

Enabling In-Network Computation in Remote Procedure Calls

Bohan Zhao*, Wenfei Wu**, Wei Xu*

**Tsinghua University, **Peking University*

NetRPC: a General INC-enabled RPC System

- **Motivation:**

In-network computation (INC) is **beneficial** to system performance but **difficult to program**

- **Goal:**

Make INC **easy** to use for normal applications with **little performance loss**

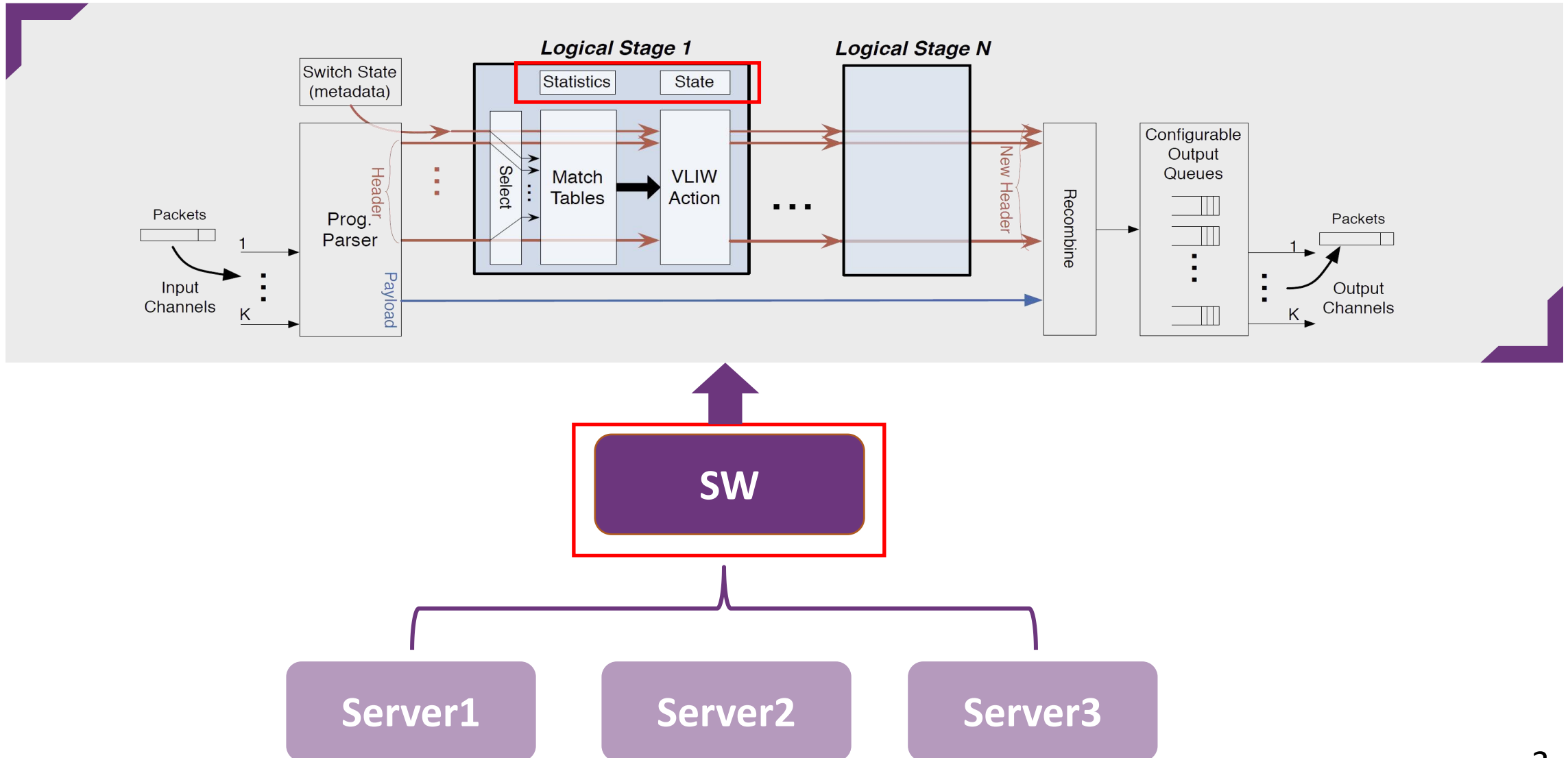
- **Metrics:**

Reduce lines of code of INC applications by up to **97%**

Support most popular INC applications

Little performance loss

INC Customizes Stateful Packet Processing



INC is Widely Used in Many Scenarios



**In-Network
Computation**



Advantages

- Server Func Offloading
- Line-rate Computation
- Network Stack Simplification

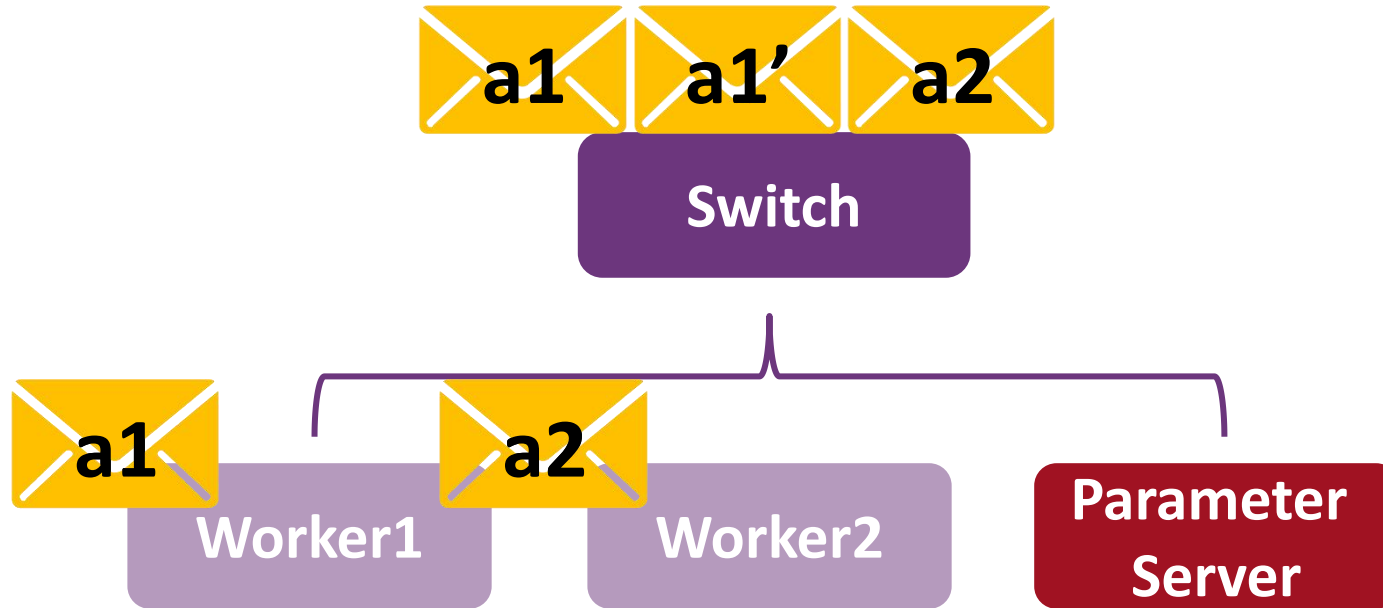


Scenario

- Synchronous Aggregation
- Asynchronous Aggregation
- Key-value Caches
- Agreement

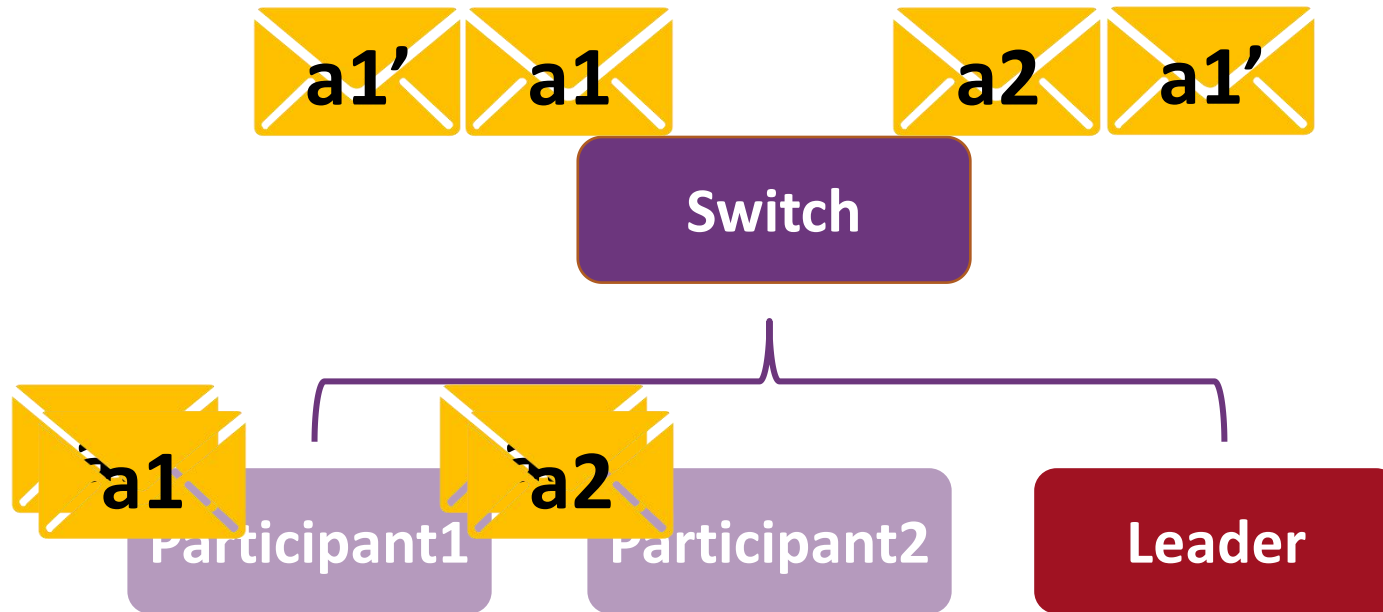
INC Provides Higher Throughput

- Eliminate incast to reduce traffic, especially for distributed training

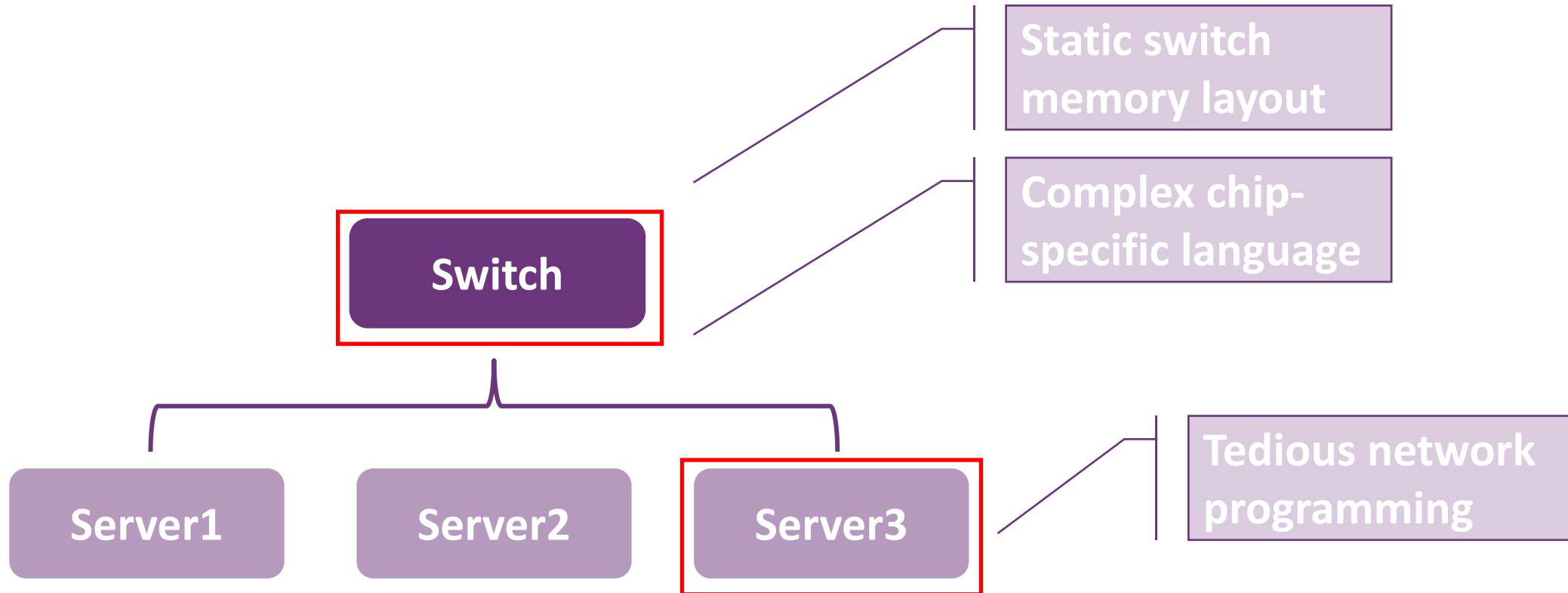


INC Provides Lower Delay

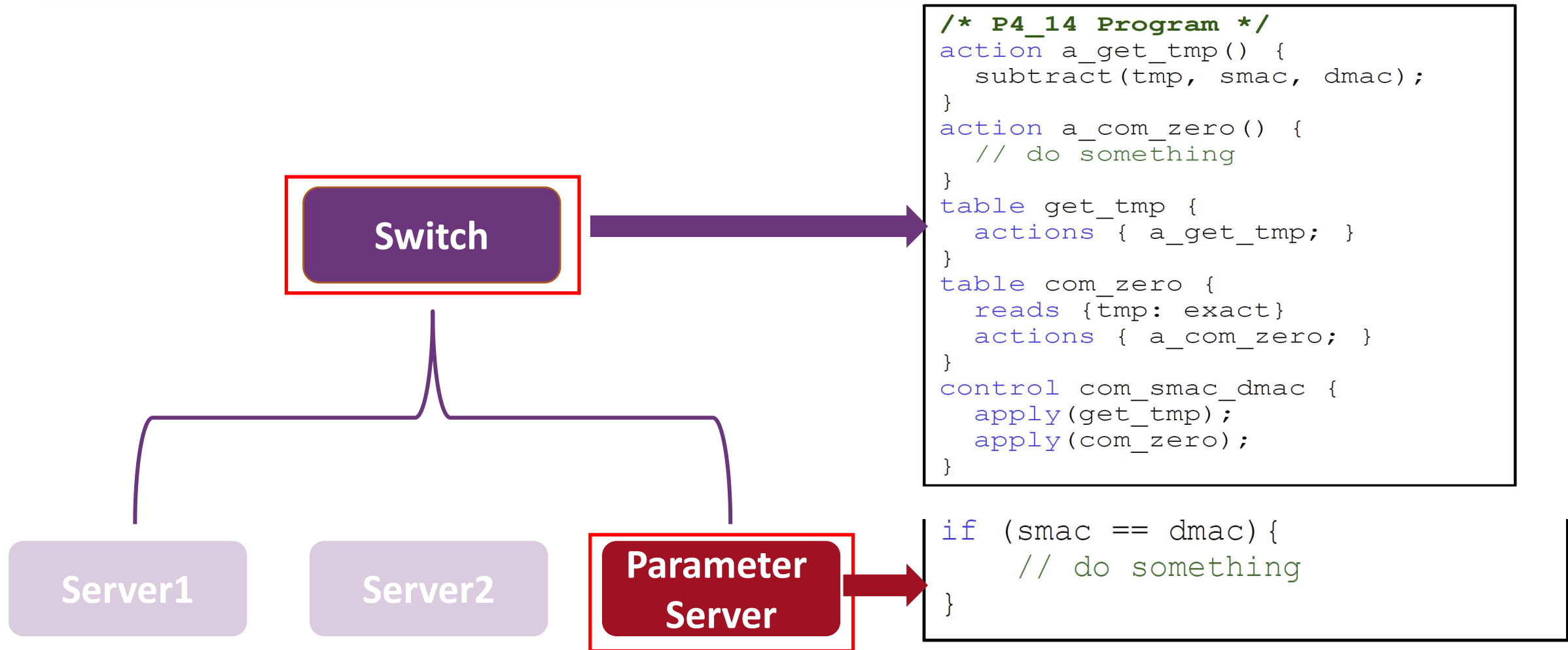
- Reduce the hops of round trip, useful for agreement applications



Challenges of Developing INC Application

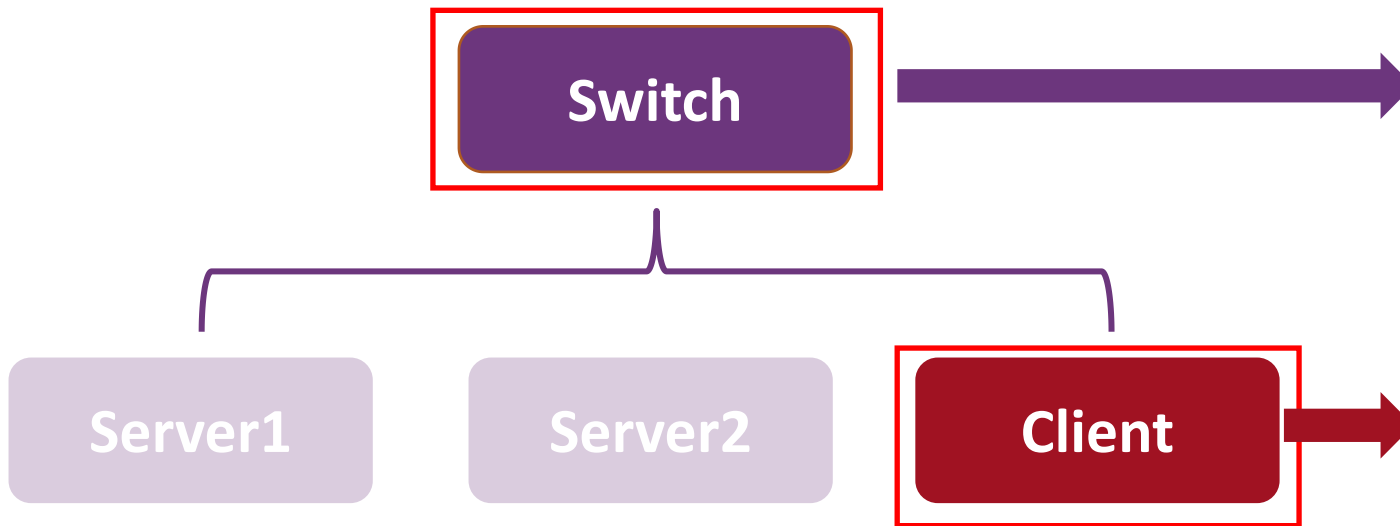


P4 Programming is Much More Complex than Software



Can We Provide a Computation-centric Programming Model to Include INC

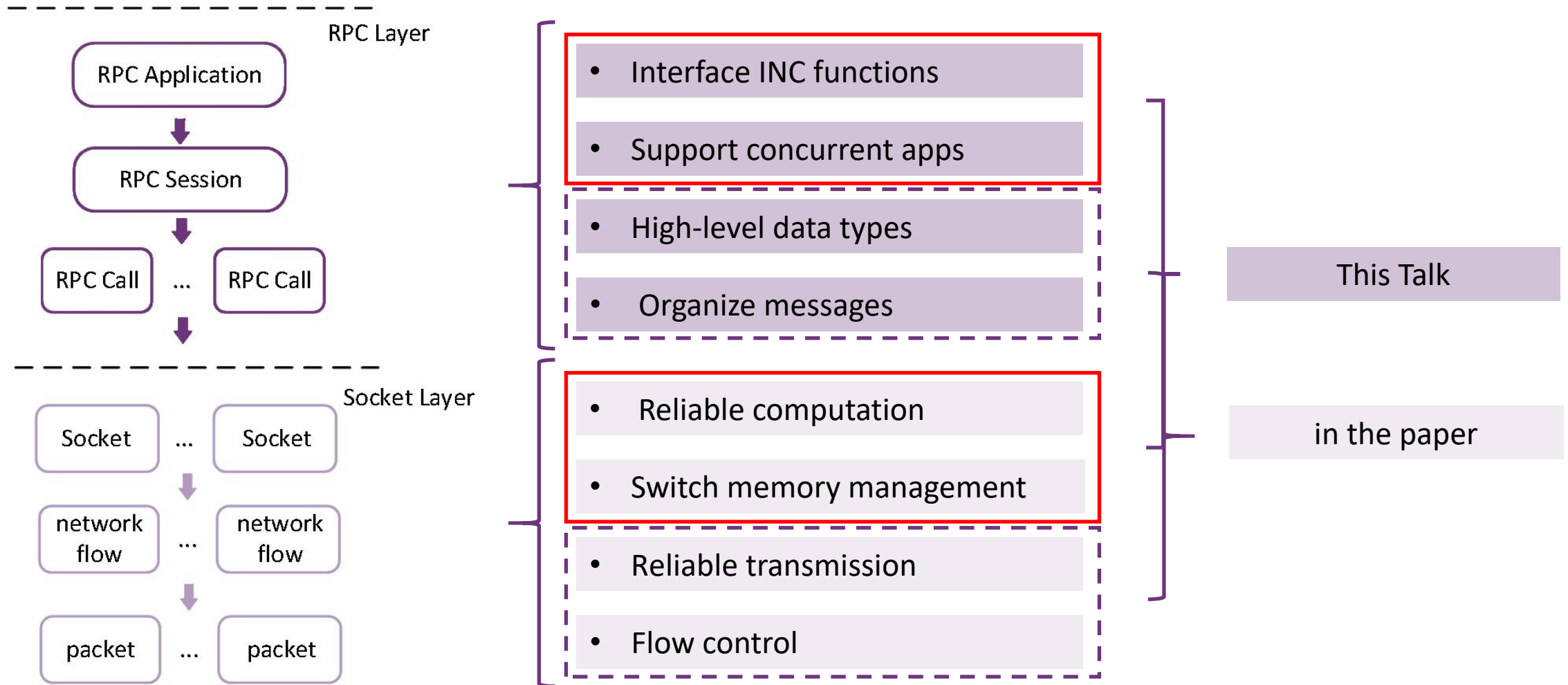
- P4 language is network-centric and focus on **communication**
- Users only take care of **computation**
- **RPC** adapts INC functions better than other models (e.g., MPI)



```
/* P4_14 Program */
action a_get_tmp() {
    subtract(tmp, smac, dmac);
}

action a_com_zero() {
void PushPull(double* data, int length) {
    NewGrad request;
    AgtrGrad reply;
    ClientContext context;
    request.mutable_tensor()->mutable_data()
        ->Add(data, data+length);
    Status status = stub_
        ->Update(&context, request, &reply);
    memcpy(data, reply.tensor().data(),
        length * sizeof(double))
    train(data);
}
```

Challenges in RPC-based INC Programming

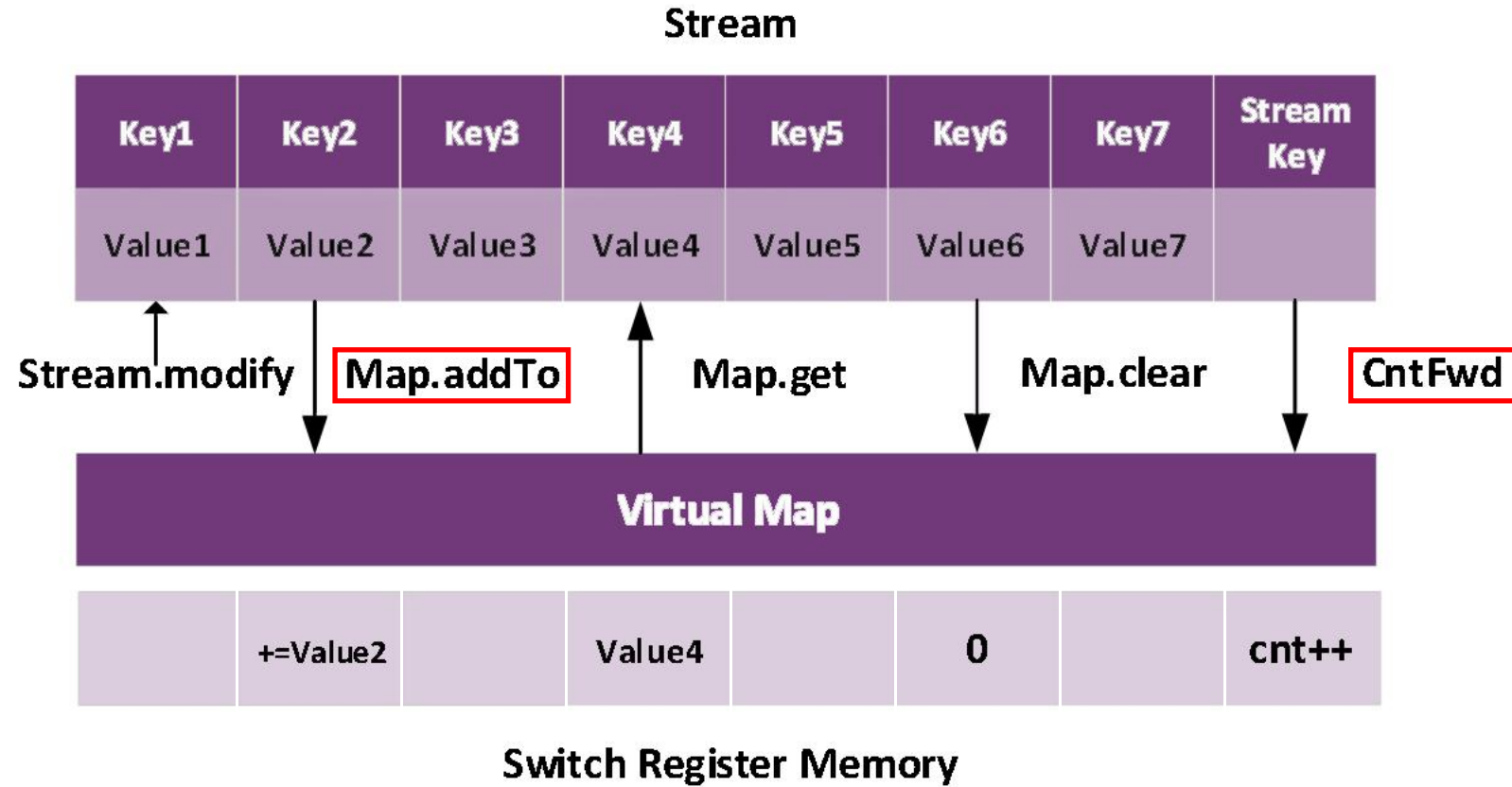


Switch Program is Complex, but We Can Provide High-level Primitives

- We identify a minimum set of **primitives** to compose INC applications, named reliable INC primitives (**RIPs**)
- We hope to use the description of INC primitives (**Netfilter**) to replace switch programs

Primitive	Args	Semantics
Map.addTo	stream	map[stream.key] += stream.value
Map.get	stream	stream.value = map[stream.key]
Map.clear	empty	map[stream.key] = 0
Stream.modify	op, para	stream.value = op(stream.value, para)
CntFwd	key, th, tgt	cnt[key]++; if cnt[key] == th then forward(tgt) else drop

We Implement RIPs Using Host and Switch Memory



NetRPC Programming Examples: Very Similar to gRPC

```
1 import "netrpc.proto"
2 message NewGrad {
3   netrpc.FPArray tensor = 1;
4 }
5 message AgtrGrad {
6   netrpc.FPArray tensor = 1;
7 }
8 service Training {
9   rpc Update(NewGrad) returns (AgtrGrad)
10  {} filter "agtr.nf"
```

Protobuf

```
1 { //agtr.nf
2   "AppName": "DT-1",
3   "Precision": 8,
4   "get": "AgtrGrad.tensor",
5   "addTo": "NewGrad.tensor",
6   "clear": "copy",
7   "modify": "nop",
8   "CntFwd": {
9     "to": "ALL",
10    "threshold": 2,
11    "key": "ClientID",
12  },
13 }
```

Netfilter

INC-enabled data types

Indicating NetFilter file name

```
1 shared_ptr<Channel> channel = CreateCustomChannel(server_ip,
2   InsecureChannelCredentials());
3 unique_ptr<Stub> stub_(NewStub(channel));
4 void PushPull(double* data, int length) {
5   NewGrad request;
6   AgtrGrad reply;
7   ClientContext context;
8   request.mutable_tensor()->mutable_data()
9   ->Add(data, data+length);
10  Status status = stub_
11    ->Update(&context, request, &reply);
12  memcpy(data, reply.tensor().data(),
13    length * sizeof(double))
14  train(data);
```

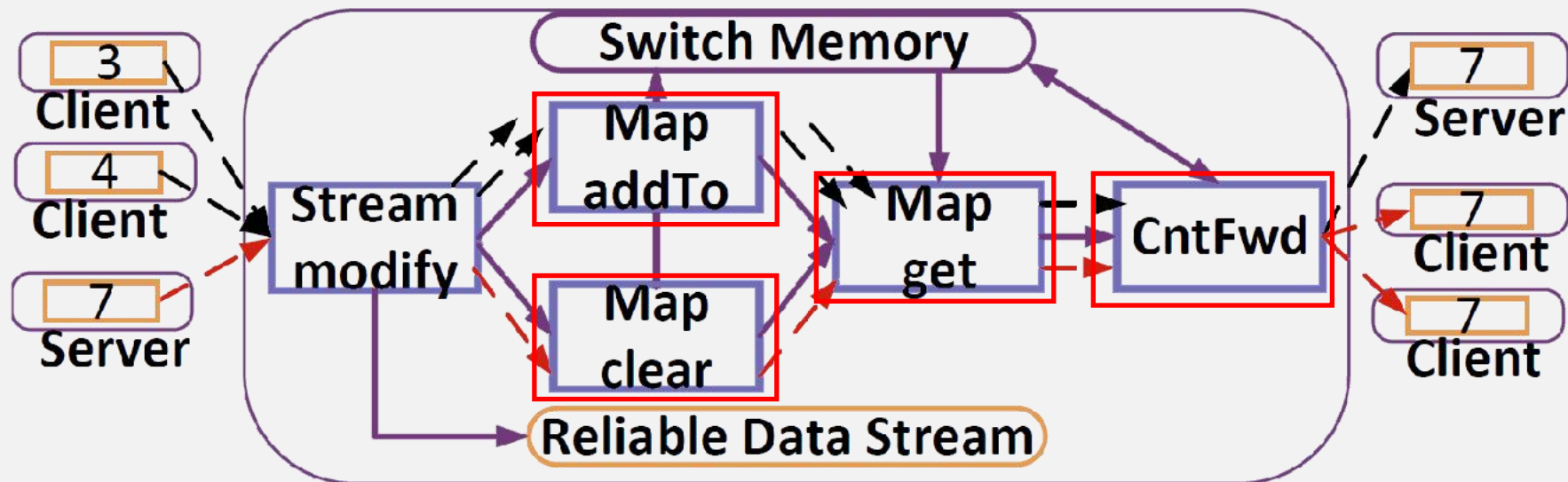
Quantization factor

Reliable INC primitives

RPC


Support Concurrent INC Applications in One Switch

- We implement RIPs on the programmable switch to support multiple applications concurrently:



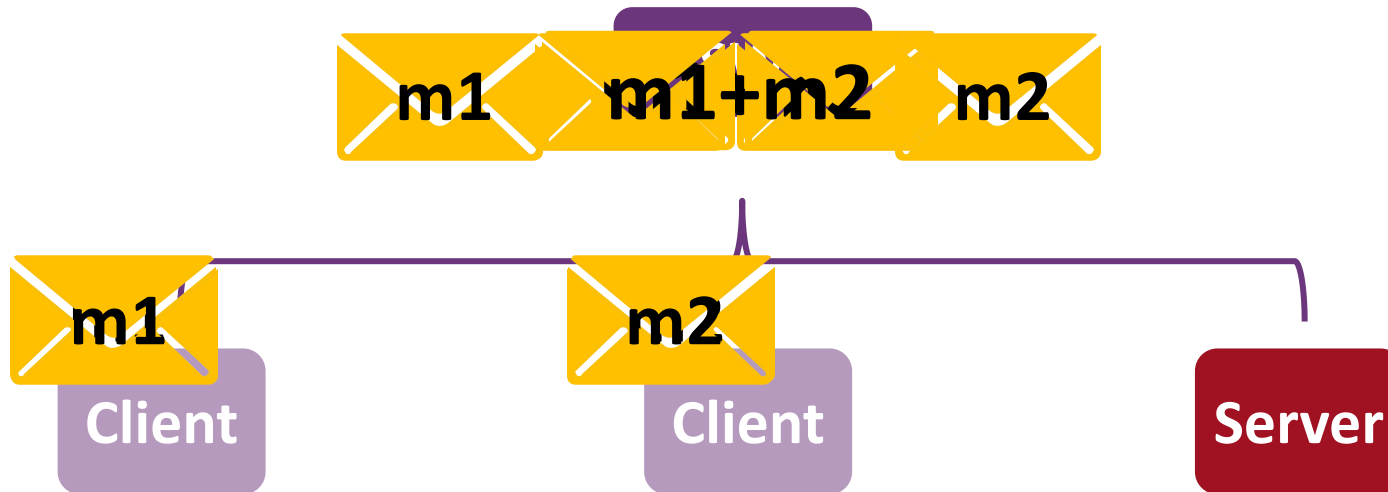
Reliable INC Requires Memory-Efficient Idempotence

- RPC calls should always **succeed eventually**, so RIPs should be same
- INC requires idempotence in addition
 - a. Sockets only guarantee at -least-once packet transmission
 - b. However, **repetitive accumulation** on the switch causes incorrect result
 - c. Normal path of some INC applications do not involve servers (**on-switch reliability**)
- We need to **detect resent packets** with **limited switch memory**

Packet	x1	x2	x3	x4	x5	x6
Flip bit	1	1	1	0	0	0
	0		1 		1	
Switch States						

Reliable INC Requires Fallback to Fit RPC Calls

- INC can **fail** due to insufficient switch memory, computation overflow, etc.
- We implement all RPs on the hosts. When INC fails, the **RPC server** can complete computation instead



Utilizing Switch Memory Efficiently Guarantees INC Benefits

- Sufficient switch memory makes INC full effect
- We need a management scheme to utilize switch resource **efficiently**
- We address switch memory in a **key-value** level by **clients**

Value Stream

1	2	3	4	5	6	7	Stream Key
Value1	Value2	Value3	Value4	Value5	Value6	Value7	

Pool-based Streaming

Value5	Value2	Value3	Value4
Value5	Value2	Value3	Value4

Utilizing Switch Memory Efficiently Guarantees INC Benefits

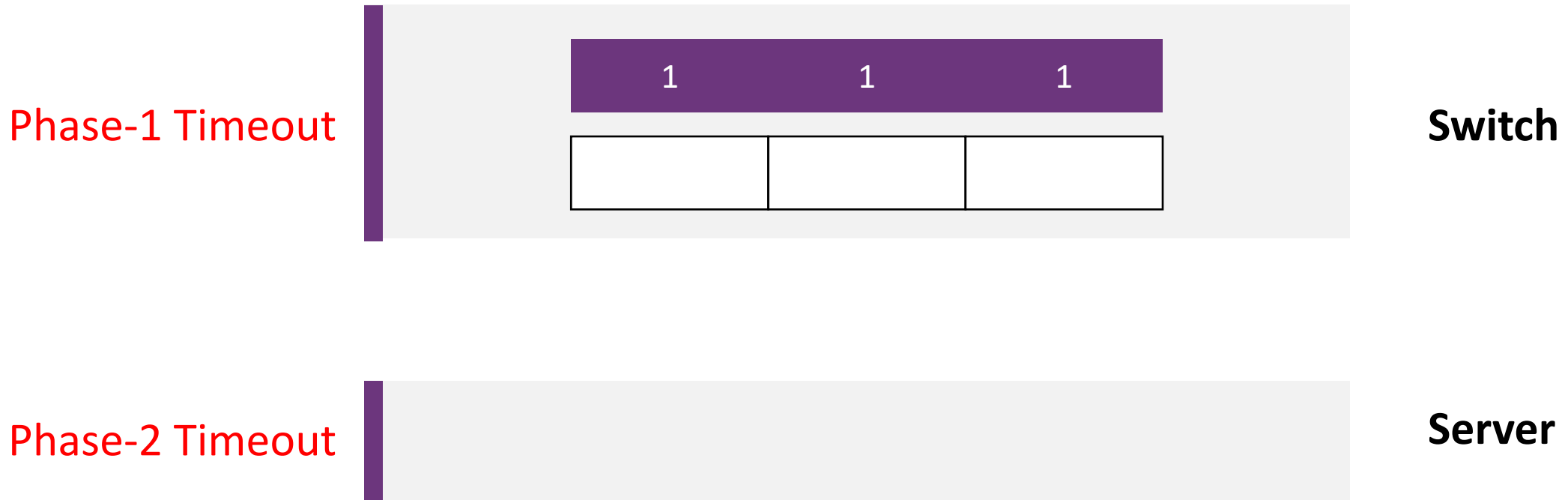
Key-value Stream	1	2	3	4	5	6	7	Stream Key
	Value1	Value2	Value3	Value4	Value5	Value6	Value7	1

On-switch Cache	Value1	Value2	Value3	Value4			
	Key1	Key2	Key3	Key4			
	Value1	Value2	Value8	Value4			
					Server		

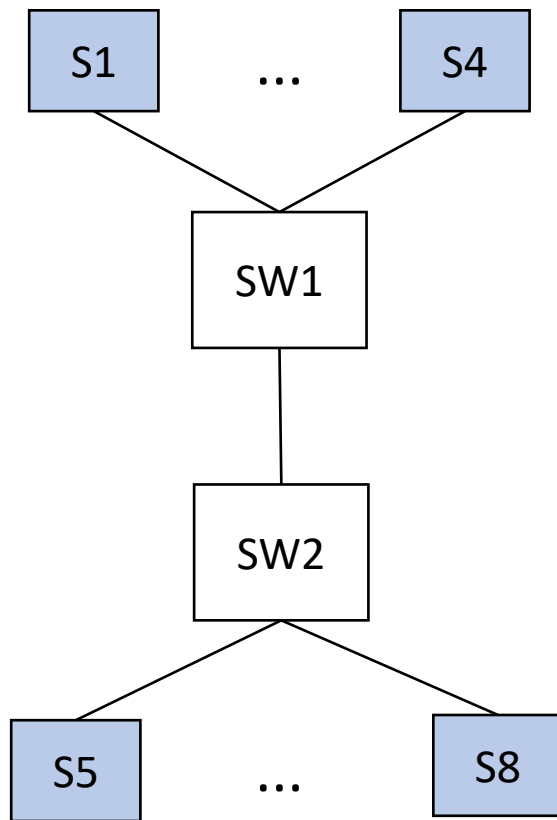
Key-value Stream	1	2	3	4	5	6	7	Stream Key
	Value1	Value2	Value3	Value4	Value5	Value6	Value7	2

On-Host Addressing Requires Handling Client Crash

- NetRPC relies on hosts to manage switch memory correctly
- **Memory leak** happens when the client crashes and loses states
- We apply a two-phase timeout to **recycle** valuable switch memory



NetRPC Evaluation: Setup

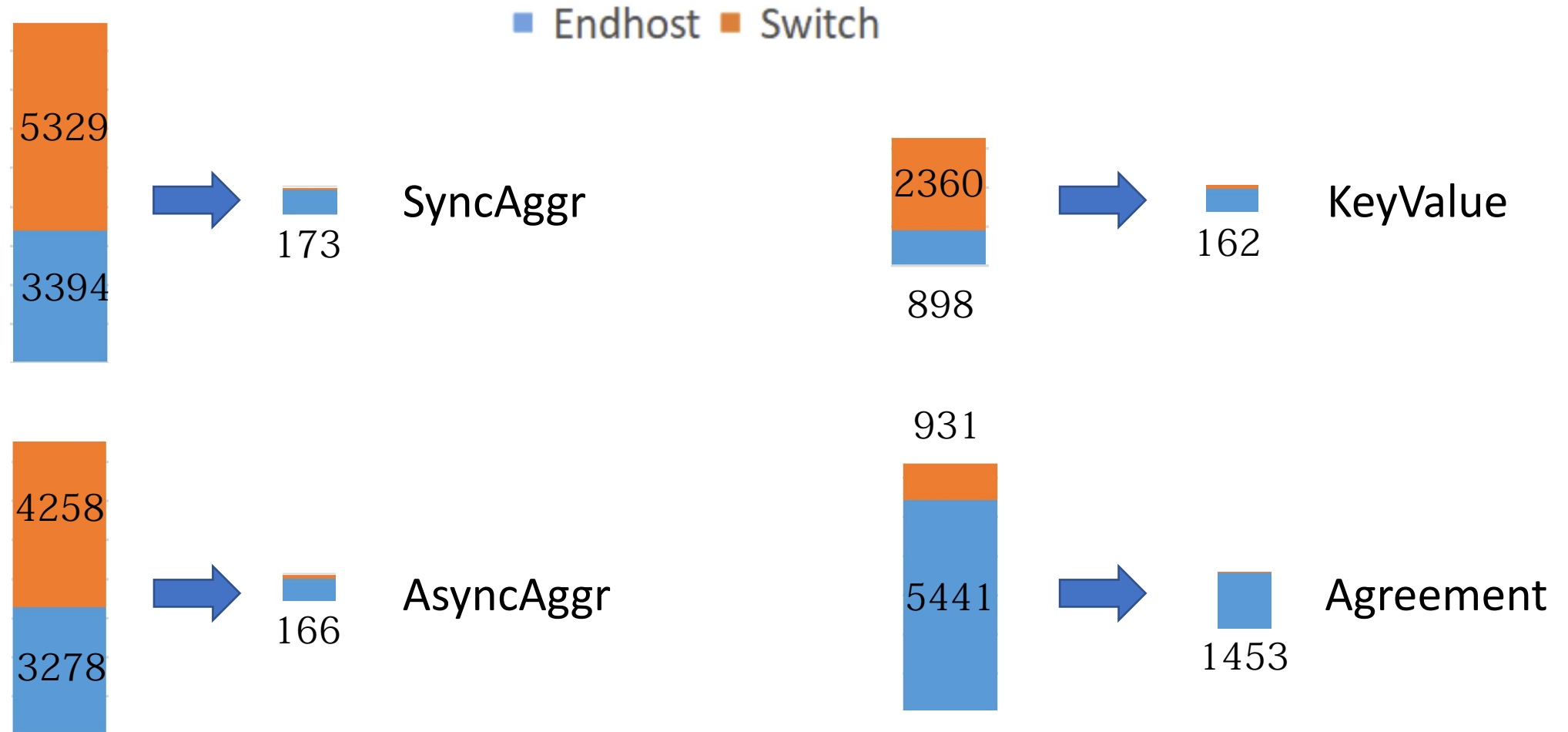


Type	Applications and Existing Systems
SyncAgtr	Distributed ML training (ATP, SHARP, SwitchML)
AsyncAgtr	MapReduce (ASK, NetAccel, Cheetah)
KeyValue	Cache (NetCache, DistCache), Monitoring (ElasticSketch)
Agreement	Synchronization (P4xos, NetChain, NetLock)

- Can NetRPC simplify INC programming?
- How does the NetRPC system perform?
- Can NetRPC support concurrent application?
- Can NetRPC guarantee reliability?

NetRPC Greatly Reduces User Code Complexity

- NetRPC reduces lines of code of INC application by up to 97%



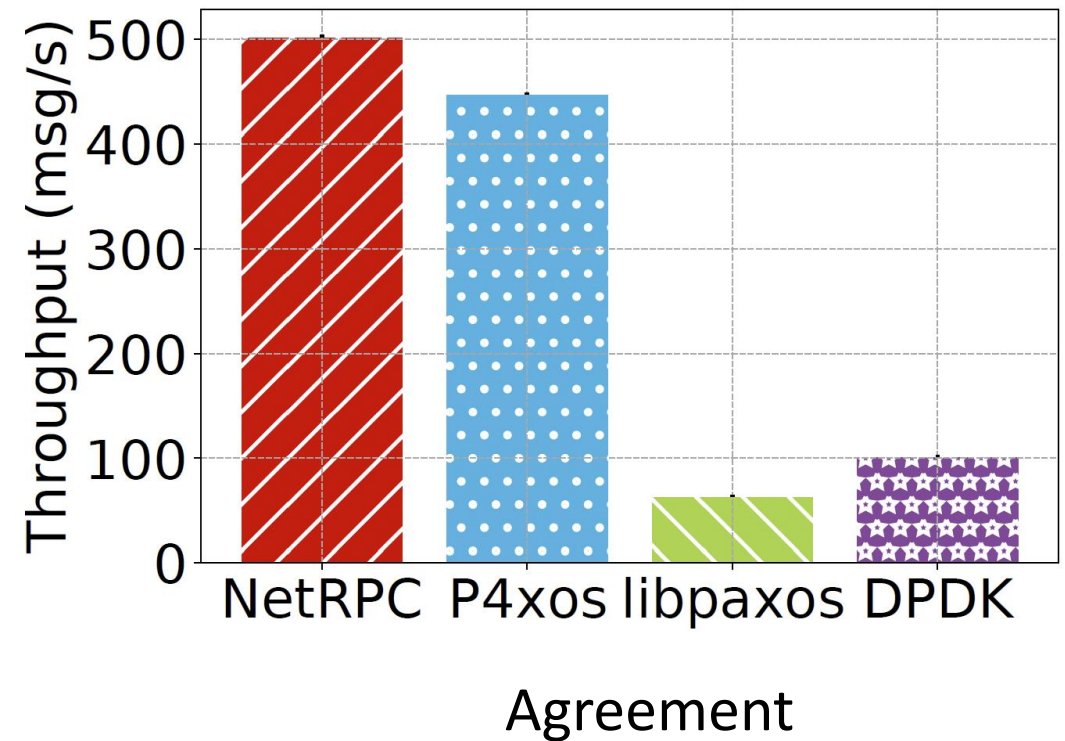
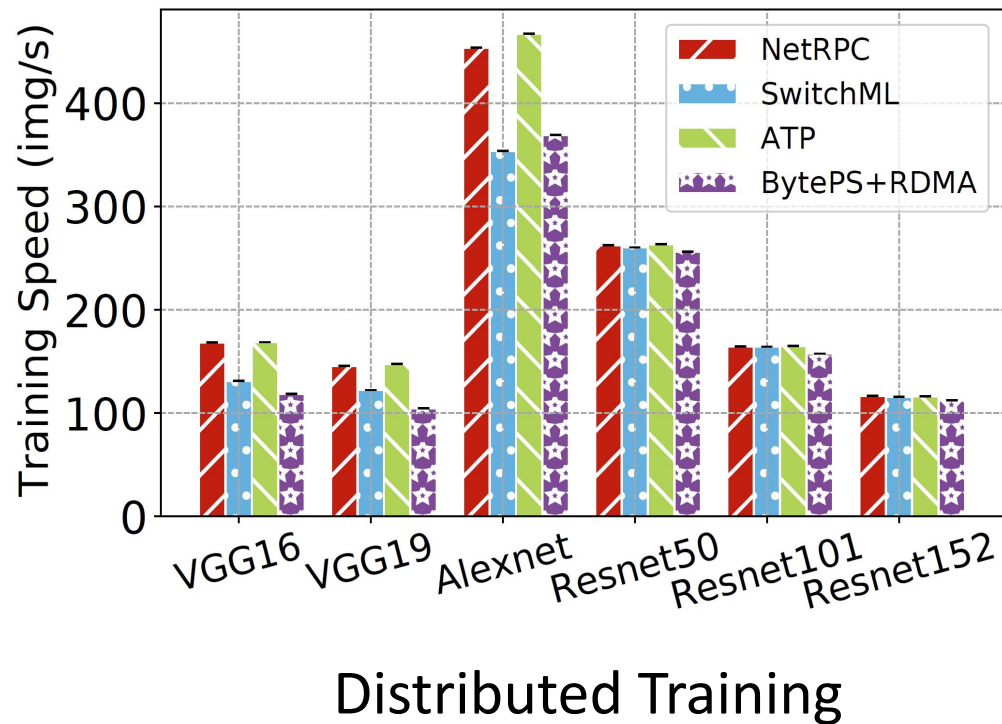
NetRPC Achieves Similar Performance to Handcrafted Code

- NetRPC achieves similar performance ($\geq 90\%$) to baselines even after programming simplification

Metrics	NetRPC	Prior Arts	DPDK
SyncAgtr Goodput(Gbps)	50.55	46.44(ATP)	40.11
AsyncAgtr Goodput(Gbps)	72.31	73.96(ASK)	45.88
Voting Delay(μ s)	20	22(P4xos)	92
Monitor Delay(ms)	3.52	3.26(ElasticSketch)	4.05

Faster than Handcrafted Code in End-to-end Application

- NetRPC achieves even better training throughput than ATP ($\geq 97\%$)
- NetRPC brings **12%** higher throughput than P4xos



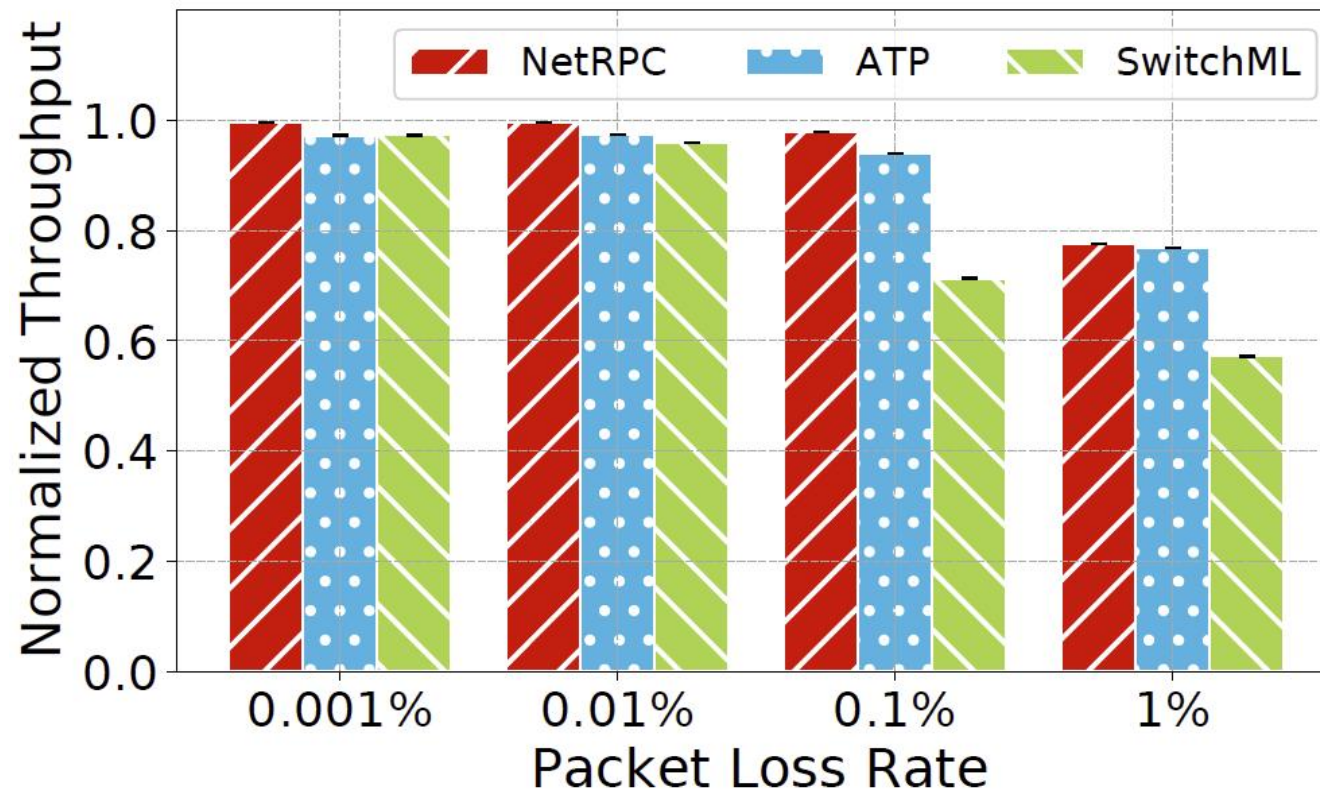
NetRPC Supports Multiple Concurrent Applications

- NetRPC can support concurrent INC applications with different types and different numbers

Metrics	1APP	4APP	4APP×5
Sync Goodput(Gbps)	50.55	24.88	24.84
Async Goodput(Gbps)	72.31	36.01	36.6
Goodput Sum(Gbps)	N/A	60.89	61.44
KeyValue Delay(ms)	3.52	3.56	3.85
AgreementDelay(μ s)	20	21	24

NetRPC is Reliable under Packet Loss

- NetRPC shows less performance degradation than prior arts with various packet loss rate.



Conclusion

- **NetRPC:**

The first framework that integrates INC into the familiar RPC programming model

- **Contribution:**

Make INC development easier and offer similar or better performance boosts than handcrafted systems

- **Future work:**

Explore scheduling policies and scale NetRPC to more complex topologies

=

Thanks!

