

Miscellaneous Thingies



Four Miscellaneous Topics

- Review of the memory model and C->Assembly flow
 - Particular weakness on the midterm, and also a perennial source of problems...
 - But now you have the knowledge to understand ***why*** it is
- Building a software pipeline in a 61C style
 - Or "How Nick pegged all cores on The Beast"
- Alternate Compute Fabrics:
 - Graphics cards, 2-D Systolic Arrays, FPGAs
 - (Not on the final but you should know...)
- Cool architecture features:
 - ARM PACs
 - RISC-V 16b ISA
 - (Not on the final but you should know...)

Remember: Everything is bits...

- We take the bits and group them into bytes
 - 8 bits to a byte
- We address memory by byte
- And memory is abstractly treated as a big *big* bucket of bytes
 - At address 0 to address **0xFFFFFFFF** (on a 32b architecture)
 - At address 0 to address **0xFFFFFFFFFFFFFFFFF** (on a 64b architecture)
- A pointer is simply an address to someplace in main memory

And Alignment...

- We group bytes into words
 - 4 byte words on a 32b architecture:
Same size as a pointer
- But we need to keep things aligned
 - We saw that when we build the load/store unit for the RISC-V project:
Bytes/characters we load anywhere
 - 16b half-words can only be stored/loaded at addresses starting with 0, 2, 4...
 - "Half-word aligned": address % 2 == 0
 - 32b words need full word alignment: address % 4 == 0
- Some architectures don't **require** alignment but...
 - Even on those, if you don't align data your performance will go down as loads become a lot slower

Alignment Continues To Larger Sizes...

- For RISC-V, we are only dealing with a 32b architecture
 - So alignment only needs to be 4B aligned
- But SIMD-vector architectures can be wider
 - EG, Intel AVX2 uses 256b (32B) vectors
- So the high performance load of 256b of data:
 - `_mm256_stream_load_si256`
 - Requires the load to be 32B aligned
 - Without this requirement the load path would be significantly more complicated

So When Creating Structures...

- The compiler needs to pad things for alignment
- The basic rule for a 32b RISC-V:
 - Characters/bytes: No padding
 - 1/2 words (e.g. shorts): **Must** start at 0, 2, 4 etc...
 - Words (e.g. int, pointers): **Must** start at 0, 4, 8 etc...
- And then add enough space at the end...
 - Have to make the total size $\% 2 == 0$ if only characters/half-words
 - Have to make the total size $\% 4 == 0$ if includes words

Live Example...

- ```
struct foo {int a;
 char *b;
 char c,
 int d};
```

# And so under-the-hood pointer arithmetic

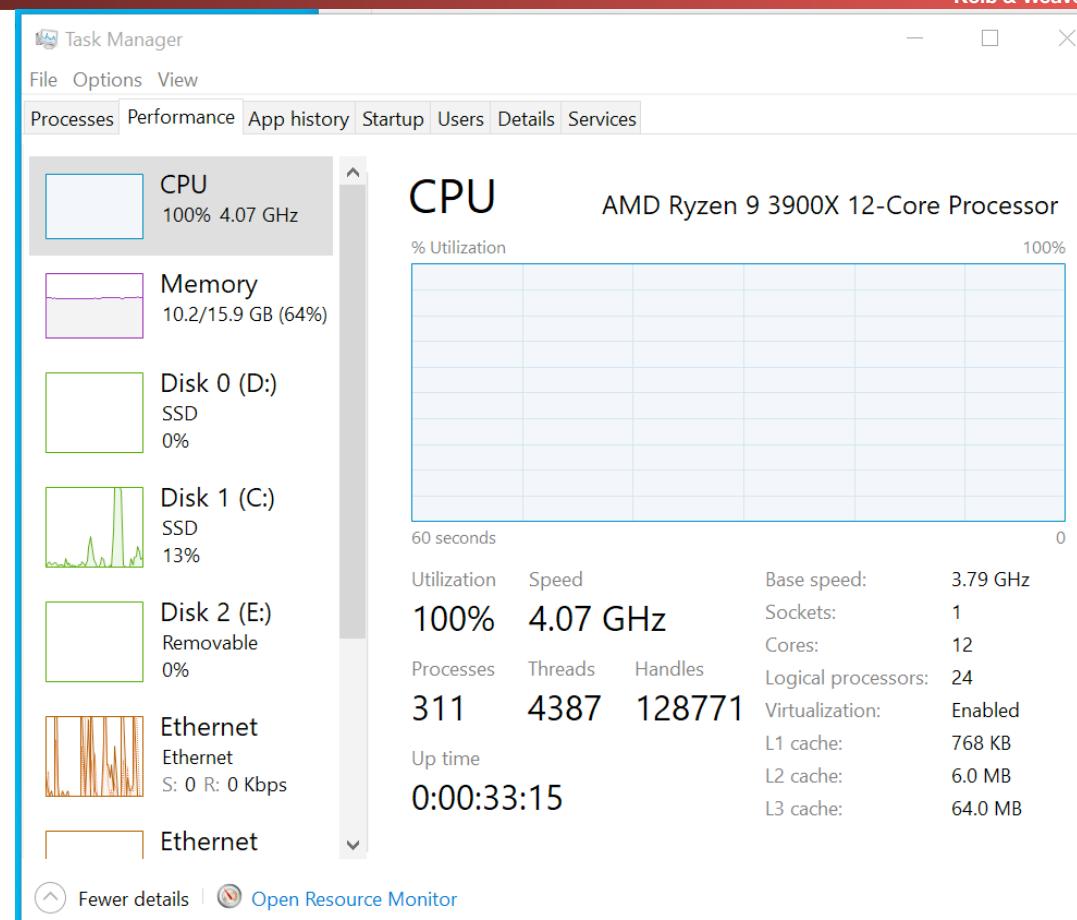
- struct foo {int a; char \*b; char c, int d};
- Remember **struct foo \*f** and **struct foo f[]** are effectively the same, and the compiler knows **sizeof(struct foo)**
- So to convert **f[i].c** to assembly:
  - Take **f** and load it in a register
  - Take **i** and multiply it by **sizeof(struct foo)**
    - Use a shift if possible
  - Add **f** and **i** together
  - Then load **(f + i) + offset-in-structure**

# So a richer example...

- ```
void bar(struct foo *f,
          int max) {
    int i;
    for(i = 0; i < max; ++i) {
        baz(f[i].c)
    }
}
```

Building A Data Pipeline... 61C Style on my Linux box

- Some bad actor stole >4 GB of data from UCOP
 - Basically everything that was on the "secure" file transfer server in December
- The bad actor released at least some of this data **publicly**
 - As a 4GB compressed archive that anyone can download
 - What information about **me** was in the archive?



What I need to know...

- I already know my social security # got breached
 - They told us that...
 - But I've got fraud alerts & freezes in place already
- But what other information?
 - Address? Phone #? Things I don't know about?
 - Tax information?
 - ***Banking information?***
 - The numbers on the bottom of a check are all an attacker needs to make fake checks

The Nature of the Dump...

- A ***lot*** of pdf files
 - PDFs are a pain to search, need to convert to text
- A ***lot*** of data tables
 - Some as comma-delimited text, some as excel spreadsheets, some in stada format
- Need to convert it to something reasonable
- Google around...
 - Nice linux OCR pipeline cobbled together:
PDF -> images -> OCR text
 - pandas can read both xlsx and stada files

Step 1: File Conversion

- Want to convert everything into text files
- Obscenely parallel problem:
 - For every PDF do X
- But with some gotchas...
 - I can't just spawn 700 PDF->txt conversion programs
 - That would grind my machine to a halt
 - And different invocations take a different amount of time
 - So I can't just spawn off tasks at a reasonable interval
- Two approaches
 - Dynamically tune based on load...
 - Or just say "F-it, and keep X jobs live"

Keeping X jobs live: Fork/join with a limiter

- Used a simple golang hack
- `capacity := make(chan bool, 10)`
`done := make(chan bool)`
- `func run(txt string) {`
 `c <- True;`
 `....`
 `<- c;`
 `done <- True;`
}
- main just calls "go func()" on every line of stdin...
and then an equal number of lines of "<- done"

A bit of tuning...

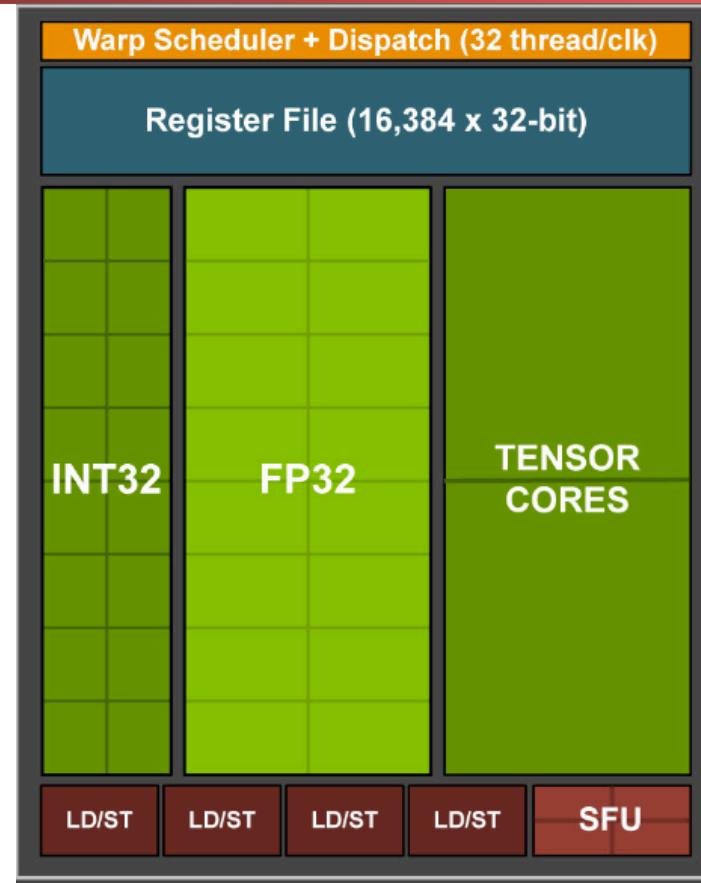
- ~10 jobs pegs all CPU cores on the OCR pipeline...
 - And it took a couple hours to PDF->txt the lot:
Driver used multiple threads for single documents
- The .xlsx and stata conversion was a lot faster...
 - Set to ~25 (since python doesn't thread, especially on this task)
 - Took ~10 minutes or so
- Not fully efficient...
 - At the end of the PDF run there was no longer pegged CPU
 - 100% CPU utilization means efficiency loss due to context switching:
Optimum would be ~95%
 - Also, really stressing the Windows virtualization...
 - I do all my work in "linux" under WSL

And Now To Search

- Just use the same pipeline with grep...
- But...
 - Ends up **not** pegging the CPUs...
Instead I'm pegging the "disk"!
- OK for just searching for me, but...
 - Want to be able to do a "for anyone who wants" service
- So to do this, parallelize on an alternate axis:
 - Don't check one person at a time, check **all** people using a single program
 - And then invoke that in parallel across all files
- Gotcha problem: Need to make sure to synchronize writes well
 - Again, golang FTW:
A channel for each user's results, the search does an atomic write to the channel

Alternate Architectures: GPUs

- The GPUs these days are semi-SIMD beasts
 - "For every element do X"
- Shortcuts that prove to be good for some compute problems
 - Single precision and ***half precision*** floating point
- New addition: TENSOR cores in addition to SIMD processing
 - Only really suitable for machine learning



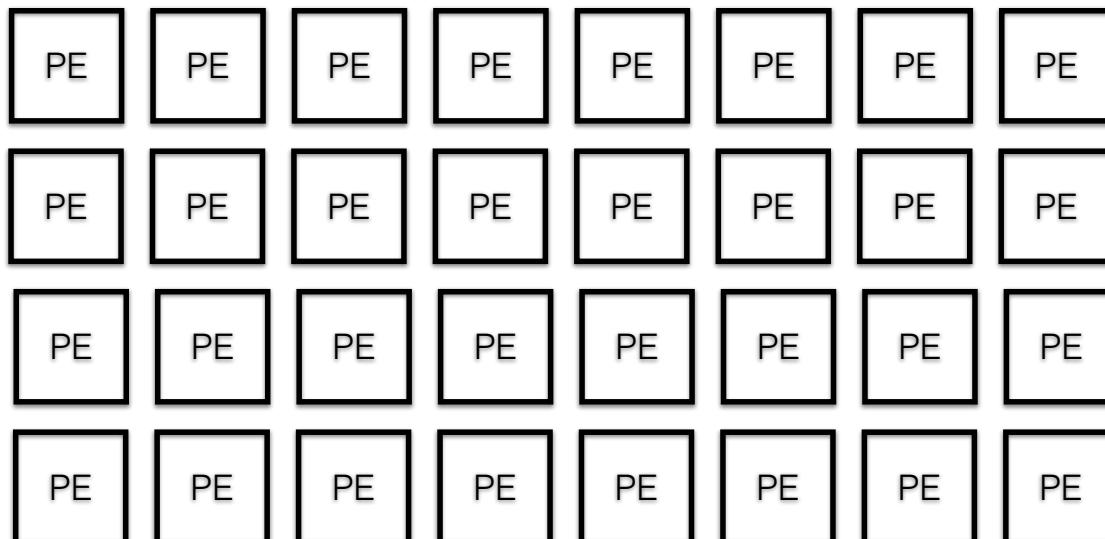
The "AI" compute problem...

- In the end, it becomes dense matrix multiplies...
 - But there is a bit of specialness...
AI only cares about **crappy** math! 16b floating point or even 8b or 4b integer math
- Even in SIMD computation we do a lot of memory loads
 - So can we have an architecture that gives us way more compute on each load?



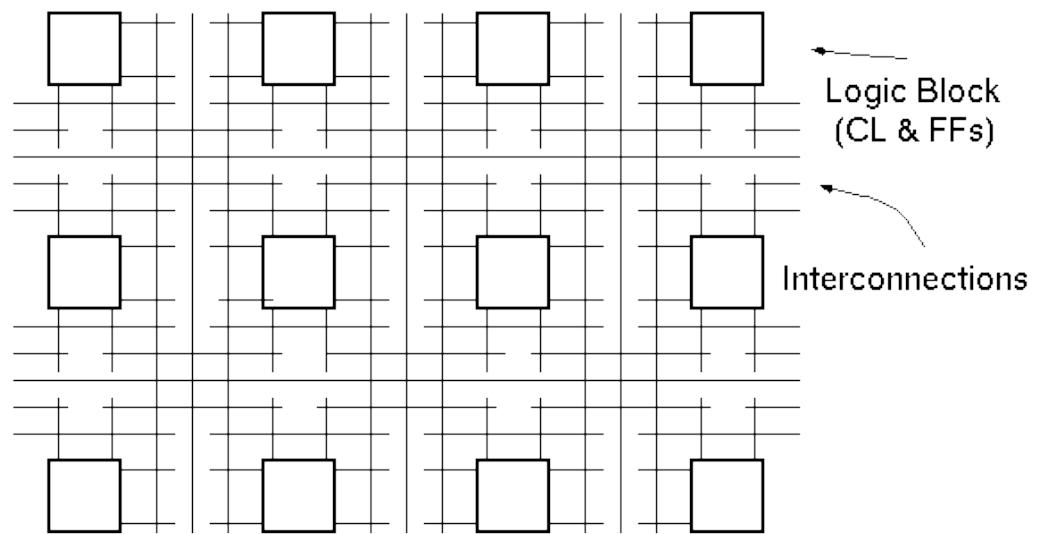
Tensor Accelerators: Specialized 2D "Systolic Arrays"

- We can pipeline and stream our loads...
 - But we are still doing a *lot* of loads
- But if we arrange compute elements in an array
 - Each one can take data, compute on it, and shift it on to the next element



FPGA Overview

- Basic idea: two-dimensional array of logic blocks and flip-flops with a means for the user to configure (program):
 - the interconnection between the logic blocks
 - the function of each block.



Simplified version of FPGA internal architecture

Why are FPGAs Interesting?

- Technical viewpoint:
 - For hardware/system-designers, like ASICs - only better: “Tape-out” new design every few minutes/hours.
 - “reconfigurability” or “reprogrammability” may offer other advantages over fixed logic?
 - In-field reprogramming? Dynamic reconfiguration? Self-modifying hardware, evolvable hardware?
- Recent trend: Datacenter FPGAs...
 - Actually not a great idea in most cases, we'll see why in a bit...

Why are FPGAs Interesting?

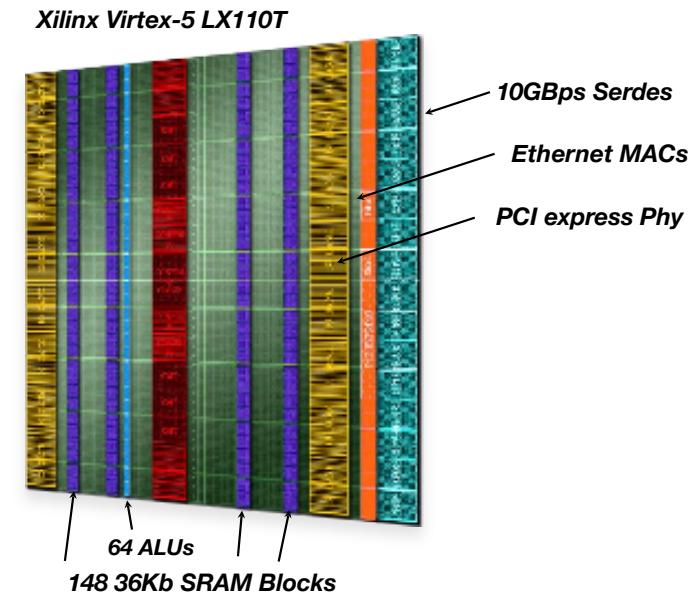
- Staggering logic capacity growth (10000x):

<i>Year Introduced</i>	<i>Device</i>	<i>Logic Cells</i>	“logic gate equivalents”
1985	XC2064	128	1024
2011	XC7V2000T	1,954,560	15,636,480

- FPGAs have tracked Moore’s Law better than any other programmable device

Why are FPGAs Interesting?

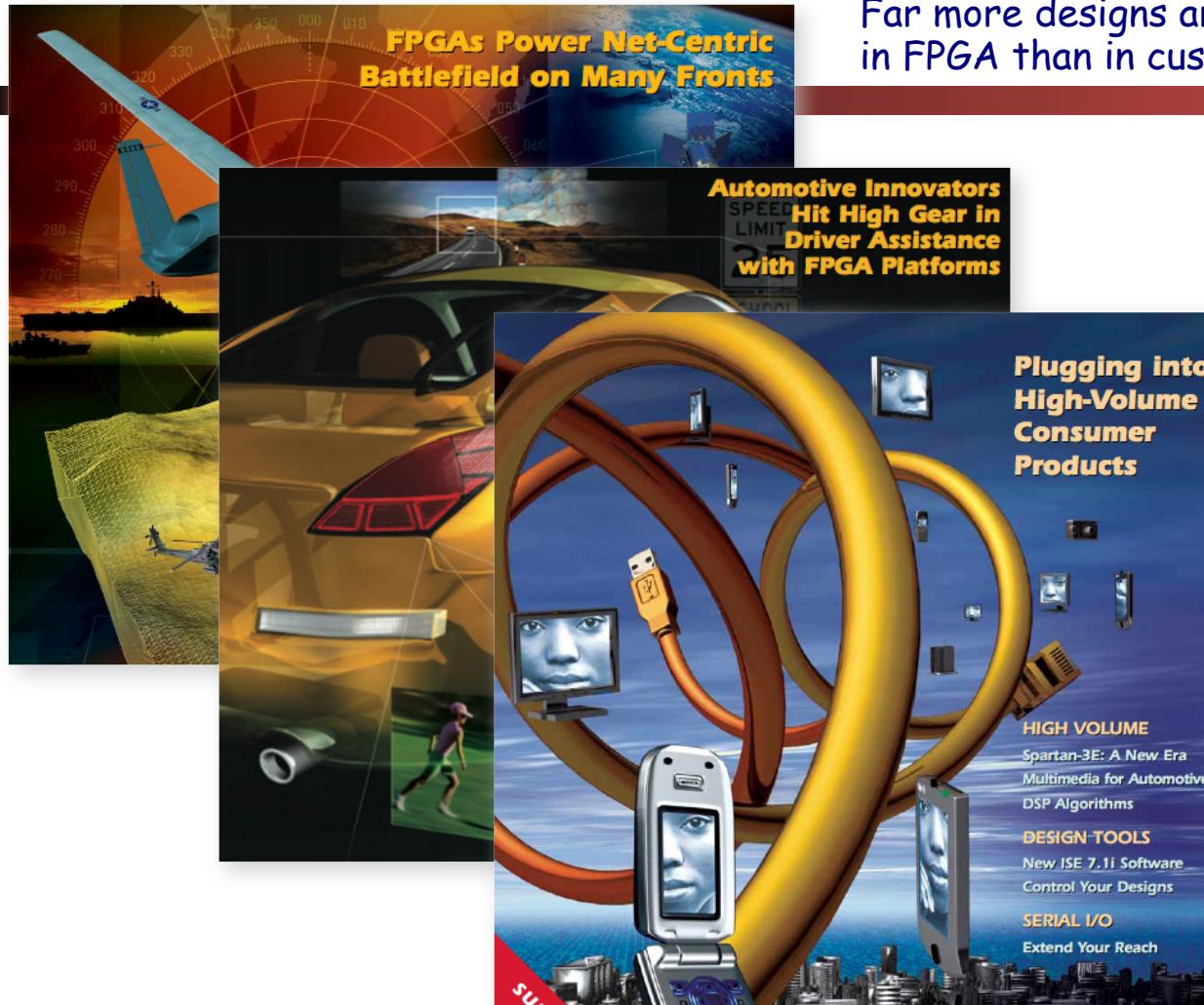
- Logic capacity now only part of the story: on-chip RAM, high-speed I/Os, “hard” function blocks, processors ...
- Modern FPGAs are “reconfigurable systems”



FPGAs are in widespread use

Computer Science 61C

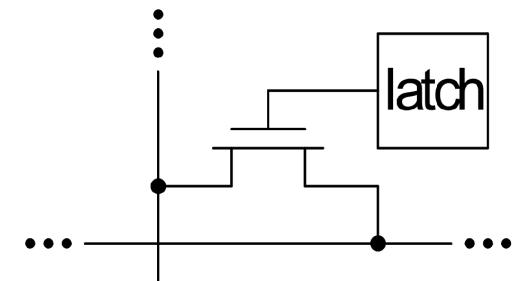
Kolb & Weaver



Far more designs are implemented in FPGA than in custom chips.

User Programmability

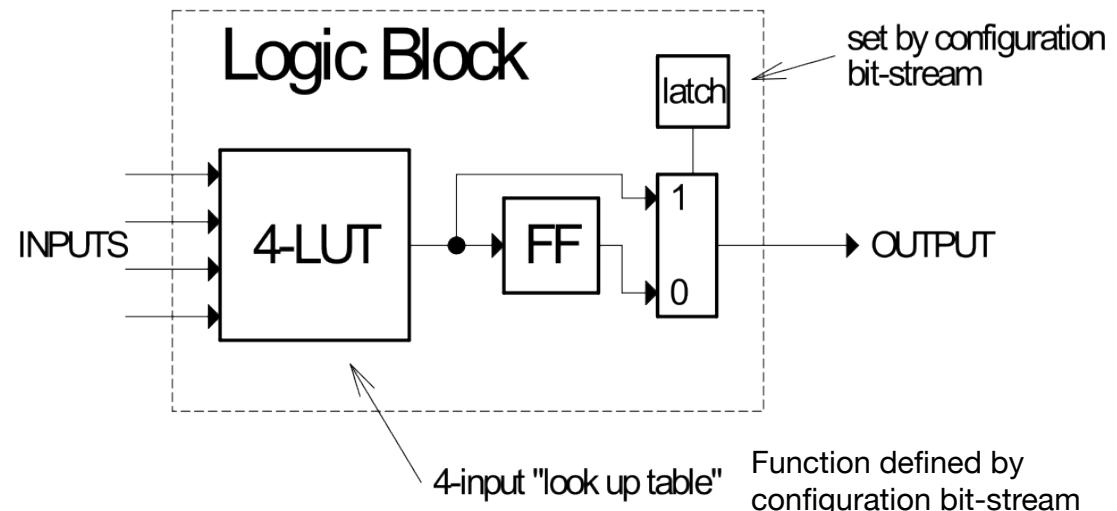
- Latches are used to:
 - control a switch to make or break cross-point connections in the interconnect
 - define the function of the logic blocks
 - set user options:
 - within the logic blocks
 - in the input/output blocks
 - global reset/clock
- “Configuration bit stream” is loaded under user control
 - *Latch-based (Xilinx, Altera, ...)*



- + *reconfigurable*
- *volatile*
- *relatively large.*

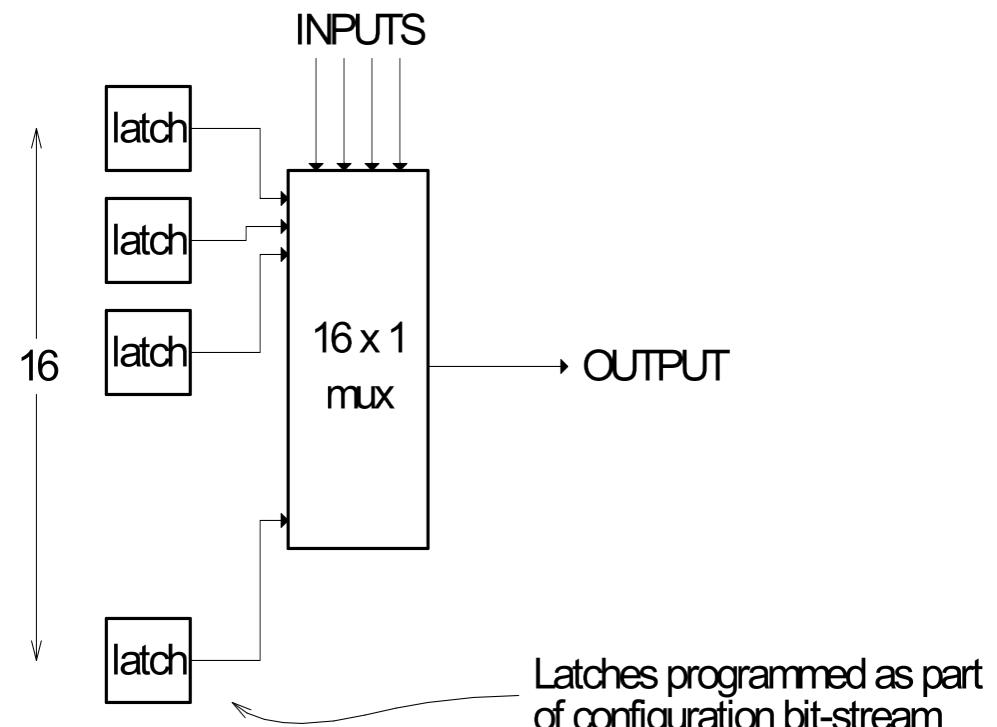
Idealized FPGA Logic Block

- 4-input look up table (LUT)
 - implements combinational logic functions
- Register
 - optionally stores output of LUT



4-LUT Implementation

- n-bit LUT is implemented as a $2^n \times 1$ memory:
 - inputs choose one of 2^n memory locations.
 - memory locations (latches) are normally loaded with values from user's configuration bit stream.
 - Inputs to mux control are the CLB inputs.
- Result is a general purpose "logic gate".
 - n-LUT can implement any function of n inputs by directly implementing the truth table



LUT as general logic gate

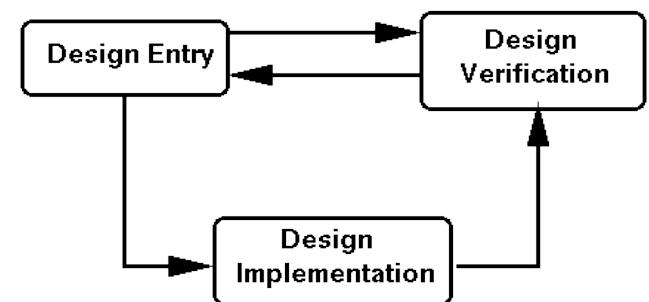
- An n-lut as a direct implementation of a function truth-table.
- Each latch location holds the value of the function corresponding to one input combination.

Example: 4-lut

INPUTS		
0000	$F(0,0,0,0)$	← store in 1st latch
0001	$F(0,0,0,1)$	← store in 2nd latch
0010	$F(0,0,1,0)$	←
0011	$F(0,0,1,1)$	←
0011		
0100		•
0101		•
0110		
0111		
1000		
1001		
1010		
1011		
1100		
1101		

FPGA Generic Design Flow

- Design Entry:
 - Create your design files using:
 - Schematic editor (like Logisim) or
 - HDL (hardware description languages: Verilog, VHDL)
- Design Implementation:
 - Logic synthesis (in case of using HDL entry) followed by,
 - Partition, place, and route to create configuration bit-stream file
- Design verification:
 - Optionally use simulator to check function,
 - Load design onto FPGA device (cable connects PC to development board), optional “logic scope” on FPGA
 - check operation at full speed in real environment.



What Are They Good For...

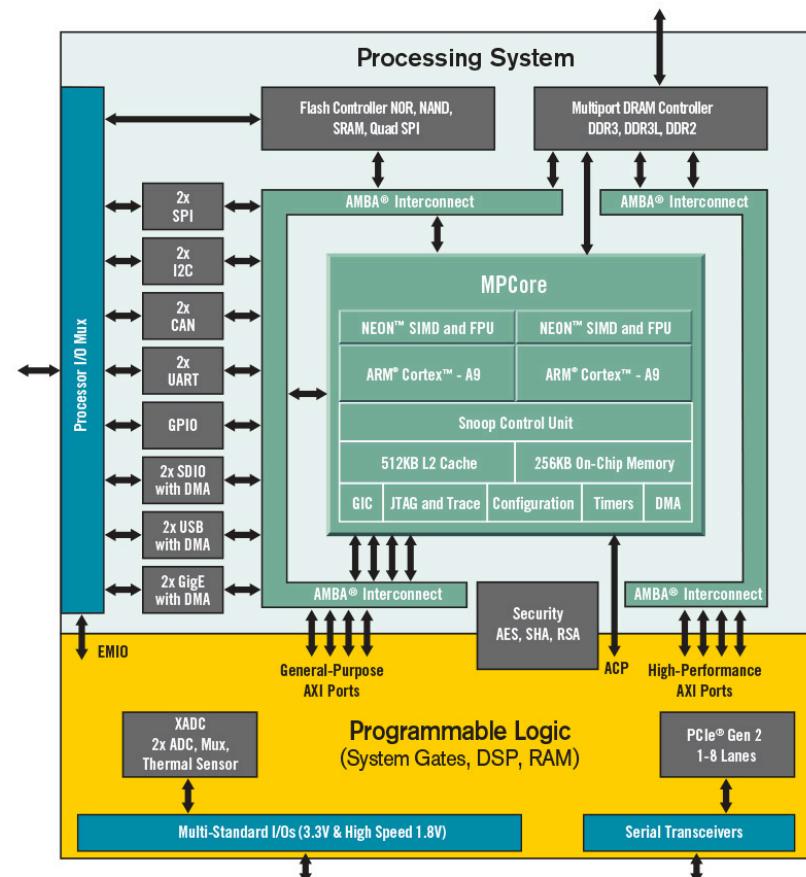
- **ASIC-lite:**
 - Custom chip tasks without the huge (\$1M+) engineering cost of a custom chip
- **"Glue" to integrate things together**
 - The problem with standards is there are so many to chose from:
Get all the devices talking together
- **"Bit-twiddling operations"**
 - Where each bit needs to be handled differently
- **Video on a per-pixel basis**
 - Cameras provide the image one pixel (or a few pixels) at a time...
 - Unlike GPU/SIMD which wants the whole image first...
delaying things by 1/30th of a second!

What They Are *Not* Good For

- Word-width operations
 - People keep trying to build word-oriented FPGAs, and just would do better using GPUs, vector machines, 2D systolic arrays, etc...
 - Instead just some word-supporting large blocks instead
- Floating point: They don't *gots* it
- Branches!
 - Conditions in a processor/GPU are designed to happen:
We understand branches. Even GPUs have some limited conditioning
 - Branches in an FPGA usually involve muxes and computing **both** options:
wastes power and computational fabric
- Why I'm down on datacenter-FPGAs
 - Only useful for a few limited apps
 - For those, efficiency/cost savings win over GPUs is minor

But Modern Coolness: FPGA/Processor Hybrids

- EG, Xilinx Zynq-7000 series
 - Low(ish) cost FPGA/processor hybrid
 - Dual core, 600 MHz ARM
 - 32b ARM, superscalar
 - Full linux software stack available
 - Coupled to a significant FPGA fabric
 - FPGA has interfaces to other devices, main memory, and L2 cache on the ARM cores
 - FPGA also has ~50-100 36kb distributed memories and an equal number of fixed-point multipliers
 - Plus a lot of pins on its own connecting to the rest of the world



Observations...

- Much worse memory bandwidth than a standard computer
 - Only 1-2 DRAM chips, so far narrower interfaces
- Decent latency
 - Normal DDR3 DRAM, good caches
- Decent CPU
 - A Raspberry Pi 3 is better, but still decent
- Good parallelism
 - Sweet spot on cost/performance superscalar, 2x cores

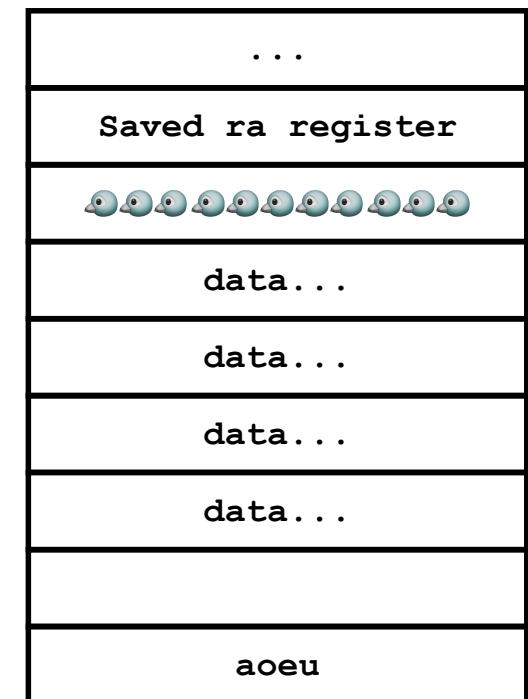
Memory Hardening... ARM Pointer Authentication Codes

- Attackers want to overwrite memory...
 - When your C code fails to check a buffer
- The classic vulnerability

```
void foo () {  
    char c[32];  
    gets(c)  
}
```
- Attacker gives you too long an input...
 - And c is stored on the stack
 - So the attacker overwrites not just c but the other stuff on the stack...
 - Such as the saved **ra**
 - The saved **ra** is overwritten to point to the attacker's code in memory

Stack Canaries...

- Goal is to protect the return pointer from being overwritten by a stack buffer...
- When the program starts up, create a ***random*** value
 - The “stack canary”
 - When starting a function, write it just below the saved frame pointer
 - When returning in a function
 - First check the canary against the stored value



Stack Canary Overhead...

- May require enabling an optional compiler flag...
 - So of course it is commonly not done!
- Requires a memory load & store on every function entrance
 - Highly cacheable so basically only 4 instructions on a typical RISC:
Load address of canary (2 instructions)
Load canary value into register
Store canary value onto stack
- Requires 2 memory loads and a (probably) not taken branch on exit
 - So 5 instructions on a typical RISC:
Load address
Load canary value
Load canary off stack
BNE (mark as probably-not-taken if you can)

So example code...

- ```
la t0 canary # Reminder, turns into two
 # instructions
lw t0 0(t0)
sw t0 x(sp) # four below where ra got stored
 # if we don't bother saving the frame pointer
```
- ```
la t0 canary
lw t0 0(t0)
lw t1 x(sp)
bne t0 t1 dead_canary
                  # Make sure this is a forward branch:
                  # So CPU assumes it won't be taken
```
- Note also generally sequential:
only parallelism present is in loading the canary from both the stack and storage

Brute Force...

- Brute force: just simply try every possibility
 - Or if its a different random # each time, just always try the same number
- Even the smallest timeout goes along way:
 - If you can try 10,000 per second, trying 2^{20} possibilities takes less than 2 minutes
 - If you can only try 10 per second, it takes a day and a half
 - And if 10 failures causes a 10 minute timeout...
Forgettaboutit!
- Exponentials matter
 - If it take 1 minute to try 2^{20} , it will take 16 hours to try 2^{30}
 - And 2 years to try 2^{40} !
 - EG, Apple added a mitigation in the latest iOS:
Crashing programs can (optionally) have an exponentially growing delay on restarting from crashes, which prevents attacks that need to repeatedly crash the service to extract information or get lucky

Pointer Protection: Modern 64b ARM 8.3 Pointer Authentication

- <https://www.qualcomm.com/media/documents/files/whitepaper-pointer-authentication-on-armv8-3.pdf>
- ARM64 uses 64b pointers
- Idea: Since our pointers are 64b but we are only using say 42b of them...
 - Lets use that upper 22b to encrypt/protect pointers of various types!
- New instructions:
 - **PAC** -> Set Pointer Authentication Code
 - Sets the upper bits with a cryptographic checksum
 - **AUT** -> Check and Remove Pointer Authentication Code
 - If the check is invalid, it will instead put an error in the checksum space:
If the pointer is dereferenced it causes an error
 - **XAUT** -> Strip PAC without checking
- Instructions are in NO-OP space if the processor doesn't support them

Plus some non-NOOP higher performance options

- When you know you will be running on a processor which supports it
 - check & return:
Check the return address has a valid PAC and if so, return
 - check & load:
Check the PAC and if so, load the pointer
 - check & branch:
Check the PAC and if so, do a jump-and-link to that pointer
- Allows the complete elimination of the overhead for checking!
 - Well, cheat: You cause it to trigger an exception on instruction committing and just assume the pointer is valid to start with...

How To Use...

- There are 5 secrets for pointer protection
 - These contain random 128b secrets that are used to authenticate the pointer:
Provided by the OS
 - Two for data (DA/DB), two for instruction (IA/IB), and one general purpose (GA)
- These are contained in processor registers,
and are ***not readable to the program itself!***
 - Key property: An information leakage vulnerability can't defeat this protection on a user-level program
 - But it could on a kernel level program:
Solution would be to also have a secret random to the CPU that is included but non readable
- Other workaround: find a vulnerability that can trick the program into authenticating new pointers it shouldn't, or be able to reuse authenticated pointers in another context

So in practice

- The PAC is a function of the pointer, an additional register (or register 0) and the hidden secret
 - **PACIA x30 sp**
AUTIA x30 sp
Protect/Authenticate x30 as a function of x30, sp, and the secret data associated with the Instruction A context (x30 is the default link register for ARM == **ra** in RISC-V)
 - Thanks to crypto-magic discussed in 161, the PAC's "look random"
 - Changing a single bit of anything should result in something looking totally different and random
 - So to guess a 22 bit PAC would be 1 in 4 million odds.

So Cheaper Stack Canaries...

- On function entry: Create the PAC for the return address
 - Using the stack pointer as the context itself:
This means the return address can't even be moved
- On function exit: Check & return as normal
- With backwards compatibility: only 2 instructions
 - **PACIA** on function start, **AUTIA** on function end
- Without backwards compatibility: only 1 instruction!
 - Just the **PACIA** on function start and a check & return on exit
 - Saves 8 instructions... Or >85%!
- Only 22 bits of entropy but...
 - If you get more than a few failures, just keep the program dead!

Or Protecting vtable pointers...

- When you allocate a new C++ object...
 - The first thing is a "vtable pointer", really a pointer to an array of pointers to functions
 - Attackers want to overwrite this with their own version
 - Protect the vtable pointer with a context and register 0:
One additional instruction when calling `new()`
 - Then have the vtable itself live in read-only space
- Now when calling a virtual function...
 - Check & Load the vtable pointer (RISC-V like pseudocode):
eg, if the object pointer is in `s0`, the vtable pointer is at the start of `s0`...
`LDRAA t0 0(s0) # Load 0 + s0, authenticated with data A`
`LW t0 x(t0) # x == the specific function to call`
`JALR t0 # Actually call it`
- Now you **can't** overwrite a C++ object's vtable pointer to something else without either being very lucky, finding a separate vulnerability, or replacing with another valid pointer that you acquire...
And the overhead is literally **nothing!**
 - Apart from you need to recompile and using the latest ARM silicon, that is

Probably the biggest benefit for Apple going to ARM

- MacOSX ARM will be able to ***assume*** PAC support!
 - Since it is Apple A12 or newer processors only
 - Latest iOS also just started really aggressively turning on PAC support
- Can therefore use the more efficient primitives:
 - Check & Branch Register, Check & Load, Check & Return which all eliminate the instruction needed in a separate check
 - Usable in both the kernel and user space:
Acts to harden both applications and the underlying OS
- x86 has nothing like this in the pipeline!
- If you have a choice of architecture for a product: ARM 8.3+
 - This gives you so much real-world security for crappy C-code

RISC-V 16b ISA

- Observation:
Although we encode instructions with 32b, a lot of the instructions follow common patterns
 - Some registers used a lot more than others
 - Immediates are often small
 - Same source and destination for 3-operand operations
- So the optional "C" instruction set
 - Look at the first two bits of an instruction can determine its type
 - 32b or 16b
 - For a C instruction, $PC \leftarrow PC + 2$ instead of $PC + 4$
 - And now normal instructions only need to be half-word aligned
- Results in ~30% smaller code
 - Roughly the same performance gain as **doubling** the icache!

The Instruction Encoding

Format	Meaning	15	14	13	12	11	10	9	8	7	6	5	4	3	2	1	0	
CR	Register	funct4			rd/rs1			rs2			op							
CI	Immediate	funct3	imm		rd/rs1			imm			op							
CSS	Stack-relative Store	funct3	imm			rs2			op									
CIW	Wide Immediate	funct3	imm			rd'			op									
CL	Load	funct3	imm		rs1'		imm		rd'		op							
CS	Store	funct3	imm		rs1'		imm		rs2'		op							
CB	Branch	funct3	offset		rs1'		offset			op								
CJ	Jump	funct3	jump target			op												

RVC Register Number
 Integer Register Number
 Integer Register ABI Name
 Floating-Point Register Number
 Floating-Point Register ABI Name

000	001	010	011	100	101	110	111
x8	x9	x10	x11	x12	x13	x14	x15
s0	s1	a0	a1	a2	a3	a4	a5
f8	f9	f10	f11	f12	f13	f14	f15
fs0	fs1	fa0	fa1	fa2	fa3	fa4	fa5

Stack Relative Load/Stores: Shrink preamble/postamble by nearly 50%

- Immediates are 0-extended
 - Because we write up from the stack pointer
- Immediates assume basic alignment (lower two bits 0 for words, three for doubles...)

	15	13	12	11		7	6		2	1	0
	funct3		imm		rd			imm		op	
	3		1		5			5		2	
C.LWSP		offset[5]			dest≠0			offset[4:2 7:6]		C2	
C.LDSP		offset[5]			dest≠0			offset[4:3 8:6]		C2	
C.LQSP		offset[5]			dest≠0			offset[4:9:6]		C2	
C.FLWSP		offset[5]			dest			offset[4:2 7:6]		C2	
C.FLDSP		offset[5]			dest			offset[4:3 8:6]		C2	

	15	13	12		7	6		2	1	0
	funct3			imm			rs2		op	
	3			6			5		2	
C.SWSP				offset[5:2 7:6]			src		C2	
C.SDSP				offset[5:3 8:6]			src		C2	
C.SQSP				offset[5:4 9:6]			src		C2	
C.FSWSP				offset[5:2 7:6]			src		C2	
C.FSDSP				offset[5:3 8:6]			src		C2	

Register Relative Loads & Stores

- Same zero-extending and alignment tricks

15	13 12	10 9	7 6	5 4	2 1	0
funct3	imm	rs1'	imm	rd'	op	
3	3	3	2	3	2	
C.LW	offset[5:3]	base	offset[2 6]	dest	C0	
C.LD	offset[5:3]	base	offset[7:6]	dest	C0	
C.LQ	offset[5 4 8]	base	offset[7:6]	dest	C0	
C.FLW	offset[5:3]	base	offset[2 6]	dest	C0	
C.FLD	offset[5:3]	base	offset[7:6]	dest	C0	

15	13 12	10 9	7 6	5 4	2 1	0
funct3	imm	rs1'	imm	rs2'	op	
3	3	3	2	3	2	
C.SW	offset[5:3]	base	offset[2 6]	src	C0	
C.SD	offset[5:3]	base	offset[7:6]	src	C0	
C.SQ	offset[5 4 8]	base	offset[7:6]	src	C0	
C.FSW	offset[5:3]	base	offset[2 6]	src	C0	
C.FSD	offset[5:3]	base	offset[7:6]	src	C0	

Jumps & Branches

15	13 12	imm	2 1	0
funct3				op
3		11		2
C.J		offset[11 4 9:8 10 6 7 3:1 5]		C1
C.JAL		offset[11 4 9:8 10 6 7 3:1 5]		C1

15	12 11	rs1	7 6	rs2	2 1	0
funct4					op	
4		5		5		2
C.JR		src ≠ 0		0		C2
C.JALR		src ≠ 0		0		C2

15	13 12	10 9	rs1'	imm	2 1	0
funct3		imm	rs1'	imm		op
3		3	3	5		2
C.BEQZ		offset[8 4:3]	src	offset[7:6 2:1 5]		C1
C.BNEZ		offset[8 4:3]	src	offset[7:6 2:1 5]		C1

And Then Assorted ALU instructions...

- Load 6 bit immediate values
 - Either to the lower 6 bits or bits 17-12
 - Useful for smaller immediates
- ADDI to self with small immediate
 - And a special form for a scaled by 16 immediate to the stack pointer
 - And another one to add a 0-extended immediate to the stack pointer to get addresses of stack-allocated variables
- Left shift more important than right shift:
 - Can left shift any register, but right shift only the encoded 8...
- Two register operations rather than 3 for the basic ALU ops:
 - Add, subtract, and, or, xor