

III. Floating Point Representation

- Floating Point numbers contain the following components
 - 1) Mantissa
 - 2) Mantissa sign (optional): Some method to allow a signed mantissa such as: a) an explicit sign bit, or b) another bit to allow a signed format with the same positive range
 - 3) Exponent
 - 4) Exponent sign (optional): Some method to allow a signed mantissa such as: a) an explicit sign bit, or b) another bit to allow a signed format with the same positive range
 - 5) Exponent base. The value is normally fixed and not explicitly provided. Common values are 2, a power-of-2, or 10 (used for financial applications).

Floating Point Basics

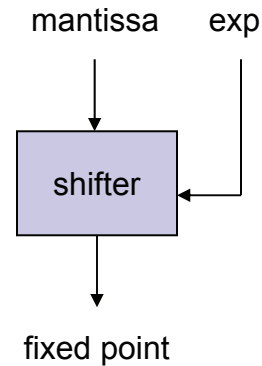
- Normally write: $\text{mantissa} * 2^{\text{exponent}}$
- Alternately, $\text{sign} * \text{mantissa} * 2^{\text{exponent}}$
- In hardware, we normally operate on, transmit, and save only the mantissa and exponent:
(MMMMMM, EEEE)
- Normalized floating point numbers contain no extra (useful) bits at the MSB of the mantissa
 - Example:

00010 * 2 ⁰	not normalized, or "denormalized"
01000 * 2 ⁽⁻²⁾	normalized, with 2's comp. mantissa
10000 * 2 ⁽⁻³⁾	normalized, with unsigned mantissa

Floating Point → Fixed Point Conversion

- If the *exp* is unsigned, the shifter shifts only to the left
- If the *exp* is signed, the shifter must shift to the left and right
- Example:

```
01011. * 22  
01011. << 2  
0101100.
```



Fixed Point → Floating Point Conversion

- Leading 0s/1s detector finds the optimum place to begin selecting bits for the mantissa
- Common pitfall: If the *mantissa* is signed, its sign bits must be maintained!

