

EE 486 : lecture 2 , the floating point numbers

M. J. Flynn

---

Computer Architecture & Arithmetic Group1Stanford University

FPN: basic idea

- $m \in M$ , machine numbers;  $r \in R$ , result numbers;  $R$  includes all  $M$
- A FPN is  $\pm 0.M \times \beta^{\pm e}$
- $\begin{matrix} \pm & 0 & . & M & \times & \beta & ^{\pm e} \\ & \uparrow & & \uparrow & \uparrow & \uparrow & \uparrow \\ & 1 & & 2 & 3 & 4 & 5 & 6 \end{matrix}$
- 6 attributes define a FPN

---

Computer Architecture & Arithmetic Group2Stanford University

FPN: basic terms

- $m \in M$ , machine numbers;  $r \in R$ , result numbers ( $R$  includes  $M$  and more)
- Expressed as six implicit or explicit parameters:
  - Sign of the number (explicit: usually  $S+M$ )
  - Radix point (implicit: fraction or significand)
  - Mantissa or significand (explicit; length implicit)
  - Radix (implicit)
  - Exponent and sign (usually excess coded)

---

Computer Architecture & Arithmetic Group3Stanford University

FPN, more on the terms

- Exponent is usually represented in form of characteristic (or an excess code)
- Characteristic =  $\text{exp} + \text{bias}$ ; where the bias is “usually”  $\frac{1}{2} \beta^n$ ,  $n$  is the number of exp digits (bits)
- Why excess code?
  - Simplifies compares
  - Representation of zero is all zeros (same as integers)
  - Zero has minimum absolute exponent

---

Computer Architecture & Arithmetic Group4Stanford University

FPN, more on terms

- Sign of the number:  $s_f$  is usually in  $S + M$  where a  $(-0)$  is not allowed
- Radix: only 2 or 16 are in use,  $\beta = 2$  provides greater representational capacity,  $\beta=16$  makes shifting easier and faster.
- Radix point: 0.1 or 1.0 with  $\beta = 2$ ; either immediately before or after the MSD of the fraction.
- Mantissa, fraction, significand synonyms

---

Computer Architecture & Arithmetic Group5Stanford University

More FPN concepts

- Maximum representable number,  $\text{max} = \beta^{\text{emax}} M_{\text{max}} = \beta^{\text{emax}} (1 - \beta^{-m})$ ,  $m$  is mantissa bits
- Minimum representable (non 0) number,  $\text{min} = \beta^{\text{emin}} M_{\text{min}}$
- Range is  $[\text{min}, \text{max}]$
- Precision is  $\beta^{-m}$  same as ulp (unit in the last place)
- Gap is  $\beta^{\text{exp}} \times (\text{ulp})$ .

---

Computer Architecture & Arithmetic Group6Stanford University

### Normalized nos., over/under flows and hidden "1"

- Normalized no. has non 0 msd; unnormalized in anything else; denormalized is unnormal with min exp
- Overflow: result > max... then either set 0 flag or set result to infinity representation.
- Underflow: min > result > 0.. Then either set *u* flag or set result to zero.
- Hidden 1: if  $\beta = 2$  then msb of mantissa = 1 so imply it (1).xxx or 0.(1)xxx

Computer Architecture & Arithmetic Group 7
Stanford University

### Rounding

- Rounding is the mapping of result number  $r \in \mathbb{R}$ , to an adjacent machine number  $m \in \mathbb{M}$ .
- Various types:
  - RZ(*r*) truncation of *r* to *m*
  - RN (*r*) round to nearest, up in a tie or round to nearest even
  - RP (*r*) round to positive infinity
  - RM (*r*) round to minus infinity

Computer Architecture & Arithmetic Group 8
Stanford University

### Four classic FPN systems

- Generic binary, classic 36b old style *fpns*, ca. 1952- 1990.
- Hex,  $\beta = 16$ , main frames
- Cray, quick and "dirty" binary, still useful for signal processing applications.
- IEEE standard binary, the most complicated and, by now, the most widely used

Computer Architecture & Arithmetic Group 9
Stanford University

### Generic binary

1 8
27

$s_f$	exp	mantissa
-------	-----	----------

36

- Bias =  $2^8/2 = 128$
- Max =  $2^{127} (1 - 2^{-27})$ , Min =  $2^{-128} (2^{-1})$
- Precision =  $2^{-27}$
- Round is RZ(*r*) only, no hidden 1
- Set 0 flag on overflow, set *r* to 0 on underflow.

Computer Architecture & Arithmetic Group 10
Stanford University

### Mainframes (S 360 ca 1963)

1 7
24

$s_f$	char	mantissa
-------	------	----------

32

Mantissa has 6 hex digits; note a normalized number may have leading digit 0001. The exponent (char) uses a binary radix.

So bias is  $2^7/2 = 64$ , emax is  $127-64 = 63$  and emin =  $-64$

Then Max =  $16^{63} (1 - 16^{-6})$  and Min =  $16^{-64} (16^{-1})$

Set 0 flag on overflows and either set *u* flag or set *r* to 0 on underflows

Computer Architecture & Arithmetic Group 11
Stanford University

### Mainframe FPN, more

- In addition to single (32b), there's double (64b) arranged 1+7+56(14digits) and quad (128b)
- Rounding is RZ(*r*) or RN (*r*)
- Implementations must use guard digit otherwise  $1 \times X$  is not =  $X$  (as 1 is 0001)
- Newer mainframes (S 390) offer both the old format and IEEE.

Computer Architecture & Arithmetic Group 12
Stanford University

