
Hyper-Threading Considered Harmful

Colin Percival

`cperciva@freebsd.org`

Hyper-Threading

- Quick introduction to Hyper-Threading:
 - Present on recent Pentium Extreme Edition, Pentium 4, Mobile Pentium 4, and Xeon processors.
 - A single processor core executes two threads simultaneously.
 - In theory, throughput is increased by taking advantage of execution units which would otherwise be idle.
 - In practice, the benefit is only significant on desktop applications.
 - Some resources are duplicated for each thread.
 - Memory caches and the external bus are shared.

Multi-level security

- Processes operating on Top Secret data aren't allowed to talk to processes without Top Secret clearance.
 - Necessary to prevent “covert declassification”.
 - Requires that security checks are performed on files, network sockets, signals, *et cetera*.
 - Robert Watson can probably explain the details far better than I can.
 - A mechanism for illicit communication is called a “Covert Channel”.
 - Terminology: The process which sends data is the “Trojan”, the process which receives data is the “Spy”.

Covert communication via paging

- Assume that the Trojan and the Spy have access to a large reference file (e.g., the Encyclopædia Britannica).
 - The Trojan reads some pages from the reference file.
 - The Spy reads the entire reference file, but measures how long each memory access takes.
 - Loading pages from disk is slow!
 - The Spy can determine which pages were accessed by the Trojan.
 - Covert channel rate: 200 - 1000 bps.
 - Having a shared reference file is essential!

Covert communication via paging

- What if we don't have any shared reference file?
 - Assume two processes each have an address space larger than half of the available memory.
 - To transmit a “1”, the Trojan reads its entire address space.
 - Some pages owned by the Spy are evicted from memory.
 - To transmit a “0”, the Trojan spends the same amount of time accessing a single page.
 - The Spy measures how long it takes to read its entire address space.
 - Covert channel rate: 0.01 - 0.1 bps.
 - Much slower, but no shared reference file is needed.

Pentium 4 cache hierarchy

- The Pentium 4 L1 data cache contains 128 lines of 64 bytes each.
 - Cache lines are divided into 32 4-way associative sets.
 - A pseudo-LRU eviction strategy is used within each cache set.
- The Pentium 4 L2 cache contains 4096 lines of 128 bytes each.
 - Cache lines are divided into 512 8-way associative sets.
 - Again, a pseudo-LRU eviction strategy is used.
- The caches are shared between Hyper-Threads.

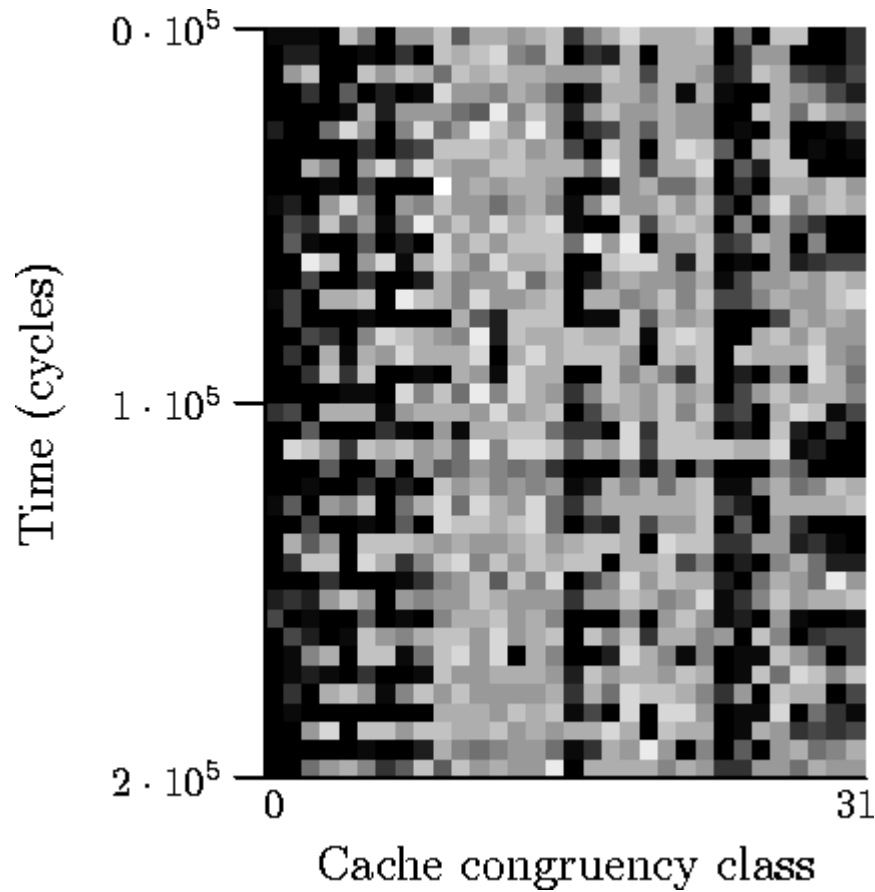
Covert communication via caching

- Each cache set behaves like a virtual memory paging system.
 - No data is shared between Hyper-Threads, but one thread can evict cache lines owned by the other thread.
 - A thread can measure how long it takes to access data, thereby determining if the data was in the L1 cache.
 - The Trojan reads memory locations determined by the data it wants to transmit.
 - The Spy repeatedly reads 4 cache lines mapping to each cache set, and measures whether the Trojan forced any of those lines to be evicted.
 - Covert channel rate: 2 - 5 Mbps (L1) or 0.5 - 1 Mbps (L2).

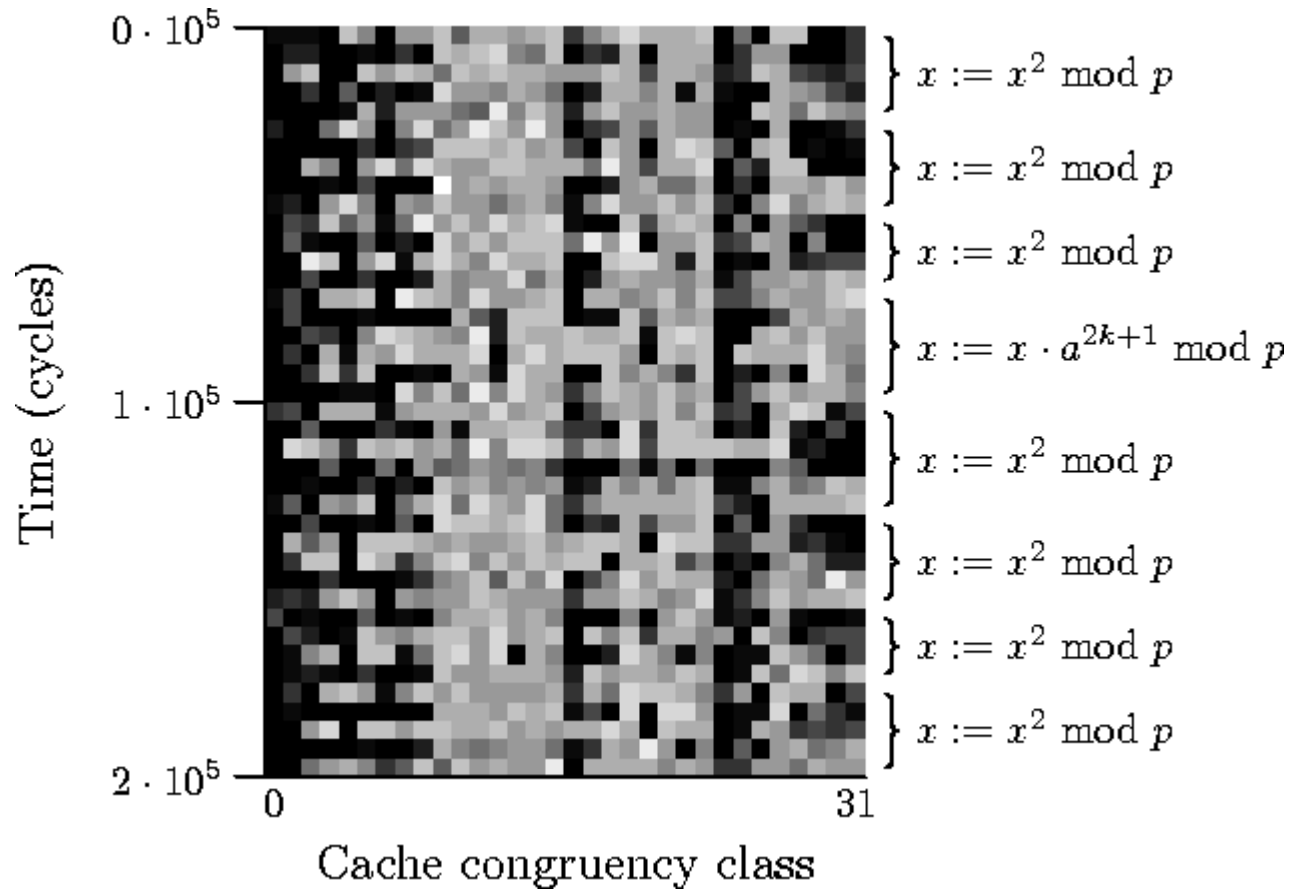
Attacking OpenSSL

- 1024-bit RSA is secure, right?
 - Only if the implementation doesn't permit any side channel attacks.
 - OpenSSL performs 1024-bit private key operations using two exponentiations modulo 512-bit primes.
 - Each exponentiation is performed using a series of squarings ($x := x^2$) and multiplications ($x := x \cdot a^{2k+1}$).
 - If we can determine those exponents, we can factor the RSA modulus and break the encryption.
- What if we use the Spy to monitor which cache sets are being accessed by OpenSSL?

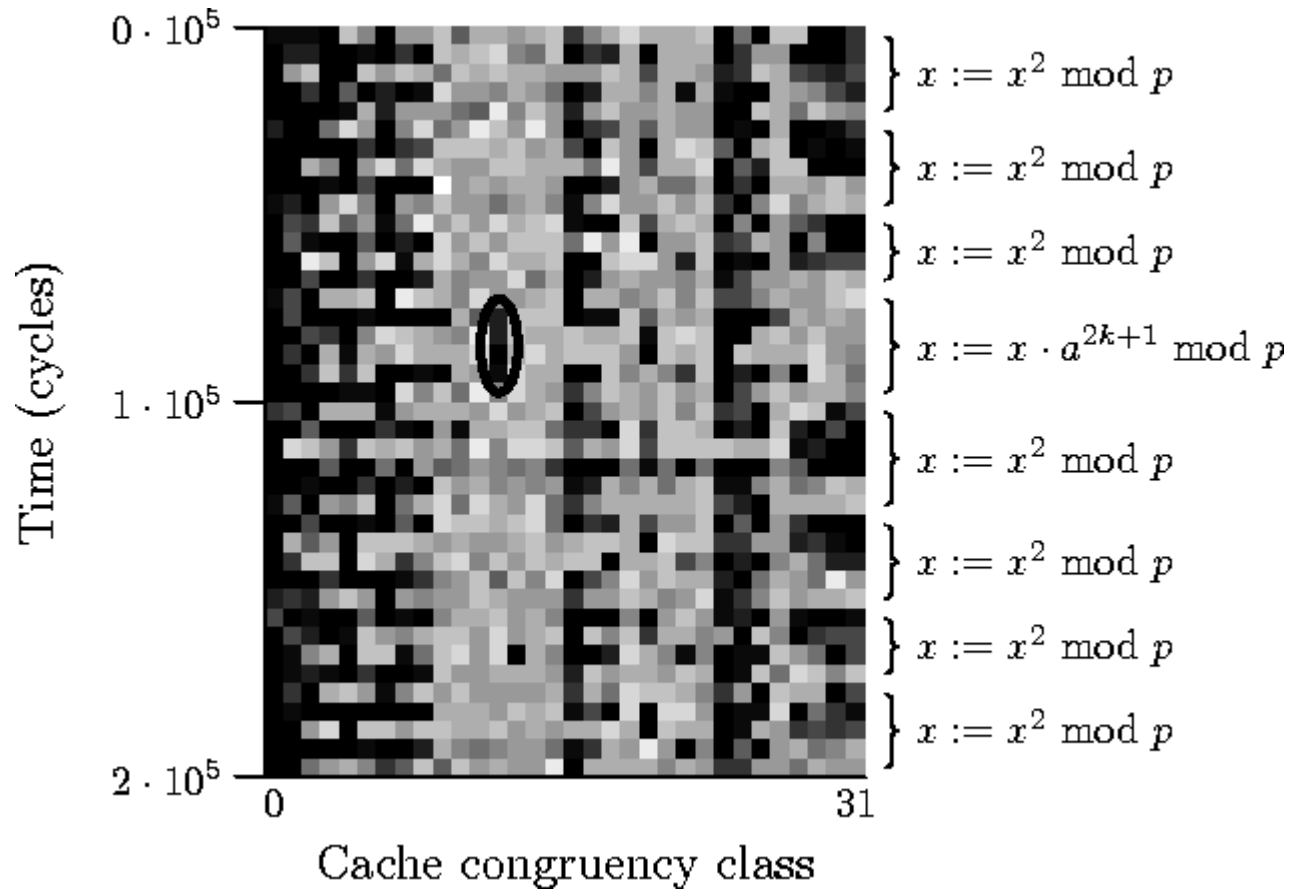
Attacking OpenSSL



Attacking OpenSSL



Attacking OpenSSL



Attacking OpenSSL

- The cache usage observed by the Spy reveals information about the exponent.
 - The sequence of squarings vs. multiplications reveals ~200 bits out of 512.
 - The locations of the multipliers reveals another ~110 bits of exponent.
 - In total we can steal ~310 bits out of 512.
- Given 256 bits from each exponent, we can factor the RSA modulus N in polynomial time.

Bad news

- If an attacker can run his code on the same processor core as your RSA operation, he can steal your key.
 - You only need to spy on OpenSSL once.
 - All Hyper-Threading servers which permit user logins over SSH are affected.
 - All shared SSL web servers which allow CGI scripts are affected.
 - Not just OpenSSL.
 - Not just RSA.

More bad news

- What if we don't have Hyper-Threading?
 - Caches aren't usually flushed across context switches.
 - There might be information left in the L1 cache after a context switch.
 - There will certainly be information left in the L2 cache after a context switch.
 - Cache evictions across context switches provides a covert channel of at least 20kbps (maybe more).
 - Can this be used as a cryptanalytic side channel?
 - Don't know for certain, but it wouldn't be very surprising.

Fixing Hyper-Threading

- Caches in Hyper-Threaded processors need to be made threading-aware.
 - Divide the cache into per-thread sub-caches.
or
 - Use a more sophisticated cache-eviction strategy.
 - Ensure that when threads compete for cache space, they each end up with their “fair share” of the cache.
 - After the first few cache evictions, threads stop evicting each others' cache lines and instead evict their own cache lines.
 - I have no idea if this is feasible in silicon.

Defenses in the OS

- The operating system could disable Hyper-Threading completely.
 - Simple, obviously correct.
 - Loss of performance, makes Intel unhappy.
- The operating system could use Hyper-Threading carefully.
 - If the threads sharing a CPU are allowed to debug each other (i.e., are not setuid and have the same uid), then Hyper-Threading is harmless.
 - Much harder to get correct.
 - Potential problems with kernel data locking.

Defenses in the Applications

- Write crypto code so that the code path and memory accesses do not depend upon the input data or key.
 - Immune to this attack, since the cache footprint will never change.
 - Even better, immune to all timing attacks.
 - (Assuming constant-time arithmetic operations.)
 - Some loss of performance is inevitable.
 - Do we really care?
 - Hard to get right; even harder to modify existing code.
 - Even if cryptographic libraries are fixed, applications will still be leaking data...

Conclusions

- Disabling Hyper-Threading is a necessary first step.
- Writing a Hyper-Threading-aware scheduler would be a good idea.
- Rewriting cryptographic libraries to function obliviously to the data and key would be a good idea.
 - Even better, write a new library from the ground up.
- Hopefully future processors will remove these channels (but I'm not optimistic).
- We still have work to do.

<http://www.daemonology.net/hyperthreading-considered-harmful/>

Questions?