# Fast and Scalable In-network Lock Management Using Lock Fission
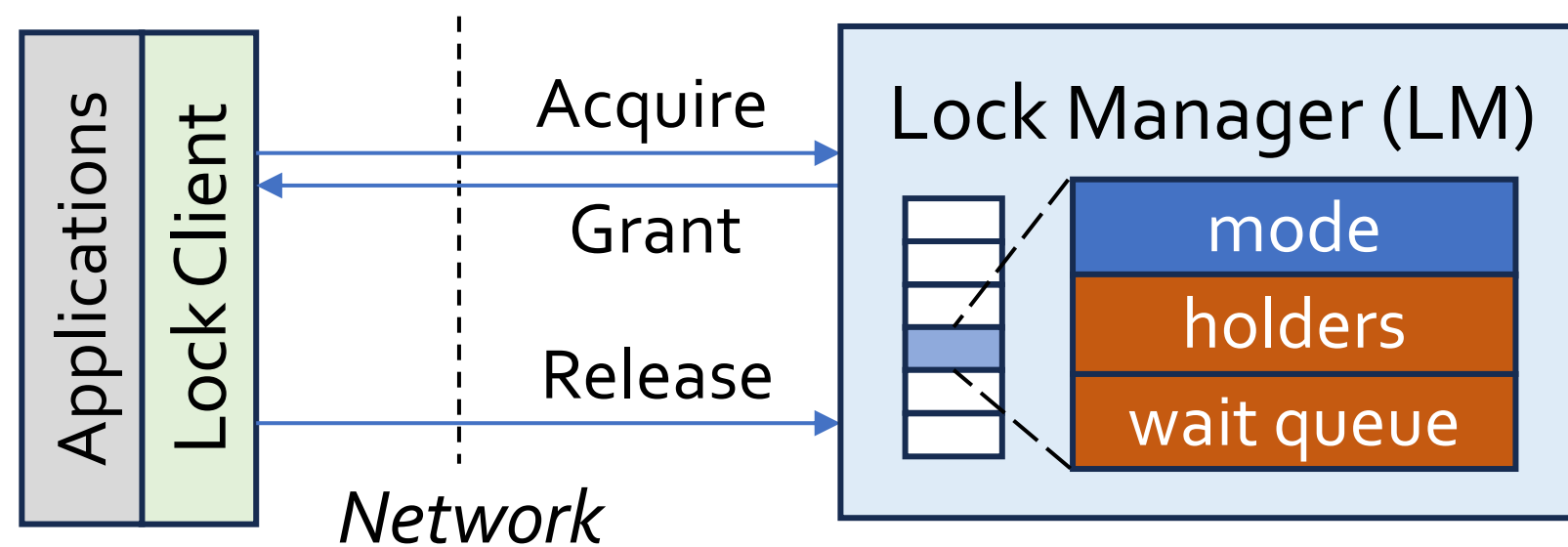
Hanze Zhang, 3$^{rd}$ year Ph.D. candidate from IPADS, SJTU

## Background

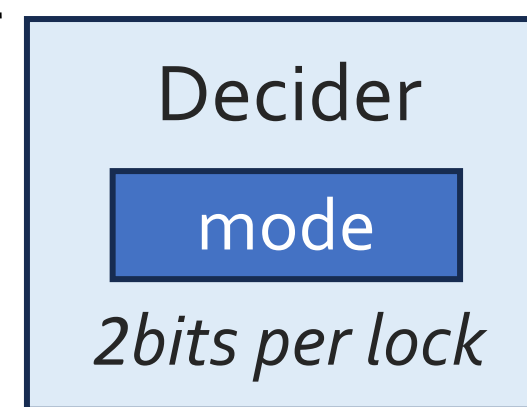Distributed workloads tends to have **low execution time** and **large data scale**:

| | | |
|---|---|---|
| Txn. Processing | 7/2.8 µs | 160M rows |
| File System | 1/10/20 µs | 10B files |
| Key-value Store | 8/15 µs | 250M keys |

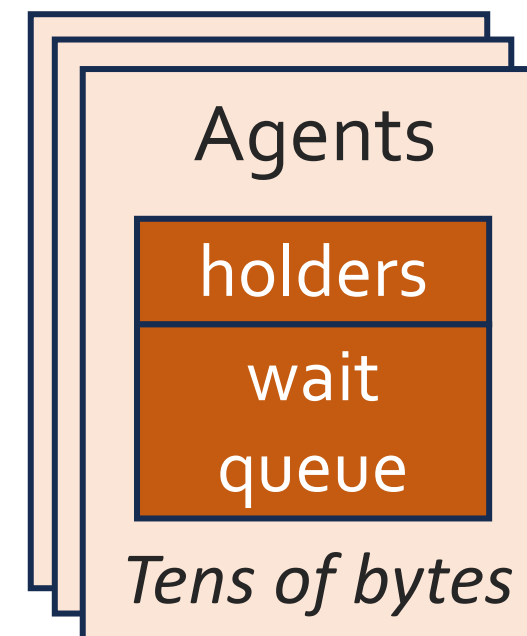Therefore, distributed lock services need to be **fast** and **scalable**.



## Key technique

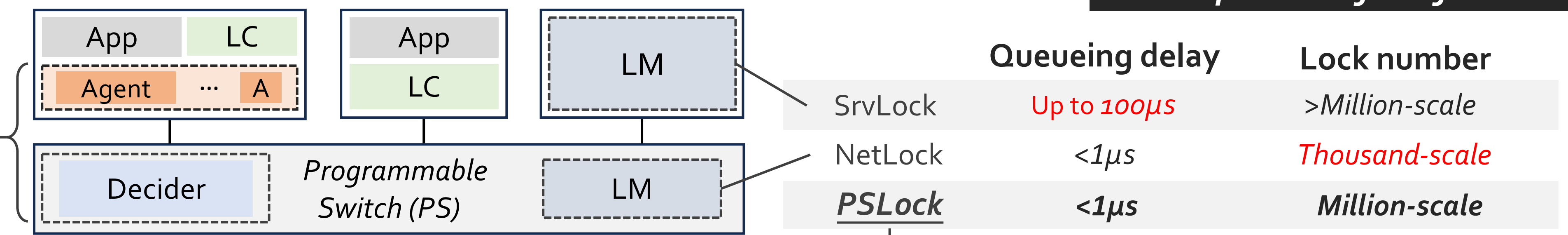We propose _**lock fission**_ that decouples lock grant **decision** and metadata **maintenance**:

**Decider**
mode
_2bits per lock_

### Synchronous decision
- _On grant critical path, must be **fast**_
- _Depends on **small, fixed-size** data_

_Suitable for programmable switch_

**Agents**
holders
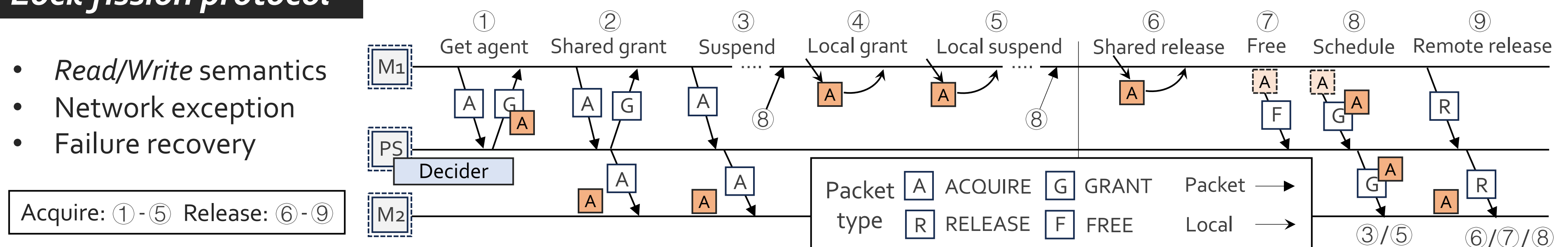wait queue
_Tens of bytes_

### Asynchronous maintenance
- _Off grant critical path, can be **slow**_
- _Depends on **large, variable-size** data_

_Suitable for servers_

- _On release critical path, require locality_

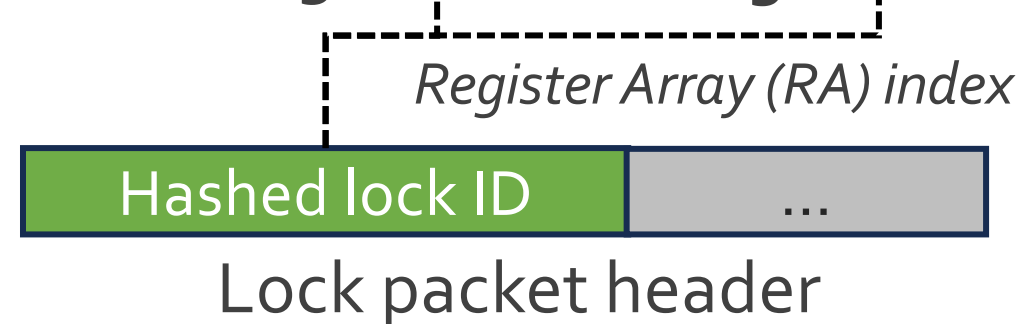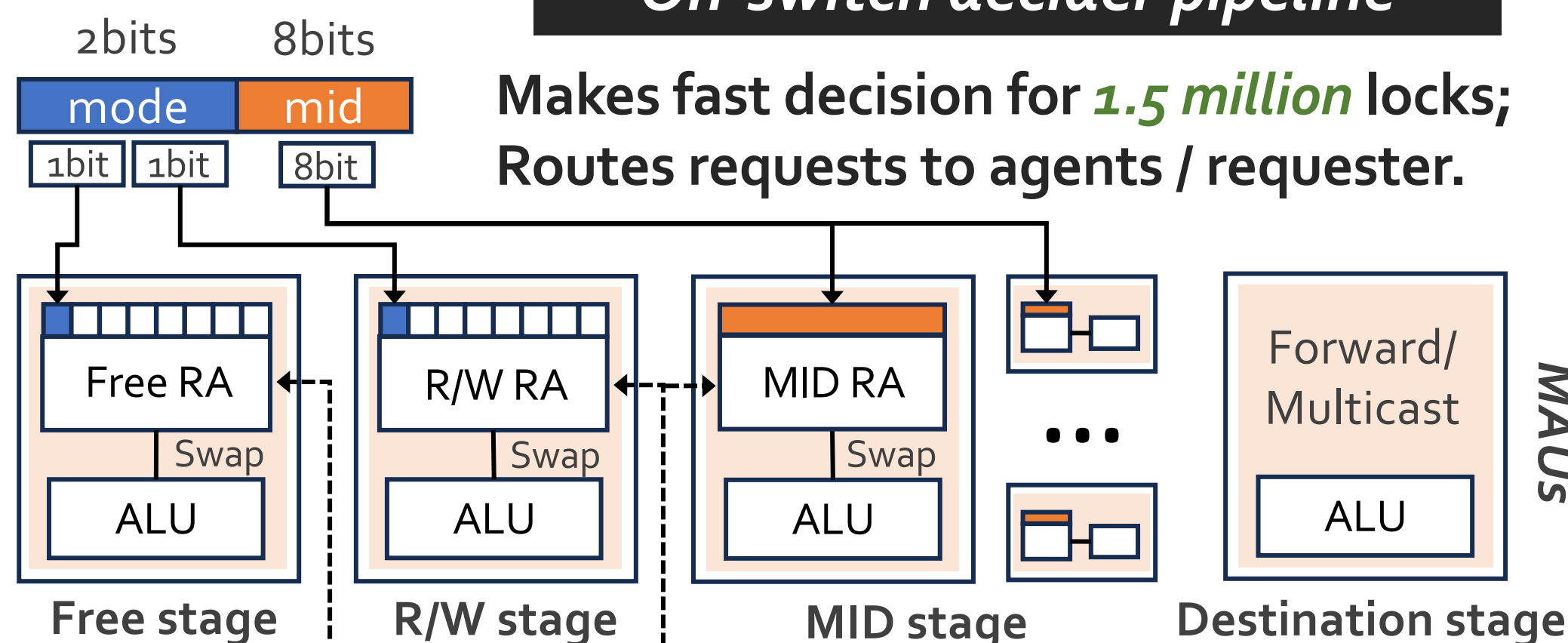_**Per-lock, migrate with lock ownership**_

_Fission_

## Comparison of LM forms



| | Queueing delay | Lock number |
|---|---|---|
| SrvLock | Up to _100µs_ | >Million-scale |
| NetLock | <1µs | Thousand-scale |
| **PSLock** | **<1µs** | **Million-scale** |

## Lock fission protocol

- _Read/Write_ semantics
- Network exception
- Failure recovery

Acquire: ①-⑤  Release: ⑥-⑨



## On-switch decider pipeline

2bits    8bits
mode    mid

Makes fast decision for _1.5 million_ locks;
Routes requests to agents / requester.



Free stage | R/W stage | MID stage | Destination stage

_Register Array (RA) index_
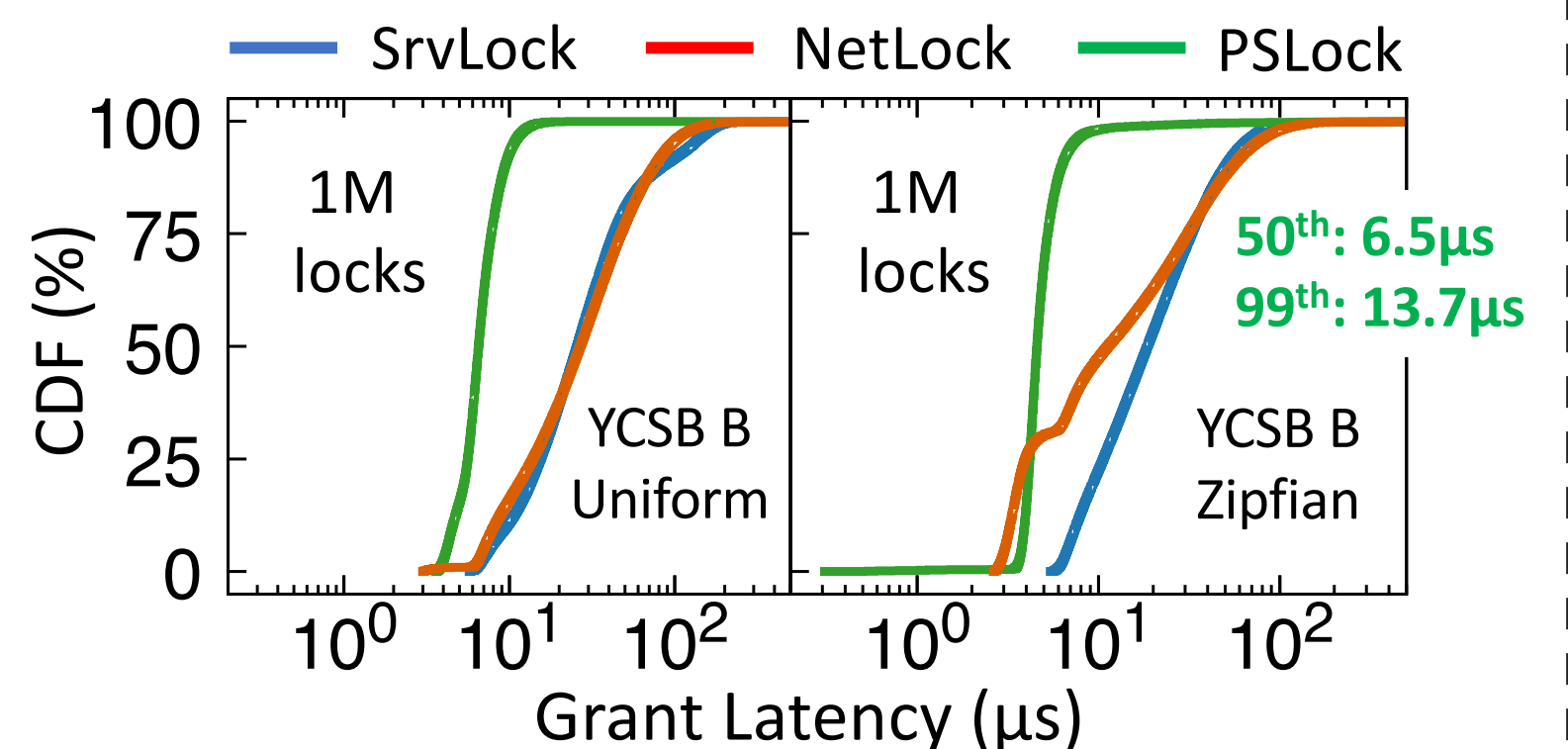
Hashed lock ID | ...
Lock packet header

### Switch memory challenges:
- Restricted layout
  (Arrays of registers)
- Limited access (Once per MAU)
- Scattered distribution
  (Limited per-MAU capacity)

## Evaluation

_Microbenchmark: Key-value store_



_Dynamic workload_



_Application benchmark: Transaction_