

# Exokernel

January 22, 2004

Key Ideas:

- o separate protection from policy: normal kernel does both
- o default policy might not be the best, but very hard to change
- o challenge: how to structure protection for flexibility \*and\* sharing among untrusted processes
- o move most of traditional OS to user-level library (a “lib OS”), to maintain compatibility

Extends the trend of better support for application control:

- o ACPM: control over paging and replacement
- o Scheduler Activations: control over threads and scheduling (also Capriccio)
- o U-net and others: control over networking

Principles:

- o Separate protection and management: use the lowest level required for protection, ideally the hardware level (e.g. disk blocks). Types of calls: allocation/revocation, sharing, tracking/changing ownership
- o Explicit allocation and revocation
- o Use physical names: avoid translation from virtual names (and its consistency problem)
- o Expose information: typically read-only access to kernel state for making decisions in the libOS

Four basic tools: [slightly different than written in the paper!]

- o software regions: like Multics segments -- sub-page protection and fault isolation
- o hierarchical capabilities, which must be passed on most system calls
- o bind hardware resources together (e.g. disk block and page frame that holds it)
- o download code (that can be verified) to control some policies. E.g. wakeup predicates determine if a process is ready to run (without waking it up)

Multiplexing Storage -- extraordinarily complicated!

- o goal: allow each app to have its own file system, but on the same disk, and share file systems when needed
- o challenge: have to figure out which blocks go with which file system, for almost arbitrary file systems (that don't trust each other!)

Some attempts that didn't work:

- o capability for each block (or extent): would work, except that you need to either put it in the block, which is a mess, or put it somewhere else have many more seeks... [would caching solve this?]
- o self-describing metadata: so that XN can figure out where the references are without knowing all the details about the file system. Hard to find the right language for this.
- o template-based metadata: use more complex language, but only once for each template. But still couldn't find the right language; i.e. flexible but simple

Actual solution: UDFs: *untrusted deterministic functions*

- o idea: don't try to understand the metadata -- just make the file system tell you in a consistent way what the references are
- o Usage: file system author writes a function for each template type,  $owns-udf_T$ , and gives it to the kernel, which verifies that it is deterministic. The kernel then uses this function to verify every proposed change to a block, so that libOS can't change references without the kernels permission.
- o  $owns-udf_T(\text{block}) \rightarrow$  list of (block #, template type) pairs  
This returns the (deterministic) set of pointers in the input block
- o Example: to add a data block,  $d$ , to an i-node,  $i$ , you call something like:  
 $XN\_allocate(i, d, \text{changes to } i)$   
XN would create  $i'$  by applying the changes to  $i$ , and would confirm that:  
 $owns-udf_T(i') = owns-udf_T(i) + d$

Access control for pages:

- o done when a block is mapped into an address space -- which means you don't need to check it on every access (bind + fast access)

Must also deal with write ordering for persistence (read the paper for this)

- o root list keeps the root block for each file system, from which XN can build the whole tree (using  $owns-udf_T$  recursively)

Buffer cache registry:

- o map of all pages and their access privileges (and dirty/clean)
- o mapped read-only into libOSs
- o any process can write back pages! (can see the data, but can force the writeback)

Fork is challenging:

- o must copy your own process while you are running!
- o idea: copy all of the pages that are not the active data/stack pages via copy-on-write, but handle the hard ones directly

Related Work:

- o microkernels: help this problem via modularity -- easier to replace one part of the OS (to change its policy)
- o faster system calls: help performance but not policy

- o “info kernel” work: make more info available to user level, and then trick the OS into doing the right thing. For example, you can change the caching policy by issuing the right set of reads (from the lib OS layer). [from Univ. of Wisconsin, SOSP 2003]

Cheetah is compelling:

- o high performance, highly tuned web server
- o precompute TCP checksum in cache, so you can send them directly over the network
- o avoid extra ack packets by knowing the protocol
- o colocate files that are part of the same HTML page to reduce seeks
- o about 4-8 x faster for small docs, maybe 2x for large docs

Some lessons:

- o changing layers has more impact than optimizing the layers as they stand
- o fast upcalls more useful than downloadable interrupt handlers (and easier)
- o download code to extend policy (wakeup predicates)
- o is it worth all this effort? how many applications will benefit? could most of the value be received via changes to Linux instead?