

# How Computers Really Do Arithmetic

Bruce Maggs

$$\begin{array}{r} 15211 \\ + 15299 \\ \hline 20510 \end{array}$$

$$\begin{array}{r} 11101101101011 \\ 11101111000011 \\ \hline 00000010101000 \\ + 111011 \leftarrow \leftarrow 0 \leftarrow 011 \end{array}$$



Why binary?

- addition
- multiplication
- division
- square roots

# "Schoolboy" Arithmetic

- carry-propagate addition

decimal

n digits

$$\begin{array}{r} \overset{1}{1} \overset{2}{9} \overset{2}{9} \overset{2}{7} \\ + 3003 \\ \hline 5000 \end{array}$$

binary

n bits

$$\begin{array}{r} \overset{1}{1} \overset{2}{0} \overset{2}{1} \overset{2}{1} \overset{2}{1} \overset{2}{1} \overset{2}{0} \overset{2}{1} \\ + 10000011 \\ \hline 1010000000 \end{array}$$

# "Schoolboy" Arithmetic

- carry-propagate addition

decimal

n digits

$$\begin{array}{r} \overset{1}{1} \overset{2}{9} \overset{3}{9} \overset{4}{7} \\ + 3003 \\ \hline 5000 \end{array}$$

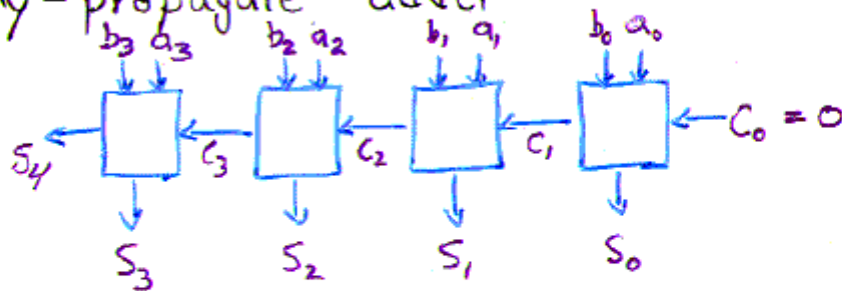
binary

n bits

$$\begin{array}{r} 1011101 \\ + 10000011 \\ \hline 101000000 \end{array}$$

- carry-propagate adder:

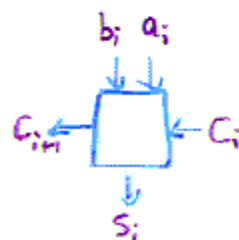
$$\begin{array}{r} a_3 \ a_2 \ a_1 \ a_0 \\ + b_3 \ b_2 \ b_1 \ b_0 \\ \hline s_4 \ s_3 \ s_2 \ s_1 \ s_0 \end{array}$$



where

$$s_i = a_i \oplus b_i \oplus c_i$$

$$c_{i+1} = a_i \cdot b_i \vee a_i \cdot c_i \vee b_i \cdot c_i$$



Time?  
n  
why?

# Carry-Look-Ahead Addition

$$\begin{array}{r} \text{To add} \quad + \quad 10111101 \\ \hline \quad \quad \quad 10000011 \end{array}$$

1) compute parity  $p_i$  at each bit position

$$p_i = a_i \oplus b_i$$

1 time step

$$\begin{array}{r} 10111101 \\ 10000011 \\ \hline 00111110 \end{array}$$

(compute all  $p_i$   
values in parallel)

# Carry-Look-Ahead Addition

$$\begin{array}{r} \text{To add} \\ + \quad 10111101 \\ \hline \end{array}$$

1) compute parity  $p_i$  at each bit position

$$p_i = a_i \oplus b_i$$

1 time step

$$\begin{array}{r} 10111101 \\ 10000011 \\ \hline 00111110 \end{array}$$

(compute all  $p_i$  values in parallel)

2) compute carry-in  $C_i$  at each bit position

$$\begin{array}{r} 10111101 \\ 10000011 \\ \hline 10 \leftarrow \leftarrow \leftarrow \leftarrow \leftarrow 1 \end{array}$$

How? What does this notation mean?

# Carry-Look-Ahead Addition

To add 
$$\begin{array}{r} 10111101 \\ + 10000011 \\ \hline \end{array}$$

1) compute parity  $p_i$  at each bit position

$$p_i = a_i \oplus b_i \quad \text{1 time step}$$

$$\begin{array}{r} 10111101 \\ 10000011 \\ \hline 00111110 \end{array} \quad \text{(compute all } p_i \text{ values in parallel)}$$

2) compute carry-in  $C_i$  at each bit position

$$\begin{array}{r} 10111101 \\ 10000011 \\ \hline 10 \leftarrow \leftarrow \leftarrow \leftarrow \leftarrow 1 \end{array} \quad \text{How? What does this notation mean?}$$

3) compute  $S_i = p_i \oplus C_i$  at each position

$$\begin{array}{r} 00111110 \\ \oplus 10111110 \\ \hline 10100000 \end{array} \quad \text{1 time step}$$

# Prefix Sums Calculation

Input: sequence  $X = X_{n-1}, X_{n-2}, \dots, X_1, X_0$   
binary associative operator  $+$   
(associative:  $(x+y)+z = x+(y+z)$ )

Output: sequence  $Y = Y_{n-1}, Y_{n-2}, \dots, Y_1, Y_0$   
where  $Y_i = \sum_{j=0}^i X_j$

Example:  $+$   $\equiv$  integer addition

$$X = 3, 2, 1, 4, 8, 2$$

$$Y = 20, 17, 15, 14, 10, 2$$

# Prefix Sums Calculation

Input: sequence  $X = X_{n-1}, X_{n-2}, \dots, X_1, X_0$   
binary associative operator  $+$   
(associative:  $(x+y)+z = x+(y+z)$ )

Output: sequence  $Y = Y_{n-1}, Y_{n-2}, \dots, Y_1, Y_0$   
where  $Y_i = \sum_{j=0}^i X_j$

Example:  $+$   $\equiv$  integer addition  
 $X = 3, 2, 1, 4, 8, 2$   
 $Y = 20, 17, 15, 14, 10, 2$

Other Associative Operators:

multiplication, min, max, and, or, left, right  
left  $(a,b) = a$  right  $(a,b) = b$



# Algorithm for Computing Prefix Sums

(assume  $n$  is a power of 2)

Base case:  $n=1$ ,  $y_0 = x_0$

Recursive case: 1) let  $z_i = x_{2i+1} + x_{2i}$ ,  $0 \leq i \leq \frac{n}{2}-1$

i.e. add pairs in  $X$ :

$x_{n-1}$   $x_{n-2}$  ...  $x_3$   $x_2$   $x_1$   $x_0$   
+  
 $z_{\frac{n}{2}-1}$   $z_1$   $z_0$

# Algorithm for Computing Prefix Sums

(assume  $n$  is a power of 2)

Base case:  $n=1$ ,  $Y_0 = X_0$

Recursive case: 1) let  $Z_i = X_{2i+1} + X_{2i}$ ,  $0 \leq i \leq \frac{n}{2}-1$

i.e. add pairs in  $X$ :

The diagram illustrates the recursive step of adding pairs of elements in the sequence  $X$ . It shows three pairs of elements:  $x_{n-1}$  and  $x_{n-2}$ ,  $x_3$  and  $x_2$ , and  $x_1$  and  $x_0$ . Each pair is summed, and the result is labeled as  $z_{\frac{n}{2}-1}$ ,  $z_1$ , and  $z_0$  respectively. The summation is shown with a plus sign and a vertical line below it.

2) recursively compute prefix sums

$W = w_{\frac{n}{2}-1}, \dots, w_1, w_0$  for sequence  $Z$

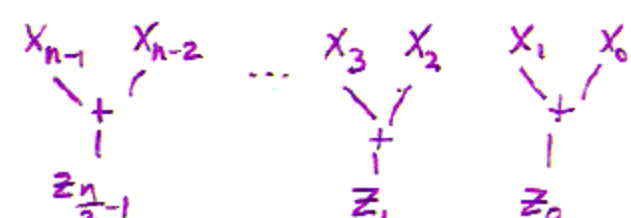
# Algorithm for Computing Prefix Sums

(assume  $n$  is a power of 2)

Base case:  $n=1$ ,  $Y_0 = X_0$

Recursive case: 1) let  $Z_i = X_{2i+1} + X_{2i}$ ,  $0 \leq i \leq \frac{n}{2}-1$

i.e. add pairs in  $X$ :



$x_{n-1}, x_{n-2}, \dots, x_3, x_2, x_1, x_0$   
 $z_{\frac{n}{2}-1}, z_i, z_0$

2) recursively compute prefix sums

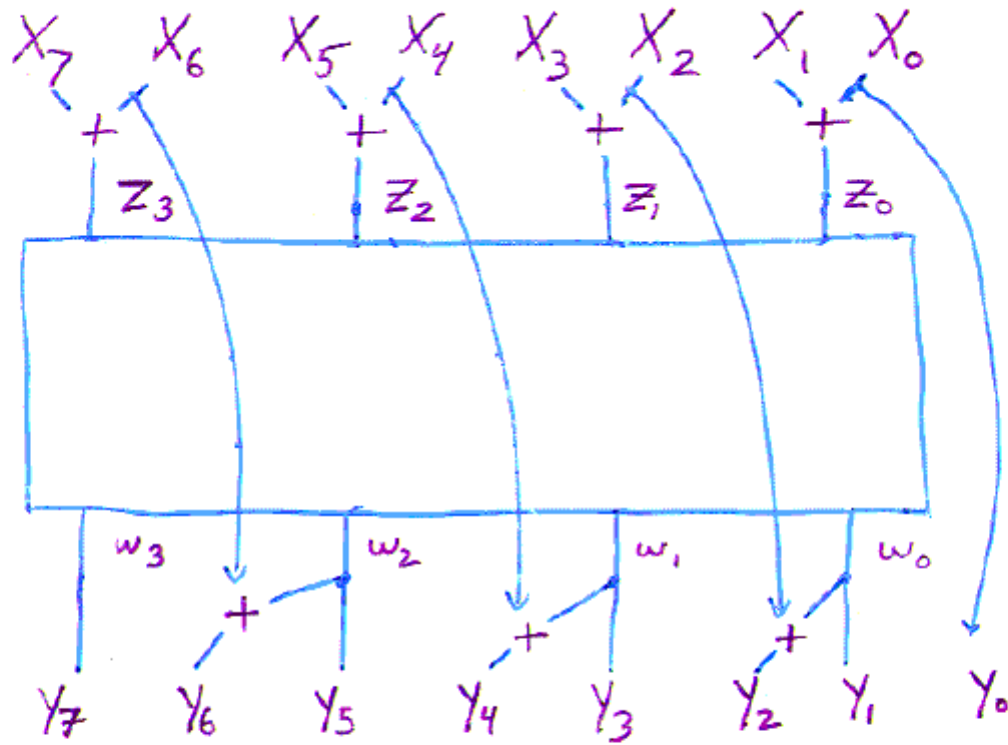
$W = w_{\frac{n}{2}-1}, \dots, w_1, w_0$  for sequence  $Z$

3) let

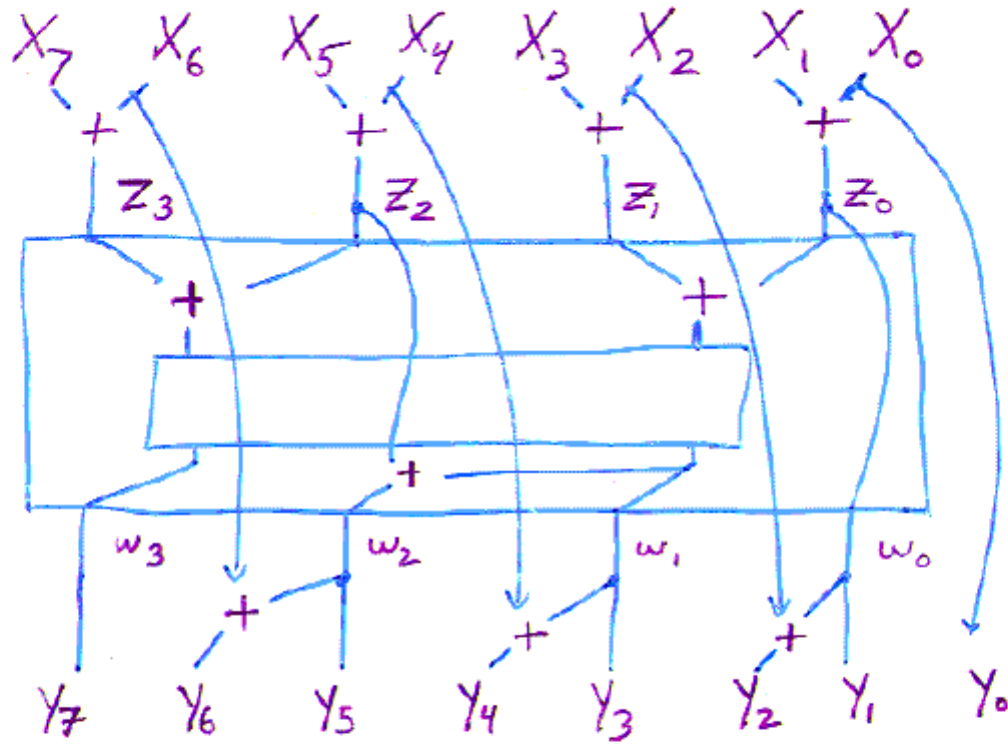
$$Y_i = \begin{cases} w_{\frac{i-1}{2}} & i \text{ odd} \\ X_i + w_{\frac{i}{2}-1} & i \text{ even, } i > 0 \\ X_0 & i = 0 \end{cases}$$

i.e.  $Y = w_{\frac{n}{2}-1}, \dots, X_1 + w_1, w_1, X_2 + w_0, w_0, X_0$

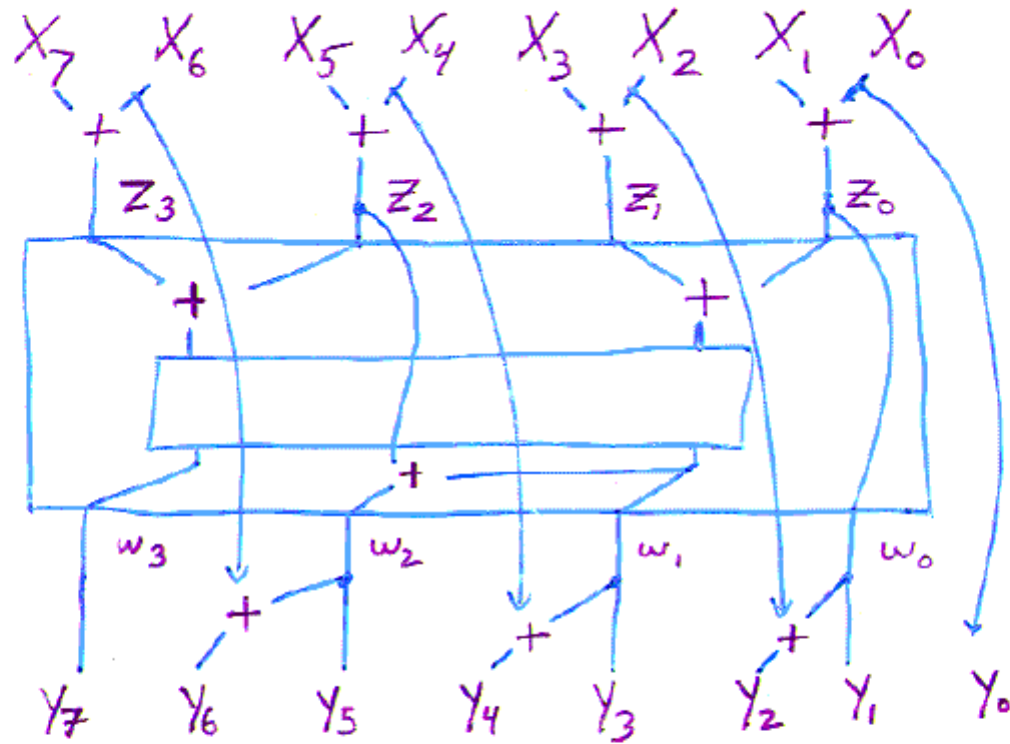
# Circuit Diagram



# Circuit Diagram



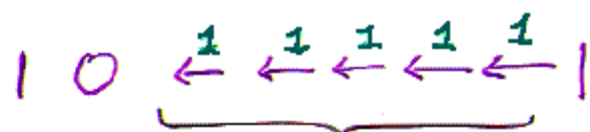
# Circuit Diagram



$$\begin{aligned} \text{Time: } T(n) &= \begin{cases} 1 & n=1 \\ 1+T(\frac{n}{2})+1 & n>1 \end{cases} \\ &= 2 \log_2 n + 1 \end{aligned}$$

Better  
than  $n$ .

# Back to carry-lookahead addition



we want to replace each  $\leftarrow$  with a 1, and not sequentially!

What's the connection to prefix sums?

define operator  $*$  as follows

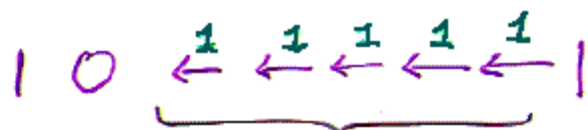
$*$	0	1	$\leftarrow$
0	0	0	0
1	1	1	1
$\leftarrow$	0	1	$\leftarrow$

now apply to

$$X = 1 \ 0 \ \leftarrow \ \leftarrow \ \leftarrow \ (\leftarrow \ (\leftarrow \ 1))$$

$$\Leftrightarrow Y = 1 \ 0 \ 1 \ 1 \ 1 \ \underline{1} \ \underline{1} \ 1$$

# Back to carry-lookahead addition



we want to replace each  $\leftarrow$  with a 1, and not sequentially!

What's the connection to prefix sums?

define operator  $*$  as follows

$*$	0	1	$\leftarrow$
0	0	0	0
1	1	1	1
$\leftarrow$	0	1	$\leftarrow$

now apply to

$$X = 1 \ 0 \ \leftarrow \leftarrow \leftarrow \left( \left( \leftarrow \leftarrow 1 \right) \right)$$

$$\Leftrightarrow Y = 1 \ 0 \ 1 \ 1 \ 1 \ 1 \ 1 \ 1$$

we should verify that

$*$  is associative. Easy if  $x=0,1$ .

$x$	$y$	$z$	$x * (y * z)$	$(x * y) * z$
0	0	0	0	0
0	0	1	0	0
0	0	$\leftarrow$	0	0
...	...	...	...	...

$x$	$y$	$z$	$x * (y * z)$	$(x * y) * z$
$\leftarrow$	0	0	0	0
$\leftarrow$	0	1	0	0
$\leftarrow$	0	$\leftarrow$	0	0
$\leftarrow$	1	0	1	1
$\leftarrow$	1	1	1	1
$\leftarrow$	1	$\leftarrow$	1	1
$\leftarrow$	$\leftarrow$	0	0	0
$\leftarrow$	$\leftarrow$	1	1	1
$\leftarrow$	$\leftarrow$	$\leftarrow$	$\leftarrow$	$\leftarrow$



## Bottom Line:

Time to add 2  $n$ -bit numbers =

$$\underbrace{2 \log_2 n + 1}_{\text{compute } C_i\text{'s}} \quad + \quad \underbrace{1}_{\text{compute } s_i = p_i \oplus C_i} = 2 \log_2 n + 2$$

<u>processor</u>	<u><math>n</math></u>	<u><math>2 \log_2 n + 2</math></u>
80186	16	10
Pentium	32	12
Alpha	64	14

# More "Schoolboy" Arithmetic

## MULTIPLICATION

decimal

$$\begin{array}{r} \overbrace{211}^n \\ \times 299 \\ \hline 1899 \\ 1899 \\ 422 \\ \hline 63089 \end{array}$$

binary

$$\begin{array}{r} \overbrace{011010011}^n \\ \times 100101001 \\ \hline 011010011 \\ 011010011 \\ 011010011 \\ \hline 011010011 \\ 011010011 \\ 011010011 \\ \hline 011010011 \\ 011010011 \\ 011010011 \\ \hline \end{array}$$

up to  $n$   
"partial products"

$a_0$   
 $a_1$   
 $a_2$   
 $\vdots$

Time?  $\approx n^2$  using one-column-at-a-time addition

# More "Schoolboy" Arithmetic

## MULTIPLICATION

decimal

$$\begin{array}{r}
 \overbrace{211}^n \\
 \times 299 \\
 \hline
 1899 \\
 1899 \\
 422 \\
 \hline
 63089
 \end{array}$$

binary

$$\begin{array}{r}
 \overbrace{011010011}^n \\
 \times 100101001 \\
 \hline
 011010011 \\
 011010011 \\
 011010011 \\
 \hline
 011010011 \\
 011010011 \\
 011010011 \\
 \hline
 011010011
 \end{array}$$

up to  $n$   
"partial products"

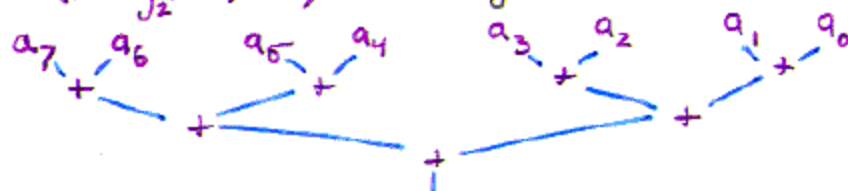
$a_0$   
 $a_1$   
 $a_2$   
 $\vdots$

Time?  $\approx n^2$  using one-column-at-a-time addition

$\approx n \cdot 2 \log_2(2n)$  computing sequentially  
 $a_0 + a_1, a_0 + a_1 + a_2, \dots$  (with carry look-ahead)

$\approx \log_2 n * (2 \log_2(2n) + 2)$  using tree

We can beat this!



Why  $2n$ ?

# Carry-Save Addition

Idea: Can convert sum of 3 numbers into a sum of 2 numbers in one step.

example

$$\begin{array}{r} 011010011 \\ 011010011000 \\ 011010011000 \\ \hline \end{array}$$

Diagram illustrating the carry-save addition of three 10-bit numbers. The numbers are aligned by their least significant bits. The sum is shown below a horizontal line. The carry bits are circled in red, and the parity bits are circled in blue. The carry bits are 1, 1, 0, 1, 0, 1, 1, 1, 0, 1. The parity bits are 0, 1, 0, 1, 0, 1, 0, 1, 0, 1.

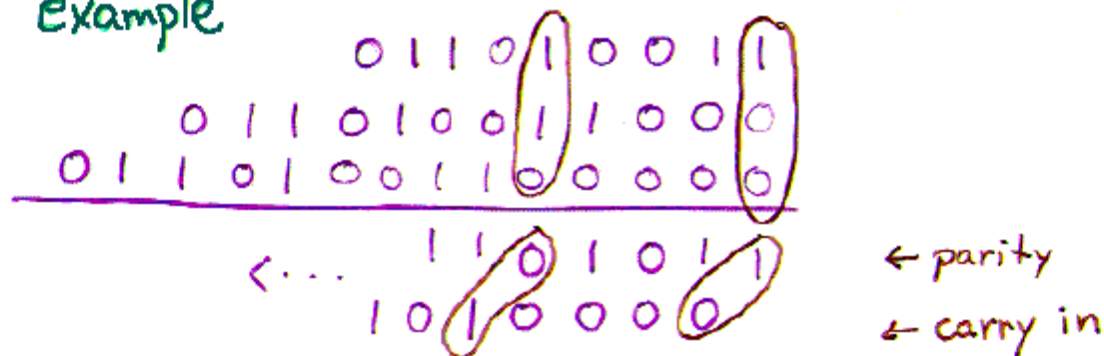
← parity

← carry in

# Carry-Save Addition

Idea: Can convert sum of 3 numbers into a sum of 2 numbers in one step.

example



Wallace Tree

depth of tree

$$\approx \log_{3/2} n$$

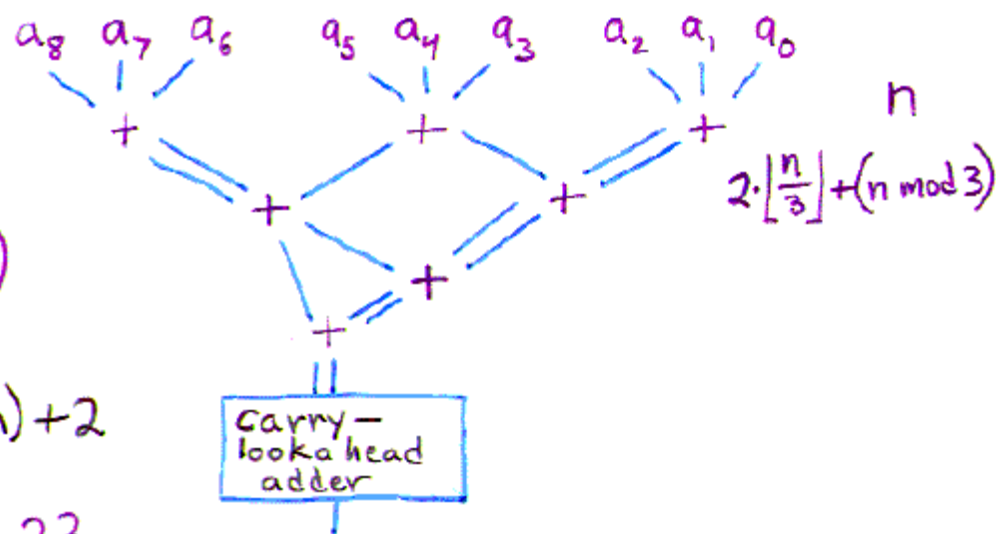
(exactly 9 for  $n=64$ )

total time

$$\approx \log_{3/2} n + 2\log(2n) + 2$$

time for  $n=64$ : 23

(vs. 14 to add!!)



# "Schoolboy" Division

$$\begin{array}{r} 1.0057 \dots \\ 15211 \overline{) 15299.} \\ \underline{-15211} \phantom{0} \\ 88\overset{9}{0}\overset{9}{0}\overset{0}{0} \\ \underline{76055} \\ 119450 \end{array}$$

$$\begin{array}{r} \textcircled{0} \\ 15211 \\ 30422 \\ 45633 \\ 60844 \\ 76055 \\ \vdots \\ \vdots \\ \vdots \end{array}$$

# "Schoolboy" Division

$$\begin{array}{r} 1.0057\dots \\ 15211 \overline{) 15299.} \\ \underline{-15211} \phantom{0} \\ 88\overset{9}{0}\overset{9}{0}0 \\ \underline{76055} \\ 119450 \end{array}$$

$$\begin{array}{r} 0 \\ 15211 \\ 30422 \\ 45633 \\ 60844 \\ 76055 \\ \vdots \\ \vdots \end{array}$$

Time to get  $n$  digits of precision:

- build table of multiples of divisor  
 $1 \times 15211, 2 \times 15211, 3 \times 15211, \dots, 9 \times 15211$

- $n$  table lookups

- $n$  subtractions ( $\approx 2 \log n + 2$  times each)

Total:  $\approx n(2 \log n + 2)$

(We can do better!!!)

## A few words about subtraction

- it never costs more than addition

Almost all machines use 2's-complement representation of signed integers.

example:

-128	+64	+32	+16	+8	+4	+2	+1	
1	1	0	0	1	1	0	0	= -51



## A few words about subtraction

- it never costs more than addition

Almost all machines use 2's-complement representation of signed integers.

example:

-128	+64	+32	+16	+8	+4	+2	+1	
1	1	0	0	1	1	0	1	= -51

Assuming no overflows, addition is unchanged:

					1				
	1	1	0	0	1	1	0	1	-51
+	0	0	1	0	1	1	0	1	+45
<hr/>									
	1	1	1	1	1	0	1	0	-6



# Redundant Representation of Integers

Idea: Allow each digit to be 0, 1, -1

examples:  $0101 = 5$

$10-1-1 = 5$

New addition algorithm:

1) add corresponding digits, no carries

$$\begin{array}{r} 0101-1011 \quad 75 \\ + 10-10-1010 \quad 90 \\ \hline 11-11-2021 \quad 165 \end{array}$$

# Redundant Representation of Integers

Idea: Allow each digit to be 0, 1, -1

examples:

$$\begin{array}{r} 0101 = 5 \\ 10-1-1 = 5 \end{array}$$

New addition algorithm:

1) add corresponding digits, no carries

$$\begin{array}{r} 0101-1011 \quad 75 \\ + 10-10-1010 \quad 90 \\ \hline 11-11-2021 \quad 165 \end{array}$$

2) remove +2's by changing to 0 passing +1 left  
" +1's " " -1 " +1 "

$$10-10-1-211-1$$

Now all digits are -2, -1, 0, 1 why?

# Redundant Representation of Integers

Idea: Allow each digit to be 0, 1, -1

examples:

$$\begin{array}{r} 0101 = 5 \\ 10-1-1 = 5 \end{array}$$

New addition algorithm:

1) add corresponding digits, no carries

$$\begin{array}{r} 0101-1011 \quad 75 \\ + 10-10-1010 \quad 90 \\ \hline 11-11-2021 \quad 165 \end{array}$$

2) remove +2's by changing to 0 passing +1 left  
" +1's " " -1 " +1 "

$$10-10-1-211-1$$

Now all digits are -2, -1, 0, 1 why?

3) remove -2's by changing to 0 passing -1 left  
" -1's " " +1 " -1 "

$$1-1+1-10010+1$$

We've added 2 n-bit numbers in 3 steps!

# SRT Division Algorithm

$$\begin{array}{r}
 \phantom{1011} \underline{111101101} \\
 1011 \overline{)111101101} \\
 \underline{-10111} \\
 1011 \\
 \underline{-10111} \\
 \phantom{1011} \underline{20} \\
 \Downarrow \\
 \phantom{1011} \underline{1001} \\
 \phantom{1011} 1011 \\
 \phantom{1011} \underline{012} \\
 \Downarrow \\
 \phantom{1011} \underline{1000} \\
 \phantom{1011} \underline{-10111} \\
 \phantom{1011} 0111
 \end{array}$$

$$\begin{array}{r}
 21r6 \\
 11 \overline{)237}
 \end{array}$$

**Rule:** Each bit of quotient is determined by comparing first bit of divisor with first bit of dividend. Easy!

$$22 r 5$$

# SRT Division Algorithm

$$\begin{array}{r}
 \phantom{1011} \overline{11110110} \phantom{11} \\
 \underline{-10111} \\
 10111 \\
 \underline{-10111} \\
 \phantom{10} 20 \\
 \downarrow \\
 \phantom{10} \underline{-1001} \\
 \phantom{10} 1011 \\
 \phantom{10} \downarrow \\
 \phantom{10} \phantom{10} 012 \\
 \phantom{10} \phantom{10} \phantom{10} \underline{1000} \\
 \phantom{10} \phantom{10} \phantom{10} \underline{-1011} \\
 \phantom{10} \phantom{10} \phantom{10} 0111
 \end{array}$$

$$\begin{array}{r}
 21r6 \\
 11 \overline{)237}
 \end{array}$$

**Rule:** Each bit of quotient is determined by comparing first bit of divisor with first bit of dividend. Easy!

$$22 \text{ r } 5$$

Time for  $n$  bits of precision in result:

$$\approx \underbrace{3 \cdot n} + \underbrace{2 \log n + 2}$$

1 addition per bit

convert to standard representation by subtracting negative bits from positive

# Intel Pentium Division Error

- used essentially the same algorithm, but computed more than one bit of result in each step. Examined several leading bits of divisor and remainder and looked in table.
- table had several bad entries
- ultimately Intel offered to replace any defective chip, estimating their loss at **\$ 475 million**



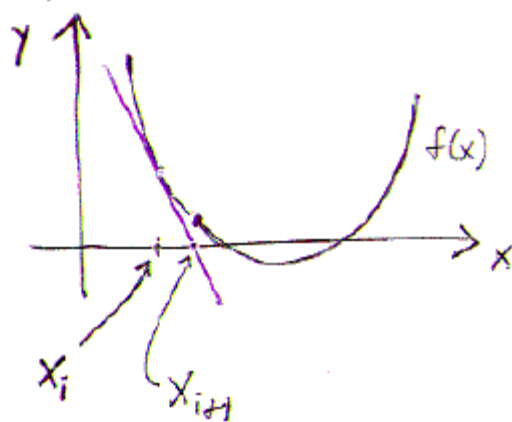
# square roots

Who remembers the "schoolboy" method?  
Fortunately, calculators became cheap before  
I managed to learn it.

Computers use another technique:

## Newton's Method

(also used in IBM 360 for division)



- finds  $x$  s.t.  $f(x) = 0$
- improves guess  $x_i$  by interpolation  $x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$
- e.g.  $f(x) = x^2 - a$   
(find  $\sqrt{a}$ )
- e.g.  $f(x) = \frac{1}{x} - a$  (find  $\frac{1}{a}$ )

## Newton's Method - Division

deriving the recurrence:

$$f(x) = \frac{1}{x} - a \quad f'(x) = -\frac{1}{x^2}$$

$$x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}$$

$$= x_i - \frac{\frac{1}{x_i} - a}{-\left(\frac{1}{x_i}\right)^2}$$

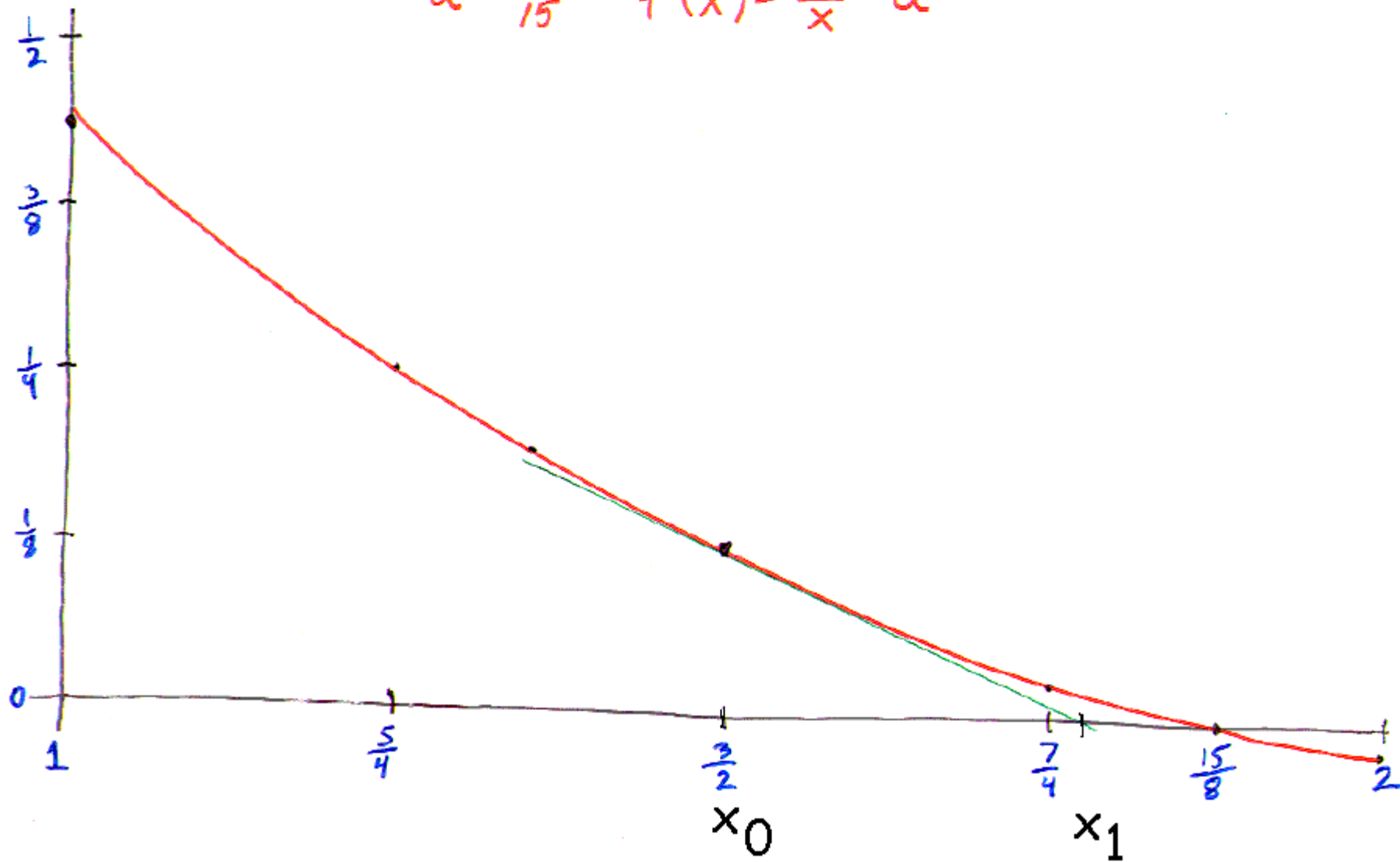
$$= 2x_i - x_i^2 a$$

Note: Some functions  $f$  s.t.  $f\left(\frac{1}{a}\right) = 0$   
don't work out. E.g.:

$$f(x) = x - \frac{1}{a} \Rightarrow x_{i+1} = \frac{1}{a}$$

$$f(x) = xa - 1 \Rightarrow x_{i+1} = \frac{1}{a}$$

$$a = \frac{8}{15} \quad f(x) = \frac{1}{x} - a$$



## Error Analysis - Division

Assume  $\frac{1}{2} < a < 1$  so that  $1 < \frac{1}{a} < 2$ .

$$\text{Let } \epsilon_i = x_i - \frac{1}{a}.$$

↑ error after  $i$  iterations

$$\text{Then } x_i = \frac{1}{a} + \epsilon_i, \text{ and}$$

$$\begin{aligned} x_{i+1} &= 2x_i - a \cdot x_i^2 \\ &= 2 \left( \frac{1}{a} + \epsilon_i \right) - a \left( \frac{1}{a} + \epsilon_i \right)^2 \\ &= \frac{1}{a} - a \epsilon_i^2, \end{aligned}$$

$$\text{so } \epsilon_{i+1} = -a \epsilon_i^2.$$

Notice that  $\epsilon_{i+1} < 0$ , and  $|\epsilon_{i+1}| < |\epsilon_i|^2$   
(for  $0 < a < 1$ ).

# Convergence - Division

$$\frac{1}{2} < a < 1 \quad 1 < \frac{1}{a} < 2$$

pick  $x_0 = \frac{3}{2}$  (initial guess)

$$\Rightarrow |\epsilon_0| < \frac{1}{2}$$

$$\Rightarrow |\epsilon_1| < \left(\frac{1}{2}\right)^2, \quad |\epsilon_2| < \left(\left(\frac{1}{2}\right)^2\right)^2$$

$$|\epsilon_i| < \frac{1}{2^{2^i}}$$

$\therefore$  After  $i$  iterations,  $x_i$  is correct to  $2^i$  bits

$$x_i = 1.\underbrace{01011011}_{2^i}101\dots$$

except for possible error smaller in magnitude than max. contribution of all bits beyond  $2^i$ .

(Adding  $|\epsilon_i|$  to  $x_i$  might change some of the leading  $2^i$  bits, though.)

Total time?  
 $O(\log^2 n)$

# Newton's Method - Square Root

$$f(x) = x^2 - a \quad \frac{1}{2} < a < 1, \quad \frac{1}{2} < \sqrt{a} < 1$$

$$\begin{aligned} \text{derivation: } X_{i+1} &= X_i - \frac{f(x_i)}{f'(x_i)} & X_0 &= \frac{3}{4} \\ &= X_i - \frac{x_i^2 - a}{2x_i} \\ &= \frac{1}{2}X_i + \frac{a}{2X_i} \end{aligned}$$

$$\text{error analysis: Let } X_i = \sqrt{a} + \epsilon_i. \quad |\epsilon_0| < \frac{1}{4}$$

$$\begin{aligned} X_{i+1} &= \frac{1}{2}(\sqrt{a} + \epsilon_i) + \frac{a}{2(\sqrt{a} + \epsilon_i)} \\ &= \frac{1}{2}(\sqrt{a} + \epsilon_i) + \frac{\sqrt{a}}{2} \left( \frac{1}{1 + \frac{\epsilon_i}{\sqrt{a}}} \right) \\ &= \frac{1}{2}(\sqrt{a} + \epsilon_i) + \frac{\sqrt{a}}{2} \left( 1 - \frac{\epsilon_i}{\sqrt{a}} + \frac{\epsilon_i^2}{a} - \dots \right) \\ &\quad \text{(provided that } \left| \frac{\epsilon_i}{\sqrt{a}} \right| < 1) \\ &= \sqrt{a} + \frac{\epsilon_i^2}{2\sqrt{a}} - \dots \end{aligned}$$

$$\therefore |\epsilon_{i+1}| < |\epsilon_i|^2$$

Total time?  $O(\log^3 n)$