

# Floating Point Operations

This presentation explores the fascinating world of floating point operations, from their fundamental concepts to advanced techniques for optimizing their performance.

 by karthik

# Introduction to Floating Point Numbers

## Representing Real Numbers

Floating point numbers are used to represent real numbers, encompassing both integers and fractional values. They offer a versatile way to express numbers with a wide range of magnitudes.

## Scientific Notation

Similar to scientific notation, floating point numbers consist of a significand (mantissa), an exponent, and a sign. This format allows for flexible representation of very large and very small numbers.

## Applications

Floating point operations are essential in numerous scientific, engineering, and computational applications. They are used in fields such as physics, finance, and computer graphics.

# IEEE 754 Standard

## Universal Standard

The IEEE 754 standard is a globally recognized standard that defines the representation and handling of floating point numbers. It ensures compatibility and consistency across different computer systems.

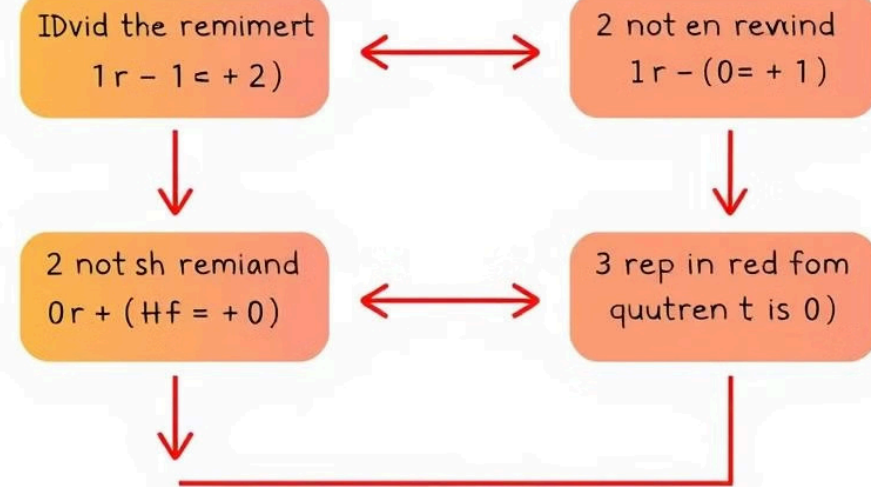
## Precision and Formats

The standard specifies various formats for floating point numbers, including single-precision (32 bits) and double-precision (64 bits), each offering a different level of precision.

## Special Values

The IEEE 754 standard includes special values like infinity (positive and negative), NaN (Not a Number), and zero (positive and negative). These values handle exceptional situations and edge cases.

Coal wibe the **num**er te  
reccimd it ncinreledn.



Read resd tshes : In in red l top:  
(1 - 11, 0) 12 = 4 139 122.4 11)  
↓  
2, + 4, 3, 1 = 3.2, 1 132 - 111)



# Representation of Floating Point Numbers



## Sign

The sign bit indicates whether the number is positive or negative. 0 represents positive, and 1 represents negative.



## Exponent

The exponent determines the magnitude of the number. It is encoded in a biased form to represent both positive and negative exponents.



## Mantissa

The mantissa stores the significant digits of the number, providing the precision of the representation.

# Floating Point Arithmetic

1

Addition and Subtraction: Floating point numbers are aligned to the same exponent before performing addition or subtraction.

2

Multiplication: Multiplying the mantissas and adding the exponents is the fundamental operation.

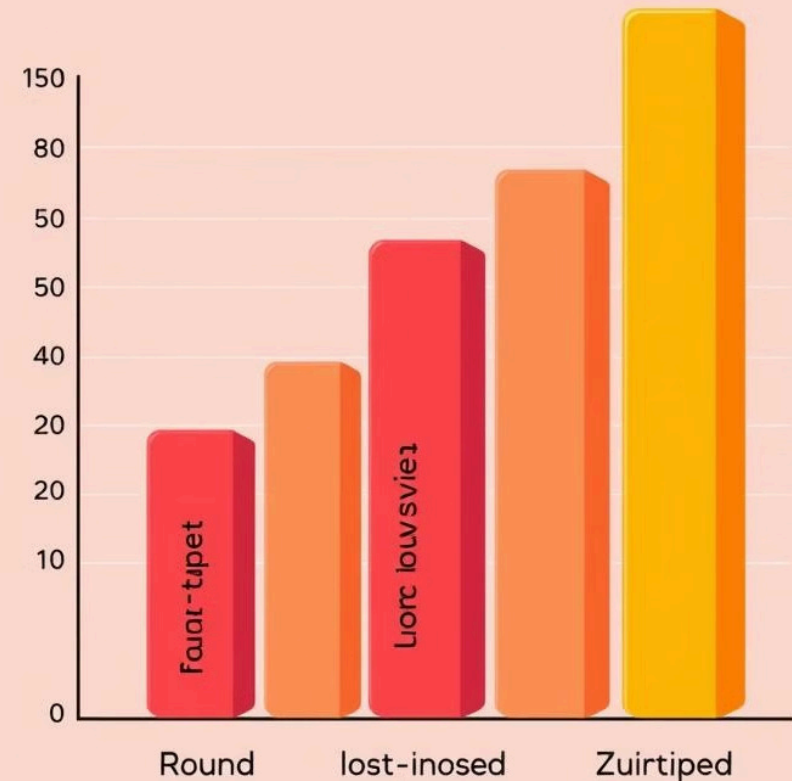
3

Division: Dividing the mantissas and subtracting the exponents is the essential operation.





# Rounding and Errors in Floating Point Computations



Rounding rounding computations used of glong them round on sures  
ronkads and load is foy thes to rof lision if floatin preccisim.

1

## Rounding

Floating point arithmetic often involves rounding to fit the limited precision of the representation.

2

## Rounding Modes

Different rounding modes, such as round-to-nearest, round-up, and round-down, exist to handle rounding decisions.

3

## Error Accumulation

Repeated operations can accumulate rounding errors, potentially leading to significant deviations from the expected result.

# Floating Point Performance Considerations

**1**

## **Computational Cost**

Floating point operations are computationally expensive, especially when dealing with complex calculations.

**2**

## **Memory Access**

Accessing floating point numbers from memory can contribute to performance overhead.

**3**

## **Processor Architecture**

The design of the processor architecture influences the efficiency of floating point operations.

# Techniques for Improving Floating Point Performance

**1**

## **Loop Unrolling**

Unrolling loops can reduce the overhead of loop control, increasing execution speed.

**2**

## **Vectorization**

Vectorization exploits the ability of modern processors to perform operations on multiple data elements simultaneously.

**3**

## **Cache Optimization**

Optimizing memory access patterns to leverage the cache hierarchy can significantly improve performance.





Horeg rafi empager doms in  
glat end for thet of essons de  
theuding assioptvies.

# Conclusion and Takeaways

## 1

### Understanding

A solid understanding of floating point operations is essential for developing accurate and efficient software.

## 2

### Precision

Be mindful of the limitations of floating point precision and potential rounding errors.

## 3

### Optimization

Employ techniques to improve floating point performance, maximizing computational efficiency.

**THANK YOU**

