# Taurus: A Data Plane Architecture for Per-Packet ML

**Tushar Swamy**

Alexander Rucker, Muhammad Shahbaz, Ishan Gaur, and Kunle Olukotun

Stanford University

----------------------------------------------------------------------

" *Our current generation — Jupiter fabrics — can deliver more than 1 Petabit/sec of total bisection bandwidth* "

----------------------------------------------------------------------

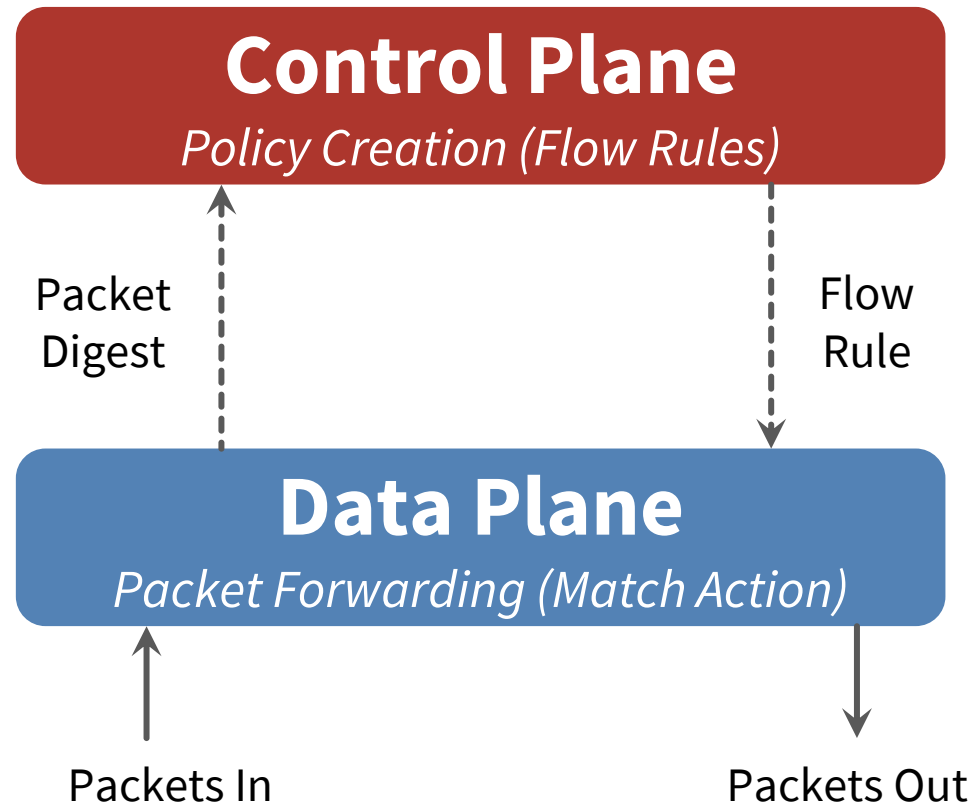— A Look Inside Google's Data Center Networks[1]

**Networks require complex management with high performance**

2

# Automate decision-making with machine learning (ML)

- Making decisions based on data ➝ ***machine learning***

- Machine learning can:

  - ***Approximate*** network functions based on data
  - ***Customize*** network functions based on data

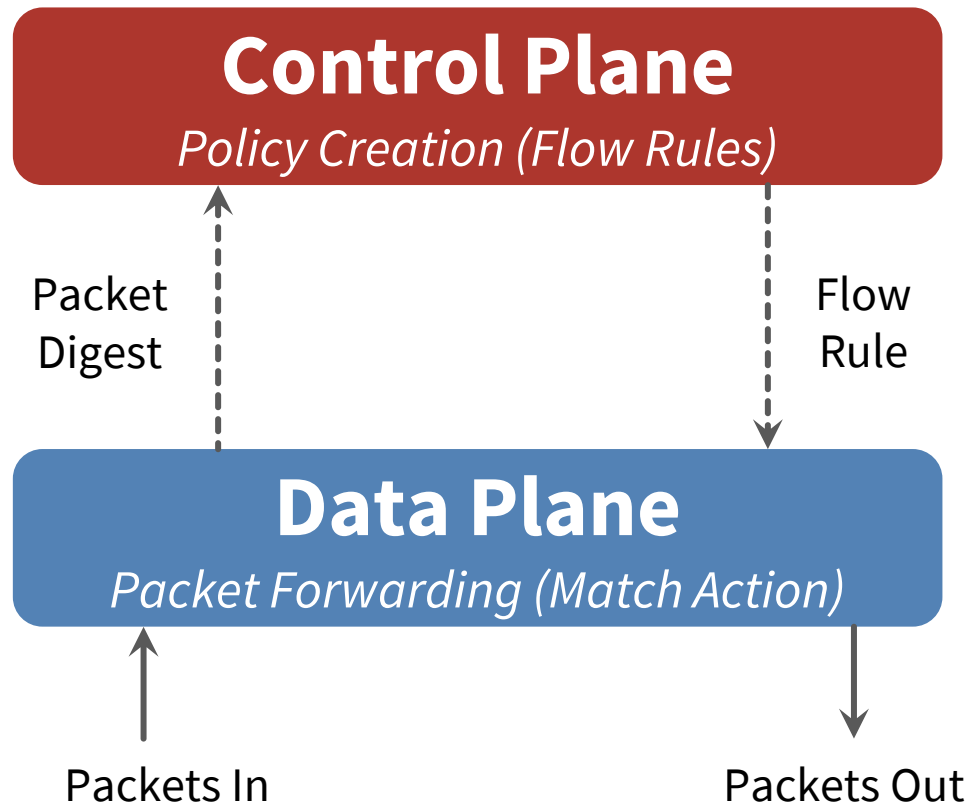- Currently, we use by hand-written heuristics in the network…

3

*Software Defined Network*

# A Taurus network introduces ML for management

**Software Defined Network**

**Control Plane**
*Policy Creation (Flow Rules)*

Packet Digest

Flow Rule

**Data Plane**
*Packet Forwarding (Match Action)*

Packets In

Packets Out

**Software Defined Network with Taurus**

**Control Plane**
*Policy Creation (Flow Rules + ML Training)*

Packet Digest

Flow Rule

ML model weights

**Data Plane**
*Packet Forwarding (Match Action) + Decision Making (ML Inference)*

Packets In

Packets Out

5

# ML inference should happen *per-packet* in the *data plane*

Processing time: **0.5 ms**
Packets missed: **1000 K**



Control Plane

Flow rule

Packet digest

Data Plane

*1.5 M Packets missed during flow rule installation time*

# Robustness and performance of the network are determined by:

→ *Quality of reaction*

→ *Speed of reaction*

**Software Defined Network with Taurus**

**ML Training is off critical path**

**Control Plane**
*Policy Creation (Flow Rules + ML Training)*

Packet Digest

Flow Rule

ML model weights

**Data Plane**
*Packet Forwarding (Match Action) + Decision Making (ML Inference)*

Packets In

Packets Out

**Software Defined Network
with Taurus**

**Control Plane**
*Policy Creation (Flow Rules + ML Training)*

Packet Digest

Flow Rule

ML model weights

***ML Inference is on critical path***

**Data Plane**
*Packet Forwarding (Match Action) + Decision Making (ML Inference)*

Packets In

Packets Out

10

*Taurus* is an architecture for per-packet ML inference in the data plane

# What do programmable switches look like?



Packets In → Packet Parser → Match-Action Tables → Traffic Manager → Packets Out
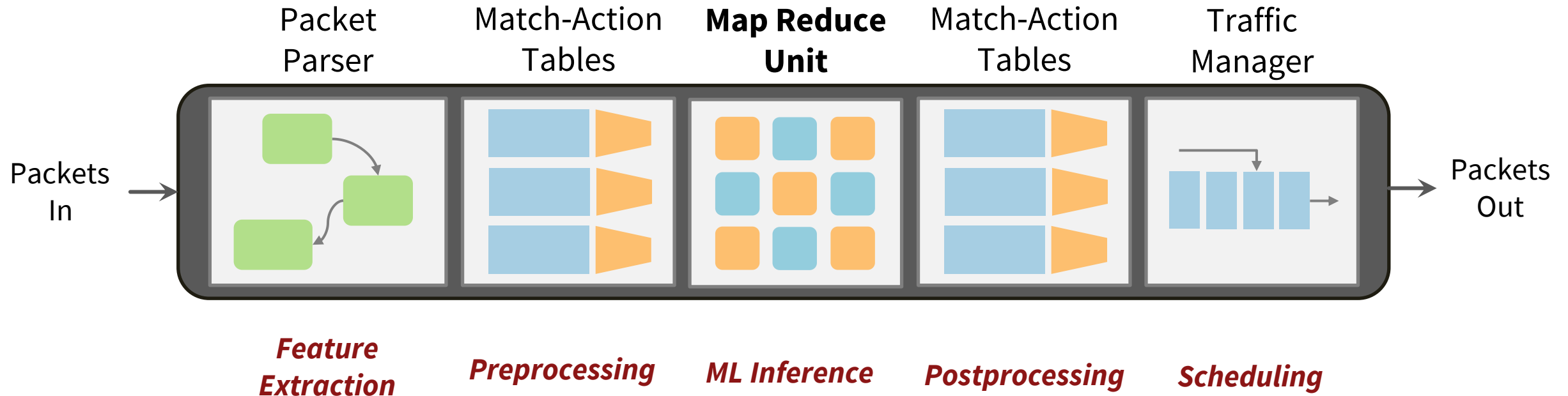
*A Protocol Independent Switch Architecture (PISA)*

# What abstraction should we use?

- ***Map-reduce*** can support linear algebra operations common in ML
  - Neural networks, SVMs, etc.

- ***SIMD Parallelism*** enables performance with minimal logic

- ***Unrolling*** patterns allows for flexibility
  - More unrolling ⟶ better performance
  - Less unrolling ⟶ less resource usage

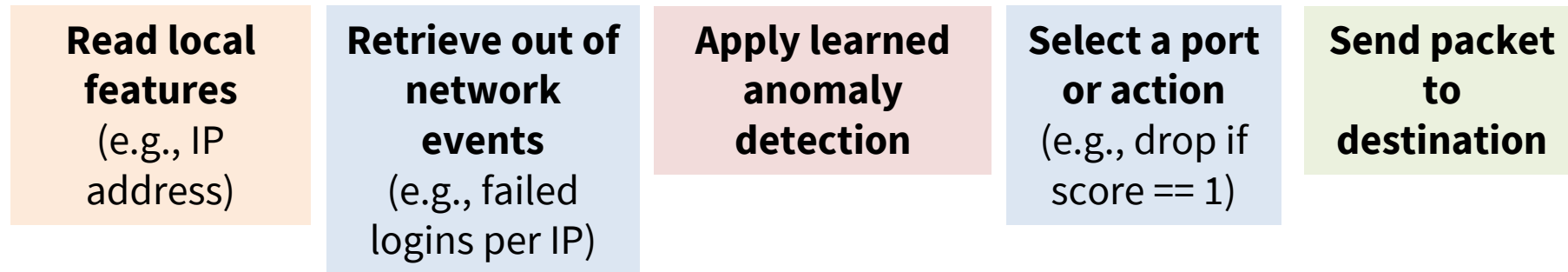# The Taurus pipeline with a Map Reduce Unit

Packet
Parser

Match-Action
Tables

**Map Reduce
Unit**

Match-Action
Tables

Traffic
Manager

Packets
In

Packets
Out

*Feature
Extraction*

*Preprocessing*

*ML Inference*

*Postprocessing*

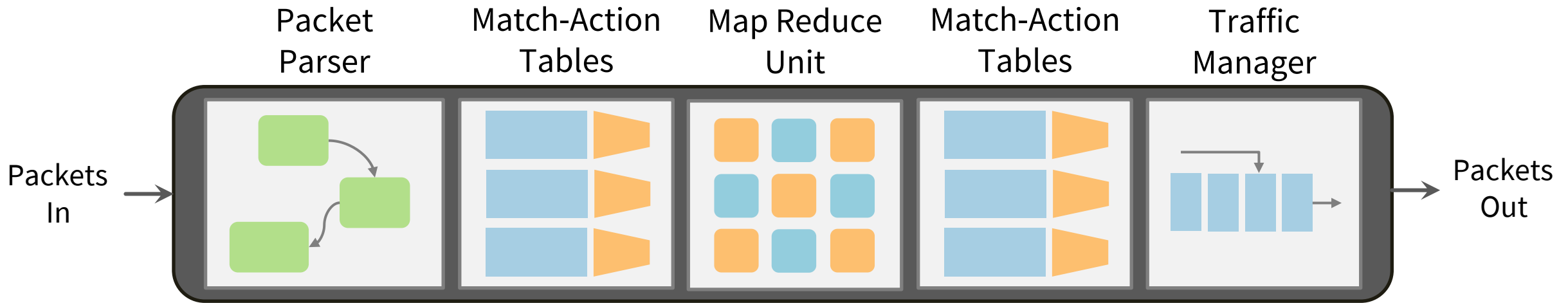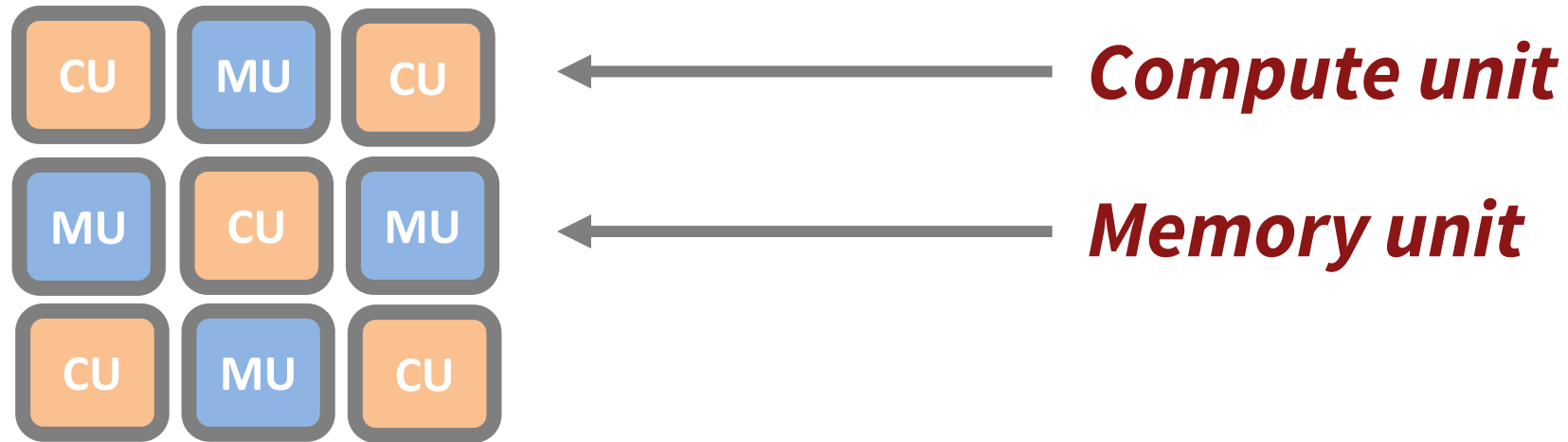*Scheduling*

- ***Map Reduce Unit*** must:
  - be reconfigurable
  - meet line rate (with a fixed clock)
  - incur minimal area and power overhead

14

# Example Application: Anomaly Detection

| Packet Parser | Match-Action Tables | Map Reduce Unit | Match-Action Tables | Traffic Manager |
|---|---|---|---|---|

Packets In → ... → Packets Out

**Packet**

| Read local features (e.g., IP address) | Retrieve out of network events (e.g., failed logins per IP) | Apply learned anomaly detection | Select a port or action (e.g., drop if score == 1) | Send packet to destination |
|---|---|---|---|---|

- Our evaluation platform is based on ***Plasticine***

- We program our map-reduce applications in the ***Spatial HDL***

| CU | MU | CU |
|----|----|----|
| MU | CU | MU |
| CU | MU | CU |

← ***Compute unit***

← ***Memory unit***

***More architectural details in full paper!***

16

# Evaluation of a Taurus ASIC
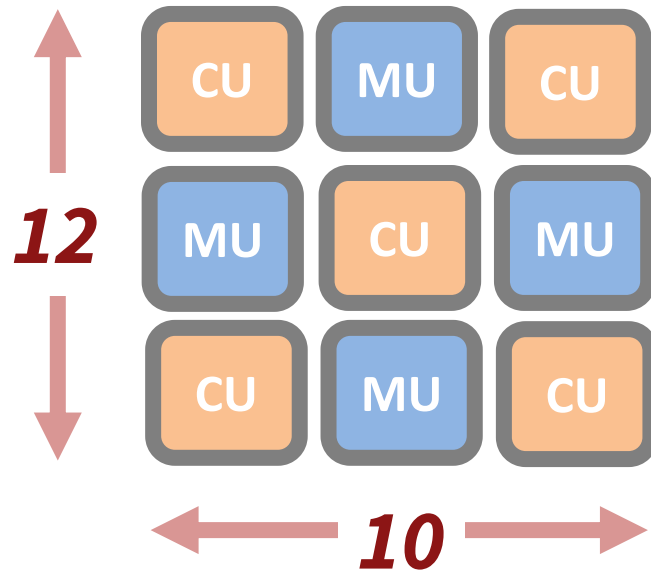
- Our evaluation platform is based on ***Plasticine***

- We program our map-reduce applications in the ***Spatial HDL***



| Hardware | Area | |
|---|---|---|
| | mm$^2$ | +% |
| 12x10 MR Grid | 4.8 x 4 | 3.8 |
| Prog. Switch | 500 | --- |

*Overheads are calculated relative to state of the art programmable switches*
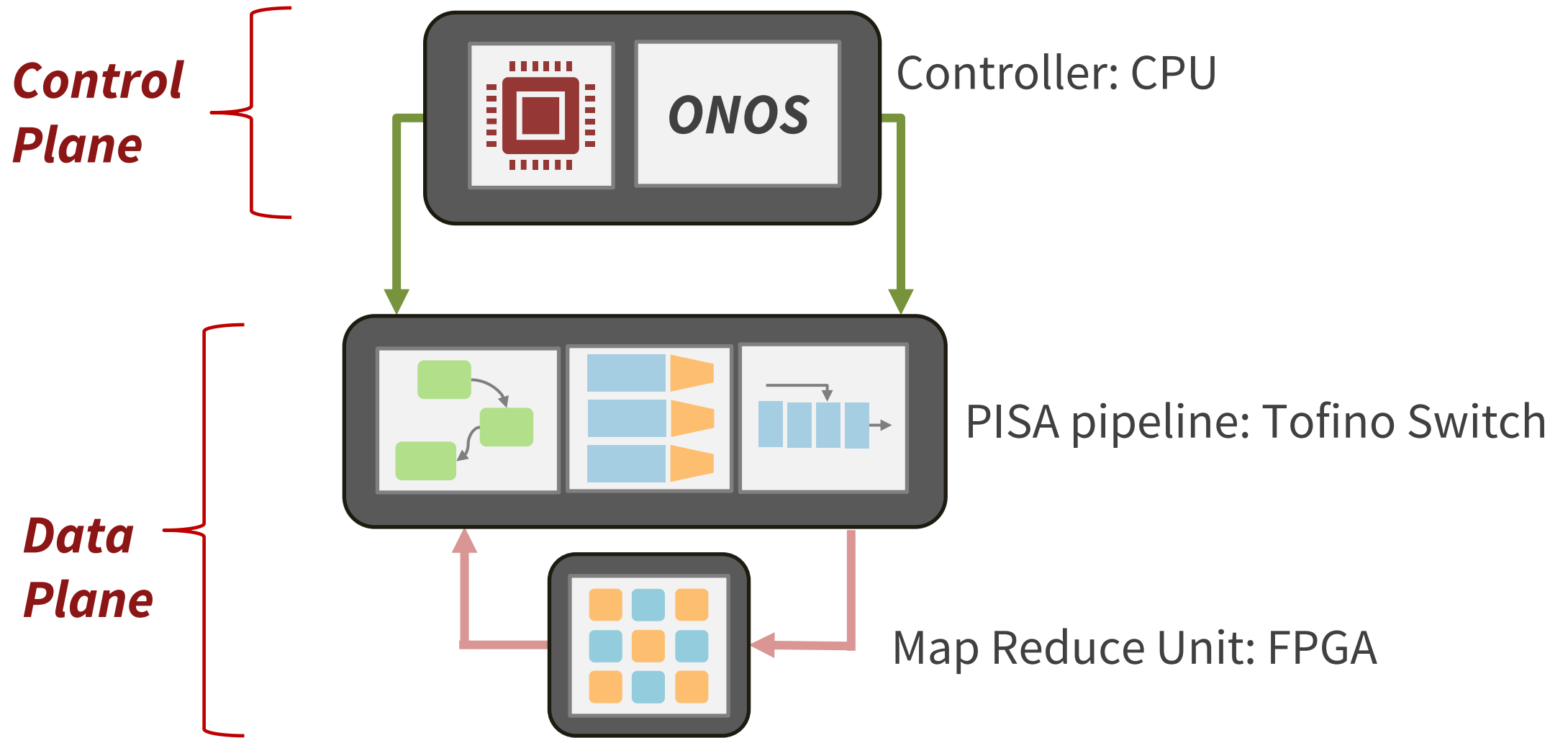
# Evaluation of an Anomaly Detection (AD) benchmark

- ***AD SVM: 8 support vectors***
- ***AD DNN: 4 layers - 12x6x3x2 neurons***

***Overhead of Map Reduce Unit***

| Model | TP (GPkt/s) | Lat (ns) | Area +% | Power +% |
|-------|-------------|----------|---------|----------|
| SVM | 1 | 83 | 0.5 | 0.6 |
| DNN | 1 | 221 | 0.8 | 1.0 |

*Overheads are calculated relative to state of the art programmable switches*

***More apps in full paper!***

**Control Plane**

Controller: CPU

**Data Plane**

PISA pipeline: Tofino Switch

Map Reduce Unit: FPGA

19

# Questions?

Try it out!

https://gitlab.com/dataplane-ai/taurus