

# On the Amplification of Rounding Errors

Erin C. Carson

Katedra numerické matematiky, Matematicko-fyzikální fakulta, Univerzita Karlova

*Advances in Numerical Linear Algebra: Celebrating the Centenary of the Birth of James H. Wilkinson*

Manchester, UK

May 29-30, 2019

This research was partially supported by OP RDE project No. CZ.02.2.69/0.0/0.0/16\_027/0008495



FACULTY  
OF MATHEMATICS  
AND PHYSICS  
Charles University



EUROPEAN UNION  
European Structural and Investment Funds  
Operational Programme Research,  
Development and Education

  
MINISTRY OF EDUCATION,  
YOUTH AND SPORTS

# Motivation

*People are awed at the prodigious speeds at which they execute primitive arithmetic operations such as addition and multiplication. Yet this speed is achieved at a price, almost every answer is wrong!*

- B. N. Parlett, *James Hardy ("Jim") Wilkinson*, ACM Turing Award site

# Motivation

*People are awed at the prodigious speeds at which they execute primitive arithmetic operations such as addition and multiplication. Yet this speed is achieved at a price, almost every answer is wrong!*

- B. N. Parlett, *James Hardy ("Jim") Wilkinson*, ACM Turing Award site

- Goal: efficient, sufficiently accurate computations in spite of rounding errors

# Motivation

*People are awed at the prodigious speeds at which they execute primitive arithmetic operations such as addition and multiplication. Yet this speed is achieved at a price, almost every answer is wrong!*

- B. N. Parlett, *James Hardy ("Jim") Wilkinson*, ACM Turing Award site

- Goal: efficient, sufficiently accurate computations in spite of rounding errors
- Accumulation versus amplification: the role of the algorithm
  - Accumulation of rounding errors: inevitable part of computation in finite precision arithmetic
  - Amplification of rounding errors: **property of the mathematical structure of the algorithm we use to transform the data**

# Example: Conjugate Gradient Algorithms

- Two algorithms for the CG method: HSCG and STCG
  - Equivalent in exact arithmetic
  - Don't look terribly different; can finite precision behavior be significantly different?

## HSCG

([Hestenes & Stiefel, 1952])

```
 $r_0 = b - Ax_0, \quad p_0 = r_0$   
for  $i = 1:nmax$   
     $\alpha_{i-1} = \frac{r_{i-1}^T r_{i-1}}{p_{i-1}^T A p_{i-1}}$   
     $x_i = x_{i-1} + \alpha_{i-1} p_{i-1}$   
     $r_i = r_{i-1} - \alpha_{i-1} A p_{i-1}$   
     $\beta_i = \frac{r_i^T r_i}{r_{i-1}^T r_{i-1}}$   
     $p_i = r_i + \beta_i p_{i-1}$   
end
```

## STCG

([Stiefel, 1952/53], [Rutishauser, 1959], [Hageman & Young, 1981])

```
 $r_0 = b - Ax_0, \quad p_0 = r_0, \quad x_{-1} = x_0,$   
 $r_{-1} = r_0, \quad e_{-1} = 0$   
for  $i = 1:nmax$   
     $q_{i-1} = \frac{(r_{i-1}, Ar_{i-1})}{(r_{i-1}, r_{i-1})} - e_{i-2}$   
     $x_i = x_{i-1} + \frac{1}{q_{i-1}} (r_{i-1} + e_{i-2} (x_{i-1} - x_{i-2}))$   
     $r_i = r_{i-1} + \frac{1}{q_{i-1}} (-Ar_{i-1} + e_{i-2} (r_{i-1} - r_{i-2}))$   
     $e_{i-1} = q_{i-1} \frac{(r_i, r_i)}{(r_{i-1}, r_{i-1})}$   
end
```

# Maximum Attainable Accuracy

- Attainable accuracy typically bounded in terms of the size of the *residual gap* (between true residual  $b - Ax_i$  and recursively updated residual  $r_i$ )

# Maximum Attainable Accuracy

- Attainable accuracy typically bounded in terms of the size of the *residual gap* (between true residual  $b - Ax_i$  and recursively updated residual  $r_i$ )
- For HSCG
  - Bound on residual gap can be written as accumulation of local errors [Greenbaum, 1997]

$$f_i \equiv b - A\hat{x}_i - \hat{r}_i = f_0 + \sum_{m=1}^i (A\delta x_m + \delta r_m)$$

# Maximum Attainable Accuracy

- Attainable accuracy typically bounded in terms of the size of the *residual gap* (between true residual  $b - Ax_i$  and recursively updated residual  $r_i$ )
- For HSCG
  - Bound on residual gap can be written as accumulation of local errors [Greenbaum, 1997]

$$f_i \equiv b - A\hat{x}_i - \hat{r}_i = f_0 + \sum_{m=1}^i (A\delta x_m + \delta r_m)$$

- For STCG
  - Attainable accuracy for STCG can be much worse than for HSCG [Gutknecht & Strakoš, 2000]

# Maximum Attainable Accuracy

- Attainable accuracy typically bounded in terms of the size of the *residual gap* (between true residual  $b - Ax_i$  and recursively updated residual  $r_i$ )
- For HSCG
  - Bound on residual gap can be written as accumulation of local errors [Greenbaum, 1997]

$$f_i \equiv b - A\hat{x}_i - \hat{r}_i = f_0 + \sum_{m=1}^i (A\delta x_m + \delta r_m)$$

- For STCG
  - Attainable accuracy for STCG can be much worse than for HSCG [Gutknecht & Strakoš, 2000]
  - Residual gap bounded by sum of local errors PLUS local errors multiplied by factors which depend on

$$\max_{0 \leq \ell < j \leq i} \frac{\|r_j\|^2}{\|r_\ell\|^2}$$

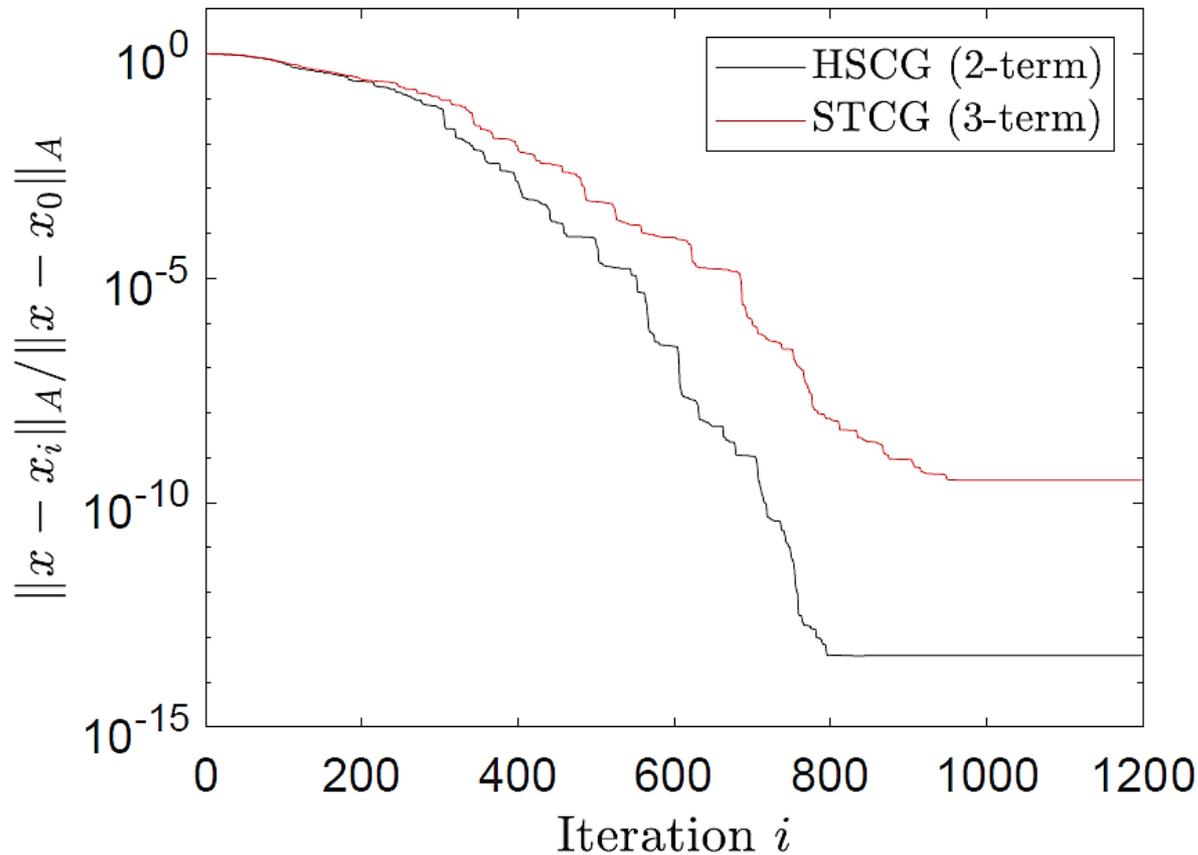
- ⇒ Large residual oscillations can cause these factors to be large!
- ⇒ Local errors can be amplified!

# Numerical Example

$A$ : bcsstk03 from SuiteSparse,

$b$ : equal components in the eigenbasis of  $A$  and  $\|b\| = 1$

$N = 112, \kappa(A) \approx 7e6$



# Algorithms Designed for HPC

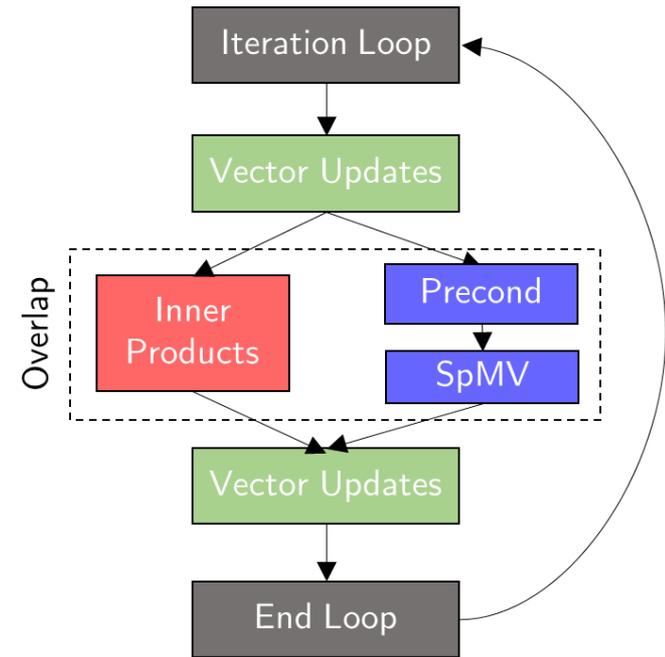
- Many other variants of CG motivated by solving large-scale problems on large-scale machines
- Example: pipelined CG [Ghysels & Vanroose, 2014]

# Algorithms Designed for HPC

- Many other variants of CG motivated by solving large-scale problems on large-scale machines
- Example: pipelined CG [Ghysels & Vanroose, 2014]
  - Main idea: add auxiliary vectors

$$s_i \equiv Ap_i, \quad w_i \equiv Ar_i, \quad z_i \equiv Aw_i$$

so that matrix-vector product and inner product computations are decoupled and can be overlapped



- How does adding auxiliary vectors effect the numerical behavior?

# Algorithms Designed for HPC

- Many other variants of CG motivated by solving large-scale problems on large-scale machines
- Example: pipelined CG [Ghysels & Vanroose, 2014]
  - Main idea: add auxiliary vectors

$$s_i \equiv Ap_i, \quad w_i \equiv Ar_i, \quad z_i \equiv Aw_i$$

so that matrix-vector product and inner product computations are decoupled and can be overlapped

$$r_0 = b - Ax_0, p_0 = r_0, s_0 = Ap_0$$

for  $i = 1:nmax$

$$\alpha_{i-1} = \frac{(r_{i-1}, r_{i-1})}{(p_{i-1}, s_{i-1})}$$

$$x_i = x_{i-1} + \alpha_{i-1} p_{i-1}$$

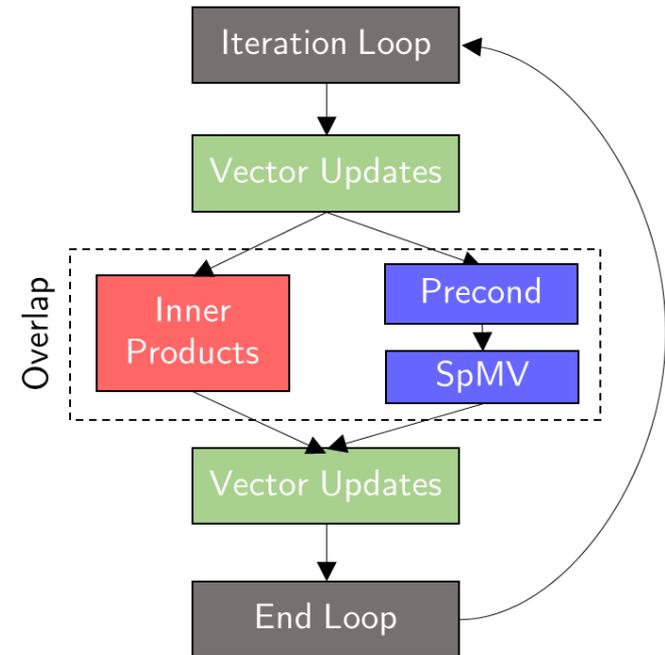
$$r_i = r_{i-1} - \alpha_{i-1} s_{i-1}$$

$$\beta_i = \frac{(r_i, r_i)}{(r_{i-1}, r_{i-1})}$$

$$p_i = r_i + \beta_i p_{i-1}$$

$$s_i = Ar_i + \beta_i s_{i-1}$$

end



- How does adding auxiliary vectors effect the numerical behavior?
- Consider simplified version, where we just add one auxiliary vector  $s_i \equiv Ap_i$  to HSCG

# Maximum Attainable Accuracy

For this simplified pipelined CG algorithm:

[C., Rozložník, Strakoš, Tichý, & Tůma, 2018]  
see also [Cools et al., 2018]

$$f_i \equiv b - A\hat{x}_i - \hat{r}_i = f_0 - \sum_{j=0}^i \hat{\alpha}_j g_j - \sum_{j=0}^i (A\delta_j^x + \delta_j^r)$$

# Maximum Attainable Accuracy

For this simplified pipelined CG algorithm:

[C., Rozložník, Strakoš, Tichý, & Tůma, 2018]  
see also [Cools et al., 2018]

$$f_i \equiv b - A\hat{x}_i - \hat{r}_i = f_0 - \sum_{j=0}^i \hat{\alpha}_j g_j - \sum_{j=0}^i (A\delta_j^x + \delta_j^r)$$

$$g_j = \left( \prod_{k=1}^j \hat{\beta}_k \right) g_0 + \sum_{k=1}^j \left( \prod_{\ell=k+1}^j \hat{\beta}_\ell \right) (A\delta_k^p - \delta_k^s)$$

# Maximum Attainable Accuracy

For this simplified pipelined CG algorithm:

[C., Rozložník, Strakoš, Tichý, & Tůma, 2018]  
see also [Cools et al., 2018]

$$f_i \equiv b - A\hat{x}_i - \hat{r}_i = f_0 - \sum_{j=0}^i \hat{\alpha}_j g_j - \sum_{j=0}^i (A\delta_j^x + \delta_j^r)$$

$$g_j = \left( \prod_{k=1}^j \hat{\beta}_k \right) g_0 + \sum_{k=1}^j \left( \prod_{\ell=k+1}^j \hat{\beta}_\ell \right) (A\delta_k^p - \delta_k^s)$$

$$\beta_\ell \beta_{\ell+1} \cdots \beta_j = \frac{\|r_j\|^2}{\|r_{\ell-1}\|^2}, \quad \ell < j$$

# Maximum Attainable Accuracy

For this simplified pipelined CG algorithm:

[C., Rozložník, Strakoš, Tichý, & Tůma, 2018]  
see also [Cools et al., 2018]

$$f_i \equiv b - A\hat{x}_i - \hat{r}_i = f_0 - \sum_{j=0}^i \hat{\alpha}_j g_j - \sum_{j=0}^i (A\delta_j^x + \delta_j^r)$$

$$g_j = \left( \prod_{k=1}^j \hat{\beta}_k \right) g_0 + \sum_{k=1}^j \left( \prod_{\ell=k+1}^j \hat{\beta}_\ell \right) (A\delta_k^p - \delta_k^s)$$

$$\beta_\ell \beta_{\ell+1} \cdots \beta_j = \frac{\|r_j\|^2}{\|r_{\ell-1}\|^2}, \quad \ell < j$$

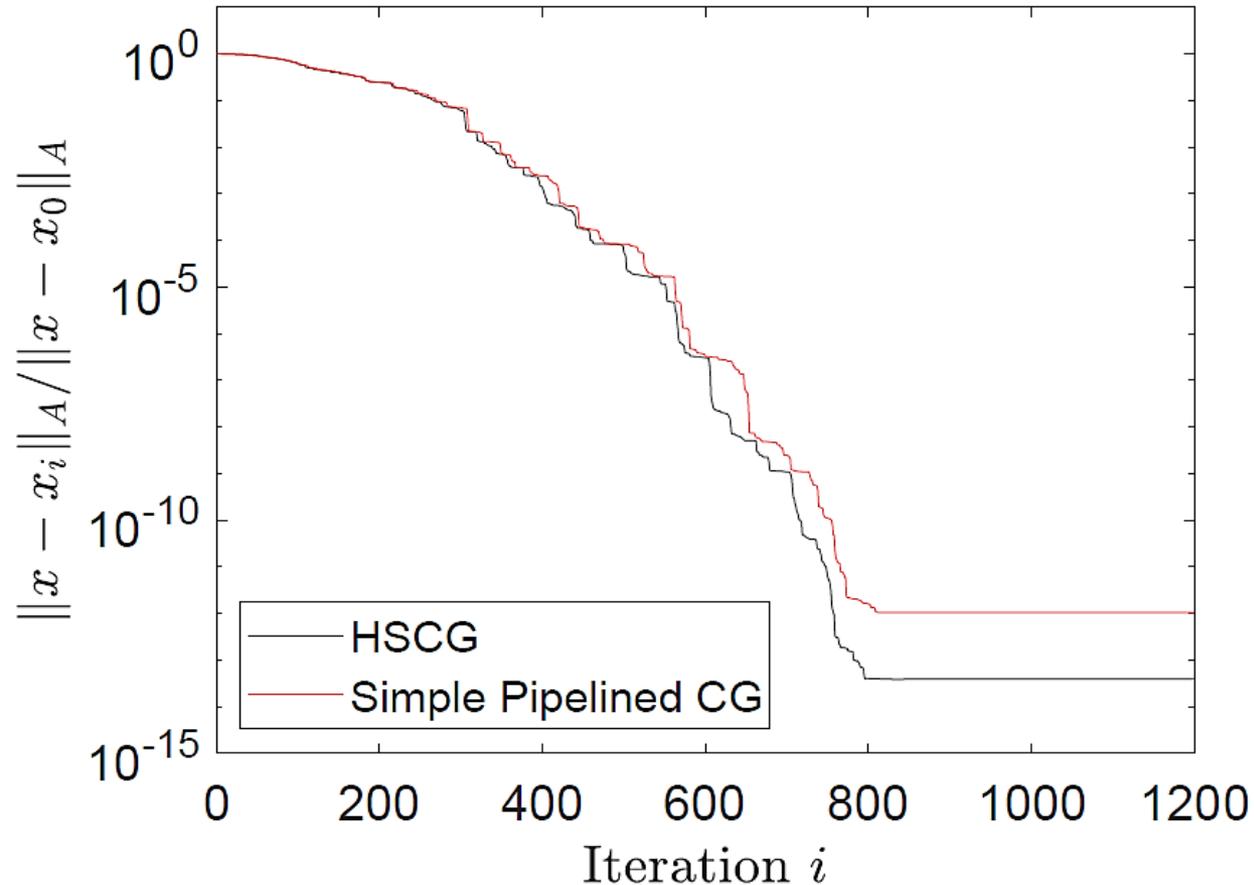
- Residual oscillations can cause these factors to be large!
- Very similar to the results for attainable accuracy in the 3-term STCG
- *Seemingly innocuous* change can cause **amplification of local rounding errors**

# Numerical Example

$A$ : bcsstk03 from SuiteSparse,

$b$ : equal components in the eigenbasis of  $A$  and  $\|b\| = 1$

$N = 112, \kappa(A) \approx 7e6$

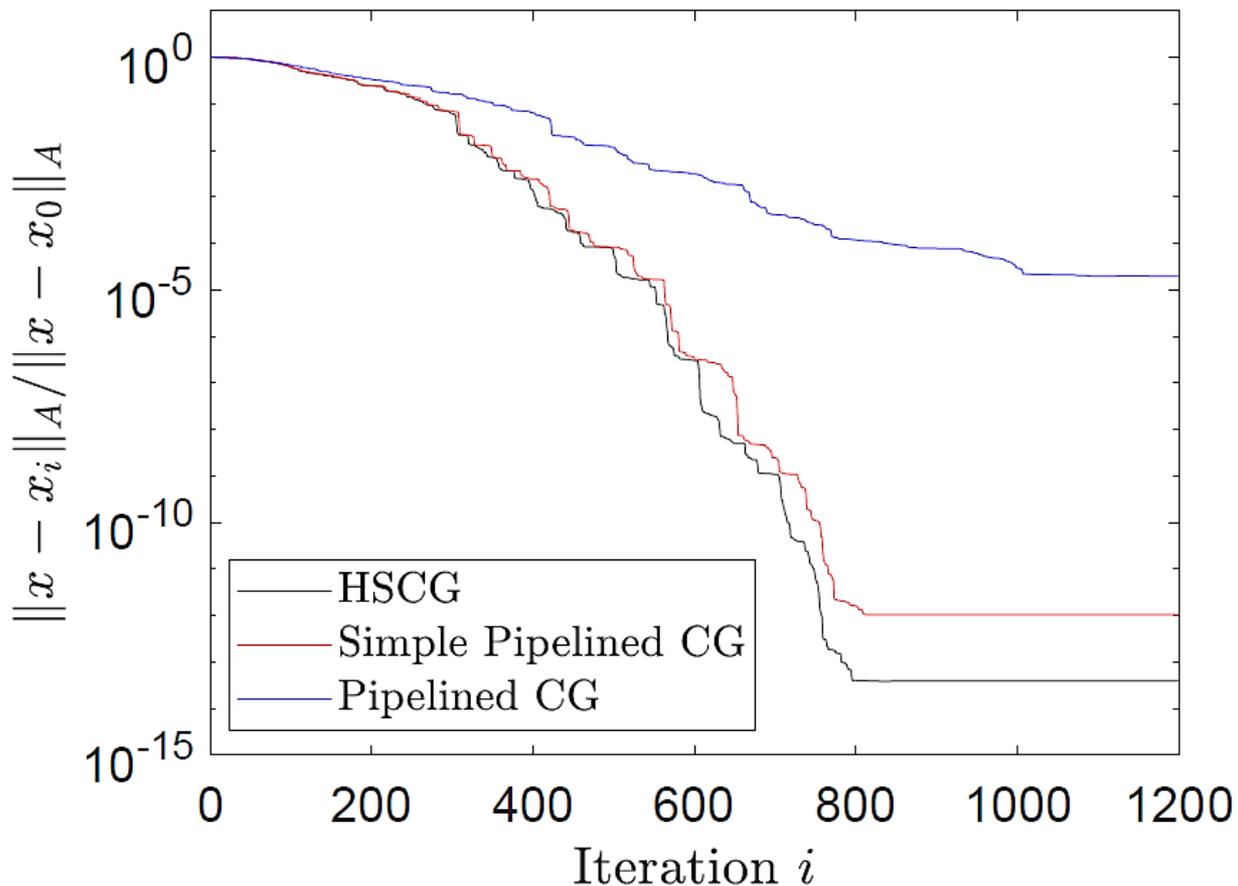


# Numerical Example

$A$ : bcsstk03 from SuiteSparse,

$b$ : equal components in the eigenbasis of  $A$  and  $\|b\| = 1$

$N = 112, \kappa(A) \approx 7e6$



# Insights from Error Analysis

- Takeaway: even a small modification to HSCG recurrences (addition of one auxiliary vector) can cause rounding errors to be amplified
  - Amplification factors depend on size of residual oscillations

# Insights from Error Analysis

- Takeaway: even a small modification to HSCG recurrences (addition of one auxiliary vector) can cause rounding errors to be amplified
  - Amplification factors depend on size of residual oscillations
- Note: bounds may be far from tight; the important thing is the insight we can obtain from the bounds

*There is still a tendency to attach too much importance to the precise error bounds obtained by an a priori error analysis. In my opinion, **the bound itself is usually the least important part of it.** The main object of such an analysis is to expose the potential instabilities, if any, of an algorithm so that, hopefully, from the insight thus obtained one might be led to improved algorithms.*

- J. H. Wilkinson, SIAM Rev. 14 (1971)

# Takeaways

- In designing new algorithms, even slight modifications of the way in which quantities are computed can cause significant changes to numerical behavior in finite precision

# Takeaways

- In designing new algorithms, even slight modifications of the way in which quantities are computed can cause significant changes to numerical behavior in finite precision
- It is critical to consider this in designing algorithms, especially in the context of HPC
  - Even if algorithms are mathematically (in infinite precision) equivalent to the classical approach, **effects of finite precision can negate any potential performance benefit**
  - Note: we only discussed maximum attainable accuracy, but convergence is also delayed due to finite precision computations
    - In all presented CG algorithms, even HSCG, *amplification* of rounding errors contributes to convergence delay

# Takeaways

- In designing new algorithms, even slight modifications of the way in which quantities are computed can cause significant changes to numerical behavior in finite precision
- It is critical to consider this in designing algorithms, especially in the context of HPC
  - Even if algorithms are mathematically (in infinite precision) equivalent to the classical approach, **effects of finite precision can negate any potential performance benefit**
  - Note: we only discussed maximum attainable accuracy, but convergence is also delayed due to finite precision computations
    - In all presented CG algorithms, even HSCG, *amplification* of rounding errors contributes to convergence delay

*It is easy to be carried away by the excitement of producing an alternative method for which convergence can be rigorously demonstrated, and to overlook the fact that this method too will suffer from the incidence of rounding errors. Attractive mathematics does not protect one from the rigors of digital computation.*

- J. H. Wilkinson, SIAM Rev. 14 (1971)

# Looking Forward

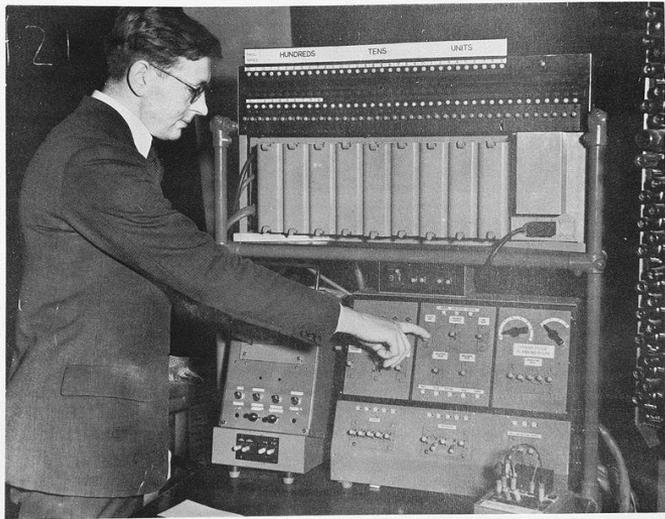
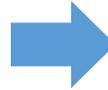


Figure 9. J. H. Wilkinson operating the console of the Pilot ACE during a press demonstration held in November 1950 (which mercifully coincided with a rare period of reliable operation). The legend HUNDREDS-TENS-UNITS was provided for the press demonstration and was not a permanent fixture. The face of the CRT monitor is obscured by Wilkinson's forearm.



[www.maths.manchester.ac.uk/~higham/photos/wilkinson/jhw\\_pilot%20ace1.htm](http://www.maths.manchester.ac.uk/~higham/photos/wilkinson/jhw_pilot%20ace1.htm)

[www.olcf.ornl.gov/summit/](http://www.olcf.ornl.gov/summit/)

# Looking Forward

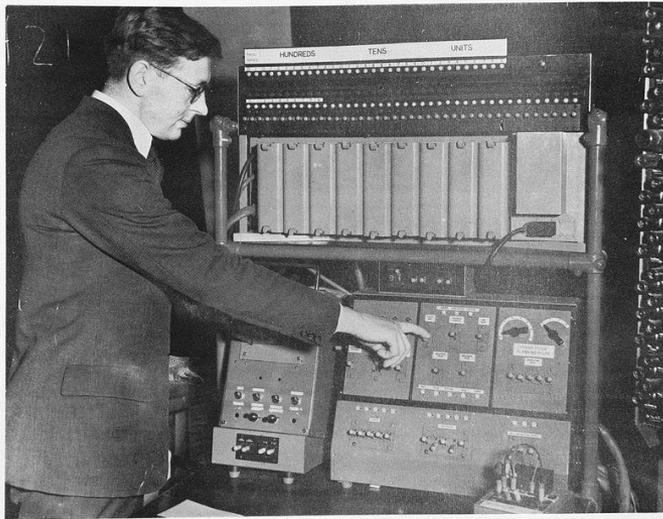
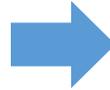


Figure 9. J. H. Wilkinson operating the console of the Pilot ACE during a press demonstration held in November 1950 (which mercifully coincided with a rare period of reliable operation). The legend HUNDREDS-TENS-UNITS was provided for the press demonstration and was not a permanent fixture. The face of the CRT monitor is obscured by Wilkinson's forearm.



[www.maths.manchester.ac.uk/~higham/photos/wilkinson/jhw\\_pilot%20ace1.htm](http://www.maths.manchester.ac.uk/~higham/photos/wilkinson/jhw_pilot%20ace1.htm)

[www.olcf.ornl.gov/summit/](http://www.olcf.ornl.gov/summit/)

- With trend of multi-precision and low-precision computation, paying attention to amplification of rounding errors becomes especially important;
  - Amplification factors that were small relative to double precision can now have a much greater affect

$$1 \cdot \varepsilon_h \approx 10^{12} \cdot \varepsilon_d$$

# Looking Forward

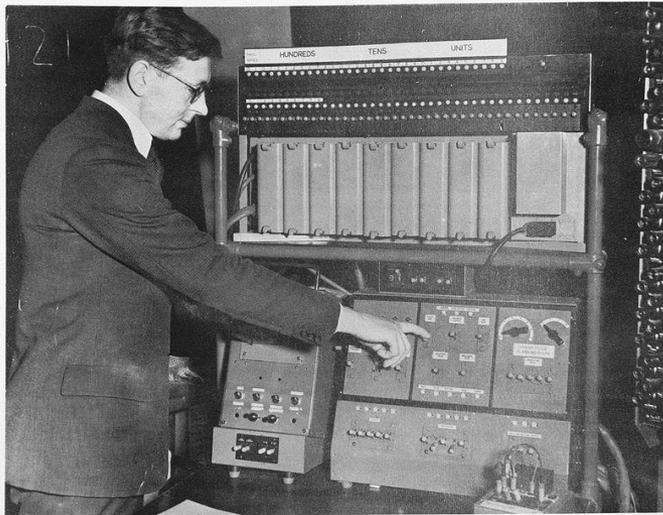
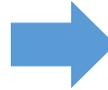


Figure 9. J. H. Wilkinson operating the console of the Pilot ACE during a press demonstration held in November 1950 (which mercifully coincided with a rare period of reliable operation). The legend HUNDREDS-TENS-UNITS was provided for the press demonstration and was not a permanent fixture. The face of the CRT monitor is obscured by Wilkinson's forearm.



[www.maths.manchester.ac.uk/~higham/photos/wilkinson/jhw\\_pilot%20ace1.htm](http://www.maths.manchester.ac.uk/~higham/photos/wilkinson/jhw_pilot%20ace1.htm)

[www.olcf.ornl.gov/summit/](http://www.olcf.ornl.gov/summit/)

- With trend of multi-precision and low-precision computation, paying attention to amplification of rounding errors becomes especially important;
  - Amplification factors that were small relative to double precision can now have a much greater affect

$$1 \cdot \varepsilon_h \approx 10^{12} \cdot \varepsilon_d$$

- Challenges: new number formats (IEEE 754 and beyond); efficient algorithms/implementations on multiprecision hardware; analysis of multiprecision algorithms; refined notions of ill-conditioning and techniques used in error analysis

# Following in Wilkinson's Footsteps

- Wilkinson's resume includes experience with applications, hardware design and construction of computers, algorithm implementation, development of backward error analysis
  - "bird's eye view" of numerical computation from the hardware to the algorithms to the application

# Following in Wilkinson's Footsteps

- Wilkinson's resume includes experience with applications, hardware design and construction of computers, algorithm implementation, development of backward error analysis
  - "bird's eye view" of numerical computation from the hardware to the algorithms to the application
- Progress in numerical mathematics and high-performance computing must be tightly **interdisciplinary** and involve close collaboration between computer engineers, software engineers, computer scientists, applied mathematicians, computational science experts, ...

# Thank you!

[carson@karlin.mff.cuni.cz](mailto:carson@karlin.mff.cuni.cz)

[www.karlin.mff.cuni.cz/~carson/](http://www.karlin.mff.cuni.cz/~carson/)

# References

- E. C. Carson, M. Rozložník, Z. Strakoš, P. Tichý, and M. Tůma. "The Numerical Stability Analysis of Pipelined Conjugate Gradient Methods: Historical Context and Methodology." *SIAM J. Sci. Comput.* 40.5 (2018): A3549-A3580.
- S. Cools, E. F. Yetkin, E. Agullo, L. Girard, and W. Vanroose, "Analyzing the effect of local rounding error propagation on the maximal attainable accuracy of the pipelined conjugate gradient method." *SIAM J. Matrix Anal. Appl.* 39.1 (2018): 426-450.
- M. R. Hestenes and E. Stiefel. *Methods of conjugate gradients for solving linear systems*. J. Research Nat. Bur. Standards, 49:409–436, 1952.
- P. Ghysels, and W. Vanroose. "Hiding global synchronization latency in the preconditioned conjugate gradient algorithm." *Parallel Comput.* 40.7 (2014): 224-238.
- A. Greenbaum, "Estimating the attainable accuracy of recursively computed residual methods." *SIAM J. Matrix Anal. Appl.* 18.3 (1997): 535-551.
- M. H. Gutknecht and Z. Strakoš. "Accuracy of two three-term and three two-term recurrences for Krylov space solvers." *SIAM J. Matrix Anal. Appl.* 22.1 (2000): 213-229
- J. H. Wilkinson, "Modern error analysis." *SIAM Review* 13.4 (1971): 548-568.