

Software Techniques for Programmable Data Plane Virtualization

David Hancock (dhancock@cs.utah.edu)



THE
UNIVERSITY
OF UTAH

Flux Research Group

Motivation

Approach

Controller (“Control Plane Hypervisor”)

HyPer4.p4 (“Data Plane Hypervisor”)

Demo

Wrap-up / Questions

Programmable data planes

- How dynamic are they? Can we insert functionality on demand without disrupting important flows?

Multi-tenant environments

Virtualization

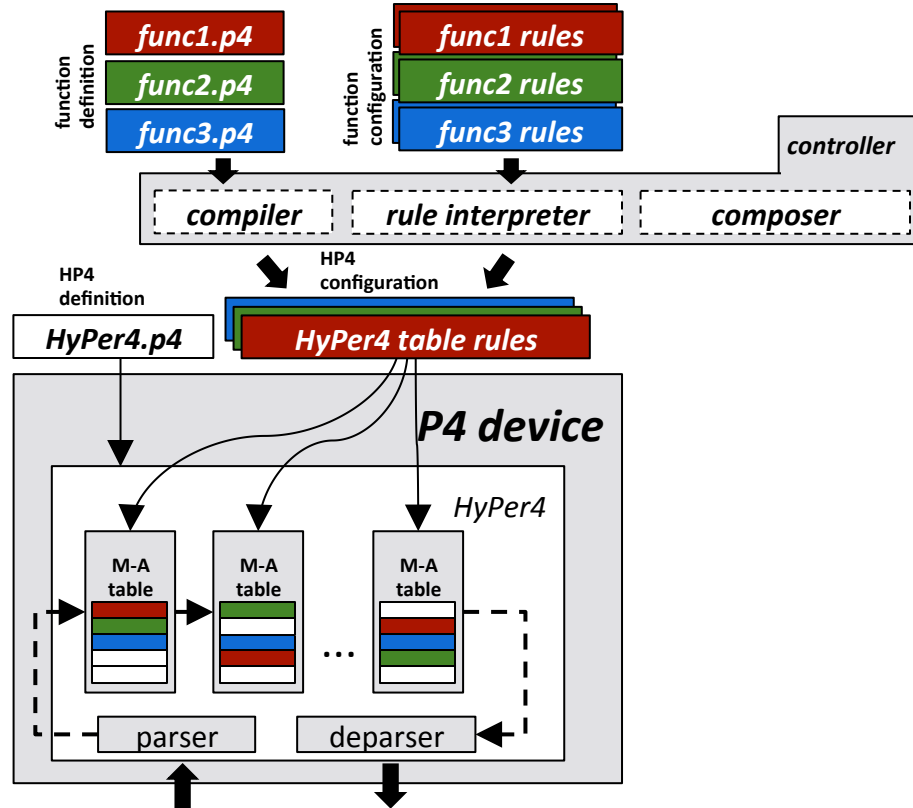
- Isolation
- Abstraction
 - Composition
 - Introspection
 - Standard functionality

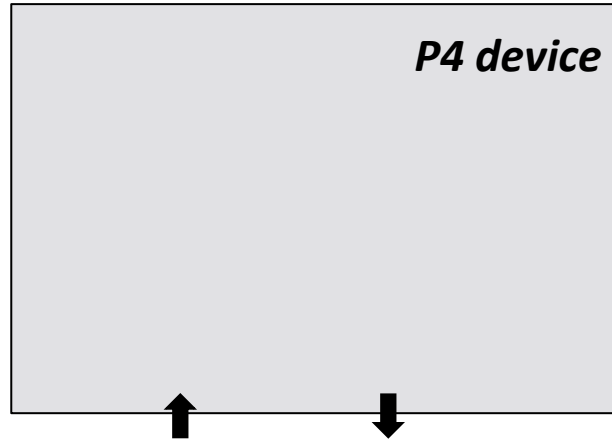
Software approach to expanding PDP device feature set

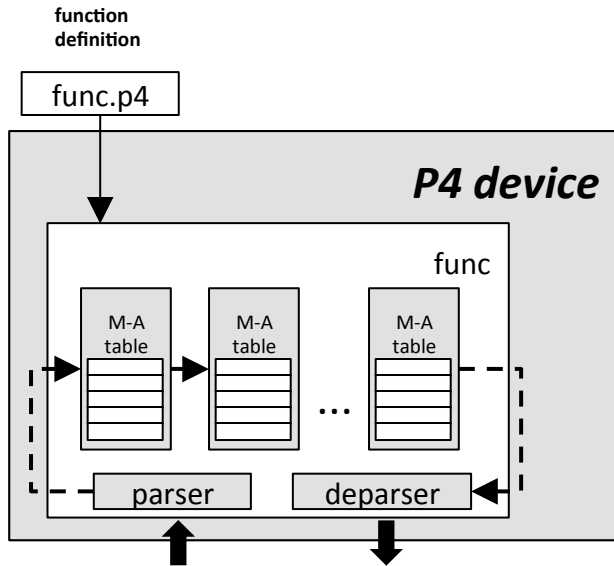
HyPer4.p4: “Data Plane Hypervisor”

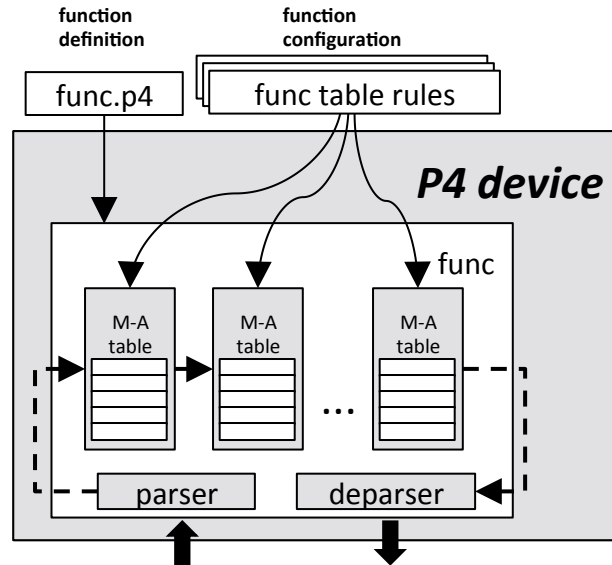
Controller: “Control Plane Hypervisor”

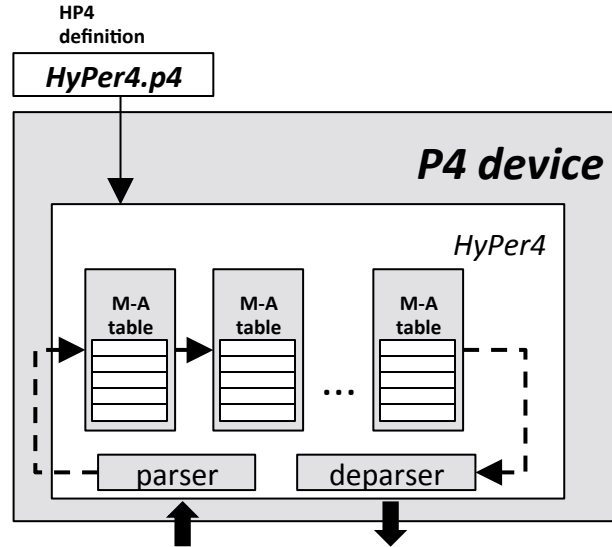
- Slice definition / provisioning
- Slice management
 - Compiler
 - Loader
 - Rule interpreter
 - Composer

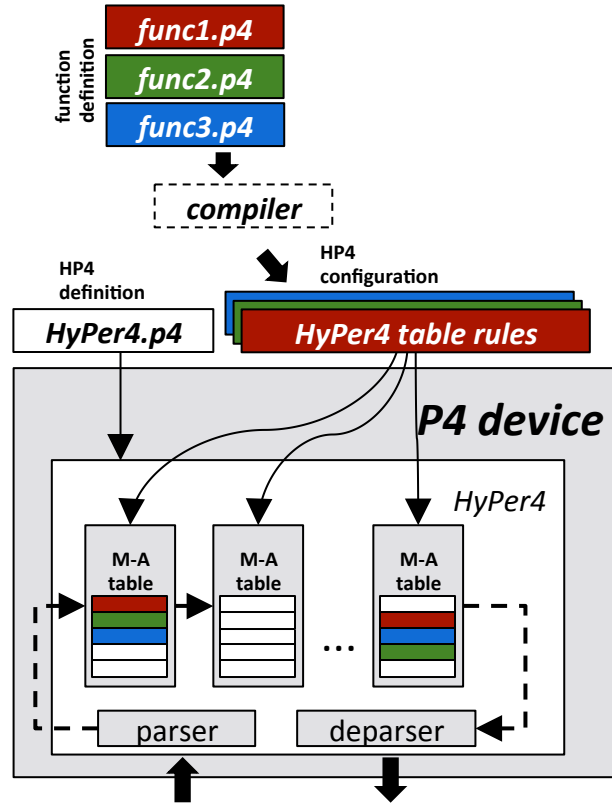


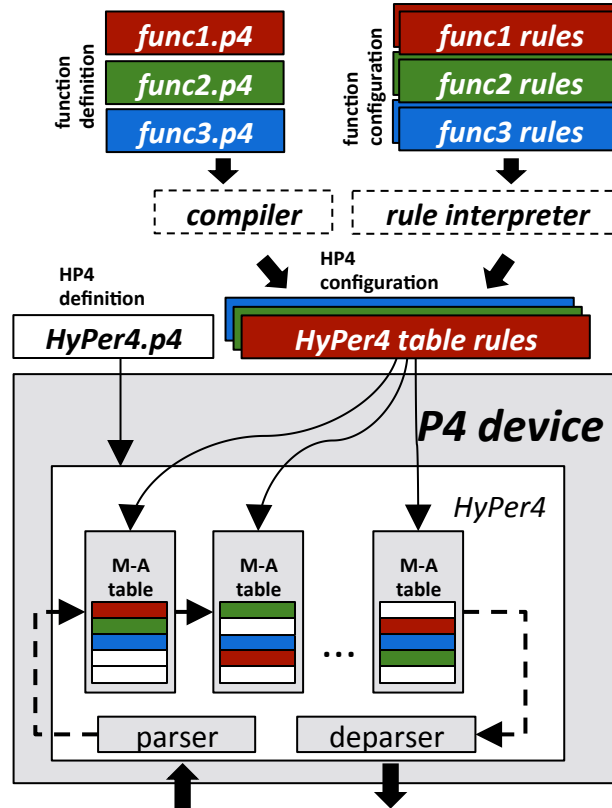


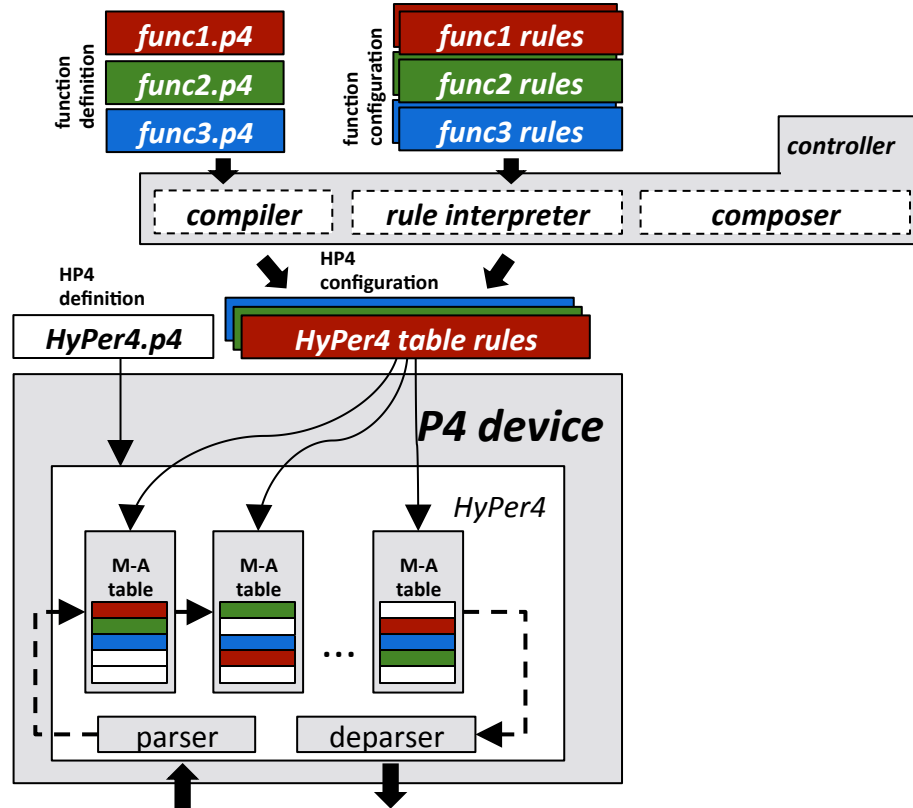










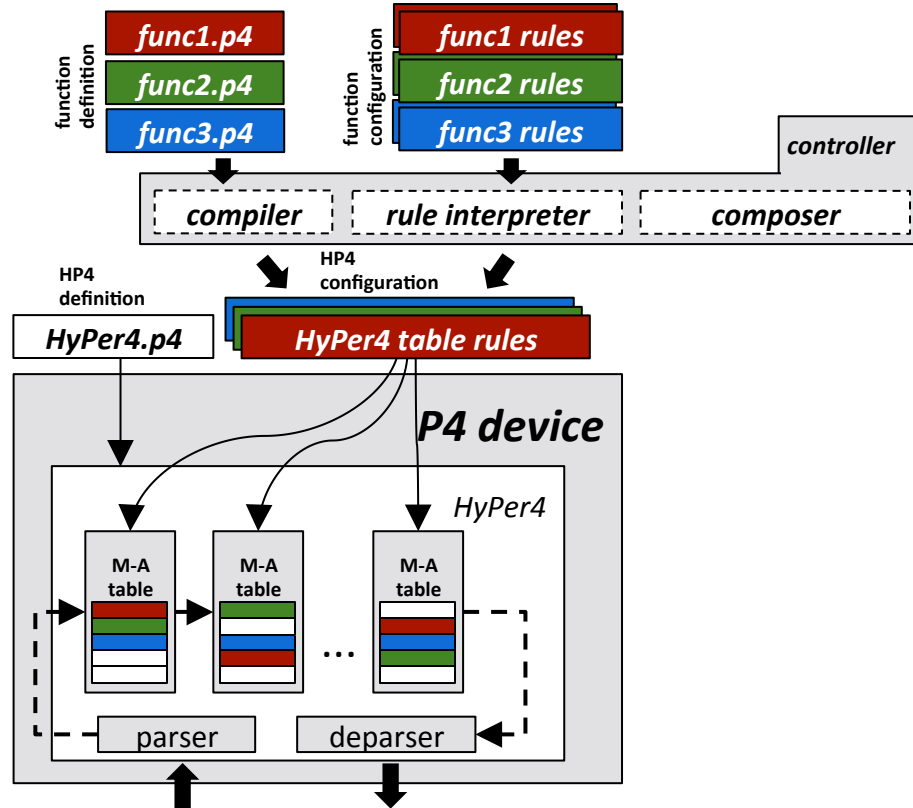


HyPer4.p4: Packet Tagging

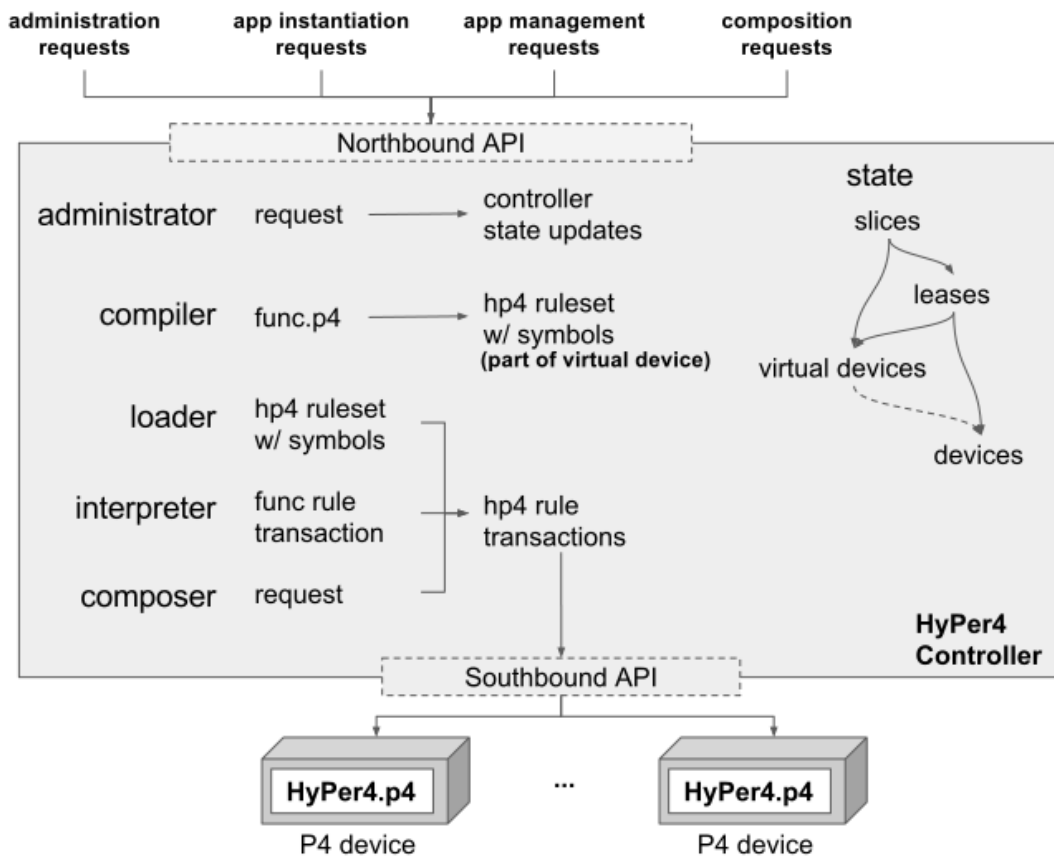
```
action a_set_context(vdev_ID, vingress) {
    modify_field(meta.vdev_ID, vdev_ID);
    modify_field(meta.vingress, vingress);
}

table tset_context {
    reads { standard_metadata.ingress_port : exact; }
    actions { a_set_context; }
}

control setup {
    if (meta.vdev_ID == 0) {
        apply(tset_context);
    }
    // ...
}
```



HyPer4 Controller



Administrator

- `create_device` `<dev name> <ip_addr> <port> \
<dev_type> <# entries> <ports>`
- `create_slice` `<slice name>`
- `grant_lease` `<slice name> <dev name> <entry limit> <ports>`
- `list_devices`
- `list_slices`
- `reset_device` `<dev name>`
- `revoke_lease` `<slice name> <dev name>`

Slice Manager

- `create_virtual_device` <p4 path> <vdev name>
- `destroy_virtual_device` <vdev name>
- `interpret` <vdev name> <API style> <command>
- `interpret_file` <vdev name> <API style> <command> \
<file>
- `lease insert | append | remove | replace` <vdev name> [args]
- `list_devs`
- `list_vdev`
- `list_vdevs`

Device class defines interface

- send_command
 - do_table_add
 - do_table_modify
 - do_table_delete
- various multicast-related methods

Device subclasses implement

- Bmv2-SSwitch
- Agilio

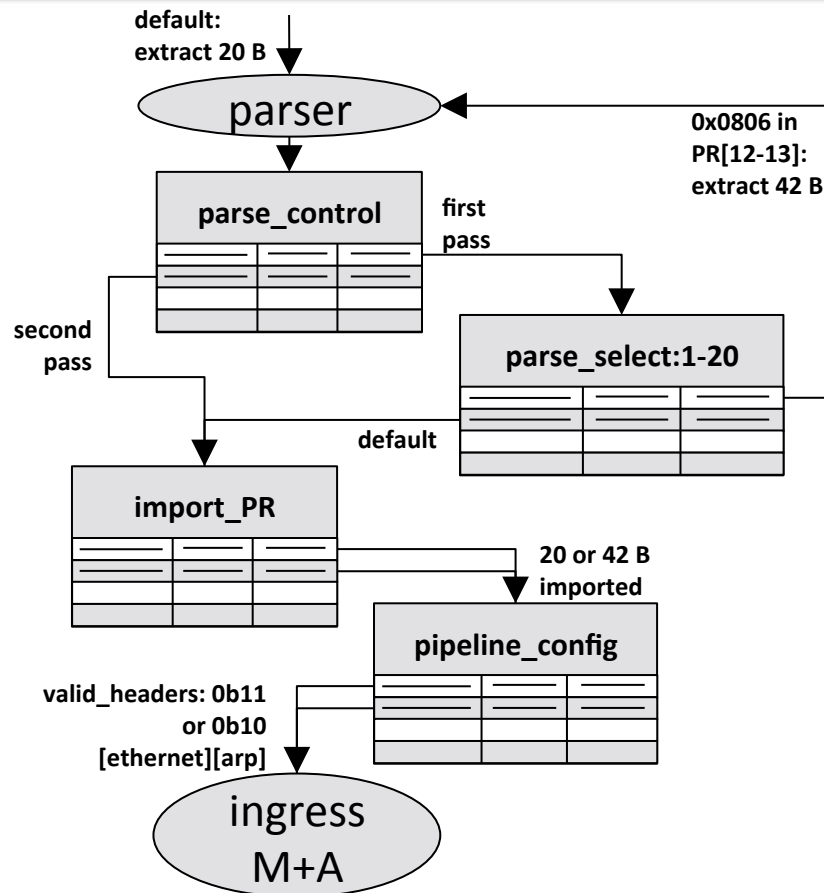
Task: P4 programs → “object-code-like” HyPer4 rulesets
→ rule interpreter guidance

Subtasks:

- metadata field addressing
 - All metadata fields consolidated into HyPer4.p4’s meta.data field
 - Reads/writes done with bitmasks and bitshifts applied to meta.data
- parsing
- matching
- actions
- control flow

Subtasks: Parsing

```
parser start {  
  extract(ethernet);  
  return select(ethernet.EtherType) {  
    0x0806: parse_arp;  
    default : ingress;  
  }  
}  
parser parse_arp {  
  extract(arp);  
  return ingress;  
}
```



Subtasks: Matching

table from myfunction.p4

```
mytable
match field : match type:
  myheader.myfield : exact
action list:
  myaction1
  myaction2
```

① compile myfunction

② manage myfunction

rule for myfunction

```
mytable
match values:
  [myfield match value]
selected action:
  myaction1
action parameters:
  [myaction1 parameters]
```

Controller

HP4C

mytable
translation
guide

rule
interpreter

vdev state

counterpart in HyPer4.p4

```
parsed_rep_[stage]_exact
match fields : match types:
  vdev_ID : exact
  parsed_rep.data : ternary
action list:
  init_action_state
```

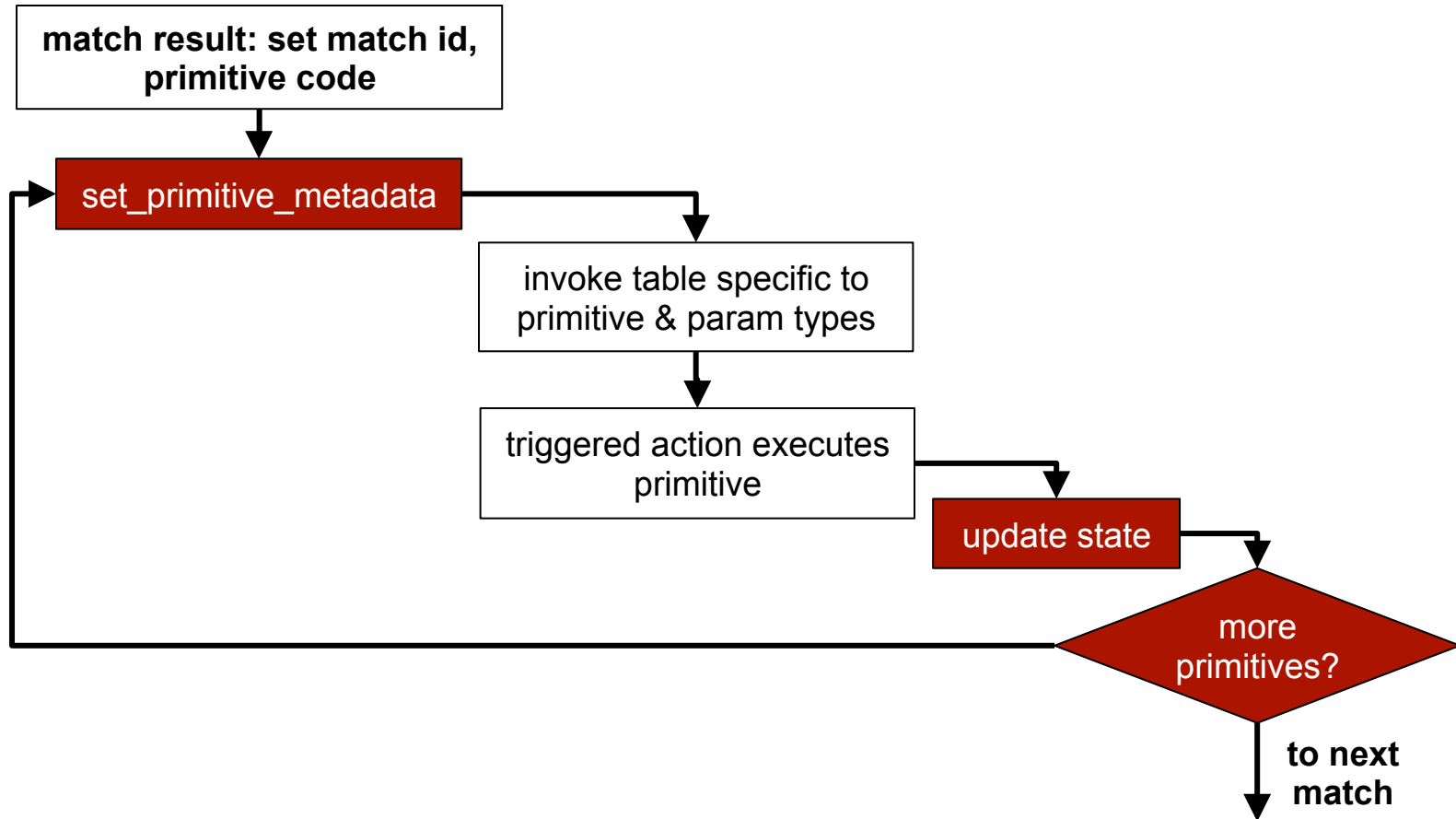
rule for HyPer4

```
[ ] parsed_rep_[stage]_exact
match values / bitmasks:
  [ ] [vdev_ID]
  [ ] [parsed_rep.data]
selected action:
  init_action_state
action parameters:
  [ ] [action ID] [ ] [match ID]
  [ ] [next table] [ ] [primitive type]
```

Exact matching vs arbitrary fields:

- Use ternary matching
 - rules include mask
- Strategy: apply mask to pr.data to isolate relevant bits

HyPer4.p4: Actions



“Object code”:

```
tset_pipeline_config a_set_pipeline :[vdev ID] 1:  
  [STDMETA_INGRESS_PORT_EXACT] 0x80[9x00s] 0
```

```
t_mod_41 mod_stdmeta_egressspec_stdmeta_ingressport :  
  [vdev ID] 11 4 0&&&0:[MAX PRIORITY]
```


“Object code”:

table (end of “setup”)	action	virtualized function ID	parse_state
tset_pipeline_config	a_set_pipeline	:[vdev ID]	1:
[STDMETA_INGRESS_PORT_EXACT]	0x80[9x00s]	0	egress handling mode
ID of 1 st match table	valid headers bitmap	(0b1000...: ethernet only)	

```
t_mod_41 mod_stdmeta_egressspec_stdmeta_ingressport :  
[vdev ID] 11 4 0&&&0:[MAX PRIORITY]
```

“Object code”:

```
tset_pipeline_config a_set_pipeline :[vdev ID] 1:  
  [STDMETA_INGRESS_PORT_EXACT] 0x80[9x00s] 0
```

table (modify_field)

action

```
t_mod_41 mod_stdmeta_egressspec_stdmeta_ingressport :  
[vdev ID] 11 4 0&&&0:[LOWEST PRIORITY]  
virtualized prim action ID match ID ternary match  
function ID subtype priority
```

“Object code”:

```
tset_pipeline_config a_set_pipeline :[vdev ID] 1:  
  [STDMETA_INGRESS_PORT_EXACT] 0x80[9x00s] 0
```

```
t_mod_41 mod_stdmeta_egressspec_stdmeta_ingressport :  
  [vdev ID] 11 4 0&&&0:[MAX PRIORITY]
```

Bmv2-SS API Commands:

```
table_add tset_pipeline_config a_set_pipeline 1 1 =>  
5 0x800000000000000000000000 0
```

```
table_add t_mod_41 mod_stdmeta_egressspec_stdmeta_ingressport  
1 11 4 0&&&0 => 2147483646
```

End of ingress

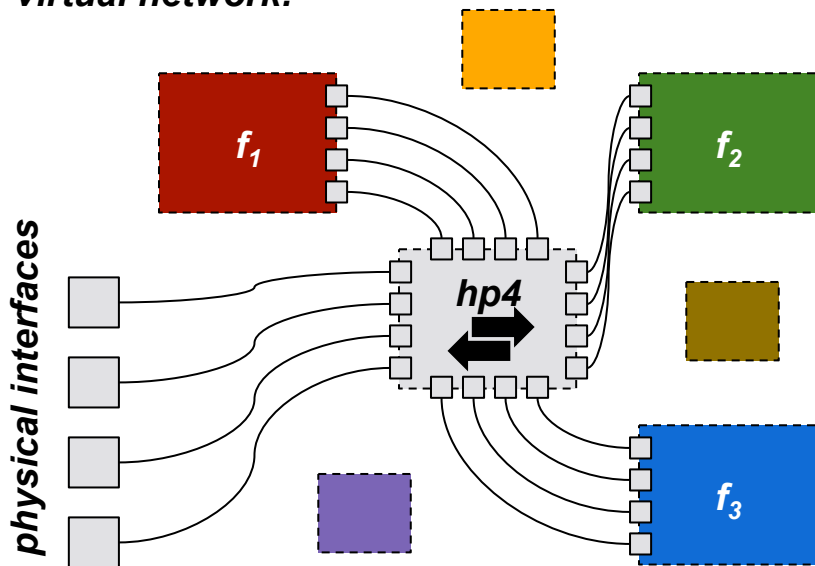
- Match on vgress
- Virtual forwarding? Set egress to ingress port, set virt forwarding flag

Egress

- Virt forwarding flag set?
 - Match on vgress
 - Set next vdev_ID, vingress
 - Recirculate with field list that includes key fields

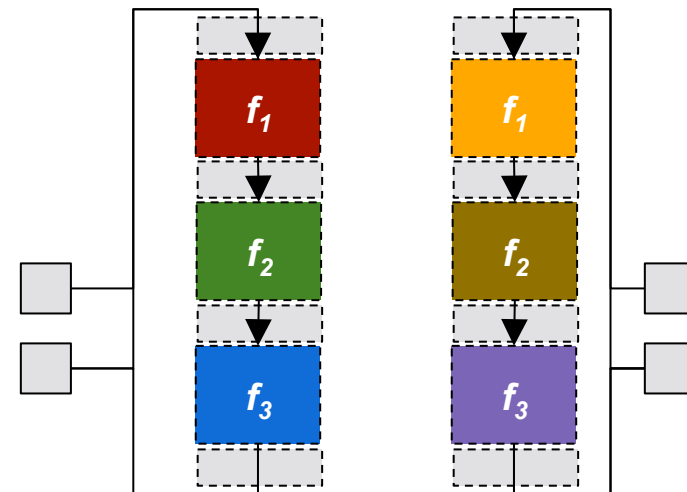
Composition Abstractions

virtual network:

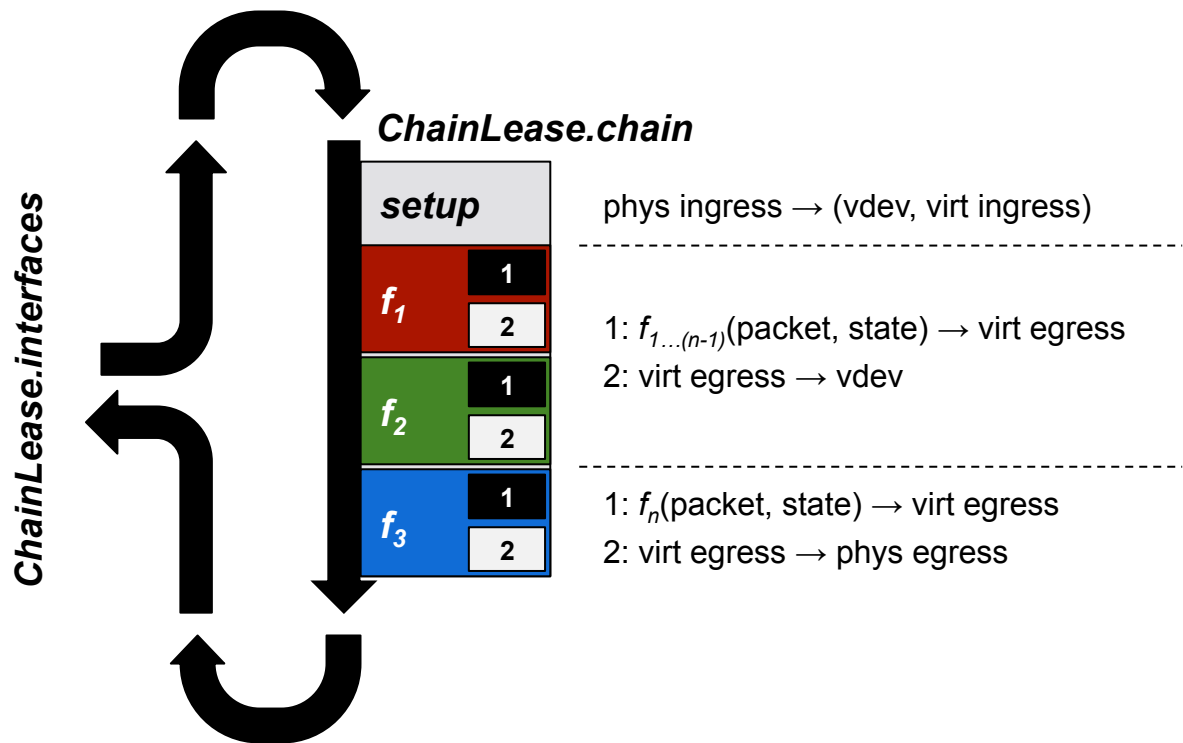


- Very flexible
- Complicated to implement and manage

vdev chains:



- Flexible
- Simpler to implement and manage



Types of Forwarding Abstractions

- return-to-sender: single port
- bump-in-the-wire [+ return-to-sender]: two ports
- switch: (packet, state) → egress: two+ ports

VDev sequence results in (sometimes unintentional) egress overrides

Resolution:

- ~~Enforce switch at the end~~
- Slice manager indicates preferences with respect to egress modifications for each vdev via controller command option:
 - always write to egress_spec
 - conditionally write to egress_spec (if void)
 - never write to egress_spec

Migration

Cloning

Maintain rulesets separately, attach to compatible vdevs on-demand

Automatically evolving functions

Linked functions (change to one begets [a]symmetrical change to another across the network)

Handle virtual networking

Egress filtering (avoid broadcast storms)

Handle checksum (e.g., IPv4)

Resize parsed representation (in case of `add_header / remove_header`)

Write `pr.data` back to parsed representation

P4 compiler / target restriction: ONE reference per table

- Want: arbitrary functionality at arbitrary points in the pipeline
- Strategy: many functionally redundant copies, differing in name only
- Want: configurable support for specified # stages, specified # primitives / stage
- Strategy: define .p4 for single stage, single primitive, add tokens directing processing tool to produce desired configuration

HyPer4.p4: Code Generation

```
[+ sloop +]
control stage[+X+] {
  match_[+X+](); // match.p4
  [+ nif +]
  if(meta_ctrl.stage_state != COMPLETE) {
    switch_primitivetype_[+X+][+Y+]();
    apply(tstg[+X+][+Y+]_update_state);
  }
  [+ endnif +]
}
[+ endsloop +]
```

```
control stage1 {
  match_1();
  if(meta_ctrl.stage_state != COMPLETE) {
    switch_primitivetype_11();
    apply(tstg11_update_state);
    if(meta_ctrl.stage_state != COMPLETE) {
      switch_primitivetype_12();
      apply(tstg12_update_state);
      // ...
    }
  }
  control stage2 {
    match_2();
    if(meta_ctrl.stage_state != COMPLETE) {
```

Performance Concerns

- May be addressed with custom targets
- May be addressed with more sensible default parsing behavior and streamlined match-action pipeline
- May be addressed by using programmable parser + fixed pipeline, or fixed parser + programmable pipeline
- Should not get in the way of exploring high level abstractions made possible by vPDPs

Motivation

Approach

Controller (“Control Plane Hypervisor”)

HyPer4.p4 (“Data Plane Hypervisor”)

- [git@gitlab.flux.utah.edu:dhankook/hp4-src.git](https://gitlab.flux.utah.edu/dhankook/hp4-src.git)
- README.md provides commit used for CoNEXT 2016 paper

Demo

Questions?

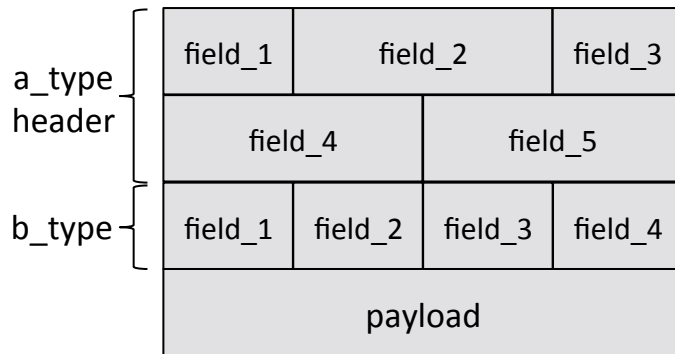
HyPer4.p4: Parsing

```
// not shown: header_type definitions
header a_t a;
header b_t b;

parser start {
  extract(a);
  return select(a.field_3) {
    0x01 : parse_b;
    default : ingress;
  }
}

parser parse_b {
  extract(b);
  return ingress;
}
```

parsed representation: a | a, b



P4 source includes references to named fields

Strategy:

- generic 8-bit header type
- ***ext[100]***: array of such headers

ext[0]	ext[1]	ext[2]	ext[3]
ext[4]	ext[5]	ext[6]	ext[7]
ext[8]	ext[9]	ext[10]	ext[11]
payload			

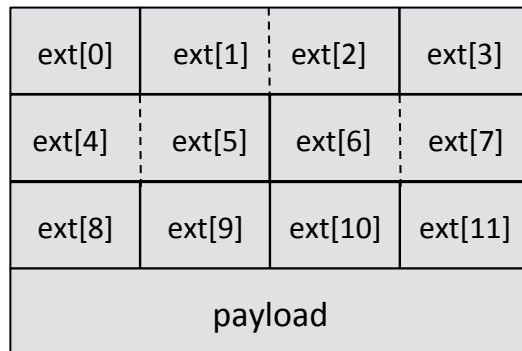
Strategy:

- generic 8-bit header type
- ***ext[100]***: array of such headers
- metadata fields: ***#bytes, parse_state***

ext[0]	ext[1]	ext[2]	ext[3]
ext[4]	ext[5]	ext[6]	ext[7]
ext[8]	ext[9]	ext[10]	ext[11]
payload			

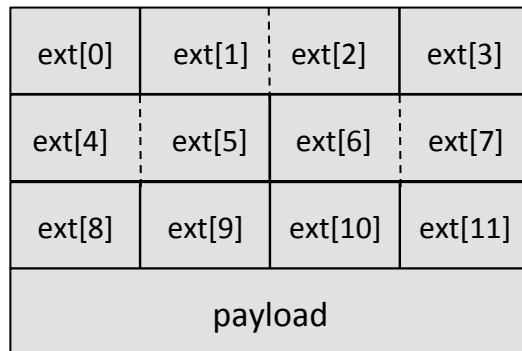
Strategy:

- generic 8-bit header type
- **ext[100]**: array of such headers
- metadata fields: **#bytes, parse_state**
- Parser: extract default -> 100 bytes
 - **#bytes** guides parser



Strategy:

- generic 8-bit header type
- **ext[100]**: array of such headers
- metadata fields: **#bytes, parse_state**
- Parser: extract default -> 100 bytes
 - **#bytes** guides parser
- Complex parsing logic moved to Ingress:
parse_state may trigger update of **#bytes** + resubmit



Strategy:

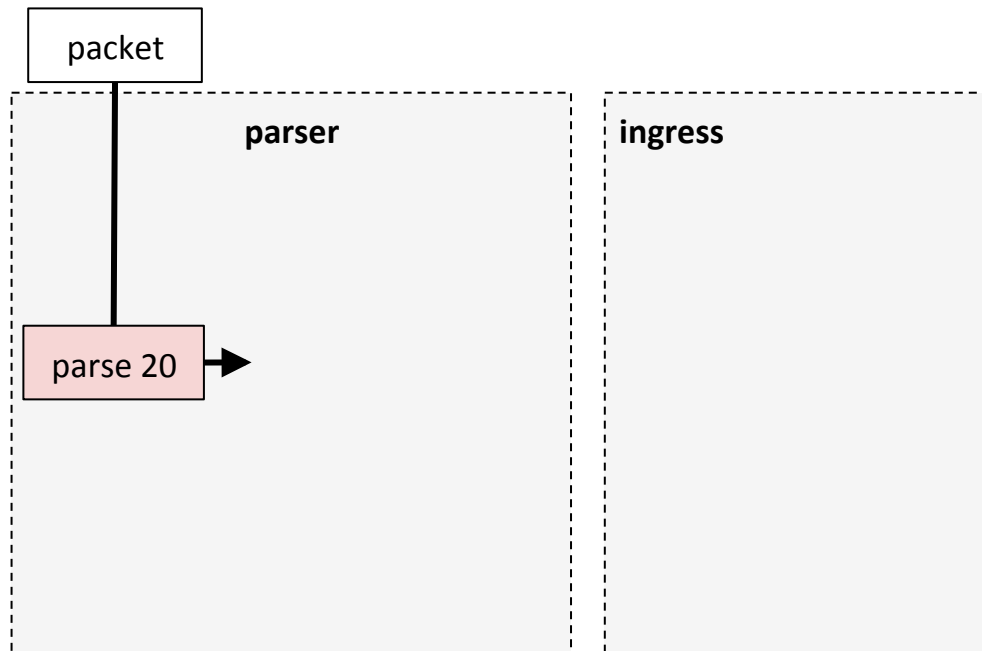
- generic 8-bit header type
- **ext[100]**: array of such headers
- metadata fields: **#bytes, parse_state**
- Parser: extract default -> 100 bytes
 - **#bytes** guides parser
- Complex parsing logic moved to Ingress:
parse_state may trigger update of **#bytes** + resubmit

ext[0]	ext[1]	ext[2]	ext[3]
ext[4]	ext[5]	ext[6]	ext[7]
ext[8]	ext[9]	ext[10]	ext[11]
payload			

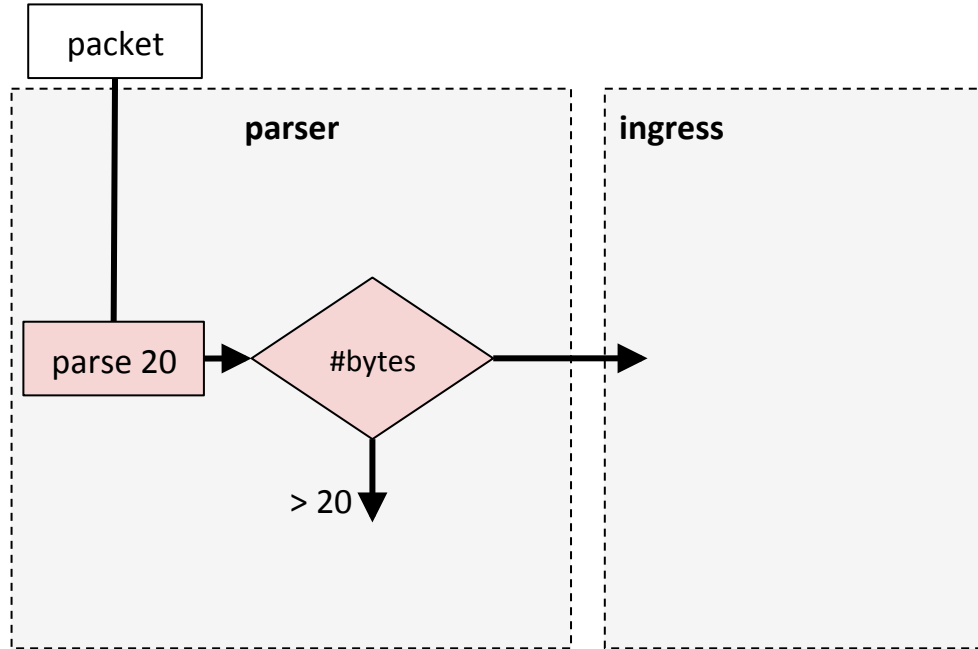
parsed representation: all valid elements of ext

- Inconvenient to work with; copy to single very wide metadata field

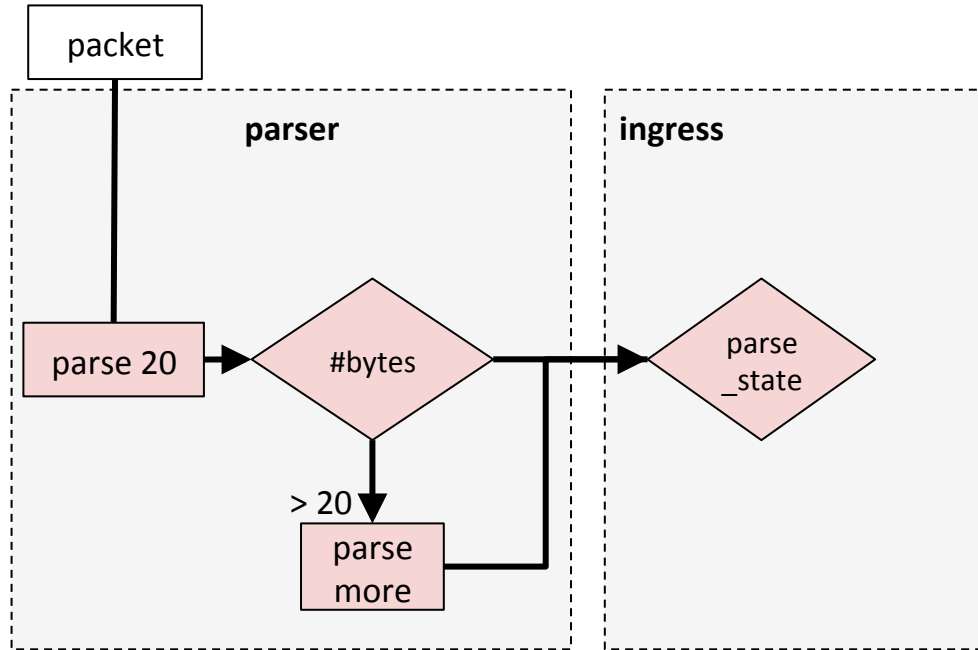
HyPer4.p4: Parsing



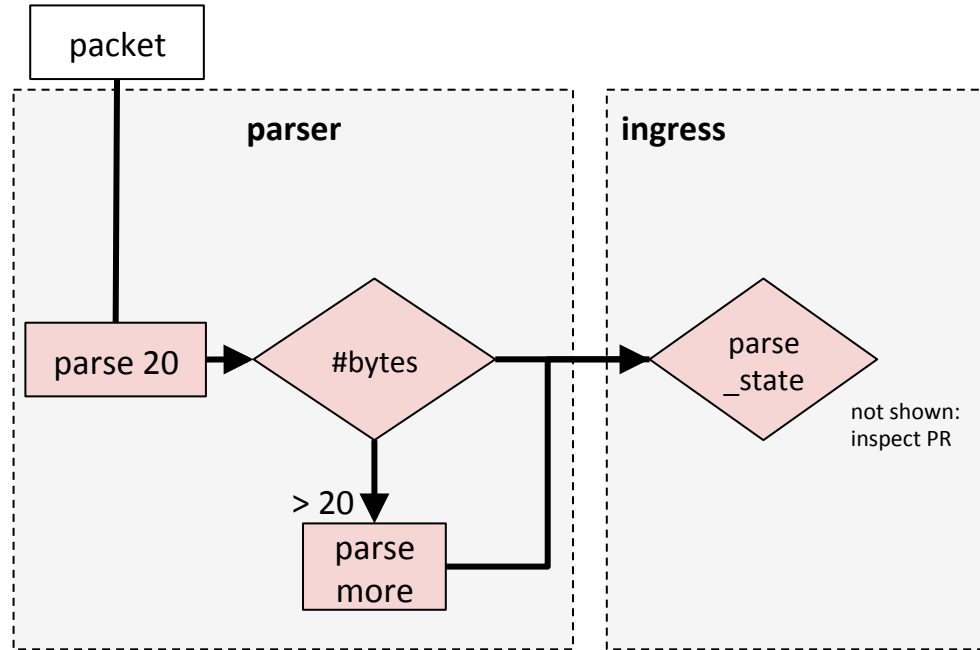
HyPer4.p4: Parsing



HyPer4.p4: Parsing



HyPer4.p4: Parsing



HyPer4.p4: Parsing

