# Machine-Verified Network Controllers

Nate Foster
Cornell University

**frenetic** »

Coq
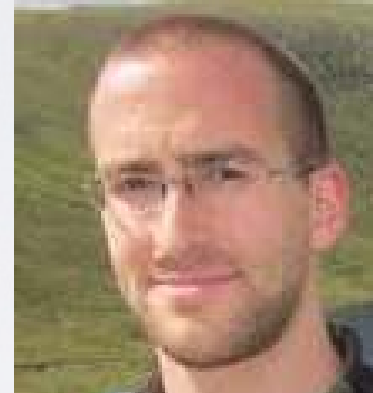
# Proof Assistants



Coq



Arjun Guha
Postdoc → UMass



Mark Reitblatt
PhD student

"The trigger for this event was a network configuration change"
—Amazon

"The service outage was due to a series of internal network events that corrupted router data tables"

—GoDaddy

# Networks in Practice

"The airline experienced a network connectivity issue..."
—United Airlines

"The airline experienced a network connectivity issue..."
—United Airlines

There are hosts...

Connected by switches...

There are also servers...

Connected by routers...

And a load balancer...

And a gateway router...

There are other ISPs...

So we need to run BGP...

# Networks in Practice

And we need a firewall to filter incoming traffic...

# Networks in Practice

There are also wireless hosts…

So we need wireless gateways...

# Networks in Practice

And yet more middleboxes for lawful intercept...

# Networks in Practice

Each color represents a different set of control plane protocols and algorithms... this is

# Software-Defined Networking

A clean-slate architecture that standardizes features and decouples forwarding from

# Software-Defined Networking

## Essential ingredients
- Decouple control and data planes
- Logically-centralized control

## Enables
- Novel functionality
- Formal reasoning

# Software-Defined Networking

## Essential ingredients

- Decouple control and data planes
- Logically-centralized control

## Enables

- Novel functionality
- Formal reasoning

# Existing Tools

There is a cottage industry in SDN configuration-checking tools...

# Existing Tools

There is a cottage industry in SDN configuration-checking tools...

- FlowChecker [SafeConfig '10]

# Existing Tools

There is a cottage industry in SDN configuration-checking tools…

- FlowChecker [SafeConfig '10]
- AntEater [SIGCOMM '11]

# Existing Tools

There is a cottage industry in SDN configuration-checking tools...

- FlowChecker [SafeConfig '10]
- AntEater [SIGCOMM '11]
- NICE [NSDI '12]

# Existing Tools

There is a cottage industry in SDN configuration-checking tools...

- FlowChecker [SafeConfig '10]
- AntEater [SIGCOMM '11]
- NICE [NSDI '12]
- Header Space Analysis [NSDI '12]

# Existing Tools

There is a cottage industry in SDN configuration-checking tools...

- FlowChecker [SafeConfig '10]
- AntEater [SIGCOMM '11]
- NICE [NSDI '12]
- Header Space Analysis [NSDI '12]
- VeriFlow [HotSDN '12]
- and many others...

# Existing Tools

There is a cottage industry in SDN configuration-checking tools...

- FlowChecker [SafeConfig '10]
- AntEater [SIGCOMM '11]
- NICE [NSDI '12]
- Header Space Analysis [NSDI '12]
- VeriFlow [HotSDN '12]
- and many others...

These are all great tools!

But they are expensive to run,
and each builds on a custom
(typically ad hoc) foundation

# Machine-Verified Controllers

## Vision

- Develop programs in a high-level language
- Reason at a high level of abstraction
- Use a compiler and run-time system to generate low-level control messages
- Machine-verified proofs of correctness

## Contributions

- NetCore compiler + optimizer
- Featherweight OpenFlow model
- General framework for establishing run-time system correctness

# OVERVIEW

# OpenFlow Switches

**Forwarding Table:** prioritized list of rules

**Rule:** pattern, actions, and counters

**Pattern:** prefix match on headers

**Action:** forward or modify

**Counters:** total bytes and packets processed

## OpenFlow

| Pattern | Action | Bytes | Packets |
|---------|--------|-------|---------|
| 1010 | Drop | 200 | 10 |
| 010* | Forward(2) | 100 | 4 |
| 011* | Controller | 0 | 0 |

Priority ↓

**NOX**

**Network Events**
- Topology changes
- **Diverted packets**
- Traffic statistics

**Control Messages**
- **Modify rules**
- Query counters

Controller

OpenFlow Switch

OpenFlow Switch

OpenFlow Switch

# Issue #1: Switch-Level Errors

What happens if...
- The controller misses a keep-alive message?
- The controller sends a malformed message?
  - Bad output port
  - Too many actions
  - Inconsistent actions
  - Unsupported actions
- The switches runs out of space for rules?

Any of these can lead to essentially arbitrary behavior

# Issue #2: Malformed Patterns

What happens if the controller sends the following message to a switch?

```
FlowMod AddFlow { match = { srcIPAddress = 10.0.1.*", ... },
                  actions = [ flood ], ... }
```

# Issue #2: Malformed Patterns

What happens if the controller sends the following message to a switch?

```
FlowMod AddFlow { match = { srcIPAddress = 10.0.1.*", ... },
                  actions = [ flood ], ... }
```

We'd expect the switch to install a rule that broadcasts all traffic from a host the given subnet...

# Issue #2: Malformed Patterns

What happens if the controller sends the following message to a switch?

```
FlowMod AddFlow { match = { srcIPAddress = 10.0.1.*", ... },
                  actions = [ flood ], ... }
```

We'd expect the switch to install a rule that broadcasts all traffic from a host the given subnet...

...but it actually installs a rule that floods *all* traffic

# Issue #2: Malformed Patterns

What happens if the controller sends the following message to a switch?

```
FlowMod AddFlow { match = { srcIPAddress = 10.0.1.*", ... },
                  actions = [ flood ], ... }
```

We'd expect the switch to install a rule that broadcasts all traffic from a host the given subnet...

...but it actually installs a rule that floods *all* traffic

Why? Switches *silently* ignore IP fields unless the Ethernet frame type is IP!

# Issue #3: Message Reordering

What happens if the controller sends the following pair of OpenFlow messages to a switch in sequence?

```
FlowMod AddFlow { match = { ethFrameType = ethTypeIP,
                            srcIPAddress =
"10.0.1.99", ... },
                  priority = 1,
                  actions = [ ] }

FlowMod AddFlow { match = { ethFrameType = ethTypeIP,
                            srcIPAddress = "10.0.1.*", ... },
                  priority = 2,
                  actions = [ flood ] }
```

The intention is to encode a negation...

# Issue #3: Message Reordering

What happens if the controller sends the following pair of OpenFlow messages to a switch in sequence?

```
FlowMod AddFlow { match = { ethFrameType = ethTypeIP,
                            srcIPAddress =
"10.0.1.99", ... },
                  priority = 1,
                  actions = [ ] }
FlowMod AddFlow { match = { ethFrameType = ethTypeIP,
                            srcIPAddress = "10.0.1.*", ... },
                  priority = 2,
                  actions = [ flood ] }
```

The intention is to encode a negation…

…but the switch may process these in either order!

# MACHINE-VERIFIED CONTROLLERS

## Syntax

```
Inductive pred : Type :=              (* Predicates *)
    | OnSwitch : Switch -> pred
    | IngressPort : Port -> pred
    | DlSrc : EthernetAddress -> pred
    | DlDst : EthernetAddress -> pred
    | DlVlan : option VLAN -> pred
    | ...
    | And : pred -> pred -> pred
    | Or : pred -> pred -> pred
    | Not : pred -> pred
    | All : pred
    | None : pred.

Inductive PseudoPort : Type :=        (* Psuedo Ports *)
    | PhysicalPort : Port -> PseudoPort
    | AllPorts : PseudoPort.

Inductive act : Type :=               (* Actions *)
    | FwdMod : Mod -> PseudoPort -> act

Inductive pol : Type :=               (* Policies *)
    | Policy : pred -> list act -> pol
    | Union : pol -> pol -> pol
    | Restrict : pol -> pred -> pol.
```

NetCore

Compiler

Flow tables — Optimizer

Run-time system

OpenFlow messages

Featherweight OpenFlow

## Semantics

$$lp = (sw, pt, pk)$$
$$lps_{out} = pol(sw, pt, pk)$$
$$S = \{\!| (T(sw, pt_{out}), pk) \mid (pt_{out}, pk) \in lps_{out} |\!\}$$
$$\overline{\{\!|lp|\!\} \uplus \{\!|lp_1 \cdots lp_n|\!\} \xrightarrow{lp} S \uplus \{\!|lp_1 \cdots lp_n|\!\}}$$

NetCore

Compiler

Flow tables → Optimizer

Run-time system

OpenFlow messages

Featherweight OpenFlow

- Models hop-by-hop forwarding behavior of the network
- Abstracts away from the underlying distributed system
- Makes it easy to reason about network-wide properties

# NetCore to Flow Tables

## Example

| Priority | Pattern | Action |
|----------|---------|--------|
| 65534 | $\mathbf{inPort} = 2, \mathbf{dlSrc} = $ `dc:ba:65:43:21` | **Fwd** 2 |
| 65533 | $\mathbf{inPort} = 2$ | **Fwd** 3 |

## NetCore compiler
- Key operation: flow table intersection
- Must restrict to "valid" patterns

## Optimizer
- Optimizer prunes (many) redundant rules
- Based on simple algebra of operations

## Correctness Theorem

NetCore ~ FlowTable

NetCore

**Compiler**

Flow tables **Optimizer**

Run-time system

OpenFlow messages

Featherweight OpenFlow

# Valid Patterns

```
Inductive ValidPattern : Pattern -> Prop :=
    | SupportedIPPatternValid : forall dlSrc dlDst dlVlan dlVlanPcp nwSrc nwDst nwTos
                                tpSrc tpDst inPort nwProto,
      In nwProto SupportedL4Protos ->
      ValidPattern (MkPattern dlSrc dlDst (WildcardExact Const_0x800)
                              dlVlan dlVlanPcp
                              nwSrc nwDst (WildcardExact nwProto)
                              nwTos tpSrc tpDst inPort)
    | UnsupportedIPPatternValid : forall dlSrc dlDst dlVlan dlVlanPcp nwSrc nwDst nwTos
                                   inPort nwProto,
      ~ In nwProto SupportedL4Protos ->
      ValidPattern (MkPattern dlSrc dlDst (WildcardExact Const_0x800)
                              dlVlan dlVlanPcp
                              nwSrc nwDst (WildcardExact nwProto)
                              nwTos WildcardAll WildcardAll inPort)
    | ARPPacketValid : forall dlSrc dlDst dlVlan dlVlanPcp nwSrc nwDst inPort,
      ValidPattern (MkPattern dlSrc dlDst (WildcardExact Const_0x806)
                              dlVlan dlVlanPcp
                              nwSrc nwDst WildcardAll
                              WildcardAll WildcardAll WildcardAll inPort)
    | UnknownDlTypPatternValid : forall dlSrc dlDst dlTyp dlVlan dlVlanPcp inPort,
      ValidPattern (MkPattern dlSrc dlDst dlTyp
                              dlVlan dlVlanPcp
                              WildcardAll WildcardAll WildcardAll
                              WildcardAll WildcardAll WildcardAll inPort)
    | EmptyPatternValid :
      ValidPattern Pattern_empty.
```

# OpenFlow Specification

42 pages...

...of informal English text

...and C struct definitions

# Featherweight OpenFlow

## Syntax

| | | | |
|---|---|---|---|
| **Devices** | Switch | $S$ | $::= \mathbb{S}(sw, pts, RT, inp.outp, inm, out$ |
| | Controller | $C$ | $::= \mathbb{C}(\sigma, f_{in}, f_{out})$ |
| | Link | $L$ | $::= \mathbb{L}(loc_{src}, pks, loc_{dst})$ |
| | OpenFlow Link to Controller | $M$ | $::= \mathbb{M}(sw, SMS, CMS)$ |
| **Packets and Locations** | Packet | $pk$ | $::= abstract$ |
| | Switch ID | $sw$ | $\in \mathbb{N}$ |
| | Port ID | $pt$ | $\in \mathbb{N}$ |
| | Location | $loc$ | $\in sw \times pt$ |
| | Located Packet | $lp$ | $\in loc \times pk$ |
| **Controller Components** | Controller state | $\sigma$ | $::= abstract$ |
| | Controller input relation | $f_{in}$ | $\in sw \times CM \times \sigma \rightsquigarrow \sigma$ |
| | Controller output relation | $f_{out}$ | $\in \sigma \rightsquigarrow sw \times SM \times \sigma$ |
| **Switch Components** | Rule table | $RT$ | $::= abstract$ |
| | Rule table Interpretation | $[\![RT]\!]$ | $\in lp \rightarrow \{\!|lp_1 \cdots lp_n|\!\} \times \{\!|CM_1 \cdots C$ |
| | Rule table modifier | $\Delta RT$ | $::= abstract$ |
| | Rule table modifier interpretation | $apply$ | $\in \Delta RT \rightarrow RT \rightarrow \Delta RT$ |
| | Ports on switch | $pts$ | $\in \{pt_1 \cdots pt_n\}$ |
| | Consumed packets | $inp$ | $\in \{\!|lp_1 \cdots lp_n|\!\}$ |
| | Produced packets | $outp$ | $\in \{\!|lp_1 \cdots lp_n|\!\}$ |
| | Messages from controller | $inm$ | $\in \{\!|SM_1 \cdots SM_n|\!\}$ |
| | Messages to controller | $outm$ | $\in \{\!|CM_1 \cdots CM_n|\!\}$ |
| **Link Components** | Endpoints | $loc_{src}, loc_{dst}$ | $\in loc$ where $loc_{src} \neq loc_{dst}$ |
| | Packets from $loc_{src}$ to $loc_{dst}$ | $pks$ | $\in [pk_1 \cdots pk_n]$ |
| **Controller Link** | Message queue from controller | $SMS$ | $\in [SM_1 \cdots SM_n]$ |
| | Message queue to controller | $CMS$ | $\in [CM_1 \cdots CM_n]$ |
| **Abstract OpenFlow Protocol** | Message from controller | $SM$ | $::= \textbf{FlowMod } \Delta RT \mid \textbf{PktOut } pt \; p$ |
| | Message to controller | $CM$ | $::= \textbf{PktIn } pt \; pk \mid \textbf{BarrierReply } n$ |

## Semantics

$$\frac{(outp', outm') = [\![RT]\!](lp)}{\mathbb{S}(sw, pts, RT, \{\!|lp|\!\} \uplus inp, outp, inm, outm) \xrightarrow{lp} \mathbb{S}(sw, pts, RT, inp, outp' \uplus outp, inm, outm' \uplus outm)} \text{ (Pkt-Process)}$$

$$\frac{}{\begin{array}{l} \mathbb{S}(sw, pts, RT, inp, \{\!|(sw, pt, pk)|\!\} \uplus outp, inm, outm) \mid \mathbb{L}((sw, pt), pks, loc') \\ \longrightarrow \mathbb{S}(sw, pts, RT, inp, outp, inm, outm) \mid \mathbb{L}((sw, pt), [pk] \mathbin{+\!\!+} pks, loc') \end{array}} \text{ (Send-Wire)}$$

$$\frac{}{\begin{array}{l} \mathbb{L}(loc, pks \mathbin{+\!\!+} [pk], (sw, pt)) \mid \mathbb{S}(sw, pts, RT, inp, outp, inm, outm) \\ \xrightarrow{(sw, pt, pk)} \mathbb{L}(loc, pks, (sw, pt)) \mid \mathbb{S}(sw, pts, RT, \{\!|(sw, pt, pk)|\!\} \uplus inp, outp, inm, outm) \end{array}} \text{ (Recv-Wire)}$$

$$\frac{RT' = apply(\Delta RT, RT)}{\mathbb{S}(sw, pts, RT, inp, outp, \{\!|\textbf{FlowMod } \Delta RT|\!\} \uplus inm, outm) \longrightarrow \mathbb{S}(sw, pts, RT', inp, outp, inm, outm)} \text{ (Switch-FlowMod)}$$

$$\frac{pt \in pts}{\mathbb{S}(sw, pts, RT, inp, outp, \{\!|\textbf{PktOut } pt \; pk|\!\} \uplus inm, outm) \longrightarrow \mathbb{S}(sw, pts, RT, inp, \{\!|(sw, pt, pk)|\!\} \uplus outp, inm, outm)} \text{ (Switch-PktOut)}$$

$$\frac{f_{out}(\sigma) \rightsquigarrow (sw, SM, \sigma')}{\mathbb{C}(\sigma, f_{in}, f_{out}) \mid \mathbb{M}(sw, SMS, CMS) \longrightarrow \mathbb{C}(\sigma', f_{in}, f_{out}) \mid \mathbb{M}(sw, [SM] \mathbin{+\!\!+} SMS, CMS)} \text{ (Ctrl-Send)}$$

$$\frac{f_{in}(sw, \sigma, CM) \rightsquigarrow \sigma'}{\mathbb{C}(\sigma, f_{in}, f_{out}) \mid \mathbb{M}(sw, SMS, CMS \mathbin{+\!\!+} [CM]) \longrightarrow \mathbb{C}(\sigma', f_{in}, f_{out}) \mid \mathbb{M}(sw, SMS, CMS)} \text{ (Ctrl-Recv)}$$

$$\frac{SM \neq \textbf{BarrierRequest } n}{\begin{array}{l} \mathbb{M}(sw, SMS \mathbin{+\!\!+} [SM], CMS) \mid \mathbb{S}(sw, pts, RT, inp, outp, inm, outm) \\ \longrightarrow \mathbb{M}(sw, SMS, CMS) \mid \mathbb{S}(sw, pts, RT, inp, outp, \{\!|SM|\!\} \uplus inm, outm) \end{array}} \text{ (Switch-Recv-Ctrl)}$$

$$\frac{}{\begin{array}{l} \mathbb{M}(sw, SMS \mathbin{+\!\!+} [\textbf{BarrierRequest } n], CMS) \mid \mathbb{S}(sw, pts, RT, inp, outp, \emptyset, outm) \\ \longrightarrow \mathbb{M}(sw, SMS, CMS) \mid \mathbb{S}(sw, pts, RT, inp, outp, \emptyset, \{\!|\textbf{BarrierReply } n|\!\} \uplus outm) \end{array}} \text{ (Switch-Recv-Barrier)}$$

$$\frac{}{\begin{array}{l} \mathbb{S}(sw, pts, RT, inp, outp, inm, \{\!|CM|\!\} \uplus outm) \mid \mathbb{M}(sw, SMS, CMS) \\ \longrightarrow \mathbb{S}(sw, pts, RT, inp, outp, inm, outm) \mid \mathbb{M}(sw, SMS, [CM] \mathbin{+\!\!+} CMS) \end{array}} \text{ (Switch-Send-Ctrl)}$$

## Key judgments:

- Controller in: $(sw, CM, \sigma) \rightsquigarrow \sigma'$
- Controller out: $\sigma \rightsquigarrow (sw, SM, \sigma')$
- Network step: $M \rightarrow M'$

Models *all* essential asynchrony

# Run-Time System

## Invariants
- Maintain a sound approximation of overall flow table each switch
- Eventually process all diverted packets

## Theorem

FlowTable ≈ Featherweight OpenFlow

## Run-time instances
- Trivial: processes all packets on controller
- Proactive: installs rules, falls back to Trivial when out of space
- Full: like Proactive, but also installs exact-match rules

NetCore

Compiler

Flow tables

Optimizer

**Run-time system**

OpenFlow messages

Featherweight OpenFlow

# Safe Wires

```
Inductive SafeWire : SF -> SF -> SF -> list CM -> Prop :=
| SafeWire_nil : forall lb ub,
    extends ub lb ->
    SafeWire lb ub lb nil
| SafeWire_cons_FlowMod : forall lb ub sf sft lst,
    SafeWire lb ub sf lst ->
    extends ub (apply_SFT sft sf) ->
    SafeWire lb ub (apply_SFT sft sf) (FlowMod sft :: lst)
| SafeWire_cons_PktOut : forall lb ub sf pt pk lst,
    SafeWire lb ub sf lst ->
    SafeWire lb ub sf (PktOut pt pk :: lst)
| SafeWire_cons_BarrierRequest : forall lb ub sf n lst,
    SafeWire lb ub sf lst ->
    SafeWire lb ub sf (BarrierRequest n :: lst).
```

# Implementation

## Source

- ~8,000 lines of Coq
- ~1,500 lines of Haskell

## Components

- NetCore compiler and optimizer
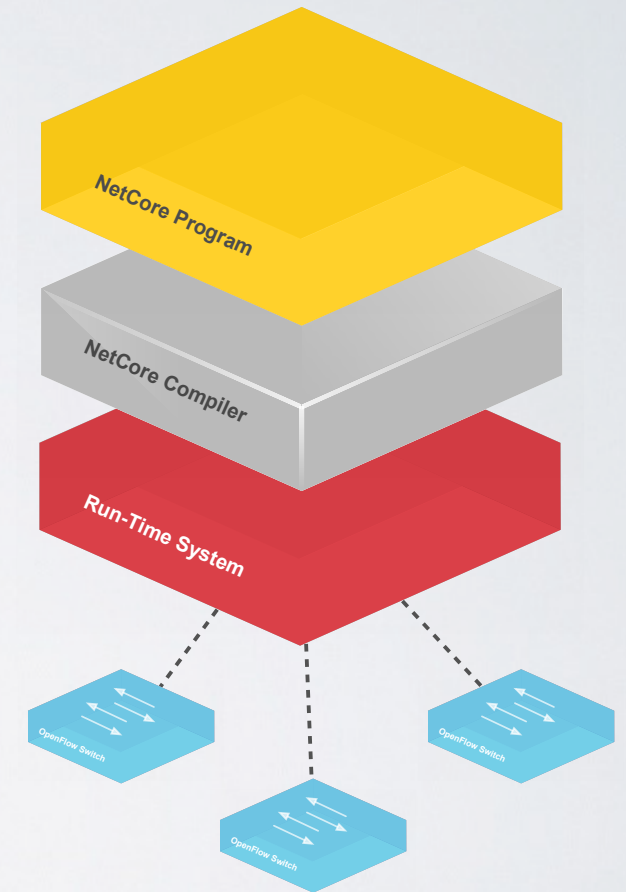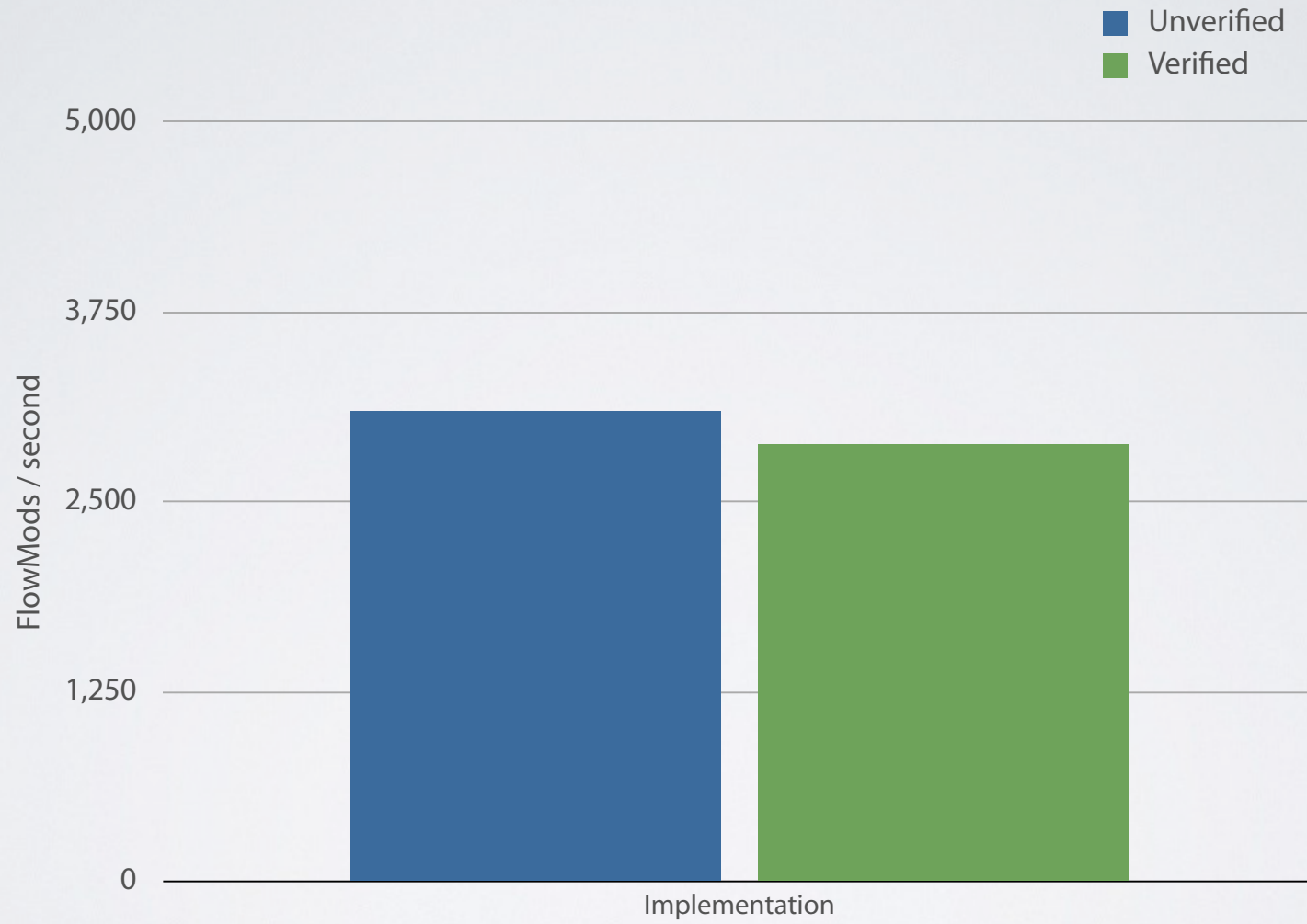- Flow tables
- Featherweight OpenFlow
- Run-time system instances
- Proofs of correctness

## Status

- Extracts to Haskell source code
- Compiles against Nettle libraries
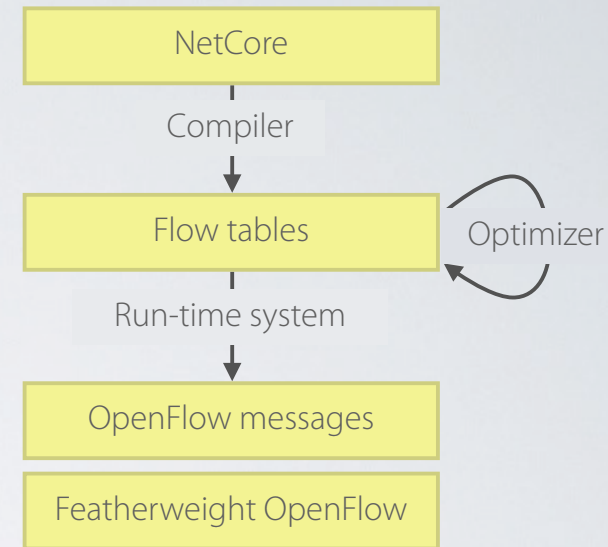- Running on "production" traffic in the lab

# Performance

# Conclusion

Networks are critical infrastructure…

…developed using 1970s-era techniques

Software-defined networks are an architecture that could be used to put networks on a solid foundation

Machine-verified controllers based on NetCore a first step in this direction

# A Grand Collaboration: Languages + Networking

**Frenetic Cornell**

Shrutarshi Basu (PhD)
Nate Foster (Faculty)
**Arjun Guha** (Postdoc)
Stephen Gutz (Undergrad)
**Mark Reitblatt** (PhD)
Robert Soulé (Postdoc)
Alec Story (Undergrad)

**Frenetic Princeton**

Chris Monsanto (PhD)
Joshua Reich (Postdoc)
Jen Rexford (Faculty)
Cole Schlesinger (PhD)
Dave Walker (Faculty)
Naga Praveen Katta (PhD)

frenetic >>

http://frenetic-lang.org