

Paper-Reading-Group

Nested Kernel: An Operating System Architecture for Intra-Kernel Privilege Separation

Nathan Dautenhahn, Theodoros Kasampalis, Will Dietz, John Criswell, and Vikram Adve

Nested Kernel - Motivation

- Monoliths have single large TCB
- How to separate into multiple protection domains?
 - Microkernels require complete redesign of kernel
 - VMMs have high performance overhead

How can we provide protection domains without the overhead using the existing code base and design principles?

Nested Kernel - Idea

- Use the MMU to isolate the MMU
 - *Nested Kernel* is small and protects MMU structures
 - *Outer Kernel* is de-privileged and only has checked access to MMU structures
- Keep the monolithic address space
- While still enabling application of intra-kernel security policies
 - Example policies: write logging, write-mediation

Nested Kernel - Design

- Separate policy from mechanism (MMU)
- OS Co-design for security policies
- MMU based privilege separation
- Fine grained resource control
- Negligible performance impact

Nested Kernel - Invariants

Invariant 1:

Active virtual-to-physical mappings for protected data are configured read-only while the outer kernel executes.

Invariant 2:

Write-protection permissions in active virtual-to-physical mappings are enforced while the outer kernel executes.

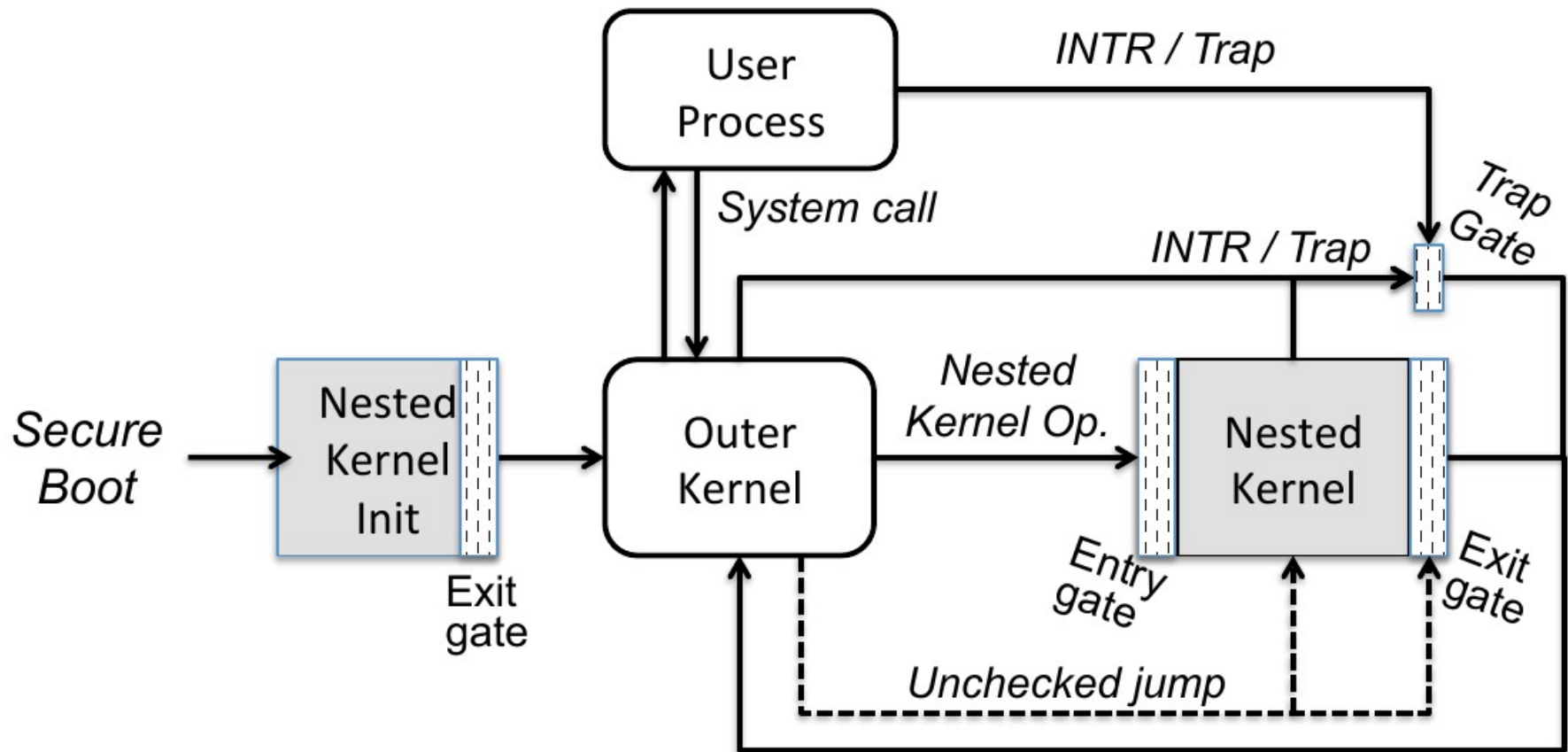
Nested Kernel – write protection

Function	Selected Arguments	Purpose
<code>nk_declare</code>	<code>mem_start, size, mediation_func</code>	Marks all pages RO; initializes an NK write descriptor <code>nk_wd</code> ; returns the <code>nk_wd</code> and the pointer to the region.
<code>nk_alloc</code>	<code>size, mediation_func, nk_wd_p</code>	Allocates memory region; invokes <code>nk_declare</code> on it; stores write descriptor in <code>nk_wd</code> ; returns <code>nk_wd</code> and pointer to the region.
<code>nk_free</code>	<code>nk_wd</code>	Deallocates memory identified by <code>nk_wd</code> . Memory must have been allocated by <code>nk_alloc</code> . Freed pages can be reused only by a future <code>nk_alloc</code> .
<code>nk_write</code>	<code>dest, src, size, nk_wd</code>	Verifies write bounds; invokes <code>mediation_func</code> , if any; then copies <code>size</code> bytes from <code>src</code> to <code>dest</code> .

PerspikuOS – Thread Model

- *Outer kernel* may be under complete attacker control
 - Can attempt to arbitrarily modify CPU state
 - Can modify outer kernel source code
 - Can modify control flow
- *Nested kernel* source code and binary are trusted
 - Including the mediation functions
- Hardware is free of vulnerabilities
- Do not protect against hardware attacks

PerspikuOS – Architecture



PerspikuOS – Invariant 1 Support

Reminder: Invariant 1:

Active virtual-to-physical mappings for protected data are configured read-only while the outer kernel executes.

Invariant 3:

Ensure that there are no unvalidated mappings prior to outer kernel execution

Invariant 4:

Only declared PTPs are used in mappings

Invariant 5:

All mappings to PTPs are marked read-only

Invariant 6:

CR3 is only loaded with a pre-declared top-level PTP.

PerspikuOS – Invariant 1 Support

Remin

Wri
whi

Invaria

Invaria

Invaria

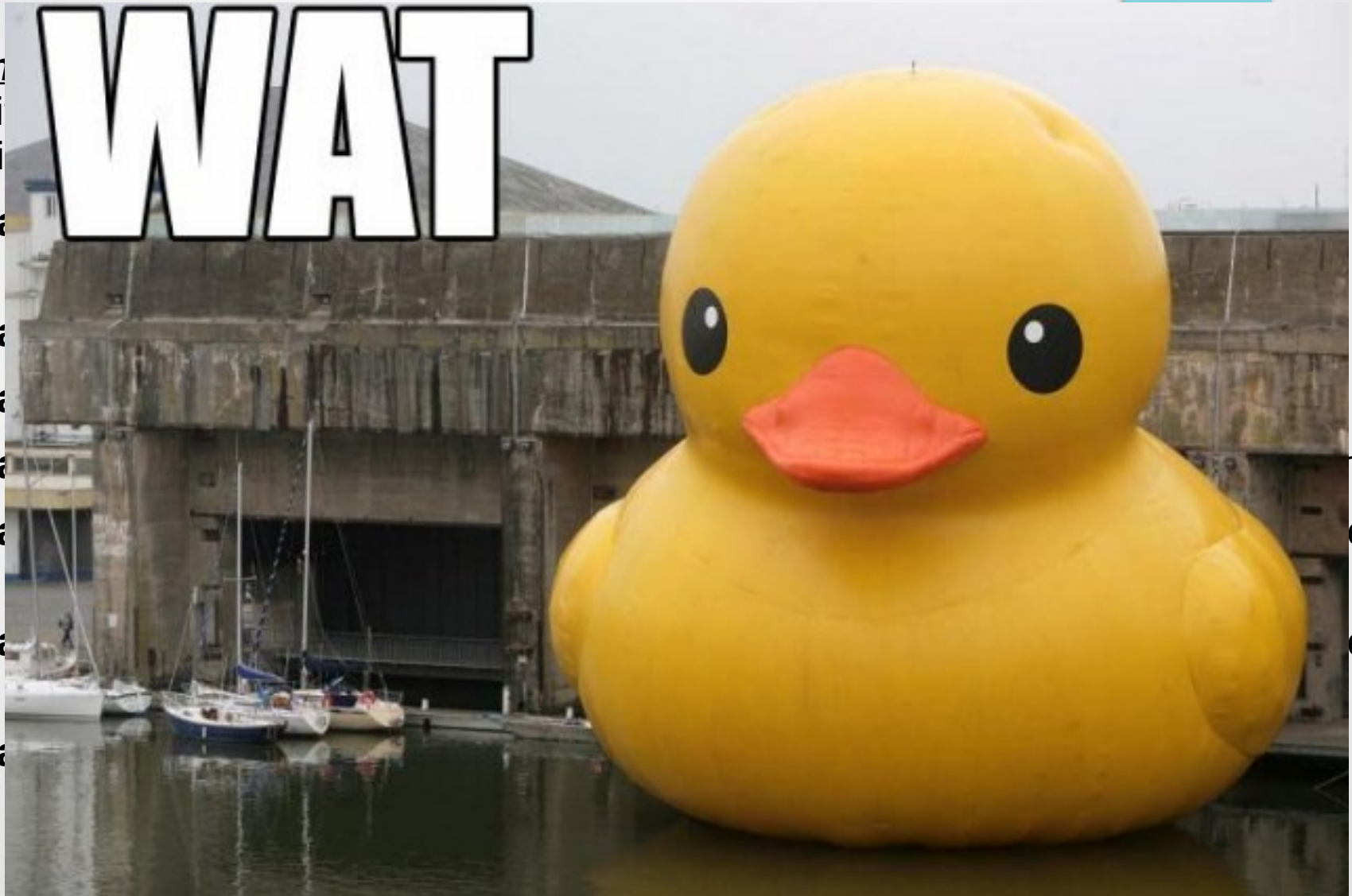
Invaria

Invaria

Invaria

Invaria

WAT



PerspikuOS – Entry & Exit

```
entry:
    pushfq                Save current flags
    cli                  Disable interrupts
    mov %rax, -8(%rsp)    Spill regs for temps
    mov %rcx, -16(%rsp)
    mov %rsp, %rcx        Save stack ptr in rcx
    mov %cr0, %rax        Get current CR0 value
    and ~CR0_WP, %rax     Clear WP bit in copy
    mov %rax, %cr0        Write back to CR0
    cli                  Disable interrupts
    mov PerCPUSecureStack, %rsp Switch to secure stack
    push %rcx             Save orig stack ptr
    mov -0x8(%rcx), %rax  Restore spilled regs
    mov -0x10(%rcx), %rcx
```

```
exit:
    mov 0(%rsp), %rsp     Restore orig stack ptr
    push %rax             Spill scratch reg
    mov %cr0, %rax        Get current CR0 value
1:
    or CR0_WP, %rax       Set WP in CR0 copy
    mov %rax, %cr0        Write back to CR0
    test CR0_WP, %eax     Ensure WP set
    je 1b                 If not, loop back
    pop %rax              Restore clobbered reg
    popfq                 Restore flags
                          (incl interrupt status)
```

PerspikuOS - Evaluation

Example Mediation Functions:

- Write-only data
- Append-only data
- Write logging

Privilege Boundary	Time (μ secs)	Time / NK Call
NK Call	0.1390	1.00x
Syscall	0.08757	0.62x
VMCALL	0.5130	3.69x

Table 3. Privilege Boundary Crossing Costs.

PerspikuOS – Evaluation (2)

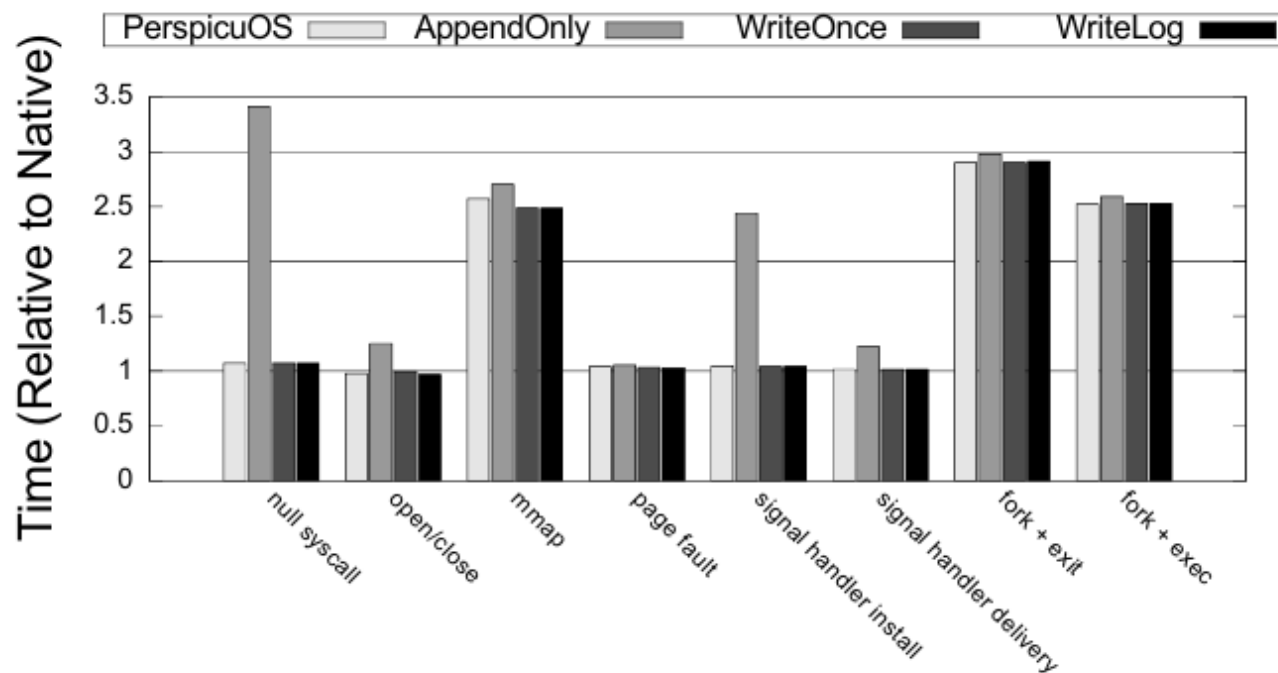


Figure 4. LMBench results.

PerspikuOS – Evaluation (3)

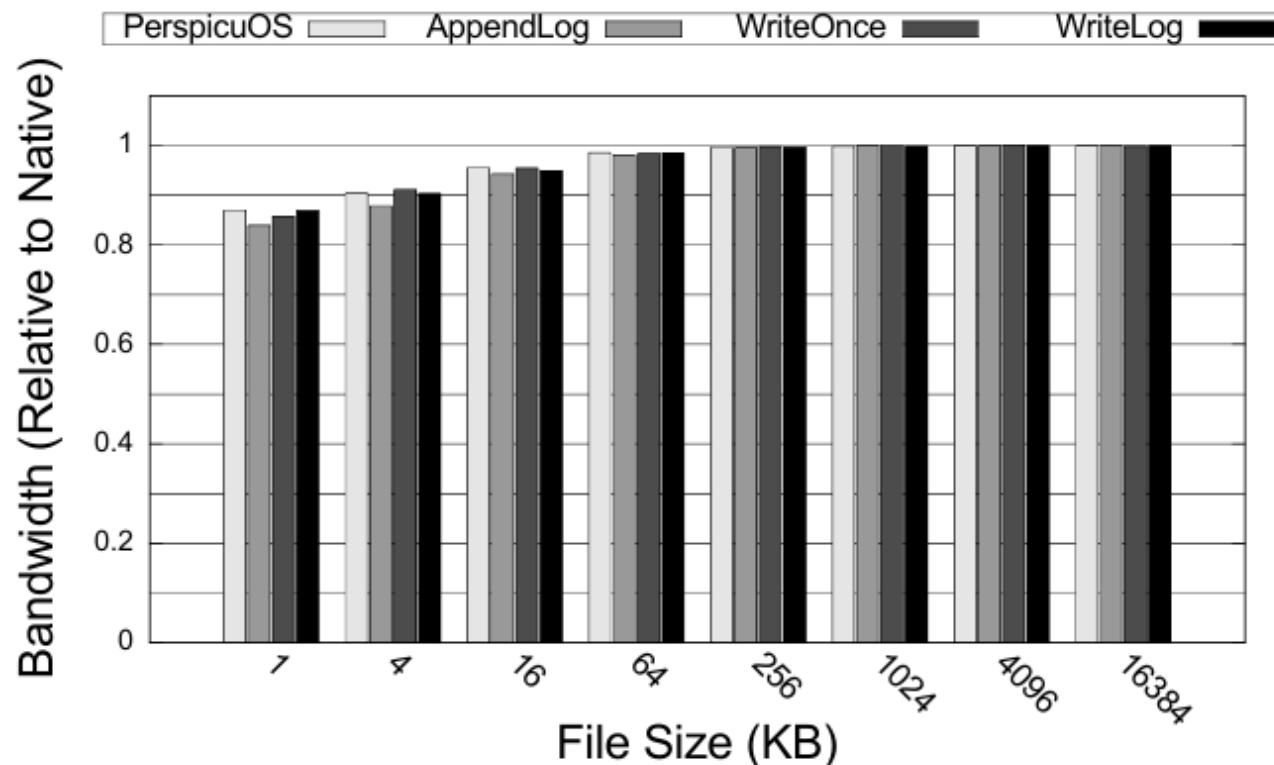


Figure 5. SSHD Average Bandwidth.

PerspikuOS – Evaluation (4)

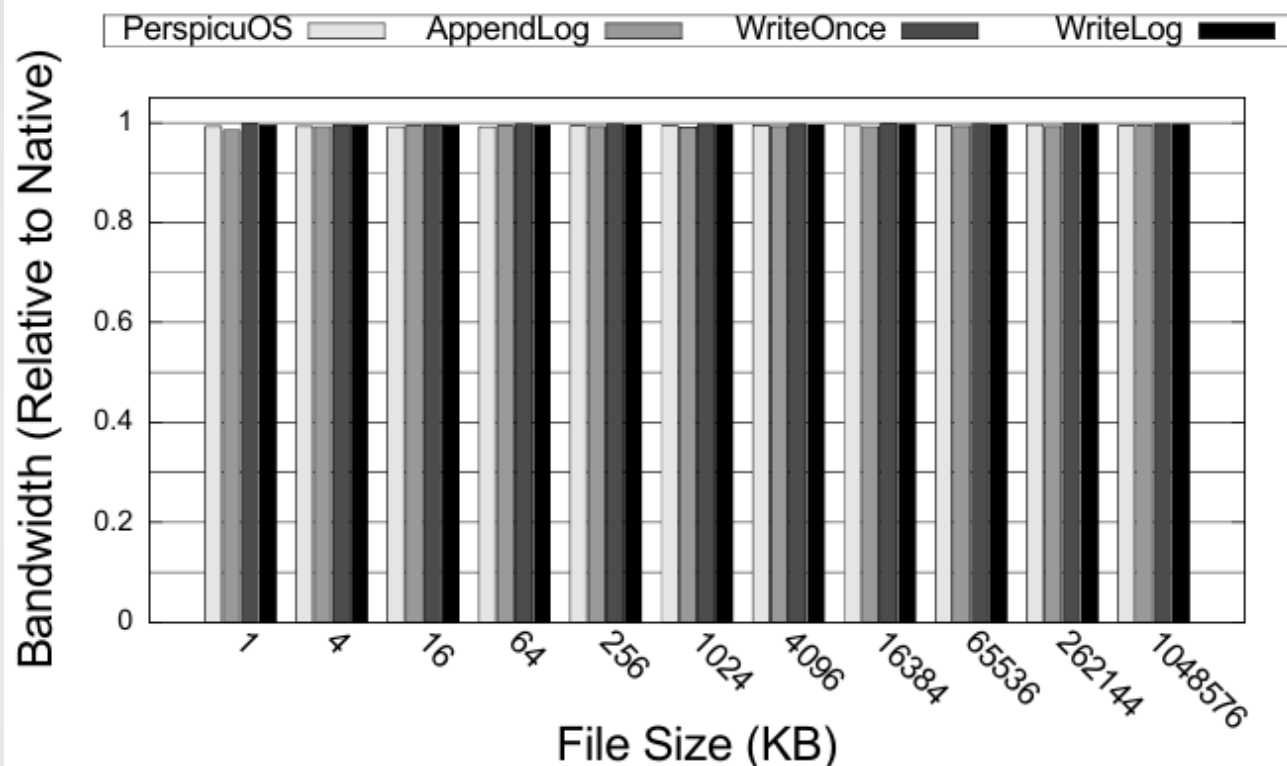


Figure 6. Apache average bandwidth.

PerspikuOS	AppendOnly	WriteOnce	WriteLog
2.6%	3.0%	2.6%	2.7%

Table 4. Kernel Build Overhead over Native.

Conclusion

- Nested kernel architecture
- Using the MMU to protect the MMU
- Write mediation policies for memory
- Protection domains without using VMM/Processes
 - Based on addresses

Discussion

Pro:

- Great if you want to log write accesses
- If the intention is rootkit detection this might be nice

Con:

- Not really suited for isolation of components
- Can't you still attack e.g. communication?
- The overhead evaluation ... sucks