

lecture 3

- BLAS complexity justified
- a word on error (absolute, relative, approximate, forward, backward, total)
- finite bit representation of numbers - introduction (big and little endian)
- IEEE floating point representation of numbers

Memory Wall Always There ...

Computation: Theoretical peak: (# cpu cores) * (flops / cycle / core) * (cycles / second)

Memory: Theoretical peak: (bus width) * (bus speed)

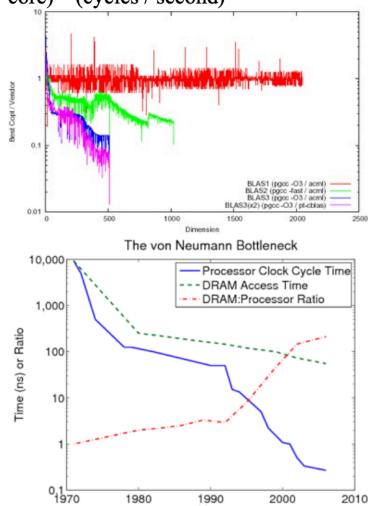
BLAS 1: O(n) operations on O(n) operands
BLAS 2: O(n^{**}2) operations on O(n^{**}2) operands
BLAS 3: O(n^{**}3) operations on O(n^{**}2) operands

$y = \alpha x + y$:
 3 loads, 1 store
 (more expensive than FP_OPs by a long shot)
 2 floating point operations (maybe 1) on 3 operands

e.g., double precision on the FY10 target platform:

$(3 \text{ operands} / 2 \text{ flop}) * (8 \text{ bytes} / \text{operand}) * 6 \text{ core} * 4 \text{ (flop / cyc /core)} * 2.6e9(\text{cyc/sec}) \sim 125 \text{ GBps}$

... We don't have this and to get it is \$\$\$... how to achieve **Sustainability??**



1/2

BLAS - Basic Linear Algebra Subprograms.

	Operations	Operands	Example
L1	$\Theta(n)$	$\Theta(n)$	$\langle x y \rangle$

Let. $x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ $y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$ $\langle x | y \rangle = (x^T, y)$

Clearly $2n$ numbers, $\Theta(n)$

$2n-1$ floating point work.

$$= \underbrace{x_1 \dots x_n}_{\text{n multiplies}} \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

$$= \sum_{i=1}^n x_i^T y_i //$$

$n-1$ adds.

$2n-1$ FP-OPs //

L2 $\Theta(n) - \Theta(n^2)$ $\Theta(n^2)$ Ax , matrix vector product.
 $\Theta(n \log n) ??$ (sparse)

if A $n \times n$, x $n \times 1$, y $n \times 1$

$$Ax = y$$

$$y_i = \sum_{j=1}^n A_{ij} x_j \quad \forall i = 1, n.$$

n multiplies
 $n-1$ adds } we do this n -times

$$n(2n-1) = 2n^2 - n \text{ FP-OPs //}$$

clearly $2n + n^2$ numbers, $\Theta(n^2)$

and $2n^2 - n$ floating point work

Operations	Time Complexity	Space Complexity	z/z
	Operations	Operands	
$\mathcal{O}(n^3)$	$\mathcal{O}(n^2)$	$\alpha A B + \beta C \rightarrow C$ general matrix multiply- (GEMM)	
$C = \alpha A B + \beta C.$			
let $\alpha, \beta \in \mathbb{R}^0$, $A, B, C \in \mathbb{R}^{n \times n}$			
use $\sum_{i=1}^n i\rangle \langle i = \hat{\mathbb{1}}_n$. where $ i\rangle = \begin{pmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \end{pmatrix}$ 1 in the i th position.			
ie let $n=3$.			
clearly $ 1\rangle \langle 1 = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$			
$ 2\rangle \langle 2 = \begin{pmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{pmatrix}$			
$ 3\rangle \langle 3 = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 1 \end{pmatrix}$			
		$\sum_{i=1}^3 i\rangle \langle i = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} = \hat{\mathbb{1}}_3$	
also, note that $\langle i j\rangle = \delta_{ij} = \begin{cases} 0 & i \neq j \\ 1 & i=j \end{cases}$			
thus.	$\langle i c_j\rangle = c_{ij}$	and $\exists n^2$ of these.	
$\langle i \alpha A B + \beta C j\rangle = \alpha \langle i A B j\rangle + \beta \langle i C j\rangle$			
so, $(2+3n^2)$ numbers is $\mathcal{O}(n^2)$ storage. //	$= \alpha \sum_{k=1}^n \langle i A(\sum_{l=1}^n l\rangle \langle l)B j\rangle + \beta C_{ij}$		
	$= \alpha \sum_{k=1}^n \sum_{l=1}^n A_{il} B_{lj} + \beta C_{ij}$	$\forall i, j = 1, n$	
	+ 1 mult, (n multiplies),	1 multiply	
		+ 1 add.	
so, $(2n+2)n^2$ floating pt operations			
	$= 2n^3 + 2n^2$ is $\mathcal{O}(n^3)$. work. //		

Matrix multiply
- a cornerstone of
many applications
including linear
algebra, PDE
solvers, and ML
(both training and
inference)

Powers of 2 (2^n)	Prefix (SI)	Numerical Range [$0, 2^n - 1$]
2^{10}	Kilo (K)	0 to 1,023
2^{20}	Mega (M)	0 to 1,048,575
2^{30}	Giga (G)	0 to 1,073,741,823
2^{40}	Tera (T)	0 to 1,099,511,627,775
2^{50}	Peta (P)	0 to 1,125,899,906,842,623
2^{60}	Exa (E)	0 to 1,152,921,504,606,846,975

1 bit : {0,1} : 2^1 combinations

8 bits = 1 byte = 1 B

1 B : 2^8 = 256 combinations

Big endian, most significant bit (2^7) stored leftmost.



Little endian, least significant bit (2^0) stored leftmost.



Big endian

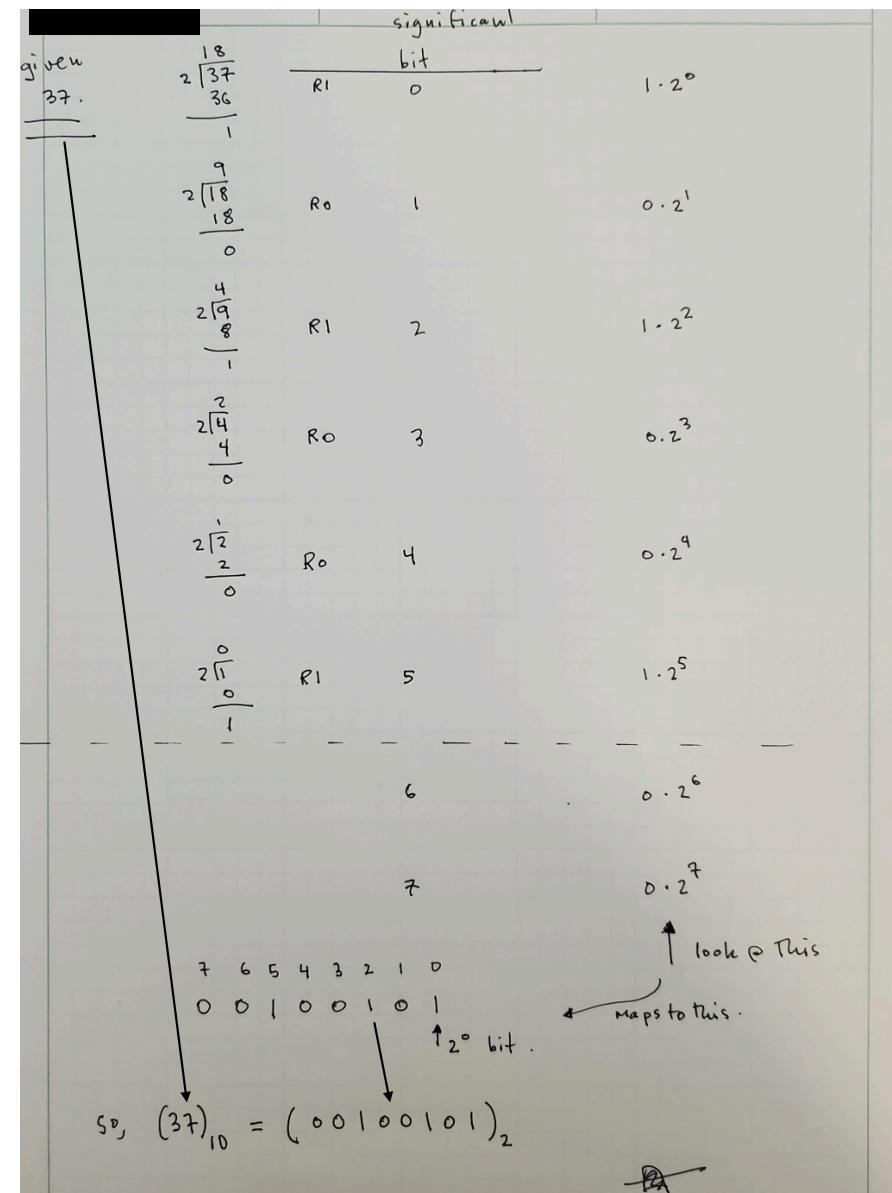
0	0	1	0	0	1	0	1
---	---	---	---	---	---	---	---

b7	b6	b5	b4	b3	b2	b1	b0
----	----	----	----	----	----	----	----

Little endian

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

b0	b1	b2	b3	b4	b5	b6	b7
----	----	----	----	----	----	----	----



A word on error

Absolute error = approximate value – true value

Relative error = $\frac{\text{absolute error}}{\text{true value}}$

10^{-p} relative error in quantity has $\sim p$ correct significant digits in its decimal representation

Approximate value = (true value) \times (1 + relative error)

Do you see any problems here?

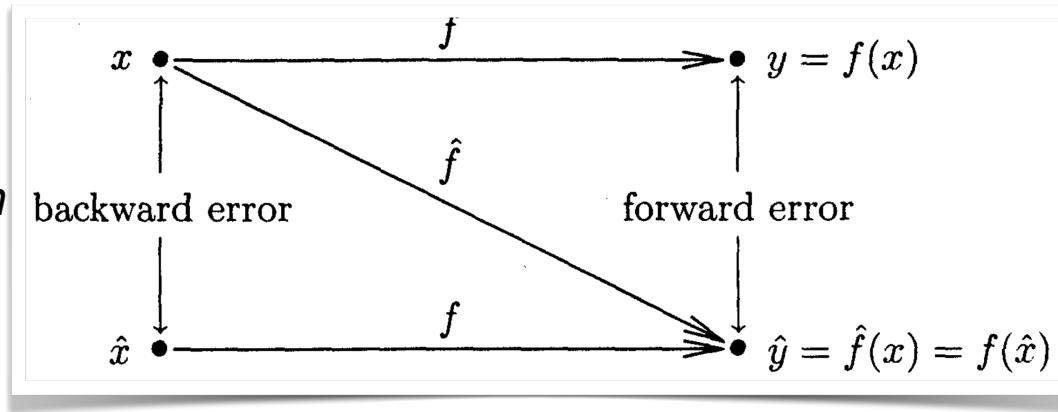
Ans. We don't usually know the *true value* or correct answer if we have committed to computing approximations. However, this is useful for composing tests where known values can be used to test fidelity of implemented algorithm.

$$\begin{aligned}
 \text{Total error} &= \hat{f}(\hat{x}) - f(x) \\
 &= (\hat{f}(\hat{x}) - f(\hat{x})) + (f(\hat{x}) - f(x)) \\
 &= \text{computational error} + \text{propagated data error}
 \end{aligned}$$

desired true result is $f(x)$

\hat{f} , is an approximation to the actual function

\hat{x} , is an approximation to the correct input



- Forward error: Measures the error in the computed function value.
- Backward error: Measures how much the input must be changed to explain the computed result.
- Total error: Includes computational error (caused by finite precision) and propagated error (caused by input errors affecting results).
- Numerical stability: An algorithm is stable if its backward error is small—meaning it computes the exact answer to a slightly different problem. (see condition number)

(in)stability

The condition number for a problem indicates the sensitivity of its computed values to perturbations in the input values.

$$\text{Condition number} = \frac{|(f(\hat{x}) - f(x))/f(x)|}{|(\hat{x} - x)/x|} = \frac{|(\hat{y} - y)/y|}{|(\hat{x} - x)/x|} = \frac{|\Delta y/y|}{|\Delta x/x|}$$

- Forward error: Measures the error in the computed function value.
- Backward error: Measures how much the input must be changed to explain the computed result.

$$f(2) = \sqrt{2}$$

The exact value (to high precision) is:

$$\sqrt{2} \approx 1.414213562373095$$

However, suppose our computer gives the approximate value:

$$\tilde{f}(2) = 1.4$$

The forward error is:

$$\begin{aligned}\text{Forward Error} &= |\tilde{f}(2) - f(2)| \\ &= |1.4 - 1.414213562373095| \\ &\approx 0.01421356\end{aligned}$$

This tells us how far off the computed result is from the true square root.

Forward and Backward Error. As an approximation to $y = \sqrt{2}$, the value $\hat{y} = 1.4$ has an absolute forward error of

$$|\Delta y| = |\hat{y} - y| = |1.4 - 1.41421\dots| \approx 0.0142,$$

or a relative forward error of about 1 percent. To determine the backward error we observe that $\sqrt{1.96} = 1.4$, so the absolute backward error is

$$|\Delta x| = |\hat{x} - x| = |1.96 - 2| = 0.04,$$

or a relative backward error of 2 percent.

$$\sqrt{\hat{x}} = 1.4$$

Squaring both sides:

$$\hat{x} = 1.4^2 = 1.96$$

The backward error is:

$$\begin{aligned}\text{Backward Error} &= |x - \hat{x}| \\ &= |2 - 1.96| = 0.04\end{aligned}$$

The sensitivity of $f(x) = \sqrt{x}$ to input errors is given by its derivative:

$$f'(x) = \frac{1}{2\sqrt{x}}$$

At $x = 2$:

$$f'(2) = \frac{1}{2 \times 1.414} \approx 0.3536$$

If the input has an error $\Delta x = 0.04$, the propagated error in the output is:

$$\text{Propagated Error} \approx 0.3536 \times 0.04 = 0.01414$$

Adding the computational error 0.01421356, the total error is:

$$\text{Total Error} \approx 0.01421356 + 0.01414 = 0.02835$$

(in)stability

The condition number for a problem indicates the sensitivity of its computed values to perturbations in the input values.

$$\text{Condition number} = \frac{|(f(\hat{x}) - f(x))/f(x)|}{|(\hat{x} - x)/x|} = \frac{|(\hat{y} - y)/y|}{|(\hat{x} - x)/x|} = \frac{|\Delta y/y|}{|\Delta x/x|}$$

Consider $\cos(x)$, $\cos(\pi/2)=0$. Let h perturb x around $\pi/2$, $\cos(x+h)$.

By definition:

$$\begin{aligned}\text{abs_error} &= \cos(x+h) - \cos(x) \sim -h \sin(x) = -h \text{ for } x \sim \pi/2 \\ \text{rel_error} &= -h \tan(x) = \text{inf (divide by zero)}\end{aligned}$$

$\pi/2 \sim 1.570796327\dots$ look at a value close to this and perturb it:

$$\begin{aligned}x &= 1.57078, \cos(x) = 1.6327e-5 \\ x+h &= 1.57079, \cos(x+h) = 6.327e-6\end{aligned}$$

$$\begin{aligned}d(\text{input}) &= 6.366e-6 \\ d(\text{output}) &= 6.125e-1 \\ d(\text{output})/d(\text{input}) &\sim 96212.7\end{aligned}$$

HIGHLY SENSITIVE

IEEE 754 fp data structures

floating-point number $x \in \mathbb{F}$

$$x = \pm \left(d_0 + \frac{d_1}{\beta} + \frac{d_2}{\beta^2} + \cdots + \frac{d_{p-1}}{\beta^{p-1}} \right) \beta^E$$

a string of p base- β digits $d_0 d_1 \cdots d_{p-1}$

mantissa
exponent

$$0 \leq d_i \leq \beta - 1, \quad i = 0, \dots, p - 1$$

sign, exponent,
mantissa stored
in distinct fields

β Base or radix

p Precision

$[L, U]$ Exponent range

$L \leq E \leq U$

System	β	p	L	U
IEEE SP	2	24	-126	127
IEEE DP	2	53	-1,022	1,023

Data structure for floating point numbers

IEEE

$$V = (-)^s \cdot M \cdot 2^E$$

s , sign , $s=0$, positive
 $s=1$, negative

sign bit

M , significand or mantissa
[$1, 2^{-\epsilon}$], or [$0, 1-\epsilon$] } n bits

$$\text{frac } f = f_{n-1} f_{n-2} \dots f_0$$

encoded value depends if $E=0$ or not .

E , exponent = $e_{k-1} e_{k-2} \dots e_1 e_0$ } k bit
exponent field

SP floating point (float in C or C++)

s , 1-bit

exp , $k = 8$ bit

frac , $n = 23$ bit

32 bit number.

IEEE

$$V = (-)^s \cdot M \cdot 2^E$$

s , sign, $s = 0$, positive
 $s = 1$, negative

M , significand or mantissa
 $[1, 2 - \epsilon]$, or $[0, 1 - \epsilon]$

$$\text{frac } f = f_{n-1} f_{n-2} \dots f_0$$

encoded value depends if $E = 0$ or not.

E , exponent = $e_{k-1} e_{k-2} \dots e_1 e_0$ { k bit exponent field }

DP floating point (double in C or C++)

s , 1-bit

exp , $k = 11$ bit

frac , $n = 52$ bit

64 bit number.

Normalized (most common)

Normalized,

when exp has not $e_i = 0 \forall i$
and not $e_i = 1 \forall i$.

$$E = e - \text{bias} , e \text{ unsigned } \# e_{k-1} \dots e_1 e_0$$

$$\text{bias} = 2^{k-1} - 1$$

$$= 2^{8-1} - 1 = 2^7 - 1 = 127 \quad \underline{\text{SP}}$$

$$= 2^{11-1} - 1 = 2^{10} - 1 = 1023 \quad \underline{\text{DP}}$$

$$\text{So, } -126 \leq E \leq 127 \quad \text{SP}$$

$$-1022 \leq E \leq 1023 \quad \text{DP}$$

Not all 0s, not all 1s

IEEE

$$V = (-)^S \cdot M \cdot 2^E$$

$S, \text{ sign}, S=0, \text{ positive}$
 $S=1, \text{ negative}$ } sign bit

$M, \text{ significand or mantissa}$ } n bits
 $[1, 2-E], \text{ or } [0, 1-E]$

$$\text{frac } f = f_{n-1} f_{n-2} \dots f_0$$

encoded value depends if $E=0$ or not.

$E, \text{ exponent} = e_{k-1} e_{k-2} \dots e_1 e_0$ } k bit exponent field

$$\text{frac}, \quad 0 \leq f < 1 \quad 0.f_{n-1} f_{n-2} \dots f_1 f_0$$

↑
binary pt to left most significant bit.

$$M = 1 + f$$

$$\text{thus, } V = (-)^S (1+f) 2^{(e-(2^k-1))}$$

IEEE

$$V = (-)^S \cdot M \cdot 2^E$$

S , sign, $S=0$, positive
 $S=1$, negative

↳ sign bit

M , significand or mantissa

$[1, 2-\epsilon]$, or $[0, 1-\epsilon]$

} n bits

$$\text{frac } f = f_{n-1} f_{n-2} \dots f_0$$

encoded value depends if $E=0$ or not.

E , exponent = $e_{k-1} e_{k-2} \dots e_1 e_0$

} k bit exponent field

**Special case: if $e=1s$ and $f \neq 0$ NaNs

De Normalized

$$e = e_{k-1} \dots e_1 e_0 \quad \text{and} \quad e_i = 0 \quad \forall i$$

"gradual
under
flow"

$$E = 1 - \text{bias}$$

$$M = f$$

$$\text{thus } V = (-)^S \cdot f \cdot 2^{1-(2^{k-1}-1)}$$

$$= (-)^S \cdot f \cdot 2^{2-2^{k-1}}$$

ie. represent 0

$$-0.0 \neq +0.0$$

↑
different
sign bit!

Denormalized

* if $e_i = 1 \neq i$ and $M=0$ $\rightarrow \pm \infty$ cases.

IEEE

$$V = (-)^s \cdot M \cdot 2^E$$

s , sign, $s=0$, positive
 $s=1$, negative

M , significand or mantissa
 $[1, 2-\epsilon]$, or $[0, 1-\epsilon]$

$$\text{frac } f = f_{n-1} f_{n-2} \dots f_0$$

encoded value depends if $E=0$ or not.

E , exponent = $e_{k-1} e_{k-2} \dots e_1 e_0$

$\left. \begin{array}{l} k \text{ bit} \\ \text{exponent field} \end{array} \right\}$

Smallest Normalized

1 in least significant bit, 0s in all other bits of exponent field.

$$M = 1$$

$$E = -2^{k-1} + 2$$

$$V = 2^{-2^{k-1} + 2}$$

Largest Normalized

sign bit 0

least significant exp. bit 0, the rest 1s

$$f = 1 - 2^{-n}$$

$$E = 2^{k-1} - 1$$

$$M = 2 - 2^{-n}$$

$$V = (2 - 2^{-n}) \cdot 2^{2^{k-1} - 1}$$

IEEE

$$V = (-)^S \cdot M \cdot 2^E$$

S , sign, $S=0$, positive
 $S=1$, negative

M , significand or mantissa
[$1, 2-\epsilon$], or [$0, 1-\epsilon$] } n bits

$$\text{frac } f = f_{n-1} f_{n-2} \dots f_0$$

encoded value depends if $E=0$ or not.

E , exponent = $e_{k-1} e_{k-2} \dots e_1 e_0$ } k bit exponent field

Smallest DeNormalized

1 in least significant bit, 0's otherwise

$$M = f = 2^{-n}$$

$$E = -2^{k-1} + 2$$

$$V = 2^{-n} \cdot 2^{-2^{k-1} + 2}$$

Largest DeNormalized

Exponent all 0s

frac, all 1s

$$M = f = 1 - 2^{-n}$$

$$E = -2^{k-1} + 2$$

$$V = (1 - 2^{-n}) \cdot 2^{-2^{k-1} + 2}$$

casting example

Example 1

cast integer as SP floating point number

$$\text{consider } (12345)_{10} = (11000000111001)_2$$

normalize. shift to right of binary pt.

$$(12345)_{10} = (1.1000000111001)_2 \times 2^{13}$$

IEEE. frac, drop leading 1, add 10 0s

$$f \sim \underline{1000000111001000000000000} \equiv M$$

exponent, add bias 127 to 13 = 140

$$(140)_{10} = (10001100)_2$$

sign bit, positive #, s=0.

Combine, ~~SEf~~ SEf

$$12345.0 = [01000110010000001110010000000000] \\ \text{SP} \quad \text{32-bit FP representation}$$

toy data structure for HW (just getting you started here)

6 bit representation	$E, k=3, n=2$	$M, f, n=2$	bias, $2^{k-1} - 1 = 3$
$s=1, k=3, n=2$	$V = \pm M \cdot 2^E$		
normalized	$E = e - \text{bias}$	$M = 1 + f$	
	$= (e_2 e_1 e_0) - \text{bias}$	$f \sim f_n, f_{no}$	
		$\cdot f_n, f_{no}$	
		/	
		$\{0, 1\} \cdot \frac{1}{2}$	$\{0, 1\} \cdot \frac{1}{2}$
e_2, e_1, e_0	e	$e - \text{bias}$	f
0 0 0	0	-2	0 0
0 0 1	1	-2	0 1
0 1 0	2	-1	1 0
0 1 1	3	0	1 1
1 0 0	4	1	$1 \cdot \frac{1}{2} + 1 \cdot \frac{1}{2^2} = \frac{3}{4}$
1 0 1	5	2	
1 1 0	6	3	
1 1 1	7	$\pm \infty$	
			$M_{00} = 0 \rightarrow 0 \cdot 2^{-2}$
0 0 0 \rightarrow denormalized, $M=f$, bias=3, $E_{000}=1-3=-2$			
1 1 1, $\pm \infty$			
$M_{00} = 1 + 0 = 1$	$1 \cdot 2^{-2} = \frac{1}{4}$	$1 \cdot 2^1 = 2$	
	$1 \cdot 2^{-1} = \frac{1}{2}$	$1 \cdot 2^2 = 4$	
	$1 \cdot 2^0 = 1$	$1 \cdot 2^3 = 8$	

End Lecture 3 (end week 1)