# Hodor: Intra-Process Isolation for High-Throughput Data Plane Libraries
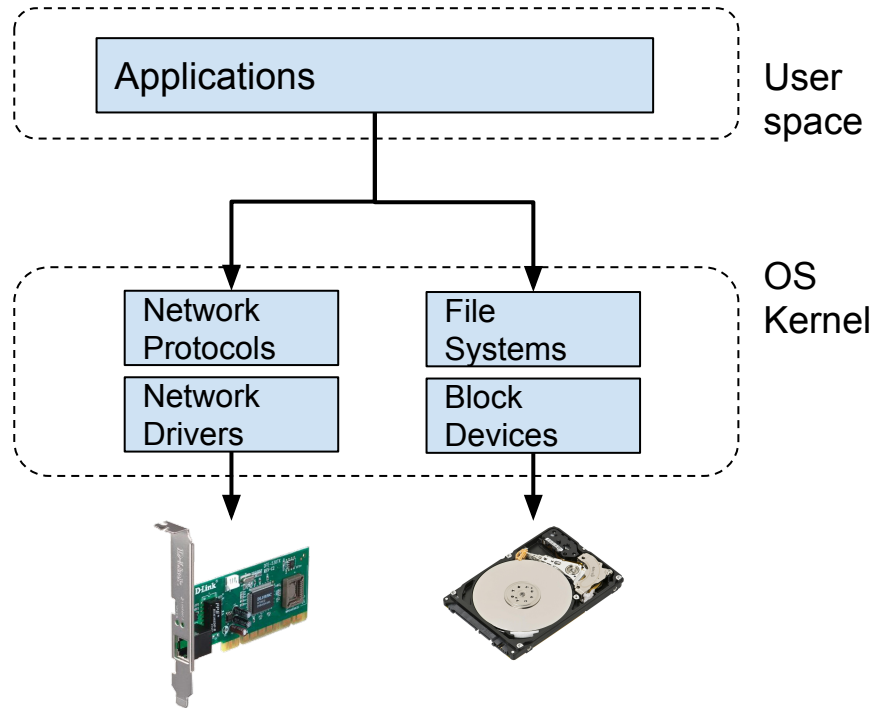
Mohammad Hedayati[1], Spyridoula Gravani[1], Ethan Johnson[1], John Criswell[1], Michael L. Scott[1], Kai Shen[2], Mike Marty[2]

[1]University of Rochester
[2]Google

UNIVERSITY *of* ROCHESTER

# Kernel I/O

- Enables sharing
- Provides guarantees
  - Fairness
  - Recovery
  - Security

- All without needing to trust
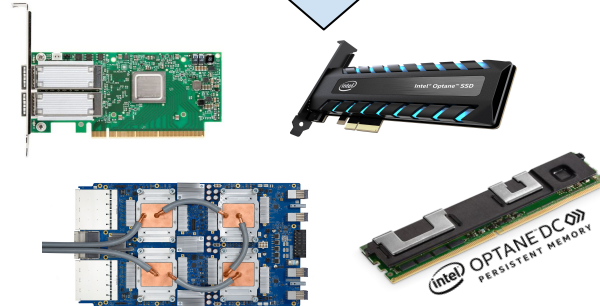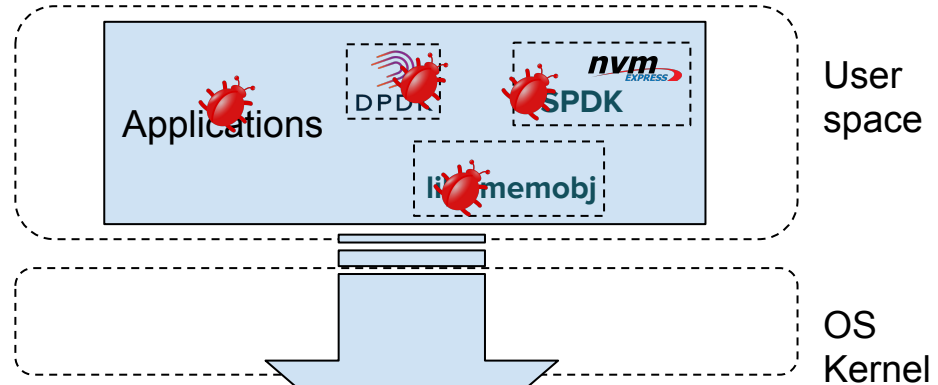
| Applications | User space |

Network Protocols

Network Drivers

File Systems

Block Devices

OS Kernel

# Kernel-bypass I/O

- **Pros:**
  - Lower latency
  - Rapid development
  - Specialization



User space

OS Kernel

- **Cons:**
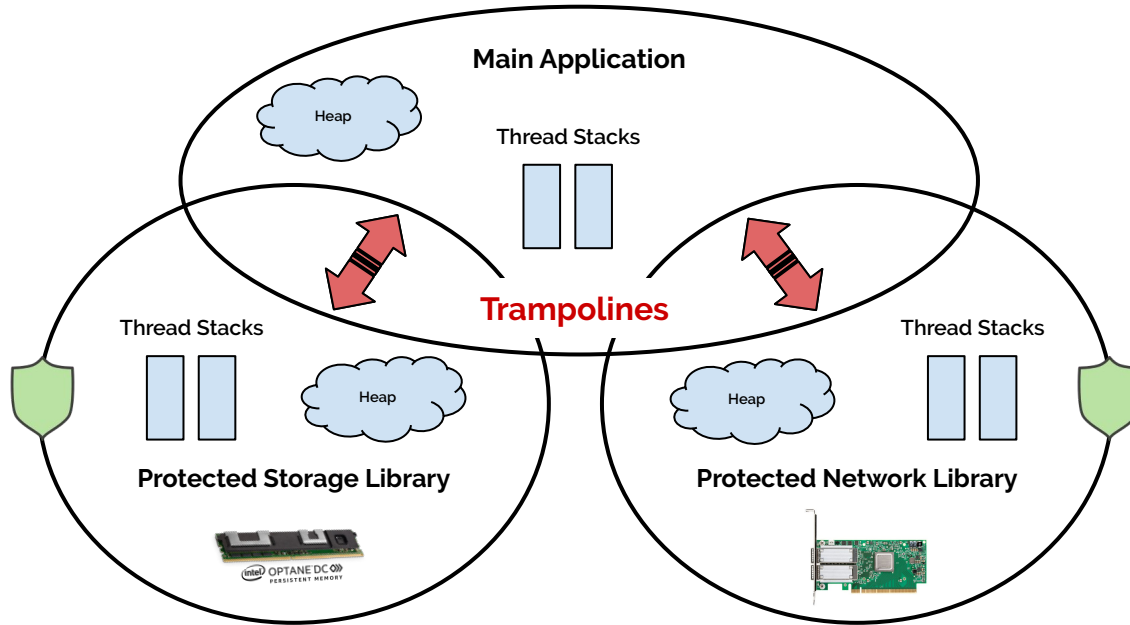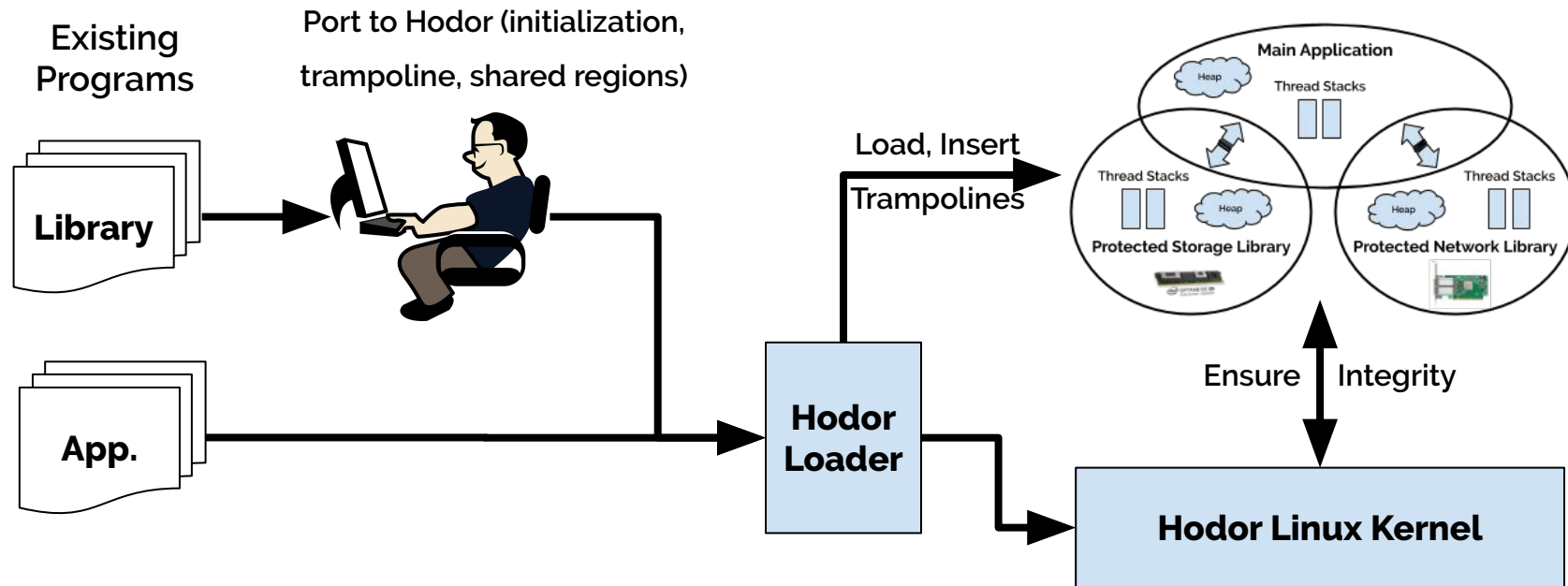  - No guarantees
  - Hard to multiplex

# Overview

- Motivation

- Design of Hodor

- Fast Memory Isolation

- Evaluation

- Conclusion

# Protected Library

# Hodor



Existing Programs

Port to Hodor (initialization, trampoline, shared regions)

**Library**

**App.**

Load, Insert Trampolines

Main Application
Heap
Thread Stacks

Thread Stacks
Heap
Protected Storage Library

Heap
Thread Stacks
Protected Network Library

**Hodor Loader**

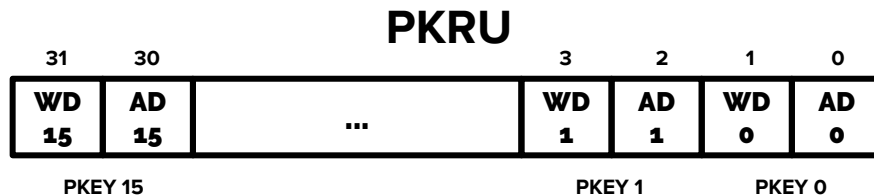Ensure Integrity

**Hodor Linux Kernel**

# Overview

- Motivation

- Design of Hodor

- Fast Memory Isolation

- Evaluation

- Conclusion

# Intel PKU

- Protection Keys for User-Space (a.k.a. MPK)

- Introduced in Skylake-SP

**PTE**

| X D | **PKEY** | *###* | U | W | P |
|-----|----------|-------|---|---|---|

Positions marked: **62** and **59** above PKEY.

- 32-bit PKRU register (Access/Write Disable)

- `WRPKRU/RDPKRU`

**PKRU**

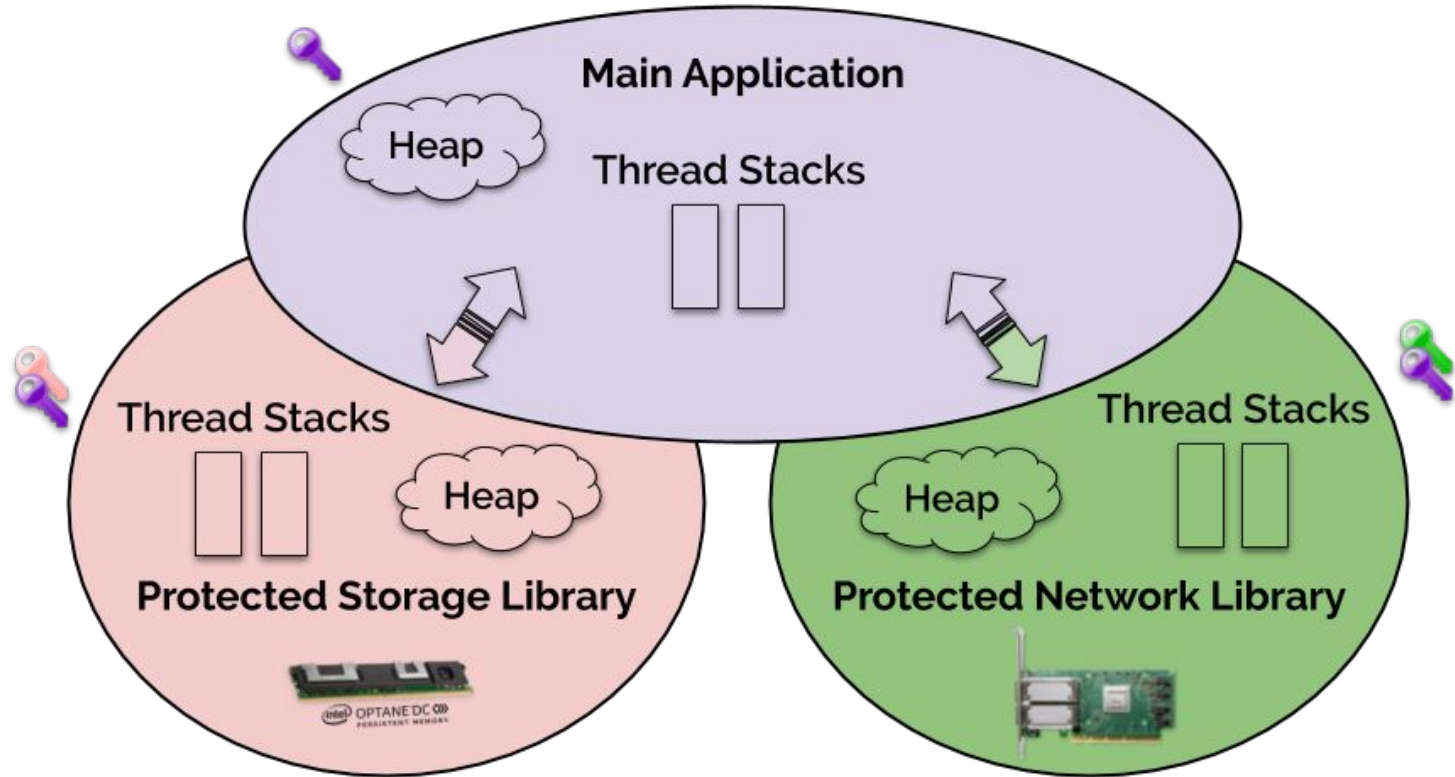| 31 | 30 | | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|----|
| WD 15 | AD 15 | ... | WD 1 | AD 1 | WD 0 | AD 0 |
| PKEY 15 | | | PKEY 1 | | PKEY 0 | |

# Intel PKU



based on a figure from: https://ubm.io/2YjGvFE
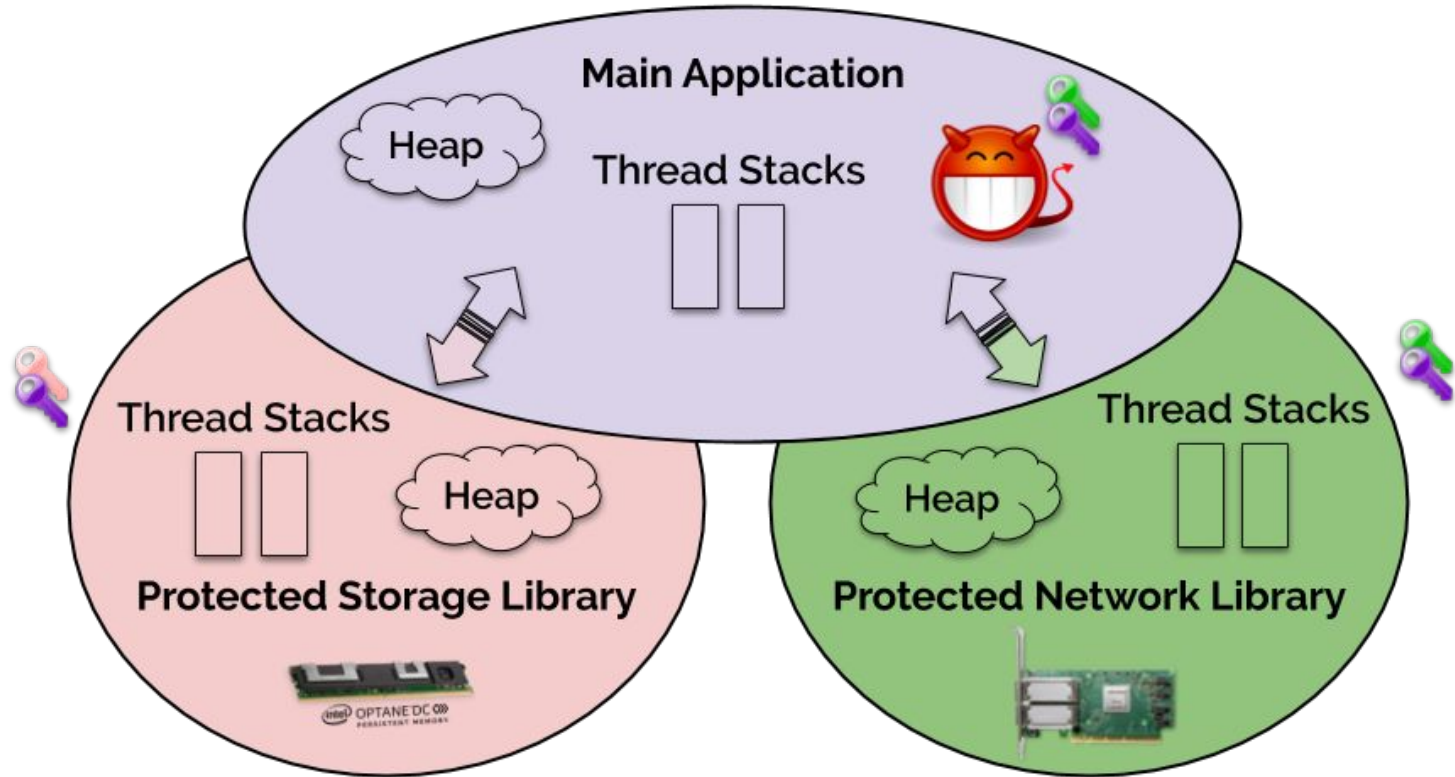
# Hodor: Memory Isolation

# Hodor: Memory Isolation

# Hodor: Vetting `WRPKRUS`

- Inspect executable regions
    - Load (by Hodor loader)
    - W→X change (by Hodor kernel at run-time)

- Look for `WRPKRU (0f 01 ef)` instances
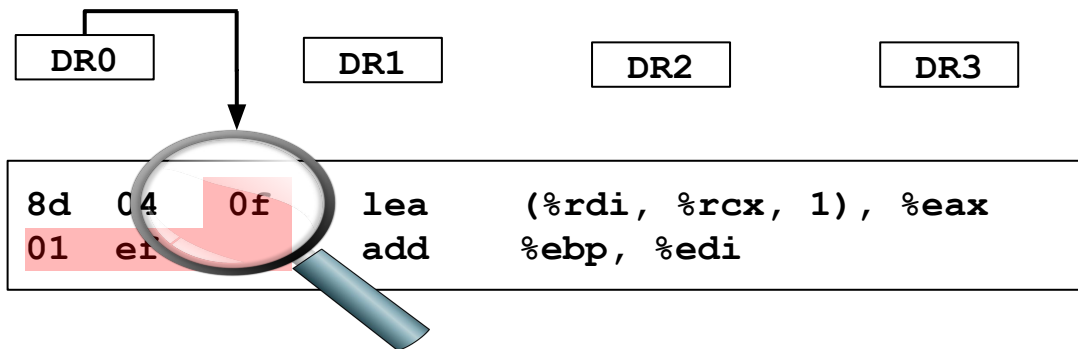
**glibc-devel-static-2.27-alt9.x86_64**

```
f7 d2          not      %edx
21 d0          and      %edx,%eax
44 89 c2       mov      %r8d,%edx
09 f0          or       %esi,%eax
0f 01 ef       wrpkru
31 c0          xor      %eax,%eax
```

**blender-2.79b-7.fc29.x86_64**

```
8d  04   0f    lea      (%rdi, %rcx, 1), %eax
01  ef         add      %ebp, %edi
```
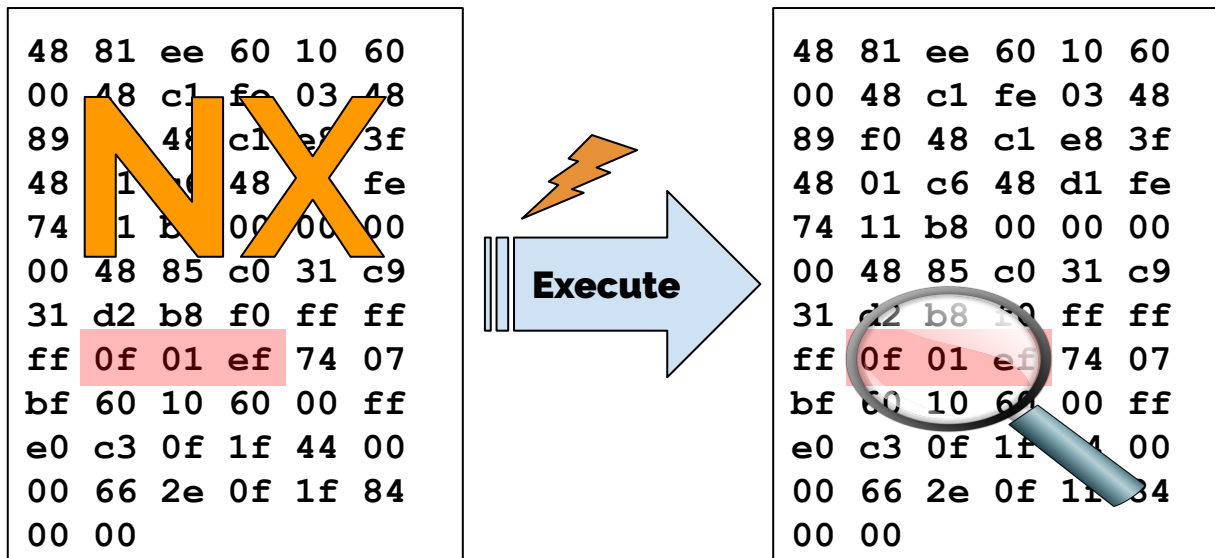
UNIVERSITY *of* ROCHESTER

# Hodor: Vetting `WRPKRUS`

- Hardware watchpoints
  - DR# registers point to the beginning of illegal byte sequence
  - No spurious traps when correctly aligned execution runs past an implicit instance

| DR0 | DR1 | DR2 | DR3 |
|-----|-----|-----|-----|

```
8d  04   0f     lea     (%rdi, %rcx, 1), %eax
01  ei          add     %ebp, %edi
```

# Hodor: Illegal WRPKRUS

- Limited hardware watchpoints
  - Only 4 on Intel Processors
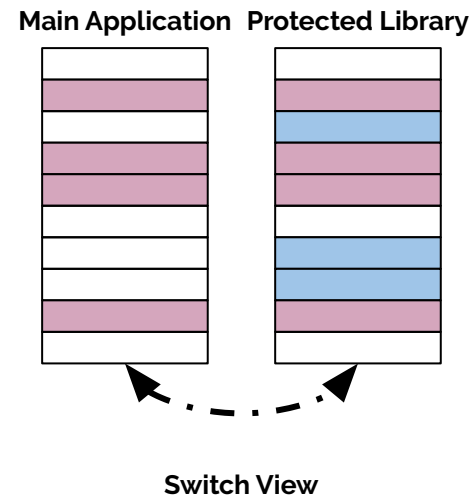  - HW watchpoints as **cache** for illegal sequences

# Hodor: Vetting `WRPKRUS`

- Vetting cost

  - Implicit instances incur no run-time overhead

  - Explicit instances should use `pkey_set()`

  - No measurable overhead as long as:

    - #hot illegal seq. fewer than #hw watchpoints

- How often?

  - **58,273** rpm packages on Fedora 29 (**108K executables**)

  - Only **123** binaries with one or more illegal byte sequences

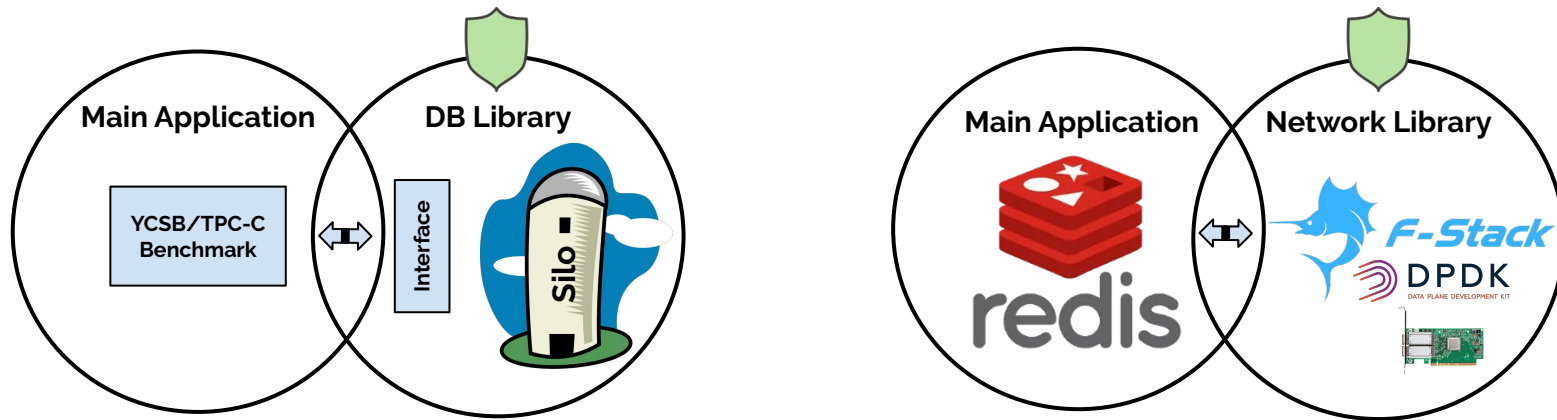    - Only **2** (**less than 0.02%**) with more than 4

# Alternative Memory Isolation

- **Per-Domain Page-Table**
  - Each mapping the view of a domain
  - Switch using system calls

- **Per-Domain Extended Page-Table**
  - Requires running virtualized (in Intel VMX)
  - Switch using VMFUNC w/o causing VMEXIT

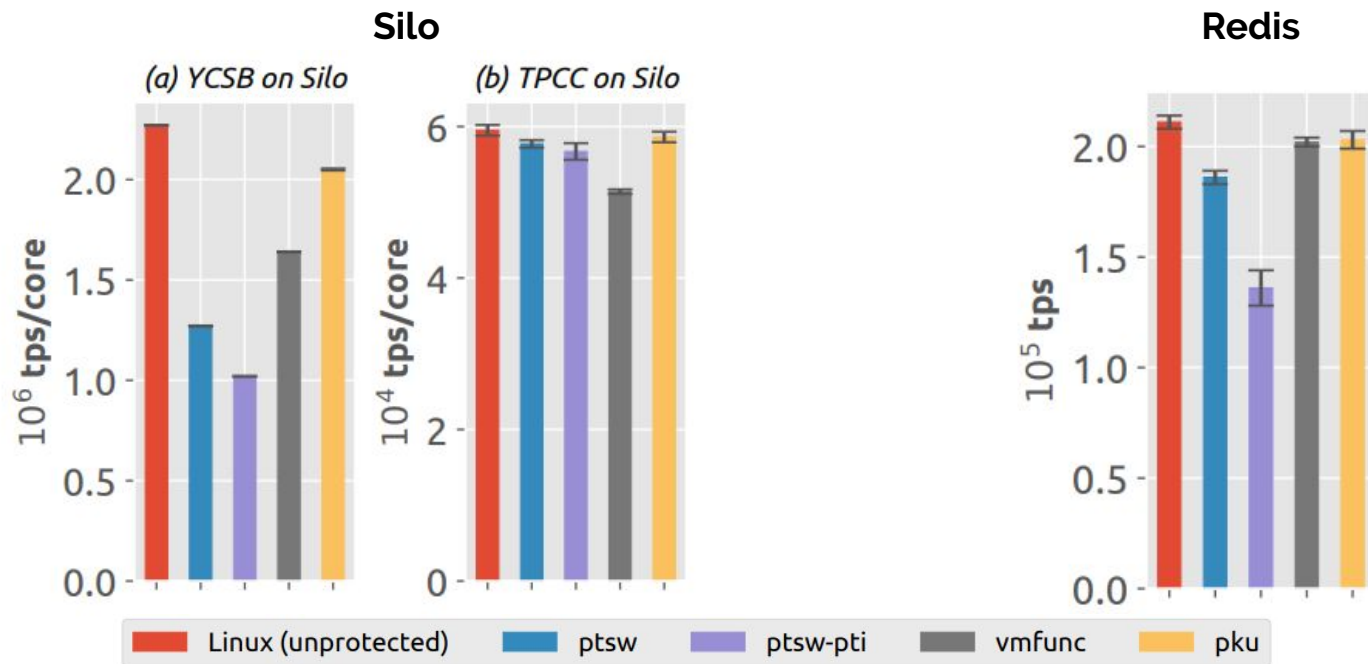**Main Application**   **Protected Library**

**Switch View**

# Evaluation: Applications

- Silo (in-memory database)

- Redis (kernel-bypass network TCP/IP stack)

- DPDK (kernel-bypass packet processing) -- in the paper

# Evaluation: Applications

**Silo**

**Redis**



(a) YCSB on Silo

(b) TPCC on Silo

Legend: Linux (unprotected), ptsw, ptsw-pti, vmfunc, pku

**ptsw**: page table switching    **ptsw-pti**: page table switching w/ KPTI,
**vmfunc**: EPT switching    **pku**: using memory protection keys

# Conclusion

Introduced:

- *Protected Libraries*: new OS abstraction for library isolation

Showed that:

- Intel PKU can be safely used to isolate protected libraries
- Doing so does not sacrifice performance
  - **90–98%** of unprotected throughput

**See the paper:** How multiple processes can share a protected library