

Sisteme de operare avansate

Exploiting Concurrency  
Vulnerabilities in  
System Call Wrappers

Robert N. M. Watson  
Computer Laboratory  
University of Cambridge  
[robert.watson@cl.cam.ac.uk](mailto:robert.watson@cl.cam.ac.uk)

# Interceptare apeluri de sistem

- Folosirea de wrappere pentru apeluri de sistem
  - nivel kernel
  - nivel userspace
- Verificări de securitate suplimentare
- Folosită în research și comercial (anti-virusi)
- Există probleme de securitate
  - datorită concurenței

# Concurența în SO moderne

- Fizică
  - multiprocesoare
- Logică
  - în userspace: multi-procese, multi-threading, async IO, semnale
  - în kernel: întreruperi, kernel thread-uri

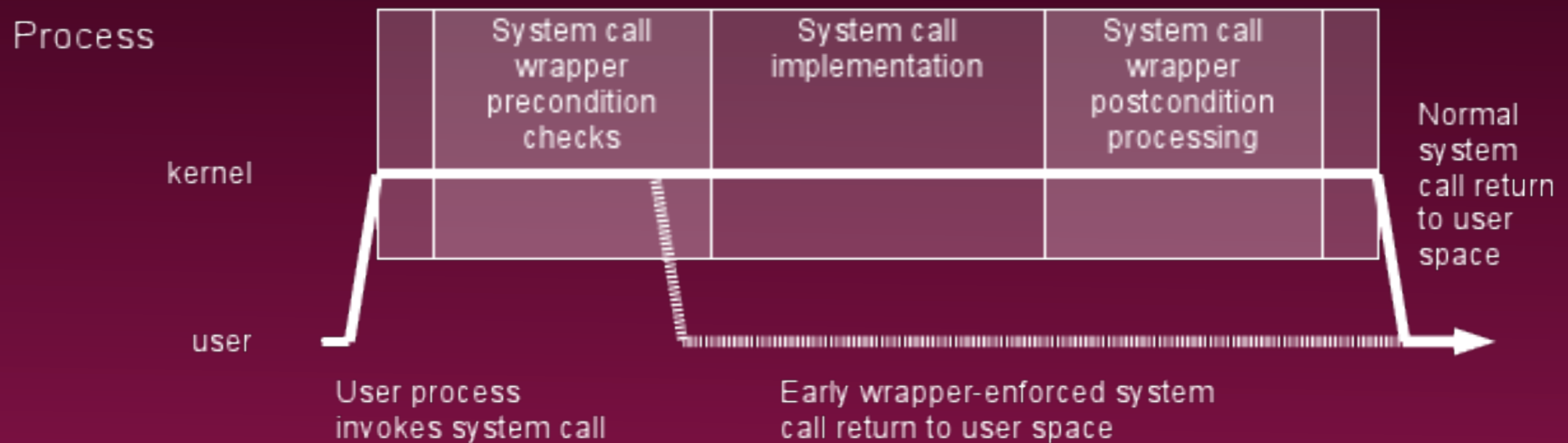
# System call wrappers

- Construite ca un reference monitor
  - tamper proof
  - always invoked
  - small enough to be tested
- System Call Wrappers
  - rezidă de obicei în kernel
  - necesită modificări minime în kernel
  - sunt invocate la fiecare apel de sistem
  - portabile accros APIs (UNIX)

# Implementare tipică în kernel

- Interceptează trap-ul apelului de system și inserează
  - verificări pre apel de sistem
    - inspectare argumente, eventual înlocuirea acestora
    - respingerea unui apel de sistem
  - verificări post apel de sistem
    - menținerea de informații pentru cazurile în care e nevoie de stare
    - logarea rezultatului apelului de sistem
    - modificarea rezultatului

# System call wrappers [2]



# Politica de securitate

- Verificarea se face în concordanța cu o politică de securitate
- Definită de utilizator
- În general statică
- Politica poate fi:
  - compilată în modulul de kernel
  - definită și încărcată din userspace

# Vulnerabilități datorate concurenței

- Politica de securitate nu este respectată datorită unor probleme de concurență
- Efecte:
  - denial of service
  - leaking of sensitive data
  - incorrect access control



# Atacuri asupra bazate pe concurență asupra SCW-relor

- Buguri de sincronizare în SCW
- Buguri ce exploatează non-atomicitatea dintre kernel și SCW
  - nesincronizare dintre SCW și kernel la copierea argumentelor (syntactic race condition)
  - nesincronizarea dintre SCW și kernel la interpretarea argumentelor (semantic race condition)

# Syntactic race conditions

- “Portable” pe mai multe kernele sau SCW-uri
- Sunte probleme fundamentale
  - nu pot fi rezolvate fără schimbări majore în arhitectură
- Time-of-audit-to-time-of-use (TOATTOU)
  - logurile nu reflectă corect acțiunile efectuate

# Syntactic race conditions [2]

- Time-of-check-to-time-of-use (TOCTTOU)
  - operațiile de verificare nu sunt atomice cu operațiile pe care le protejează
- Time-of-replacement-to-time-of-use (TORTTOU)
  - argumentele apelurilor de sistem sunt înlocuite între verificare în SCW și folosirea lor în kernel

# Tehnici de exploatare

- Identificarea unor resurse partajate între user, SCW și kernel
- Argumentele directe sunt pasate în registri dar argumentele indirecte...
- .. sunt pasate prin Userspace memory

# Argumente indirecte

- Exemple:
  - file paths
  - adrese de socket
  - resource limits
- Aceste argumente sunt copiate de două ori
  - de SCW pentru a face verificările
  - de kernel pentru a executa apelul de sistem

# Condiții de exploatare race-uri

- Concurență
  - se poate obține cu procese multiple, thread-uri, semnale, async IO
- Accesul la memorie de către cel puțin două entități
  - pentru procese multiple: memorie partajată prin moștenire (minherit, rfork, clone)
  - thread-uri

# Exploatare race-uri pe sisteme uni-procesor

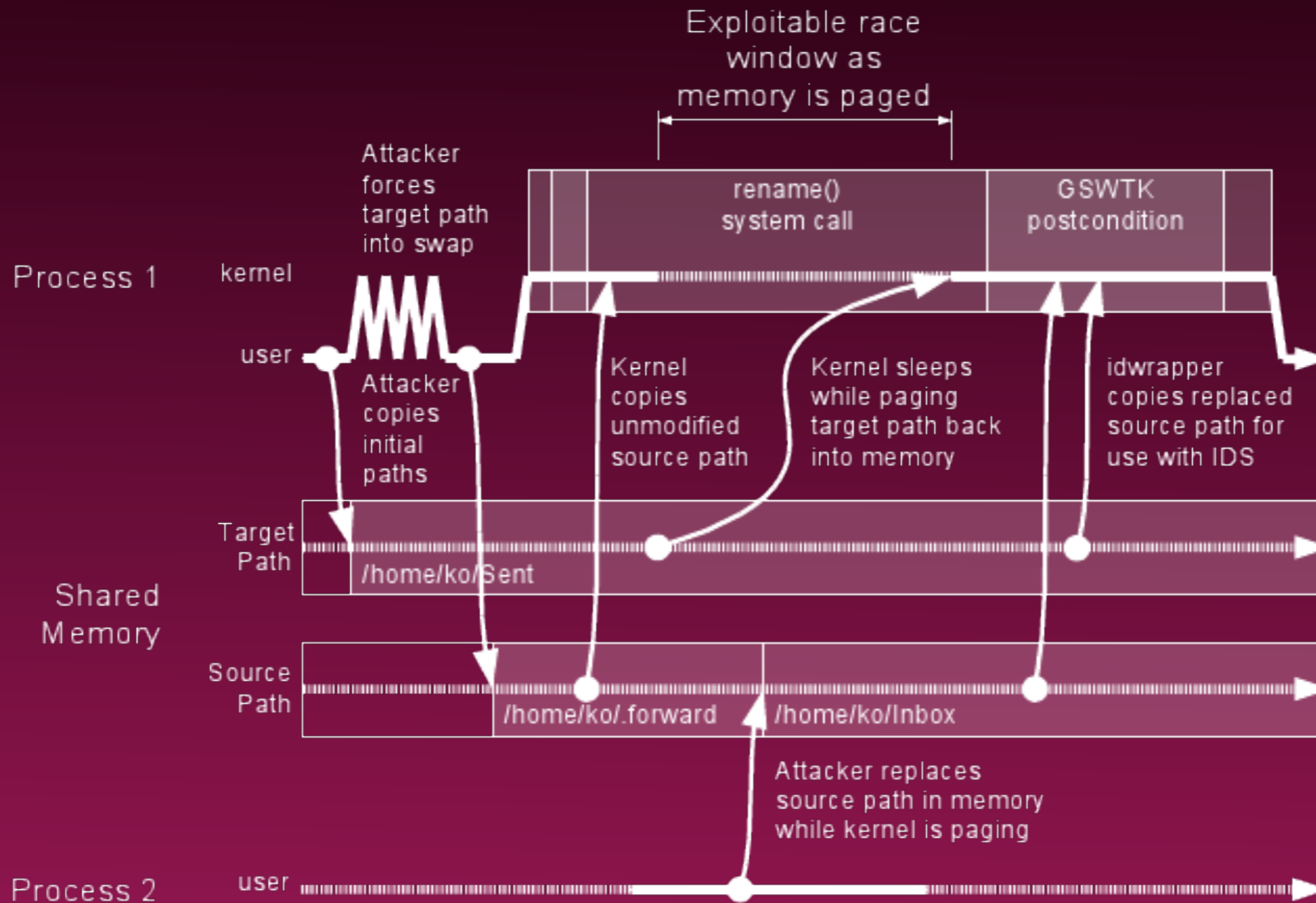
- Atacatorul trebuie sa determine kernel-ul să cedeze procesorul între pre-condition și kernel sau între kernel și post-conditions
  - voluntară: blocking IO pe un socket sau disk
  - involuntară: page fault
    - daca pagina cu argumente este în swap, deschide o fereastră de câteva ms în care se poate schimba argumentul

# Atac cu page fault-uri

- rename(from, to)
- from este in memorie, to este swapat
- SCW va accessa from și verifica accesul după care va incerca să acceseze to, se va genera un page fault și se va porni operația de swap-in
- Între timp, controlul ajunge la un helper thread care modifică from



# Atac cu page fault-uri [2]



# Cazul unui singur argument

- Argumentul se plasează între două pagini una în memorie alta în swap
- Se suprascrive doar bucata din prima pagină
  - în ideea ca deja a fost citită de SCW

# Atac cu cedare voluntară

- connect asteaptă primirea unui SYN-ACK
- În acest timp atacatorul poate schimba argumentele (de exemplu adresa la care se face conectarea)
- Post-condiția va evalua o adresa falsă

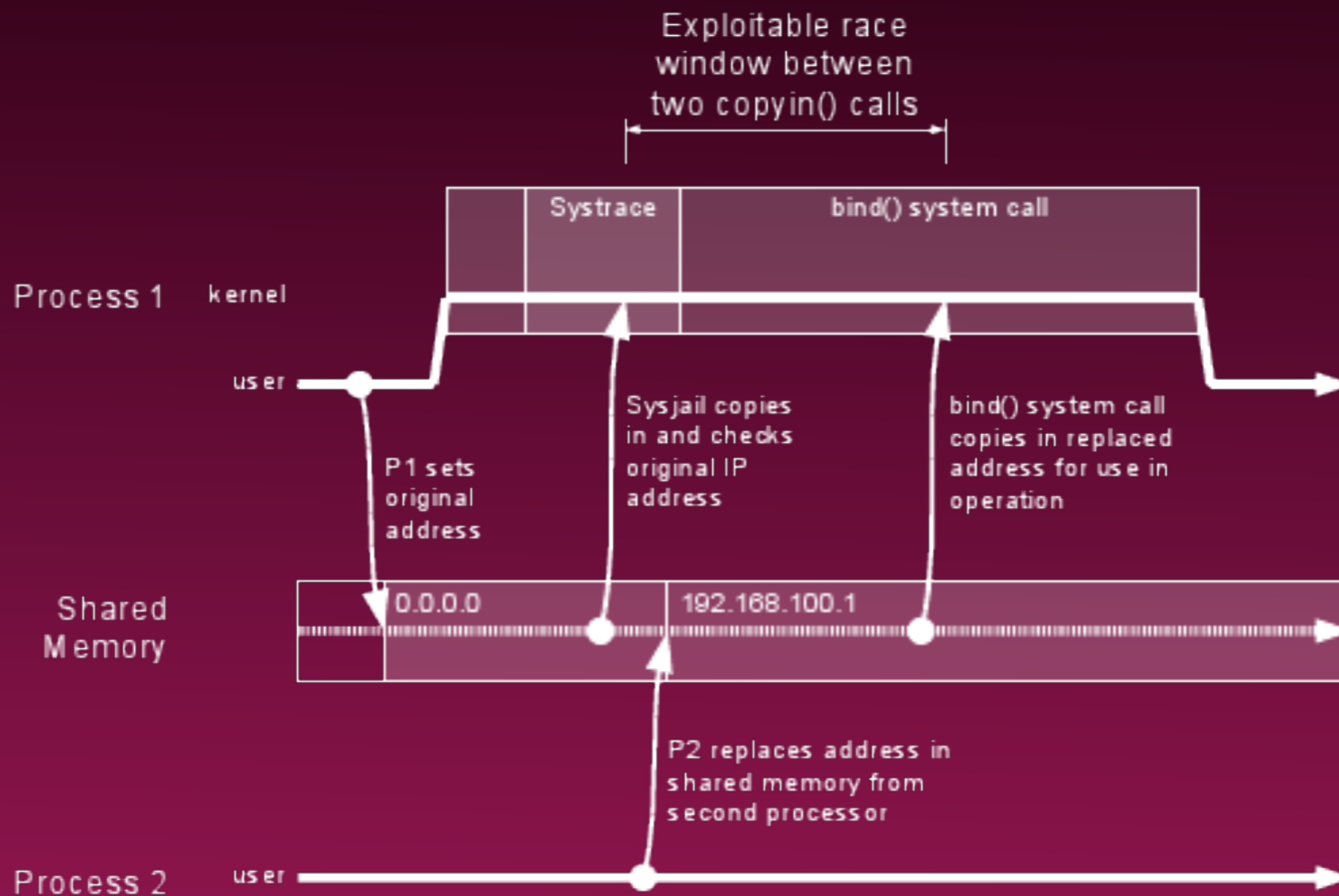
# Exploatare race-uri pe sisteme multiprocesor

- Se folosește de paralelismul oferit de hardware
- Trebuie determinată fereastra de timp în care se poate face înlocuirea
  - se folosește un timer de mare precizie (TSC) pentru a temporiza înlocuirea + binary search pentru a determina fereastra de vulnerabilitate
  - la argumentele care sunt înlocuite de SWC, urmărirea în buclă a faptului că argumentul e modificat

# Exploatare race-uri pe sisteme multiprocesor [2]

- Rezultatele experimentale indică o fereastră de la 5K-15K cicluri (GSWTK) până la 100K (systrace)
- Mai mult decât suficientă pentru atacuri

# Exemplu



# Generic Software Wrapper Toolkit

- Sistem generic, bazat de wrappere configurabile
- Wrapperele sunt scrise folosind un limbaj gen C extins cu suport SQL
- Rulează pe Solaris, FreeBSD, BSD/OS, și Linux
- Există o multitudine de wrappere, de la sisteme de control al accesului până la sisteme de detecție a intruziunilor



# Din 23 de wrapper inspectate 16 au fost identificate ca vulnerabile și exploate cu succes

Wrapper	Description	Vulnerabilities
<code>callcount</code>	Count system calls	None: relies on the system call number.
<code>conwatch</code>	Track IP connections by processes.	Postcondition TOATTOU race on <code>connect()</code> and <code>bind()</code> allows masking address/port used.
<code>dbfencrypt</code>	Encrypt files with '\$' in their names; prevent rename so that encryption policy on a file cannot be changed.	Postcondition TOCTTOU race allows incorrect name in policy check; precondition TORTTOU races on I/O write unencrypted data and bypass rename checks.
<code>dbexec</code>	Authorize execution of programs based on a pathname database.	Precondition TOCTTOU race allows bypass by substituting the name during the wrapper check.
<code>dbsynthetic</code>	Synthetic file system sandbox substituting pathnames from a database.	Precondition TORTTOU race bypasses path replacement; postcondition TORTTOU race leaks true paths
<code>life</code>	Prints the process life cycle.	Precondition TOATTOU race replaces <code>exec()</code> paths.
<code>noadmin</code>	Deny all privileged operations.	None: relies on the kernel's process credential.
<code>aks.wr</code>	Audit file operations	Pre/postcondition TOATTOU races avoid audit.
<code>seq-kernel.wr</code>	Sequence-based intrusion detection	None: relies on the system call number.
<code>imapd.wr</code>	Detect anomalous access by <code>imapd</code> .	Postcondition TOATTOU path races prevent alerts.



# Systrace

- Sistem prin care un proces poate controla un target process prin inspectarea și modificarea argumentelor apelurilor de sistem
- OpenBSD, NetBSD; ports pentru Linux, Mac OS X, FreeBSD
- Au fost testate module sudo monitor mode și sysjail

# sudo monitor mode

- se interceptează execuțiile și se face audit la comenzile executate
- Politica este menținută în user-space
  - sunt necesare context switch-uri pentru a citi și acționa conform politicii
- Pe sisteme MP fereastra de vulnerabilitate a fost determinată la 430K cicluri
- Atacul a reușit înlocuirea argumentelor, mascând astfel acțiunile în log

# sysjail

- rulează procesele ce sunt necesare a rula ca root
  - pentru a minimiza impactul asupra unei posibile vulnerabilități în aplicație
- se atașează la toate procesele ce rulează în jail
- validează sau în anumite cazuri rescrie argumentele apelurilor de sistem

# Exemplu bind

- adresa la bind trebuie să fie
  - egală cu cea configurată pe jail
  - `INADDR_ANY`, caz în care jail-ul o va înlocui cu cea configurată pe jail
- S-a reușit exploatare pe un sistem MP prin înlocuirea adresei de bind

# CernNG

- Similar cu GSWTK
- S-au exploatat vulnerabilități de tipul TOATTOU și TOCTTOU
- În ciuda unor metode de protecție targetate împotriva unor astfel de atacuri (bazate pe memorie virtuală)

# Cum se pot preveni aceste race-uri?

- Tehnici de mitigare
  - se modifică SCW-urile pentru a ne proteja la aceste tipuri de atacuri
- Modificarea sistemului de operare dar păstrarea SCW-urilor
- Renunțarea la SCW-uri

# Tehnici de mitigare

- Detectarea modificării argumentelor la post-condiție și roll-back
  - rollback-ul nu este practic pentru apeluri de sistem complexe
  - detectarea modificării argumentelor poate fi și ea exploatată cu aceleași tehnici

# Tehnici de mitigare [2]

- Marcarea paginilor care mențin argumentele ca RO pe durată apelului de sistem
  - se poate face doar la nivel de pagina; dacă alte thread-uri stochează date în acea pagină...
  - trebuie proiectate toate paginile mapate asociate cu acea pagina fizică
  - memory mapped files: trebuie avute în vedere și operația de write()
  - trebuie avut în vedere și faptul că aplicația poate face remap cu RW



# Tehnici de mitigare [3]

- Stack gap
  - argumentele indirecte se copiază într-un zonă rezervată, pe stivă
- Probleme
  - de implementare: zona este accesibilă din alte thread-uri și poate fi mapată chiar și în alte procese
  - copieri în plus

# Tehnici de mitigare [4]

- Concluzii
  - greu de implementat
  - probleme fundamentale
  - În practică ineficiente: atât “stack gap”-ul cât și sistemul de VM protection oferit de CerbNG au putut fi ocolite

# Message Passing Systems

- Argumentele apelurilor de sistem sunt trimise sistemului de operare într-un bundle
- Elimină race-urile sintactice, nu și cele semantice
- Trap-ul ce se ocupa cu tratarea apelurilor de sistem trebuie să aibă știe formatul argumentelor

# Renunțarea la SCW

- Integrarea unui framework de securitate în kernel
- Modelul cel mai larg adoptat în curent de comunitatea “OS security”: LSM, SeLinux, SEBSD, SEDarwin