NEU CS 3650 Computer Systems

Instructor: Dr. Ziming Zhao

Process

- The concept of process
- ELF and other executable file format
- How is an ELF loaded to memory?
- Virtualizing CPU with time-sharing
- The mechanism to switch processes

Diving into the Operating Systems

- So far, we have been preparing for our further exploration:
 - CPU, Memory, Assembly, C

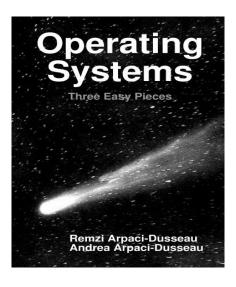
- Today we will dive into the OS itself. What we learned will be useful
 - Registers and instruction concepts
 - Memory as a linear array and ways to work with memory addresses
 - C is at the core of many common OSes

OS: Virtualization + Abstraction

- The OS is a (software) land of magic and illusions
- OS makes a computer "easy" to use
- OS hides overwhelming complexities of hardware behind an API
 - This is abstraction
- OS creates the illusion of an ideal, general, and powerful machine
 - This is virtualization
- We will start by looking at how OSes provides CPU virtualization at the process level.
- Other abstractions the OS uses

Required Reading

- The OSTEP book: up to Ch. 3-6
- Online: https://pages.cs.wisc.edu/~remzi/OSTEP/
- Hard copy: Lulu or Amazon



Virtualization

- 3 <u>Dialogue</u>
- 4 Processes
- 5 Process API code
- 6 Direct Execution
- 7 CPU Scheduling
- 8 Multi-level Feedback
- 9 Lottery Scheduling code
- 10 Multi-CPU Scheduling
- 11 <u>Summary</u>

Process

• Basic function of an OS is to execute and manage code dynamically:

for example,

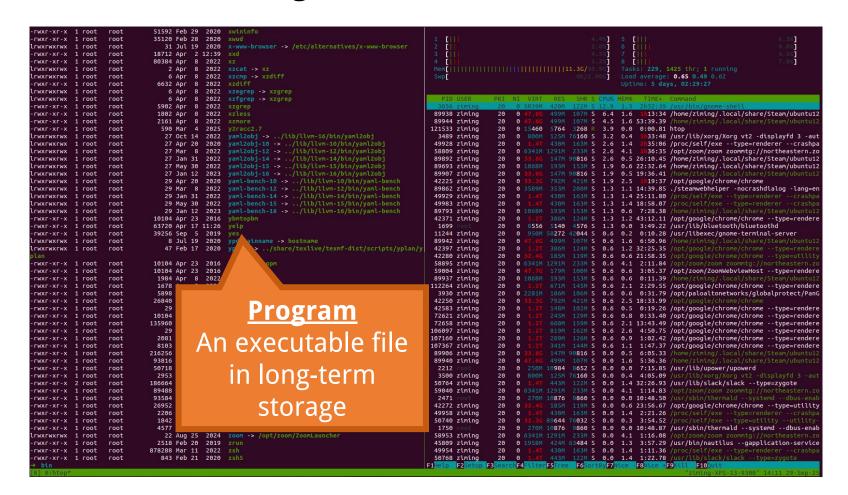
- A command issued at a command line terminal
- An icon double clicked from the desktop
- Jobs/tasks run as part of a batch system

• A process is the basic unit of a program in execution. On other words, process is a running program.

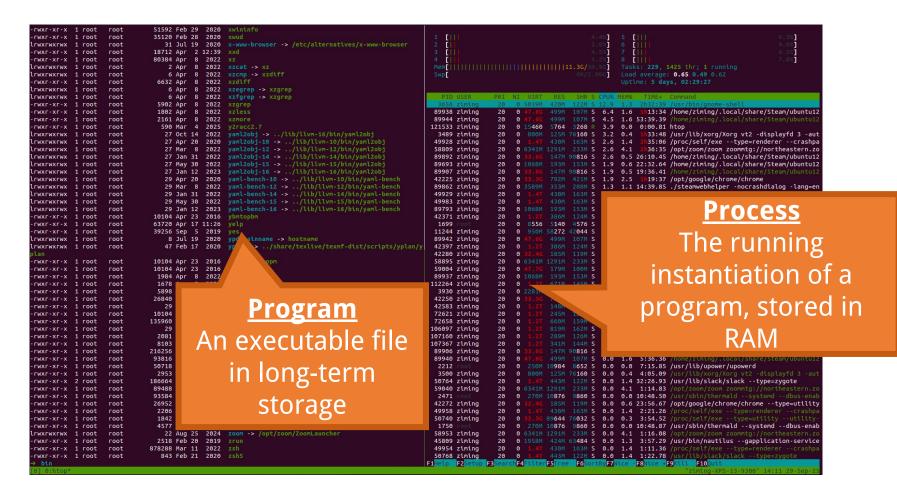
Programs and Processes (Commands: Is ps htop)

The Content of the	- FWXF-XF-X	1 cont	root	51592 Feb 29 2	2020	vwininfo								
Tenting Foot Foot							1 [11]					4%]	5 [11]	6.3%]
Total Continue 1												exi		9.0%
MARY 1 Foot Foot 1 1 1 1 1 1 1 1 1												5%		4.5%
Information Front														7.8%
Tenthalman Front Foot									11111	0.000	1.3G/	961		
Text														
Interviews Foot F												-		
**************************************	lrwxrwxrwx	1 root		6 Apr 8 2	2022	xzegrep -> xzgrep								
Part Frank	lrwxrwxrwx	1 root	root	6 Apr 8 2	2022	xzfgrep -> xzgrep	PID USER	PRI	NI	VIRT RES		CPU%	MEM% TIME+	
### STRICT FOOT 2216 Apr 8 2022 2026	- FWXF-XF-X	1 root	root	5902 Apr 8 2	2022	xzgrep	3656 ziming	20	0 5	6039M 420M	122M S	12.9	1.3 2h32:39	/usr/bin/gnome-shell
12533 trining	- FWXF-XF-X	1 root	root	1802 Apr 8 2	2022		89938 ziming	20	0			6.4	1.6 1h13:34	/home/ziming/.local/share/Steam/ubuntu
INTERNATION 1 cot	- FWXF-XF-X	1 root	root				89944 ziming	20	0 4			4.5	1.6 53:39.39	
International Continues Troot Foot 27 Apr 20 2000 ymallobj-10 111/10 111/10 121/10 111/10 12	- FWXF-XF-X	1 root	root	590 Mar 4 2	2025		121533 ziming	20	0 1	15460 5764	3268 R	3.9	0.0 0:00.81	htop
International Continues Front Fr	lrwxrwxrwx	1 root	root	27 Oct 14 2	2022	<pre>yaml2obj ->/lib/llvm-16/bin/yaml2obj</pre>	3489 ziming	20			76160 S	3.2	0.4 1h33:48	/usr/lib/xorg/Xorg vt2 -displayfd 3 -a
Interview is root	lrwxrwxrwx	1 root	root					20						/proc/self/exetype=renderercrash
Internative 1 root	lrwxrwxrwx	1 root	root						0 (opt/zoom/zoom zoommtg://northeastern.
Invariance Toot T														
Invariance Toot	lrwxrwxrwx	1 root	root	27 May 30 2										
Instruction Tool														
Instruction Tool														
Therefore														
Impurence 1														
-mwr.xr.x root os 3720 Apr 17 11:25 109 mst 20 0 129 380M 124M 5 1.3 1.2 43:12.11 pot/google/chrone/chrome -fwr.xr.x root root 3720 Apr 17 11:25 109 mst 20 0 8556 3140 5756 5 1.3 0.0 31:49.22 usr/llbjuetoothly bluetoothly blueto														
1 109 100 109														
Track - x - x Toot Toot Toot Sull 19 200 Sull														
Transfrance														
Transformer Troot Foot														
1918 Apr. 23 2916 yusplittoppn 5889 zining 20 6.341.19M 5.06 6.6 21:58.35 Opt./Google/chrome/chrome - trwxr-xr-x 1 root root 1918 Apr. 23 2916 yuvtoppn 5989 zining 20 6.3411.19M 2919 2318 5.06 6.6 21:58.35 Opt./Google/chrome/chrome - trwxr-xr-x 1 root root 1984 Apr. 23 2916 yuvtoppn 5989 zining 20 6.3411.19M 2919 2318 5.06 6.6 3:05.37 Opt./zoom/ZoomwebviewHost - trwxr-xr-x 1 root root 1678 Apr. 8 2922 zcap 112264 zining 20 0 1068M 1994 1538 5.0 6.6 3:05.37 Opt./zoom/ZoomwebviewHost - trwxr-xr-xr-x 1 root root 5988 Apr. 8 2922 zcap 112264 zining 20 0 1068M 1994 5389 5389 78 2922 zcap 112264 zining 20 0 1068M 1994 50 6.0 6.5 21:31.39 Opt./google/chrome/chrome - trwxr-xr-xr-xr-xr-xr-xr-xr-xr-xr-xr-xr-xr-x														
Track - x - x root root 10104 App 23 2016 yuvsplittoppm 58895 zining 20 0.3418 12918 2338 5 0.6 4.1 2111.84 /ojt //zoon/zoom/com/edv/edv/eds - t-type rack - x root root 10104 App 23 2016 yuvsplittoppm 59004 zining 20 0.47.75 7799 1006 3.150.37 /ojt //zoon/zoom/com/edv/edv/eds - t-type rack - x root root 1078 App 8 2022 zcat 89937 zining 20 0.1088 1938 1538 5 0.6 0.6 0.11.39 /home/zining / rack - x root root 5076 App 8 2022 zcat 2204 2		1 root	root	47 Feb 17 2	2020	<pre>yplan ->/share/texlive/texmr-dist/scripts/yplan/y</pre>								
FWKY-KY-X root 10194 Apr 23 2016 vuvoppn 10194 Apr 23 2016 vuvoppn 10194 Apr 23 2016 2016				40404 4 33	2046									
PROF. NRT. X Tool														
12264 zining 20 1.27 671M 145% 5 6 2.1 2:29.55 popt/goolejchrone/chrometymer-x-r-x root root 26840 May 26 04:69 2dump 42250 zining 20 20.281M 180M 5 6 6 6 6.1 13.179 popt/goolejchrone/chrometymer-x-r-x root root 26840 May 26 04:69 2dump 42250 zining 20 20.281M 180M 10.65 6 6 6 1.371 2dump 20.281M 180M 10.65 6 6 6 6 1.379 popt/goolejchrone/chrometymer-x-r-x root root 26840 May 26 04:69 2dump 42250 zining 20 0 33.36 792M 4211 5 6 6 5 511:33.99 popt/goolejchrone/chrometymer-x-r-x root root 10104 Apr 23 2016 zeisstopnn 72621 zining 20 0 1.271 2dsM 102M 5 6 6 5 511:33.99 popt/goolejchrone/chrometymer-x-r-x root root 10104 Apr 23 2016 zeisstopnn 72621 zining 20 0 1.271 2dsM 102M 2dsM 102M 2dsM														
Track - xr - x 1 root 5898 Apr 8 2022 2diff 20840 Apr 8 2022 2durp 4250 2durh 20 0 22811 186M 186M S 0.6 0.6 181.79 90tt/placaltionetworks/globalpro-rewr-xr-x 1 root 700t 29 Apr 8 2022 2grep 42583 2durh 20 0 1.27 148M 180M S 0.6 0.5 1813.39 90tt/placaltionetworks/globalpro-rewr-xr-x 1 root 700t														
-rwar-xr-x 1 root root 26840 May 26 94:09 zdump														
FYRKT-KT-X TOOL FOOL 29 Apr 8 2022 zegrep 42583 ziming 20 0 1.27 1488 102N S 0.6 0.5 0.19.26 Opt/google/chrome/chrometyp FYRKT-KT-X TOOL FOOL 135966 Feb 27 2020 zentty 72658 ziming 20 0 1.27 245N 1.29N S 0.6 0.8 0.13.43.49 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72658 ziming 20 0 1.27 245N 1.29N S 0.6 2.6 4150.75 Opt/google/chrome/chrometyp 72678 ziming 20 1.27 245N 245N S 245														
Transfer														
Track-xr-x 1 root root 135966 feb 27 2020 zentty 72658 ziming 20 0 1.27 6608 159M S 0.6 2.1 134.34 9 70pt/google/chrome/chrometyp 106097 ziming 20 0 1.27 819M 102M S 0.6 2.1 134.34 9 70pt/google/chrome/chrometyp 107097 ziming 20 0 1.27 819M 102M S 0.6 2.6 4159.75 0pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.9 1102.42 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.9 1102.42 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.9 1102.42 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.9 1102.42 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.9 1102.42 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.1 1147.37 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.1 1147.37 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.9 1102.42 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.1 1147.37 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 0.6 0.1 1147.37 70pt/google/chrome/chrometyp 107367 ziming 20 0 1.27 819M 102M S 200.00 20														
Trigory No.														
-rwxr-xr-x 1 root root 2081 Apr 8 2022 zforce 107160 ziming 20 0 1.27 2098 1260 S 0.6 0.9 11:04.2 /opt/google/chrome/chrometyp -rwxr-xr-x 1 root root 216256 Apr 21 2017 ztp 89906 ziming 20 0 33.06 1478 90816 S 0.0 0.5 6:05.33 /home/ziming/local/share/Steam-rwxr-xr-x 1 root root 216256 Apr 21 2017 ztp 89906 ziming 20 0 33.06 1478 90816 S 0.0 0.5 6:05.33 /home/ziming/local/share/Steam-rwxr-xr-x 1 root root 99316 Apr 21 2017 ztpcloak 89906 ziming 20 0 33.06 1478 90816 S 0.0 0.5 6:05.33 /home/ziming/local/share/Steam-rwxr-xr-x 1 root root 50718 Nov 23 2023 ztpdetatls 2212 root 20 0 250M 10984 8652 S 0.0 0.5 6:05.33 /home/ziming/local/share/Steam-rwxr-xr-x 1 root root 2953 Feb 1 2024 ztpgrep 3500 ziming 20 0 800M 1258 70160 S 0.0 0.4 4:05.09 /usr/llb/slack/slack-rtype=Zug-rwxr-xr-x 1 root root 89488 Apr 21 2017 ztpnote 59040 ziming 20 0 6.3418 12218 S 0.0 1.4 32:2639 usr/llb/slack/slack-rtype=Zug-rwxr-xr-x 1 root root 2095 Apr 20 2023 zisdecode 42272 ziming 20 0 6.3418 12918 2338 S 0.0 4.1 1:14.83 /opt/zoom/zoom/zoom/zoom/tor-rwxr-xr-x 1 root root 2095 Apr 20 2023 zisdecode 42272 ziming 20 0 32.46 1858 1198 S 0.0 0.6 23:56.67 /opt/google/chrome/chrometyp 2007 xr-xr-xr-x 1 root root 1842 Apr 8 2022 znore 50740 ziming 20 0 32.48 0544 76032 S 0.0 0.1 42:21.25														
Transfer														
-rwxr-xr-x 1 root root 216256 Apr 21 2017 zip -rwxr-xr-x 1 root root 93816 Apr 21 2017 zipclask -rwxr-xr-x 1 root root 50718 Nov 23 2023 zipdetails -rwxr-xr-x 1 root root 50718 Nov 23 2023 zipdetails -rwxr-xr-x 1 root root 50718 Nov 23 2023 zipdetails -rwxr-xr-x 1 root root 50718 Nov 23 2023 zipdetails -rwxr-xr-x 1 root root 2053 Feb 1 2024 zipprep -rwxr-xr-x 2 root root 186604 Feb 1 2024 zipprep -rwxr-xr-x 1 root root 89488 Apr 21 2017 zipnote -rwxr-xr-x 1 root root 89488 Apr 21 2017 zipnote -rwxr-xr-x 1 root root 2055 Mar 20 2023 zipdetails -rwxr-xr-x 1 root root 89488 Apr 21 2017 zipnote -rwxr-xr-x 1 root root 2055 Mar 20 2023 zipdetails -rwxr-xr-x 1 root root 2055 Mar 20 2023 zipdetails -rwxr-xr-x 1 root root 89488 Apr 21 2017 zipnote -rwxr-xr-x 1 root root 2055 Mar 20 2023 zipdetails -rwxr-xr-x 1 root root 2055 Mar 20 2023 zipdetails -rwxr-xr-x 1 root root 2055 Mar 20 2023 zipdetails -rwxr-xr-x 1 root root 2055 Mar 20 2023 zipdetails -rwxr-xr-x 1 root root 2055 Mar 20 2023 zipdetails -rwxr-xr-x 1 root root 2056 Mar 8 2022 zeles -rwxr-xr-x 1 root root 1842 Apr 8 2022 zonore -rwxr-xr-x 1 root root 2206 Apr 8 2022 zeles -rwxr-xr-x 1 root root 2206 Apr 8 2022 zeles -rwxr-xr-x 1 root root 2218 Feb 20 2019 zrun -rwxr-xr-x 1 root root 2218 Feb 20 2019 zrun -rwxr-xr-x 1 root root 2218 Feb 20 2019 zrun -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 -rwxr-														
-rwxr-xr-x 1 root root 93816 Apr 21 2017 zipcloak 89940 ziming 20 0 47.66 499N 167M S 0.0 1.6 5:36.36 /home/ziming/.local/share/Steam-rwxr-xr-x 1 root root 50718 Nov 23 2023 zipdetails 2212 root root 50718 Nov 23 2023 zipdetails 2212 root root 50718 Nov 23 2023 zipdetails 2212 root root root 2953 Feb 1 2024 zipgrep 3500 ziming 20 0 800M 125M 76160 S 0.0 0.0 4 4:05.09 /usr/lib/yorg/korg vtz -dtsplay -rwxr-xr-x 1 root root 186664 Feb 1 2024 zipgrep 3500 ziming 20 0 800M 125M 76160 S 0.0 0.0 4 4:05.09 /usr/lib/yorg/korg vtz -dtsplay -rwxr-xr-x 1 root root 186664 Feb 1 2024 zipgrep 3500 ziming 20 0 1.47 443M 122M S 0.0 1.4 3:226.93 /usr/lib/yorg/korg vtz -dtsplay -rwxr-xr-x 1 root root 93584 Apr 21 2017 zipsplit 20 0 6341M 129M 233M S 0.0 4.1 1:14.85 /usr/lib/slack/slacktype=zygrep 20 0 270M 10876 9860 S 0.0 0.0 10:48.50 /usr/sbin/thermaldsystemdrwxr-xr-x 1 root root 20552 Mar 20 2023 zipse 49958 ziming 20 0 32.46 185M 119M S 0.0 0.6 23:56.67 /porc/setif/exetype=renderer -rwxr-xr-x 1 root root 1842 Apr 8 2022 znew 1750 root 206 Apr 1870 200 0 32.30 85644 76032 S 0.0 0.3 3:54.52 /porc/setif/exetype=renderer -rwxr-xr-x 1 root root 22 202 Zibs 2 zibs 2 200M 2507 Ziming 20 0 32.30 85644 76032 S 0.0 0.3 3:54.52 /porc/setif/exetype=renderer -rwxr-xr-x 1 root root 22 202 Zibs 2 zibs 2 200M 2507 Ziming 20 0 32.30 85644 76032 S 0.0 0.3 3:54.52 /porc/setif/exetype=renderer -rwxr-xr-x 1 root root 8777 Apr 8 2022 znew 1750 root 22 200 Zibs 200M 2507 Ziming 20 0 32.30 85644 76032 S 0.0 0.1 3 3:54.52 /porc/setif/exetype=renderer -rwxr-xr-x 1 root root 878288 Mar 11 2022 zib 49958 ziming 20 0 32.47 430M 103M S 0.0 1.4 1:11.36 /porc/setif/exetype=renderer -rwxr-xr-x 1 root root 878288 Mar 11 2022 zib 4958 ziming 20 0 1.47 430M 103M S 0.0 1.4 1:11.36 /porc/setif/exetype=renderer -rwxr-xr-x 1 root root 883 Feb 21 2020 zibs 500 S 0.0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0														
-rwxr-xr-x 1 root root 50718 Nov 23 2023 zipdetalls 2212 root 20 0 2508 10984 8652 S 0.0 0.0 7:15.85 /usr/lib/pupower/upowerd 3500 ziming 20 0 8008 1258 76160 S 0.0 0.4 4:05.09 xir/lib/sorcy/xorg vtz disslay -rwxr-xr-x 1 root root 18664 Feb 1 2024 zipinfo 50764 ziming 20 0 8.081 1258 76160 S 0.0 0.4 4:05.09 xir/lib/sorcy/xorg vtz disslay -rwxr-xr-x 1 root root 89488 Apr 21 2017 zipnote 50764 ziming 20 0 8.041 1228 S 0.0 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 0.1 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 0.1 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 0.1 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 0.1 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 0.1 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 0.1 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 0.1 1.4 32:26.93 /usr/lib/slack/slacktype=zyg 59040 ziming 20 0 8.041 1228 S 0.0 0.1 1.4 1221.26 /opt/google/chrome/chrometype 222 ziming 20 8.041 1228 S 0.0 0.1 1.4 1212.26 /usr/lib/slack/slacktype=zyg 59040 ziming 20 8.041 1228 S 0.0 0.1 1.4 1212.26 /usr/lib/slack/slacktype=zyg 59040 ziming 20 8.041 1228 S 0.0 0.1 1.4 1212.26 /usr/lib/slack/slacktype=zyg 59040 ziming 20 8.041 1228 S 0.0 0.1 1.4 1212.26 /usr/lib/slack/slacktype=renderer -rwxr-xr-x 1 root root 873288 Mar 11 2022 zimi 49954 ziming 20 8.044 6348 S 0.0 1.3 3:57.29 /usr/bin/nautilustype=renderer -rwxr-xr-x 1 root root 873288 Mar 11 2022 zimi 49954 ziming 20 8.044 838 1228 S 0.0 1.4 1:12.278 /usr/lib/slack/slacktype=renderer -rwxr-xr-x 1 root root 843 Feb 21 2020 zish5 50768 ziming 20 8.044 1438 1228 S 0.0 1.4 1:12.278 /usr/lib/slack/slacktype=zyg 50768 ziming 20 8.044 1438 1228 S 0.0 1.4 1:12.278 /usr/lib/slack/slacktype=zyg 50768 ziming 20 8.044 1438														
-rmar-xr-x 1 root root 2953 Feb 1 2024 ztpgrep 3500 ztming 20 0 800M 125% 76160 S 0.0 0.4 4:05.09 /usr/llb/xorg/Xorg vtz -display														
-rwxr-xr-x 2 root root 186664 Feb 1 2024 zipinfo 59764 ziming 20 0 1.47 443M 122M S 0.0 1.4 32:26.93 /wsr/lib/slack/slacktype=zyg-rwxr-xr-x 1 root root 89488 Apr 21 2017 zipnote 59040 ziming 20 0 6341M 1291M 233M S 0.0 1.4 32:26.93 /wsr/lib/slack/slacktype=zyg-rwxr-xr-x 1 root root 93584 Apr 21 2017 zipnote 2471 rows 20 0 270M 10876 9860 S 0.0 0 4.1 1:14.83 /poz/zoom/zoom zoommtg://horther-rwxr-xr-x 1 root root 20952 Mar 20 2023 zjsdecode 42272 ziming 20 0 32:46 189M 119M S 0.0 0.6 23:56.67 /opt/google/chrome/chrometype-rwxr-xr-x 1 root root 22060 Apr 8 2022 zers 49958 ziming 20 0 32:46 189M 119M S 0.0 0.6 23:56.67 /opt/google/chrome/chrometype-rwxr-xr-x 1 root root 1842 Apr 8 2022 zmore 50740 ziming 20 0 32:36 8064 76032 S 0.0 0.3 3:54.52 /proc/self/exetype-rendererrwxr-xr-x 1 root root 4577 Apr 8 2022 zmew 1750 root 1750 root 22 Aug 25 2024 zoom -> /opt/zoom/zoomLauncher 59953 ziming 20 0 6341M 1291M 233M S 0.0 4.1 1:16.08 /opt/zoom/zoom zoommtg://horther-rwxr-xr-x 1 root root 22 Aug 25 2024 zoom -> /opt/zoom/zoomLauncher 45809 ziming 20 0 1.61 430M 103M S 0.0 1.3 3:57.29 /usr/bin/nauttlusgapplicaticrwxr-xr-x 1 root root 843 Feb 21 2020 zsh5 49954 ziming 20 0 1.61 4348 122M S 0.0 1.4 1:11.36 /proc/self/exetype-rendererrwxr-xr-x 1 root root 843 Feb 21 2020 zsh5 59768 ziming 20 0 1.61 4348 122M S 0.0 1.4 1:12.78 /usr/bin/slack/slacktype-zyg-yom-bin														
-rwxr-xr-x 1 root root 89488 Apr 21 2017 zipnote 59940 ziming 20 0 6341M 1291W 233M S 0.0 4.1 1:14.83 /ppt/zoom/zoommtos//horther-rwxr-xr-x 1 root root 2952 Mar 20 22 zipsole 2211 root 20 0 zoom/soom/soommtos/rhorther-rwxr-xr-x 1 root root 2952 Mar 20 222 zipsole 42272 ziming 20 0 zining 20 0 32.46 185H 119M S 0.0 0.5 zipsole/chrome/chrometyxr-xr-x 1 root root 2206 Apr 8 2022 zipsole 42272 ziming 20 0 zining 20														
-rwxr-xr-x 1 root root 93584 Apr 21 2017 zipsplit 2471 root 20 0 270M 10876 9860 S 0.0 0.0 10:48.50 /usr/ship/thermaldsystendcycr-xr-x 1 root root 20952 Mar 20 2023 zjsdecode 42272 ziming 20 0 32.46 185M 119M S 0.0 0.6 23:56.67 /opt/google/chrome/chrometycr-xr-x 1 root root 2206 Apr 8 2022 zless 49958 zlming 20 0 1.45 430M 163M S 0.0 1.4 2:21.26 /proc/setf/exetype=rendererrwxr-xr-x 1 root root 1842 Apr 8 2022 znew 50740 zlming 20 0 32.36 89644 76032 S 0.0 0.3 3:54.52 /proc/setf/exetype=rendererrwxr-xr-x 1 root root 4577 Apr 8 2022 znew 1750 root 22 Aug 25 2024 zoom -> /opt/zoom/ZoomLauncher 58953 zlming 20 0 6.34M 1291M 233M S 0.0 4.1 1:16.08 /opt/zoom/zoom.zoomnto://norther-rwxr-xr-x 1 root root 878288 Mar 11 2022 zsh 49954 zlming 20 0 1.47 430M 163M S 0.0 1.4 1:11.36 /proc/setf/exetype=rendererrwxr-xr-x 1 root root 8843 Feb 21 2020 zsh5 50768 zlming 20 0 1.47 430M 163M S 0.0 1.4 1:12.78 /usr/bln/nautilusgapplicationrwxr-xr-x 1 root root 843 Feb 21 2020 zsh5 50768 zlming 20 0 1.47 430M 163M S 0.0 1.4 1:12.78 /usr/bln/slack/slacktype=zyg-yb-bin														
-rwxr-xr-x 1 root root 26952 Mar 20 2023 zjsdecode 4227z ztming 20 0 32-46 185% 119M S 0.0 0.6 23:56.67 /opt/google/chrome/chrometymer-xr-x 1 root root 2206 Apr 8 2022 zers 49958 ztming 20 0 32-46 185% 119M S 0.0 0.6 23:56.67 /opt/google/chrome/chrometymer-xr-x 1 root root 1842 Apr 8 2022 zmore 50740 ztming 20 0 32-46 186% 14 300M 103M S 0.0 0.6 23:56.56 /opt/google/chrome/chrometymer-xr-x 1 root root 700 1842 Apr 8 2022 zmore 50740 ztming 20 0 32-46 186% 14 300M 103M S 0.0 0.6 23:56.56 /opt/google/chrome/chrometymer-xr-x 1 root root 700 4577 Apr 8 2022 zmore 50740 ztming 20 0 270M 10876 9860 S 0.0 0.0 1048.87 /usr/sbin/thermaldsystemdtwxr-xr-x 1 root root 22 Aug 25 2024 zoom -> /opt/zoom/zoomLauncher 58953 ztming 20 0 3341M 1291M 233M S 0.0 0.4 1 11:16.08 /opt/zoom/zoomzoomtsj://norther-rwxr-xr-x 1 root root 700 2518 Feb 20 2019 zrun 45809 ztming 20 0 1958M 424M 63484 S 0.0 1.3 3:57.29 /usr/bln/nauttlusgapplicatic 49954 ztming 20 0 1.4 430M 120M 30M S 0.0 1.4 1:11:2.78 /usr/bls/sack/slacktype=zype-btm														
-rwxr-xr-x 1 root root 2206 Apr 8 2022 zless 49958 zlming 20 0 3.4T 430M 163M 5 0.0 1.4 2:21.26 /proc/setf/exetype=renderer 50740 zlming 20 0 32.36 89644 76032 5 0.0 0.3 3:54.52 /proc/setf/exetype=renderer 50740 zlming 20 0 32.36 89644 76032 5 0.0 0.3 3:54.52 /proc/setf/exetype=renderer 50740 zlming 20 0 270M 108To 9860 5 0.0 0.0 10:48.37 /pur/sbln/thermaidsystemd trwxr-xr-x 1 root root 22 Aug 25 2024 zoom -> /opt/zoom/zoomLauncher 58953 zlming 20 0 6341M 1291M 233M 5 0.0 4.1 1:10.08 /opt/zoom/zoom zoomntg://norther -rwxr-xr-x 1 root root 878288 Mar 11 2022 zsh 49954 zlming 20 0 1.4T 430M 163M 5 0.0 1.4 1:11.36 /proc/setf/exetype=renderer -rwxr-xr-x 1 root root 878288 Mar 11 2022 zsh 49954 zlming 20 0 1.4T 443M 122M 5 0.0 1.4 1:22.78 /psr/tbl/slack/slacktype=zype-bin														
-rwxr-xr-x 1 root root 1842 Apr 8 2022 zmore 59740 ztming 20 0 32:36 89644 76032 S 0.0 0.3 3:54.52 /proc/setf/exetype-utilityrwxr-xr-x 1 root root 4577 Apr 8 2022 zmew 1759 root 20 0 27040 18876 9860 S 0.0 0.0 10:48.87 /yors/sbin/thermaidsystemdtwxr-xr-x 1 root root 22 Aug 25 2024 zoom -> /opt/zoom/zoomLauncher 58953 ztming 20 0 6341M 1291M 233M S 0.0 4.1 1:16.08 /opt/zoom/zoom zoomntg://norther-rwxr-xx 1 root root 82518 Feb 20 2019 zrun 45809 ztming 20 0 1958M 424M 63484 S 0.0 1.3 3:57.29 /usr/bin/nautilusspapilcastor 49954 ztming 20 0 1958M 424M 6348 S 0.0 1.4 1:1:2.78 /vsr/bin/reketype-renderertwxr-xr-x 1 root root 843 Feb 21 2020 zsh5 50768 ztming 20 0 1.41 438 122M S 0.0 1.4 1:1:2.78 /usr/bin/slack/slacktype-zype-bin														
-rwxr-xr-x 1 root root 4577 Apr 8 2022 znew 1759 root 22 Aug 25 2024 zoon -> /opt/zoom/ZoomLauncher 58953 ziming 20 0.6341M 1291M 2331 S 0.0 4.1 1:10.08 /opt/zoom/zoom/to://nort-rwxr-xr-x 1 root root 2518 Feb 20 2019 zrun 45809 ziming 20 0.1958M 424M 63484 S 0.0 4.1 1:10.08 /opt/zoom/zoom/to://nort-rwxr-xr-x 1 root root 878288 Mar 11 2022 zsh 49954 ziming 20 0.1958M 424M 63484 S 0.0 1.3 3:57.29 /usr/bin/nautilusgapplication-rwxr-xr-x 1 root root 878288 Mar 11 2022 zsh 49954 ziming 20 0.1958M 424M 63484 S 0.0 1.4 1:11.36 /proc/self/exetype=rendererrwxr-xr-x 1 root root 843 Feb 21 2020 zsh 5 50768 ziming 20 0.147 433M 1221 S 0.0 1.4 1:22.78 /usr/lib/slack/slacktype=zyg-y-bin														
Trixkrikkrikk														
-rwxr-xr-x 1 root root 2518 Feb 20 2019 zrun 45809 ziming 20 0 1958M 424M 63484 S 0.0 1.3 3:57.29 /usr/bin/nautilusgapplicatio -rwxr-xr-x 1 root root 878288 Mar 11 2022 zsh 49954 ziming 20 0 1.47 430M 163M S 0.0 1.4 1:11.36 /proc/self/exetype=renderer -rwxr-xr-x 1 root root 843 Feb 21 2020 zsh5 50768 ziming 20 0 1.47 433M 122M S 0.0 1.4 1:22.78 /usr/lib/slack/slacktype=zyg + bin F1Help F2Setup F3_SearchF4+IllerF5_wree F0_SortsyF7_Wcce F8_Hcce F9_Kllce_F9_Kllc			root			zoom -> /opt/zoom/ZoomLauncher	58953 ziming	20						
-rwxr-xr-x 1 root root 878288 Mar 11 2022 zsh 49954 ztming 20 0 1.47 430% 163M 5 0.0 1.4 1:11.36 /proc/seif/exetype=genderer -rwxr-xr-x 1 root root 843 Feb 21 2020 zsh5 50768 ztming 20 0 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slack/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slacktype=zype b 1.47 443M 122M 5 0.0 1.4 1:22.78 /usr/tlb/slacktype=z			root											
-rwxr-xr-x 1 root root 843 Feb 21 2020 zsh5 50768 ziming 20 0 1.4T 443M 122M S 0.0 1.4 1:22.78 /usr/lib/slack/slacktype=zyg → bin F1Help F2 Setup F3 SearchF4Filter F5 ree F6 SortByF7Nice F8 Mice F9 Kill F10 Jutt	- FWXF-XF-X	1 root	root					20						
→ bin F1/elp F2/setup F3/search F4/Filter F5/Tree F6/sort8/F7/Vice -F8/Vice +F9/Vill F10/Vit	- rwxr-xr-x	1 root	root	843 Feb 21 2	2020			20	0					
	→ bin					7-96	F1Help F2Setup	F3Sear	chF4					9Kill <mark>F10</mark> Quit
[0] 0:ntop: "ztming-XPS-13-9300" 14:11	[0] 0:htop*													"ziming-XPS-13-9300" 14:11 29-Sep-

Programs and Processes



Programs and Processes



Programs and Processes

1:13 /usr/lib/slack/slack --type=zy

32:27 /usr/lib/slack/slack --type=zygote

4:07 /proc/self/exe --type=utility --utilit

```
ziming
                 0.0 0.1 34136220 61784 ?
                                                      Sep25
                                                              0:06 /proc/self/exe --type=utility --utilit
ziming
                 0.0 0.1 1067900 57816 ?
                                                      Sep25
                                                              0:01 /usr/bin/qnome-calendar --gapplication
ziming
                 0.0 0.0
                           7708 1968 ?
                                                      Sep25
                                                              0:00 /usr/bin/zoom zoommtg://northeastern.;
ziming
                                                 SLL
                                                      Sep25 156:39 /opt/zoom/zoom zoommtg://northeastern.
                 2.9 4.0 6493848 1322692 ?
ziming
                 0.0 0.6 2400544 209692
                                                sl
                                                      Sep25
                                                              0:12 /opt/zoom/ZoomWebviewHost
ziming
                 0.0 0.2 34173764 71188 3
                                                      Sep25
                                                              0:00 /opt/zoom/ZoomWebviewHost --type=zygot
ziming
                 0.0 0.2 34173752 68888 3
                                                      Sep25
                                                              0:00 /opt/zoom/ZoomWebviewHost --type=zygot
ziming
                                                      Sep25
           58833
                 0.0 0.0 34173776 17440 ?
                                                              0:00 /opt/zoom/ZoomWebviewHost --type=zygot
ziming
                 0.0 0.3 35602488 119860 ?
                                                sl
                                                      Sep25
                                                              0:11 /opt/zoom/ZoomWebviewHost --type=gpu-p
zimina
           58883
                 0.0 0.3 1095732 98028 ?
                                                sl
                                                      Sep25
                                                              0:01 /opt/zoom/ZoomWebviewHost --tvpe=utili
ziming
                 0.0 0.1 34545108 45116 3
                                                sl
                                                      Sep25
                                                              0:00 /opt/zoom/ZoomWebviewHost --type=utili
ziming
                 0.0 0.4 49749360 160980 ?
                                                      Sep25
                                                              0:02 /opt/zoom/ZoomWebviewHost --type=rende
ziming
                 0.0 0.5 49977456 183792
                                                      Sep25
                                                              3:05 /opt/zoom/ZoomWebviewHost --type=rende
ziming
                                                sl
                                                              0:04 /opt/zoom/ZoomWebviewHost --type=rende
                 0.0 0.6 49751972 195264 ?
                                                      Sep25
zimina
                 0.0 0.5 1283861032 184084 ?
                                                sl
                                                      Sep26
                                                              0:23 /opt/google/chrome/chrome --type=rende
ziming
                 0.0 0.6 1285596768 201684 ?
                                                sl
                                                      Sep26
                                                              0:18 /opt/google/chrome/chrome --type=rende
ziming
           72621 0.0 0.7 1283877664 251512 ?
                                                      Sep26
                                                              0:33 /opt/google/chrome/chrome --type=rende
ziming
                 0.0 0.4 1283855656 142184 ?
                                                sl
                                                      Sep26
                                                              0:13 /opt/google/chrome/chrome --type=rende
zimina
                                                sl
                                                      Sep26
                 0.0 0.4 1283855580 134700 ?
                                                              0:12 /opt/google/chrome/chrome --type=rende
zimina
                                                      Sep26
                                                             13:43 /opt/google/chrome/chrome --type=rende
                 0.3 2.1 1289104872 690636 ?
ziming
                 0.0 0.1 471484 46088 ?
                                                 SLl
                                                      Sep26
                                                              0:01 /usr/bin/seahorse --gapplication-servi
                                                      Sep26
ziming
                 0.0 0.0 12972
                                                              0:00 bash /home/ziming/.local/share/Steam/s
ziming
                                                              0:00 /home/ziming/.local/share/Steam/ubuntu
          89535 0.0 0.0 29684
                                                      Sep26
ziming
                 0.5 0.6 1094244 198576 ?
                                                      Sep26
                                                             22:35 /home/ziming/.local/share/Steam/ubuntu
           89693
ziming
                 0.0 0.0
                           4632
                                   676 ?
                                                 Ss
                                                      Sep26
                                                              0:00 /home/ziming/.local/share/Steam/steam/
ziming
          89729 0.0 0.0 27404
                                  3068 ?
                                                      Sep26
                                                              0:00 /home/ziming/.local/share/Steam/ubuntu
ziming
           89817 0.0 0.0 322500
                                                sl
                                                      Sep26
                                                              0:00 steam-runtime-launcher-service --along
ziming
                                                      Sep26
                                                              0:00 /usr/lib/pressure-vessel/from-host/lib
                 0.0 0.0 32276
                                                 Ss
ziming
                 0.3 1.1 3671232 362312 3
                                                      Sep26
                                                             14:41 ./steamwebhelper -nocrashdialog -lang=
                                                              0:00 /home/ziming/.local/share/Steam/ubuntu
zimina
           89865
                 0.0 0.1 419800 55896 ?
                                                      Sep26
                                                              0:00 /home/ziming/.local/share/Steam/ubuntu
ziming
                 0.0 0.2 34190312 67572
                                                      Sep26
ziming
                                                              0:00 /home/ziming/.local/share/Steam/ubuntu
                 0.0 0.2 34190300 67452
                                                      Sep26
ziming
                 0.0 0.0 34190324 16936 ?
                                                      Sep26
                                                              0:00 /home/ziming/.local/share/Steam/ubuntu
ziming
           89892
                 0.6 0.4 35278576 150892 ?
                                                sl
                                                      Sep26
                                                             26:13 /home/ziming/.local/share/Steam/ubuntu
zimina
           89917
                 0.0 0.2 1040808 92332 ?
                                                sl
                                                      Sep26
                                                              0:01 /proc/self/exe --type=utility --utilit
ziming
                 0.0 0.1 34562112 41236 ?
                                                sl
                                                      Sep26
                                                              0:00 /home/ziming/.local/share/Steam/ubuntu
ziming
                      1.5 49949696 507136 ?
                                                sl
                                                             73:43 /home/ziming/.local/share/Steam/ubuntu
ziming
                 0.0 0.2 49758268 87720 ?
                                                sl
                                                      Sep26
                                                              0:00 /home/ziming/.local/share/Steam/ubuntu
zimina
                 0.0 0.0 43996 6576 ?
                                                 Ss
                                                      Sep28
                                                              0:00 /usr/lib/bluetooth/obexd
          105305
zimina
                 0.0 0.6 1283876196 209712 ?
                                                      Sep28
                                                              0:12 /opt/google/chrome/chrome --type
ziming
                      0.4 1283872056 150432 ?
                                                sl
                                                      Sep28
                                                              0:07 /opt/google/chrome/chrome --type=renu
ziming
                      0.5 1283873736 172316 ?
                                                      Sep28
                                                              0:05 /opt/google/chrome/chrome --type=rende
                                                     Sep28
                                                              4:50 /opt/google/chrome/chrome --type=rende
ziming
          106097
                 0.2 2.5 1285679540 837956 ?
ziming
                 0.0 0.6 1285596948 213200 ?
                                                sl
                                                      Sep28
                                                              0:21 /opt/google/chrome/chrome --type=rende
          106178
ziming
                 0.0 0.4 1283863960 144060 ?
                                                sl
                                                      Sep28
                                                              0:05 /opt/google/chrome/chrome --type=rende
ziming
          106265
                      1.0 1295892516 356152 ?
                                                sl
                                                      Sep28
                                                              0:41 /opt/google/chrome/chrome --type=rende
                                                Sl
                                                      Sep28
ziming
                 0.0 0.5 1283860776 180696 ?
                                                              0:04 /opt/google/chrome/chrome --type=rende
                                                      Sep28
                                                              0:03 /opt/sublime text/sublime text --detac
ziming
          106653
                 0.0 0.3 776052 121148 ?
          106664 0.0 0.0 152684 1140 ?
                                                      Sep28
                                                              0:00 /opt/sublime_text/crash_handler --no-n
zimina
ziming
          106707 0.0 0.0 57032 16188 ?
                                                 SI
                                                      Sep28
                                                              0:00 /opt/sublime text/plugin host-3.3 1066
          106710 0.0 0.0 73760 25668 3
                                                              0:00 /opt/sublime text/plugin host-3.8 1066
```

Sep25

Sep25

50731 0.0 0.4 34498188 146668 ?

0.5 1.4 1461505468 462016 ?

50734 0.0 0.2 33900548 89724 ?

ziming

ziming

One-to-many relationship between program and processes

Why so many Chrome processes?

To shield the browser from attacks, Google Chrome developers eyed three key problems.

BY CHARLES REIS, ADAM BARTH, AND CARLOS PIZANO

Browser Security: Lessons from Google Chrome

https://dl.acm.org/doi/pdf/10.1145/1536616.1536634

THE WEB HAS become one of the primary ways people interact with their computers, connecting people with a diverse landscape of content, services, and applications. Users can find new and interesting content on the Web easily, but this presents a security challenge: malicious Web site operators can attack

Browsers face the challenge of keeping their users safe while providing a rich having their computers compromised. platform for Web applications.

a wide network-visible interface. Historically, every browser at some point has contained a bug that let a malicious Web site operator circumvent the browser's security policy and compromise the user's computer. Even after these vulnerabilities are patched, many users continue to run older, vulnerable | dering engine.

users through their Web browsers, | versions,5 When these users visit malicious Web sites, they run the risk of

Generally speaking, the danger Browsers are an appealing target for posed to users comes from three facattackers because they have a large and tors, and browser vendors can help complex trusted computing base with keep their users safe by addressing each of these factors:

▶ The severity of vulnerabilities. By sandboxing their rendering engine, browsers can reduce the severity of vulnerabilities. Sandboxes limit the damage that can be caused by an attacker who exploits a vulnerability in the ren-

How to Run a Program?

• How does the OS turn an executable file into a process?

What information must an executable file contain to run as a program?

Program Formats

- Programs obey specific file formats
 - Unix/Linux: Executable and Linkable Format (ELF)
 - Mac OSX: Mach object file format (Mach-O)
 - Windows Portable Executable (PE, PE32+) (*.exe)
 - Modified version of Unix COFF executable format
 - PE files start with an MZ header.
 - CP/M (control program monitor) and DOS (disk operating system)
 COM executables (*.com)
 - DOS: MZ executables (*.exe)
 - Named after Mark Zbikowski, a DOS developer

Program Formats

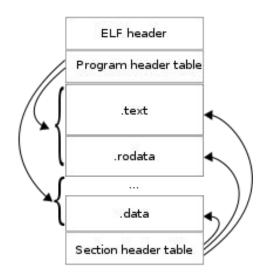
Format	Platform	Spec availability	Ownership
ELF	Linux, BSD, Unix	Open System V ABI	Community / Linux Foundation
PE	Windows	Public (Microsoft PE/COFF spec)	Microsoft
Mach-O	macOS, iOS	Public headers/docs by Apple	Apple

"System V" (often written SysV) was a line of commercial UNIX from AT&T in the 1980s. System V ABI is a set of open specifications defining: (1) the Executable and Linkable Format (ELF) for object files and executables. (2) Calling conventions, symbol tables, relocation formats, dynamic linking rules, etc.

The docs were published and made freely available, not proprietary or locked. That's why Linux, BSD, Solaris, and many other Unix-like OSes could adopt ELF with no licensing restrictions.

ELF File Format

- Spec: https://refspecs.linuxfoundation.org/elf/elf.pdf
- ELF Header
 - Contains compatibility info
 - Entry point of the executable code
- Program header table
 - Lists all the segments in the file
 - Used to load and execute the program
- Section header table
 - Used by the linker



ELF Header Example

```
@zzm7000 → /workspaces/assignment-3-zzm7000 (main) $ readelf -a ./a.out
ELF Header:
  Magic:
          7f 45 4c 46 02 01 01 00 00 00 00 00 00 00 00 00
  Class:
                                      ELF64
  Data:
                                      2's complement, little endian
  Version:
                                      1 (current)
  OS/ABI:
                                      UNIX - System V
  ABT Version:
                                      0
  Type:
                                      DYN (Position-Independent Executable file)
  Machine:
                                      Advanced Micro Devices X86-64
  Version:
                                      0x1
  Entry point address:
                                      0x1040
  Start of program headers:
                                      64 (bytes into file)
  Start of section headers:
                                      13928 (bytes into file)
  Flags:
                                      0x0
  Size of this header:
                                     64 (bytes)
  Size of program headers:
                                     56 (bytes)
  Number of program headers:
                                      13
  Size of section headers:
                                      64 (bytes)
  Number of section headers:
  Section header string table index: 28
Section Headers:
  [Nr] Name
                         Type
                                           Address
                                                             Offset
       Size
                         EntSize
                                           Flags Link Info Align
  [ 0]
                         NULL
                                           00000000000000000
                                                             00000000
       00000000000000000
                         00000000000000000
                                                                 0
                         PROGBITS
                                           00000000000000318
                                                             00000318
  [ 1] .interp
       000000000000001c 0000000000000000
```

00000000000000000

00000000000000000

NOTE

0000000000000338

0000000000000368

0

000000000000038c 0000038c

Α

00000338

8 00000368

4

[2] .note.gnu.pr[...] NOTE

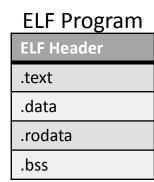
[3] .note.gnu.bu[...] NOTE

00000000000000024

[4] .note.ABI-tag

The Program Loader

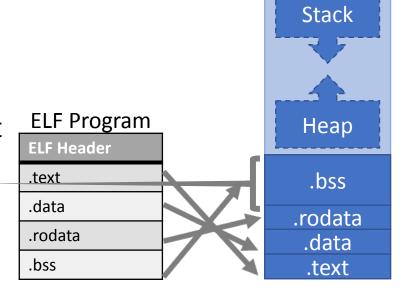
- OS functionality that loads programs into memory, creates processes
 - Places segments into memory
 - Loads necessary dynamic libraries
 - Performs relocation
 - Allocated the initial stack frame
 - Sets EIP to the programs entry point
- Process is a live program execution context or basic unit of execution



Memory

The Program Loader

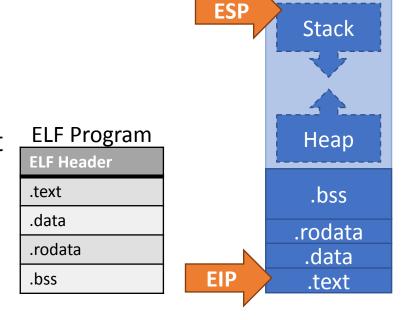
- OS functionality that loads programs into memory, creates processes
 - Places segments into memory
 - Loads necessary dynamic libraries
 - Performs relocation
 - Allocated the initial stack frame
 - Sets EIP to the programs entry point
- Process is a live program execution context or basic unit of execution



Memory

The Program Loader

- OS functionality that loads programs into memory, creates processes
 - Places segments into memory
 - Loads necessary dynamic libraries
 - Performs relocation
 - Allocated the initial stack frame
 - Sets EIP to the programs entry point
- Process is a live program execution context or basic unit of execution



Memory

xv6 as an example

Background: xv6 programs are statically linked. All the user programs you compile for xv6 (like ls, sh, etc.) are **fully statically linked** into one ELF file — no shared libraries. Kernel itself is the loader.

In exec.c (exec() in kernel/exec.c):

- The kernel opens the ELF file from the file system.
- Reads the ELF header (struct elfhdr) and each program header (struct proghdr).
- For each loadable segment, it allocates memory, maps it into the new process's address space, and copies the data from the file.
- Sets up the stack and the initial instruction pointer to the ELF's entry address.

https://github.com/mit-pdos/xv6-public/blob/master/exec.c

xv6 as an example

myproc() (kernel/proc.c)

Returns a pointer to the **current process** (struct proc *) running on this CPU. It looks up the per-CPU structure and fetches its proc field.

setupkvm() (kernel/vm.c)

Builds a **fresh kernel page table** (kernel virtual memory only): maps kernel text/data, devices, I/O space, etc. Used when creating a new address space for a process (before loading user pages).

allocuvm(pagetable, oldsz, newsz) (kernel/vm.c)

Grows a process's user address space from oldsz to newsz. Allocates physical pages and maps them into pagetable with user permissions. Returns the new size (or 0 on failure).

loaduvm(pagetable, va, ip, offset, sz) (kernel/vm.c / kernel/exec.c)

Loads program bytes from disk into user memory: reads sz bytes from inode ip at file offset and copies them into the user virtual region starting at virtual address va (assumes those pages are already allocated/mapped). Used by exec() to load ELF segments.

switchuvm(p) (kernel/vm.c / kernel/proc.c)

Switches the hardware MMU to a process's page table (and does any CPU/TSS setup on x86). Called when the scheduler is about to run process p, or when a process's memory image has just changed (e.g., after exec).

freevm(pagetable) (kernel/vm.c)

Frees a process's user page table and all user pages it points to (tears down the address space). Used on process exit or on exec failure paths.

Linux as an example

How it starts:

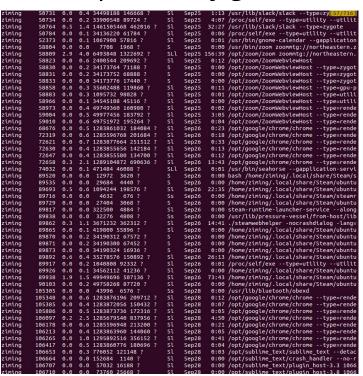
- You run an ELF file => kernel's execve() system call reads its ELF headers.
- 2. If the file has an INTERP section, the kernel:
 - Loads that interpreter binary (e.g., Id-linux-x86-64.so.2) into memory.
 - Jumps to its entry point in user space.
- 3. The interpreter (now running in ring 3) maps your program's segments and shared libraries, relocates symbols, and finally transfers control to your program's _start.

```
ldd /usr/bin/ls
linux-vdso.so.1 (0x00007ffc3ed1f000)
libselinux.so.1 => /lib/x86_64-linux-gnu/libselinux.so.1 (0x00007f9ec75ce000)
libc.so.6 => /lib/x86_64-linux-gnu/libc.so.6 (0x00007f9ec73dc000)
libpcre2-8.so.0 => /lib/x86_64-linux-gnu/libpcre2-8.so.0 (0x00007f9ec734b000)
libdl.so.2 => /lib/x86_64-linux-gnu/libdl.so.2 (0x00007f9ec7345000)
/lib64/ld-linux-x86-64.so.2 (0x00007f9ec764c000)
libpthread.so.0 => /lib/x86_64-linux-gnu/libpthread.so.0 (0x00007f9ec7322000)
```

Warmup

- How many processes do you have open at any given time?
 - 10s, 100s? More!?:)

ps -aux | wc -l



First: Instruction Execution

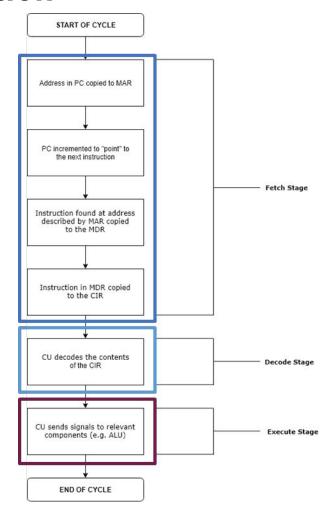
- Code in an executable is a sequence of instructions
- A CPU/core runs an instruction at a time
- This is done in a fetch-decode-execute cycle
- If you have 4 cores, your processor can do
 4 FDE cycles at a time
- But how do we see ~100s of programs running on 4 cores?
- What about a single core CPU?

Internal CPU registers that implement the fetch-decode-execute cycle:

Memory Address Register: holds address of current instruction

Memory Data Register: holds contents of address in MAR

Current Instruction Register: stores current instruction, so not overwritten by additional fetches to MBR/MDR



From the warm up

Many programs are running, but only 8 CPUs that

do the work

Iscpu

```
Architecture:
                                      x86 64
CPU op-mode(s):
                                      32-bit. 64-bit
Byte Order:
                                      Little Endian
Address sizes:
                                      39 bits physical, 48 bits virtual
CPU(s):
On-line CPU(s) list:
                                      0-7
Thread(s) per core:
Core(s) per socket:
Socket(s):
NUMA node(s):
Vendor ID:
                                      GenuineIntel
CPU family:
Model:
Model name:
                                      Intel(R) Core(TM) i7-1065G7 CPU @ 1.30GHz
Stepping:
CPU MHz:
                                      1500.000
CPU max MHz:
                                      3900.0000
CPU min MHz:
                                      400.0000
BogoMIPS:
                                      2995.20
Virtualization:
                                      VT-x
L1d cache:
                                      192 KiB
L1i cache:
                                      128 KiB
L2 cache:
                                      2 MiB
L3 cache:
                                      8 MiB
NUMA node0 CPU(s):
                                      0-7
Vulnerability Gather data sampling:
                                      Mitigation; Microcode
Vulnerability Itlb multihit:
                                      KVM: Mitigation: VMX disabled
Vulnerability L1tf:
                                      Not affected
Vulnerability Mds:
                                      Not affected
Vulnerability Meltdown:
                                      Not affected
Vulnerability Mmio stale data:
                                      Mitigation; Clear CPU buffers; SMT vulnerable
Vulnerability Reg file data sampling: Not affected
Vulnerability Retbleed:
                                      Mitigation; Enhanced IBRS
Vulnerability Spec rstack overflow:
                                      Not affected
Vulnerability Spec store bypass:
                                      Mitigation: Speculative Store Bypass disabled via prctl and seccomp
Vulnerability Spectre v1:
                                      Mitigation; usercopy/swapgs barriers and user pointer sanitizatio
```

From the warm up

 Many programs are running, but only 8 CPUs that do the work

<u>The Problem:</u> So how does our Operating System provide the illusion of hundreds of processes running at once?

Virtualization with time sharing

- If one CPU, the Operating System runs one process at a time,
 - That executes one instruction a time
 - After some amount of time the process stops or finishes
 - Then the OS starts another process
 - Eventually the same process will run again and continue where it left off
 - Repeat

Time	Process ₀	$Process_1$	Notes
1	Running	Ready	
2	Running	Ready	
3	Running	Ready	
4	Running	Ready	Process ₀ now done
5	_	Running	
6	_	Running	
7	<u></u>	Running	
8	_	Running	Process ₁ now done

Virtualization with time sharing

- If one CPU, the Operating System runs one process at a time,
 - That executes one instruction a time
 - After some amount of time the process stops or finishes
 - Then the OS starts another process
 - Eventually the same process will run again and continue where it left off

D +	,							
• Repeat	Time	Process ₀	$Process_1$	Notes				
	1	Running	Ready					
	2	Running	Ready					
	3	Running	Ready					
	4	Running	Ready	Process ₀ now done				
	5	_	Running					
	6	_	Running					
	7	_	Running					
	8	_	Running	Process ₁ now done				

- This concept is known as time sharing
- Are the two states, Running and Ready, enough?

• What if the process needs to read/write to disk or perform a network request? Any problems?

- What if the process needs to read/write to disk or perform a network request? Any problems?
 - These operations take (comparatively) long to complete

- What if the process needs to read/write to disk or perform a network request? Any problems?
 - These operations take (comparatively) long to complete
 - Keeping process state to Running?
 - Keeping process state to Ready?

- What if the process needs to read/write to disk or perform a network request? Any problems?
 - These operations take (comparatively) long to complete
 - Keeping process state to Running?
 - Hogs the CPU just waiting for disk/network access to complete
 - Keeping process state to Ready?

- What if the process needs to read/write to disk or perform a network request? Any problems?
 - These operations take (comparatively) long to complete
 - Keeping process state to Running?
 - Hogs the CPU just waiting for disk/network access to complete
 - Keeping process state to Ready?
 - Might not be ready to run when its turn comes
 - Asking it to run may be waste of time

- What if the process needs to read/write to disk or perform a network request? Any problems?
 - These operations take (comparatively) long to complete
 - Keeping process state to Running?
 - Hogs the CPU just waiting for disk/network access to complete
 - Keeping process state to Ready?
 - Might not be ready to run when its turn comes
 - Asking it to run may be waste of time
- Solution?

- What if the process needs to read/write to disk or perform a network request? Any problems?
 - These operations take (comparatively) long to complete
 - Keeping process state to Running?
 - Hogs the CPU just waiting for disk/network access to complete
 - Keeping process state to Ready?
 - Might not be ready to run when its turn comes
 - Asking it to run may be waste of time
- Solution?
 - Introduce a 3rd state, Blocked
 - Meaning: the process requested some I/O operation and cannot run until that operation is completed

- Each process can be in one of several states
- The OS schedules the state the process is in
- Typically, these are:
 - Running: the process is executing on the CPU
 - Ready: the process is ready to execute, but the OS did not choose to run it
 - Blocked the process issued some blocking operation
 - I/O is a common blocking operation

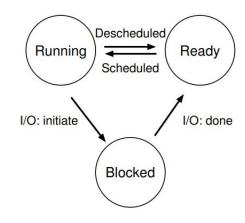


Figure 4.2: Process: State Transitions

xv6 as an example

- **UNUSED** empty slot in the process table; no process here.
- **EMBRYO** (x86) / **USED** (RISC-V) slot allocated and being initialized (after allocproc()), but not runnable yet.
- **SLEEPING** blocked, waiting on some event (a "channel"), e.g., disk I/O to finish.
- RUNNABLE ready to run; waiting for CPU.
- **RUNNING** currently executing on a CPU.
- **ZOMBIE** finished execution (exit() called), but parent hasn't wait()ed yet; holds exit status/resources to be reaped.

https://github.com/mit-pdos/xv6-public/blob/master/proc.h

Why I Don't Feel Apps Lagging Even Though They Time-Share the CPU?

Then how does OS switch processes?

OS Challenges to Virtualization

- Performance
 - How to implement virtualization without excessive overhead?

- Control
 - How to run multiple processes without losing control over the CPU?
 - Without OS control, a process
 - could occupy the CPU and run forever
 - access memory it does not have access impacting safety and security

Switching between processes

• Switching between processes is a challenge, because

If the CPU is running a program, then the OS is not running

Switching between processes

• Switching between processes is a challenge, because

If the CPU is running a program, then the OS is not running

- If OS is not running, then how can it switch out/in processes?
 - Think about how you would design the OS!

When Do You Switch Processes?

- To share CPU between multiple processes, control must eventually return to the OS
 - When should this happen?
 - What mechanisms implements the switch from user process back to the OS?

Four approaches:

When Do You Switch Processes?

- To share CPU between multiple processes, control must eventually return to the OS
 - When should this happen?
 - What mechanisms implements the switch from user process back to the OS?

- Four approaches:
 - 1. Voluntary yielding
 - 2. Switch during API calls to the OS
 - 3. Switch on I/O
 - 4. Switch based on a timer interrupt

Voluntary Yielding

 Idea: processes must voluntary give up control by calling an OS API, e.g. thread_yield()

Voluntary Yielding

 Idea: processes must voluntary give up control by calling an OS API, e.g. thread_yield()

Problems?

Voluntary Yielding

 Idea: processes must voluntary give up control by calling an OS API, e.g. thread_yield()

• Problems?

- Misbehaving or buggy apps may never yield
 e.g., while (1) { //do something without yielding }
- No guarantee that apps will yield in a reasonable amount of time
- Waste of CPU resources, i.e. what if a process is idle-waiting on I/O?

Interjection on OS APIs

- Idea: whenever a process calls an OS API, this gives the OS an opportunity to context switch
 - E.g. printf(), fopen(), socket(), etc...

- The original Apple Macintosh used this approach
 - Cooperative multi-tasking

Interjection on OS APIs

- Idea: whenever a process calls an OS API, this gives the OS an opportunity to context switch
 - E.g. printf(), fopen(), socket(), etc...

- The original Apple Macintosh used this approach
 - Cooperative multi-tasking

Problems?

Interjection on OS APIs

- Idea: whenever a process calls an OS API, this gives the OS an opportunity to context switch
 - E.g. printf(), fopen(), socket(), etc...

- The original Apple Macintosh used this approach
 - Cooperative multi-tasking

- Problems?
 - Misbehaving or buggy apps may never call OS APIs
 - Some normal apps don't use OS APIs for long periods of time
 - E.g. a long, CPU intensive matrix calculation

Switching on I/O

- Idea: when one process is waiting on I/O, switch to another process
 - I/O APIs already go through the OS, so context switching is easy

Switching on I/O

- Idea: when one process is waiting on I/O, switch to another process
 - I/O APIs already go through the OS, so context switching is easy

• Problems?

Switching on I/O

- Idea: when one process is waiting on I/O, switch to another process
 - I/O APIs already go through the OS, so context switching is easy

- Problems?
 - Some apps don't have any I/O for long periods of time

Preemptive Switching

- So far, processes will not switch to another until an action is taken
 - e.g. an API call or an I/O interrupt
- Idea: use a timer to force context switching at set intervals
 - Timer is running at a fixed frequency to measure how long a process has been running
 - If it's been running for some max duration (scheduling quantum),
 the handler switches to the next process

Preemptive Switching

- So far, processes will not switch to another until an action is taken
 - e.g. an API call or an I/O interrupt
- Idea: use a timer to force context switching at set intervals
 - Timer is running at a fixed frequency to measure how long a process has been running
 - If it's been running for some max duration (scheduling quantum), the handler switches to the next process

Problems? Who will trigger the timer

Preemptive Switching

- So far, processes will not switch to another until an action is taken
 - e.g. an API call or an I/O interrupt

- Idea: use a timer to force context switching at set intervals
 - Timer is running at a fixed frequency to measure how long a process has been running
 - If it's been running for some max duration (scheduling quantum), the handler switches to the next process

- Problems? Who will trigger the timer
 - Requires hardware support (a programmable timer)
 - Thankfully, this is built-in to most modern CPUs

Why I Don't Feel Apps Lagging Even Though They Time-Share the CPU?

Linux as an Example

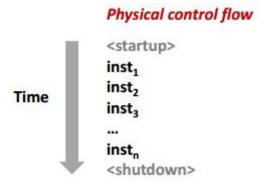
On most desktop Linux systems, the default scheduler "tick" runs **250 Hz or 1000 Hz** — meaning the kernel timer fires every **4 ms or 1 ms**. This timer tick alone lets the scheduler run up to **1000 times per second per CPU**, checking whether to switch processes. On top of this, Linux can also reschedule when events happen (I/O completes, a higher-priority process wakes up), so the system can respond even faster when needed.

Mechanisms for switching:

Exceptional Control Flow

Normal Control Flow or Regular Execution

- Computers only really do one thing, they execute one instruction one after another
 - This is based on the execution in your program.
 - Your programs follow some control flow based on jumps and branches (and calls and returns)
 - This is based on your programs state.



Reasons to break normal control flow

However, sometimes we want to break normal control flow

- Multitasking. Stop one so others can run.
- The program needs to take care of user input. E.g., you hit Ctrl+C on the keyboard in your terminal and execution stops.
- Program can make a system call.

Exceptional Control Flow Mechanisms

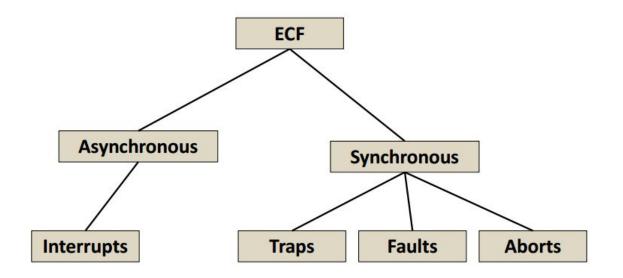
- Low level mechanism
 - Exceptions
 - Change in control flow in response to a system event.
 - This is implemented in hardware and OS software

Exceptional Control Flow Mechanisms

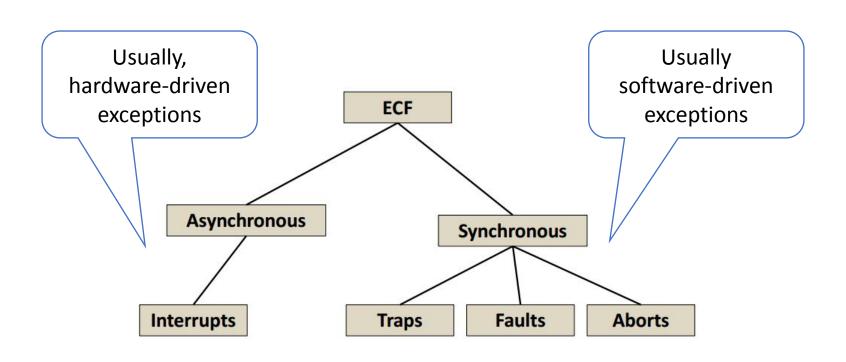
- High level mechanisms
 - Process context switch
 - e.g. It appears that multiple programs are running at once on your OS, but remember only one instruction at a time.
 - Context switches provide this illusion
 - Signals
 - Implemented by OS software

Exceptional Control Flow Taxonomy

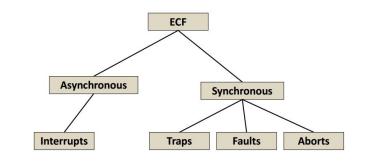
Exception (broad sense): any event that causes the CPU to stop its current normal execution path and transfer control to special handler code.



Exceptional Control Flow Taxonomy



Asynchronous Exceptions (Interrupts)

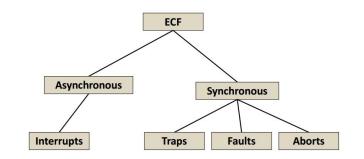


- Caused by events external to processor
 - I.e., not from the result of an instruction the user wrote
 - E.g.
 - Timer interrupts scheduled to happen every few milliseconds
 - A kernel can use this to take back control from a program/user
 - Some network data arrives (I/O)
 - A nice example is while reading from disk
 - The processor can start reading, then hop over and perform some other tasks until memory is actually fetched.

→ Cact	→ CactiLab.github.io git:(master) cat /proc/interrupts										
	CPU0	CPU1	CPU2	CPU3	CPU4	CPU5	CPU6	CPU7			
1:	0	0	Θ	Θ	0	0	0	130731	IR-IO-APIC 1-edge i8042		
8:	0	0	0	0	0	0	0	0	IR-IO-APIC 8-edge rtc0		
9: 12:	0	4232726	0	0	0	0	0	0	IR-IO-APIC 9-fasteoù acpù IR-IO-APIC 12-edge ù8042		
14:	0	9	261613	0	0	0	140	0	IR-IO-APIC 12-edge		
16:	0	0	201013	2043581	0	0	0	0	IR-IO-APIC 14-Tasteot intel_ish_ipc, idma64.0, i801_smbus, i2c_designware.0		
17:	ŏ	ō	1514	114950	16686610	12166	ŏ	ŏ	IR-IO-APIC 17-fasteoi idma64.1, i2c designware.1		
20:	0	Θ	Θ	Θ	0	0	0	Θ	IR-IO-APIC 20-fasteoi idma64.2		
120:	0	0	0	0	0	0	0	0	DMAR-MSI 0-edge dmar0		
121:	0	0	0	0	0	0	0	0	DMAR-MSI 1-edge dmar1		
122:	0	0	0	0	0	0	0	0	IR-PCI-MSI 114688-edge PCIe PME, pciehp		
123: 124:	0	9	9	8	9	0	0	0	IR-PCI-MSI 118784-edge PCIe PME, pciehp IR-PCI-MSI 475136-edge PCIe PME		
125:	0	0	ĕ	0	0	0	0	0	IR-PCI-MSI 489472-edge PCIe PME, pciehp		
126:	Ö	ō	ŏ	Ö	Ö	733	ŏ	ŏ	IR-PCI-MSI 217088-edge thunderbolt		
127:	0	0	Θ	Θ	0	0	1017	0	IR-PCI-MSI 217089-edge thunderbolt		
142:	20318	286966	625	28130	78172	2829		3437	IR-PCI-MSI 212992-edge xhci_hcd		
143:	0	0	0	0	0	0	0	106	IR-PCI-MSI 46137344-edge rtsx_pci		
144:	2233079 1	6412488	2024743	13565716 0	88127	1357648	2781621	385438	IR-PCI-MSI 327680-edge xhci_hcd	Onlinuv	
145: 146:	452940	28 0	0	0	3	0	72 0	0	IR-PCI-MSI 45613056-edge nvme0q0 IR-PCI-MSI 45613057-edge nvme0q1	On Linux	
147:	432340	460347	0	0	0	0	0	0	IR-PCI-MSI 45613058-edge		
148:	Ö	0	435105	Ö	Ö	0	ŏ	Ö	IR-PCI-MSI 45613059-edge nvme0q3		
149:	0	0	0	452592	Ō	0	0	0	IR-PCI-MSI 45613060-edge nvme0q4		
150:	0	0	Θ	Θ	448577	0	0	Θ	IR-PCI-MSI 45613061-edge nvme0q5		
151:	0	0	0	0	0	477954	0	0	IR-PCI-MSI 45613062-edge nvme0q6	cat /proc/interrupts	
152: 153:	0	0	0	0	0	0	462865	452517	IR-PCI-MSI 45613063-edge nvme0q7	cat / proc/interrupts	
153:	0	0	261613	0	9	0	0	452517	IR-PCI-MSI 45613064-edge nvmeθq8 INT3455:00 35 DLL096D:01		
155:	5064963	17310	3206231	220145	4504975	5283354	58368619	499374	IR-PCI-MSI 32768-edge i915		
156:	0	0	0	0	0	733	0	0	IR-PCI-MSI 219136-edge thunderbolt		
157:	0	0	0	0	0	0	1025	0	IR-PCI-MSI 219137-edge thunderbolt	Gran IACONITIC 117-1	
172:	48	0	Θ	Θ	0	0	0	Θ	IR-PCI-MSI 360448-edge mei_me	grep '^CONFIG_HZ='	
173:	322800	1008195	1919846	1051926	1934430	241138	6366	1054598	IR-PCI-MSI 333824-edge iwlwifi:default_queue		
174: 175:	77855	68860 48889	52230 168899	149625	85835 96700	89668	194793 44309	115748	IR-PCI-MSI 333825-edge iwlwifi:queue_1	/boot/config-\$(uname -r)	
176:	55587 109035	48889	42471	79574 586984	93652	157103 115404	39459	82157 80553	IR-PCI-MSI 333826-edge iwlwifi:queue_2 IR-PCI-MSI 333827-edge iwlwifi:queue 3	/boot/coming-a(uname-i)	
177:	50653	64542	36201	78944	81029	89647	47000	86400	IR-PCI-MSI 333828-edge iwlwifi:queue_4		
178:	74801	37065	36703	88397	190921	84557	37344	192138	IR-PCI-MSI 333829-edge iwlwifi:queue 5		
179:	52167	76304	391103	169174	130707	429014	49444	73248	IR-PCI-MSI 333830-edge iwlwifi:queue_6		
180:	78602	57998	97479	112693	69675	80179	86634	321517	IR-PCI-MSI 333831-edge iwlwifi:queue_7	CE CO4 4EE/OEO ~ OCO EOO	
181:	79443	83061	177047	98436	76296	110730	44851	71342	IR-PCI-MSI 333832-edge iwlwifi:queue_8	65,631,155/250 ≈ 262,500	
182: 183:	1	3 0	0	0	1	0	0	47 0	IR-PCI-MSI 333833-edge iwlwifi:exception als-dev0 als_consumer0		
185:	406	0	9	0	0	0	9	0	ats_devo ats_consumero IR-PCI-MSI 514048-edge snd_hda_intel:card0	econds ≈ 3 days of active ticks	
NMI:	0	0	0	9	Ö	0	0	0	Non-maskable interrupts	FLUTIUS ~ 3 days of active ticks	
LOC:	65631155	64596281	63676227	64238073	64059082	64735579	67949340	63804742	Local timer interrupts	1	
SPU:	0	0	Θ	Θ	Θ	0	0	0	Spurious interrupts		
PMI:	0	0	0	0	0	0	0	0	Performance monitoring interrupts		
IWI:	3121548	610104	2096845	704034	2787303	3280945	30305026	917879	IRQ work interrupts		
RTR: RES:	933823	907409	925545	917141	956591	955761	1183767	943247	APIC ICR read retries Rescheduling interrupts		
CAL:	14480409	13290398	12253091	11837881	11574765	11343450	11072357	11242028	Function call interrupts		
TLB:	7137555	7250981	7202742	7093647	7180035	7198626	6894040	7193845	TLB shootdowns		
TRM:	978	978	978	978	978	978	978	978	Thermal event interrupts		
THR:	0	0	0	Θ	Θ	0	O	0	Threshold APIC interrupts		
DFR:	0	0	0	0	0	0	0	0	Deferred Error APIC interrupts		
MCE: MCP:	0 858	859	859	859	859	0 859	0 859	859	Machine check exceptions Machine check polls		
ERR:	858	629	659	659	655	659	659	655	machine check potes		
MIS:	ő										
PIN:	0	0	0	Θ	0	0	0	0	Posted-interrupt notification event		
NPI:	0	0	0	0	0	0	0	0	Nested posted-interrupt event		
PIW:	0	0	_0	Θ	0	0	0	0	Posted-interrupt wakeup event		

Synchronous Exceptions

- Events caused by executing an instruction
 - Traps
 - Intentionally done by the user
 - e.g. system calls, breakpoints (like in gdb)
 - Returns control to the next instruction
 - Faults
 - Unintentional, but possibly recoverable
 - e.g. <u>page faults</u> (we'll learn more about soon), floating point exceptions
 - Handled by re-executing current instruction or aborting execution
 - Aborts
 - Unintentional and unrecoverable
 - e.g. illegal instruction executed, <u>parity error</u>



```
→ CactiLab.github.io git:(master) grep pgfault /proc/vmstat
pgfault 1432038924
→ CactiLab.github.io git:(master) grep pgmajfault /proc/vmstat
pgmajfault 34497
```

What is a page fault?

The CPU tried to access a virtual page not currently mapped.

Two types

- **Minor fault:** Page not mapped yet but data already in RAM; kernel just sets up page tables (fast).
- Major fault: Page must be read from disk or swap (slow).

Why so many?

- Linux uses **demand paging** memory isn't backed until first use.
- Shared libraries and file mappings are faulted in on demand.
- Growing stacks and heap also cause minor faults.

Key metrics

- pgfault total page faults (mostly minor).
- pgmajfault major faults (performance cost).

Software breakpoint — int3 / #BP

The debugger overwrite the first byte of the target instruction with the opcode 0xCC (INT3).

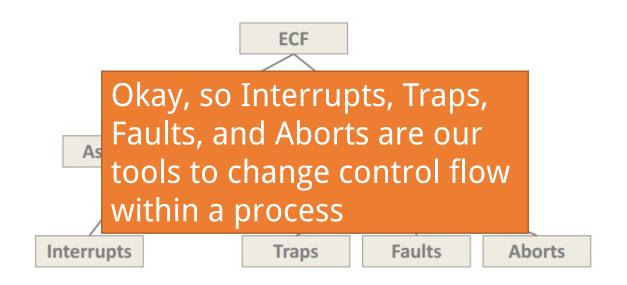
When the CPU fetches and executes that byte, it **generates #BP (Breakpoint Exception, vector 3)**.

The OS delivers this exception back to the debugger (e.g., via ptrace on Linux).

After stopping, the debugger typically restores the original byte so the instruction can run if you continue.

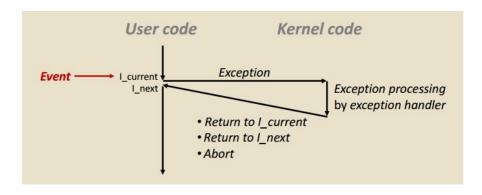
No debug registers involved.

Exceptional Control Flow Taxonomy



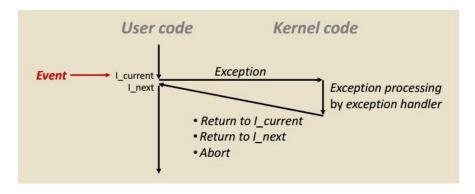
Exceptions

- An exception is a transfer of control to the OS kernel
 - The kernel is the memory-resident part of the OS
 - Meaning OS lives in memory forever: we do not modify this!
- Examples of exceptions we may be familiar with:
 - Divide by 0, arithmetic overflow, or typing Ctrl+C



Exceptions

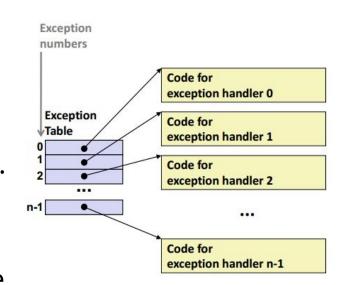
- An exception is a transfer of control to the OS kernel
 - The kernel is the memory-resident part of the OS
 - Meaning OS lives in memory forever: we do not modify this!
- Examples of exceptions we may be familiar with:
 - Divide by 0, arithmetic overflow, or typing Ctrl+C



How does the OS know how to handle the exception?

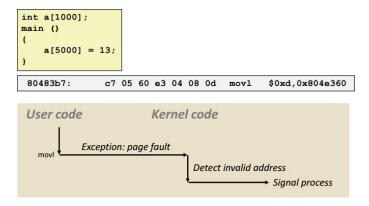
Exception Tables

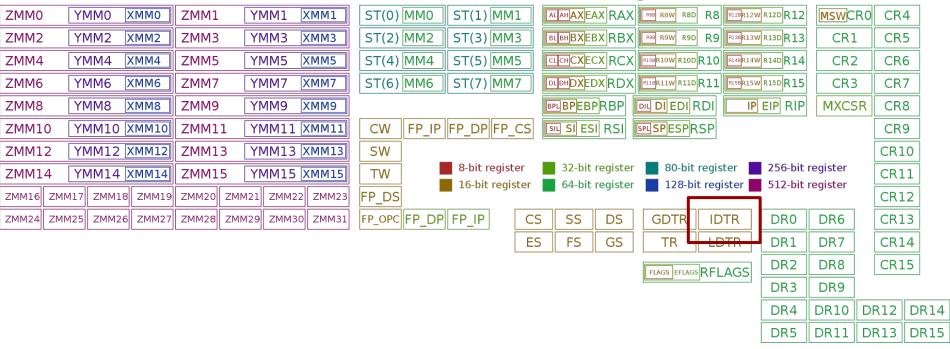
- Somewhere in the OS, a table exists with different exceptions.
 - Think of it like a giant switch or many if else-if statements.
- This is part of a kernel that you cannot modify.
 - This code is in a "protected region" of memory
- For each exception, there is one way to handle it
 - (We call these "exception handlers")



Our favorite: Invalid Memory Reference

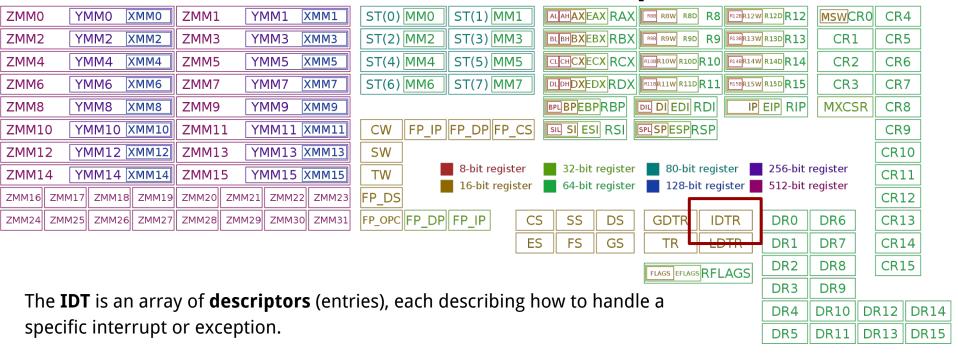
- That is, the segmentation fault
 - OS sends signal SIGSEGV to our user process
 - This time the program gets terminated.





In x86 processors, the Interrupt Descriptor

Table (IDT) is a special data structure the CPU uses
to look up what code to run when an interrupt,
exception, or trap happens.



Each entry tells the CPU:

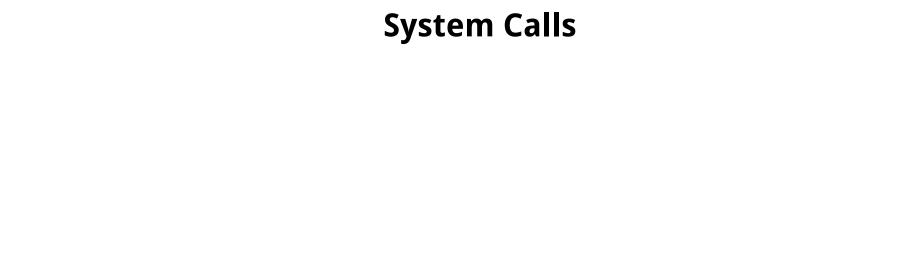
- The **address** of the interrupt/trap handler function (in kernel space).
- The **segment selector** to use (usually the kernel code segment).
- Attributes such as privilege level and gate type.

7.15	EXCEPTION AND INTERRUPT REFERENCE	7-23					
	Interrupt 0—Divide Error Exception (#DE)	7-24					
	Interrupt 1—Debug Exception (#DB)	7-25					
	Interrupt 2—NMI Interrupt	7-27					
	Interrupt 3—Breakpoint Exception (#BP)	7-28					
	Interrupt 4—Overflow Exception (#OF)	7-29					
	Interrupt 5—BOUND Range Exceeded Exception (#BR)	7-30					
	Interrupt 6—Invalid Opcode Exception (#UD)	7-31					
	Interrupt 7—Device Not Available Exception (#NM)						
	Interrupt 8—Double Fault Exception (#DF)	7-33					
	Interrupt 9—Coprocessor Segment Overrun	7-35					
	Interrupt 10—Invalid TSS Exception (#TS)	7-36					
	Interrupt 11—Segment Not Present (#NP)	7-38					
	Interrupt 12—Stack Fault Exception (#SS)						
	Interrupt 13—General Protection Exception (#GP)						
	Interrupt 14—Page-Fault Exception (#PF)	7-44					
	Interrupt 16—x87 FPU Floating-Point Error (#MF)	7-48					
	Interrupt 17—Alignment Check Exception (#AC)	7-50					
	Interrupt 18—Machine-Check Exception (#MC)	7-52					
	Interrupt 19—SIMD Floating-Point Exception (#XM)	7-53					
	Interrupt 20—Virtualization Exception (#VE)	7-55					
	Interrupt 21—Control Protection Exception (#CP)	7-56					
	Interrupts 32 to 255—User Defined Interrupts	7-58					

https://cdrdv2.intel.com/v1/dl/getContent/671447

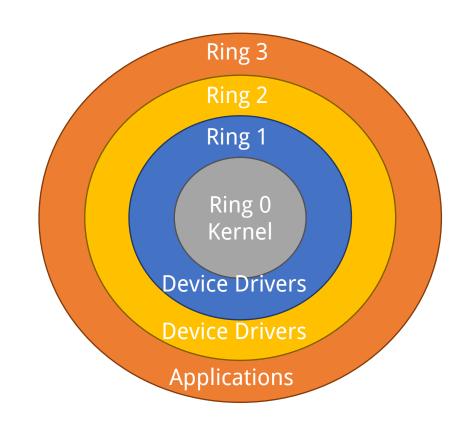
```
// Interrupt descriptor table (shared by all CPUs).
struct gatedesc idt[256];
extern uint vectors[]; // in vectors.S: array of 256 entry pointers
struct spinlock tickslock;
uint ticks;
void
tvinit(void)
  int i;
  for(i = 0; i < 256; i++)
   SETGATE(idt[i], 0, SEG_KCODE<<3, vectors[i], 0);</pre>
  SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3, vectors[T_SYSCALL], DPL_USER);</pre>
  initlock(&tickslock, "time");
void
idtinit(void)
  lidt(idt, sizeof(idt));
```

https://github.com/mit-pdos/xv6-public/blob/master/trap.c



Different privilege levels

- Most modern CPUs support protected mode
- x86 CPUs support three rings with different privileges
 - Ring 0: OS kernel. Most privileged. Full hardware access.
 - Ring 1, 2: device drivers
 - Ring 3: userland. Restricted.
- Most OSes only use rings 0 and 3



Dual-Mode Operation

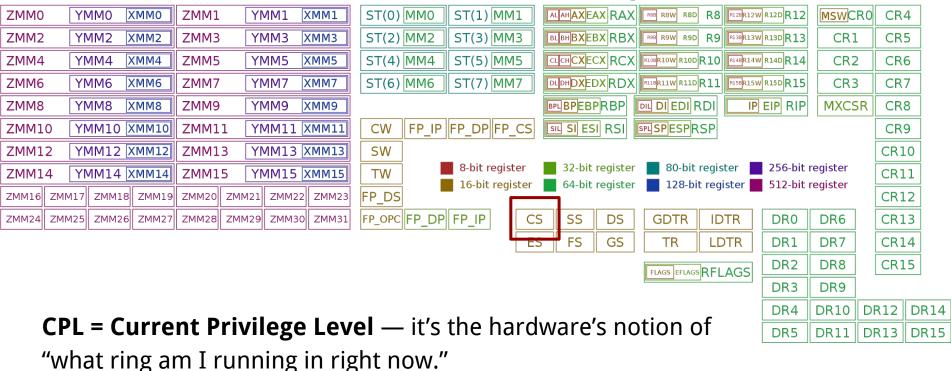
- Ring 0: kernel/supervisor mode
 - Execution with the full privileges of the hardware
 - Read/write to any memory, access any I/O device, read/write any disk sector, send/read any packet

- Ring 3: user mode or "userland"
 - Limited privileges
 - Only those granted by the operating system kernel

Protected Features

- What system features are impacted by protection?
 - Privileged instructions
 - Only available to the kernel
 - Limits on memory accesses
 - Prevents user code from overwriting the kernel
 - Access to hardware
 - Only the kernel may directly interact with peripherals
 - Programmable Timer Interrupt
 - May only be set by the kernel
 - Used to force context switches between processes

x86-64 CPL - Current Privilege Level



CPL is stored in the **lowest two bits of the CS (Code Segment) selector register** on x86/x86-64.

Key Differences: Ring 0 vs Ring 3

Feature	Ring 0 (Kernel)	Ring 3 (User)
Privilege level	Highest (CPL=0)	Lowest (CPL=3)
Access to hardware	Direct (I/O ports, control registers)	No direct hardware access
Memory access	Can map & modify page tables	Limited to user space virtual memory
Instructions allowed	All x86 privileged instructions	Non-privileged subset only
System calls	Not needed (already privileged)	Must trap to kernel via syscalls / int ?? /syscall

What are System Calls?

When a process needs to invoke a kernel service, it invokes a procedure call in the operating system interface using special instructions (not a **call** instruction in x86; but **int** or **syscall**). Such a procedure is called a system call.

The system call enters the kernel; the kernel performs the service and returns. Thus a process alternates between executing in user space and kernel space.

System calls are generally not invoked directly by a program, but rather via wrapper functions in glibc (or perhaps some other library).

System Calls

- Syscall is the lowest level of interaction with an operating system from a C programmer
- A user program can ask the OS for services that the OS manages
 - You may have used '_exit' in your assignment
 - Anything else you can think of?

System Calls

- Syscall is the lowest level of interaction with an operating system from a C programmer
- A user program can ask the OS for services that the OS manages
 - You may have used '_exit' in your assignment
 - Anything else you can think of?

Number	Name	Description
0	read	Read file
1	write	Write file
2	open	Open file
3	close	Close file
4	stat	Get info about file
57	fork	Create process
59	execve	Execute a program
60	_exit	Terminate process
62	kill	Send signal to process

Popular System Call

On Unix, Unix-like and other POSIX-compliant operating systems, popular system calls are open, read, write, close, wait, exec, fork, exit, and kill.

Many modern operating systems have hundreds of system calls. For example, Linux and OpenBSD each have over 300 different calls, FreeBSD has over 500, Windows 7 has close to 700.

Glibc interfaces

Often, but not always, the name of the wrapper function is the same as the name of the system call that it invokes.

For example, glibc contains a function chdir() which invokes the underlying "chdir" system call.

System Calls change the CPU Mode (CPL)

- Applications often need to access the OS
 - i.e. system calls
 - Writing files, displaying on the screen, receiving data from the network, etc...
- But the OS is ring 0, and apps are ring 3
- How do apps get access to the OS?
 - Apps invoke system calls with an interrupt
 - E.g. int 0x80
 - int causes a mode transfer from ring 3 to ring 0

Tools: strace

```
paper-auto-glitching git:(main) strace ls
execve("/usr/bin/ls", ["ls"], 0x7fff83ec4e70 /* 54 vars */) = 0
brk(NULL)
                                    = 0x55b54b6a0006
arch prctl(0x3001 /* ARCH ??? */, 0x7ffe8074a7d0) = -1 EINVAL (Invalid argument)
access("/etc/ld.so.preload", R_OK)
                                   = -1 ENOENT (No such file or directory)
openat(AT FDCWD, "/etc/ld.so.cache", O RDONLY|O CLOEXEC) = 3
fstat(3, {st mode=S IFREG|0644, st size=179304, ...}) = 0
mmap(NULL, 179304, PROT READ, MAP PRIVATE, 3, 0) = 0x7f7c7fba9000
close(3)
openat(AT_FDCWD, "/lib/x86_64-linux-gnu/libselinux.so.1", O_RDONLY|O_CLOEXEC) = 3
fstat(3, {st mode=S IFREG|0644, st size=163200, ...}) = 0
mmap(NULL, 8192, PROT READ|PROT WRITE, MAP PRIVATE|MAP ANONYMOUS, -1, 0) = 0x7f7c7fba7000
mmap(NULL, 174600, PROT_READ, MAP_PRIVATE|MAP_DENYWRITE, 3, 0) = 0x7f7c7fb7c000
mprotect(0x7f7c7fb82000, 135168, PROT NONE) = 0
mmap(0x7f7c7fb82000, 102400, PROT READ|PROT EXEC, MAP PRIVATE|MAP FIXED|MAP DENYWRITE, 3, 0x6000) = 0x7f7c7fb82000
mmap(0x7f7c7fb9b000, 28672, PROT READ, MAP PRIVATE|MAP FIXED|MAP DENYWRITE, 3, 0x1f000) = 0x7f7c7fb9b000
mmap(0x7f7c7fba3000, 8192, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3, 0x26000) = 0x7f7c7fba3000
mmap(0x7f7c7fba5000, 6664, PROT_READ|PROT_WRITE, MAP_PRIVATE|MAP_FIXED|MAP_ANONYMOUS, -1, 0) = 0x7f7c7fba5000
openat(AT FDCWD, "/lib/x86 64-linux-gnu/libc.so.6", 0 RDONLY|O CLOEXEC) = 3
read(3, "\177ELF\2\1\1\3\0\0\0\0\0\0\0\1\0\0\300A\2\0\0\0\0\0\"..., 832) = 832
pread64(3, "\4\0\0\0\20\0\0\0\5\0\0\0NU\0\2\0\0\300\4\0\0\0\3\0\0\0\0\0\0", 32, 848) = 32
pread64(3, "\4\0\0\0\24\0\0\0\3\0\0\0W\222s/x1X\306o\264\363udX\312$"..., 68, 880) = 68
fstat(3, {st mode=S IFREG|0755, st size=2029592, ...}) = 0
pread64(3, "\4\0\0\0\24\0\0\0\3\0\0\0W\222s/x1X\306o\264\363udX\312$"..., 68, 880) = 68
mmap(NULL, 2037344, PROT READ, MAP PRIVATE|MAP DENYWRITE, 3, 0) = 0x7f7c7f98a000
mmap(0x7f7c7f9ac000, 1540096, PROT_READ|PROT_EXEC, MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3, 0x22000) = 0x7f7c7f9ac000
mmap(0x7f7c7fb24000, 319488, PROT_READ, MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3, 0x19a000) = 0x7f7c7fb24000
mmap(0x7f7c7fb72000, 24576, PROT READ|PROT WRITE, MAP PRIVATE|MAP FIXED|MAP DENYWRITE, 3, 0x1e7000) = 0x7f7c7fb72000
mmap(0x7f7c7fb78000, 13920, PROT READ|PROT WRITE, MAP PRIVATE|MAP FIXED|MAP ANONYMOUS, -1, 0) = 0x7f7c7fb78000
openat(AT FDCWD, "/lib/x86 64-linux-qnu/libpcre2-8.so.0", O RDONLY|O CLOEXEC) = 3
read(3, "\177ELF\2\1\1\0\0\0\0\0\0\0\0\0\1\0\0\0\340\"\0\0\0\0\0\0\0\0"..., 832) = 832
fstat(3, {st mode=S IFREG|0644, st size=588488, ...}) = 0
mmap(NULL, 590632, PROT_READ, MAP_PRIVATE|MAP_DENYWRITE, 3, 0) = 0x7f7c7f8f9000
mmap(0x7f7c7f8fb000, 413696, PROT_READ|PROT_EXEC, MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3, 0x2000) = 0x7f7c7f8fb000
mmap(0x7f7c7f960000, 163840, PROT READ, MAP PRIVATE|MAP FIXED|MAP DENYWRITE, 3, 0x67000) = 0x7f7c7f960000
mmap(0x7f7c7f988000, 8192, PROT READ|PROT WRITE, MAP PRIVATE|MAP FIXED|MAP DENYWRITE, 3, 0x8e000) = 0x7f7c7f988000
close(3)
openat(AT FDCWD, "/lib/x86 64-linux-gnu/libdl.so.2", O RDONLY|O CLOEXEC) = 3
read(3, "\177ELF\2\1\1\0\0\0\0\0\0\0\0\0\1\0\0\0 \22\0\0\0\0\0\0"..., 832) = 832
fstat(3, {st mode=S IFREG|0644, st size=18848, ...}) = 0
mmap(NULL, 20752, PROT READ, MAP PRIVATE|MAP DENYWRITE, 3, 0) = 0x7f7c7f8f3000
mmap(0x7f7c7f8f4000, 8192, PROT_READ|PROT_EXEC, MAP_PRIVATE|MAP_FIXED|MAP_DENYWRITE, 3, 0x1000) = 0x7f7c7f8f4000
mmap(0x7f7c7f8f6000, 4096, PROT READ, MAP PRIVATE|MAP FIXED|MAP DENYWRITE, 3, 0x3000) = 0x7f7c7f8f6000
mmap(0x7f7c7f8f7000, 8192, PROT READ|PROT WRITE, MAP PRIVATE|MAP FIXED|MAP DENYWRITE, 3, 0x3000) = 0x7f7c7f8f7000
openat(AT_FDCWD, "/lib/x86_64-linux-gnu/libpthread.so.0", O_RDONLY|O_CLOEXEC) = 3
read(3, "\177ELF\2\1\1\0\0\0\0\0\0\0\0\1\0\0\220q\0\0\0\0\0\0"..., 832) = 832
pread64(3, "\4\0\0\0\24\0\0\0\3\0\0GNU\0\227Sr\5\2W;\227\333\254Y[a\375r\302"..., 68, 824) = 68
fstat(3, {st mode=S IFREG|0755, st size=157224, ...}) = 0
```

strace Is

Use "man 2 syscall_name" to check out its usage

On x86/x86-64, most system calls rely on the software interrupt.

A software interrupt is caused either by an exceptional condition in the processor itself, or a special instruction (the **int 0x80** instruction or **syscall** instruction).

For example: a divide-by-zero exception will be thrown if the processor's arithmetic logic unit is commanded to divide a number by zero as this instruction is in error and impossible.

Making a System Call in x86 Assembly (INT 0x80)

x86 (32-bit)

Compiled from Linux 4.14.0 headers.

NR	syscall name	references	%eax	arg0 (%ebx)	arg1 (%ecx)	arg2 (%edx)	arg3 (%esi)	arg4 (%edi)	arg5 (%ebp)
0	restart_syscall	man/ cs/	0x00	-	-	100	-	9 - 0	-
1	exit	man/ cs/	0x01	int error_code	-	s-	4 - 5	j	-
2	fork	man/ cs/	0x02	÷	5	-	-	-	-
3	read	man/ cs/	0x03	unsigned int fd	char *buf	size_t count	100	CE1	2
4	write	man/ cs/	0x04	unsigned int fd	const char *buf	size_t count	-	-	-
5	open	man/ cs/	0x05	const char *filename	int flags	umode_t mode	-	-	-
6	close	man/ cs/	0x06	unsigned int fd	-	-	-	-	-
7	waitpid	man/ cs/	0x07	pid_t pid	int *stat_addr	int options		(*)	-
8	creat	man/ cs/	0x08	const char *pathname	umode_t mode		-	-	-
9	link	man/ cs/	0x09	const char *oldname	const char *newname	-	1.50	100	-
10	unlink	man/ cs/	0x0a	const char *pathname	=		S.T.S.	l=e	es .
11	execve	man/ cs/	0x0b	const char *filename	const char *const *argv	const char *const *envp	-	-	-
12	chdir	man/ cs/	0x0c	const char *filename	±1		621	526	=
13	time	man/ cs/	0x0d	time_t *tloc	=	N2	14		2
14	mknod	man/ cs/	0x0e	const char *filename	umode_t mode	unsigned dev	-	~	-
15	chmod	man/ cs/	0x0f	const char *filename	umode_t mode		-	5=0	2
10	lahaum	man! cal	0.10	const show	a +a.v	and t areas			

https://chromium.googlesource.com/chromiumos/docs/+/master/constants/syscalls.md#x86-32_bit

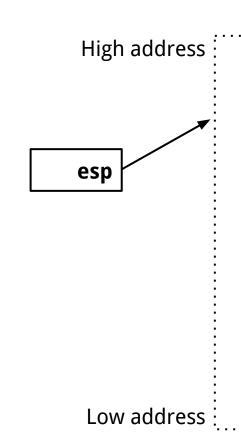
xor eax,eax push eax 0x68732f2f push 0x6e69622f push ebx,esp mov push eax push ebx mov ecx,esp mov al,0xb 08x0 int

```
Dec Hx Oct Char
                                      Dec Hx Oct Html Chr
                                                          Dec Hx Oct Html Chr Dec Hx Oct Html Chr
                                                            64 40 100 @ 0
 0 0 000 NUL (null)
                                      32 20 040   Space
                                                                               96 60 140 6#96;
                                      33 21 041 6#33; !
                                                            65 41 101 A A
    1 001 SOH (start of heading)
                                                                               97 61 141 6#97; @
                                      34 22 042 6#34; "
                                                            66 42 102 B B
                                                                               98 62 142 6#98; b
    2 002 STX (start of text)
    3 003 ETX (end of text)
                                      35 23 043 6#35; #
                                                            67 43 103 C C
                                                                               99 63 143 4#99; 0
    4 004 EOT (end of transmission)
                                      36 24 044 6#36; $
                                                            68 44 104 D D
                                                                              100 64 144 d d
                                                                              101 65 145 6#101; 6
 5 5 005 ENQ (enquiry)
                                      37 25 045 6#37; %
                                                            69 45 105 E E
                                      38 26 046 4#38; 4
                                                            70 46 106 @#70; F
                                                                              102 66 146 @#102; f
    6 006 ACK (acknowledge)
    7 007 BEL (bell)
                                      39 27 047 6#39; 1
                                                            71 47 107 @#71; G
                                                                              103 67 147 @#103; g
                                                            72 48 110 @#72; H
                                                                              104 68 150 6#104; h
    8 010 BS
              (backspace)
                                      40 28 050 6#40;
                                                            73 49 111 6#73; I
                                                                             105 69 151 6#105; 1
    9 011 TAB (horizontal tab)
                                      41 29 051 6#41; )
   A 012 LF
              (NL line feed, new line)
                                      42 2A 052 6#42; *
                                                            74 4A 112 6#74; J
                                                                              106 6A 152 @#106; j
                                      43 2B 053 6#43; +
                                                            75 4B 113 6#75; K
                                                                              107 6B 153 k k
11 B 013 VT
              (vertical tab)
   C 014 FF
              (NP form feed, new page)
                                      44 2C 054 ,
                                                            76 4C 114 6#76; L
                                                                             108 6C 154 6#108; 1
13 D 015 CR
                                      45 2D 055 6#45; -
                                                            77 4D 115 @#77; M
                                                                              109 6D 155 m m
              (carriage return)
14 E 016 SO
              (shift out)
                                      46 2E 056 . .
                                                            78 4E 116 @#78; N
                                                                              110 6E 156 n n
15 F 017 SI
             (shift in)
                                      47 2F 057 @#47; /
                                                                             111 6F 157 @#111; 0
                                                            79 4F 117 6#79: 0
                                      48 30 060 6#48; 0
                                                            80 50 120 6#80; P
                                                                             112 70 160 @#112; p
16 10 020 DLE (data link escape)
                                      49 31 061 6#49; 1
                                                            81 51 121 6#81; 0
17 11 021 DC1 (device control 1)
                                                                             113 71 161 4#113; 9
                                      50 32 062 6#50; 2
                                                            82 52 122 @#82; R
                                                                             114 72 162 @#114; r
18 12 022 DC2 (device control 2)
19 13 023 DC3 (device control 3)
                                      51 33 063 4#51; 3
                                                            83 53 123 6#83; $
                                                                             115 73 163 4#115; 8
                                                            84 54 124 6#84; T
                                                                             116 74 164 @#116; t
20 14 024 DC4 (device control 4)
                                      52 34 064 6#52; 4
                                      53 35 065 4#53; 5
                                                            85 55 125 6#85; U
                                                                             117 75 165 6#117; u
21 15 025 NAK (negative acknowledge)
                                                                              118 76 166 4#118; 7
22 16 026 SYN (synchronous idle)
                                      54 36 066 6#54; 6
                                                            86 56 126 V V
23 17 027 ETB (end of trans. block)
                                      55 37 067 6#55; 7
                                                            87 57 127 6#87; W
                                                                             119 77 167 w W
24 18 030 CAN (cancel)
                                      56 38 070 6#56; 8
                                                            88 58 130 6#88; X
                                                                             120 78 170 x X
                                                                              121 79 171 @#121; 7
25 19 031 EM
             (end of medium)
                                      57 39 071 6#57; 9
                                                            89 59 131 Y Y
26 1A 032 SUB (substitute)
                                       58 3A 072 : :
                                                            90 5A 132 Z Z
                                                                              122 7A 172 z Z
27 1B 033 ESC (escape)
                                      59 3B 073 4#59; ;
                                                            91 5B 133 6#91; [
                                                                             123 7B 173 6#123;
28 1C 034 FS
              (file separator)
                                      60 3C 074 < <
                                                            92 5C 134 @#92; \
                                                                              124 7C 174 @#124;
                                                            93 5D 135 6#93; ]
                                                                              125 7D 175 } }
29 1D 035 GS
              (group separator)
                                      61 3D 075 = =
                                                                              126 7E 176 ~ ~
                                      62 3E 076 > >
                                                            94 5E 136 @#94; ^
30 1E 036 RS
              (record separator)
31 1F 037 US
              (unit separator)
                                      63 3F 077 ? ?
                                                            95 5F 137 @#95;
                                                                           127 7F 177  DEL
```

Source: www.LookupTables.com

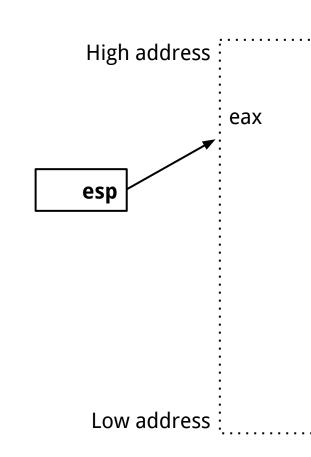
stack

eax,eax xor push eax 0x68732f2f push 0x6e69622f push ebx,esp mov push eax push ebx mov ecx,esp al,0xb mov int 0x80



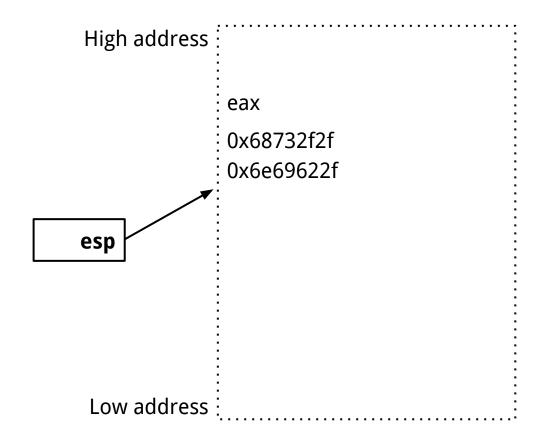
stack

eax,eax xor push eax push 0x68732f2f 0x6e69622f push ebx,esp mov push eax push ebx mov ecx,esp al,0xb mov int 0x80



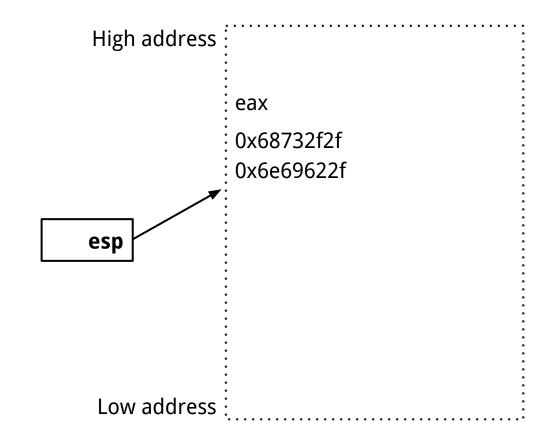
stack

eax,eax xor push eax 0x68732f2f push 0x6e69622f push ebx,esp mov push eax ebx push mov ecx,esp al,0xb mov int 0x80



stack

eax,eax xor push eax 0x68732f2f push 0x6e69622f push mov ebx,esp push eax push ebx mov ecx,esp al,0xb mov int 0x80



```
EXECVE(2)
                                   Linux Programmer's Manual
NAME
       execve - execute program
SYNOPSIS
       #include <unistd.h>
       int execve(const char *filename, char *const argv[],
                   char *const envp[]);
       /bin/sh, 0x0
                              0x00000000
                                              Address of /bin/sh, 0x00000000
           EBX
                                  EDX
                                                         ECX
```

execve("/bin/sh", address of string "/bin/sh", 0)

Making a System Call in x86_64 (64-bit) Assembly

x86_64 (64-bit)

Compiled from Linux 4.14.0 headers.

NR	syscall name	references	%rax	arg0 (%rdi)	arg1 (%rsi)	arg2 (%rdx)	arg3 (%r10)	arg4 (%r8)	arg5 (%r9)
0	read	man/ cs/	0x00	unsigned int fd	char *buf	size_t count	-		
1	write	man/ cs/	0x01	unsigned int fd	const char *buf	size_t count	-	i.e.	-
2	open	man/ cs/	0x02	const char *filename	int flags	umode_t mode	5	i.i.	ā
3	close	man/ cs/	0x03	unsigned int fd		₹.	-	8.50	-
4	stat	man/ cs/	0x04	const char *filename	struct old_kernel_stat *statbuf	B	8	(8)	-
5	fstat	man/ cs/	0x05	unsigned int fd	struct old_kernel_stat *statbuf	-	-	1.00	-
6	Istat	man/ cs/	0x06	const char *filename	struct old_kernel_stat *statbuf	E	8	(8)	*
7	poll	man/ cs/	0x07	struct pollfd *ufds	unsigned int nfds	int timeout	-		-
8	lseek	man/ cs/	0x08	unsigned int fd	off_t offset	unsigned int whence	-	1.5	
9	mmap	man/ cs/	0x09	?	?	?	?	?	?
10	mprotect	man/ cs/	0x0a	unsigned long start	size_t len	unsigned long prot	-	15	
11	munmap	man/ cs/	0x0b	unsigned long addr	size_t len		-		
12	brk	man/ cs/	0x0c	unsigned long brk		-	-	(e)	-
13	rt_sigaction	man/ cs/	0x0d	int	const struct sigaction *	struct sigaction *	size_t		-

https://chromium.googlesource.com/chromiumos/docs/+/master/constants/syscalls.md#x86-32_bit

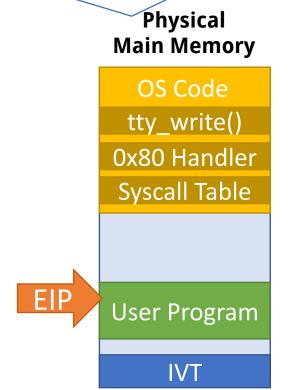
Making a System Call in x86_64 (64-bit) Assembly

NR	syscall name	references	%rax	arg0 (%rdi)	arg1 (%rsi)	arg2 (%rdx)	arg3 (%r10)	arg4 (%r8)	arg5 (%r9)
59	execve	man/ cs/	0x3b	const char *filename	const char *const *argv	const char *const *envp	-		18

push rax xor rdx, rdx xor rsi, rsi mov rbx, '/bin//sh' push rbx push rsp pop rdi mov al, 59 syscall

- 1. Software executes int 0x80
 - Pushes EIP, CS, and EFLAGS

Note: this shows a physical memory layout. The user program thinks it owns the entire memory space (the diagram that we saw in previous lectures).



- 1. Software executes int 0x80
 - Pushes EIP, CS, and EFLAGS
- CPU transfers execution to the OS handler
 - Look up the handler in the Interrupt Vector Table (IVT)
 - Switch from ring 3 to 0

Note: this shows a physical memory layout. The user program thinks it owns the entire memory space (the diagram that we saw in previous lectures).

Physical Main Memory

OS Code

tty_write()
0x80 Handler

Syscall Table

EIP

User Program

- 1. Software executes int 0x80
 - Pushes EIP, CS, and EFLAGS
- CPU transfers execution to the OS handler
 - Look up the handler in the Interrupt Vector Table (IVT)
 - Switch from ring 3 to 0

Note: this shows a physical memory layout. The user program thinks it owns the entire memory space (the diagram that we saw in previous lectures).

Physical Main Memory

OS Code tty_write()

0x80 Handler

Syscall Table

User Program

Note: this shows a physical memory layout. The user program thinks it owns the entire memory space (the diagram that we saw in previous lectures).

- 1. Software executes int 0x80
 - Pushes EIP, CS, and EFLAGS
- 2. CPU transfers execution to the OS handler
 - Look up the handler in the Interrupt Vector Table (IVT)
 - Switch from ring 3 to 0
- 3. OS executes the system call
 - Save the processes state
 - Use EAX to locate the system call
 - Execute the system call
 - Restore the processes state
 - Put the return value in EAX

Physical
Main Memory

OS Code tty_write()

0x80 Handler

Syscall Table

User Program

1. Software executes int 0x80

Pushes EIP, CS, and EFLAGS

CPU transfers execution to the OS handler

- Look up the handler in the Interrupt Vector Table (IVT)
- Switch from ring 3 to 0
- 3. OS executes the system call
 - Save the processes state
 - Use EAX to locate the system call
 - Execute the system call
 - Restore the processes state
 - Put the return value in EAX

Note: this shows a physical memory layout. The user program thinks it owns the entire memory space (the diagram that we saw in previous lectures).

Physical Main Memory

OS Code tty_write()

0x80 Handler

Syscall Table

User Program

Note: this shows a physical memory layout. The user program thinks it owns the entire memory space (the diagram that we saw in previous lectures).

- 1. Software executes int 0x80
 - Pushes EIP, CS, and EFLAGS
- 2. CPU transfers execution to the OS handler
 - Look up the handler in the Interrupt Vector Table (IVT)
 - Switch from ring 3 to 0
- 3. OS executes the system call
 - Save the processes state
 - Use EAX to locate the system call
 - Execute the system call
 - Restore the processes state
 - Put the return value in EAX

Physical
Main Memory

OS Code tty_write()

0x80 Handler

Syscall Table

User Program

Note: this shows a physical memory layout. The user program thinks it owns the entire memory space (the diagram that we saw in previous lectures).

- 1. Software executes int 0x80
 - Pushes EIP, CS, and EFLAGS
- 2. CPU transfers execution to the OS handler
 - Look up the handler in the Interrupt Vector Table (IVT)
 - Switch from ring 3 to 0
- 3. OS executes the system call
 - Save the processes state
 - Use EAX to locate the system call
 - Execute the system call
 - Restore the processes state
 - Put the return value in EAX
- 4. Return to the process with iret
 - Pops EIP, CS, and EFLAGS
 - Switches from ring 0 to 3

Physical Main Memory

OS Code tty_write()

0x80 Handler

Svscall Table

EIP

User Program

System Calls and arguments

- Helpful webpage with syscalls and arguments
 - https://filippo.io/linux-syscall-table/

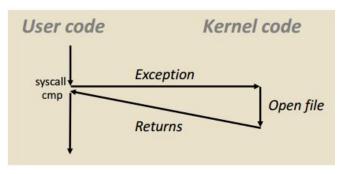
	F	-7r	No. 1 State
8	lseek	sys_lseek	fs/read_write.c
9	mmap	sys_mmap	arch/x86/kernel/sys_x86_64.c
10	mprotect	sys_mprotect	mm/mprotect.c
11	munmap	sys_munmap	mm/mmap.c
12	brk	sys_brk	mm/mmap.c

Opening a File

- rax holds the system call # that we want to pass.
 - Other arguments accessed as follows

%rax	Name	Entry point		Implementation
0	read	sys_read		fs/read_write.c
1	write	sys_write		fs/read_write.c
2	open	sys_open		fs/open.c
%rdi		%rsi	%rdx	
const char user * filename			int flags	umode t mode

Opening a File | Illustration



x86 and xv6 as an example

```
#include "traps.h"
 8
       #include "spinlock.h"
 9
10
       // Interrupt descriptor table (shared by all CPUs).
11
12
       struct gatedesc idt[256];
       extern uint vectors[]; // in vectors.S: array of 256 entry pointers
13
       struct spinlock tickslock;
14
       uint ticks;
15
16
       void
17
       tvinit(void)
18
19
         int i;
20
21
         for(i = 0; i < 256; i++)
22
           SETGATE(idt[i], 0, SEG_KCODE<<3, vectors[i], 0);</pre>
23
24
         SETGATE(idt[T_SYSCALL], 1, SEG_KCODE<<3, vectors[T_SYSCALL], DPL_USER);</pre>
25
         initlock(&tickslock, "time");
26
27
28
```

https://github.com/mit-pdos/xv6-public/blob/master/trap.c https://github.com/mit-pdos/xv6-public/blob/master/syscall.c

Piping and Redirection

Channels of Communication for Linux Process

Every process in Linux has three initial, standard channels of communication:

- Standard Input (stdin, fd=0) is the channel through which the process takes input. For example, your shell uses Standard Input to read the commands that you input.
- Standard Output (stdout, fd=1) is the channel through which processes output normal data, such as the flag when it is printed to you in previous challenges or the output of utilities such as *ls*.
- Standard Error (stderr, fd=2) is the channel through which processes output error details. For example, if you mistype a command, the shell will output, over standard error, that this command does not exist.

Examples

```
Redirecting output > or 1>
echo hi > asdf echo hi 1> asdf
ls > files.txt
```

Appending output >> echo hi >> asdf

Redirecting errors 2> /challenge/run 2> errors.log

```
Redirecting input <
rev < messagefile
sort < names.txt</pre>
```

Pipe

The | (pipe) operator. Standard output from the command to the left of the pipe will be connected to (piped into) the standard input of the command to the right of the pipe.

echo hello-world | wc -c

How to use the fork syscall?

What it does

- Creates a new process by duplicating the calling process.
- The new process is called the child; the original is the parent.

Key details

- Both processes continue execution from the point after the fork() call.
- Return values:
 - Parent gets the child's PID (positive number).
 - Child gets 0.
 - If fork fails, parent gets -1.
- Child initially gets a copy of:
 - Parent's address space (code, data, stack).
 - File descriptors.
 - Environment and signal dispositions.

How to use the fork syscall?

```
#include <stdio.h>
#include <unistd.h>
int main() {
  pid_t pid = fork();
  if (pid == 0) {
    // Child process
    printf("I am the child! PID = %d\n", getpid());
  } else if (pid > 0) {
    // Parent process
    printf("I am the parent! Child PID = %d\n", pid);
  } else {
    // Error
    perror("fork failed");
  return 0;
```

strace -f ./a.out

How many processes?

```
#include <stdio.h>
#include <unistd.h>
int main() {
  fork(); // First fork
  fork(); // Second fork
  fork(); // Third fork
  printf("Hello from PID
%d\n", getpid());
  return 0;
```

How to use the dup and dup2 syscall?

Both: create a duplicate of an existing file descriptor (FD)

An FD is a handle the OS gives to open files, pipes, sockets, etc.

dup(oldfd)

- Creates a new FD that refers to the same file/pipe/socket as oldfd.
- Returns **lowest-numbered free FD** (you don't choose the number).
- Leaves oldfd unchanged.

int fd2 = dup(fd1); // fd2 is some free number

dup2(oldfd, newfd)

- Makes newfd refer to the same thing as oldfd.
- If newfd is already open, it's closed first.
- Guarantees the new FD number is exactly newfd.

dup2(fd1, STDOUT_FILENO); // redirect stdout to fd1

How to use the wait syscall?

Purpose

- Allows a parent process to wait for one of its child processes to finish.
- Collects the child's exit status to avoid creating a zombie process.
 pid_t wait(int *status);

Behavior:

- Suspends the calling (parent) process until one of its children exits.
- Returns the PID of the terminated child.
- If status is not NULL, stores the child's exit code.

Returns **-1** if no children or on error.

How to use the wait syscall?

```
#include <sys/wait.h>
#include <unistd.h>
#include <stdio.h>
int main() {
  pid_t pid = fork();
  if (pid == 0) {
    printf("Child: exiting\n");
    _exit(42);
  } else {
    int status;
     pid_t child = wait(&status);
     printf("Parent: child %d exited with code %d\n",
         child, WEXITSTATUS(status));
  return 0;
```

What does this program do?

```
pid_t pid = fork();
  if (pid == 0) {
    // Child process
    int fd = open("out.txt", O_CREAT | O_WRONLY | O_TRUNC, 0644);
    if (fd < 0) { perror("open"); exit(1); }
    // Redirect stdout to the file
    dup2(fd, STDOUT FILENO);
    close(fd);
    // Run command
    execlp("ls", "ls", NULL);
    perror("execlp"); // only runs if execlp fails
    exit(1);
  } else if (pid > 0) {
    wait(NULL);
```

What does this program do?

```
pid = fork();
  if (pid == 0) {
    int fd = open("in.txt", O_RDONLY);
    if (fd < 0) { perror("open"); exit(1); }
    // Redirect stdin from the file
    dup2(fd, STDIN_FILENO);
    close(fd);
    execlp("sort", "sort", NULL);
    perror("execlp");
    exit(1);
  } else if (pid > 0) {
    wait(NULL);
  return 0;
```

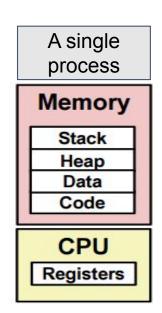
Announcements

- 1. The class scheduled for Tuesday will be delivered online. I will post the exact meeting time on Canvas later. Attendance is optional, and students are welcome to join if they wish. The recording of the session will be uploaded to YouTube afterward for anyone who prefers to watch it later.
- 2. Midterm exam on Thursday. If you've done the readings and the assignments for the first 4 weeks, then you should be fine. Open book/notes, no electronic devices, no substantial coding questions.



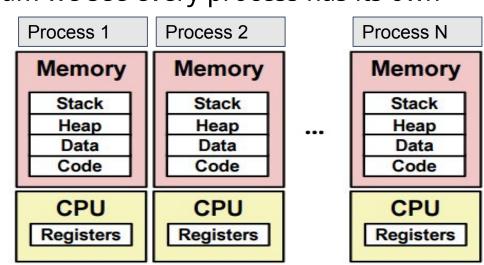
The Process

- A process is alive, a program is dead. A program is just the code.
- Processes are organized by the OS using two key abstractions
 - Logical Control Flow
 - Programs "appear" to have exclusive control over the CPU
 - Done by "context switching"
 - Private Address Space
 - Each program "appears" to have exclusive use of main memory
 - Provided by mechanism called virtual memory



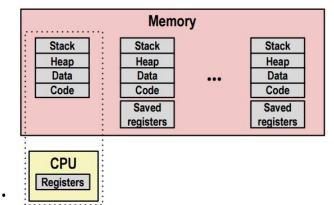
Multiprocessing: Illusion

- When running processes, it appears that we are running many different tasks at the same time
- It also appears that our memory is neatly organized.
 - Note from this diagram we see every process has its own
 - stack
 - heap
 - data
 - code
 - registers

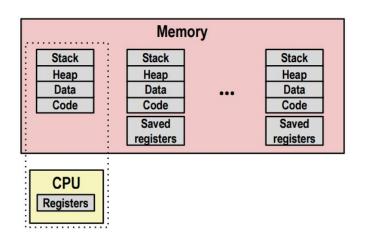


Multiprocessing: Reality

- Remember, at any time, only one processor is really running code
- Program execution is interleaved
- OS manages memory addresses in virtual memory
- OS stores the saved registers for different programs.
 - (At some point in this class, you probably figured 16 registers is not enough for all of the processes that you were running.)
- When we switch which process is executing: this is a context switch



Context switch: a high-level view

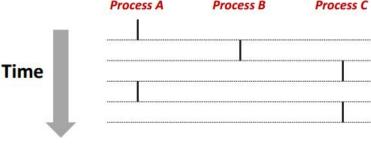


Memory Stack Stack Stack Heap Heap Heap Data Data Data Code Code Code Saved Saved registers registers CPU Registers Context Switch

- Save register values to memory
- Move on to the next process
 - Point to the stack of the next process
 - Restore saved register values
- Start running executing the next process

- Each process running has its own control flow
- If they overlap in their lifetime, then they are running concurrently
 - otherwise they are sequential
- Remember only 1 process at a time can execute
 - On a single core, which processes here are concurrent to each other?

 Process A Process B Process C
 - Concurrent:
 - Which are sequential?
 - Sequential:



- Each process running has its own control flow
- If they overlap in their lifetime, then they are running concurrently
 - otherwise they are sequential
- Remember only 1 process at a time can execute
 - On a single core, which processes here are concurrent to each other?

Time

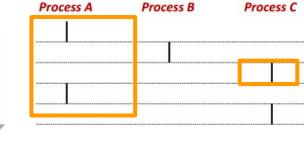
- Concurrent: A&B
- Which are sequential?
 - Sequential:



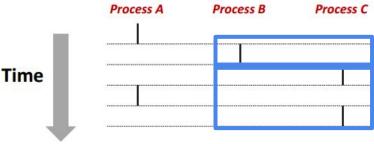
- Each process running has its own control flow
- If they overlap in their lifetime, then they are running concurrently
 - otherwise they are sequential
- Remember only 1 process at a time can execute
 - On a single core, which processes here are concurrent to each other?

Time

- Concurrent: A&B, A&C
- Which are sequential?
 - Sequential:

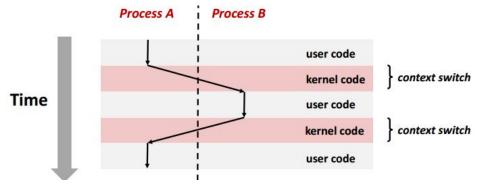


- Each process running has its own control flow
- If they overlap in their lifetime, then they are running concurrently
 - otherwise they are sequential
- Remember only 1 process at a time can execute
 - On a single core, which processes here are concurrent to each other?
 - Concurrent: A&B, A&C
 - Which are sequential?
 - Sequential: B &C



Context Switching Illustration

- Processes are managed by a shared chunk of memory-resident OS code called the <u>kernel</u>
 - The kernel is not a separate process itself, but runs as part of other existing processes
- Context Switches pass the control flow from one process to another
 - Note how going from A to B (and B to A) requires some kernel code to be executed



Process 1's Code

a = b + 1;EIP switch();

Process 2's Code

puts(my_str); switch(); $my_str[0] = '\n';$ i = strlen(my_str); An example of a context switch: there can be different implementations **OS Code** <switch>:

push eax push edx mov [saved esp], esp edx pop ebx pop pop

OS Memory Saved ESP for Process 2

Process 1's Stack

Top Frame

Top Frame

Process 2's Stack

Return addr Saved FAX

ESP

Saved EDX

Process 1's Code

a = b + 1; switch(); EIP b--;

Process 2's Code

puts(my_str); switch(); $my_str[0] = '\n';$ i = strlen(my_str); An example of a context switch: there can be different implementations **OS Code** <switch>: push eax

> push edx mov [saved esp], esp edx pop

> > ebx

pop

pop

ESP

Return addr

Process 1's Stack

Top Frame

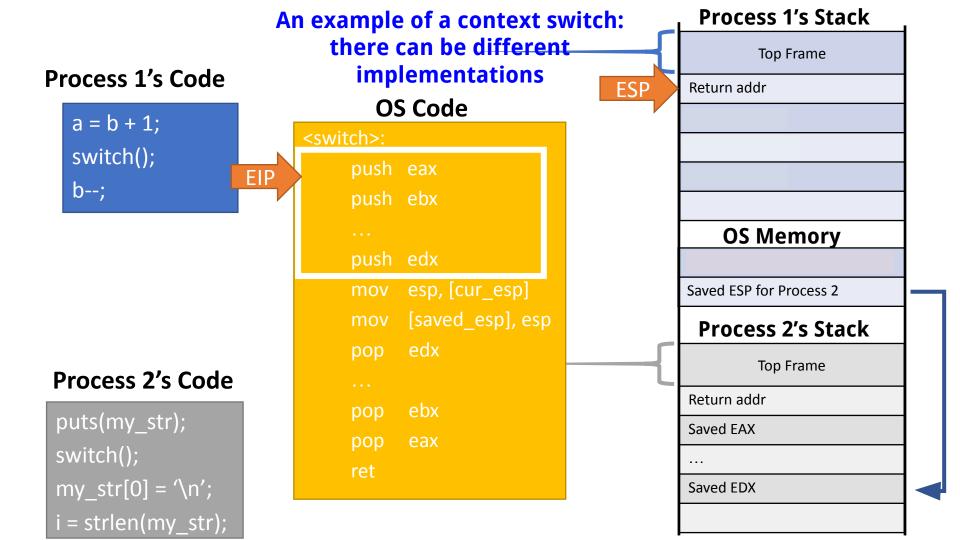
OS Memory

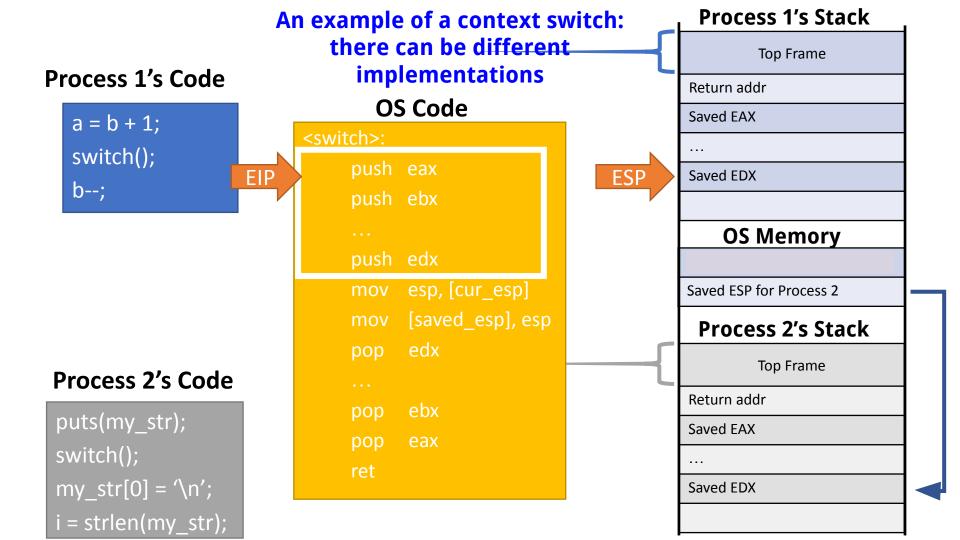
Saved ESP for Process 2 **Process 2's Stack**

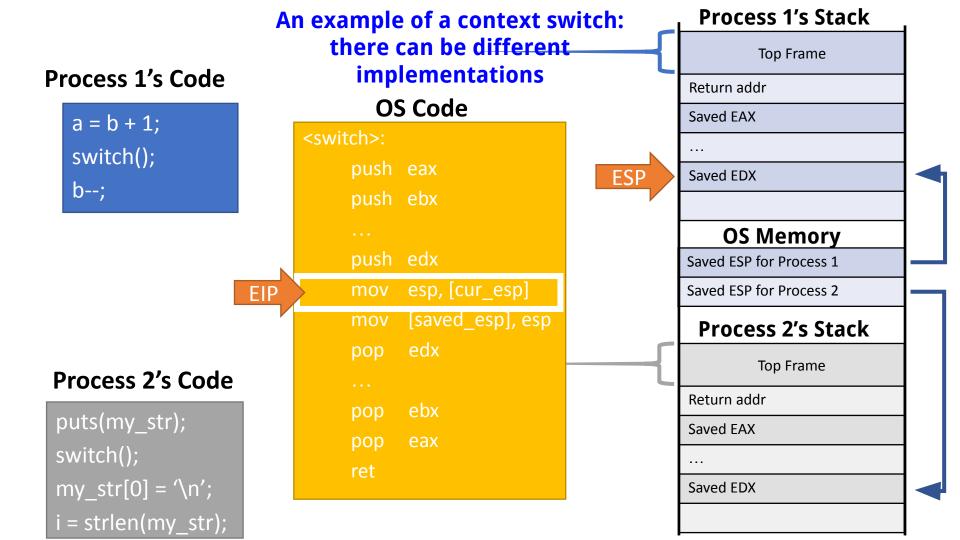
Top Frame Return addr

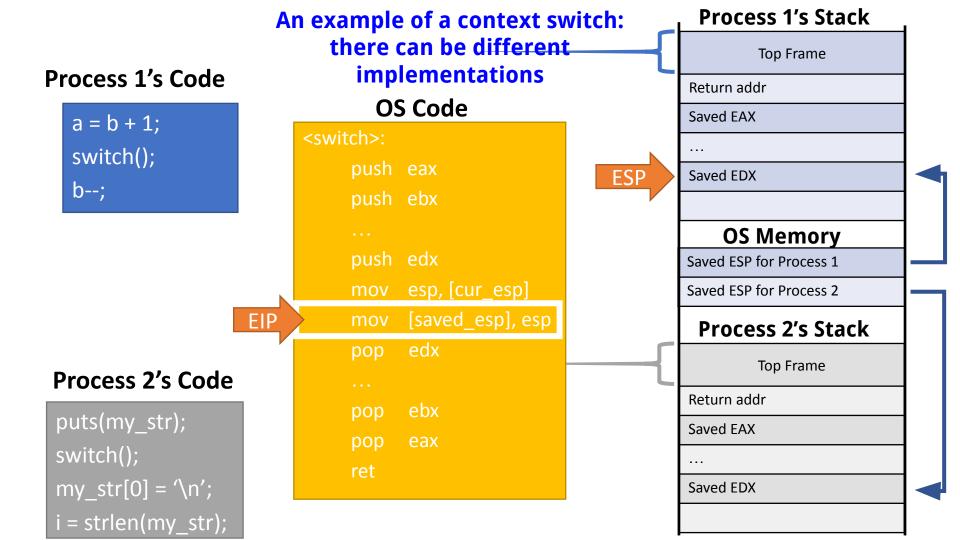
Saved FAX

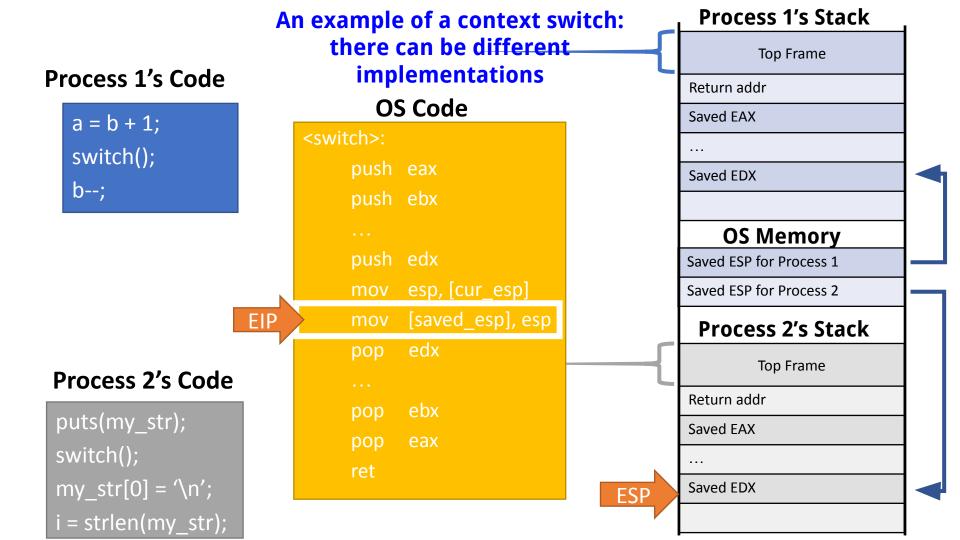
Saved EDX

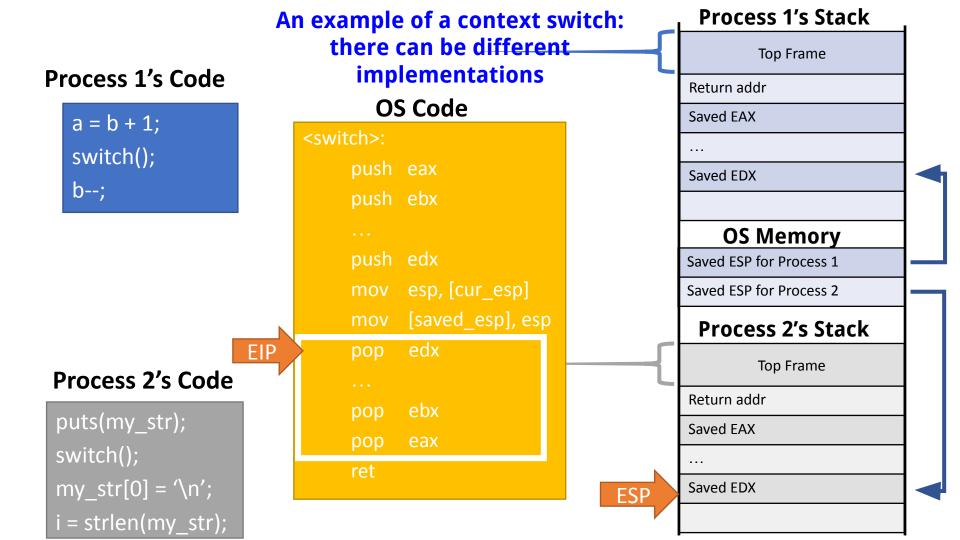


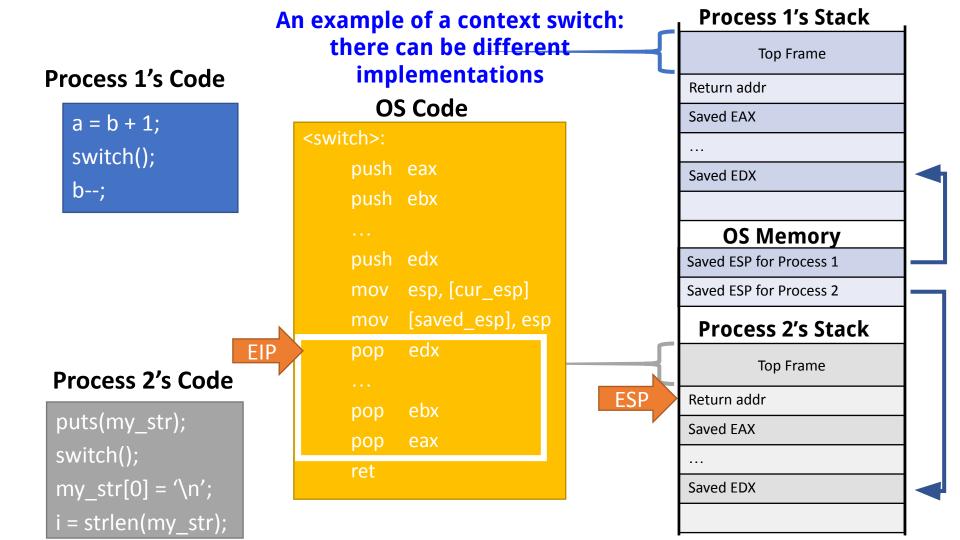


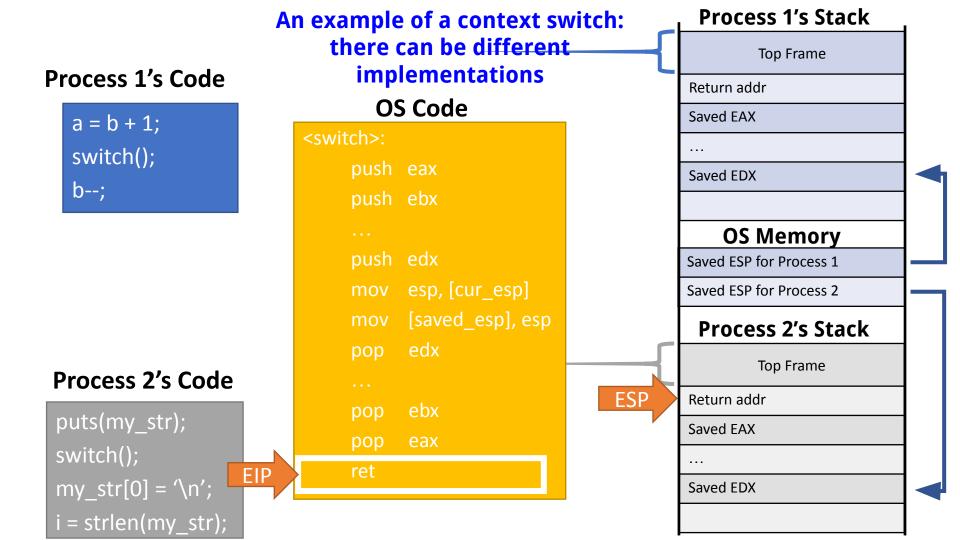












a = b + 1; switch(); b--;

Process 2's Code

puts(my_str);
switch();
my_str[0] = '\n';
i = strlen(my_str);

An example of a context switch: there can be different implementations OS Code

```
<switch>:
    push eax
    push edx
          edx
    pop
          ebx
    pop
    pop
```

Top Frame Return addr Saved EAX Saved EDX **OS Memory** Saved ESP for Process 1 **Process 2's Stack** Top Frame

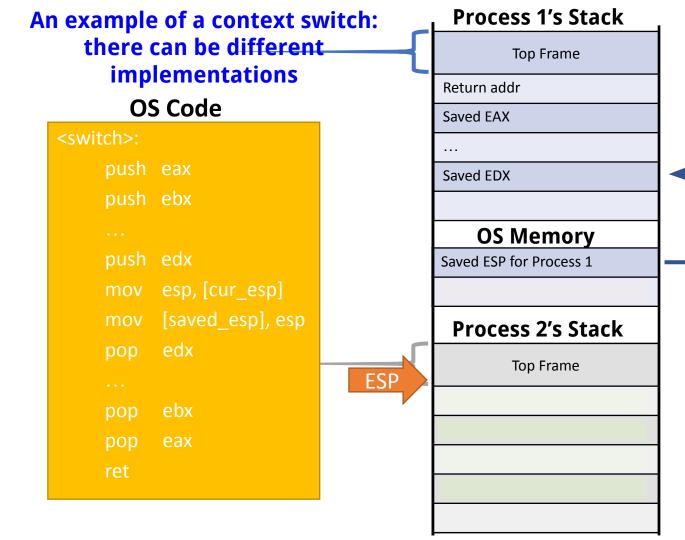
ESP

Process 1's Stack

a = b + 1; switch(); b--;

Process 2's Code

puts(my_str);
switch();
my_str[0] = '\n';
i = strlen(my_str);



Process 2's Code

puts(my_str); switch(); $my_str[0] = '\n';$ i = strlen(my_str); An example of a context switch: there can be different implementations

OS Code

<switch>:

pop

EIP

```
push eax
push edx
mov [saved esp], esp
     edx
pop
     ebx
pop
```

Return addr

Process 1's Stack

Top Frame

Saved EAX

Saved EDX

OS Memory

Saved ESP for Process 1

Process 2's Stack

Top Frame

Return addr Saved FAX

Saved EDX

ESP

Process 1's Stack An example of a context switch: there can be different Top Frame implementations **Process 1's Code** Return addr **OS Code** Saved EAX <switch>: push eax Saved EDX **OS Memory** push edx Saved ESP for Process 1 Saved ESP for Process 2 mov [saved esp], esp **Process 2's Stack** edx pop EIP Top Frame **Process 2's Code** Return addr ebx Saved FAX pop $my_str[0] = '\n';$ Saved EDX i = strlen(my_str);

a = b + 1;

switch();

puts(my_str);

switch();

b--:

Process 1's Stack An example of a context switch: there can be different Top Frame implementations ESP Return addr **OS Code** Saved EAX <switch>: push eax Saved EDX **OS Memory** push edx Saved ESP for Process 1 Saved ESP for Process 2 mov [saved esp], esp **Process 2's Stack** edx pop Top Frame Return addr ebx pop Saved FAX EIP Saved EDX

Process 1's Code

Process 2's Code

puts(my_str);

 $my_str[0] = '\n';$

i = strlen(my_str);

switch();

a = b + 1;

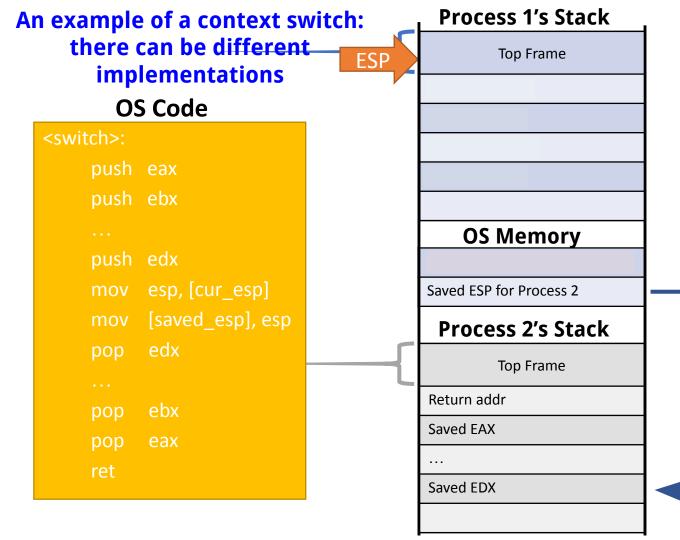
switch();

b--:

a = b + 1; switch(); b--;

Process 2's Code

puts(my_str);
switch();
my_str[0] = '\n';
i = strlen(my_str);



xv6 as an example

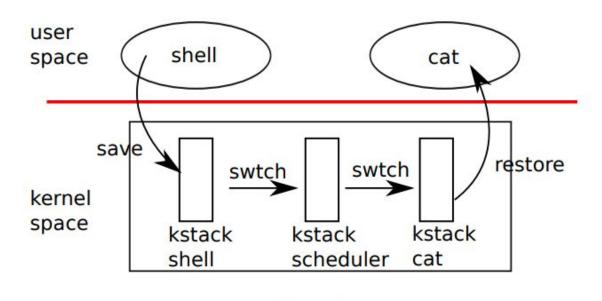


Figure 5-1. Switching from one user process to another. In this example, xv6 runs with one CPU (and thus one scheduler thread).

Kernel

xv6 as an example: process A running in user mode

```
// Per-process state
••• 37
        struct proc {
  39
          uint sz;
                                // Size of process memory (bytes)
          pde_t* pgdir;
                               // Page table
  40
          41
         enum procstate state; // Process state
  42
  43
          int pid;
                              // Process ID
          struct proc *parent; // Parent process
  44
          struct trapframe *tf; // Trap frame for current syscall
  45
  46
          struct context *context; // swtch() here to run process
  47
          void *chan;
                        // If non-zero, sleeping on chan
                     // If non-zero, have been killed
          int killed;
  48
          struct file *ofile[NOFILE]; // Open files
  49
         struct inode *cwd; // Current directory
  50
          char name[16]; // Process name (debugging)
  51
        };
  52
  53
```

https://github.com/mit-pdos/xv6-public/blob/master/proc.h

xv6 as an example: timer interrupt fires

```
#include "mmu.h"
 2
         # vectors.S sends all traps here.
       .globl alltraps
       alltraps:
         # Build trap frame.
         pushl %ds
         pushl %es
         pushl %fs
         pushl %qs
11
         pushal
12
13
         # Set up data segments.
         movw $(SEG_KDATA<<3), %ax
14
15
         movw %ax, %ds
         movw %ax, %es
17
         # Call trap(tf), where tf=%esp
         pushl %esp
         call trap
21
         addl $4, %esp
22
         # Return falls through to trapret...
       .globl trapret
24
       trapret:
         popal
         popl %qs
         popl %fs
         popl %es
         popl %ds
31
         addl $0x8, %esp # trapno and errcode
         iret
```

Hardware (on privilege change) switches to **A's kstack** (from TSS.esp0)

Hardware **pushes**: SS, ESP, EFLAGS, CS, EIP (old user context) onto A's **kstack**

alltraps then saves the rest:

pushal # eax,ecx,edx,ebx,esp,ebp,esi,edi

https://github.com/mit-pdos/xv6-public/blob/master/trapasm.S

xv6 as an example: timer interrupt fires

trapframe

xv6 as an example: trap()

```
//PAGEBREAK: 41
35
       void
36
       trap(struct trapframe *tf)
37
38
         if(tf->trapno == T_SYSCALL){
39
           if(myproc()->killed)
40
             exit();
41
           myproc()->tf = tf;
42
           syscall();
43
44
           if(myproc()->killed)
45
             exit();
46
           return;
47
48
         switch(tf->trapno){
49
50
         case T_IRO0 + IRO_TIMER:
51
           if(cpuid() == 0){
52
             acquire(&tickslock);
53
             ticks++;
             wakeup(&ticks);
54
             release(&tickslock);
55
56
           lapiceoi();
57
           break;
58
```

```
103
         // Force process to give up CPU on clock tick.
         // If interrupts were on while locks held, would need to check nlock.
104
105
         if(myproc() && myproc()->state == RUNNING &&
106
             tf->trapno == T IROO+IRO TIMER)
           yield();
107
```

xv6 as an example: yield() from A to scheduler

```
// Give up the CPU for one scheduling round.
384
385
       void
386
       yield(void)
387
          acquire(&ptable.lock); //DOC: yieldlock
388
          myproc()->state = RUNNABLE;
389
390
          sched();
391
          release(&ptable.lock);
392
```

xv6 as an example: sched()

```
// Enter scheduler. Must hold only ptable.lock
358
       // and have changed proc->state. Saves and restores
359
360
       // intena because intena is a property of this
       // kernel thread, not this CPU. It should
361
       // be proc->intena and proc->ncli, but that would
362
363
       // break in the few places where a lock is held but
       // there's no process.
364
365
        void
        sched(void)
366
367
368
          int intena;
         struct proc *p = myproc();
369
370
371
         if(!holding(&ptable.lock))
            panic("sched ptable.lock");
372
         if(mycpu()->ncli != 1)
373
            panic("sched locks");
374
         if(p->state == RUNNING)
375
376
            panic("sched running");
         if(readeflags()&FL_IF)
377
            panic("sched interruptible");
378
379
          intena = mycpu()->intena;
          swtch(&p->context, mycpu()->scheduler);
380
         mycpu()->intena = intena;
381
382
```

sched() calls swtch(&A->context,
cpu->scheduler):

- saves A's kernel context (callee-saved regs & return EIP) into A->context
- loads the CPU's scheduler context
- jumps to it

xv6 as an example: swtch

```
# Context switch
       # void swtch(struct context **old, struct context *new);
       # Save the current registers on the stack, creating
       # a struct context, and save its address in *old.
      # Switch stacks to new and pop previously-saved registers.
       .globl swtch
       swtch:
10
         movl 4(%esp), %eax
11
12
         movl 8(%esp), %edx
13
14
         # Save old callee-saved registers
15
         pushl %ebp
         pushl %ebx
16
17
         pushl %esi
18
         pushl %edi
19
20
         # Switch stacks
         movl %esp, (%eax)
21
22
         movl %edx, %esp
23
24
         # Load new callee-saved registers
         popl %edi
25
26
         popl %esi
         popl %ebx
28
         popl %ebp
29
         ret
```

the kernel context of the **CPU's scheduler** (the code that runs
scheduler()

Save the current process's kernel registers into p->context,

and restore the scheduler's registers from cpu->scheduler.

xv6 as an example

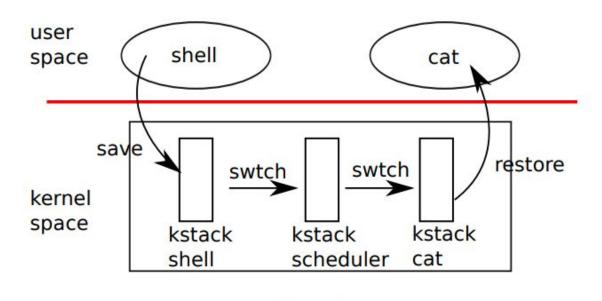


Figure 5-1. Switching from one user process to another. In this example, xv6 runs with one CPU (and thus one scheduler thread).

Kernel

xv6 as an example: CPU's scheduler loop picks B

```
//PAGEBREAK: 42
314
315
        // Per-CPU process scheduler.
        // Each CPU calls scheduler() after setting itself up.
317
        // Scheduler never returns. It loops, doing:
        // - choose a process to run
318
319
        // - swtch to start running that process
320
        // - eventually that process transfers control
321
                via swtch back to the scheduler.
322
        void
323 · · · scheduler(void)
324
325
          struct proc *p;
326
          struct cpu *c = mycpu();
327
          c->proc = 0;
328
329
          for(;;){
            // Enable interrupts on this processor.
330
331
            sti();
332
333
            // Loop over process table looking for process to run.
            acquire(&ptable.lock);
334
            for(p = ptable.proc; p < &ptable.proc[NPROC]; p++){</pre>
335
336
              if(p->state != RUNNABLE)
337
                continue;
338
339
              // Switch to chosen process. It is the process's job
340
              // to release ptable.lock and then reacquire it
341
              // before jumping back to us.
342
              c->proc = p;
343
              switchuvm(p);
344
              p->state = RUNNING;
345
346
              swtch(&(c->scheduler), p->context);
347
              switchkvm();
348
349
              // Process is done running for now.
350
              // It should have changed its p->state before coming back.
```

The scheduler scans for a **RUNNABLE** process (here, **B**).

switchuvm(B):

- loads CR3 with B's page directory (B's address space)
- sets the CPU's TSS.esp0 = top of B's kstack (for future user→kernel entries)

swtch(&cpu->scheduler, B->context):

- saves scheduler's context
- restores B's kernel context
- jumps to B->context->eip (resuming B's kernel thread where it last yielded/slept) - ret instruction in swtch