

# Rethinking the SDN Abstraction

Chengchen Hu

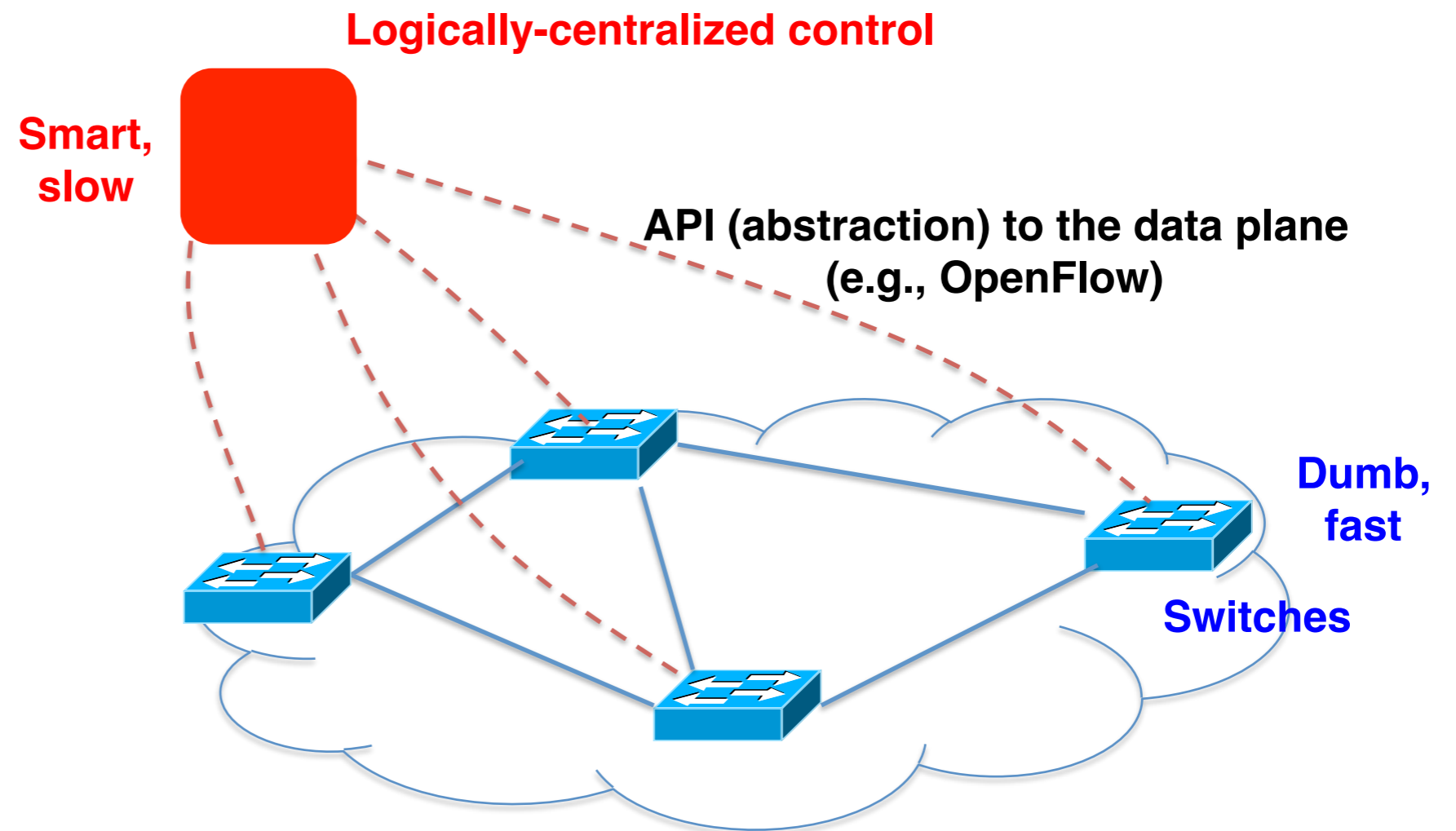
XJTU

Oct. 15, 2016

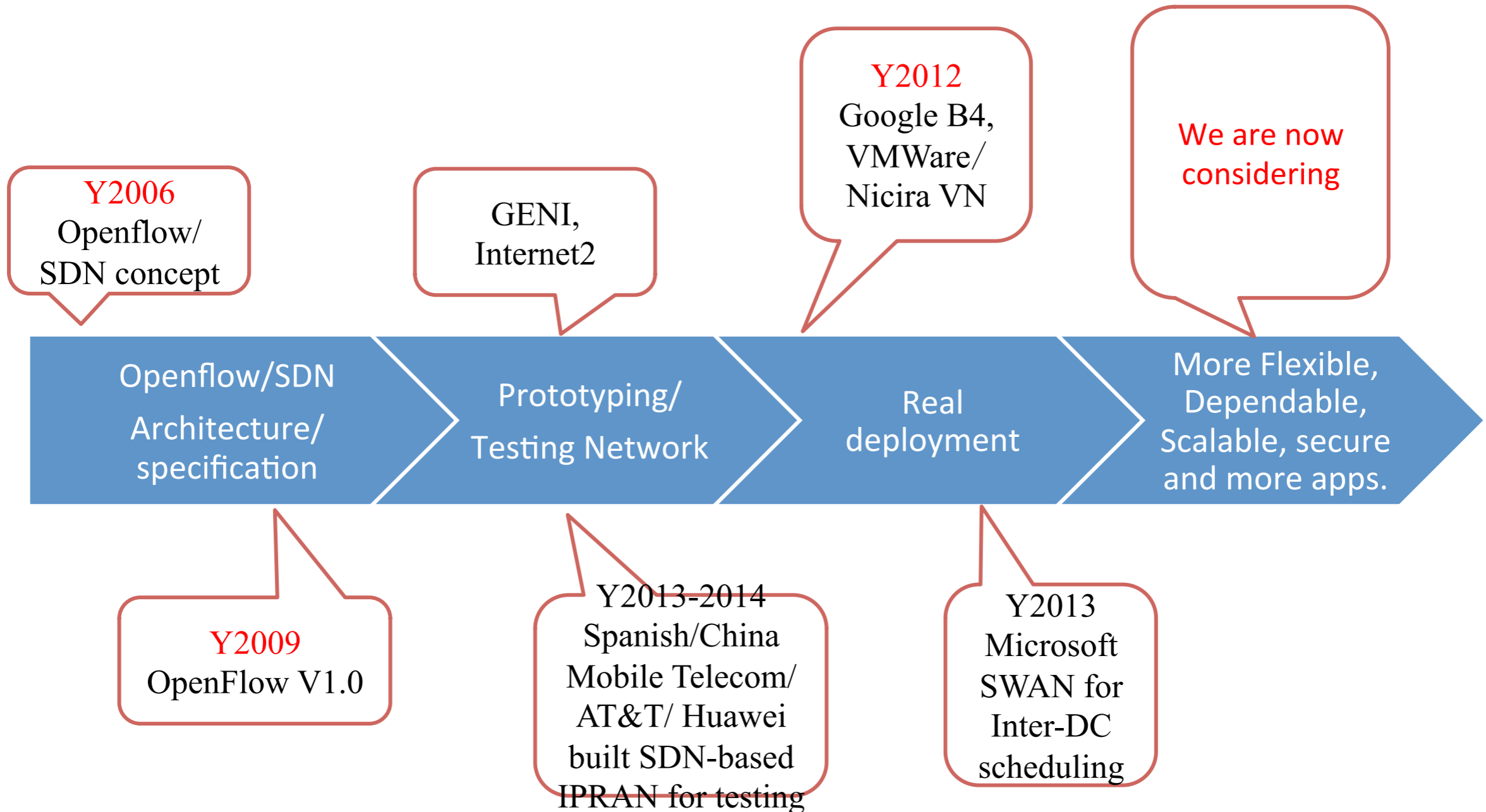


西安交通大学  
XIAN JIAOTONG UNIVERSITY

- Separate Control and Data
- Abstraction
- Global view



# History





OPEN NETWORKING  
FOUNDATION

# OF-PI: A Protocol Independent Layer

Version 1.1

September 5, 2014

ONF TR-505

# POF: Protocol Oblivious Forwarding

- Table search keys are Match defined as {offset, length} tuples
- Instructions/Actions access packet data or metadata using {offset, length} tuples
- Include other math, logic, move, branching, and jump instructions

Match ~40 matching header fields defined yet still **many** uncovered protocols/headers

Action

OFFPAT\_COPY\_TTL\_OUT  
OFFPAT\_COPY\_TTL\_IN  
OFFPAT\_SET\_MPLS\_TTL  
OFFPAT\_DEC\_MPLS\_TTL  
OFFPAT\_PUSH\_VLAN  
OFFPAT\_POP\_VLAN  
OFFPAT\_PUSH\_MPLS  
OFFPAT\_POP\_MPLS  
OFFPAT\_SET\_NW\_TTL  
OFFPAT\_DEC\_NW\_TTL  
OFFPAT\_PUSH\_PBB  
OFFPAT\_POP\_PBB  
**and on and on and on ...**



{offset, length} covers **any** frame based formats



POFAT\_SET\_FIELD  
POFAT\_ADD\_FIELD  
POFAT\_DELETE\_FIELD  
POFAT\_MOD\_FIELD  
**Period.**

Current OpenFlow

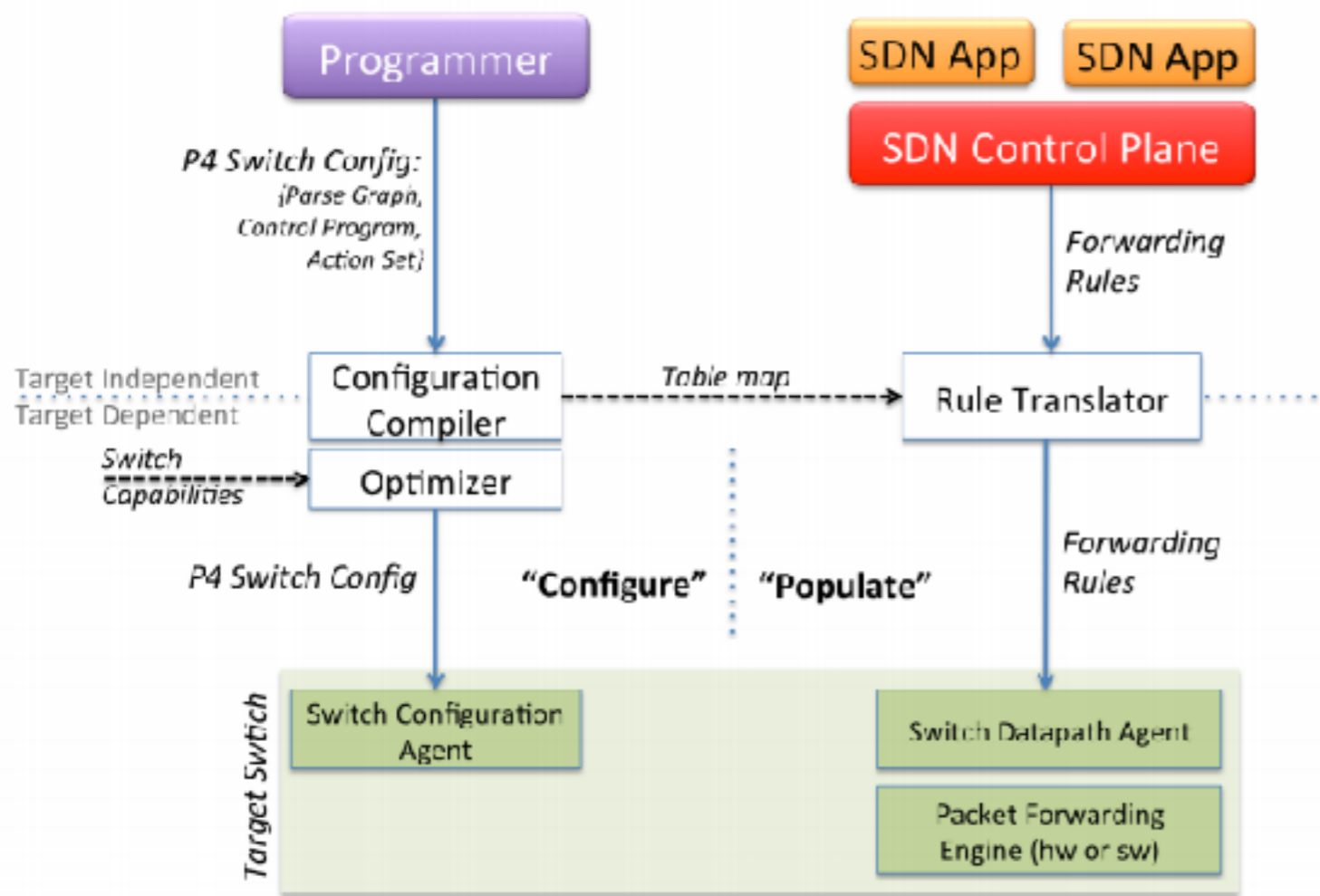
POF

**Packet field parsing and handling are abstracted as generic instructions to enable flexible and future proof forwarding elements. This is simple yet has profound implications to SDN.**

# P4: Programming Protocol-Independent Packet Processors

Proposed by

- Nick McKeown
- Jennifer Rexford
- Amin Vahdat
- George Varghese

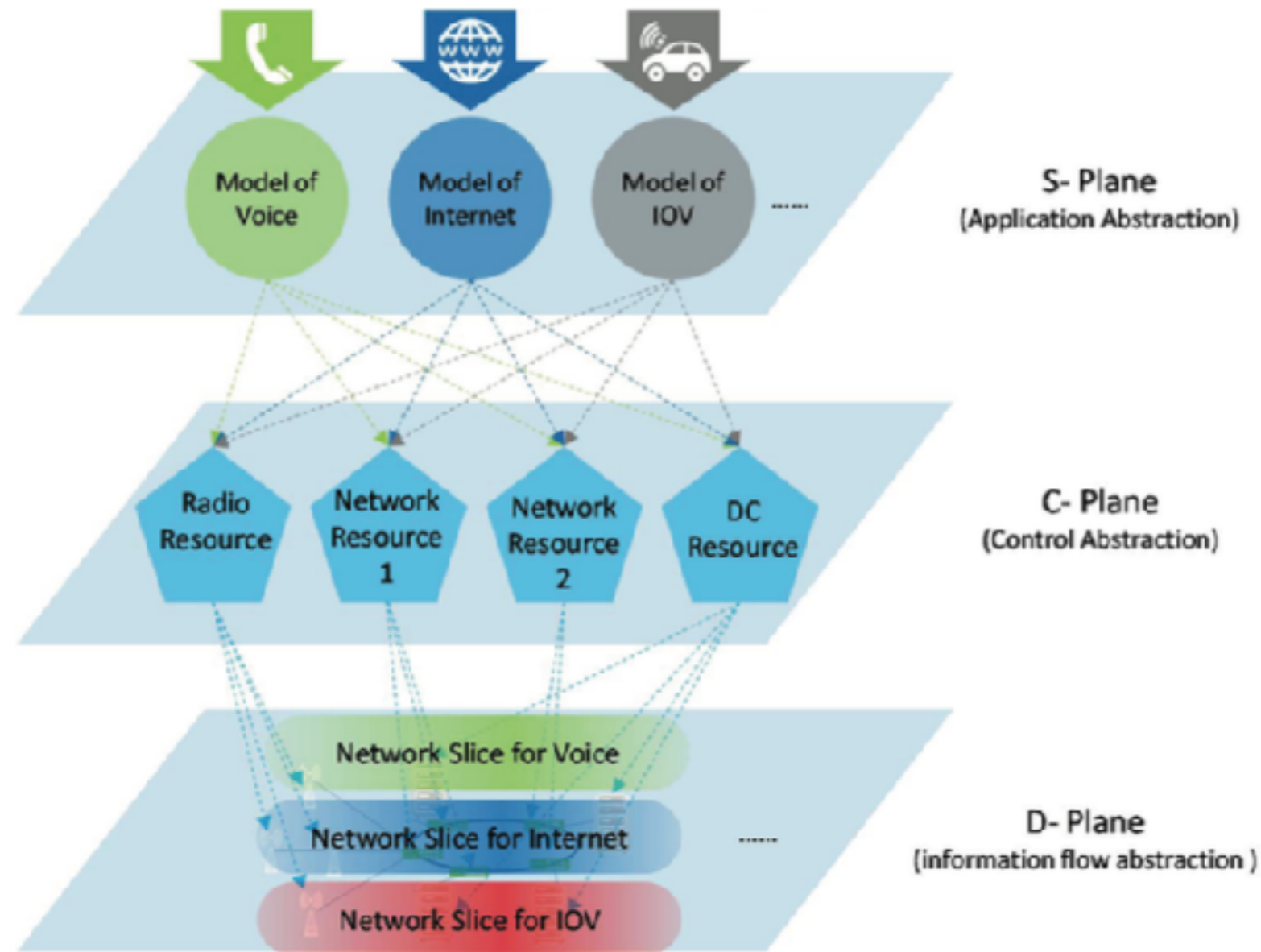


Goals

- Protocol independence
- Target independence
- Reconfigurability

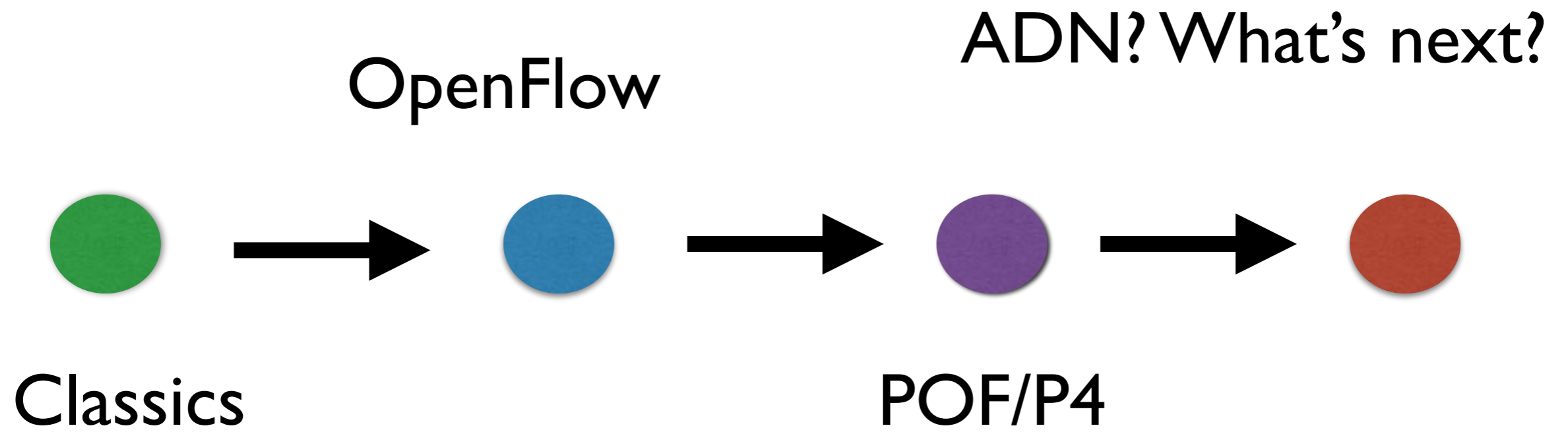
"P4: Programming Protocol-Independent Packet Processors," *ACM Sigcomm Computer Communications Review (CCR)*. Volume 44, Issue #3 (July, 2014)

# Application-Driven Network



Application Driven Network: providing On-Demand Services for Applications. Demo&Poster @ **SIGCOMM '16** (co-work with Huawei).

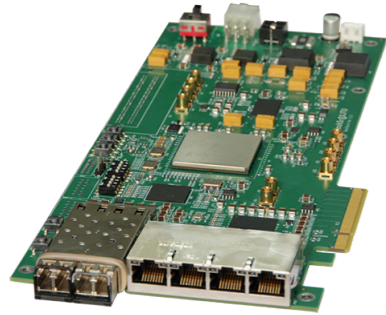




POF/P4 is flexible enough?



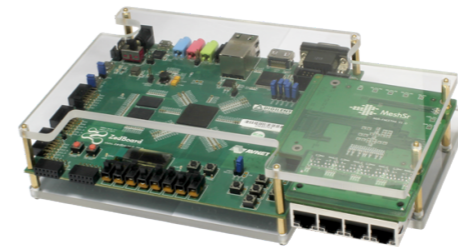
# Case 0: programmable SDN Data Plane



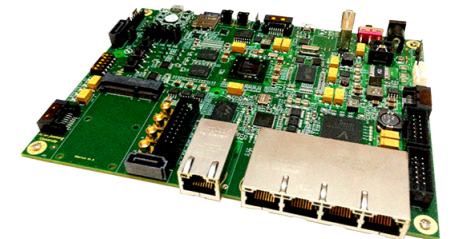
ONetCard  
2012 Aug  
PCIe Card



ONetSwitch 45  
4\*10G, 4\*GE, wifi  
2013 Aug



ONetSwitch 20  
4\*GE, with ZEDboard  
2013 Dec



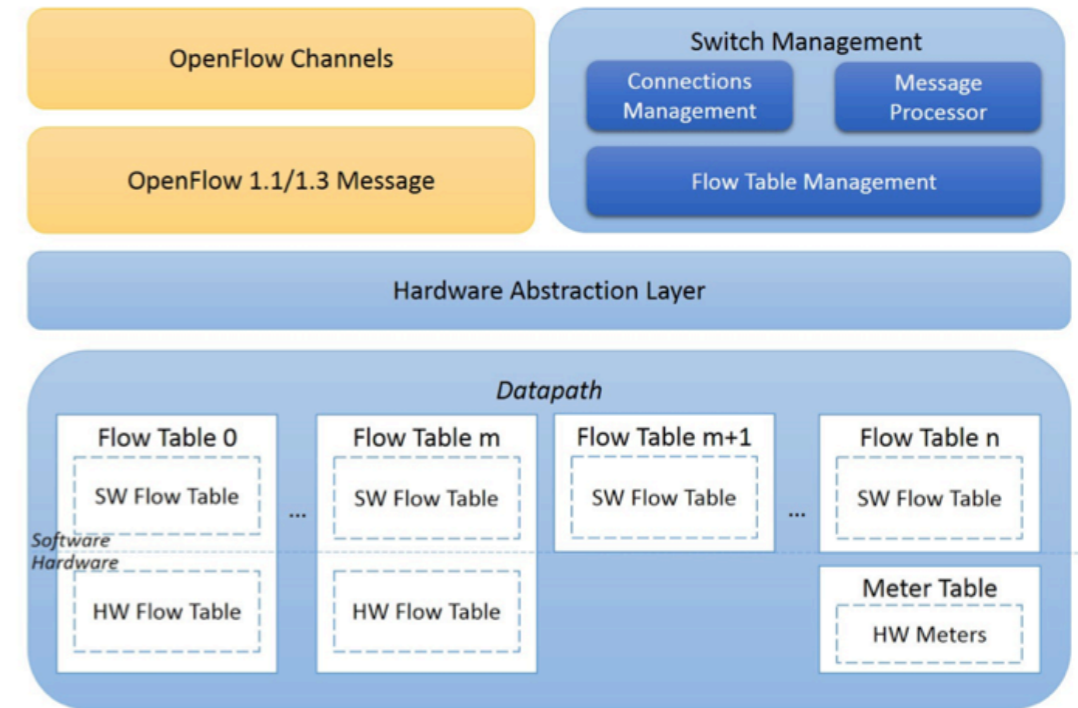
ONetSwitch 30  
wifi/storage, 5\*GE  
2014 Dec.

**Sponsored by Xilinx**

<http://onetswitch.org>

# ONetSwitch: All programmable SDN Switch

Chengchen Hu, Ji Yang, Hongbo Zhao, and Jiahua Lu.  
“Design of all programmable innovation platform for software defined networking”. **Open Networking Summit (ONS) 2014**, Santa Clara, CA, US, 2014



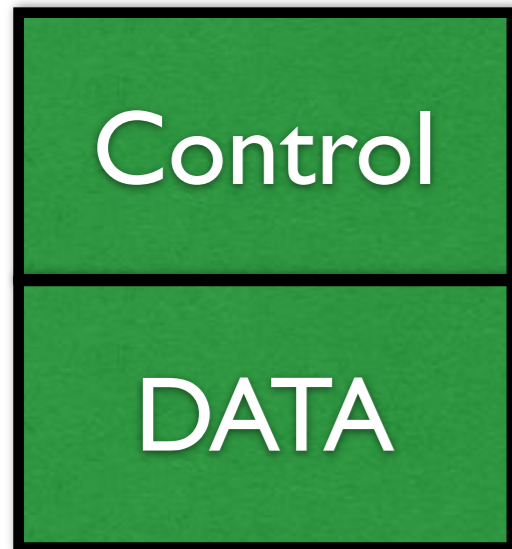
Chengchen Hu, Ji Yang, Zhimin Gong, Shuoling Deng, Hongbo Zhao. “DesktopDC: Setting All Programmable Data Center Networking Testbed on Desk”, **Poster&Demo at SIGCOMM 2014**, Chicago, IL, US, 2014

# Users



400+ ONetSwitches deployed in China, Europe, US

# Case I: Frequent protocols



## Case 1: ARP Gateway

Physical Gateway usually doesn't exist in SDN  
Controller reply to server end hosts  
Data plane queries controller repeatedly

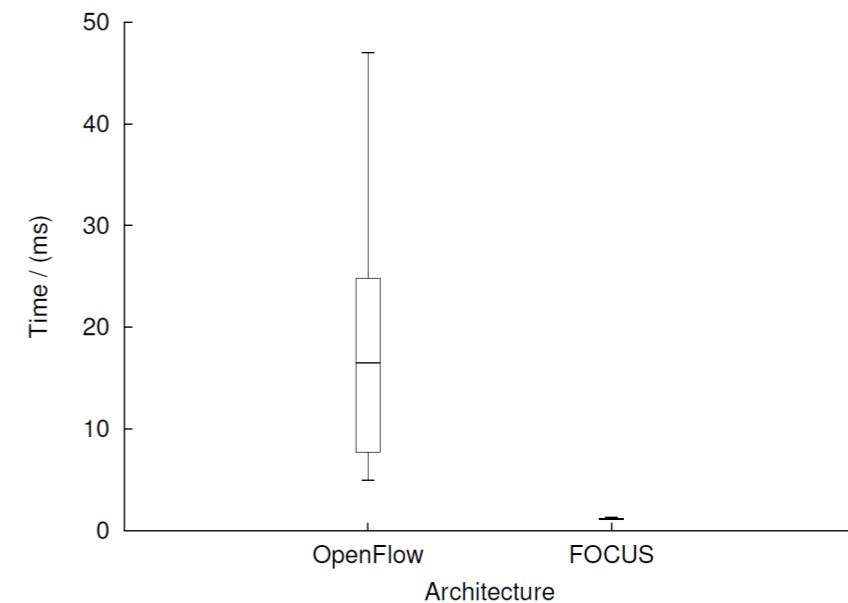
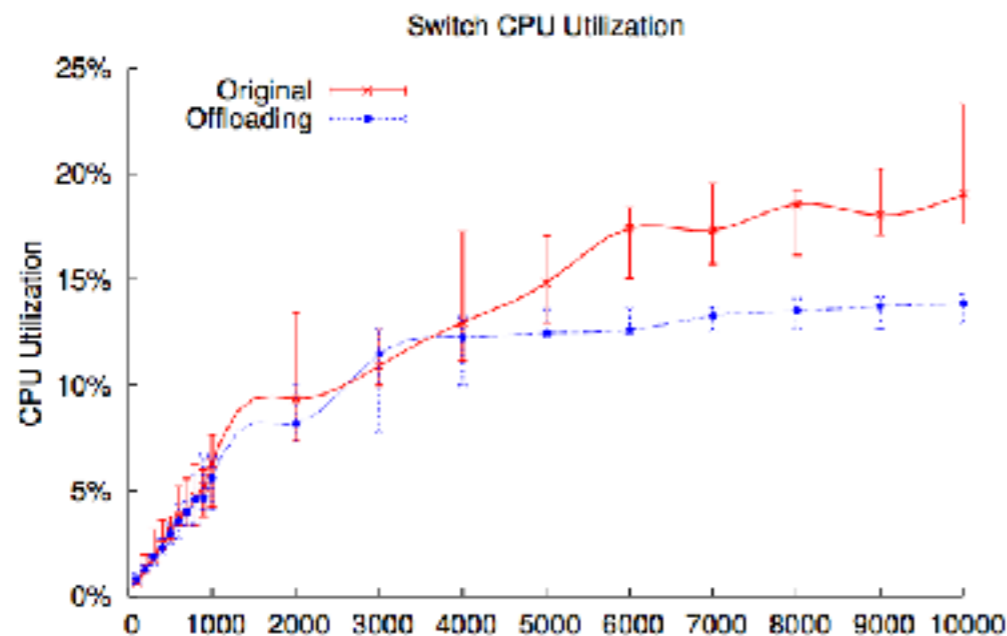
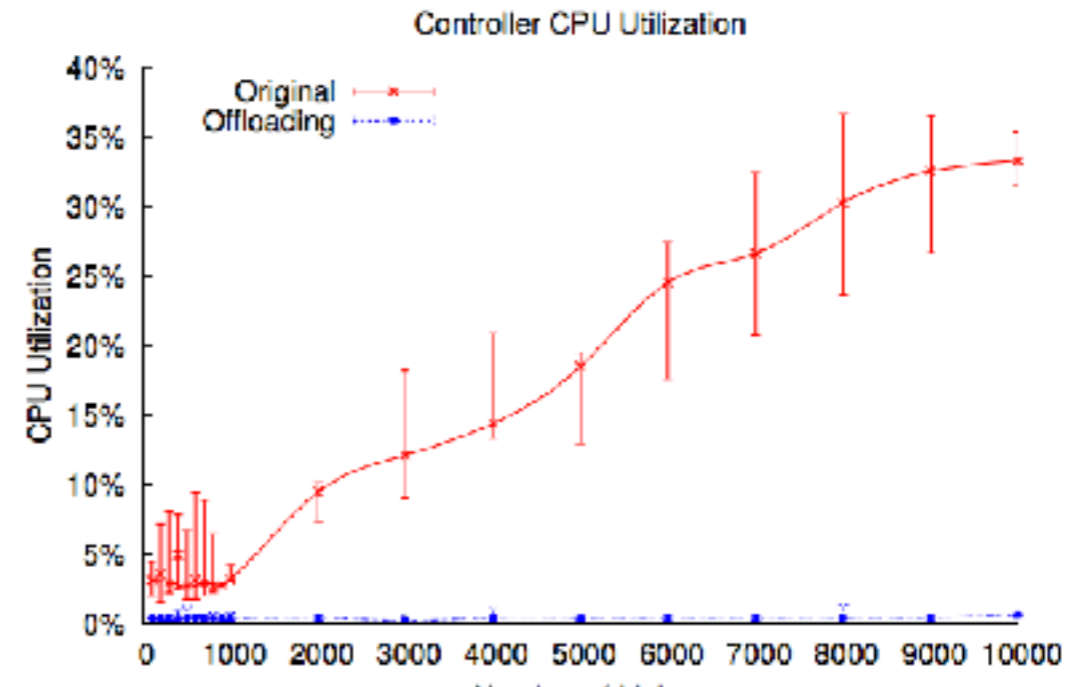
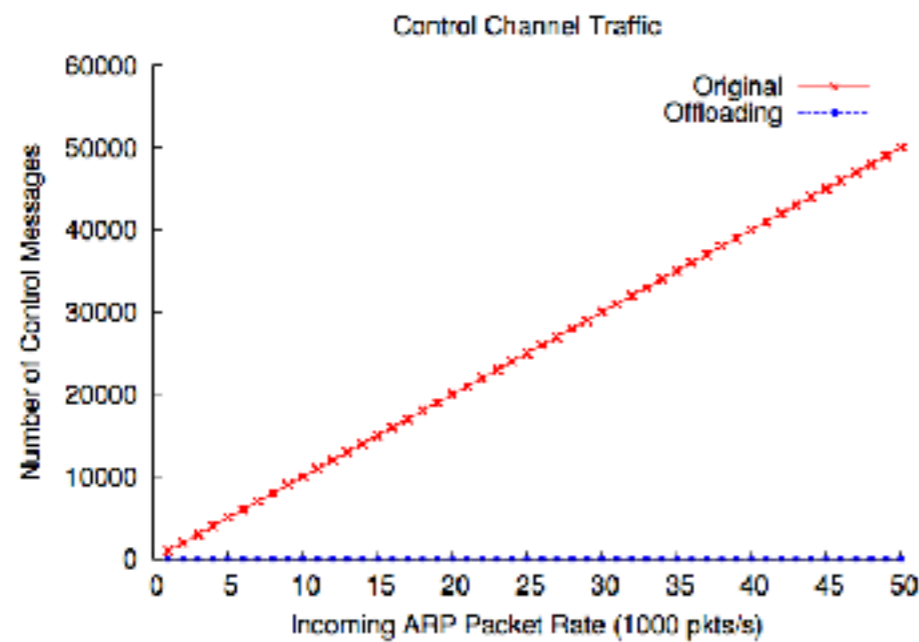


## Case 2: LLDP (Link Layer Discovery Protocol)

Get neighbor info. then keep watching  
Data plane packet goes to controller repeatedly even topology is stable  
Usually 2-5s per port per packet

## Case 3: LACP (Link aggregation Control Protocol)

State maintenance  
Repeat the same work again and again after topology is stable  
Usually 8 packets to converge  
Usually 1/30s per port per packet after convergence



Communication overhead reduced: 50%-100%

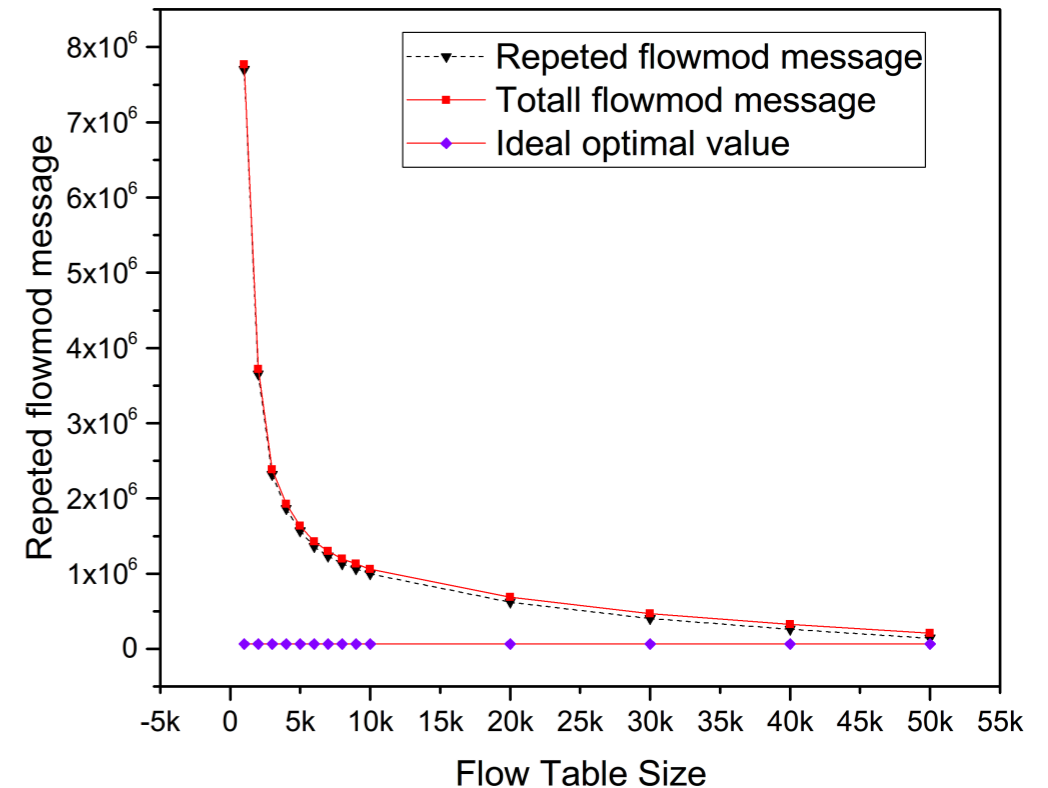
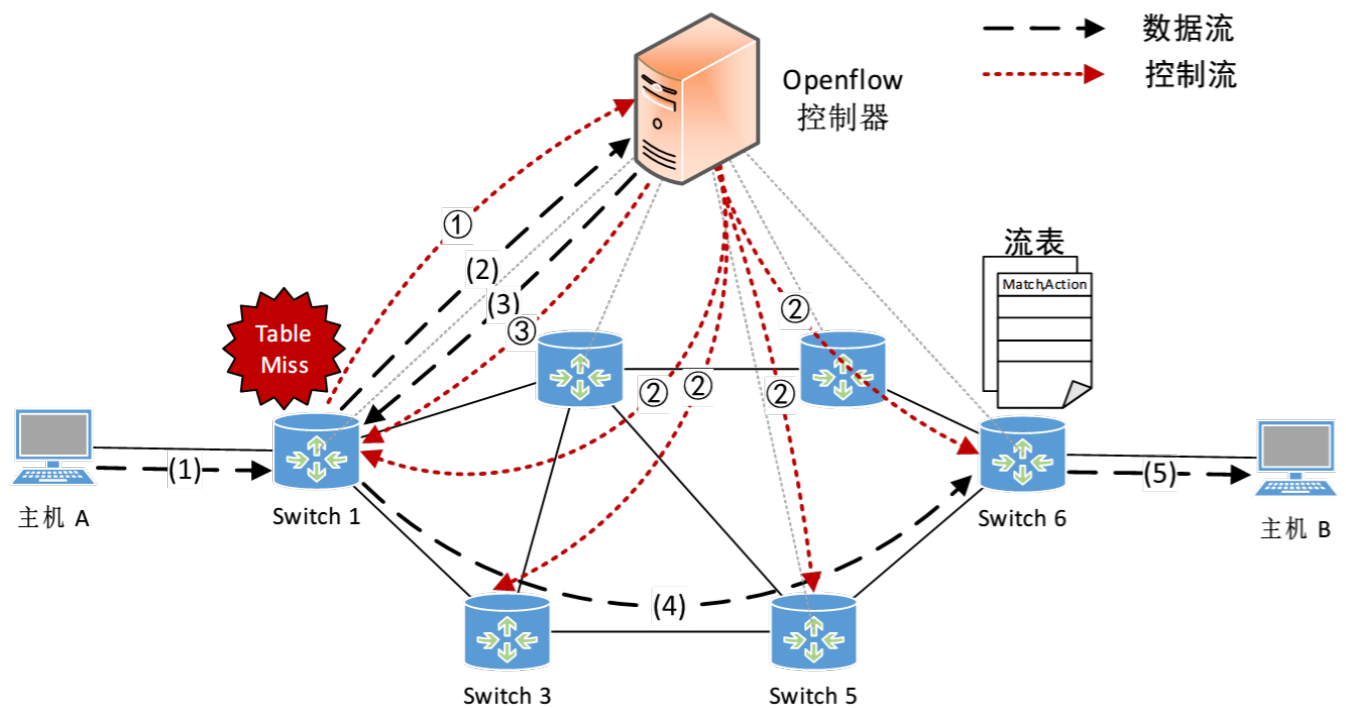
Controller CPU work load reduced: 80%-98%

ARP response time: from 10+ms to us

FOCUS: Function Offloading from a Controller to Utilize Switch Power.

**Poster @ NSDI'16 and @ SDN/NFV workshop'16**

# Case2: Table-miss



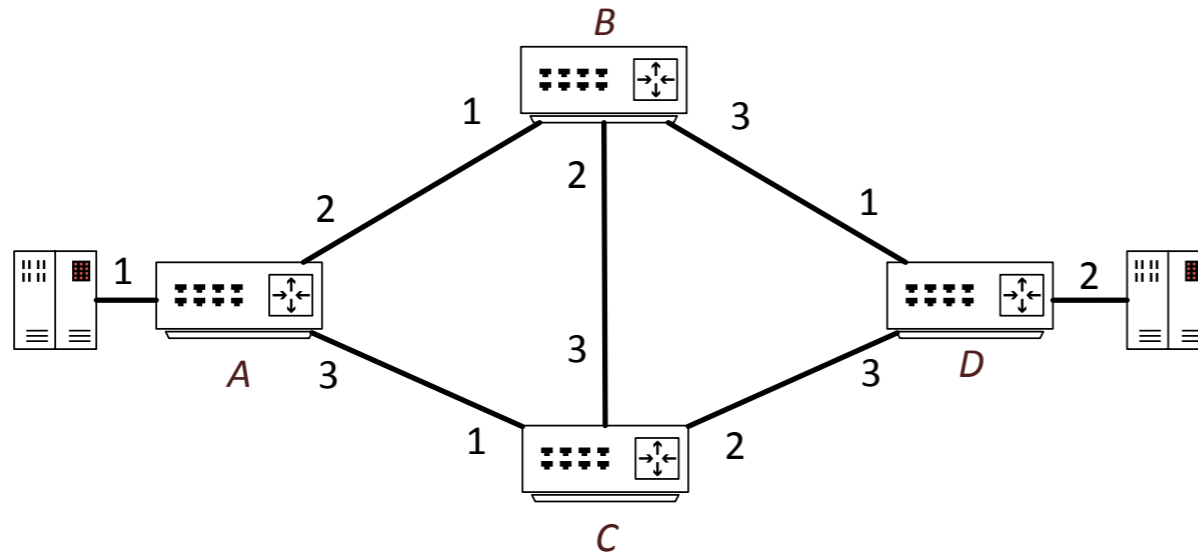
Repeated flowmod message: 68%-98% when flow table with 1k-50k entries.  
Eat either fast path memory or bandwidth between switch and controller  
Becomes bottleneck between Slow path and fast path in switch  
Small flows make the problem worse.

CoSwitch: A Cooperative Switching Design for Software Defined Data Center Networking  
**@HotData, 2014** ( best paper award)

Co-Work with IBM Research Lab

Taming the Flow Table Overflow in OpenFlow Switch. **Poster at SIGCOMM '16.**

# Case 3: Rules Conflicts



**A** DstIP 10.0.0.2, FORWARD 2  
**B** DstIP 10.0.0.2, FORWARD 3  
**D** DstIP 10.0.0.2, FORWARD 2

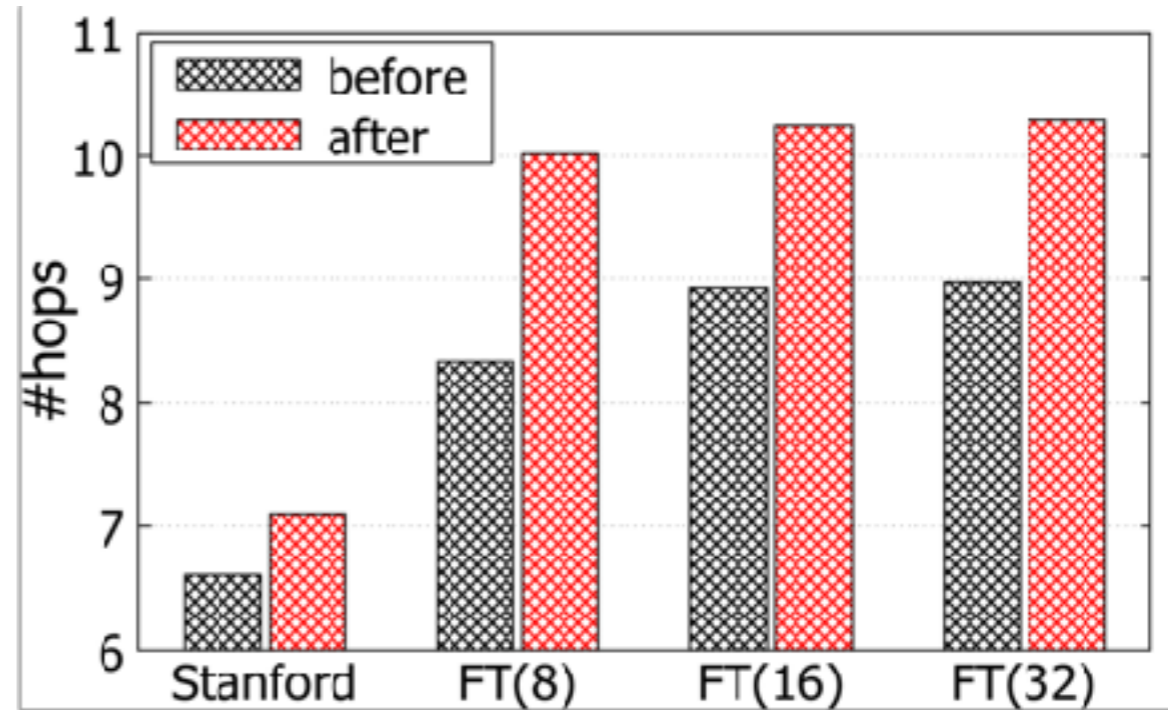
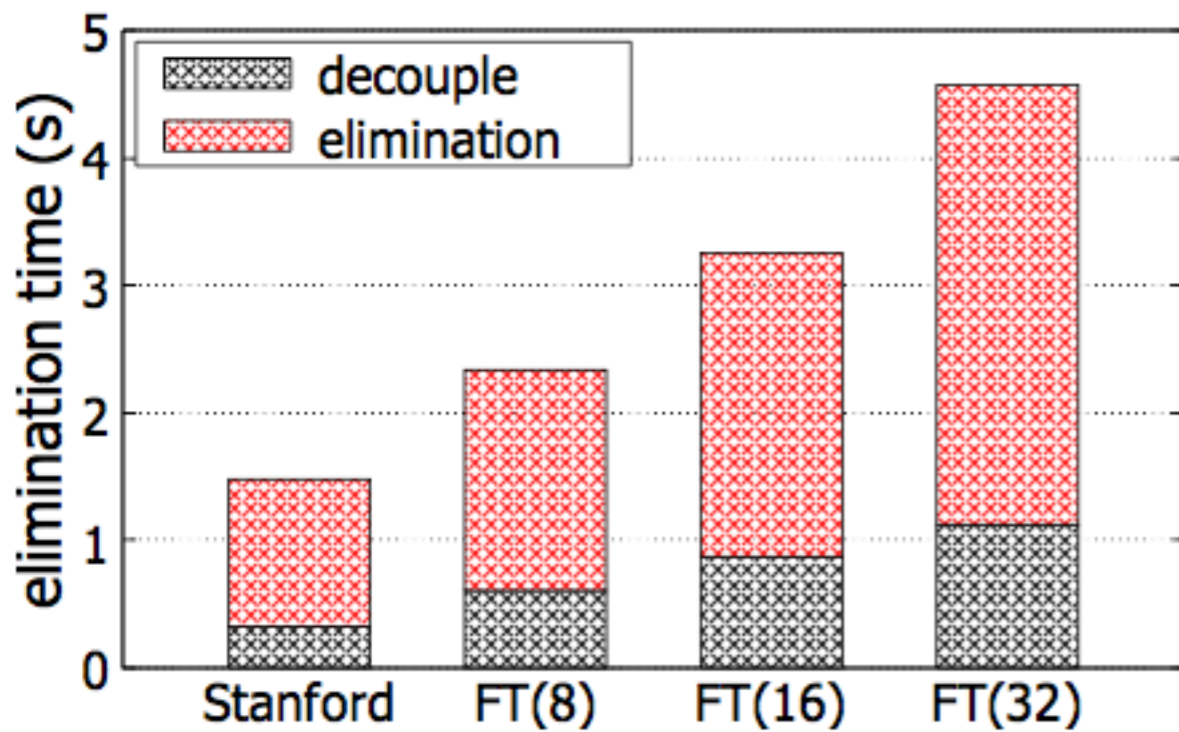
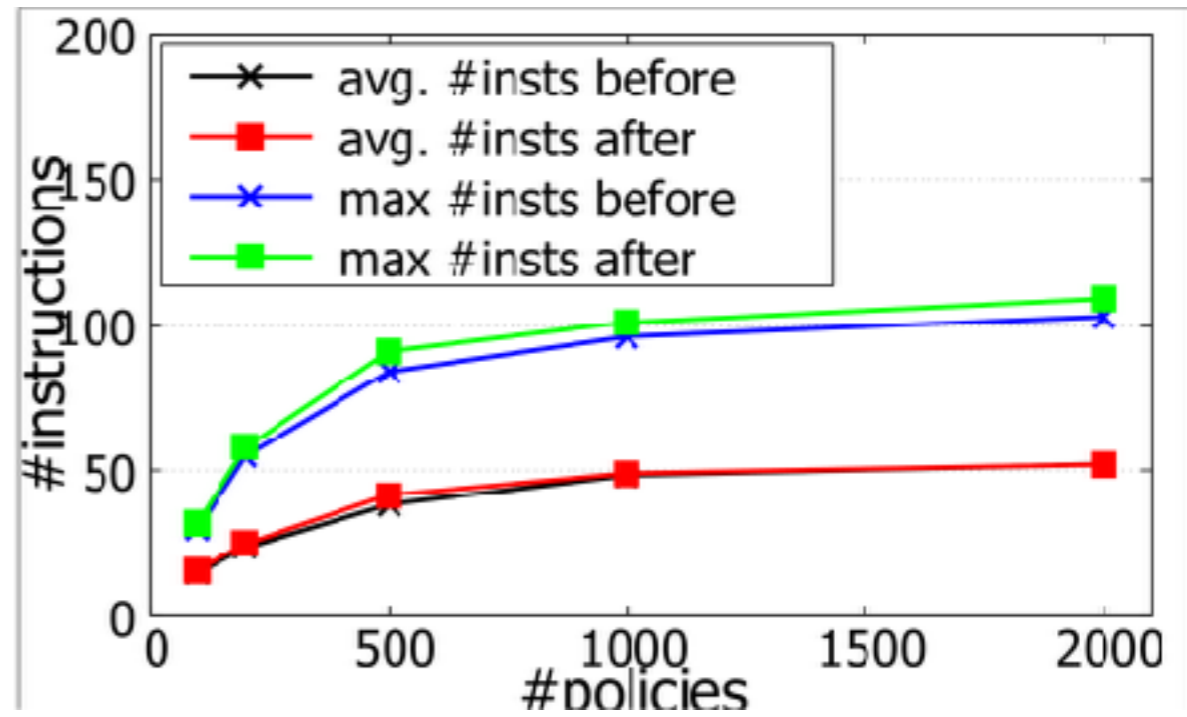
**A** DstIP 10.0.0.2, FORWARD 3  
**C** DstIP 10.0.0.2, FORWARD 2  
**D** DstIP 10.0.0.2, FORWARD 2

Hard to eliminate the conflicts  
without the high-level intents

- blue and orange policies are only to connect A and D → random overwriting
- blue and orange policies critically specify the current path → unsolvable
- blue/orange policy is to connect A and with B/C in path → A to B to C to D

Modular SDN Compiler Design with Intermediate Representation.  
**poster @ SIGCOMM '16.**

{A, DstIP 10.0.0.2}, FORWARD	towards B $\cap$ forbid A
{A, DstIP 10.0.0.2}, FORWARD	towards C $\cap$ forbid A
{B, DstIP 10.0.0.2}, FORWARD	towards D $\cap$ forbid B
{C, DstIP 10.0.0.2}, FORWARD	towards D $\cap$ forbid C
{D, DstIP 10.0.0.2}, FORWARD	fixed_forward 2
{D, DstIP 10.0.0.2}, FORWARD	fixed_forward 2





# Case 4: Rule update

- Flowtable update bottleneck
  - 10s to 100s of rule edits per second
  - Full refresh of 5K entries takes minutes

Pattern	Priority
<1, 2>	3
<*, 2>	2
<*, *>	1

Old





Pattern	Priority
<1, 2>	5
<2, *>	4
<1, *>	3
<*, 2>	3
<3, *>	2
<*, *>	1

New

Priority Updates

3 rule adds + 2 priority updates

-  Unmodified fields
-  Modified fields

# Try to minimize the update

	Predicate
1:	00*
2:	**0
3:	0*1
4:	**1
5:	***

(a) An example of rule update.

	Predicate	Prio
1:	00*	20
6:	0*0	17
2:	**0	15
3:	0*1	15
4:	**1	10
5:	***	5

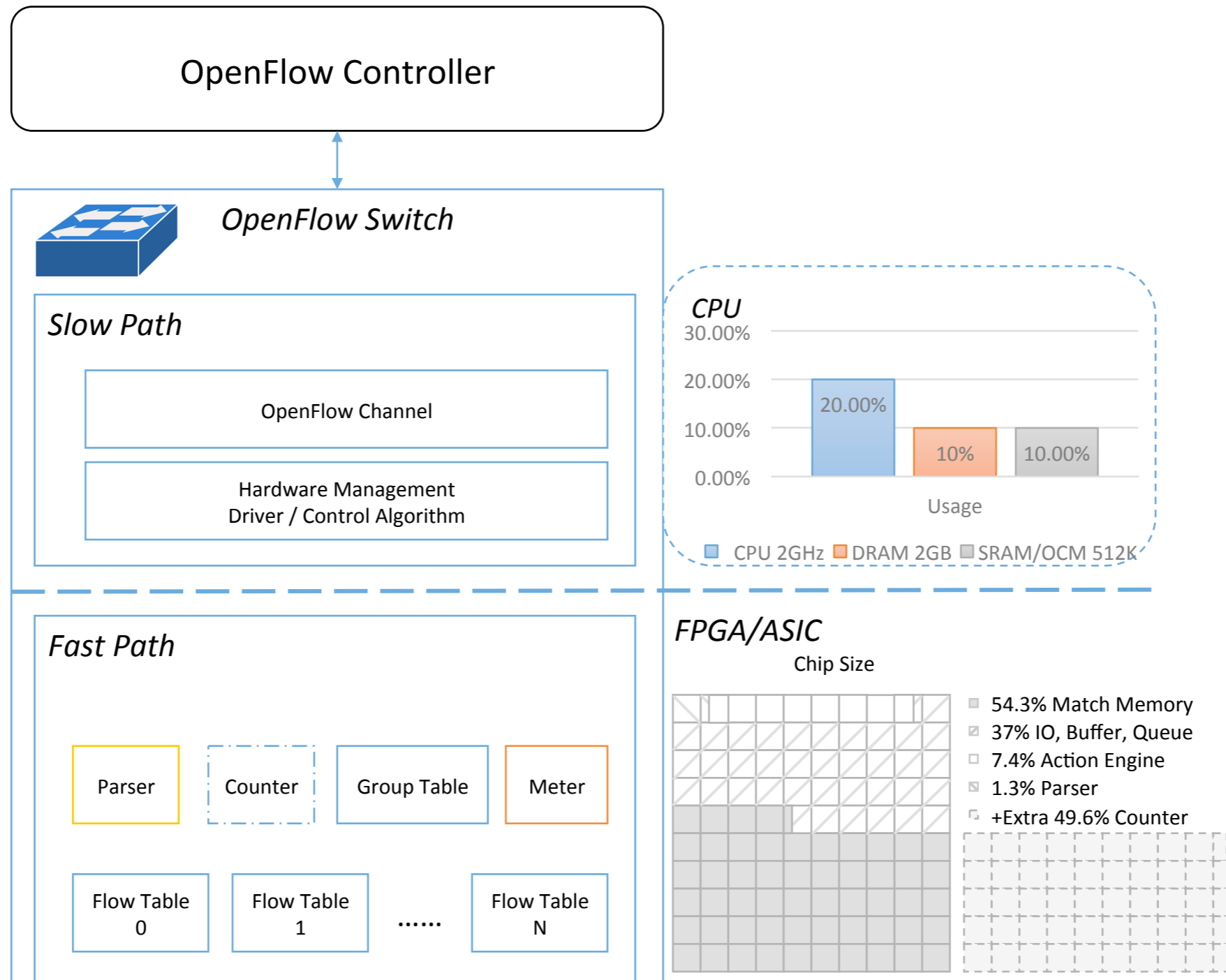
(b) Update with priority.

	Predicate	Dependency
1:	00*	1
6:	0*0	6
3:	0*1	3
4:	**1	4
2:	**0	2
5:	***	5

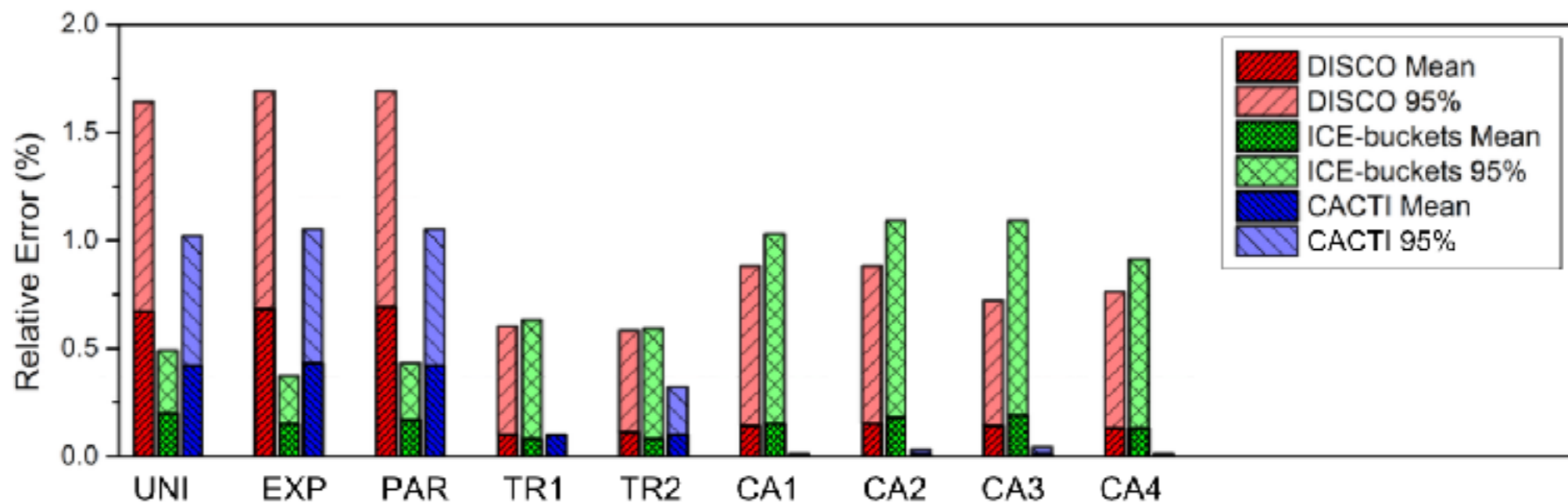
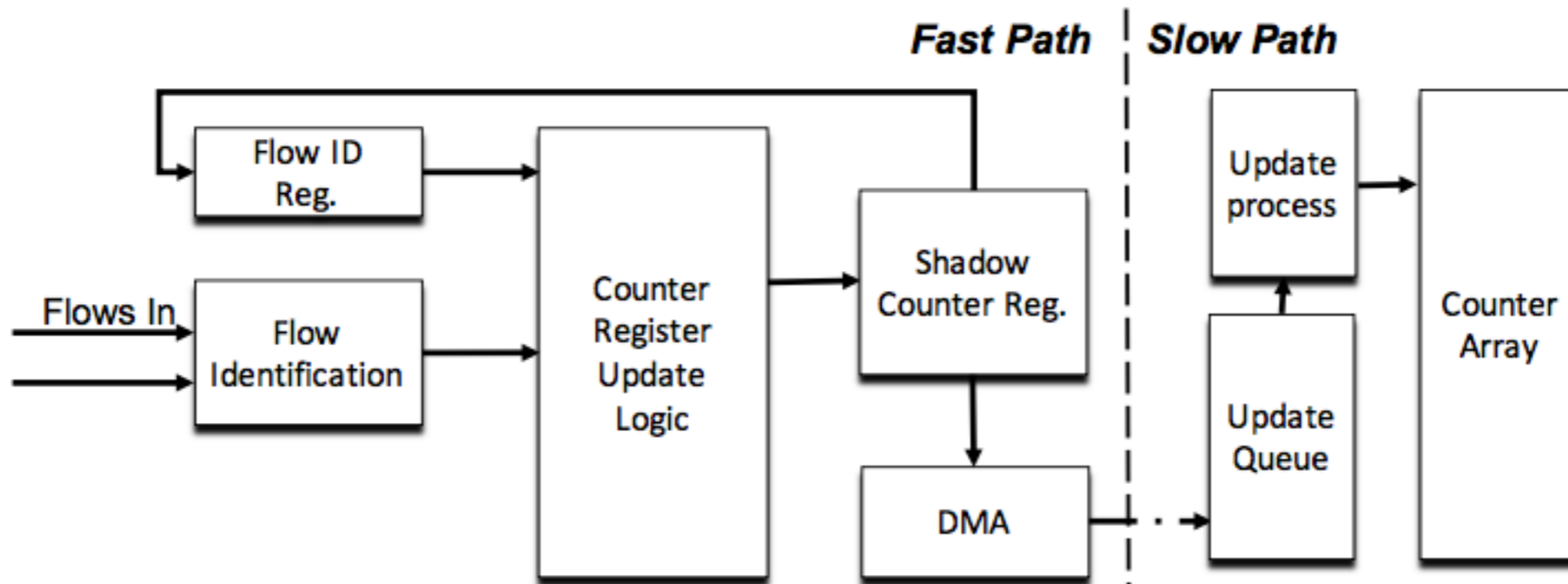
(c) Update with dependency graph.

Xitao Wen, Bo Yang, Yan Chen, Li Erran Li, Kai Bu, Peng Zheng, Yang Yang, Chengchen Hu, RuleTris: Minimizing Rule Update Latency for TCAM-based SDN Switches, **ICDCS 2016**, Nara, Japan, June 27 – June 30, 2016.

# Case 5: SDN Counters



# CACTI: CAche CounTing

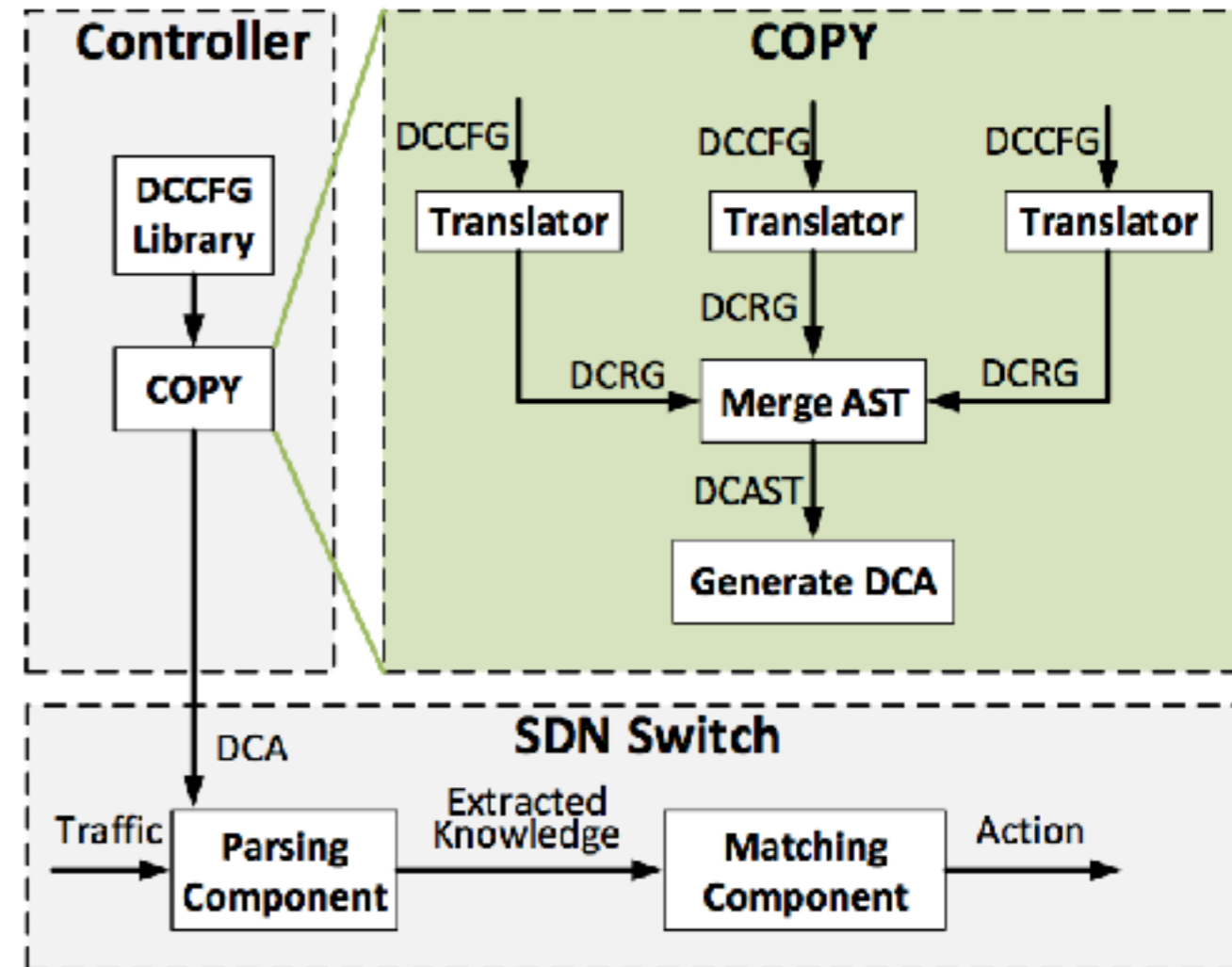


# Case 6: Abstract L7

We propose a new specification form DCCFG, which can be formulated as a six-tuple,  $\Gamma = (\mathbb{N}, \Sigma, \mathbb{C}, \mathbb{R}, S, \mathbb{E})$ , where  $\mathbb{N}$ ,  $\Sigma$ ,  $\mathbb{C}$ ,  $\mathbb{R}$ ,  $S$ ,  $\mathbb{E}$  are the finite set of non-terminals, terminals, counters, production rules, start non-terminal, and extraction tokens, respectively. The non-terminals are the symbols where the terminals can be derived. The terminals can be a single character or a RegEx. The production rules can be described as  $\langle \text{guard} \rangle : \langle \text{non-terminal} \rangle \rightarrow \langle \text{body} \rangle \langle \text{action} \rangle$ .

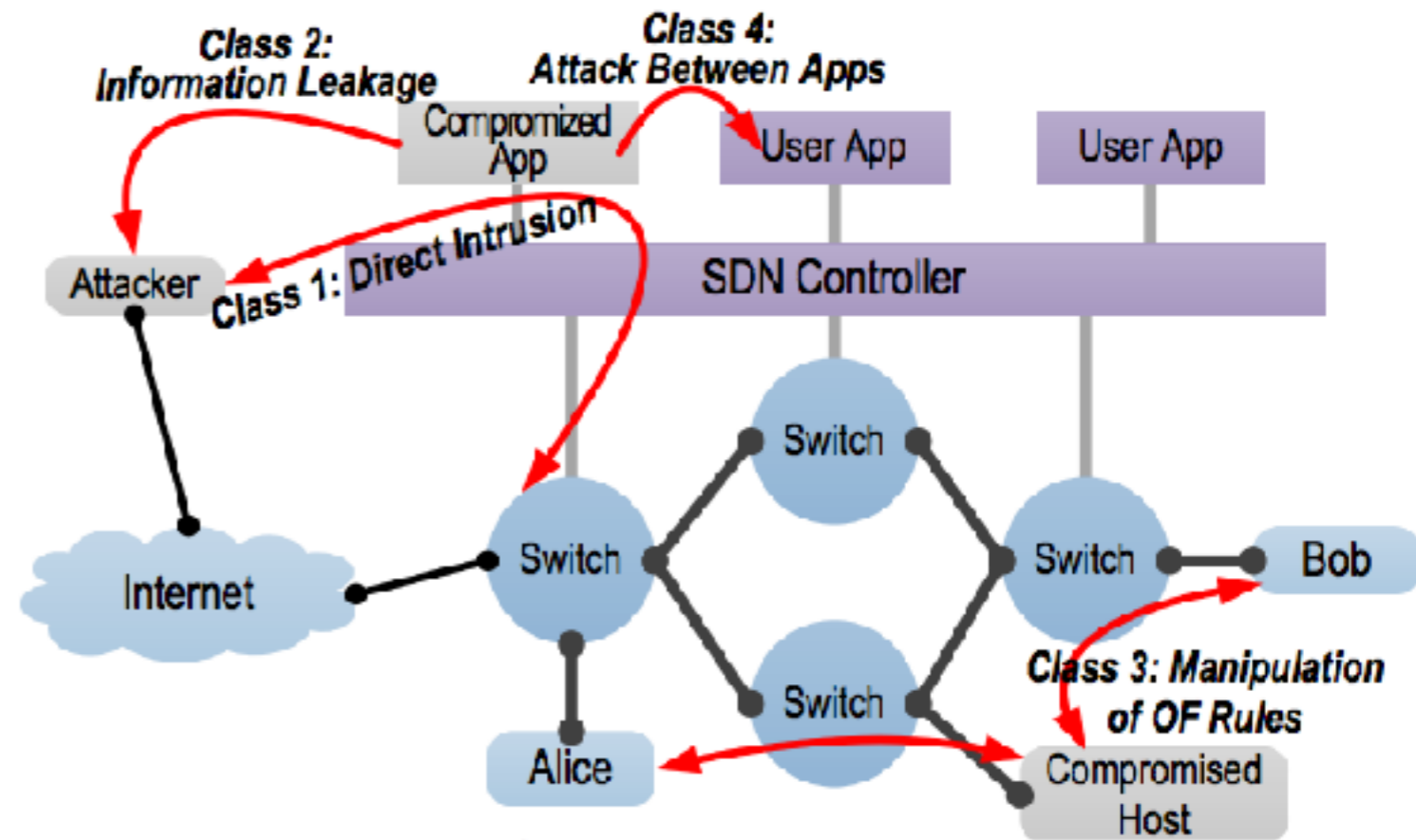
1		S	→	T L V
2		T	→	$\langle \text{key}, 0, "[A-Z]" \rangle$
3		L	→	$"[0-9]" \text{ [len=getnum()]}$
4		$[\text{len}>0]V$	→	$\langle \text{value}, 0, "[A-Za-z]*" \rangle \text{ [len=reducc()]}$
5		$[\text{len}=0]V$	→	$\epsilon$

Figure 3: A TLV specification  $\Gamma$  in DCCFG.



# Case 7: Vulnerabilities

- Flow table attack
- Secure Channel attack
- Session Hijacking
- Policy bypass
- ....



On Denial of Service Attacks in Software Defined Networks @ **IEEE Network**  
SDNShield: Reconciling Configurable Application Permissions for SDN App Markets @ **DSN 2016**

Mind the Gap: Monitoring the Control-Data Plane Consistency in Software Defined Networks @ **CoNext 2016**

# Open Questions

- How to make data plane programmable ?
- Flow-action abstraction?
- Forwarding—Switch, Control—Controller?
- Fully Centralization?
- Anything between high level intents and low level rules?
- How to co-design Fast Path & Slow Path in Switch

**Thank you**