

# Overview

- This set of slides is a tour through the Ethos source code
- Ethos obtains significant advantage by implementing on top of Xen
- Much simpler I/O device interface
- Availability of a paired OS for performing OS services such as file system.

# Architectural Components

- An operating system can be targeted to different architectures
- Ethos is designed for 64-bit architectures
- Which means its designed for architectures with considerable resources, especially memory
- Ethos current implementation is 32-bit
- Ethos will probably have two ultimate architectures, x86 and ARM.

< 67 ►

▲ 臣 ▶ ▲ 臣 ▶ = = • • • • • •

# Architectural Components

- An OS is implemented in two layers
- The base abstractions using machine dependent layer
- On top of that a far larger machine independent layer is built
- To port to a different architecture, need only implement a machine dependent layer
- This assumes that the layering is properly done
- And that takes a while, so we'll concentrate on a single architecture for now

# Architectural dependent components

Part II

# on A. Solworth Secure OS Design and Implementation Ethos Tour

・ロト ・四ト ・ヨト ・ヨト

31

3

# <</td> □ > <</td> ≥ ≥ <</td> ≥ <</td> ≥ <</td> ≥ ≥ ≥ ≥ ≥ <</td> ≥ > ≥ ≥ ≥ ≥ > ≥ ≥ > ≥ > ≥ > ≥ > ≥ > ≥ > ≥ > ≥ > <

Ethos Tou

▲ 臣 ▶ ▲ 臣 ▶ ○ 臣 ○ � � �

# Architectural dependent components

The architectural dependant components include:

- Types are used to express size requirements of different components of the OS
- The compiler generates machine language for the architecture (only non-privileged instructions)
- The machine dependent OS components include:
  - The privileged instructions
  - Architectural dependent data structures
  - Bits of glue logic (such as for system calls)
- Since Ethos is built on top of Xen, it does not have any privileged instructions

# Architectural dependent source files

 $\operatorname{arch}/\times 86$  contains the code for both 32-bit and 64-bit  $\times 86$  architectures

- The 32-bit code works in non-PAE mode
- The PAE mode and 64-bit code are residuals from Mini-OS
- debug.c: performs a register dump on crash
- archPageTable.c multilevel page table operations (exports a flat page table with R/W permissions)
- minios-x86\_32.1ds: 32-bit load script for the OS
- sched.c scheduling and process code.
- setup.c various bits of code for Ethos start up
- time.c time using Xen's time facilities
- trap.c setting up traps and interrupts using Xen events
- x86\_32.S assembly glue code for traps, etc.

# Architectural dependent include files

include/ethos/x86 contains the architecture dependent files

- archMemory.h exports a flat page table info to architecture independent portion of OS
- archPageTable.h is the internal hierarchical page table
- arch\_debug.h interface to dump registers
- arch\_elf.h interface to elf object loader for user space
- arch\_sched.h interface to scheduling and process primitives
- bits.h low-level locking primitives
- cpu.h low-level CPU definitions, such as interrupt numbers, etc.
- ldsyms.h symbols set by the linker script
- synchro.h various types of memory barriers defined

Jon A. Solworth

• x86\_32/hypercall-x86\_32.h hypercalls for the x86 32-bit architecture

```
Part III
Architecture-independent components
```

# Architecture-independent components

# Primary directories

console Xen console devices. This is the code that processes Xen error messages

- net Xen pseudo ethernet driver code.
- file Filesystem implementation
- mm Memory management
- process process support, including process scheduling
  - rpc remote procedure call support for communicating with the shadowdaemon
- syscall system call interface

userspace contains the shadow daemon

# Part IV

# Xen Components

< 17 →

-2

・ロト ・ 日 ・ ・ 日 ・ ・ 日 ・ うへの

Ethos Tour

Secure OS Design and Implementation

• = •

- 2

・ ロ と ・ ( 雪 と ・ ( 目 と

Secure OS Design and Implementation

3

Ethos Tour

Ethos is not intended to be run stand alone, and hence is dependent on Xen.

- Parts of the Xen-interfacing components in the architecturally dependent components
- Others are in the Xen directory
- The include files for Xen are in include/xen

Jon A. Solworth

• The Xen include files are currently from Xen version 3.2, and of course follow Xen coding style

・ロト ・四ト ・ヨト ・ヨト

Secure OS Design and Implementation

3

Ethos Tour

# Xen

The **xen** directory includes:

- xenPageTable.c includes procedures to:
  - Switch the current page table being used (this is part of a context switch)
  - To pin and unpin pages. Pinned pages are maintained in the Xen hypervisor to minimize checking of page reference.
- xenSchedule.c schedule the Ethos OS vs. other OSs running under Xen.
- xenGrant.c grant pages are pages which are shared between multiple OSs. Xen provides mechanism by which pages can be shared and sharing changing. This is used for device drivers (in particular, ethernet pseudo devices).
- xen\_hypervisor.c Hypervisor commands which should be moved to other files.
- xenbus.c provides Xenbus/Xenstore interface to store and retrieve name-value pairs.

Jon A. Solworth Secure OS Design and Implementation Ethos Tour



# Include directories

include/ethos include files which are used only within the Ethos kernel include/userspace this are include files which are used both within the Ethos kernel and by Dom0 and/or Ethos user space include/xen include files from Xen



# KernelInitialize

Steps to initialize (boot start up code):

- arch\_init Architectural specific initialization
- ② trap\_init Initialize interrupt vector
- sen\_event\_init initialize Xen events
- Interrupts
- onsoleInit initialize the console
- 6 setup\_xen\_features
- @ memoryInit initialize virtual memory
- sen\_grant\_init initialize grant memory for pseudo-devices
- Init\_time initialize wall clock time computation

# Kernellnitialize (cont'd)

- o eventInit initialize Ethos events
- xenbus\_init
- netInit initialize the network interface
- pairedConnection set up connection to DomO shadowdaemon
- fileSystemInit initialize file system
- initstor\_init initialize the initstor
- debug\_init read in debugging info for the Ethos kernel
- scheduleInit creates the init process, after which Ethos is up and running.

▲母▶▲臣▶▲臣▶ 臣 のQの

< 67 ►

▲ 臣 ▶ ▲ 臣 ▶ ▲ 臣 → ���





- Ethos deals with typed data
  - In the file system,
  - for inter-process communication, and
  - across the network.
- The most flexible way of implementing this is RPC
- RPC right now is used between the Ethos kernel and Dom0 shadowdaemon
- That means that the RPC code needs to be compiled in both environments

▲ 臣 ▶ ▲ 臣 ▶ ▲ 臣 → ���

< A

# **RPC** components

- netInterface.c provides the interface to network device. This is the "logical" interface, the "physical" interface is either Xen's net front or UNIX sockets.
- tunnel.c provides a logical conduit between two network interfaces
- connection.c is a logical connection. Each connection is associated with a tunnel, a tunnel can support an arbitrary number of connections.
- procedure.c stores and looks up the RPC specification
- rpcInterface.c describes how to marshall and unmarshall RPC types
- **rpcType.c** describes the primitive types and vectors of the primitive types.
- packet.c connections are implemented with a sequence of packets

Secure OS Design and Implementation

Jon A. Solworth

Ethos Tour



# Procedure-dependent components

- RPC allows an interface to be defined consisting of a set of procedure
- ethosRpc.c describes the remote procedures. It consists of
  two types:
  - procedureDispatcher invokes a procedure. It has a case statement with an entry for each procedure.
  - **rpcInit** describes the remote procedures in a table format.
- It is possible to automatically generate ethosRpc.c from a short description, but to do that, we would need an IDL compiler. Unfortunately, we haven't gotten around to writing one yet.

## Jon A. Solworth Secure OS Design and Implementation Ethos Tour



# Net

In console/console.c

• Provides printk, the kernel printf equivalent

# in console/xencons\_ring.c

- Contains the interface to a 2 kilobyte Xen-based console buffer
- Not invoked directly, but only through console.c

Jon A. Solworth

There are three externally facing functions in net/net.c

- **netInit** initialize network interfaces
- netSend send a packet to the network
- InetIncoming process incoming packets

net/xenNetInterface.c are the Xen interface. There needs to be more code moved from net.c to xenNetInterface.c.

# Memory Management

- kmap.c these are reserved mapping in the kernel so that Ethos doesn't run out of memory at inopportune times.
- memoryInit.c initializes the memory subsystem
- pageAllocator.c allocate pages using a buddy allocator. At kernel initialization, all unallocated pages are put into the buddy allocator.
- pageFault.c handles page faults
- pageTable.c exports a flat page table abstraction to rest of OS

# $\mathsf{Part}\ \mathsf{X}$

# Memory management

・ロト ・聞 ト ・ヨト ・ヨト

Secure OS Design and Implementation

E 990

Ethos Tour

< 注 → < 注 → □ 注 □

・ロト ・ 日本・ ・ 日本・

Secure OS Design and Implementation

- 34

Ethos Tour

# Memory Management

- pmem.c this manages memory for a single process.
- slob.c sub-page allocator
- vaddr.c checks whether virtual addresses are mapped into the virtual address space.
- vma.c manages process regions
- vmalloc.c a second buddy allocator

Jon A. Solworth



# File system

- directory.c directory operations to open, create, and delete directory entries.
- file.c file operations to read and write
- fileSystem.c create initial in memory data structures to access file system.
- mode.c permission string to/from integer values for access type
- resourceClose.c close a resource descriptor
- terminal.c terminal read and write



< 17 →

▲ 臣 ▶ ▲ 臣 ▶ ● 臣 ● � � �

・ロト ・四ト ・ヨト ・ヨト

Secure OS Design and Implementation

∃ 990

Ethos Tour

- process.c manages processes, including creating, deleting, and blocking processes.
- fileDescriptor.c is the per-process file descriptors
- schedule.c is scheduling for processes, which happens when a process blocks or is coming out of a (completed) system call.
- elf.c loads elf file into a process on exec or creation of init process.

# <section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><section-header><image>

# Syscalls

- System calls are the invocation of OS facilities from user space
- Ethos syscalls pass the values in the registers

Jon A. Solworth

- Syscall number
- Pointer to input parameter structure
- Pointer to output parameter structure
- System calls must check that parameters are valid
- Including that pointers are to allocated memory in user space
- In addition, Ethos returns values whose size is not known when the call is made.
- Note that syscall/syscall.c is the kernel side of system
  calls, the user side is in sub-directories userspace/ethos/
  called libsyscall

# Syscall macros

getArgs(\_var) gets the argument structure and puts the result in \_var. This retrieves the fixed sized arguments getMemStruct(\_memStruct) gets a variable size MemStruct, which is a pair of size and pointer. The memStruct is one of the arguments retrieved by getArgs. Allocates space in the kernel for the variable sized structure.

- putReturn(\_var) put the fixed sized components of the value to be returned to user space.

These macros branch to done if the macros fail. The variable size putMemStruct is mostly used in **retire**.

・ロト ・聞 ト ・ヨト ・ヨト

Secure OS Design and Implementation

E 990

Ethos Tour

▲■ ▶ ▲ 臣 ▶ ▲ 臣 ▶ ○ 臣 ○ の Q @

