

- N bits \rightarrow represent 2^N things
- doesn't always have to be value
- converting between bases, think buckets filled left \rightarrow right
- 4 bits = nibble; 8 bits = byte
- Signed Decimal Representation
 - sign and magnitude - never used, 2 0's
 - 2's complement
 - $(-x) = x + 1$ // to make negative
 - msb is $-(2^M)$
 - sign extension pretty easy to make more bits
 - 1 zero only, a little mean, widely used today
 - 1's complement
 - $-x = x$ // just flip bits to make negative
 - max value is half, ex. 0111 \rightarrow 1000
 - 2 0's, not used anymore
- Bias Encoding
 - # = bit rep (unsigned) + bias
 - ex...00 most negative, 11...11 most positive
- C
 - string - array of chars, null terminated '\0'
 - `malloc(strlen(s)+1)` // the +1 terminator
 - `int strlen(char* c)` // length of c, not '\0'
 - `int strcmp(char* s1, char* s2)` // 0 if same
 - `int char* strcpy(char* dst, char* src)` // src \rightarrow dst
 - returns pointer to dst string
 - pointers - function unsigned int basically, need type for size of
 - `int* a, b, c;` // make ptr, 2 vars
 - add integer, subtract, compare, compare to null, nothing else
 - `*ptr++ = *(ptr+1)`
 - C knows size of pointers, so can increment right amount
 - casting always fine
 - array - n length, index 0, 1, ..., n-1
 - C defines 1 element past end of array must be valid, to check for end of array, compare to that extra address
 - calculation with all pointers not assignable
 - but array name is read only pointer to start of array
 - struct - data structure made of smaller ones, need ; at end
 - can declare more at end of struct declaration if want
 - ; accessor
 - pointer to struct - use \rightarrow dereference operator
 - `p \rightarrow x == (*p).x;`
 - typedef - typedef <name> <alias>;
 - good for structs - typedef struct name { ... } alias;
 - good for linking pointers - typedef char* strings;
 - Bit Error - anomalies condition on bus
 - Segmentation Fault - access memory not allocated to it
 - `a = b` is assignment (`char* a`) vs. `a == b`
 - switch (var) { case <>: break; ... }
 - if no break; continues on
 - Storage Classes
 - auto, register - never used
 - extern - this default elsewhere, default to 0
 - declaration vs. definition
 - function must be declared before main()
 - `int x = 1;` // when outside func, loaded before program starts
 - `int x;` // when inside function, or definition, can have multiple
 - `extern int x;` // uses x which is defined elsewhere
 - static - default to 0
 - scope - where var available, external vs internal linkage
 - extent - how long var lasts
 - inside procedure - local scope, infinite extent
 - outside procedure - same file scope, infinite extent
 - function - \rightarrow
 - var initialized once, value persists between calls

- `int main(int argc, char** argv);`
 - can't do `char* argv[0]`
 - `argc[0]` is program name
 - little endian - smallest byte on right
 - big endian - big byte right
- #define <stuff> // preprocessor
- preprocessor replaces directly
- `getchar()` and `putchar()`
- `int putchar(c)` // extra character
- `!=` before `=`
- `int vs. +;`
- `if` before `||`, left to right
- `if, else if, else`
- local var - 'auto' be fixed by default
 - garbage initially
- extern
 - access by name by any func
 - defined once, declared everywhere
 - usually place at beginning of file
 - declaration - `static`, no storage
 - definition - `static`, storage allocated
 - by default functions have extern
 - `extern int x;` // considered a definition, also
- variable names - begin w/ letter
 - letter, - digits no "
- constants - letter dot L for long, 0 for unsigned, leading 0 for int
 - default double if decimal point, F for float, L for long double
- enum - allocate to itself
 - very similar to struct
 - first value is 0 unless otherwise stated
 - `enum names { ... } vars;`
- arithmetic operators = +, -, /, %, *
 - relational operators >, >=, <=, <
 - `==` and `!=` better in practice
 - all lower than arithmetic operators
- if, ||, short circuit
- binary operators - &, |, ^, <<, >>, ~
 - higher than arithmetic
- Control Flow
 - don't count indent
 - break; / next loop
 - continue; // back to top
 - `do { } while();`
 - `goto - write`
- static - limit scope to source file
 - can only init w/ constant literals
- array if use { }, sets extra to 0
- printf: %d - decimal; %u - unsigned
- void* - exist
- pointer to function
 - `void (*func)(int, int) = name;`
 - array don't know size & arr meaningless
 - handle - pointer to pointer address of a pointer
 - always free() or malloc() ?
 - size of() in bytes, know size of array
 - malloc() // need a NULL check
 - ternary operator - `eval? c1? : c2;`
 - const - don't change it
 - should cast malloc()

C Memory Management

- Aside
 - structure declaration don't allocate memory
 - var variable declaration allocate it
- `char* c = "hi";` // static
- `char* c = "hi";` // stack

- code - loaded when program starts, doesn't change
- static - var declared outside main, does not grow or shrink; global/static
 - global variables, permanent (compiler time contents only)
- heap - space requested for pointer, resize dynamically, grows upward
 - dynamic null storage, data live until deallocated
- stack - local vars, grows down, include main()
 - local vars
- Stack Frame - function calls LIFO, new return address from here
- stack pointer kills where top stack frame is

- Heap
 - malloc() calls - only one granular contiguous block, separate case special
 - want to avoid fragmentation
 - want minimal memory overhead
 - stored as circular-linked list - size of block, pointer to next block
 - choosing a block
 - best fit - smallest block big enough results in others
 - first fit - first block big enough, many small block at beginning
 - next fit - first fit but remember where finished searching
 - best fit - best fit at beginning, better than first-fit

IEEE Prefixes

- Ki - 2^{10}
- Mi - 2^{20} *increases by power of 10*
- Gi, Ti, Pi, Ei, Zi, Yi
- Gi, Ti, Pi, Ei, Zi, Yi
- SI
 - K, M, G, T, P, E, Z, Y
 - Ki, Mi, Gi, Ti, Pi, Ei, Zi, Yi
 - 1/2 bit / 0 bit trip value

Floating Point

- Fixed Point - set "binary" point, math: $a \cdot 10^b$; $a = mantissa$, $b = exponent$
- format: $1.x \dots x \cdot 2^{exp}$
- 32 bit float (fp32) - 1 bit sign, 8 bit exponent, 23 bit significand
 - $2 \cdot 10^{-38}$ to $2 \cdot 10^{38}$ // 0 is negative
 - 2 zeros tho
 - overflow - exponent too big
 - underflow - number too small
 - counting 0 \rightarrow max \rightarrow 0 \rightarrow -max
 - IEEE 754, base of -127 for exponent
 - $(-1)^s \cdot (1 + significand) \cdot 2^{exp-127}$
 - special cases
 - all zero = 0, right can make 0
 - 1 on highest exponent, significant all 0
 - NaN; highest exponent, not 0 significand
 - Denorm - 0 exponent, 0.x, # = $0.x \cdot 2^{denorm}$
 - all 1 < min
 - $2^{39} + 1$ is smallest number bigger can't represent
 - adding tiny num to big num just returns big number bc of precision & control

Matthew Tran

CS61C

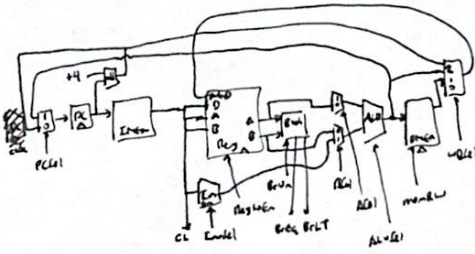
PS-1

c61c-a20

- RISC-V is an assembly language
 - RISC - let software do the rest
 - open source license free spec
 - 32 registers x 0...x31
 - 32 bit each, x0 == 0 always
 - Assembly
 - each line the instruction, # comment
 - memory address
 - 4 bytes = word
 - little endian - smallest bytes smallest addr
 - lw rd, imm(rs1) # rd = MEM[rs1+imm]
 - sw rs2, imm(rs1) # MEM[rs1+imm] = rs2
 - lb, sb equivalents, remember little endian, sign extended
 - lb - unsigned, 0 extend
 - bge rs1, rs2, label # goto label if rs1 >= rs2
 - jal rd, label # rd = PC+4, goto label & PC = PC+imm
 - jalr rd, rs1 # rd = PC+4, goto rs1 addr; PC = rs1+imm
 - all, srl, sra - sra sign extend, all have immediate versions
 - and, or, xor; not = xor w/ FF...
 - pseudo-instructions - j, li, la, mv, nop, not, ret
 - Call Convention
 - keep eye on who saves what, use Green Card
 - if calling function, make sure to store ra
 - use sp, addi sp, 4 to reserve 4 bytes
 - Instruction Formats
 - R - register-register ops
 - opcode - distinguishes operation, partially
 - funct3 + funct7 - fully specifies
 - I - register-immediate ops
 - imm is 12-bit signed, extended to 32 bit in hardware
 - 2048 in 2047
 - shift by imm ignores the upper bits
 - S - stores ~~store~~ aka SB
 - B - branches
 - imm - 4096 to 4094 as 2 byte offsets
 - 2 not 4 bit compressed instruction rd
 - drop the 0 bit by don't need it
 - U - upper-immediate stuff, for huge imm
 - imm (31:12), upper 20 bits
 - lui rd, imm
 - clear imm 12 bits, then imm
 - use addi to set lower 12 bits
 - bit addi does sign extension so need to prefill 1
 - j: handle all 4 bits
 - auipc - rd = PC + imm
 - good for storing PC
 - J - jumps, aka UJ
 - like B, drops 0 bit, 2¹⁴ 2 byte len, 2¹⁸ instr
 - jump to any where
 - auipc rd, chi 20 bits
 - jalr x0, xl, clow 12 bits # careful sign extension

- CALL
 - Compiler
 - inputs: high-level language
 - outputs: assembly language, may have pseudo-instructions
 - Assembler
 - inputs: assembly, for .s
 - outputs: object code, for .o, information tables (true assembly only)
 - Assembler Directives
 - .text - assembly
 - .data - put in .data segment (static)
 - .global sym - let sym access from other files
 - .string str - store str and null terminate
 - .word val... - store n 32-bit words in successive locs
 - Once no more pseudo-instructions - can find actual jump distance
 - "Forward reference" problem - need two passes to find everything
 - PIC - position independent code
 - what about static data? make table to tell linker
 - Symbol Table - items used by other files, aka labels, .data
 - Relocation Table - specifies file needs, the exact locations
 - ex. external labels, static data
 - Object File Format
 - object file header - specifies # of other pieces
 - text segment - machine code
 - data segment - binary rep of static data
 - relocation info - how that needs fixing
 - symbol table - list of labels and static data that can be ref.
 - debugging info
 - usually ELF file format

- Synchronous Digital Systems (SDS)
 - clock - we use rising edge
 - D-type flip flop
 - D is "data", Q is "output"
 - Q becomes D on rising edge
 - t_{setup} - time stable before clock pulse
 - t_{hold} - how long after clock pulse is stable
 - t_{clock} - how long for Q to become D
 - undefined if violate anything
 - max delay = t_{setup} + t_{clock} + longest CL
 - determine max clock speed
 - max hold = t_{clock} + shortest CL
 - determine register logic logic
- FSM (Finite State Machine)
 - finite # of states
 - use PS (present state) and input to determine next state (NS)
 - use truth table
 - usually inferring
- Boolean Algebra (+, *)
 - Complementarity - x + x' = 1, x * x' = 0
 - law of 0s and 1s - x + 0 = x, x + 1 = 1, x * 0 = 0, x * 1 = x
 - idempotent - x + x = x, x * x = x
 - idempotent - x + x = x, x * x = x
 - commutative - x + y = y + x, x * y = y * x
 - associative - (x + y) + z = x + (y + z), (x * y) * z = x * (y * z)
 - distributive - x(y + z) = xy + xz, x(y * z) = (x * y) * z
 - uniting theorem - xy + xz = x(y + z), x(y * z) = (x * y) * z
 - Demorgan's - (x + y)' = x' * y', (x * y)' = x' + y'
 - x + x' = 1, x * x' = 0



- Linker
 - inputs: object files, info tables
 - outputs: executable code
 - "link" - read .o files, allow separate compilation
 - steps
 - put all .text together
 - put all .data together, add to end of text
 - resolve references - use relocation table
 - fill in absolute addresses
 - 4 address types
 - PC-relative (bge, bne, jal, auipc, addi)
 - PIC, no need relocate
 - absolute function addr (auipc/jalr)
 - external function ref (auipc/jalr) } always relocate
 - static data ref (lui/addi)
 - base and store to static, relocate global pointer (got)
 - Relocation References
 - assume first word of text at 0x10000 for RV32
 - search in "user" symbol table
 - not found - search library files
 - fill in machine code appropriately
 - Static vs Dynamically Linked Library
 - OEL allow don't have to link new copy of compilation
 - less space, less dependencies, less memory, file level
 - but more time overhead to do link
 - Loader
 - load executable into memory and run, usually at 0x0
 - bunch of stuff
 - reserve enough memory

- Single Cycle Datapath and Control
 - ALU (Arithmetic Logic Unit)
 - adder
 - half adder - sum-xor, carry-and
 - full adder - sum-xor, carry-majority
 - overflow?
 - unsigned - highest C out
 - signed - C_n ^ C_{n-1}, out of highest bits
 - sbs - xor out put in carry in, basically 2 complement
 - Stages
 - IF - instruction fetch 200ps, Typ speed 1.25 GHz
 - IO - instruction decode 100ps, normal bc. flow use all 5 them
 - EX - execute (ALU) 200ps
 - MEM - memory access 200ps
 - WB - write back assume 200ps
 - RISC-V has 47 instructions, 37 enough for any C code
 - RTL (Register Transfer Level) - describe machinery as C. Verilog
 - Control Logic - use * when data care, simplify
 - ROM - popular, early programmable, fix errors
 - CL - use gates, count truth tables
 - BrEq - if equal
 - BrLT - less than
 - BrUN - branch unsigned
 - PCSel - ALU or PC+4
 - ImmSel - which immel format
 - ASel - Reg or PC
 - BSel - Reg or Imm
 - ALUSel - add/sub/shift
 - MemRW - write to MEM
 - RegWRen - write to register
 - WBSel - when to write
 - Formulae of Processor Performance
 - Time Program = Instructions Program / Instructions / Cycle = Cycles * Time Cycle
 - Energy Program = Instructions Program / Instructions * Energy Instruction
 - Energy "Iron Law" = Time Program * Cycles * Time Cycle
 - Time = Power (J/s) * Cycles * Time Cycle

Matthew Tran
CS61C
Pg. 2
cs61c-a20

Pipelining

- Add registers between stages
- shared clock period over many of stages
- can't get 5x speedup bc reasons
- Hazards
 - Structural
 - resource busy, needed in multiple stages
 - solution
 - 1) Stall
 - 2) add more hardware (always ok)
 - Dataflow
 - access hit in IO and cache
 - get's base register and write ports
 - double pumping - for some registers, need to write in 1st half, read in 2nd half
 - a pipeline = read/write
 - Memory
 - write-back IF and MEM stage
 - M1 has separate DPORT and I/MEM
 - actually just use 2 caches

Data

- processor hasn't finished writing yet
- Stalling - add NOP as compiler optimization
- forwarding - ex. the ALU, value of register
- get extra forwarding
- write-back to stall

Control

- caused by jump and branch
- could stall, but slow
- use branch prediction
- if wrong, flush the pipeline, insert nops

Other Performance Tweaks

- clock rate - limited by heat and power dissipation
- pipelining - good, but has costs, CPI > 1
- Superscalar (Pipeliner)
 - instruction-level parallelism
 - not all instructions, but separate pipelines
 - short multiple instructions per cycle
 - CPI < 1
 - "out of order" execution - reduce hazards

Caches

- DRAM slow, use SRAM which is cheap, smaller
- what to store?
 - Temporal locality - used it before, prob need it again
 - Spatial locality - prob need neighbors
- Direct Mapped Cache
 - each address mapped to one block in cache
 - block - unit of transfer between cache and memory
 - bigger block = better spatial locality = more locality
- TID - each address split
 - Tag: 1 bit response, EOC each block
 - Index: which mem if cache to check, which block
 - Offset: specifies byte within a block
- Access Read
 - processor checks cache
 - cache checks
 - hit: ✓ read data to processor
 - miss: ✗ get from memory, replace, send new value
 - valid bit = says if something there
- Cache Writing
 - write-through - update both, slow
 - write-back - dirty bit, write when read by the cache
- Cache Misses "Three C's"
 - Compulsory - can't have it never asked for it
 - Conflict - had before, but got replaced by no room
 - Capacity - miss because not stored

4th C

- coherence - race condition, pretty big and
- Cache Structure
 - Direct-mapped - one block per index (if data = 8 bits)
 - No-way set associative - block in any one of N places
 - Fully associative - whole compares to check if in cache
- Block Replacement Policy
 - for cache and associativity
 - LRU (least recently used) - kick least used out
 - FIFO - queue
 - Random - good for low temporal locality
- Average Memory Access Time (AMAT)
 - AMAT = hit time + miss rate * miss penalty
 - cache hit - hit time
 - multibank caches = increases hit

Virtual Memory

- give each process illusion of "main" memory
- provides protection/security
- all programs start at 0
- memory split into "pages"
- each process has own page table in memory
- Page Table Base Register (PTBR) - holds page table's address
- when OS runs the process force page table's data to disk, then dump in new page table
- invalid bit in TLB
- context switching
- Terms
 - Virtual Address (VA)
 - Virtual Memory (VM)
 - Virtual Page Number (VPN)
 - Physical Address (PA)
 - Physical Memory (PM)
 - Physical Page Number (PPN)
 - Split VA = VPN | offset
 - PA = PPN | offset
- Page Table Entry (PTE) - valid bit, access rights, PPN
- Translation Lookaside Buffer (TLB)
 - cache for the page table
 - usually highly associative (8B)
 - VPN split into tag and index (make sure fully address)
- Memory Read
 - check TLB
 - TLB Hit ✓ get PPN
 - TLB Miss - check page table (PT)
 - PT Hit - look up PTE in TLB, return a cache entry
 - PT Miss - Page Fault, fetch data to memory, update PTE, flush TLB
 - exception - create a new mapping
 - no entries - write new mapping
 - check cache w/ PA
 - partition fault - wrong access bit
 - Thrashing - constant disk memory swap

Parallelism

- Flynn's Taxonomy
 - SISD - single instruction, single data stream
 - SIMD - not used
 - MISD - not used
 - MIMD - single person multiple data
 - SMD - single person multiple data
 - Data Level Parallelism (DLP)
 - MMX, SSE, AVX
 - C Interleave - assembly in C
 - Multiprocessor Execution Model
 - each core runs its own instructions
 - separate and shared resources
 - Threads
 - list of instructions for one task
 - low PC, registers, access to shared memory
 - each physical core has N threads
 - OS manages them
 - multiplex software threads into hardware threads
 - switch out kernel threads (cache miss, interrupt, etc)
 - switch in when
 - resume software threads by interrupt mechanism and PC and registers
 - load w/ by load register & PC, jump to it
 - Multi-threading
 - multiprocessor runs really quickly
 - Hardware Threaded Multiprocessing
 - hyperthreading, both threads share full hardware SMT
 - 2 copies of PC and registers
 - 50% mechanical ~ 2x performance
 - Logical Thread
 - pretty much what can handle a software thread
 - OpenMP
 - C library handle parallelism
 - bit control of data dependencies
 - #pragma omp parallel for // how code with // line in each thread
 - omp_get_num_threads()
 - omp_get_thread_num()
 - #pragma omp critical // only one thread at a time
 - #pragma omp reduction (exp: +vars) // handle primitive and that
 - // can combine threads
 - MapReduce and Sparks - map function, then reduce, heavily parallelized
 - Relaxed Ordered Outlets (ROO)
 - map OP, flatMap (F), reduce by Key (F), reduce (F)
 - abstract OS level

Parallelism (cont.)

- only dependent on processor speed increase
- word level bus
- Locks
 - variable that says being used
 - issue is check value before I subtract
 - possible 2 processor access same thing
 - solution:
 - Hardware Synchronization
 - atomic read/write - 1 instr
 - must use shared memory
 - usually using register and memory
 - semaphore
 - RISC-V Atomic Memory Operations (AMO)
- Deadlock - everything depends on each other, cycle
- race? - measure elapsed time
- Multicore Multiprocessor
 - Shared Memory Symmetric Multiprocessor (SMP)
 - all arithmetic CPU taking memory
 - 2 identical cores
 - single shared memory (each with private own cache)
 - communicate via shared variables
 - Multiprocessor Cache
 - How long coherence?
 - if miss, write, kill others, invalid bit optic
 - MOESI - cache block status
 - Invalid - up to date, memory wrong, no cache copy
 - Owner - up to date, other has copy (in main cache)
 - Exclusive - up to date, no cache copy, reading up to date
 - Shared - up to date, other cache may have copy
 - Invalid - nothing
 - owner - want speed to let other read owner's data
 - exclusive - want to be modified to avoid write to memory in cache
 - False sharing - write block size, avoid reading of data that is modified, ex. F.x and F.y in same block
- Amdahl's Law
 - speedup = $\frac{1}{(1-p) + p/s}$
 - p: proportion that parallel
 - s: #x speedup for parallel
- Request Level Parallelism (RLP) - many independent requests
- Data Level Parallelism (DLP) - data in memory / disk
- Cloud Computing
 - Warehouse Scale Computers (WSC)
 - 10k to 100k copies
 - 41 hour/yr 49,999\$
 - cost of ownership >> cost of operation equipment
 - need high efficiency Power Usage Efficiency (PUE) = $\frac{\text{total power}}{\text{useful power}}$
 - cooling log issue - air flow, water, water, climate
 - also need monitor high CPU usage
 - I/O
 - process need to input output bytes to internet
 - can use special I/O instructions and hardware
 - on-chip Memory Mapped I/O (MMIO) - very common
 - control register, data register
 - Billing
 - constantly ask you for more, but for HDD
 - Interrupts
 - device & CPU try handle PC section Machine Exception PC (MEPC)
 - store registers, handle code, cache block related
 - Interrupt - external, asynchronous
 - Exception - caused by program
 - Trap - act as "software" an interrupt or exception
 - prevent traps, aware every multi-processor to trap done, more after done
 - don't need to know pipeline just jump back
 - handled like a pipeline hazard
 - 1) compute instr before trap
 - 2) flush cache, write-back
 - 3) optionally store exception in cache register
 - 4) handle exception to trap handler
 - Golden Age of Computer Architecture
 - software advances can impact architecture innovation
 - reusing hardware interface opportunities for innovation
 - multiplexed cache, cache coherence, cache
 - control vs. datapath
 - microprocessors - logic unit to design control
 - Demand calling - grow further from a processor, more control
 - Security challenges
 - Domain Specific Architecture (DSA) - do the job well
 - Tensor Processing Unit (TPU) - makes unit - multiply accumulate 6570 MAC 750 MAC 868 DDG 440 multiply 3200
 - domain specific, > 25x MAC = GPU, 250x a CPU
 - small size, memory all off board, no cache
 - GPU
 - lots of dataflow architectures
 - many cores, TLB, software coherence, low freq, high latency
 - many registers, limited bus bandwidth, in-order, many execution

- Extra
 - even parity - make bits even
 - >> not always work as flow control but
 - n-bit machine usually n-bit VM
 - to point memory, access edge offset