Daniele De Sensi, Salvatore Di Girolamo, Saleh Ashkboos, Shigang Li, Torsten Hoefler
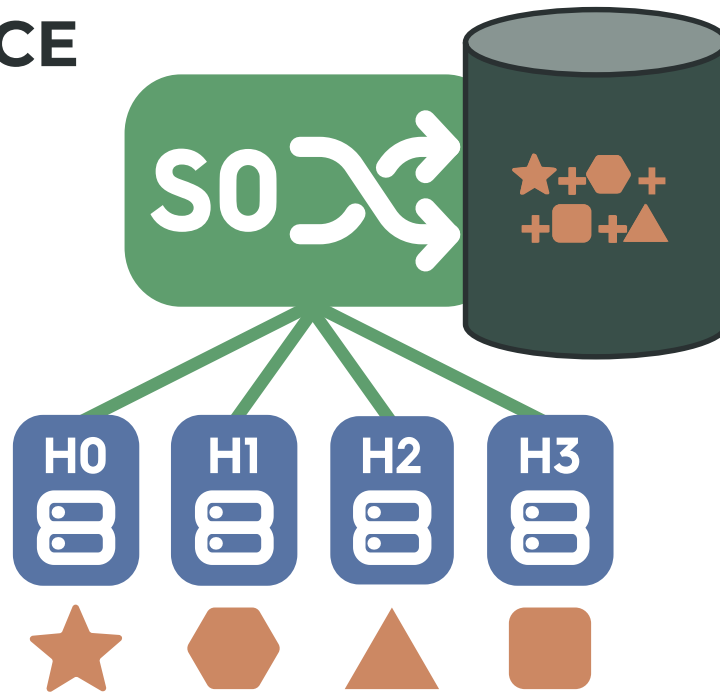
# Flare: Flexible In-Network Allreduce

# IN-NETWORK ALLREDUCE
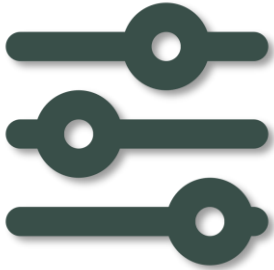
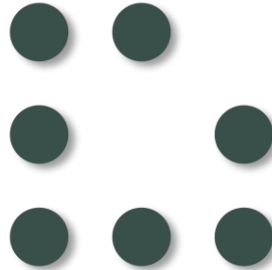**2x traffic reduction** compared to host-based allreduce

**2x bandwidth improvement**

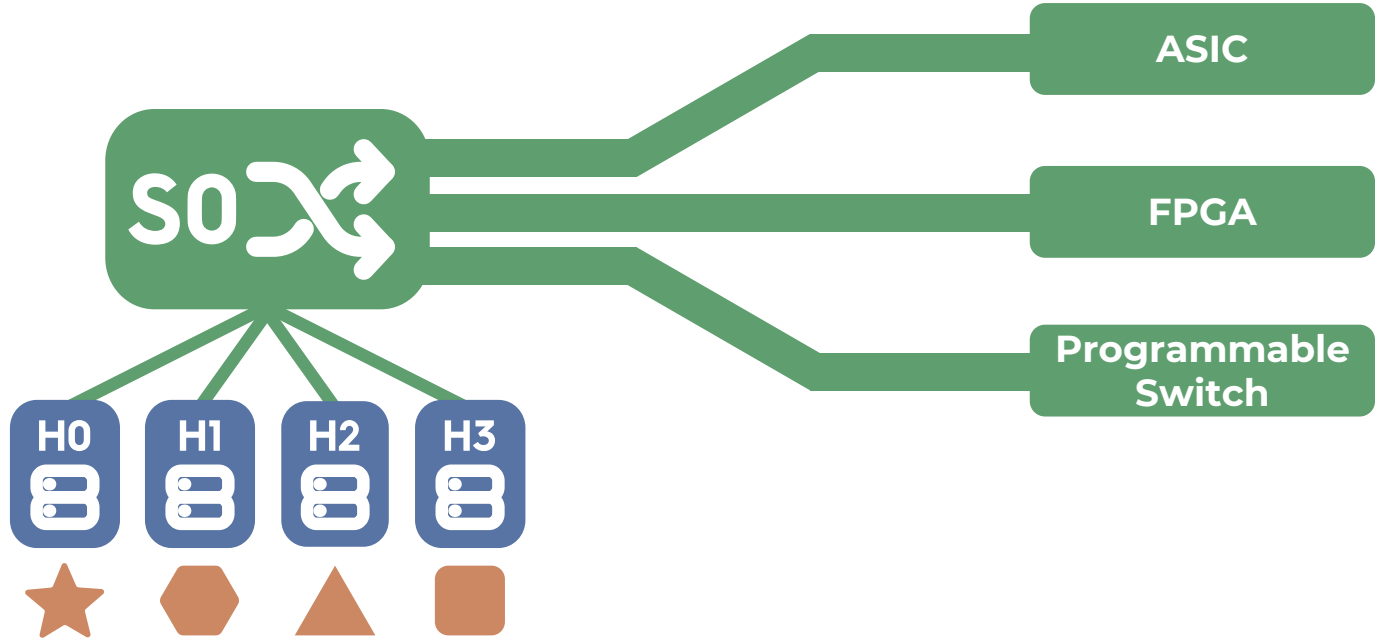# MISSING FEATURES

Custom
**operators** and
**datatypes**

Support for
**sparse data**

**Reproducibility**

# EXISTING SWITCHES ARCHITECTURES

# FLARE

Programmable switch **architecture**
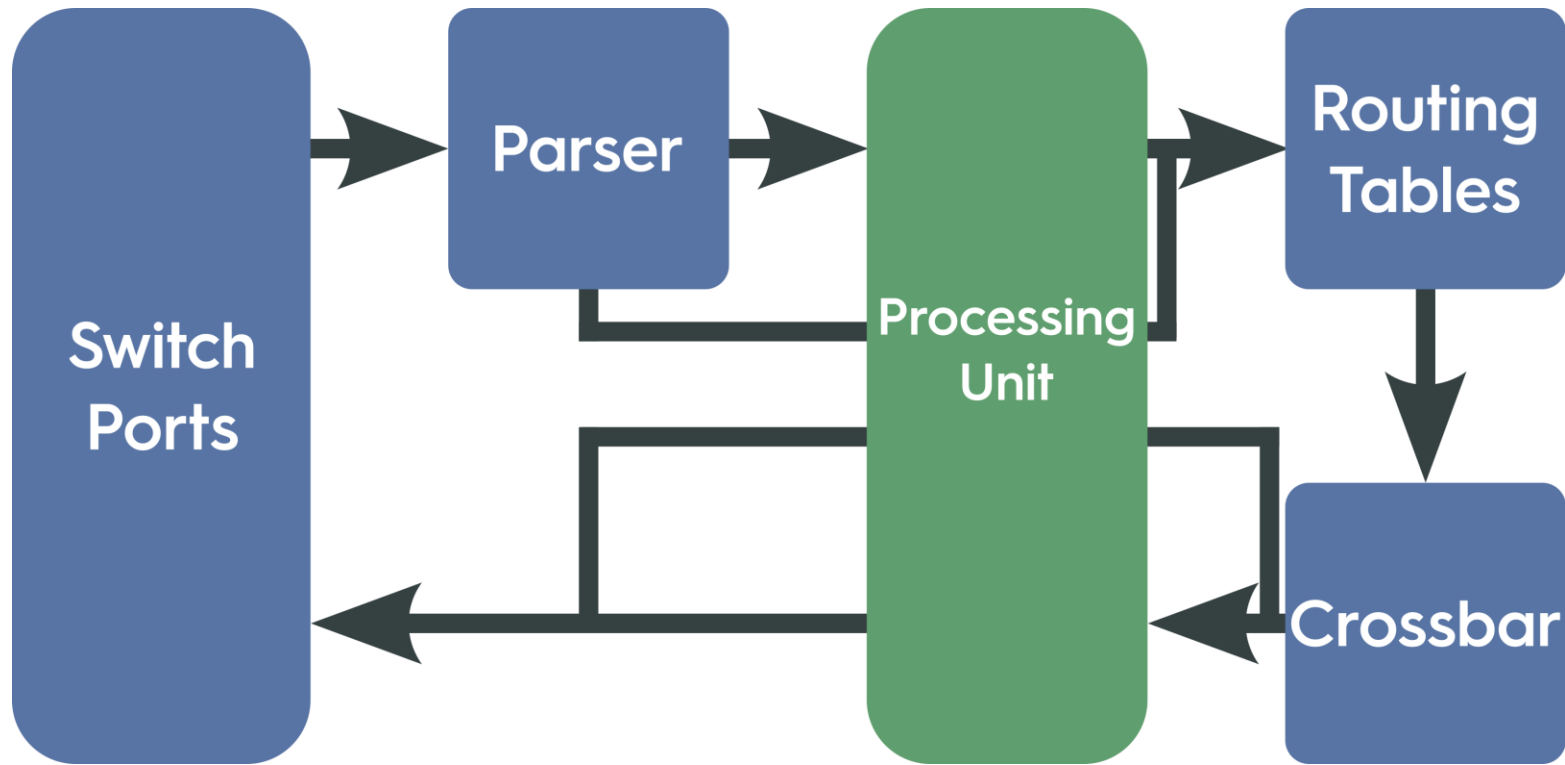
Set of **algorithms**
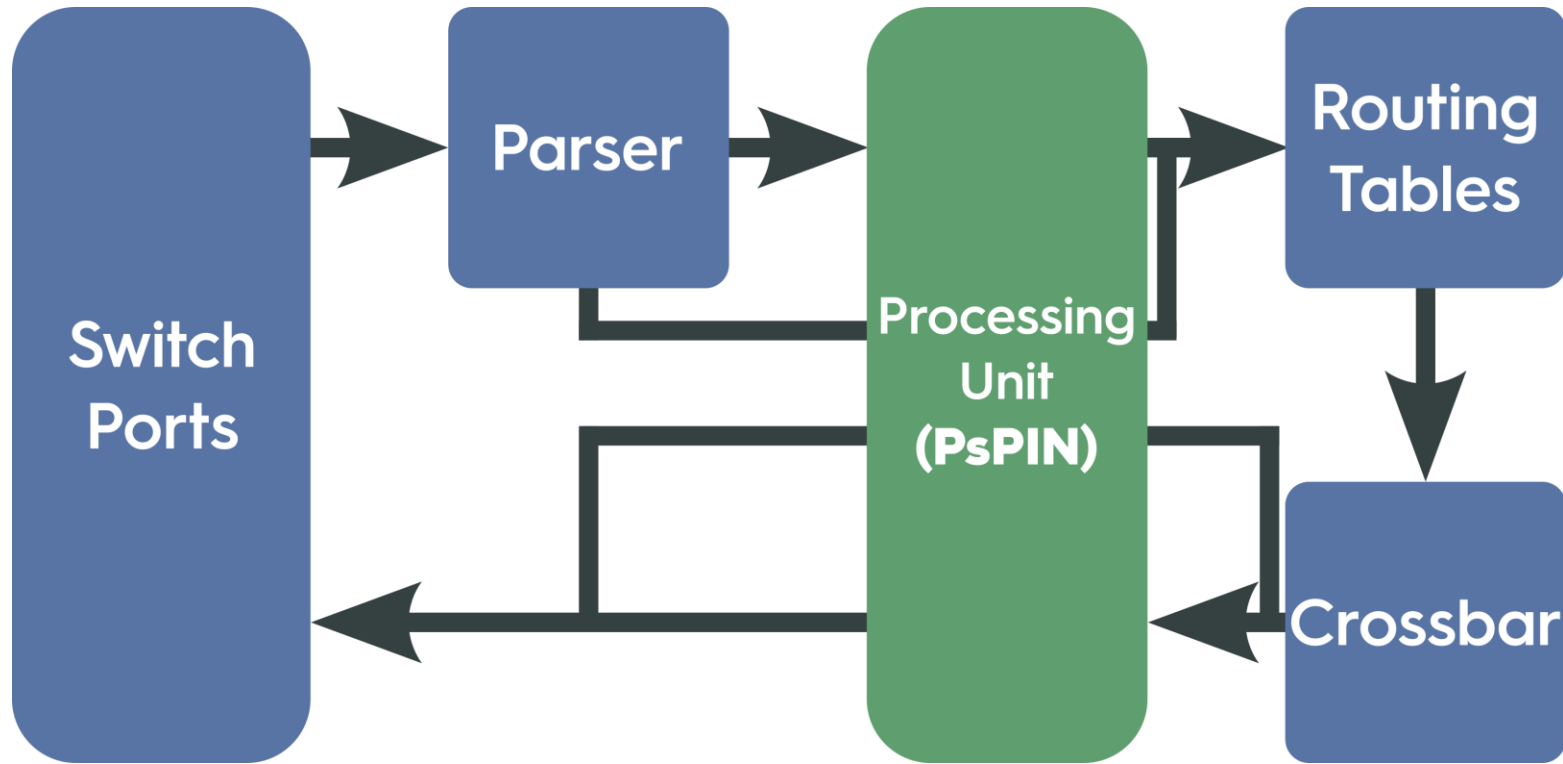
Performance and memory occupancy **models**
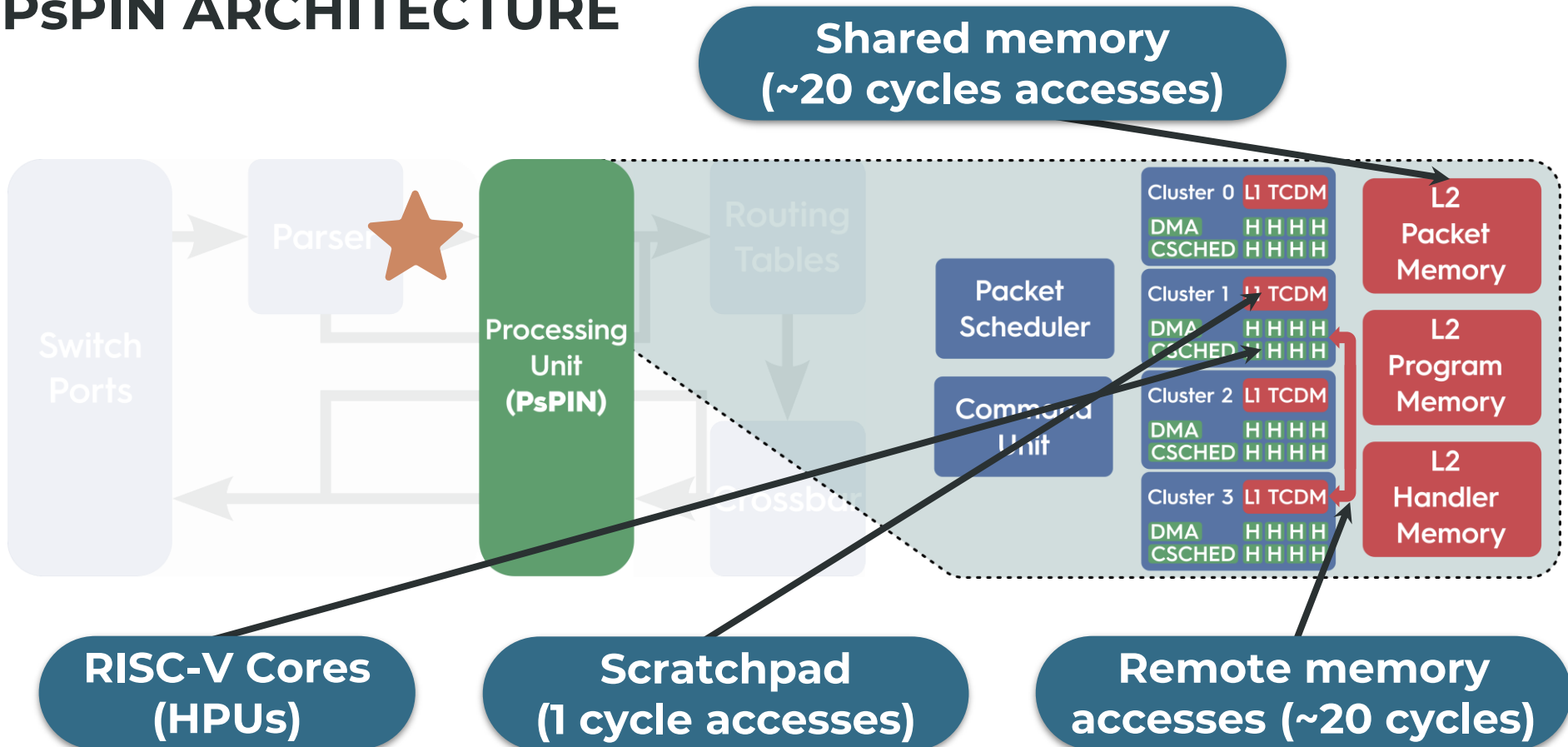
Programmable switch **architecture**

# SWITCH ARCHITECTURE

# SWITCH ARCHITECTURE



*PsPIN: A high-performance low-power architecture for flexible in-network compute* - **S. Di Girolamo et al. - ISCA'21**
*sPIN: High-performance streaming Processing in the Network* – **T. Hoefler et al. – SC'17**

# PsPIN ARCHITECTURE

*PsPIN: A high-performance low-power architecture for flexible in-network compute* - **S. Di Girolamo et al. - ISCA'21**

# ADVANTAGES

Processing functions specified as **C kernels**

Coverage of **more use cases**

We can fit **512 cores + memory**

# Algorithms

# KEY FEATURES

**Where to store**
the data

**How to access**
the data

How to manage
**sparse data**

# KEY FEATURES



**Where to store**
the data

**How to access**
the data

How to manage
**sparse data**

# DATA TRANSMISSION

Host 0 | Packet 0,0 | Packet 0,1 | Packet 0,2 | Packet 0,3 | Packet 0,4

# PACKET SCHEDULING

# KEY FEATURES

**Where to store**
the data

**How to access**
the data

How to manage
**sparse data**

# SHARED BUFFER CONTENTION

# SHARED BUFFER CONTENTION

# SHARED BUFFER CONTENTION



**Staggered sending**

20

# HOW TO REDUCE CONTENTION

**Staggered sending**

| **Single buffer** | **Multiple buffers** | **Tree** |
|---|---|---|

👍 **Minimal memory occupancy**

👎 **High contention**

👉 **Higher memory occupancy**

👉 **Lower contention**

👎 **Highest memory occupancy**

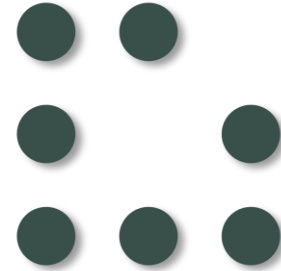👍 **No contention**

👍 **Reproducible**

22

# KEY FEATURES

**Where to store**
the data

**How to access**
the data

How to manage
**sparse data**

# DATA FILL-IN



26

# TWO-LEVEL APPROACH



ARRAY
(OR DENSE)

HASH TABLE
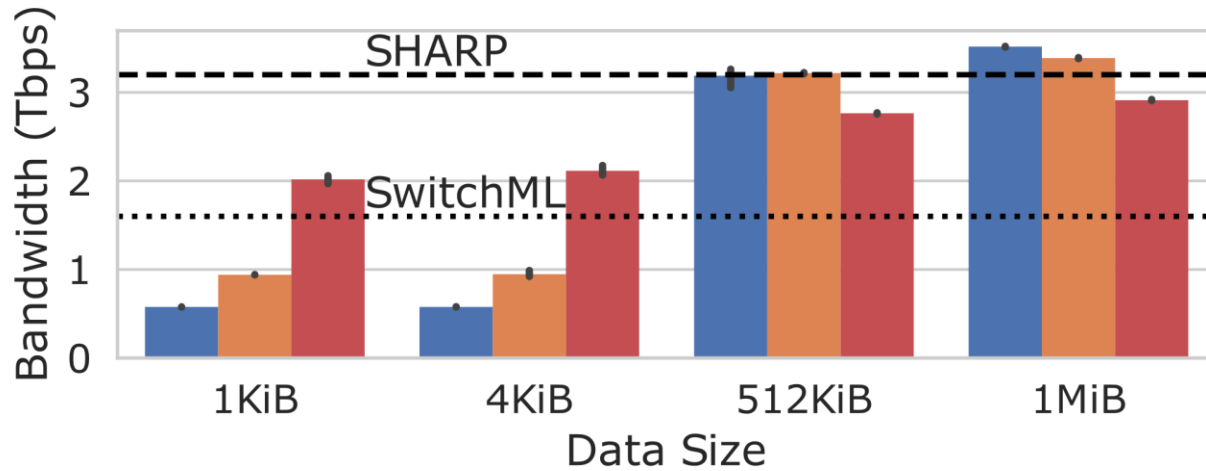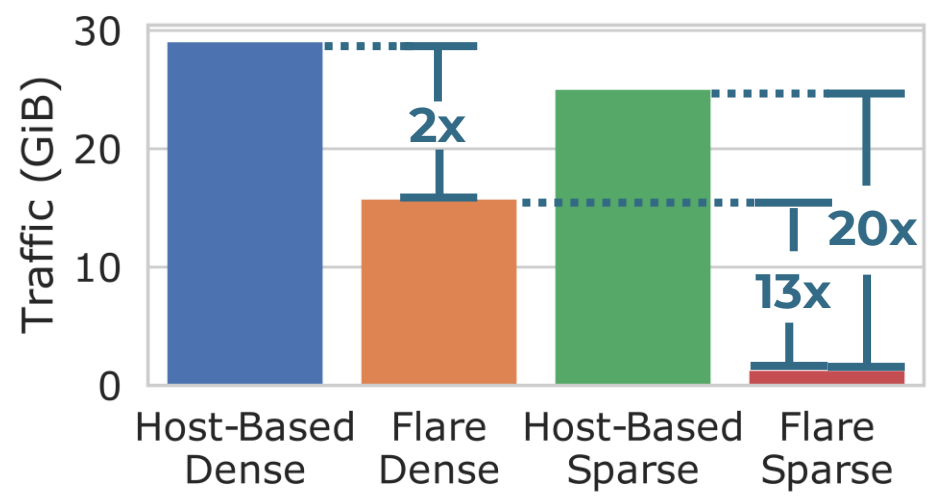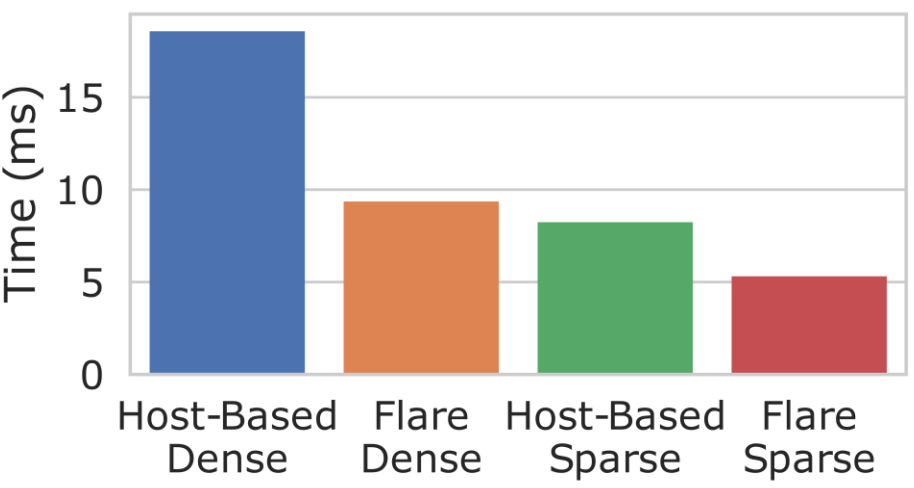
# Evaluation

# RESULTS – SINGLE SWITCH

# RESULTS - 64 NODES, 2-LEVELS FAT TREE



**Communication time** of a **ResNet50** iteration
with **sparsified gradients** (0.2% density)

# CONCLUSIONS



STATE OF THE ART LIMITATIONS

No custom **operators** and **datatypes**

No support for **sparse data**

No **reproducibility** guarantees



SWITCH ARCHITECTURE

PsPIN: A high-performance low-power architecture for flexible in-network compute - S. Di Girolamo et al. - ISCA'21



CONTENTION – STAGGERED SENDING



RESULTS - 64 NODES, 2-LEVELS FAT TREE

**Communication time** of a **ResNet50** iteration with **sparsified gradients**