3.3. Floating point addition without rounding
– while aligning, keeps all
the bits of the RShifted mantissa

$$X \pm Y = (X_n \pm Y_n \cdot 2^{(Y_e - X_e)}) \cdot 2^{X_e}$$
$$\text{if } X_e \geq Y_e$$

m+2 mantissas
1 sign
1 underflow (separate bit)
m front part of signif.

adder / subtracter 1

declare registers $A[(m+1):0]$, $M[(m+1):0]$, $E_1[(e-1):0]$,
$\quad E_2[(e-1):0]$, $E[(e-1):0]$, A_COUT, ERROR;

declare bus INBUS$[(m+e+1):0]$, OUTBUS$[(m+e+1):0]$;

BEGIN:    A_COUT := 0, ERROR := 0,

INPUT:    $E_1$ := INBUS($X_e$), $A$ := INBUS($X_n$);
$\quad\quad E_2$ := INBUS($Y_e$), $M$ := INBUS($Y_n$);

COMPARE:    $E := E_1 - E_2$;

ALIGN:    if $E < 0$ then $A$ := RShift($A$), $E := E+1$, goto ALIGN;
   if $E > 0$ then $M$ := RShift($M$), $E := E-1$, goto ALIGN;

ADD/SUBTRACT:    $A := A \pm M$, $E := \max(E_1, E_2)$;

OVERFLOW:    if A_COUT == 1 then begin
     if $E == E_{MAX}$ then goto ERROR,
     $A$ := RShift($A$), $E := E+1$, goto END;   ← unconditional jump
   end

ZERO:    if $A == 0$ then $E := 0$, goto END;   ← conditional jump

NORMALISE:    if isNormalized($A$) == 1 then goto END,

UNDERFLOW:    if $E > E_{min}$ then
     $A$ := LShift($A$), $E := E-1$, goto NORMALISE;

ERROR:    ERROR := 1;
END:

Pseudolanguage:
1) declare register → name
     → width
   – concatenation: A_COUT . A    {A_COUT, A} in Verilog.

2) declare bus → name
     → width   – $2m + 2e + 4$ bits.

   – unified IOBUS $m + e + 2$ bit

3) Synchronous execution
- non-conflicting operations: executed concurrently, separated by , (same cycle)
- sequential operations: executed sequentially separated by ; (different cycle)

4) := assign operator
- hardwired operations: $\underline{Max}$, $\underline{RShift}$, $\underline{LShift}$, is Normalized

5) flow control
- unconditional jump: go to $\overline{OUT}$
- conditional jump: if $A = 0$ then $E := 0$, go to $\overline{OUT}$

6) simultaneous read/write to registers/buses:
$$A[7] := \Pi[7] \text{ and } \overline{Q[0]}, \overline{Q[0]} := 0;$$

Comments regarding the f.p. algorithm:

A) operands are on $m + e + 2$ bits
- 1 sign
- $e$ exponent } $m+2$ of sections
- $m+1$ significand
  - 1 hidden
  - in fractal part
- received in packed format from the INBUS

B) registers $A$ and $M$
- store significant + significand ($\equiv$ mantissa).
- RShift capabilities for alignment.
- A has LShift capabilities for normalisation.
- A is extended by the A_cout bit (flag)



c) register $E$:
- has $++$ / $--$ capabilities
- loaded with the maximum of $E_1, E_2$.

D) flag OVERR
- indicates exceptions: OVERFLOW, UNDERFLOW

## 3.4. Rounding and normalisation rules for f.p. addition.

Consider $X_n = \underbrace{1}_{\text{hidden bit}} . \overbrace{X_{m-2}}^{2^{-1}} \overbrace{X_{m-3}}^{2^{-2}} \cdots X_i \cdots X_1 \overbrace{X_0}^{2^{-m+1}}$   $\quad X_E \geq Y_E$
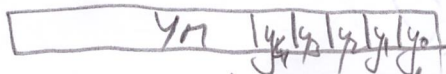
$Y_n = 1 . Y_{m-2} Y_{m-3} \cdots Y_i \cdots Y_1 Y_0$

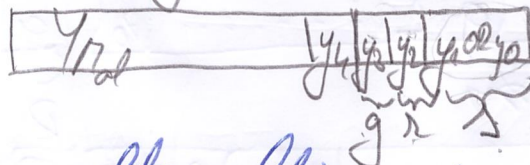alligment of $Y_n$ by RShifting with $d$ bits    $d = |X_E - Y_E|$

$Y_n'$ alignment $\Big\{$

with infinite precision: keeps all bits of $Y_n$ including bits with weight $< 2^{-m+1}$

with finite precision: keep only **3 bits out** of all bits having weight $< 2^{-m+1}$

the 3 bits : the sticky bits.

$\Big\{$

$g$ : guard bit, weight of $2^{-m}$
  guards against loss of precision after $Y_n$ align.

$r$ : round bit, weight of $2^{-m-1}$
  used for result rounding.

$s$ : sticky bit, weight of $2^{-m-2}$
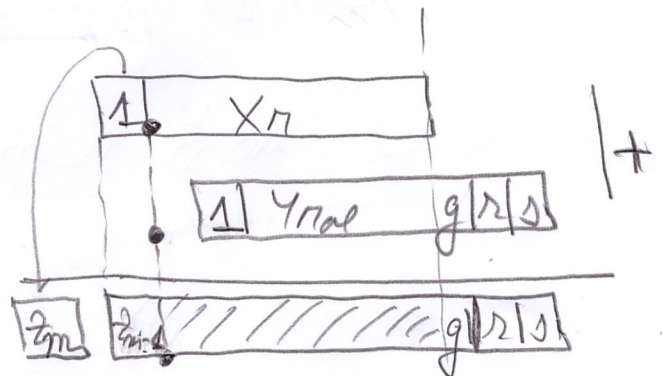  — obtained as a logic OR of all other less' significant bits that were RShifted out of $Y_n$, except $g$ and $r$



$Y_n$ alligment by 4 bit RShift.



$Y_{nal} = Y_n$ after alligment

$\longrightarrow Z_n = X_n + Y_{nal}$

Consider $z_n$ to be

$$z_n = z_m z_{m-1} . z_{m-2} z_{m-3} \cdots z_1 z_0 | g\ r\ s$$

Normalization of $z_n \rightarrow z_{nn}$

only need 2 bits for implementing the rounding modes

$$z_{nn} = 1 . z_{m-2_n} z_{m-3_n} \cdots z_{1_n} z_{0_n} | R\ S$$

cases.

1) $z_m = 1$
⟹ 1-bit R Shift of $z_n$   $z_e ++;$

2) $z_m = 0, z_{m-1} = 1$
⟹ $z_n$ is already normalized

3) $z_m = 0, z_{m-1} = 0, z_{m-2} = 1$
⟹ 1-bit L Shift of $z_n$   $z_e --$

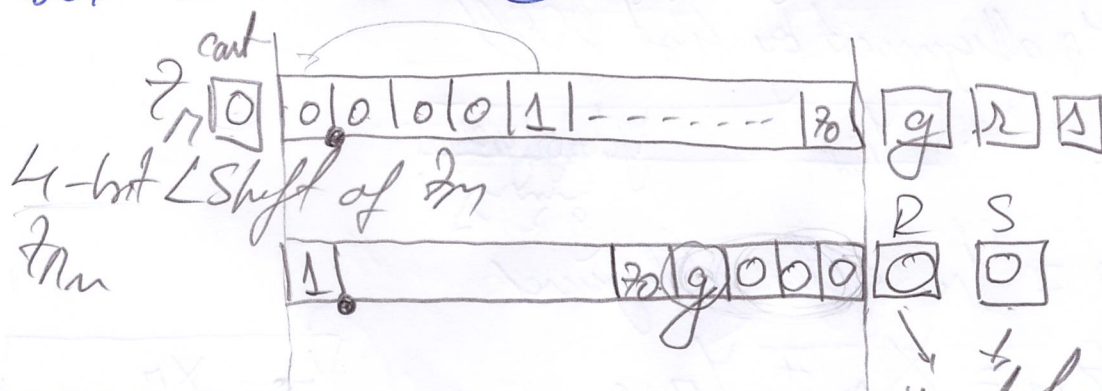4) $z_m = 0, z_{m-1} = 0, z_{m-2} = 0,$
$z_{m-3} = 1$
⟹ 2-bit L Shift of $z_n$   $z_e -= 2;$

| 1. | | | | | | | R | S |
|---|---|---|---|---|---|---|---|---|
| 1. | $z_{m-1}$ | $z_{m-2}$ | ---- | $z_2$ | $z_1$ | $z_0$ | | (g or r or s) |
| 1. | $z_{m-2}$ | $z_{m-3}$ | — | $z_1$ | $z_0$ | g | | (r or s) |
| 1. | $z_{m-3}$ | $z_{m-2}$ | — | $z_0$ | g | r | s | |
| 1. | $z_{m-4}$ | $z_{m-5}$ | -- | g | ⓪ | ⓪ | ⓪ | |

— if a 2-bit or more bits L Shift is required for $z_n$'s normalization

 — append g bit to $z_n$, after $z_0$

 — complete all remaining positions with 0s

 — set R = S = 0

$z_n$ [cout 0] [ 0 | 0 | 0 | 0 | 0 | 1 ------------ $z_0$ ] [ g ] [ r ] [ s ]

4-bit L Shift of $z_n$

$z_{nn}$ [ 1 ] [ ........ $z_0$ g 0 0 0 ] [R 0] [S 0]

used for implementing the rounding modes

Rounding of $z_{m}$ : uses R and S bits.
— use 2 "fractional" bits "X" $= x_{m-1} x_{m-2} \cdots x_1 x_0 . RS$
to 0 discard RS

| Rounding mode | $z_{m} > 0$ | $z_{m} < 0$ | |
|---|---|---|---|
| to 0 | $\overline{\text{(discard R and S)}}$ | $\overline{\text{(discards R and S)}}$ | |
| towards $-\infty$ | $\overline{\text{(discard R and S)}}$ | if (R or S) then $z_{m} - 1$ | $-1.0005$ $\to -2$ |
| towards $\infty$ | if (R or S) then $z_{m} + 1$ | $\overline{\text{(discard R and S)}}$ | $+7.0000001$ $\to +8$ |
| to nearest even | if (R and (S or $z_{0m}$)) then $z_{m} + 1$ | if (R and (S or $z_{0m}$)) then $z_{m} - 1$ | |

$z_{m} = \underbrace{1 . z_{m-2m} \cdots z_{0m} z_{0m}}_{} \mid \underbrace{RS}_{}$

round to nearest even for positives:

a) if fractional part $< \frac{1}{2}$ $\Rightarrow$ discard R and S

$\quad 4.(49) \to 4$

b) if fractional part $= \frac{1}{2}$ and the integer part is an even no $\Rightarrow$ discard R and S

$\quad 4.(5) \to 4$

c) if fractional part $= \frac{1}{2}$ and the integer part is an odd no $\Rightarrow +1$ to int. part

$\quad 3.5 \to 4$

d) if fractional part $> \frac{1}{2}$ $\Rightarrow +1$ to integer part

$\quad 4.55 \to 5$

in c) and d) = odd & unit. using 2 bits of fr. part. when to $+1$ to $z_{m}$: $\Rightarrow R = S = 1 (.75)$
$\searrow R = 1, S = 0 (.5)$
and $z_{0m} = 1$

| | R | S hit part |
|---|---|---|
| R | 0 | 0 (.0) |
| S | 0 | 1 (.25) |
| if integer part is odd | 1 | 0 (.5) |
| $\boxed{+1}$ | 1 | 1 (.75) |

condition for +1 to $z_{rn}$ when rounding to nearest even.

$$R \cdot S + R \cdot \bar{S} \cdot 7_{on} = R \cdot (S + 7_{on})$$

$(1 \quad t = 0.75)$

$(1 \quad 0 = 0.5)$

! if rounding generates a carry out
$\Rightarrow$ post normalisation.