

# CAPSICUM

Practical capabilities for UNIX

19<sup>th</sup> USENIX Security Symposium  
11 August 2010 - Washington, DC

Robert N. M. Watson  
Jonathan Anderson  
Ben Laurie  
Kris Kennaway

Google UK Ltd  
FreeBSD Project  
University of Cambridge



# Introduction

- Capsicum: hybrid UNIX/capability operating system
  - Requirements of complex, security-aware applications
  - Why MAC isn't quite what we want
  - Capsicum's *Capability Mode* and *Capabilities*
  - Interactions between applications and sandboxing
- Building on Capsicum

# Paradigm shift

... change is coming here

- Multi-user machines ➔ multi-machine users
- “Applications” frame competing interests
- Thin client one point of confluence
- DAC/MAC-centric access control ➔ sandboxing
- Application security rather than OS security
- Primitives for mapping distributed to local security domains

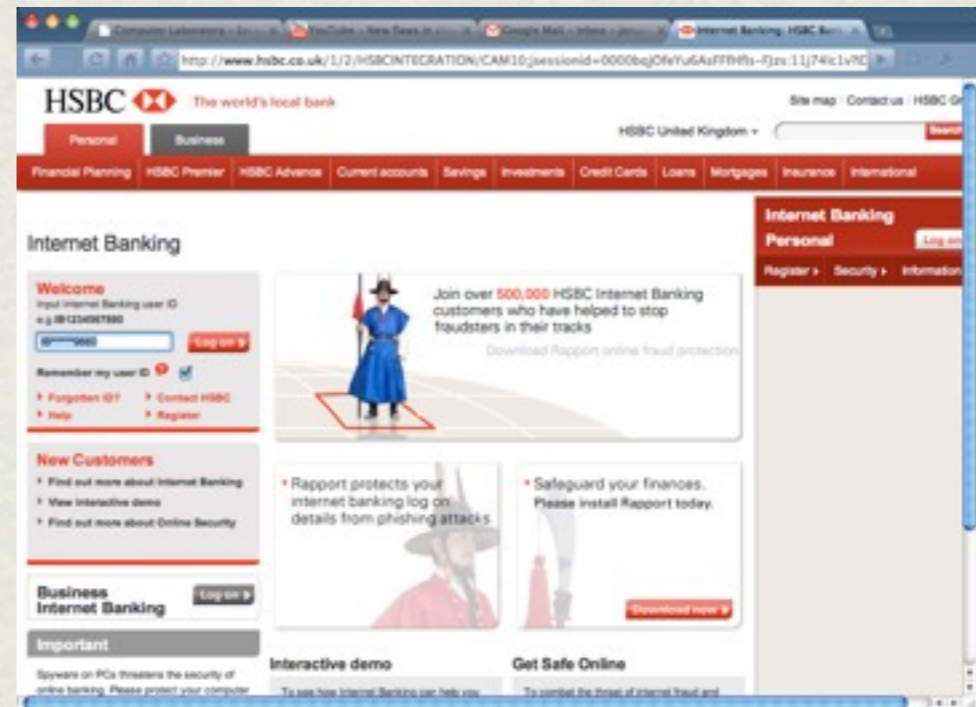
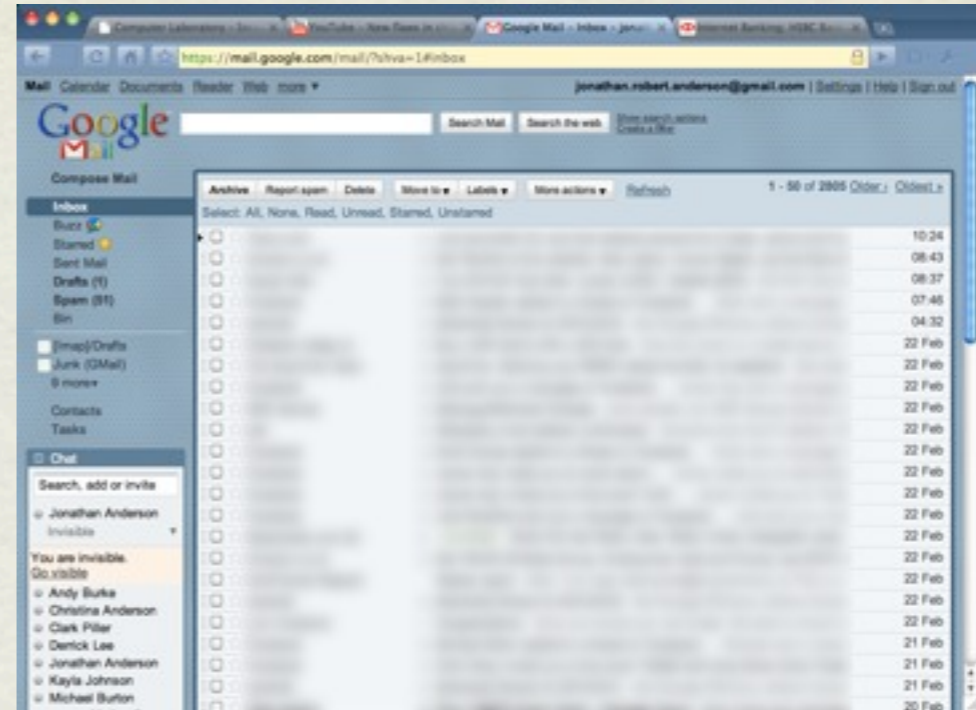
Upcoming Deadlines

Deadlines	Event	Date(s)	Location	
<b>Papers</b>				
2010-02-18	PETS	21 - 23 Jul 2010	Berlin, DE	
2010-02-18	CRYPTO	15 - 19 Aug 2010	Santa Barbara, CA, US	
2010-02-22	WEIS	7 - 8 Jun 2010	Cambridge, MA, US	
2010-02-25	USENIX-LBET	27 Apr 2010	San Jose, CA, US	
2010-02-18	WOSN	22 Jun 2010	Boston, MA, US	
2010-02-25				
2010-02-26	SECSI	26 - 27 Apr 2010	Cologne, DE	
2010-03-01	CHES	18 - 20 Aug 2010	Santa Barbara, CA, US	
2010-03-05	2010-05-28	SOUPS	14 - 16 Jul 2010	Redmond, WA, US
2010-02-18	2010-03-15	USENIX-OTPS	27 Apr 2010	San Jose, CA, US
2010-03-19	BCS-HCI	6 - 10 Sep 2010	Dundee, Scotland, UK	
2010-03-26	CEAS	13 - 14 Jul 2010	Redmond, WA, US	
2010-04-01	ESORICS	20 - 22 Sep 2010	Athens, GR	
2010-04-05	ASA	21 Jul 2010	Edinburgh, Scotland, UK	
2010-04-16	NSPW	21 - 23 Sep 2010	Concord, MA, US	
2010-04-17	ACH-CCS	4 - 8 Oct 2010	Chicago, IL, US	
2010-02-22	2010-05-18	ACH-SIGCOMM	30 Aug - 3 Sep 2010	New Delhi, IN

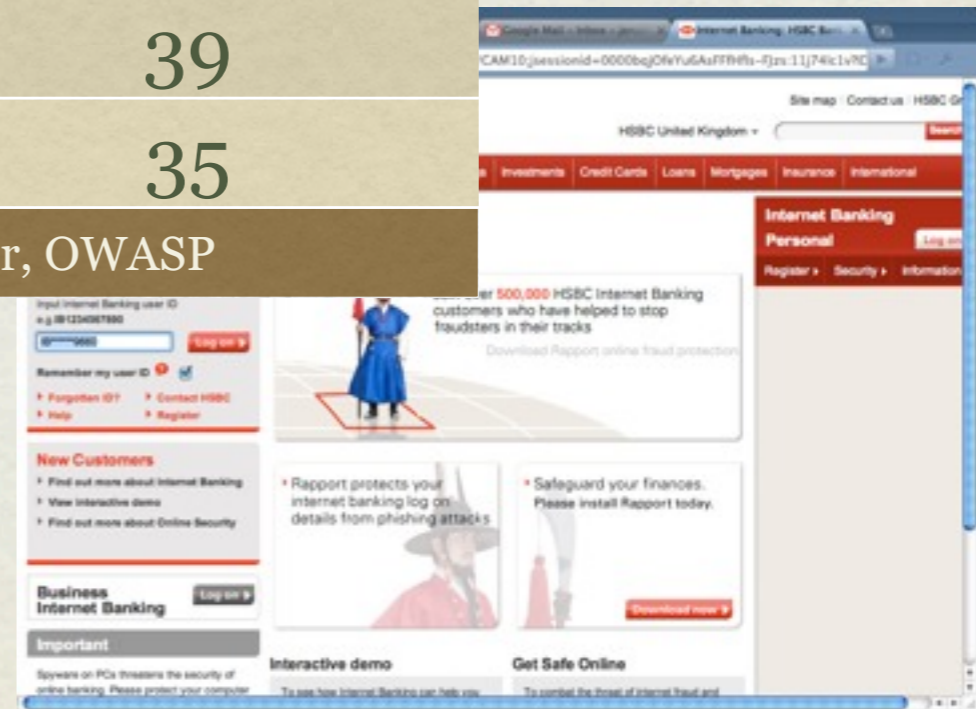
Upcoming Conferences

Event	Date(s)	Location	URLs
NOSS	28 Feb - 3 Mar 2010	San Diego, CA, US	permalink
CT-ISA	1 - 5 Mar 2010	San Francisco, CA, US	permalink
ACH-SAC	22 - 26 Mar 2010	Leuven, CH	permalink



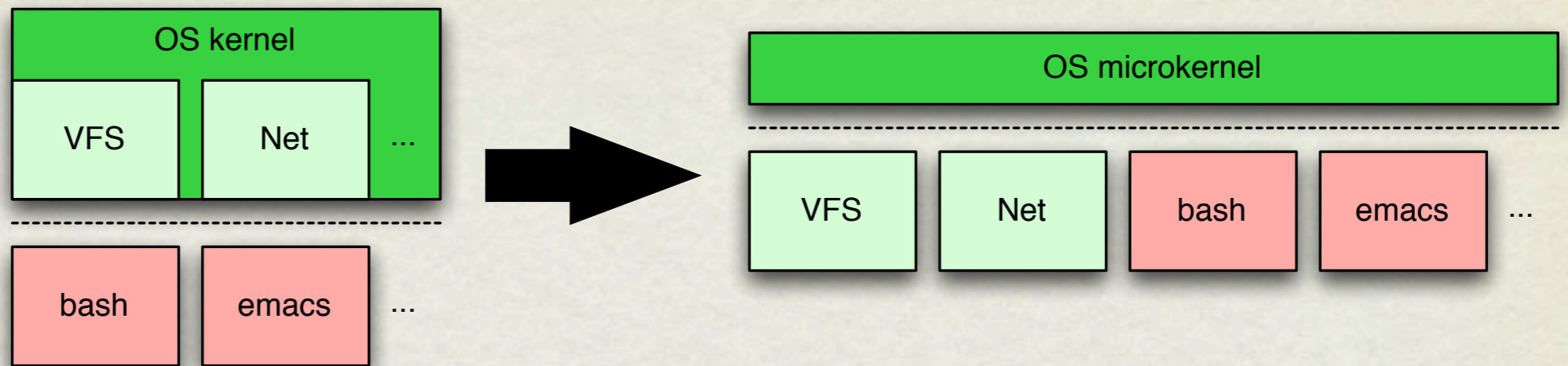


CVEs in Jan-Aug 2009	
Firefox	85
Safari	59
IE	48
Chrome	39
Flash	35
source; Justin Foster, OWASP	



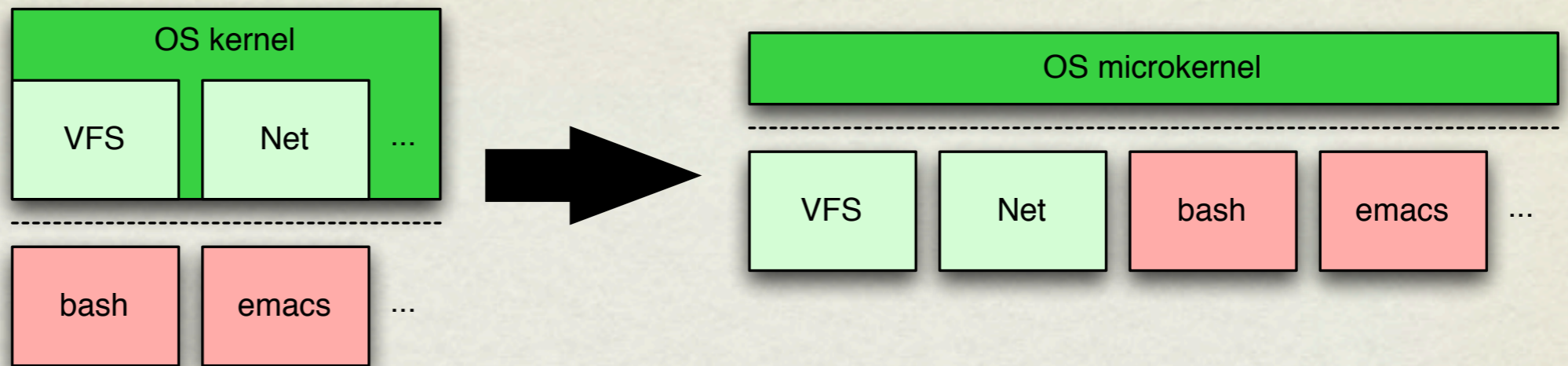
# Microkernels to compartmentalisation

1980's

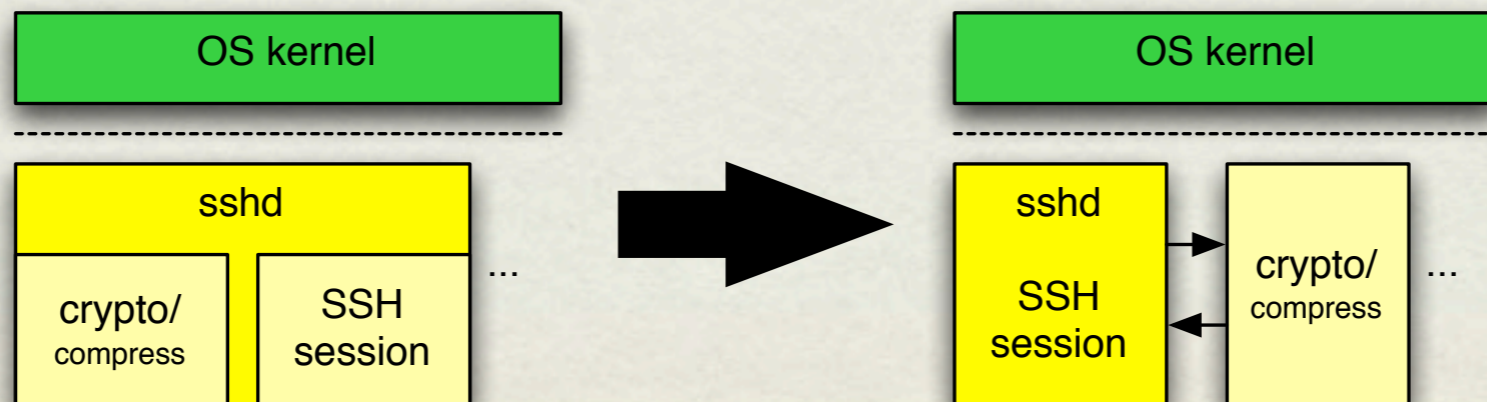


# Microkernels to compartmentalisation

1980's



2000's



# What about MAC?

	Type Enforcement (TE)	What we need
Interests of	Administrator	User or application
Sandbox creation	Administrator modifies global policy	On demand without using privilege
Policy source	Access control rules in global policy files	Embedded in applications, from UI



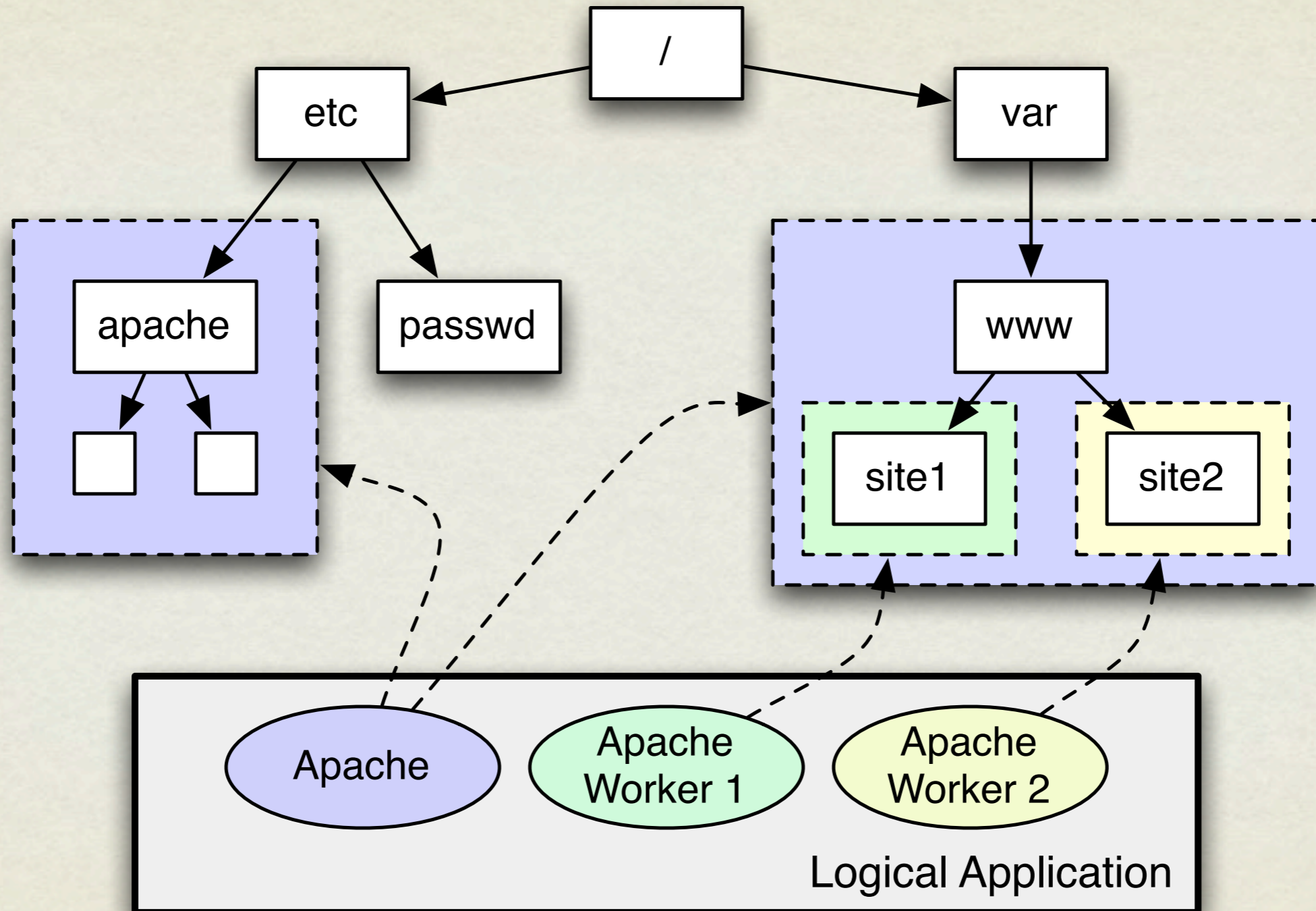
# What about MAC?

	Type Enforcement (TE)	What we need
Interests of	Administrator	User or application
Sandbox creation	Administrator modifies global policy	On demand without using privilege
Policy source	Access control rules in global policy files	Embedded in applications, from UI

# What about MAC?

	Type Enforcement (TE)	What we need
Interests of	Administrator	User or application
Sandbox creation	Administrator modifies global policy	On demand without using privilege
Policy source	Access control rules in global policy files	Embedded in applications, from UI

# Application-driven rights delegation



# Capability systems



*A capability* is an unforgeable token of authority.

Supports delegation-centric access control.

# Where to start?

## Production monolithic systems

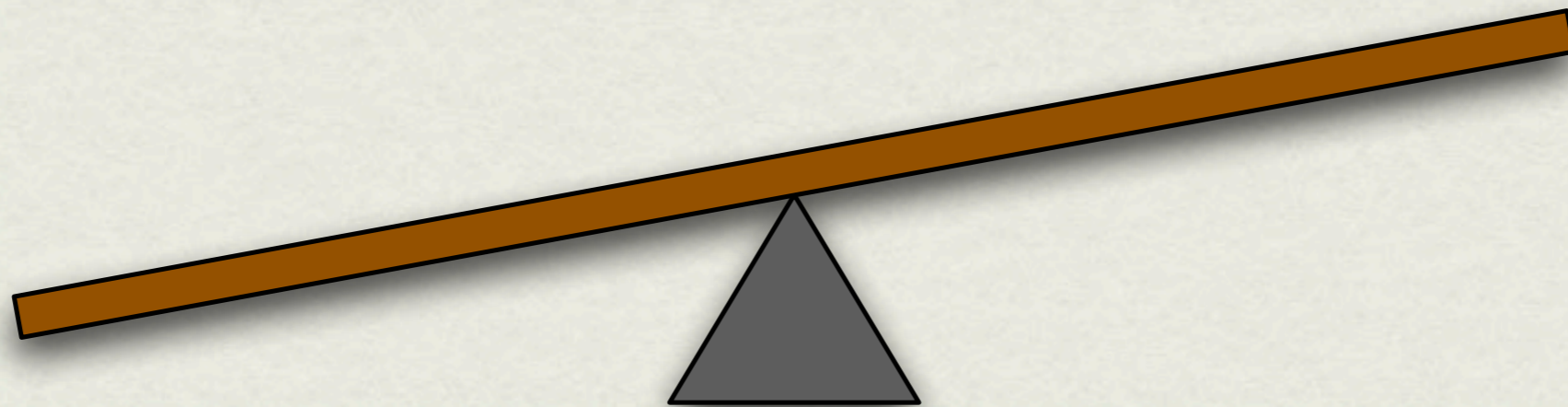
UNIX, Linux, Windows, Mac OS X

- ✗ Monolithic kernel security model
- ✓ Real application stack today

## Research capability systems

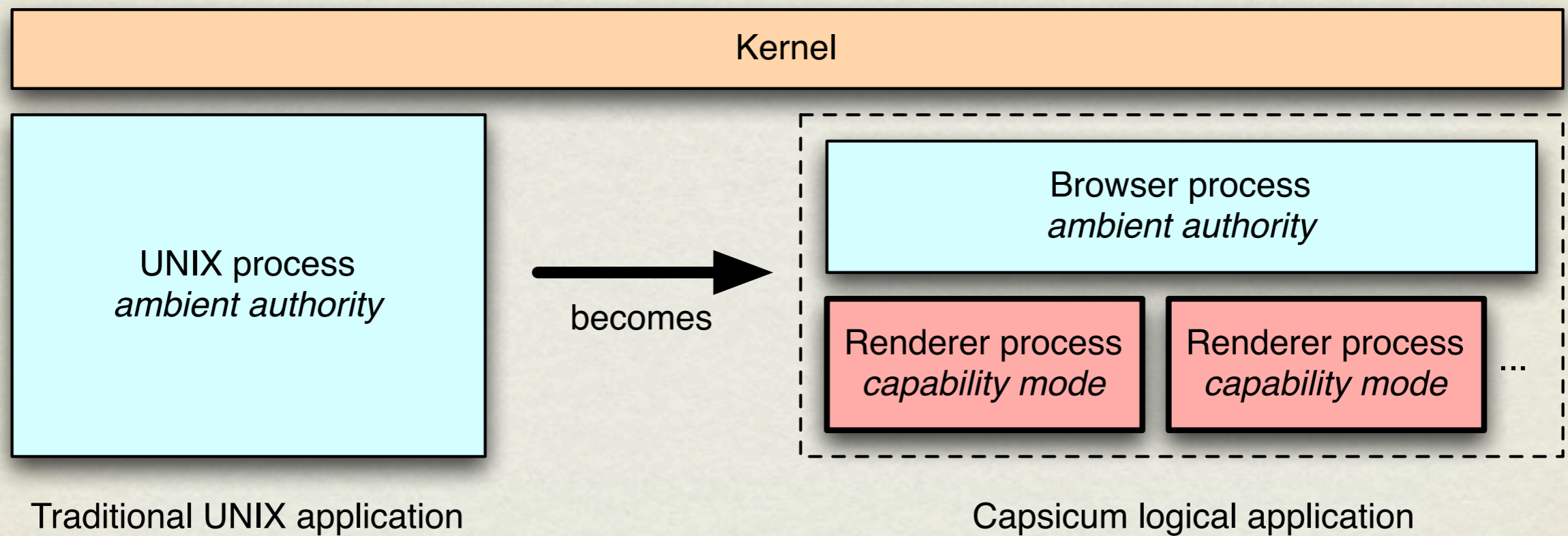
EROS (CAPROS), CoyoteOS

- ✓ Least privilege design
- ✗ No extant application stack



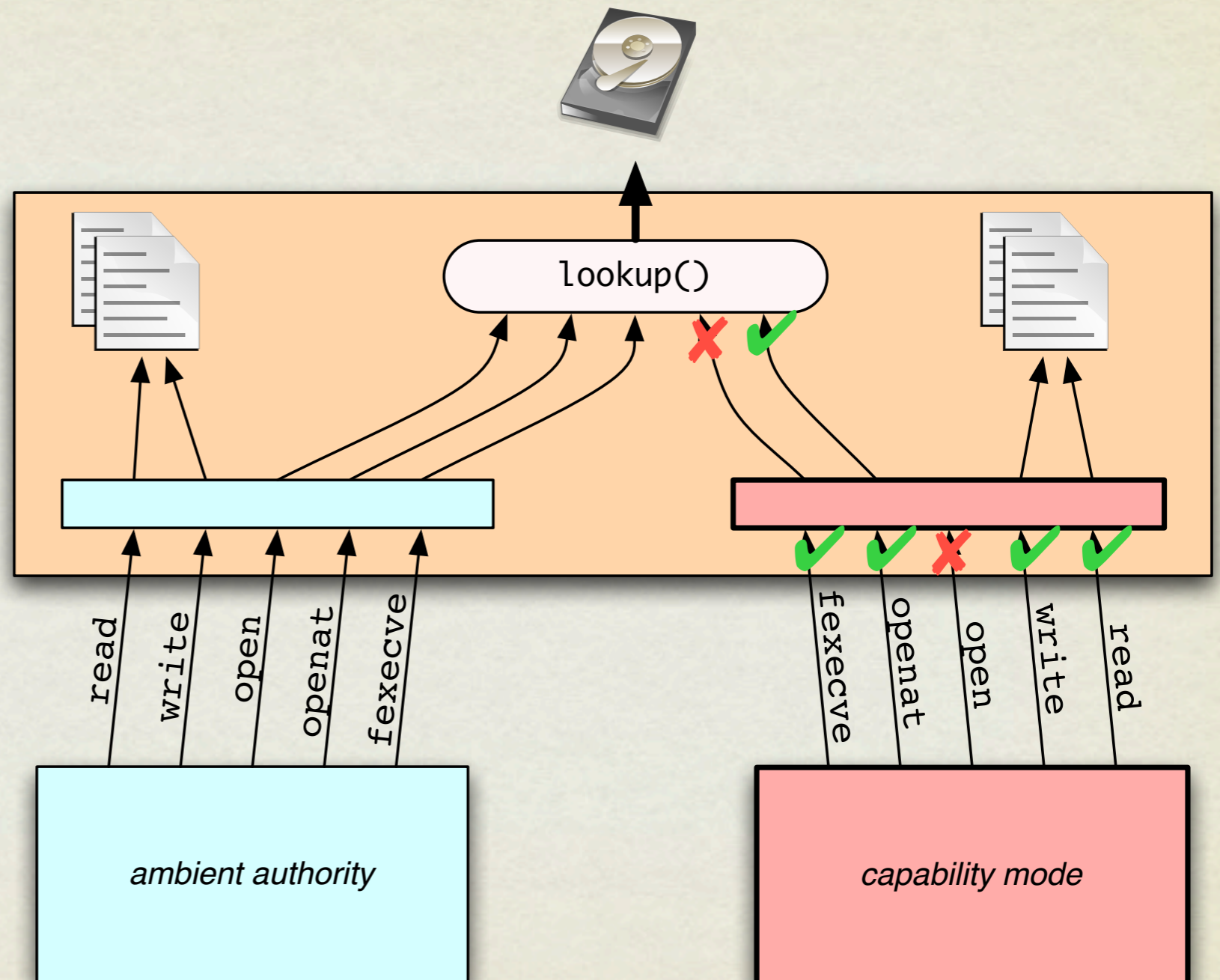
Hybrid approach: **immediate security benefits** with a **long-term capability system vision**

# Logical applications in Capsicum

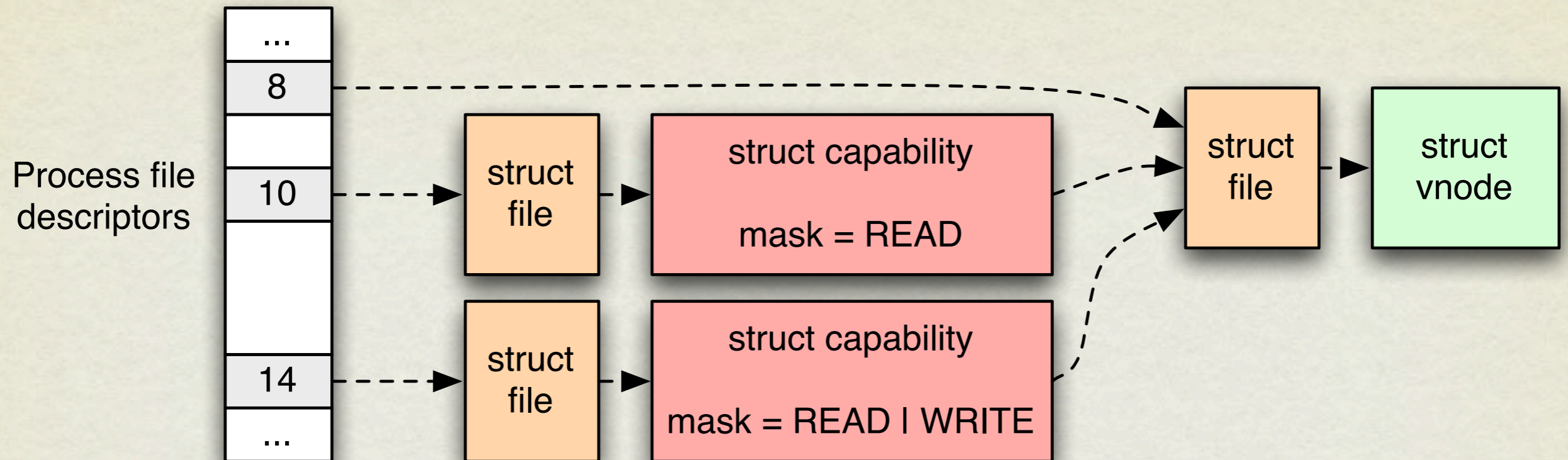


# Capability mode

- New system call `cap_enter` sets inherited credential flag
- Global OS name spaces restricted: only delegated rights available
- Interface thinning and other constraints on system calls



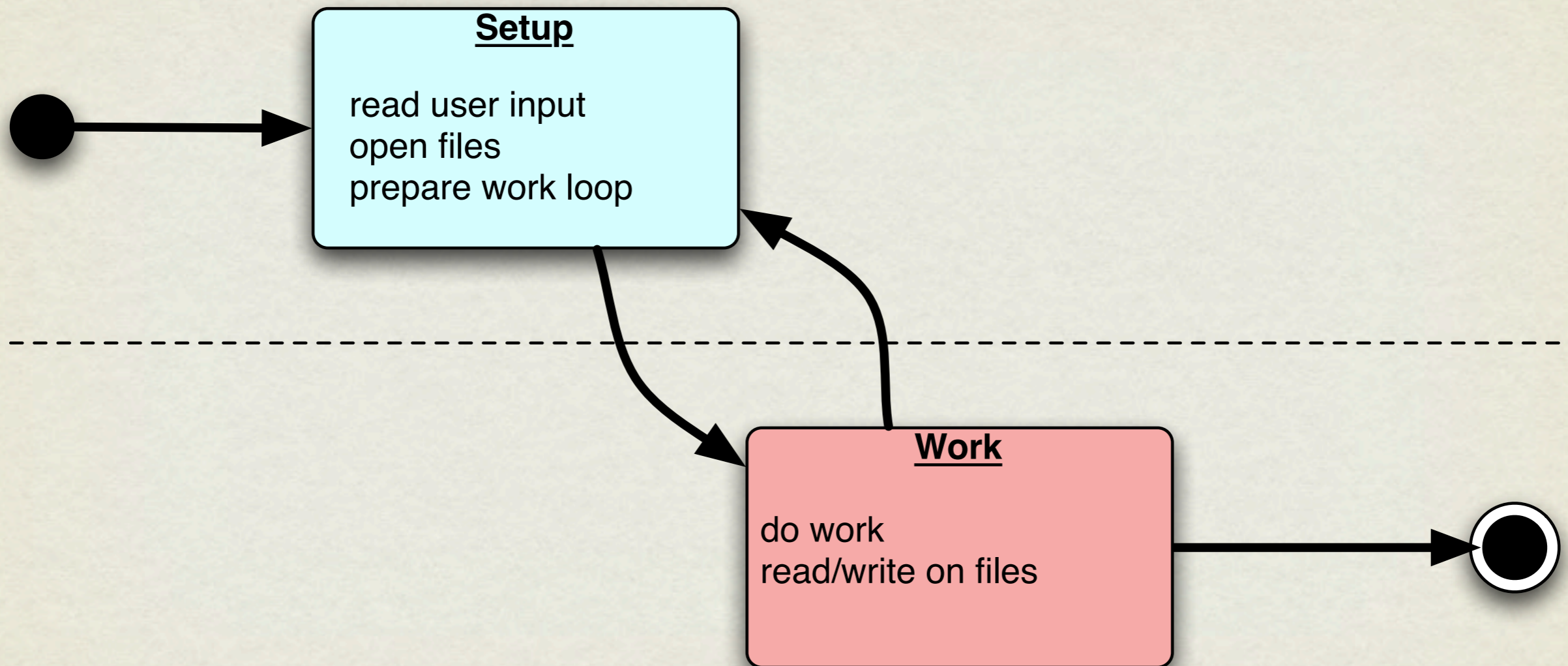
# Capabilities



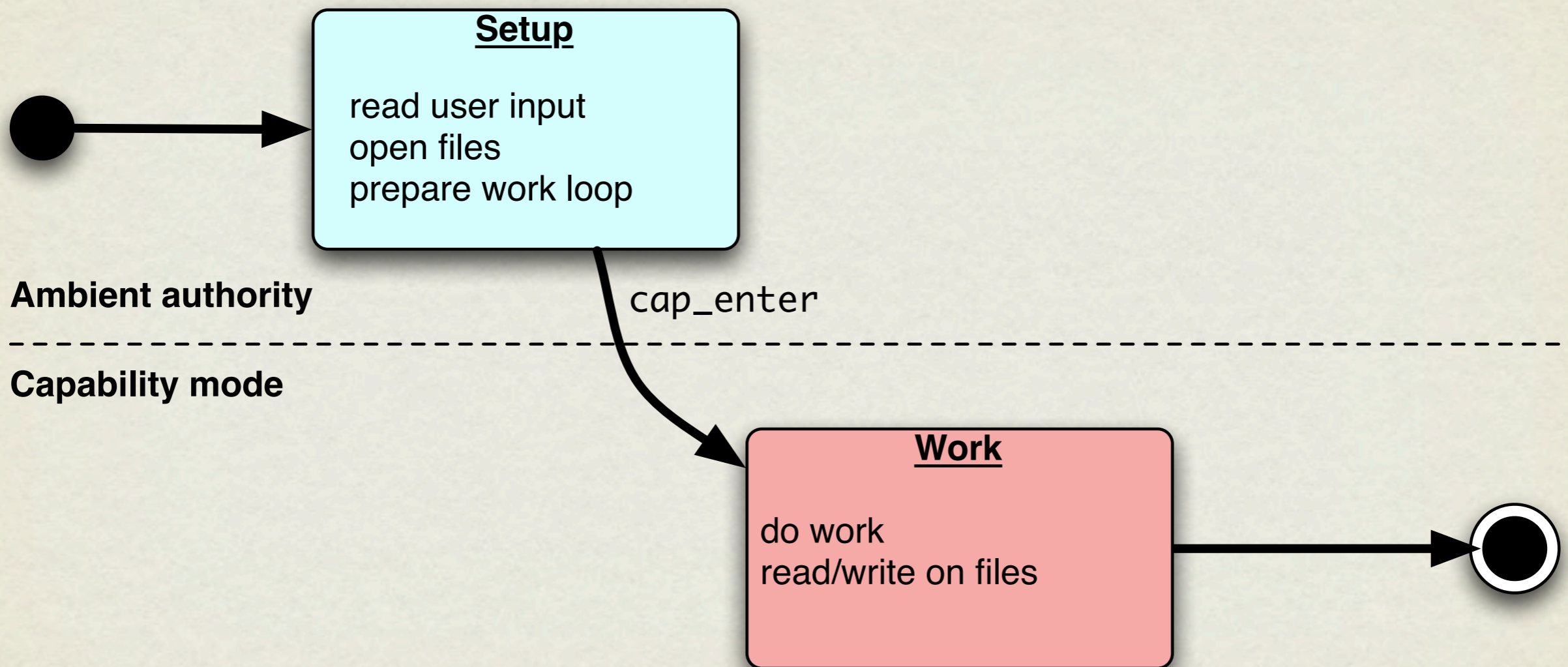
- Capabilities refine *open flags* on file descriptors
- `cap_new` on a capability further restricts access; no chains
- Inherited across `fork/exec` or passed via sockets
- Directory capabilities allow subtree delegation



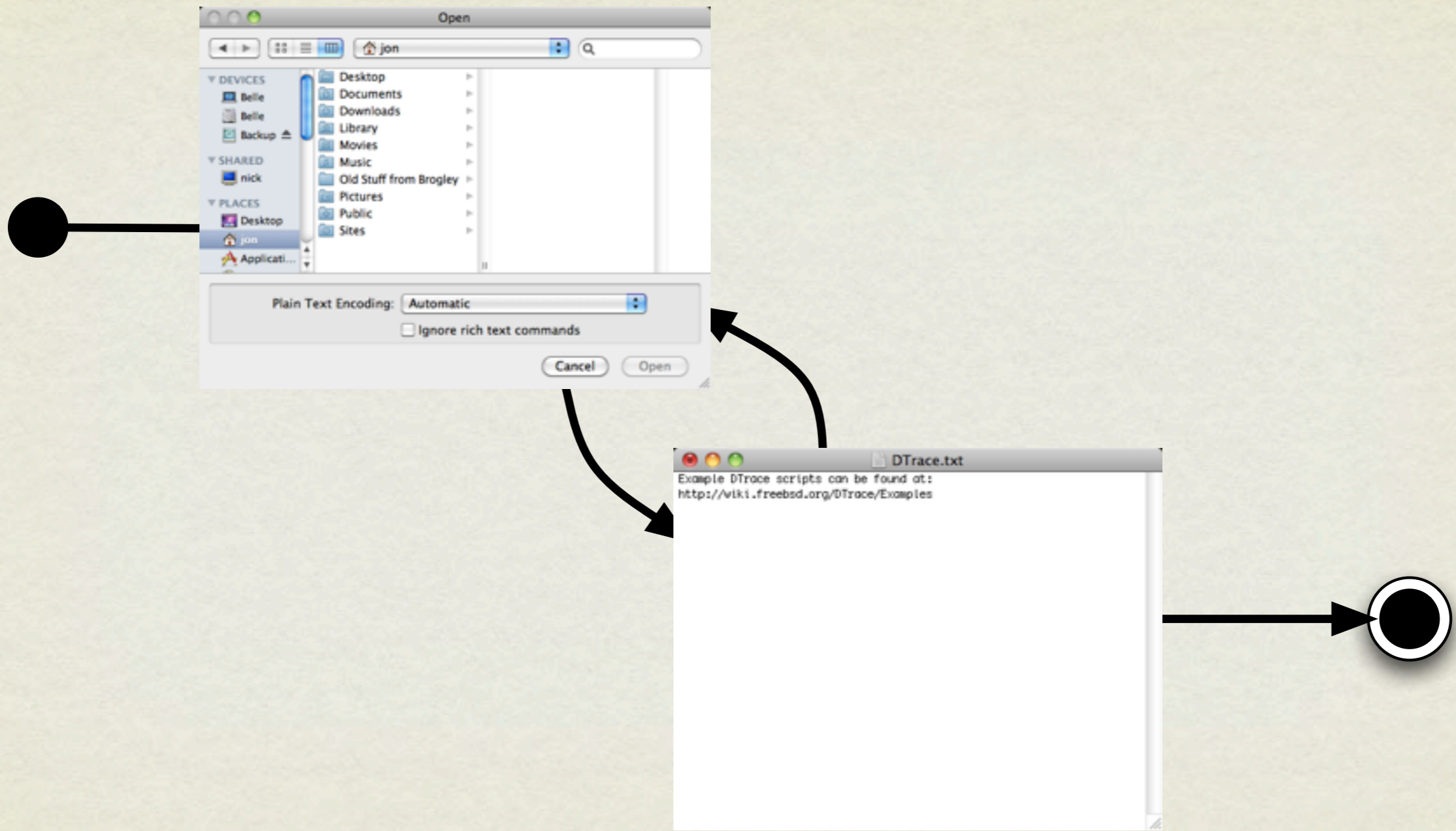
# Possible application



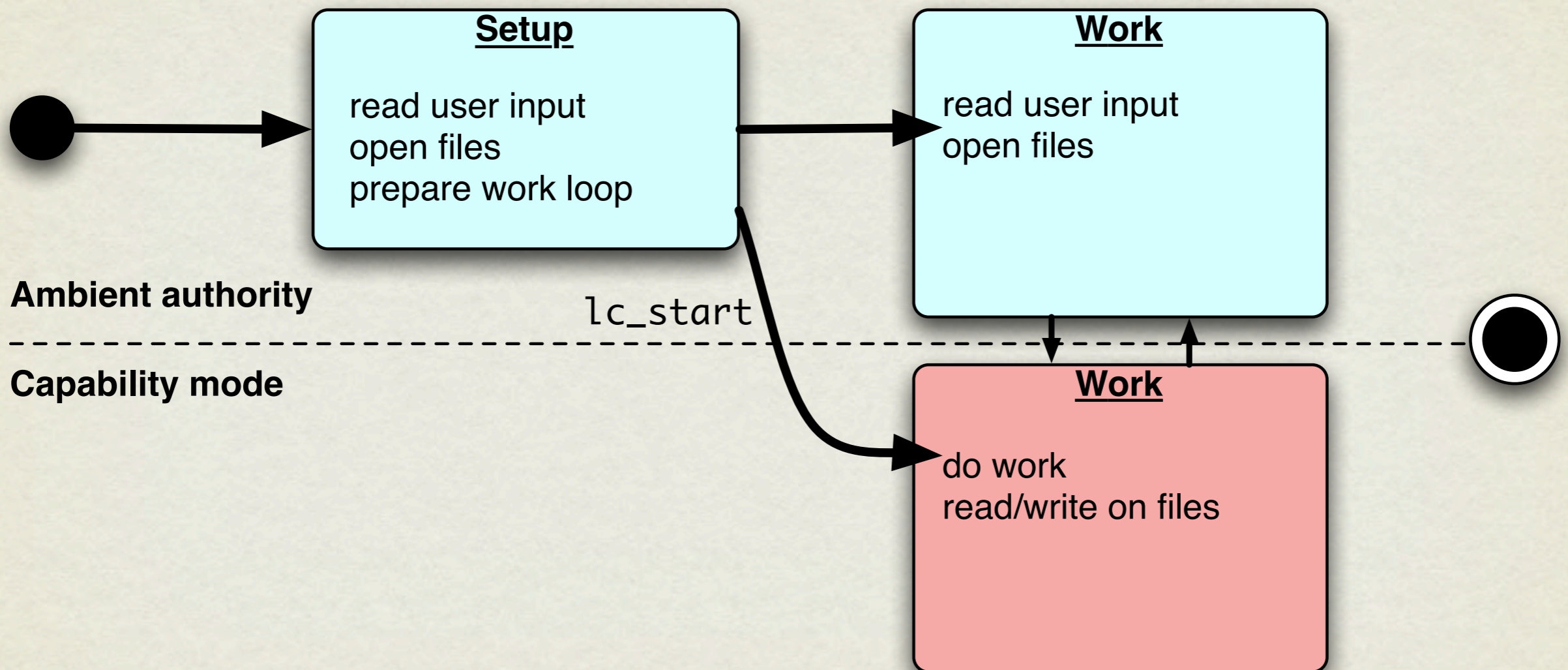
# System call API



# Interactive applications



# libcapsicum API



# Adapted applications

Program	Approach	Changes
tcpdump	<code>cap_enter</code>	Enter for parse/render work loop
dhclient	<code>cap_enter</code>	Reinforce existing <code>chroot/setuid</code> privilege separation
gzip	<code>libcapsicum</code>	Open files with ambient authority, pass capabilities to sandbox
Chromium	<code>cap_enter</code>	Sandbox Javascript and HTML processing in renderer processes

# tcpdump

```
@@ -1197,6 +1199,14 @@
        (void) fflush(stderr);
    }
#endif /* WIN32 */
+   if (lc_limitfd(STDIN_FILENO, CAP_FSTAT) < 0)
+       error("lc_limitfd: unable to limit STDIN_FILENO");
+   if (lc_limitfd(STDOUT_FILENO, CAP_FSTAT | CAP_SEEK | CAP_WRITE) < 0)
+       error("lc_limitfd: unable to limit STDIN_FILENO");
+   if (lc_limitfd(STDERR_FILENO, CAP_FSTAT | CAP_SEEK | CAP_WRITE) < 0)
+       error("lc_limitfd: unable to limit STDERR_FILENO");
+   if (cap_enter() < 0)
+       error("cap_enter: %s", pcap_strerror(errno));
status = pcap_loop(pd, cnt, callback, pcap_userdata);
if (WFileName == NULL) {
```

# Chromium sandboxing

	OS	Sandbox	LoC	FS	IPC	NET	S≠S'	Priv
DAC	Windows	DAC ACLs	22,350	⚠	⚠	✗	✗	✓
	Linux	chroot()	600	✓	✗	✗	✓	✗
MAC	Mac OS X	Sandbox	560	✓	⚠	✓	✓	✓
	Linux	SELinux	200	✓	⚠	✓	✗	✗
Cap	Linux	seccomp	11,300	⚠	✓	✓	✓	✓
	FreeBSD	Capsicum	100	✓	✓	✓	✓	✓

# Chromium sandboxing

	OS	Sandbox	LoC	FS	IPC	NET	S≠S'	Priv
DAC	Windows	DAC ACLs	<b>22,350</b>	⚠	⚠	✗	✗	✓
	Linux	chroot()	600	✓	✗	✗	✓	✗
MAC	Mac OS X	Sandbox	560	✓	⚠	✓	✓	✓
	Linux	SELinux	200	✓	⚠	✓	✗	✗
Cap	Linux	seccomp	<b>11,300</b>	⚠	✓	✓	✓	✓
	FreeBSD	Capsicum	100	✓	✓	✓	✓	✓



# Chromium sandboxing

	OS	Sandbox	LoC	FS	IPC	NET	S≠S'	Priv
DAC	Windows	DAC ACLs	22,350	⚠	⚠	✗	✗	✓
	Linux	chroot()	600	✓	✗	✗	✓	✗
MAC	Mac OS X	Sandbox	560	✓	⚠	✓	✓	✓
	Linux	SELinux	200	✓	⚠	✓	✗	✗
Cap	Linux	seccomp	11,300	⚠	✓	✓	✓	✓
	FreeBSD	Capsicum	100	✓	✓	✓	✓	✓

# Chromium sandboxing

	OS	Sandbox	LoC	FS	IPC	NET	S≠S'	Priv
DAC	Windows	DAC ACLs	22,350	⚠	⚠	✗	✗	✓
	Linux	chroot()	600	✓	✗	✗	✓	✗
MAC	Mac OS X	Sandbox	560	✓	⚠	✓	✓	✓
	Linux	SELinux	200	✓	⚠	✓	✗	✗
Cap	Linux	seccomp	11,300	⚠	✓	✓	✓	✓
	FreeBSD	Capsicum	100	✓	✓	✓	✓	✓

# Chromium sandboxing

	OS	Sandbox	LoC	FS	IPC	NET	S≠S'	Priv
DAC	Windows	DAC ACLs	22,350	⚠	⚠	✗	✗	✓
	Linux	chroot()	600	✓	✗	✗	✓	✗
MAC	Mac OS X	Sandbox	560	✓	⚠	✓	✓	✓
	Linux	SELinux	200	✓	⚠	✓	✗	✗
Cap	Linux	seccomp	11,300	⚠	✓	✓	✓	✓
	FreeBSD	Capsicum	100	✓	✓	✓	✓	✓

# Chromium sandboxing

	OS	Sandbox	LoC	FS	IPC	NET	S≠S'	Priv
DAC	Windows	DAC ACLs	22,350	⚠	⚠	✗	✗	✓
	Linux	chroot()	600	✓	✗	✗	✓	✗
MAC	Mac OS X	Sandbox	560	✓	⚠	✓	✓	✓
	Linux	SELinux	200	✓	⚠	✓	✗	✗
Cap	Linux	seccomp	11,300	⚠	✓	✓	✓	✓
	FreeBSD	Capsicum	100	✓	✓	✓	✓	✓

# Building on Capsicum

- Assisted compartmentalisation (static, dynamic analysis)
- Critical network services: routing daemon, etc.
- Monolithic applications: OpenOffice.org, KDE..
- Distributed domains ➔ local domains: browsers, databases...
- Gesture-Based Access Control (GBAC)
  - Power boxes, “Drag and drop” ➔ assign capabilities

# Conclusion

- Multi-user security ➔ compartmentalised applications
- Capsicum APIs **faster, cleaner, and more secure**
  - Delegation-centric approach to granular policy
  - Avoid policy dual-coding, no privilege requirement
- **Supplement** rather than replace DAC and MAC
- API/semantics + prototype on FreeBSD 9.x, 8.x backport
- Linux/ChromeOS port in progress at Google

# Questions?



\* <http://www.cl.cam.ac.uk/research/security/capsicum/>