

Adapting HW to SW for the Cloud

Han Dong

Mass Open Cloud Workshop

10/30/2018



<http://sesa.github.io/>

Addressing Performance of Datacenter Workloads

- Kernel Bypass (Unikernels, User-level network stacks, DPDK, FaRM, ...)
- Hardware Offload (Near-data Processing, Accelerators, FPGAs, Disaggregated Datacenters, ...)

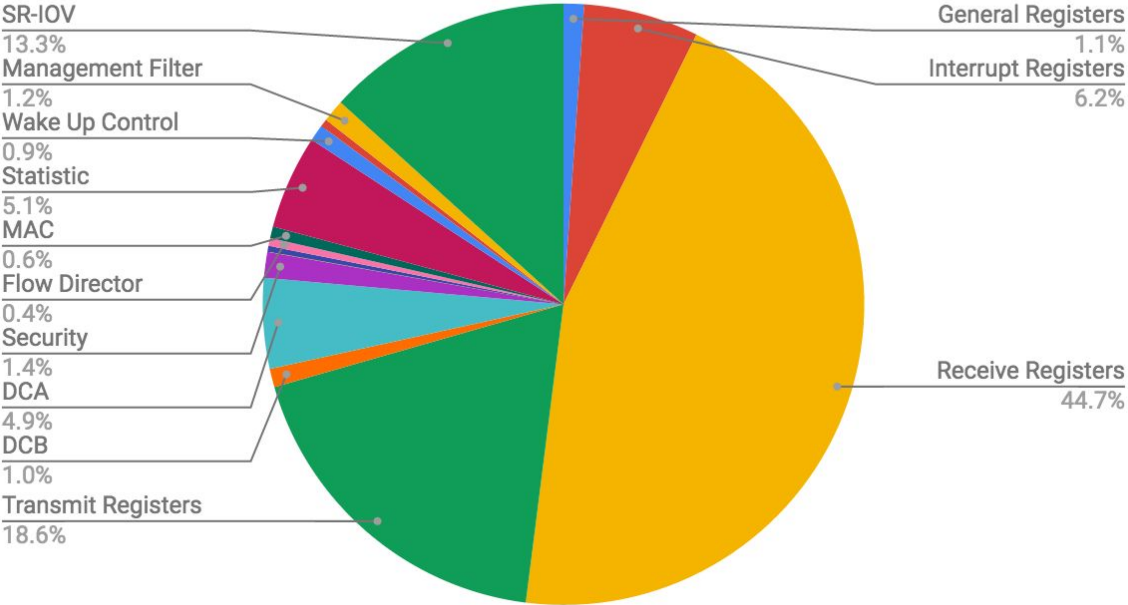
Opportunity to dynamically adapt the HW
to the SW

Intel 82599 Datasheet

Total Registers:
~5630

Linux Device Register
Initialization: ~1360

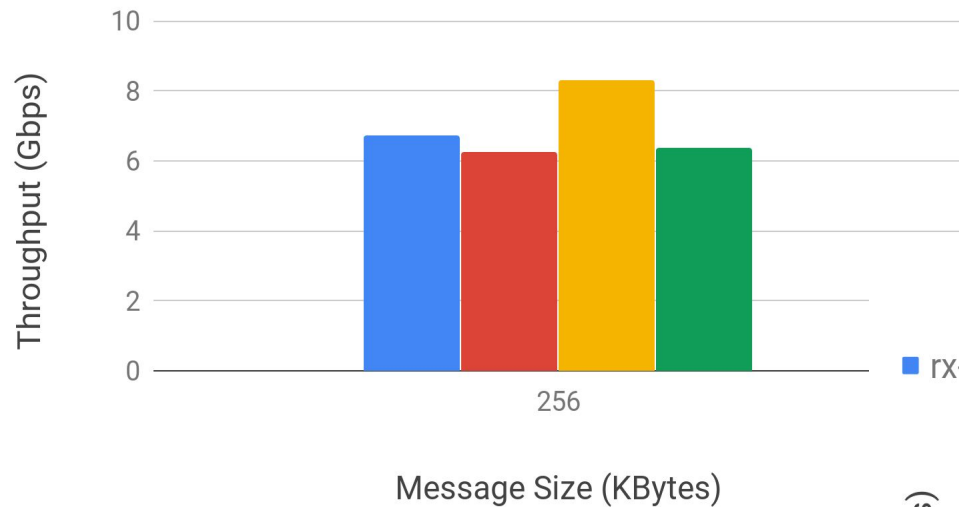
Intel 82599 NIC Register Breakdown



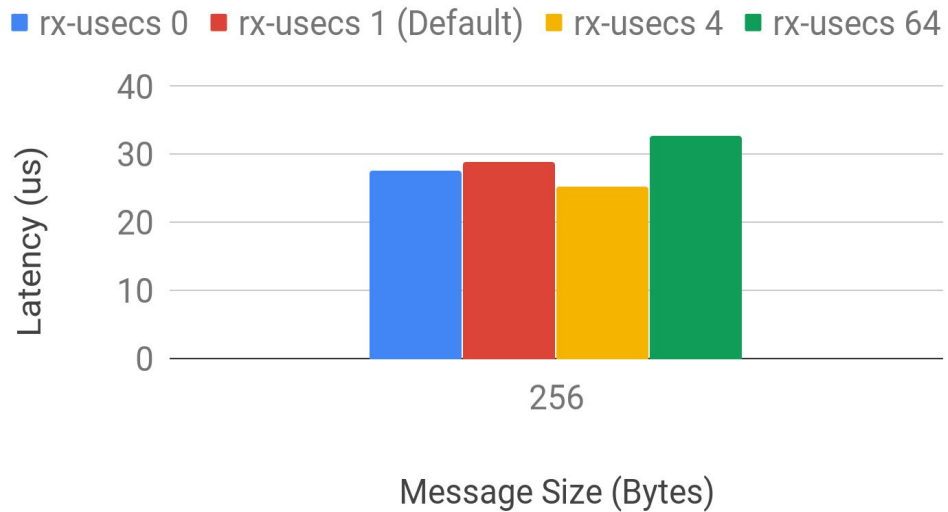
Experimental Setup

1. Two MOC machines:
 - a. CPU: 2 * 8 Core Intel(R) Xeon(R) CPU E5-2650
 - b. NIC: Intel Corporation Ethernet X520 10GbE
2. Linux ethtool:
 - a. rx-usecs - how many usecs to delay an RX interrupt after a packet arrives
 - i. 0us, 1us (Linux Default), 4us, and 64us
3. Benchmarks:
 - a. iperf, netpipe, memcache (Facebook ETC workload)

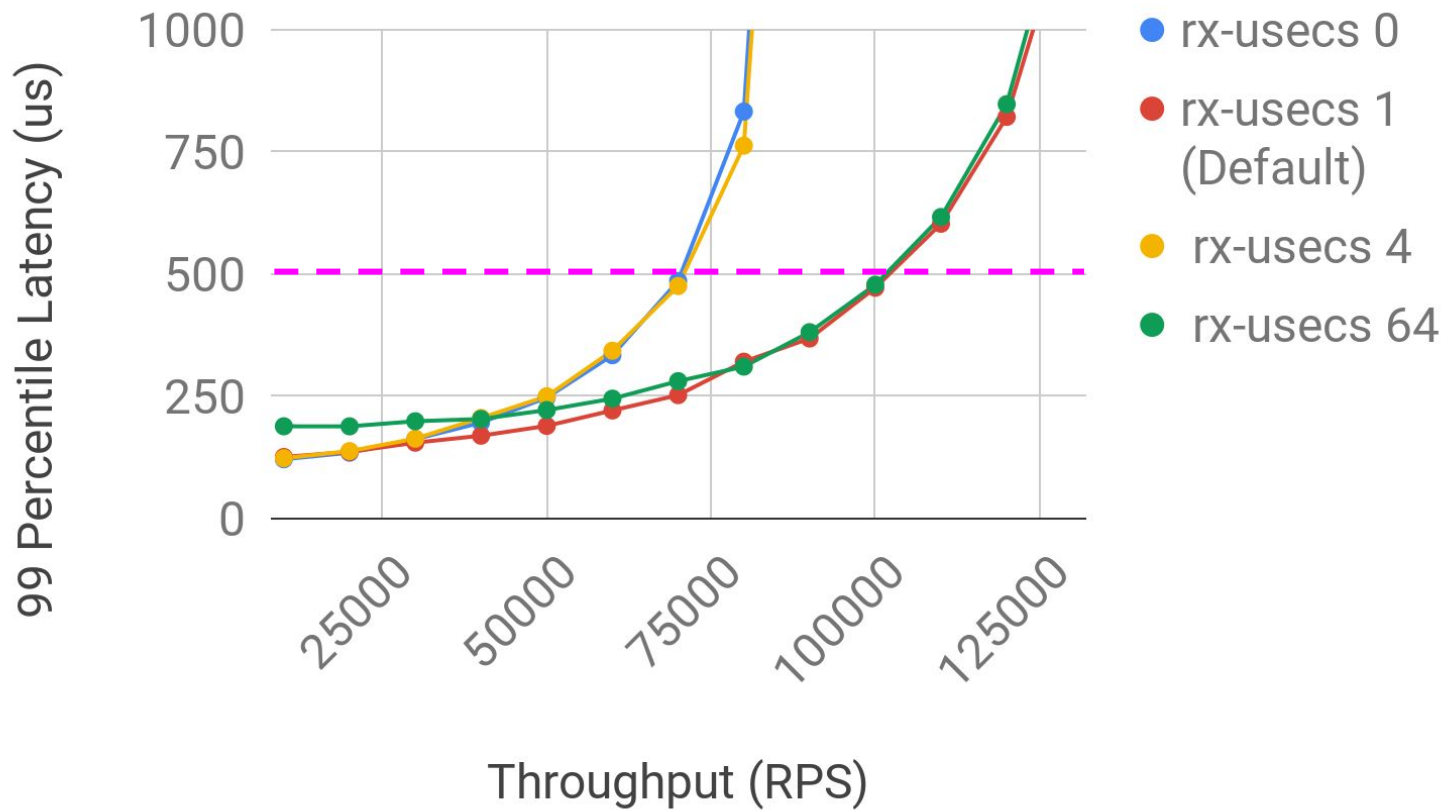
iperf



netpipe



memcache



Questions

- No single configuration works for all workloads.
- Which registers have an impact on network benchmarks?
 - Which registers have more “weight”?
- How to begin discovery of said registers?
- Is there a way to expose settings up to application?
- What about other devices?

Future Directions

1. Frame configuration problem as a learning problem:
 - a. NIC contains ~5k unique registers, ~80 statistic registers
 - b. ethtool contains ~100 device features that are tunable in Linux
2. Dynamically modify device features as workload changes:
 - a. Cache injection, interrupts, packet filter
3. Per application optimization of HW and SW via unikernels or kernel bypass