

***Computer Architecture:
Fundamentals, Tradeoffs, Challenges***

Chapter 11: Floating Point Arithmetic

Yale Patt

The University of Texas at Austin

***Austin, Texas
Spring, 2021***

Scientific Notation and its representation in binary

- **Avogadro's Number: 6.022×10^{23}**
- **Where is the Decimal Point**
 - **What determines where the decimal point is?**
 - **Ergo, "floating point"**
- **Bits for range vs bits for precision**
 - **Von Neumann's Quip**
- **Binary Representations**
 - **32 bits: 1 bit sign, 8 bits exponent, 23 bits of fraction**
 - **64 bits (the default): 1 bit sign, 11 bits exponent, 52 bits fraction**
 - **16 bits (recently for graphics): 1 bit sign, 5 bits exp, 10 bits fra**

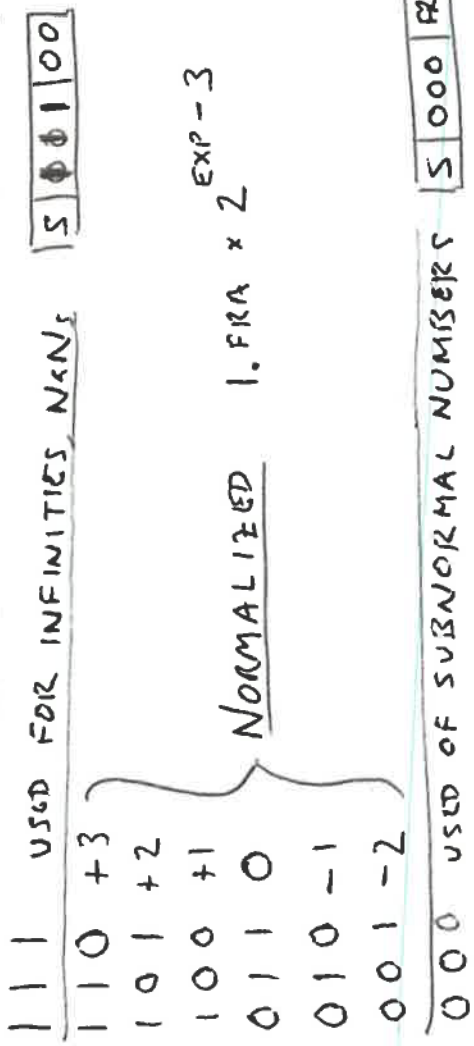
Outline

- **Scientific Notation and its representations in binary**
 - Avogadro's Number
 - Format (Precision vs Range)
 - Common floating point data types
 - Why it is called floating point
- **My 6-bit floating point data type**
 - Redundant most significant bit
 - Excess code for exponents (*infinity, subnormals*)
- **Rounding**
 - Four rounding modes
 - Guard, Round, Sticky bits
 - Wobble
- **Infinities, NaNs**
- **Subnormal numbers (gradual underflow)**
- **Exceptions: Invalid, DivBy0, Underflow, Overflow, Inexact**
 - Quiet vs Signaling

My 6-bit floating point data type



- Excess code for range (excess is 011, i.e. 3)

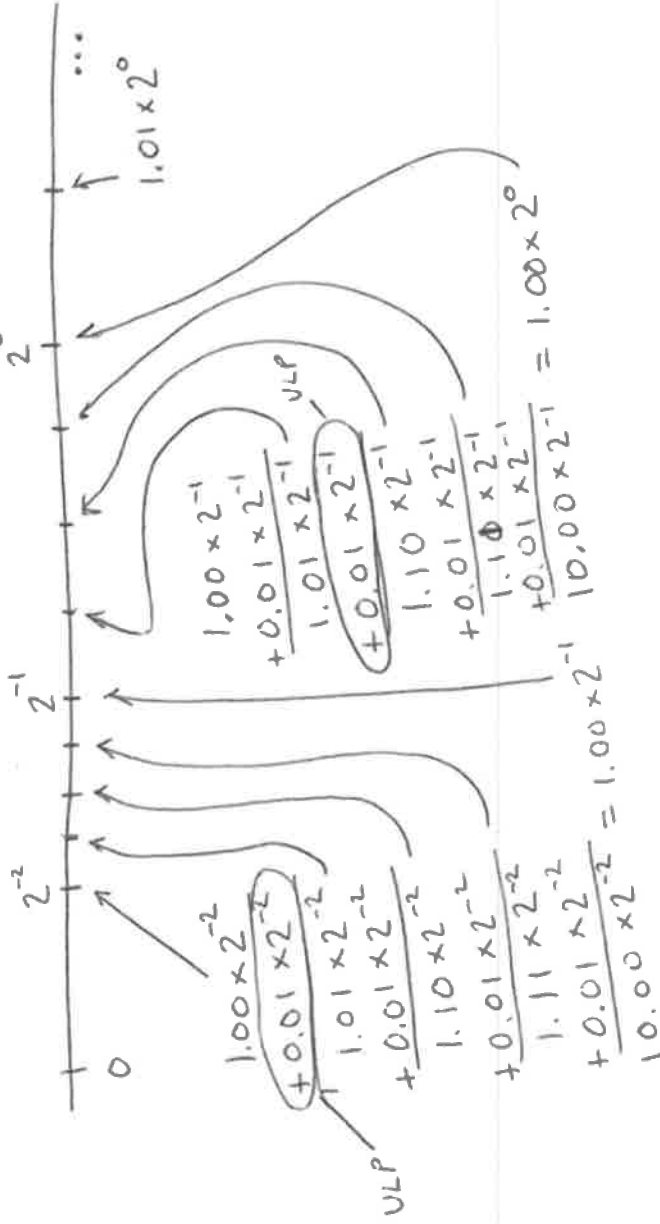


- Signed-magnitude for precision
 - Redundant most significant bit (because radix = 2)
- An example: 6 5/8. 110.101 (Unnormalized)
 - Normalize it: 1.10101 $\times 2^2$

C	1	0	1	0	1	0
---	---	---	---	---	---	---
 - Store in memory: Subsequent read from memory: $1.10 \times 2^2 = 6$
 - We lost 5/8 because we had 2 fraction bits, but we needed 5 fraction bits

My 6-bit floating point data type

- Exact representations on the real line



- Maximum value:

$$\boxed{011011} = 1.11 \times 2^3 = 14$$

- Minimum normalized value:

$$\boxed{01001000} = 1.00 \times 2^{-2} = \frac{1}{4}$$

$$0.01 \times 2^{-2}$$

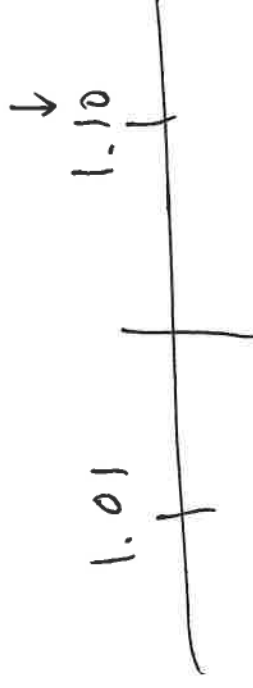
Rounding

- **Why, When**
 - **Why:** when a value can not be represented exactly
 - **When:** Often, most values can not be represented exactly
 - **Example** (with our 6 bit floating point): **1.011 (1 3/8)**

- **Rounding Modes**

- **Round up:** **1.10 (1 1/2)**
- **Round down:** **1.01 (1 1/4)**
- **Round to zero:** **1.01 (1 1/4)**
- **Unbiased Nearest:** **1.10 (1 1/2)**

- Why not 1.01 (1 1/4)?
- 1 1/4 is just as “near” to 1 3/8 as 1 1/2 is
- Unbiased → when equal, round to the value with 0 in the ULP



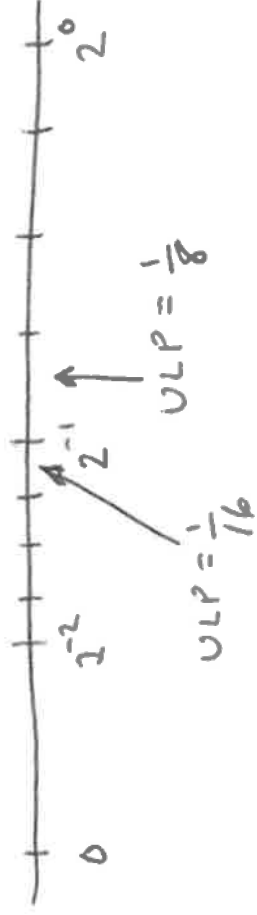
Rounding

- Guard, Round, Sticky bits
 - Extra bits carried along to help in rounding (1 bit or 2 bits?)
 - Example: Add $15 + 2^{1/2}$ (This example uses 3 bits of fraction!)

$15: 1.111 \times 2^3$ $2^{1/2}: 1.010 \times 2^1$

① ALIGN EXP	15: 1.111 00 $\times 2^3$
(NOTE 2 EXTRA BITS)	$2^{1/2}: 0.010 10 \times 2^3$
② ADD	$\begin{array}{r} 17\frac{1}{2}: 10.001 10 \times 2^3 \\ 1.000 11 \times 2^4 \\ \hline \end{array}$
③ NORMALIZE	EXTRA BITS
④ ROUND	

- Wobble (since $1 \text{ ULP} < 2^k$ is half the size of $1 \text{ ULP} > 2^k$)
 - Best case, because radix = 2



Infinity

$$5 \overline{) 11100} + \infty$$

- **Exact vs Overflow**
 - (finite operands) → infinite results
 - Examples: $\tan(90 \text{ degrees})$, 5 divided by 0
- **Infinity is NOT the same as undefined**
 - Example: continued fraction expansion
 - Simple examples:
 - $\text{infinity} + 7 = \text{infinity}$
 - $\text{infinity} + \text{infinity} = \text{infinity}$
 - 5 divided by $\text{infinity} = 0$
 - Code example
 - $X=5$
 - $Y=0$
 - $Z=X/Y$
 - $W= \text{arctan}(Z)$

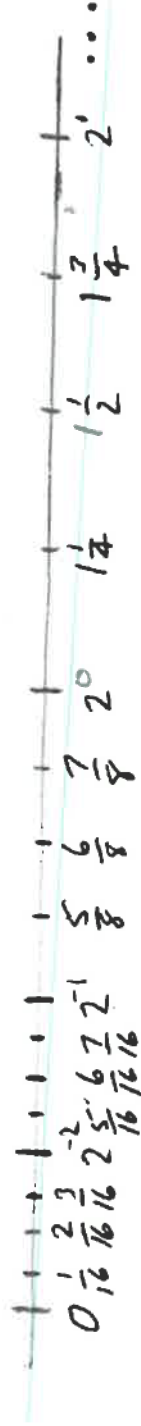
Subnormals

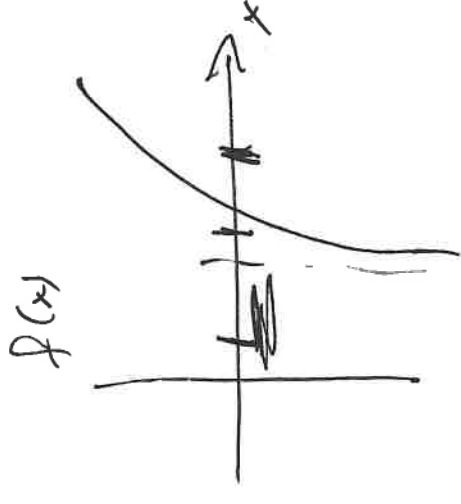
- **Why?**
 - Underflow vs inexact discrepancy was unacceptable
 - 1 divided by underflow produces infinity
- **Tradeoff**
 - Subnormals provide gradual underflow
 - Underflow is no worse than inexact
 - Cost is loss of precision

- **Without subnormals: Store zero**



- **With subnormals: Gradual underflow**





Not a Number (NaN)

- **Examples**
 - *Infinity minus infinity, infinity divided by infinity, 0 divided by 0*
 - *arcsin(2), sqrt(2) (negative number)*
 - *A function that asymptotes*
- **It was here before IEEE Floating Point**
 - *Supercomputers had them, for example*
- **The difference:**
 - *IEEE Floating Pt allows exception handlers to be involved*
 - *Allows correction of the problem and continue processing*

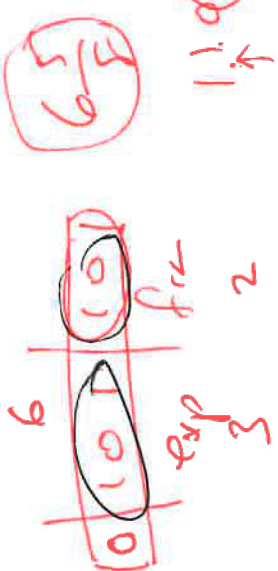
Five Floating Point Exceptions

- **What are they?**
 - **Overflow:** too large to represent in normalized form)
 - **Underflow:** too small to represent as a subnormal number
 - **Inexact:** not a value that can be represented exactly)
 - **Divide by zero:** function (finite arguments) \rightarrow infinity
 - **Invalid:** creation of a NaN
- **Quiet vs Signaling**
 - **Quiet:** sets a sticky bit, handled under program control
 - For example, usual way to deal with “inexact”
 - **Signaling:** takes an exception to deal with the problem
 - For example, usual way to deal with “NaNs”

Arigato!

ROUNDING (S) (S2)

$1.09 \times 2^{(exp - bias)}$



655

1.09101



1.09101 x 2²



WOBBLE

1.09×2^0
 0.01×2^0
 1.01×2^0

EXCESS CODE

3

BIAS

BINARY

1.01×2^0
 $+ 0.01 \times 2^0$
 1.10×2^0
 0.01×2^0
 1.11×2^0
 0.01×2^0
 1.000×2^0

1.00×2^1
 1.10×2^2

ULP

4 ROUNDING

111.7

11.0

1.10 x 2²

110 x 2¹

6

11.08

110 + 3

101 + 2

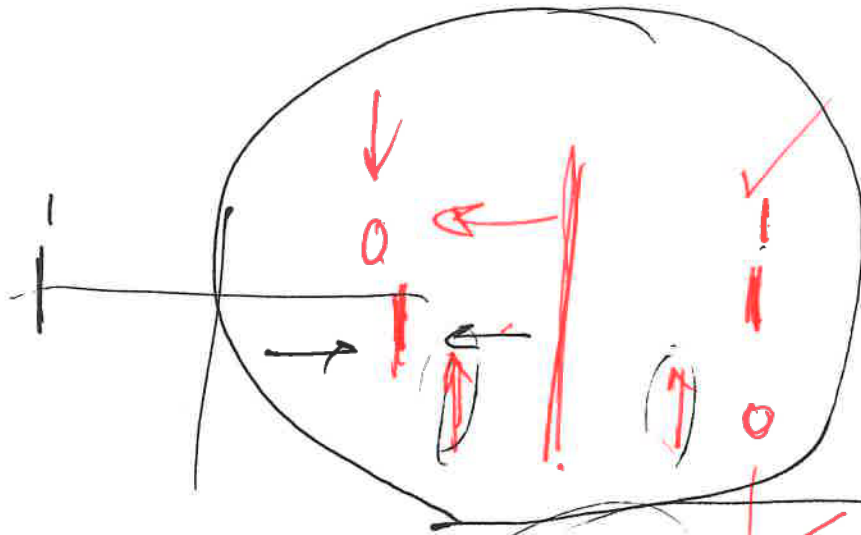
100 + 1

011

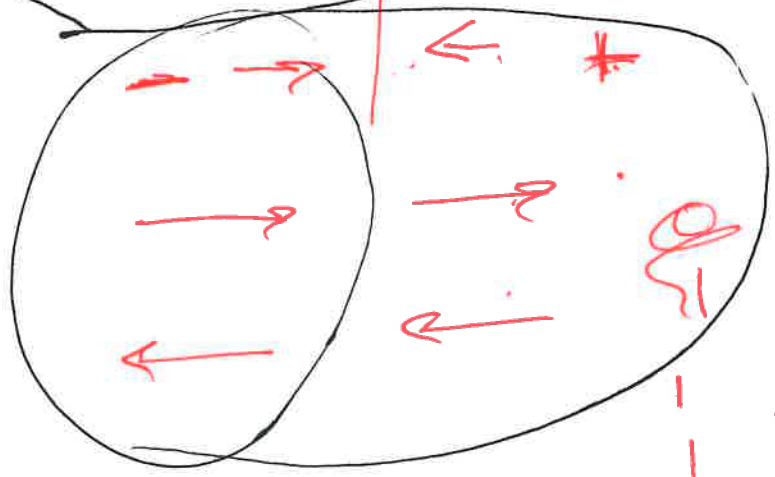
010 - 1

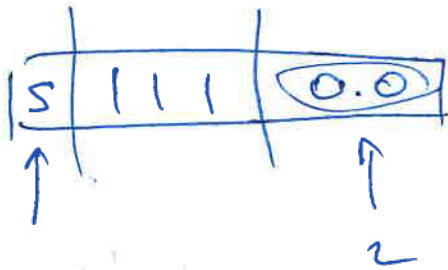
001 - 2

000



NEARBY





INFINITY

NaN

$$\tan\left(\frac{\pi}{2}\right) = \infty$$

$x = 2$
 $y = \arcsin(x)$

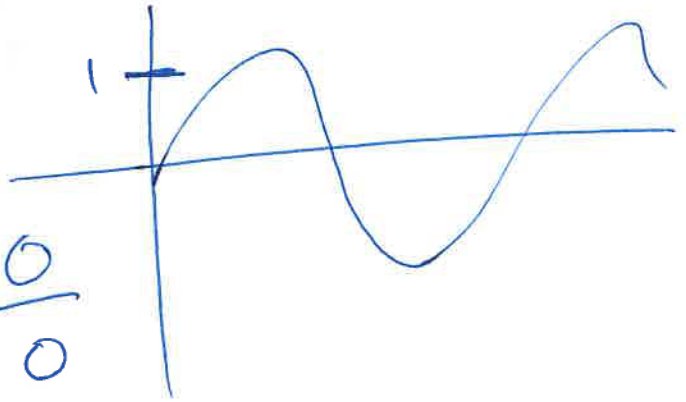
$\frac{3}{4}$

$\frac{5}{0}$

$$x = \frac{3}{4 + \frac{2}{7 + \frac{3}{0}}}$$

$\frac{0}{0}$

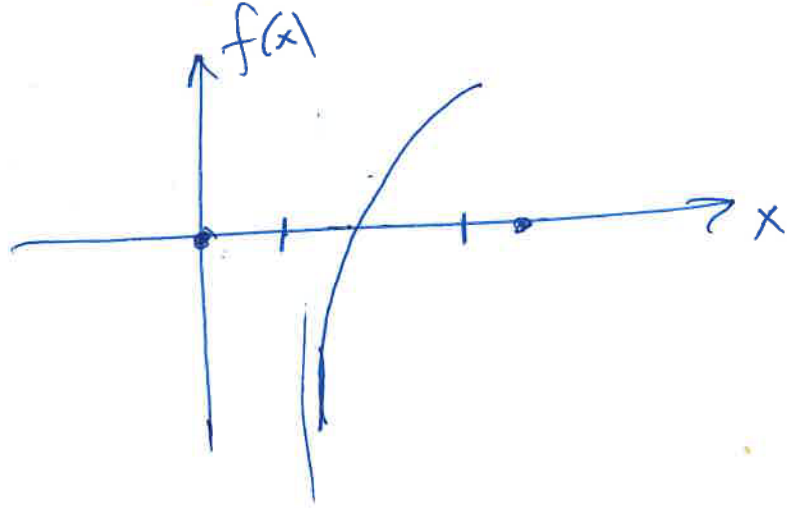
$\frac{\infty}{\infty}$



CONTINUOUS
 FRACTION
 EXPANSION

$$\infty + \infty = \infty$$

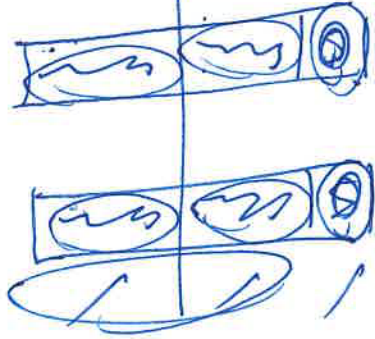
~~$\infty = \infty$~~



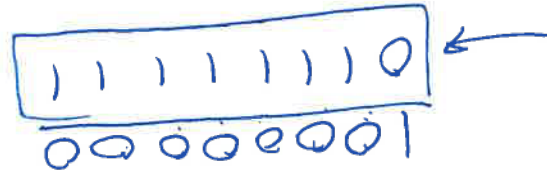
6

$+ 1.10 \times 2$

~~$\frac{2}{3}$~~

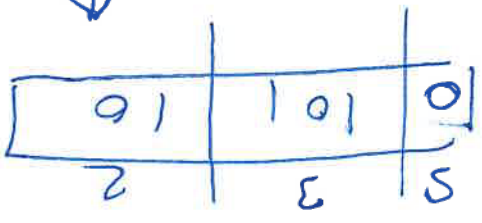


000	000
001	2
010	1
011	0
100	1
101	2
110	3
111	3



ROUNDING

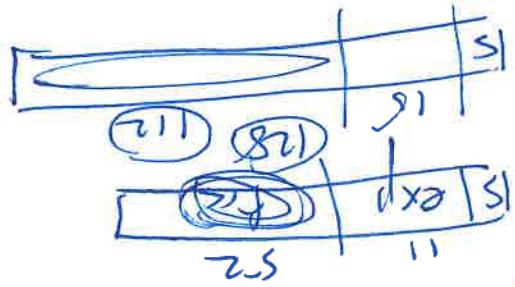
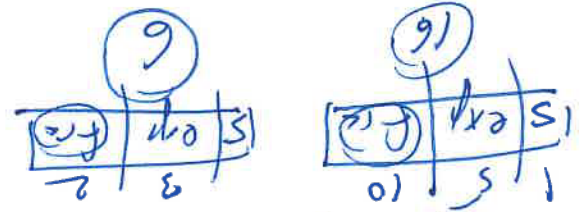
010
110



~~4.9~~
 4.9×10^{-1}

$+ 6 \frac{8}{5}$

32

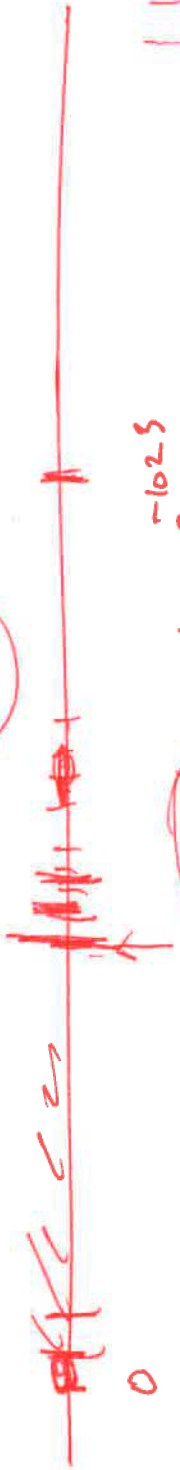


64

$(.025) \times 10^{(2)}$

AFFINE

$\frac{1}{252}$



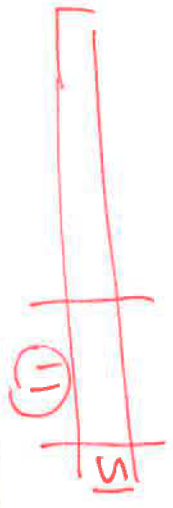
1.0×2^{-1025}

1.0×2^{-1024}

Smaller



Smaller



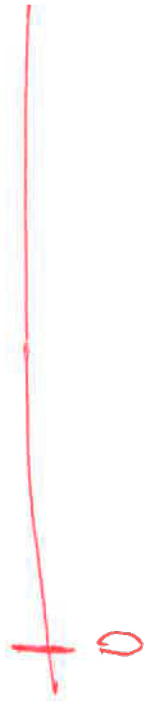
111

100
010 0
001 -1
000 -2

111111111

1000000000
0111111111
-1
-2

-1024
0000 — 0



0	001
---	-----

$$1.0 \times 2^{-2}$$

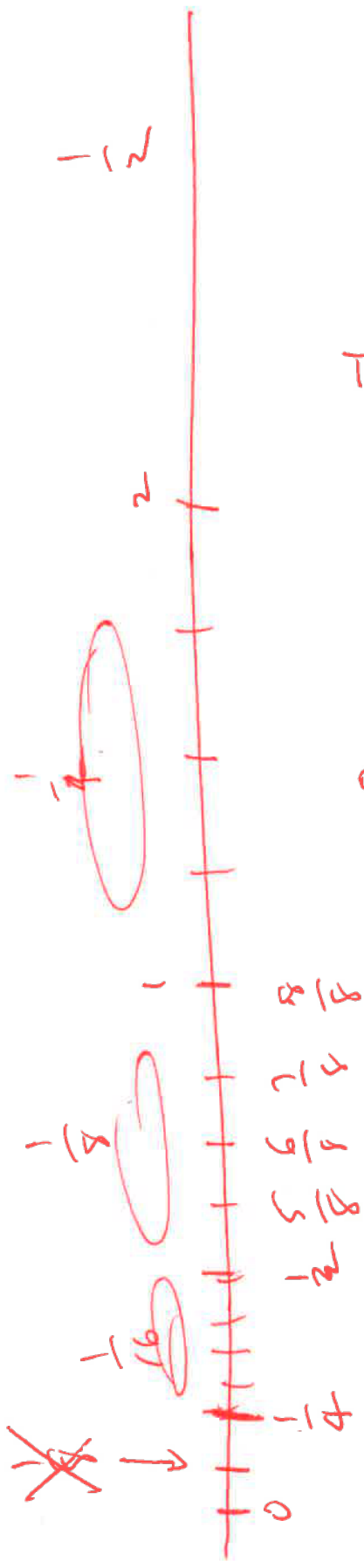
$$\frac{001}{001}$$

0	000
---	-----

$$0.0 \times 2^{-2}$$

$$0.00 \times 2$$

0.01
0.110
0.111
1.000



$$\rightarrow \frac{1.00 \times 2^{-1}}{0.01 \times 2^{-1}}$$

$$\frac{1.01 \times 2^{-1}}{0.1 \times 2^{-1}} = \boxed{\frac{5}{8}}$$

$$1.00 \times 2^{-2}$$

$$1.00 \times 2^{-3}$$

$$\frac{0.000}{0.00}$$

$$0.01 \times 2^{-2} = \frac{1}{4}$$

$$\frac{1.11 \times 2^{-1}}{0.01 \times 2^{-1}} = \frac{10,000 \times 2^{-1}}{10,000 \times 2^{-1}}$$

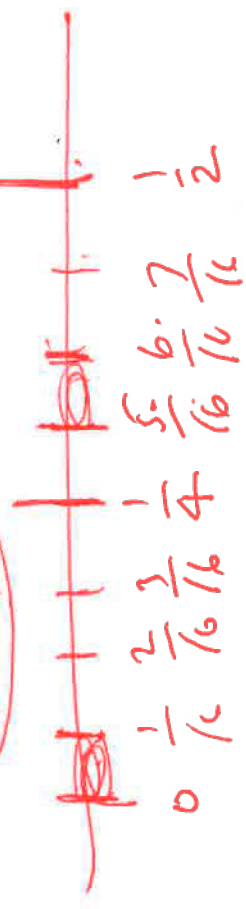
$$\frac{3}{2} \times \frac{1}{2} = \frac{6}{8}$$

$$1.01 \times 2^{-2}$$

$$0.01 \times 2^{-2}$$

NOT A 01

NO 01



$$1.00 \times 2^{-2}$$

$$0.00 \times 2^{-2}$$