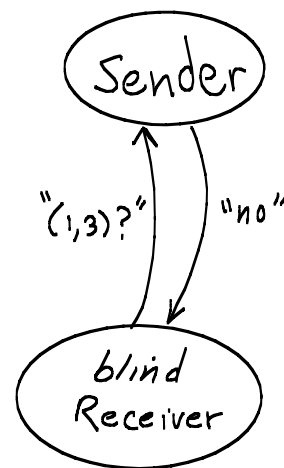


where's the ball?

	1	2	3	4
1				
2				
3		●		
4				



• 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)?

$$P(\text{Hit 1st}) = 1/16$$

$$\Rightarrow P(\text{Hit 2nd}) = P(\text{Hit 2nd} \mid \text{Miss 1st})P(\text{Miss 1st}) = (1/15)(15/16) = 1/16$$

$$P(\text{Hit 3rd}) = (1/14) * P(\text{Miss 2nd and 1st}) = (1/14)(14/15)(15/16) = 1/16$$

$$E(n) = 1*(1/16) + 2*(1/16) + \dots + 15*(1/16) + 15(1/16) = (1+2+3+\dots+15+15) / 16 = 8 \frac{1}{2} - 1/16$$

• Different set of question?

"in (1,*) or (2,*)" , "in (1,*)" , "in (2,1) or (2,2)" ...

\Rightarrow Each Q reduces space by $\frac{1}{2} \Rightarrow$ Exactly 4 questions

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

Measure info content of answer?

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

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

$$= -1 \quad (\text{hmm, make +?})$$

$$\Rightarrow 1 \text{ bit}$$

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

$$\text{Prob}(\text{yes}) = \frac{1}{4} \quad \left. \vphantom{\text{Prob}(\text{yes})} \right\} \text{assume}$$

$$\text{Prob}(\text{no}) = \frac{3}{4}$$

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

↑
Q = "Left of here?"

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

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


Avg info of answer?

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

$$= \frac{1}{4}(2) + \frac{3}{4}\left(\frac{1}{2}\right) = \frac{1}{2} + \frac{3}{8} = \frac{7}{8} \text{ bit}$$

Extreme: $\text{Prob}(\text{yes}) = \frac{1}{2^{10}}$ $\text{Prob}(\text{no}) = \frac{(2^{10}-1)}{2^{10}} \approx 1$

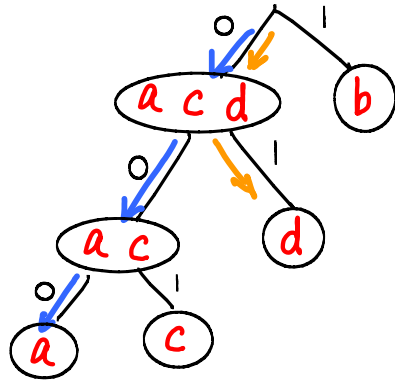
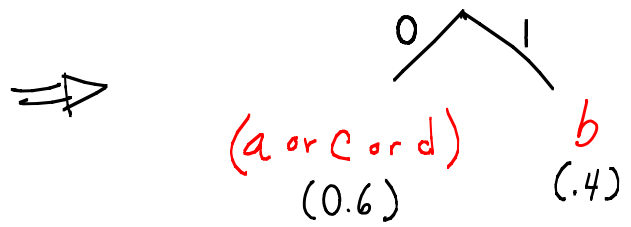
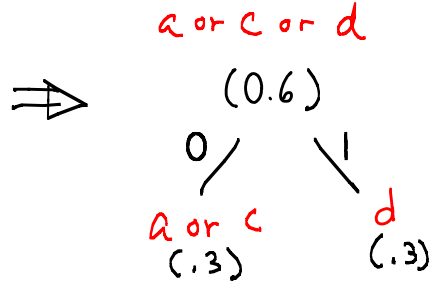
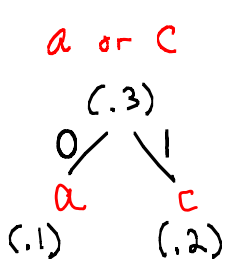
$$\text{avg} = \frac{1}{2^{10}}(10 \text{ bits}) + (1) \log(1) \Rightarrow \frac{1}{100} \text{ bit}$$



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

How to encode to get 50-50? \Rightarrow (Pair least Probable) \sim equally likely.

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



Huffman Code	msg
000	a
010	c
011	d
1	b

Shannon information = (Expected (Avg.) bits of information, per message) = $E(-\log \text{Prob})$

$$\begin{aligned}
 (\text{Entropy, } H) &= - \left[\underset{a}{0.1 \log(0.1)} + \underset{b}{0.4 \log(0.4)} + \underset{c}{0.2 \log(0.2)} + \underset{d}{0.3 \log(0.3)} \right] \\
 &= \left[0.33 + 0.53 + 0.46 + 0.52 \right] = 1.84 \text{ bits per message.}
 \end{aligned}$$

How'd we do?

Avg. #bits sent, using our code.	=	$(0.1) \underset{a}{3}$	+	$(0.4) \underset{b}{1}$	+	$(0.2) \underset{c}{3}$	+	$(0.3) \underset{d}{2}$	
	=	0.3	+	0.4	+	0.6	+	0.6	= 1.9 bits per message.

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: $\text{Prob}(\text{message-}i) = (1/n)$ for any message- i .

Expected information per message is,

$$\text{Sum} \left[- (1/n) \log \left[1/n \right] \right] = - n \left(1/n \log \left[1/n \right] \right) = -1 \log \left[2^{-k} \right] = -1(-k) = k \text{ bits per message.}$$

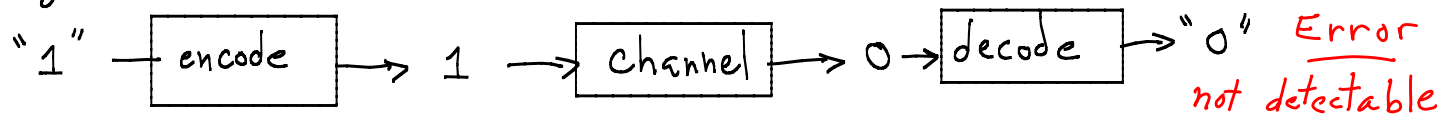
If we use a k-bit code for our messages, we will be 100% compressed. (k-bit integers? Are they equally likely?)

Error Detection / Correction

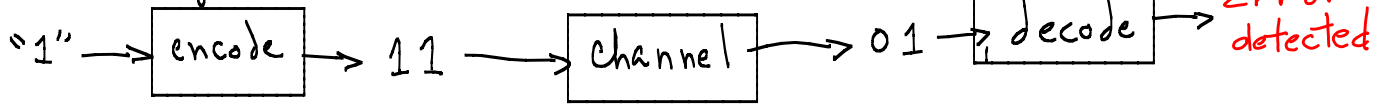
"message" could be a bit, a string of bits, a character, a page of characters, ...



message coded in bits:



2-bit encoding



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"

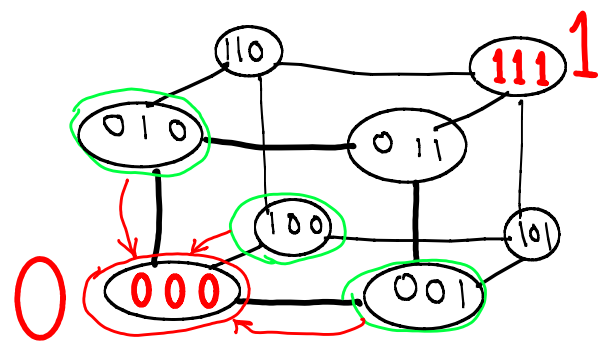
011 ==> "1"

010 ==> "0"

101 ==> "1"

100 ==> "0"

110 ==> "1"



1-bit Correction, 2-bit Detection

-- odd parity: 1-bit error corrected

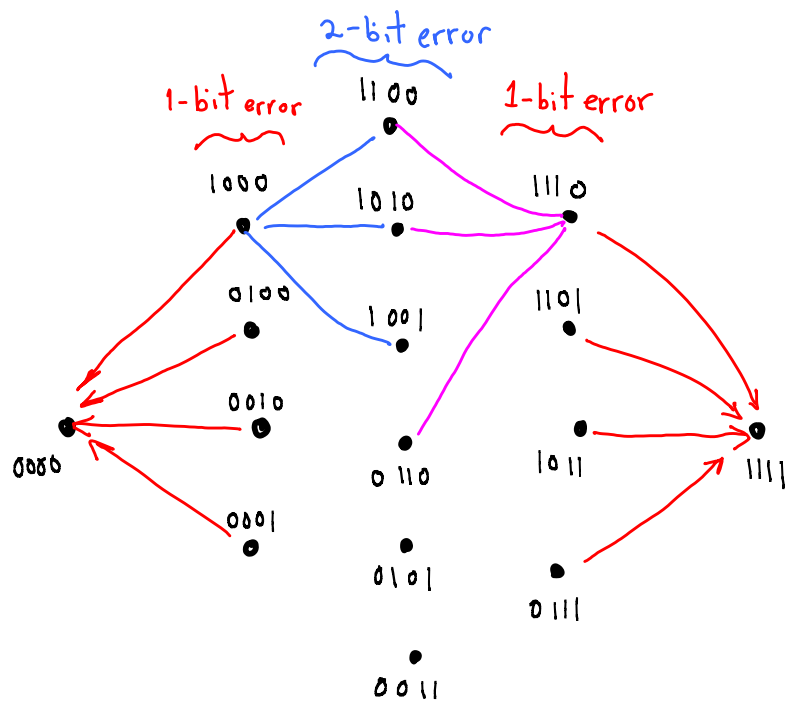
-- exactly two 1's: 2-bit error detected

-- otherwise: no error

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

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

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

unsigned 3-bit binary
binary:
--- digits = {"0", "1"}
--- base = 2

$$\begin{aligned} 000 &\rightarrow value \triangleq 0 \cdot 2^2 + 0 \cdot 2^1 + 0 \cdot 2^0 = 0 \\ 001 &\rightarrow value \triangleq 0 \cdot 2^2 + 0 \cdot 2^1 + 1 \cdot 2^0 = 1 \\ \dots &\dots \dots \dots \dots \dots \dots \\ 111 &\rightarrow value \triangleq 1 \cdot 2^2 + 1 \cdot 2^1 + 1 \cdot 2^0 = 7 \end{aligned}$$

oops, more symbols?

Let's do some arithmetic.

ADD:

$$\begin{array}{r} 0 \\ +0 \\ \hline 0 \end{array} \quad \begin{array}{r} 0 \\ +1 \\ \hline 1 \end{array} \quad \begin{array}{r} 1 \\ +0 \\ \hline 0 \end{array} \quad \begin{array}{r} 1 \\ +1 \\ \hline 10 \end{array}$$

↑ carry

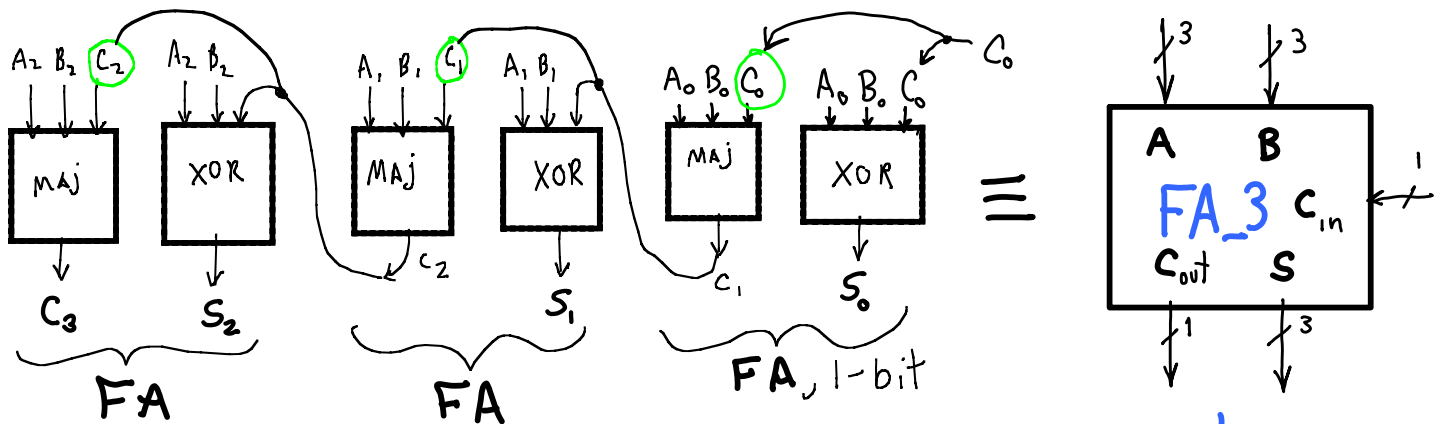
$$\begin{array}{c|c|c} A & B & S \\ \hline 0 & 0 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{array} \quad \left. \vphantom{\begin{array}{c|c|c} A & B & S \\ \hline 0 & 0 & 0 \\ 0 & 1 & 1 \\ 1 & 0 & 1 \\ 1 & 1 & 0 \end{array}} \right\} \text{XOR}$$

w/ Carry In:

$$\begin{array}{r} 1 \\ +0 \\ \hline 1 \end{array} \quad \begin{array}{r} 1 \\ +1 \\ \hline 10 \end{array} \quad \begin{array}{r} 1 \\ +0 \\ \hline 10 \end{array} \quad \begin{array}{r} 1 \\ +1 \\ \hline 11 \end{array}$$

$$\begin{array}{c|c|c} A & B & S \\ \hline 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \end{array} \quad \left. \vphantom{\begin{array}{c|c|c} A & B & S \\ \hline 0 & 0 & 1 \\ 0 & 1 & 0 \\ 1 & 0 & 0 \\ 1 & 1 & 1 \end{array}} \right\} \overline{\text{XOR}}$$

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



3-bit FULL Adder

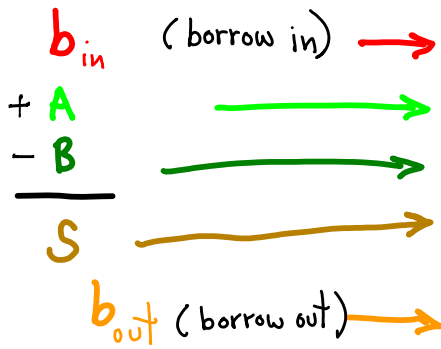
Let's try
SUBTRACTION

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

possible borrow

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

All possible 1-bit subtractions

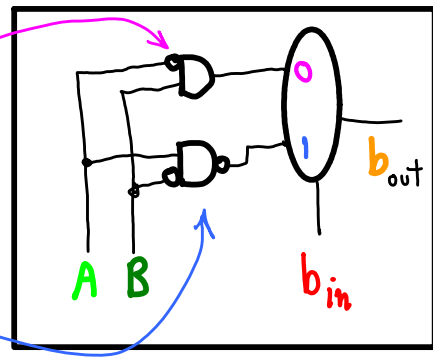


0	0	0	0	-1	-1	-1	-1
0	0	1	1	0	0	1	1
0	-1	0	-1	0	-1	0	-1
0	1	1	0	1	0	0	1
0	1	0	0	1	1	0	1

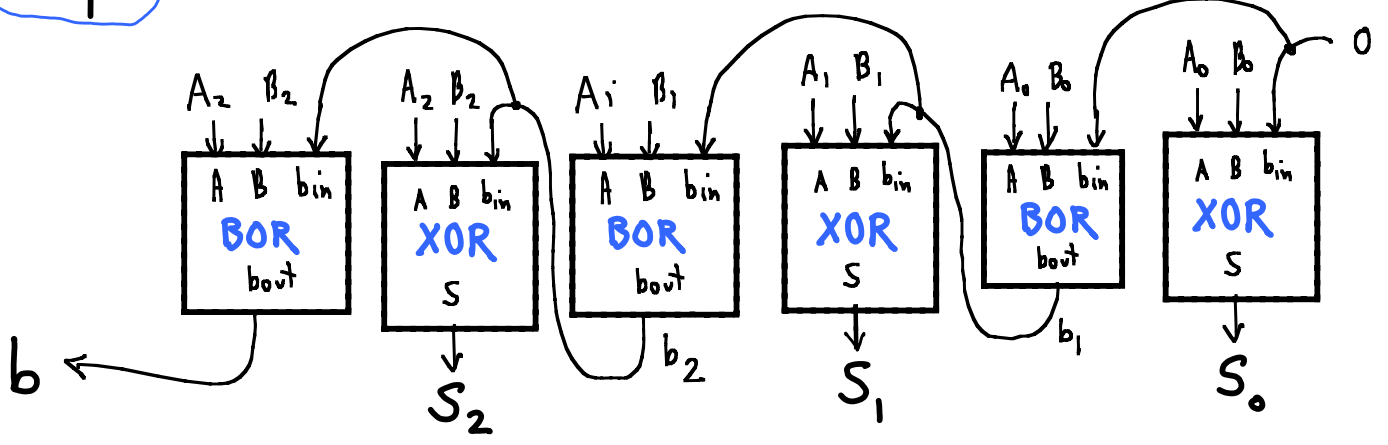
$= \text{XOR}(A, B, b_{in})$
 $= \text{BOR}(A, B, b_{in})$

b_{in}	A	B	b_{out}
0	0	0	0
0	0	1	1
0	1	0	0
0	1	1	0
1	0	0	1
1	0	1	1
1	1	0	0
1	1	1	1

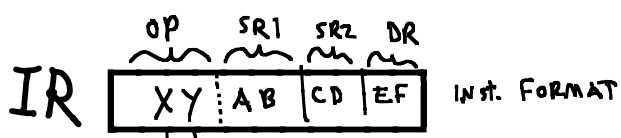
$(\bar{A}_i \cdot B_i)$
 $(A_i \cdot \bar{B}_i)$



BOR

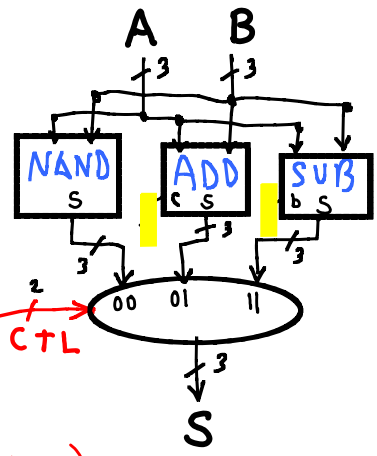


ALU

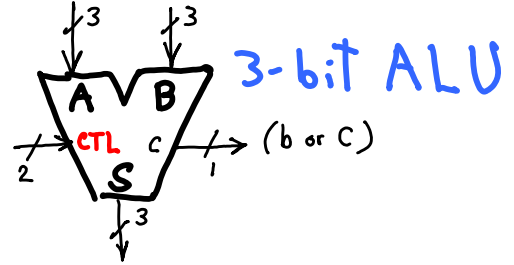
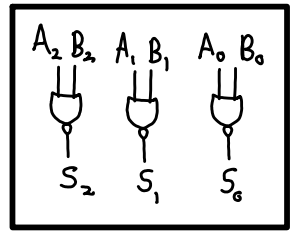


OP CODES

00	NAND
01	ADD
11	SUB



3-bit, bit-wise NAND

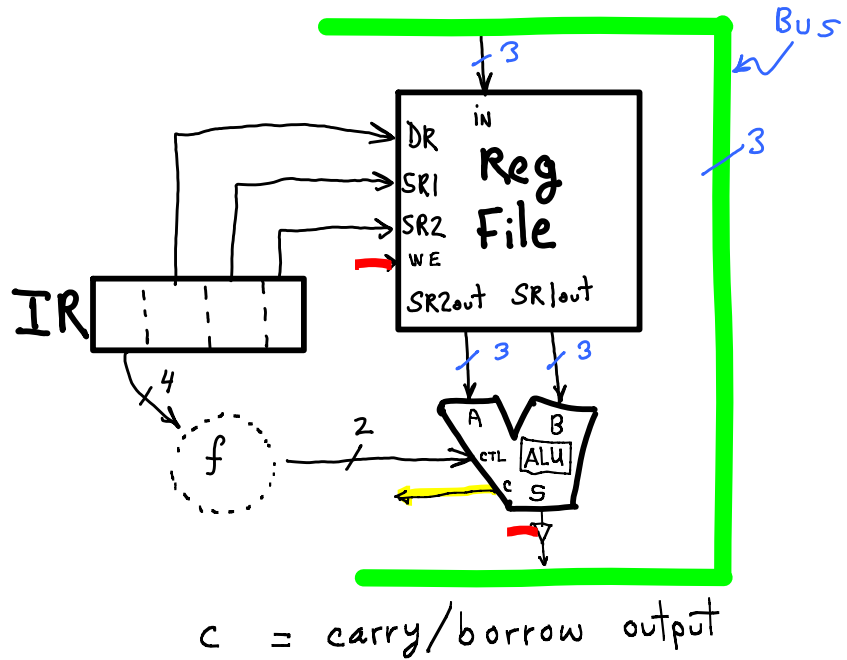


(simplified instruction, only shows ALU ops)

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 control bits in ROM.



Unsigned Arith. Errors

3-bit numbers

$$A + B > 7 \implies c = 1$$

$$A - B < 0 \implies b = 1$$

$$S = (A + B) \bmod 2^3 \leq 7$$

$$S = (A - B) \bmod 2^3 \geq 0$$

$$\left(\underline{c_3} \cdot 2^3 + s_2 \cdot 2^2 + s_1 \cdot 2^1 + s_0 \cdot 2^0 \right) \bmod 2^3 = s_2 s_1 s_0$$

$$\left(\underline{b_3}(-1) \cdot 2^3 + s_2 \cdot 2^2 + s_1 \cdot 2^1 + s_0 \cdot 2^0 \right) \bmod 2^3 = s_2 s_1 s_0$$

$$b = c = 1 \implies \text{Overflow Error}$$

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

Choose an interpretation.

3-bit Code	interpretation as value
000	0
001	1
010	2
011	3
100	4
101	5
110	6
111	7

NB--We represent the value using another encoding: base 10!

What other number values are we interested in?

Are other encodings useful?

Scaled Numbers

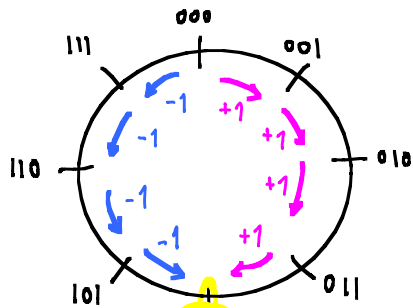
3-bit Code	interpretation as value
000	0
001	8
010	16
011	24
100	32
101	40
110	48
111	56

Sign-magnitude

3-bit Code	interpretation as value
000	+0
001	+1
010	+2
011	+3
100	-0
101	-1
110	-2
111	-3

2's-Complement Encoding:
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

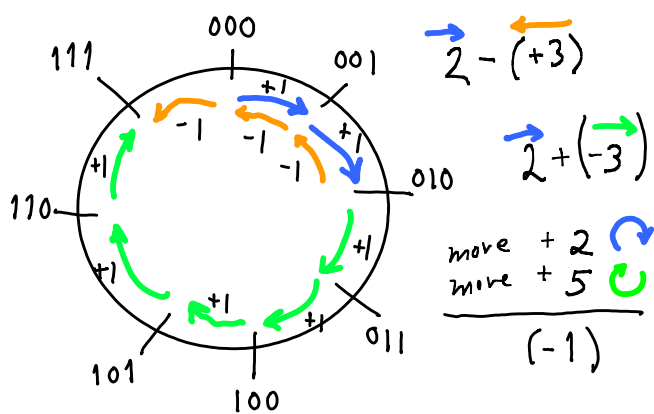


moving +4 \equiv moving -4
Which value makes sense?

sanity check: 0 - 1?

$$\begin{array}{r}
 \begin{array}{ccc}
 \overset{+10}{-1} & \overset{+10}{-1} & \overset{+10}{-1} \\
 0 & 0 & 0 \\
 -0 & 0 & 1 \\
 \hline
 (-1) & 1 & 1 & 1
 \end{array}
 \end{array}$$

hmm, kind of makes sense.



-k can be moving:

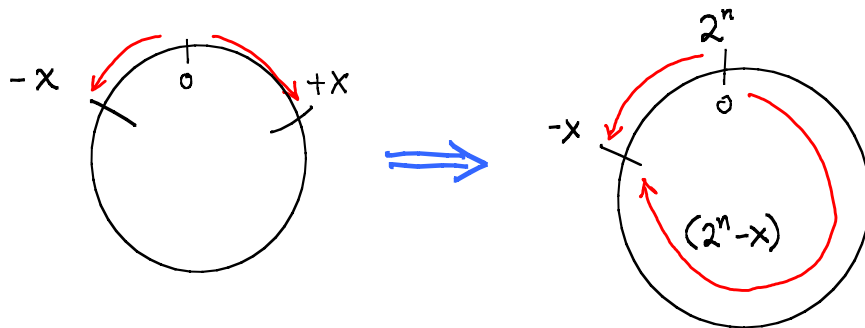
- \circlearrowright k steps
- OR \circlearrowleft $2^n - k$ steps:

-3 \Rightarrow \circlearrowright 3 steps

\Rightarrow \circlearrowleft $2^3 - 3 = 5$ steps

n-bit 2's Complement

$-x \rightarrow 2^n - x$



Sanity check $-(-x)$?

$$\begin{aligned}
 -x &\Rightarrow (2^n - x) \\
 -(-x) &\Rightarrow 2^n - (2^n - x) \\
 &= x \\
 -(-x) &= x \\
 &\text{in 2's comp.}
 \end{aligned}$$

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

$(-3)_{2's\ comp} \rightarrow 2^n - 3 = 8 - 3 = 5$

$(-(-3)_{2's\ comp})_{2's\ comp} \rightarrow 2^n - 5 = 8 - 5 = +3$

Converting To 2's Comp?

n-bit

How do we do this simply, in general?

$$\begin{array}{r}
 2^n = 1 \ 0 \ 0 \ \dots \ 0 \ 0 \ 0 \ 0 \ \dots \ 0 \\
 -x = -x_{n-1} \ x_{n-2} \ \dots \ x_{n-j} \ 1 \ 0 \ 0 \ \dots \ 0 \\
 \hline
 S = S_{n-1} \ S_{n-2} \ \dots \ S_{n-j} \ 1 \ 0 \ 0 \ \dots \ 0
 \end{array}$$

borrow (with arrows pointing from right to left)

Note: columns with borrows give a bit flip.

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

 sum bit = 1

$$\begin{array}{r}
 1 \\
 -x_k \\
 \hline
 \bar{x}_k
 \end{array}$$

$$\begin{array}{r}
 2^n = 0 \ 1 \ 1 \ \dots \ 1 \ 0 \ 0 \ 0 \ \dots \ 0 \\
 -x = -x_{n-1} \ x_{n-2} \ \dots \ x_{n-j} \ 1 \ 0 \ 0 \ \dots \ 0
 \end{array}$$

$$\begin{array}{r}
 (2^n - x) = \bar{x}_{n-1} \ \bar{x}_{n-2} \ \dots \ \bar{x}_{n-j} \ 1 \ 0 \ 0 \ \dots \ 0 \\
 (2^n - x) - 1 = \bar{x}_{n-1} \ \bar{x}_{n-2} \ \dots \ \bar{x}_{n-j} \ 0 \ 1 \ 1 \ \dots \ 1
 \end{array}$$

-1 (green arrow pointing to the 1 in the second row)
+1 (blue arrow pointing to the 1 in the second row)

NB--These are the **negated bits** of x .

To Get 2sComp(x):
 Negate bits,
 add 1.

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

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's Comp ($1x_3x_2x_1x_0$)

a neg. number in 2's comp.

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

flip bits, add 1

	<u>Least neg.</u>		<u>between</u>		<u>most neg.</u>
	1111		1abc		1000
flip =>	0000		0 $\bar{a}\bar{b}\bar{c}$		0111 <i>← flip</i>
+1 =>	0001		0xyz		1000 <i>← +1</i>

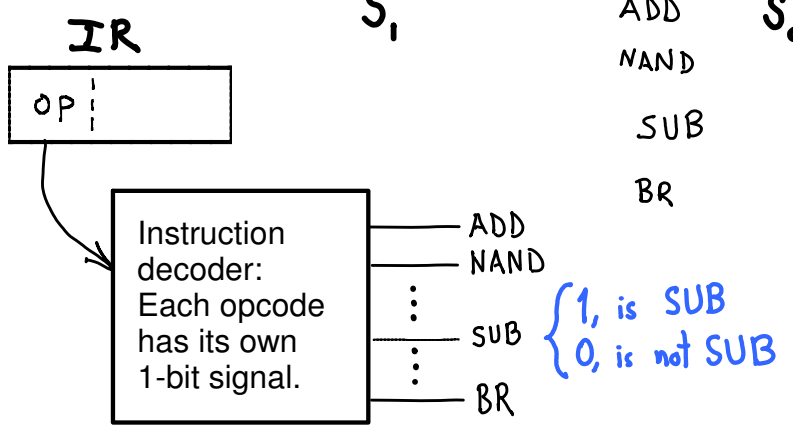
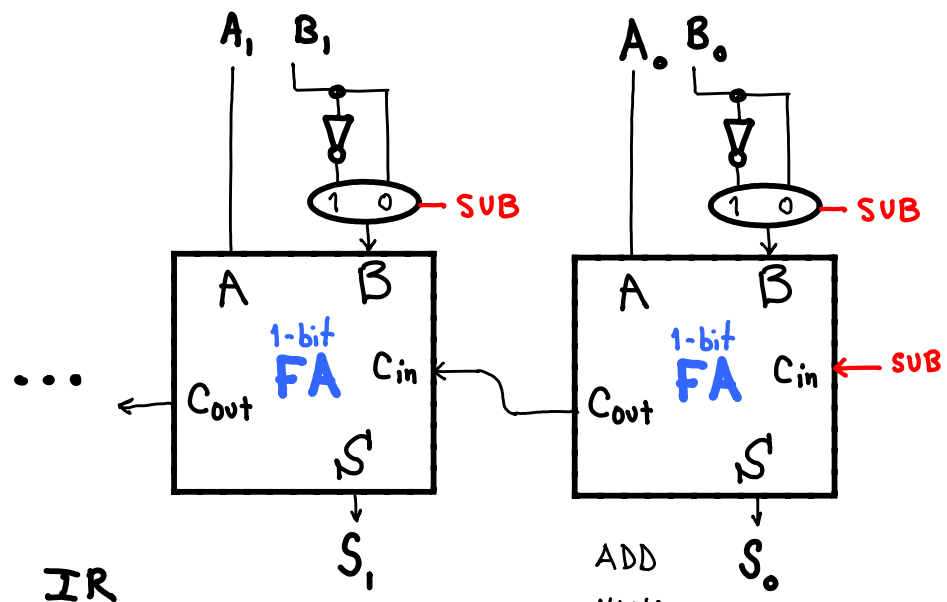
Extreme case **✓ok**

in-between case **✓ok**

Extreme case
~~oops! what's wrong?~~

$$A + 2s\text{Comp}(B)$$

$$\text{sub} = \begin{cases} 0, & \text{ADD} \\ 1, & \text{SUB} \end{cases}$$

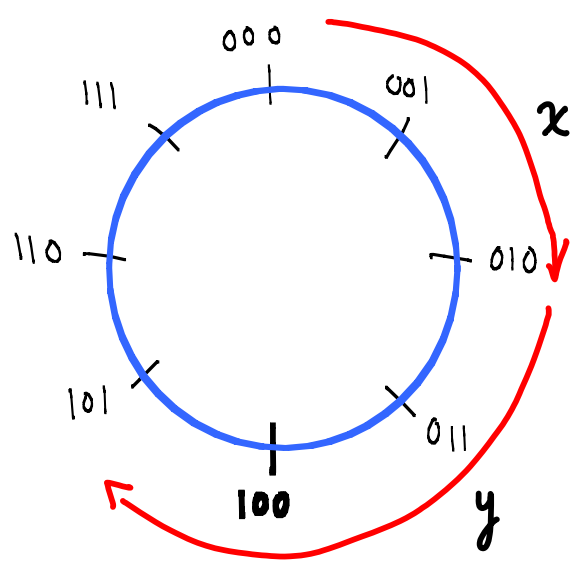


e.g., 3-bit (3-1)

A = 3 \Rightarrow 011
B = 1 \Rightarrow 001
 $2s\text{Comp}(B) \Rightarrow 110 + 1 = 111$

$$\begin{array}{r} \overset{\text{carries}}{1} \overset{1}{\uparrow} \overset{1}{\uparrow} \overset{1}{\uparrow} \\ 011 \text{ (A)} \\ + 111 \text{ (-B)} \\ \hline (1) 010 \\ \text{Pos} \quad 2 \end{array}$$

2's complement Arith., Overflow

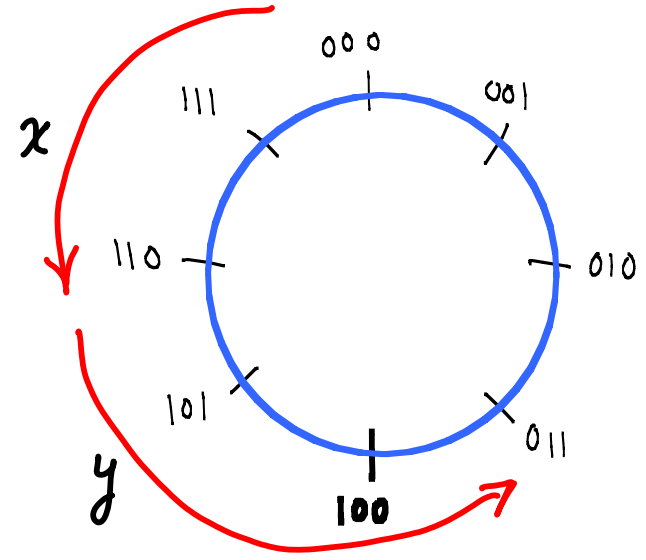


$x > 0, y > 0$
and $x+y < 0$

$$\begin{array}{r} 0xxx \\ 0yyy \\ \hline 1sss \end{array}$$

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

$$\begin{array}{r} 0xxx \\ - 1yyy \Rightarrow 0zzz^* \\ \hline 1sss \end{array}$$



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

$$\begin{array}{r} 1xxx \\ 1yyy \\ \hline 0sss \end{array}$$

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

$$\begin{array}{r} 1xxx \\ - 0yyy \Rightarrow 1zzz \\ \hline 0sss \end{array}$$

2's comp

2's comp

* could be 1000

ERROR = Same signs in, diff. out

Multiply

does $2x$
= Left shift?

$$\frac{1}{x \times 10} = \frac{10}{100}$$

$$\frac{10}{x \times 10} = \frac{100}{1000}$$

$$\frac{100}{x \times 10} = \frac{1000}{10000}$$

$$2x = x + x \Rightarrow$$

$$\begin{array}{r} x_n x_{n-1} \dots x_k \quad \overset{1^{st}}{1} 00 \dots 0 \\ + x_n x_{n-1} \dots x_k \quad 1 00 \dots 0 \\ \hline C_{n+1} S_n S_{n-1} \dots S_k \quad 0 0 0 \dots 0 \end{array}$$

Col. j

$x_j = 0$

$x_j = 1$

$$\begin{array}{r} C_j \\ x_j \\ + x_j \\ \hline C_{j+1} S_j \end{array}$$

$$\begin{array}{r} C_j \\ 0 \\ + 0 \\ \hline C_{j+1} S_j \end{array}$$

$$\begin{array}{r} C_j \\ 1 \\ + 1 \\ \hline C_{j+1} S_j \end{array}$$

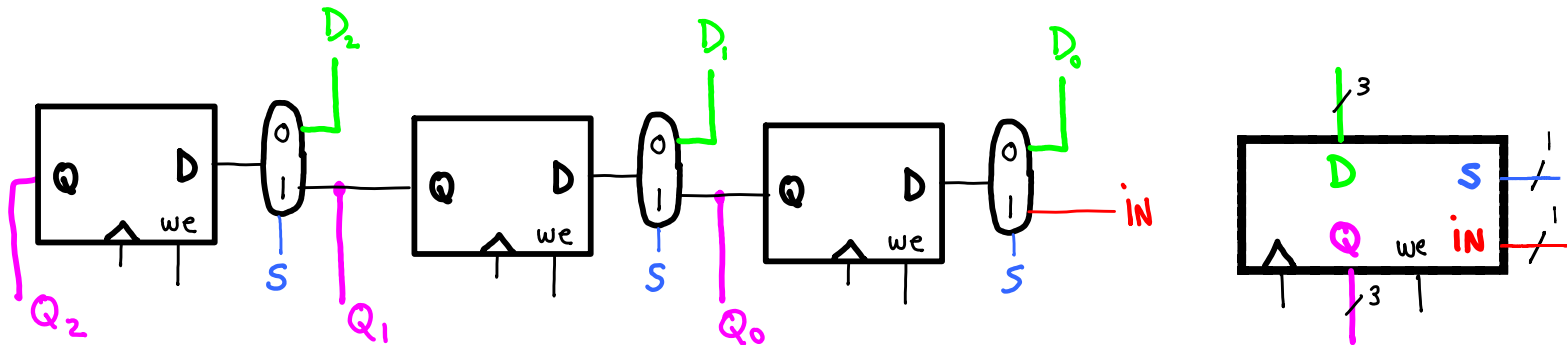
$$\left\{ \begin{array}{l} S_j = C_j \\ C_{j+1} = x_j \end{array} \right\}$$

$$S_{j+1} = C_{j+1} = x_j$$

left shift

$x == 0$: $S ==$ carry in; carry out $== 0$.
 $x == 1$: $S ==$ carry in; carry out $== 1$.

3-bit, parallel load, left-shift



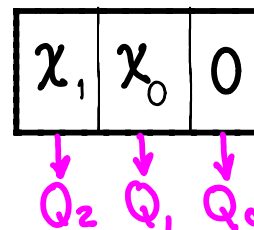
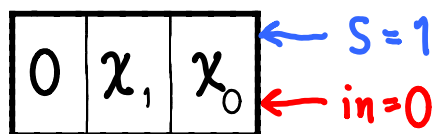
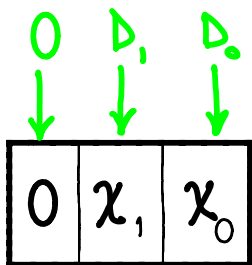
Parallel write/load : $we=1$ and $S=0$: $Q[2:0] \leftarrow D[2:0]$

after next clock tick.

Shift Left : $we=1$ and $S=1$: $Q[2:0] \leftarrow \{ Q[1:0], IN \}$

after next tick.

Multiplier (by 2)



What about signed numbers?

Convert to unsigned.

Multiply.

Convert back.

$$x \times 7 == x_n x_{n-1} \dots x_1 x_0 \times (4 + 2 + 1)$$

==

$$\begin{aligned}
 & x_n x_{n-1} \dots x_1 x_0 \times (1) \leftarrow 0 \text{ left-shift} \\
 + & x_n x_{n-1} \dots x_1 x_0 \times (2) \leftarrow 1 \text{ left-shift} \\
 + & x_n x_{n-1} \dots x_1 x_0 \times (4) \leftarrow 2 \text{ left-shifts}
 \end{aligned}$$

==

$$\begin{aligned}
 & x_n x_{n-1} \dots x_1 x_0 \\
 + & x_n x_{n-1} \dots x_1 x_0 0 \\
 + & x_n x_{n-1} \dots x_1 x_0 0 0
 \end{aligned}$$

MULTIPLY:

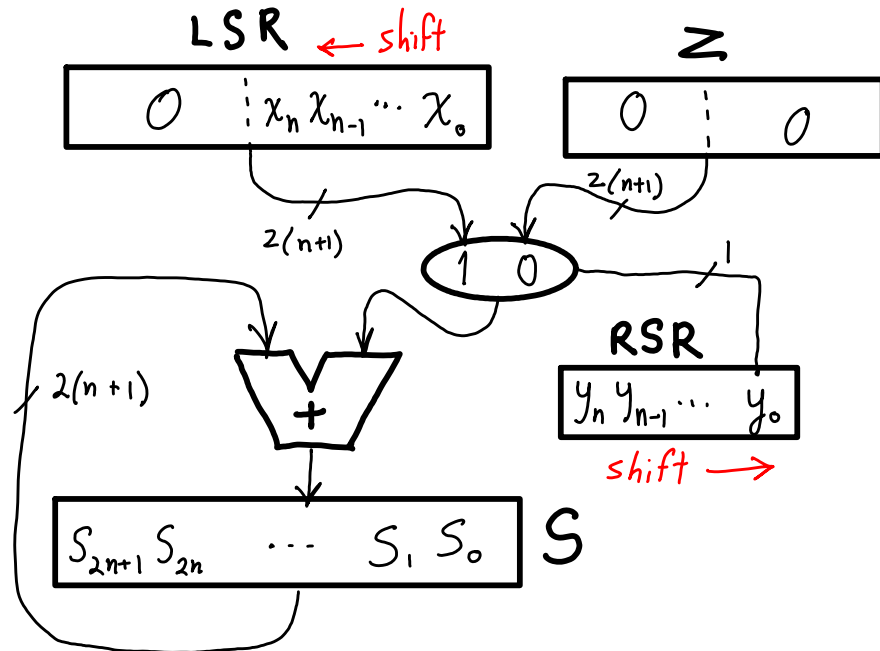
LSR: partial products, initially x.

S: partial sum, initially 0.

RSR: initially y.

Z: all 0s

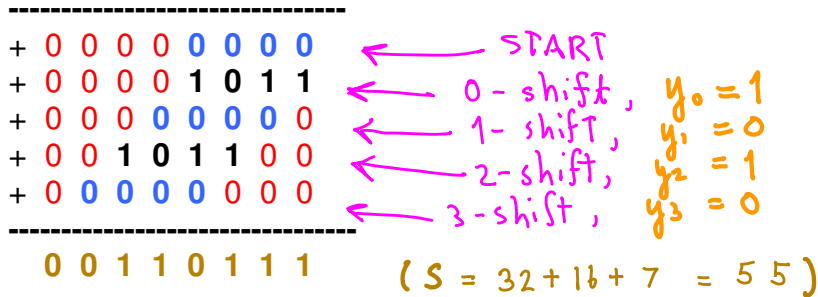
RSR's low-bit MUXes Z or LSR to adder.



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

e.g. 1011 × 0101

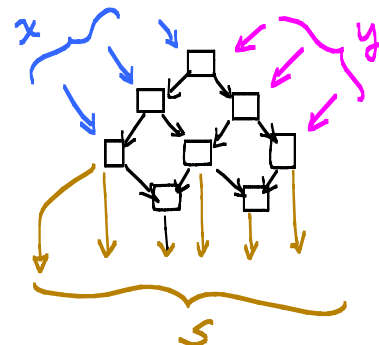
1 0 1 1 (x = 11)
 x 0 1 0 1 (y = 5)



Can we simplify multiplier?

- Get rid of zero register and mux.
- Use y_i to write-enable S register.

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

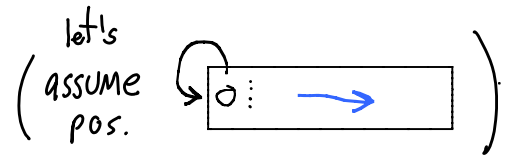
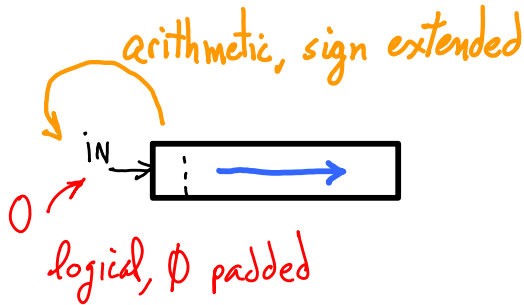


left-shift(y) == y X 2 ==>>> right-shift(y) == y / 2

Ok, for division by a power of 2.

Div by 2

R-shift:



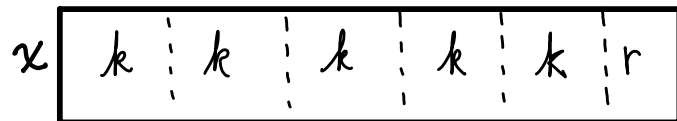
Integer Division = drop remainder

011 R-shift → 001
3 ÷ 2 = 1

0111 (R-shift)² → 0001
7 ÷ 4 = 1

$$x = k \cdot q + r \quad \left\{ \begin{array}{l} k \text{ is divisor} \\ q \text{ is quotient} \end{array} \right.$$

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



R-shift(n) = divide-by-2ⁿ

If divisor is not power of 2?

divBySubtraction(x, k)

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

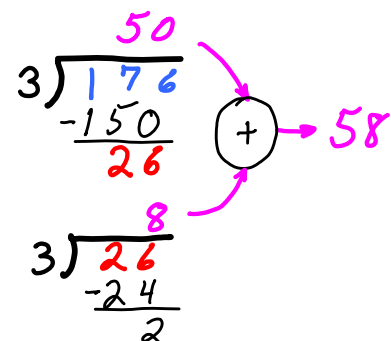
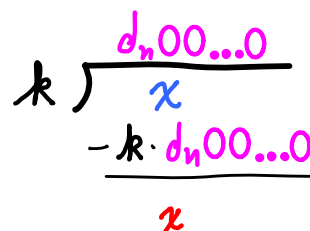
divByAddition(x, k)

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

$$\text{time} = O(q)$$

We'd like $O(\log q) = \# \text{ bits of } q \implies \text{long division}$

1. Try n -th power of 10, $d_n 00 \dots 0$
subtract: $x - k \times d_n 00 \dots 0$
2. result x non-negative?
yes: save $d_n 00 \dots 0$, $x \lll x$
3. next power of 10: $n \lll n-1$



INTEGER (unsigned) DIVISION

$$x = kq + r \quad k = \text{divisor}, \quad q = \text{quotient}, \quad r = \text{remainder (ignore for now)}. \quad \text{FIND } q.$$

$$x = kq = kq_n 2^n + kq_{n-1} 2^{n-1} \dots + kq_0 2^0$$

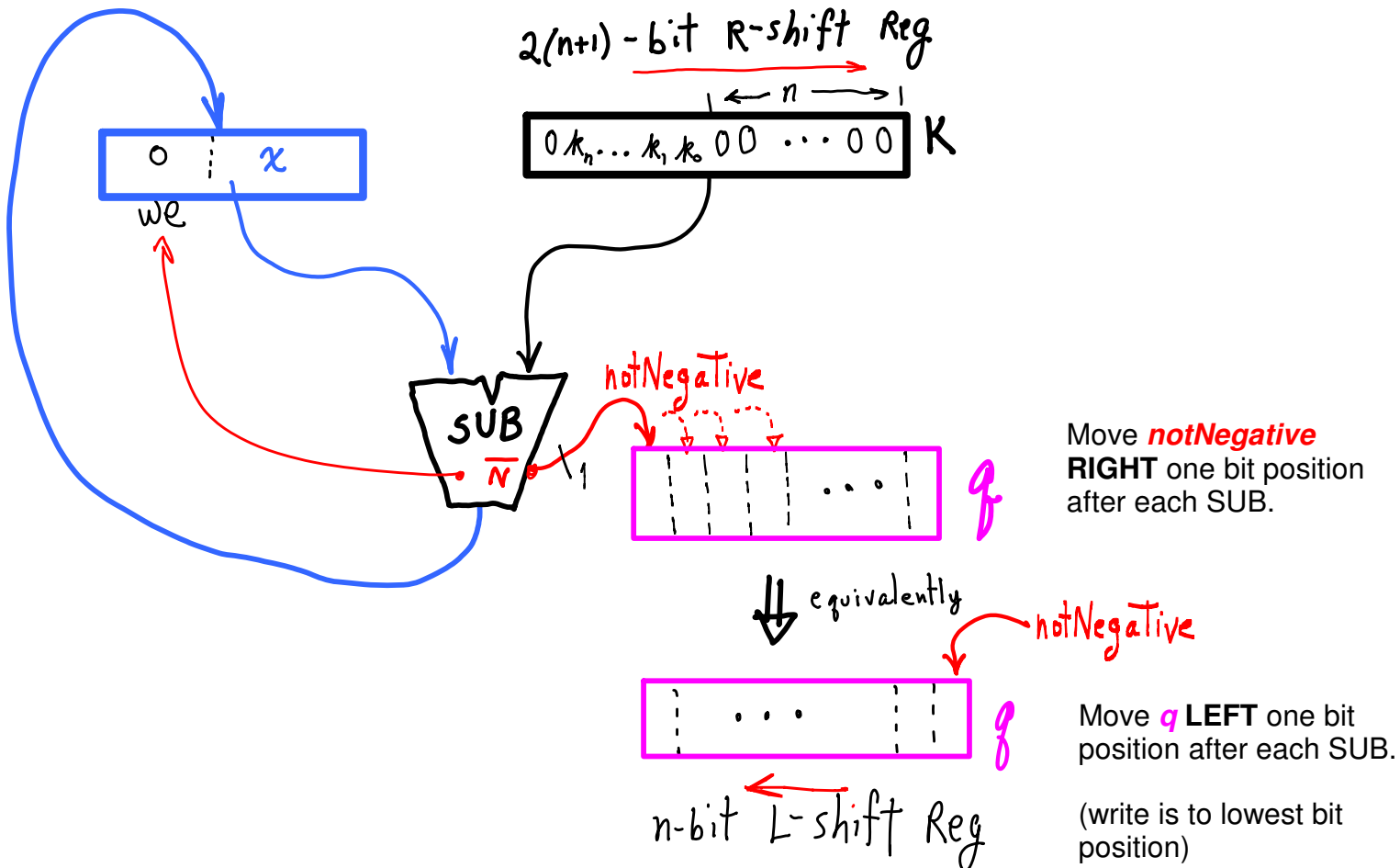
try $q_n = 1$

$$(x - k \cdot 1 \cdot 2^n) \geq 0 ? \quad \begin{cases} \text{yes: } q_n = 1 \\ \text{no: } q_n = 0 \end{cases}$$

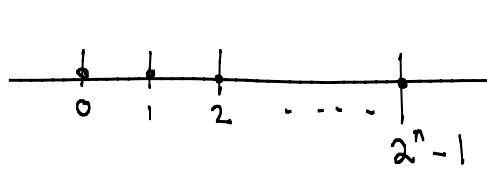
$$x \leftarrow kq_{n-1} 2^{n-1} + \dots + kq_0 2^0 \iff x \leftarrow x - (kq_n 2^n)$$

try $q_{n-1} = 1$

$$(x - k \cdot 1 \cdot 2^{n-1}) \geq 0 ? \quad \begin{cases} \text{yes: } 1 \\ \text{no: } 0 \end{cases}$$

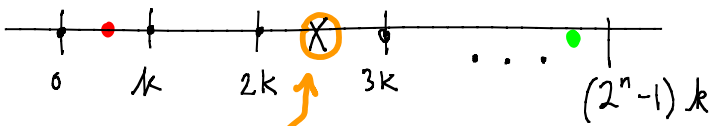


Floating Point



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

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

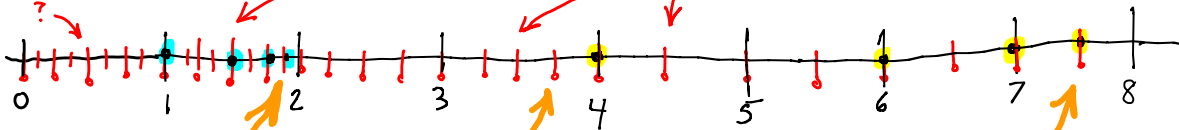


this can't be represented, error $\approx k/2$

discretization error $\approx k/2$
Near $\bullet k$, % error $\Rightarrow \frac{k/2}{k} = 50\%$
Near $\bullet k(2^n - 1)$ $\Rightarrow \frac{k/2}{k(2^n - 1)} \approx \frac{1}{2^{n+1}}$

FP, exponential scaling

$$2^m(1.XYZ)$$



Can we get more consistent errors?

$$2^0(1.111)$$

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

$$2^2(1.111)$$

$$2^0(1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8})$$

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

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

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

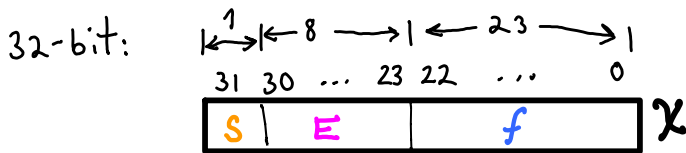
The part inside "(...)" is essentially integer.
The exponent determines the scaling.

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

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

====> geometrical-progression scaled integers

FP Format, single float



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

$s = 0 : +$
 $s = 1 : -$

E pos.?
E neg.?

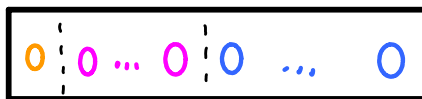
Use 2's comp for E?

not STORED

$$+ 2^2 \times 1.11010\dots 0$$

$$\Rightarrow + 2^2 \times (1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{16})$$

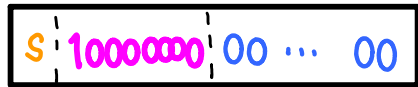
Represent 0?



$$+ 2^0 \times 1.0 \dots 0 \Rightarrow 1?$$

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

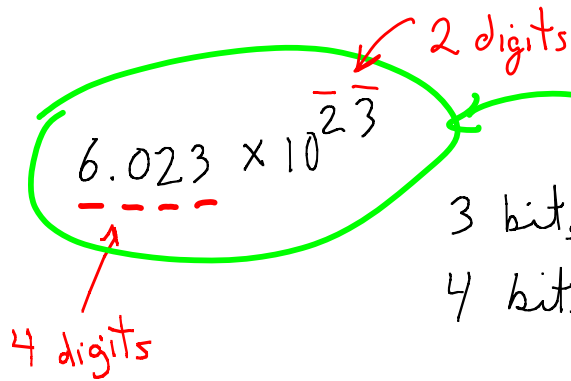
$$2^{-128} \times 1.00 \dots 0 \quad \text{that's small. Do we need it?}$$



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?



a nice number

$$\left. \begin{array}{l} 3 \text{ bits per digit } (2^3 \rightarrow 0..7) \\ 4 \text{ bits per digit } (2^4 \rightarrow 0..15) \end{array} \right\} \approx 3.5 \text{ bits per digit}$$

$$\left. \begin{array}{l} 4 \text{ dig.} \times \frac{3 \text{ bit}}{\text{digit}} = 12 \text{ bits} \\ 4 \text{ dig.} \times \frac{4 \text{ bit}}{\text{digit}} = 16 \text{ bits} \end{array} \right\} \Rightarrow 12 \leq \left(\begin{array}{c} \text{Precision} \\ \text{of} \\ f \end{array} \right) \leq 16$$

(23 bits are enough)

range of E?
(2 dec. digits)

$$\left. \begin{array}{l} 2 \text{ dig.} \times \frac{3 \text{ bit}}{\text{dig.}} = 6 \text{ bits} \\ 2 \text{ dig.} \times \frac{4 \text{ bit}}{\text{dig.}} = 8 \text{ bits} \end{array} \right\} 6 \leq \left(\begin{array}{c} \text{bits} \\ \text{of} \\ E \end{array} \right) \leq 8$$

Let's check

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

← 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?

x 0

y 1

← sign bit looks ok,
 \Rightarrow Looks like + 2's-comp
 \Rightarrow Looks like - 2's-comp
 \Rightarrow ALL (+) floats are greater than all (-) floats, as 2's-comp.

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

x 0 A 0 0

y 0 B 0 0

$A = 011$
 $B = 111$

3-bit, 2's-comp exp

	0	1	1
+exp	0	1	0
	0	0	1
	0	0	0
-exp	1	1	1
	1	0	1
	1	0	0

} these are larger as unsigned, darn.

$x = +2^3 \cdot (1.00 \dots 0)$
 $y = +2^{-1} \cdot (1.0 \dots 0)$

\Rightarrow

$(0$	011	$00 \dots 0)$	x
$(0$	111	$00 \dots 0)$	y

y looks like bigger 2's comp. number than x .
OOPS!

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

Suppose $\text{sign}(x) = \text{sign}(y)$

(Mag. comparison)



or



(Reverse result for neg.)

(for signs \neq result is obvious)

let $e_i = \text{value}(E_i)$

Note: $e_1 > e_2 \Rightarrow 2^{e_1}(1.f_1) > 2^{e_2}(1.f_2)$
regardless of the fractional parts.

Let's check

Suppose $e_2 = e_1 + 1$

$2^{e_2} \cdot (1.0)$ $2^{e_1} (2 - \epsilon) = 2^{e_1+1} (1 - \epsilon/2)$
smallest possible x Largest possible y

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

Y is less than X

So, make E_1 look bigger than E_2

what we have so far

8-bit exponent (single float)

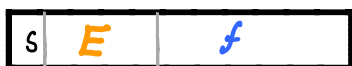
2's complement

{	0111 1111 $(2^7 - 1 = 127)$
	\vdots
	0000 0000 (0)
	\vdots
	1000 0001 (-127, not used, reserved for signal)
	1000 0000 (-128, not used, signals 0?)

Let's
fix Σ_{exp}

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

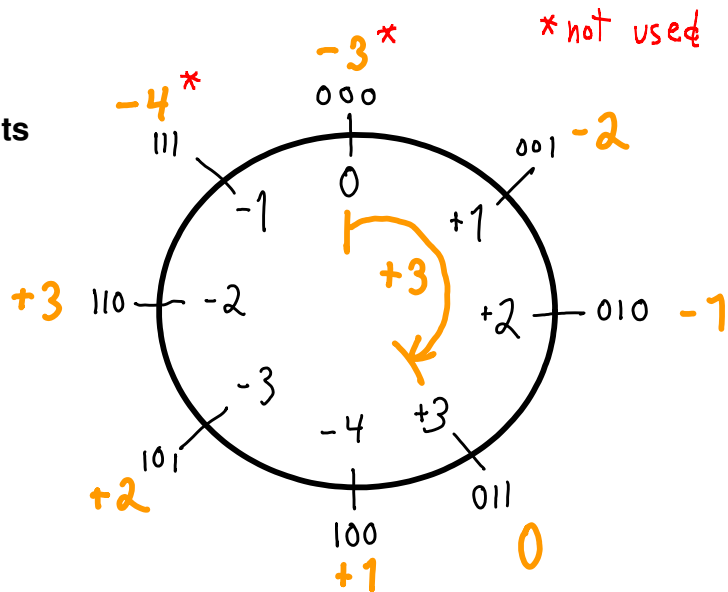
Negative exponents look smaller than positive exponents AS unsigned ints.



$$E = e + 011$$

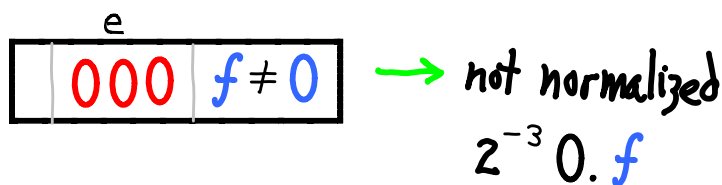
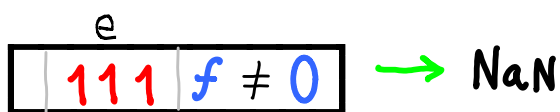
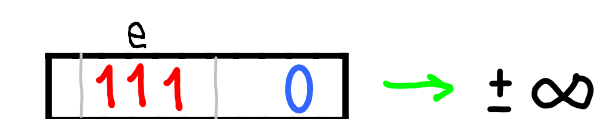
value	e in 2s-comp	E in excess-3
+3	011	+011 ==> 110
+2	010	+011 ==> 101
+1	001	+011 ==> 100
0	000	+011 ==> 011
-1	111	+011 ==> 010
-2	110	+011 ==> 001
-3	101	+011 ==> 000 *
-4	100	+011 ==> 111 *

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

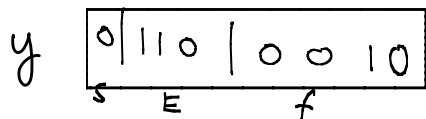
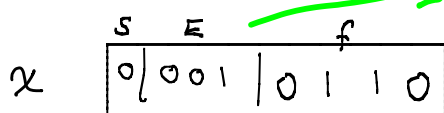


Rotate 2's comp so that +3 becomes largest number available.

How to represent 0?



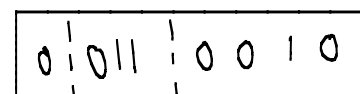
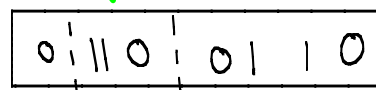
8-bit FP, ADD



$$001 - 011 = 110 \rightarrow -2$$

Convert from excess-3 to 2's comp.

$$110 - 011 = 011 = +3$$



$$2^3 \times 0.0000010110$$

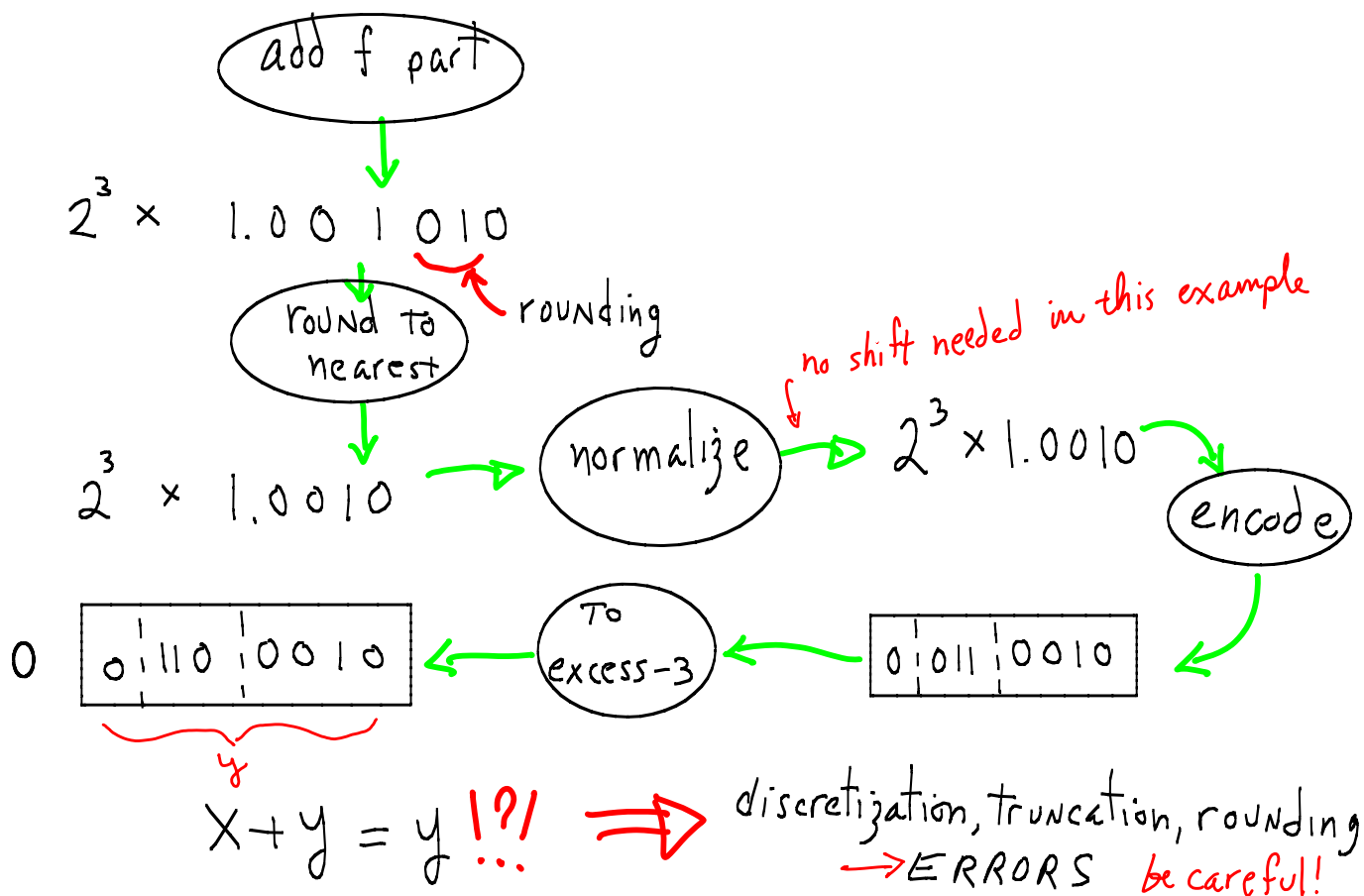
$$2^3 \times 1.0010$$

Shift/align exponents

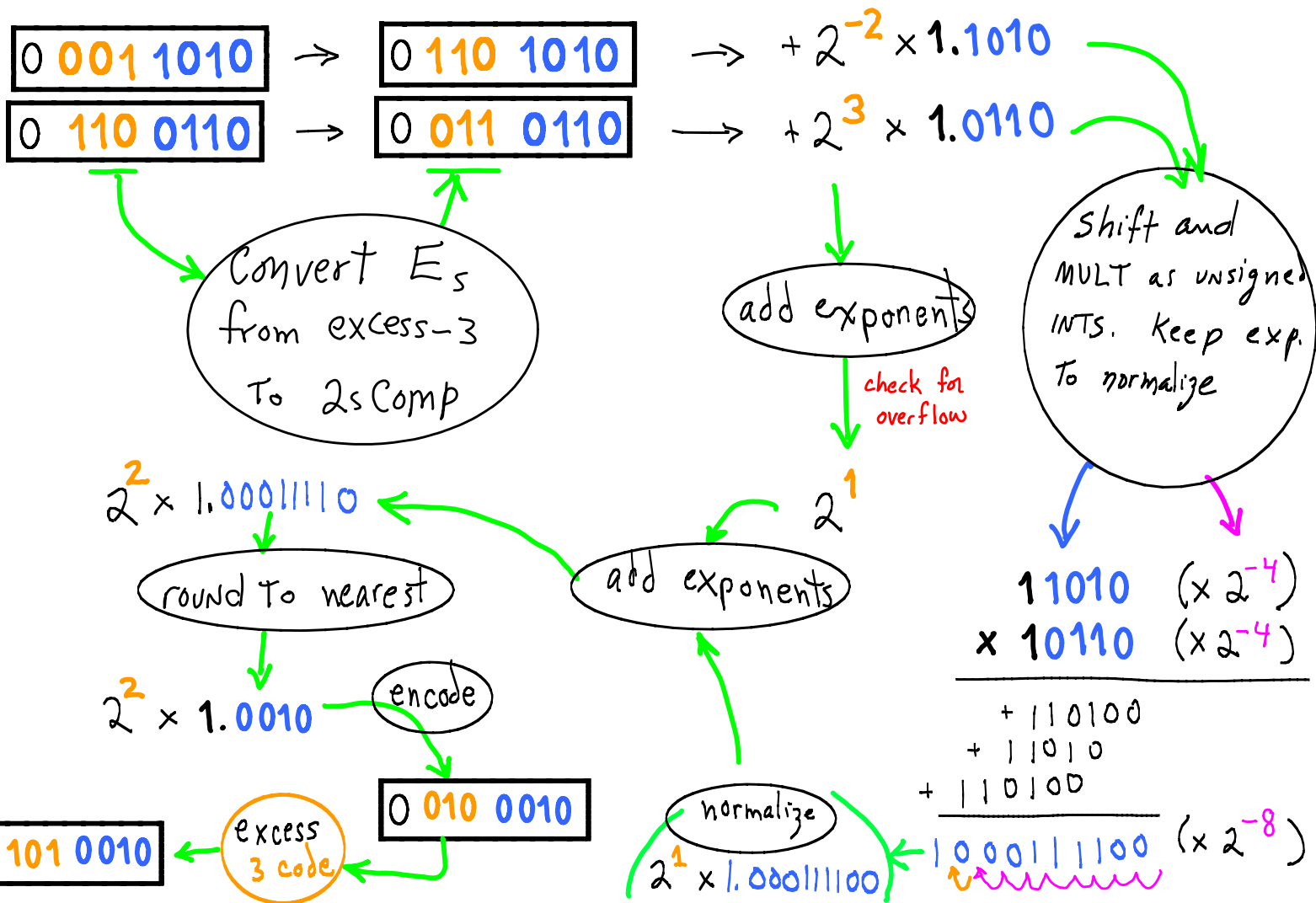
$$+2^{-2} \times 1.0110$$

$$+2^3 \times 1.0010$$

TRUNCATION?



8-bit FP, MULT



Convert to 32-bit FP 28

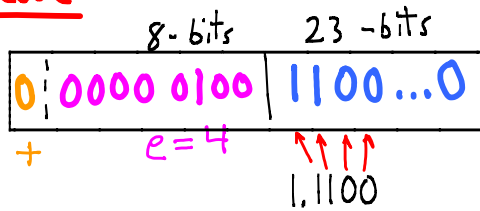
1. Convert to binary

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

2. Normalize

$$\begin{array}{c}
 11100. \\
 \uparrow \uparrow \uparrow \uparrow \\
 e=4
 \end{array}
 \Rightarrow 2^4 \times 1.1100$$

3. Encode

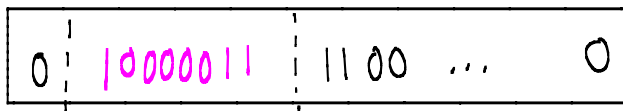


4. Convert e to excess ($2^{n-1}-1$)

excess ($2^{8-1}-1 = 127$)

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

5. Replace e with E



Convert back

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

2. convert E

$$\begin{array}{r}
 \leftarrow \text{borrows} \\
 -11111 \\
 1000011 \\
 -0111111 \\
 \hline
 0000100 = 4 = e
 \end{array}$$

$$\Rightarrow 2^4 \times 1.11 \quad (\text{convert } f) \Rightarrow 1.1100 = 11100$$

(mult. by 2^e) $\uparrow \uparrow \uparrow \uparrow$
 $\times 2^4$

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

- Problem with both truncation and rounding
 - They cause errors to **accumulate**
 - E.g., if always round up, result will gradually "crawl" upwards
- One solution: **round to nearest even**
 - If un-rounded LSB is 1 → round up (011 → 10)
 - If un-rounded LSB is 0 → round down (001 → 00)
 - Round up half the time, down other half → overall error is stable
- Another solution: **multiple intermediate precision bits**
 - IEEE 754 defines 3: guard + round + sticky
 - Guard and round are shifted by de-normalization as usual
 - Sticky is 1 if any shifted out bits are 1
 - Round up if 101 or higher, round down if 011 or lower
 - Round to nearest even if 100

- Suppose you added two numbers and came up with
 - 0 101 **11111 101**
 - What happens when you round?
 - Number becomes denormalized... arrrrgggghh
- FP adder actually has more than three steps...
 - Align exponents
 - Add/subtract significands
 - Re-normalize
 - **Round**
 - **Potentially re-normalize again**
 - **Potentially round again**

- Latency in cycles of common arithmetic operations
- Source: *Software Optimization Guide for AMD Family 10h Processors, Dec 2007*
 - Intel "Core 2" chips similar

	Int 32	Int 64	Fp 32	Fp 64
Add/Subtract	1	1	4	4
Multiply	3	5	4	4
Divide	14 to 40	23 to 87	16	20

- Floating point divide faster than integer divide?
 - Why?

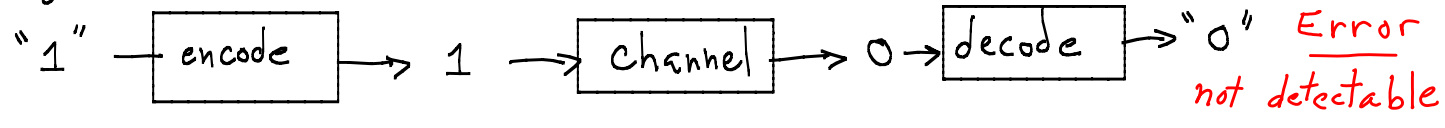
So much for encoding data. We could go on to audio, video, ... But, back to noise and errors.

Error Detection/Correction

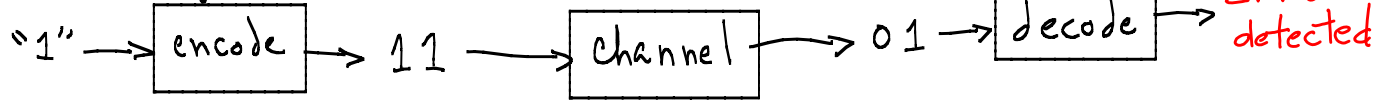
"message" could be a bit, a string of bits, a character, a page of characters, ...



message coded in bits:



2-bit encoding



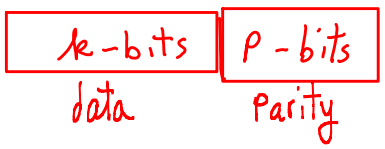
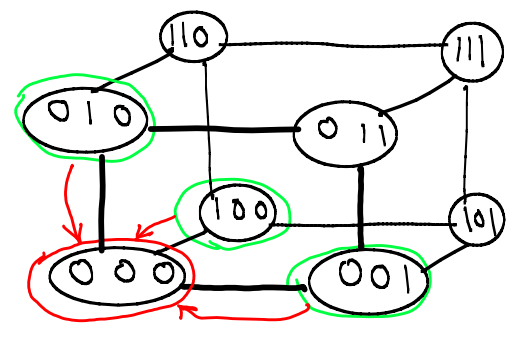
Code words 00 and 11 are good data, 10 and 01 indicate 1-bit errors. Last bit is "parity" bit, odd parity codeword indicates error. Works for k-bit messages w/ 1 parity bit (if 2-bit errors very unlikely).

Parity scheme data 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"



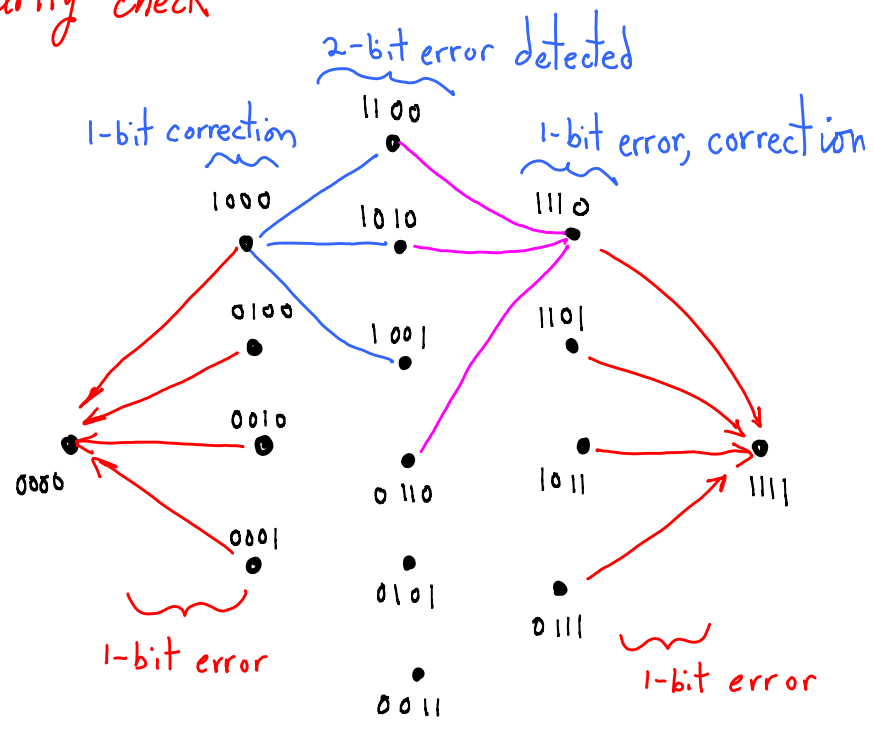
P-bit parity check

1-bit Correction, 2-bit Detection

- odd parity: 1-bit error corrected
- exactly two 1's: 2-bit error detected
- otherwise: no error

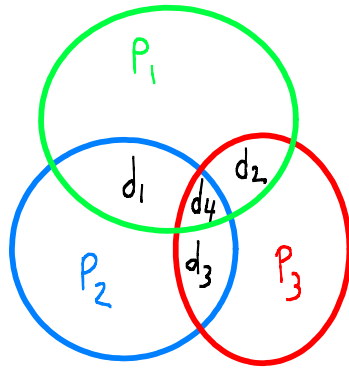
How many extra bits, at minimum?
Depends on noise in channel:
Shannon Noisy Coding Theorem.

We use 4 bits, 1-bit data.



Hamming (7,4) code (Single Error Detection / Single Error Correction)

7 bits per code word
 4 data bits
 3 parity bits

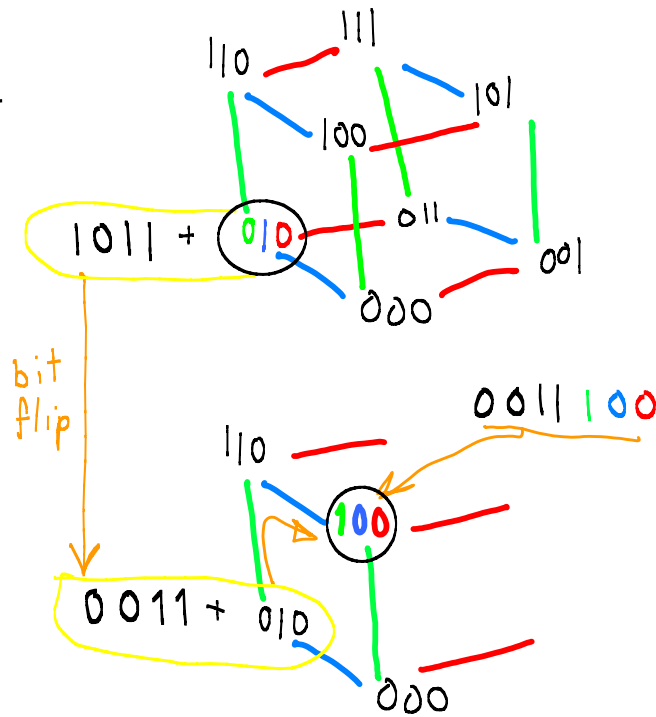
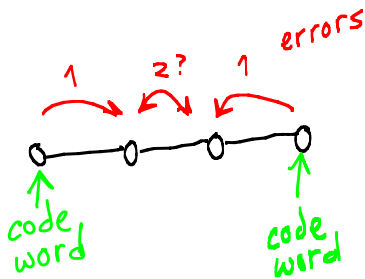


$d_1, d_2, d_3, d_4, P_1, P_2, P_3$

$P_1 = \text{parity}(d_1, d_2, d_4)$
 $P_2 = \text{parity}(d_1, d_3, d_4)$
 $P_3 = \text{parity}(d_2, d_3, d_4)$

3 steps to next codeword

1-bit error: can detect and correct
 2-bit error: cannot detect



other neighbor code words

- 0001 111
- 0010 011
- 0111 001

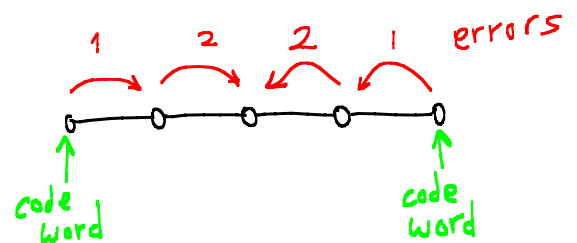
What can we do about 2-bit errors?
 Add another parity bit.

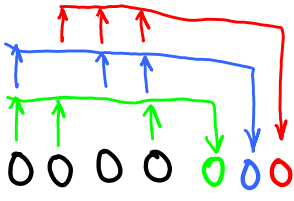
$$P_4 = \text{parity}(d_1, d_2, d_3, d_4, P_1, P_2, P_3)$$

$$\text{code word} = d_1, d_2, d_3, d_4, P_1, P_2, P_3, P_4$$

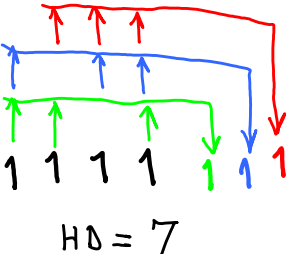
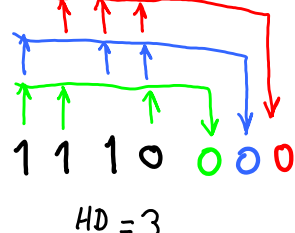
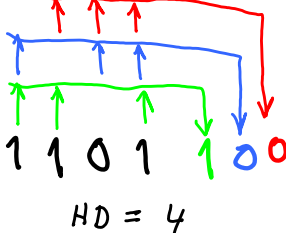
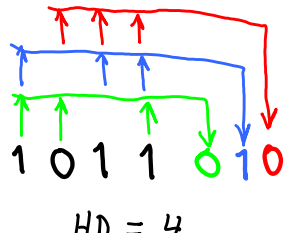
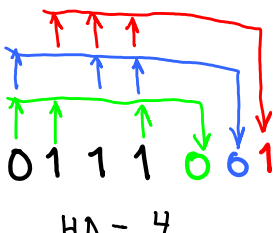
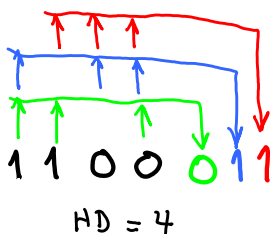
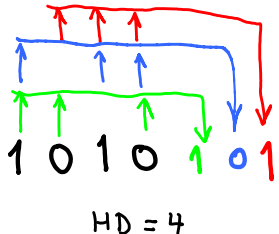
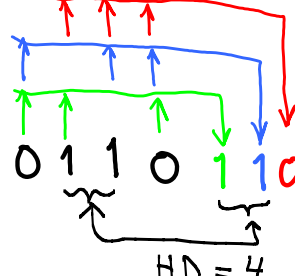
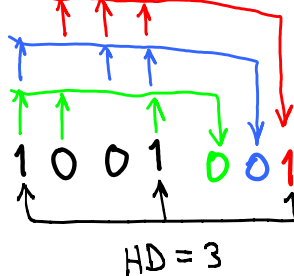
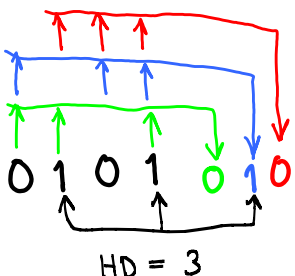
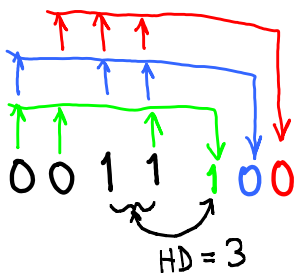
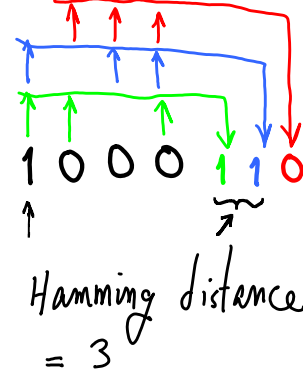
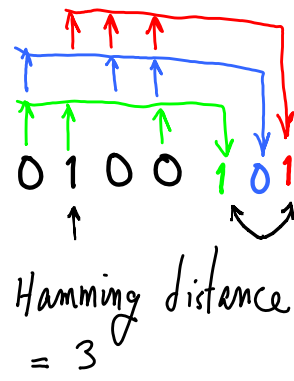
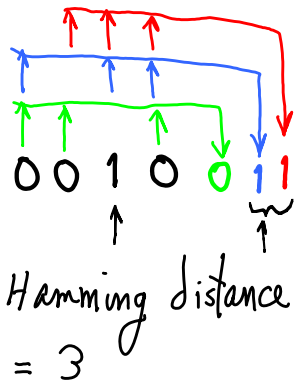
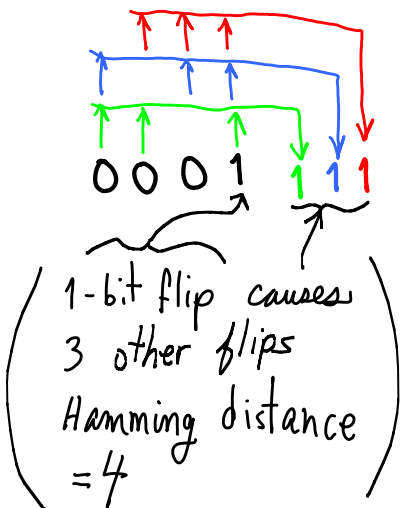
⇒ 4 steps min to next code word

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





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

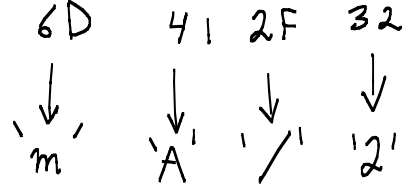


ASCII (See back cover of PP)

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?



Byte Addressable

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

low bit (red arrow pointing to bit 0 of address 0)
high bit (red arrow pointing to bit 7 of address 3)

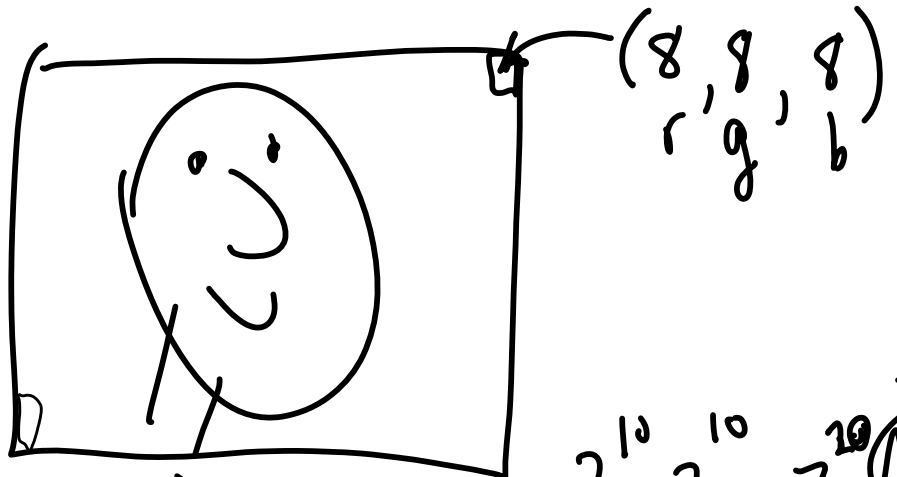
What to Print	Starting Memory Address	What is displayed (left-to-right)
4-byte number (in hex notation)	0	6D412F32
two 2-byte numbers (in hex)	0	2F32 6D41
four 1-byte numbers (in hex)	0	32 2F 41 6D
one 4-byte string	0	2 / A m

(see "od" in unix)

24 Mb
3 MB

~~1K~~
2¹⁰

2¹⁰ 1K ~~1M~~



2¹⁰ 2¹⁰ = 2²⁰ (M)
~~2²⁰ (0.2) = 10⁶~~

run length

