

Ma221: Matrix Computations

Lecture 28:

Communication-Avoiding Sparse Iterative Methods

James Demmel
(and many collaborators)

www.cs.berkeley.edu/~demmel

Motivation

- Iterative Methods (both Splitting Methods and Krylov Subspace Methods) rely on Sparse-Matrix-Vector-Multiplication (SpMV)
 - If A too large to fit in fast memory, need to reread A at each iteration, do only 2 flops per matrix entry, so limited by communication
 - On a parallel machine, typically need to do an operation like a dot product (eg to compute $\|Ax-b\|_2$) after each iteration, also a bottleneck

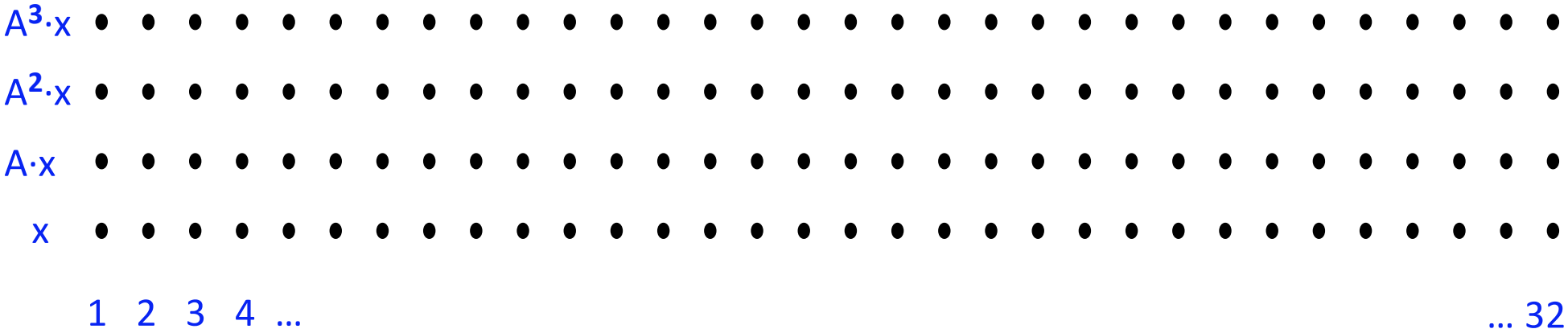
Goal

- Reorganize iterative algorithms to take k steps for the communication cost of 1 step
 - k will be a tuning parameter to optimize speedup
- Challenges:
 - A needs to be sparse enough
 - Tradeoff avoiding communication with redundant flops
 - Need to maintain numerical stability
 - Simplest approaches not stable
 - Preconditioning

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$

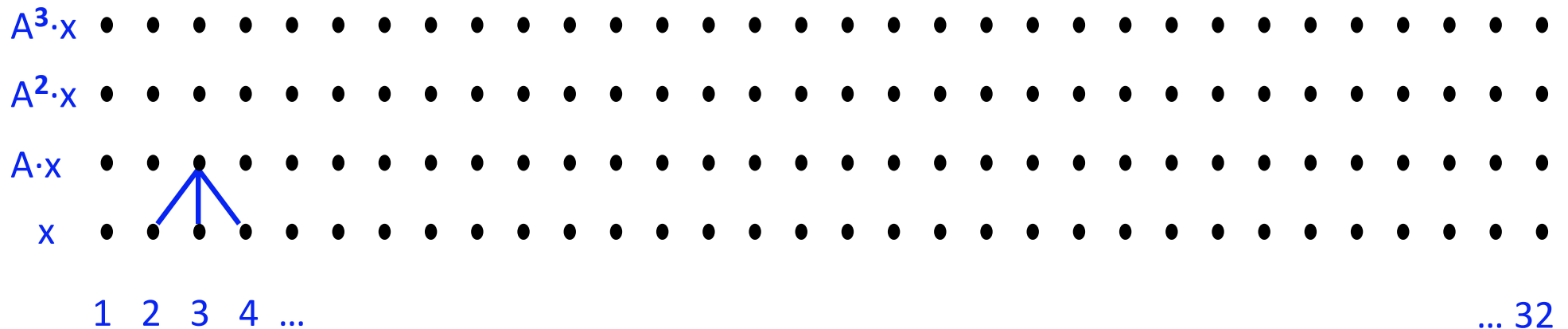


- Example: A tridiagonal, $n=32$, $k=3$
- Works for any “well-partitioned” A

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$

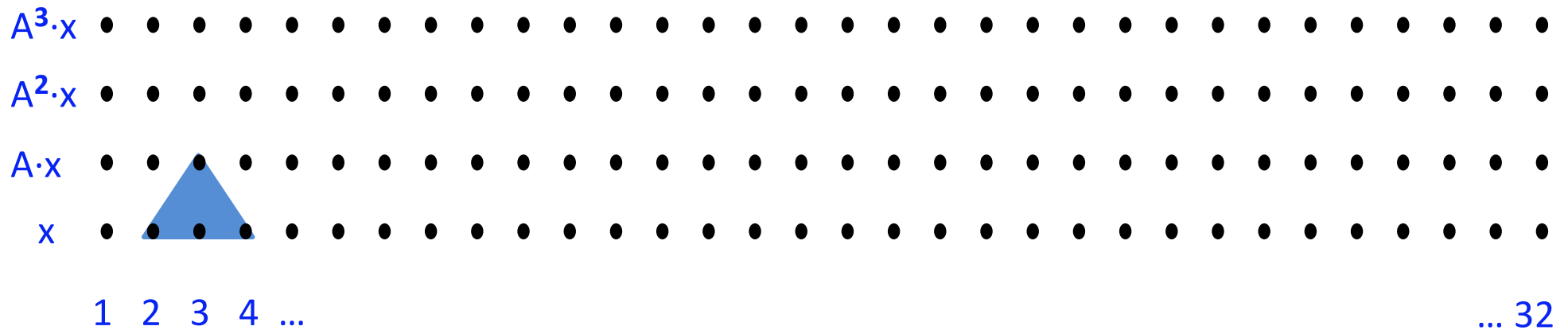


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$

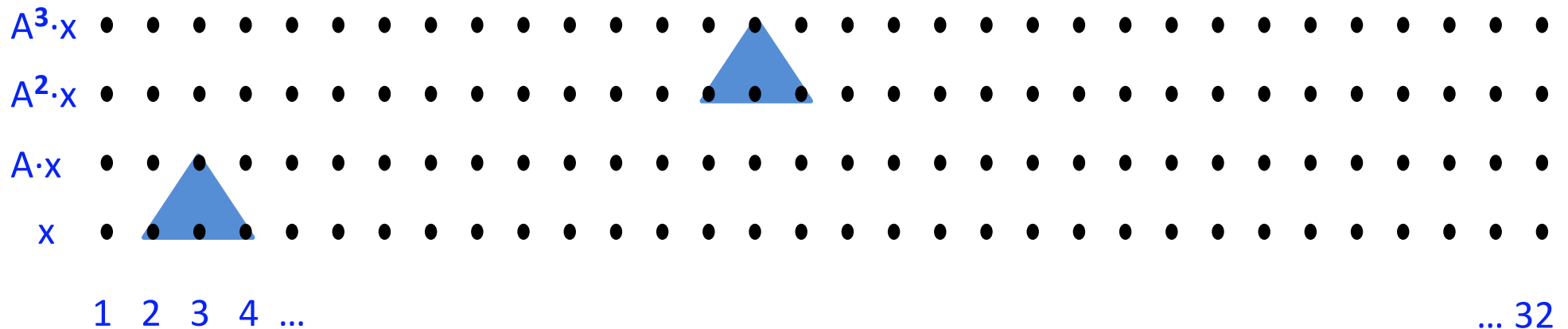


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$

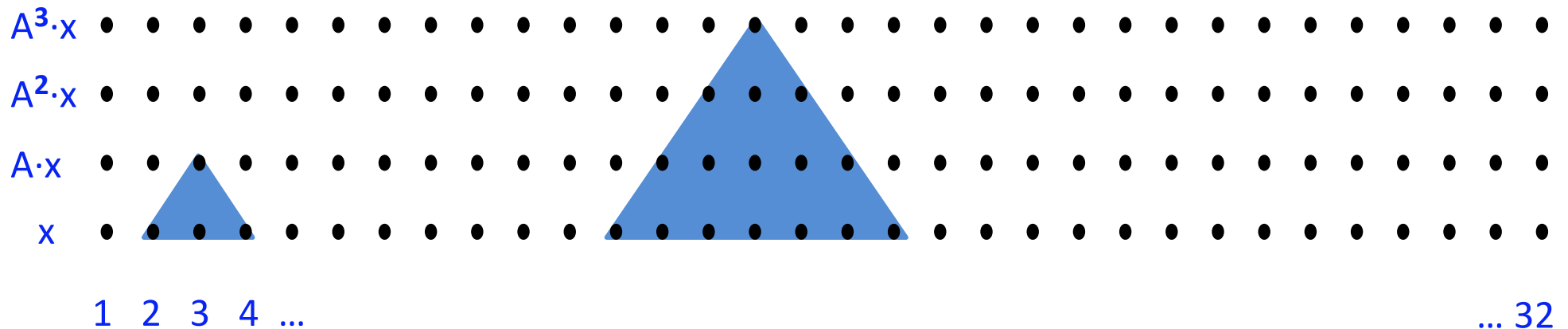


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$

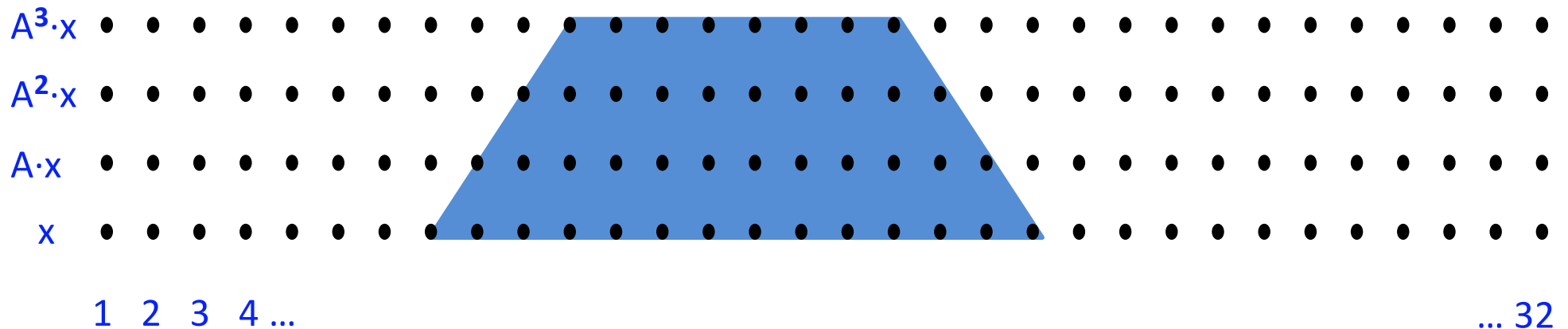


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$

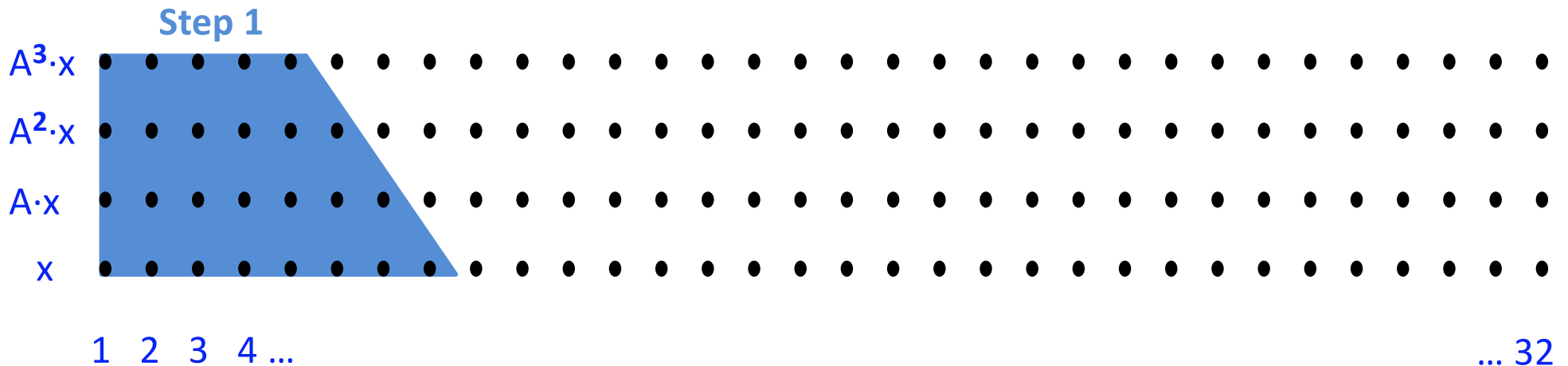


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$
- Sequential Algorithm

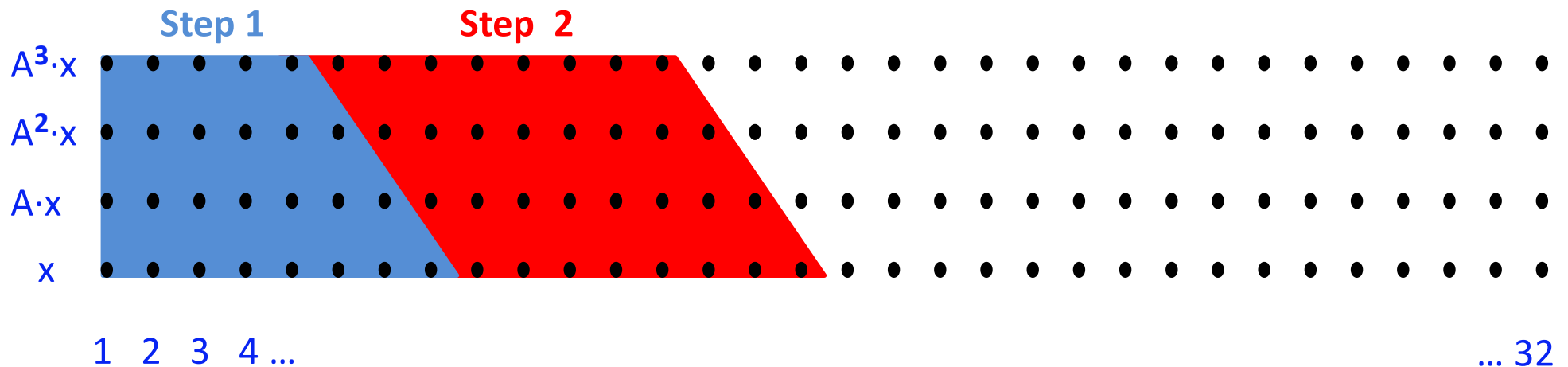


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$
- Sequential Algorithm

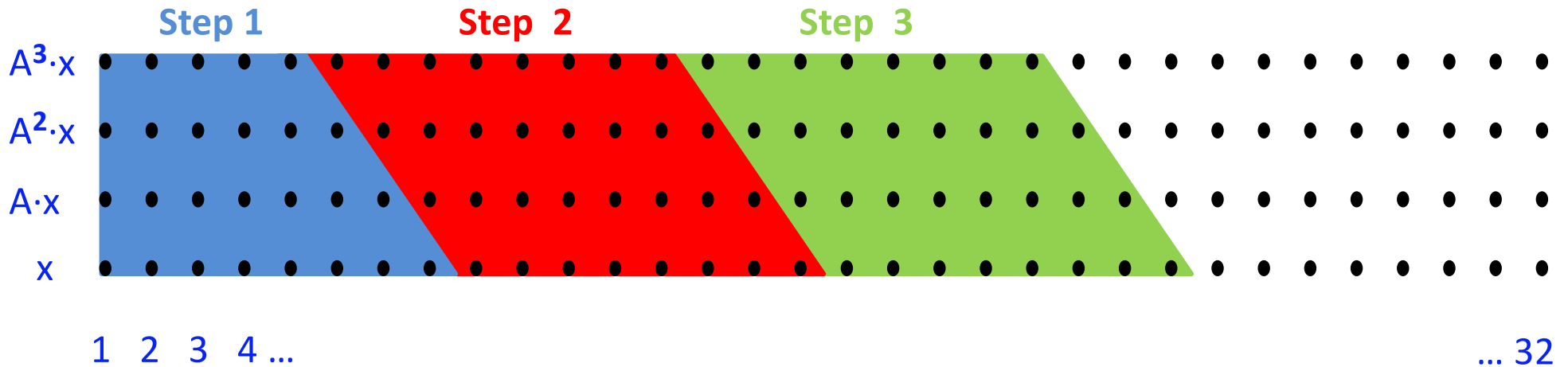


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$
- Sequential Algorithm

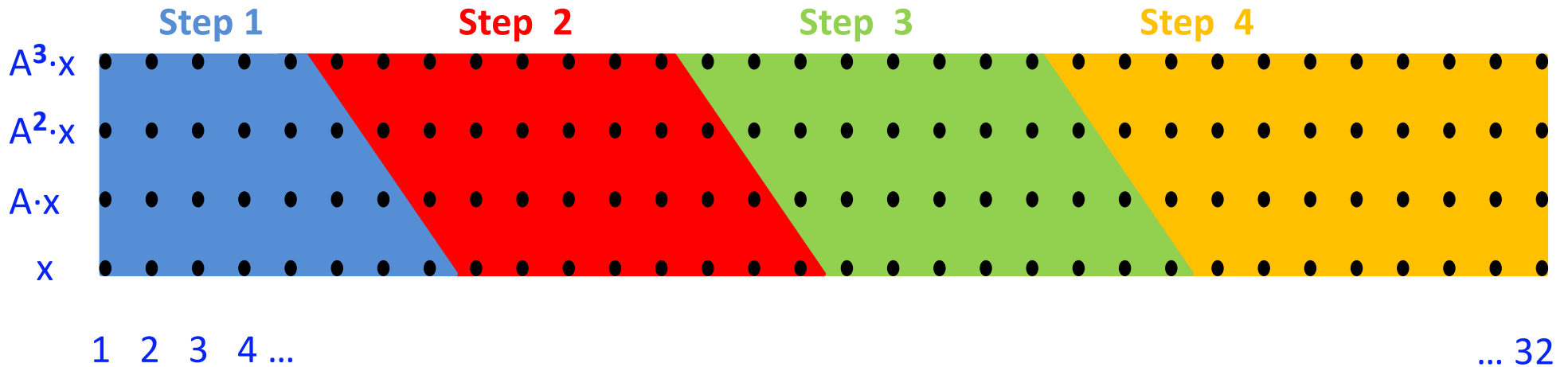


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$
- Sequential Algorithm

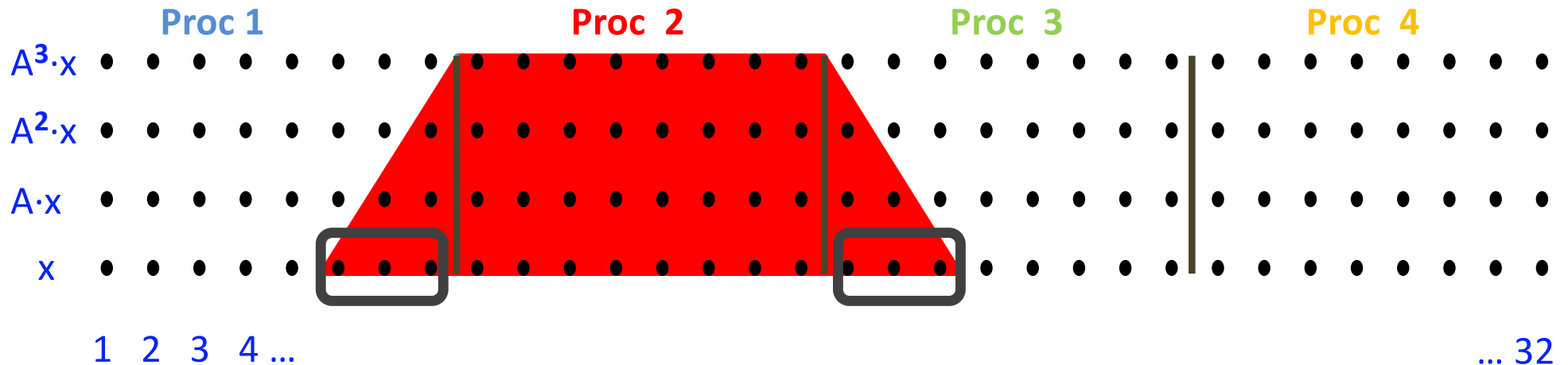


- Example: A tridiagonal, $n=32$, $k=3$

Communication Avoiding Kernels:

The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$
- Parallel Algorithm

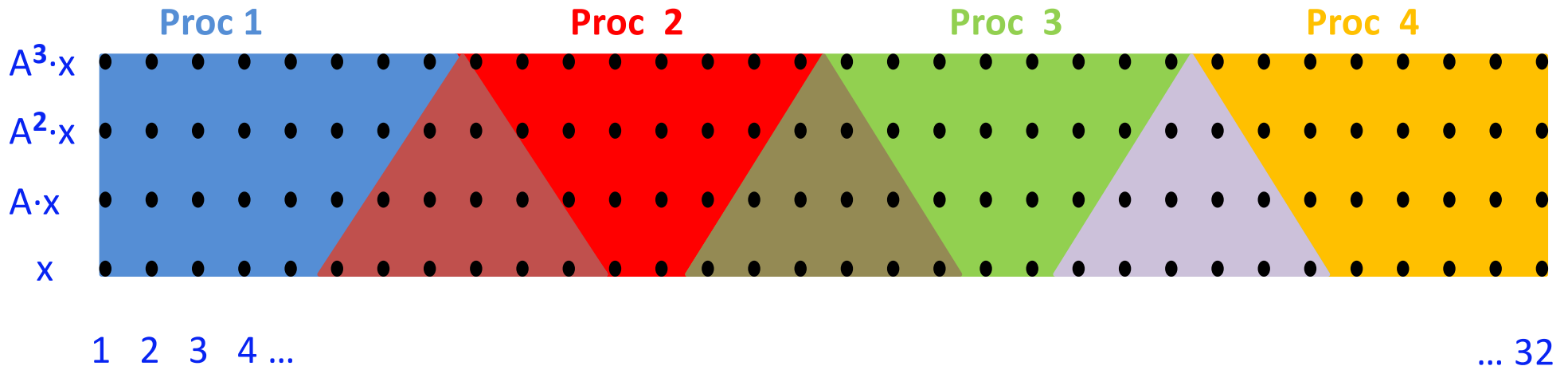


- Example: A tridiagonal, $n=32$, $k=3$
- Each processor communicates once with neighbors

Communication Avoiding Kernels:

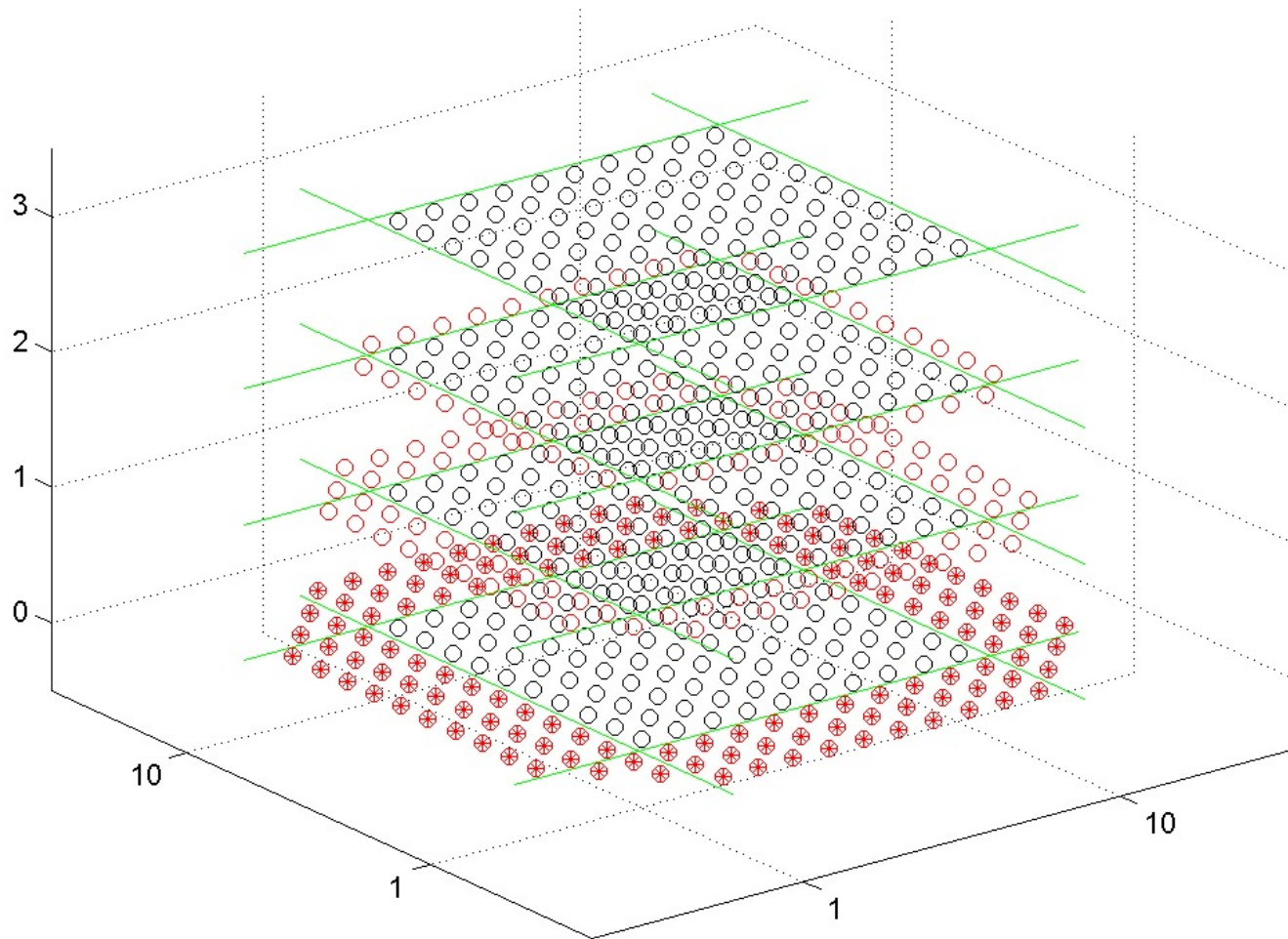
The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$

- Replace k iterations of $y = A \cdot x$ with $[Ax, A^2x, \dots, A^kx]$
- Parallel Algorithm

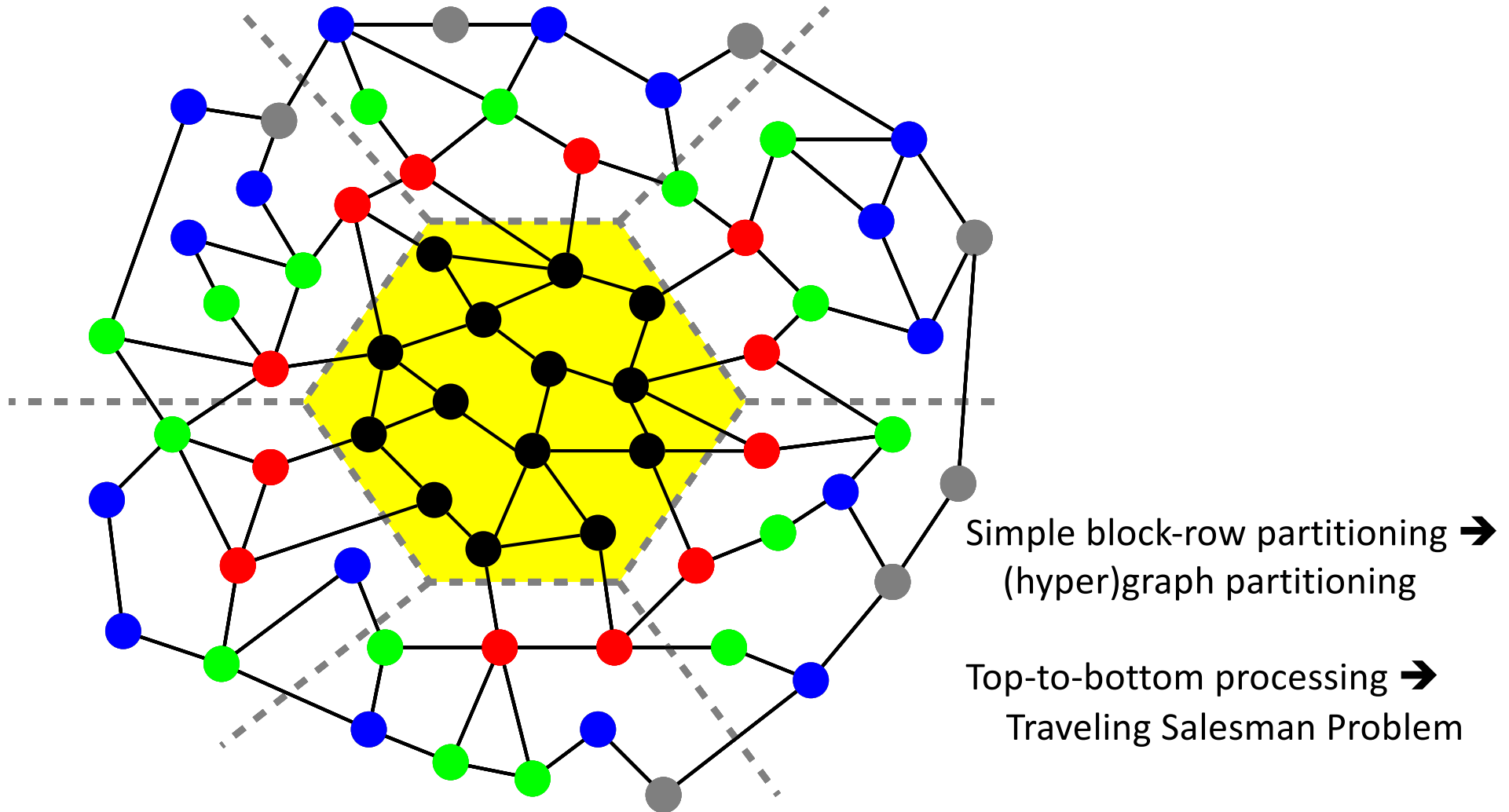


- Example: A tridiagonal, $n=32$, $k=3$
- Each processor works on (overlapping) trapezoid

Matrix Power Kernel $[x, Ax, A^2x, A^3x]$ for 2D Poisson



The Matrix Powers Kernel : $[Ax, A^2x, \dots, A^kx]$ on a general matrix (nearest k neighbors on a graph)



Same idea for general sparse matrices: k -wide neighboring region

Minimizing Communication of GMRES to solve $Ax=b$

- GMRES: find x in $\text{span}\{b, Ab, \dots, A^k b\}$ minimizing $\|Ax-b\|_2$

Standard GMRES

for $i=1$ to k

$w = A \cdot v(i-1)$... *SpMV*

MGS($w, v(0), \dots, v(i-1)$)

update $v(i), H$

endfor

solve LSQ problem with H

Communication-avoiding GMRES

$W = [v, Av, A^2v, \dots, A^k v]$

$[Q, R] = \text{TSQR}(W)$

... *“Tall Skinny QR”*

build H from R

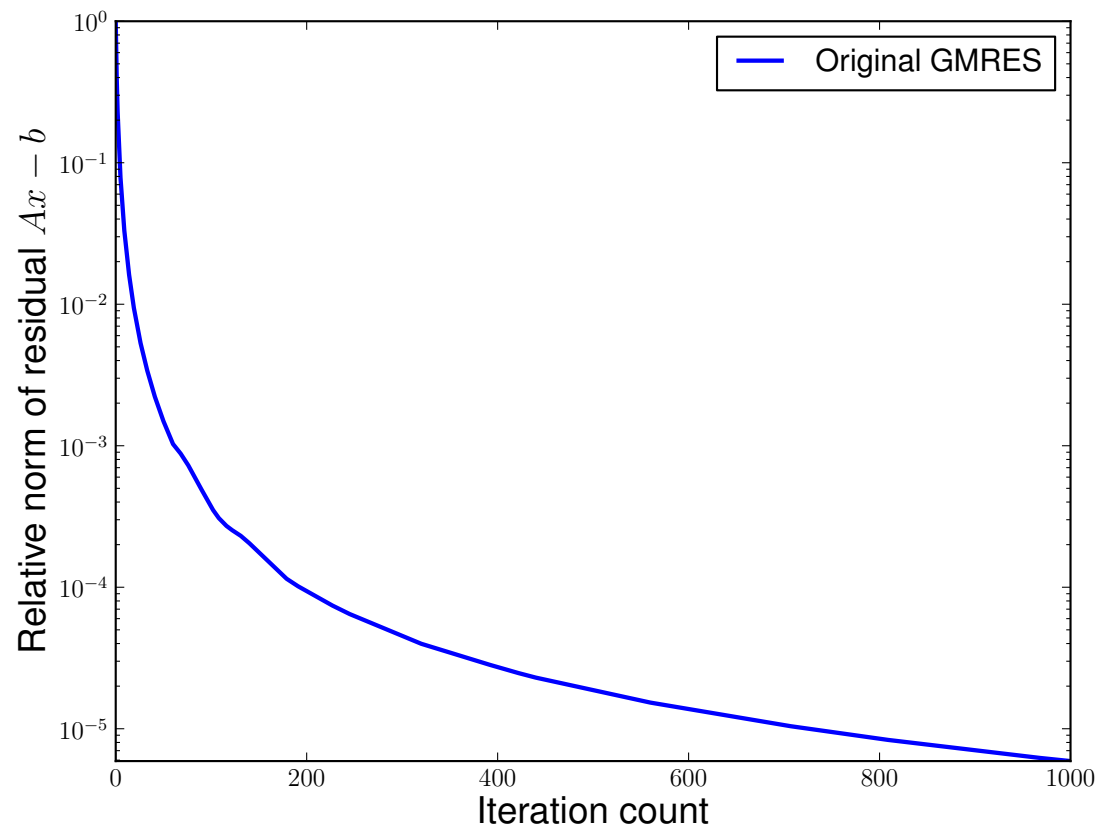
solve LSQ problem with H

Sequential case: #words moved decreases by a factor of k

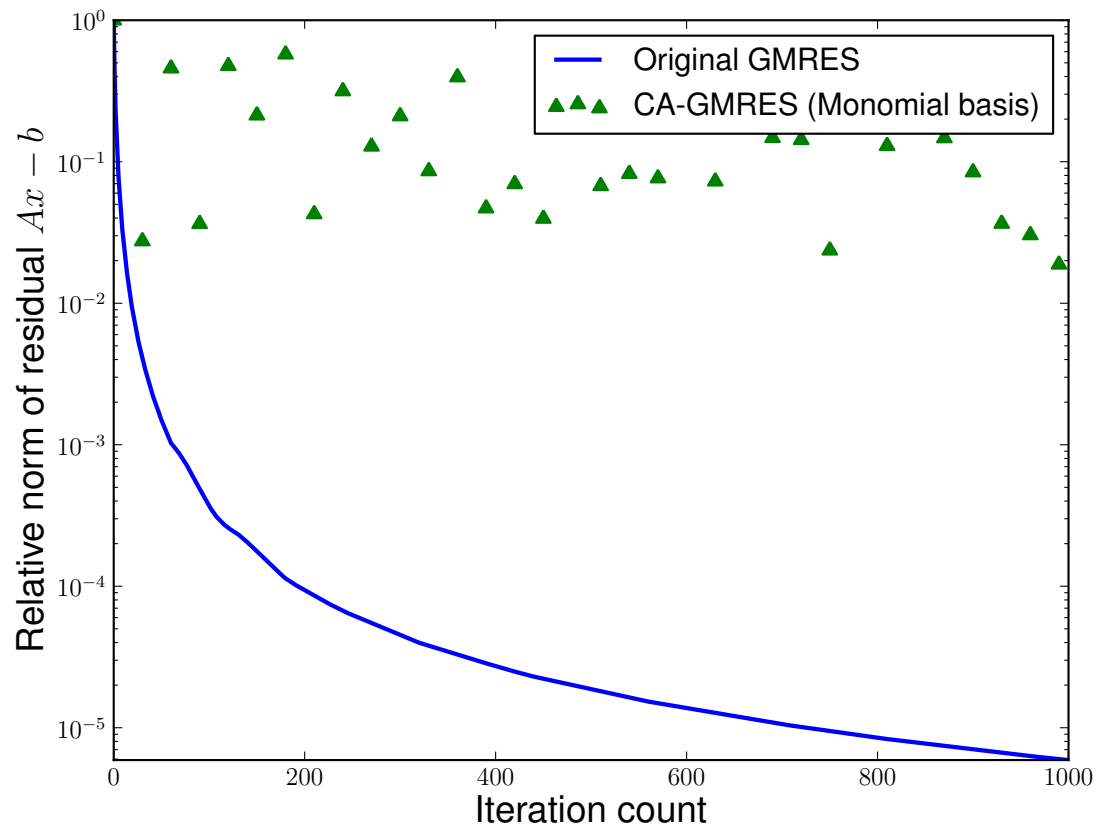
Parallel case: #messages decreases by a factor of k

- **Oops – W from power method, precision lost!**
- **Fix: replace W by $[v, p_1(A)v, p_2(A)v, \dots, p_k(A)v]$**

Convergence of GMRES

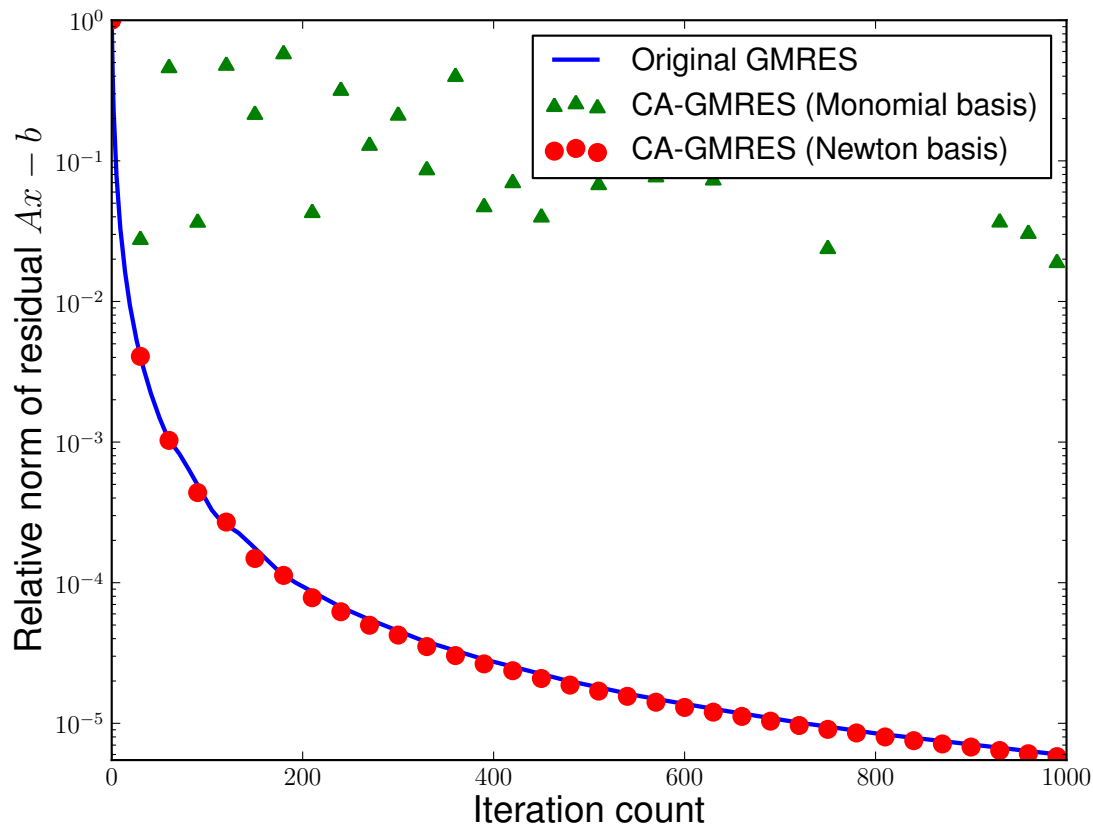


Convergence of GMRES



Convergence of GMRES

Up to **2.3x** speedup on 8 core Intel Clovertown (k=5)



Compute $r_0 = b - Ax_0$. Choose r_0^* arbitrary.

Set $p_0 = r_0$, $q_{-1} = 0_{N \times 1}$.

For $k = 0, 1, \dots$, until convergence, Do

$$P = [p_{sk}, Ap_{sk}, \dots, A^s p_{sk}]$$

$$Q = [q_{sk-1}, Aq_{sk-1}, \dots, A^s q_{sk-1}]$$

$$R = [r_{sk}, Ar_{sk}, \dots, A^s r_{sk}]$$

//Compute the $1 \times (3s+3)$ Gram vector.

$$g = (r_0^*)^T [P, Q, R]$$

//Compute the $(3s+3) \times (3s+3)$ Gram matrix

$$G = \begin{bmatrix} P^T \\ Q^T \\ R^T \end{bmatrix} [P \quad Q \quad R]$$

For $\ell = 0$ to s ,

$$b_{sk}^\ell = [B_1(:, \ell)^T, 0_{s+1}^T, 0_{s+1}^T]^T$$

$$c_{sk-1}^\ell = [0_{s+1}^T, B_2(:, \ell)^T, 0_{s+1}^T]^T$$

$$d_{sk}^\ell = [0_{s+1}^T, 0_{s+1}^T, B_3(:, \ell)^T]^T$$

1. Compute $r_0 := b - Ax_0$; r_0^* arbitrary;
2. $p_0 := r_0$.
3. For $j = 0, 1, \dots$, until convergence Do:
4. $\alpha_j := (r_j, r_0^*) / (Ap_j, r_0^*)$
5. $s_j := r_j - \alpha_j Ap_j$
6. $\omega_j := (As_j, s_j) / (As_j, As_j)$
7. $x_{j+1} := x_j + \alpha_j p_j + \omega_j s_j$
8. $r_{j+1} := s_j - \omega_j As_j$
9. $\beta_j := \frac{(r_{j+1}, r_0^*)}{(r_j, r_0^*)} \otimes \frac{\alpha_j}{\omega_j}$
10. $p_{j+1} := r_{j+1} + \beta_j (p_j - \omega_j Ap_j)$
11. EndDo

CA-BiCGStab

Up to **4.2x** speedup for BiCGStab on 24K core Cray XE6 (used as coarse grid solver in multigrid)

For $j = 0$ to $\lfloor \frac{s}{2} \rfloor - 1$, Do

$$\alpha_{sk+j} = \frac{\langle g, d_{sk+j}^0 \rangle}{\langle g, b_{sk+j}^1 \rangle}$$

$$q_{sk+j} = r_{sk+j} - \alpha_{sk+j} [P, Q, R] b_{sk+j}^1$$

For $\ell = 0$ to $s - 2j + 1$, Do

$$c_{sk+j}^\ell = d_{sk+j}^\ell - \alpha_{sk+j} b_{sk+j-1}^{\ell+1}$$

//such that $[P, Q, R] c_{sk+j}^\ell = A^\ell q_{sk+j}$

$$\omega_{sk+j} = \frac{\langle c_{sk+j+1}^1, Gc_{sk+j+1}^0 \rangle}{\langle c_{sk+j+1}^1, Gc_{sk+j+1}^1 \rangle}$$

$$x_{sk+j+1} = x_{sk+j} + \alpha_{sk+j} p_{sk+j} + \omega_{sk+j} q_{sk+j}$$

$$r_{sk+j+1} = q_{sk+j} - \omega_{sk+j} [P, Q, R] c_{sk+j+1}^1$$

For $\ell = 0$ to $s - 2j$, Do

$$d_{sk+j+1}^\ell = c_{sk+j+1}^\ell - \omega_{sk+j} c_{sk+j+1}^{\ell+1}$$

//such that $[P, Q, R] d_{sk+j+1}^\ell = A^\ell r_{sk+j+1}$

$$\beta_{sk+j} = \frac{\langle g, d_{sk+j+1}^0 \rangle}{\langle g, d_{sk+j}^0 \rangle} \times \frac{\alpha}{\omega}$$

$$p_{sk+j+1} = r_{sk+j+1} + \beta_{sk+j} p_{sk+j} - \beta_{sk+j} \omega_{sk+j} [P, Q, R] b_{sk+j}^1$$

For $\ell = 0$ to $s - 2j$, Do

$$b_{sk+j+1}^\ell = d_{sk+j+1}^\ell + \beta_{sk+j} b_{sk+j}^\ell - \beta_{sk+j} \omega_{sk+j} b_{sk+j}^{\ell+1}$$

//such that $[P, Q, R] b_{sk+j+1}^\ell = A^\ell p_{sk+j+1}$.

EndDo

EndDo

Conclusions

- Communication is much more expensive than arithmetic, so avoid it, even if you need to do (a little) more arithmetic!
- For more details, see links on the class webpage under Reference for Communication-Avoiding Algorithms