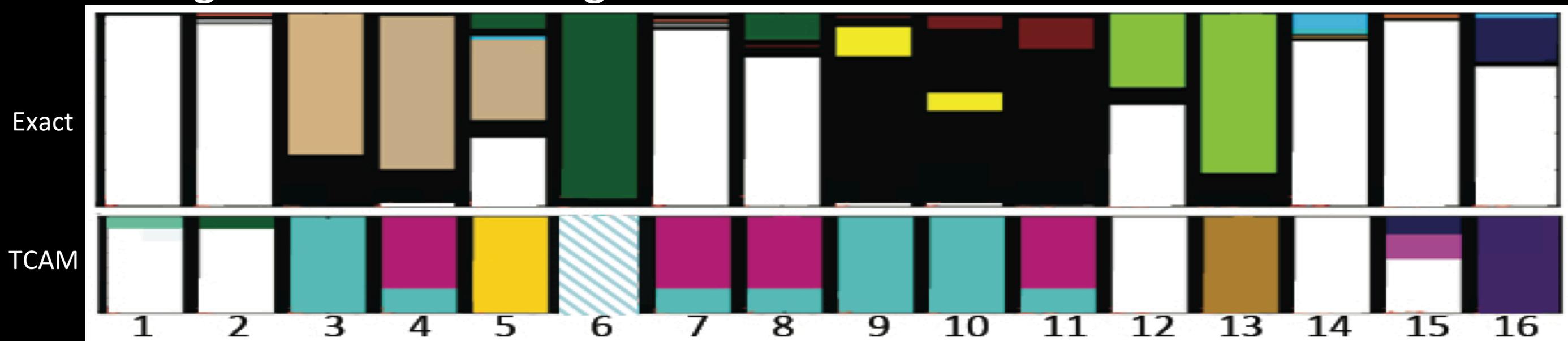
Switch Pipeline



Example Use Case: Typical TDG



Configuration for 16-stage PISA



Mapping Techniques

[NSDI 2015]

Compare: Greedy Algorithm versus Integer Linear Programming (ILP)

Greedy Algorithm runs 100-times faster

ILP Algorithm uses 30% fewer stages

Recommendations:

1. If enough time, use ILP
2. Else, run ILP offline to find best parameters for Greedy algorithm

PISCES: Protocol Independent Software Hypervisor Switch

Mohammad Shahbaz*, Sean Choi, Jen Rexford*, Nick Feamster*, Ben Pfaff, NM

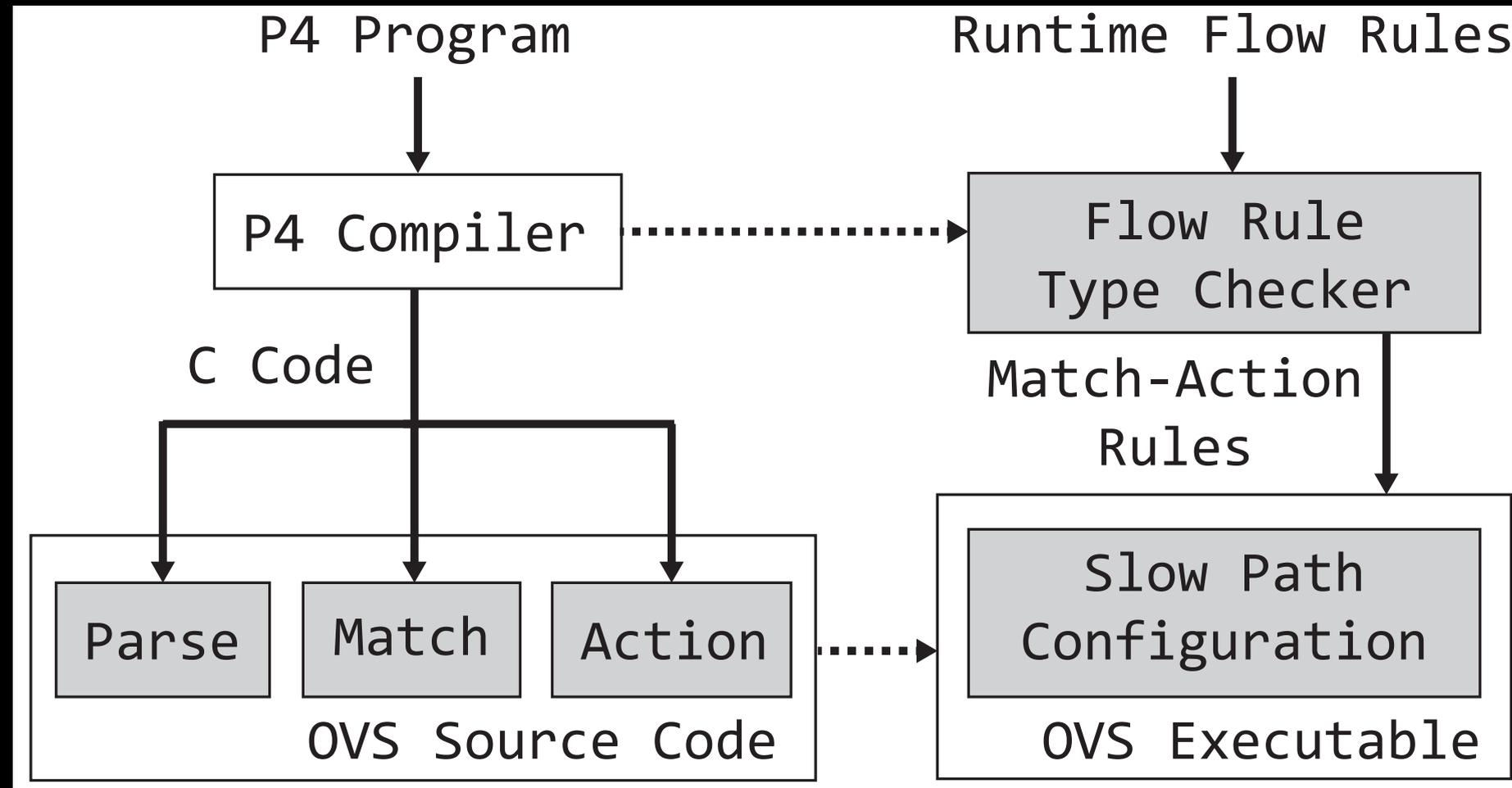
Problem: Adding new protocol feature to OVS is complicated

- Requires domain expertise in kernel programming *and* networking
- Many modules affected
- Long QA and deployment cycle: typically 9 months

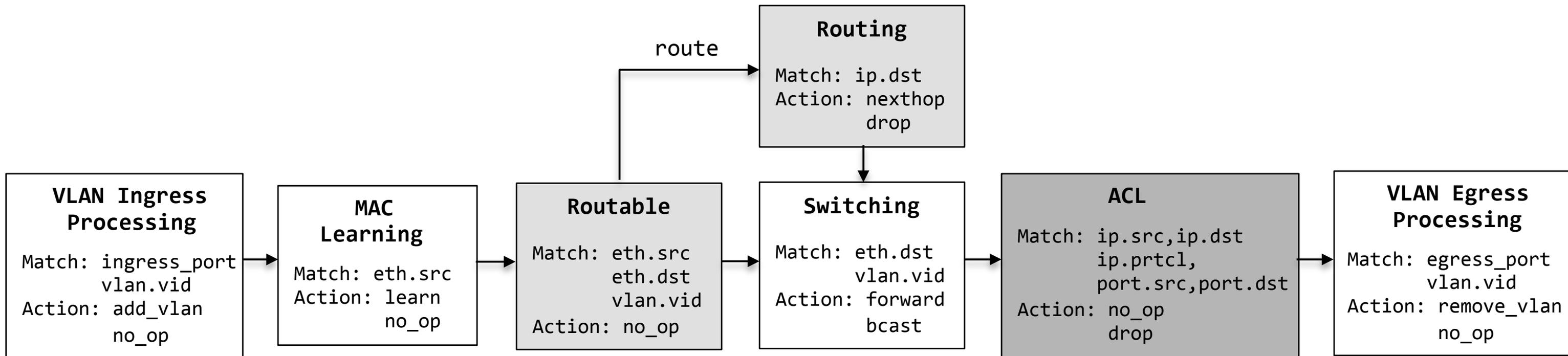
Approach: Specify forwarding behavior in P4; compile to modify OVS

Question: How does the PISCES switch performance compare to OVS?

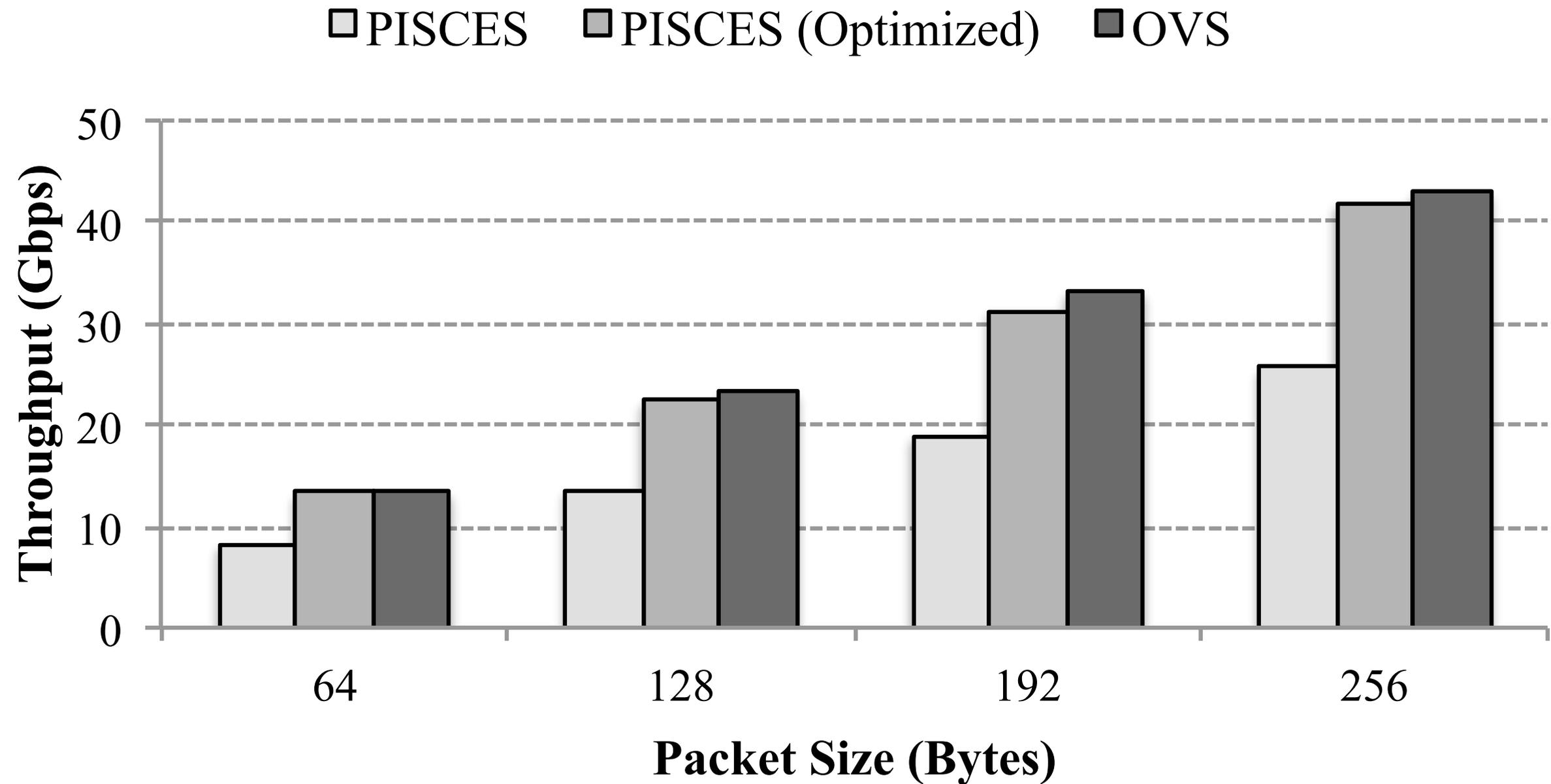
PISCES Architecture



Native OVS expressed in P4



PISCES vs Native OVS



Complexity Comparison

	LOC	Methods	Method Size
Native OVS	14,535	106	137.13
ovs.p4	341	40	8.53

40x reduction in LOC
20x reduction in method size

		Files Changed	Lines Changed
Connection Label	OVS	28	411
	ovs.p4	1	5
Tunnel OAM Flag	OVS	18	170
	ovs.p4	1	6
TCP Flags	OVS	20	370
	ovs.p4	1	4

Code mastery no longer needed

Next Steps

1. Make PISCES available as open-source (May 2016)
2. Accumulate experience, measure reduction in deployment time
3. Develop P4-to-eBPF compiler for kernel forwarding

PERC: Proactive Explicit Rate Control

Lavanya Jose, Stephen Ibanez, Mohammad Alizadeh, George Varghese, Sachin Katti, NM

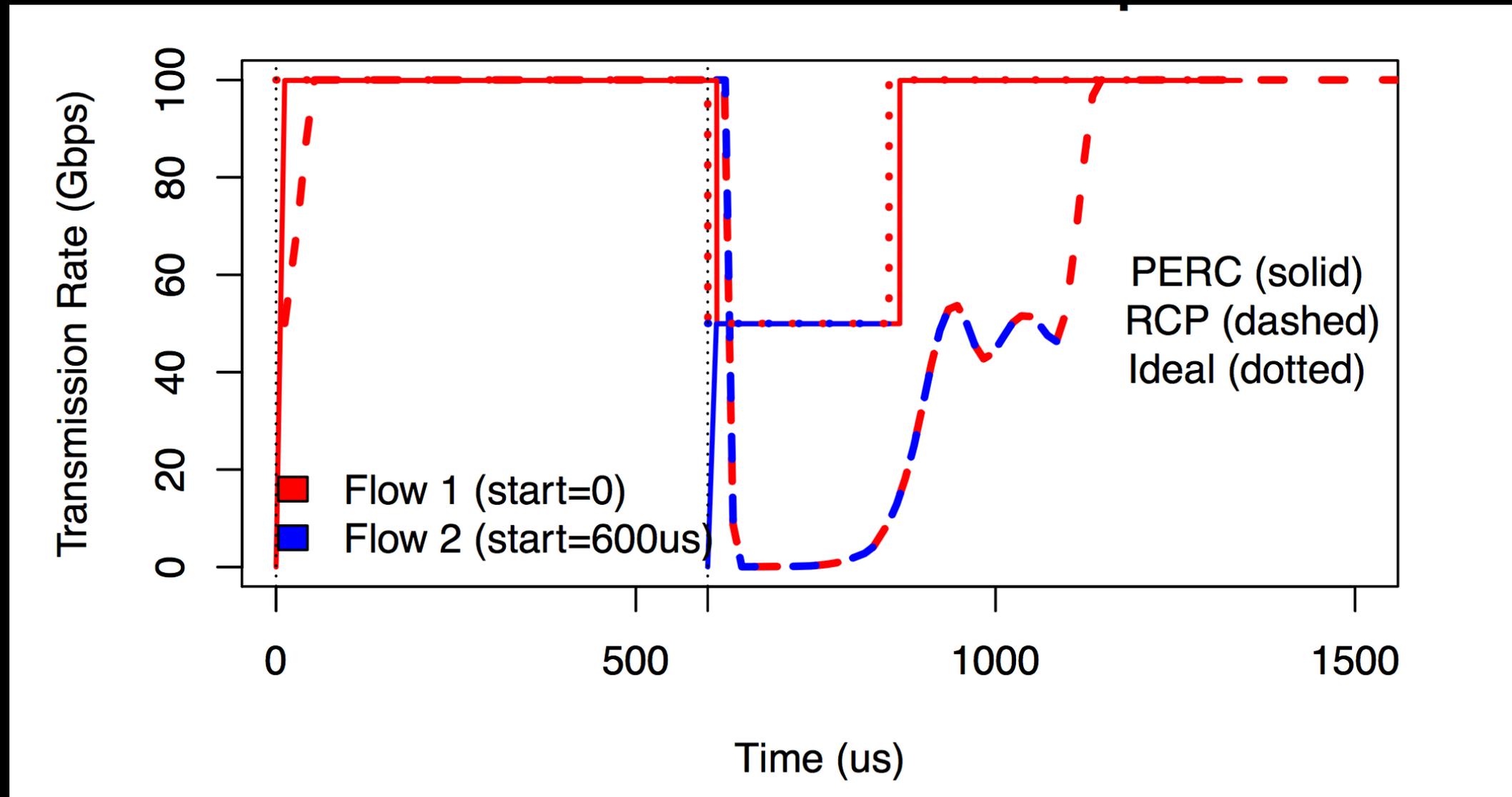
Problem: Congestion control algorithms in DCs are “reactive”

- Typically takes 100 RTTs to converge to fair-share rates (e.g. TCP, RCP, DCTCP)
- The algorithm it doesn't know the answer; it uses successive approximation

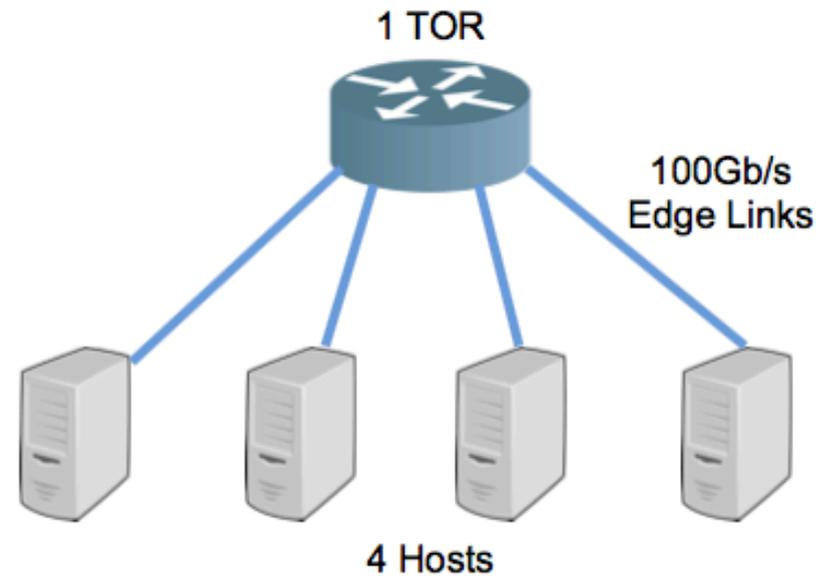
Approach: Explicitly calculate the fair-share rates in the forwarding plane

Question: Does it converge much faster? Is it practical?

Reactive vs Proactive Algorithms



Performance Results



	RCP	PERC
Median	14 RTTs	4 RTTs
Tail (99 th)	71 RTTs	10 RTTs

Convergence time
determined by
dependency chain

Next Steps

Convergence time

- Proof that convergence time equals length of dependency chain
- Reduce measured time to provable minimum

Develop practical algorithm

- Resilient to imperfect and lost update information
- Calculated in PISA-style forwarding plane

<The End>