

Do Work: Allocates + frees memory  
on per-thread heaps

Is This Better as one Process or two?

# Virtual Memory

A Huge Step Backwards

Daniel Bittman (he/him)

HPTS 2024

# Virtual Memory is Evil

Semantics

Performance

Security

Flat, Uniform,

Contiguous, infinite,

Isolated, Fast

Flat, Uniform,

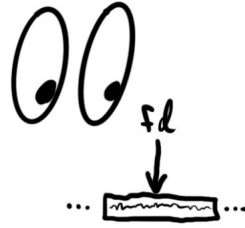
non-uniform Access latency  
Page faults, discontinuous allocation,  
holes, Varying memory coherence,

Contiguous, infinite,

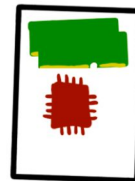
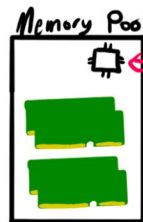
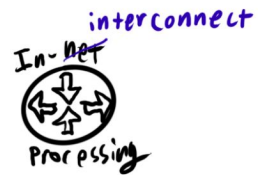
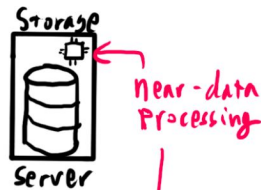
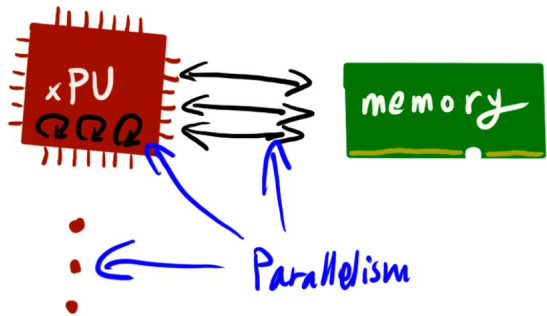
infallible malloc, Spatial safety, temporal  
safety, buffer overflow, ROP,

Isolated, Fast

TLB miss, shutdown, meltdown,  
Cross-machine coherence, local execution  
context, Page table walking

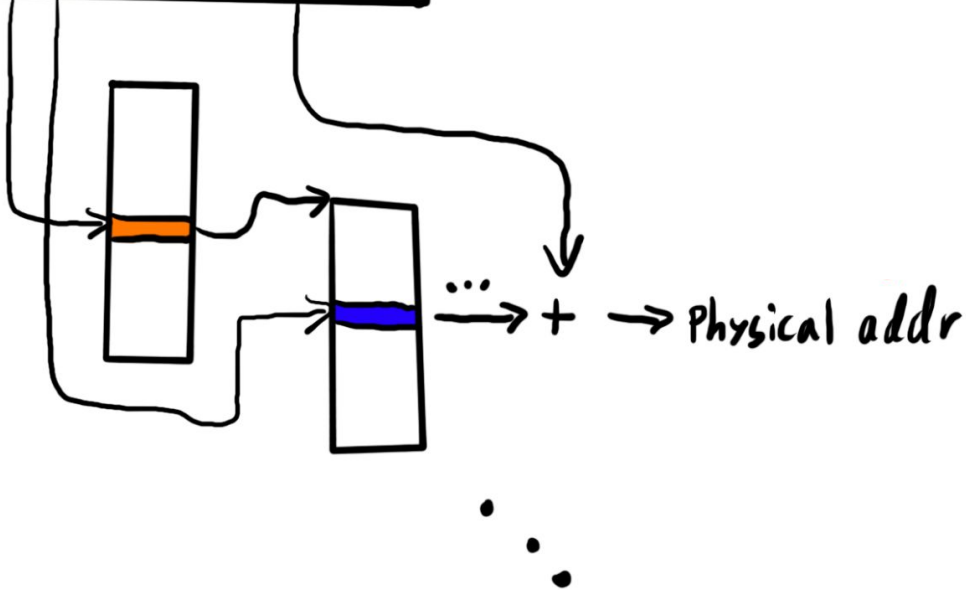
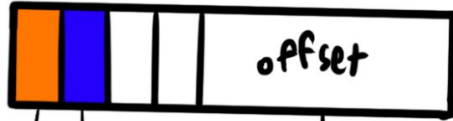


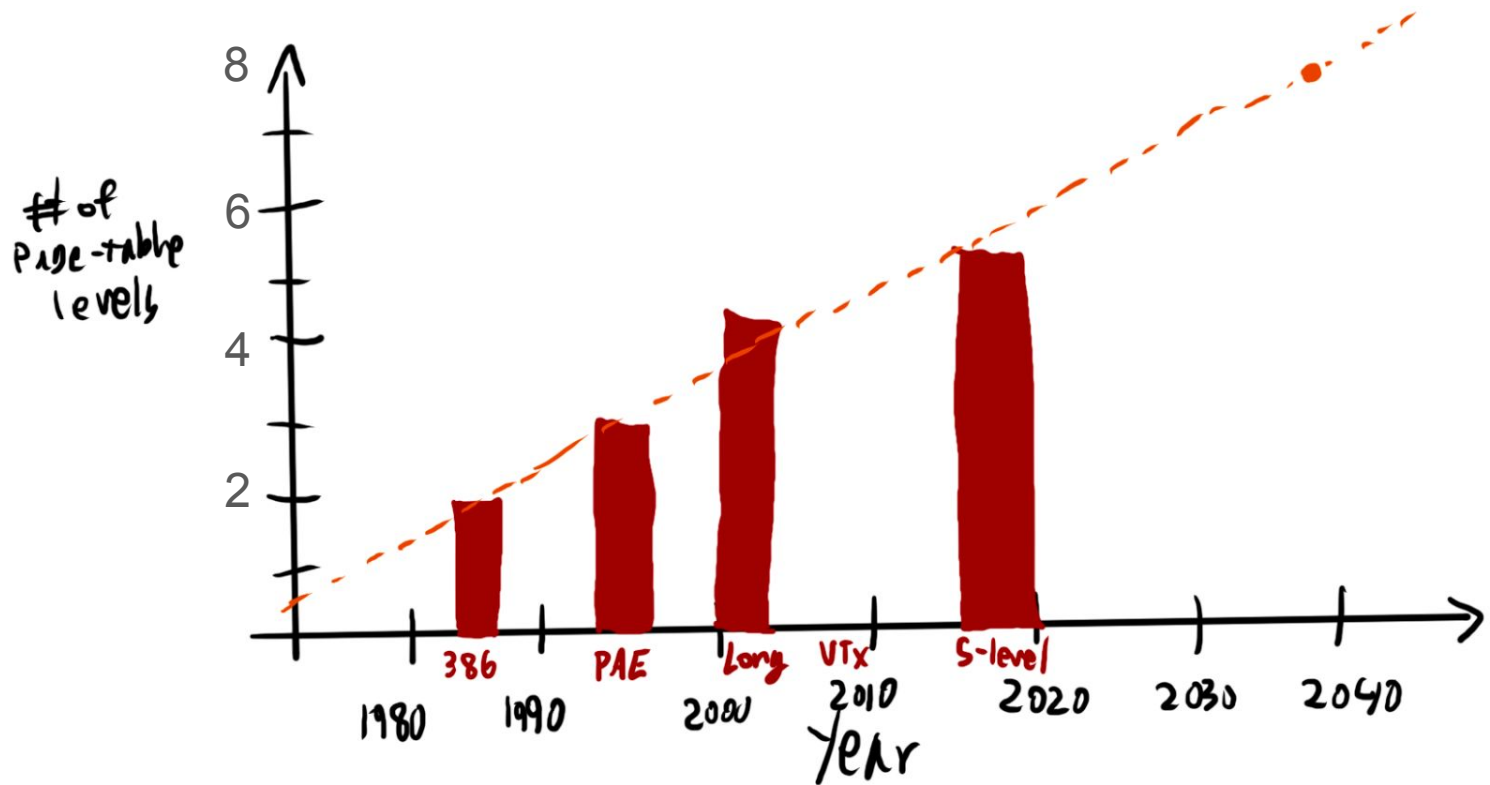


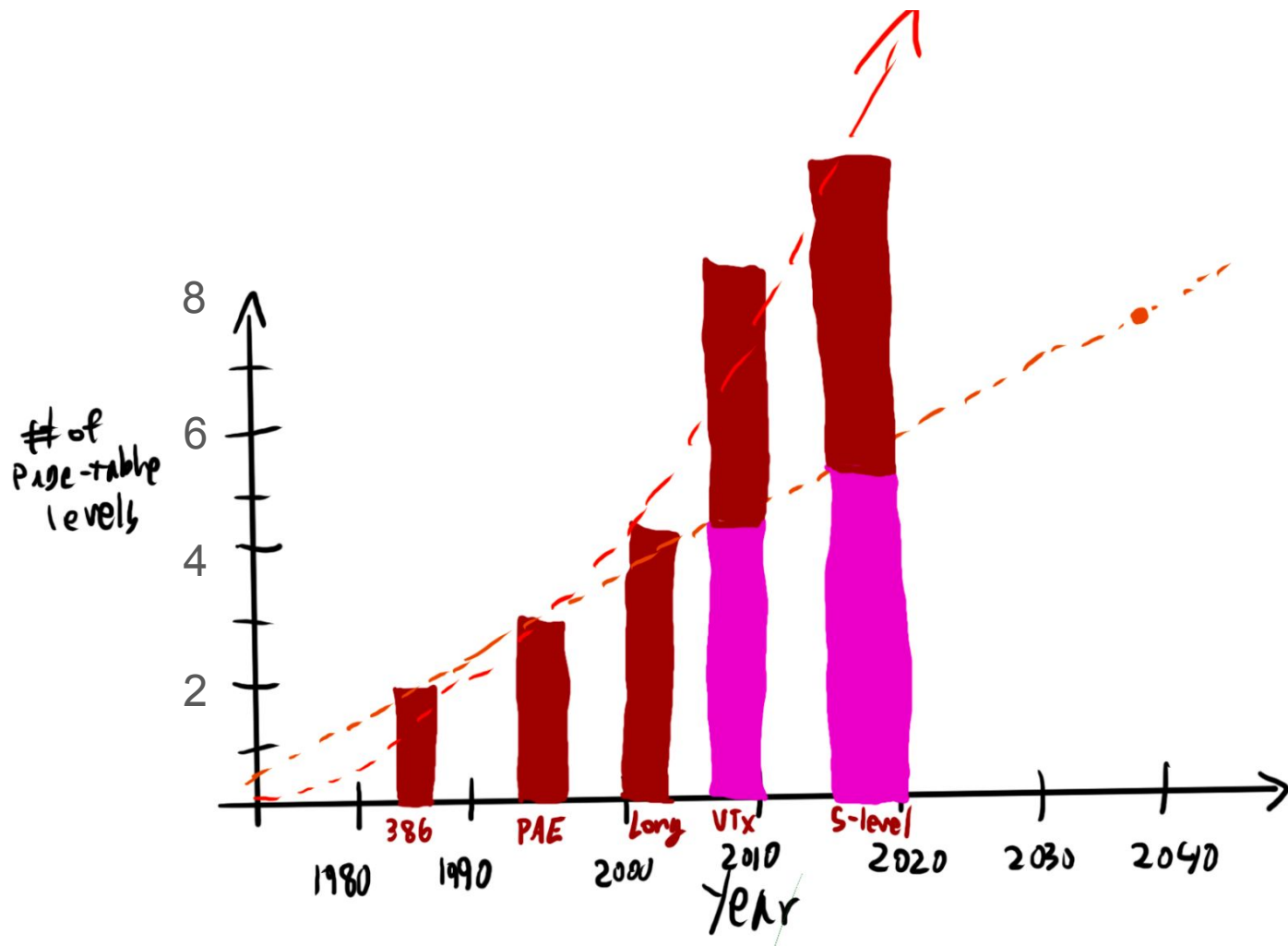




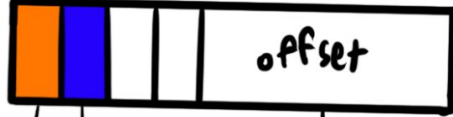
Virtual Addr



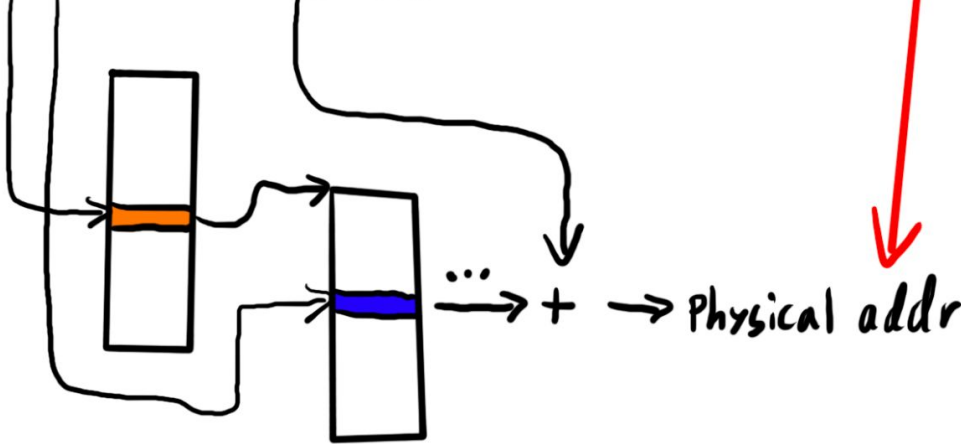




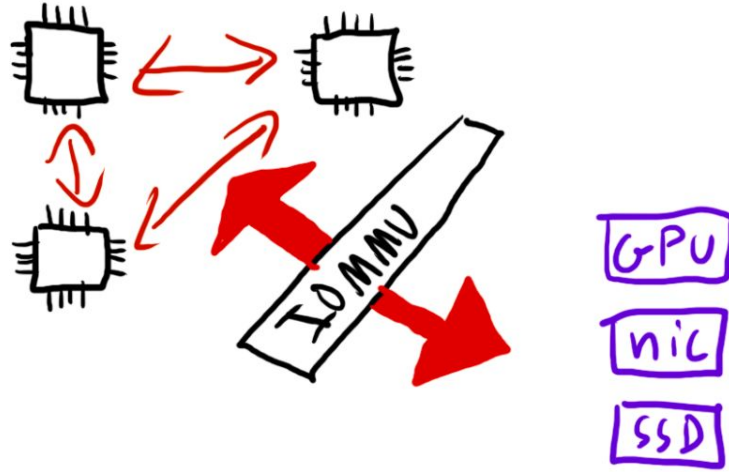
Virtual Addr



translate va  
Page Walk



...



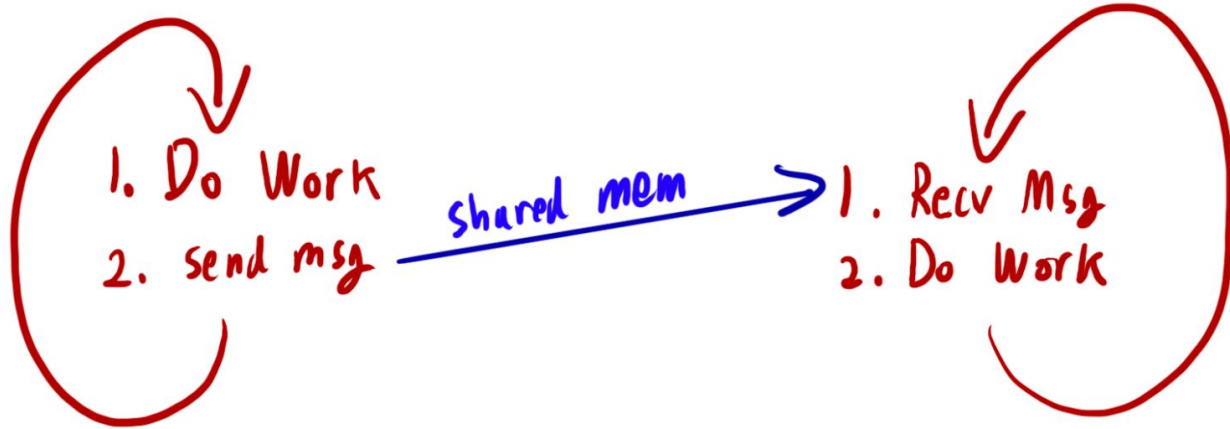
THIS

SUCKS

# TLB Coherence

Linux TLB ~~optimizations~~ <sup>Complications</sup>  
>10k LoC

CPU 0 ← 2 Threads → CPU 1

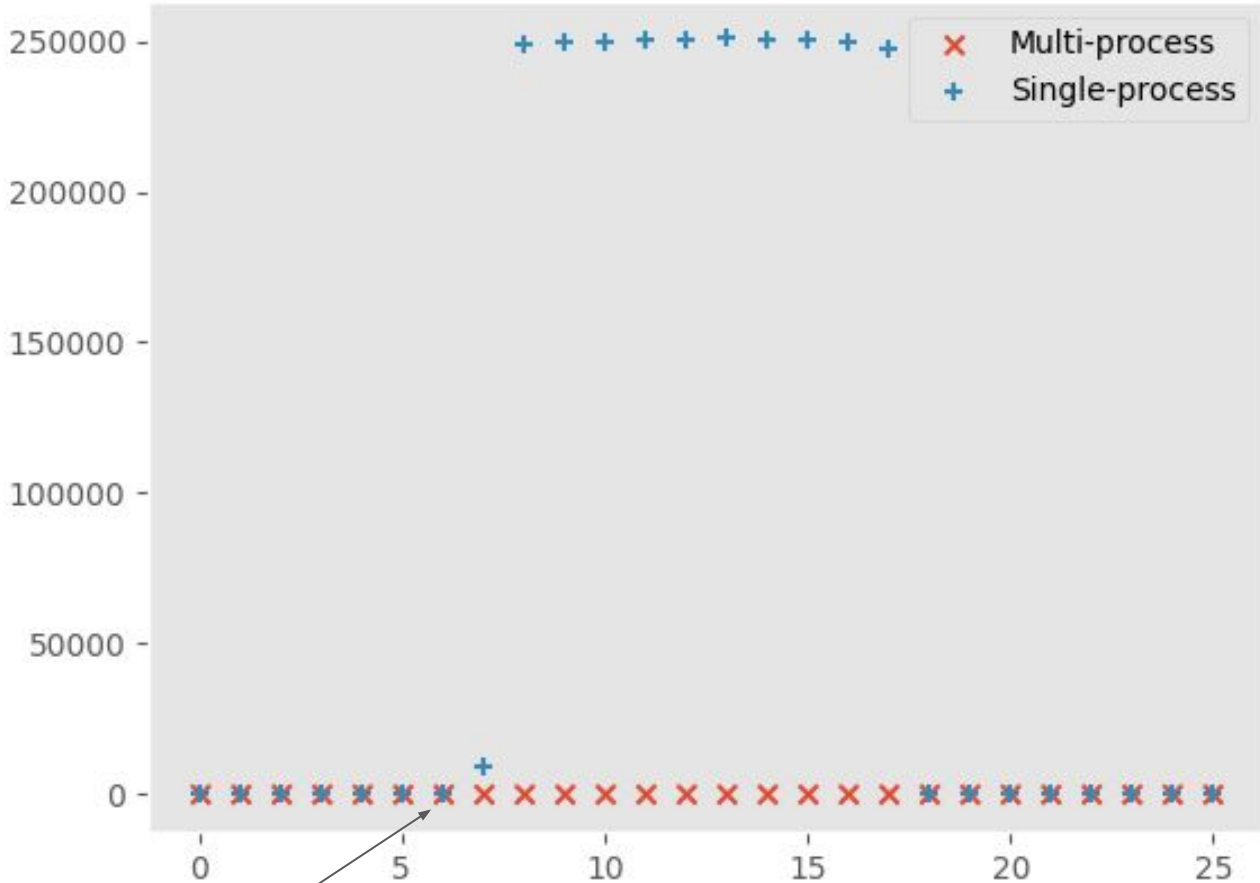


Do Work: Allocates + frees memory on per-thread heaps  
calls munmap (2)

Is This Better as one Process or two?



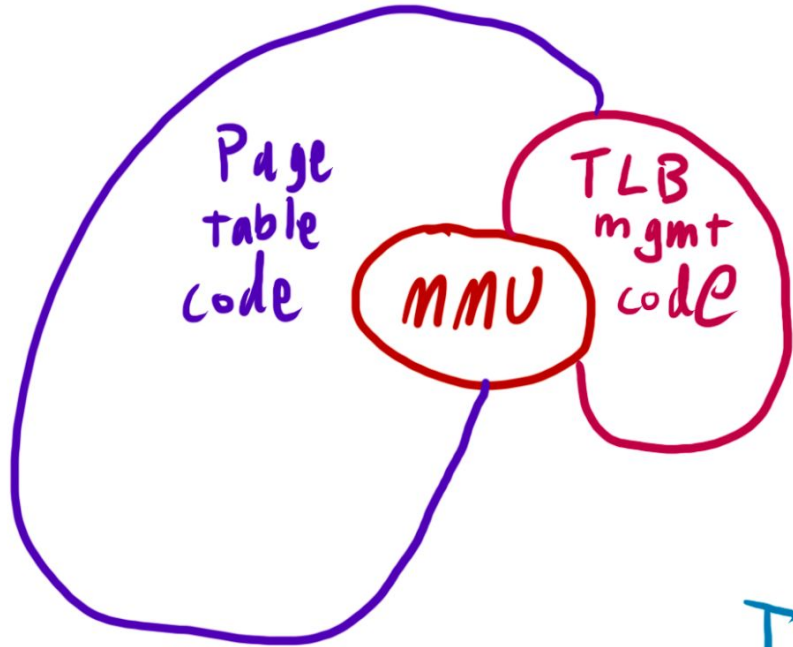
TLB Shootdowns /s



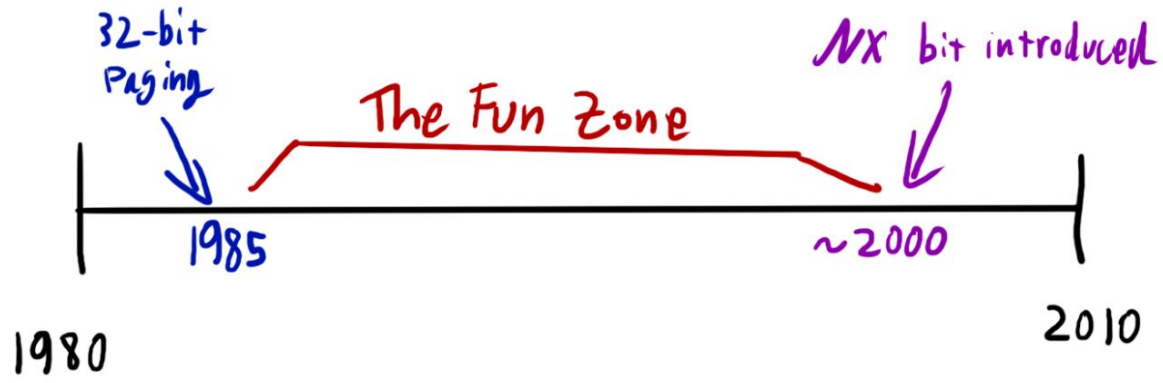
runs start

time (s)

# Attack Surface



~~Io MMU~~



Infinite Memory  
and  
Infallible Malloc

Infinite Memory  
and  
Infallible Malloc

vm.overcommit\_memory = thoughts and prayers

Infinite Memory  
and

Infallible Malloc

How many libc functions  
allocate memory?

```
char* buffer = NULL;
while(buffer == NULL){
    buffer = malloc(sizeof(char) * 100);
}
```

```
char* buffer = NULL;
while(buffer == NULL){
    buffer = malloc(sizeof(char) * 100);
}
```

→ malloc never fails

→ oh, but it's UB not to check

→ wait it totally does fail

∴ Don't worry about it, it's impossible to recover





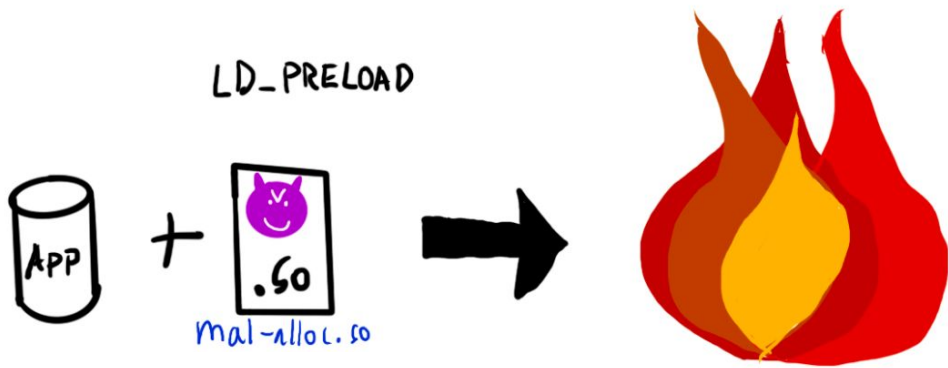
+



mal-alloc.so

LD\_PRELOAD





```
void *malloc(size_t size) {  
    if (should_i_fail lol (size)) return NULL;  
    // call real malloc  
}
```

	Doesn't Lose Data	Doesn't SIGSEGV
Redis		
Postgres		
Mongo		
SQLite		
RustC		
Is		
Python		
SSH		

	Doesn't Lose Data	Doesn't SIGSEGV	Doesn't Exit	Stays Available
Redis				
Postgres				
Mongo				
SQLite				
RustC				
Is				
Python				
SSH				

	Doesn't Lose Data	Doesn't SIGSEGV	Doesn't Exit	Stays Available
Redis	✓			
Postgres	✓			
Mongo	✓			
SQLite	✓			
RustC	?			
Is	N/A			
Python	?			
SSH	✓			

	Doesn't Lose Data	Doesn't SIGSEGV	Doesn't Exit	Stays Available
Redis	✓	X	X	X
Postgres	✓	X	X*	X
Mongo	✓	✓	X	X
SQLite	✓	✓	✓	✓
Rust C	?	✓	X	X
Is	N/A	✓	X	X
Python	?	X	X	X
SSH	✓	✓	X	X

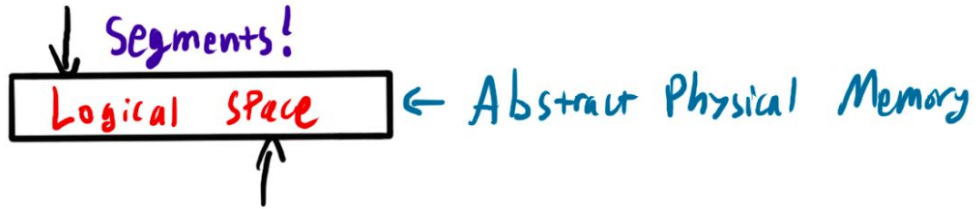
LD\_PRELOAD = libmal-alloc.so sh

( ✓ ✓ ✗ ✗ )

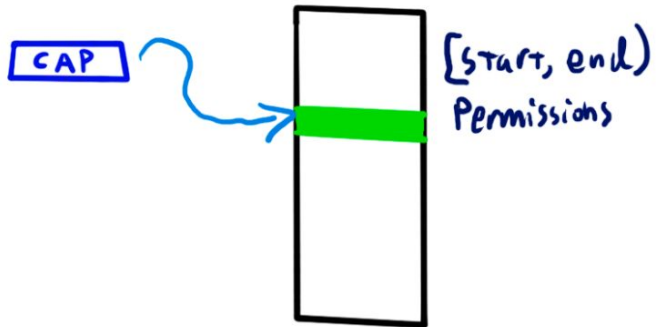








CHERI



# Allocating Memory SUCKS

my\_poor\_api () → Result<..., AllocError>

So,

a\_more\_correct\_expression\_of\_the\_api

Dont.

Twizzler

Qmut T  
&T

This reasoning is needed for  
C, C++ too

Depending on appending

Precise  
Application  
Semantics



Immutability changes everything



Safe, fast APIs for shared memory

→ Virtual Memory has bad Security + Semantics  
Hidden Complexity is still complexity!

→ Virtual Memory has huge Performance overhead  
And sometimes surprising overhead!

→ Allocating memory can fail  
And Freeing memory can send messages!

→ Intel Page table depth followed Moore's Law ~2008

→ The OS + DB communities should go after the CPU designers

→ Just run SPDK on baremetal x86 segment mode.



Thank you,

Any Questions?



Twizzler.io  
dbittman.com

Daniel Bittman (he/him)