Powers of 2

| 10 | 9 | 8 | 7 | 6 | 5 | 4 | 3 | 2 | 1 | 0 | -1 | -2 | -3 | -4 | -5 | -6 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1024 | 512 | 256 | 128 | 64 | 32 | 16 | 8 | 4 | 2 | 1 | 0.5 | 0.25 | 0.125 | 0.0625 | 0.3125 | 0.015625 |

von Neumann Architecture
- stores program & data in memory

$b \wedge b = 0$
$b \wedge 0 = b$
$a \wedge (b \wedge b) = a \wedge 0 = a$  } for swapping int

## C programming

int : 4 bytes
float : 4 bytes   } depends on CPU.
double : 4 bytes

char : always 1 byte

8 bits = 1 byte

(int) float → truncates

$32 \Rightarrow (32)_{10}$
$032 \Rightarrow (32)_8$   ← MIPS no octal input
$0 \times 32 \Rightarrow (32)_{16}$

Variable has:
 - name
 - data type
 - value
 - address

Everything in C is
pass-by-value not reference.

↙ address
scanf ("string", $x_1$, , $x_2$, ...)
printf ("string")
printf ("string", $x_1$, $x_2$, ...)
puts (str)  // w/ newline

%c  char
%d  int
%f  float  (double uses %f for output)
%lf  double  (for input)
%p  address ← different execution, different address
%zu  size_t

## Base-R to decimal

$1101.101_2 = 1 \times 2^3 + 1 \times 2^2 + 0 \times 2^1 + 1 \times 2^0 + 1 \times 2^{-1} + 0 \times 2^{-2} + 1 \times 2^{-3}$
   3 2 1 0 -1 -2 -3

## Decimal to Base-R

Whole numbers : repeated division by R
 → division ∵ $R^k \Rightarrow R^{k-1}$

| 2 | 43 |
| 2 | 21 rem 1 ← LSB |
| 2 | 10 rem 1 |
| 2 | 5 rem 0 |
| 2 | 2 rem 1 |
| 2 | 1 rem 0 |
|   | 0 rem 1 ← MSB |

$43 = (101011)_2$

↙ repeat til 0

Fractional portion : repeated multiplication
 → multiplication ∵ $R^{1-k} \leftarrow R^{-k}$

take this ↓
0.3125 × 2 = 0.625  ← MSB
0.625 × 2 = 1.25
0.25 × 2 = 0.5
0.5 × 2 = 1  ← LSB
repeat til 0
( might not terminate
  but will have pattern )

$0.3125 = (0.0101)_2$

$R1 \Rightarrow R2 \equiv R1 \Rightarrow decimal \Rightarrow R2$

but if $R1^n = R2$
 → partition into groups of n

$(10 \ 111 \ 011 \ 001 . 101 \ 110)_2 = (2731.56)_8$
  2   7   3   1    5   6

n ←.→ n

## Fixed-Point Representation

- fixed number of bits for integer & fractional components

## Sign and Magnitude

□ [ magnitude ]
sign

$0 \to$ +ve
$1 \to$ -ve

range: $[-(2^{n-1}-1), 2^{n-1}-1]$

negation: flip sign bit

## 1's complement : diminished radix complement

$-X = 2^n - X - 1$ $\qquad X = r^n - r^{-m} - X$
$\qquad\qquad\qquad\qquad\qquad\qquad\quad$ fractional

MSB is sign bit

range: $[-(2^{n-1}-1), 2^{n-1}-1]$

negation: flip <u>all</u> bits $\qquad$ e.g. FFFF - 0001 = FFFE

## 2's complement : radix complement $\quad$ e.g. 10000 - 0001 = FFFF

$-X = 2^n - X$

MSB is sign bit $\qquad$ no duplicate 0

range: $[-2^{n-1}, 2^{n-1}-1]$

negation: flip <u>all</u> bits, +1 to LSB (incl. fractional)

$$A - B = A + (-B)$$

1. Binary addition on $a + b$
2. If MSB carry out, add 1 to result
3. Check overflow.

$\longrightarrow$

Overflow if MSB 1+1 = 10

```
  1 1 0 1
+ 1 0 1 0
---------
1 0 1 1 1
```

1. Binary addition on $a + b$
2. Ignore MSB carry out
3. Check overflow.

$\longrightarrow$

## Floats   IEEE 754

[ Sign ] [ exponent ] [ mantissa ]

← normalised w/ implicit leading bit & pad w/ 0s to the right

defn: fractional component

**Single-precision** (32-bits)
— 1-bit
— 8-bit $\to$ excess-127
— 23-bit

**double-precision**
- 1-bit
- 11-bit $\to$ excess-1023
- 52-bit

all 0s is a "denormal form" that represents 0.
$\to$ cannot be represented precisely otherwise

1. Convert $|f|$ to binary
   $\to$ repeated mul/div
2. Normalise to $\pm 1.$mantissa $\times 2^{exponent}$
3. Sign: 0/1 for +ve/-ve
4. Merge

Multiply : 1. Multiply mantissa
2. Add exponents
3. Normalise

Addition : 1. Equalise exponents
2. Add new mantissa
3. Normalise

## Excess - N  on M bits

start at -N

$(x)_{binary} + N = (x)_{Excess-N}$

$\to$ easy to compare values

| Excess -4 | Value |
|-----------|-------|
| 000 | -4 |
| 001 | -3 |
| 010 | -2 |
| 011 | -1 |
| 100 | 0 |
| 101 | 1 |
| 110 | 2 |
| 111 | 3 |

## Pointers → ++ and -- based on sizeof(type)

```
int a = 123;
&a ⇒ addr of a                    void swap (int *a, int *b)
int *a_ptr = &a;                            ↖ ↘ pointers

    *a_ptr = 123 = a
```

## Arrays : homogeneous collection of data
— contiguous in memory

values.length > size ⇒ warning

```
int x[4] = {1};  →  x[0] = 1
                    x[1] = 0
                    :          } default
x = &x[0]           x[3] = 0     value

arr[i] = *(arr + i)

arr2 = arr1;  ✗ CANNOT do assignment
```

## Structures
— can be nested
— allows heterogeneous members
— allows assignment (unlike arrays)



```
typedef struct {              void func (box_t box) {...}
    int a,b;                  func (&box);
    float f;                   ↗
} box_t;                      use address to use "reference".

box_t box = {0,1,0.3};
box.f; → 0.3                  (*box_ptr).b  ≡  box_ptr → b
    ↑                                              ↑
    dot                                          arrow
```

## Strings
size = string length + 1 ↖ '\0'

ends with '\0' → ascii of 0.

```
char str[] = "abc";  ≡  char str[] = {'a','b','c','\0'}

            ↙ stdin
fgets (str, size, file)                    scanf ("%s", str)
    reads up to size OR newline                — until whitespace
    ↘  | e | a | t | \n | \0 |
                        ↑
              should change to  \0
```

## < string.h >
1. strlen (s)
    → no. of char

2. strcmp (s1, s2)
    → s1 lexicographical comparison

3. strncmp (s1, s2, n)
    → up to \0 or n

4. strcpy (s1, s2)
5. strncpy (s1, s2, n)
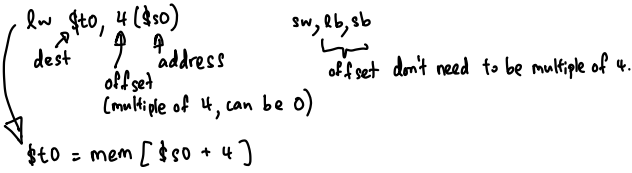
## MIPS

add \$s0, \$t1, \$t2  ≡  \$s0 = \$t1 + \$t2

immediate: $[-2^{15}, 2^{15}-1]$  ← 16-bit 2's complement
 ↳ zero-extension

for large constants (> 16 bits)
 lui \$t0, 0xAAAA
 ori \$t0, \$t0, 0xBBBB   # \$t0 = 0xAAAABBBB

lw \$t0, 4(\$s0)         sw, lb, sb
 dest ↗ ↑ ↖ address    ↳ offset don't need to be multiple of 4.
   offset
   (multiple of 4, can be 0)
\$t0 = mem[\$s0 + 4]

beq
bne }  \$r1, \$r2, Label    ← relative

j Label                   ← absolute

slt \$t0, \$s0, \$s1   if \$s0 < \$s1,
                         \$t0 = 1
                      else
                         \$t0 = 0

Amdahl's Law
→ make frequently used instructions FAST!

### Instruction Set Architecture (ISA)

① Complex instruction set computer (CISC)  ← x86_32 (IA32)

② Reduced instruction set computer (RISC)
     ↑ MIPS, ARM

### Storage Architectures

① Stack
② Accumulator
   - one operand implicitly on accumulator
③ General purpose
   - register-memory
   - register-register (MIPS)
④ Memory-memory
   - go through BUS (slow)
   - small instruction set

### Instruction Length

Variable vs Fixed vs Hybrid (≈ variable)
 → OPCODE (must have!)
 → operands
     → usually $2^n$ sizes
Work w/ most constrained instruction type.
            ↑
           most operands

### R-format

| OPCODE = 0 | rs | rt | rd (dest) | shamt (shift) | funct |
|---|---|---|---|---|---|
| 6 | 5 | 5 | 5 | 5 | 6 |

add \$rd, \$rs, \$rt
sll/srl \$rd, \$rt, shamt   ← \$rs = 0

### I-format

sign-extended except andi, ori, xori

| OPCODE | rs | rt (dest) | IMMEDIATE |
|---|---|---|---|
| 6 | 5 | 5 | 16 |

beq/bne \$rs, \$rt, immd   → up to $\pm 2^{17}$ bytes
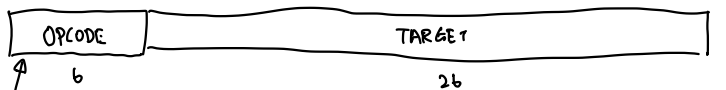addi \$rt, \$rs, immd       if jump: PC = PC + 4 + (immd × 4)
lw \$rt, immd(\$rs)
                            ① beq ..., end ↙②
                            ②            ℓ₀
                            ③            ℓ₁
                            ④ end: hi    ℓ₂

### J-format

| OPCODE | TARGET |
|---|---|
| 6 | 26 |

$(2)_{10} = (000010)_2$

256MB boundary
$2^{26} - 1$ - (instructions between)

if at 0xAFFFFFFE boundary,
cannot use j.

| 4 MSB of next PC | + target + | 00 |
|---|---|---|

### Endians

**Big Endian**
→ MSB in lowest address

0x DE AD BE EF

| 0 | DE |
| 1 | AD |
| 2 | BE |
| 3 | EF |

**Little Endian**
→ LSB in lowest address
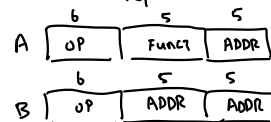
0x DE AD BE EF
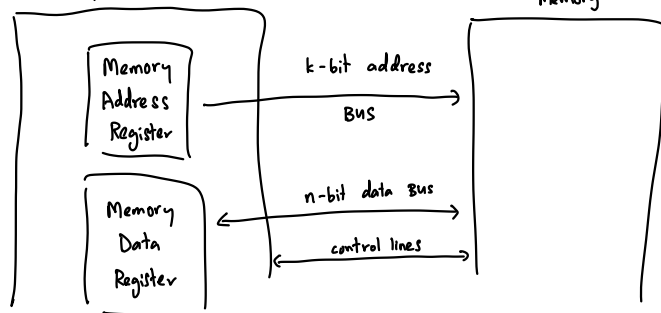
| 0 | EF |
| 1 | BE |
| 2 | AD |
| 3 | DE |

for $a_1 < a_2 < \dots < a_n$,
$$\min: 2^{a_1} + \sum_{k=1}^{n-1} 2^{(a_{k+1} - a_k)} - n + 1$$
$$\max: 2^{a_n} - \sum_{k=1}^{n-1} 2^{(a_n - a_k)} + n - 1$$

largest opcode
Max: $2^M - 2^A + 1$
         ↳ repeat

| A | OP (6) | FUNCT (5) | ADDR (5) |
|---|---|---|---|

| B | OP (6) | ADDR (5) | ADDR (5) |
|---|---|---|---|

$2^{11} - 2^5 + 1$
= 2017 //

Processor — Memory

| Memory Address Register | → k-bit address BUS → |
|---|---|

| Memory Data Register | ← n-bit data Bus → |
|---|---|
|  | ← control lines → |

\* up to $2^k$ addressable locations

## Control & Datapath

1. Fetch
   - get instruction from memory
   - get address from PC register
2. Decode
   - find which operation
3. Operand fetch
   - get operands needed for operation
4. Execute
5. Writeback
   - store results in register

PC read during first half of clock period
& updated at rising clock edge

### RegDst

0: Write register = inst $[20:16]$ ← I-type
1: Write register = inst $[15:11]$ ← R-type

### RegWrite

0: no register write

1: write to register

### ALUSrc

0: Operand 2 = RD2
1: Operand 2 = SignExt( inst $[15:0]$) ← lw/sw

### Mem Read

0: no memory read
1: memory read using Address (ALU result)

### MemWrite

0: no memory write
1: mem $[Address]$ ← write data (RD2)

### MemToReg

1: WD = memory read data

0: WD = ALU result

### PCSrc / Branch   ← PCSrc = Branch · isZero

0: next PC = PC + 4

1 & isZero: PC + 4 + SignExt( inst $[15:0]$) << 2

### ALU Control

| | |
|---|---|
| 0000 | AND |
| 0001 | OR |
| 0010 | ADD |
| 0110 | SUB → A − B = A + (−B) |
| 0111 | SLT  = A + ~B + 1 ← $C_{in}$ |
| 1100 | NOR |

↑ ↖
Ainvert  Binvert



A
B

Ainvert
Binvert

result

00
01
10
11 is undefined

| | RegDst | ALUSrc | MemToReg | RegWrite | MemRead | MemWrite | Branch | ALUOp |
|---|---|---|---|---|---|---|---|---|
| R-type | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 10 |
| lw | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 00 |
| sw | X | 1 | X | 0 | 0 | 1 | 0 | 00 |
| beq | X | 0 | X | 0 | 0 | 0 | 1 | 01 |

# Boolean Algebra

## Identity

$$A + 0 = 0 + A = A \qquad A \cdot 1 = 1 \cdot A = A$$

## Inverse / Complement

$$A + A' = A' + A = 1 \qquad A \cdot A' = A' \cdot A = 0$$

## Commutative

$$A + B = B + A \qquad A \cdot B = B \cdot A$$

## Associative

$$A + (B + C) = (A + B) + C \qquad A \cdot (B \cdot C) = (A \cdot B) \cdot C$$

## Distributive

$$A \cdot (B + C) = (A \cdot B) + (A \cdot C) \qquad A + (B \cdot C) = (A + B) \cdot (A + C)$$

## Idempotency

$$X + X = X \qquad X \cdot X = X$$

## One element / Zero element

$$X + 1 = 1 + X = 1 \qquad X \cdot 0 = 0 \cdot X = 0$$

## Involution

$$(X')' = X$$

## Absorption 1

$$X + X \cdot Y = X \qquad X \cdot (X + Y) = X$$

## Absorption 2

$$X + X' \cdot Y = X + Y \qquad X \cdot (X' + Y) = X \cdot Y$$

## De Morgan → generalisable to >2 variables

$$(X + Y)' = X' \cdot Y' \qquad (X \cdot Y)' = X' + Y'$$

## Consensus

$$X \cdot Y + X' \cdot Z + Y \cdot Z = X \cdot Y + X' \cdot Z \qquad (X + Y) \cdot (X' + Z) \cdot (Y + Z) = (X + Y) \cdot (X' + Z)$$

## Duality

$$X + 1 = 1 \iff X \cdot 0 = 0$$

inverts

| + and · |
| 0 and 1 |

---

must have all n literals

Minterm of n variables : $X_1 \cdot X_2 \cdot \cdots \cdot X_n$
Maxterm of n variables : $X_1 + X_2 + \cdots + X_n$

up to $2^n$ possibilities

$$m_X = (M_X)'$$

don't compress

| x | y | Minterms | | Maxterms | |
|---|---|----------|----|----------|----|
| 0 | 0 | $x' \cdot y'$ | m0 | $x + y$ | M0 |
| 0 | 1 | $x' \cdot y$ | m1 | $x + y'$ | M1 |
| 1 | 0 | $x \cdot y'$ | m2 | $x' + y$ | M2 |
| 1 | 1 | $x \cdot y$ | m3 | $x' + y'$ | M3 |

sum of minterms = canonical sum of products (SOP)
product of maxterms = canonical product of sums (POS)

$$F = m1 + m2 = \Sigma m(1,2) = \Sigma m(1-2)$$
$$G = M0 + M4 = \Pi M(0,4)$$
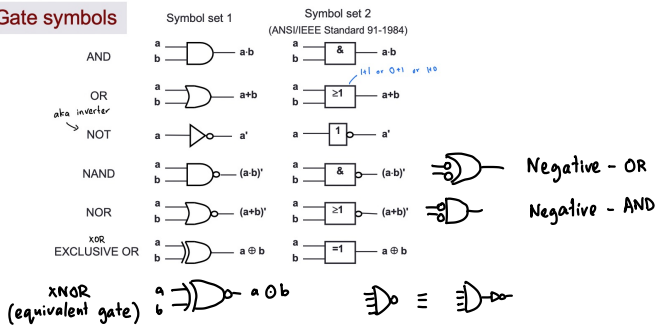
$\Sigma m$ and $\Pi M$ can convert with De Morgan's

# Logic Circuits

## Gate symbols

| A | B | AND | OR | XOR | NAND | NOR | XNOR |
|---|---|-----|----|-----|------|-----|------|
| 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 |
| 1 | 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 |

| | Symbol set 1 | Symbol set 2 (ANSI/IEEE Standard 91-1984) | |
|---|---|---|---|
| AND | a·b | & a·b | |
| OR (aka inverter) → NOT | a+b | ≥1 a+b | |
| NOT | a' | 1 a' | |
| NAND | (a·b)' | & (a·b)' | Negative - OR |
| NOR | (a+b)' | ≥1 (a+b)' | Negative - AND |
| XOR EXCLUSIVE OR | a ⊕ b | =1 a ⊕ b | |
| XNOR (equivalent gate) | a ⊙ b | | |

fan – in : number of inputs of a gate

* assume complemented literals
  don't exist unless stated for CS2100

## Universal Gates

Complete sets of logic: {AND, OR, NOT}, {NAND}, {NOR}, {AND, NOT}, {OR, NOT}, ... (a lot)

Others still useful
- XOR : parity bit generation
- economical (save on gates)

### NAND



$X \to x'$

$x, y \to (x \cdot y)' \to x \cdot y$

$x \to x'$, $y \to y'$, $\to (x' \cdot y')' = x+y$

### NOR



$x \to x'$

$x \to x'$, $y \to y'$, $\to (x' + y')' = x \cdot y$

$x, y \to (x+y)' \to x+y$

SOP can easily be implemented with:
1) 2-level AND-OR circuit
2) 2-level NAND circuit
   ↑ invertor not considered

POS can easily be implemented with:
1) 2-level OR-AND circuit
2) 2-level NOR circuit
   ↑ invertor not considered

e.g. $F = A \cdot B + C' \cdot D + E$

AND-OR to NAND



## Programming Language Array (PLA)



input ≡ AND

OR ≡ output

— may not be able to implement every mapping

## Read Only Memory (ROM)

—Similar to PLA
— fully decoded : able to implement any mappings

## Simplification
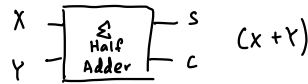
- minimise literals
- minimise terms used

## Gray Code — aka reflected binary code

binary $\overset{\nearrow}{\rightarrow}$

- unweighted (opposite of arithmetic code)
- only 1-bit change $\left(max \overset{loops}{\underset{back}{\Rightarrow}} 0\right)$
- good for error detection
- n bits $\Rightarrow 2^n$ values

| Decimal | Binary | Gray Code | Decimal | Binary | Gray code |
|---------|--------|-----------|---------|--------|-----------|
| 0 | 0000 | 0000 | 8 | 1000 | 1100 |
| 1 | 0001 | 0001 | 9 | 1001 | 1101 |
| 2 | 0010 | 0011 | 10 | 1010 | 1111 |
| 3 | 0011 | 0010 | 11 | 1011 | 1110 |
| 4 | 0100 | 0110 | 12 | 1100 | 1010 |
| 5 | 0101 | 0111 | 13 | 1101 | 1011 |
| 6 | 0110 | 0101 | 14 | 1110 | 1001 |
| 7 | 0111 | 0100 | 15 | 1111 | 1000 |

standard 4-bit gray code
$\nearrow$
many possible implementations

## Half-Adder



$(x + y)$

| X | Y | C | S |
|---|---|---|---|
| 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 |
| 1 | 0 | 0 | 1 |
| 1 | 1 | 1 | 0 |

$C = X \cdot Y = m3$
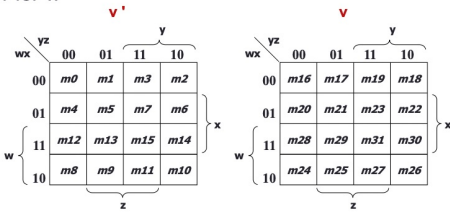$S = X' \cdot Y + X \cdot Y = m1 + m2$
$\quad = X \oplus Y$



## Karnaugh - Maps

- each square represents a minterm
- 2 adjacent squares $\Rightarrow$ differ by one literal
- wrap around exists
- n - variable $\Rightarrow$ n neighbours

- $2^n$ cells
- for $2^x$ adjacent minterms, n - x literals



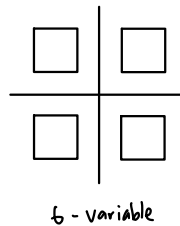2 - variable

$0 = b'$
$1 = b$

3 - variable

4 - variable

▪ Organised as two 4-variable K-maps. One for v' and the other for v.



Corresponding squares of each map are adjacent.
Can visualise this as *one 4-variable K-map* being on TOP *of the other 4-variable K-map.*

5 - variable

6 - variable

## Groupings in K-Maps

based on $A + A' = 1$ (Unifying Theorem)

1. group as many as possible
   $\rightarrow$ PI
   $\rightarrow$ results in less terms

2. ignore redundant groups
   $\rightarrow$ for EPI

✳ remember wrap arounds

implicant: product term that can be used to cover minterms of the function

EPI: PI that includes $\geqslant 1$ minterm __not covered__ by any other PI.

$\rightarrow$ does NOT need to cover all 1s

## Don't Care

- denote with X or d
- use $\Sigma d$ or $\Pi D$

$\rightarrow$ include in PI/EPI only if the PI/EPI contains 1.

finding SOP
$\rightarrow$ all EPIs + remaining PIs not covered by any EPIs
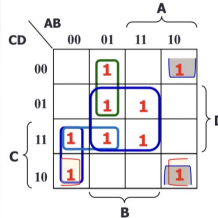
## Quine - McCluskey

- similar to K-map but not limited to ~6 variables

Steps:

1. list out all minterms in groups with same number of 1s

2. Combine codes that differ by 1 bit into bigger group, write combined with " - "

3. Repeat step 2 and continue combining

Identify EPIs (columns containing single tick)

| | 2 | 3 | 4 | 5 | 7 | 8 | 10 | 13 | 15 |
|---|---|---|---|---|---|---|---|---|---|
| 2,3: 001- (A'.B'.C) | ✓ | ✓ | | | | | | | |
| 2,10: -101 (B.C'.D) | ✓ | | | | | | ✓ | | |
| EPI 4,5: 010- (A'.B.C') | | | ✓ | ✓ | | | | | |
| EPI 8,10: 10-0 (A.B'.D') | | | | | | ✓ | ✓ | | |
| 3,7: 0-11 (A'.C.D) | | ✓ | | | ✓ | | | | |
| EPI 5,7,13,15: -1-1 (B.D) | | | | ✓ | ✓ | | | ✓ | ✓ |

Eliminate EPI rows

Answer = EPIs + remaining ticks

## Combinational Circuit

Each output depends entirely on the immediate inputs
→ like mathematical functions
→ no internal state

To analyse:

1. Label inputs/outputs

2. Obtain functions of intermediate points & outputs ← after gates/blocks

3. Draw truth table $^{xx}$

4. Deduce/guess functionality
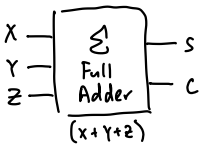
## Gate-level Design

- with logic gates
  OR, AND, NOT, XOR,...

- analyse w/ K-map or truth table to get SOP (& simplify) to build circuit

## Block-level Design

- with functional blocks
- using simpler blocks to build complex blocks
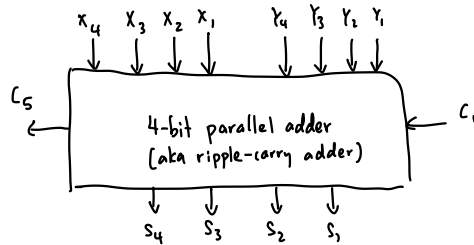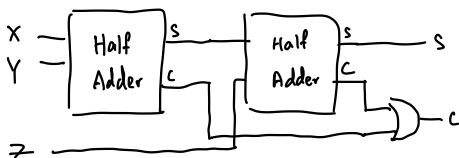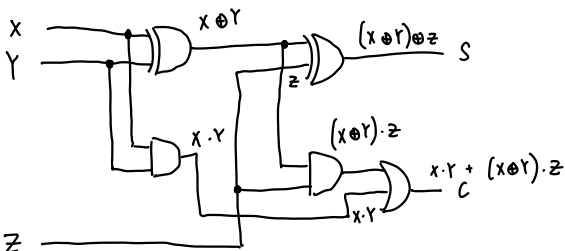
- find common patterns & intermediate states

- Reduce cost
- Increase speed
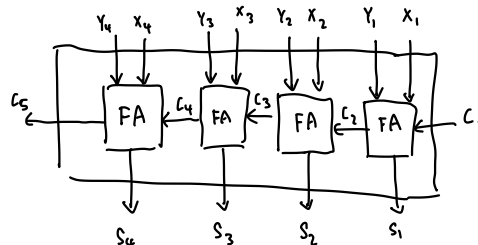- Design simplicity (re-use if possible)

## Common Blocks

$$S = X \oplus (Y \oplus Z)$$
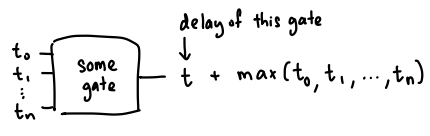$$C = X \cdot Y + (X \oplus Y) \cdot Z$$

gate-level too complicated to analyse
→ $2^9$ possible inputs

→ can build 8, 16, ... bit adders

## Circuit Delays



$$t + \max(t_0, t_1, \ldots, t_n)$$

delay of this gate

→ have to wait for all inputs to be stable

## MSI Components

← medium-scale integrations

- An integrated chip (IC) is a set of electronics on one small flat piece/chip of semiconductor material.
- Scale of integration: # of components fitted on a standard IC chip

## n:m Decoder   or n×m or n-to-m

n: # of input bits

m: # of possible outputs $\left[ m \leq 2^n \right]$ ← usually, $m = 2^n$

→ input order matters (not always but sometimes)

### 2:4 decoder

| E | X | Y | $F_0$ | $F_1$ | $F_2$ | $F_3$ |
|---|---|---|-------|-------|-------|-------|
| 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 1 | 0 | 0 |
| 1 | 1 | 0 | 0 | 0 | 1 | 0 |
| 1 | 1 | 1 | 0 | 0 | 0 | 1 |
| 0 | d | d | 0 | 0 | 0 | 0 |

↑
1-enabled

minterms
∴ can use to implement SOP functions with OR gates



← only 1 high/low

$F_0 \quad x' \cdot Y'$
$F_1 \quad x' \cdot Y$
$F_2 \quad x \cdot Y'$
$F_3 \quad x \cdot Y$

AND gates only (+ 1 inverter)

0-enabled usually denoted with $\bar{E}$ or $E'$

→ more commonly used in MSI decoders

→ truth table just swap 0 & 1

active-high (normal output)
  → 1 if true
active-low (negated output)
  → 0 if true

can build $(n+1):2^{n+1}$ decoder with 2 $n:2^n$ decoders



use MSB
∴ easier

74138 (3:8 decoder)

→ very common

→ G1 = 1, G2A = 0, G2B = 0
  to enable

→ active-low output

## m:n Encoders

→ converse of decoders

→ only 1 input bit is high

(or in the case of priority encoders, highest priority is chosen)

m: # of input bits $\left[ m \leq 2^n \right]$

n: # of outputs

(all 0 is invalid!)

### 8:3 encoder



$x = D_4 + D_5 + D_6 + D_7$

$y = D_2 + D_3 + D_6 + D_7$

$z = D_1 + D_3 + D_5 + D_7$

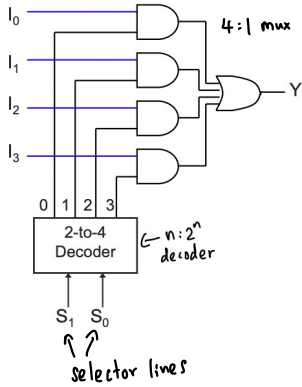use truth table / binary

$1 : 2^n$

**Demultiplexer** or demux

→ take input & a set of selection lines

→ directs input to one of the output lines

→ identical circuit as decoder with enable

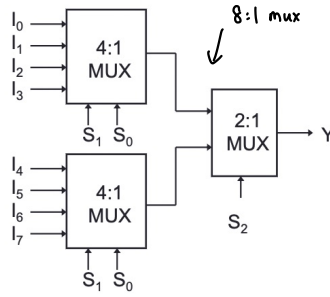$2^n : 1$

**Multiplexer** or mux or data selector

→ many inputs & a set of selection lines

→ chooses one of the inputs as output

output $= I_0 \cdot m_0 + I_1 \cdot m_1 + \cdots + I_{n-1} \cdot m_{n-1}$

→ use $n : 2^n$ decoder & AND with the input bits (as above)



4:1 mux

2-to-4 Decoder ← $n : 2^n$ decoder

$S_1$ $S_0$

selector lines

**Merging mux**



8:1 mux

→ group by common term in truth table

| $S_2$ | $S_1$ | $S_0$ | Y |
|---|---|---|---|
| 0 | 0 | 0 | $I_0$ |
| 0 | 0 | 1 | $I_1$ |
| 0 | 1 | 0 | $I_2$ |
| 0 | 1 | 1 | $I_3$ |
| 1 | 0 | 0 | $I_4$ |
| 1 | 0 | 1 | $I_5$ |
| 1 | 1 | 0 | $I_6$ |
| 1 | 1 | 1 | $I_7$ |

to implement SOP function { put 1 if it is a minterm of the function else 0
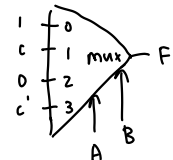
& variables = select

BUT we can do BETTER!!

→ use $2^{n-1} : 1$ mux is sufficient

e.g. $F(A,B,C) = \Sigma m(0,1,3,6)$

| A | B | C | F | mux |
|---|---|---|---|---|
| 0 | 0 | 0 | 1 | 1 |
| 0 | 0 | 1 | 1 | |
| 0 | 1 | 0 | 0 | C |
| 0 | 1 | 1 | 1 | |
| 1 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | |
| 1 | 1 | 0 | 1 | C' |
| 1 | 1 | 1 | 0 | |

group by AB ∴ using C

can use any



74151A 8-to-1 mux

→ 0-enabled (strobe = enable)

→ both active-high & active-low

| | A B C D | AB | AC | AD | BC | BD | CD |
|---|---|---|---|---|---|---|---|
| 0 | 0000 | 00 | 00 | 00 | 00 | 00 | 00 |
| 1 | 0001 | 00 | 00 | 01 | 00 | 01 | 01 |
| 2 | 0010 | 00 | 01 | 00 | 01 | 00 | 10 |
| 3 | 0011 | 00 | 01 | 01 | 01 | 01 | 11 |
| 4 | 0100 | 01 | 00 | 00 | 10 | 10 | 00 |
| 5 | 0101 | 01 | 00 | 01 | 10 | 11 | 01 |
| 6 | 0110 | 01 | 01 | 00 | 11 | 10 | 10 |
| 7 | 0111 | 01 | 01 | 01 | 11 | 11 | 11 |
| 8 | 1000 | 10 | 10 | 10 | 00 | 00 | 00 |
| 9 | 1001 | 10 | 10 | 11 | 00 | 01 | 01 |
| 10 | 1010 | 10 | 11 | 10 | 01 | 00 | 10 |
| 11 | 1011 | 10 | 11 | 11 | 01 | 01 | 11 |
| 12 | 1100 | 11 | 10 | 10 | 10 | 10 | 00 |
| 13 | 1101 | 11 | 10 | 11 | 10 | 11 | 01 |
| 14 | 1110 | 11 | 11 | 10 | 11 | 10 | 10 |
| 15 | 1111 | 11 | 11 | 11 | 11 | 11 | 11 |

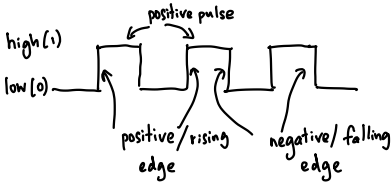## Sequential Circuit

→ has internal state/memory

memory element: device that remembers a value indefinitely, or change value on command from inputs

↖ 1-bit data

$Q(t)$ or $Q$: current state
$Q(t+1)$ or $Q^+$: next state

## Clock 🕐 − square wave



high (1)
low (0)

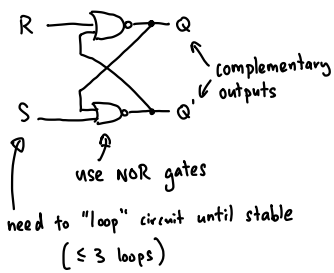positive pulse

positive/rising edge    negative/falling edge

pulse-triggered
− latches

ON = 1, OFF = 0

edge-triggered
− flip-flops
− positive edge-triggered: ON = rising edge, OFF = other time
− negative edge-triggered: ON = falling edge, OFF = other time

## S-R Latch   or R-S

(active-high)



R — Q
S — Q'

complementary outputs

use NOR gates

need to "loop" circuit until stable
(≤ 3 loops)

**Characteristic Table**

| S | R | Q | Q' | |
|---|---|---|----|---|
| 0 | 0 | NC | NC | |
| 0 | 1 | 0 | 1 | Latch SET |
| 1 | 0 | 1 | 0 | Latch RESET |
| 1 | 1 | 0 | 0 | Invalid input |

| S | R | $Q^+$ | |
|---|---|-------|---|
| 0 | 0 | Q | |
| 0 | 1 | 0 | |
| 1 | 0 | 1 | |
| 1 | 1 | N/A | ← breaks circuit |

→ $Q^+ = S + R' \cdot Q$

active-low variant
− use NAND or negative-or instead of NOR

− sometimes, S and R are labelled S' and R'

| S | R | $Q^+$ |
|---|---|-------|
| 0 | 0 | N/A |
| 0 | 1 | 1 |
| 1 | 0 | 0 |
| 1 | 1 | Q |

## Gated S-R Latch   ← gated = enable (EN)



S
EN
R

1-enabled

S
EN
R
Q
Q'

## Gated D Latch   aka data latch



D
EN
Q
Q'

1-enabled

D
EN
Q
Q'

→ avoids having invalid input

| EN | D | $Q^+$ | |
|----|---|-------|---|
| 1 | 0 | 0 | Reset |
| 1 | 1 | 1 | Set |
| 0 | x | Q | No change |

When EN = 1, $Q^+ = D$

↑ if non-gated, D latch pretty useless

Flip Flops ← use this over latch (only difference is the source of enable)

– synchronous (by clock)

– bistable

– output changes state at a specified point based on a clock

- rising edge
  or
- falling edge



positive edge-triggered
S-R flip-flop

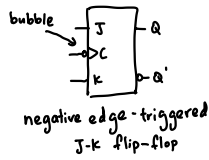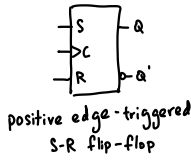negative edge-triggered
J-K flip-flop



almost instant

* use edge ∵ period of time is <u>very short</u> compared to period of high/low

* choosing positive or negative edge-triggered is arbitrary

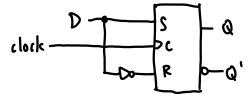* flip-flops below are positive edge-triggered

---

S-R flip-flop

X – don't care
↑ – rising edge

| S | R | Clock | $Q^+$ | |
|---|---|-------|-------|---|
| 0 | 0 | X | Q | no change |
| 0 | 1 | ↑ | 0 | reset |
| 1 | 0 | ↑ | 1 | set |
| 1 | 1 | ↑ | N/A | invalid |

→ $Q^+ = S + R' \cdot Q$

---

D flip-flop

| D | Clock | $Q^+$ | |
|---|-------|-------|---|
| 1 | ↑ | 1 | set |
| 0 | ↑ | 0 | reset |

→ $Q^+ = D$

→ used for parallel data transfer



build D flip-flop using S-R flip-flop

---

J-K flip-flop

| J | k | Clock | $Q^+$ | |
|---|---|-------|-------|---|
| 0 | 0 | X/↑ | Q | no change |
| 0 | 1 | ↑ | 0 | reset |
| 1 | 0 | ↑ | 1 | set |
| 1 | 1 | ↑ | Q' | toggle |

→ $Q^+ = J \cdot Q' + k' \cdot Q$



PTD = Pulse transition detector

---

T flip-flop

| T | Clock | $Q^+$ | |
|---|-------|-------|---|
| 0 | ↑ | Q | no change |
| 1 | ↑ | Q' | toggle |

→ $Q^+ = T' \cdot Q + T \cdot Q'$
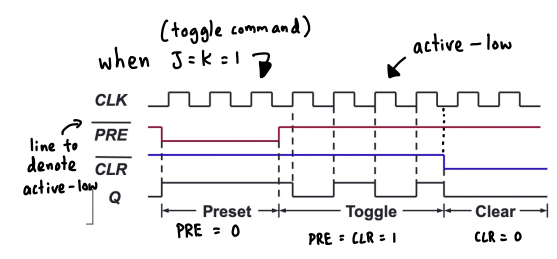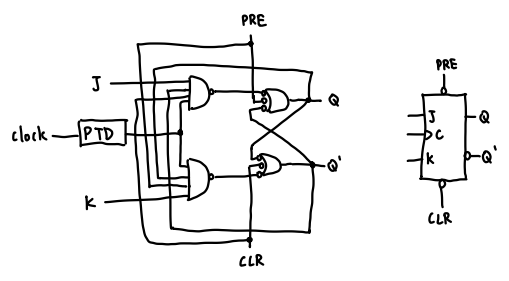
= $T \oplus Q$





build T flip-flop using J-K flip-flop

---

Asynchronous Inputs

PRE (preset) = active , Q is immediately set to 1 ⎫
                                                    ⎬ ignores clock
CLR (clear) = active , Q is immediately set to 0  ⎭

flip-flop works normally when PRE = CLR = not active ← can use to initialise flip-flop

→ if PRE = CLR = active, invalid input!

* active-high: active = 1

* active-low: active = 0



when J = k = 1 (toggle command)

active-low

line to denote active-low



|← Preset →|← Toggle →|← Clear →|
PRE = 0    PRE = CLR = 1   CLR = 0

m flip-flops + n inputs → $2^{m+n}$ rows

| Present State | | Input | Next State | | Output |
|---|---|---|---|---|---|
| A | B | x | $A^+$ | $B^+$ | y |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

Compact →

| Present State | Next State | | Output | |
|---|---|---|---|---|
| | x=0 | x=1 | x=0 | x=1 |
| AB | $A^+B^+$ | $A^+B^+$ | y | y |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

## State Diagram

state → circle

arrow → transition (with a label of I/o)    ↙ input/output

m flip-flops → ≤ $2^m$ states

— circuit output functions
  → y = ...

— flip-flop input functions
  → JA = ...  ⎫ 1$^{st}$ letter: flip-flop input
     KA = ...  ⎭ 2$^{nd}$ letter: flip-flop name

— state equations
  $A^+$ = ...

## Excitation Table

— used in design

— given transition, find flip-flop inputs

| Q | $Q^+$ | J | K | S | R | D | T |
|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | x | 0 | x | 0 | 0 |
| 0 | 1 | 1 | x | 1 | 0 | 1 | 1 |
| 1 | 0 | x | 1 | 0 | 1 | 0 | 1 |
| 1 | 1 | x | 0 | x | 0 | 1 | 0 |
| x | x | x | x | x | x | x | x |

## Analysis

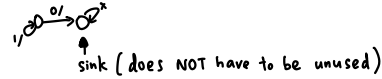| Present State | | Input | Next State | | flip-flop inputs | | | |
|---|---|---|---|---|---|---|---|---|
| A | B | x | $A^+$ | $B^+$ | JA | KA | JB | KB |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

Start from circuit diagram, obtain state table or state diagram

if unused states, "x" for everything

draw K-map / truth table to get boolean functions



Sink: once the circuit enters this state, never exits



sink (does NOT have to be unused)

self-correcting circuit: if circuit (somehow) enters an unused state,
                    will exit the unused state in a finite number of transitions
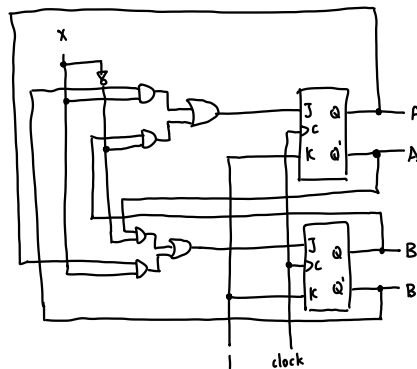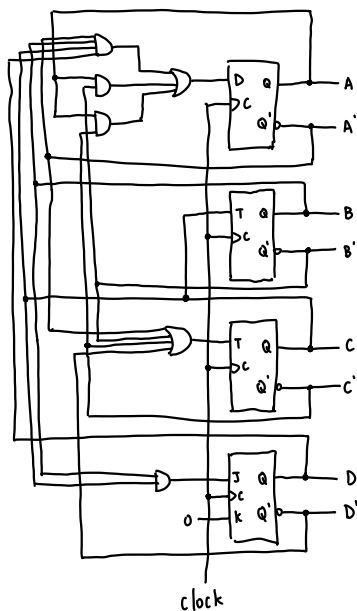


unused state        * draw slash even if only input/output present



clock



1  clock

## Memory

→ stores program & data

1 byte = 8 bits
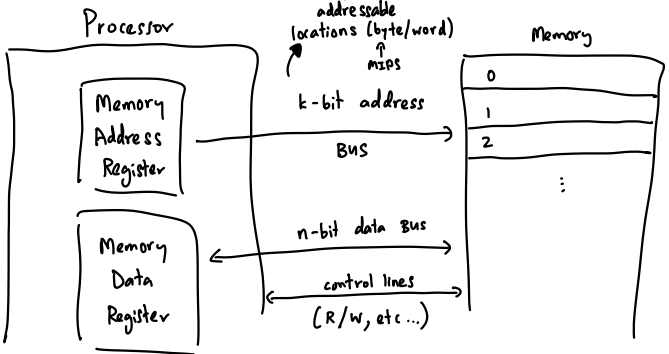
1 word: multiple of bytes, for transfer main memory ⇆ register
usually size of register

1 KB : $2^{10}$ bytes ⎫
1 MB : $2^{20}$ bytes ⎬ powers
1 GB : $2^{30}$ bytes ⎪ of 2
1 TB : $2^{40}$ bytes ⎭

### desirable properties
— fast access
— large capacity     ⎫ but most
— economical cost    ⎬ don't have
— non-volatile (saved even when   ⎭ all these
  no power)

### memory hierarchy

registers          ⎫ fast,
main memory        ⎪ expensive (small numbers),
disk storage       ⎬ volatile
magnetic tapes     ⎪
                   ⎪ slow,
                   ⎭ cheap (large numbers),
                     non-volatile

*  up to $2^k$ addressable locations (byte/word)
↑ mips

**Processor**    **Memory**

Memory Address Register

$k$-bit address — Bus

0
1
2
⋮

Memory Data Register

$n$-bit data Bus

control lines (R/W, etc ...)

data input ↓ $n$

address — $k$ → memory unit $2^k$ words $n$ bits per word

read/write control →

↓ $n$

data output

| memory enable | read/write | |
|---|---|---|
| 0 | x | none |
| 1 | 0 | write |
| 1 | 1 | read |

Memory unit stores binary information in groups of bits (words)

data consists of n lines (for n-bit words)

data input lines carry data to be written

data output lines carry data to be read

address consists of k lines

control line read/write specifies direction of transfer

**Write**
- transfer address
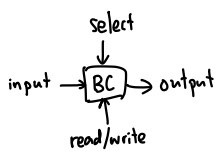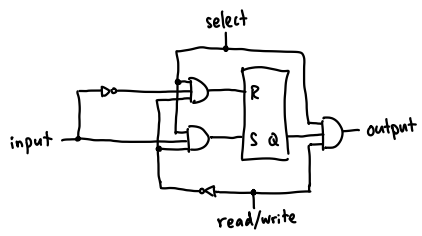- transfer data to be written
- set read/write = 0

**Read**
- transfer address
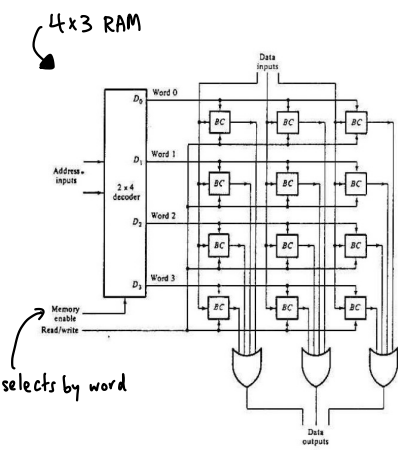- set read/write = 1

---

## RAM : Random Access Memory

### Static RAM
- use flip-flops as memory cells

select

input    R / S Q → output

read/write

select
input → BC → output
read/write

BC : binary/bit cell
[single bit]

### Dynamic RAM
- use capacitor charges to represent data
- simpler in circuitry BUT need to be constantly refreshed

4×3 RAM

Data inputs

$D_0$  Word 0
$D_1$  Word 1
Address inputs  2×4 decoder
$D_2$  Word 2
$D_3$  Word 3

Memory enable
Read/write

BC

Data outputs

selects by word

→ keep combining blocks of RAM
to make BIGGER blocks of RAM

word size
↓
1K × 8 : 1024 bits $(2^{10})$

RAM 1K × 8

input data — $8$ → DATA (8)    (8) $8$ → output data
address — $10$ → ADDR (10)
chip select — CS
read/write — RW

4K × 8 : 4096 bits $(2^{12}$ address$)$

2M × 32 : $2^{21}$ address, 32 bit word size
  └ use 512k × 8 memory chips

$2^{21}$ ∴ $2^{21}$ > 2 million

# Pipelining

- does not help latency of single task
- helps **throughput** of entire workload
  **multiple** tasks operating **simultaneously** using different resources

limited by:
  - slowest pipelining stage
  - stall for dependencies

## Pipelined Implementation

- one cycle per pipeline stage
- data required for each stage needs to be stored separately because previous component might be used for something else already

Data used by subsequent instructions are stored in programmer-visible state elements: PC, register file, memory
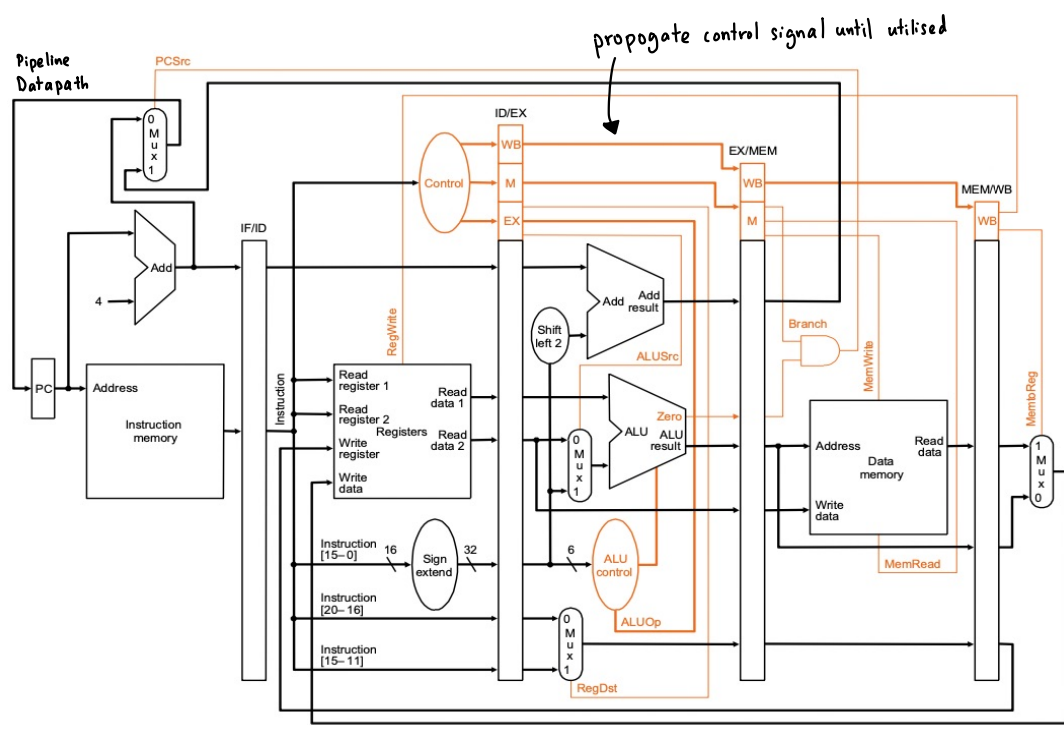
## Pipeline Registers — not 32 bits (not really registers, its a collection of things)

| IF/ID | → | ID/EX | → | Ex/MEM | → | MEM/WB | → | ~ end ~ |
|---|---|---|---|---|---|---|---|---|
| · RR1 and RR2 | | · RD1 and RD2 | | · (PC + 4) + (immediate × 4) | | · ALU result | | · result is written back |
| · 16-bit offset (to be sign-extended to 32-bits) | | · 32-bit immediate value | | · ALU result | | · memory read data | | to register file (if applicable) |
| · PC + 4 | | · PC + 4 | | · isZero? signal | | · write register | | using the write register stored in MEM/WB |
| | | · write register | | · RD2 | | | | |
| | | | | · write register | | | | |

\* need to "pass" write register along the pipeline registers

## MIPS Pipeline Stages

IF: instruction fetch
ID: instruction decode & register file read
Ex: execute/address calculation
MEM: memory access
WB: write back

group control signals based on pipeline stage  →
[ IF & ID no control signals ]

| | EX Stage | | | MEM Stage | | | WB Stage | |
|---|---|---|---|---|---|---|---|---|
| | RegDst | ALUSrc | ALUOp | MemRead | MemWrite | Branch | MemToReg | RegWrite |
| R-type | 1 | 0 | 10 | 0 | 0 | 0 | 0 | 1 |
| lw | 0 | 1 | 00 | 1 | 0 | 0 | 1 | 1 |
| sw | X | 1 | 00 | 0 | 1 | 0 | X | 0 |
| beq | X | 0 | 01 | 0 | 0 | 1 | X | 0 |

propogate control signal until utilised



Pipeline Datapath
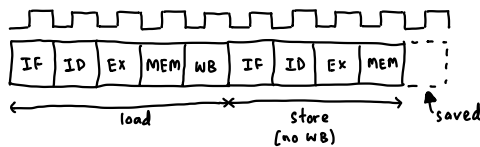
## Single-cycle Implementation

- update all state elements at the end of clock cycle
  (PC, register file, data memory)

- every instruction takes up 1 cycle

- cycle length based on slowest instruction



```
| load | store | //// |
                 ↑ unused
```

Cycle time: $CT_{seq} = max\left(\mathcal{E}_{k=1}^{N} T_k\right)$ ←— bounded by slowest instruction

For j instructions: $Time_{seq} = j \times CT_{seq}$


## Multi-cycle Implementation

- every stage takes up 1 cycle

- cycle length based on slowest stage

- instructions: n cycles ⟷ n stages



```
| IF | ID | EX | MEM | WB | IF | ID | EX | MEM |  ⌐ ⌐
←——————— load ———————×——————— store ————————→  ↑ saved
                                  (no WB)
```

Cycle time: $CT_{multi} = max(T_k)$ ←— bounded by slowest stage

For j instructions: $Time_{multi} = \underbrace{j \times average\ CPI}_{total\ cycles} \times CT_{multi}$


## Pipelining Implementation

Cycle time: $CT_{pipeline} = max(T_k) + T_d$ ←— overhead / latency
∴ pipeline register etc

For j instructions: $j + N - 1$ cycles ← after first, each additional only adds 1 more

$$Time_{pipeline} = (j + N - 1) \times CT_{pipeline}$$

assumptions → | Ideal case | → Speedup $= \dfrac{Time_{seq}}{Time_{pipeline}}$ ∵ ratio

- every stage take equal time

- no overhead, $T_d = 0$

$\qquad = \dfrac{j \times \mathcal{E}_{k=1}^{N} T_k}{(j+N-1) \times (max(T_k) + T_d)}$

- $j \gg N$

$\qquad = \dfrac{j \times N \times T_1}{(j+N-1) \times T_1}$  ∵ $T_1 = T_2 = \cdots = T_N$

N: number of pipeline stages

$\qquad \approx \dfrac{j \times N \times T_1}{j \times T_1}$  ∵ $j \gg N$
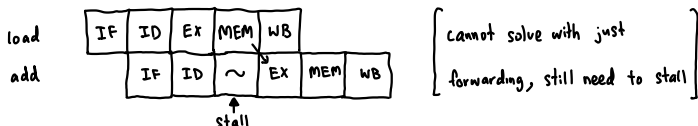
$\qquad = N$ ☆


## Pipeline Hazards

Structural hazard

  - simultaneous usage of hardware

Data hazard / dependency

  - read / write data in same register

Control hazard / dependency

  - to branch or not

} instruction dependencies


## Read-After-Write (RAW) aka true data dependency

  add $1, $2, $3

  sub $4, $1, $2 ← still uses stale/old value of $1

Soln: Forward the data to next instructions

  Bypass data read from register file

```
add | IF | ID | EX | MEM | WB |
add      | IF | ID | EX | MEM | WB |
```

```
load | IF | ID | EX | MEM | WB |
add       | IF | ID | ~ | EX | MEM | WB |
                      ↑
                    stall
```
[ cannot solve with just forwarding, still need to stall ]

WAR and WAW also exist but do <u>not</u> cause problems
unless executed out of order


## Solutions

① Stall / delay pipeline

② Separate Data and Instruction memory

③ Split register cycle into half

  - write then read

| pipeline stage | | hardware |
|---|---|---|
| IF | → | IM (instruction memory) |
| ID | → | REG |
| MEM | → | ALU |
| EX | → | DM (data memory) |
| WB | → | REG |

] separate to avoid structural hazard

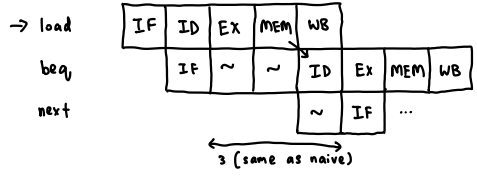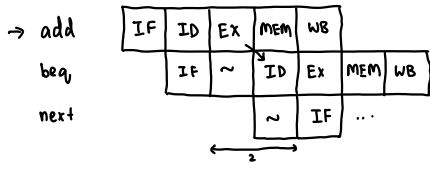∴ conditional branching

Control Dependency          * can use all these methods in conjunction

① Stall pipeline
- waiting for branch outcome wastes 3 clock cycles [next instruction starts after MEM of beq/bne]
- branching happens very often so, not acceptable (BAD!)

② early branch resolution          magnitude comparator
- move computation of isZero? earlier [ID stage instead of MEM]
- reduced from 3 to 1 cycle delay
- but if register modified previously, still need to stall (RAW problem)

→ add

| IF | ID | EX | MEM | WB |

beq

| IF | ~ | ID | EX | MEM | WB |

next

| ~ | IF | ... |

├──2──┤

→ load

| IF | ID | EX | MEM | WB |

beq

| IF | ~ | ~ | ID | EX | MEM | WB |

next

| ~ | IF | ... |

├──3 (same as naive)──┤

③ branch prediction
- assume not branching & start computing      } simple/naive way
- flush successor instructions from pipeline
  (if branch) → terminate immediately
- flush happens before anything is written to register/memory
  so no problem
- if correct → no stall needed
- if fail → equivalent to not doing branch prediction in the first place

- software can optimise dynamic prediction based on previous branches
  (out of scope)

④ delayed branching (software/compiler solution)
- shift order of execution
- compute some non-control dependent instruction
  [known as branch-delay slot]
        ↳ for MIPS (w/ early branch), 1 slot          } done by compiler

if no suitable instruction, use a NOP (no-op) instruction
~50% of the time can find

*easier with early branch resolution
  (less branch-delay slots to fill)

          not tested
⑤ Multiple Issue Processor
- multiple instructions in each stage

Static multiple issue
- compiler specifies set of instructions to be executed together
- simple hardware, complex compiler

    Explicitly Parallel Instruction Compiler (EPIC)
    Very Long Instruction Word (VLIW) — IA64

dynamic static issue
- hardware decides
- complex hardware, simple compiler

    Superscalar processor — most modern processors

          accessing memory takes
          ~50x the clock cycles
Memory

          uses flip-flop    main
                           memory

|  | Registers | SRAM | DRAM | Hard Drive |
|---|---|---|---|---|
| transistors/memory cell | | 6 | 1 | |
| density | | low | high | |
| access latency | very fast 20ps | fast 0.5-5ns | slow 50-70ns | very slow 5-20ms |
| cost | $$$$ | $$$ | $ | ¢ |
| capacity | 100s bytes | 100s kB | 100s MB | 100s GB |

ideal:
    - small but fast memory near CPU
    - large but small memory further away

DDR SDRAM
Double Data Rate Synchronous Dynamic RAM
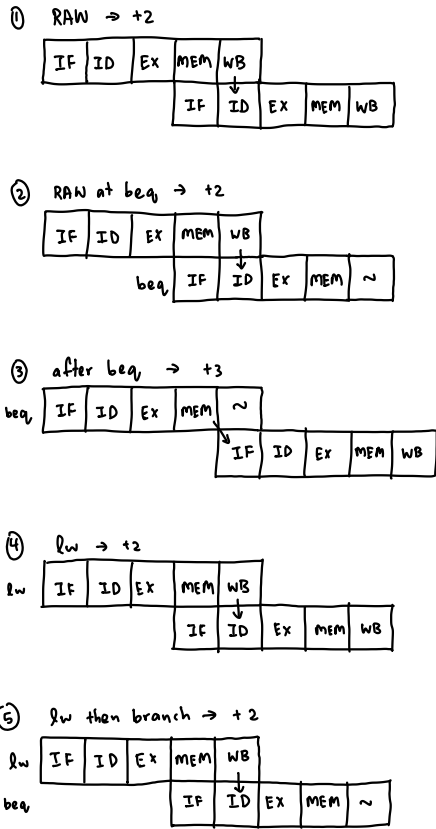↳ both positive & negative clock edges used

CS2100                          CS2106
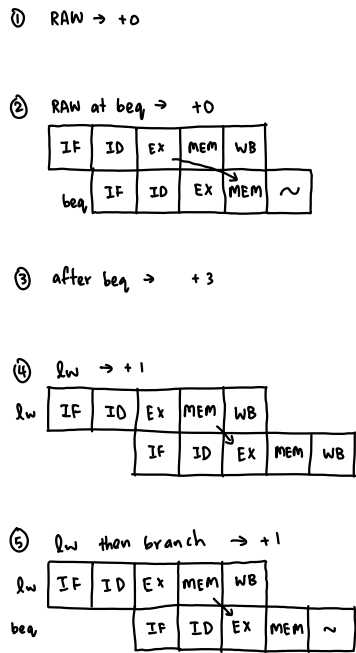make SLOW memory look FAST       make SMALL memory look BIG
- cache                          - virtual memory
- hardware managed               - OS managed

              transparent
              to programmer

## w/o forwarding & MEM stage

① RAW → +2

IF | ID | EX | MEM | WB
              IF | ID | EX | MEM | WB

② RAW at beq → +2

IF | ID | EX | MEM | WB
beq: IF | ID | EX | MEM | ~

③ after beq → +3

beq: IF | ID | EX | MEM | ~
              IF | ID | EX | MEM | WB

④ lw → +2

lw: IF | ID | EX | MEM | WB
         IF | ID | EX | MEM | WB

⑤ lw then branch → +2

lw: IF | ID | EX | MEM | WB
beq: IF | ID | EX | MEM | ~

## w/o forwarding & ID stage

① RAW → +2

② RAW at beq → +2

③ after beq → +1

beq: IF | ID | EX | MEM | ~
        IF | ID | EX | MEM | WB

④ lw → +2

⑤ lw then branch → +2

## w/ forwarding & MEM stage

① RAW → +0

② RAW at beq → +0

IF | ID | EX | MEM | WB
beq: IF | ID | EX | MEM | ~

③ after beq → +3

④ lw → +1

lw: IF | ID | EX | MEM | WB
        IF | ID | EX | MEM | WB

⑤ lw then branch → +1

lw: IF | ID | EX | MEM | WB
beq: IF | ID | EX | MEM | ~

## w/ forwarding & ID stage

① RAW → +0

② RAW at beq → +1

IF | ID | EX | MEM | WB
beq: IF | ID | EX | MEM | ~

③ after beq → +1

④ lw → +1

⑤ lw then branch → +2

## branch prediction

① if correct, after beq → +0

② if wrong, same as no branch prediction

## jump instruction resolution

① at ID → +1

② at MEM → +3

☆ RAW can last for multiple instruction (w/o forwarding)

**Cache** ← uses fast SRAM (usually)     * also usable for instruction memory

Keep the frequently and recently used data in smaller but faster memory!

Principle of locality: Program accesses only a <u>small portion</u> of the memory address space within a <u>small time interval.</u>

<u>Temporal Locality</u>
If an item is referenced,
tend to be referenced
soon <u>again</u>.

<u>Spatial Locality</u>
If an item is referenced,
<u>nearby items</u> tend
to be referenced soon.

Instructions & data are in different localities!

Average Access Time =     Hit Rate × Hit Time
                  + (1 - Hit Rate) × Miss Penalty
                    ↳ time to replace cache block + hit time    ∴ miss penalty < hit time

<u>larger</u> block size
  ++ spatial locality
  -- larger miss penalty ∵ larger block to fill
  -- lesser cache blocks → miss rate go up
  —→ need balance to maximise average access time

miss rate

block size

block size > word size
to exploit locality

Cache contains:
1. data block
2. Tag of memory block
3. Valid bit (initial: 0, after storing data: 1)

cache hit = valid bit is TRUE
          && tag[index] == tag in cache

<u>Compulsory Miss</u>
  — on first access
  — valid == FALSE
  — aka cold start miss,
    first reference miss

<u>Conflict Miss</u>
  — tag[index] ≠ tag in cache
  — aka collision/interference miss
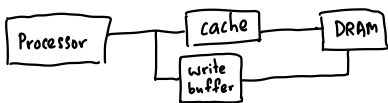  — FA cache no conflict misses

<u>Capacity Miss</u>
  — blocks discarded from cache
    ∵ cache cannot contain all needed blocks
  — occurs in FA Cache

→ on read miss:
    1. load data into cache
    2. load from there to register

**Writing with a Cache**

1. Write-through cache
   — write data both to cache and main memory
   — use write buffer to avoid slow write

2. Write-back cache
   — only write to cache
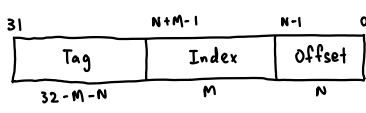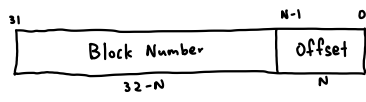   — only write to main memory when cache block is replaced/evicted

   — use extra bit (dirty bit) to store if a write occurred
     ↳ only write to main memory if dirty == TRUE

**Write Miss**

1. Write allocate
   — load the entire block into cache
   — change the word
   — write back depending on write policy

2. Write around
   — write directly to memory

## Direct Mapping Cache

```
31          N-1      0
┌──────────────┬────────┐
│ Block Number │ Offset │
└──────────────┴────────┘
      32-N         N
```

```
31      N+M-1    N-1     0
┌──────┬────────┬────────┐
│ Tag  │ Index  │ Offset │
└──────┴────────┴────────┘
 32-M-N    M        N
```

→ essentially 1-way SA cache

Number of blocks in memory $= 2^{32-N}$

total bytes (cache) = words × 4 (MIPS) $= 2^{N+M}$ bytes $(\because 2^N \times 2^M)$

total bytes (block) $= 2^N$ bytes $(N = offset)$

Number of cache blocks $= 2^M$ $(M = index)$

Tag $= 32 - (N+M)$ bits
↳ $2^{32-(N+M)}$ share same spot

### Cache trashing

- 0% hit rate when A[i] and B[i] mapped to same cache location so constantly re-write to same block
- common reason: multiple of cache size

## Set Associative (SA) Cache  ← minimise conflict misses

N-way, N>1, N = number of cache blocks in a set

1. find set (w/ set index)

2. place in any block in the set (prioritise left to right for CS2100)

∴ need to search all the cache blocks in the set (search is done concurrently)



Miss rate of N-sized direct-mapped cache
$\approx$ Miss rate of $\frac{N}{2}$-sized 2-way SA cache

memory addressing same as direct mapping

set index = block number % num of cache sets
↳ $= \frac{\text{num of cache blocks}}{N}$

number of cache sets $= 2^M$   ← NOT blocks!

$(M = set index)$

## Fully Associative (FA) Cache

place memory block ANYWHERE (no more mapping function)

++ can be placed anywhere (no conflict miss)

-- need to search ALL cache blocks

- same cold miss

- capacity miss happens & decreases as cache size increases

```
31          N-1      0
┌──────────────┬────────┐
│ Block Number │ Offset │
└──────────────┴────────┘
      32-N         N
```

block number = tag

```
31          N-1      0
┌──────────────┬────────┐
│     Tag      │ Offset │
└──────────────┴────────┘
      32-N         N
```

→ essentially all cache blocks in same set of SA cache

total bytes (block) $= 2^N$ bytes
$(N = offset)$

Tag $= 32 - N$ bits
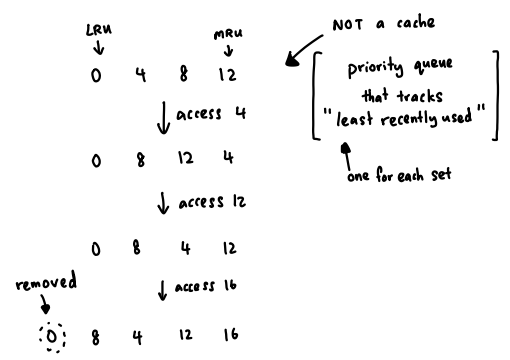
## Block Replacement Policy

① Least Recently Used (LRU)

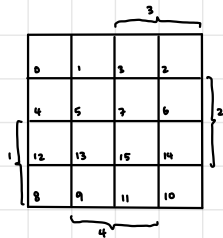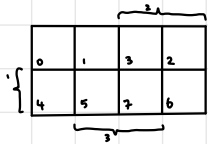- for cache hit, record the cache block that was accessed

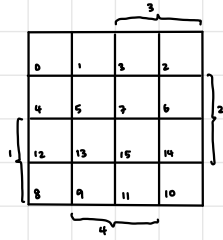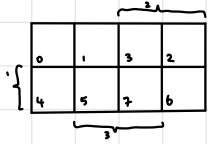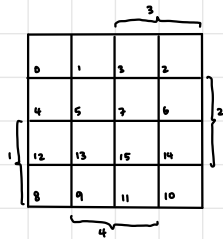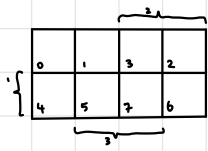- when replacing, replace the one that has not been used for the longest time

- good ∵ temporal locality

- but can be hard to keep track of

```
LRU              MRU
 ↓                ↓
 0    4    8    12          ┌─────────────────┐
                            │ NOT a cache     │
          ↓ access 4        │ priority queue  │
                            │ that tracks     │
 0    8    12   4           │ "least recently │
                            │  used"          │
          ↓ access 12       └─────────────────┘
                              ↑ one for each set
 0    8    4    12

          ↓ access 16
removed
 ↓        ↓ access 16
:0:  8    4    12   16
```

② FIFO queue

③ Random replacement (RR)

④ Least frequently used (LFU)

## Multilevel Cache

- use multiple caches (w/ different associativities) in tandem!

- Modern CPUs have 3 levels: L1, L2, L3
                              ↑ Intel/AMD

Fast boolean algebra / Karnaugh map worksheet

**Truth table (4-bit)**

| | ABCD |
|---|---|
| 0 | 0000 |
| 1 | 0001 |
| 2 | 0010 |
| 3 | 0011 |
| 4 | 0100 |
| 5 | 0101 |
| 6 | 0110 |
| 7 | 0111 |
| 8 | 1000 |
| 9 | 1001 |
| 10 | 1010 |
| 11 | 1011 |
| 12 | 1100 |
| 13 | 1101 |
| 14 | 1110 |
| 15 | 1111 |

**Pair table**

| | ABCD | AB | AC | AD | BC | BD | CD |
|---|---|---|---|---|---|---|---|
| 0 | 0000 | 00 | 00 | 00 | 00 | 00 | 00 |
| 1 | 0001 | 00 | 00 | 01 | 00 | 01 | 01 |
| 2 | 0010 | 00 | 01 | 00 | 01 | 00 | 10 |
| 3 | 0011 | 00 | 01 | 01 | 01 | 01 | 11 |
| 4 | 0100 | 01 | 00 | 00 | 10 | 10 | 00 |
| 5 | 0101 | 01 | 00 | 01 | 10 | 11 | 01 |
| 6 | 0110 | 01 | 01 | 00 | 11 | 10 | 10 |
| 7 | 0111 | 01 | 01 | 01 | 11 | 11 | 11 |
| 8 | 1000 | 10 | 10 | 10 | 00 | 00 | 00 |
| 9 | 1001 | 10 | 10 | 11 | 00 | 01 | 01 |
| 10 | 1010 | 10 | 11 | 10 | 01 | 00 | 10 |
| 11 | 1011 | 10 | 11 | 11 | 01 | 01 | 11 |
| 12 | 1100 | 11 | 10 | 10 | 10 | 10 | 00 |
| 13 | 1101 | 11 | 10 | 11 | 10 | 11 | 01 |
| 14 | 1110 | 11 | 11 | 10 | 11 | 10 | 10 |
| 15 | 1111 | 11 | 11 | 11 | 11 | 11 | 11 |

**Bottom-left truth table**

| | | | | |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 0 | 0 | 1 |
| 2 | 0 | 0 | 1 | 0 |
| 3 | 0 | 0 | 1 | 1 |
| 4 | 0 | 1 | 0 | 0 |
| 5 | 0 | 1 | 0 | 1 |
| 6 | 0 | 1 | 1 | 0 |
| 7 | 0 | 1 | 1 | 1 |
| 8 | 1 | 0 | 0 | 0 |
| 9 | 1 | 0 | 0 | 1 |
| 10 | 1 | 0 | 1 | 0 |
| 11 | 1 | 0 | 1 | 1 |
| 12 | 1 | 1 | 0 | 0 |
| 13 | 1 | 1 | 0 | 1 |
| 14 | 1 | 1 | 1 | 0 |
| 15 | 1 | 1 | 1 | 1 |

unused = X

**Fast boolean algebra**

2:  3, 5
3:  F, 33, 55
4:  FF, F0F, 3333, 5555
5:  FFFF, FF00FF, F0F0F0F, 3333 3333, 5555 5555