



Programmers Guide

Craig Riecke

with

Others

Draft of April 18, 2016

Contents

1	Quick Start	4
1.1	Installation	4
1.2	What Do You Get With Frenetic VM?	5
1.3	An Attempt at Hello World	5
1.4	A Repeater	7
1.5	Running The Repeater Application	9
1.6	Summary	11
2	NetKAT	12
2.1	Introduction to OpenFlow	13
2.2	OpenFlow is Difficult	16
2.3	Predicates	17
2.4	Policies	21
2.5	Commands and Hooks	22
2.5.1	The packet.in Hook	24
2.5.2	The pkt.out Command	27
2.5.3	Buffering	29
2.6	OpenFlow Constructs Not Supported by NetKAT	30
3	NetKAT Principles	31
3.1	Efficient SDN	31
3.2	Combining NetKAT Policies	33
3.3	Keeping It Stateless	38
3.4	Summary	46
4	Learning Switch	47
4.1	Design	47
4.2	A First Pass	49
4.3	A More Efficient Switch	51
4.4	Timeouts and Moves	53

5	Proxy ARP	54
5.1	The ARP Protocol	54
5.2	Design	54
5.3	Snooping on ARP Requests and Replies	54
5.4	Replying	54
5.5	Maintaining State Across Restarts	54
6	Handling Vlans	55
6.1	Vlan Uses	55
6.2	Design	55
6.3	Maintaining Separate Vlan Tables	55
6.4	Tagging and Untagging	55
7	Gathering Statistics	56
8	SDN Development Tools	57
8.1	Tmux	57
8.2	Open VSwitch Utilities	59
8.3	TCPDump	59
8.4	Mininet Network Modelling	59
9	Network Address Translation	60
9.1	Why Do We Need NAT?	60
9.2	Design	60
10	Spanning Tree Alternatives	61
10.1	What's Wrong with STP Protocols?	61
10.2	Design	61
10.3	Calculating Shortest Paths	61
10.4	Core/Edge Separation	61
11	Routing	62
11.1	Design	62
11.2	Configuring Route Tables	62
11.3	Using Link State	62
12	Modularization	63
12.1	Sharing Actions in Rules and Packet Outs	63
12.2	Subclassing	63
12.3	Multiple Network Apps	63
12.4	Clustering	63

13 Frenetic REST API	64
13.1 REST URL's	64
13.2 Incoming, Northbound Events	64
13.3 Messages	64
14 Frenetic/NetKAT Reference	65
14.1 Predicates	65
14.2 Policies	65
14.3 Event Hooks	65
14.4 Compiler Directives	65
14.5 Frenetic Command Line	65
15 Productionalizing	66
15.1 Installing Frenetic on Bare Metal Linux	66
15.2 Control Scripts	66
15.3 Logging	66

Chapter 1

Quick Start

In this book, you will use Frenetic to create a full-programmable network. For the moment, let's assume you're familiar with Software Defined Networking and the OpenFlow protocol, and just dive right in. (If you're not, don't worry! We'll introduce some bedrock concepts in the next chapter and explain everything that happened here.)

1.1 Installation

There are several ways to get started with Frenetic, but the easiest is to use Frenetic VM. Frenetic itself only runs on Linux, but the Frenetic VM will run on any host system that supports VirtualBox, including Windows, Mac OS X and practically any version of Linux itself. Keeping Frenetic in its own VM will keep your own system clean and neat. Later on, if you want to install Frenetic on a bare metal Ubuntu Linux server or network device, you can use the instructions in 15.1.

- Install VirtualBox from <https://www.virtualbox.org/wiki/Downloads>. Use the latest version platform package appropriate for your system.
- Install Vagrant at <http://www.vagrantup.com/downloads>. Vagrant automates the process of building VM's from scratch, and Frenetic VM uses it to build its own environment. This is more reliable than downloading a multi-gigabyte VM file.
- Install Frenetic VM from <https://github.com/frenetic-lang/frenetic-vm>. You can simply use the Download Zip button and unzip to an appropriate directory on your system, like `frenetic-vm`. Then from a terminal or command prompt:

```
$ cd /path/to/frenetic-vm
$ vagrant up
... lots of text
```

The build process may take 15 minutes to an hour, depending on the speed of your system and Internet connection.

1.2 What Do You Get With Frenetic VM?

At the end of the process you will have a working copy of Frenetic with lots of useful open source infrastructure:

Mininet software builds a test network inside of Linux. It can simulate a topology with many switches and hosts. Writing a network application and throwing it into production is ... well, pretty risky, but running it on Mininet first can be a good test for how it works beforehand. We'll use it throughout this book.

Wireshark captures and analyzes network traffic. It's a great debugging tool, and very necessary for sifting through large amounts of data.

Frenetic . This layer provides an easy-to-use programmable layer on top of ODL. Its main job is to shuttle OpenFlow messages between ODL and your application, and to translate the language NetKAT into OpenFlow flow tables. We'll see the differences between the two as we go.

Hmmm, that's a lot of software - what do *you* bring to the table? You write your network application in Python, using the Frenetic framework. As you'll see, it's quite easy to build a network device from scratch, and easy to grow it organically to fit your requirements. Python is fairly popular, and knowing it will give you a head start into Frenetic programming. But if you're a Python novice that's OK. As long as you know one object-oriented language fairly well, you should be able to follow the concepts. We'll introduce you to useful Python features, like list comprehensions, as we go.

1.3 An Attempt at Hello World

So let's dive right in. We'll set up a Mininet network with one switch and two hosts. First you should work from the directory where you installed Frenetic VM.

```
$ cd /path/to/frenetic-vm
```

Then start up the VM:

```
$ vagrant up
ringing machine 'default' up with 'virtualbox' provider...
==> default: Clearing any previously set forwarded ports...
==> default: Clearing any previously set network interfaces...
==> default: Preparing network interfaces based on configuration...
default: Adapter 1: nat
==> default: Forwarding ports...
default: 22 => 2222 (adapter 1)
==> default: Running 'pre-boot' VM customizations...
```

```
==> default: Booting VM...
==> default: Waiting for machine to boot. This may take a few minutes...
    default: SSH address: 127.0.0.1:2222
    default: SSH username: vagrant
    default: SSH auth method: private key
    default: Warning: Connection timeout. Retrying...
==> default: Machine booted and ready!
==> default: Checking for guest additions in VM...
==> default: Setting hostname...
==> default: Mounting shared folders...
    default: /vagrant => /Users/cr396/frenetic-vm
    default: /home/vagrant/src => /Users/cr396/frenetic-vm/src
==> default: Machine already provisioned. Run `vagrant provision` or...
==> default: to force provisioning. Provisioners marked to run always...
```

Then log in to the VM. At this point your command prompt will change to `vagrant@frenetic` to distinguish it from your host machine.

```
$ vagrant ssh
Welcome to Ubuntu 14.04.2 LTS (GNU/Linux 3.16.0-30-generic x86_64)

 * Documentation:  https://help.ubuntu.com/
Last login: Tue Oct  6 10:35:06 2015 from 10.0.2.2
vagrant@frenetic:~$
```

So you are now working inside an Ubuntu-based VM. You don't really need to know Ubuntu, but just know that Mac OS and Windows commands won't necessarily work here.

Let's start up a Mininet network with one switch and two nodes.

```
vagrant@frenetic:~$ sudo mn --topo=single,2 --controller=remote
*** Creating network
*** Adding controller
Unable to contact the remote controller at 127.0.0.1:6633
*** Adding hosts:
h1 h2
*** Adding switches:
s1
*** Adding links:
(h1, s1) (h2, s1)
*** Configuring hosts
h1 h2
*** Starting controller
c0
*** Starting 1 switches
s1 ...
*** Starting CLI:
mininet>
```

The prompt changes to `mininet>` to show your working in Mininet. The error message `Unable to contact controller at 127.0.0.1:6633` looks a little ominous, but not fatal.

You now have an experimental network with two hosts named `h1` and `h2`. To see if there's connectivity between them, use the command `consoleh1 ping h2` which means "On host `h1`, ping the host `h2`."

```
mininet> h1 ping h2
PING 10.0.0.2 (10.0.0.2) 56(84) bytes of data.
From 10.0.0.1 icmp_seq=1 Destination Host Unreachable
From 10.0.0.1 icmp_seq=2 Destination Host Unreachable
From 10.0.0.1 icmp_seq=3 Destination Host Unreachable
^C
--- 10.0.0.2 ping statistics ---
6 packets transmitted, 0 received, +3 errors, 100\% packet loss, time 5014ms
pipe 3
```

The ping gets executed over and over again, but it's clearly not working. So we press CTRL-C to stop and quit out of Mininet:

```
mininet> quit
```

So by default, hosts can't talk over the network to each other. We're going to fix that by writing a *network application*. Frenetic will act as the controller on the network, and the network application tells Frenetic how to act.

1.4 A Repeater

You write your network application in Python, using the Frenetic framework. Mininet is currently running in our VM under its own terminal window, and we can leave it like that. We'll do our programming in another window, so start up another one and log into our VM:

```
$ cd /path/to/frenetic-vm
$ vagrant ssh
Welcome to Ubuntu 14.04.2 LTS (GNU/Linux 3.16.0-30-generic x86_64)

 * Documentation:  https://help.ubuntu.com/
Last login: Tue Oct  6 10:35:06 2015 from 10.0.2.2
vagrant@frenetic:~$
```

Get the sample code from the Frenetic Github repository:

```
vagrant@frenetic:~$ mkdir tutorial
vagrant@frenetic:~$ cd tutorial
```

Now we'll write our first network application. Or rather ...you can follow along in the sample code for this app.

```
vagrant@frenetic:~$ git clone https://github.com/frenetic-lang/manual
Cloning into 'manual'...
remote: Counting objects: 102, done.
remote: Total 102 (delta 0), reused 0 (delta 0), pack-reused 102
Receiving objects: 100% (102/102), 1.46 MiB | 728.00 KiB/s, done.
Resolving deltas: 100% (47/47), done.
Checking connectivity... done.
vagrant@frenetic:~$ cd manual/programmers_guide
vagrant@frenetic:~/manual/programmers_guide$ cd code/quick_start
```

The following code is in `code/quick_start/repeater.py`:

```
1 import sys
2 sys.path.append('/home/vagrant/src/frenetic/lang/python')
3 import frenetic
4 from frenetic.syntax import *
5
6 class RepeaterApp(frenetic.App):
7
8     client_id = "quick_start"
9
10    def connected(self):
11        self.update( id >> SendToController("repeater_app") )
12
13    def packet_in(self, dpid, port_id, payload):
14        out_port_id = 2 if port_id == 1 else 1
15        self.pkt_out(dpid, payload, [ Output(Physical(out_port_id)) ] )
16
17 app = RepeaterApp()
18 app.start_event_loop()
```

Lines 1-4 are pretty much the same in every Frenetic network application. Similarly, lines 17-18 are similar in most cases. The meat of the application is an object class named `RepeaterApp`, whose base class is `frenetic.App`. A frenetic application can hook code into different points of the network event cycle. In our `Repeater` network app, the only two events we're interested in here are `connected`, which is fired when a switch connects for the first time to Frenetic, and `packet_in`, which is fired every time a packet bound for a controller arrives.

The code in `connected` merely directs the switch to send all packets to our application. The interesting code is in `packet_in` and implements a *repeater*. A repeater is the oldest

type of network device, and is sometimes called a *hub*. In a repeater, if a packet enters on port 1, it should get copied out to port 2. Conversely, if a packet enters on port 2, it should get copied out to port 1. If there were more ports in our switch, we'd write a more sophisticated repeater – one that outputs the packet to all ports except the one on which it arrived (called the *ingress port*).

`pkt_out` is a method provided by Frenetic to actually send the packet out the switch. It takes three parameters: a switch, a packet, and a policy. Here the policy sends the packet out to port `out_port_id`.

1.5 Running The Repeater Application

So let's get this running in a lab setup. Three programs need to be running: Mininet, Frenetic, and our new Repeater application. For now, we'll run them in three separate command lines, having typed `vagrant ssh` in each to login to the VM.

In the first terminal window, we'll start up Frenetic:

```
vagrant@frenetic:~src/frenetic$ ./frenetic.native http-controller \  
> --verbosity debug  
[INFO] Calling create!  
[INFO] Current uid: 1000  
[INFO] Successfully launched OpenFlow controller with pid 3062  
[INFO] Connecting to first OpenFlow server socket  
[INFO] Failed to open socket to OpenFlow server: (Unix.Unix_error...  
[INFO] Retrying in 1 second  
[INFO] Successfully connected to first OpenFlow server socket  
[INFO] Connecting to second OpenFlow server socket  
[INFO] Successfully connected to second OpenFlow server socket
```

In the second, we'll start up Mininet with the same configuration as before:

```
vagrant@frenetic:~$ sudo mn --topo=single,2 --controller=remote  
*** Creating network  
*** Adding controller
```

The following will appear in your Frenetic window to show a connection has been made:

```
[INFO] switch 1 connected  
[DEBUG] Setting up flow table  
+-----+  
| 1 | Pattern | Action |  
|-----|  
|           |           |  
+-----+
```

And in the third, we'll start our repeater application:

```
vagrant@frenetic:~/manual/code/quick_start$ python repeater.py
Starting the tornado event loop (does not return).
```

The following will appear in the Frenetic window.

```
[INFO] GET /version
[INFO] POST /quick_start/update_json
[INFO] GET /quick_start/event
[INFO] New client quick_start
[DEBUG] Installing policy
drop | port := pipe(repeater_app) | port := pipe(repeater_app)
[DEBUG] Setting up flow table
+-----+
| 1 | Pattern | Action |
|-----|
|           | Output(Controller(128)) |
+-----+
```

And finally, we'll pop over to the Mininet window and try our connection test once more:

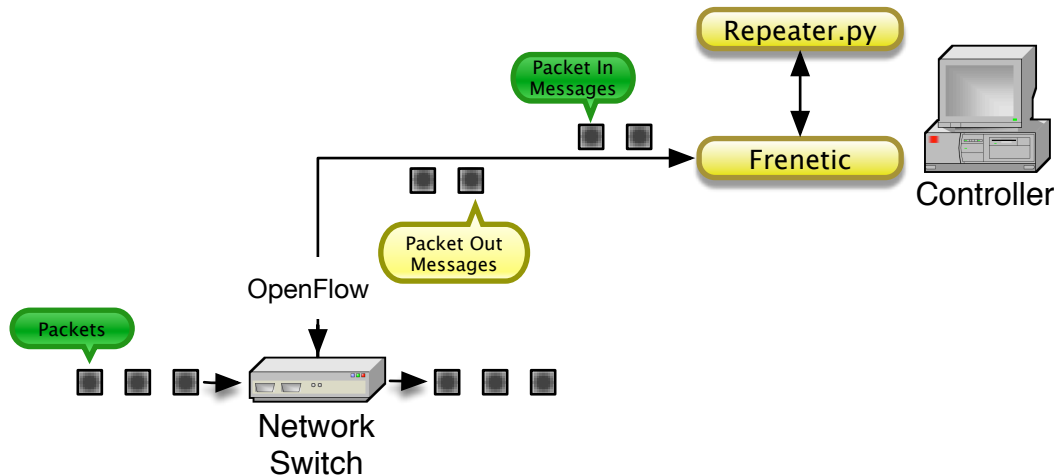
```
mininet> h1 ping h2
PING 10.0.0.2 (10.0.0.2) 56(84) bytes of data.
64 bytes from 10.0.0.2: icmp_seq=1 ttl=64 time=149 ms
64 bytes from 10.0.0.2: icmp_seq=2 ttl=64 time=97.2 ms
64 bytes from 10.0.0.2: icmp_seq=3 ttl=64 time=88.7 ms
```

Ah, much better! Our pings are getting through. You can see evidence of this in the Frenetic window:

```
[INFO] GET /quick_start/event
[INFO] POST /pkt_out
[DEBUG] SENDING PKT_OUT
[INFO] GET /quick_start/event
[INFO] POST /pkt_out
[DEBUG] SENDING PKT_OUT
[INFO] GET /quick_start/event
[INFO] POST /pkt_out
[DEBUG] SENDING PKT_OUT
```

1.6 Summary

You now have a working SDN, or Software Defined Network! Like much software, it works in layers:



1. At the bottom is your switches and wires. In our lab setup, Mininet and OpenVSwitch is a substitute for this layer.
2. In the middle is Frenetic. It talks the OpenFlow protocol to the switches (or to Mininet) – this is called the Southbound interface. It also accepts its own language called NetKAT from network applications – this is called the Northbound interface.
3. At the very top is your network application, which you write in Python. It defines how packets are dealt with.

We wrote a very simple network application that emulates a network repeater. It responds to the `packet_in` event coming from the switches through Frenetic when a packet arrives at the switch. And it sends the `pkt_out` message to send the packet back out through Frenetic to the switch. Frenetic-vm makes installing and testing all the pieces straightforward. When you're done, your network application can be deployed to a real production network.

Obviously you can do much more than just simple repeating with SDN! We'll cover that next with some background on OpenFlow and NetKAT, the underlying language of Frenetic.

Chapter 2

NetKAT

Software Defined Networking, or SDN, is a huge paradigm shift in the computing world. Traditional networking involves expensive, proprietary “boxes” from major vendors, plugging them in, configuring them, and hoping they meet your needs. But traditional networking suffers from these maladies:

- The devices are flexible only within narrow configuration parameters. Special requirements, like preventing certain kinds of devices from mobility, or configuring the spanning tree to prefer certain paths, are either impossible or expensive.
- While the devices are powered by software, there’s no workable way to examine the underlying code or prove it’s correct.
- Upgrades tend to be the forklift-variety, since mixing and matching old and new hardware is a dicey proposition ... not to mention mixing hardware from different vendors.
- Configuration is not easily automated. Solutions are often proprietary, require special programming languages, and are not interchangeable. Because of this, modern data center virtualization is more difficult.
- Adding support for new protocols is slow and expensive.

With SDN, data centers can step off the proprietary network treadmill. It’s a shift similar to the personal computer revolution of the 1980’s. Before that, IBM and similar mainframes controlled the computer space, and required the same kinds of forklift upgrades networks do. The IBM Personal Computer opened up the architecture to competitors, who could then design and build extensions that made it more useful. This created a snowball effect, with hardware extensions like the mouse and Ethernet cards opening the way for new software like Microsoft Windows and the Netscape browser.

SDN opens network devices in a similar manner, allowing them to be extended and manipulated in interesting, user-defined ways. Although the term SDN has been hijacked to mean many things, it most often refers to OpenFlow-enabled network software and

devices. OpenFlow is an open protocol defined by the Open Network Foundation that manipulates the control plane of a network intermediary.

Frenetic is an OpenFlow controller, meaning it talks the OpenFlow protocol to network intermediaries. In turn, it exposes an API that can be used to write network programs easily. It works with any network intermediary that understands the OpenFlow 1.0 protocol – both hardware devices like the HP 2920 and software “devices” like Open vSwitch. So let’s take a brief look at OpenFlow itself.

2.1 Introduction to OpenFlow

Every network device – from the lowliest repeater, to firewalls and load balancers, all the way up to the most complex router – has two conceptual layers:

The data plane performs actions on packets. It manipulates headers, copies packets to outgoing (or egress) ports, or drops packets. It consults control tables - MAC address tables, ARP caches, OSPF tables, etc. - to guide it.

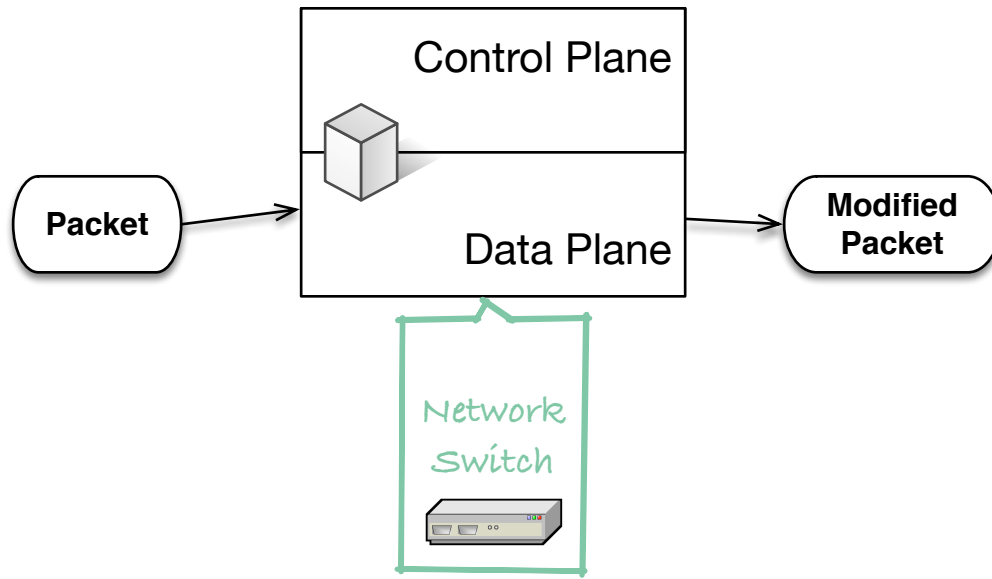
The control plane manipulates the control tables themselves. People, in turn, may manipulate the control plane through configuration software. Or packets may do it: specialized ones like OSPF neighbor exchange, ARP requests, or just examining plain ol’ packets themselves. But they never actually touch the packets.

This separation is only conceptual. You’d be hard pressed to open a network device, point to a chip and say, “That’s the data plane.” It helps in understanding a device, though, because they have different goals:

The data plane ’s job is to move data as quickly as possible. It relies more on fast table lookups than complex algorithms.

The control plane works more flexibly, yet conservatively. Control tables should not change often, and when they do, appropriate checks and balances should be applied to ensure the data plane keeps working. For example, the Spanning Tree Protocol (or STP) ensures packets are routed along the shortest path with no loops. Calculating the spanning tree is the control plane’s job, as its complex. But once that’s calculated, the data plane can use it to forward packets quickly.

A traditional network device looks like this. The control plane is closed and contained fully within the box.



The OpenFlow protocol makes the control plane *programmable*. Rather than relying on the entire program being inside the box, you write a program that advises the control plane and runs outside the box. It's like an advisor that makes arbitrarily complex control table manipulations. The programmable piece runs on any standard computer, and is collectively called the *controller*.

The controller can be written in any language and run on any computer ... the only requirement is it must speak the OpenFlow protocol to the device. You can think of these two pieces working in a tandem through an OpenFlow conversation:

Device: I just got a packet, but I don't know what to do with it. It came in port 3 from Ethernet mac address 10:23:10:59:12:fb and it's going to mac address 5c:fb:12:59:10:23.

Controller: OK. I'll memorize that the 10:23:10:59:12:fb address is on port 3. But I don't know which port has a device with address 5c:fb:12:59:10:23. So just send it out all ports on the switch except port 3.

Device: OK. ...Oops, here's another packet I don't know what to do with. It came in port 5 from Ethernet mac address 5c:fb:12:59:10:23 and it's going to mac address 10:23:10:59:12:fb.

Controller: Oh yeah. That looks like a reply. I'll memorize that the 5c:fb:12:59:10:23 address is on port 5. Meanwhile, I know the destination is on port 3. Install a rule so all packets going to that mac address go out port 3, then forward this packet out port 3 as well.

Device: OK!

Controller: How many packets have went out port 3, by the way?

Device: 82,120.

Device: (To itself) I just saw a packet destined for Ethernet mac address 10:23:10:59:12:fb:5c, but I have a rule for dealing with it. I'm gonna send it out port 3.

OpenFlow boils down control plane functionality to a common core set of actions. A list of rules and actions that can be handled in the device itself are kept in a *flow table*. Any complex decisions that can't be handled independently by the control plane may be offloaded to the controller. In a well-designed Software Defined Network, the controller gets involved only when necessary. After all, the conversation between the device and the controller takes time, and anything that can eliminate this conversation makes the packets go faster. So in the example above, a rule for any packets going to 10:23:10:59:12:fb:5c to be output on port 5 keeps all the processing on the switch, and out of the controller. That makes it really fast.

So, central to the OpenFlow model is the *flow table*. Flow tables have *entries*, sometimes called *flow rules*, that are consulted for making decisions. A sample flow table might look like this:

Match	Actions	Priority
dl_src = 10:23:10:59:12:fb:5c, dl_type = 0x806	OFPAT_OUTPUT(9)	100
nw_src = 128.56.0.0/16, dl_type = 0x800	OFPAT_SET_DL_DST(5c:fb:12:59:10:23), OFPAT_OUTPUT(1)	90
Wildcard	OFPAT_OUTPUT(Controller)	1

The main components of a flow entry are:

A match specifies patterns of packet header and metadata values. OpenFlow 1.3 defines 40 different fields to match: some popular ones are the Input Port, the Ethernet Destination mac address, and the TCP destination port. The match values can either be exact (like 10:23:10:59:12:fb:5c above) or wild carded (like 128.56.0.0/16 for a particular Ip subnet).

Actions tell what to do if the match occurs. Instructions can apply actions (send a packet out a port, or write some header information, or send a packet to the controller), invoke groups (like a function call in a programming language), or set variables.

A priority defines the order that matches are consulted. When more than one entry matches a particular packet, the entry with the highest priority wins.

In our example above, the controller installed a flow entry matching Ethernet Destination 10:23:10:59:12:fb:5c, and an instruction applying the action "Output it through Port 3".

OpenFlow's flow table model is *abstract*. An OpenFlow device is not necessarily going to find a RAM chip with the matches, instructions and priorities ... although a pure

software switch like Open vSwitch might mirror it quite closely. Instead, the controller asks the device to install entries, and the device accommodates it by placing entries in its own tables. For example, a real network device might have an L3 table that matches subnets with the various ports that have IP gateways. A programmer, knowing this table exists, can write instructions that match on those ports, and place it directly in the table accordingly.

Suppose you wanted to write your own controller from scratch. You could do that just by talking the OpenFlow protocol. Let's say you wrote this program in Node.js and placed it on the server "controller.example.com", listening on TCP port 6653. Then you'd just point your OpenFlow network device at controller.example.com:6653. Then your program could install flow table entries into the network device over the OpenFlow protocol.

Hmmm. Sounds pretty easy, but ...

2.2 OpenFlow is Difficult

From a programmer's perspective, a table looks awfully primitive. "That's not code, that's data," you might say. But a table is very easy for switch hardware to interpret, and the faster they can interpret and carry out rules, the faster the packets travel. It's like machine language, where the CPU interprets simple instructions very quickly, and even in parallel.

You don't program in machine language most of the time, though, and you shouldn't have to program directly in OpenFlow tables. Why not?

- The set of matches and actions is very limited
- They are difficult to modularize and compose
- They are difficult to prove correct.

Programming OpenFlow tables directly, you begin to find out the subtle missing details:

- You can only match packets with $=$. There's no \neq .
- There is an implicit And in all match rows and an implicit Or between all rules, but you can't be more flexible than that.
- Matching against a set of values requires you to write one rule per value.

Because tables often have thousands of rules, they are difficult to construct and debug. In the programming world, *modularization* aids both of these problems since smaller units of code are easier to understand.

OpenFlow tables have no inherent grouping mechanism, but we could simply modularize them by constructing small tables that do target packet processing. Smoosh them together into one big OpenFlow table when we're done, right?

But as the paper Foster et al. [2013] points out in section IIA, even simple modules can be difficult to *compose*. Suppose your SDN switch needed to do two things: repeat all traffic, but drop all HTTP packets coming from port 2 (a makeshift firewall). The repeater table might look something like this:

Match	Actions	Priority	
in_port=1	OFPAT_OUTPUT(2)	200	
in_port=2	OFPAT_OUTPUT(1)	100	

And the firewall table might look like this:

Match	Actions	Priority	
in_port = 2, tp_src_port = 80, dl.type = 0x800, nw_proto = 0x1	None	100	

If we simply smooshed the two tables together, the firewall rule would never fire because the first rule in the repeater table overshadows it. In this case, reordering the priorities might work, but it's impossible to do this correctly without a spec to guide it.

Finally, it's difficult to reason about OpenFlow tables. While it's true that the set of possible packets is a finite set, it's still a large set. One could loop over all header values (200 bits worth in an OpenFlow 1.0 structured packet) and give the corresponding actions. But it's tough to actually enumerate all these cases.

Frenetic obeys mathematically-defined rules about packets, and its algorithms are provably correct, which you can see in the paper Smolka et al. [2015] And as outlined in Foster et al. [2015], you can prove properties like loop-freeness and connectivity about NetKAT programs.

2.3 Predicates

No matter what controller you use, underneath it all, you still have OpenFlow tables. How does the Frenetic controller improve things?

Other SDN controllers like OpenDaylight, RYU, and Beacon force you to manipulate OpenFlow tables directly. Frenetic works at a higher abstraction level. Here, instead of writing programs that directly call OpenFlow primitives, you write programs in the NetKAT language. These programs are called *network applications* or more succinctly *net apps*.

Frenetic's main job is to compile NetKAT predicates and policies into OpenFlow flow tables. It directly communicates with the switch hardware (southbound) and to your net apps (northbound). Net apps talk to Frenetic with good ol' fashioned HTTP, REST, and JSON. The JSON-based NetKAT dialect is available to anyone, and any programming language that can talk HTTP and JSON can talk to Frenetic. In this manual, we use the Python bindings for NetKAT because they're easy to use and extend, and they come bundled with Frenetic. This saves you from dealing with the esoterica of HTTP communication and JSON formatting.

So let's look at NetKAT predicates first. A *predicate* is a clause used to match packets. The base language is pretty straightforward:

<code>SwitchEq(<i>n</i>)</code>	Matches packets that arrive on switch <i>n</i> , where <i>n</i> is the Datapath ID of the switch.
<code>PortEq(<i>n</i>)</code>	Matches packets that arrive on port <i>n</i> . Generally ports are numbered 1- <i>m</i> , where <i>m</i> is the number of interfaces, but they don't need to be consecutive.
<code>EthSrcEq(<i>mac</i>)</code>	Matches packets whose Ethernet Mac source address is <i>mac</i> , which is a string in the standard form <i>nn : nn : nn : nn : nn</i> where the <i>n</i> 's are lowercase hexadecimal digits.
<code>EthDstEq(<i>mac</i>)</code>	Matches packets whose Ethernet Mac destination address is <i>mac</i> .
<code>VlanEq(<i>vlan</i>)</code>	Matches packets whose VLAN is <i>vlan</i> , and integer from 1-4096. Packets without a VLAN are never matched by this predicate.
<code>VlanPcpEq(<i>p</i>)</code>	Matches packets whose VLAN Priority Code Point is <i>p</i> . Packets without a VLAN are never matched by this predicate.
<code>EthTypeEq(<i>t</i>)</code>	Matches packets whose Ethernet Type is <i>t</i> , where <i>t</i> is a 32 bit integer. Popular values of <i>t</i> are 0x800 for IP and 0x806 for ARP.
<code>IPProtoEq(<i>p</i>)</code>	Matches packets whose IP Protocol is <i>p</i> , a number from 0-255. Popular values are 1 for ICMP, 6 for TCP and 17 for UDP. This match only makes sense when <code>EthTypeEq(0x800, 0x806)</code> (IP or ARP).
<code>IPSrcEq(<i>addr</i>, <i>mask</i>)</code>	Matches packets whose IP source address is <i>addr</i> . If <i>mask</i> is provided, a number from 1-32, this matches all hosts on a network whose first <i>mask</i> bits match the host. If it's omitted, the entire address is matched – i.e. only one IP host. This match only makes sense when <code>EthTypeEq(0x800, 0x806)</code> (IP or ARP).
<code>IPDstEq(<i>addr</i>, <i>mask</i>)</code>	Matches packets whose IP destination address is <i>addr</i> . Follows same rules as <code>IPSrcEq</code> .
<code>TCPSrcPortEq(<i>p</i>)</code>	Matches packets whose TCP source port is <i>p</i> , an integer from 0-65535.
<code>TCPDstPortEq(<i>p</i>)</code>	Matches packets whose TCP destination port is <i>p</i> , an integer from 0-65535. Popular values are 80 for HTTP, 22 for SSH, and 443 for SSL.

If you're familiar with OpenFlow, this list should look familiar to you – it's the same list of fields you can use in an OpenFlow flow table. One special case is `SwitchEq`, matching a switch id, which we'll talk about in a second.

For each rule, OpenFlow allows only one kind of boolean operator: AND. If the Open-

Flow rule specifies match criteria p_1, p_2, \dots, p_n , the packet must match p_1 AND p_2 AND \dots AND p_n . NetKAT is much more expressive, allowing the following boolean operators between predicates:

$p_1 \ \& \ p_2$	Matches packets that satisfy p_1 AND p_2
<code>And([p_1, p_2, \dots, p_n])</code>	Matches packets that satisfy all the predicates: p_1 AND p_2 AND \dots AND p_n
$p_1 \ \ p_2$	Matches packets that satisfy p_1 OR p_2
<code>Or([p_1, p_2, \dots, p_n])</code>	Matches packets that satisfy one of the predicates: p_1 OR p_2 OR \dots OR p_n
$\sim p_1$	Matches packets that DO NOT satisfy p_1
<code>Not(p_1)</code>	Synonym for $\sim p_1$

The precedence is the same as for most Boolean operators in normal programming languages: Not, then And, then Or.

As a shortcut, each of the field-matching predicates `FieldEq` has an analagous `FieldNotEq` predicate which matches all values *except* the given ones.

A bunch of predicates, stitched together with Boolean operators, can be used wherever a simple predicate is used. Furthermore, predicates can be assigned to Python variables. So here are some examples in Python code:

```
import sys
sys.path.append('../src/frenetic/lang/python')
import frenetic
from frenetic.syntax import *

# Note this program doesn't actually do anything

# Match packets from a particular mac address
src_match = EthSrc("10:23:10:59:12:fb:5c")

# Match packets from a particular port that are either IP or ARP packets
port_ip_arp_match = PortEq(9) & EthType(0x800, 0x806)

# Matches packets from a particular port on switches 2 or 3 only
port_switch_match = PortEq(8) & SwitchEq(2, 3)

# Matches packets from all ports except 2 or 3
non_router_match = PortNotEq(2, 3)

# Matches broadcast packets or packets from a particular port
# or packets with a particular Vlan
all_criteria = [ EthSrc("ff:ff:ff:ff:ff:ff"), PortEq(1), VlanEq(2345) ]
brd_or_port_match = Or( all_criteria )
```

One predicate requires some explanation: `SwitchEq`. An OpenFlow flow table belongs to one and only one switch, but a NetKAT program belongs to every switch connected

to that controller. So a predicate tagged with `SwitchEq` will limit a particular match to a particular switch. Any predicates that don't have a `SwitchEq` predicate will apply to *all* switches in the network.

Finally, there are a few special predicates:

<code>id</code>	Matches all packets
<code>drop</code>	Match no packets

Why would you need these? They're useful for "catch all" rules that appear last in a list. A good example is our repeater, where we had an `id` rule that matched all packets and forwarded them to the controller.

2.4 Policies

NetKAT predicates are useless by themselves. To make them work, you need to form them into NetKAT *policies*. A policy is like a command. Often they are compiled down to OpenFlow actions, and in fact there's a lot of overlap between the two concepts. But just as NetKAT predicates are more powerful than OpenFlow matches, NetKAT policies are more powerful than OpenFlow action lists.

<code>Filter (p)</code>	Select packets that match NetKAT predicate <i>p</i> , and quietly forget the rest
<code>SetPort (n)</code>	Set the output port for the packet to port <i>n</i> .
<code>SendToController (tag)</code>	After all actions have been performed, send packet to controller with tag <i>tag</i>
<code>SetEthSrc (mac)</code>	Set Ethernet Mac source address to <i>mac</i>
<code>SetEthDst (mac)</code>	Set Ethernet Mac destination address to <i>mac</i> .
<code>SetVlan (vlan)</code>	Set packet VLAN to <i>vlan</i> . Note this is not a Vlan push - it overwrites whatever Vlan is in the packet (if there is one).
<code>SetVlanPcp (p)</code>	Set VLAN Priority Code Point to <i>p</i> .
<code>SetEthType (t)</code>	Set Ethernet Type to <i>t</i> , where <i>t</i> is a 32 bit integer.
<code>SetIPProto (p)</code>	Set IP Protocol to <i>p</i> . This action is only done when <code>EthTypeEq(0x800, 0x806)</code> (IP or ARP).
<code>SetIPSrc (addr)</code>	Set IP source address to <i>addr</i> . Note there is no mask here, as in the equivalent predicate. This action is only done when <code>EthTypeEq(0x800, 0x806)</code> (IP or ARP).
<code>SetIPDst (addr)</code>	Set IP destination address to <i>addr</i> . Follows same rules as <code>SetIpSrc</code> .
<code>SetTCPSrcPort (p)</code>	Sets TCP source port to <i>p</i> .
<code>SetTCPDstPort (p)</code>	Sets TCP destination port to <i>p</i> .

Note that `Set` policies mirror each `Eq` predicate, so for example the predicate `VlanEq(vlan)` has a matching `SetVlan(vlan)`. The exception is `Switch`. This field is not physical fields in a packet header - rather, it is metadata that OpenFlow fills in “invisibly” in a packet. So you can’t set it to a particular value, like you can other header fields. To send a packet to a different switch, you also use `Send`, but you are restricted to switches that are directly connected to the current switch, and you must know out which port to send it (e.g. you must know the topology). We’ll cover strategies for dealing with this in 10.

And just as you can combine predicates with Boolean operators, you can combine policies with NetKAT policy operators:

$pol1 \mid pol2$	Copy the packet and apply both $pol1$ and $pol2$ to it. Also known as <i>parallel composition</i> .
$Union([pol1, pol2, \dots, poln])$	Copy the packet n times and apply policy $pol[i]$ to copy i . Also known as <i>parallel composition</i> .
$pol1 \gg pol2$	Apply the policy $pol1$ to the packet, then apply $pol2$ to it. Also known as <i>sequential composition</i> .
$Seq([pol1, pol2, \dots, poln])$	Apply each of the policies $pol1, pol2, \dots, poln$ to the packet, in order
$IfThenElse(p, pol1, pol2)$	If packet matches predicate p , then apply policy $pol1$, or else apply $pol2$. Either $pol1$ or $pol2$ is applied, but never both. Equivalent to $Filter(p) \gg pol1 \mid Filter(\sim p) \gg pol2$

The `>>` should look familiar to C++ programmers. Like in C++, the `>>` operator changes a piece of data, then forwards it to the next step in the chain, one after the other. It’s especially helpful in I/O, where you build a string from pieces, then send it to the output device (file or screen) as the last step. The `|` symbol is somewhat like the equivalent in UNIX shell programming: the components actually run in parallel. However, unlike `|`, in NetKAT you are actually running separate copies of each policy without any connections between them. In other words, you don’t send packets from the output of one into the input of another.

Some examples will clarify.

2.5 Commands and Hooks

The truth is, just like NetKAT predicates, NetKAT policies don’t do anything by themselves. And so we come to the last level of a net app: commands and hooks. Commands are instructions from the net app to the switches (via Frenetic). OpenFlow calls

these *controller-to-switch messages*. Hooks are instructions from the switches to the net app (again, via Frenetic), and OpenFlow calls these *switch-to-controller messages*.

The commands are:

<code>pkt_out(sw, payload, plist, inport)</code>	Send a packet out switch with DPID <i>sw</i> . We'll describe this in detail below.
<code>update(policy)</code>	Update all switches with the given NetKAT policy. This is the equivalent of setting the OpenFlow flow tables all in one shot.
<code>port_stats(sw, port)</code>	Get statistics for a particular switch and port.
<code>query(label)</code>	Get user-defined statistics associated with a particular label. We'll cover this in 7.
<code>current_switches()</code>	Gets a list of DPID's of all current, operating switches, and the operating ports for each. This is most useful in the <code>connected()</code> hook.
<code>config(compiler_options)</code>	Set Frenetic compiler options. This is an instruction to Frenetic, not the switch.

You call commands in your net app through plain ol' Python method calls, e.g. `self.update(policy)`. The hooks are:

<code>connected()</code>	Called when Frenetic has finished startup and some (perhaps not all) switches have connected to it.
<code>packet_in(sw, port, payload)</code>	A packet has arrived on DPID <i>sw</i> , port <i>port</i> and either the matching Policy had a Send-ToController policy, or the packet matched no rules and the TableMiss rule forwards TableMiss packets to the controller. This is described in detail below.
<code>switch_up(sw)</code>	Called when a switch has been initialized and is ready for commands. Some switches send this message once every 5 minutes or some user-defined interval, and it doesn't necessarily mean it was down before this.
<code>switch_down(sw)</code>	Called when a switch has been powered-down gracefully.
<code>port_up(sw, port)</code>	Called when a port has been activated.
<code>port_down(sw, port)</code>	Called when a port has been de-activated - most of the time, that means the link status is down, the network cord has been unplugged, or the host connected to that port has been powered-off.

Your net app may be interested in one or more of these hooks. To add code to it, you write a *handler* which implements the hook's signature. Following Python conventions, it must be named exactly the same as the hook. However, you don't need to provide handlers for every hook. If you don't provide one, Frenetic uses its own default handler – in the case of `connected()`, for example, it merely logs a message to the console.

We'll see most of these commands and hooks used in the next few chapters. But since `pkt_out()` and `packet_in()` are crucial to net apps, we'll describe them first here.

2.5.1 The `packet_in` Hook

`packet_in()` is used to inspect network packets. Note that *not all packets* coming in through all switches arrive here – indeed, if they did, your controller would be horribly slow (as our Repeater example is, but we'll see how to improve it in a bit.) Packets arrive here in one of two ways:

- It may match an OpenFlow rule with the action `OFFPAT_OUTPUT (OFPP_CONTROLLER)` also known in NetKAT as `SendToController`.
- It may not match any OpenFlow rule. In OpenFlow 1.0, the default action is to send the packet to the controller. In OpenFlow 1.3, the default action is to drop the packet, but it can be configured to send it to the controller instead.

Once at the controller, Frenetic will deliver the packet to `packet_in()`. The default handler will simply log a message and drop the packet, but of course that's not very interesting.

The handler is called with three parameters:

<code>sw</code>	The DPID of the switch where the packet arrived.
<code>port</code>	The port number of the port where the packet arrived.
<code>payload</code>	The raw packet data.

There are two formats for the packet data: *buffered* or *unbuffered*. Most switches will only send unbuffered data, meaning the entire network packet - header and data - will be transferred to the controller. Unbuffered data can be held at the controller, changed and resent, dropped, multiplied into many different packets, or any arbitrary action without penalty.

Buffered data, on the other hand, exchanges flexibility for some efficiency. Only a subset of buffered packet data is sent to the controller - the amount is configurable, but defaults to 128 bytes which is usually enough for the Ethernet and IP headers. If packets are large, or the channel between the switch and controller is slow, this can save much precious time. The buffered packet will wait at the switch until it either gets changed and sent (through `pkt_out` in the case of Frenetic), or explicitly dropped, or dropped due to timeout.

Although buffering data sounds like a good idea, it requires greater care on the net apps part to prevent duplicate packets, race conditions, and so on. It's also possible that truncated packets become unparseable - for example, if they get cut off in the middle of the IP header. Because of this, very few OpenFlow switches support buffering at all. If capacity is an issue, it's better to send as few packets to the controller as possible. A judicious use of NetKAT policies in the flow table will ensure this.

If you don't need to examine any packet data, you can simply drop *packet*, or pass it directly to `pkt_out`, as we did in our Repeater app:

```
...
class RepeaterApp(frenetic.App):

...
    def packet_in(self, dpid, port_id, payload):
        out_port = 2 if port_id == 1 else 1
        self.pkt_out(dpid, payload, [ Output(Physical(out_port_id)) ] )

...
```

Here, we merely send the payload out unchanged. If it came in as buffered data, it will be sent out as buffered data. If it came in unbuffered, it will be sent out unbuffered. Pretty simple. And in the cases where all you want to decide in the controller which ports to send out a packet, this is all the infrastructure you need.

If you need to examine the payload, you'll need some assistance. The payload is the raw network packet data, not parsed or translated at all. So for the Ethernet Source address, for example, you'd need to examine bytes 14 through 19 (byte ordering starts at 0). Working at this low-level is exceedingly error prone.

So Frenetic leverages the RYU Packet library. RYU is an open source project spearheaded by NTT (Nippon Telephone and Telegraph) and its Python packet parsing library is very solid and complete.

Frenetic provides a simple API on top of RYU Packet.

```
# You must import the proper protocol from the RYU packet library to use it
from ryu.lib.packet import ethernet, arp

...
    def packet_in(self, dpid, port_id, payload):
        ethernet_packet = self.packet(payload, 'ethernet')
        src_mac = ethernet_packet.src
        if ethernet_packet.ethertype == 0x806:
            arp_packet = self.packet(payload, 'arp')
            src_ip = arp_packet.src_ip

...
```

`p = packet(payload, protocol)` turns the raw *payload* into a parsed packet of type *protocol*. From there, the fields of *p* are the parsed values of that protocol. You can

think of *p* as a view into *payload* with the *protocol* lens attached. In the example above, for example, *payload* is both an Ethernet packet and an ARP packet, and the variables *ethernet_packet* and *arp_packet* provide respective views into it. If *payload* is not parseable into that protocol, `None` is returned.

A complete reference for RYU packets is available at <http://ryu-zhdoc.readthedocs.org/en/latest/>. The following is a list of the most popular protocols and fields:

ethernet

<code>dst</code>	Destination mac address string, formatted as "08:60:6e:7f:74:e7"
<code>src</code>	Source mac address string
<code>ethertype</code>	Ethernet frame type. Popular values are 0x800 for IP version 4, or 0x806 for ARP.

Constructor: `ethernet.ethernet(dst, src, ethertype)`

vlan

<code>vid</code>	VLAN id
<code>pcp</code>	Priority Code Point
<code>ethertype</code>	Ethernet frame type. The outer Ethernet packet has ethertype 0x8100, so this is the ethertype of the packet's data.

Constructor: `ethernet.vlan(pcp, vid, ethertype)`

ipv4

<code>proto</code>	The IP version 4 protocol. Popular values are 6 for TCP and 17 for UDP.
<code>src</code>	Source address, a 32 bit integer
<code>dst</code>	Destination address, a 32 bit integer

Constructor: `ipv4.ipv4(proto, src, dst)`

arp

<code>opcode</code>	Popular values are <code>arp.ARP_REQUEST</code> and <code>arp.ARP_REPLY</code> .
<code>src_mac</code>	Source mac address string, formatted as "08:60:6e:7f:74:e7"
<code>src_ip</code>	Source IP address, a 32 bit integer
<code>dst_mac</code>	Destination mac address string
<code>dst_ip</code>	Destination IP address

Constructor: `arp.arp_ip(opcode, src_mac, src_ip, dst_mac, dst_ip)`

2.5.2 The `pkt.out` Command

`pkt.out` is used to send out packets from the switch. Most of the time, the packets you send are packets you received through `packet.in`. But there's nothing stopping from you sending arbitrarily-constructed packets from here as well.

The command takes the following parameters:

<code>sw</code>	The DPID of the switch from where the packet should be sent.
<code>payload</code>	The raw packet data, wrapped in a <code>Buffered</code> or <code>Unbuffered</code> object type.
<code>policy_list</code>	Python list of NetKAT policies to apply to the packet data.
<code>in_port</code>	Port ID from which to send it. This parameter is optional, and only applies to buffered packets presumably sitting at a particular port on the switch waiting to be released.

The `policy_list` effectively tells the switch how to act on the packet. Most NetKAT policies are usable here including `SetIPSrc` and `SendToController`. The exceptions are:

- `Filter` is not usable. To optionally send or not send packets, use the packet library from RYU to make decisions.
- `SetPort` is not available, but `Output(Physical(p))` can be substituted. The difference is subtle. `SetPort` can be used anywhere in a policy sequence, and can be followed by more packet modifications. `Output` sends the packet out immediately, and according to OpenFlow it is always the last action executed if present.
- Only simple policies are doable, so you can't use policy operators like `Union`, `Seq` and `IfThenElse`. Instead, you can send multiple policies in a list, where there's an implied `Seq` operator between them.

What if you want to modify a packet before sending it out? There are actually two ways to do it, each appropriate for a particular use case:

Direct Modification where you set particular data in the packet itself, reserialize it through RYU's packet library API, and send it in `payload`. This is the only way to modify data that's not accessible to a NetKAT policy - e.g. the IPv4 source address is settable by NetKAT policy `SetIPv4`, but there's no equivalent policy for the ARP opcode. Direct modification is only available for unbuffered packets.

Policy Modification is achieved through the Policy list, and is limited to modification through NetKAT policies. It can be done on buffered or unbuffered packets.

In general, Policy Modification is preferable since it works on all packets, and saves you the costly step of parsing and reserializing the packet in the controller. For example, a common routing function is to change the destination MAC address for a next hop router. Here's an example of doing this with Policy Modification:

```
def packet_in(self, dpid, port_id, payload):
    (next_hop_mac, next_port) = calculate_next_hop(payload)
    self.pkt_out(dpid, new_payload,
        [ SetEthDst(next_hop_mac), Output(Physical(next_port)) ]
    )
```

For those times when you need Direct Modification, here's an example of how to do it:

```
from ryu.lib.packet import ethernet, arp
...

def payload(self, e, pkt):
    p = packet.Packet()
    p.add_protocol(e)
    p.add_protocol(pkt)
    p.serialize()
    return NotBuffered(binascii.a2b_base64(binascii.b2a_base64(p.data)))

def packet_in(self, dpid, port_id, payload):
    ethernet_packet = self.packet(payload, 'ethernet')
    arp_packet = self.packet(payload, 'arp')
    # Flip a request into a reply and vice versa
    arp_packet.opcode = \
        arp.ARP_REPLY if arp_packet.opcode == arp.ARP_REQUEST \
        else arp.ARP_REQUEST
    # Reserialize it back into a raw packet
    new_payload = self.payload(ethernet_packet, arp_packet)
    self.pkt_out(dpid, new_payload, [ Output(Physical(1)) ])
...
```

You can also create packets from scratch with a variation of Direct Modification. First create the packet views with standard RYU Packet library constructors. Then just call `self.payload()` on them to create the raw packet and send it.

```
from ryu.lib.packet import ethernet, arp
...

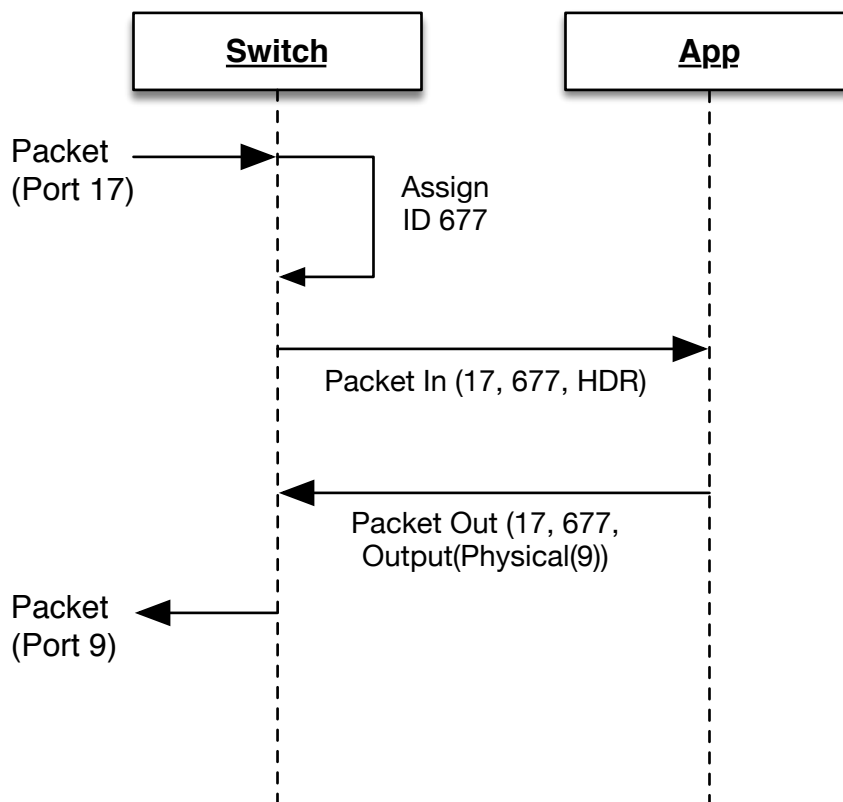
def send_arp_reply(port, src_mac, dst_mac, src_ip, dst_ip):
    # Immediately send an ARP reply when a port comes up.
    e = ethernet.ethernet(dst=dst_mac, src=src_mac, ethertype=0x806)
    a = arp.arp_ip(arp.ARP_REPLY, src_mac, src_ip, dst_mac, dst_ip)
    # Serialize it into a raw packet
```

```
new_payload = self.payload(ethernet_packet, arp_packet)
self.pkt_out(dpid, new_payload, [ Output(Physical(port)) ])
...
```

2.5.3 Buffering

Packet buffering is often a settable option in a switch's OpenFlow configuration. It's especially useful when packets are large, as in Jumbo Frames. Most switching decisions made in the controller look only at the packet headers, not the data, so why should you have to send it? By buffering the entire packet, the switch only sends the headers to the controller. The controller only does Policy Modifications to the buffered packet, saving the trouble of sending all of it back.

A buffered Packet Out *must always* be preceded by a buffered Packet In. That's because you need to send back two things: the incoming port id, and the buffer id. The switch actually has separate buffer list for every port on the switch, and sending it back with the proper port helps it match it to the correct buffer. This sequence diagram shows how it commonly works:



Here, the HDR is the first 128 bytes of the packet, enough to hold the Ethernet and IP header information, typically. The app doesn't send any of this header information back to the switch - just the instruction `Output(Physical(9))` to direct the switch's data

plane.

What if we *never* send back the Packet Out? Buffer contents are not held indefinitely – they timeout after a certain period. At that point, the buffered packet will drop out. If we send a Packet Out for that buffer after the timeout period, the Packet Out will generate an error (which Frenetic ignores). A similar fate awaits Packet Outs that fabricate a random buffer id, or send it to the wrong port.

So where is the buffer id in the `PacketOut` call? It's embedded in the `payload` object. When a `PacketIn` delivers a buffered payload, you can simply it send it back out the `PacketOut` call, which sends the buffer id along with it.

2.6 OpenFlow Constructs Not Supported by NetKAT

NetKAT supports the most popular features of OpenFlow 1.0. But there a few things it can't do:

- It can't output to special pseudo-ports NORMAL, FLOOD, etc. The semantics of sending to these ports depends on the spanning tree, VLAN's, and other settings that NetKAT is not aware of, so it cannot reason about them.
- It can't set the size for buffered packet data in `SendToController`. All buffered packets get sent with the default 128 bytes, usually enough to encapsulate the header information.
- It can't retrieve flow table counters. Since one NetKAT policy may expand to many OpenFlow rules, or even be optimized to OpenFlow Rules, there is no good way to map and retrieve this data.
- Frenetic must talk to switches through an unsecured channel. TLS is not supported.
- Most controller-to-switch messages like Features, Configuration or Barrier are not supported.
- Some switch-to-controller messages like Error, Flow Removed and Vendor are ignored by Frenetic.
- Some OpenFlow actions are not supported, like Enqueue.

We haven't discussed some implemented features like Statistics yet, but those will be described in chapter 7.

Chapter 3

NetKAT Principles

3.1 Efficient SDN

In a nutshell, the `packet_in()` hook receives network packets and the `pkt_out()` command sends network packets. In theory, you could use these two to implement arbitrarily-complex network clients and servers. You could build switches and routers, but also HTTP servers, Email servers, Database servers, or any other network server.

That said, you probably wouldn't want to. OpenFlow and Frenetic are optimized for small, very selective packet inspections and creations. The more packets you inspect through `packet_in()`, the slower your controller will be, and the more likely that packets will be dropped or sent out of sequence.

Principle 1 *Keep as much traffic out of the controller as possible. Instead, program NetKAT policies to make most of the decisions inside the switch.*

So let's go back to our Repeater application:

```
import sys
sys.path.append('/home/vagrant/src/frenetic/lang/python')
import frenetic
from frenetic.syntax import *

class RepeaterApp(frenetic.App):

    client_id = "quick_start"

    def connected(self):
        self.update( id >> SendToController("repeater_app") )

    def packet_in(self, dpid, port_id, payload):
        out_port_id = 2 if port_id == 1 else 1
        self.pkt_out(dpid, payload, [ Output(Physical(out_port_id)) ] )
```



```
app = RepeaterApp()
app.start_event_loop()
```

Here, *every single packet* goes from the switch to Frenetic to the net app and back out. That's horribly inefficient, and unnecessarily so since all the decisions can be made inside the switch. So let's write a more efficient one. The following code is in `code/netkat_principles/re`

```
import sys
sys.path.append('/home/vagrant/src/frenetic/lang/python')
import frenetic
from frenetic.syntax import *

class RepeaterApp2(frenetic.App):

    client_id = "repeater"

    def connected(self):
        rule_port_one = Filter(PortEq(1)) >> SetPort(2)
        rule_port_two = Filter(PortEq(2)) >> SetPort(1)
        self.update( rule_port_one | rule_port_two )

app = RepeaterApp2()
app.start_event_loop()
```

This program takes principle 1 very seriously, to the point where *no* packets arrive at the controller. All of the configuration of the switch is done up front.

The `Filter(PortEq(1)) >> SetPort(2)` policy is a pretty common pattern in NetKAT. You first whittle down the incoming flood of packets to a certain subset with `Filter` and a predicate. Then you apply a policy or series of policies, `SetPort` being the most popular. We'll look at combining policies in the section 3.2.

If you've worked with OpenFlow, you might wonder how the NetKAT rules get translated to OpenFlow rules. In this example, it's fairly straightforward. You get two OpenFlow rules in the rule table, which you can see in the Frenetic debug window:

```
[INFO] POST /repeater/update_json
[INFO] GET /repeater/event
[INFO] New client repeater
[DEBUG] Installing policy
drop | (filter port = 1; port := 2 | filter port = 2; port := 1)
[DEBUG] Setting up flow table
+-----+
| 1 | Pattern | Action   |
|-----|
| InPort = 1 | Output(2) |
|-----|
| InPort = 2 | Output(1) |
```



But this is not true in general. One NetKAT rule may expand into many, many OpenFlow rules. And it may go the opposite direction to: where different NetKAT rules are combined to create one OpenFlow rule. It's the same thing that happens with most compiled languages – the rules that govern the compiled code are non-trivial. If they were easy, you wouldn't need a compiler to do it!

There are two problems with RepeaterApp2:

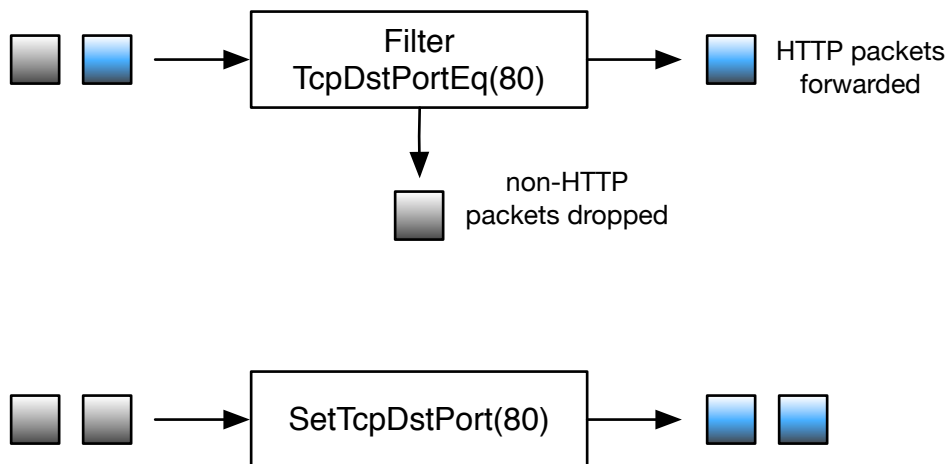
- It works on a two port switch, but not anything bigger. And the ports absolutely have to be numbered 1 and 2 ... otherwise, the whole program doesn't work. And those ports need to be functioning.
- More subtly, this program can drop packets. There is a short lag in between when the switches come up and the `self.update()` installs the policies. During this lag, packets will arrive at the controller by default and get dropped by the default `packet_in` handler in Frenetic.

We will correct both of these problems in section 3.3

3.2 Combining NetKAT Policies

In our Repeater2 network app, the two rules have the `Seq` operator `>>` in them, then the two rules are joined together with `Union` or `|`. But what is the difference between the two? When do you use one or the other?

To illustrate this, let's go back to NetKAT basics. At the lowest level there are two types of building blocks: filters and modifications.

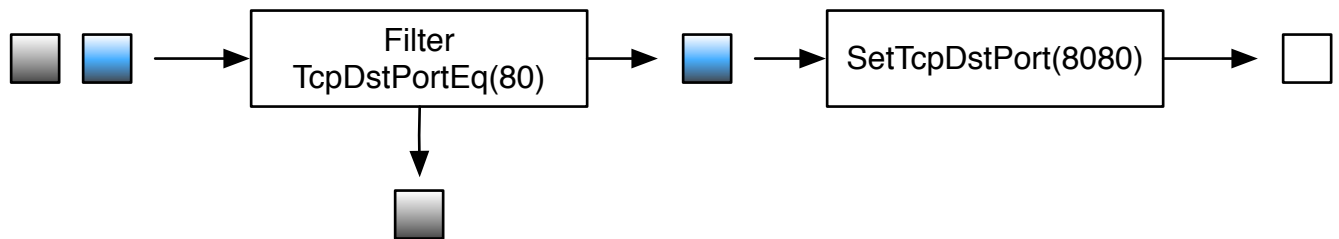


A filter, shown at the top, takes as input a stream of packets. Packets that match its criteria are forwarded, packets that don't match are dropped. Two special filters, `id` and `drop`, let through all or no packets respectively.

A modification changes packet header values. In the bottom example, the `TcpSrcPort` header gets changed to 80. This modification, but rather happens on a logical copy of the packet. That means subsequent `Filters` match against the original value in the header, not the one we just changed it to. But most of the time, you can ignore this subtlety.

Modifications to different headers can generally be specified in any order. If two modifications to the same header value occur one after the other, the last one wins.

You combine these building blocks with `Seq` and `Union`.



Here is a `Seq` of two policies, a filter and a modification. In sequences, policies are chained together in a straight line, and packets travel through each of one of the policies in order. Combining a filter plus modifications is very common in NetKAT programs, and we generally use `Seq` to combine them into a *rule*.

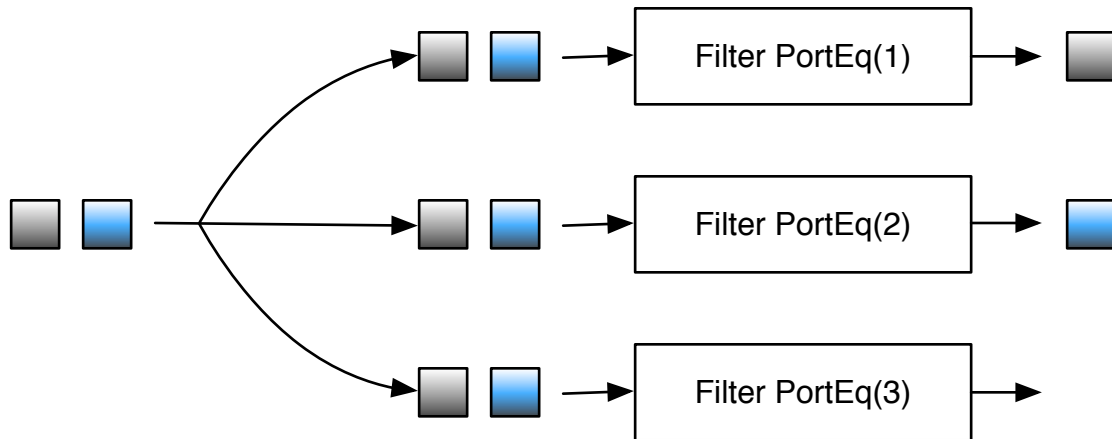
A rule is like an OpenFlow flow table entry, but with more options. If you think about it, a switch receives a huge firehose blast of packets, and you want the switch to perform a targeted action on a very small subset of those packets. In other words, you want a *filter* and some *actions*. It's like the MapReduce paradigm, but in reverse order - you first *reduce* the stream to a manageable level, then *map* the result by performing actions on it.

With `Seq`, order matters. So the rule `drop >> Filter(EthTypeEq(0x806))` drops ALL packets, no matter the type. The second filter is never reached. In general, putting all the filters before the actions in a sequence chain is the clearest way to write it.

The following principle is a good guideline:

Principle 2 Use `Seq` or `>>` between filters and all actions except multiple `SetPorts`.

`SetPort` is an exception to the rule that we'll see in a minute.



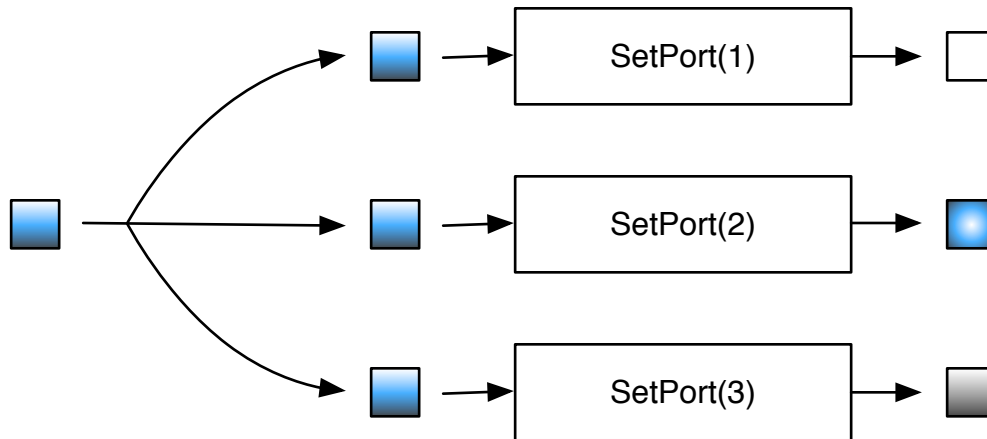
A `Union` of policies makes several logical copies of policies and sends the copy of each packet through each policy in parallel. In the bottom example, the two incoming packets are copied three times, one for each rule. In this case, the top filter only forwards a copy of the first packet. The middle filter only forwards a copy of the second packet. The bottom filter doesn't forward any packets at all.

Principle 3 Use `Union` or `|` between multiple `SetPorts` and rules that **DO NOT** overlap. Use `IfThenElse` to combine rules that **DO** overlap.

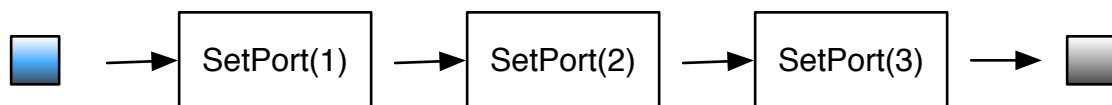
We'll explain `IfThenElse` in a little bit.

You might think, "All that packet copying must be really tough on the switch." In fact, all the packet copying is conceptual, it doesn't necessarily happen in the switch. Frenetic turns NetKAT programs into flow tables, and these flow tables generally work through sophisticated TCAM's, Hash tables, and pipelines that eliminate actual packet copying. So don't worry about stringing 20,000 policies with a `Union`. Your switch can handle it.

Now `SetPort` is a little different. Other modifications like `SetTcpDstPort` act on a single packet at a time. A single `SetPort` is like this - we merely want to change the egress port of the packet itself. But if you want to send it out multiple ports at once, as in flooding, you need to use multiple `SetPort` actions.



A Union of Three SetPorts Sends Three Packets



A Seq of Three SetPorts Sends One Packet

If you string these `SetPorts` together in a `Seq`, the single packet will be assigned each output port in turn, and the last assigned egress port will “win”. This is probably not what you want. On the other hand, stringing these together with `Union` makes more sense. You actually WANT to make extra copies of the packet, and send each copy out a different port.

Now let’s look at `Union` and `IfThenElse`. As we have seen in Chapter 2, `Union` is parallel composition: we effectively make copies of each packet and apply the rules to each packet. So in our repeater application, the rules are:

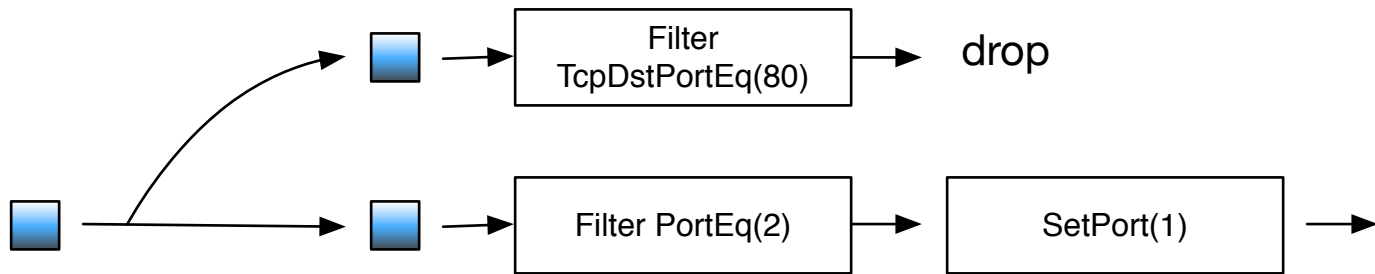
- `Filter(PortEq(1)) >> SetPort(2)` Sends traffic from port 1 out port 2
- `Filter(PortEq(2)) >> SetPort(1)` Sends traffic from port 2 out port 1

A packet cannot match both rules – in other words, there is no overlap. So a `Union` between these two rules is appropriate. You make a copy of each packet, send it through each of the rules. One or the other will match and send it out the appropriate port, the other will simply drop it.

Suppose you have two rules:

- `Filter(TcpDstPortEq(80)) >> drop` Drops non-encrypted HTTP traffic

- `Filter(PortEq(2)) >> SetPort(1)` Sends port 2 traffic out port one.

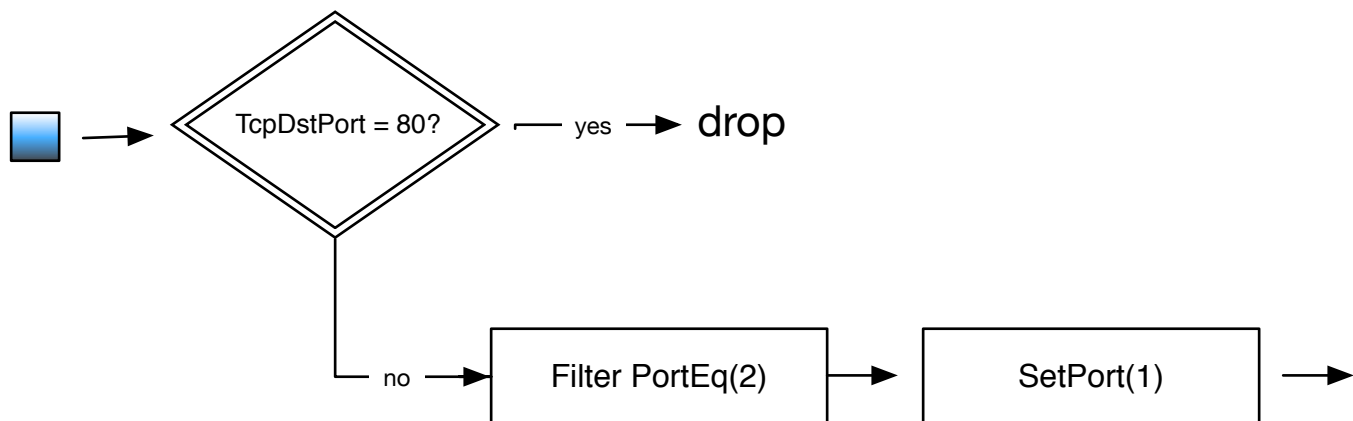


These two rules overlap. A packet can both have a TCP destination port of 80 and arrive on switch 1, port 2. What happens if we combine them with `Union`, and a packet like that arrives at the switch? The packet will be copied twice, then:

- `Filter(TcpDstPortEq(80)) >> drop` will drop the first copy of the packet
- `Filter(PortEq(2)) >> SetPort(1)` will send out the second copy to port 1

But that's probably not what you intended. You probably want the HTTP rule to take precedence over the sending rule. In such cases, `IfThenElse` makes the precedence crystal clear:

```
IfThenElse(TcpDstPortEq(80), drop, Filter(PortEq(2)) >> SetPort(1))
```



The predicate is tested, and if it matches, the first policy is performed, otherwise the second is performed.

That's pretty powerful, and in fact we can write the repeater app to use `IfThenElse` instead of `Union`:

```
IfThenElse(PortEq(1), SetPort(2), Filter(PortEq(2)) >> SetPort(1))
```

So why not just use `IfThenElse` for combining rules all the time? In a real switch, you might have hundreds of non-overlapping rules. Writing them as:

```
IfThenElse(predicate1), action1, IfThenElse(predicate2, action2,
IfThenElse...
```

is doable but a lot more verbose. It also hides the fact that the rules are non-overlapping, and this is useful information when rummaging through an SDN app.

IfThenElse is really syntactic sugar. The statement:

```
IfThenElse(pred, true_policy, false_policy)
```

is just a shortcut for:

```
Filter(pred) && true_policy — Filter(! pred) && false_policy
```

but it's a lot shorter and easier to understand, especially for programmers accustomed to seeing If-Then-Else constructs in other programming languages.

3.3 Keeping It Stateless

So let's kick it up a notch. Our repeater is fast, but pretty static – it only works with two ports, and only ports that are numbered 1 and 2. So let's extend our repeater to do two things:

1. On connection, read the list of currently available ports and build the initial NetKAT policy accordingly.
2. When ports go up or down, adjust the policy dynamically.

The first task requires us to inquire which ports a switch has. Frenetic provides a `current_switches()` method that does exactly that. The best time to call it is on the `connected()` hook since, at that point, Frenetic knows the switches and ports out there.

The following code is in `code/netkat_principles/repeater3.py`:

```
import sys
sys.path.append('/home/vagrant/src/frenetic/lang/python')
import frenetic
from frenetic.syntax import *

class RepeaterApp3(frenetic.App):

    client_id = "repeater"

    def connected(self):
        def handle_current_switches(switches):
            logging.info("Connected to Frenetic - Switches: "+str(switches))
            self.current_switches(callback=handle_current_switches)

logging.basicConfig(stream = sys.stderr, \
    format='%(asctime)s [%(levelname)s] %(message)s', level=logging.INFO \
)
app = RepeaterApp3()
app.start_event_loop()
```

`current_switches()` is an asynchronous call, meaning you don't get the results immediately. Instead, we pass a callback procedure named `handle_current_switches`, which Frenetic calls when the `current_switches()` call is complete and has results. In our first pass, the callback is just a one line procedure that prints the results to the log (the screen by default). Although we could define `handle_current_switches()` outside the `connected()` method, placing it inside clarifies the connection between the two.

So let's start up Mininet with 4 hosts this time:

```
vagrant@frenetic: ~/workspace$ sudo mn --topo=single,4 --controller=remote
*** Creating network
*** Adding controller
Unable to contact the remote controller at 127.0.0.1:6633
*** Adding hosts:
h1 h2 h3 h4
*** Adding switches:
s1
*** Adding links:
(h1, s1) (h2, s1) (h3, s1) (h4, s1)
*** Configuring hosts
h1 h2 h3 h4
*** Starting controller
c0
*** Starting 1 switches
s1 ...
*** Starting CLI:
mininet>
```

Start up Frenetic:

```
vagrant@frenetic: ~/workspace$ ~/src/frenetic/frenetic.native http-controller
[INFO] Calling create!
[INFO] Current uid: 1000
[INFO] Successfully launched OpenFlow controller with pid 2035
[INFO] Connecting to first OpenFlow server socket
[INFO] Failed to open socket to OpenFlow server: (Unix.Unix_error "Connection refused" co
[INFO] Retrying in 1 second
[INFO] Successfully connected to first OpenFlow server socket
[INFO] Connecting to second OpenFlow server socket
[INFO] Successfully connected to second OpenFlow server socket
[INFO] switch 1 connected
```

And start up our application:

```
vagrant@frenetic: ~/workspace$ python repeater3.py
No client_id specified. Using 0323e4bc31c94e81939d9f51640f1f5b
Starting the tornado event loop (does not return).
2016-04-12 09:51:14,578 [INFO] Connected to Frenetic - Switches: {1: [4, 2, 1, 3]}
```

The `handle_current_switches` callback gets called with a Python dictionary. The keys in this dictionary are the datapath ID's (dpid's) of each switch – in our case, we only have one switch with a dpid of 1. The value associated with the key is a list of ports that the switch has operational. In this case, we have four ports labelled 1 through 4 (they will be in some random order in the list, as you see above).

Great! Once we have the list of ports, we can construct policies. If a packet comes in on port 1, it should be repeated to all the ports that are not 1, e.g. 2, 3 and 4, and so on. The policy, written out manually, will look like this:

```
Filter(PortEq(1)) >> (SetPort(2) | SetPort(3) | SetPort(4)) |
Filter(PortEq(2)) >> (SetPort(1) | SetPort(3) | SetPort(4)) |
Filter(PortEq(3)) >> (SetPort(1) | SetPort(2) | SetPort(4)) |
Filter(PortEq(4)) >> (SetPort(1) | SetPort(2) | SetPort(3))
```

A combination of Frenetic and Python lists makes this. Suppose you have the list named `sw [1; 2; 3; 4]` from the callback. The following Python list comprehension:

```
SetPort(p) for p in sw
```

returns the list:

```
[ SetPort(1); SetPort(2); SetPort(3); SetPort(4) ]
```

Now suppose we're looking only at port 1. We want all of the `SetPort`'s here except `SetPort(1)`. A little tweak to the list comprehension:

```
SetPort(p) for p in sw if p != in_port
```

Removes the input port from the list, leaving:

```
[ SetPort(2); SetPort(3); SetPort(4) ]
```

And now you can pass this to the NetKAT Union operator:

```
Union( SetPort(p) for p in sw if p != in_port )
```

... which creates a parallel composition out of the entire list:

```
SetPort(2) | SetPort(3) | SetPort(4)
```

You can pass a Python list of policies to either `Union` or `Seq`, making it easy to build very large policies from component parts very quickly.

So here is our repeater that installs an initial configuration. The following code is in `code/netkat_principles/repeater4.py`:

```
import sys
sys.path.append('/home/vagrant/src/frenetic/lang/python')
import frenetic
from frenetic.syntax import *

class RepeaterApp4(frenetic.App):

    client_id = "repeater"

    def port_policy(self, in_port, all_ports):
        return \
            Filter(PortEq(in_port)) >> \
            Union( SetPort(p) for p in all_ports if p != in_port )

    def all_ports_policy(self, all_ports):
        return Union( self.port_policy(p, all_ports) for p in all_ports )

    def policy(self, switches):
        # We take advantage of the fact there's only one switch
        dpid = switches.keys()[0]
        return self.all_ports_policy(switches[dpid])

    def connected(self):
        def handle_current_switches(switches):
            logging.info("Connected to Frenetic - Switches: "+str(switches))
            self.update(self.policy(switches))
            self.current_switches(callback=handle_current_switches)

logging.basicConfig(stream = sys.stderr, \
    format='%(asctime)s [%(levelname)s] %(message)s', level=logging.INFO \
)

app = RepeaterApp4()
app.start_event_loop()
```

Now it's just a hop, skip and jump to a fully dynamic repeater. First we need to capture packets from ports that we haven't seen yet. We could write an extremely long filter, filtering out every port that's not on the list, but since port numbers can be 32-bits long, that's gonna be pretty huge:

```
Filter(PortEq(5, 6, 7, ...) >> SendToController("repeater_app"))
```

It's easier just to write an overlapping rule. Here `id` is a filter that matches all packets:

```
id >> SendToController("repeater_app")
```

But obviously, this overlaps with the rules for known ports in our repeater. So we use an `IfThenElse` to resolve the overlap:

```
def known_ports_pred(self, all_ports):
    return PortEq(all_ports)

def policy(self):
    return IfThenElse( \
        self.known_ports_pred(self.all_ports), \
        self.all_ports_policy(self.all_ports), \
        SendToController("repeater_app") \
    )

def connected(self):
    def handle_current_switches(switches):
        logging.info("Connected to Frenetic - Switches: "+str(switches))
        # We take advantage of the fact there's only one switch
        dpid = switches.keys()[0]
        self.all_ports = switches[dpid]
        self.update(self.policy())
    self.current_switches(callback=handle_current_switches)
```

We take this opportunity to save the port list in an instance variable `self.all_ports`. This instance variable is the beginning of a Network Information Base, or NIB – it encapsulates the known state of the network at this time. We’re going to see the NIB a lot in future apps. You can think of our app as a function with the NIB as input and a NetKAT program as output. Funky!

Now how we do learn about new ports? A packet arriving at `pkt_in` will signal we’re seeing a new port. So suppose we’re seeing a packet on port 40. If you think about, the entire NetKAT program will change from:

```
Filter(PortEq(1)) >> (SetPort(2) | SetPort(3) | SetPort(4)) |
Filter(PortEq(2)) >> (SetPort(1) | SetPort(3) | SetPort(4)) |
Filter(PortEq(3)) >> (SetPort(1) | SetPort(2) | SetPort(4)) |
Filter(PortEq(4)) >> (SetPort(1) | SetPort(2) | SetPort(3))
```

to:

```
Filter(PortEq(1)) >> (SetPort(2) | SetPort(3) | SetPort(4) | SetPort(40)) |
Filter(PortEq(2)) >> (SetPort(1) | SetPort(3) | SetPort(4) | SetPort(40)) |
Filter(PortEq(3)) >> (SetPort(1) | SetPort(2) | SetPort(4) | SetPort(40)) |
Filter(PortEq(4)) >> (SetPort(1) | SetPort(2) | SetPort(3) | SetPort(40)) |
Filter(PortEq(40)) >> (SetPort(1) | SetPort(2) | SetPort(3) | SetPort(4))
```

Fortunately, we have all the logic we need in `self.all_ports_policy` already. We just need to pass it an amended list of known ports. Not a problem!

```
def packet_in(self, dpid, port_id, payload):  
    self.all_ports.add(port_id)  
    self.update(self.policy())
```

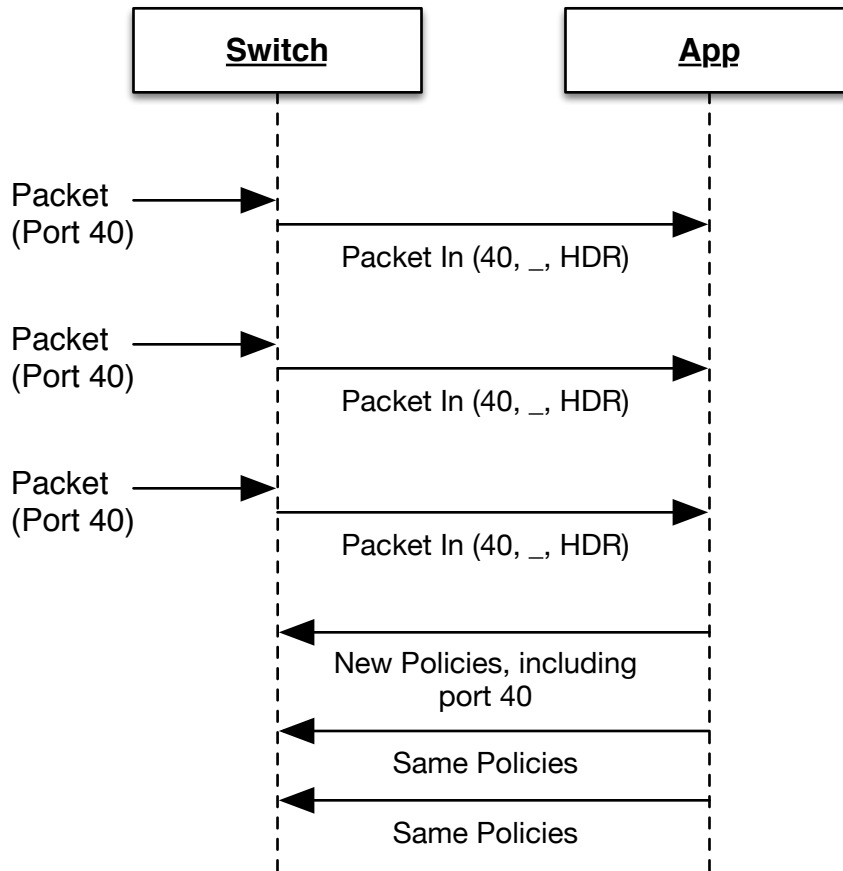
What do we do with the packet we just received? We could just drop it and hope the host is smart enough to resend it. But that's not very hospitable.

What if we just do a `pkt_out` immediately afterward? After all, we just installed a rule that will deal with appropriately, right?

```
def packet_in(self, dpid, port_id, payload):  
    self.all_ports.add(port_id)  
    self.update(self.policy())  
    self.pkt_out(dpid, payload, [ ??? ] )
```

There are two problems. First, it's unclear what the action should be on `pkt_out`. If we send an empty list of actions, the OpenFlow switch will interpret it as "drop the packet."

Second, there's a timing problem. Even though we just sent a `self.update`, the call is asynchronous with no response when it's done updating the OpenFlow flow table. It could take a microsecond, it could take an hour ... we just don't know. This timing problem actually can cause more problems. What if, before the new rules are installed, the switch receives 100 more packets? `pkt_in` will be called 100 times, and each time the policy will be recalculated and sent to the switch. That could cause major problems since installing switch rules can be a CPU-intensive process on the switch.



Hence the following principle:

Principle 4 *When you install a new switch policy, do not assume it'll be installed right away.*

We can solve this with two quick fixes. One is to make a quick check that port has not actually be learned yet. The other is to add actions to `pkt_out` that emulate the new rule.

Here's our final repeater program in `code/netkat_principles/repeater5.py`:

```

import sys
sys.path.append('/home/vagrant/src/frenetic/lang/python')
import frenetic
from frenetic.syntax import *

class RepeaterApp5(frenetic.App):

    client_id = "repeater"

    def port_policy(self, in_port, all_ports):
        return \
            Filter(PortEq(in_port)) >> \
            Union( SetPort(p) for p in all_ports if p != in_port )
  
```

```

def all_ports_policy(self, all_ports):
    return Union( self.port_policy(p, all_ports) for p in all_ports )

def known_ports_pred(self, all_ports):
    return PortEq(all_ports)

def policy(self):
    return IfThenElse( \
        self.known_ports_pred(self.all_ports), \
        self.all_ports_policy(self.all_ports), \
        SendToController("repeater_app") \
    )

def connected(self):
    def handle_current_switches(switches):
        logging.info("Connected to Frenetic - Switches: "+str(switches))
        # We take advantage of the fact there's only one switch
        dpid = switches.keys()[0]
        self.all_ports = switches[dpid]
        self.update(self.policy())
    self.current_switches(callback=handle_current_switches)

def packet_in(self, dpid, port_id, payload):
    if port_id not in self.all_ports:
        self.all_ports.append(port_id)
        self.update(self.policy())
    flood_actions = [ Output(Physical(p)) for p in self.all_ports if p != port_id ]
    self.pkt_out(dpid, payload, flood_actions )

logging.basicConfig(\
    stream = sys.stderr, \
    format='%(asctime)s [%(levelname)s] %(message)s', level=logging.INFO \
)
app = RepeaterApp5()
app.start_event_loop()

```

In mininet with a simple 4-host topology, we can test our repeater by adding a 5th host and port:

```

mininet> py net.addHost('h5')
<Host h5: pid=3187>
mininet> py net.addLink(s1, net.get('h5'))
<mininet.link.Link object at 0x7fd432ace810>
mininet> py s1.attach('s1-eth5')
mininet> py net.get('h5').cmd('ifconfig h5-eth0 10.5')

```

The new port causes the switch to send an OpenFlow portUp message to our app. We don't have a handler for the port_up hook, so the default message appears:

```
port_up(switch_id=1, port_id=5)
```

Back in Mininet, let's try to ping from our new host:

```
mininet> h5 ping -c 1 h1
PING 10.0.0.1 (10.0.0.1) 56(84) bytes of data.
64 bytes from 10.0.0.1: icmp_seq=1 ttl=64 time=996 ms

--- 10.0.0.1 ping statistics ---
1 packets transmitted, 1 received, 0% packet loss, time 0ms
rtt min/avg/max/mdev = 996.738/996.738/996.738/0.000 ms
```

Looks good! A combination of the new packet rules and the `pkt_out` is sending the packets to the correct place.

3.4 Summary

In building a dynamic repeater, we have learned the following NetKAT principles:

Principle 1 *Keep as much traffic out of the controller as possible. Instead, program NetKAT policies to make most of the decisions inside the switch.*

Principle 2 *Use `Seq` or `>>` between filters and all actions except multiple `SetPorts`.*

Principle 3 *Use `Union` or `|` between multiple `SetPorts` and rules that DO NOT overlap. Use `IfThenElse` to combine rules that DO overlap.*

Principle 4 *When you install a new switch policy, do not assume it'll be installed right away.*

And we'll learn a fifth in the next chapter:

Principle 5 *Do not rely on network state. Always assume the current packet is the first one you're seeing.*

These principles are meant as guidelines, not straitjackets. As we build network applications over the next few chapters, we may find good reasons to violate them to achieve shorter code or better modularity. That's OK.

Now we have a robust repeater that acts mostly correctly. There is still a small timing problem. If a packet arrives from port 1 bound for port 40 in that short interval before the rules for port 40 are installed, it will be copied to ports 2, 3, and 4 but not 40. There's little we can do in the context of our simple repeater.

But of course, modern network equipment doesn't act like a repeater. That's 1980's technology! Repeaters do a lot of unnecessary packet copying, flooding hosts with packets that are clearly destined for them. So our next project will be building an L2 learning switch. In that context, we'll correct some of the remaining timing flaws. And the result will be much more efficient.

Chapter 4

Learning Switch

4.1 Design

Layer 2, or L2, switching revolutionized networking in the 1990's. As LAN traffic grew, hub performance rapidly degraded as collisions became more and more frequent. L2 switches effectively cut the LAN into logical segments, performing the forwarding between them. This dramatically reduced the number of collisions, and also cut down on the traffic that individual NIC's had to filter out. Just like the evolution from party lines to direct lines in the Plain Old Telephone Network, the evolution from Hubs to L2 switches improved security, speed and quality.

Of course, L2 switches were more sophisticated than hubs. They required a processor, memory, and some form of programming. In order to know which segments to forward traffic, they needed to watch the Ethernet MAC addresses of traffic and remember their respective ports.

We can simulate the L2 switch with Frenetic. By doing so, as we'll see in 4.4, we can add features to the switch with just a little extra programming effort.

At a high-level, you can think of a Frenetic network application as:

$$netkat = f(nib, env)$$

Where f is your application, nib is the Network Information Base – the information you have dynamically determined in your network through packets received by `pkt_in` – and env is other information (fixed configuration files, out-of-band network measurements, or whatever you want). The output, $netkat$ is the NetKAT program.

Naturally, designing nib is critical to a good design. You don't need to see every packet, and you don't need to record every detail about the packets you see. In an L2 switch, we are really only interested in three pieces of data:

- The source MAC address
- The destination MAC address

- The switch port connected to the host with a particular MAC address

Here's the way we want the switch to behave, in pseudocode:

```
if port_for_mac(EthSrc) == None:
    learn(EthSrc, Port)
if port_for_mac(EthDst) != None:
    pkt_out(payload, port)
else
    pkt_out(payload, all_ports_except(port))
```

Admittedly this is pretty sketchy, but it covers the interesting cases. In particular, it covers Ethernet broadcasts to MAC address ff:ff:ff:ff:ff:ff just by the fact that a source MAC will never equal ff:ff:ff:ff:ff:ff. And flooding is exactly what you want to do in that case.

So our NIB must maintain at least a list of MAC-to-port mappings. In our Repeater app, our NIB was a single instance variable in the application itself: `self.ports`, which held a list of connected ports on the switch. Now we'll evolve a little. In what will become a standard part of our network apps, we'll model the NIB as a separate object class in Python. The following code is in `code/l2/_learning/_switch/network/_information/_base.py`:

```
class NetworkInformationBase():

    # hosts is a dictionary of MAC addresses to ports
    # { "11:11:11:11:11:11": 2, ...}
    hosts = {}

    # ports on switch
    ports = []

    def __init__(self, logger):
        self.logger = logger

    def learn(self, mac, port):
        # Do not learn a mac twice
        if mac in self.hosts:
            return

        self.hosts[mac] = port
        self.logger.info(
            "Learning: "+mac+" attached to ( "+str(port)+" )"
        )

    def port_for_mac(self, mac):
        if mac in self.hosts:
            return self.hosts[mac]
        else:
```

```

        return None

    def all_mac_port_pairs(self):
        return [ (mac, self.hosts[mac]) for mac in self.hosts.keys() ]

    def all_learned_macs(self):
        return self.hosts.keys()

    def set_ports(self, list_p):
        self.ports = list_p

    def add_port(self, p):
        if p not in ports:
            self.ports.append(p)

    def delete_port(self, p):
        if p in ports:
            self.ports.remove(p)

    def all_ports_except(self, in_port):
        return [p for p in self.ports if p != in_port]

```

That encapsulates the state in one place, making it easy to change underlying data structures later. It also separates the NIB details from the NetKAT details, making it easier to reuse the code in other applications later.

4.2 A First Pass

One of the problems with our switch pseudocode design is it doesn't fit our notions of NetKAT very well. NetKAT programs do not have variables, so they can't remember MAC-to-port mappings on their own. So it appears that every packet must pass through the controller so we can make decisions. Processing every single packet through the controller clearly violates the First NetKAT principle, but we can leave that aside for now. It'll be instructive to build an easy but inefficient L2 switch first.

We can use elements from our Repeater application. The following code is in `code/l2/_learning/_sw`

```

1  import sys, logging
2  sys.path.append('/home/vagrant/src/frenetic/lang/python')
3  import frenetic
4  from frenetic.syntax import *
5  from ryu.lib.packet import ethernet
6  from network_information_base import *
7
8  class LearningApp1(frenetic.App):
9
10     client_id = "l2_learning"
11

```

```

12 def __init__(self):
13     frenetic.App.__init__(self)
14     self.nib = NetworkInformationBase(logging)
15
16 def connected(self):
17     def handle_current_switches(switches):
18         logging.info("Connected to Frenetic - Switches: "+str(switches))
19         dpid = switches.keys()[0]
20         self.nib.set_ports( switches[dpid] )
21         self.update( id >> SendToController("learning_app") )
22         self.current_switches(callback=handle_current_switches)
23
24 def packet_in(self, dpid, port_id, payload):
25     nib = self.nib
26
27     # Parse the interesting stuff from the packet
28     ethernet_packet = self.packet(payload, "ethernet")
29     src_mac = ethernet_packet.src
30     dst_mac = ethernet_packet.dst
31
32     # If we haven't learned the source mac, do so
33     if nib.port_for_mac( src_mac ) == None:
34         nib.learn( src_mac, port_id)
35
36     # Look up the destination mac and output it through the
37     # learned port, or flood if we haven't seen it yet.
38     dst_port = nib.port_for_mac( dst_mac )
39     if dst_port != None:
40         actions = [ Output(Physical(dst_port)) ]
41     else:
42         actions = [ Output(Physical(p)) for p in nib.all_ports_except(port_id) ]
43     self.pkt_out(dpid, payload, actions )
44
45 if __name__ == '__main__':
46     logging.basicConfig(\
47         stream = sys.stderr, \
48         format='%(asctime)s %(levelname)s %(message)s', level=logging.INFO \
49     )
50     app = LearningApp1()
51     app.start_event_loop()

```

There are a couple of new details to note:

- The `__init__` constructor must call the superclass constructor to properly initialize.
- Because we are writing in classes, we now distinguish the main loop of this application with a check on `__main__`.
- We are using the RYU packet parsing routines discussed in 2.5.1

- The `packet._in` looks almost exactly like our pseudocode design

Starting up Mininet, Frenetic and our application respectively, we try a `pingall` in Mininet and see the following on the console:

```
vagrant@frenetic:~/manual/programmers_guide/code/12_learning_switch$ python learning1.py
Starting the tornado event loop (does not return).
2016-04-14 12:49:17,228 [INFO] Connected to Frenetic - Switches: {1: [4, 2, 1, 3]}
2016-04-14 12:49:17,229 [INFO] Learning: 9a:0f:ec:39:54:f5 attached to ( 1 )
2016-04-14 12:49:17,258 [INFO] Learning: be:3f:5a:90:8a:ac attached to ( 2 )
2016-04-14 12:49:17,303 [INFO] Learning: 3a:a4:6b:e6:24:25 attached to ( 3 )
2016-04-14 12:49:17,343 [INFO] Learning: f2:a7:c0:cb:90:23 attached to ( 4 )
```

The switch works perfectly! But it's a huge violation of Principle 1.

4.3 A More Efficient Switch

Once we've learned a MAC-to-port mapping, we shouldn't have to go to the controller for packets destined for that MAC. The switch should handle it by itself.

This is actually pretty straightforward. If we know that MAC 11:11:11:11:11:11 is on port 2, we can handle it with the following NetKAT program:

```
Filter(EthDstEq("11:11:11:11:11:11")) >> SetPort(2)
```

And we just need one of these rules for each MAC we've learned. But all of these rules are non-overlapping because they involve different values for `EthDst`. So we just Union them all together and that's our entire NetKAT program.

So let's write some methods for calculating the policies We'll add this code to `learning1` (listed in `code/12/_learning/_switch/learning2.py`:

```
def policy_for_dest(self, mac_port):
    (mac, port) = mac_port
    return Filter(EthDstEq(mac)) >> SetPort(port)

def policies_for_dest(self, all_mac_ports):
    return [ self.policy_for_dest(mp) for mp in all_mac_ports ]

def policy(self):
    return Union( self.policies_for_dest(self.nib.all_mac_port_pairs()) )
```

When shall we install these rules? We could install them on every incoming packet, but that's a little overkill. We really only need to recalculate them when we see a newly learned MAC and port. So we add them to that conditional:

```
# If we haven't learned the source mac, do so
if nib.port_for_mac( src_mac ) == None:
    nib.learn( src_mac, port_id)
    self.update(self.policy())
```

Now run it and try a pingall from Mininet:

```
vagrant@frenetic:~/manual/programmers_guide/code/12_learning_switch$ python learning1.py
Starting the tornado event loop (does not return).
2016-04-14 13:33:22,965 [INFO] Connected to Frenetic - Switches: {1: [2, 4, 1, 3]}
2016-04-14 13:33:26,447 [INFO] Learning: 86:d8:df:f0:95:75 attached to ( 1 )
2016-04-14 13:33:26,453 [INFO] Learning: 4a:1c:9e:9b:50:7c attached to ( 2 )
... STOP
```

Uh oh. Why did we only learn the first two ports? Let's look at the Frenetic console for a clue:

```
[DEBUG] Installing policy
drop |
(filter ethDst = 4a:1c:9e:9b:50:7c; port := 2 |
 filter ethDst = 86:d8:df:f0:95:75; port := 1)
[DEBUG] Setting up flow table
+-----+
| 1 | Pattern                | Action      |
+-----+
| EthDst = 4a:1c:9e:9b:50:7c | Output(2)   |
+-----+
| EthDst = 86:d8:df:f0:95:75 | Output(1)   |
+-----+
|                               |              |
+-----+
```

Can you see the problem? There's no longer a rule to send packets to the controller. If packets are destined for the first two MAC addresses, that's not a problem, but if they're not, it is definitely a problem.

One thing that definitely *won't* work is to add the following rule with a Union:

```
id >> SendToController("learning_app")
```

The `id` filter matches all packets, and therefore overlaps every other rule. Even if we place this rule as the last rule in a set of Unions, *that does not guarantee it'll be fired last*. Frenetic does not guarantee the OpenFlow rules will follow the order of the NetKAT rules.

There are a few ways to solve this problem, but we'll try an easy one first. In Chapter 2, we mentioned briefly that for every `FieldEq` NetKAT predicate, there is a corresponding `FieldNotEq` predicate. We can use that in our policy, as we see in `learning3.py`:

```
def policy(self):
    return \
        IfThenElse(
            EthDstNotEq( self.nib.all_learned_macs() ),
            SendToController("learning_app"),
            Union( self.policies_for_dest(self.nib.all_mac_port_pairs()) )
        )
```

TODO: The first ping from h1 to h2 fails the first time, even though all the subsequent ones work, and the second pingall works 100 percent.

4.4 Timeouts and Moves

Our learning switch works fine if MAC to port assignments never change. But a network is usually more fluid than that:

- People unplug a host from one port and plug it into another. In our application, packets will continue to go to the old port.
- People replace one host and MAC with another host and MAC in the same port. In our application, the old MAC will continue to take up rule space on the switch.

TODO: These can be handled with portup and portdown hooks. But it's unclear that they fully solve the problem. And it's not clear they solve malicious users from forging an existing MAC and hijacking traffic.

Chapter 5

Proxy ARP

5.1 The ARP Protocol

5.2 Design

5.3 Snooping on ARP Requests and Replies

5.4 Replying

5.5 Maintaining State Across Restarts

Chapter 6

Handling Vlan

6.1 Vlan Uses

6.2 Design

6.3 Maintaining Separate Vlan Tables

6.4 Tagging and Untagging

Chapter 7

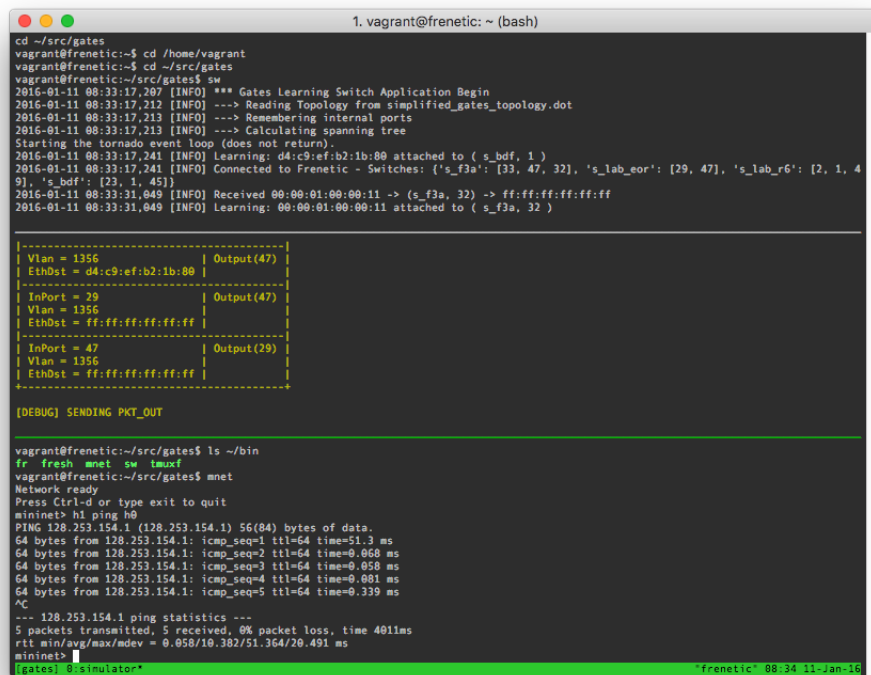
Gathering Statistics

Chapter 8

SDN Development Tools

8.1 Tmux

Tmux stands for *terminal multiplexor*, and it's indispensable for all kinds of multi-program development. It's a Window Manager for the command line, of sorts. With tmux you can split the screen into *panes*, each of which can run a different shell.



```
1. vagrant@frenetic: ~ (bash)
cd ~/src/gates
vagrant@frenetic:~$ cd /home/vagrant
vagrant@frenetic:~$ cd ~/src/gates
vagrant@frenetic:~/src/gates$ sw
2016-01-11 08:33:17,207 [INFO] *** Gates Learning Switch Application Begin
2016-01-11 08:33:17,212 [INFO] ---> Reading Topology from simplified_gates_topology.dot
2016-01-11 08:33:17,213 [INFO] ---> Remembering internal ports
2016-01-11 08:33:17,213 [INFO] ---> Calculating spanning tree
Starting the tornado event loop (does not return).
2016-01-11 08:33:17,241 [INFO] Learning: d4:c9:ef:b2:1b:80 attached to ( s_bdf, 1 )
2016-01-11 08:33:17,241 [INFO] Connected to Frenetic - Switches: {'s_f3a': [33, 47, 32], 's_lab_eor': [29, 47], 's_lab_r6': [2, 1, 4
9], 's_bdf': [23, 1, 45]}
2016-01-11 08:33:31,049 [INFO] Received 00:00:01:00:00:11 -> ( s_f3a, 32 ) -> ff:ff:ff:ff:ff:ff
2016-01-11 08:33:31,049 [INFO] Learning: 00:00:01:00:00:11 attached to ( s_f3a, 32 )

-----
| Vlan = 1356 | Output(47) |
| EthDst = d4:c9:ef:b2:1b:80 | |
-----
| InPort = 29 | Output(47) |
| Vlan = 1356 | |
| EthDst = ff:ff:ff:ff:ff:ff | |
-----
| InPort = 47 | Output(29) |
| Vlan = 1356 | |
| EthDst = ff:ff:ff:ff:ff:ff | |
-----

[DEBUG] SENDING PKT_OUT

vagrant@frenetic:~/src/gates$ ls ~/bin
fr fresh mnet sw tmuxf
vagrant@frenetic:~/src/gates$ mnet
Network ready
Press Ctrl-d or type exit to quit
mininet> h1 ping h0
PING 128.253.154.1 (128.253.154.1) 56(84) bytes of data:
64 bytes from 128.253.154.1: icmp_seq=1 ttl=64 time=51.3 ms
64 bytes from 128.253.154.1: icmp_seq=2 ttl=64 time=0.068 ms
64 bytes from 128.253.154.1: icmp_seq=3 ttl=64 time=0.050 ms
64 bytes from 128.253.154.1: icmp_seq=4 ttl=64 time=0.001 ms
64 bytes from 128.253.154.1: icmp_seq=5 ttl=64 time=0.339 ms
^C
--- 128.253.154.1 ping statistics ---
5 packets transmitted, 5 received, 0% packet loss, time 401ms
rtt min/avg/max/mdev = 0.050/10.382/51.364/20.491 ms
mininet>
[gates] 0:simulator* "frenetic" 08:34 11-Jan-16
```

In Frenetic development, it's helpful to run tmux with at least 3 panes, which you can see in the example above:

1. One pane runs your network application, e.g. `python repeater.py`

2. One pane runs Frenetic, e.g. `frenetic http-server --verbosity debug`
3. One pane runs Mininet

Tmux is very personalizable and configurable. But if you haven't used tmux before, here is a good setup to get you started.

1. Install tmux on Frenetic-vm with `sudo apt-get install tmux`
2. Create a file in your home directory, `./tmux.conf` with the line `set -g prefix C-a`. This maps the tmux prefix key to `Ctrl` + `A`, which is much easier to reach than the default `Ctrl` + `B`
3. Type `[style=BashInputStyle]tmux` to start.

`Ctrl` + `A` is called the *prefix key*, and we'll denote it as `Prefix` below. Because you will be typing the prefix a lot, it's really helpful to map your `Caps Lock` key to `Ctrl`, if you haven't already done so. On a Mac, for example, you can go to Apple Menu → System Preferences → Keyboard → Keyboard Tab → Modifier Keys and select Control for Caps Lock.





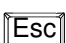
Once there, you can use the following key combinations:

- `Prefix` `=` splits the pane at the cursor into two panes: one above the cursor and one below. You can split existing panes as many times as you want, all the way down to panes with one line (which are probably not very useful).
- `Prefix` `↑` moves the cursor to the pane above. (If you're on the top pane already, the cursor moves to the bottom-most pane.)
- `Prefix` `↓` moves the cursor to the pane below. (If you're on the bottom pane already, the cursor moves to the top-most pane.)
- `Prefix` `Z` zooms the current pane, so that it takes up the entire window. The other panes continue to run, even though they're not visible. Pressing `Prefix` `Z` again unzooms the window.
- `Prefix` `D` detaches from the Tmux session. You can start it up again later, even after having logged off the Frenetic VM, by using `tmux attach`.

Because tmux operates outside the normal window manager realms, you can no longer scroll up or down in a pane using scroll bars. But tmux has a scrolling mechanism inside itself which scrolls panes independently.

- `Prefix` `I` enters scroll mode. You can see a cursor position status at the top right hand corner of the pane: 67/900 means you're on line 67 of 900 lines in the pane.

- once in scroll mode:

-  moves the cursor one line up.
-  moves the cursor one line down.
-  moves the cursor one page up.
-  moves the cursor one page down.
-  leaves scroll mode and scrolls all the way down to the bottom

If you end up doing the same keystrokes each time you start up an SDN session, you can automate it with tmuxinator software.

The book Hogan [2012] is a great introduction and reference to tmux.

8.2 Open VSwitch Utilities

8.3 TCPCDump

8.4 Mininet Network Modelling

Chapter 9

Network Address Translation

9.1 Why Do We Need NAT?

9.2 Design

Chapter 10

Spanning Tree Alternatives

10.1 What's Wrong with STP Protocols?

10.2 Design

10.3 Calculating Shortest Paths

10.4 Core/Edge Separation

Chapter 11

Routing

11.1 Design

11.2 Configuring Route Tables

11.3 Using Link State

Chapter 12

Modularization

12.1 Sharing Actions in Rules and Packet Outs

12.2 Subclassing

12.3 Multiple Network Apps

12.4 Clustering

Chapter 13

Frenetic REST API

13.1 REST URL's

13.2 Incoming, Northbound Events

13.3 Messages

Chapter 14

Frenetic/NetKAT Reference

14.1 Predicates

14.2 Policies

14.3 Event Hooks

14.4 Compiler Directives

14.5 Frenetic Command Line

Chapter 15

Productionalizing

15.1 Installing Frenetic on Bare Metal Linux

15.2 Control Scripts

15.3 Logging

Bibliography

Nate Foster, Michael J. Freedman, Arjun Guha, Rob Harrison, Naga Praveen Katta, Christopher Monsanto, Joshua Reich, Mark Reitblatt, Jennifer Rexford, Cole Schlesinger, Alec Story, and David Walker. Languages for software-defined networks. *IEEE Communications*, 51(2):128–134, Feb 2013.

Nate Foster, Dexter Kozen, Matthew Milano, Alexandra Silva, and Laure Thompson. A coalgebraic decision procedure for netkat. *SIGPLAN Not.*, 50(1): 343–355, Jan 2015. ISSN 0362-1340. doi: 10.1145/2775051.2677011. URL <http://doi.acm.org/10.1145/2775051.2677011>.

Brian P. Hogan. *tmux: Productive Mouse-Free Development*. The Pragmatic Programmers, LLC, Raleigh, NC, and Dallas, TX, 2012. ISBN 978-1-93435-696-8. URL <http://www.pragprog.com/titles/bhtmux/tmux>.

Steffen Smolka, Spiridon Aristides Eliopoulos, Nate Foster, and Arjun Guha. A fast compiler for netkat. *CoRR*, abs/1506.06378, 2015. URL <http://arxiv.org/abs/1506.06378>.