

# CS 450: Numerical Analysis<sup>1</sup>

## Introduction to Scientific Computing

University of Illinois at Urbana-Champaign

---

<sup>1</sup>*These slides have been drafted by Edgar Solomonik as lecture templates and supplementary material for the book “Scientific Computing: An Introductory Survey” by Michael T. Heath ([slides](#)).*

# Scientific Computing Applications and Context

- ▶ **Mathematical modelling for computational science** *Typical scientific computing problems are numerical solutions to PDEs*
  - ▶ *Newtonian dynamics: simulating particle systems in time*
  - ▶ *Fluid and air flow models for engineering*
  - ▶ *PDE-constrained numerical optimization: finding optimal configurations (used in engineering of control systems)*
  - ▶ *Quantum chemistry (electronic structure calculations): many-electron Schrödinger equation*
- ▶ **Linear algebra and computation**
  - ▶ *Linear algebra and numerical optimization are building blocks for machine learning methods and data analysis*
  - ▶ *Computer architecture, compilers, and parallel computing use numerical algorithms (matrix multiplication, Gaussian elimination) as benchmarks*

## Example: Mechanics<sup>2</sup>

- ▶ Newton's laws provide incomplete particle-centric picture
- ▶ Physical systems can be described in terms of *degrees of freedom* (DoFs)
  - ▶ A piston moving up and down requires \_\_\_\_\_ DoFs
  - ▶ 1-particle system requires \_\_\_\_\_ DoFs
  - ▶ 2-particle system requires \_\_\_\_\_ DoFs
  - ▶ 2-particles at a fixed distance require \_\_\_\_\_ DoFs
- ▶  $N$ -particle system *configuration* described by  $3N$  DoFs

---

<sup>2</sup>*Variational Principles of Mechanics*, Cornelius Lanczos, Dover Books on Physics, 1949.



# Numerical Analysis

- ▶ **Numerical Problems involving Continuous Phenomena:**

- ▶ **Error Analysis:**

## Sources of Error

- ▶ **Representation of Numbers:**

- ▶ **Propagated Data Error:**

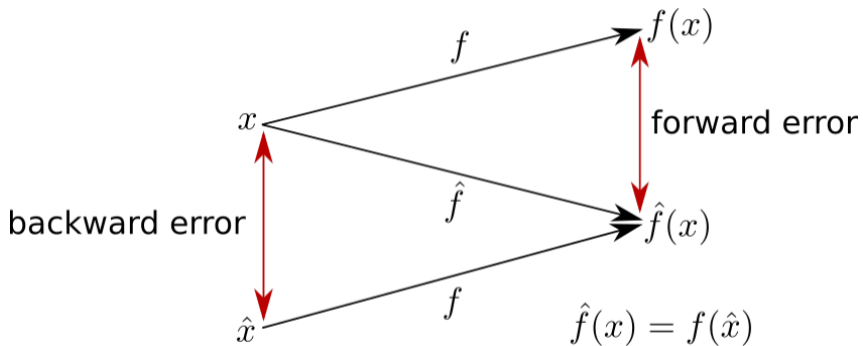
- ▶ **Computational Error =  $\hat{f}(x) - f(x)$  = Truncation Error + Rounding Error**

# Error Analysis

- ▶ **Forward Error:**

- ▶ **Backward Error:**

## Visualization of Forward and Backward Error





## Conditioning

▶ **Absolute Condition Number:**

▶ **(Relative) Condition Number:**

## Posedness and Conditioning

- ▶ **What is the condition number of an ill-posed problem?**

# Stability and Accuracy

▶ **Accuracy:**

▶ **Stability:**

## Error and Conditioning

- ▶ Two major sources of error: *roundoff* and *truncation* error.
  - ▶ roundoff error concerns floating point error due to finite precision
  - ▶ truncation error concerns error incurred due to algorithmic approximation, e.g. the representation of a function by a finite Taylor series
  
- ▶ To study the propagation of roundoff error in arithmetic we can use the notion of conditioning.

# Floating Point Numbers

*Demo: Picking apart a floating point number*

*Demo: Density of Floating Point Numbers*

- ▶ **Scientific Notation**

- ▶ **Significand (Mantissa) and Exponent** Given  $x$  with  $s$  leading bits  $x_0, \dots, x_{s-1}$

# Rounding Error

*Demo: Floating point and the Harmonic Series*

*Demo: Floating Point and the Series for the Exponential Function*

## ► Maximum Relative Representation Error (Machine Epsilon)

$$1.\underbrace{01101}_{\text{significantand}} \times 2^{-3}$$

$$\epsilon = \operatorname{argmin}_{\epsilon > 0} f(1 + \epsilon) = 1 + \epsilon$$

$$\begin{array}{r} 1.0000 \times 2^0 \\ + 0.0001 \times 2^{-1} \\ \hline 1.0001 \end{array}$$

# Rounding Error in Operations (I) *Activity: Cancellation in Standard Deviation Computation*

► Addition and Subtraction

Cancellation - loss of sig. digits

$$\begin{array}{r} 1.247 \\ - 1.233 \\ \hline \end{array}$$

$$\begin{array}{l} (1 + \epsilon) \\ \uparrow \\ (1 + \epsilon) \end{array}$$

$$= \underline{.014} \Rightarrow \underline{1.4} \times 10^{-2} + \epsilon \cdot 1.247 + \epsilon \cdot 1.233$$

$$\kappa_{\text{add}}(x, y) = \lim_{\delta x \rightarrow 0} \left| \frac{(x + \delta x) + y - (x + y)}{x + y} \right|$$

$\delta x \quad | \quad 1 \times 1$

$$= \lim_{\delta x \rightarrow 0} \frac{\left| \frac{\delta x}{|x+y|} \right|}{|\delta x|/|x|} = \frac{|x|}{|x+y|}$$

$$f(x) = x + y$$

$$y \approx -x$$

→ ∞

$$K_f(x) = \frac{|f'(x)| \cdot |x|}{|f(x)|} = \frac{|x|}{|x+y|}$$

$$K_{abs}(x) = |f'(x)|$$

(relative)  $K_f(x) = K_{abs}(x) \cdot \frac{|f(x)|}{|x|}$



## Rounding Error in Operations (II)

- ▶ **Multiplication and Division**

# Exceptional and Subnormal Numbers

## ▶ Exceptional Numbers

not normalized numbers  $(1.\underbrace{0110\dots}_{\text{significant}} \times 2^E)$

$0, -0, \underline{0}, -0, \underline{NaN} = 0/0 \rightarrow \underline{0} - \underline{0}$

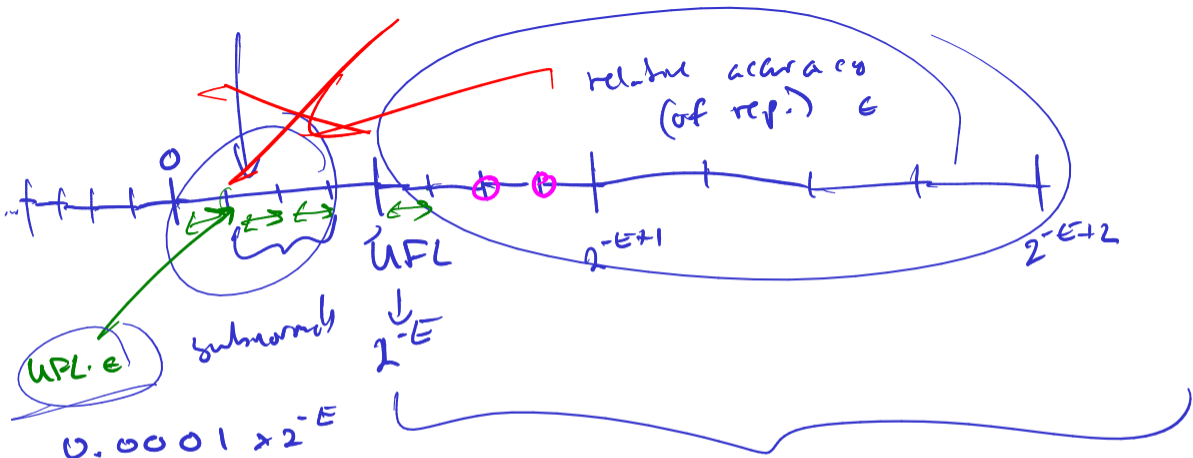
$\{-1023, 1024\}$   
E

## ▶ Subnormal (Denormal) Number Range

$0.\underbrace{\hspace{2cm}}_{\text{significant}} \times 10^{-1023}$

UFL  $1.0000\dots 0 \times 10^{-1023}$   
3 smallest possible exponent

## ▶ Gradual Underflow: Avoiding underflow in addition



$$f(x) - f(y) \neq 0$$

$$\Leftrightarrow f(x) \neq f(y)$$

because of subnormal

normalized numbers

# Floating Point Number Line

