

• Worse is Better: (rarely)

} MIT/Stanford philosophy

" New Jersey "

Worse is better \rightarrow simplicity of implementation is

main concern

consistency & consistency for the other one

(Point the PC back to user mode)

Buy a user prog invokes sys. routine. sys. routine

should point PC to the routine.

(because not simple)

New Jersey doesn't deal with this while MIT/Stanford

wrong.

Because in high demand, get 90% functionality.

Def. of complexity: sys is a set of interacting parts

collectively have a behavior. Complex \rightarrow large # of

components, many ppl req. to maintain

Problem: many emergent properties

- Scaling: not same as designing for users

or for billions.

- ~~point~~ if one is wrong, others could be

Hierarchically layering \rightarrow reducing connection

Therac-25: (every 2 years).

Follows Therac-20 and other machines

an prob: 2 modes - Electron therapy

Energy	$\frac{E \cdot T}{5-25 MeV}$	X-Ray therapy mode X-Ray (Photon therapy): ^{has} because X-ray flat error
Current	X	100x

where E.T is active & current set to X-ray

dangerous

before hardware locks preventing dangerous state

Therac-25 only relied on software. There were

cryptic error msg - maljon comes up, operator has

choices → restart

→ Press (P)

Since error msgs were cryptic & operator would press (P)

detecting was tough

simplified by ← copy, if happens quickly enough

can enter bad state

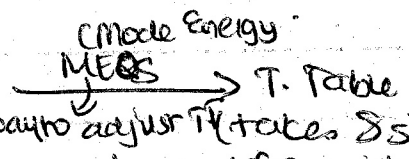
Keyboard handler

Apple IIish type

↓

Data entry complete flag

MEOS



If operator enters DE. and modifies within 8s, not

identified by T.T.

① Reliability vs safety.

if seemed to do nothing wrong doesn't mean safe.

② lack of defensive design. → self checks

Not making cryptic error checks.

③ lack of two locks → single pt. of failure.

~~① single pt of failure~~

~~② ref. uni~~

-unrealistic risk assessment, company uses fault tree analysis. Unrealistic → prob. of software failing → have to test all parts. Some parts difficult to assign value to.

~~Safe~~ Friendly UI vs safe UI.

Assigned 0 probability to ~~these~~ software errors ← Problem.

Not about locks being acquired in wrong order. Problem of not having locks.

T.T. sensors were there. ^{did not} ATCL responded to Therac 20's complaints but to Therac 25 ← false.

Client/server

- call D

RPC: H

at-least

from

at-least

send more

more than

Exactly on

X-windows

In general

win

X

Want to draw

20090.

If overlapp

send requ

Client/Server.

- call DUNCS on other machine, look like local call.

RPC: How would want to work on RPC:

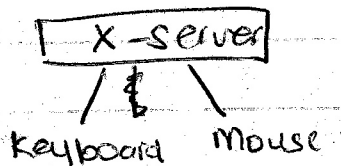
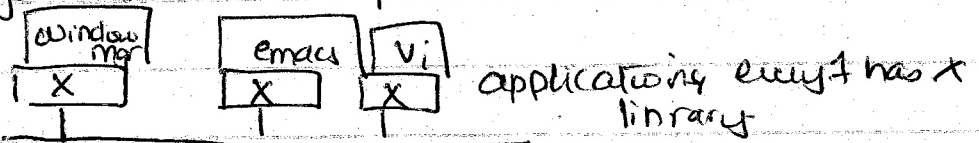
at-least^{least} once: Resend packages until get something from client.

at-most^{most} once: Send it once, if don't receive don't send more than once. Amazon: don't get credit card more than once, but might not get anything. CI

Exactly once: Hard to do.

X-windows: Client/Server (just switched).

In general, server is your machine.



Want to draw line/shape on your screen.

20090.

If overlapping windows, move window out of way. Do you send request? Usually not.

Naming in file sys:

- Provides abstraction eg. DNS.
- Provides user friendly way.

UNIX file system:

Block layer ~~→ File layer~~ → Inode # → Filename
→ Path → Abs. Path → Symbolic path

file → Inode #
inode

One to one relation.

Hard links vs Symbolic link

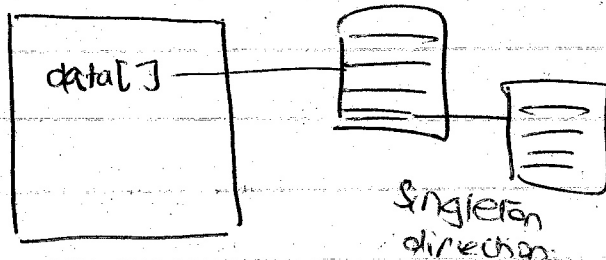
↓
ref to actual
inode,
pt. to same
inode

No simple way to go from inode to path name.

Questions about how look up.

2006 4 - good set of questions (Q5, 6).

UNIX v6



Can't store large files. Scheme in 2006 pitfalls,
try to bump size to 64 bit.

what happens if change block addresses to 4 bytes, how much increase file size & look at 2008 & iv gain idea of math experts

Lec 5.6,7 (UNIX, OS, Eraser)

OS structure: ~~Abstract devices~~

an OS {

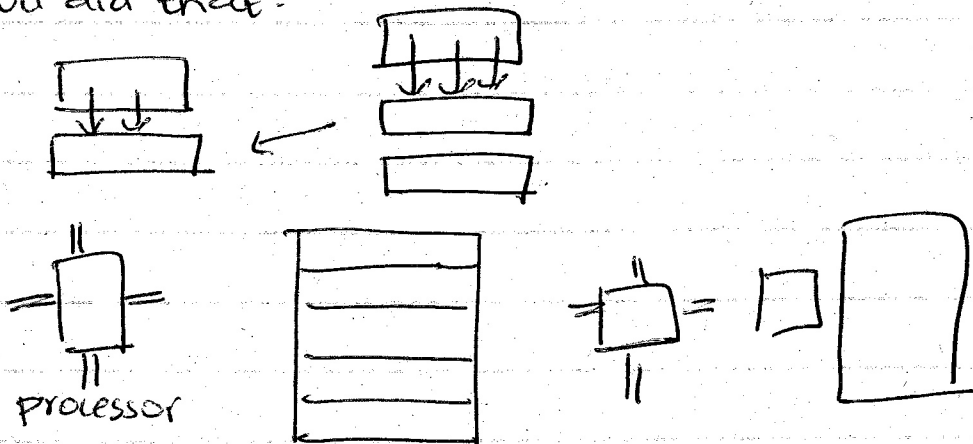
- abstract devices
- run many prog. together/multi program

Techniques {

- virtualization
- abstraction

2 layers communicate in certain defined interface.

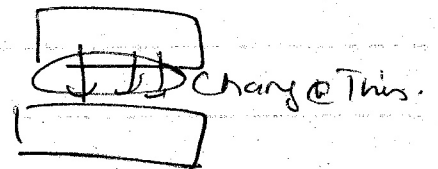
Virtualization put 3rd layer w/o telling upper layer that you did that.



read 0x01

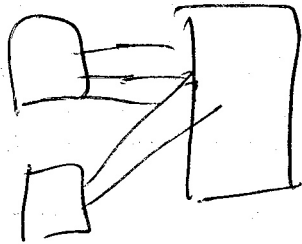
Virtual memory → put something in it

Abstraction: take lower layer, change lang. you use



Virtual memory

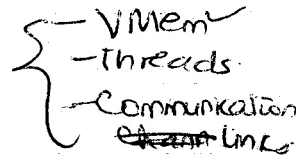
can have many processors share same memory.



No matter how messed up your code is never writing to what other is writing.

Another visualization:

Threads:



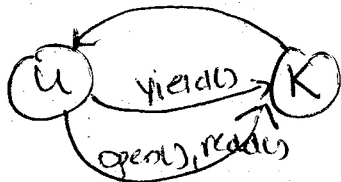
Virtualizing processor

Basic abstraction. For threads:

- yield(): stop current thread go to next one

yield() is kind of ^{OS} procedure called system call.

It goes in kernel. If prog. your own sys. call won't be able to. ~~for~~ Since your bit will be in user mode not kernel, can't write your own sys. call.



Why want to control how go from U to K?

Security, and reliability. Simplify the interface.

Very limited way of going from U to K.

Instead of complete freedom, gives limited freedom to user

State where user ~~has~~ should not touch but Kernel has to touch.

e.g. in virtual memory, each processor knows where its mapping table is. To do that looks at that table. Has to be changed but done in restricted way. That is what a kernel can do, but user can not.

Communication b/w processors:

PA	PB.
X=0	X=0

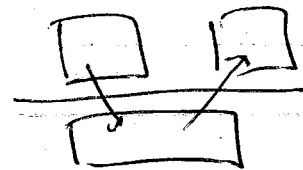
X defined ⁱⁿ context of the program. Talking about diff. X both of them.

But what if want to communicate?

SHARED BUFFER → OS simulate the network.

Links:

↳ Bounded Buffer: when 2 prog want to communicate, make sys. call SEND(), RECEIVE().
don't know when to send or receive.



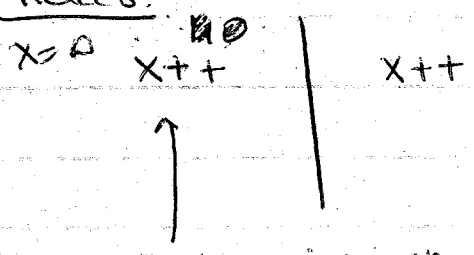
when call send(), kernel takes

care of it. Keeps an array. When somebody sends msg, put in queue, receive → take out of queue.

Problem → Concurrency Issues.

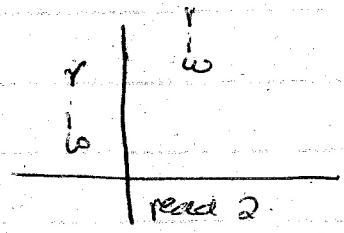
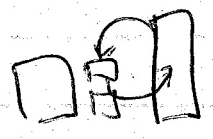
When someone calls kernel, another kernel code may be using it at same time.
 What could go wrong? (Race condition).
 Why race cond. a problem?

Races:

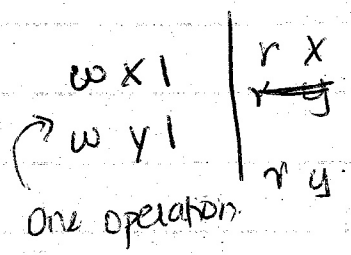
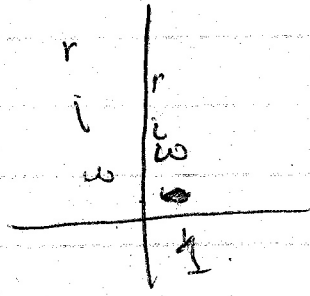


abstraction breaks internally ~~read~~

x++ → { read x to
 inc to
 write to x

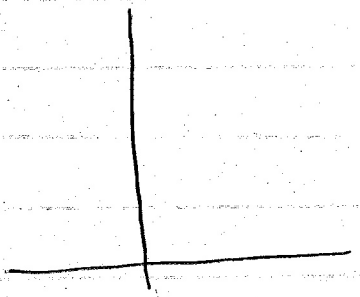


Not all races are bugs, but most of them are.



x	y
0	0
0	1
1	0
1	1

could be diff. values



When ppl reading thing \rightarrow NOT a problem.

Problems occur when at least 1 is writing.

when you write \rightarrow change state.

Soln: \textcircled{D} Use locks. (not only way of dealing it).

mechanism to control interleavings. For us, they are abstractions.

lock(l)	lock
x = 1	read x
y = 1	read y
unlock(l)	

what happens? not solved problem.

lock(l)	lock(l)	2 possibilities: read 0 is or read 1 is.
x = 1	read x	
y = 1	read y	
unlock(l)	unlock(l)	

Simon's Paper: (2005).

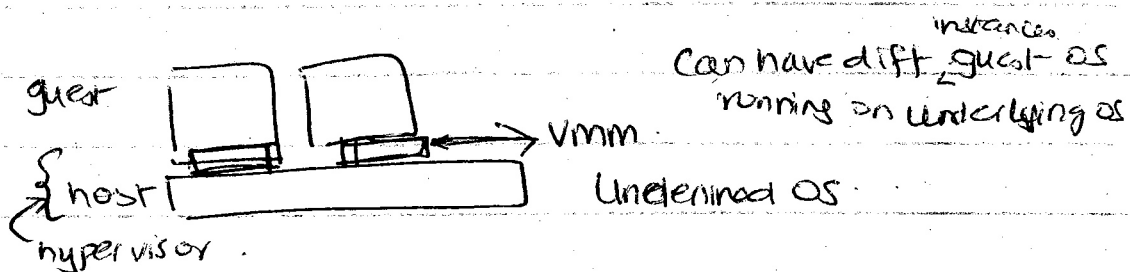
- Use hierarchies.

- Hierarchy \rightarrow collection of components ^{organized by strength of interaction} (see paper).

- Evolve naturally because stable system.

Virtual machines:

1. Terminology:



VMM: responsible for all virtualization.

Takes care of CPU virtualization & memory virtualization.

CPU virtualization:

making it seem like sharing single processors among different processes simultaneously.

① Trap & emulate:

(VMM + host = hypervisor)

take all guest code run on processor.

Caveat: run in user mode.

Not prob. for applications. But for guest kernel mode \rightarrow ~~it~~ has to run in kernel code. Results in traps.

At this pt. VMM comes & emulates it.

It doesn't work in x86 because of some instructions like popf.

fails silently.

then run popf in user mode (doesn't change flag, but doesn't trap)
in kernel mode (changes flags / ~~interrupt~~ update
sets eflags, interrupt) in kernel flag

The signalling mechanism is not present.

Binary Translation:

- ① Rewrite problem instructions like "popf"
- ② All others copy them identically. (IDENT)
- ③ Rewrite control flow instructions.

Translated code lives in ^{Region of memory} Translated cache that belongs to VMM.

TC → only guest kernel mode instructions are translated.

This is what processor executes.

Control flow e.g. jmp

xor
popf
jmp 24

In previous step rewriting

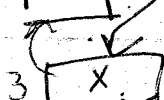
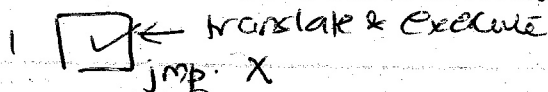
xor ✓

popf == 3 popf translated to 3 times

now line 24 doesn't correspond to it

24 "x"

Thus you write jmp to point to correct address. Why dynamically?



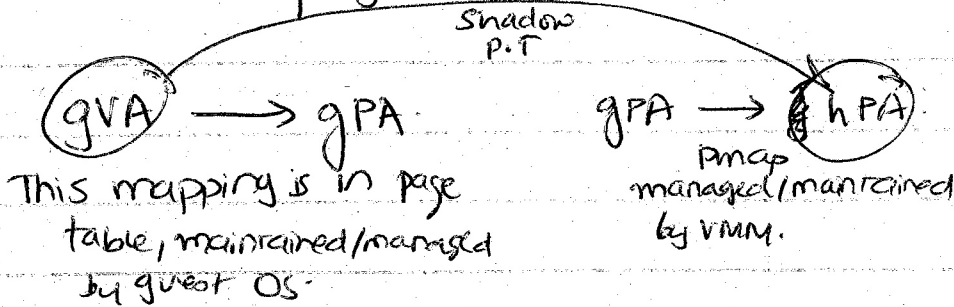
You translate 1, 3, 2.

HW Support: User, Kernel bit + VMM mode bit

2 levels of two page tables instead of one

Memory Virtualization:

Use shadow page tables.



Shadow PT created by VMM, actual ^{HW} pg. Tb ptr points to it.

in TLB, Shadow pg table will be cached.

Pg. Table: = Guest phy. page "X"

GVA	1	7	GPA	GPA	Pmap	HPA
	2	3		7	—	11
	0	"X"		6	—	5
				3	—	4
				X	—	9

Shadow Pg. Table

1	—	11
2	—	4
0	—	9

1 bit

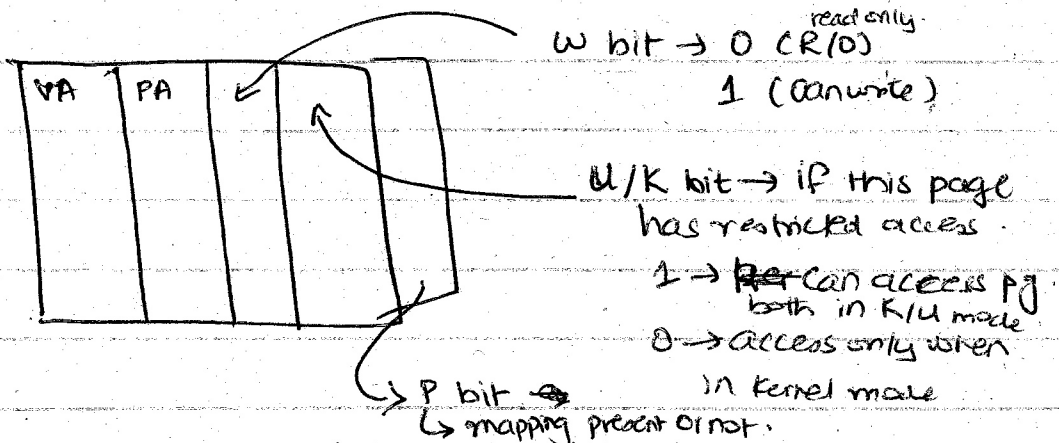
Problem: If guest changes 1-11 to 1-6
mapping in STP also needs to be updated. 1-5.

We use traces for that ← Signalling mechanism for
when pages need to be updated. Marks the
entry corresponding to the page where page table lies in.

mark 0-y as \textcircled{R} takes it as a trap.

mark entry for page that contains page table.

i.e. 0 → y. everytime try to write to y which is x which
is where there is page table, update it.

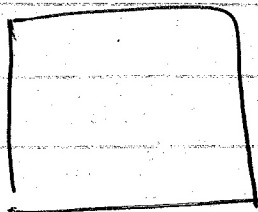


Guest OS believes that it has U/K bit.

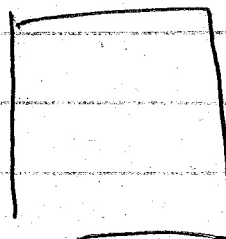
HW doesn't know abt it. That's why have

2 pg tables

User STP



Kernel STP



2011 Q5.

Guest pg. table = pg. 4

gVA	gPA	P	w	V
0	5	1	1	1
1	9	0	1	1
2	11	1	0	1
3	4	1	1	0
4	7	1	1	0

maintained by guest OS

pmap:

gPA	hPA
4	2
5	6
6	4
7	3
11	8

maintained by host OS.

when virtual bit set to user → user mode STP

else kernel STP.

User Shadow PT: → present bit

gVA	hPA	P	w	u
0	6	1	1	1
1	-	0	-	-
3	-	0	-	-

(no mapping)

if present bit is 0, user bit set to 0 so when

VA
9
3
4
for Any acc
the host
is never
In any

Kernel SPT

g VA	h PA	P	W	U
9	8	1	0	1
3	2	1	0	1
4	3	1	1	1

acquired
by guest OS

lives in host OS, guest kernel is running in user mode & it is to host OS.

can't set to 1

as 3 maps to

4. but 4 is the pg. table entry so using trace set to 0.

For Any access is a user access to the host OS. For kernel, user bit is never 0.

In any SPT, user bit is 1.