

# Building Compilers for Reconfigurable Switches

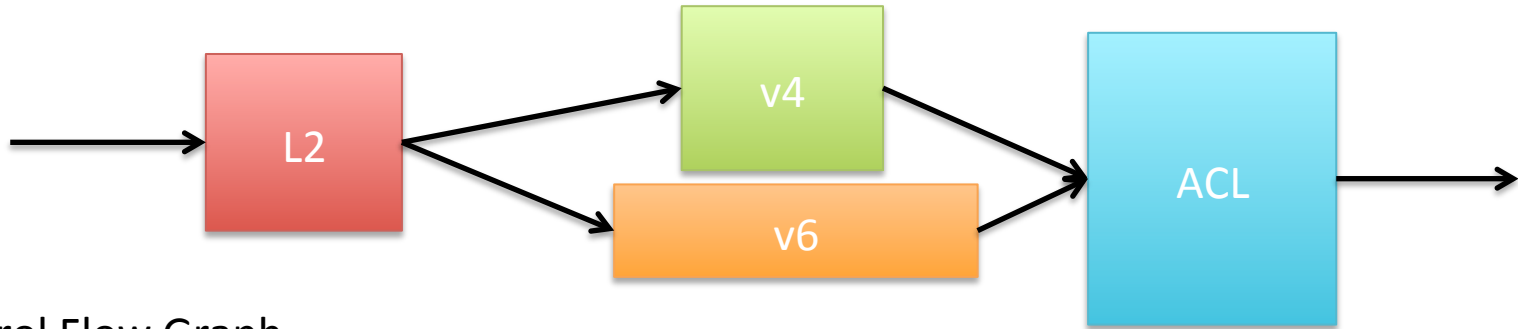
Lavanya Jose, *Lisa Yan*,  
Nick McKeown, and George Varghese

Research funded by AT&T, Intel, Open Networking Research Center.

# In the next 20 minutes

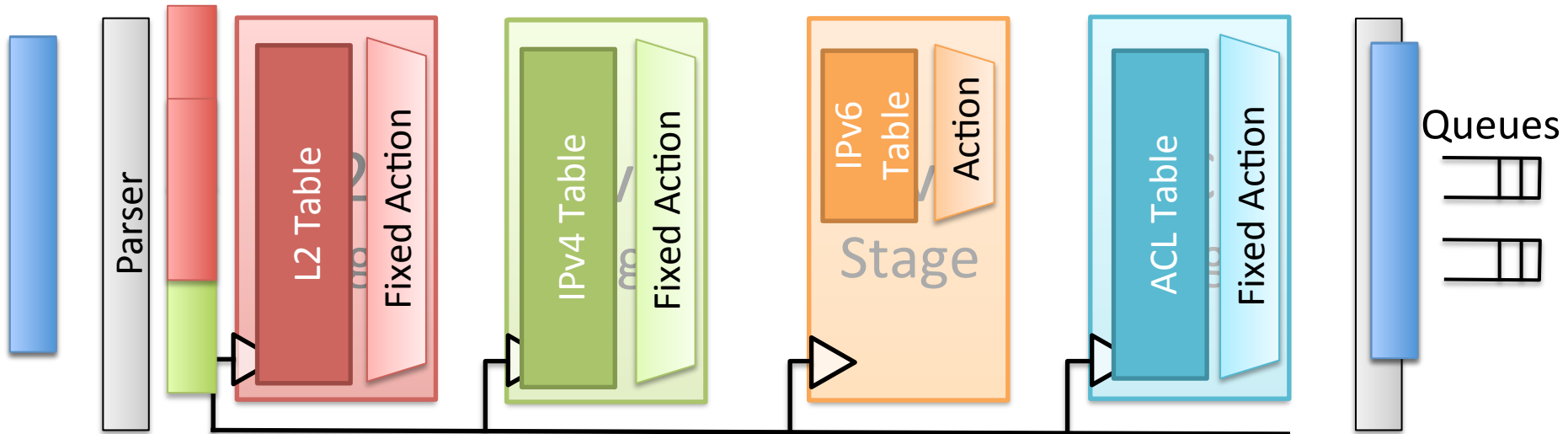
- Reconfigurable switch chips provide functionality and efficiency
- We will program them using languages like P4
- We need a compiler to compile P4 programs to reconfigurable switch chips.

# Fixed-Function Switch Chips



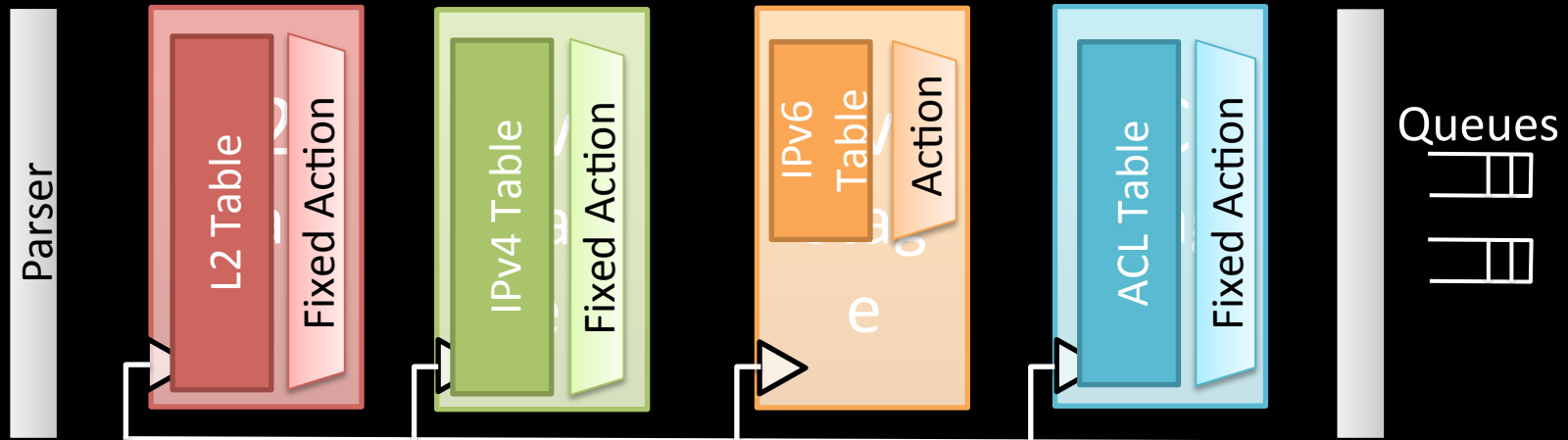
Control Flow Graph

Switch Pipeline

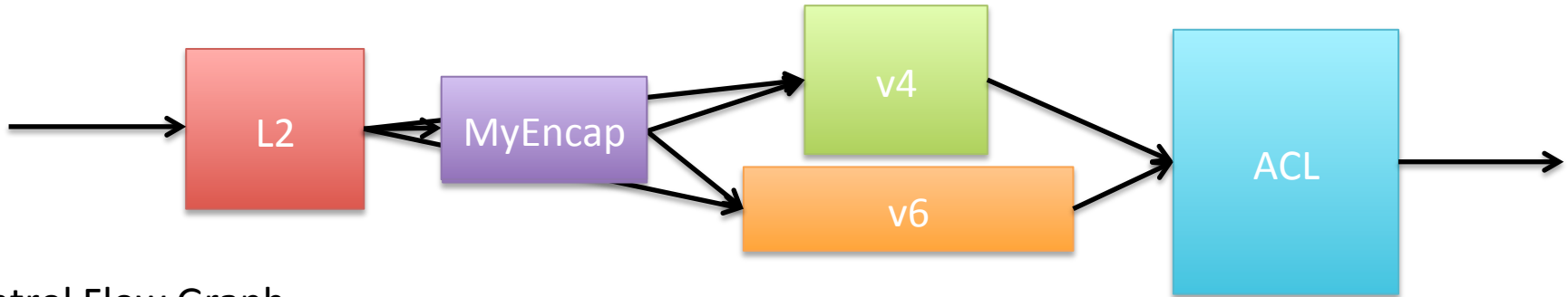


# Fixed-Function Switch Chips Are Limited

1. Can't add new forwarding functionality
2. Can't add new monitoring functionality
3. Can't move resources between functions

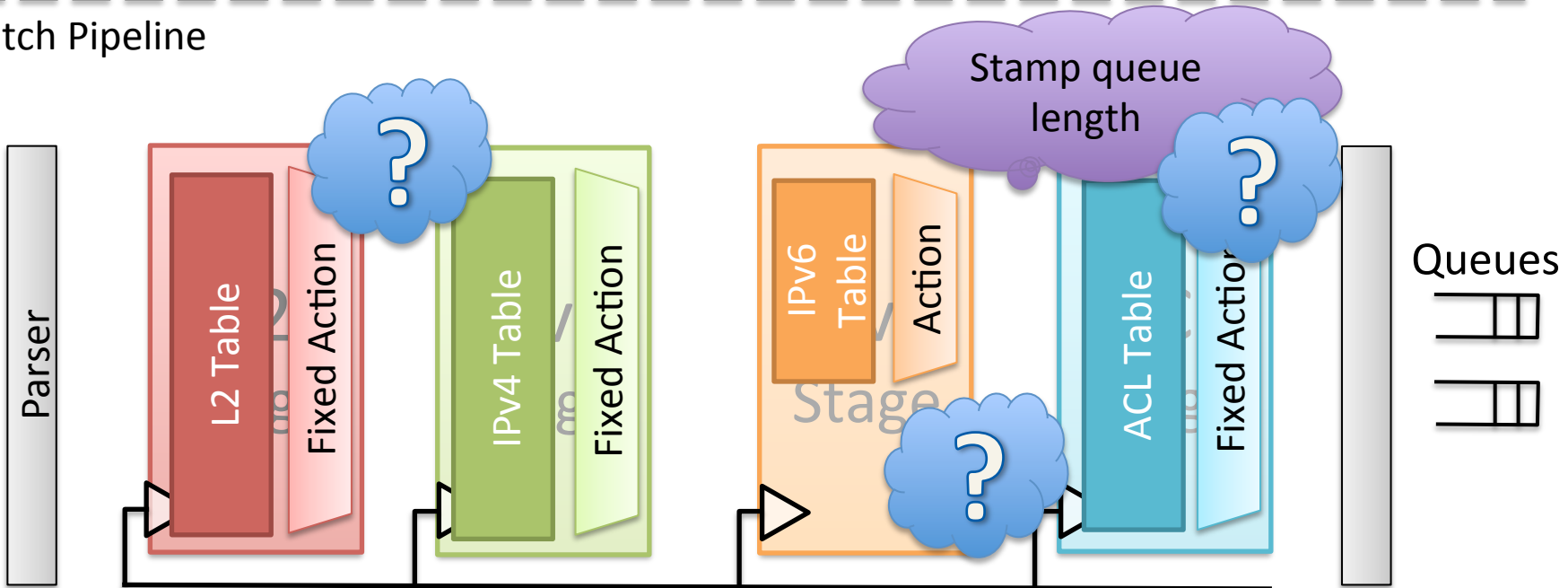


# Fixed-Function Switch Chips

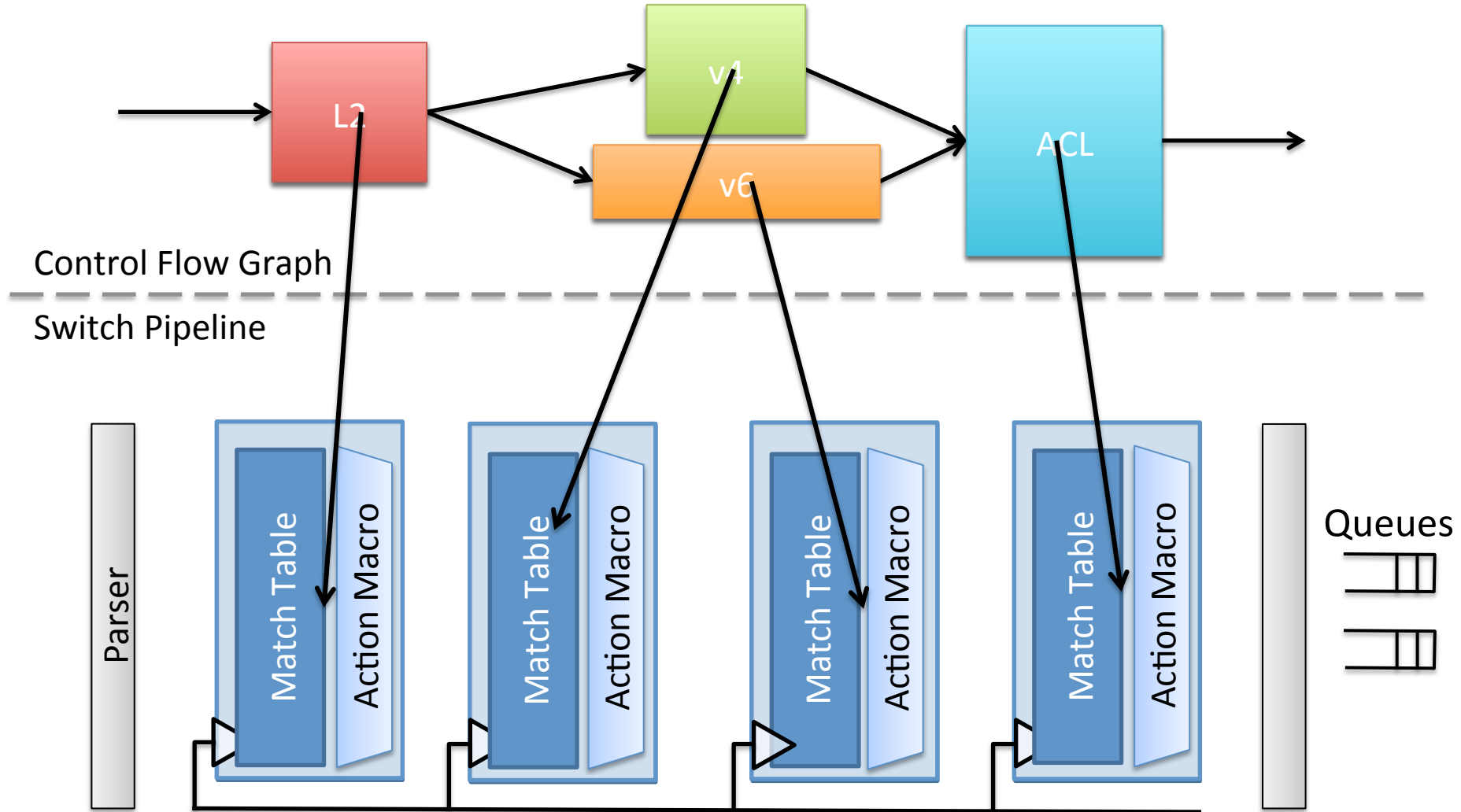


Control Flow Graph

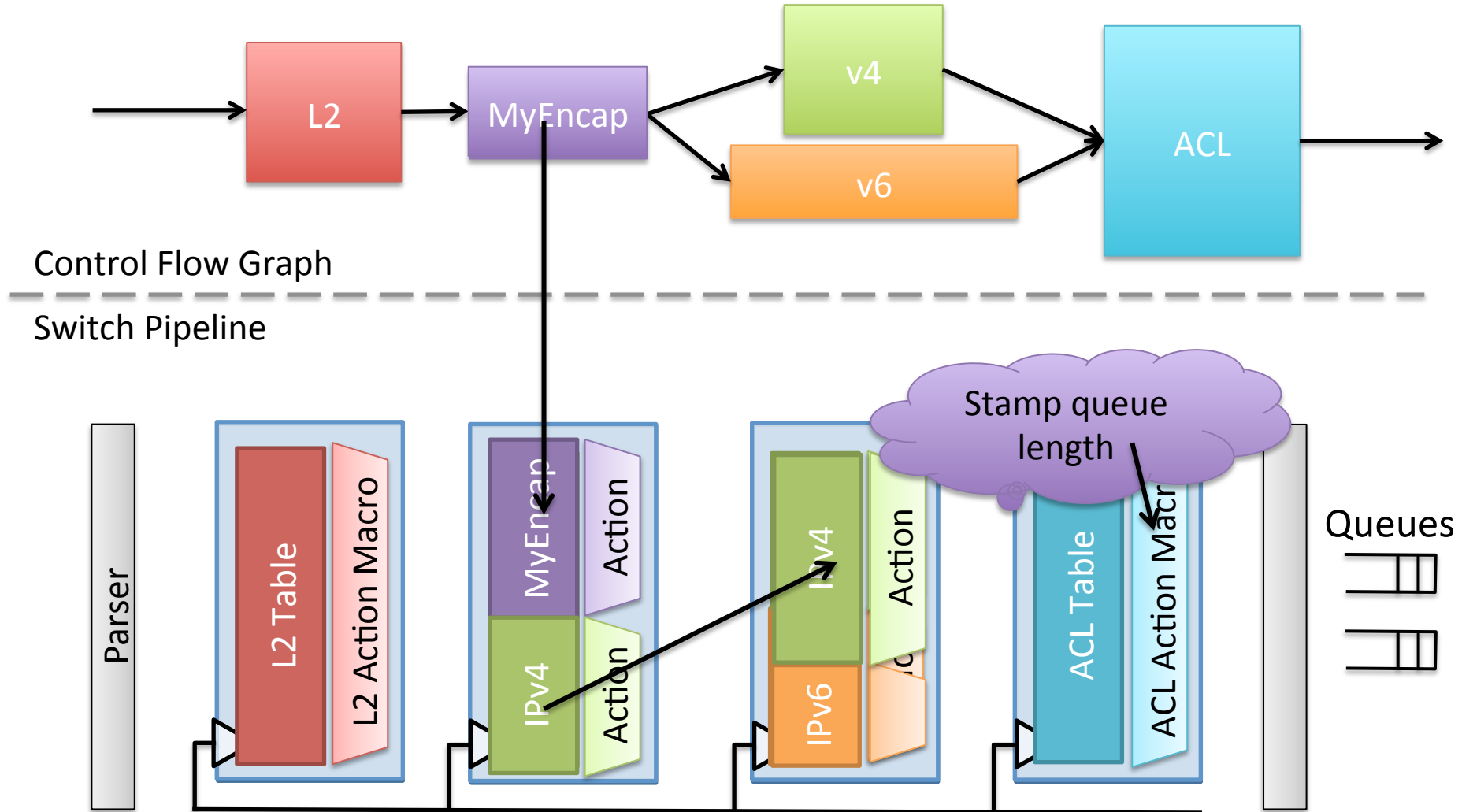
Switch Pipeline

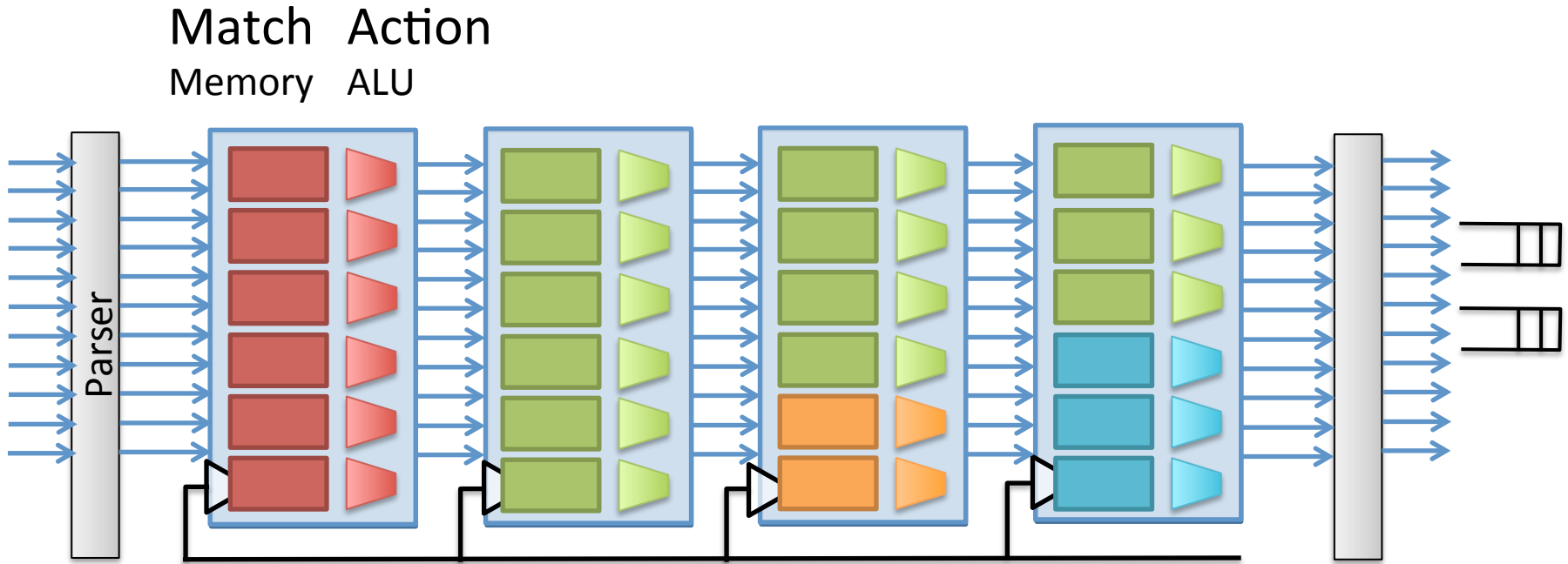


# Reconfigurable Switch Chips



# Reconfigurable Switch Chips





Protocol Independent Switch Architecture (PISA)



# Match + Action Processor: pipelined and in-parallel

Press

## Fixed Function Broadcom Tomahawk: 3.2 Tbps Reconfigurable Cavium Xpliant: 3.2 Tbps

Broadcom Delivers Industry's First High-Density 25/100 Gigabit Ethernet Switch for Cloud-Scale Networks

Customers, New StrataXGS® Tomahawk™ Series Delivers 3.2 Tbps  
SDN Control and Visibility Features



PRODUCTS SUPPORT NEWS & EVENTS SALES SOFTWARE PARTNERS CORPORATE INVESTORS

> News & Events

News&Events

Press Releases

Cavium Networks in the News

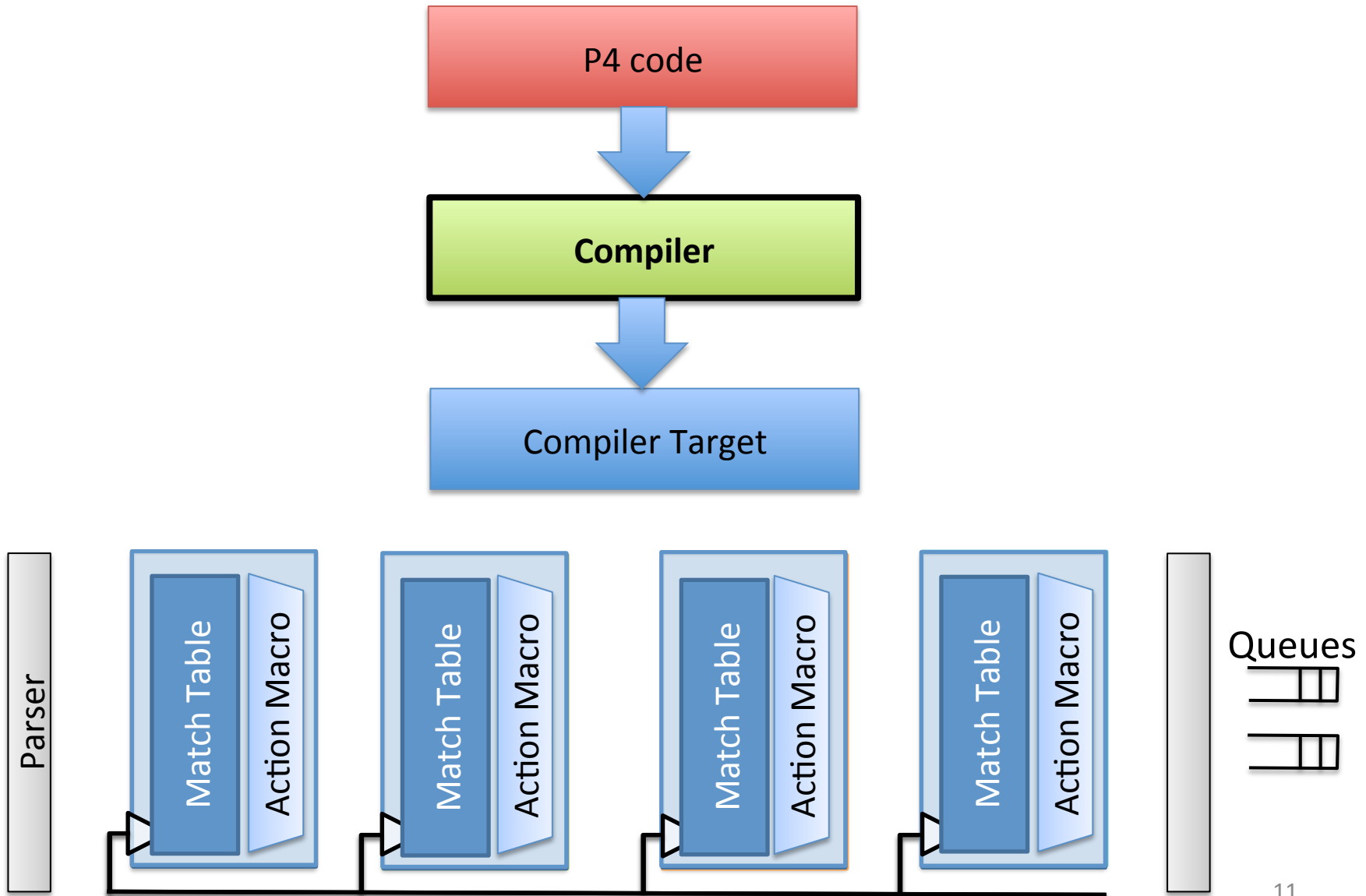
News & Events > Press Releases

**Cavium and XPliant Introduce a Fully Programmable Switch Silicon Family Scaling to 3.2 Terabits per Second**

**XPliant™ Packet Architecture (XPA™) Enables Data Plane Flexibility for SDN Without Compromising Raw Speed**



# Configuring Switch Chips



# P4 and PISA

## Parser (ANCS'13)

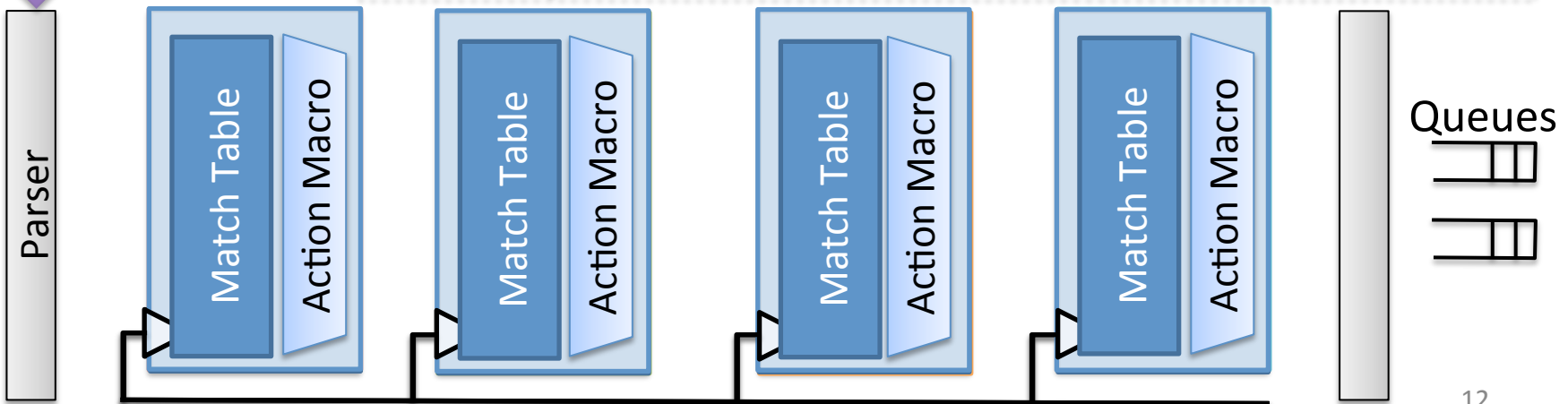
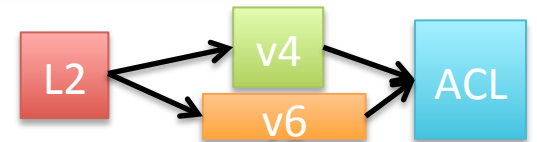
```
parser parse_ethernet {  
  extract(ethernet);  
  select(latest.etherType)  
  {  
    0x800 : parse_ipv4;  
    0x86DD : parse_ipv6;  
  }  
}
```

## Match Action Tables

```
table ipv4_lpm {  
  reads {  
    ipv4.dstAddr :  
      lpm;  
  }  
  actions {  
    set_next_hop;  
    drop;  
  }  
}
```

## Control Flow Graph

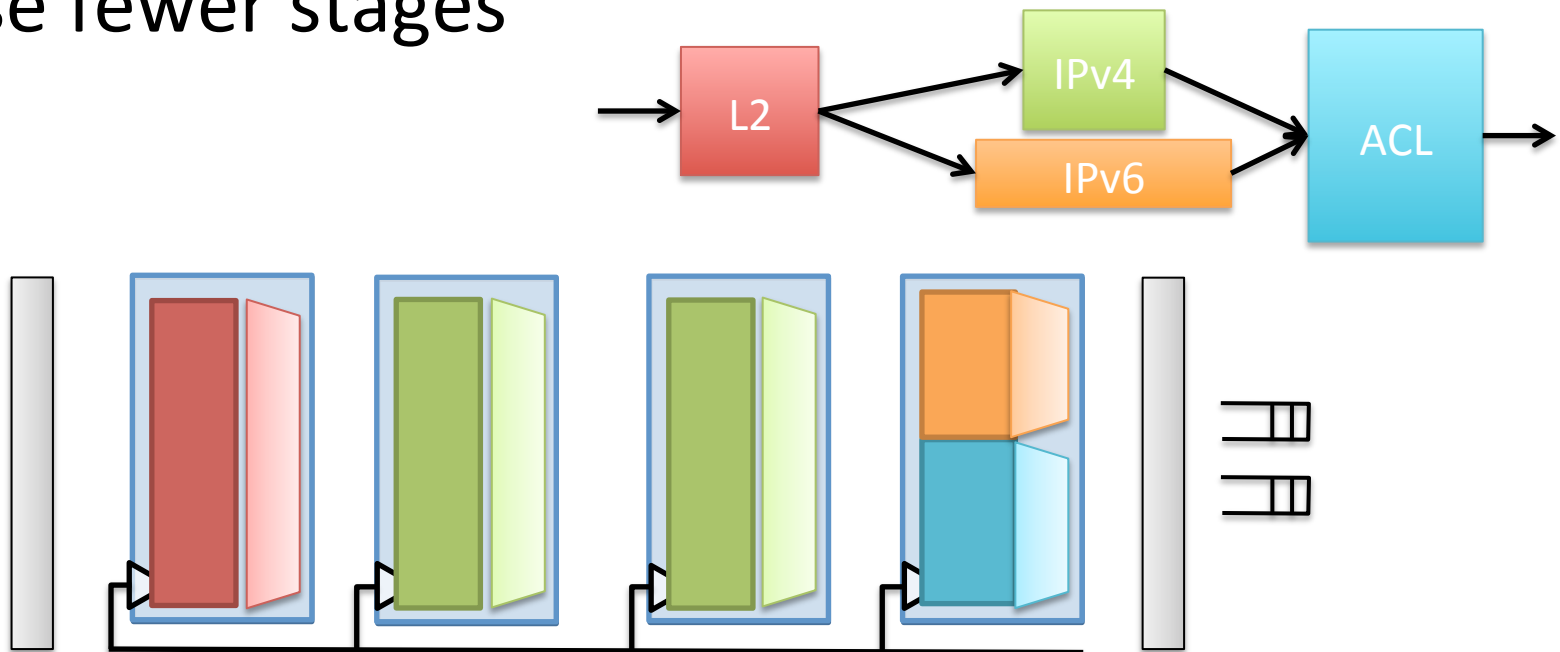
```
control ingress  
{  
  apply(l2_table);  
  if (valid(ipv4)) {  
    apply(ipv4_table);  
  }  
  if (valid(ipv6)) {  
    apply(ipv6_table);  
  }  
  apply (acl);  
}
```



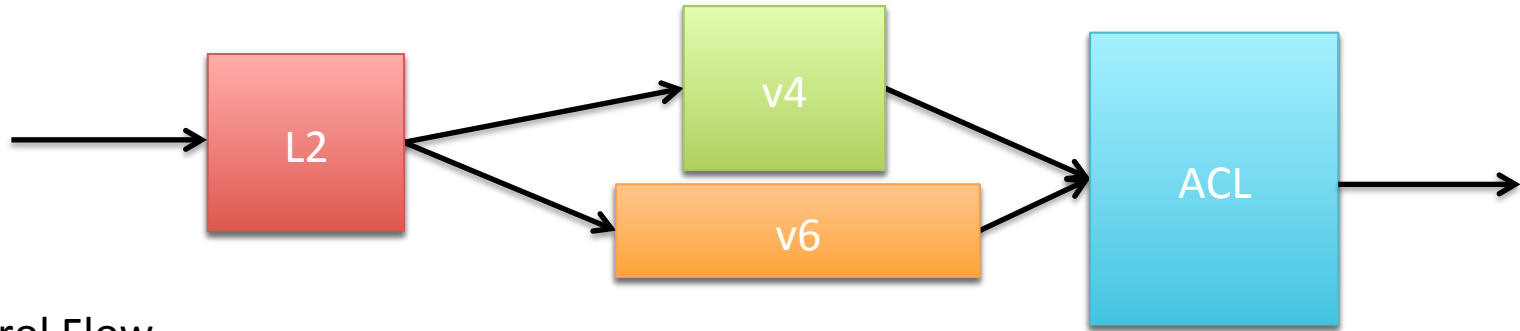
What does reconfigurability  
buy us?

# Optimizing with Reconfigurability

- Use resources efficiently
  - Multiple tables per stage
  - Big table in multiple stages
- Use fewer stages

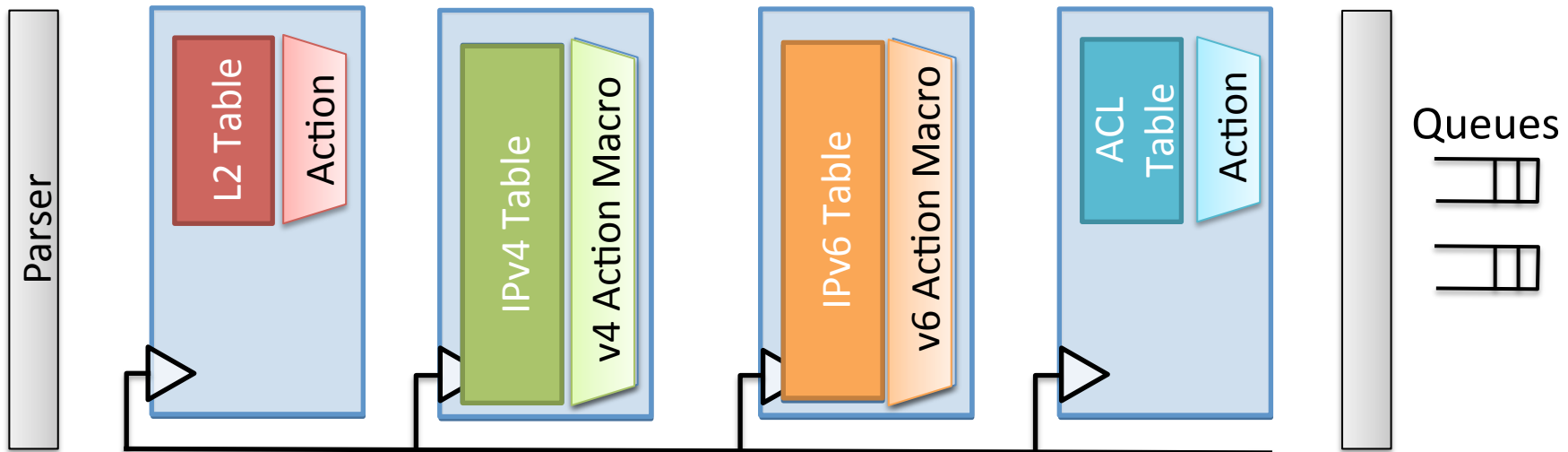


# Naïve Mapping: Control Flow Graph

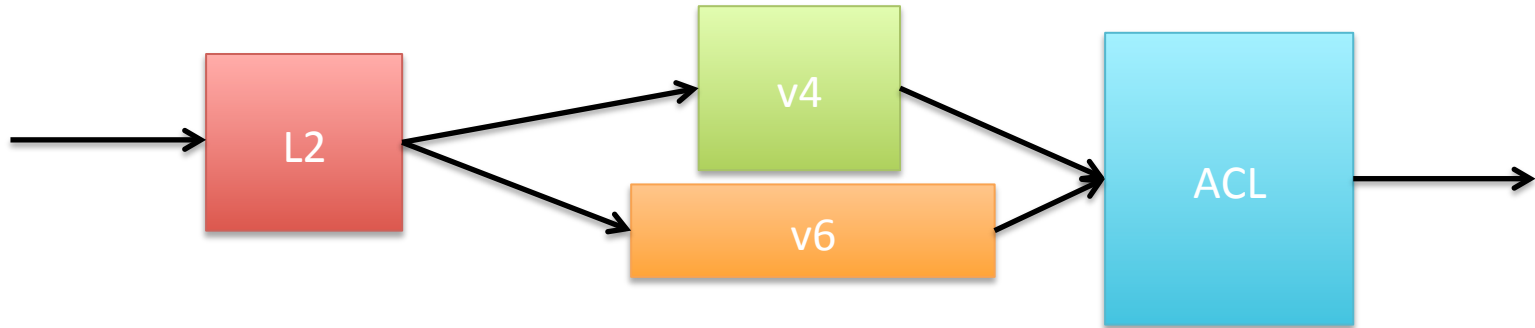


Control Flow

Switch Pipeline

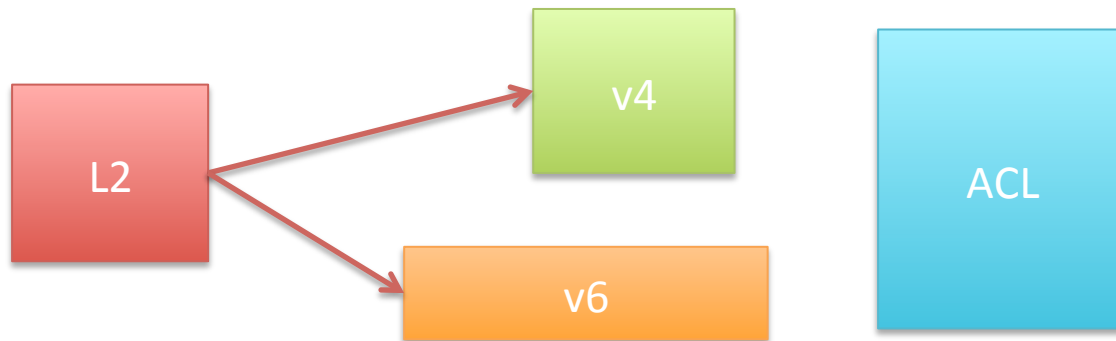


# Table Dependency Graph (TDG)



Control Flow Graph

Table Dependency Graph





# Optimizing the Mapping: TDG

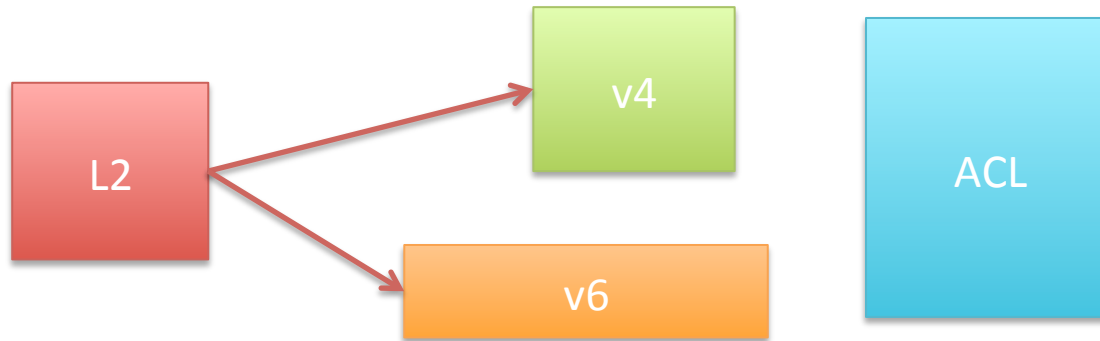
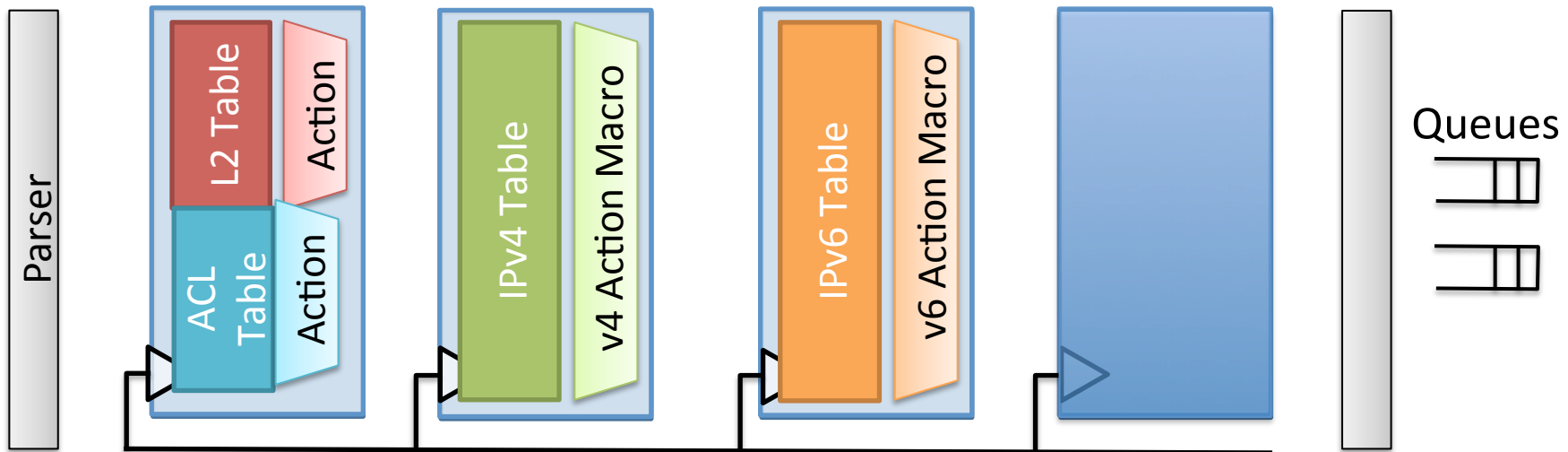
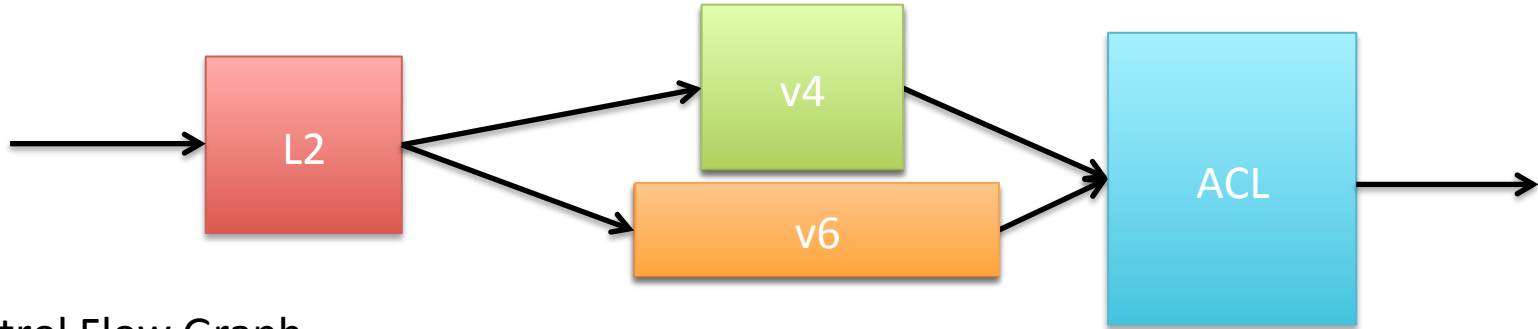


Table Dependency Graph

Switch Pipeline

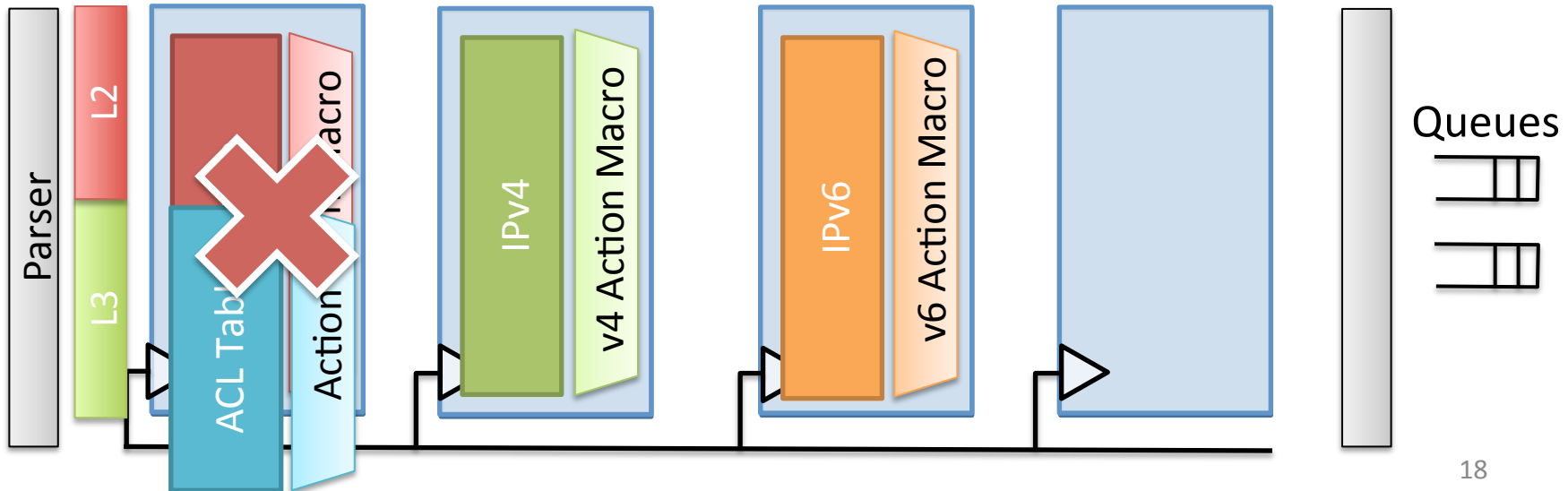


# Resource constraints



Control Flow Graph

Switch Pipeline



# More resource constraints

*Table parallelism*

*Action Memory*

*Memory Type*

*Action ALU input*

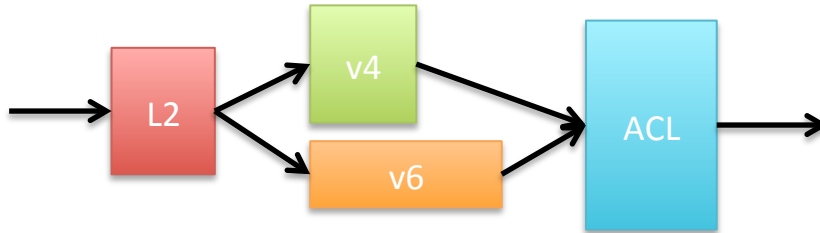
*Header widths*

# The Compiler Problem

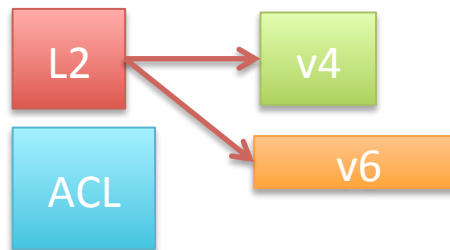
Map match action tables in a TDG to a switch pipeline while respecting dependency and resource constraints.

## Step 1: P4 Program

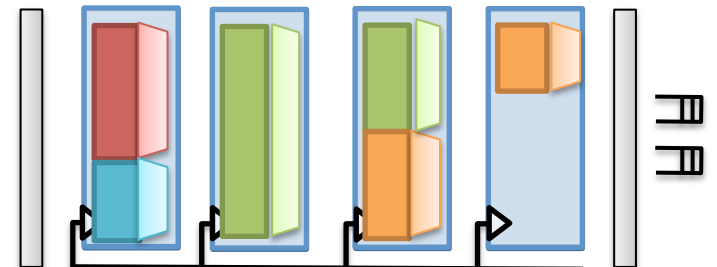
### Step 2: Control Flow Graph



### Step 3: Table Dependency Graph

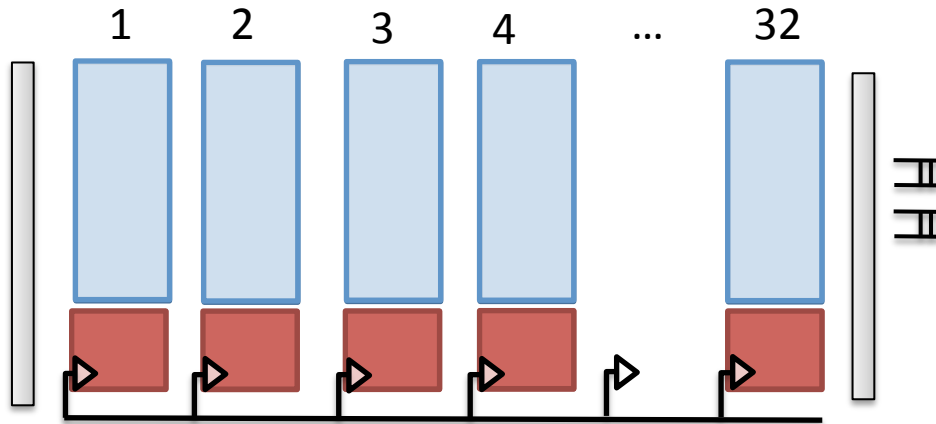


### Step 4: Table Configuration



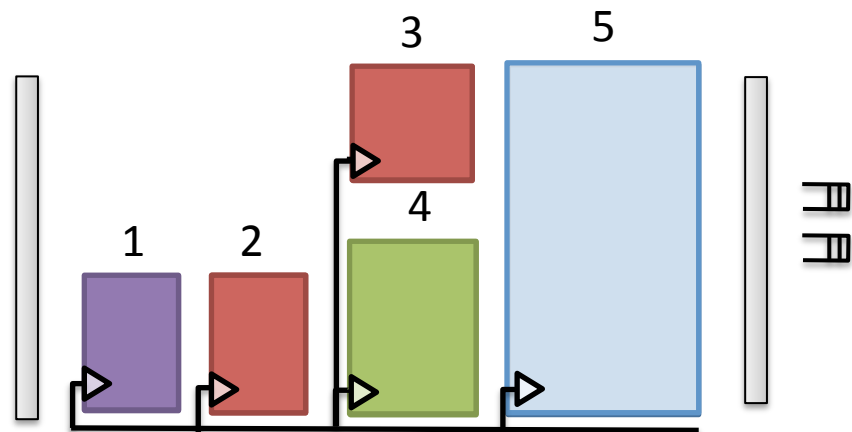
Is that it?

# Two Switches We Studied



RMT  
32 Stages  
(SIGCOMM 2013)

FlexPipe  
5 Stages  
(Intel FM6000)



# Additional switch features

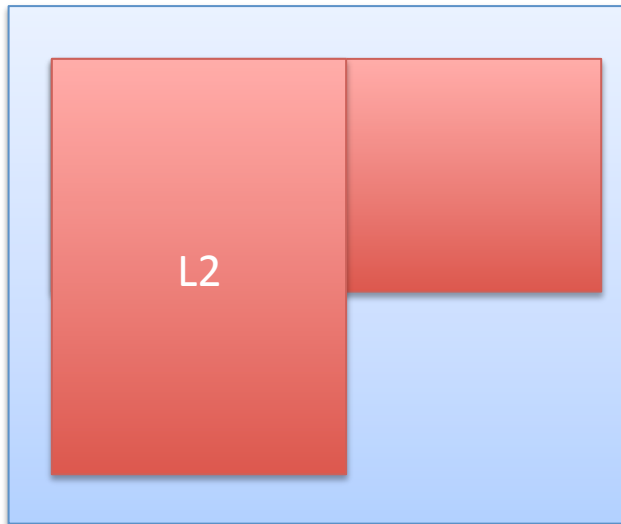


Table shaping in RMT

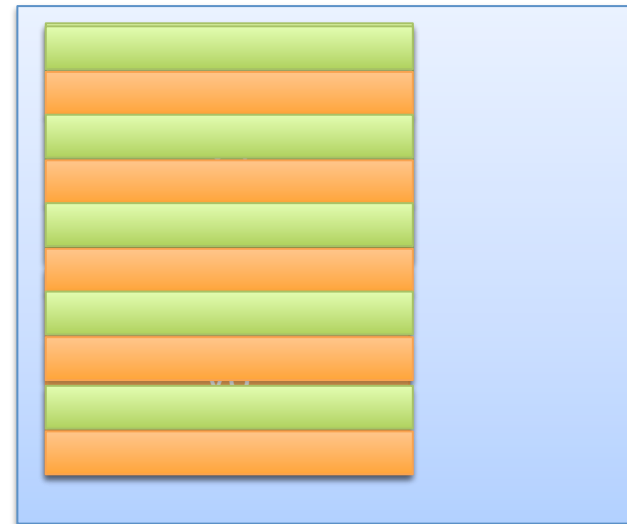


Table sharing in FlexPipe



*Action Memory*

*Table parallelism*

*Memory Type*

# The Compiler Problem

*Header widths*

*Action ALU input*

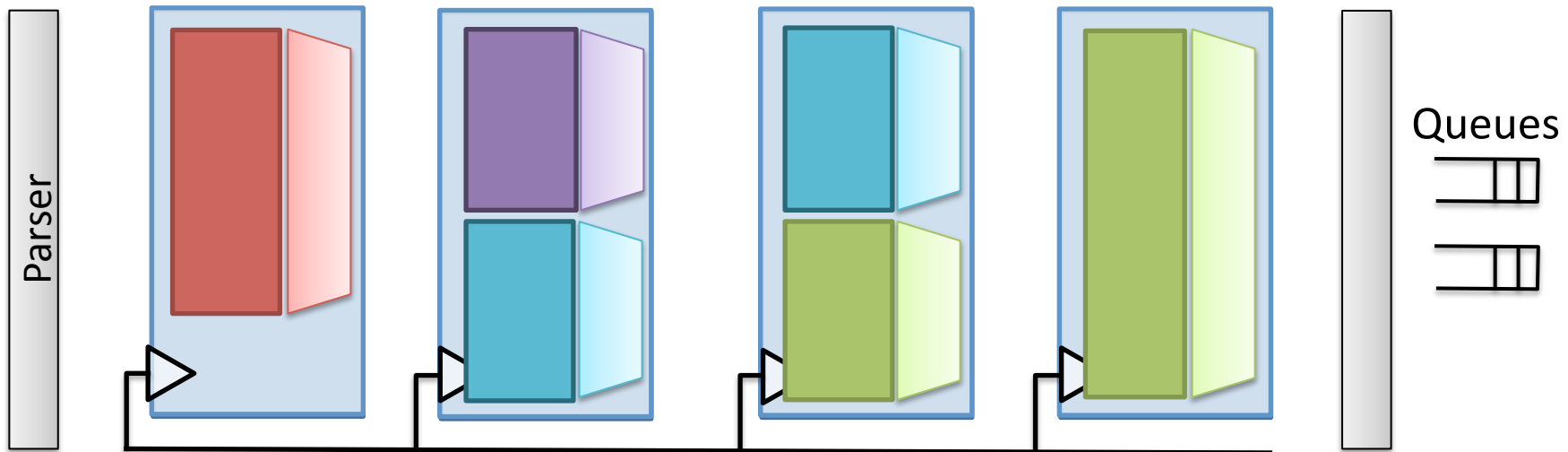
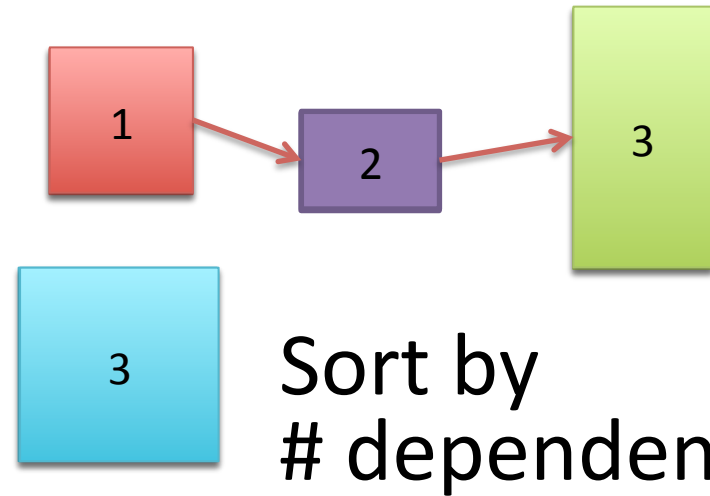
Map match action tables in a TDG to a switch pipeline while respecting dependency and resource constraints.

***Table shaping***

***Table sharing***

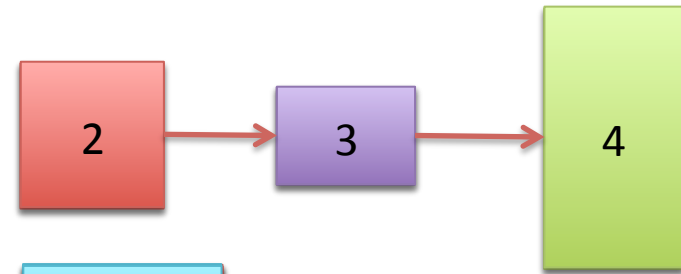
# First approach: Greedy

- Prioritize one constraint
- Sort tables
- Map tables one at a time

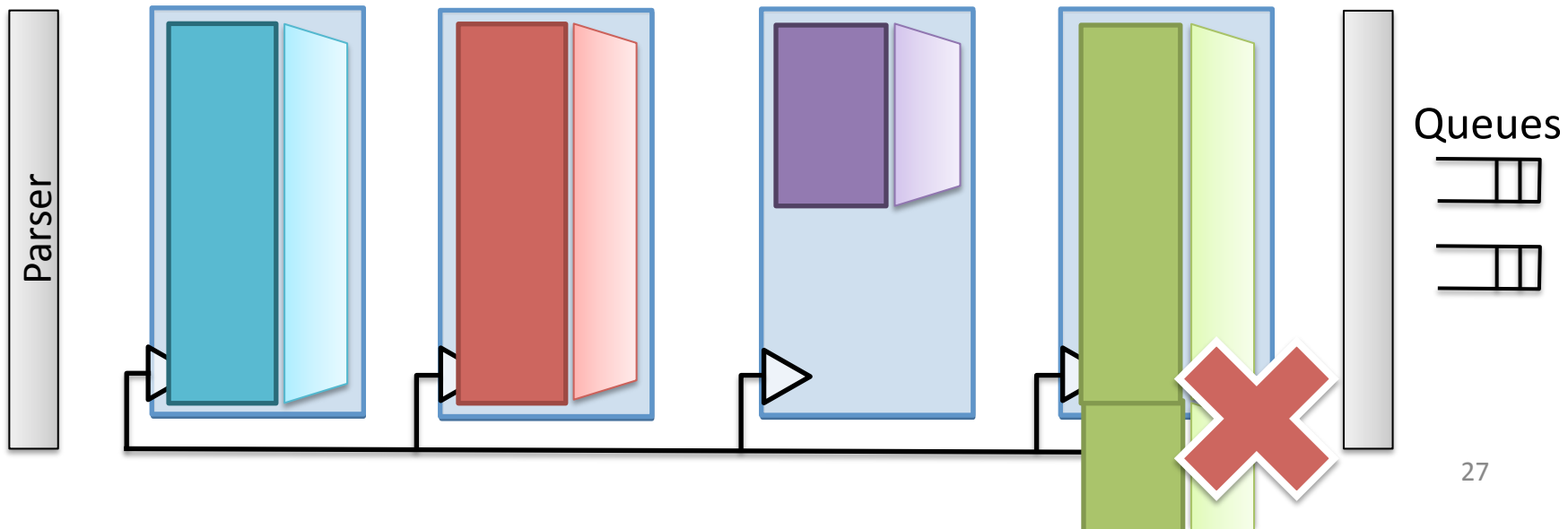


# First approach: Greedy

- Prioritize one constraint
- Sort tables
- Map tables one at a time



Sort by  
match width



# Too many constraints for Greedy

- Any greedy must sort tables based on a metric that is a *fixed* function of constraints.
- As the number of constraints gets larger, it's harder for a fixed function to represent the interplay between all constraints.
- Can we do better than greedy?

# Second approach: Integer Linear Programming (ILP)

Find an optimal mapping.

## Pros:

- Takes in all constraints
- Different objectives
- Solvers exist (CPLEX)

## Cons:

- Blackbox solver
- Encoding is an art
- Slow

# ILP Setup

min # stages

subject to:

table sizes  
assigned

$\geq$

table sizes  
specified

memories  
assigned

$\leq$

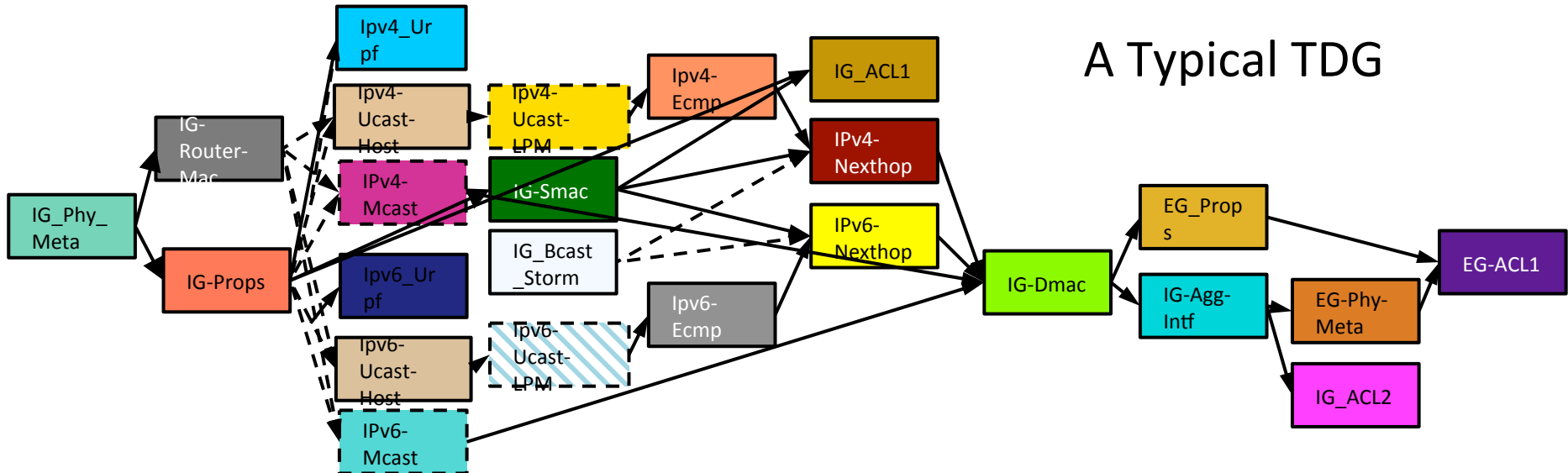
memories in  
physical stage

dependency constraints

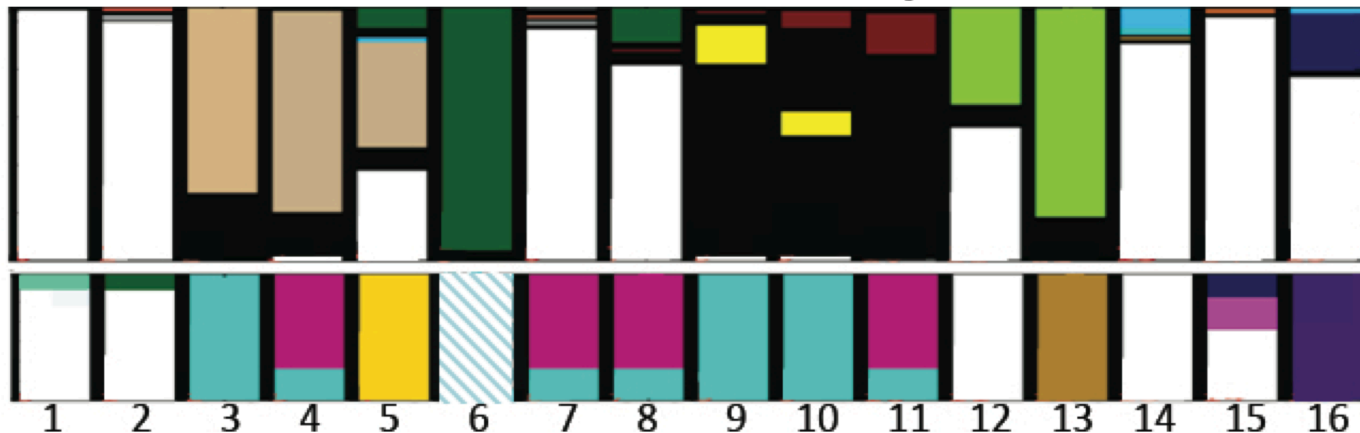
# Experiment Setup

- 4 datacenter use cases from Intel, Barefoot
- Differ in tables, table sizes, and dependencies

# Example Use Case



Configuration for RMT





# Setup: Greedy vs ILP

1. Ability to fit: FlexPipe
  - Variants of use cases in 5-stage pipeline.
2. Optimality: RMT
  - Minimum stage, pipeline latency, power
3. Runtime: both switches

# Results: Greedy vs ILP

1. Can Greedy fit my program?
  - Yes, if resources aplenty (RMT, 32 stages)
  - No, if resources constrained (FlexPipe, 5 stages),  
Can't fit 25% of programs .
2. How close to optimal is Greedy?
  - 30% more time for packet to get through RMT pipeline.
3. Hmm.. looks like I need ILP. How slow is it?
  - 100x slower than Greedy
  - Reasonable if programs don't change often.

If we have time,  
we should run ILP.

# Use ILP to suggest best Greedy for program type.

## Critical constraints

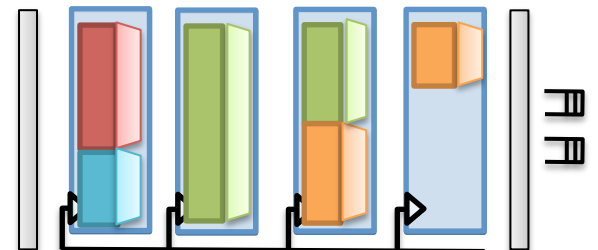
- Dependency critical: 16  $\rightarrow$  13 stages
- Additional resource constraints less important

## Critical resources

- TCAM memories critical: 16  $\rightarrow$  14 stages
  - Results for one of our datacenter L2/L3 use cases

# Conclusion

- **Challenge:** Parallelism and constraints in reconfigurable chips makes compiling difficult.
- **TDG:** highlights parallelism in program.
- **ILP:** better if enough time, fitting is critical, or objectives are complicated.
- **Best Greedy:** ILP can choose via notion of *critical* constraints and *critical* resources.



# Thank you!

NSDI '15: Compiling Packet Programs to  
Reconfigurable Switches

Research funded by AT&T, Intel, Open Networking Research Center.

Back up slides

# Greedy vs ILP

- Greedy often, but not always, can fit.
  - 17 v/s 16 stages for RMT (increased latency)
  - fits 93% of FlexPipe use cases
- Greedy is suboptimal for complicated objectives
  - Latency: 130 cycles v/s 104 cycles
- ILP can be slower than greedy
  - 3.8 min for Latency v/s 1s for Greedy
  - Longest ILP run time: FlexPipe program 45 min.
  - Reasonable if programs to change every few months



# ILP Run time

- Number of constraints? Not obvious. For RMT:
  - Min. stage: few secs.
  - Min. power: few secs.
  - Min. pipeline latency 10x slower
- Number of variables? How fine-grained is the resource assignment? For FlexPipe:
  - One match entry at a time: many days..
  - 100-500 match entries at a time: < 1 hr

# Use Cases

| Name             | Switch   | $N$ | Dependencies |        |       |
|------------------|----------|-----|--------------|--------|-------|
|                  |          |     | Match        | Action | Other |
| L2L3<br>-Complex | RMT      | 24  | 23           | 2      | 10    |
| L2L3<br>-Simple  | RMT      | 16  | 4            | 0      | 15    |
|                  | FlexPipe | 13  | 12           | 0      | 4     |
| L2L3<br>-Mtag    | RMT      | 19  | 6            | 1      | 16    |
|                  | FlexPipe | 11  | 9            | 1      | 3     |
| L3DC             | RMT      | 13  | 7            | 3      | 1     |