

Block Cyclic Distribution of Data in pbdR and its Effects on Computational Efficiency

UMBC REU Site: Interdisciplinary Program in High Performance Computing

Team members: Matthew Bachmann¹, Ashley Dyas², Shelby Kilmer³, and Julian Sass⁴

Graduate assistant: Andrew Raim⁴, Faculty mentor: Nagaraj K. Neerchal⁴, Clients: George Ostrouchov⁵ and Ian F. Thorpe⁶

¹Northeast Lakeview College, ²Contra Costa College, ³Bucknell University, ⁴Math. and Stat., UMBC, ⁵Oak Ridge Nat. Lab., ⁶Chemistry, UMBC

Summary

- pbdR is an R package used to implement high-performance statistical computing on very large data sets.
- Block cyclic arrangement is used in pbdR (through ScaLAPACK) for distributed operations on large, dense matrices amongst parallel processes.
- Selection of block size and grid processor layout greatly influence computational efficiency.

In this work, we explore block cyclic distribution by implementing the statistical method PCA. We illustrate a large-scale PCA with an application to protein movements, focusing on motional correlations.

Methods

PCA

Principal component analysis (PCA) determines the directions of maximal variability in large dimensional data and is widely used for dimensional reduction. PCA is implemented on the sample covariance matrix:

$$\mathbf{S} = \frac{1}{n-1} \mathbf{Y}' \left(\mathbf{I} - \frac{1}{n} \mathbf{J} \right) \mathbf{Y}, \quad \text{where } \mathbf{J} = \mathbf{1}\mathbf{1}' \text{ and } \mathbf{Y} = \begin{bmatrix} y_{11} & \cdots & y_{1k} \\ \vdots & \ddots & \vdots \\ y_{n1} & \cdots & y_{nk} \end{bmatrix}$$

\mathbf{S} is positive semi-definite and therefore can be diagonalized: $\mathbf{\Lambda} = \mathbf{O}'\mathbf{S}\mathbf{O}$, where \mathbf{O} is an orthogonal matrix. Each row of \mathbf{O} is an eigenvector of \mathbf{S} . These eigenvectors are called principal components. $\mathbf{\Lambda}$ is the diagonal matrix of eigenvalues of \mathbf{S} . The number of significant eigenvalues determines the true underlying dimensionality of the data.

We analyze the effects of several factors of the PCA algorithm on computational speed:

- Vary the size of the $n \times k$ data matrix
- Vary block size and grid layout on a single node for constant n and k
- Extend study of block size and grid layout to multiple nodes
- Demonstrate real-world applicability

Background

Block Cyclic Distribution

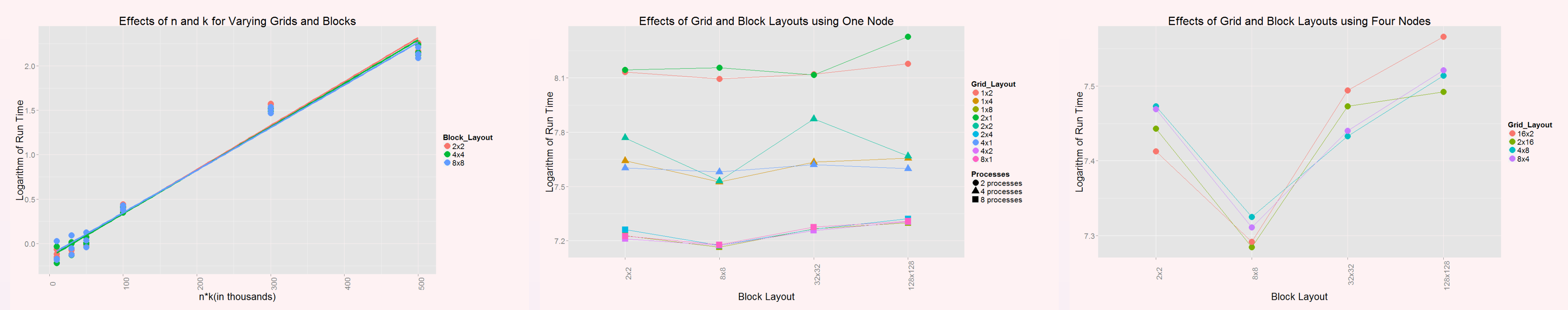
$$\mathbf{x} = \begin{bmatrix} x_{11} & x_{12} & x_{13} & x_{14} & x_{15} & x_{16} & x_{17} & x_{18} & x_{19} \\ x_{21} & x_{22} & x_{23} & x_{24} & x_{25} & x_{26} & x_{27} & x_{28} & x_{29} \\ x_{31} & x_{32} & x_{33} & x_{34} & x_{35} & x_{36} & x_{37} & x_{38} & x_{39} \\ x_{41} & x_{42} & x_{43} & x_{44} & x_{45} & x_{46} & x_{47} & x_{48} & x_{49} \\ x_{51} & x_{52} & x_{53} & x_{54} & x_{55} & x_{56} & x_{57} & x_{58} & x_{59} \\ x_{61} & x_{62} & x_{63} & x_{64} & x_{65} & x_{66} & x_{67} & x_{68} & x_{69} \\ x_{71} & x_{72} & x_{73} & x_{74} & x_{75} & x_{76} & x_{77} & x_{78} & x_{79} \\ x_{81} & x_{82} & x_{83} & x_{84} & x_{85} & x_{86} & x_{87} & x_{88} & x_{89} \\ x_{91} & x_{92} & x_{93} & x_{94} & x_{95} & x_{96} & x_{97} & x_{98} & x_{99} \end{bmatrix}$$

$$\text{Processor Grid Layout} = \begin{bmatrix} 0 & 1 & 2 \\ 3 & 4 & 5 \end{bmatrix}$$

Here a 9×9 matrix is block cyclically divided with block size of 2×2 and processor grid layout of 2×3 . Each colored 2×2 partition will be distributed to the process in the 2×3 grid with the corresponding color.

SNP and its Applications SNPs are variations that occur in a DNA sequence when a nucleotide differs between paired chromosomes, changing the amino acid sequence. These can be analyzed through motional correlations.

Results

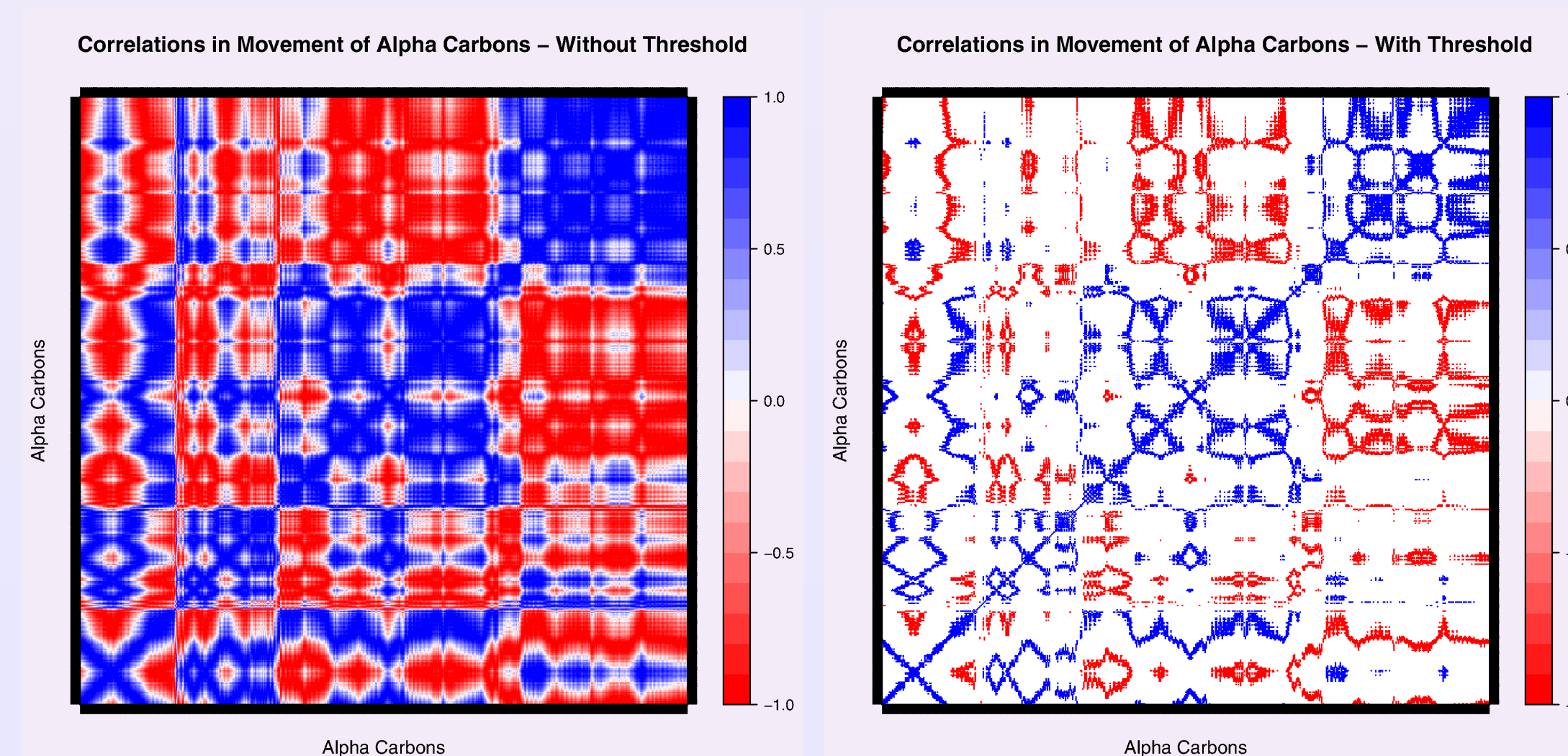


The left figure shows the predictable nature of the effects n and k have on efficiency. The center and right figures show the effects of grid layout and block size on computational speed for a constant n and k .

ANOVA Results:

Source	DF	One Node Study			Four Node Study		
		SS	MS	F	SS	MS	F
Grid	3	0.0008	0.0003	3.0000	0.0007	0.0002	0.2593
Block	3	0.0393	0.0131	131.0000	0.1040	0.0347	38.5556
Residuals	9	0.0011	0.0001		0.0080	0.0009	

These tables show ANOVA results for one node 8 processes on the left, and for four nodes 8 processes on the right. Results show that block size has a clear effect on computational speed. 8×8 block sizes are the most efficient for this n and k .



These figures show the level plot of Dr. Thorpe's data before and after PCA and greying out all statistically non-significant correlations. Note that relatively few correlations are statistically significant.

References and Acknowledgments

References at www.umbc.edu/hpcf > Publications:

- Full Technical Report: HPCF-2013-11
- A. M. Raim, pbdR Tutorial, HPCF-2013-2.

Acknowledgments:

- REU Site: www.umbc.edu/hpcreu
- NSF, NSA, UMBC, HPCF, CIRC