This paper describes objective technical results and analysis. Any subjective views or opinions that might be expressed in the paper do not necessarily represent the views of the U.S. Department of Energy or the United States Government.

SAND2022-7865C

# Mixed Precision *s*-step Conjugate Gradient with Residual Replacements on (NVIDIA) GPUs

Approved for public release

Ichitaro Yamazaki*, Erin Carson T, Brian Kelley*

*Sandia National Laboratories, New Mexico, USA
TCharles University, Prague, Czech Republic

U.S. DEPARTMENT OF ENERGY | Office of Science

# Aims with "Mixed Precision *s*-step Conjugate Gradient with Residual Replacements on (NVIDIA) GPUs"

- improve the stability of the solver, using higher-precision arithmetic
  with two main aims:

  1. obtain higher accuracy (converging to lower residual norm)

  2. improve performance (converging with a fewer iterations, but without significant increase
     in the per-iteration time = faster time-to-solution)
     - Careful design and implementation
       - using higher precision only at the critical parts of the algorithms
       - optimizing the underlying kernels for particular properties

# Conjugate Gradient (CG)

**Require:** SPD matrix $A$, RHS vector $\mathbf{b}$, and initial approximate solution vector $\mathbf{x}_1$

1: $\mathbf{r}_1 := \mathbf{b} - A\mathbf{x}_1$, $\delta_1 := \mathbf{r}_1^T\mathbf{r}_1$, $\mathbf{p}_1 := \mathbf{r}_1$
2: **for** $j = 1, 2, \ldots$ **do**
3:     // SpMV with P2P communication
4:     $\mathbf{w}_j := A\mathbf{p}_j$
5:     // dot-product with global-reduce
6:     $\gamma_j := \mathbf{p}_j^T\mathbf{w}_j$
7:     $\alpha_j := \delta_j/\gamma_j$
8:     // update solution and residual vectors
9:     $\mathbf{x}_{j+1} := \mathbf{x}_j + \alpha_j\mathbf{p}_j$
10:     $\mathbf{r}_{j+1} := \mathbf{r}_j - \alpha_j\mathbf{w}_j$
11:     // dot-product with global-reduce
12:     $\delta_{j+1} := \mathbf{r}_{j+1}^T\mathbf{r}_{j+1}$
13:     **if** Converged **then**
14:         break
15:     **else**
16:         // compute next search direction
17:         $\beta_{j+1} := \delta_{j+1}/\delta_j$
18:         $\mathbf{p}_{j+1} := \mathbf{r}_{j+1} + \beta_j\mathbf{p}_j$
19:     **end if**
20: **end for**

- CG is a popular iterative method for symmetric positive definite (SPD) linear system, $Ax = b$.

- It relies on two types of kernels
  - Matrix Vector multiply (SpMV)
    - for generating Krylov subspace = span($p$, $Ap$, $A^2p$, …)
    - typically combined with precondiner to improve convergence
    - used as a black box, provided by users, for supporting a wide range of applications

  - BLAS-1 operations (focus of the paper)
    - for computing search direction, to update solution and residual vectors
    - two **dot**'s (with global all-reduce) and three **axpy**'s

- CG iteration relies on efficient short-term recurrence, but underlying BLAS-1 kernels are latency bound with low performance
  - could become significant in the iteration time (e.g., at large scale)

## *s*-step Conjugate Gradient [Chronopoulos, Gear '89]

- *s*-step CG generates a set of *s* basis vectors at a time
  - Potential reduction in communication cost by a factor of *s*
    - reducing latency cost (one synchronization per *s* steps)
    - exposing more parallelism and data reuse (BLAS-3 instead of BLAS-1)
    - Require $O(1)$ communication for generating $O(s)$ basis vectors.

- Two challenges for practical use
  1. Computational overheads
  2. Numerical stability

```
1: for j = 1, 1 · s + 1, 2 · s + 1, ... do
2:     // Matrix Powers Kernel to generate Krylov space
3:     for k = 1, 2, ..., s do
4:         [r_{j+k}, p_{j+k+1}] := A [r_{j+k-1}, p_{j+k}]
5:     end for
6:     // Compute new solution and residual vector
7:     G := V^T V  with V = [P_{j:j+s}, R_{j:j+s-1}]
8:     ...
9:     [r_{j+s}, x_{j+1}, p_{j+2}] := V[y, t, c]
10: end for
```

# Computational overhead (first challenge)

- In order to reduce communication, it requires additional computation
  - If underlying kernels are optimized for multiple vectors, performance may be improved



**Require:** SPD matrix $A$, RHS vector $\mathbf{b}$, and initial approximate solution vector $\mathbf{x}_1$
1: $\mathbf{r}_1 := \mathbf{b} - A\mathbf{x}_1,\ \delta_1 := \mathbf{r}_1^T\mathbf{r}_1,\ \mathbf{p}_1 := \mathbf{r}_1$
2: **for** $j = 1, 2, \ldots$ **do**
3:     *// SpMV with P2P communication*
4:     $\mathbf{w}_j := A\mathbf{p}_j$
5:     *// dot-product with global-reduce*
6:     $\gamma_j := \mathbf{p}_j^T\mathbf{w}_j$
7:     $\alpha_j := \delta_j/\gamma_j$
8:     *// update solution and residual vectors*
9:     $\mathbf{x}_{j+1} := \mathbf{x}_j + \alpha_j\mathbf{p}_j$
10:    $\mathbf{r}_{j+1} := \mathbf{r}_j - \alpha_j\mathbf{w}_j$
11:    *// dot-product with global-reduce*
12:    $\delta_{j+1} := \mathbf{r}_{j+1}^T\mathbf{r}_{j+1}$
13:    **if** Converged **then**
14:       break
15:    **else**
16:       *// compute next search direction*
17:       $\beta_{j+1} := \delta_{j+1}/\delta_j$
18:       $\mathbf{p}_{j+1} := \mathbf{r}_{j+1} + \beta_j\mathbf{p}_j$
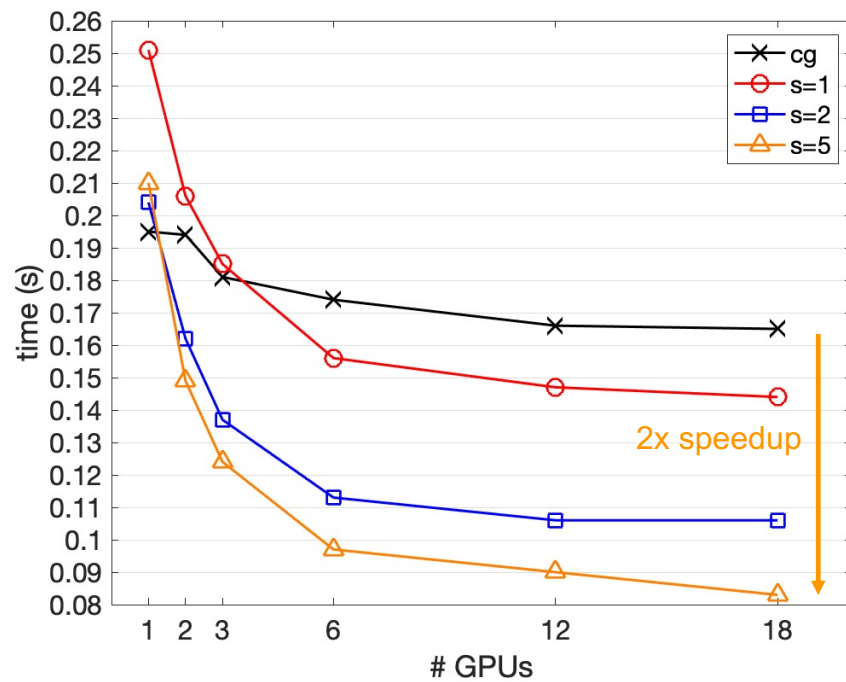19:    **end if**
20: **end for**

1: **for** $j = 1,\ 1\cdot s+1,\ 2\cdot s+1, \ldots$ **do**
2:     *// Matrix Powers Kernel to generate Krylov space*
3:     **for** $k = 1, 2, \ldots, s$ **do**
4:       $[\mathbf{r}_{j+k}, \mathbf{p}_{j+k+1}] := A\,[\mathbf{r}_{j+k-1}, \mathbf{p}_{j+k}]$
5:     **end for**
6:     *// Compute new solution and residual vector*
7:     $G := V^TV$ with $V = [P_{j:j+s}, R_{j:j+s-1}]$
8:     ...
9:     $[\mathbf{r}_{j+s}, \mathbf{x}_{j+1}, \mathbf{p}_{j+2}] := V[\mathbf{y}, \mathbf{t}, \mathbf{c}]$
10: **end for**

"$s$" **SpMM**s with two vectors = **2× flops**, but two vectors at a time

One dot-products with "2s+1" vectors = **2s× flops**, but with one **SYRK**

3 **GEMV**s with "$2s+1$" columns = **2× flops**, but with one **GEMM**
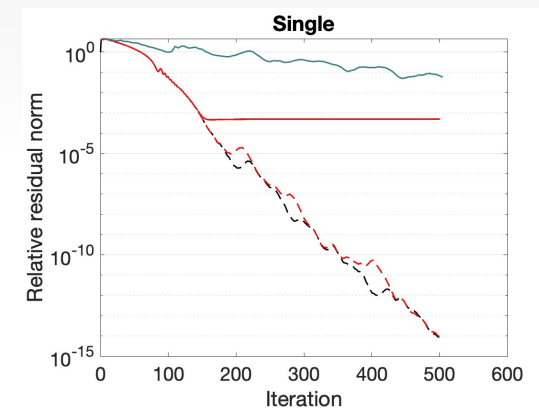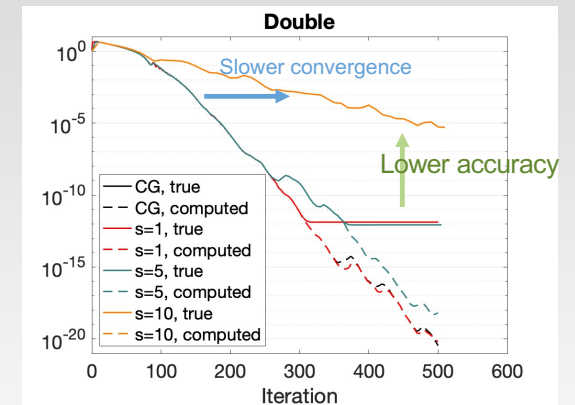
# Strong-scaling results on Summit (V100 GPUs)



- 500 CG iteration time with 7-pts Laplace 3D ($n_x=100$)

- When communication (e.g., **latency**) becomes significant, *s*-step CG may reduce iteration time, even with the computational overhead

# Potential numerical instability with *s*-step CG (second challenge)

- **Both convergence rate and attainable accuracy can deteriorate**

- Potentially very ill-conditioned *s*-step basis vectors $V_{2s+1}$
  - Condition number $\kappa(V_{2s+1})$ can grow exponentially with *s*
  - Orthogonality errors can grow quadratically to the condition number
    - The Gram matrix $G$ has the squared condition number

# Mixed-precision *s*-step CG with **residual replacement** on GPUs
for improving convergence and accuracy

- Higher-precision to improve convergence behavior [Carson, Gergelits, '21]
  - form the Gram matrix $G$ in **double** the working precision
  - orthogonality error depends linearly, instead of quadratically, to condition number of *s*-step basis vectors

- Residual replacement to improve solution accuracy [Carson, Demmel '14]
  - replace computed residual vector with true residual vector at "selected" iterations
  - the selection requires the computation of $\bar{G} = |V_{(j-1)/s}|^T |V_{(j-1)/s}|$, in the working precision
  - *s*-step CG obtains the same residual norm bound as standard CG, $O(\epsilon)\|A\|\|x\|$.

```
1:  r_1 := b − Ax_1, δ_1 := r_1^T r_1, p_1 := r_1
2:  z = 0
3:  for j+= ŝ do
4:      // MPK (on GPUs) with P2P communication
5:      P_j := p_j, R_j := r_j
6:      P_{j+1} := AP_j
7:      for k = 1, 2, . . . , s − 1 do
8:          [R_{j+k}, P_{j+k+1}] := A[R_{j+k−1}, P_{j+k}]
9:      end for
10:     // dot-products (on GPUs) with global-reduce
11:     G := V_ℓ^T V_ℓ,
12:         where V_ℓ := [P_{j:j+s}, R_{j:j+s−1}] and ℓ := (j−1)/s
13:     if Converged then
14:         break
15:     end if
16:     // update coefficients (redundantly on each host)
17:     c_1 := e_1, t_1 := e_{s+1}, δ_k := t_k G t_k
18:     for k = 1, 2, . . . , s do
19:         d_k := Bc_k, γ_k := c_k G d_k, α_k := δ_k/γ_k
20:         y_{k+1} := y_k + α_k d_k
21:         t_{k+1} := t_k − α_k d_k
22:         δ_{k+1} := t_{k+1}^T G t_{k+1}, β_k := δ_{k+1}/δ_k
23:         c_{k+1} := t_{k+1} + β_k c_k
24:         if time to replace residual vector then
25:             x_{j+k} := x_j + [P_{j:j+s} R_{j:j+s−1}]y_{k+1}
26:             p_{j+k} := [P_{j:j+s} R_{j:j+s−1}]c_{k+1}
27:             z = z + x_{j+k}
28:             r_{j+k} := b − Az_{j+k}
29:             x_{j+k} := 0
30:             ŝ := k
31:             break
32:         end if
33:     end for
34:     if not replaced then
35:         // update vectors (on GPUs)
36:         x_{j+s} := x_j + [P_{j:j+s} R_{j:j+s−1}]y_{s+1}
37:         r_{j+s} := [P_{j:j+s} R_{j:j+s−1}]t_{s+1}
38:         p_{j+s} := [P_{j:j+s} R_{j:j+s−1}]c_{s+1}
39:         ŝ := s
40:     end if
41: end for
42: x := x_end + z
```

# Numerical results with mixed-precision *s*-step CG with RR
using 3D Laplace ($n = 100^3$)



**Uniform-precision**

Legend:
- CG
- CG+RR
- s=5
- s=5, RR
- s=10
- s=10, RR

mixed

RR

**Mixed-precision**
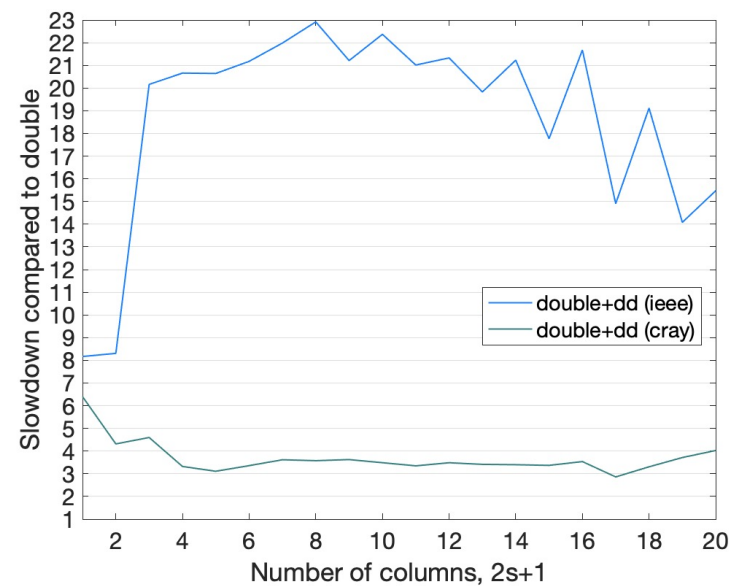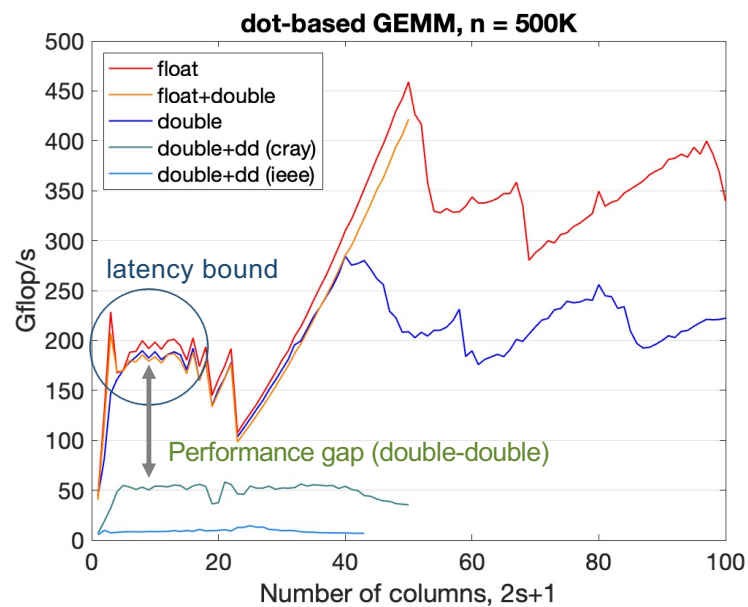
Relative residual norm vs Iteration

Both mixed-precision and RR are needed to obtain convergence similar to standard CG

- Higher-precision improves the convergence
- RR improves the accuracy

# Implementation for performance study on a GPU cluster (Summit)

- MPI (cuda-aware) for data exchange between GPUs

- Kokkos for portable performance on different manycore architectures
  - we only show performance on NVIDIA GPUs,

- Mixed-precision dot-products to compute $G$
  - It reads "big" tall-skinny $V$ in working precision, but internally use higher precision to compute "small" $G$
  - It is latency bound, hopefully with a small overhead (of computing and writing $G$ in higher precision)

- **double** or **single** precision as our working precision
  - double working precision
    - typical for scientific and engineering application
    - may require software-emulated higher precision (double-double in our experiments on V100 GPUs)
  - single working precision
    - experiments where higher precision is implemented by hardware,
    - practical use of single-precision CG exists, e.g., mixed-precision reliable updates and iterative refinements

# Performance of mixed-precision dot-products using Kokkos on one NVIDIA V100 GPU
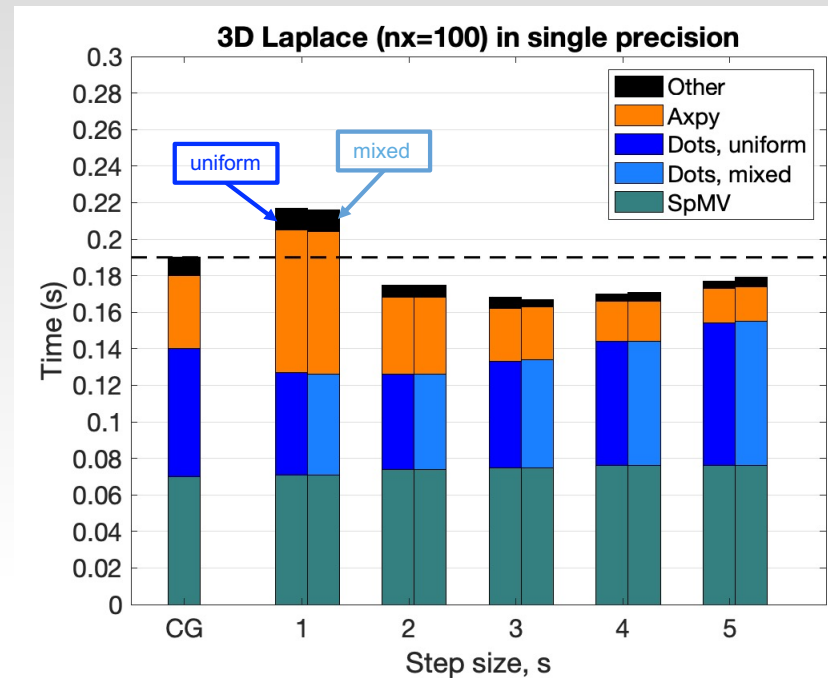


- Dot-products is often latency bound, and virtually no overhead using higher precision, when implemented by hardware
  - For mixed single+double or uniform double, vs uniform single
- Mixed-precision dot-products with Cray-style double-double requires 17× more flops, while IEEE variant requires 21× more flops
- The overhead tends to become smaller on multiple GPUs.

# Iteration time breakdown for single working precision on <u>one</u> NVIDIA V100 GPU
## hardware-implemented higher precision

- Uniform precision

  - SpMM with two vectors is as fast as SpMV with one vector

  - Vector-updates with multiple vectors continue to improve the performance with a larger $s$

  - Dot-products improves the performance for a small $s$, but the computational overhead ($2s\times$) becomes significant for a large $s$

- Mixed-precision, float + double, dot-product

  - Input vectors are read in single precision

  - No overhead in performing multiply-add and write back $G$, in double precision

  - Any reduction in the iteration count has direct impact to time-to-solution, no overhead even if not reduction in the iteration count



3D Laplace (nx=100) in single precision

Kernel performance study on <u>one</u> GPU (up to 1.1x), while larger speedups on multiple GPUs

# Time-to-solution with mixed-precision single+double *s*-step CG on <u>six</u> NVIDIA V100 GPUs
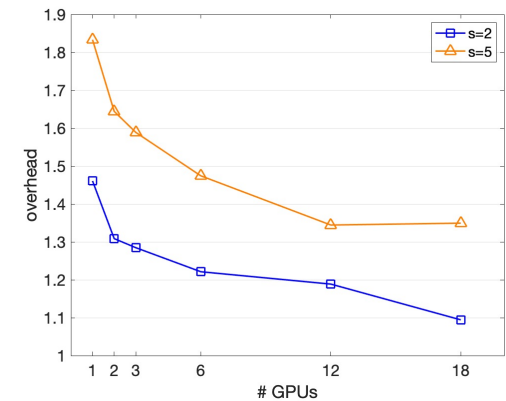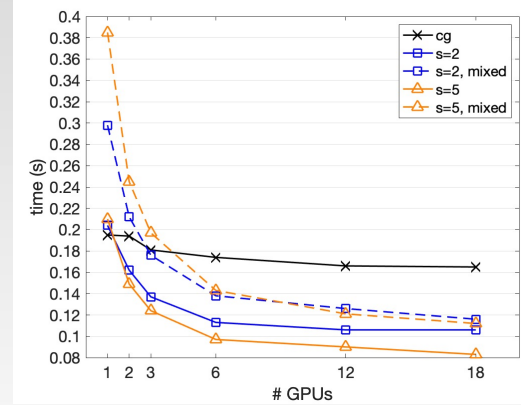hardware-implemented higher precision



- *s*-step reduces the time-per-iteration, but can suffer from numerical instability

- mixed-precision improves the stability with virtually no overhead, leading to faster to solution

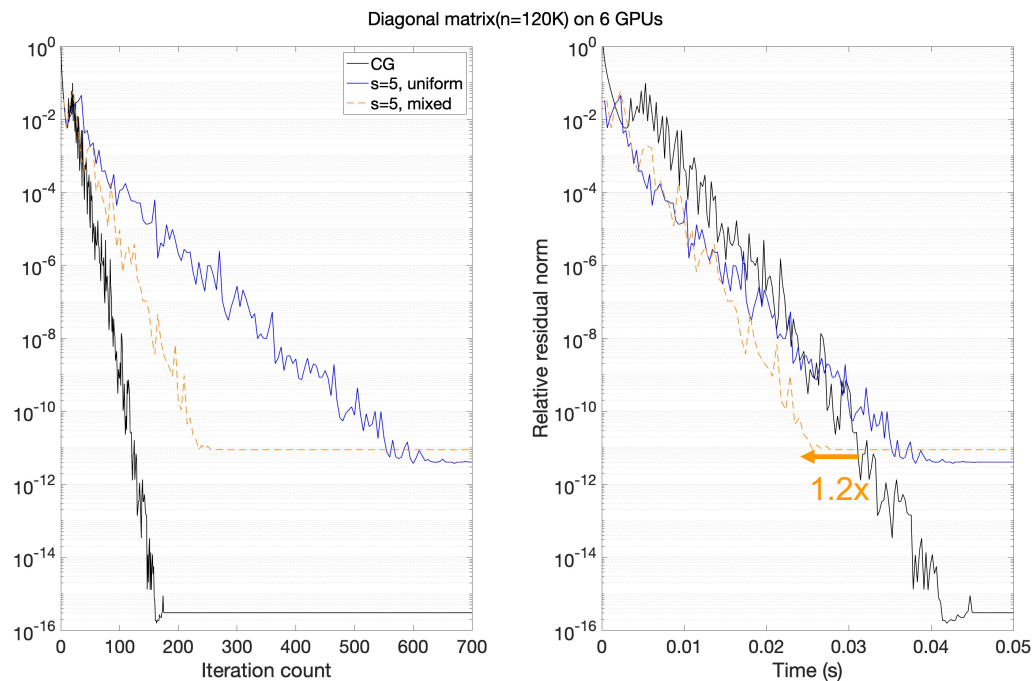# Iteration time breakdown with double-precision on an NVIDIA V100 GPU
## software-emulated higher precision



3D Laplace (nx=100) in double precision





- mixed-precision dot-products has significant overhead especially with a large $s$
  - Overhead becomes smaller as latency become more significant on multiple GPUs

# Time-to-solution with mixed double + double-double $s$-step CG
 software-emulated higher precision



Diagonal matrix(n=120K) on 6 GPUs

- Using diagonal matrix,

  $$a_{i,i} = \lambda_1 + (i-1)/(n-1)(\lambda_n - \lambda_1)\rho^{n-i}.$$

  - Overhead of dot-products is more significant
  - Allows controlling the conditioning of the matrix

- Even with software-emulated higher precision, there are cases where mixed-precision variant reduces the time-to-solution

  - when $s$-step CG suffers from instability with a small step size, $s$

# Mixed-precision *s*-step CG with **residual replacement**

```
1:  r₁ := b − Ax₁, δ₁ := r₁ᵀr₁, p₁ := r₁
2:  z = 0
3:  for j+ = ŝ do
4:      // MPK (on GPUs) with P2P communication
5:      Pⱼ := pⱼ, Rⱼ := rⱼ
6:      Pⱼ₊₁ := APⱼ
7:      for k = 1, 2, . . . , s − 1 do
8:          [Rⱼ₊ₖ, Pⱼ₊ₖ₊₁] := A[Rⱼ₊ₖ₋₁, Pⱼ₊ₖ]
9:      end for
10:     // dot-products (on GPUs) with global-reduce
11:     G := VₗᵀVₗ,
12:         where Vₗ := [Pⱼ:ⱼ₊ₛ, Rⱼ:ⱼ₊ₛ₋₁] and ℓ := (j−1)/s
13:     if Converged then
14:         break
15:     end if
16:     // update coefficients (redundantly on each host)
17:     c₁ := e₁, t₁ := e_{s+1}, δ₁ := t₁Gt₁
18:     for k = 1, 2, . . . , s do
19:         dₖ := Bcₖ, γₖ := cₖGdₖ, αₖ := δₖ/γₖ
20:         yₖ₊₁ := yₖ + αₖcₖ
21:         tₖ₊₁ := tₖ − αₖdₖ
22:         δₖ₊₁ := tₖ₊₁ᵀGtₖ₊₁, βₖ := δₖ₊₁/δₖ
23:         cₖ₊₁ := tₖ₊₁ + βₖcₖ
24:         if time to replace residual vector then
25:             xⱼ₊ₖ := xⱼ + [Pⱼ:ⱼ₊ₛRⱼ:ⱼ₊ₛ₋₁]yₖ₊₁
26:             pⱼ₊ₖ := [Pⱼ:ⱼ₊ₛRⱼ:ⱼ₊ₛ₋₁]cₖ₊₁
27:             z = z + xⱼ₊ₖ
28:             rⱼ₊ₖ := b − Azⱼ₊ₖ
29:             xⱼ₊ₖ := 0
30:             ŝ := k
31:             break
32:         end if
33:     end for
34:     if not replaced then
35:         // update vectors (on GPUs)
36:         xⱼ₊ₛ := xⱼ + [Pⱼ:ⱼ₊ₛRⱼ:ⱼ₊ₛ₋₁]y_{s+1}
37:         rⱼ₊ₛ := [Pⱼ:ⱼ₊ₛRⱼ:ⱼ₊ₛ₋₁]t_{s+1}
38:         pⱼ₊ₛ := [Pⱼ:ⱼ₊ₛRⱼ:ⱼ₊ₛ₋₁]c_{s+1}
39:         ŝ := s
40:     end if
41: end for
42: x := x_end + z
```

- **residual replacement** to improve the attainable accuracy
  - With residual replacement, *s*-step CG obtains the same residual norm bound as standard CG, $O(\epsilon)\|A\|\|x\|$, [Carson, Demmel '14]
  - The detection require the computation of $\bar{G} = |V_{(j-1)/s}|^T |V_{(j-1)/s}|$, in the working precision
  - If the residual needs to be replaced before *s*-th step, we waste some computation needed to form $G$, and also take a step smaller than *s*

- dot-products is latency or bandwidth limited
  improve convergence (a fewer iterations) without significant increase in iteration time

# Time-to-solution with mixed *s*-step CG with RR



- RR adds overhead, but improves the attainable solution accuracy
- More results in paper

## Final remarks

- We studied mixed-precision $s$-step CG with residual replacements on GPUs
  - When the higher-precision is supported by hardware, it improves the stability with virtually no overhead, and hence reduces time-to-solution, or if not, no overhead.
  - If the higher-precision requires software-emulation, the overhead becomes significant. It may still help when $s$-step CG becomes unstable with a small step size, and the latency becomes significant in the iteration time.

- We are planning on some extensions/variations of the algorithm
  - We have only looked at monomial basis. Combining these techniques with more stable basis (e.g., Newton, Chebychev) may further improve the stability, and practicability, of $s$-step CG
  - We only looked at two-term recurrence variant of s-step CG. There is also three-term recurrence variant, where relative cost of dot-products is smaller in iteration cost, and hence the mixed-precision may be more attractive.

Thank you!!

## Acknowledgements