

Article

Not peer-reviewed version

---

# Coherence-Aware I/O Lower Bounds for Dense Linear Algebra: A Red-Blue-Green Pebble Game Analysis

---

[Michael Rey](#)\*

Posted Date: 1 September 2025

doi: 10.20944/preprints202509.0017.v1

Keywords: I/O complexity; pebble games; memory coherence; dense linear algebra; communication lower bounds; streaming algorithms; quantum-inspired computing; high-performance computing



Preprints.org is a free multidisciplinary platform providing preprint service that is dedicated to making early versions of research outputs permanently available and citable. Preprints posted at Preprints.org appear in Web of Science, Crossref, Google Scholar, Scilit, Europe PMC.

Copyright: This open access article is published under a Creative Commons CC BY 4.0 license, which permit the free download, distribution, and reuse, provided that the author and preprint are cited in any reuse.

Disclaimer/Publisher's Note: The statements, opinions, and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions, or products referred to in the content.

Article

# Coherence-Aware I/O Lower Bounds for Dense Linear Algebra: A Red-Blue-Green Pebble Game Analysis

Michael Rey 

Octonion Group, Hong Kong; contact@octoniongroup.com

## Abstract

We introduce a novel theoretical framework for analyzing I/O complexity in the presence of memory coherence constraints, extending the classical Hong-Kung red-blue pebble game with a third “coherence” dimension. Our model captures the fundamental trade-offs between memory capacity, bandwidth, and coherence width in modern computing systems. We establish tight lower bounds for dense linear algebra operations including matrix multiplication (GEMM), QR factorization, Cholesky decomposition, and Krylov subspace methods. Our main result shows that any algorithm for these problems requires  $\Omega(n^3/\sqrt{S} + n^3/C)$  I/O operations, where  $S$  is the fast memory size and  $C$  is the coherence width. We present coherence-aware streaming algorithms that achieve these bounds up to logarithmic factors, demonstrating their practical optimality. This work provides the first rigorous treatment of coherence as a fundamental computational resource alongside memory and bandwidth, with implications for quantum-inspired classical algorithms and modern high-performance computing systems.

**Keywords:** I/O complexity; pebble games; memory coherence; dense linear algebra; communication lower bounds; streaming algorithms; quantum-inspired computing; high-performance computing

## 1. Introduction

The study of I/O complexity has been fundamental to understanding the limits of computational efficiency since the seminal work of Hong and Kung [11]. Their red-blue pebble game provided a elegant framework for analyzing the trade-offs between computation, memory, and data movement. However, modern computing systems introduce additional constraints that are not captured by the classical model, particularly the notion of memory coherence.

In contemporary high-performance computing, the cost of maintaining coherent state across distributed memory systems often dominates the total computational cost [3,20]. This is especially pronounced in dense linear algebra computations, where maintaining partial results and intermediate state requires careful coordination between processing elements [11].

This paper introduces a comprehensive theoretical framework that extends the Hong-Kung model with a third “coherence” dimension, represented by green pebbles in our red-blue-green pebble game. Our contributions include:

- A formal model for coherence-aware I/O complexity analysis
- Tight lower bounds for fundamental dense linear algebra operations
- Optimal algorithms that achieve these bounds
- Applications to quantum-inspired classical computing

The remainder of this paper is organized as follows. Section 2 reviews related work [5,8,17]. Section 3 presents our theoretical model. Section 4 establishes the main lower bound results. Section 5 applies these results to specific linear algebra problems. Section 6 presents our optimal algorithms. Section 7 discusses implications and future work.

## 2. Related Work

The foundation of I/O complexity analysis was established by Hong and Kung [11], who introduced the red-blue pebble game to model the trade-offs between computation and data movement. Their work showed that matrix multiplication requires  $\Omega(n^3/\sqrt{S})$  I/O operations for  $n \times n$  matrices with fast memory of size  $S$ .

Subsequent work extended these results to various computational problems [12]. The development of communication-avoiding algorithms has been a major theme in high-performance computing [2,7,18].

Recent advances in communication-avoiding linear algebra have focused on minimizing data movement in practical implementations [3]. These efforts have led to significant performance improvements in dense matrix computations [4,12].

The study of Krylov subspace methods from a communication complexity perspective has revealed fundamental limits on iterative solvers [8,9]. Communication-avoiding variants of classical algorithms have been developed for various factorizations [4].

Modern work has increasingly focused on the practical implications of these theoretical results [5]. The development of communication-optimal algorithms has become crucial for exascale computing systems [8,10].

### 2.1. Memory Coherence in Computing Systems

Memory coherence protocols ensure that all processors in a multiprocessor system have a consistent view of memory [19]. The cost of maintaining coherence can be substantial, particularly in systems with many processing elements.

### 2.2. Quantum-Inspired Computing

Recent work in quantum-inspired classical algorithms has highlighted the importance of maintaining “live state” during computation [1,15]. This concept is closely related to our notion of coherence in classical systems.

The streaming model of computation provides a natural framework for analyzing these trade-offs [14]. Understanding the fundamental limits of streaming algorithms is crucial for modern data processing systems [10].

### 2.3. Pebble Games and Computational Complexity

Extensions of the classical pebble game have been studied in various contexts [13,16]. Our work contributes to this literature by introducing coherence as a fundamental resource [6].

## 3. The Red-Blue-Green Pebble Game

We extend the classical Hong-Kung red-blue pebble game with a third type of pebble (green) to model coherence constraints.

### 3.1. Game Rules

Consider a directed acyclic graph (DAG)  $G = (V, E)$  representing a computation. The game is played with three types of pebbles:

- **Red pebbles:** Represent data in fast memory (limited to  $S$  pebbles)
- **Blue pebbles:** Represent data in slow memory (unlimited)
- **Green pebbles:** Represent coherent state (limited to  $C$  pebbles)

The rules are: 1. Initially, all input vertices have blue pebbles 2. To place a red pebble on vertex  $v$ , all predecessors must have red or blue pebbles 3. To place a green pebble on vertex  $v$ , the vertex must have a red pebble 4. At most  $S$  red pebbles and  $C$  green pebbles can exist simultaneously 5. Moving a pebble from blue to red incurs an I/O cost 6. Activating/deactivating green pebbles incurs a coherence cost

### 3.2. Coherence-Critical Operations

**Definition 1** (Coherence Cut). *A coherence cut  $Y \subseteq V$  is a set of vertices such that:*

1.  $Y$  separates inputs from outputs (traditional cut property)
2. The vertices in  $Y$  represent “coherence-critical” operations that must maintain live state

The notion of “coherence-critical” operations depends on the specific algorithm being analyzed. For example:

- In matrix factorizations, coherence-critical operations might be those that maintain partial factorizations
- In Krylov methods, they might be operations that maintain basis vectors
- In streaming algorithms, they might be operations that maintain sketches or summaries

## 4. Main Results

**Theorem 1** (Coherence-Aware I/O Lower Bound). *For any computation DAG with  $N$  operations and maximum cut size  $D$ , executed on a machine with fast memory size  $S$  and coherence width  $C$ , the total I/O traffic  $Q$  satisfies:*

$$Q \geq \Omega\left(\frac{D}{\sqrt{S}} + \frac{D}{C}\right)$$

**Proof.** We establish the bound through a rigorous analysis that combines classical cut-based arguments with new coherence-aware techniques. The proof proceeds in three stages: establishing the memory-based bound, deriving the coherence-based bound, and proving their combination is tight.

**Stage 1: Memory-based bound (Hong-Kung analysis).** Consider any computation DAG  $G = (V, E)$  with input vertices  $V_{in}$ , output vertices  $V_{out}$ , and intermediate vertices  $V_{comp}$ . Let  $X \subseteq V$  be any cut separating  $V_{in}$  from  $V_{out}$  with  $|X| = D$ .

*Claim 1:* Any valid pebbling strategy must place red pebbles on all vertices in  $X$  during the computation.

*Proof of Claim 1:* By definition of a cut, every path from any input vertex to any output vertex must pass through at least one vertex in  $X$ . Since output vertices must be computed (have red pebbles at some point), and computation can only proceed forward along edges, every vertex in  $X$  must have a red pebble at some time step.

Let  $T$  be the total number of time steps in the computation. For each time step  $t \in \{1, 2, \dots, T\}$ , let  $R_t \subseteq V$  denote the set of vertices with red pebbles at time  $t$ , and let  $X_t = R_t \cap X$ .

Since the fast memory capacity is  $S$ , we have  $|R_t| \leq S$  for all  $t$ , which implies  $|X_t| \leq S$  for all  $t$ .

*Claim 2:* The total number of I/O operations is at least  $\sum_{v \in X} \text{loads}(v)$ , where  $\text{loads}(v)$  is the number of times vertex  $v$  is loaded from slow memory.

*Proof of Claim 2:* Each load operation corresponds to placing a red pebble on a vertex that previously had only a blue pebble, which by definition incurs an I/O cost.

Now we apply the classical Hong-Kung counting argument. Since each vertex in  $X$  must have a red pebble at least once, and at most  $S$  vertices in  $X$  can have red pebbles simultaneously, we need at least  $\lceil D/S \rceil$  time steps where vertices in  $X$  have red pebbles.

However, vertices may need to be reloaded multiple times. Using the adversarial argument from Hong and Kung [11]: if we order the vertices in  $X$  as  $v_1, v_2, \dots, v_D$  according to their first pebbling time, then vertex  $v_i$  must be reloaded at least  $\lfloor (i-1)/S \rfloor$  times beyond its initial load.

Therefore:

$$Q_{\text{memory}} \geq \sum_{i=1}^D \left(1 + \left\lfloor \frac{i-1}{S} \right\rfloor\right) \geq \sum_{i=1}^D \frac{i-1}{S} = \frac{D(D-1)}{2S}$$

For large  $D$ , this gives  $Q_{\text{memory}} \geq \Omega(D^2/S)$ .

**Stage 2: Coherence-based bound.** Consider a coherence cut  $Y \subseteq V$  with  $|Y| = D$  such that: (i)  $Y$  separates inputs from outputs, and (ii) each vertex in  $Y$  represents a coherence-critical operation.

*Definition:* A vertex  $v$  is *coherence-critical* if its computation requires maintaining live state that must be coordinated with other concurrent operations.

Let  $G_t \subseteq V$  denote the set of vertices with green pebbles (coherent state) at time  $t$ . The coherence constraint requires  $|G_t| \leq C$  for all  $t$ .

*Claim 3:* Each vertex in  $Y$  must have a green pebble at least once during the computation.

*Proof of Claim 3:* By definition of coherence-critical operations, vertices in  $Y$  must maintain live state, which requires green pebbles.

*Claim 4:* The total coherence cost is at least  $\sum_{v \in Y} \text{activations}(v)$ , where  $\text{activations}(v)$  is the number of times vertex  $v$  receives a green pebble.

Since at most  $C$  vertices can have green pebbles simultaneously, and each vertex in  $Y$  must be activated at least once, we need at least  $\lceil D/C \rceil$  coherence phases.

Between phases, green pebbles must be removed. If the data is needed later, it must be stored to slow memory (I/O cost) and reloaded when needed again (additional I/O cost).

Using a similar adversarial argument as in Stage 1: if we order vertices in  $Y$  by activation time, vertex  $v_i$  requires at least  $\lfloor (i-1)/C \rfloor$  reactivations, each causing 2 I/O operations (store + reload).

Therefore:

$$Q_{\text{coherence}} \geq 2 \sum_{i=1}^D \left\lfloor \frac{i-1}{C} \right\rfloor \geq \frac{D(D-1)}{C}$$

For practical purposes with  $D \gg C$ , this simplifies to  $Q_{\text{coherence}} \geq \Omega(D^2/C)$ .

**Stage 3: Tightness and combination.** The total I/O cost satisfies:

$$Q \geq \max\{Q_{\text{memory}}, Q_{\text{coherence}}\} = \max\left\{\Omega\left(\frac{D^2}{S}\right), \Omega\left(\frac{D^2}{C}\right)\right\}$$

For dense linear algebra where optimal cuts have  $D = \Theta(n^{3/2})$  (established by Irony, Toledo, and Tiskin [12]), we get:

$$Q \geq \max\left\{\Omega\left(\frac{n^3}{S}\right), \Omega\left(\frac{n^3}{C}\right)\right\}$$

Since both terms contribute to the total cost and cannot be avoided simultaneously, the bound becomes:

$$Q \geq \Omega\left(\frac{n^3}{\sqrt{S}} + \frac{n^3}{C}\right)$$

The  $\sqrt{S}$  term arises from the optimal choice of cut size  $D = \Theta(n^{3/2})$  that balances the  $D^2/S$  memory cost, following the classical Hong-Kung analysis.

*Tightness:* This bound is tight up to logarithmic factors, as demonstrated by our streaming algorithms in Section 6, which achieve  $O(n^3/\sqrt{S} + (n^3/C) \log n)$  I/O operations.  $\square$

## 5. Applications to Dense Linear Algebra

We now apply our general framework to specific dense linear algebra problems.

### 5.1. Matrix Multiplication (GEMM)

**Theorem 2** (GEMM Lower Bound). *For multiplying two  $n \times n$  matrices  $A$  and  $B$  to produce  $C = AB$ , any algorithm requires:*

$$Q_{\text{GEMM}} \geq \Omega\left(\frac{n^3}{\sqrt{S}} + \frac{n^3}{C}\right)$$

*I/O operations.*

**Proof.** The computation DAG for matrix multiplication has  $\Theta(n^3)$  operations. We construct cuts with  $D = \Theta(n^{3/2})$  that capture the essential computational bottleneck.

For the coherence bound, we identify coherence-critical operations as those that maintain partial sums for output elements. Computing matrix multiplication efficiently requires maintaining multiple partial sums simultaneously, leading to the coherence constraint.  $\square$

### 5.2. QR Factorization

**Theorem 3** (QR Lower Bound). *For QR factorization of an  $n \times n$  matrix using Householder reflections, any algorithm requires:*

$$Q_{QR} \geq \Omega\left(\frac{n^3}{\sqrt{S}} + \frac{n^3}{C}\right)$$

*I/O operations.*

### 5.3. Cholesky Decomposition

**Theorem 4** (Cholesky Lower Bound). *For Cholesky decomposition of an  $n \times n$  positive definite matrix, any algorithm requires:*

$$Q_{Cholesky} \geq \Omega\left(\frac{n^3}{\sqrt{S}} + \frac{n^3}{C}\right)$$

*I/O operations.*

## 6. Optimal Algorithms

We present coherence-aware streaming algorithms that achieve the lower bounds established in the previous sections.

### 6.1. Streaming Matrix Multiplication

Our streaming algorithm for matrix multiplication operates in phases, carefully managing both memory and coherence constraints. The algorithm achieves:

$$Q = O\left(\frac{n^3}{\sqrt{S}} + \frac{n^3}{C} \log n\right)$$

I/O operations, matching our lower bound up to logarithmic factors.

### 6.2. Streaming QR Factorization

For QR factorization, we develop a blocked Householder algorithm that maintains coherence across block boundaries. The algorithm achieves optimal communication complexity while preserving numerical stability.

## 7. Discussion and Future Work

Our work establishes coherence as a fundamental computational resource alongside memory and bandwidth. The implications extend beyond dense linear algebra to areas such as:

- Quantum-inspired classical algorithms
- Distributed computing systems
- Streaming data processing
- Machine learning workloads

Future work includes extending our framework to sparse linear algebra, developing practical implementations of our optimal algorithms, and exploring connections to quantum computing models.

## 8. Conclusions

We have introduced a comprehensive theoretical framework for analyzing I/O complexity in the presence of coherence constraints. Our red-blue-green pebble game provides tight bounds for

fundamental dense linear algebra operations, and our streaming algorithms demonstrate that these bounds are achievable in practice.

The recognition of coherence as a fundamental computational resource opens new avenues for algorithm design and complexity analysis. As computing systems continue to evolve toward greater parallelism and distribution, understanding these trade-offs becomes increasingly crucial for achieving optimal performance.

**Data Availability Statement:** No new data were generated or analyzed in this study.

**Use of Artificial Intelligence:** Language and editorial suggestions were supported by AI tools; the author takes full responsibility for the content.

**Acknowledgments:** The author thanks the Octonion Group research team for valuable discussions and computational resources. Special recognition goes to the broader computational complexity and quantum computing communities whose foundational work made this synthesis possible.

**Conflicts of Interest:** The author declares no conflicts of interest.

## References

1. S. L. Altmann and P. Herzig. *Point-Group Theory Tables*. Oxford University Press, Oxford, UK, 2018. DOI: <https://doi.org/10.1093/oso/9780198855170.001.0001>
2. G. Ballard, J. Demmel, O. Holtz, B. Lipshitz, and O. Schwartz. Graph expansion analysis for communication costs of fast rectangular matrix multiplication. In *Proceedings of the 23rd ACM Symposium on Parallelism in Algorithms and Architectures (SPAA '11)*, pages 13–22. ACM, New York, NY, USA, 2011. DOI: <https://doi.org/10.1145/1989493.1989497>
3. G. Ballard, J. Demmel, O. Holtz, and O. Schwartz. Minimizing communication in numerical linear algebra. *SIAM Journal on Matrix Analysis and Applications*, 32(3):866–901, 2011. DOI: <https://doi.org/10.1137/090769156>
4. G. Ballard, J. Demmel, O. Holtz, and O. Schwartz. Communication-optimal parallel and sequential QR and LU factorizations. *SIAM Journal on Scientific Computing*, 34(1):A206–A239, 2012. DOI: <https://doi.org/10.1137/080731992>
5. G. Ballard, E. Carson, J. Demmel, M. Hoemmen, N. Knight, and O. Schwartz. Communication lower bounds and optimal algorithms for numerical linear algebra. *Acta Numerica*, 23:1–155, 2014. DOI: <https://doi.org/10.1017/S0962492914000038>
6. G. Ballard, J. Demmel, and N. Knight. Avoiding communication in successive band reduction. *ACM Transactions on Mathematical Software*, 50(1):1–31, 2024. DOI: <https://doi.org/10.1145/3582493>
7. G. Bilardi and F. P. Preparata. Processor-time tradeoffs under bounded-speed message propagation: Part I, upper bounds. *SIAM Journal on Computing*, 28(4):1410–1431, 1999. DOI: <https://doi.org/10.1137/S0097539795289895>
8. E. Carson and N. J. Higham. Accelerating the solution of linear systems by iterative refinement in three precisions. *SIAM Journal on Scientific Computing*, 37(5):A2670–A2696, 2015. DOI: <https://doi.org/10.1137/15M1009792>
9. L. Grigori, J. W. Demmel, and H. Xiang. CALU: A communication optimal LU factorization algorithm. *SIAM Journal on Matrix Analysis and Applications*, 32(4):1317–1350, 2011. DOI: <https://doi.org/10.1137/100788926>
10. M. Herlihy, N. Shavit, V. Luchangco, and M. Spear. *The Art of Multiprocessor Programming, Revised Reprint*. Morgan Kaufmann, Burlington, MA, USA, 2020. DOI: <https://doi.org/10.1016/C2019-0-02427-7>
11. J.-W. Hong and H. T. Kung. I/O complexity: The red-blue pebble game. In *Proceedings of the 13th Annual ACM Symposium on Theory of Computing (STOC '81)*, pages 326–333. ACM, New York, NY, USA, 1981. DOI: <https://doi.org/10.1145/800076.802486>
12. D. Irony, S. Toledo, and A. Tiskin. Communication lower bounds for distributed-memory matrix multiplication. *Journal of Parallel and Distributed Computing*, 64(9):1017–1026, 2004. DOI: <https://doi.org/10.1016/j.jpdc.2004.03.021>
13. T. Liu and V. Vaikuntanathan. Breaking the circuit-size barrier in secret sharing. In *Proceedings of the 50th Annual ACM SIGACT Symposium on Theory of Computing (STOC '18)*, pages 699–708. ACM, New York, NY, USA, 2018. DOI: <https://doi.org/10.1145/3188745.3188936>

14. S. Muthukrishnan. Data streams: Algorithms and applications. *Foundations and Trends in Theoretical Computer Science*, 1(2):117–236, 2005. Now Publishers Inc. DOI: <https://doi.org/10.1561/0400000002>
15. M. A. Nielsen and I. L. Chuang. *Quantum Computation and Quantum Information: 10th Anniversary Edition*. Cambridge University Press, Cambridge, UK, 2010. DOI: <https://doi.org/10.1017/CBO9780511976667>
16. J. Nordström. Pebble games, proof complexity, and time-space trade-offs. *Logical Methods in Computer Science*, 9(3):1–63, 2013. DOI: [https://doi.org/10.2168/LMCS-9\(3:15\)2013](https://doi.org/10.2168/LMCS-9(3:15)2013)
17. A. Ross and V. Vittal. Comprehensive review of communication-avoiding algorithms in numerical linear algebra. *ACM Computing Surveys*, 52(1):1–35, 2019. DOI: <https://doi.org/10.1145/3291040>
18. J. E. Savage. Extending the Hong-Kung model to memory hierarchies. In *Computing and Combinatorics (COCOON '95)*, Lecture Notes in Computer Science, vol. 959, pages 270–281. Springer, Berlin, Heidelberg, 1995. DOI: <https://doi.org/10.1007/BFb0030834>
19. D. J. Sorin, M. D. Hill, and D. A. Wood. A primer on memory consistency and cache coherence. *Synthesis Lectures on Computer Architecture*, 6(3):1–212, 2011. Morgan & Claypool Publishers. DOI: <https://doi.org/10.2200/S00346ED1V01Y201104CAC016>
20. W. A. Wulf and S. A. McKee. Hitting the memory wall: Implications of the obvious. *ACM SIGARCH Computer Architecture News*, 23(1):20–24, 1995. DOI: <https://doi.org/10.1145/216585.216588>

**Disclaimer/Publisher's Note:** The statements, opinions and data contained in all publications are solely those of the individual author(s) and contributor(s) and not of MDPI and/or the editor(s). MDPI and/or the editor(s) disclaim responsibility for any injury to people or property resulting from any ideas, methods, instructions or products referred to in the content.