

Lens: Experiencing Multi-level Page Tables at Close Quarters



Arun KP Rohit Singh Debadatta Mishra
Indian Institute of Technology, Kanpur



Importance of Operating Systems



Importance of Operating Systems

- OS is a building block in today's digital world.



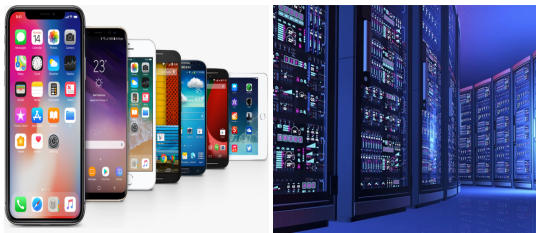
Importance of Operating Systems



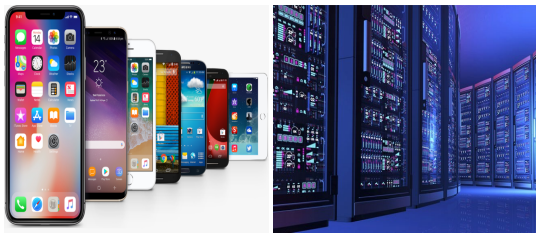
Importance of Operating Systems



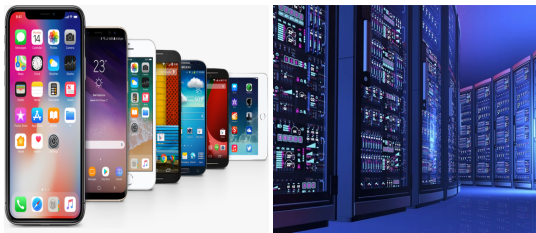
Importance of Operating Systems



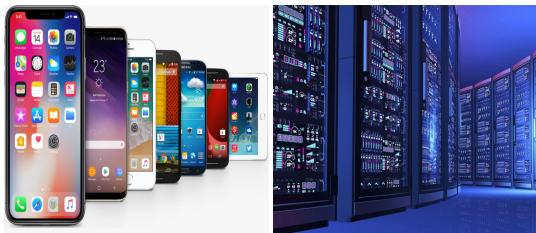
Importance of Operating Systems



Importance of Operating Systems

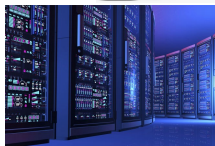


Importance of Operating Systems



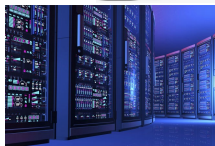
Importance of Operating Systems

- OS is a building block in today's digital world.
- OS helps us to perform our day-to-day digital activities better.



Importance of Operating Systems

- OS is a building block in today's digital world.
- OS helps us to perform our day-to-day digital activities better.
- OS encompasses multiple subsystems such as File system, Virtual Memory, and the interplay between them.



The Virtual Memory Sub-system



Understanding the VM Sub-system

- Virtual memory allows inter-process isolation.



Understanding the VM Sub-system

- Virtual memory allows inter-process isolation.
- Each process has an independent view of memory.



Understanding the VM Sub-system

- Virtual memory allows inter-process isolation.
- Each process has an independent view of memory.
- Virtual memory makes a programmer's life easy.



Understanding the VM Sub-system

- Virtual memory allows inter-process isolation.
- Each process has an independent view of memory.
- Virtual memory makes a programmer's life easy.
- Virtual memory facilitates memory access permission.



Understanding the VM Sub-system

Text Books

- 3 Use simple examples to explain concepts

Educational OS

Debugging Interfaces



Understanding the VM Sub-system

Text Books

- 3 Use simple examples to explain concepts
- 7 Do not capture real-world complexities

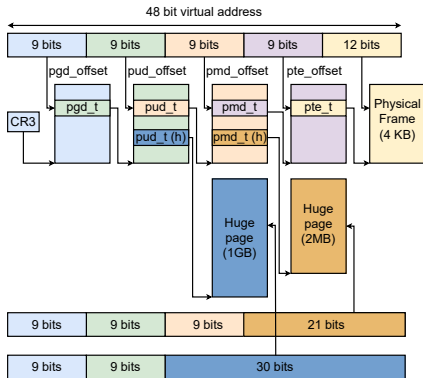
Educational OS

Debugging Interfaces



Page table in Intel x86-64 processor

- 4 level paging translates 48 bit virtual address to physical address.
- Allows translation into 4 KB, 2MB, and 1 GB pages.
- Each entry provides details about the next level page.



Ref: [Guide(2011)]

Page table entry in Intel x86-64 system

- PTE captures access permissions, access mode and other details.
- Intel software manual details page table structures and different configurations.
- Challenging for a beginner to correlate page table details with a program's virtual memory state.

63	62..59	58..52	51... M	M-1	12	11..9	8	7	6	5	4	3	2	1	0
X D	PK	AVL	Reserved (0)	Bits 12 - (M-1) of address	AVL	G	P A T	D	A	P C D	U / S	U / S	R / W	R / W	P

P: Present	G: Global
R/W: Read/Write	AVL: Available
U/S: User/Supervisor	PAT: Page Attribute Table
PWT: Write-Through	M: Maximum Physical Address Bit
PCD: Cache Disable	PK: Protection Key
A: Accessed	XD: Execute Disable
D: Dirty	

Ref: [Guide(2011)]



Understanding the VM Sub-system

Text Books

Educational OS

3 Profiling the code with debug statements

Debugging Interfaces



Understanding the VM Sub-system

Text Books

Educational OS

3 Profiling the code with debug statements

7 Requires understanding the code base

Debugging Interfaces



Understanding the VM Sub-system

Text Books

Educational OS

Debugging Interfaces

3 Reading the interface for fetching predefined virtual memory details



Understanding the VM Sub-system

Text Books

Educational OS

Debugging Interfaces

- 3 Reading the interface for fetching predefined virtual memory details
- 7 Gives limited information for education



VM interfaces in Linux

- `#lfc WdIX# Udg` shows currently mapped memory regions.
- `#lfc WdIX# Udg` is a pseudo file maintained under proc file system for each process.
- It is not straightforward to correlate a program variable address with information from `#lfc WdIX# Udg`

```
55a4c3600000-55a4c3601000 r-xp 00000000 08:02 25958718 program_code
55a4c3800000-55a4c3801000 r--p 00000000 08:02 25958718 program_code
55a4c3801000-55a4c3802000 rw-p 00001000 08:02 25958718 program_code
55a4c55fa000-55a4c561b000 rw-p 00000000 00:00 0 [heap]
7f4407200000-7f44072e7000 r-xp 00000000 08:02 36177285 libc-2.27.so
7f44072e7000-7f44072e7000 ---p 001c7000 08:02 36177285 libc-2.27.so
7f44072e7000-7f44072e7000 r--p 001e7000 08:02 36177285 libc-2.27.so
7f44072e7000-7f44072e7000 r-wp 001eb000 08:02 36177285 libc-2.27.so
7f44072e7000-7f44072e7000 r-wp 00000000 00:00 0
7f4407600000-7f4407629000 r-xp 00000000 08:02 36177281 ld-2.27.so
7f4407629000-7f440762a000 r--p 00029000 08:02 36177281 ld-2.27.so
7f440762a000-7f440762b000 rw-p 0002a000 08:02 36177281 ld-2.27.so
7f440762b000-7f440762c000 rw-p 00000000 00:00 0
7f4407929000-7f440792b000 rw-p 00000000 00:00 0
7ff9934b0000-7ff9934b0000 r-wp 00000000 00:00 0 [stack]
7ff9934b0000-7ff9934f0000 r--p 00000000 00:00 0 [vvar]
7ff9934f0000-7ff9934f0000 r--p 00000000 00:00 0 [vdso]
ffffffff600000-ffffffff601000 --xp 00000000 00:00 0 [vsyscall]
```



VM interfaces in Linux

- `#cat /proc/$pid/pagemap` shows virtual to physical page mapping.
- It requires calculating the file offset corresponding to a program variable in the pagemap file to get its physical page.
- Only gives information in the last level of a page table (PTE).

63	62	61	60..57	56	55	54	53	4	3	2	1	0
R	S	MT	0	E	SD	Page Frame Number							

SD: Soft Dirty
E: Exclusively Mapped Page
MT: Mapping Type
 (file-mapped page or a shared anonymous page)
S: Page is in swap space
R: Page is in RAM



Tool for Learning Virtual Memory



Features of an educational tool for VM

- Augments concept learning with practical exposure.



Features of an educational tool for VM

- Augments concept learning with practical exposure.
- Provides details from a real system without missing actual complexities.



Features of an educational tool for VM

- Augments concept learning with practical exposure.
- Provides details from a real system without missing actual complexities.
- Presents an interface that is easy-to-use.



Features of an educational tool for VM

- Augments concept learning with practical exposure.
- Provides details from a real system without missing actual complexities.
- Presents an interface that is easy-to-use.
- Expects only basic C programming knowledge from learners.



Features of an educational tool for VM

- Augments concept learning with practical exposure.
- Provides details from a real system without missing actual complexities.
- Presents an interface that is easy-to-use.
- Expects only basic C programming knowledge from learners.
- Allows learners to correlate program execution with memory state.



Features of an educational tool for VM

- Augments concept learning with practical exposure.
- Provides details from a real system without missing actual complexities.
- Presents an interface that is easy-to-use.
- Expects only basic C programming knowledge from learners.
- Allows learners to correlate program execution with memory state.
- Visualizes details for easy understanding by learners.



Features of an educational tool for VM

- Augments concept learning with practical exposure.
- Provides details from a real system without missing actual complexities.
- Presents an interface that is easy-to-use.
- Expects only basic C programming knowledge from learners.
- Allows learners to correlate program execution with memory state.
- Visualizes details for easy understanding by learners.
- Shows changes in real-time by updating visualization in sync with program changes.



Lens: An Education Tool for VM



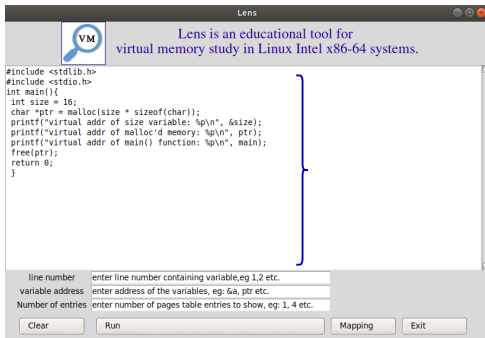
What does Lens offer?

- Provides a practical exposure to virtual memory in Linux.



What does Lens offer?

- Provides a practical exposure to virtual memory in Linux.
- Provides an interface to
 - Write C programs.



What does Lens offer?

- Provides a practical exposure to virtual memory in Linux.
- Provides an interface to
 - Write C programs.
 - Correlate program execution with memory state.

The screenshot shows the Lens application window. At the top, there is a title bar with the text "Lens" and a logo with "VM" inside a magnifying glass. Below the title bar, a banner reads "Lens is an educational tool for virtual memory study in Linux Intel x86-64 systems." The main area contains a C program:

```
#include <stdlib.h>
#include <stdio.h>
int main(){
    int size = 16;
    char *ptr = malloc(size * sizeof(char));
    printf("virtual addr of size variable: %p\n", &size);
    printf("virtual addr of malloc'd memory: %p\n", ptr);
    printf("virtual addr of main() function: %p\n", main);
    free(ptr);
    return 0;
}
```

Below the code, there is a table with the following fields:

line number	4
variable address	&size
Number of entries	1

Below the table are buttons for "Clear", "Run", "Mapping", and "Exit".

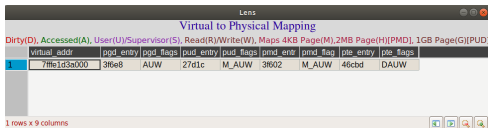
Below the buttons, the "output:" section shows the following text:

```
virtual addr of size variable: 0x7ffe1d3ae1c
virtual addr of malloc'd memory: 0x555ce50ee6d0
virtual addr of main() function: 0x555ce4ed9349
```



What does Lens offer?

- Provides a practical exposure to virtual memory in Linux.
- Provides an interface to
 - Write C programs.
 - Correlate program execution with memory state.
 - Visualize virtual to physical address translation.



The screenshot shows a window titled 'Lens' with the subtitle 'Virtual to Physical Mapping'. Below the subtitle, there is a header line with various memory-related attributes: 'Dirty(D), Accessed(A), User(U)/Supervisor(S), Read(R)/Write(W), Maps 4KB Page(M), 2MB Page(H)[PMD], 1GB Page(G)[PUD]'. Below this is a table with 9 columns: 'virtual_addr', 'pgd_entry', 'pgd_flags', 'pud_entry', 'pud_flags', 'pmd_entry', 'pmd_flag', 'pte_entry', and 'pte_flags'. The first row of data shows: '7ffe1d3a000', '3f6e8', 'AUW', '27d1c', 'M_AUW', '3f602', 'M_AUW', '46cbd', and 'DAUW'. At the bottom left of the window, it says '1 rows x 9 columns'. At the bottom right, there are three small icons: a magnifying glass, a refresh icon, and a close icon.

virtual_addr	pgd_entry	pgd_flags	pud_entry	pud_flags	pmd_entry	pmd_flag	pte_entry	pte_flags
7ffe1d3a000	3f6e8	AUW	27d1c	M_AUW	3f602	M_AUW	46cbd	DAUW

Start Video

COMPUTER SCIENCE & ENGINEERING

IIT KANPUR



What does Lens offer?

- Provides a practical exposure to virtual memory in Linux.
- Provides an interface to
 - Write C programs.
 - Correlate program execution with memory state.
 - Visualize virtual to physical address translation.
- Covers underlying hardware-OS complexities.



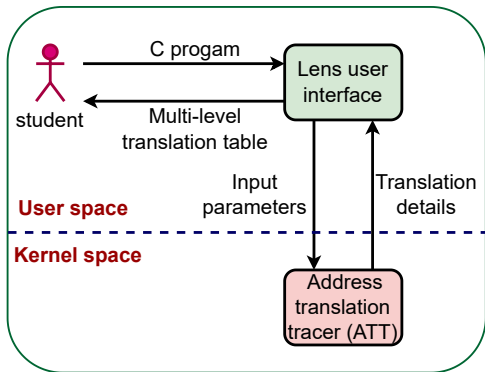
What does Lens offer?

- Provides a practical exposure to virtual memory in Linux.
- Provides an interface to
 - Write C programs.
 - Correlate program execution with memory state.
 - Visualize virtual to physical address translation.
- Covers underlying hardware-OS complexities.
- Shows address translation table with latest updates in program.



High level view of Lens

- Student interacts with graphical interface.
- ATT handles underlying hardware-OS complexities.



Example use cases with Lens

- Lazy page allocation in Linux.
- Linux allocates physical pages on first access.



Inspecting lazy page allocation

```
1 #define TWOMB 2097152
2 int main(int argc, char* argv[]){
3     char * ptr = (char*)mmap(NULL, TWOMB, PROT_READ|PROT_WRITE, MAP_PRIVATE|
4     MAP_ANONYMOUS,0,0);
5     ptr[0] = 'A';
6     printf("ptr:%p\n", ptr);
7     ptr[4096] = 'B';
8     mprotect(&ptr[4096],4096,PROT_READ);
9     munmap(ptr, TWOMB);
10    return 0;}
```

Inspecting at line number 4



Inspecting lazy page allocation

```
1 #define TWOMB 2097152
2 int main(int argc, char* argv[]){
3     char * ptr = (char*)mmap(NULL, TWOMB, PROT_READ|PROT_WRITE, MAP_PRIVATE|
4     MAP_ANONYMOUS,0,0);
5     ptr[0] = 'A';
6     printf("ptr:%p\n", ptr);
7     ptr[4096] = 'B';
8     mprotect(&ptr[4096],4096,PROT_READ);
9     munmap(ptr, TWOMB);
10    return 0;}
```

Inspecting at line number 4

pud_entry	pud_flags	pmd_entry	pmd_flags	pte_entry	pte_flags
36d4b	M_AUW	ffa2	M_AUW	1c8e4	DAUW
36d4b	M_AUW	ffa2	M_AUW	0	---



Example use cases with Lens

- Lazy page allocation in Linux.
- Linux allocates physical pages on first access.
- Changing memory access permission.
- Linux provides a *dfchmm* sys call to change access protection of pages.



Changing memory access permission

```
1 #define TWOMB 2097152
2 int main(int argc, char* argv){
3     char * ptr = (char*)mmap(NULL, TWOMB, PROT_READ|PROT_WRITE, MAP_PRIVATE|
4     MAP_ANONYMOUS,0,0);
5     ptr[0] = 'A';
6     printf("ptr:%p\n", ptr);
7     ptr[4096] = 'B';
8     mprotect(&ptr[4096],4096,PROT_READ);
9     munmap(ptr, TWOMB);
10    return 0;}
```

Inspecting at line number 7



Changing memory access permission

```
1 #define TWOMB 2097152
2 int main(int argc, char* argv[]){
3     char * ptr = (char*)mmap(NULL, TWOMB, PROT_READ|PROT_WRITE, MAP_PRIVATE|
4     MAP_ANONYMOUS,0,0);
5     ptr[0] = 'A';
6     printf("ptr:%p\n", ptr);
7     ptr[4096] = 'B';
8     mprotect(&ptr[4096],4096,PROT_READ);
9     munmap(ptr, TWOMB);
10    return 0;}
```

Inspecting at line number 7

pud_entry	pud_flags	pmd_entry	pmd_flags	pte_entry	pte_flags
36d4b	M_AUW	ffa2	M_AUW	1c8e5	DAUR



Example use cases with Lens

- Lazy page allocation in Linux.
- Linux allocates physical pages on first access.
- Changing memory access permission.
- Linux provides a `mprotect` sys call to change access protection of pages.
- Lens can also provide insights to other virtual memory concepts like hugepages.



Future Directions

- Use Lens in an OS course for collecting feedback.
- Show address translation table as a radix tree and make the interface interactive for students to zoom-in.
- Visualize more virtual memory changes such addition, deletion and merging of VM areas, stack/heap growth and shrink with changes in program.
- Visualize more OS activities related to virtual memory, like page faults.



Conclusion

- Practical understanding of virtual to physical address translation is essential in learning virtual memory.
- Students require hands-on experience with virtual memory concepts to enhance their learning.
- Lens provides an easy interface for students to write C code and correlate program state with changes in virtual memory.
- Lens helps students in gaining practical exposure to virtual memory concepts.
- Lens is available at <https://github.com/arunkp1986/Lens.git>



Thank You!


Questions?

Arun KP

kparun@cse.iitk.ac.in



References

-  Part Guide.
Intel® 64 and ia-32 architectures software developer's manual.
Jc ĩ a Y' 6. GghYá dfc[fUá a]b[; i]XZDUh, 2:5, 2011.

