# Taurus: An Intelligent Data Plane

## Muhammad Shahbaz

Tushar Swamy, Alex Rucker, and Kunle Olukotun

# Taurus: An Intelligent Data Plane

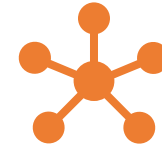Programmable Data Plane

fused with

Machine Intelligence

Deep Learning

# Managing Networks is Hard!
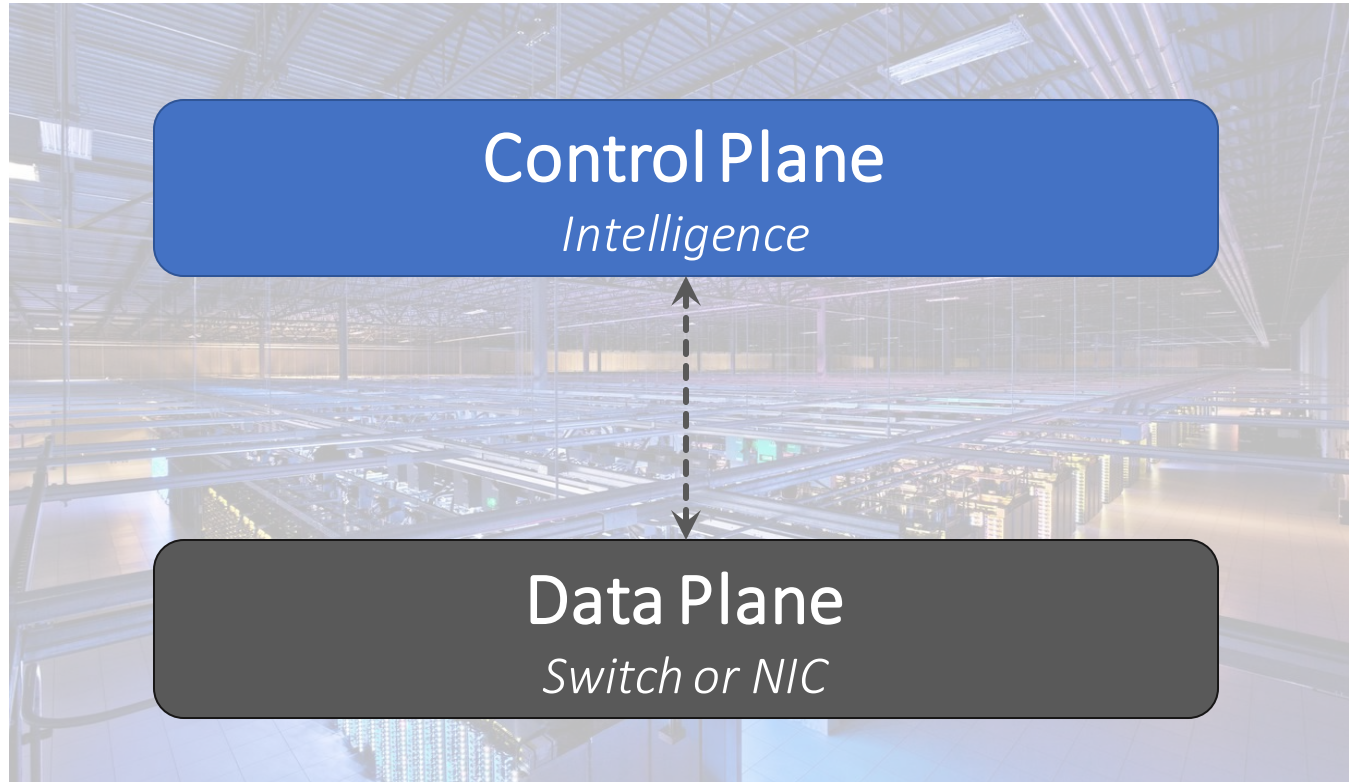


Cloud Computing

Internet of Things (IoT)

Augmented and Virtual Reality (AR/VR)

# Approaches to Manage Networks are …

**Control Plane**
*Intelligence*

**Data Plane**
*Switch or NIC*

*Slow* but *intelligent*
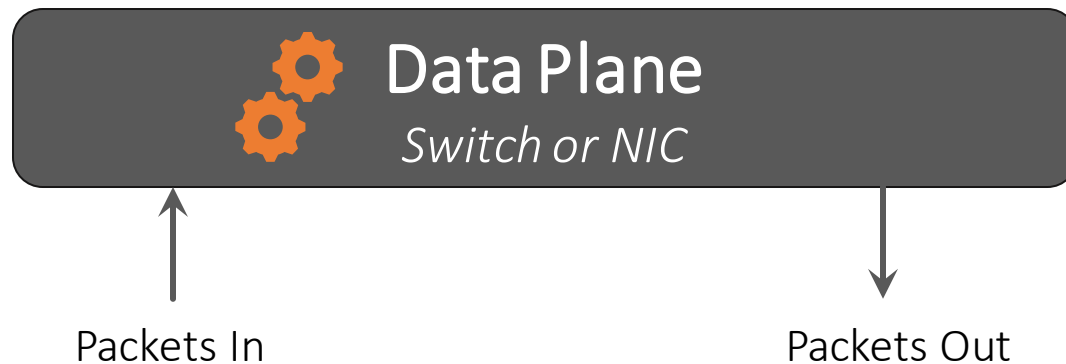
*Fast* yet *dumb*

amazon    Google    Microsoft

# Approaches to Manage Networks are …

Examples:
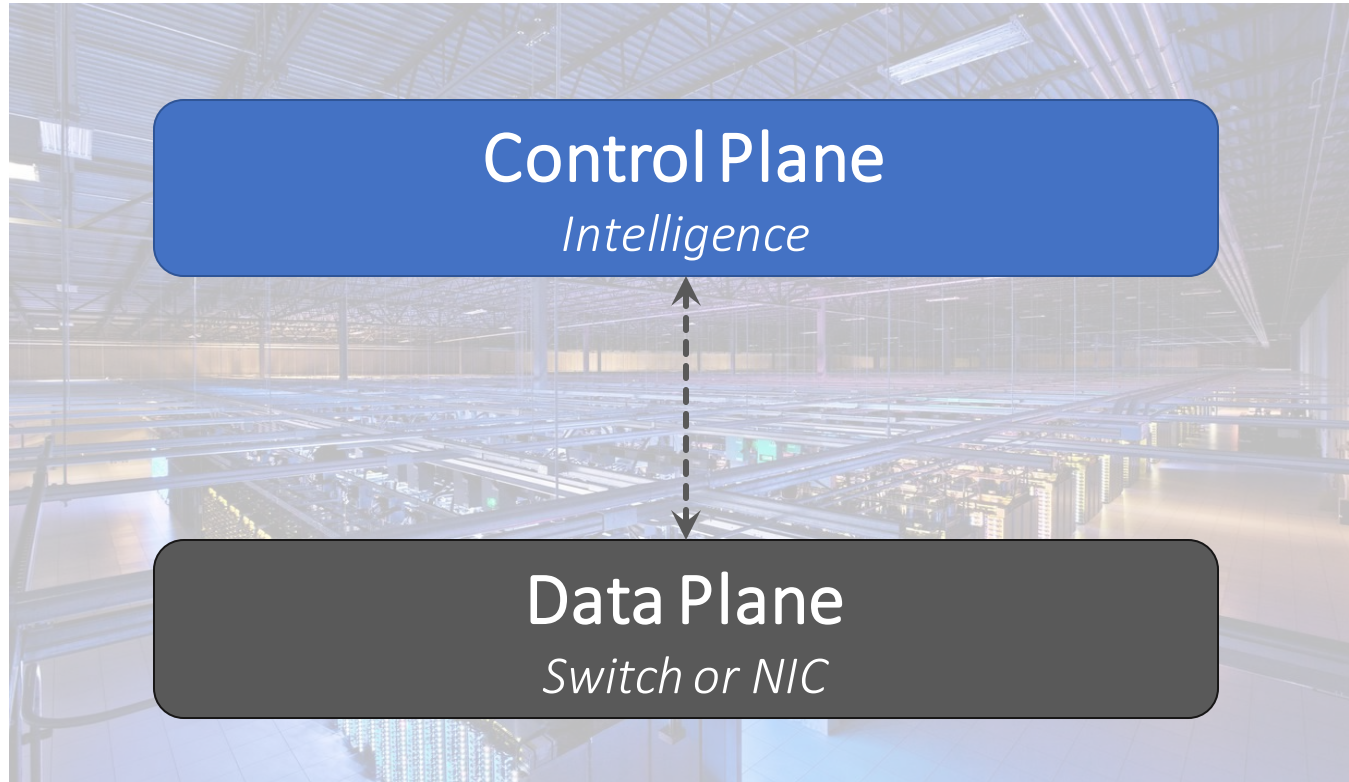- Congestion control
- Load balancing (ECMP, RSS)
- Queue scheduling
- and more

Characteristics:
- Operates on packets or flowlets (*i.e.,* bursts of packets)
- Uses heuristics … hash, etc.
- Low latency … ≤ sub µs
- High throughput … Tbps

**Data Plane**
*Switch or NIC*

Packets In

Packets Out

*Fast* yet *dumb*

# Approaches to Manage Networks are ...

**Control Plane**
*Intelligence*

**Data Plane**
*Switch or NIC*

*Slow* but *intelligent*

*Fast* yet *dumb*

amazon    Google    Microsoft
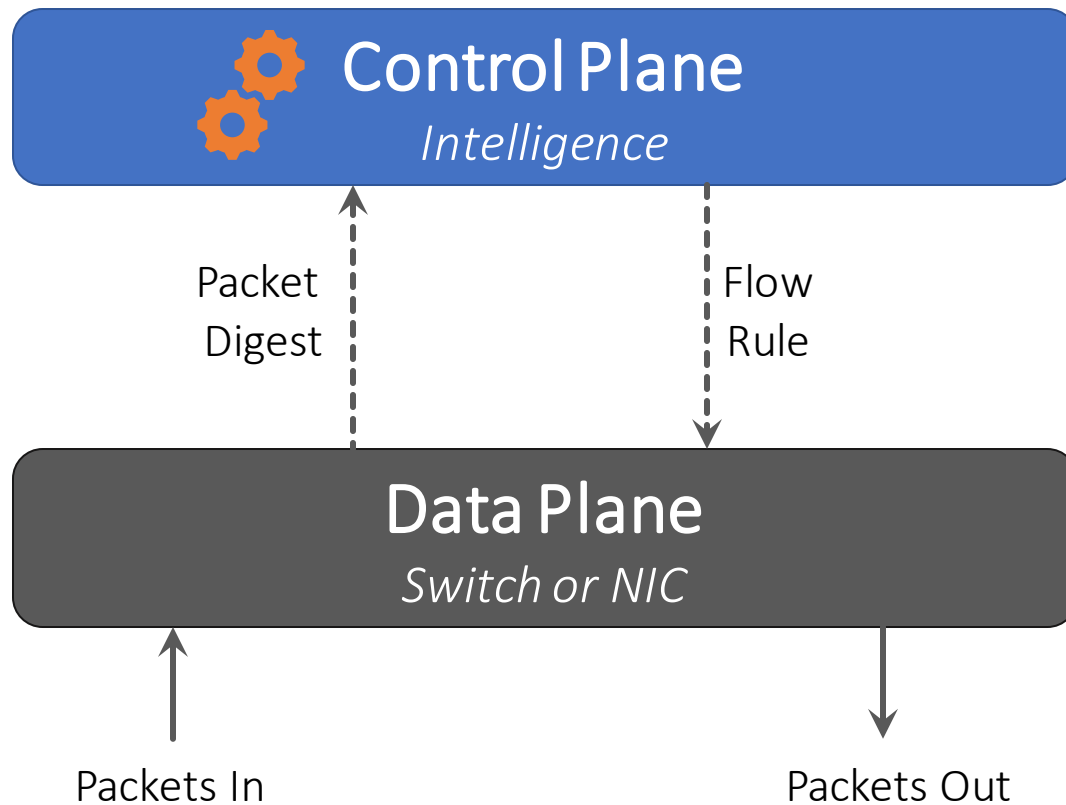
# Approaches to Manage Networks are …



Slow but *intelligent*

Examples:
- Anomaly detection
- Automation
- Recommendation

Characteristics:
- Operates on flows
- Performs complicated tasks
- Sub-second latency
- Low throughput

Control Plane
*Intelligence*

Packet
Digest

Flow
Rule

Data Plane
*Switch or NIC*

Packets In

Packets Out

# Approaches to Manage Networks are ...

**Control Plane**
*Intelligence*

*Slow* but *intelligent*

**Data Plane**
*Switch or NIC*

*Fast* yet *dumb*

# Approaches to Manage Networks are ...

**Control Plane**
*Intelligence*

~~Slow~~ but *intelligent*

**Data Plane**
*Switch or NIC*

*Fast* yet ~~dumb~~

# Approaches to Manage Networks are …

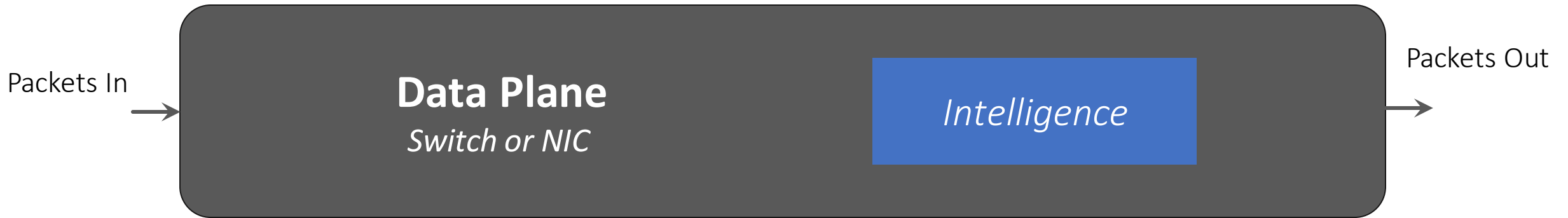**Control Plane**
*Intelligence*

*intelligent*

**Data Plane**
*Switch or NIC*   Intelligence

*Fast* and *intelligent*

# Taurus: An Intelligent Data Plane

# What does "intelligence" mean?

- Networks are becoming autonomous, *Self-Driving Networks*.

- **Machine learning (ML)** will play a key role in the future of networks.
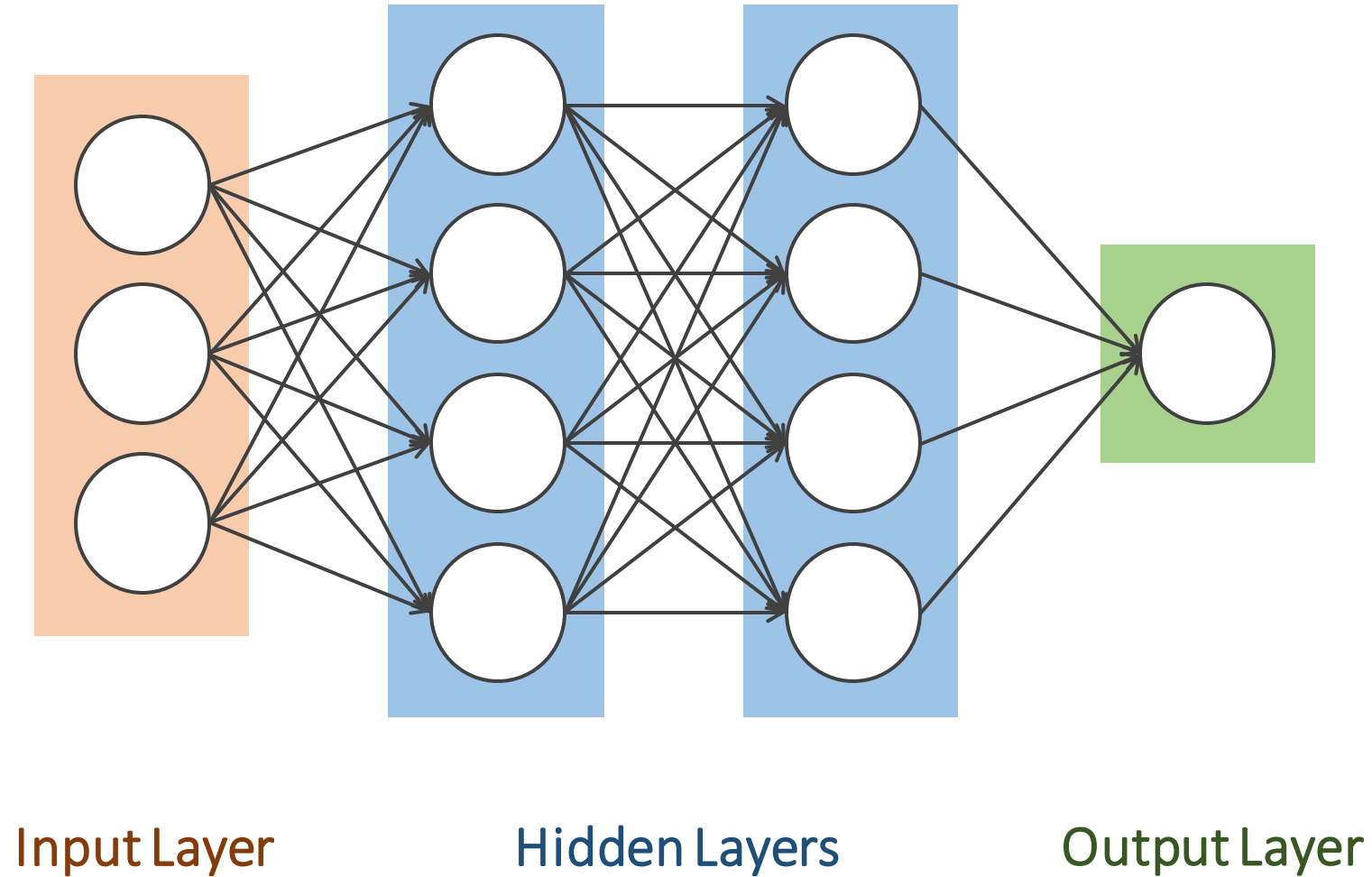
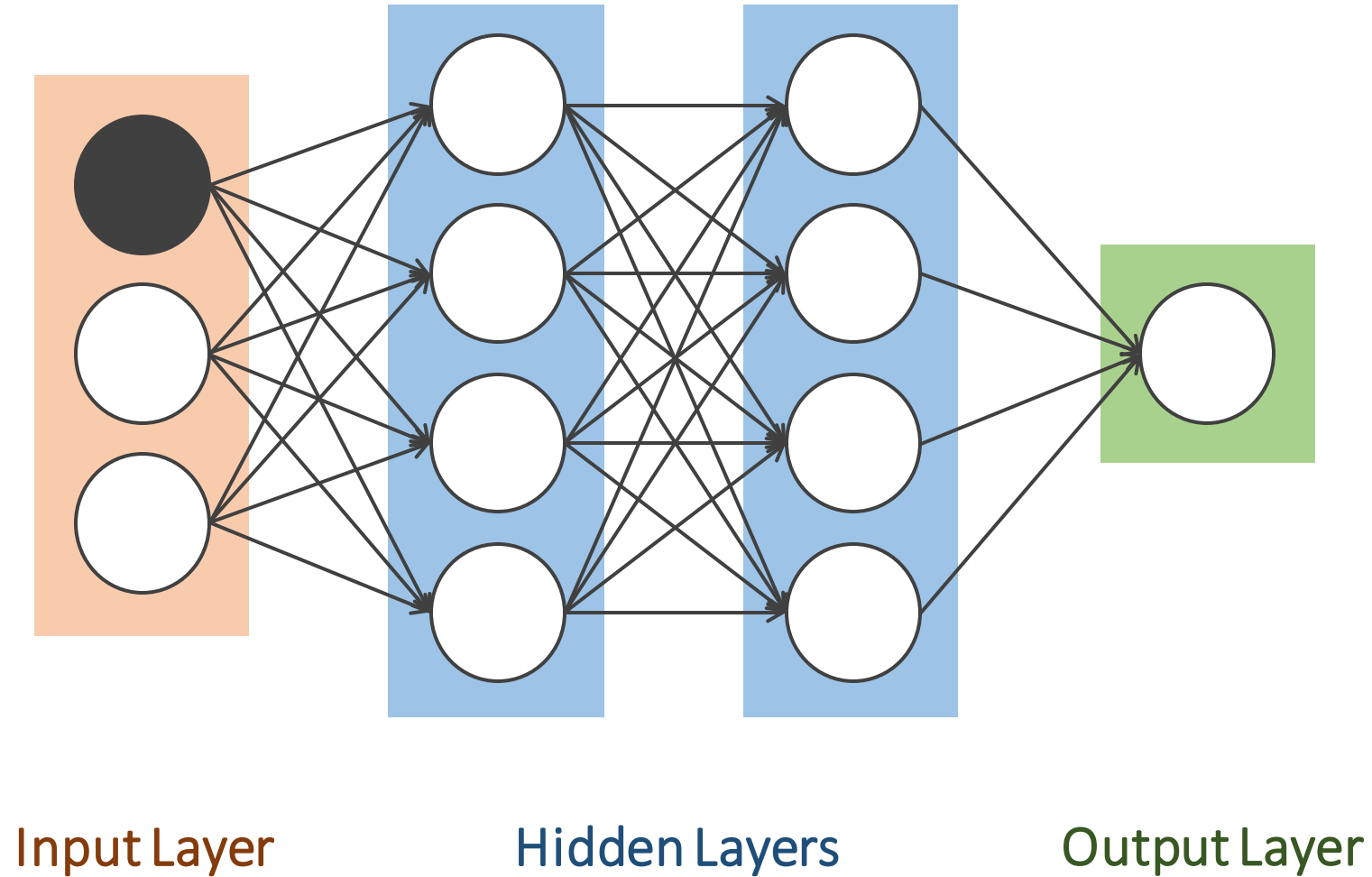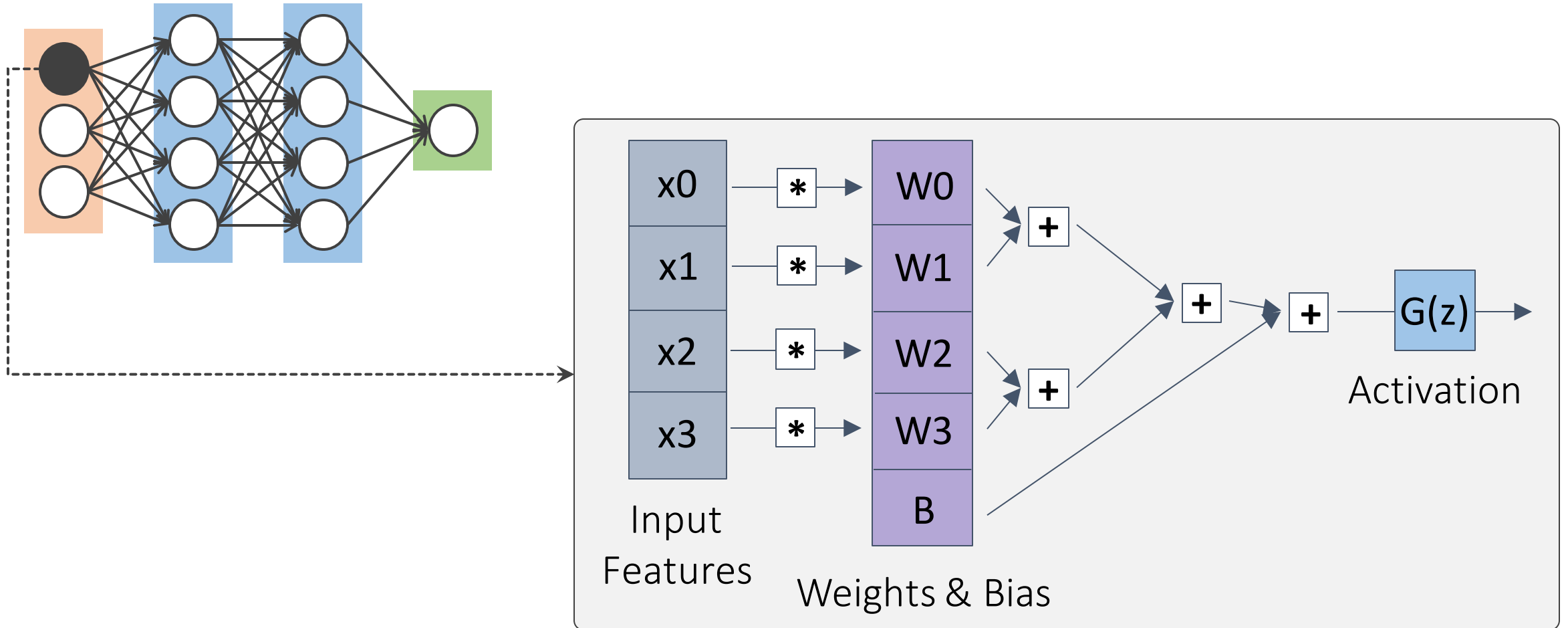**Security**                    **Control**                    **Analytics**

# ML Inference: Neural Networks

Input Layer       Hidden Layers       Output Layer

# ML Inference: Neural Networks



Input Layer       Hidden Layers       Output Layer

# ML Inference: Neural Networks

# Modern Network Data Plane

Packet
Parser

Match-Action
Tables

Traffic
Manager

Packets In
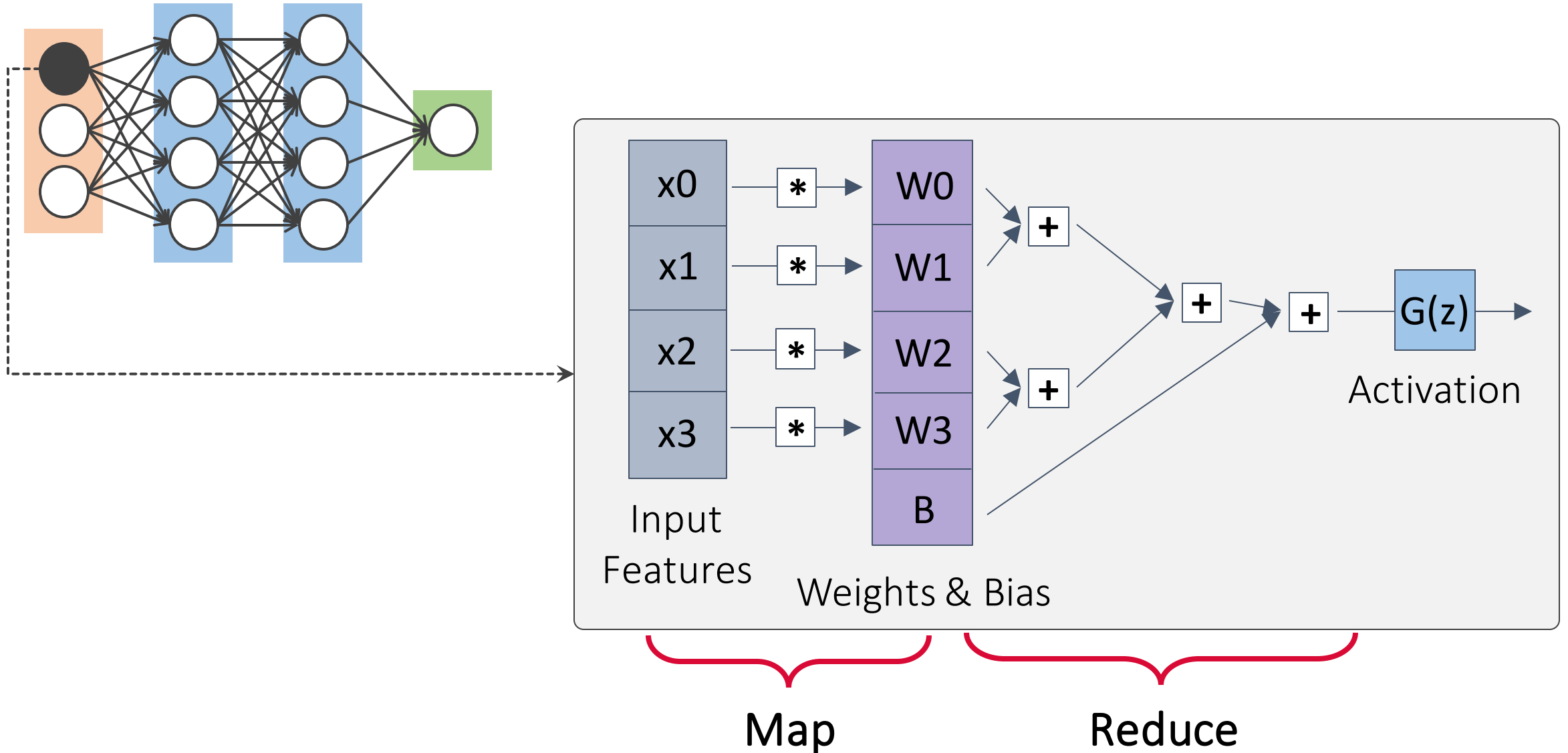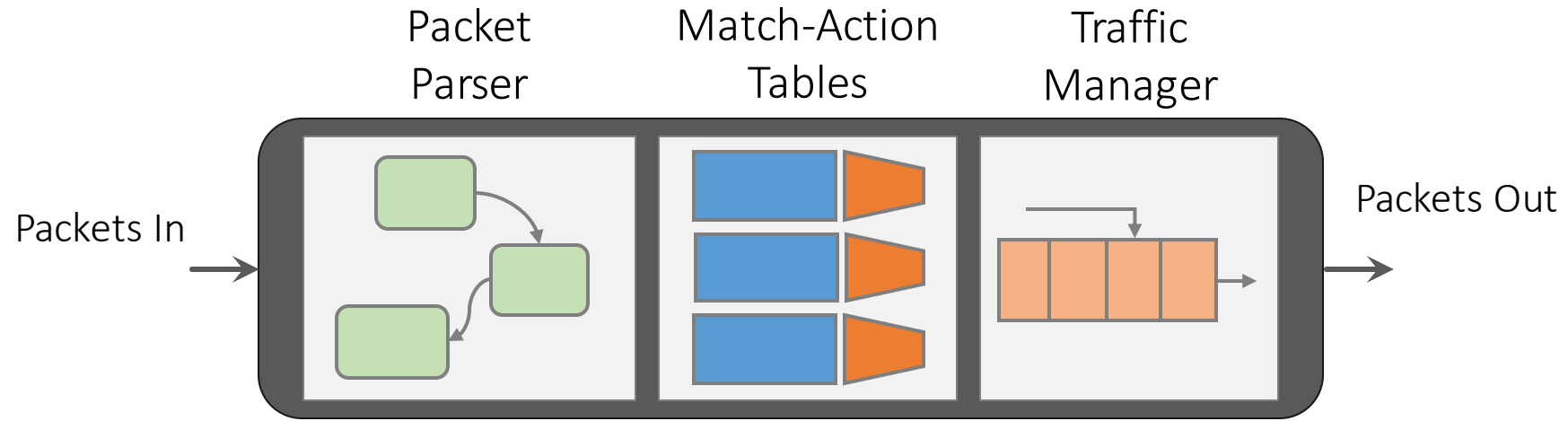
Packets Out

# Modern Network Data Plane



- It's **not suitable** for **learned operations:**
  - Arithmetic intensity is too low to perform ML operations
  - Not enough intermediate storage to carry feature and state
  - and more …

# ML Inference: Neural Networks

# Modern Network Data Plane
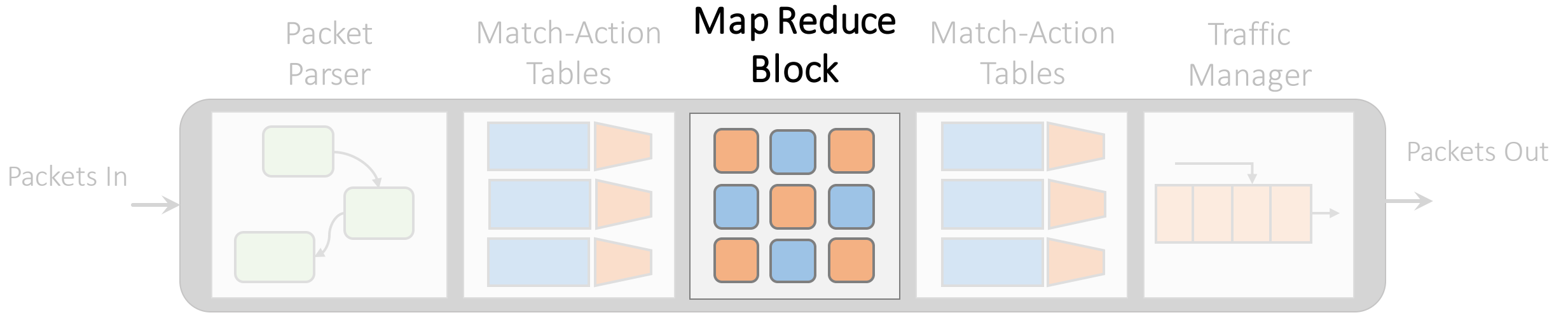
Packet
Parser

Match-Action
Tables

Traffic
Manager

Packets In

Packets Out

# Taurus: An Intelligent Data Plane

Packet Parser  Match-Action Tables  **Map Reduce Block**  Match-Action Tables  Traffic Manager

Packets In →

Packets Out →

# Taurus: An Intelligent Data Plane

Packet Parser — Match-Action Tables — **Map Reduce Block** — Match-Action Tables — Traffic Manager
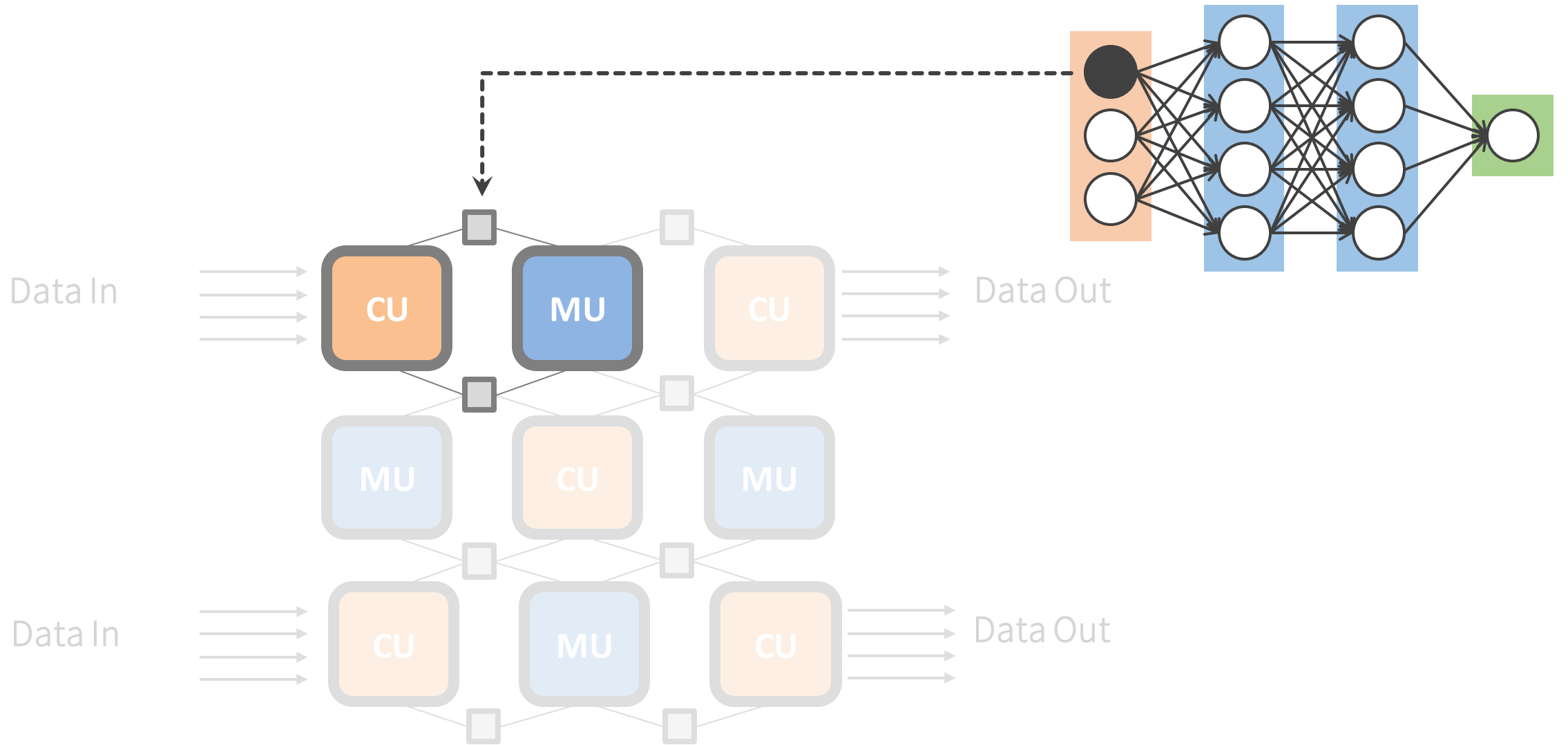
Packets In → → Packets Out

- Implements a **spatial SIMD architecture**

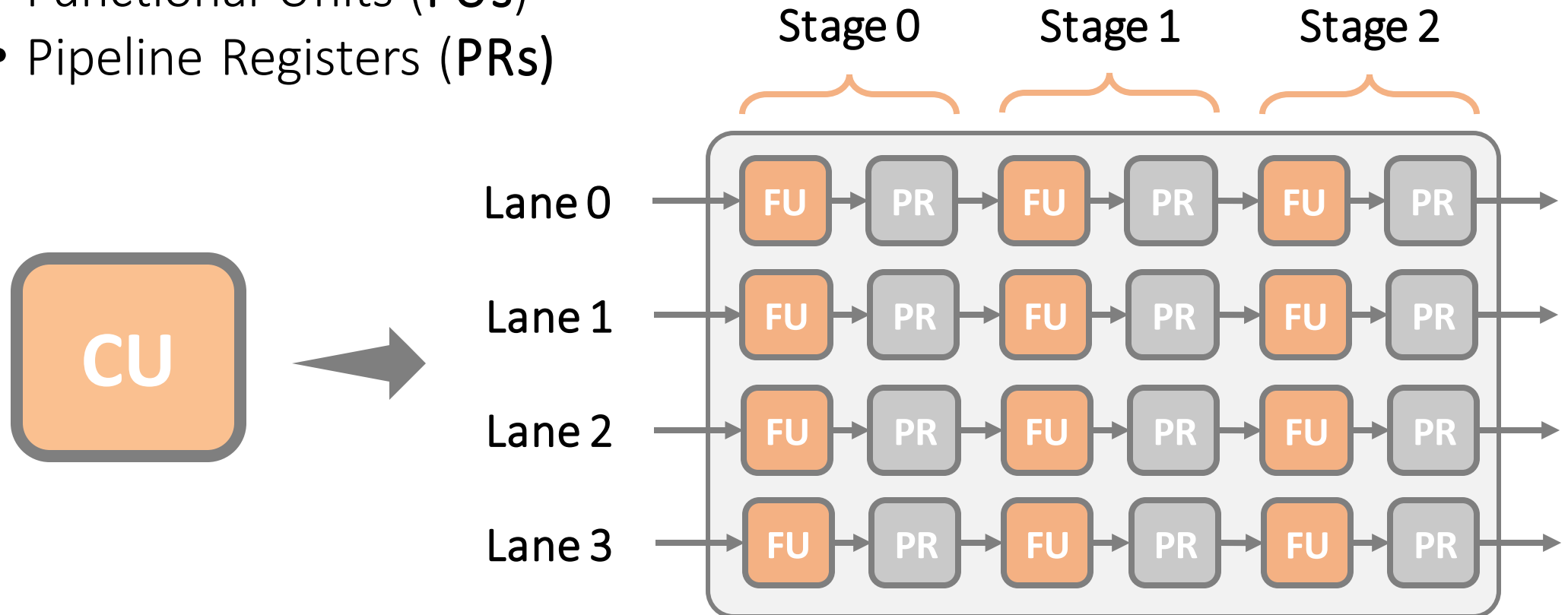# Taurus: Map Reduce Block

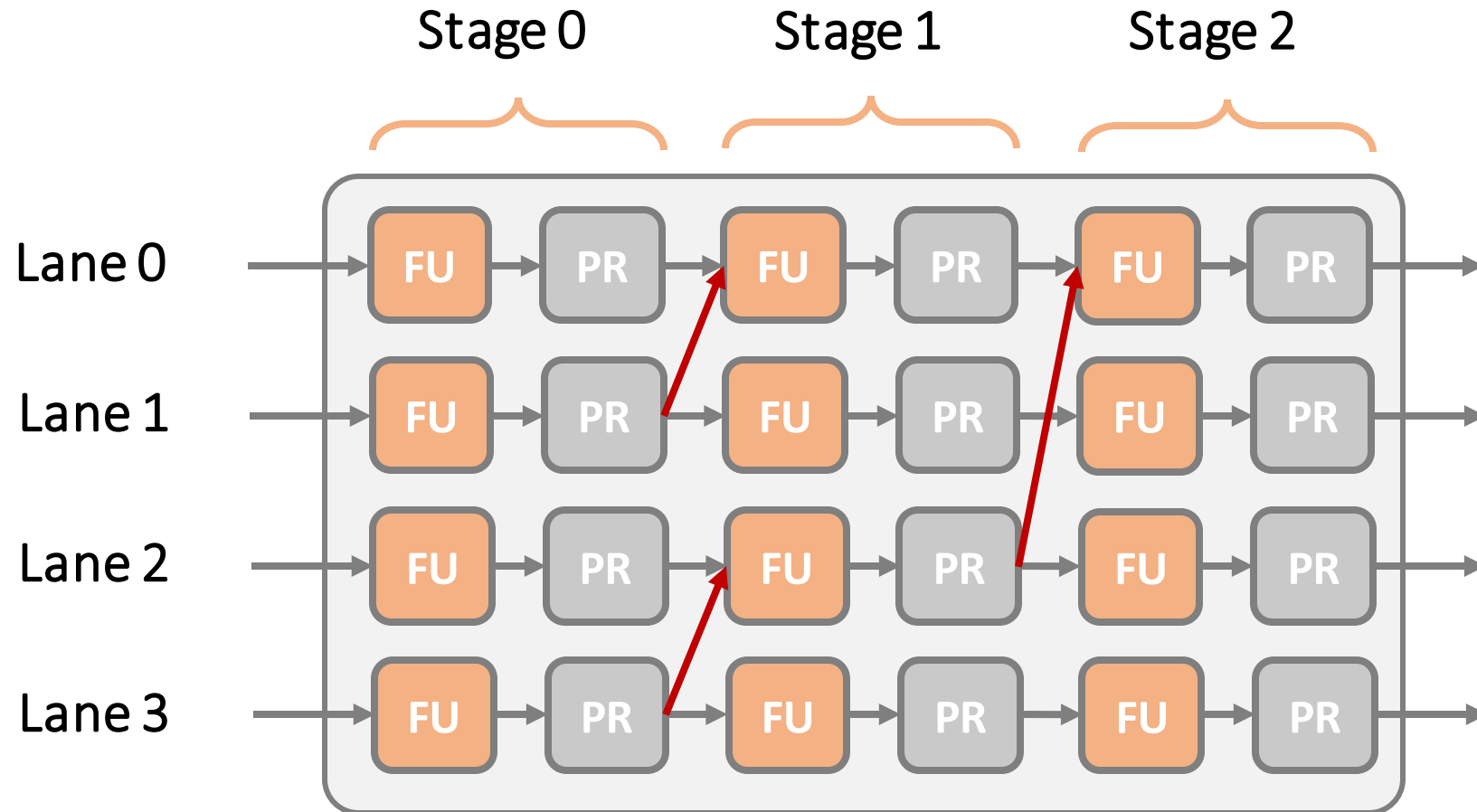# Taurus: Map Reduce Block

# Map Reduce Block: Compute Unit (CU)

- Taurus **CUs** are array-based:
  - Functional Units (**FUs**)
  - Pipeline Registers (**PRs**)
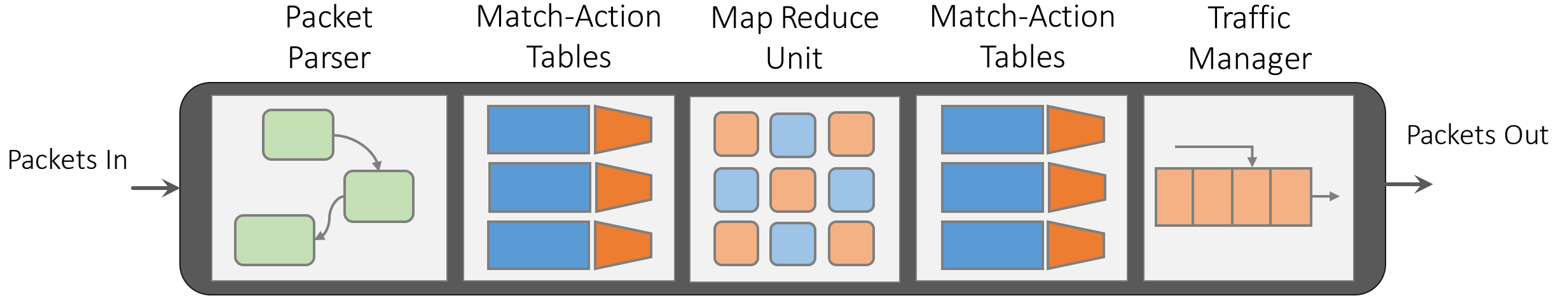
# Map Reduce Block: Compute Unit (CU)



Reduction network condenses vectors to scalars

# Taurus: An Intelligent Data Plane

Packet
Parser

Match-Action
Tables

Map Reduce
Unit

Match-Action
Tables

Traffic
Manager

Packets In

Packets Out

# Example: Anomaly Detection



Packets In → | Packet Parser | Match-Action Tables | Map Reduce Unit | Match-Action Tables | Traffic Manager | → Packets Out

Parse packets and **read local features** (e.g., IP address)

Retrieve **out of network events** (e.g., failed logins per IP)

Apply **learned functions** to mark anomalous packets

**Select a port** or **action** (drop if score == 1)

**Send packets** out the selected port

# Evaluation: Anomaly Detection in Switches

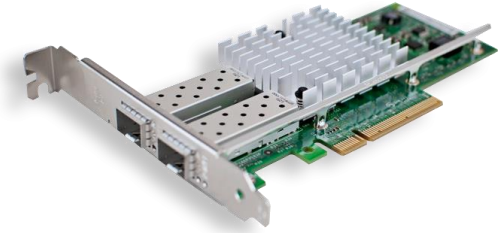- Taurus examines **every** packet at **line rate**

- Added **latency** is **less** than **port-to-port** latency

| Model | Throughput | Latency | Area +% | Power +% |
|:---:|:---:|:---:|:---:|:---:|
| SVM | 1 GPkt/s | **68** ns | 6.1 | 1.1 |
| DNN | 1 GPkt/s | **362** ns | 11.7 | 2.0 |

*\*Overheads are calculated relative to a 300 mm$^2$ chip with 4 reconfigurable pipelines, each drawing an estimated 25 W*

# Evaluation: Congestion Control at the NICs



- **Indigo** is a **congestion control** LSTM network

- Taurus updates **every 12.5 ns** (software updates every 10 ms)

| Model | Throughput | Latency | Area +% | Power +% |
|---|---|---|---|---|
| LSTM | 0.08 GPkt/s | 380 ns | 23.6 | 4.1 |

*Overheads are calculated relative to a 300 mm$^2$ chip with 4 reconfigurable pipelines, each drawing an estimated 25 W

# Conclusion

**Taurus**

**Data Plane**
*Switch or NIC*
*Intelligence*

*Fast* and *intelligent*

- Designed to **run machine-learning inference** inside a data plane
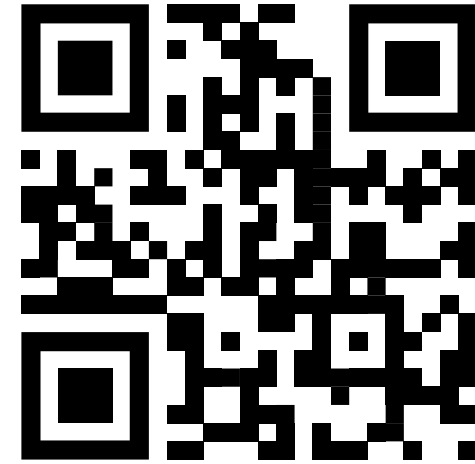- Provides **orders of magnitude improvement** over existing approaches

# Conclusion

Taurus

Data Plane
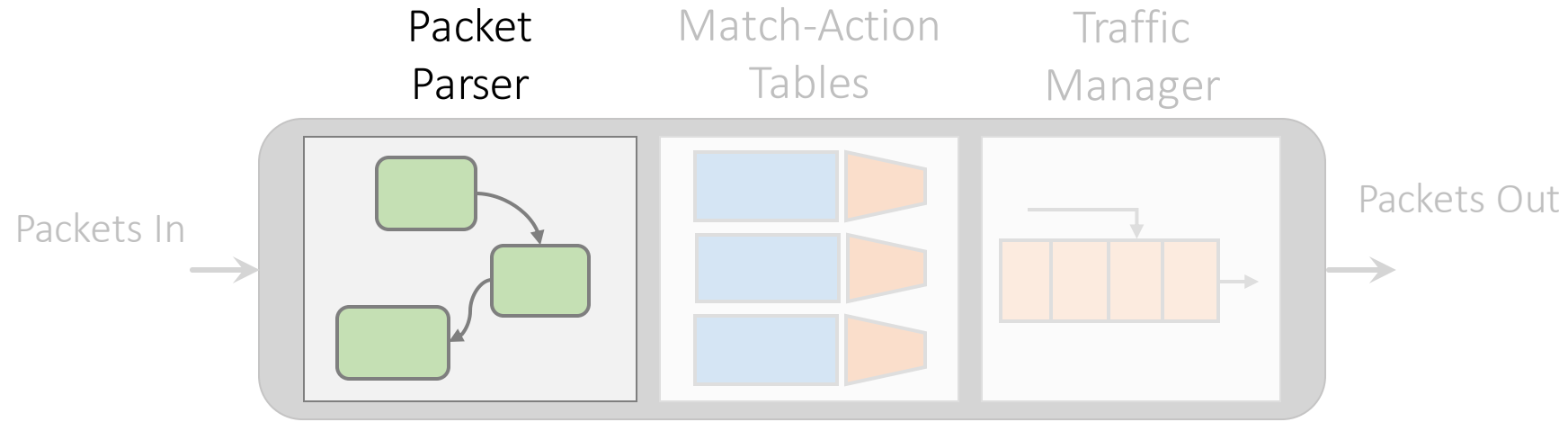*Switch or NIC*          *Intelligence*

*Fast* and *intelligent*    **?**

dataplane.ai

Muhammad **Shahbaz**

http://cs.stanford.edu/~mshahbaz

# Backup slides ...

# Modern Network Data Plane



Packet Parser

Match-Action Tables

Traffic Manager

Packets In

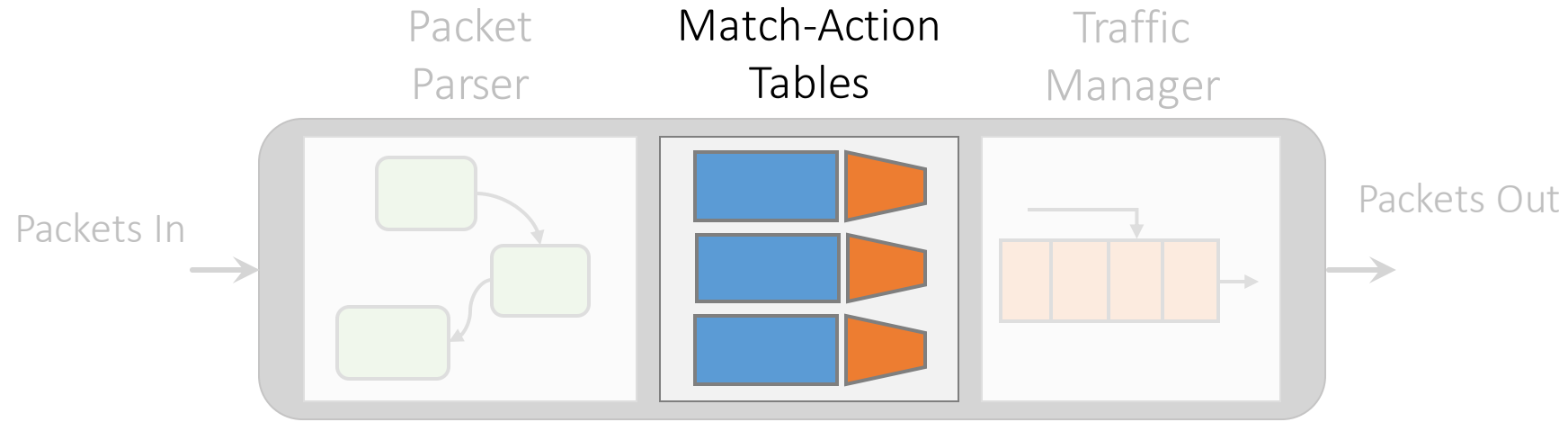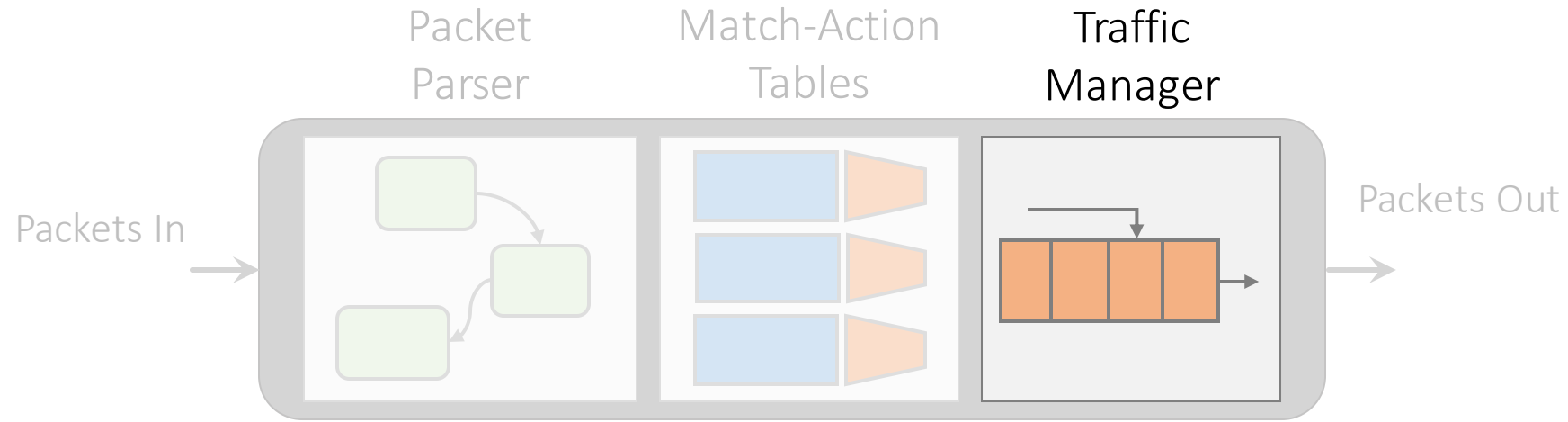Packets Out

- Implements a **finite state machine (FSM)** that operates on a user-defined **parse graph**

- Converts the **incoming packet bit stream into vectors**, *e.g.,*
  - headers (IP or TCP)

# Modern Network Data Plane

Packet
Parser

Match-Action
Tables

Traffic
Manager

Packets In

Packets Out

Memory | ALU

- A **match-action table:**
  - Memory for **exact** (SRAM) and **ternary** (TCAM) match
  - ALU for basic single-cycle **VLIW operations** (no loops or multiplication)

# Modern Network Data Plane



Packet Parser

Match-Action Tables

Traffic Manager

Packets In

Packets Out

- Responsible for **storing** and **forwarding** packets off of the chip:
  - **Queuing**: buffer incoming packet
  - **Replication**: clone packets across multiple egress ports (*e.g.*, multicast)
  - **Scheduling**: forward packets based on a queuing discipline (*e.g.*, PIFO) or instructions from the match-action tables