**CODING and INFORMATION**
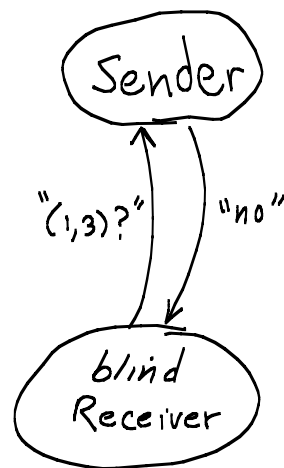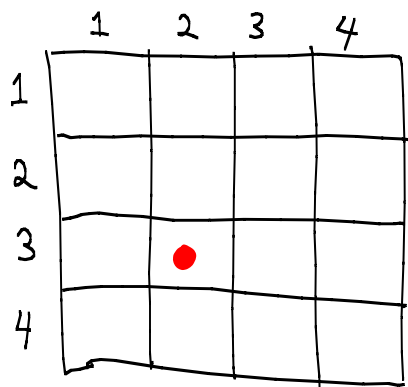We need encodings for data.

where's the ball?



- Ask yes/no questions.

- Are all questions-answers equally informative?

- Min number of questions?

  How many questions must be asked to be certain where the ball is. (cases: avg, worst, best)

- How much do you learn from each yes/no?

- Is that a Good series of questions, "In (1,3)"?

Avg number of questions needed (assume equally likely boxes)?

P( Hit 1st ) = 1/16

$\Rightarrow$ P( Hit 2nd ) = P( Hit 2nd | Miss 1st )P( Miss 1st ) = (1/15)(15/16)     = 1/16

P( Hit 3rd ) = (1/14) * P( Miss 2nd and 1st ) = (1/14)(14/15)(15/16)   = 1/16

E( n ) = 1*(1/16) + 2*(1/16) + ... + 15*(1/16) + 15(1/16)   =   (1+2+3+...+15+15) / 16  = 8 1/2 - 1/16

- Different set of question?
  "in (1,*) or (2,*)", "in (1,*)", "in (2,1) or (2,2)" ...

  $\Rightarrow$ Each Q reduces space by 1/2  $\Rightarrow$ Exactly 4 questions

$\Rightarrow$ Most information if each Q-A splits possibilities 50-50.

Measure info content of answer?

prob(yes) = prob(no)  $\Rightarrow$  $\log_2(\text{Prob(yes)})$

$= \log_2(2^{-1})$

$= -1$   (hmm, make +?)

$\Rightarrow$ 1 bit

$-\log_2(\text{Prob})$ = **info measure**

$\text{Prob}(\text{yes}) = \frac{1}{4}$ $\Big\}$ assume
$\text{Prob}(\text{no}) = \frac{3}{4}$

yes: $\log(2^{-2}) = 2$ bits

no: $\log(\frac{3}{4}) \approx -\log(.7) \approx -\log(\frac{1}{\sqrt{2}}) = \frac{1}{2}$ bit

| yes | no |
|-----|-----|
| $\frac{1}{4}$ | $\frac{3}{4}$ |

$Q = $ "Left of here?"

### Avg info of answer?

$$= \text{prob}(\text{yes}) \cdot (2 \text{ bits}) + \text{prob}(\text{no}) \cdot (\frac{1}{2} \text{ bit})$$

$$= \frac{1}{4}(2) + \frac{3}{4}(\frac{1}{2}) = \frac{1}{2} + \frac{3}{8} = \frac{7}{8} \text{ bit}$$
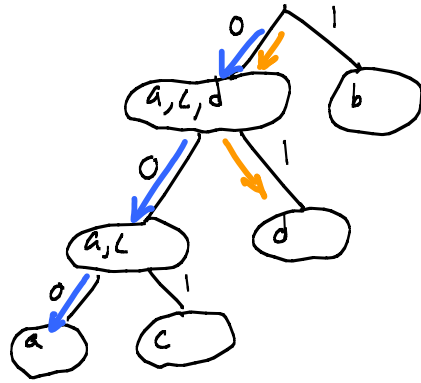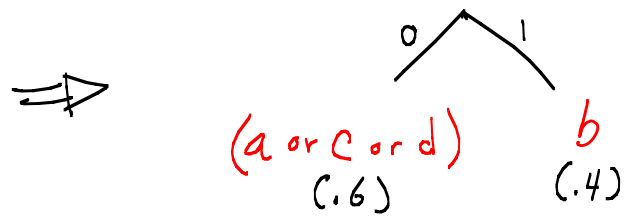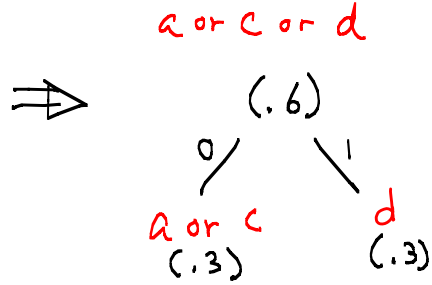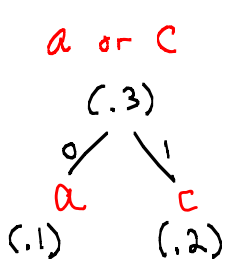
**Extreme:** $\text{Prob}(\text{yes}) = \frac{1}{2^{10}}$   $\text{Prob}(\text{no}) = \frac{(2^{10}-1)}{2^{10}} \approx 1$

$\text{avg} = \frac{1}{2^{10}}(10 \text{ bits}) + (1)\underbrace{\log(1)}_{\longrightarrow 0 \text{ bits}}$   $\Rightarrow \frac{1}{100}$ bit

**Thm** max avg. if $P_i = \frac{1}{n}$   for $n$ possibilities.
$(\sum P_i = 1)$

How to encode to get 50-50? $\Rightarrow \left(\begin{array}{c}\text{Pair} \\ \text{least} \\ \text{Probable}\end{array}\right) \sim$ equally likely.

$(a, b, c, d)$ w/ prob. $(.1, .4, .2, .3)$

a or c
(.3)

$\overset{0}{\diagup}\overset{1}{\diagdown}$

a     c
(.1)   (.2)

$\Rightarrow$

a or c or d
(.6)

$\overset{0}{\diagup}\overset{1}{\diagdown}$

a or c    d
(.3)    (.3)

$\Rightarrow$

$\overset{0}{\diagup}\overset{1}{\diagdown}$

(a or c or d)    b
(.6)     (.4)

Huffman Code | msg

0 0 0 | a
0 0 1 | c
0 1 | d
1 | b

Shannon Info Theory = Expected # bits = $E[(-\log(P_r))]$

Avg. content =
(Entropy, H)
$$-\left[0.1\log(0.1) + 0.4\log(0.4) + 0.2\log(0.2) + 0.3\log(0.3)\right]$$
a         b         c         d

$$= 0.33 + 0.53 + 0.46 + 0.52 = 1.84$$
bits per message.

How'd we do?

avg. # bits = $(0.1)3 + (0.4)1 + (0.2)3 + (0.3)2$
                   a        b        c        d

$$= 0.3 + 0.4 + 0.6 + 0.6 = 1.9 \text{ bits}$$

We are sending more bits than information content, but we are very close.

MIN-Length code ==> MAX compression ==> most info bits in least number of communicated bits.

Suppose n different "messages" to send, n = 2^k.
Maximum entropy => equally likely:  Prob( message-i ) = (1/n)   for any message-i.

Expected information per message is,

 Sum[ - (1/n) log[ 1/n ] ]   =  - n ( 1/n log[ 1/n ] )   =  -1 log[ 2^-k ]   =  -1 (-k)   =  k  bits per message.  If we
use a k-bit code for our messages, we will be 100% compressed. (k-bit integers? Are they equally likely?)
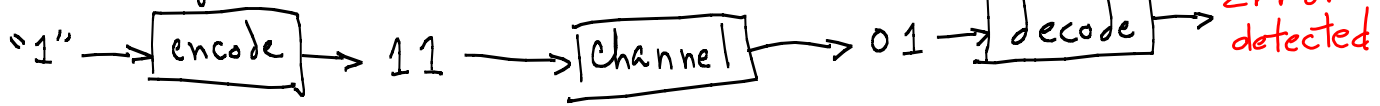
# Error Detection / Correction

"message" could be a bit, a string of bits, a character, a page of characters, ...

"message" → Communication Channel → "mfssage"

message coded in bits:

"1" → encode → 1 → Channel → 0 → decode → "0"   1-bit Error — not detectable

2-bit encoding

"1" → encode → 11 → Channel → 01 → decode → Error detected

Code words: 00 and 11   ---   "0" and "1"
Code words: 10 and 01   ---   1-bit errors: odd parity codeword indicates error.
Works for k-bit messages w/ 1 parity bit, but only if 2-bit errors very unlikely (never occur?).

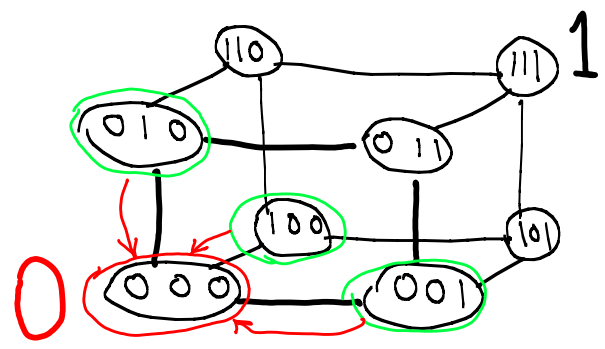1-bit Error Correction w/ 3-bit code words:
"0" ==> 000
"1" ==> 111

| | |
|---|---|
| 001 ==> "0" | 011 ==> "1" |
| 010 ==> "0" | 101 ==> "1" |
| 100 ==> "0" | 110 ==> "1" |



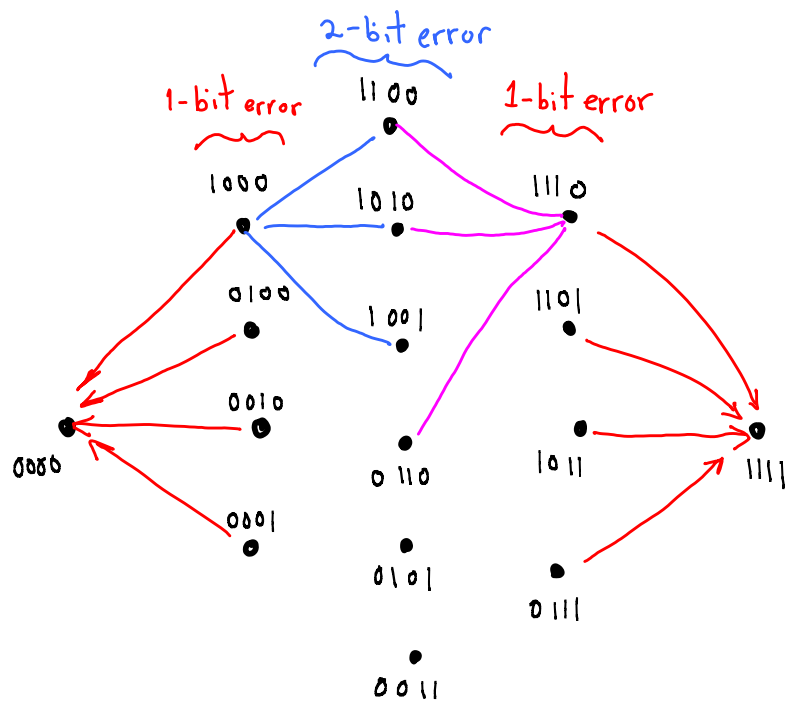1-bit Correction, 2-bit Detection

-- odd parity: 1-bit error corrected

-- exactly two 1's: 2-bit error detected

-- otherwise: no error

How many extra bits are needed at minimum? Depends on noise in channel: Shannon Noisey Coding Theorem.

Can you think of a scheme like the parity-bit scheme that uses as few bits as possible? (See Reed-Solomon codes, for instance.)

More bits, higher error probability.

# ALU, numbers

d_i is a "digit", a symbol for a value: value( "d_i")

b is a value, the "base" of the number notation.

**Positional notation for numbers**

There is a rule to find the value, given the symbols.

$$value(\text{"}d_n d_{n-1} \cdots d_1 d_0\text{"}) = d_n \cdot b^n + d_{n-1} \cdot b^{n-1} + \cdots \; d_1 \cdot b^1 + d_0 \cdot b^0$$

**unsigned 3-bit binary**

binary:
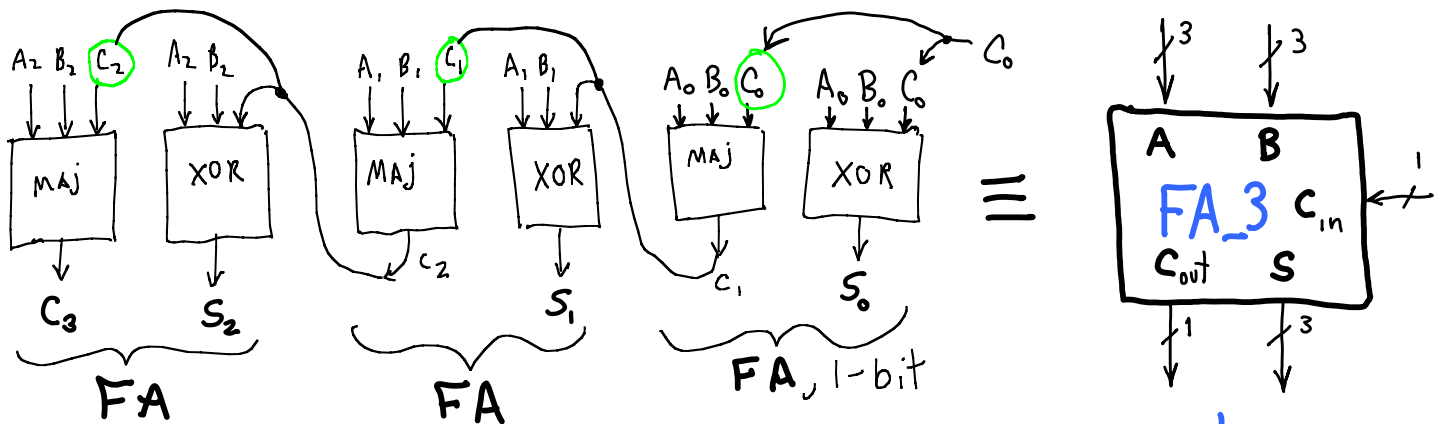--- digits = { "0", "1" }
--- base = 2

$$000 \longrightarrow value \triangleq 0 \cdot 2^2 + 0 \cdot 2^1 + 0 \cdot 2^0 = 0$$
$$001 \longrightarrow value \triangleq 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 = 1$$
$$\cdots$$
$$111 \longrightarrow value \triangleq 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = 7$$

*oops, more symbols?*

Let's do some 3-bit arithmetic.
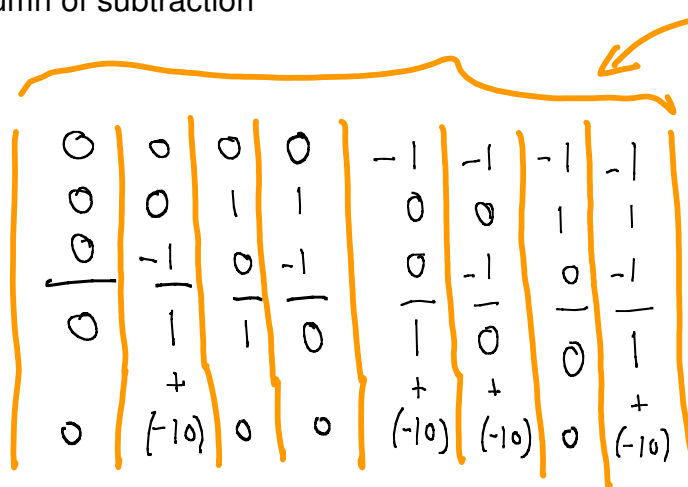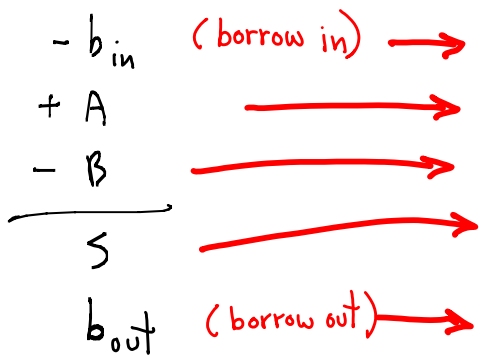ADD:

$$A_2 A_1 A_0 + B_2 B_1 B_0 = C_3 S_2 S_1 S_0$$



**3-bit Full Adder**

Let's try
SUBTRACTION

$$A_2 A_1 A_0 - B_2 B_1 B_0 = b_3 \; S_2 S_1 S_0$$

*← possible borrow*

But first, let's look at a single column of subtraction

$-b_{in}$ (borrow in) →
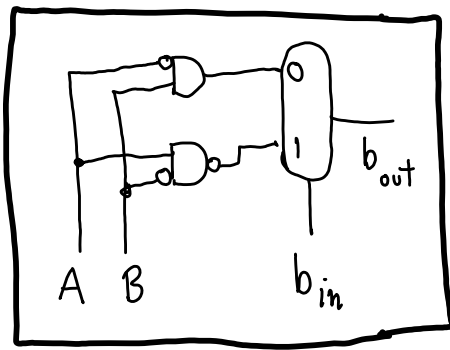$+A$ →
$-B$ →
_____
$S$ →
$b_{out}$ (borrow out) →



| 0 | 0 | 0 | 0 | -1 | -1 | -1 | -1 |
| 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 |
|___|___|___|___|___|___|___|___|
| 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 |
| 0 | (-10) | 0 | 0 | (-10) | (-10) | 0 | (-10) |

$= XOR(A, B, b_{in})$

$(b=1, A=1, B=1)$
one of the possible bit combinations in a single col. of SUB

| $b_{in}$ | A | B | $b_{out}$ |
|----|----|----|----|
| 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 | → $\overline{A_i} \cdot B_i$
| 0 | 1 | 0 | 0 |
| 0 | 1 | 1 | 0 |
| 1 | 0 | 0 | 1 |
| 1 | 0 | 1 | 1 |
| 1 | 1 | 0 | 0 | → $A_i \cdot \overline{B_i}$
| 1 | 1 | 1 | 1 |



BOR



$b ←$   $S_2$   $b_2$   $S_1$   $b_1$   $S_0$

# ALU

A   B



3-bit, bit-wise NAND



IR  | XY | AB | CD | EF |  INST. FORMAT

OP  SR1  SR2  DR

OP CODES
_____
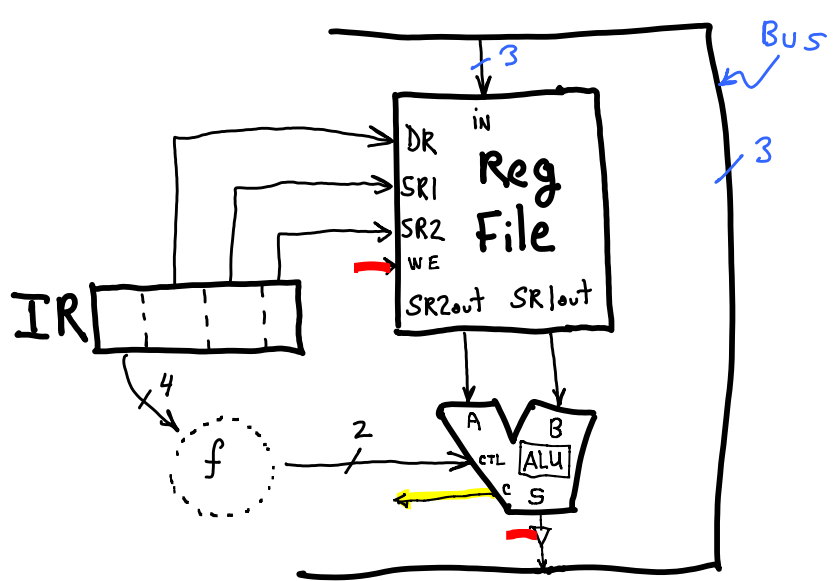00 NAND
01 ADD
10 SUB

CTL[1]
CTL[0]  CTL

S

3-bit ALU

(b or C)

(simplified instruction, only shows ALU ops)

It's almost this simple in the LC3.

This is a 3-bit version of LC3 (sort of).

Some sort of function $f$ converts 4-bit opcode to 2-bit ALU.ctl. In uCoded control, this function is implemented as 0/1 control bits in ROM.



## UNsigned errors

$c$ = carry/borrow output

$$A + B > 7 \implies c = 1 \quad, \quad S = (A+B) \bmod 2^3$$
$$A - B < 0 \implies c = 1 \quad, \quad S = (A-B) \bmod 2^3$$

$$C_3 S_2 S_1 S_0 \implies \left( \underline{C_3 \cdot 2^3} + S_2 \cdot 2^2 + S_1 \cdot 2^1 + S_0 \cdot 2^0 \right)_{\bmod 2^3} = S_2 \cdot 2^2 + S_1 \cdot 2^1 + S_0 \cdot 2^0$$

$$b_3 S_2 S_1 S_0 \implies \left( \underline{b_3 (-1) \cdot 2^3} + S_2 \cdot 2^2 + S_1 \cdot 2^1 + S_0 \cdot 2^0 \right)_{\bmod 2^3} = S_2 \cdot 2^2 + S_1 \cdot 2^1 + S_0 \cdot 2^0$$

$$c = 1 \implies Overflow$$

We have 8 possible 3-bit patterns (or symbols).

How else might we assign an interpretation to them?

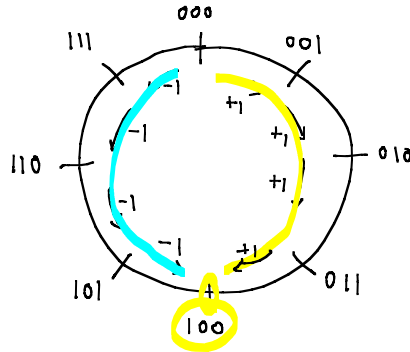What else might we want as number values?

| 3-bit Code | interpretation as value |
|---|---|
| 0 0 0 | 0 |
| 0 0 1 | 1 |
| 0 1 0 | 2 |
| 0 1 1 | 3 |
| 1 0 0 | 4 |
| 1 0 1 | 5 |
| 1 1 0 | 6 |
| 1 1 1 | 7 |

( NB — We are representing the values with
  an alternate representation : base 10
     Is there no end to this madness !?!  )

Two's-Complement Encoding:
We can represent POSTIVE and NEGATIVE

| CODE | Value |
|------|-------|
| 0 1 1 | +3 |
| 0 1 0 | +2 |
| 0 0 1 | +1 |
| **0 0 0** | 0 |
| 1 1 1 | -1 |
| 1 1 0 | -2 |
| 1 0 1 | -3 |
| 1 0 0 | -4 |

+1 +1 +1 (between top rows)
-1 -1 -1 -1 (between bottom rows)

sanity check: 0 - 1?

$$
\begin{array}{ccc}
 & 0 & 0 & 0 \\
- & 0 & 0 & 1 \\
\hline
(-1) & 1 & 1 & 1
\end{array}
$$

+10 +10 +10
-1 -1 -1

hmm, kind of makes sense.

moving +4 ≡ moving -4
Which value makes sense?

$\vec{2} - (\overleftarrow{+3})$

$\vec{2} + (\overrightarrow{-3})$

move +2 ↻
move +5 ↺
———
(-1)

-k can be moving:

↺ k steps

OR

↻ $2^n - k$ steps:

-3 ⟹ ↺ 3 steps

⟹ ↻ $2^3 - 3$ = 5 steps

n-bit Two's Comp (X)
———————————
$-X \implies 2^n - X$

-x    0    +x

⟹    $2^n$    0    -x    $(2^n - x)$

Sanity check  $-(-x)$ ?
———————————
$-x \implies (2^n - x)$

$-(-x) \implies 2^n - (2^n - x)$

$= x$

$-(-X) = X$ in 2's comp.

Try  n=3  2's comp : $2^n = 2^3 = 8$  (- (-3))

$2^3 - (2^3 - 3)$
$2^3 - (8 - 3)$
$2^3 - (5)$
$8 - 5$
$3$
⟹ 011   +3 (in 3-bit 2's comp.)

# Converting To 2's Comp? 
## n-bit

How do we do this simply, in general?

$$2^n = 1\ 0\ 0\ \cdots\ 0\ 0\ 0\ 0\ \cdots\ 0$$
(borrow: 10, 10, ... 10, 10; -1 -1 -1 ... -1 -1)
$$-x = -\ x_{n-1}\ x_{n-2}\ \cdots\ x_{n-j}\ 1\ 0\ 0\ \cdots\ 0$$
$$\overline{\hspace{6cm}}$$
$$S = \quad S_{n-1}\ S_{n-2}\ \cdots\ S_{n-j}\ 1\ 0\ 0\ \cdots\ 0$$

Notice: The 1st non-zero bit of $x$ gets copied to sum S:
```
        borrow   = 10
     subtract bit = - 1
     ----------------------
        sum bit  =    1
```

Note: columns with borrows give a bit flip.

$$\dfrac{1}{\dfrac{-x_k}{\overline{x_k}}}$$

$$2^n = 0\ 1\ 1\ \cdots\ 1\ 0\ 0\ 0\ \cdots\ 0 \quad (10)$$
$$-x = -\ x_{n-1}\ x_{n-2}\ \cdots\ x_{n-j}\ 1\ 0\ 0\ \cdots\ 0$$
$$\overline{\hspace{6cm}}$$
$$(2^n - x) = \overline{x}_{n-1}\ \overline{x}_{n-2}\ \cdots\ \overline{x}_{n-j}\ 1\ 0\ 0\ \cdots\ 0$$
$$(2^n - x) - 1 = \overline{x}_{n-1}\ \overline{x}_{n-2}\ \cdots\ \overline{x}_{n-j}\ 0\ 1\ 1\ \cdots\ 1 \quad -1 \quad +1$$

negate

$$x = x_{n-1}\ x_{n-2}\ \cdots\ x_{n-j}\ \cdots\ x_2\ x_1\ x_0$$

Notice: These are the negated bits of x.

To Get 2sComp( $x$ ):

Negate bits, add 1.

Produce $-x$ in 2's Complement (regardless of whether $x$ is + or -):

Negate bits (aka 1's Complement), then add 1.

Simple logic: inverter on each bit, carry in to lowest FA set to 1.

==> We can use adder for signed subtraction

Let's try
2's Comp of neg. number (expressed) in 2's comp.

$$2'sComp\ (\,1\ x_3\ x_2\ x_1\ x_0\,)$$
— a neg. number in 2's comp.

$$\rightarrow 0\ \overline{x}_3\ \overline{x}_2\ \overline{x}_1\ \overline{x}_0 + 1$$

flip bits, add 1

| Least neg. | between | most neg. |
|---|---|---|
| 1111 | 1 a b c | 1000 |
| flip ⟹ 0000 | 0 $\overline{a}$ $\overline{b}$ $\overline{c}$ | 0111 ⟸ flip |
| +1 ⟹ 0001 | 0 x y z | 1000 ⟸ +1 |
| | | ? |
| Extreme case ✓ok | in-between case ✓ok | Extreme case ✗ oops! what's wrong? |

$$A + 2sComp(B)$$

ADD: sub = 0
SUB: sub = 1



**IR**

| OP | |

Instruction decoder:
Each opcode
has its own
1-bit signal.

ADD
NAND
⋮ SUB $\begin{cases} 1, \text{ is SUB} \\ 0, \text{ is not SUB} \end{cases}$
⋮ BR

e.g., 3-bit (3-1)

$A = 3 \Rightarrow 011$
$B = 1 \Rightarrow 001$
$2sComp(B) \Rightarrow 110 + 1 = 111$

carries

$A \Rightarrow 1\ 1\ 1$
$1\ 0\ 1\ 1$
$+ (-B) \quad + 1\ 1\ 1$
$(1)\ 0\ 1\ 0$

carry ignored

Pos 2

# 2's complement Arith., Overflow



x > 0, y > 0
and x+y < 0

0xxx
0yyy
---------
1sss

x > 0, y < 0
and x-y < 0

2's comp

0xxx
- 1yyy => 0zzz *
----------
1sss

* could be 1000



x < 0, y < 0
and x+y > 0

1xxx
1yyy
----------
0sss

x < 0, y > 0
and x-y > 0

2's comp

1xxx
- 0yyy => 1zzz
----------
0sss

ERROR = Same signs in, diff. out

# Hex, Oct

binary: base = 2

digits = $\{0,1\}$

"$b_i b_{i-1} \cdots b_0$" $\xrightarrow{means}$ $b_i 2^i + b_{i-1} 2^{i-1} + \cdots + b_0 2^0$ $= n$

$100_2$ $\xrightarrow{means}$ $(1)\cdot 2^2 + (0)\cdot 2^1 + (0)\cdot 2^2 \Rightarrow 4_{10}$

octal: base = 8 ($= 2^3$)

"$d_i d_{i-1} \cdots d_0$" $\xrightarrow{means}$ $d_i (2^3)^i + d_{i-1}(2^3)^{i-1} + \cdots d_0 \cdot 1$

digits = $\{0,1,2,3,4,5,6,7\}$

$301_8$ $\xrightarrow{means}$ $(3)\cdot(2^3)^2 + (0)\cdot(2^3)^1 + (1)\cdot 1$

$3\cdot 64 + 0 + 1 \Rightarrow 193_{10}$

## Octal $\leftrightarrow$ binary

$501_8 \longrightarrow$ $(101_2)\cdot (2^3)^2 + 0\cdot(2^3)^1 + (001_2)(2^3)^0$

$\to$ $1\cdot 2^2 + 0\cdot 2^1 + 1\cdot 2^0 \ (2^3)^2 + 0\cdot(2^3)^1 + 1\cdot 2^0 (2^3)^0$

$\to$ $1\cdot 2^{2+6} + 0\cdot 2^{1+6} + 1\cdot 2^6 + 0\cdot 2^3 + 1\cdot 2^0$

$\to$ $1\cdot 2^8 + 0\cdot 2^7 + 1\cdot 2^6 + 0\cdot 2^5 + 0\cdot 2^4 + 0\cdot 2^3 + 0\cdot 2^2 + 0\cdot 2^1 + 1\cdot 2^0$

$\to$ $(\underline{1 \quad 0 \quad 1} \quad \underline{0 \quad 0 \quad 0} \quad \underline{0 \quad 0 \quad 1})_2$

$\to$ $(\quad 5 \qquad\qquad 0 \qquad\qquad 1 \quad)_8$

$\Rightarrow$ Octal digit $\leftrightarrow$ 3-bit binary representation of digit

## Hex

hexadecimal: base = 16 = $(2^4)$

| hex digits | bin rep. | hex digits | bin rep. |
|---|---|---|---|
| 0 | 0 0 0 0 | 8 | 1 0 0 0 |
| 1 | 0 0 0 1 | 9 | 1 0 0 1 |
| 2 | 0 0 1 0 | A | 1 0 1 0 |
| 3 | 0 0 1 1 | B | 1 0 1 1 |
| 4 | 0 1 0 0 | C | 1 1 0 0 |
| 5 | 0 1 0 1 | D | 1 1 0 1 |
| 6 | 0 1 1 0 | E | 1 1 1 0 |
| 7 | 0 1 1 1 | F | 1 1 1 1 |

$\Rightarrow$ Hex digit $\leftrightarrow$ 4-bit binary representation of digit

# Multiply

$$\begin{array}{r} 1 \\ \times 1 \\ \hline 1 \end{array} \qquad \begin{array}{r} 1 \\ \times 10 \\ \hline 10 \end{array} \qquad \begin{array}{r} 10 \\ \times 10 \\ \hline 100 \end{array} \qquad \begin{array}{r} 100 \\ \times 10 \\ \hline 1000 \end{array}$$

(does $2x$ = Left shift?)

$$2x = x + x \implies$$

$$\begin{array}{r} X_n \; X_{n-1} \cdots X_k \; 1 \; 0 \; 0 \ldots 0 \\ + \; X_n \; X_{n-1} \cdots X_k \; 1 \; 0 \; 0 \ldots 0 \\ \hline C_{n+1} \; S_n \; S_{n-1} \cdots S_k \; 0 \; 0 \; 0 \ldots 0 \end{array}$$

← 1ST 1

**Per Col.**

$$\begin{bmatrix} C_j \\ X_j \\ + X_j \\ \hline S_j \\ C_{j+1} \end{bmatrix} \implies \begin{bmatrix} C_j \\ 0 \\ + 0 \\ \hline S_j \\ C_{j+1} \end{bmatrix} \text{ OR } \begin{bmatrix} C_j \\ 1 \\ + 1 \\ \hline S_j \\ C_{j+1} \end{bmatrix} \longrightarrow \begin{matrix} S_j = C_j \\ C_{j+1} = X_j \end{matrix} \implies \begin{bmatrix} S_{j+1} = C_{j+1} \\ C_{j+1} = X_j \\ S_{j+1} = X_j \end{bmatrix} \checkmark$$

Left shift

In each column, if x = 0 then S is equal to the carry in.
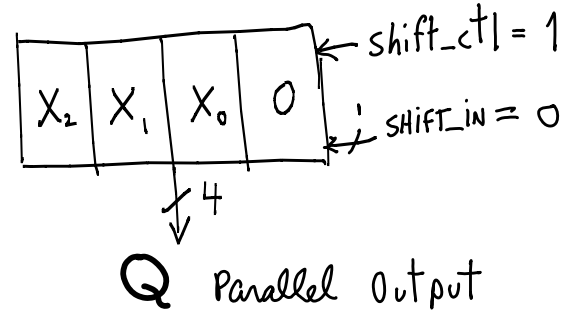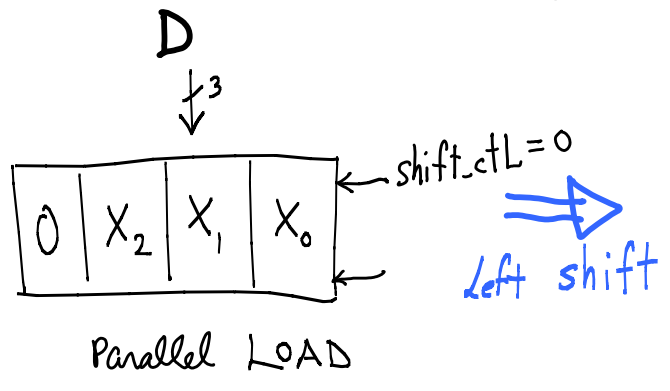If x = 1, there is a carry out and S is equal to the carry in.

## 3-bit, parallel load, left-shift register



Parallel write/load: we=1 and S=0:   Q[2:0]  <==   D[2:0]         after next clock tick.

Shift Left: we=1 and S=1:         Q[2:0]  <==  { Q[1:0], IN }      after next tick.

## 3-bit Doubler [using 4-bit LSR] (unsigned)

D
↓3

| 0 | $X_2$ | $X_1$ | $X_0$ |
|---|---|---|---|

← shift_ctL = 0

**Parallel LOAD**

⟹ *Left shift*

| $X_2$ | $X_1$ | $X_0$ | 0 |
|---|---|---|---|

← shift_ctl = 1
← SHIFT_IN = 0
↓4

**Q** Parallel output

What about signed numbers?

Convert to unsigned.

Multiply.

Convert back.

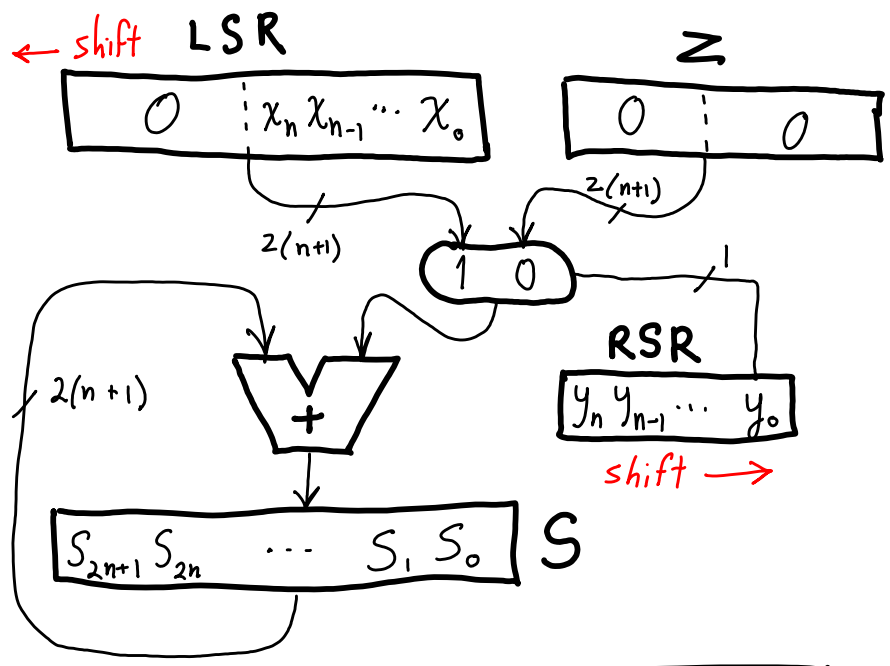$$x \cdot 7 = x(4 + 2 + 1) = 4x + 2x + x$$

$y = 00\ldots0111$

2 shifts (C)    1 shift (B)    0 shifts (A)

MULTIPLY:
LSR: partial products, initially x.
S:    partial sum, initially 0.
RSR: initially y.
Z:    all 0s

RSR's low-bit MUXes Z or LSR to adder.

$$
\begin{array}{r}
0 \quad 0 \quad \cdots \quad 0 \quad 0 \\
+ \quad x_n \; x_{n-1} \cdots x_1 \; x_0 \\
\hline
S'_n \; S'_{n-1} \cdots S'_1 \; S'_0
\end{array}
\qquad
\begin{array}{l}
= S^0 \\
+ (A) \\
\hline
S'
\end{array}
$$

shift

$$
\begin{array}{r}
+ \quad x_n \; x_{n-1} \cdots x_1 \; x_0 \; 0 \\
\hline
S''_{n+1} \; S''_n \cdots S''_2 \; S''_1 \; S''_0
\end{array}
\qquad
\begin{array}{l}
+ (B) \\
\hline
S''
\end{array}
$$

shift

$$
\begin{array}{r}
+ \quad x_n \; x_{n-1} \cdots x_1 \; x_0 \; 0 \; 0 \\
\hline
S'''_{n+2} \; S'''_{n+1} \cdots S'''_3 \; S'''_2 \; S'''_1 \; S'''_0
\end{array}
\qquad
\begin{array}{l}
+ (c) \\
\hline
S'''
\end{array}
$$

What if y has a 0 bit? Then add 0 instead of shifted x: e.g., y = 0...101 add 0, not **B**.



← shift    **LSR**     **Z**

LSR: $0 \mid x_n \; x_{n-1} \cdots x_0$    Z: $0 \mid 0$

$2(n+1)$

$1 \quad 0$    $2(n+1)$

$2(n+1)$    **+**    **RSR**

RSR: $y_n \; y_{n-1} \cdots y_0$

shift →

S: $S_{2n+1} \; S_{2n} \cdots S_1 \; S_0$

---

e.g.

$$
\begin{array}{r}
1\;0\;1\;1 \\
\times \; 0\;1\;0\;1 \\
\hline
\cdots \; 1\;0\;1\;1 \\
+ \; .\; 1\;0\;1\;1 \; . \; . \\
\hline
0 \; 0\;1\;1\;0\;1\;1\;1
\end{array}
$$

↗ (possible carry)

rewrite ⇒

multiplier bit

$$
\begin{array}{r}
1\;0\;1\;1 \\
\times \; 0\;1\;0\;1 \\
\hline
0\;0\;0\;0\;0\;0\;0\;0 \\
+ \; 0\;0\;0\;0\;1\;0\;1\;1 \\
+ \; 0\;0\;0\;0\;0\;0\;0\;0 \\
+ \; 0\;0\;1\;0\;1\;1\;0\;0 \\
+ \; 0\;0\;0\;0\;0\;0\;0\;0 \\
\hline
= 0\;0\;1\;1\;0\;1\;1\;1
\end{array}
$$

← start
← 0-shift (1011) ← $y_0 = 1$
← 1-shift ( ) ← $y_1 = 0$
← 2-shift ( ) ← $y_2 = 1$
← 3-shift ( ) ← $y_3 = 0$

⇒ We shift left (1011) every time, but add either the shifted (1011) or all zeroes, depending on whether $y_i$ is 1 or 0.

Can we simply multiplier?
-- Get rid of zero register and mux.
-- Use $y_i$ to write enable write-enable S register.

Can we speed up multiply? We currently iterate n times to multiply n-bit numbers. Add more hardware? How?
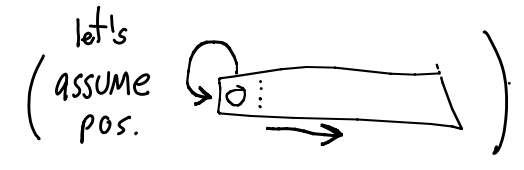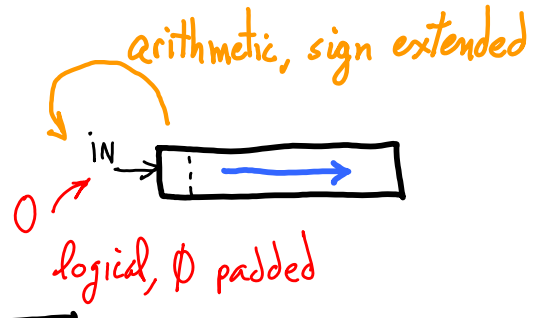


x      y

S

DIVIDE:   left-shift( y ) == y X 2    ===>    right-shift( y ) == y / 2          Ok, for division by a power of 2.

# Div by 2

### R-shift :

arithmetic, sign extended

iN

$0$

logical, $\emptyset$ padded

( let's assume pos. )

---

Integer Division
=
drop remainder

$011$  R-shift $\to 001$
$3$     $\div 2$  $= 1$

$0111$  (R-shift)$^2 \to 0001$
$7$       $\div 4$   $= 1$

---

R-shift(n) = divide-by-2^n

But, if divisor is not power of 2?

$$x = k \cdot q + r \begin{cases} k \text{ is divisor} \\ q \text{ is quotient} \end{cases}$$

$$q = \# k_s \text{ in } x$$

$$x \boxed{\; k \;\vdots\; k \;\vdots\; k \;\vdots\; k \;\vdots\; k \;\vdots\; r \;}$$

```
divBySubtraction( x, k )
  q = 0;
  while( x >= k )
    q++
    x = x - k
  endWhile
```

```
divByAddition( x, k )
  q = 0;   y = 0
  while ( x - y >= k )
    q++
    y = y + k
  endWhile
```

$$time = O(q)$$

---

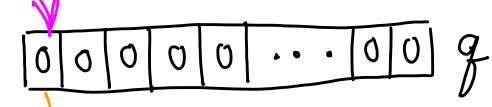We'd like $O(\log(q))$ = # bits of $q$

$\implies$ long division

1. Try largest power of 10, subtract from x.
2. If partial sum is non-negative, save digit.
3. Try next smaller power of 10, subtract from x or from remainder depending on prior result.
etc.

$$k \overline{) \begin{array}{l} d_n 00\cdots 0 \\ \overline{x} \\ -k \cdot d_n 00\cdots 0 \\ \hline x' \end{array}}$$

$$3 \overline{)\begin{array}{l} 50 \\ 176 \\ -150 \\ \hline 26 \end{array}} \quad (+) \to 58$$

$$3 \overline{)\begin{array}{l} 8 \\ 26 \\ -24 \\ \hline 2 \end{array}}$$

INTEGER (unsigned) DIVISON
x = kq +r      k = divisor,  q = quotient,  r = remainder  (let's ignore r for now).   FIND q.

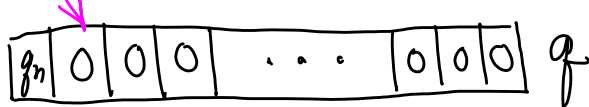$$x = k \cdot q = k \, q_n \, 2^n + k \, q_{n-1} \, 2^{n-1} \cdots + k \, q_0 \, 2^0$$

try $q_n = 1$

$$\left( x - k \cdot 1 \cdot 2^n \right) \geq 0 \quad \text{then} \quad q_n = 1$$

$k$  L-shifted n

$$x \leftarrow k \, q_{n-1} \, 2^{n-1} + \cdots + k \, q_0 \, 2^0 \qquad x \leftarrow (x - k \, q_n \, 2^n)$$

0 or 1

Try $q_{n-1} = 1$

$$\left( x - k \cdot 1 \cdot 2^{n-1} \right) \geq 0 \quad \text{then} \quad q_{n-1} = 1$$

$k$  L-shifted (n-1)

$2(n+1)$ - bit R-shift Reg

$\mid \leftarrow n \rightarrow \mid$

$0 \, k_n \cdots k_1 \, k_0 \, 0 \, 0 \cdots 0 \, 0$  K

we

SUB
$\widetilde{N}$    1

$q$

Move *notNegative*
RIGHT one bit position
after each SUB.

notNegative

$q$

n-bit L-shift Reg

Move **q** LEFT one bit
position after each SUB.

(write is to lowest bit
position)

# Floating Point



n-bit integers, range $= 2^n$
no discretization error

K-scaled integers: n-bit integer **x** represents k*x,  range = k*2^n.



this can't be represented, error $\approx k/2$

discretization error $\approx k/2$

Near 🔴 $k$, % error $\Rightarrow \frac{k/2}{k} = 50\%$

Near 🟢 $k(2^n-1)$ $\Rightarrow \frac{k/2}{k(2^n-1)} \approx 1/2^{n+1}$

---

# FP, exponential scaling
## $2^m(1.XYZ)$

Can we get more consistent errors?

values we can represent



$2^0(1.111)$

$2^0(1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8})$
$\Rightarrow \underbrace{1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8}}_{\approx 2} \pm (\frac{1}{16})$

$\Rightarrow$ error $\approx \frac{(\frac{1}{16})}{2} = \frac{1}{32}$

$2^1(1.111) = (2 + 1 + \frac{1}{2} + \frac{1}{4})$

We can't represent every number.
We choose what type of errors to live with.

The part inside "( ... )" is essentially integer.
The exponent determines the scaling.

===> geometrical-progression scaled integers

$2^2(1.111)$

$\underbrace{4 + 2 + 1 + \frac{1}{2}}_{\approx 8} \pm (\frac{1}{4})$

error $\approx \frac{(\frac{1}{4})}{8} = \frac{1}{32}$

---

# FP Format, single Float

32-bit:



$\text{value}(x) = S \, 2^E \times 1.f$

$S = 0 : +$
$S = 1 : -$

E pos.?
E neg.?

Use 2's comp for E?

not STORED

$+ \, 2^2 \times 1.11010...0$ $\Rightarrow + 2^2 \times (1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{16})$
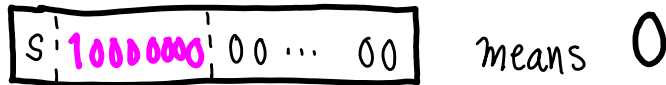
# Represent 0?



$$+2^0 \times 1.0 \cdots 0 \Rightarrow 1\,?$$

If we used 2s-Comp for E, what's the smallest: 8-bit, $1000\,0000 = -128$

$$2^{-128} \times 1.00\ldots0 \qquad \text{that's small. Do we need it?}$$

$$S\;\underline{1000\,0000}\;00\cdots00 \qquad \text{means } 0$$

Well, maybe we can live with that?
We have to stop somewhere.

---

How many bits do we need for 4 decimal digits of precision?

$$6.023 \times 10^{\overline{2}\,\overline{3}}$$

2 digits

a nice number

4 digits

3 bits per digit $(2^3 \to 0..7)$
4 bits per digit $(2^4 \to 0..15)$ $\quad\Big\}\quad \approx 3.5$ bits per digit

$$4 \text{ dig.} \times \frac{3\,\text{bit}}{\text{digit}} = 12 \text{ bits}$$
$$4 \text{ dig.} \times \frac{4\,\text{bit}}{\text{digit}} = 16 \text{ bits}$$

$$\Rightarrow \quad 12 \le \binom{\text{Precision of } f}{} \le 16$$

(23 bits are enough)

---

range of E?
(2 dec. digits)

$$2 \text{ dig.} \times \frac{3\,\text{bit}}{\text{dig}} = 6 \text{ bits}$$
$$2 \text{ dig.} \times \frac{4\,\text{bit}}{\text{dig.}} = 8 \text{ bits}$$

$$6 \le \binom{\text{bits of } E}{} \le 8$$

## Let's check

$$10^{23} \approx 8^{23} = (2^3)^{23} = 2^{\boxed{69}} \leftarrow E$$

how many bits
do we need for E?

$E = 69$, how
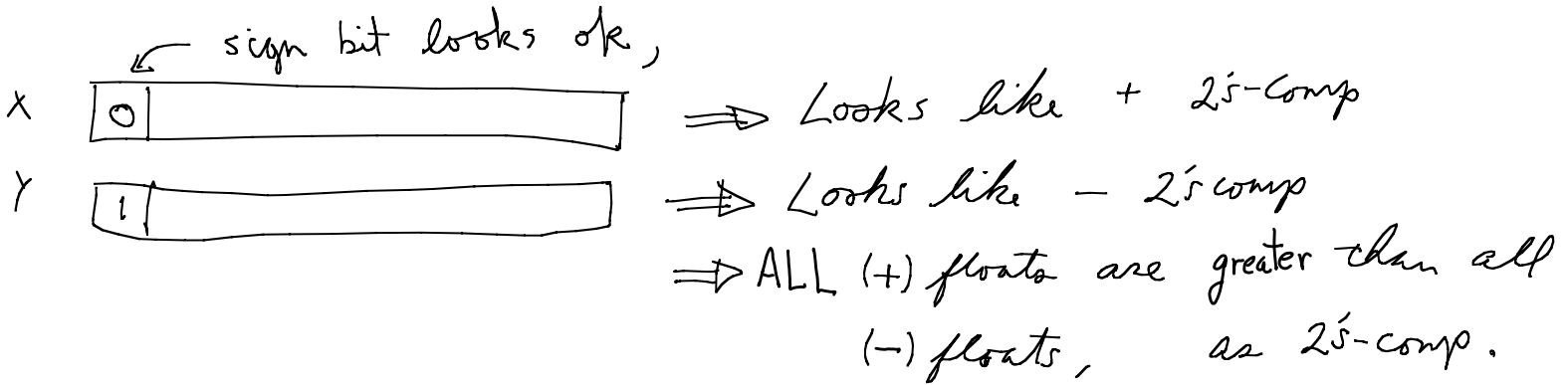many bits needed?

$$Log(69) \approx 1 + Log(64) = 7$$

Sorting is most common operation for numerical data
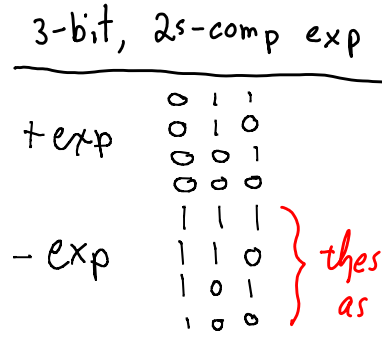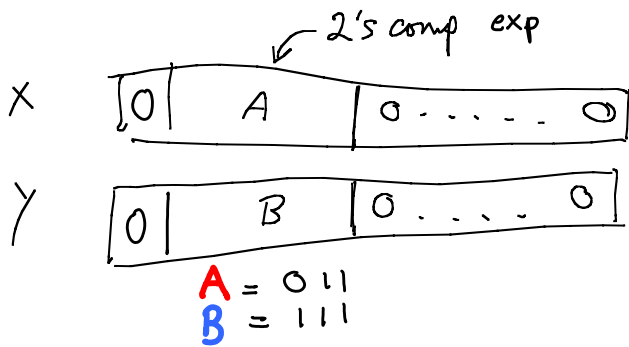
Checking x > y seems hard for floats.

Checking n > m for ints: do ( n - m ) and check sign bit, if 0 then True.

Can we check x > y using integer hardware?

That is, can we treat x and y as if they were integers, and do integer subtraction?

sign bit looks ok,

X [0 _____ ]   $\Longrightarrow$ Looks like + 2's-comp

Y [1 _____ ]   $\Longrightarrow$ Looks like − 2's comp

$\Longrightarrow$ ALL (+) floats are greater than all
(−) floats,      as 2's-comp.

How about the exponent part? x−y as 2's-comp

2's comp exp

X [0 | A | 0 · · · · · · 0 ]

Y [0 | B | 0 · · · · · 0 ]

A = 011
B = 111

3-bit, 2's-comp exp

+ exp
```
0 1 1
0 1 0
0 0 1
0 0 0
```

− exp
```
1 1 1
1 1 0   } these are larger
1 0 1   } as unsigned, darn.
1 0 0
```

$X = +2^{3} \cdot (1.00 \cdots 0)$
$y = +2^{-1} \cdot (1.0 \cdots 0)$  $\Big\} \Rightarrow$  $(0 \ 011 \ 00 \cdots 0) \ x$
$(0 \ 111 \ 00 \cdots 0) \ y$

y looks like a
bigger 2's comp. number than X.
**O O P S !**

$\longrightarrow$ Let's see if we can patch this up.
Recall, our only problem is if both x and y
have the same sign.

suppose sign(x) = sign(y)
_____
(mag. comparison)

X:  | 0 | $E_2$ | $f_2$ |

Y:  | 0 | $E_1$ | $f_1$ |

or

X:  | 1 | $E_2$ | $f_2$ |

Y:  | 1 | $E_1$ | $f_1$ |

(Reverse result for neg.)
(for signs ≠ result is obvious)

let $e_i$ = value($E_i$)

Note: $e_1 > e_2 \implies 2^{e_1}(1.f_1) > 2^{e_2}(1.f_2)$

regardless of the fractional parts.
_____

Let's check

Suppose $e_2 = e_1 + 1$                 $1.11\cdots 1$

$\underbrace{2^{e_2} \cdot (1.0)}_{\substack{\text{Smallest} \\ \text{possible} \\ x}}$     $\underbrace{2^{e_1}(2-\varepsilon)}_{\substack{\text{Largest} \\ \text{possible } y}} = 2^{e_1+1}(1-\varepsilon/2)$

$= 2^{e_2}(1-\varepsilon/2)$

Y is less than X

So, make $E_1$ look bigger than $E_2$

_____

what we have so far

8-bit exponent (single float)
_____

2's complement ⟶  ⟹ { 
0111 1111   $(2^7 - 1 = 127)$
     ⋮
0000 0000  (0)
1111 1111
1111 1110
     ⋮
1000 0001  (-127, Not used, reserved for signal)
1000 0000  (-128, Not used, signals 0 ?)

Let's
fix Exp

E.G., 3-bit exponents in 2s-complement

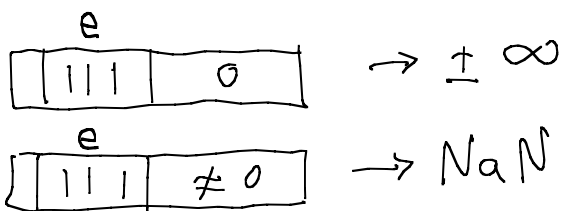goal: make all negative exponents look smaller than all positive exponents AS unsigned ints.

```
| S | E |   f   |
```

$$E = e + 011$$

| value | e in 2s-comp |   | E in excess-3 |   |
|-------|--------------|---|---------------|---|
| +3 | 011 | +011 ==> | 110 |   |
| +2 | 010 | +011 ==> | 101 |   |
| +1 | 001 | +011 ==> | 100 |   |
| 0 | 000 | +011 ==> | 011 |   |
| -1 | 111 | +011 ==> | 010 |   |
| -2 | 110 | +011 ==> | 001 |   |
| -3 | 101 | +011 ==> | 000 | * |
| -4 | 100 | +011 ==> | 111 | * |

\* These codes are reserved for special uses.
  The exponent values -3 and -4 are not allowed.



Rotate 2's comp so that +3 becomes largest number available.

How to represent 0?

```
e
| 111 |   0   |  →  ± ∞
e
| 111 |  ≠ 0  |  →  NaN
```

```
e
| 0 0 0 |   0   |  →  0   (±0)
e
| 0 0 0 |  ≠ 0  |  →  not normalized
                      →  2⁻³ × 0.f
```

$$2^{-3} \times 0.f$$

---

## 8-bit FP, ADD

X
```
| S | E    | f       |
| 0 | 0 0 1 | 0 1 1 0 |
```

y
```
| 0 | 1 1 0 | 0 0 1 0 |
|   S | E   | f       |
```

$001 - 011 = 110 \rightarrow -2$

Convert from excess-3 to 2's comp.

$110 - 011 = 011 \Rightarrow +3$

```
| 0 | 1 1 0 | 0 1 | 1 0 |
```

```
| 0 | 0 1 1 | 0 0 1 0 |
```

$+2^{-2} \times 1.0110$
$+2^{3} \times 1.0010$

shift/align exponents

$2^{3} \times 0.00001 0110$
$2^{3} \times 1.0010$

TRUNCATION?

add f part

$$2^3 \times 1.001010$$

round to nearest ← rounding

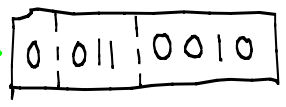$$2^3 \times 1.0010 \rightarrow \text{normalize} \rightarrow 2^3 \times 1.0010 \rightarrow \text{encode}$$

no shift needed in this example

| 0 | 110 | 0010 | ← To excess-3 ← | 0 | 011 | 0010 |

$\underbrace{\phantom{0110 0010}}_{y}$

$$x + y = y \ !?! \ \cdots \Rightarrow \text{discretization, truncation, rounding}$$
$$\rightarrow ERRORS \quad \text{be careful!}$$

# 8-bit FP, MULT

| 0 | 001 | 1010 | → | 0 | 110 | 1010 | → $+2^{-2} \times 1.1010$

| 0 | 110 | 0000 | → | 0 | 011 | 0110 | → $+2^{3} \times 1.0110$

Convert $E_s$ from excess-3 To 2sComp

Shift and MULT as unsigned INTs. Keep exp. To normalize

add exponents
check for overflow

$2^1$

add exponents

$$2^2 \times 1.00011110 \leftarrow$$

round to nearest

$$2^2 \times 1.0010 \rightarrow \text{encode}$$

normalize

$$2^1 \times 1.00011100$$

| 0 | 010 | 0010 |

excess 3 code

| 0 | 101 | 0010 |

$11010 \ (\times 2^{-4})$
$\times 10110 \ (\times 2^{-4})$
$\underline{\phantom{xxxxxxx}}$
$+ \ 110100$
$+ \ 11010$
$+ \ 110100$
$\underline{\phantom{xxxxxxx}}$
$1 00011 1100 \ (\times 2^{-8})$

# Convert to 32-bit FP ⟩ 28

## 1. Convert to binary

$$
\begin{array}{r}
28 \\
-16 \\ \hline
12 \\
-8 \\ \hline
4
\end{array}
\quad\longrightarrow\quad
\begin{aligned}
1 \cdot 2^4 &= 10000 \\
+ \\
1 \cdot 2^3 &= 1000 \\
+ \\
1 \cdot 2^2 &= 100 \\ \hline
&\phantom{=}11100
\end{aligned}
$$

## 2. Normalize

$$
11100. \quad\Rightarrow\quad 2^4 \times 1.1100
$$

$e = 4$

## 3. Encode

8-bits   23-bits

$$
\boxed{0 \mid 0000\,0100 \mid 1100...0}
$$

$+$    $e = 4$    $1.1100$

## 4. Convert $e$ to excess $(2^{n-1}-1)$

$$\text{excess}\left(2^{8-1}-1 = 127\right)$$

$$
\begin{array}{r}
00000100 \;= e \\
+\; 01111111 \;= 127 \\ \hline
10000011 \;= E
\end{array}
$$

## 5. Replace $e$ with $E$

$$
\boxed{0 \mid 10000011 \mid 1100 \; ... \; 0}
$$

## Convert back

1. decode:  $+\, 2^{10000011} \times 1.1100...0$

2. convert $E$

← borrows

$$
\begin{array}{r}
\llap{^{-1\,1\,1\,1\,1\,1}}1000011 \\
-\; 0111111 \\ \hline
0000100 \;= 4 = e
\end{array}
$$

$$
\Rightarrow \quad 2^4 \times 1.11
\quad
\left(\begin{array}{l}\text{convert } f \\ \text{mult. by } 2^e\end{array}\right)
\Rightarrow \quad 1.1100 = 11100
$$

$\times 2^4$

(convert to dec.) $\Rightarrow$   $16 + 8 + 4 = 28$

So much for encoding data. We could go on to audio, video, ...  But, back to noise and errors.

# Error Detection / Correction

"message" could be a bit, a string of bits, a character, a page of characters, ...

"message" → | Communication Channel | → "mfssage"
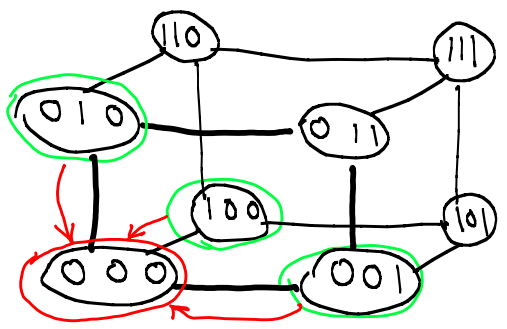
message coded in bits:

"1" → | encode | → 1 → | Channel | → 0 → | decode | → "0"

**1-bit Error** — not detectable

2-bit encoding

"1" → | encode | → 11 → | Channel | → 01 → | decode | →

**Error detected**

Code words 00 and 11 are good data, 10 and 01 indicate 1-bit errors. Last bit is "parity" bit, odd parity codeword indicates error. Works for k-bit messages w/ 1 parity bit (if 2-bit errors very unlikely).

## Parity scheme | data | p | Parity bit

1-bit Error Correction w/ 3-bit code words:
"0" ==> 000
"1" ==> 111

| | |
|---|---|
| 001 ==> "0" | 011 ==> "1" |
| 010 ==> "0" | 101 ==> "1" |
| 100 ==> "0" | 110 ==> "1" |

| k-bits | P-bits |
|--------|--------|
| data | Parity |

P-bit parity check

1-bit Correction, 2-bit Detection

-- odd parity: 1-bit error corrected

-- exactly two 1's: 2-bit error detected

-- otherwise: no error

How many extra bits, at minimum?
Depends on noise in channel:
Shannon Noisey Coding Theorem.

We use 4 bits, 1-bit data.

2-bit error detected
1100

1-bit correction
1000
1010
1-bit error, correction
1110

0100
1001
1101

0010
0110
1011
1111

0000
0001
0101
0111

0011

1-bit error

1-bit error

# Hamming (7,4) Code (Single Error Detection / Single Error Correction)

### 7 bits per code word
4 data bits
3 parity bits



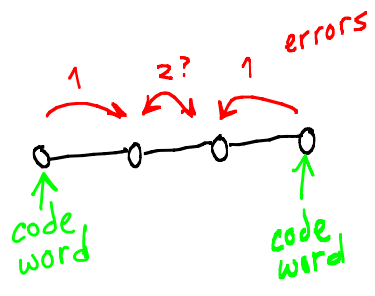$d_1 \, d_2 \, d_3 \, d_4$ $P_1$ $P_2$ $P_3$

$P_1 = parity(d_1 \, d_2 \, d_4)$

$P_2 = parity(d_1 \, d_3 \, d_4)$

$P_3 = parity(d_2 \, d_3 \, d_4)$

## 3 steps to next codeword

1-bit error: can detect and correct

2-bit error: cannot detect



What can we do about 2-bit errors?
Add another parity bit.

$P_4 = parity(d_1 \, d_2 \, d_3 \, d_4 \, p_1 \, p_2 \, p_3)$

Code word $= d_1 \, d_2 \, d_3 \, d_4 \, p_1 \, p_2 \, p_3 \, p_4$

$\Rightarrow$ 4 steps min to next code word

1-bit error: detect + correct
2-bit error: detect

other neighbor code words

0001 111
0010 011
0111 001

Hamming 7,4 code:
Find distances to all other code words.
**GREEN-PARITY: Bits[ 3, 2,    0 ]**
**BLUE-PARYT:    Bits[ 3,    1, 0 ]**
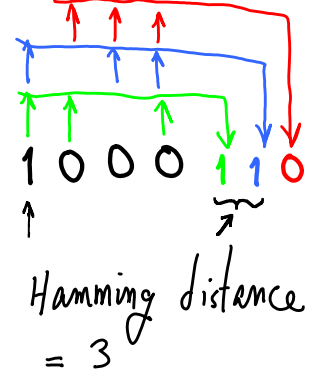**RED-PARITY:    Bits[    2, 1, 0 ]**

O O O O O O

O O O 1 1 1
( 1-bit flip causes
3 other flips
Hamming distance
= 4 )

O O 1 O O 1 1
Hamming distance
= 3

O 1 O O 1 O 1
Hamming distance
= 3

1 O O O 1 1 O
Hamming distance
= 3

O O 1 1 1 O O
HD = 3

O 1 O 1 O 1 O
HD = 3

1 O O 1 O O 1
HD = 3

O 1 1 O 1 1 O
HD = 4

1 O 1 O 1 O 1
HD = 4

1 1 O O O 1 1
HD = 4

O 1 1 1 O O 1
HD = 4

1 O 1 1 O 1 O
HD = 4

1 1 O 1 1 O O
HD = 4

1 1 1 O O O O
HD = 3

1 1 1 1 1 1 1
HD = 7

ASCII (See back cover of PP)

| HEX CODE | MEANING | Printable? |
|----------|---------|------------|
| 00 | NUL | no |
| 01 | SOH | no |
| ... | ... | |
| 20 | space | yes |
| ... | ... | |
| 30 | "0" | yes |
| 31 | "1" | yes |
| ... | ... | |
| 41 | "A" | yes |
| 42 | "B" | yes |
| ... | ... | |
| 61 | "a" | yes |
| 62 | "b" | yes |
| ... | ... | |
| 7A | "z" | yes |
| ... | ... | |

(other stuff, non-standard)

data
} Communications
Control signals

## Who's on first?

6D   41   2F   32

↓    ↓    ↓    ↓

'm'  'A'  '/'  '2'

## Byte Addressable

print order

| addr | memory bits | |
|------|-------------|---|
| 0 | 0011 0010 | → 32 ("2") |
| 1 | 0010 1111 | → 2F ("/") |
| 2 | 0100 0001 | → 41 ("A") |
| 3 | 0110 1101 | → 6D ("m") |

← low bit

high bit →

| What to Print | Starting Memory Address | What is displayed (left-to-right) |
|---------------|-------------------------|-----------------------------------|
| 4-byte number (in hex notation) | 0 | 6D412F32 |
| two 2-byte numbers (in hex) | 0 | 2F32   6D41 |
| four 1-byte numbers (in hex) | 0 | 32  2F  41  6D |
| one 4-byte string | 0 | 2  /  A  m |

(see "od" in unix)