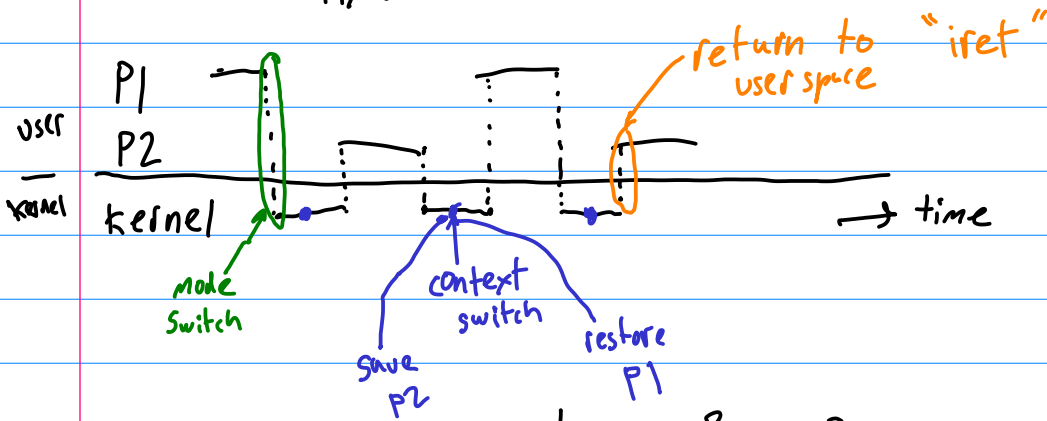
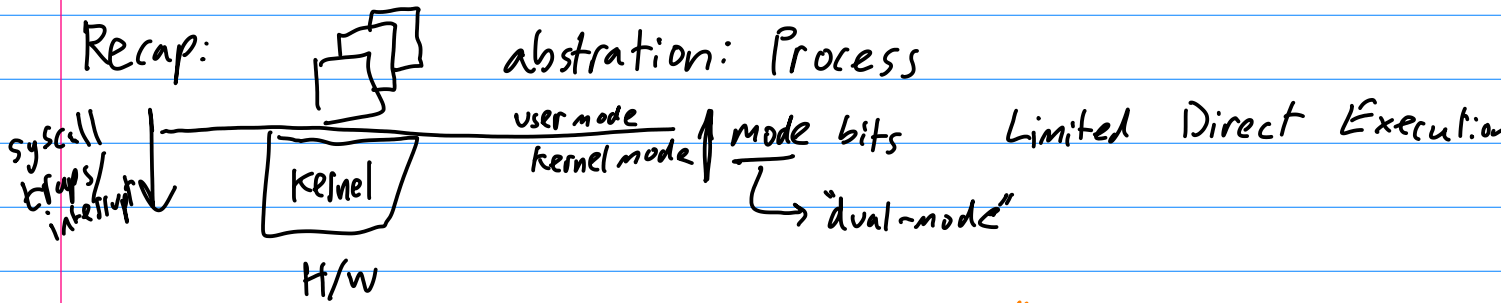


# CS 3214 Lecture #3 "Unicode"

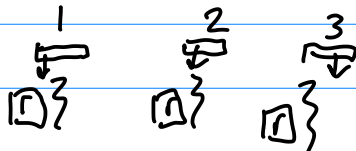
[courses.cs.vt.edu/cs3214/williams/](http://courses.cs.vt.edu/cs3214/williams/)

Office Hours: T/Th 1-2pm on Zoom

Recap:  abstraction: Process



user-threading  
stacks?



"scheduling"

reg?

[address space  
is shared]

1. cooperative scheduling: "yield"

preemption

Coming up: Process states: (runnable?)  
Management: (create / destroy) // next time

Detail about interacting w/ the system: ExO  
character sets & unicode  
"appreciation of details"

UNIX philosophy: small programs do one thing well connect w/ universal text  
 cat foo | grep bar | wc

when do you need unicode/char sets?

- user input (web forms/web servers)
- text file processing (copying / conversion/validation/display)
- i18n internationalization kubernetes: k8s

byte - unit of digital info

octet: 8-bits

0... 255

0x00 ... 0xff

History: byte didn't have to be 8

C: uint-8 foo;  
 signed char foo;

CHAR\_WIDTH

POSIX = 8

ints, shorts, longs

endianness

0x0102  
 1x256 + 2 = 258  
 1 + 256\*2 = 513

Characters: abstract entities from some "alphabet"

				0	
"	"	"	"	"	"
0	1	2	3	4	5 <del>XX</del> 2 <sup>(3)</sup> 36bits

variable-length encoding

00   01   10   1100   1101   1110

[ 00 01 10 1100 ]  
 K   T   F   ♥

~~00 11 11~~

Real world characters

- used to be many ASCII, ISO-8859-1, ISO-8859-2...
- Unicode ← one character set to rule them all

characters "code points"

↓  
1, 114, 112

U+0041

"A"

U+00C4

"Ä"

U+1F385



Santa

← glypheme

but... U+0308 unilaut over preceding char

U+0041 U+0308 "Ä"

how to encode

- in prog mem.?
- on disk?

ASCII

Tradeoff: ease of processing vs. efficiency of storage

UTF-32

- 32 bits for each char

- santa 0x0001F385

Linux wchar\_t

C++ char32\_t

endianness

BOM "byte order mark"

S[10]

4x overhead over ASCII

UTF-16

- 16 bits for each char

(some have 4 bytes)

0xD83C

0xDF85

- sometimes can index S[10]

- still wasteful, error-prone

UTF-8

- 8 bits variable-length 1, 2, 3 or 4 bytes

- Santa 0xFD 0x9F 0x8E 0x85

Python's  
Swift

Java  
JavaScript

UTF-8 (cont'd) ← most common when unicode transmitted/stored

Rust,  
Go

- ASCII ← 1 byte      0x7A "z"
- space efficient
- indexing is hard/impossible

"crate"  
"packages"

Given a seq. of bytes, can you tell the encoding?

- No

- can say something is invalid

- metadata is needed HTTP response:

Content-Type: text/html; charset=UTF-8

- files? operator/program