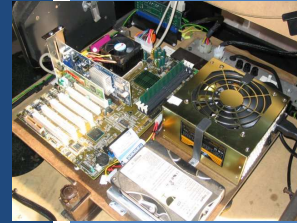


TCSS 422: OPERATING SYSTEMS

Memory Virtualization, Segmentation, Memory Paging



Wes J. Lloyd
School of Engineering and Technology,
University of Washington - Tacoma

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

FEEDBACK FROM 11/14


- How to kill all child threads with a `pthread_cond_broadcast()` ?
- At end of the program, some threads (producers or consumers) may be asleep waiting on a signal.
- For consumers, there are no more matrices being produced, so there is no signal for “consumption”
- Need some way to shutdown/end the program
- Can leverage when producer threads finish their work
- Producers last “signal” can be a “broadcast” to awaken all consumers to evaluate special “end of program” state variable.

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.2

OBJECTIVES		
<ul style="list-style-type: none">▪ Program 2▪ Program 3▪ Memory Virtualization<ul style="list-style-type: none">▪ Chapter 14 – The Memory API▪ Chapter 15 – Address Translation▪ Segments<ul style="list-style-type: none">▪ Chapter 16 – Segmentation▪ Chapter 17 – Free Space Management▪ Paging<ul style="list-style-type: none">▪ Chapter 18 – Introduction to Paging▪ Chapter 19 – Translation Lookaside Buffer		
November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.3

<h1>CHAPTER 14: THE MEMORY API</h1>		
November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.4

MALLOC

```
#include <stdlib.h>

void* malloc(size_t size)
```

- Allocates memory on the heap
- `size_t` unsigned integer (must be +)
- `size` size of memory allocation in bytes

- Returns
 - SUCCESS: A void * to a memory address
 - FAIL: NULL

- `sizeof()` often used to ask the system how large a given datatype or struct is

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.5
-------------------	---	-------

sizeof()

- Not safe to assume data type sizes using different compilers, systems
- Dynamic array of 10 ints
- Static array of 10 ints

```
int *x = malloc(10 * sizeof(int));
printf("%d\n", sizeof(x));
```

4

```
int x[10];
printf("%d\n", sizeof(x));
```

40

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.6
-------------------	---	-------

FREE()

```
#include <stdlib.h>
void free(void* ptr)
```

- Free memory allocated with malloc()
- Provide: (void *) ptr to malloc'd memory

- Returns: nothing

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.7

```
#include<stdio.h>
```

```
int * set_magic_number_a()
{
    int a =53247;
    return &a;
}

void set_magic_number_b()
{
    int b = 11111;
}

int main()
{
    int * x = NULL;
    x = set_magic_number_a();
    printf("The magic number is=%d\n",*x);
    set_magic_number_b();
    printf("The magic number is=%d\n",*x);
    return 0;
}
```

What will this code do?

8

```
#include<stdio.h>

int * set_magic_number_a()
{
    int a =53247;
    return &a;
}

void set_magic_number_b()
{
    int b = 11111;
}

int main()
{
    int * x = NULL;
    x = set_magic_number_a();
    printf("The magic number is=%d\n",*x);
    set_magic_number_b();
    printf("The magic number is=%d\n",*x);
    return 0;
}
```

What will this code do?

Output:
\$./pointer_error
The magic number is=53247
The magic number is=11111

We have not changed *x but the value has changed!!

Why?

9

DANGLING POINTER (1/2)

- Dangling pointers arise when a variable referred (a) goes “out of scope”, and it’s memory is destroyed/overwritten (by b) without modifying the value of the pointer (*x).
- The pointer still points to the original memory location of the deallocated memory (a), which has now been reclaimed for (b).

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.10
-------------------	---	--------

DANGLING POINTER (2/2)

- Fortunately in the case, a compiler warning is generated:

```
$ g++ -o pointer_error -std=c++0x pointer_error.cpp
```

```
pointer_error.cpp: In function 'int*  
set_magic_number_a()':  
pointer_error.cpp:6:7: warning: address of local  
variable 'a' returned [enabled by default]
```

- This is a common mistake - - -
accidentally referring to addresses that have
gone “out of scope”

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.11

CALLOC()

```
#include <stdlib.h>  
  
void *calloc(size_t num, size_t size)
```

- Allocate “C”lear memory on the heap
- Calloc wipes memory in advance of use...
- `size_t num` : number of blocks to allocate
- `size_t size` : size of each block(in bytes)

- Calloc() prevents...

```
char *dest = malloc(20);  
printf("dest string=%s\n", dest);
```

dest string=??F

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.12

REALLOC()

```
#include <stdlib.h>

void *realloc(void *ptr, size_t size)
```

- Resize an existing memory allocation
- Returned pointer may be same address, or a new address
 - New if memory allocation must move
- `void *ptr`: Pointer to memory block allocated with `malloc`, `calloc`, or `realloc`
- `size_t size`: New size for the memory block(in bytes)
- EXAMPLE: `realloc.c`
- EXAMPLE: `nom.c`

November 19, 2018

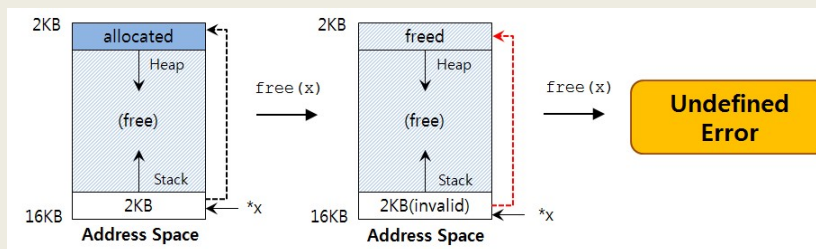
TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.13

DOUBLE FREE

```
int *x = (int *)malloc(sizeof(int)); // allocated
free(x); // free memory
free(x); // free repeatedly
```

- Can't deallocate twice
- Second call core dumps



November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.14

SYSTEM CALLS


- **brk(), sbrk()**
 - Used to change data segment size (the end of the heap)
 - Don't use these

- **Mmap(), munmap()**
 - Can be used to create an extra independent “heap” of memory for a user program

- See man page

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.15
-------------------	---	--------

CHAPTER 15: ADDRESS TRANSLATION



November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.16
-------------------	---	--------

OBJECTIVES

- Address translation
- Base and bounds
- HW and OS Support
- Memory segments
- Memory fragmentation

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.17
-------------------	---	--------

ADDRESS TRANSLATION

- 64KB Address space example
- Translation: mapping virtual to physical

The diagram illustrates the mapping of a 64KB virtual address space to physical memory. On the left, the 'Virtual mapping' shows a 64KB 'Address Space' divided into segments: Program Code (0KB-16KB), Heap (16KB-32KB), heap (free) (32KB-48KB), stack (48KB-64KB), and Stack (64KB-16KB). On the right, 'Physical Memory' shows the OS (0KB-16KB), a gap (16KB-32KB), the 'Relocated Process' (32KB-48KB) containing Code, Heap, and Stack, and another gap (48KB-64KB). Dotted lines indicate the mapping from virtual to physical memory.

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.18
-------------------	---	--------

BASE AND BOUNDS

- Dynamic relocation
- Two registers base & bounds: **on the CPU**
- OS places program in memory
- Sets base register

$$physical\ address = virtual\ address + base$$

- Bounds register
 - Stores size of program address space (16KB)
- OS verifies that every address:

$$0 \leq virtual\ address < bounds$$

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.19
-------------------	---	--------

INSTRUCTION EXAMPLE

`128 : movl 0x0(%ebx), %eax`

- Base = 32768
- Bounds = 16384
- Fetch instruction at 128 (virt addr) ↑
 - Phy addr = virt addr + base reg
 - 32896 = 128 + 32768 (base)
- Execute instruction
 - Load from address (var x is @ 15kb=15360)
 - 48128 = 15360 + 32768 (base) -- found x...
- Bounds register: terminate process if
 - **ACCESS VIOLATION:** Virtual address > bounds reg

$$physical\ address = virtual\ address + base$$

0KB	128	movl 0x0(%ebx), %eax
	132	Addl 0x03,%eax
1KB	135	movl %eax,0x0(%ebx)
Program Code		
Heap		
2KB		
3KB		
4KB		
↓ heap		
(free)		
↑ stack		
14KB		
15KB	3000	Int x
16KB		Stack

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.20
-------------------	---	--------

MEMORY MANAGEMENT UNIT

- MMU
 - Portion of the CPU dedicated to address translation
 - Contains base & bounds registers
- Base & Bounds Example:
 - Consider address translation
 - 4 KB (4096 bytes) address space, loaded at 16 KB physical location

Virtual Address	Physical Address
0	16384
1024	17408
3000	19384
FAULT 4400	20784 (out of bounds)

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.21
-------------------	---	--------

DYNAMIC RELOCATION OF PROGRAMS

- Hardware requirements:

Requirements	HW support
Privileged mode	CPU modes: kernel, user
Base / bounds registers	Registers to support address translation
Translate virtual addr; check if in bounds	Translation circuitry, check limits
Privileged instruction(s) to update base / bounds regs	Instructions for modifying base/bound registers
Privileged instruction(s) to register exception handlers	Set code pointers to OS code to handle faults
Ability to raise exceptions	For out-of-bounds memory access, or attempts to access privileged instr.

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.22
-------------------	---	--------

OS SUPPORT FOR MEMORY VIRTUALIZATION

- For base and bounds OS support required
 - When process starts running
 - Allocate address space in physical memory
 - When a process is terminated
 - Reclaiming memory for use
 - When context switch occurs
 - Saving and storing the base-bounds pair
 - Exception handlers
 - Function pointers set at OS boot time

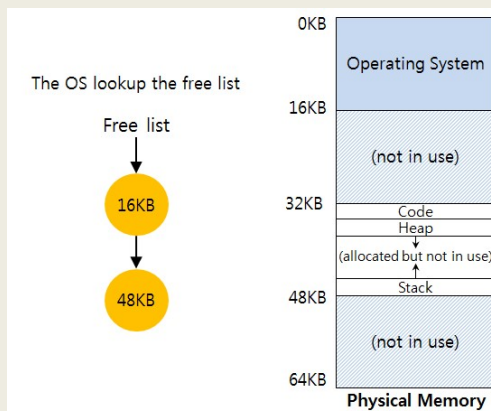
November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.23

OS: WHEN PROCESS STARTS RUNNING

- OS searches for free space for new process
 - Free list: data structure that tracks available memory slots



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.24

OS: WHEN PROCESS IS TERMINATED

- OS places memory back on the free list

Physical Memory **Free list** **Physical Memory**

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.25
-------------------	---	--------

OS: WHEN CONTEXT SWITCH OCCURS

- OS must save base and bounds registers
 - Saved to the Process Control Block PCB (task_struct in Linux)

Physical Memory Context Switching **Physical Memory**

Process A PCB
...
base : 32KB
bounds : 48KB
...

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.26
-------------------	---	--------

DYNAMIC RELOCATION

- OS can move process data when not running
 1. OS deschedules process from scheduler
 2. OS copies address space from current to new location
 3. OS updates PCB (base and bounds registers)
 4. OS reschedules process

- When process runs new base register is restored to CPU

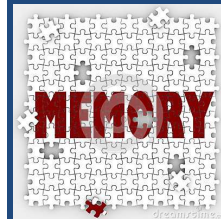
- **Process doesn't know it was even moved!**

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.27

CHAPTER 16: SEGMENTATION



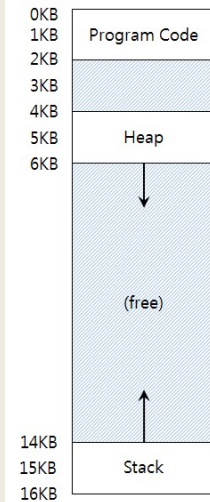
November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.28

BASE AND BOUNDS INEFFICIENCIES

- **Address space**
 - Contains significant unused memory
 - Is relatively large
 - Preallocates space to handle stack/heap growth
- **Large address spaces**
 - Hard to fit in memory
- **How can these issues be addressed?**



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.29

MULTIPLE SEGMENTS

- **Memory segmentation**
- **Address space has (3) segments**
 - Contiguous portions of address space
 - Logically separate segments for: code, stack, heap
- **Each segment can be placed separately**
- **Track base and bounds for each segment (registers)**

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.30

SEGMENTS IN MEMORY

■ Consider 3 segments:

Segment	Base	Size
Code	32K	2K
Heap	34K	2K
Stack	28K	2K

Much smaller
↓

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.31

ADDRESS TRANSLATION: CODE SEGMENT

$physical\ address = offset + base$

■ Code segment - physically starts at 32KB (base)
 ■ Starts at "0" in virtual address space

November 19, 2018

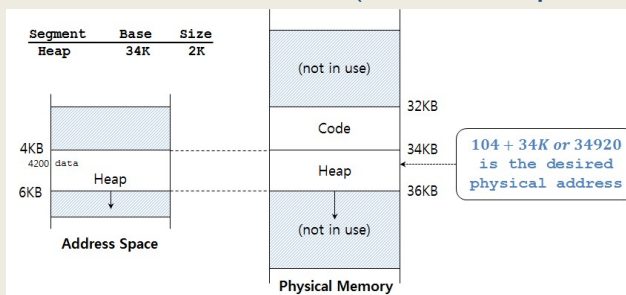
TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.32

ADDRESS TRANSLATION: HEAP

Virtual address + base is not the correct physical address.

- Heap starts at virtual address 4096
- The data is at 4200
- Offset= $4200 - 4096 = 104$ (virt addr - virt heap start)
- Physical address = $104 + 34816$ (offset + heap base)



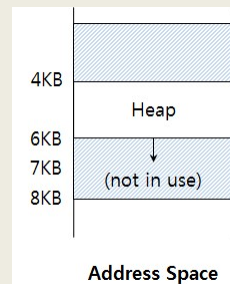
November 19, 2018

TCSS422: Operating Systems [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma

L14.33

SEGMENTATION FAULT

- Access beyond the address space
- Heap starts at virtual address: 4096
- Data pointer is to 7KB (7168)
- Is data pointer valid?
- Heap starts at $4096 + 2 \text{ KB seg size} = 6144$
- Offset= $7168 > 4096 + 2048 (6144)$



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma

L14.34

SEGMENT REGISTERS

- Used to dereference memory during translation

13	12	11	10	9	8	7	6	5	4	3	2	1	0
Segment		Offset											

- First two bits identify segment type
- Remaining bits identify memory offset
- Example: virtual heap address 4200 (01000001101000)

13	12	11	10	9	8	7	6	5	4	3	2	1	0
0	1	0	0	0	0	0	1	1	0	1	0	0	0
Segment		Offset											

Segment	bits
Code	00
Heap	01
Stack	10
-	11

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.35

SEGMENTATION DEREFERENCE

```

1 // get top 2 bits of 14-bit VA
2 Segment = (VirtualAddress & SEG_MASK) >> SEG_SHIFT
3 // now get offset
4 Offset = VirtualAddress & OFFSET_MASK
5 if (Offset >= Bounds[Segment])
6     RaiseException(PROTECTION_FAULT)
7 else
8     PhysAddr = Base[Segment] + Offset
9     Register = AccessMemory(PhysAddr)
    
```

- VIRTUAL ADDRESS = 01000001101000 (on heap)
- SEG_MASK = 0x3000 (11000000000000)
- SEG_SHIFT = 01 → *heap* (mask gives us segment code)
- OFFSET_MASK = 0xFFF (00111111111111)
- OFFSET = 000001101000 = 104 (isolates segment offset)
- OFFSET < BOUNDS : 104 < 2048

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.36

STACK SEGMENT

- Stack grows backwards (FILO)
- Requires hardware support:
- Direction bit: tracks direction segment grows

Segment Register(with Negative-Growth Support)

Segment	Base	Size	Grows	Positive?
Code	32K	2K	1	1
Heap	34K	2K	1	1
Stack	28K	2K	0	0

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.37

SHARED CODE SEGMENTS

- Code sharing: enabled with HW support
- Supports storing shared libraries in memory only once
- DLL: dynamic linked library
- .so (linux): shared object in Linux (under /usr/lib)
- Many programs can access them
- Protection bits: track permissions to segment

Segment Register Values(with Protection)

Segment	Base	Size	Grows	Positive?	Protection
Code	32K	2K	1	1	Read-Execute
Heap	34K	2K	1	1	Read-Write
Stack	28K	2K	0	0	Read-Write

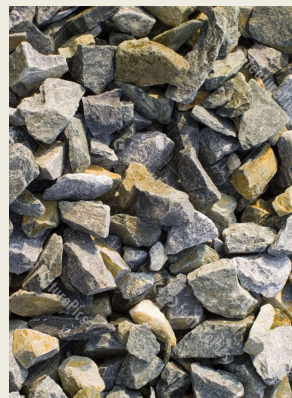
November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.38

SEGMENTATION GRANULARITY

- Coarse-grained
- Manage memory as large purpose based segments:
 - Code segment
 - Heap segment
 - Stack segment



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.39

SEGMENTATION GRANULARITY - 2

- Fine-grained
- Manage memory as list of segments
- Code, heap, stack segments composed of multiple smaller segments
- Segment table
 - On early systems
 - Stored in memory
 - Tracked large number of segments



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.40

MEMORY FRAGMENTATION

- Consider how much free space?
- We'll say about 24 KB

- Request arrives to allocate a 20 KB heap segment

- Can we fulfil the request for 20 KB of contiguous memory?

Not compacted

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.41
-------------------	---	--------

COMPACTION

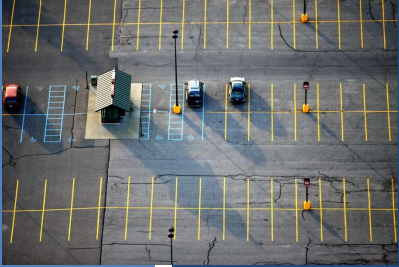
- Supports rearranging memory
- Can we fulfil the request for 20 KB of contiguous memory?

- **Drawback: Compaction is slow**
 - Rearranging memory is time consuming
 - 64KB is fast
 - 4GB+ ... slow

- **Algorithms:**
 - Best fit: keep list of free spaces, allocate the most snug segment for the request
 - Others: worst fit, first fit... (in future chapters)

Compacted

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.42
-------------------	---	--------



CHAPTER 17: FREE SPACE MANAGEMENT

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.43

FREE SPACE MANAGEMENT

- Management of memory using
 - Only fixed-sized units
 - Easy: keep a list
 - Memory request → return first free entry
 - Simple search
 - With variable sized units
 - More challenging
 - Results from variable sized malloc requests
 - Leads to fragmentation

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.44
-------------------	---	--------

FRAGMENTATION

- Consider a 30-byte heap

30-byte heap:

free	used	free
0	10	20
0	10	30

- Request for 15-bytes

free list: head →

addr:0 len:10	→	addr:20 len:10	→	NULL
------------------	---	-------------------	---	------

- Free space: 20 bytes
- No available contiguous chunk → return NULL

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.45
-------------------	---	--------

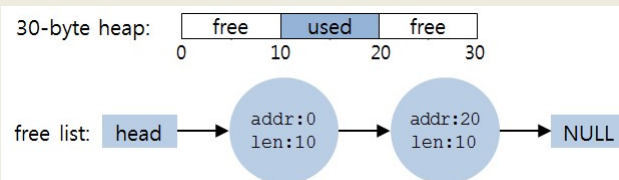
FRAGMENTATION - 2

- **External:** OS can compact
 - Example: Client asks for 100 bytes: malloc(100)
 - OS: No 100 byte contiguous chunk is available: returns NULL
 - Memory is externally fragmented - - Compaction can fix!
- **Internal:** lost space - OS can't compact
 - OS returns memory units that are too large
 - Example: Client asks for 100 bytes: malloc(100)
 - OS: Returns 125 byte chunk
 - Fragmentation is *in* the allocated chunk
 - Memory is lost, and unaccounted for - can't compact

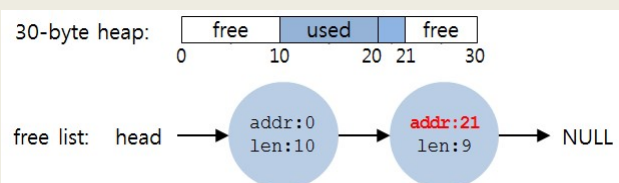
November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.46
-------------------	---	--------

ALLOCATION STRATEGY: SPLITTING

- Request for 1 byte of memory: `malloc(1)`



- OS locates a free chunk to satisfy request
- Splits chunk into two, returns first chunk



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.47

ALLOCATION STRATEGY: COALESCING

- Consider 30-byte heap
- `Free()` frees all 10 bytes segments (*list of 3-free 10-byte chunks*)



- Request arrives: `malloc(30)`
- SPLIT DOES NOT WORK** - no contiguous 30-byte chunk exists!
- Coalescing regroups chunks into contiguous chunk



- Allocation can now proceed
- Coalescing is defragmentation of the free space list

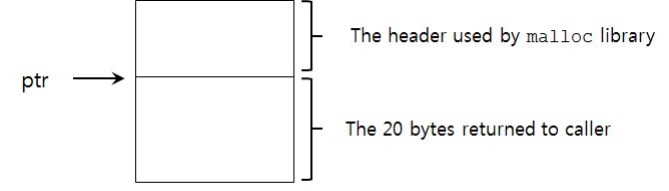
November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.48

MEMORY HEADERS

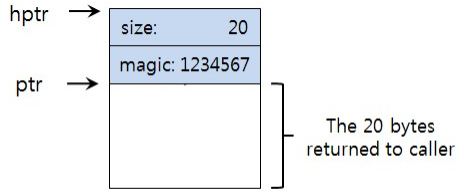
- `free(void *ptr)`: Does not require a size parameter
- How does the OS know how much memory to free?
- Header block
 - Small descriptive block of memory at start of chunk



An Allocated Region Plus Header

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.49
-------------------	---	--------

MEMORY HEADERS - 2



Specific Contents Of The Header

```
typedef struct __header_t {  
    int size;  
    int magic;  
} header_t;
```

A Simple Header

- Contains size
- Pointers: for faster memory access
- Magic number: integrity checking

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.50
-------------------	---	--------

MEMORY HEADERS - 3

- Size of memory chunk is:
 - Header size + user malloc size
 - N bytes + sizeof(header)
- Easy to determine address of header

```
void free(void *ptr) {  
    header_t *hptr = (void *)ptr - sizeof(header_t);  
}
```

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.51

THE FREE LIST

- Simple free list struct
- Use mmap to create free list
- 4kb heap, 4 byte header, one contiguous free chunk

```
typedef struct __node_t {  
    int size;  
    struct __node_t *next;  
} node_t;
```

```
// mmap() returns a pointer to a chunk of free space  
node_t *head = mmap(NULL, 4096, PROT_READ|PROT_WRITE,  
                    MAP_ANON|MAP_PRIVATE, -1, 0);  
head->size = 4096 - sizeof(node_t);  
head->next = NULL;
```

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.52

FREE LIST - 2

- Create and initialize free-list “heap”

```
// mmap() returns a pointer to a chunk of free space
node_t *head = mmap(NULL, 4096, PROT_READ|PROT_WRITE,
                    MAP_ANON|MAP_PRIVATE, -1, 0);
head->size = 4096 - sizeof(node_t);
head->next = NULL;
```

- Heap layout:

head →

size:	4088
next:	0
...	

[virtual address: 16KB]
header: size field

header: next field(NULL is 0)

the rest of the 4KB chunk

November 19, 2018TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - TacomaL14.53

FREE LIST: MALLOC() CALL

- Consider a request for a 100 bytes: `malloc(100)`
- Header block requires 8 bytes
 - 4 bytes for size, 4 bytes for magic number
- Split the heap – header goes with each block

A 4KB Heap With One Free Chunk

head →

size:	4088
next:	0
...	

the rest of the 4KB chunk

A Heap : After One Allocation

ptr →

size:	100
magic:	1234567

First block is used

head →

size:	3980
next:	0
...	

the 100 bytes now allocated

the free 3980 byte chunk

November 19, 2018TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - TacomaL14.54

FREE LIST: FREE() CALL

- Addresses of chunks
- Start=16384
 - + 108 (end of 1st chunk)
 - + 108 (end of 2nd chunk)
 - + 108 (end of 3rd chunk)
 - = 16708

Free Space With Three Chunks Allocated

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.55

FREE LIST: FREE() CHUNK #2

- Free(sptr)
- Our 3 chunks start at 16 KB (@ 16,384 bytes)
- Free chunk #2 - sptr
- Sptr = 16500
 - addr - sizeof(node_t)
- Actual start of chunk #2
 - 16492

The free 3764-byte chunk

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.56

FREE LIST- FREE ALL CHUNKS

- Now free remaining chunks:
- Free(16392)
- Free(16608)
- Walk back 8 bytes for actual start of chunk
- External fragmentation
- Free chunk pointers out of order
- Coalescing of next pointers is needed

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.57

GROWING THE HEAP

- Start with small sized heap
- Request more memory when full
- sbrk(), brk()

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.58

MEMORY ALLOCATION STRATEGIES

- **Best fit**
 - Traverse free list
 - Identify all candidate free chunks
 - Note which is smallest (has best fit)
 - When splitting, “leftover” pieces are small (and potentially less useful -- fragmented)
- **Worst fit**
 - Traverse free list
 - Identify largest free chunk
 - Split largest free chunk, leaving a still large free chunk

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.59

EXAMPLES

- Allocation request for 15 bytes



- Result of Best Fit



- Result of Worst Fit



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.60

MEMORY ALLOCATION STRATEGIES - 2

- **First fit**
 - Start search at beginning of free list
 - Find first chunk large enough for request
 - Split chunk, returning a “fit” chunk, saving the remainder
 - Avoids full free list traversal of best and worst fit

- **Next fit**
 - Similar to first fit, but start search at last search location
 - Maintain a pointer that “cycles” through the list
 - Helps balance chunk distribution vs. first fit
 - Find first chunk, that is large enough for the request, and split
 - Avoids full free list traversal

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.61

SEGREGATED LISTS

- For popular sized requests
e.g. for kernel objects such as locks, inodes, etc.
- Manage as segregated free lists
- Provide object caches: stores pre-initialized objects

- How much memory should be dedicated for specialized requests (object caches)?

- If a given cache is low in memory, can request “slabs” of memory from the general allocator for caches.
- General allocator will reclaim slabs when not used

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.62

BUDDY ALLOCATION

- Binary buddy allocation
 - Divides free space by two to find a block that is big enough to accommodate the request; the next split is too small...
- Consider a 7KB request

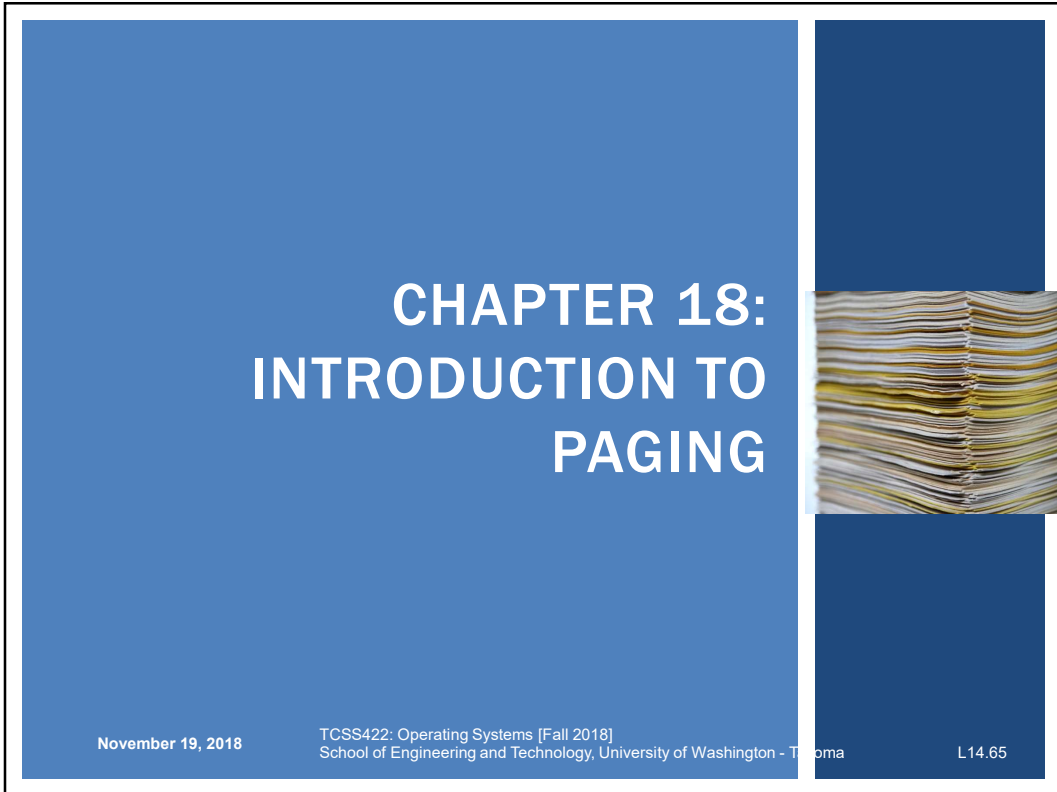
64KB free space for 7KB request

November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.63
-------------------	---	--------

BUDDY ALLOCATION - 2

- Buddy allocation: suffers from internal fragmentation
- Allocated fragments, typically too large
- Coalescing is simple
 - Two adjacent blocks are promoted up

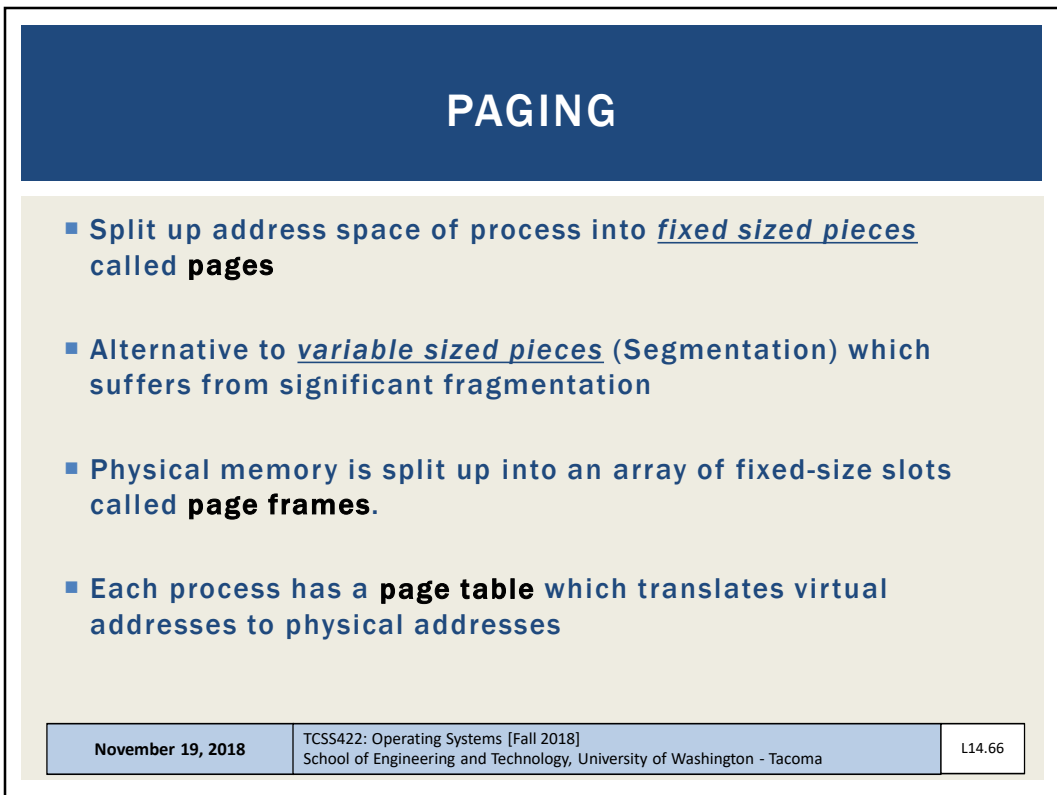
November 19, 2018	TCSS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.64
-------------------	---	--------



CHAPTER 18:
INTRODUCTION TO
PAGING

November 19, 2018 TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma L14.65

This slide features a blue background with a stack of papers on the right side. The title 'CHAPTER 18: INTRODUCTION TO PAGING' is centered in white. At the bottom, there is a footer with the date 'November 19, 2018', the course name 'TCSS422: Operating Systems [Fall 2018]', the school name 'School of Engineering and Technology, University of Washington - Tacoma', and the slide number 'L14.65'.



PAGING

- Split up address space of process into *fixed sized pieces* called **pages**
- Alternative to *variable sized pieces* (Segmentation) which suffers from significant fragmentation
- Physical memory is split up into an array of fixed-size slots called **page frames**.
- Each process has a **page table** which translates virtual addresses to physical addresses

November 19, 2018 TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma L14.66

This slide has a dark blue header with the word 'PAGING' in white. The main content is on a light beige background and consists of four bullet points. The footer is a light blue bar with the date 'November 19, 2018', the course name 'TCSS422: Operating Systems [Fall 2018]', the school name 'School of Engineering and Technology, University of Washington - Tacoma', and the slide number 'L14.66'.

ADVANTAGES OF PAGING

- Flexibility
 - Abstracts the process address space into pages
 - No need to track direction of HEAP / STACK growth
 - *Just add more pages...*
 - No need to store unused space
 - *As with segments...*

- Simplicity
 - Pages and page frames are the same size
 - Easy to allocate and keep a free list of pages

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.67
-------------------	---	--------

PAGING: EXAMPLE

Page Table:
 VP0 → PF3
 VP1 → PF7
 VP2 → PF5
 VP3 → PF2

- Consider a 128 byte address space with 16-byte pages
- Consider a 64-byte program address space

0	(page 0 of the address space) (page 1) (page 2) (page 3)
16	
32	
48	
64	

A Simple 64-byte Address Space

0	reserved for OS	page frame 0 of physical memory
16	(unused)	page frame 1
32	page 3 of AS	page frame 2
48	page 0 of AS	page frame 3
64	(unused)	page frame 4
80	page 2 of AS	page frame 5
96	(unused)	page frame 6
112	page 1 of AS	page frame 7
128		

64-Byte Address Space Placed In Physical Memory

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.68
-------------------	---	--------

PAGING: ADDRESS TRANSLATION

- **PAGE:** Has two address components
 - **VPN:** Virtual Page Number
 - **Offset:** Offset within a Page

VPN		offset			
Va5	Va4	Va3	Va2	Va1	Va0

- **Example:**
 Page Size: 16-bytes, Address Space: 64-bytes

VPN		offset			
0	1	0	1	0	1

Here there are just four pages...

November 19, 2018
TCSS422: Operating Systems [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma
L14.69

EXAMPLE: PAGING ADDRESS TRANSLATION

- Consider a 64-byte program address space (4 pages)
- Stored in 128-byte physical memory (8 frames)
- Offset is preserved
- VPN is looked up

Page Table:

VP0 → PF3

VP1 → PF7

VP2 → PF5

VP3 → PF2

Virtual Address

VPN		offset			
0	1	0	1	0	1

Address Translation					
----------------------------	--	--	--	--	--

Physical Address

1	1	1	0	1	0	1
PFN			offset			

November 19, 2018
TCSS422: Operating Systems [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma
L14.70

PAGING DESIGN QUESTIONS

- (1) Where are page tables stored?
- (2) What are the typical contents of the page table?
- (3) How big are page tables?
- (4) Does paging make the system too slow?

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.71

(1) WHERE ARE PAGE TABLES STORED?

- Example:
 - Consider a 32-bit process address space (up to 4GB)
 - With 4 KB pages
 - 20 bits for VPN (2^{20} pages)
 - 12 bits for the page offset (2^{12} unique bytes in a page)
- Page tables for each process are stored in RAM
 - Support potential storage of 2^{20} translations
= 1,048,576 pages per process
 - Each page has a page table entry size of 4 bytes

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.72

PAGE TABLE EXAMPLE

- With 2^{20} slots in our page table for a single process
- Each slot dereferences a VPN
- Provides physical frame number
- Each slot requires 4 bytes (32 bits)
 - 20 for the PFN on a 4GB system with 4KB pages
 - 12 for the offset which is preserved
 - (note we have no status bits, so this is unrealistically small)
- How much memory to store page table for 1 process?
 - 4,194,304 bytes (or 4MB) to index one process

VPN ₀
VPN ₁
VPN ₂
...
...
VPN ₁₀₄₈₅₇₆

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.73

NOW FOR AN ENTIRE OS

- If 4 MB is required to store one process
 - Consider how much memory is required for an entire OS?
 - With for example 100 processes...
 - Page table memory requirement is now 4MB x 100 = 400MB
 - If computer has 4GB memory (maximum for 32-bits), the page table consumes 10% of memory
- 400 MB / 4000 GB
- Is this efficient?

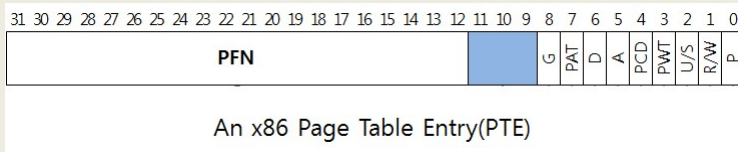
November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.74

(2) WHAT'S ACTUALLY IN THE PAGE TABLE

- Page table is data structure used to map virtual page numbers (VPN) to the physical address (Physical Frame Number PFN)
 - Linear page table → simple array
- Page-table entry
 - 32 bits for capturing state



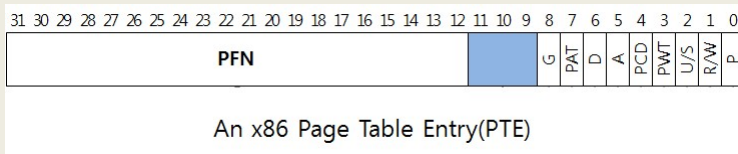
November 19, 2018

TCSS422: Operating Systems [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma

L14.75

PAGE TABLE ENTRY

- P: present
- R/W: read/write bit
- U/S: supervisor
- A: accessed bit
- D: dirty bit
- PFN: the page frame number



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma

L14.76

PAGE TABLE ENTRY - 2

- **Common flags:**
- **Valid Bit:** Indicating whether the particular translation is valid.
- **Protection Bit:** Indicating whether the page could be read from, written to, or executed from
- **Present Bit:** Indicating whether this page is in physical memory or on disk(swapped out)
- **Dirty Bit:** Indicating whether the page has been modified since it was brought into memory
- **Reference Bit(Accessed Bit):** Indicating that a page has been accessed

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.77

(3) HOW BIG ARE PAGE TABLES?

- Page tables are too big to store on the CPU
- Page tables are stored using physical memory
- Paging supports efficiently storing a sparsely populated address space
 - Reduced memory requirement
Compared to base and bounds, and segments

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.78

(4) DOES PAGING MAKE THE SYSTEM TOO SLOW?

- Translation

- **Issue #1:** Starting location of the page table is needed

- HW Support: Page-table base register
 - stores active process
 - Facilitates translation

Page Table:

VP0 → PF3
VP1 → PF7
VP2 → PF5
VP3 → PF2

Stored in RAM →

- **Issue #2:** Each memory address translation for paging requires an extra memory reference

- HW Support: TLBs (Chapter 19)

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.79

PAGING MEMORY ACCESS

```
1. // Extract the VPN from the virtual address
2. VPN = (VirtualAddress & VPN_MASK) >> SHIFT
3.
4. // Form the address of the page-table entry (PTE)
5. PTEAddr = PTBR + (VPN * sizeof(PTE))
6.
7. // Fetch the PTE
8. PTE = AccessMemory(PTEAddr)
9.
10. // Check if process can access the page
11. if (PTE.Valid == False)
12.     RaiseException(SEGMENTATION_FAULT)
13. else if (CanAccess(PTE.ProtectBits) == False)
14.     RaiseException(PROTECTION_FAULT)
15. else
16.     // Access is OK: form physical address and fetch it
17.     offset = VirtualAddress & OFFSET_MASK
18.     PhysAddr = (PTE.PFN << PFN_SHIFT) | offset
19.     Register = AccessMemory(PhysAddr)
```

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.80

COUNTING MEMORY ACCESSES

- **Example: Use this Array initialization Code**

```

int array[1000];
...
for (i = 0; i < 1000; i++)
    array[i] = 0;
    
```

- **Assembly equivalent:**

```

0x1024 movl $0x0, (%edi,%eax,4)
0x1028 incl %eax
0x102c cmpl $0x03e8,%eax
0x1030 jne 0x1024
    
```

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.81
-------------------	---	--------

VISUALIZING MEMORY ACCESSES: FOR THE FIRST 5 LOOP ITERATIONS

- **Locations:**
 - Page table
 - Array
 - Code
- **50 accesses for 5 loop iterations**

November 19, 2018	TCCS422: Operating Systems [Fall 2018] School of Engineering and Technology, University of Washington - Tacoma	L14.82
-------------------	---	--------


PAGING SYSTEM EXAMPLE

- Consider a 4GB Computer:
 - With a 4096-byte page size (4KB)
 - How many pages would fit in physical memory?
- Now consider a page table:
 - For the page table entry, how many bits are required for the VPN?
 - If we assume the use of 4-byte (32 bit) page table entries, how many bits are available for status bits?
 - How much space does this page table require?
Page Table Entries x Number of pages
 - How many page tables (for user processes) would fill the entire 4GB of memory?

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.83



CHAPTER 19: TRANSLATION LOOKASIDE BUFFER (TLB)

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.84

OBJECTIVES

- Chapter 19
 - TLB Algorithm
 - TLB Tradeoffs
 - TLB Context Switch

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.85

TRANSLATION LOOKASIDE BUFFER

- Legacy name...
- Better name, “Address Translation Cache”
- TLB is an on CPU cache of address translations
 - virtual → physical memory

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.86

TRANSLATION LOOKASIDE BUFFER - 2

- **Goal:**
Reduce access to the page tables
- **Example:**
50 RAM accesses for first 5 for-loop iterations
- **Move lookups** from RAM to TLB by caching page table entries

The figure consists of three vertically stacked graphs sharing a common x-axis labeled 'Memory Access' from 0 to 50.

- Top Graph (Page Table):** Shows 'Page Table[39]' and 'Page Table[1]' as single points at the top of the y-axis (1174-1224 PA). The rest of the page table entries are clustered at the bottom (1024-1074 PA).
- Middle Graph (Array):** Shows 'Array(VA)' with values 40000, 40050, 40100. The y-axis ranges from 7232 to 7132 PA.
- Bottom Graph (Code):** Shows 'Code(VA)' with values 1024, 1074, 1124. The y-axis ranges from 4096 to 4196 PA. Operations 'mov', 'inc', 'cmp', and 'jne' are labeled.

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.87

TRANSLATION LOOKASIDE BUFFER (TLB)

- Part of the CPU's Memory Management Unit (MMU)
- Address translation cache

The diagram illustrates the flow of address translation:

- CPU** provides a **Logical Address**.
- The **MMU** performs a **TLB Lookup**.
- If **TLB Hit**, the **MMU** outputs a **Physical Address**.
- If **TLB Miss**, the **MMU** consults the **Page Table** (all v to p entries) to find the **Physical Address**.
- The **Physical Address** is used to access **Physical Memory**, which contains **Page 0**, **Page 1**, **Page 2**, ..., **Page n**.

Address Translation with MMU

November 19, 2018

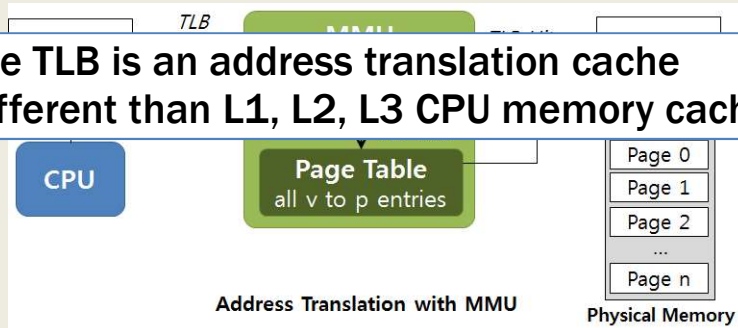
TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.88

TRANSLATION LOOKASIDE BUFFER (TLB)

- Part of the CPU's Memory Management Unit (MMU)
- Address translation cache

The TLB is an address translation cache
Different than L1, L2, L3 CPU memory caches



November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.89

TLB BASIC ALGORITHM

- For: array based page table
- Hardware managed TLB

```
1: VPN = (VirtualAddress & VPN_MASK ) >> SHIFT
2: (Success , TlbEntry) = TLB_Lookup(VPN)
3:   if(Success == True){ // TLB Hit
4:     if(CanAccess(TlbEntry.ProtectBits) == True ){
5:       Offset = VirtualAddress & OFFSET_MASK
6:       PhysAddr = (TlbEntry.PFN << SHIFT) | Offset
7:       AccessMemory( PhysAddr )
8:     }else RaiseException( PROTECTION_ERROR)
```

Generate the physical address to access memory

November 19, 2018

TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.90

TLB BASIC ALGORITHM - 2

```
11:     else{ //TLB Miss
12:         PTEAddr = PTBR + (VPN * sizeof(PTE))
13:         ➡ PTE = AccessMemory(PTEAddr)
14:         (...) // Check for, and raise exceptions...
15:
16:         ➡ TLB_Insert( VPN , PTE.PFN , PTE.ProtectBits)
17:         ➡ RetryInstruction()
18:     }
19: }
```

Retry the instruction... (requery the TLB)

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.91

TLB – ADDRESS TRANSLATION CACHE

- Key detail:
- For a TLB miss, we first access the page table in RAM to populate the TLB... we then requery the TLB
- All address translations go through the TLB

November 19, 2018

TCCS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma

L14.92

TLB EXAMPLE

```

0:   int sum = 0 ;
1:   for( i=0; i<10; i++){
2:       sum+=a[i];
3:   }
    
```

- **Example:**
- **Program address space: 256-byte**
 - Addressable using 8 total bits (2^8)
 - 4 bits for the VPN (16 total pages)
- **Page size: 16 bytes**
 - Offset is addressable using 4-bits
- **Store an array: of (10) 4-byte integers**

	OFFSET				
	00	04	08	12	16
VPN = 00					
VPN = 01					
VPN = 03					
VPN = 04					
VPN = 05					
VPN = 06		a[0]	a[1]	a[2]	
VPN = 07	a[3]	a[4]	a[5]	a[6]	
VPN = 08	a[7]	a[8]	a[9]		
VPN = 09					
VPN = 10					
VPN = 11					
VPN = 12					
VPN = 13					
VPN = 14					
VPN = 15					

November 19, 2018
TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma
L14.93

TLB EXAMPLE - 2

```

0:   int sum = 0 ;
1:   for( i=0; i<10; i++){
2:       sum+=a[i];
3:   }
    
```

- **Consider the code above:**
- **Initially the TLB does not know where a[] is**
- **Consider the accesses:**
 - a[0], a[1], a[2], a[3], a[4], a[5], a[6], a[7], a[8], a[9]
- **How many pages are accessed?**
- **What happens when accessing a page not in the TLB?**

	OFFSET				
	00	04	08	12	16
VPN = 00					
VPN = 01					
VPN = 03					
VPN = 04					
VPN = 05					
VPN = 06		a[0]	a[1]	a[2]	
VPN = 07	a[3]	a[4]	a[5]	a[6]	
VPN = 08	a[7]	a[8]	a[9]		
VPN = 09					
VPN = 10					
VPN = 11					
VPN = 12					
VPN = 13					
VPN = 14					
VPN = 15					

November 19, 2018
TCSS422: Operating Systems [Fall 2018]
School of Engineering and Technology, University of Washington - Tacoma
L14.94

TLB EXAMPLE - 3

```

0:   int sum = 0 ;
1:   for( i=0; i<10; i++){
2:       sum+=a[i];
3:   }
    
```

- For the accesses: a[0], a[1], a[2], a[3], a[4],
- a[5], a[6], a[7], a[8], a[9]

- How many are hits?
- How many are misses?
- What is the hit rate? (%)
 - 70% (3 misses one for each VP, 7 hits)

	OFFSET				
	00	04	08	12	16
VPN = 00					
VPN = 01					
VPN = 03					
VPN = 04					
VPN = 05					
VPN = 06		a[0]	a[1]	a[2]	
VPN = 07	a[3]	a[4]	a[5]	a[6]	
VPN = 08	a[7]	a[8]	a[9]		
VPN = 09					
VPN = 10					
VPN = 11					
VPN = 12					
VPN = 13					
VPN = 14					
VPN = 15					

November 19, 2018
TCSS422: Operating Systems [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma
L14.95

TLB EXAMPLE - 4

```

0:   int sum = 0 ;
1:   for( i=0; i<10; i++){
2:       sum+=a[i];
3:   }
    
```

- What factors affect the hit/miss rate?
 - Page size
 - Data locality
 - Temporal locality

	OFFSET				
	00	04	08	12	16
VPN = 00					
VPN = 01					
VPN = 03					
VPN = 04					
VPN = 05					
VPN = 06		a[0]	a[1]	a[2]	
VPN = 07	a[3]	a[4]	a[5]	a[6]	
VPN = 08	a[7]	a[8]	a[9]		
VPN = 09					
VPN = 10					
VPN = 11					
VPN = 12					
VPN = 13					
VPN = 14					
VPN = 15					

November 19, 2018
TCSS422: Operating Systems [Fall 2018]
 School of Engineering and Technology, University of Washington - Tacoma
L14.96

