



# Accelerating Sparse Iterative Solvers and Preconditioners Using RACE

Christie Alappat, Georg Hager, Gerhard Wellein

Erlangen National High Performance Computing Center (NHR@FAU)

SIAM PP24

# Matrix power kernel (MPK)

- Calculate:  $y = A^p x$
- Repeatedly perform back to back SpMVs

```
for k=1:p; do
  y[k] = SpMV(A, y[k-1])
done
```

Matrix A loaded p times from main memory.

\_\_\_\_\_

Can I cache matrix A?

Mohiyuddin et al., 2009. Minimizing communication in sparse matrix solvers. In Proceedings of the SC'09. <u>https://doi.org/10.1145/1654059.1654096</u> But requires "ghosting". Indirect accesses or redundant copies of the matrix entries, which typically increases with number of threads.



#### Can I avoid these overhead?



## Matrix power



Alappat et al, "Level-Based Blocking for Sparse Matrices: Sparse Matrix-Power-Vector Multiplication," in *IEEE Transactions on Parallel and Distributed Systems*, 2023, doi: 10.1109/TPDS.2022.3223512.

### Matrix power kernel: Performance



#### Matrix power kernel: Performance



### Matrix power kernel: Performance







#### **Iterative solvers**

#### Solve for $\overline{x: Ax = b}$



### Application: s-step Krylov schemes

#### Basic computation kernel

#### GMRES scheme: main kernel

```
for j=0:1:m; do
  v[j+1] = SpMV(A, v[j])
  Orthogonalize v[j+1] against v[0:j]
done
```

Also called CA-GMRES

Available with Trilinos framework.

RACE integrated with Trilinos to accelerate MPK.

Mark Hoemmen, 2010, Communication-avoiding Krylov subspace methods, PhD Thesis, <a href="https://www2.eecs.berkeley.edu/Pubs/TechRpts/2010/EECS-2010-37.pdf">https://www2.eecs.berkeley.edu/Pubs/TechRpts/2010/EECS-2010-37.pdf</a>

RACE, SIAM PP24, Baltimore

s-step GMRES scheme: main kernel





#### s-step GMRES



*s*-step GMRES: Belos (TPETRA) library settings: *s*=4, restart length (*m*)=50

RACE, SIAM PP24, Baltimore





## **MPK and preconditioners**

## Solve for *x*: $AM^{-1}u = b$ , $M^{-1}u = x$



#### Preconditioning s-step GMRES



## Polynomial preconditioner

- Use matrix polynomials  $(c_1Ax + c_2A^2x + ... + c_dA^dx)$  to approximate  $A^{-1}$ MPK
- Higher degree (d) of polynomial  $\rightarrow$  Better approximation
- Can use preconditioners like Jacobi, ILU on top of it
- Available in Trilinos\*

\*J. Loe, H. Thornquist, E. Boman, 2020, Polynomial Preconditioned GMRES in Trilinos: Practical Considerations for High-Performance Computing, <u>https://doi.org/10.1137/1.9781611976137.4</u>

#### Polynomial preconditioner with GMRES

Baseline:MPK+PreconBaseline:OrthoBaseline:MiscRACE:MPK+PreconRACE:OrthoRACE:Misc



Can be combined with any Krylov solvers. Does not strictly need *s*-step solvers.

High powers in MPK→ RACE very effective.

Ortho cost becomes negligible.

GMRES polynomial precon: Belos

Poly+Jacobi precon, degree (d)=80

RACE, SIAM PP24, Baltimore





## Algebraic multigrid (AMG)



#### AMG

Smoothers involve back-to-back application of the same operation.

→ Use RACE to cache block the smoother.

Currently RACE applied only to the finest level.







## **Case study: Nalu-Wind**



Wind turbine simulation using Navier-Stokes equation on unstructured grid.

Picture taken from: Sprague et al., ``ExaWind: A multifidelity modeling and simulation environment for wind energy", Journal of Physics, 2020, 10.1088/1742-6596/1452/1/012071

#### Case study: Nalu-Wind



In traditional setting polynomial preconditioner is not beneficial.

But RACE can change the picture.

→ Adds another dimension to solver choice.

RACE can accelerate multiple sweeps of the same preconditioner <sup>(3)</sup>

GMRES solver runtime for solving momentum equation of dimension  $N_r \approx 12 M$ ,  $N_{nz} = 300 M$ 

- MPK kernel's performance can be improved by level-based cache blocking using RACE.
- Speedups up to 5x possible.
- Benefits iterative solvers and its components: s-step Krylov solvers, polynomial preconditioners, AMG, …
- RACE adds another dimension to solver selection/tuning.

## Outlook

- Solvers like *s*-step GMRES and polynomial preconditioners are easy to parallelize and have lower communication overheads → promising for large-scale solvers.
- MPI parallel version of RACE already in deployment.
- Cache blocking on GPU produces first successful results.





# Thank you Questions



#### https://github.com/RRZE-HPC/RACE

C. Alappat, G. Hager, O. Schenk and G. Wellein, "Level-Based Blocking for Sparse Matrices: Sparse Matrix-Power-Vector Multiplication," in IEEE Transactions on Parallel and Distributed Systems, 2023, doi: 10.1109/TPDS.2022.3223512.

C. Alappat, A. Basermann, A.R. Bishop, H. Fehske, G. Hager, O. Schenk, J. Thies, and G. Wellein. 2020. A Recursive Algebraic Coloring Technique for Hardware-efficient Symmetric Sparse Matrix-vector Multiplication. ACM Trans. Parallel Comput., 2020. https://doi.org/10.1145/3399732