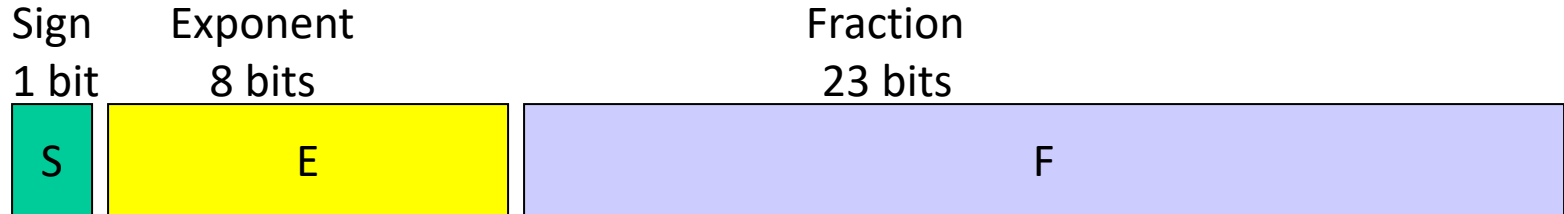


Floating Point

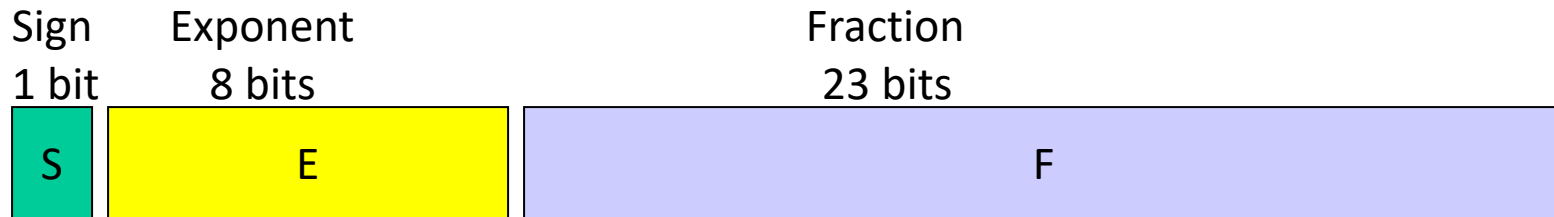
- Normalized scientific notation: single non-zero digit to the left of the decimal (binary) point – example: 3.5×10^9
- $1.010001 \times 2^{-5}_{\text{two}} = (1 + 0 \times 2^{-1} + 1 \times 2^{-2} + \dots + 1 \times 2^{-6}) \times 2^{-5}_{\text{ten}}$
- A standard notation enables easy exchange of data between machines and simplifies hardware algorithms – the IEEE 754 standard defines how floating point numbers are represented

Sign and Magnitude Representation



- More exponent bits → wider range of numbers (not necessarily more numbers – recall there are infinite real numbers)
- More fraction bits → higher precision
- Register value = $(-1)^S \times F \times 2^E$
- Since we are only representing normalized numbers, we are guaranteed that the number is of the form 1.xxxx..
Hence, in IEEE 754 standard, the 1 is implicit
Register value = $(-1)^S \times (1 + F) \times 2^E$

Sign and Magnitude Representation

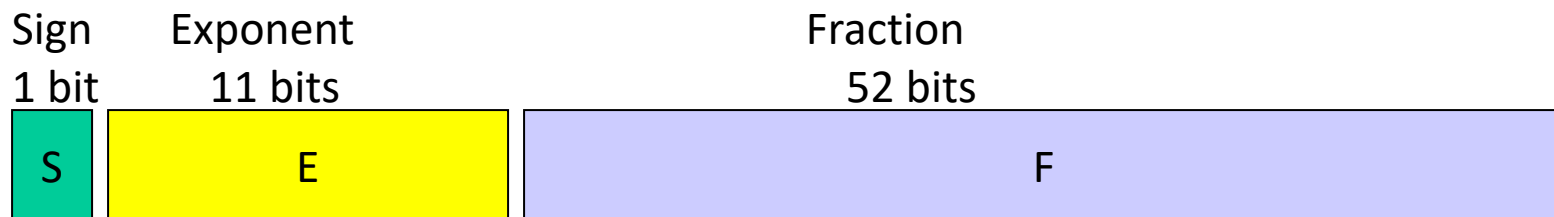


- Largest number that can be represented: $2.0 \times 2^{128} = 2.0 \times 10^{38}$
(not really – see upcoming details)
- Smallest number that can be represented: $1.0 \times 2^{-127} = 2.0 \times 10^{-38}$
(not really – see upcoming details)
- Overflow: when representing a number larger than the max;
Underflow: when representing a number smaller than the min

- Double precision format: occupies two 32-bit registers:

Largest:

Smallest:



Details

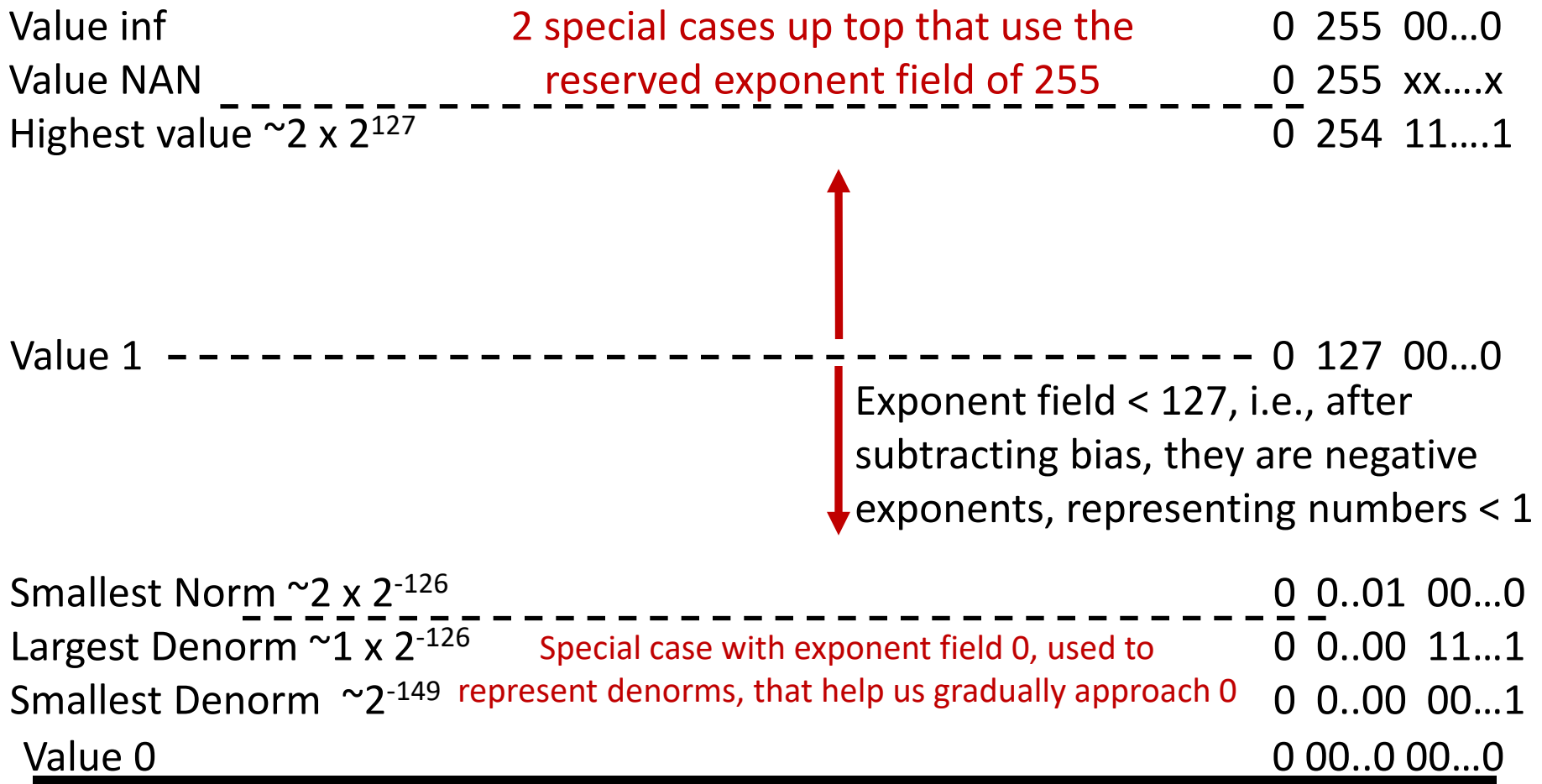
- The number “0” has a special code so that the implicit 1 does not get added: the code is all 0s
(it may seem that this takes up the representation for 1.0, but given how the exponent is represented, that’s not the case)
(see discussion of denorms in the textbook)
- The largest exponent value (with zero fraction) represents +/- infinity
- The largest exponent value (with non-zero fraction) represents NaN (not a number) – for the result of 0/0 or (infinity minus infinity)
- Note that these choices impact the smallest and largest numbers that can be represented

Exponent Representation

- To simplify sort, sign was placed as the first bit
- For a similar reason, the representation of the exponent is also modified: in order to use integer compares, it would be preferable to have the smallest exponent as 00...0 and the largest exponent as 11...1
- This is the biased notation, where a bias is subtracted from the exponent field to yield the true exponent
- IEEE 754 single-precision uses a bias of 127 (since the exponent must have values between -127 and 128)...double precision uses a bias of 1023

Final representation: $(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$





2 special cases up top that use the reserved exponent field of 255

Special case with exponent field 0, used to represent denorms, that help us gradually approach 0

Same rules as above, but the sign bit is 1
Same magnitudes as above, but negative numbers

Examples

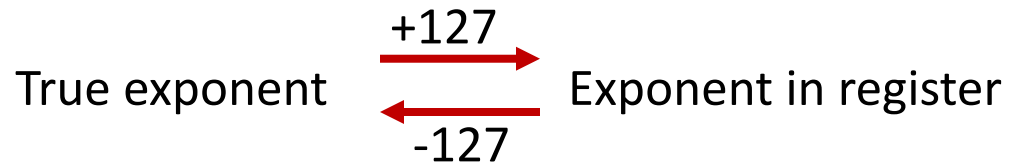
Final representation: $(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$

- Represent -0.75_{ten} in single and double-precision formats

Single: $(1 + 8 + 23)$

Double: $(1 + 11 + 52)$

Remember:



- What decimal number is represented by the following single-precision number?

1 1000 0001 01000...0000

Examples

Final representation: $(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$

- Represent -0.75_{ten} in single and double-precision formats

Single: (1 + 8 + 23)

1 0111 1110 1000...000

Double: (1 + 11 + 52)

1 0111 1111 110 1000...000

- What decimal number is represented by the following single-precision number?

1 1000 0001 01000...0000

-5.0

Example 2

Final representation: $(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$

- Represent 36.90625_{ten} in single-precision format

$$36 / 2 = 18 \text{ rem } 0$$


$$18 / 2 = 9 \text{ rem } 0$$

$$9 / 2 = 4 \text{ rem } 1$$

$$4 / 2 = 2 \text{ rem } 0$$

$$2 / 2 = 1 \text{ rem } 0$$

$$1 / 2 = 0 \text{ rem } 1$$


36 is 100100

$$0.90625 \times 2 = 1.81250$$


$$0.8125 \times 2 = 1.6250$$

$$0.625 \times 2 = 1.250$$

$$0.25 \times 2 = 0.50$$

$$0.5 \times 2 = 1.00$$

$$0.0 \times 2 = 0.0$$


0.90625 is 0.1110100...0

Example 2

Final representation: $(-1)^S \times (1 + \text{Fraction}) \times 2^{(\text{Exponent} - \text{Bias})}$

We've calculated that $36.90625_{\text{ten}} = 100100.1110100\dots0$ in binary

Normalized form = $1.001001110100\dots0 \times 2^5$

(had to shift 5 places to get only one bit left of the point)

The sign bit is 0 (positive number)

The fraction field is 001001110100...0 (the 23 bits after the point)

The exponent field is $5 + 127$ (have to add the bias) = 132,
which in binary is 10000100

The IEEE 754 format is 0 10000100 001001110100.....0
 sign exponent 23 fraction bits

FP Addition

- Consider the following decimal example (can maintain only 4 decimal digits and 2 exponent digits)

$$9.999 \times 10^1 + 1.610 \times 10^{-1}$$

Convert to the larger exponent:

$$9.999 \times 10^1 + 0.016 \times 10^1$$

Add

$$10.015 \times 10^1$$

Normalize

$$1.0015 \times 10^2$$

Check for overflow/underflow

Round

$$1.002 \times 10^2$$

Re-normalize

FP Addition

- Consider the following decimal example (can maintain only 4 decimal digits and 2 exponent digits)

$$9.999 \times 10^1 + 1.610 \times 10^{-1}$$

Convert to the larger exponent:

$$9.999 \times 10^1 + 0.016 \times 10^1$$

Add

$$10.015 \times 10^1$$

Normalize

$$1.0015 \times 10^2$$

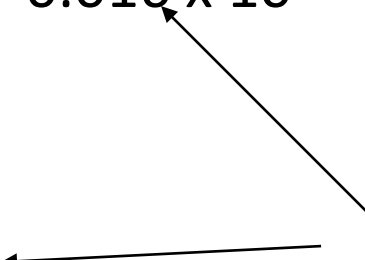
Check for overflow/underflow

Round

$$1.002 \times 10^2$$

Re-normalize

If we had more fraction bits,
these errors would be minimized



FP Addition – Binary Example

- Consider the following binary example

$$1.010 \times 2^1 + 1.100 \times 2^3$$

Convert to the larger exponent:

$$0.0101 \times 2^3 + 1.1000 \times 2^3$$

Add

$$1.1101 \times 2^3$$

Normalize

$$1.1101 \times 2^3$$

Check for overflow/underflow

Round

Re-normalize

IEEE 754 format: 0 10000010 110100000000000000000000

FP Multiplication

- Similar steps:
 - Compute exponent (careful!)
 - Multiply significands (set the binary point correctly)
 - Normalize
 - Round (potentially re-normalize)
 - Assign sign

MIPS Instructions

- The usual add.s, add.d, sub, mul, div
- Comparison instructions: c.eq.s, c.neq.s, c.lt.s....
These comparisons set an internal bit in hardware that is then inspected by branch instructions: bc1t, bc1f
- Separate register file \$f0 - \$f31 : a double-precision value is stored in (say) \$f4-\$f5 and is referred to by \$f4
- Load/store instructions (lwc1, swc1) must still use integer registers for address computation

Code Example

```
float f2c (float fahr)
{
    return ((5.0/9.0) * (fahr - 32.0));
}
```

(argument fahr is stored in \$f12)

```
lwc1 $f16, const5
lwc1 $f18, const9
div.s $f16, $f16, $f18
lwc1 $f18, const32
sub.s $f18, $f12, $f18
mul.s $f0, $f16, $f18
jr    $ra
```

Fixed Point

- FP operations are much slower than integer ops
- Fixed point arithmetic uses integers, but assumes that every number is multiplied by the same factor
- Example: with a factor of $1/1000$, the fixed-point representations for 1.46, 1.7198, and 5624 are respectively 1460, 1720, and 5624000
- More programming effort and possibly lower precision for higher performance

Subword Parallelism

- ALUs are typically designed to perform 64-bit or 128-bit arithmetic
- Some data types are much smaller, e.g., bytes for pixel RGB values, half-words for audio samples
- Partitioning the carry-chains within the ALU can convert the 64-bit adder into 4 16-bit adders or 8 8-bit adders
- A single load can fetch multiple values, and a single add instruction can perform multiple parallel additions, referred to as subword parallelism