

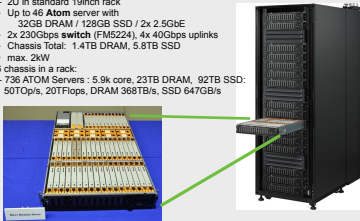
# RAMCloud on an Atom Server

Satoshi Matsushita  
Stanford University / NEC

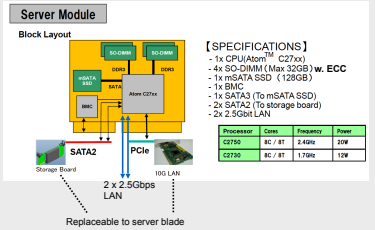


## Overview

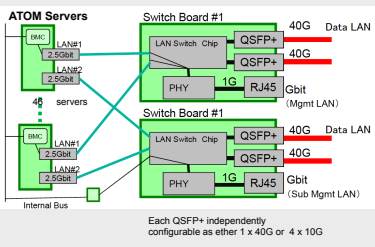
- NEC Micro Modular Server**
- Globally announced on May 20, 2014: ( Press release : [NEC raises the bar for high density IT solution platforms for the public and private cloud](#) )
  - Chassis: Redundancy (power supply, Networks, Fans) + Hot Swap
  - ~ 2U in standard 19inch rack
  - Up to 46 Atom server with
    - 32GB DRAM + 128GB SSD + 2x 2.5GbE
    - 2x 230Gbps switch (FM5224), 4x 40Gbps uplinks
    - Chassis Total: 1.4TB DRAM, 5.8TB SSD
  - max. 2kW
  - 16 chassis in a rack:
  - 736 ATOM Servers : 5.9k core, 23TB DRAM, 92TB SSD, 50TOP/s, 20TFlops, DRAM 368TB/s, SSD 64.7GB/s



## ATOM Server Blade

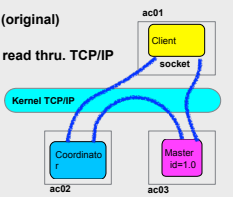


## Connection in a Chassis



## Base Evaluation

- Disable replication (backup) and collocation of entity
- CentOS 6.5
  - Ping **120-150 us** (original)
- Ported RAMCloud
  - 67.8 us** for 100B read thru. TCP/IP (tuned)



## User Space Driver

### Improvement with User Space Driver

- User space driver only for critical path, i.e. Master-Client data path
  - No modification in RAMCloud code, changing startup parameter.
  - Developed user mode driver for NIC2 based on Intel DPDK (Data Plane Development Kit)
- 
- ```

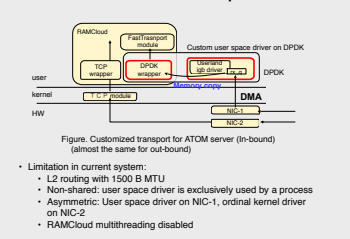
Command Line:
$ coordinator -C tcp:host=192.168.100.31,port=12246
$ server -C tcp:host=192.168.100.31,port=12246 -L fast+dpdk:host=192.168.101.29,mac=94:DE:86:AB:81:79,port=12247 + 0
$ ClusterPerf -C tcp:host=192.168.100.31,port=12246 --numClients 1 basic
    
```

### Development Platforms for User Space Driver

| PACKAGE                               | Summary                                     | Performance                                                              | License                     | Comment                                                                   |
|---------------------------------------|---------------------------------------------|--------------------------------------------------------------------------|-----------------------------|---------------------------------------------------------------------------|
| PACKET_MMAP                           | Implementation on the standard linux kernel | At least one buffer copy needed because a device buffer cannot be mapped | GPL                         |                                                                           |
| netmap                                |                                             |                                                                          | GPL/ BSD                    | Higher safety because user land code cannot access NIC registers directly |
| PF_RING / DMA (Direct NIC Access)     | Possible to map device queue to user space  | Feasible to realize zero-copy in user space driver                       | GPL/ BSD                    |                                                                           |
| Intel DPDK (Data Plane Developer Kit) |                                             |                                                                          | BSD (GPL for kernel module) | Rather widely used                                                        |

Our choice

## RAMCloud with User Space Driver



## Current Performance

Clusterperf.py basic, 30B key, Store and forward LAN switch  
Average and best/worst in 100 ms period. (7000 samples in 100B read)  
Room for tuning: [long tail](#) (Max), [slow write](#).

| Type        | Ave.     | Min.  | Max.  | Bandwidth  | Ave.     | Bandwidth  |
|-------------|----------|-------|-------|------------|----------|------------|
| 100B read   | 13.8 us  | 13.3  | 32.7  | 6.9 MB/s   | 5.1 us   | 18.7 MB/s  |
| 1KB read    | 20.7 us  | 20.0  | 37.7  | 46.1 MB/s  | 6.9 us   | 137.6 MB/s |
| 10KB read   | 52.8 us  | 52.1  | 68.6  | 180.8 MB/s | 10.4 us  | 914.1 MB/s |
| 100KB read  | 373.2 us | 371.3 | 379.0 | 255.5 MB/s | 47.2 us  | 2.0 GB/s   |
| 1MB read    | 3.9 ms   | 3.8   | 3.9   | 247.2 MB/s | 420.8 us | 2.2 GB/s   |
| 100B write  | 18.2 us  | 17.4  | 43.6  | 5.2 MB/s   | 15.7 us  | 6.1 MB/s   |
| 1KB write   | 25.6 us  | 24.7  | 64.1  | 37.2 MB/s  | 19.9 us  | 48.0 MB/s  |
| 10KB write  | 64.2 us  | 62.5  | 95.5  | 148.6 MB/s | 38.5 us  | 247.7 MB/s |
| 100KB write | 431.4 us | 423.2 | 463.0 | 221.0 MB/s | 235.3 us | 405.3 MB/s |
| 1MB write   | 4.7 ms   | 4.6   | 4.8   | 204.6 MB/s | 2.2 ms   | 436.0 MB/s |

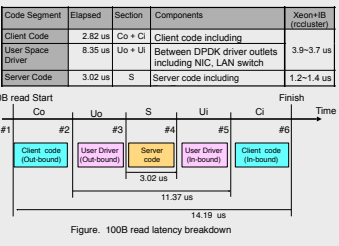
Backup Enabled

## Latency Analysis

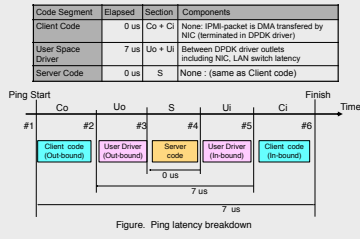
### Analysis

- Considerable gap between min. and max. implies room for improvement
- Latency breakdown
  - Comparison against ping with DPDK
  - Switch mode effect: store-and-forward vs. cut-through

### Latency Breakdown: 100B-Read



### Latency Breakdown: Ping



### Cut-through

Slight improvement for larger object size (due to 1500B MTU)  
Clusterperf.py basic, 30B key  
Average and best/worst in 100 ms period. (7000 samples in 100B read)

| LAN SW      | Store-and-Forward |       |       |            | Cut-Through |       |       |            |
|-------------|-------------------|-------|-------|------------|-------------|-------|-------|------------|
|             | Ave.              | Min.  | Max.  | Bandwidth  | Ave.        | Min.  | Max.  | Bandwidth  |
| 100B read   | 13.8 us           | 13.3  | 32.7  | 6.9 MB/s   | 13.8 us     | 13.3  | 32.7  | 6.9 MB/s   |
| 1KB read    | 20.7 us           | 20.0  | 37.7  | 46.1 MB/s  | 19.9 us     | 19.9  | 37.7  | 46.1 MB/s  |
| 10KB read   | 52.8 us           | 52.1  | 68.6  | 180.8 MB/s | 48.6 us     | 47.9  | 55.9  | 196.4 MB/s |
| 100KB read  | 373.2 us          | 371.3 | 379.0 | 255.5 MB/s | 609.6 us    | 567.3 | 576.1 | 258.4 MB/s |
| 1MB read    | 3.9 ms            | 3.8   | 3.9   | 247.2 MB/s | 3.4 ms      | 3.8   | 3.8   | 251.4 MB/s |
| 100B write  | 18.2 us           | 17.4  | 43.6  | 5.2 MB/s   | 18.2 us     | 17.4  | 43.6  | 5.2 MB/s   |
| 1KB write   | 25.6 us           | 24.7  | 64.1  | 37.2 MB/s  | 25.6 us     | 24.7  | 64.1  | 37.2 MB/s  |
| 10KB write  | 64.2 us           | 62.5  | 95.5  | 148.6 MB/s | 60.9 us     | 58.2  | 100.3 | 150.6 MB/s |
| 100KB write | 431.4 us          | 423.2 | 463.0 | 221.0 MB/s | 429.3 us    | 418.9 | 470.9 | 222.7 MB/s |
| 1MB write   | 4.7 ms            | 4.6   | 4.8   | 204.6 MB/s | 4.6 ms      | 4.7   | 4.7   | 206.8 MB/s |

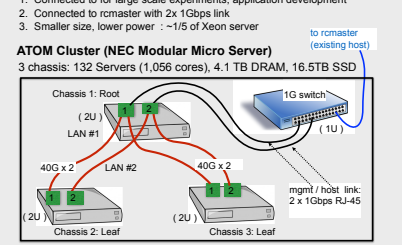
Improved Degraded

## Conclusions

### Considerations

- Large latency in user space driver (DPDK):
  - ~8.35 us with 100B read, 7 us with ping
- Copy overhead would be negligible:
  - ~0.4us for 100B (~1Kbit) transfer at 2.5Gbps
- Slight improvement with Cut-through mode
  - Negligible time for 100B memcopy (50 ns for 1KB copy on 2.4GHz Xeon)
- To tune DPDK driver:
  - Further latency breakdown
  - DPDK parameter tuning
  - Cache footprint optimization??

### Spine-switch-less cluster at Stanford



### Conclusions

- Initial performance evaluation:
  - 13.8 us for 100B-read with custom user space driver
  - on ATOM server through chassis switch (1 hop)
- Further analysis and tuning
- Functional enhancement:
  - Symmetric driver and link aggregation with two NICs
  - Providing a turn-key-solution
  - with job/network/storage/VM management tools
  - on a standardized hardware platform
- Further evaluation on a larger scale system
  - On a new ATOM cluster at Stanford
  - Application development and evaluation
- Very welcome for feedback to improve the Micro modular server and future systems