

P51: High Performance Networking

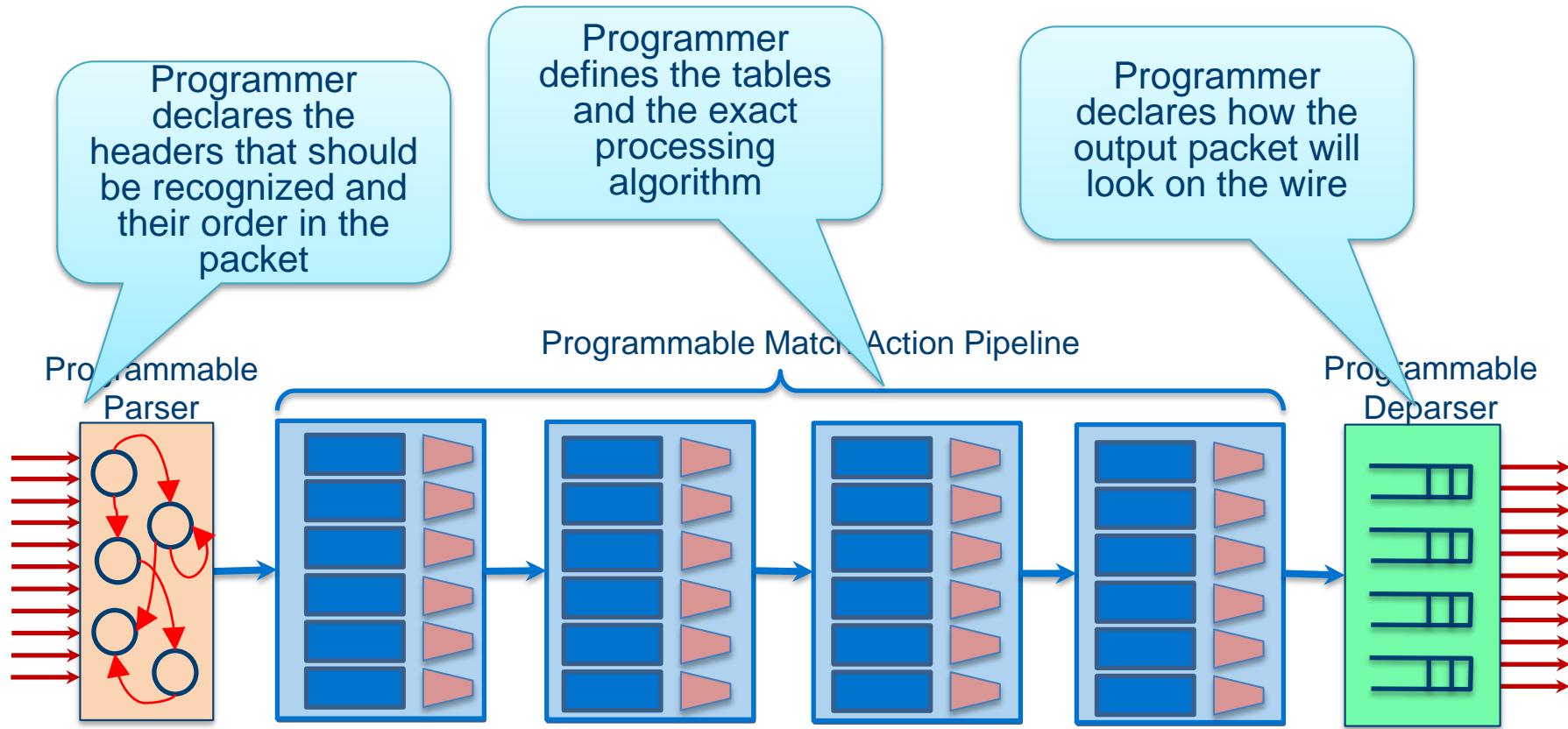
Introduction to P4 and P4-NetFPGA

Noa Zilberman
noa.zilberman@cl.cam.ac.uk

Lent 2019/20

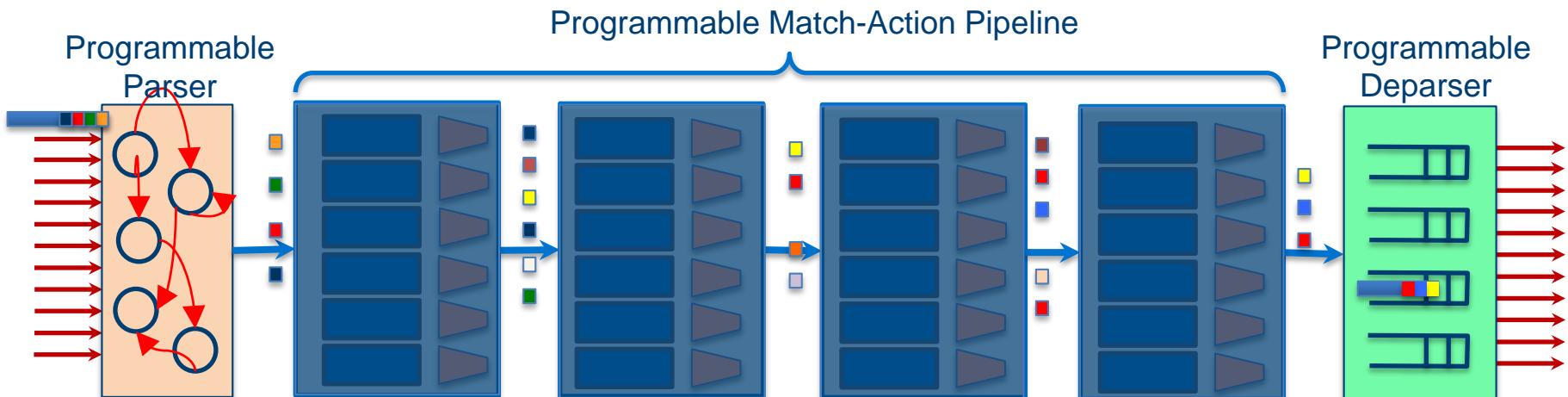
P4

PISA: Protocol-Independent Switch Architecture

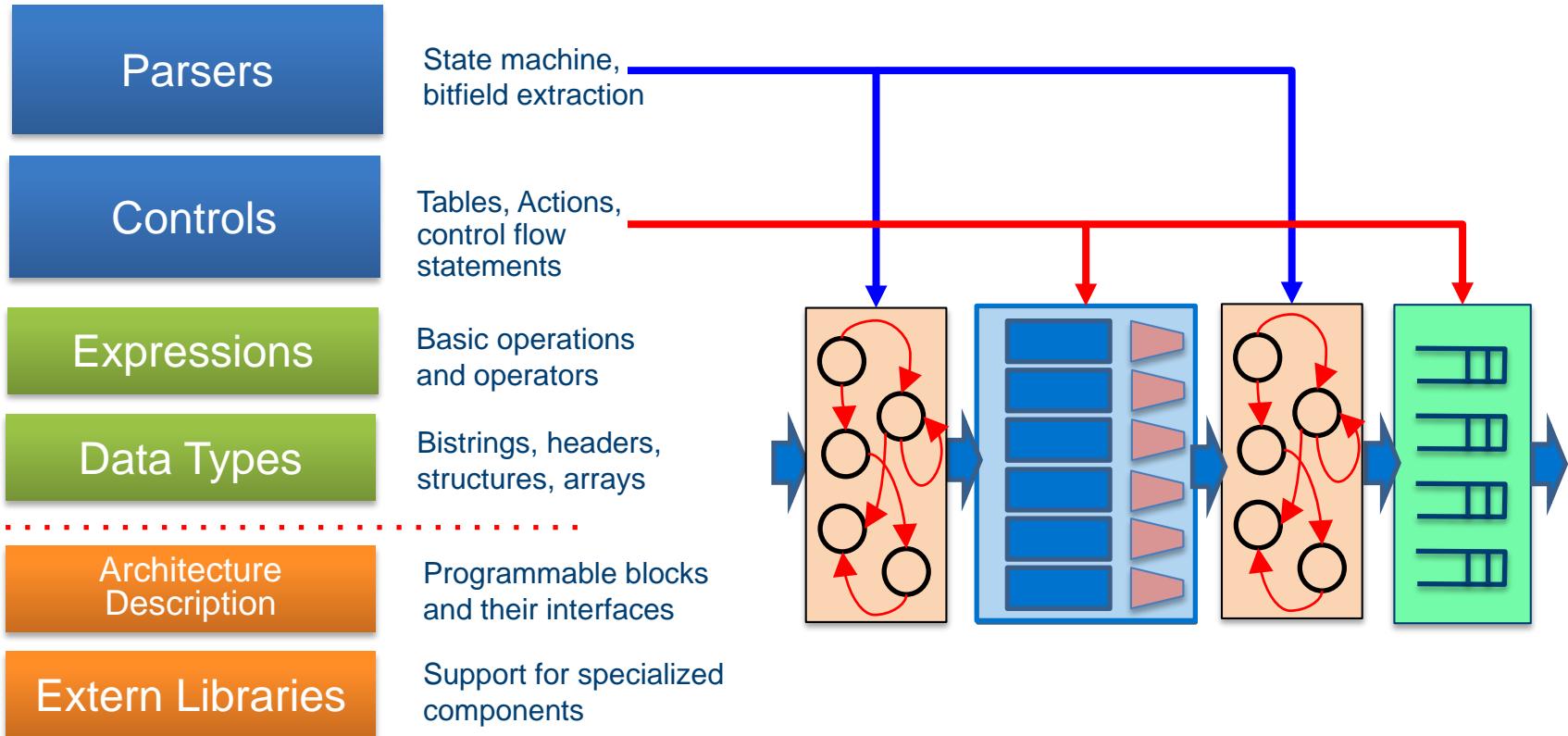


PISA in Action

- Packet is parsed into individual headers (parsed representation)
- Headers and intermediate results can be used for matching and actions
- Headers can be modified, added or removed
- Packet is deparsed (serialized)



P4₁₆ Language Elements

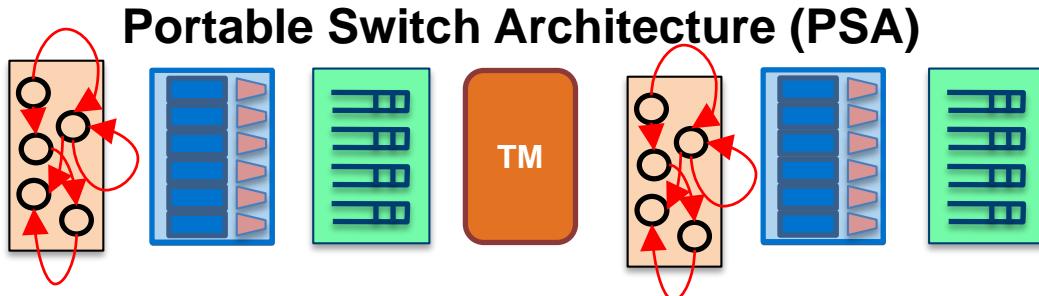
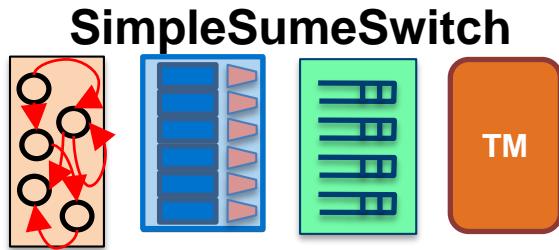
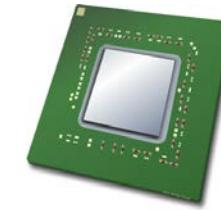
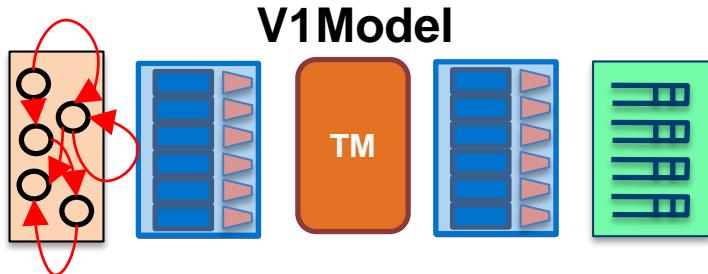


P4_16 Approach

Term	Explanation
P4 Target	An embodiment of a specific hardware implementation
P4 Architecture	A specific set of P4-programmable components, externs, fixed components and their interfaces available to the P4 programmer
P4 Platform	P4 Architecture implemented on a given P4 Target

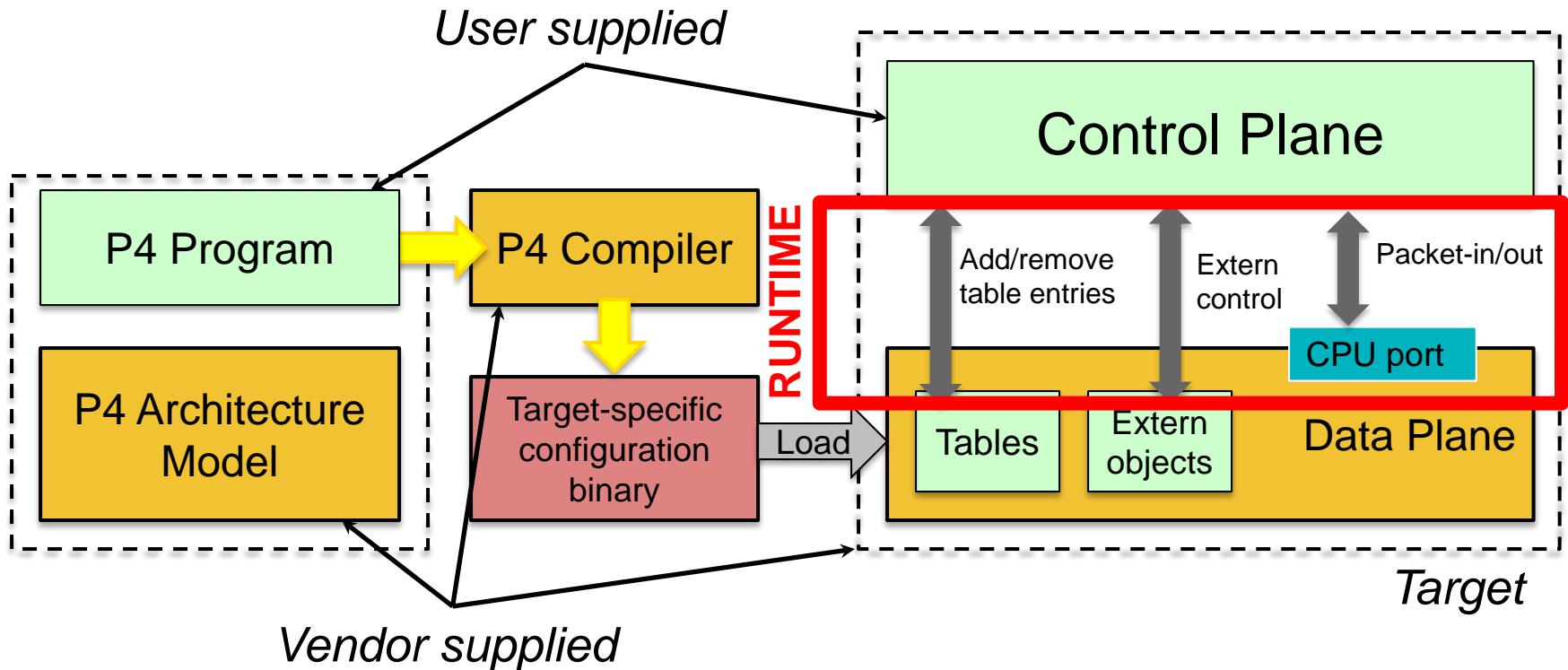


Example Architectures and Targets

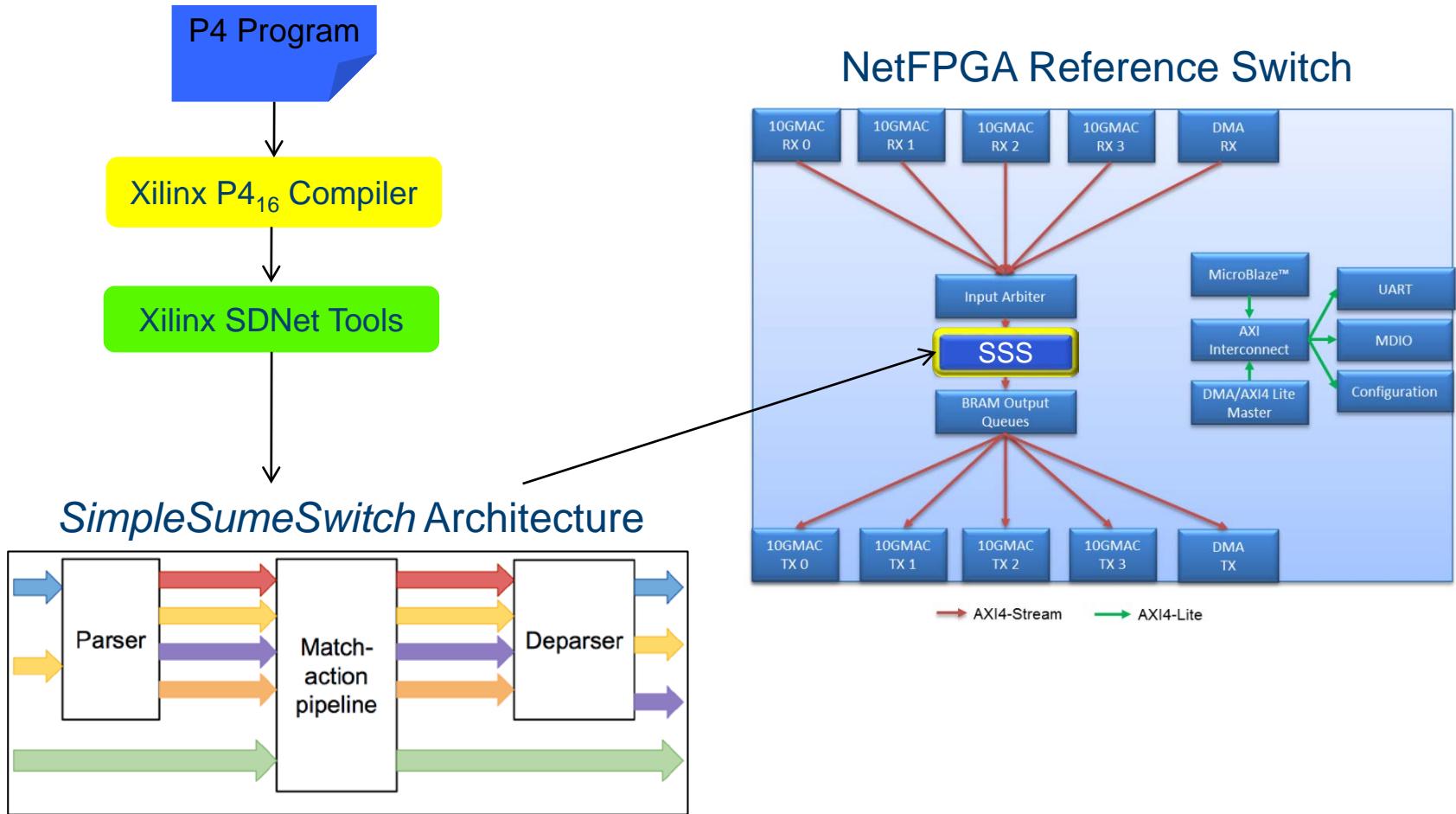


Anything

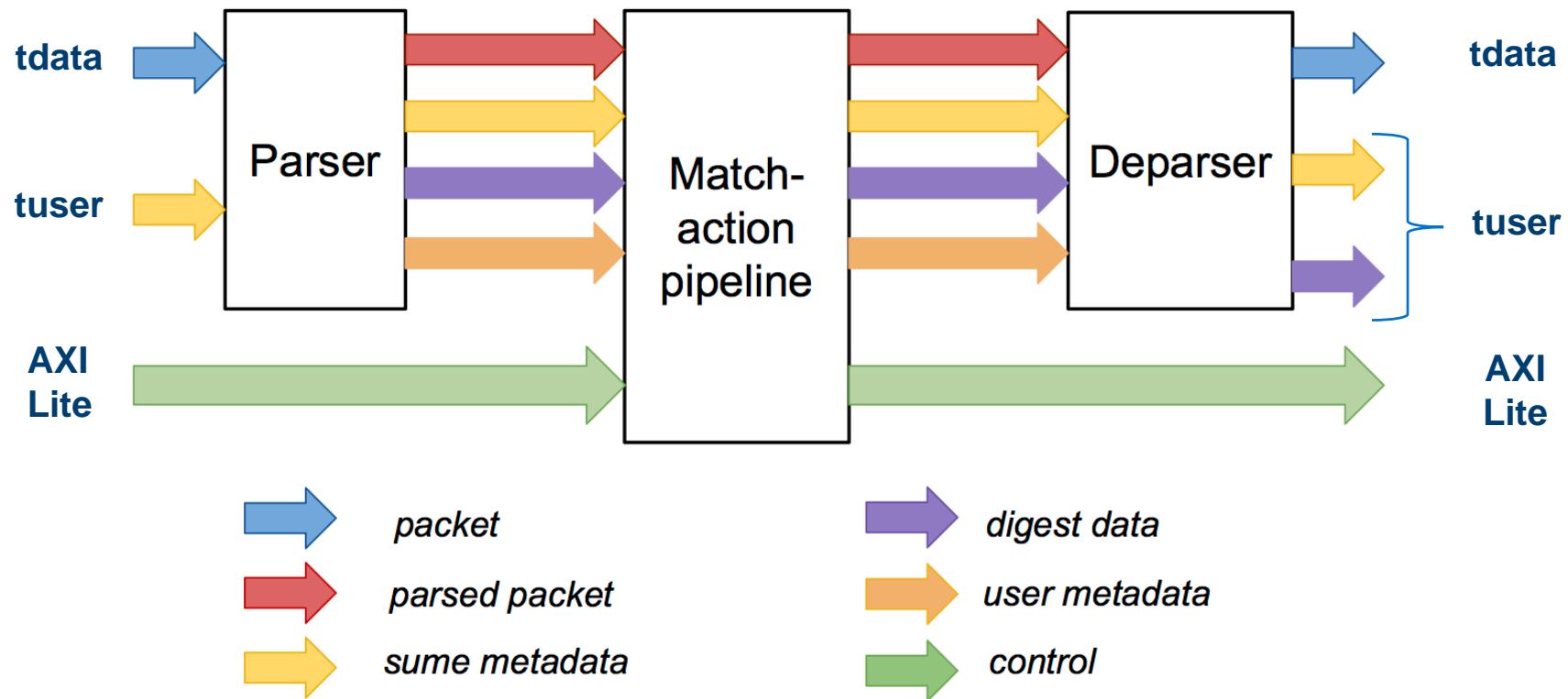
Programming a P4 Target



Reminder: P4 on NetFPGA (P4-NetFPGA)



SimpleSumeSwitch Architecture Model



P4 used to describe parser, match-action pipeline, and deparser

V1Model Standard Metadata

```
struct standard_metadata_t {  
    bit<9> ingress_port;  
    bit<9> egress_spec;  
    bit<9> egress_port;  
    bit<32> clone_spec;  
    bit<32> instance_type;  
    bit<1> drop;  
    bit<16> recirculate_port;  
    bit<32> packet_length;  
    bit<32> enq_timestamp;  
    bit<19> enq_qdepth;  
    bit<32> deq_timedelta;  
    bit<19> deq_qdepth;  
    bit<48> ingress_global_timestamp;  
    bit<32> lf_field_list;  
    bit<16> mcast_grp;  
    bit<1> resubmit_flag;  
    bit<16> egress_rid;  
    bit<1> checksum_error;  
}
```

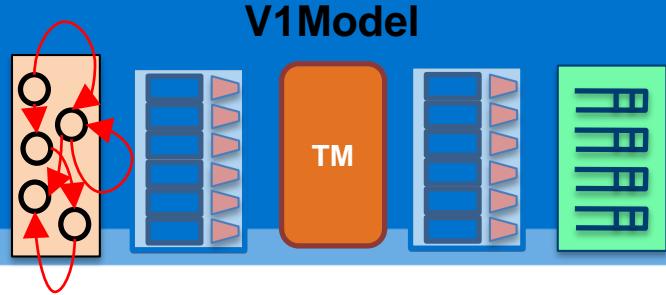
- **ingress_port** - the port on which the packet arrived
- **egress_spec** - the port to which the packet should be sent to
- **egress_port** - the port on which the packet is departing from (read only in egress pipeline)

Standard Metadata in SimpleSumSwitch Architecture

```
struct sume_metadata_t {
    bit<16> dma_q_size;
    bit<16> nf3_q_size;
    bit<16> nf2_q_size;
    bit<16> nf1_q_size;
    bit<16> nf0_q_size;
    bit<8> send_dig_to_cpu; // send digest_data to CPU
    bit<8> dst_port; // one-hot encoded
    bit<8> src_port; // one-hot encoded
    bit<16> pkt_len; // unsigned int
}
```

- *_q_size – size of each output queue, measured 32B words
Taken when the packet starts being processed by the P4 program
- src_port/dst_port – one-hot encoded, easy to do multicast
- user_metadata/digest_data – structs defined by **the user**

P4₁₆ Program Template (V1 Model)



```
#include <core.p4>
#include <v1model.p4>

/* HEADERS */
struct metadata { ... }
struct headers {
    ethernet_t    ethernet;
    ipv4_t        ipv4;
}

/* PARSER */
parser MyParser(packet_in packet,
                 out headers hdr,
                 inout metadata meta,
                 inout standard_metadata_t smeta) {
    ...
}

/* CHECKSUM VERIFICATION */
control MyVerifyChecksum(in headers hdr,
                         inout metadata meta) {
    ...
}

/* INGRESS PROCESSING */
control MyIngress(inout headers hdr,
                  inout metadata meta,
                  inout standard_metadata_t std_meta) {
    ...
}
```

```
/* EGRESS PROCESSING */
control MyEgress(inout headers hdr,
                  inout metadata meta,
                  inout standard_metadata_t std_meta) {
    ...
}

/* CHECKSUM UPDATE */
control MyComputeChecksum(inout headers hdr,
                          inout metadata meta) {
    ...
}

/* DEPARSER */
control MyDeparser(inout headers hdr,
                    inout metadata meta) {
    ...
}

/* SWITCH */
V1Switch(
    MyParser(),
    MyVerifyChecksum(),
    MyIngress(),
    MyEgress(),
    MyComputeChecksum(),
    MyDeparser()
) main;
```

P4₁₆ Program Template (SimpleSumeSwitch)

```
#include <core.p4>
#include <sume switch.p4>
/* HEADERS */
struct user_metadata_t { ... }
struct digest_data_t { ... }
struct headers {
    ethernet_t   ethernet;
    ipv4_t        ipv4;
}
/* PARSER */
parser TopParser (packet_in b,
                  out Parsed_packet p,
                  out user_metadata_t user_metadata,
                  out digest_data_t digest_data,
                  inout sume_metadata_t sume_metadata) {
    ...
}
/* PROCESSING */
control TopPipe (inout Parsed_packet p,
                  inout user_metadata_t user_metadata,
                  inout digest_data_t digest_data,
                  inout sume_metadata_t sume_metadata) {
    ...
}
```

```
/* DEPARSER */
control TopDeparser (packet_out b,
                      in Parsed_packet p,
                      in user_metadata_t user_metadata,
                      inout digest_data_t digest_data,
                      inout sume_metadata_t sume_metadata) {
    ...
}

/* SWITCH */
SimpleSumeSwitch (
    TopParser(),
    TopPipe(),
    TopDeparser()
) main;
```

P4₁₆ Hello World (V1Model)

```
#include <core.p4>
#include <v1model.p4>
struct metadata {}
struct headers {}

parser MyParser(packet_in packet,
    out headers hdr,
    inout metadata meta,
    inout standard_metadata_t standard_metadata) {
    state start { transition accept; }

control MyVerifyChecksum(inout headers hdr, inout metadata
meta) { apply { } }

control MyIngress(inout headers hdr,
    inout metadata meta,
    inout standard_metadata_t standard_metadata) {
apply {
    if (standard_metadata.ingress_port == 1) {
        standard_metadata.egress_spec = 2;
    } else if (standard_metadata.ingress_port == 2) {
        standard_metadata.egress_spec = 1;
    }
}
```

```
control MyEgress(inout headers hdr,
    inout metadata meta,
    inout standard_metadata_t standard_metadata) {
    apply { }
}

control MyComputeChecksum(inout headers hdr, inout metadata
meta) {
    apply { }
}

control MyDeparser(packet_out packet, in headers hdr) {
    apply { }
}

V1Switch(
    MyParser(),
    MyVerifyChecksum(),
    MyIngress(),
    MyEgress(),
    MyComputeChecksum(),
    MyDeparser()
) main;
```

P4₁₆ Hello World (V1Model)

```
#include <core.p4>
#include <v1model.p4>
struct metadata {}
struct headers {}

parser MyParser(packet_in packet, out headers hdr,
    inout metadata meta,
    inout standard_metadata_t standard_metadata) {
    state start { transition accept; }
}

control MyIngress(inout headers hdr, inout metadata meta,
    inout standard_metadata_t standard_metadata) {
    action set_egress_spec(bit<9> port) {
        standard_metadata.egress_spec = port;
    }
}






```

```
control MyEgress(inout headers hdr,
    inout metadata meta,
    inout standard_metadata_t standard_metadata) {
    apply {    }
}

control MyVerifyChecksum(inout headers hdr, inout metadata meta) {    apply {    }    }

control MyComputeChecksum(inout headers hdr, inout metadata meta) {    apply {    }    }

control MyDeparser(packet_out packet, in headers hdr) {
    apply {    }
}

V1Switch( MyParser(), MyVerifyChecksum(), MyIngress(),
MyEgress(), MyComputeChecksum(), MyDeparser() ) main;
```

Key	Action Name	Action Data
1	set_egress_spec	2
2	set_egress_spec	1

P4₁₆ Types (Basic and Header Types)

```
typedef bit<48> macAddr_t;
typedef bit<32> ip4Addr_t;
header ethernet_t {
    macAddr_t dstAddr;
    macAddr_t srcAddr;
    bit<16> etherType;
}
header ipv4_t {
    bit<4> version;
    bit<4> ihl;
    bit<8> diffserv;
    bit<16> totalLen;
    bit<16> identification;
    bit<3> flags;
    bit<13> fragOffset;
    bit<8> ttl;
    bit<8> protocol;
    bit<16> hdrChecksum;
    ip4Addr_t srcAddr;
    ip4Addr_t dstAddr;
}
```

Basic Types

- **bit<n>**: Unsigned integer (bitstring) of size n
- **bit** is the same as **bit<1>**
- **int<n>**: Signed integer of size n (≥ 2)
- **varbit<n>**: Variable-length bitstring

Header Types: Ordered collection of members

- Can contain **bit<n>**, **int<n>**, and **varbit<n>**
- Byte-aligned
- Can be valid or invalid
- Provides several operations to test and set validity bit:
isValid(), **setValid()**, and **setInvalid()**

Typedef: Alternative name for a type

P4₁₆ Types (Other Types)

```
/* Architecture */
struct standard_metadata_t {
    bit<9> ingress_port;
    bit<9> egress_spec;
    bit<9> egress_port;
    bit<32> clone_spec;
    bit<32> instance_type;
    bit<1> drop;
    bit<16> recirculate_port;
    bit<32> packet_length;
    ...
}

/* User program */
struct metadata {
    ...
}
struct headers {
    ethernet_t   ethernet;
    ipv4_t        ipv4;
}
```

Other useful types

- **Struct:** Unordered collection of members (with no alignment restrictions)
- **Header Stack:** array of headers
- **Header Union:** one of several headers

Intrinsic Metadata

- The data that a P4-programmable components can use to interface with the rest of the system
- Defined in the files supplied by the vendor

Declaring and Initializing Variables

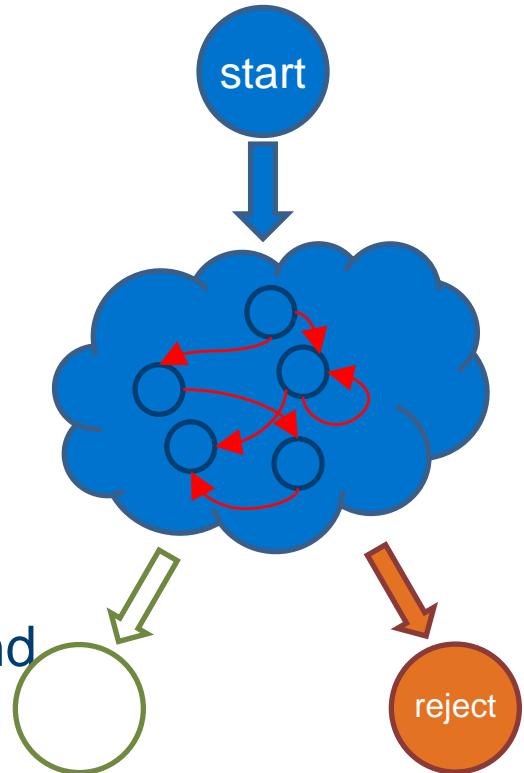
```
bit<16>      my_var;  
bit<8>       another_var = 5;  
  
const bit<16> ETHERTYPE_IPV4 = 0x0800;  
const bit<16> ETHERTYPE_IPV6 = 0x86DD;  
  
ethernet_t    eth;  
vlan_tag_t   vtag = {3w2, 0, 12w13, 16w0x8847};
```

Safe constants with
explicit widths

- In P4₁₆ you can instantiate variables of both base and derived types
- Variables can be initialized
 - Including the composite types
- Constant declarations make for safer code
- Infinite width and explicit width constants

P4₁₆ Parsers

- Parsers are functions that map packets into headers and metadata,
 - written in a state machine style
- Every parser has three predefined states
 - start
 - accept
 - reject
- Other states may be defined by the programmer
- In each state, execute zero or more statements, and then transition to another state (loops are OK)



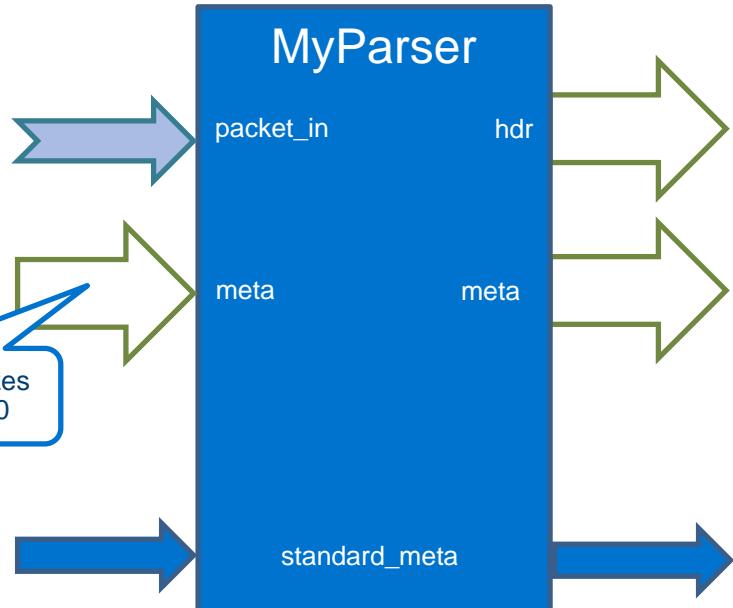
Parsers (V1Model)

```
/* From core.p4 */
extern packet_in {
    void extract<T>(out T hdr);
    void extract<T>(out T variableSizeHeader,
                    in bit<32> variableFieldSizeInBits);
    T lookahead<T>();
    void advance(in bit<32> sizeInBits);
    bit<32> length();
}

/* User Program */
parser MyParser(packet_in packet,
               out headers hdr,
               inout metadata meta,
               inout standard_metadata_t std_meta) {

    state start {
        packet.extract(hdr.ethernet);
        transition accept;
    }
}
```

The platform Initializes User Metadata to 0



Select Statement

```
state start {  
    transition parse_ethernet;  
}  
  
state parse_ethernet {  
    packet.extract(hdr.ethernet);  
    transition select(hdr.ethernet.etherType) {  
        0x800: parse_ipv4;  
        default: accept;  
    }  
}
```

P4₁₆ has a select statement that can be used to branch in a parser

Similar to case statements in C or Java, but without “fall-through behavior”—i.e., break statements are not needed

In parsers it is often necessary to branch based on some of the bits just parsed

For example, etherType determines the format of the rest of the packet

Match patterns can either be literals or simple computations such as masks

P4₁₆ Controls

- **Similar to C functions (without loops)**
- **Can declare variables, create tables, instantiate externs, etc.**
- **Functionality specified by code in apply statement**
- **Represent all kinds of processing that are expressible as DAG:**
 - Match-Action Pipelines
 - Deparsers
 - Additional forms of packet processing (updating checksums)
- **Interfaces with other blocks are governed by user- and architecture-specified types (typically headers and metadata)**

Example: Reflector (V1 Model)

```
control MyIngress(inout headers hdr,
                  inout metadata meta,
                  inout standard_metadata_t std_meta) {
    /* Declarations region */
    bit<48> tmp;

    apply {
        /* Control Flow */
        tmp = hdr.ethernet.dstAddr;
        hdr.ethernet.dstAddr = hdr.ethernet.srcAddr;
        hdr.ethernet.srcAddr = tmp;
        std_meta.egress_spec = std_meta.ingress_port;
    }
}
```

Desired Behavior:

- Swap source and destination MAC addresses
- Bounce the packet back out on the physical port that it came into the switch on

Example: Simple Actions

```
control MyIngress(inout headers hdr,
                  inout metadata meta,
                  inout standard_metadata_t std_meta) {

    action swap_mac(inout bit<48> src,
                    inout bit<48> dst) {
        bit<48> tmp = src;
        src = dst;
        dst = tmp;
    }

    apply {
        swap_mac(hdr.ethernet.srcAddr,
                  hdr.ethernet.dstAddr);
        std_meta.egress_spec = std_meta.ingress_port;
    }
}
```

- **Very similar to C functions**
- **Can be declared inside a control or globally**
- **Parameters have type and direction**
- **Variables can be instantiated inside**
- **Many standard arithmetic and logical operations are supported**
 - +, -, *
 - ~, &, |, ^, >>, <<
 - ==, !=, >, >=, <, <=
 - No division/modulo
- **Non-standard operations:**
 - Bit-slicing: [m:l] (works as l-value too)
 - Bit Concatenation: ++

Operating on Headers

```
action decap_ip_ip() {
    hdr.ipv4 = hdr.inner_ipv4;
    hdr.inner_ipv4.setInvalid();
}

action pop_mpls_label() {
    hdr.mpls.pop_front(1);
}

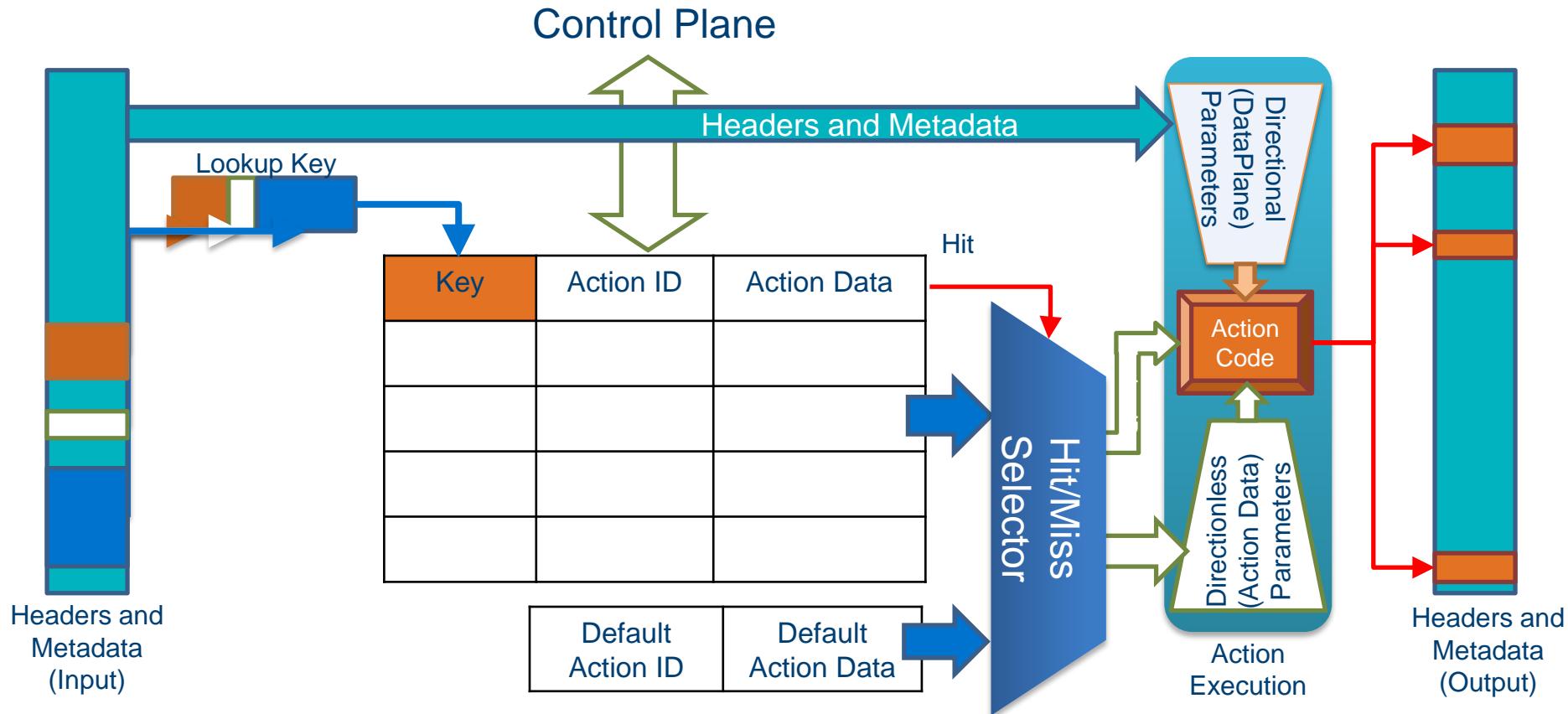
action push_mpls_label(in bit<20> label, in
bit<3> exp) {
    hdr.mpls.push_front(1);
    hdr.mpls[0].setValid();
    hdr.mpls[0] = { label, exp, 0, 64 };
}
```

- Header Validity bit manipulation:
 - header.setValid() – add_header
 - header.setInvalid() – remove_header
 - header.isValid()
- Header Assignment
 - header = { f1, f2, ..., fn }
 - header1 = header2
- Special operations on Header Stacks
 - In the parsers
 - header_stack.next
 - header_stack.last
 - header_stack.lastIndex
 - In the controls
 - header_stack[i]
 - header_stack.size
 - header_stack.push_front(int count)
 - header_stack.pop_front(int count)

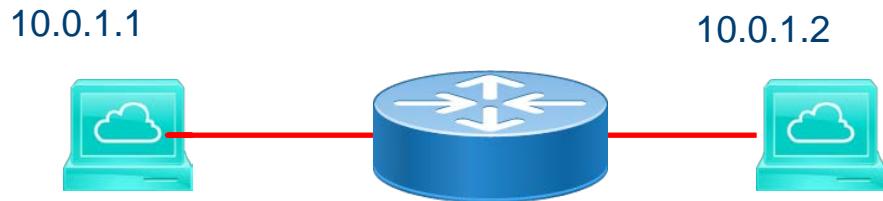
P4₁₆ Tables

- **The fundamental unit of a Match-Action Pipeline**
 - Specifies what data to match on and match kind
 - Specifies a list of *possible* actions
 - Optionally specifies a number of table **properties**
 - Size
 - Default action
 - Static entries
 - etc.
- **Each table contains one or more entries (rules)**
- **An entry contains:**
 - A specific key to match on
 - A **single** action that is executed when a packet matches the entry
 - Action data (possibly empty)

Tables: Match-Action Processing



Example: IPv4_LPM Table



Key	Action	Action Data
10.0.1.1/32	ipv4_forward	dstAddr=00:00:00:00:01:01 port=1
10.0.1.2/32	drop	
*	NoAction	

- **Data Plane (P4) Program**
 - Defines the format of the table
 - Key Fields
 - Actions
 - Action Data
 - Performs the lookup
 - Executes the chosen action
- **Control Plane (IP stack, Routing protocols)**
 - Populates table entries with specific information
 - Based on the configuration
 - Based on automatic discovery
 - Based on protocol calculations

IPv4_LPM Table

```
table ipv4_lpm {
    key = {
        hdr.ipv4.dstAddr: lpm;
    }
    actions = {
        ipv4_forward;
        drop;
        NoAction;
    }
    size = 1024;
    default_action = NoAction();
}
```

Match Kinds

```
/* core.p4 */
match_kind {
    exact,
    ternary,
    lpm
}

/* v1model.p4 */
match_kind {
    range,
    selector
}

/* Some other architecture */
match_kind {
    regexp,
    fuzzy
}
```

- **The type `match_kind` is special in P4**
- **The standard library (`core.p4`) defines three standard match kinds**
 - Exact match
 - Ternary match
 - LPM match
- **The architecture (`v1model.p4`) defines two additional match kinds:**
 - range
 - selector
- **Other architectures may define (and provide implementation for) additional match kinds**
- **P4-SDNet supports** exact, ternary, lpm and direct.

Defining Actions for L3 forwarding

```
/* core.p4 */
action NoAction() {

}

/* basic.p4 */
action drop() {
    mark_to_drop();
}

/* basic.p4 */
action ipv4_forward(macAddr_t dstAddr,
                     bit<9> port) {
    ...
}
```

- **Actions can have two different types of parameters**
 - Directional (from the Data Plane)
 - Directionless (from the Control Plane)
- **Actions that are called directly:**
 - Only use directional parameters
- **Actions used in tables:**
 - Typically use directionless parameters
 - May sometimes use directional parameters too



Applying Tables in Controls

```
control MyIngress(inout headers hdr,
                  inout metadata meta,
                  inout standard_metadata_t standard_metadata) {
    table ipv4_lpm {
        ...
    }
    apply {
        ...
        ipv4_lpm.apply();
        ...
    }
}
```

Table Initialization

```
control MyIngress(inout headers hdr,
                  inout metadata meta,
                  inout standard_metadata_t standard_metadata) {
    table ipv4_lpm {
        ...
    }
    apply {
        ...
    }
    const entries = {
        /*{ dstAddr, srcAddr, vlan_tag[0].isValid(), vlan_tag[1].isValid() } : action([action_data])*/
        { 48w000000000000, _, _, _ } : malformed_ethernet(ETHERNET_ZERO_DA);
        { _, 48w000000000000, _, _ } : malformed_ethernet(ETHERNET_ZERO_SA);
        { _, 48w010000000000 && 48w010000000000, _, _ } : malformed_ethernet(ETHERNET_MCAST_SA);
        { 48wFFFFFFFFFFFF, _, 0, _ } : broadcast_untagged();
        { 48wFFFFFFFFFFFF, _, 1, 0 } : broadcast_single_tagged();
        { 48wFFFFFFFFFFFF, _, 1, 1 } : broadcast_double_tagged();
        { 48w010000000000 && 48w010000000000, _, 0, _ } : multicast_untagged();
        { _, _, 0, _ } : unicast_untagged();
        { _, _, 1, 0 } : unicast_single_tagged();
    }
}
```

P4₁₆ Deparsing

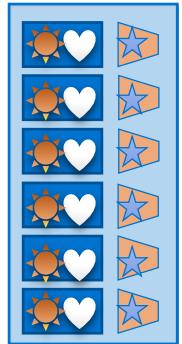
```
/* From core.p4 */
extern packet_out {
    void emit<T>(in T hdr);
}

/* User Program */
control DeparserImpl(packet_out packet,
                      in headers hdr) {
    apply {
        ...
        packet.emit(hdr.ethernet);
        ...
    }
}
```

- **Assembles the headers back into a well-formed packet**
- **Expressed as a control function**
 - No need for another construct!
- **packet_out extern is defined in core.p4:** emit(hdr): serializes header if it is valid
- **Advantages:**
 - Makes deparsing explicit...
 - ...but decouples from parsing

The Need for Externs

- Most platforms contain specialized facilities
 - Differ from vendor to vendor
 - Can't be expressed in the core language
 - Might have control-plane accessible state or configuration
- The language should stay the same
 - In P4₁₄ almost 1/3 of all the constructs were dedicated to specialized processing
 - In P4₁₆ all specialized objects use the same interface
- Objects can be used even if their implementation is hidden
 - Through instantiation and method calling



Stateless and Stateful Objects

- Stateless Objects: Reinitialized for each packet
 - Variables (metadata), packet headers, packet_in, packet_out
- Stateful Objects: Keep their state between packets
 - Tables
 - Externs
 - V1 architecture: Counters, Meters, Registers, Parser Value Sets, Selectors, etc.

P4-NetFPGA Extern Function library

- Implement platform specific functions
 - Black box to P4 program
- Implemented in HDL
- Stateless – reinitialized for each packet
- Stateful – keep state between packets
- Xilinx Annotations
 - `@Xilinx_MaxLatency()` – maximum number of clock cycles an extern function needs to complete
 - `@Xilinx_ControlWidth()` – size in bits of the address space to allocate to an extern function

P4-NetFPGA Extern Function library

- HDL modules invoked from within P4 programs

- Stateful atomic operations:

Extern	Description
R/W	Read or write state
RAW	Read, add to, or overwrite state
PRAW	Predicated version of RAW
ifElseRAW	Two RAWs, one each for when predicate is true or false
Sub	IfElseRAW with stateful subtraction capability

- Stateless Externs

Extern	Description
IP Checksum	Given an IP header, compute IP checksum
LRC	Longitudinal redundancy check, simple hash function
timestamp	Generate timestamp (granularity of 5 ns)

Using Externs in P4 – Resetting Counter

Packet processing pseudo code:

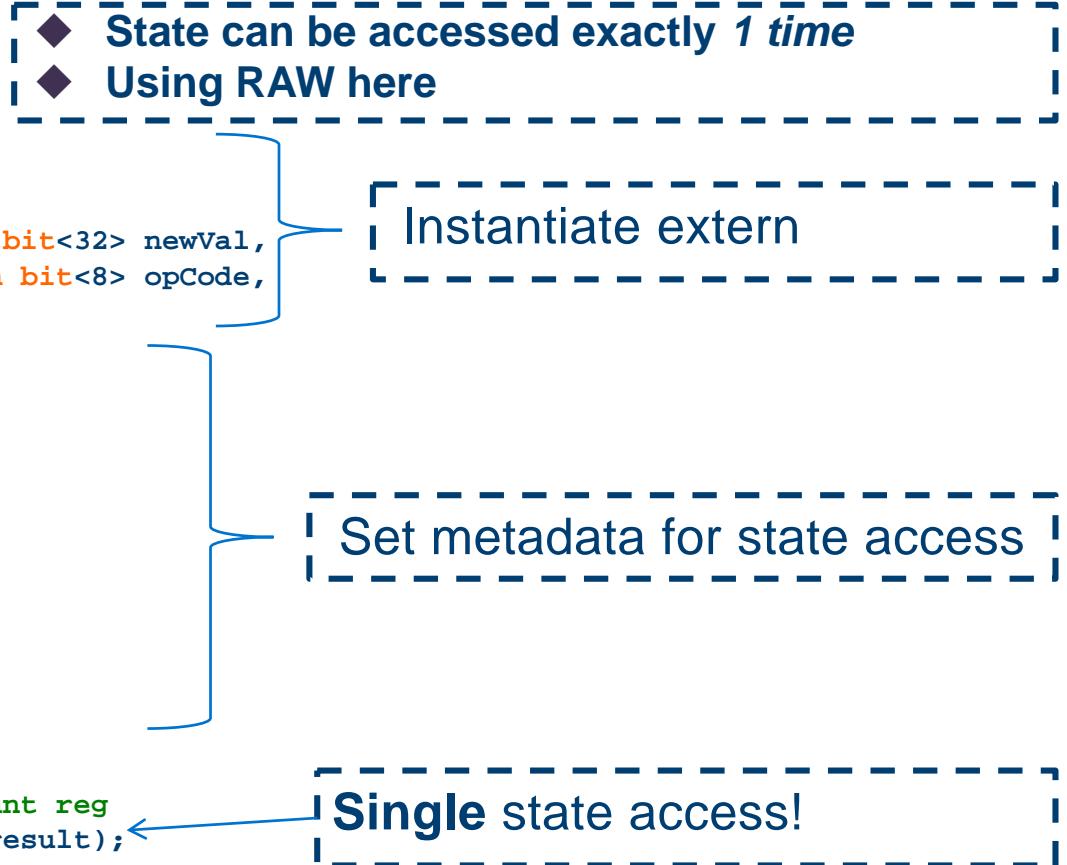
```
count[NUM_ENTRIES];  
  
if (pkt.hdr.reset == 1):  
    count[pkt.hdr.index] = 0  
else:  
    count[pkt.hdr.index]++
```

Using Externs in P4 – Resetting Counter

```
#define REG_READ    0
#define REG_WRITE   1
#define REG_ADD     2
// count register
@Xilinx_MaxLatency(1)
@Xilinx_ControlWidth(3)
extern void count_reg_raw(in bit<3> index, in bit<32> newVal,
                          in bit<32> incVal,in bit<8> opCode,
                          out bit<32> result);

bit<16> index = pkt.hdr.index;
bit<32> newVal; bit<32> incVal; bit<8> opCode;
if(pkt.hdr.reset == 1) {
    newVal = 0;
    incVal = 0; // not used
    opCode = REG_WRITE;
} else {
    newVal = 0; // not used
    incVal = 1;
    opCode = REG_ADD;
}

bit<32> result; // the new value stored in count reg
count_reg_raw(index, newVal, incVal, opCode, result);
```



FAQs

- **Can I apply a table multiple times in my P4 Program?**
 - No (except via resubmit / recirculate)
- **Can I modify table entries from my P4 Program?**
 - No (except for direct counters)
- **What happens upon reaching the reject state of the parser?**
 - Architecture dependent
 - Currently not supported by P4-SDNet
- **How much of the packet can I parse?**
 - Architecture dependent