

# Intel® Programmable Integrated Unified Memory Architecture (PIUMA)

*Hardware for Faster and Deeper Insights from Large Scale Graph*

Nikhil Deshpande Ph.D.

[Nikhil.m.Deshpande@intel.com](mailto:Nikhil.m.Deshpande@intel.com)

Product Director, AI and HPC Innovations

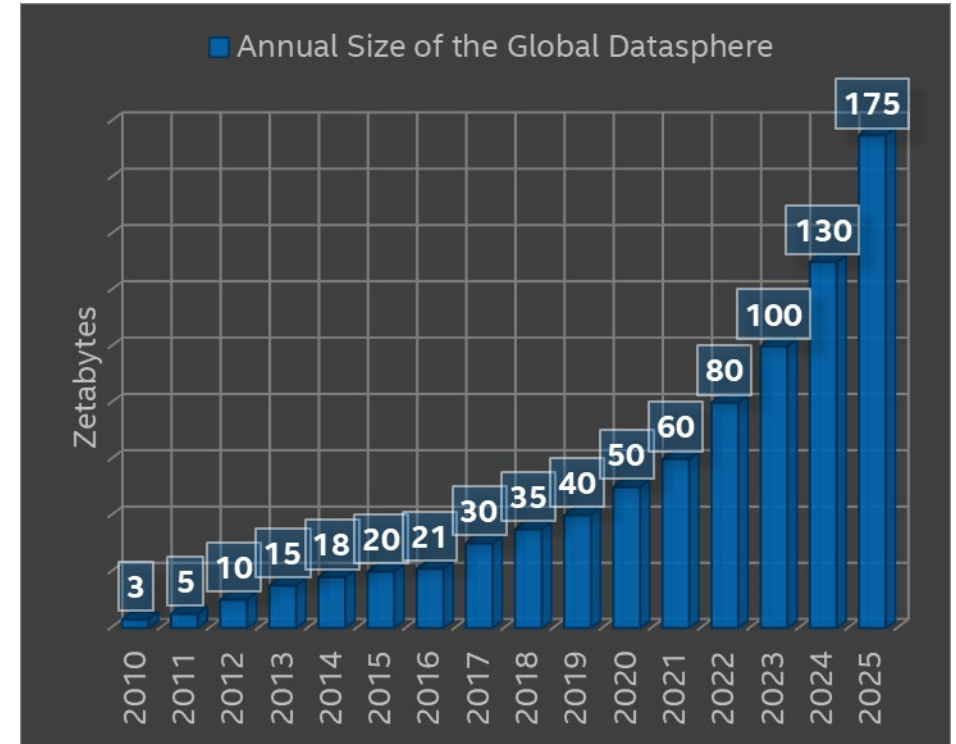
Intel Corporation

# Today

- Environment
- Graph Usages and Scalability Issues
- Graph and Traditional Compute
- Rethinking Hardware - Intel® PIUMA
- Programming Intel® PIUMA
- Summary

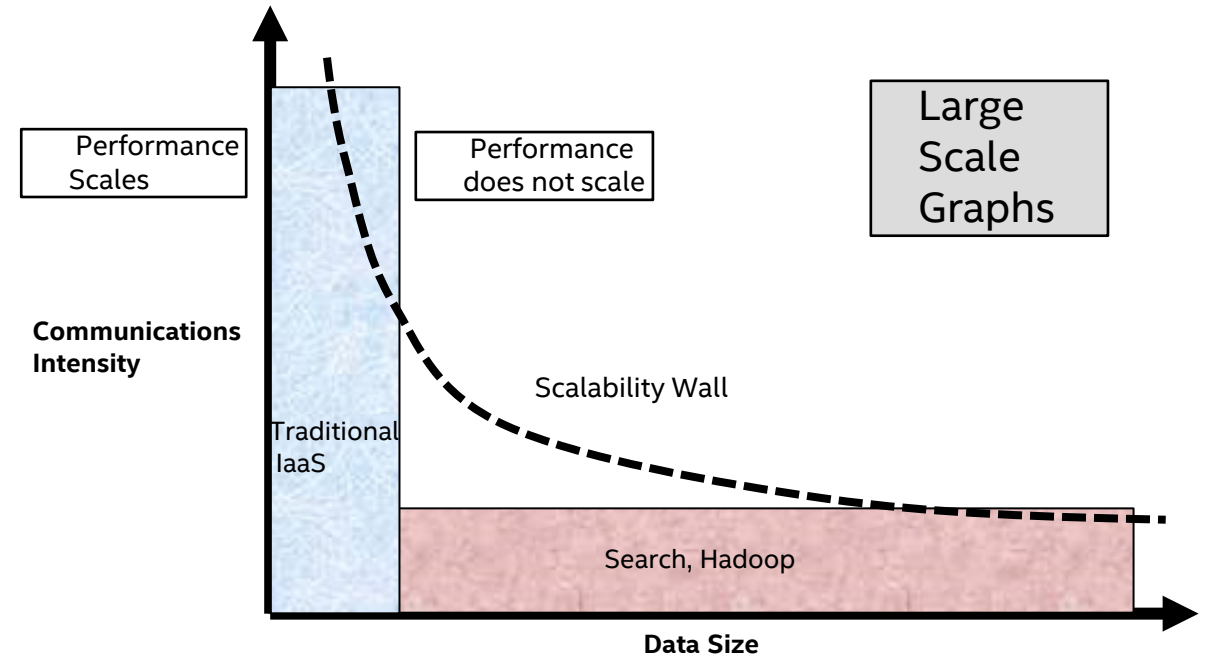
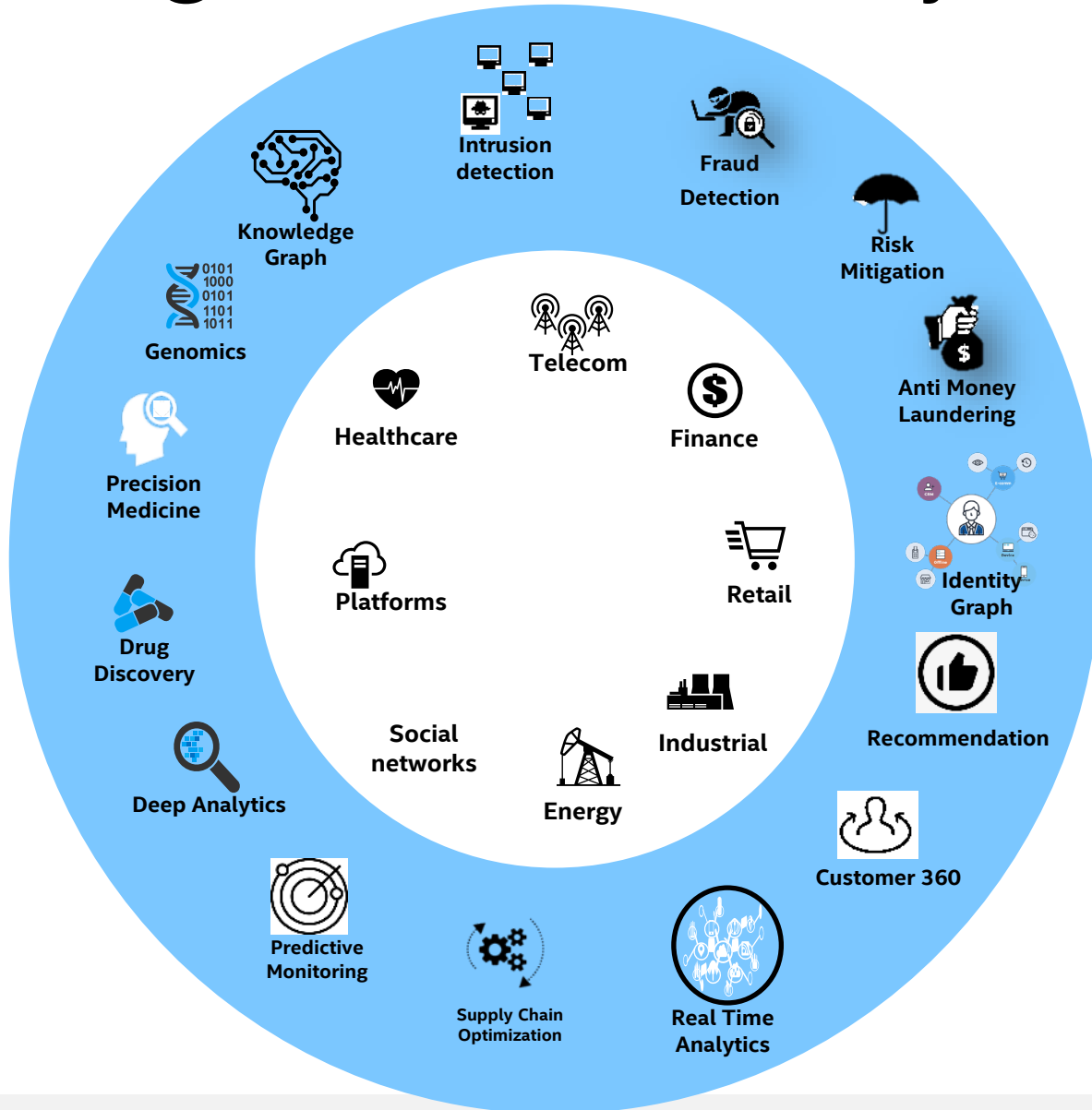
# Environment

- Data Deluge
- What are we after: Data, Knowledge or Insights?
- Is Data still the oil? Refining perhaps?
- Graphs: Best Representation for Capturing Relationships/Insights
- Faster and Deeper



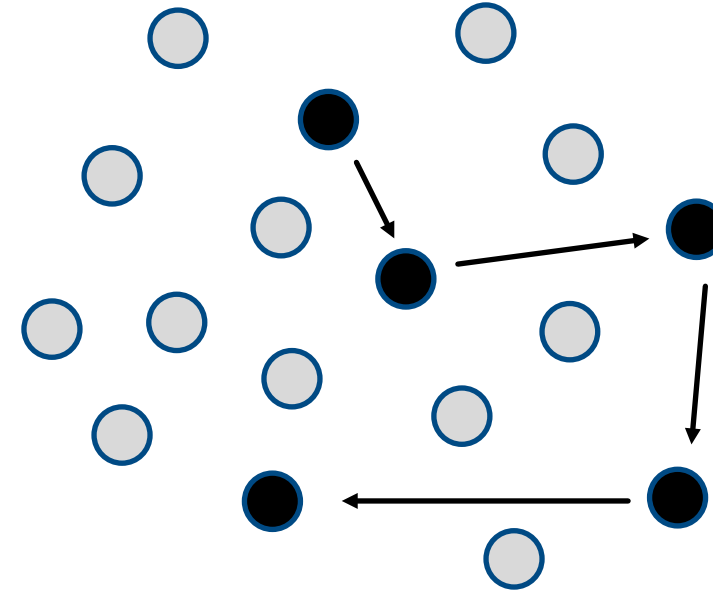
Source: Data Age 2025, Sponsored by Seagate with Data from IDC Global Datasphere, Nov. 2018

# Usages and Scalability Wall



# New Workload Behavior Needs New Thinking

	Regular	Graph
<b>Branch Prediction</b>	<i>Branches have regular pattern</i>	<i>Branch outcome is data dependent ("Pointer hopping")</i>
<b>Locality</b>	<i>Same or neighboring data likely used</i>	<i>Data is randomly scattered</i>
<b>Data Access</b>	<i>Same operation on neighboring data</i>	<i>Operations on scattered data</i>

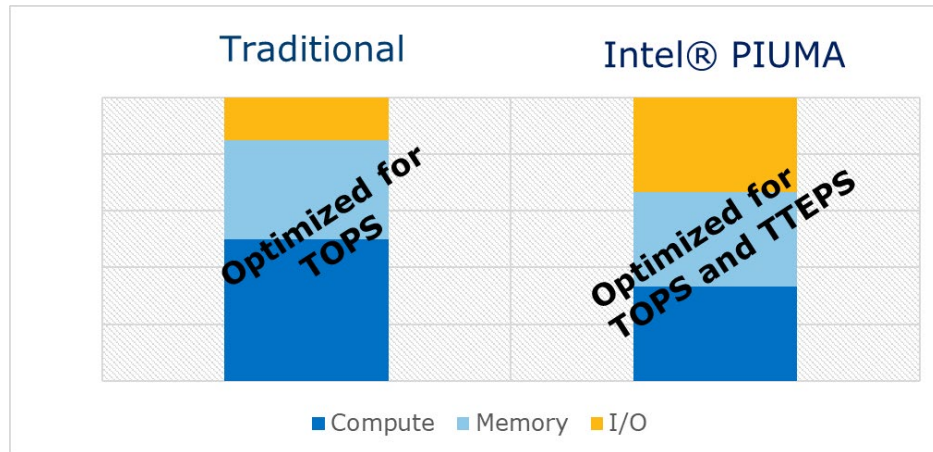


Graph Traversal is all about Pointer Chasing  
... with lowest latency!

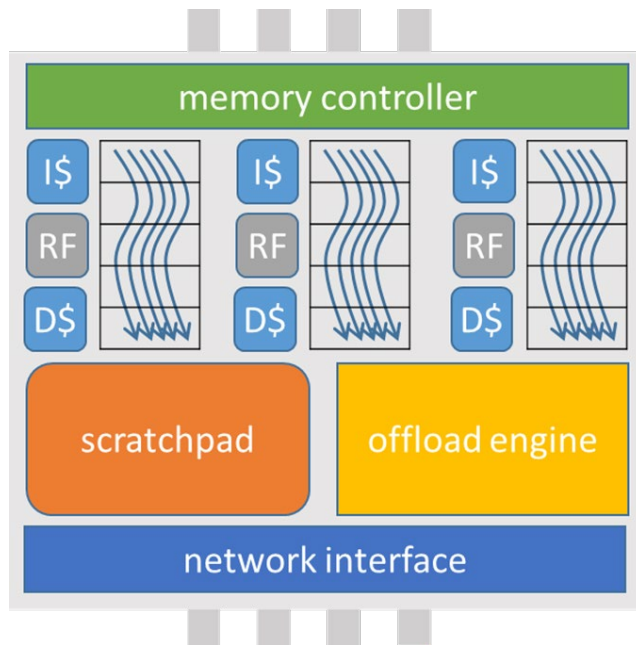
Today: Caching, Large message sizes, FLOPs...

Needs: Granular access, Small message sizes, +Traversed Edges Per Second!

# Intel® PIUMA Technology



Balance I/O, Memory, Compute and Prioritize in that order



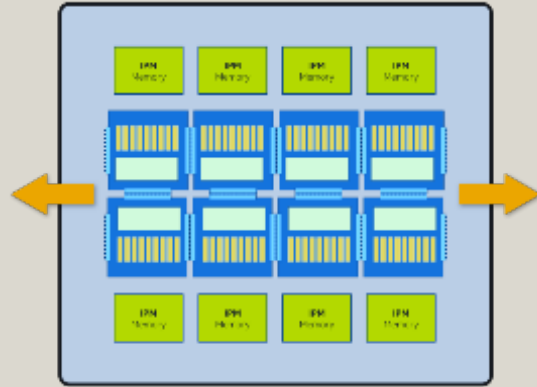
- 1 PIUMA core = multiple processor pipelines with multiple threads
- Novel 64-bit RISC instruction set, graph optimized instructions
- Global address space (GAS), accessible for all other cores
- Memory controller with 8B granularity
- Network interface with 8B packets
- Tiles → Nodes → Systems → Millions of threads...

TOPS: Tera-Operations Per Second  
TTEPS: Tera-Traversed Edges Per Second

# Intel's Approach: Intel® PIUMA

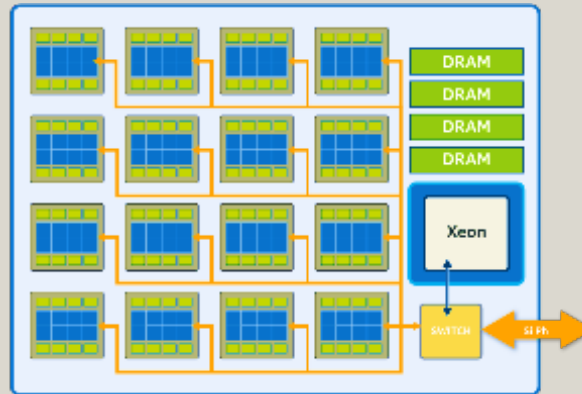
A Programmable Integrated Unified Memory Architecture

## Re-imagined Architecture



CPU Support for Small,  
Irregular Memory Accesses  
Near-Memory Atomics

## Fully Integrated



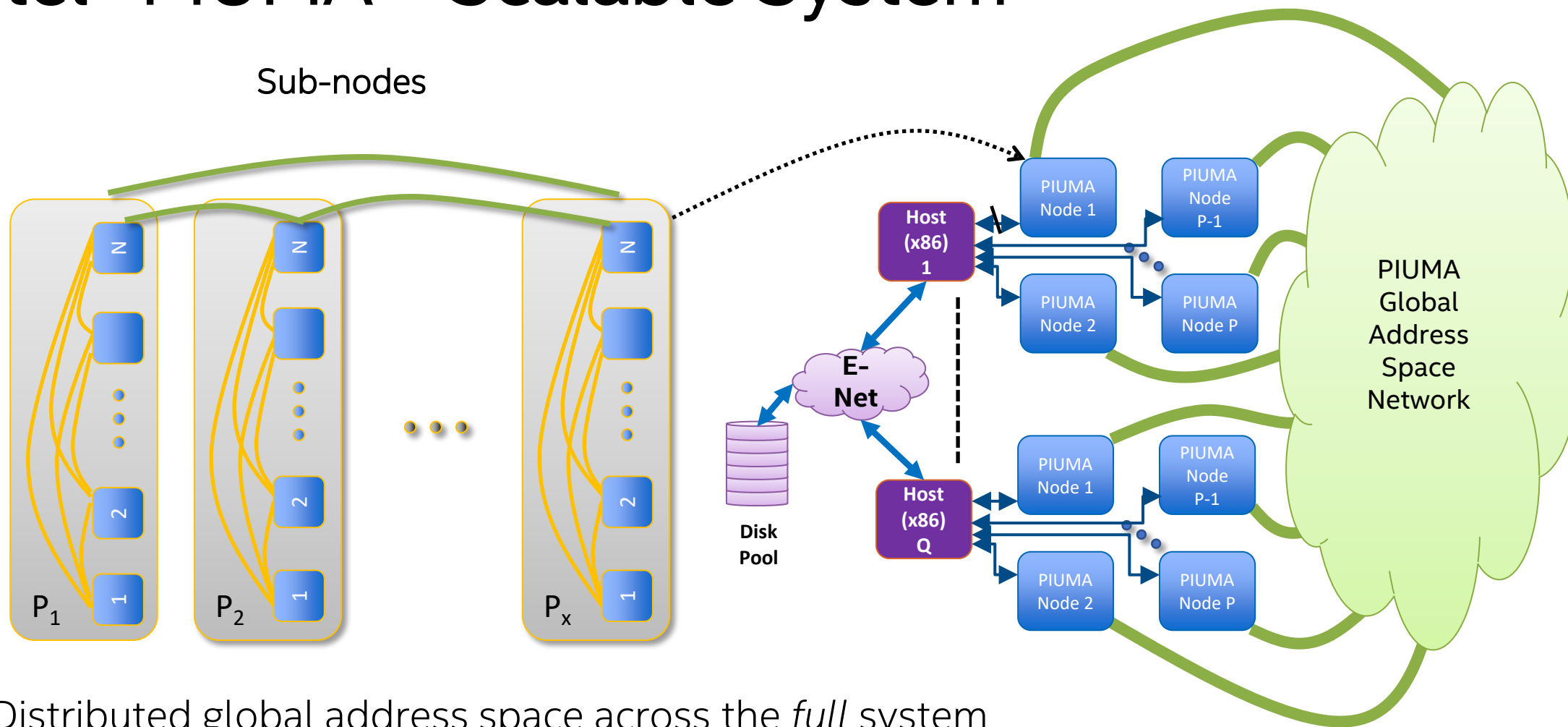
Global Memory Model  
Packaging for High IO &  
Memory bandwidth

## Architected to scale



Network as 1<sup>st</sup>-class Citizen  
Flatten Latency Hierarchy  
Point-to-Point Messages

# Intel® PIUMA – Scalable System



- Distributed global address space across the *full* system
- Hierarchical topology with all-to-all connections → low diameter
- High radix fabric between tiles and between nodes for high bandwidth and low latency



# Speedup: Intel® PIUMA versus 1 Intel® Xeon® node

Application	Intel® PIUMA 1 node	Intel® PIUMA 16 nodes
Application Classification	6.9 x	111 x
Random Walks	279 x	2,606 x
Graph Search	34 x	544 x
Louvain Community	41 x	555 x
TIES Sampler	93 x	419 x
Graph2Vec	42 x	178 x
GraphSAGE	3.1 x	46 x

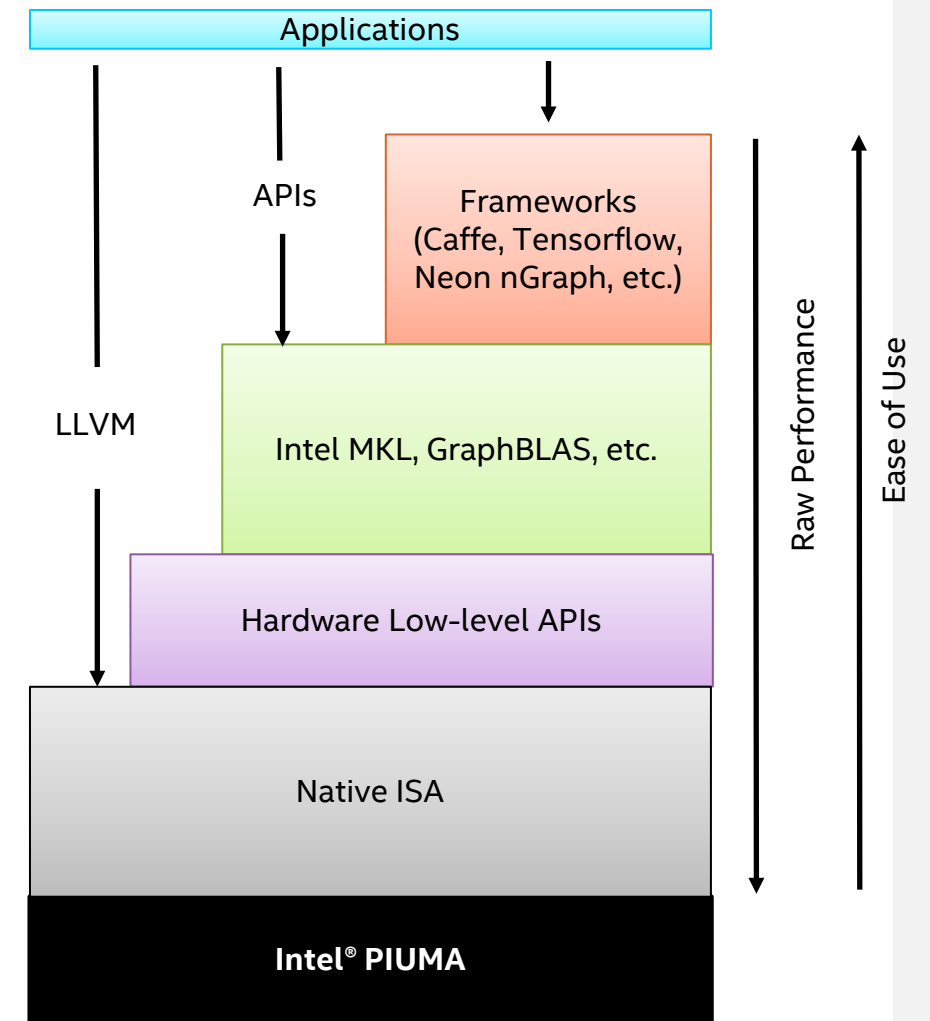
Application	Intel® PIUMA 1 node	Intel® PIUMA 16 nodes
Graph Wave	8.0 x	125 x
Parallel Decoding FST	6.8 x	109 x
Geolocation	15 x	243 x
SpMV	29 x	467 x
SpMSPV	111 x	1,387 x
Breadth-First Search	7.5 x	117 x

\*Results have been estimated or simulated.

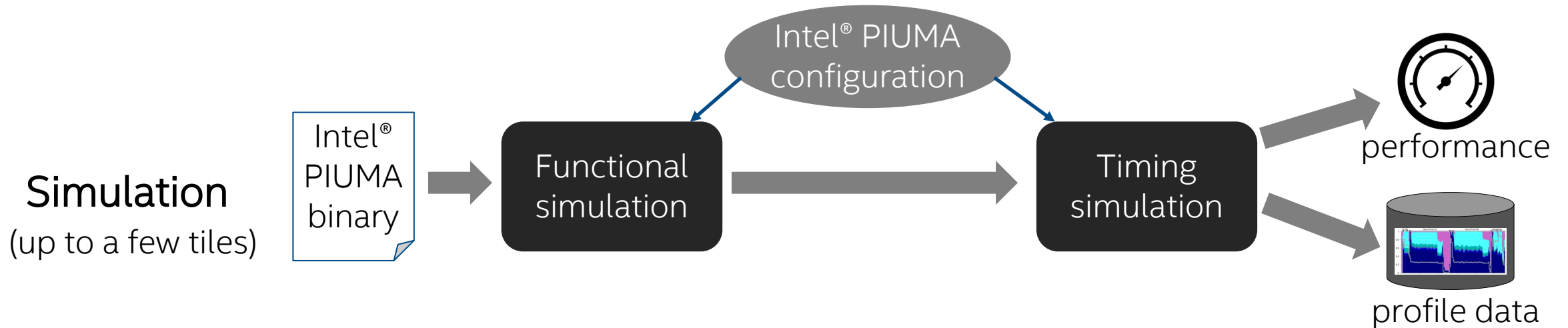
For more complete information about performance and benchmark results, visit [www.intel.com/benchmarks](http://www.intel.com/benchmarks).

# PIUMA Software Stack

- Performance, Performance, Performance!
- Currently support for C++, pthreads, OpenMP
- Under development: graph libraries (GraphBLAS), other programming languages (OpenCL), Python frontend
- Customer Choice and Options



# How to Engage?



- Learn how your workload kernels would perform on Intel® PIUMA
- Establish software readiness with your choice of productivity suite
- Be ready to take advantage of hardware!

Credit: Eyerman et.al. FODSEM 2020

# Summary

- Real-time Insights from large scale Data via Graph Analytics require new hardware thinking
- Intel® PIUMA is a programmable instruction set processor optimized for sparse graph applications
- Main features: high-bandwidth system-wide shared memory, small granularity memory accesses, massive multithreading, configurable caching, offload engines
- Programming Intel® PIUMA via C/C++, MPI, OpenMP plus other productivity suite
- Collaborate on workloads and datasets; Co-design with customer input. Join us.

**Thank you!**