

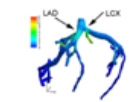




Using Computational Patterns to Understand Heterogeneity

David Sheffield and Kurt Keutzer
and the PALLAS team

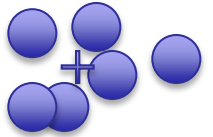
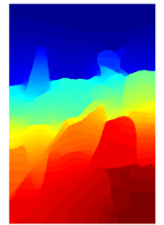
Michael Anderson, Bryan Catanzaro (→ Nvidia),
Katya Gonina, Chao-Yue Lai, Mark Murphy,
Bor-Yiing Su, Narayanan Sundaram

- Future generations of microprocessors are certain to have heterogeneous elements
- But determining the precise mix of heterogeneity is difficult:
 - Different microarchitectures to support the same ISA
 - Large and small cores
 - Significant coprocessor cores (eg, GPUs)
 - Special purpose execution units
 - ISA additions (eg Tensilica TIE)
 - Autonomous execution units (eg next-route lookup in network processors)
 - Reconfigurable logic
- Computational patterns give a new perspective on identifying heterogeneity to create new architectures

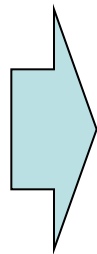
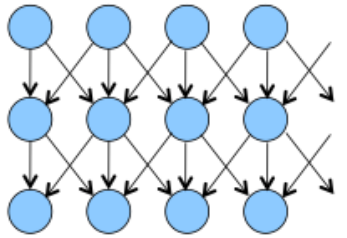
- We will review the computational patterns
- We will identify key program features to best support the computational patterns
- We will identify key micro-architectural elements that best support these program features
- Show implications of these micro-architectural elements on microprocessor-level heterogeneity
- Caveat: workload dependent: your mileage may vary

Apps Dwarves	Embed	SPEC	DB	Games	ML	HPC	CAD	 Health	 Image	 Speech	 Music	 Browser
	Graph Algorithms	Red	Yellow	Yellow	Yellow	Red	Light Blue	Red	Red	Green	Red	Green
Graphical Models	Light Blue	Light Blue	Yellow	Green	Red	Light Blue	Light Blue	Light Blue	Green	Red	Red	Light Blue
Backtrack / B&B	Light Blue	Light Blue	Yellow	Green	Red	Light Blue	Red	Light Blue	Light Blue	Light Blue	Yellow	Light Blue
Finite State Mach.	Red	Red	Red	Yellow	Yellow	Light Blue	Yellow	Light Blue	Light Blue	Light Blue	Light Blue	Red
Circuits	Red	Light Blue	Green	Light Blue	Green	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Light Blue	Red
Dynamic Prog.	Yellow	Light Blue	Red	Light Blue	Red	Light Blue	Yellow	Light Blue	Light Blue	Yellow	Light Blue	Red
Structured Grid	Red	Red	Light Blue	Yellow	Light Blue	Red	Light Blue	Light Blue	Red	Light Blue	Light Blue	Light Blue
Dense Matrix	Red	Red	Yellow	Red	Red	Light Blue	Yellow	Light Blue	Red	Red	Red	Light Blue
Sparse Matrix	Yellow	Yellow	Light Blue	Red	Red	Light Blue	Yellow	Red	Light Blue	Light Blue	Red	Light Blue
Spectral (FFT)	Yellow	Light Blue	Light Blue	Yellow	Yellow	Red	Light Blue	Light Blue	Green	Red	Red	Red
Monte Carlo	Light Blue	Light Blue	Light Blue	Yellow	Light Blue	Red	Light Blue	Yellow	Light Blue	Light Blue	Light Blue	Light Blue
N-Body	Light Blue	Yellow	Light Blue	Yellow	Light Blue	Red	Light Blue	Green	Light Blue	Light Blue	Light Blue	Light Blue

Applications

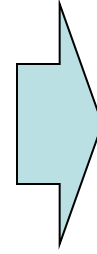


Pixel Dependencies

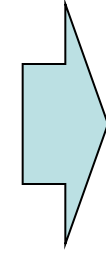
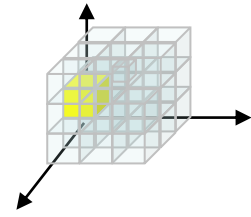
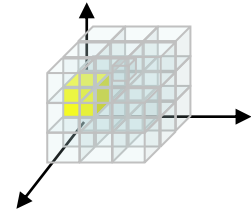
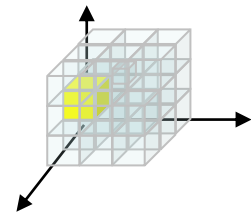


Computational Patterns

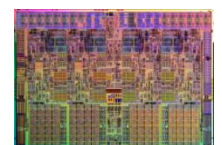
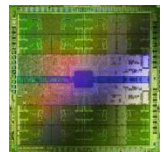
Computational patterns



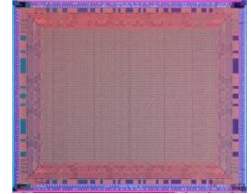
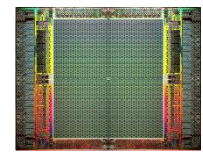
Program features



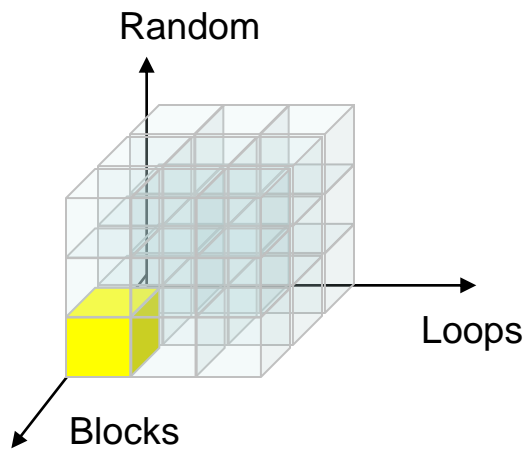
Architectures



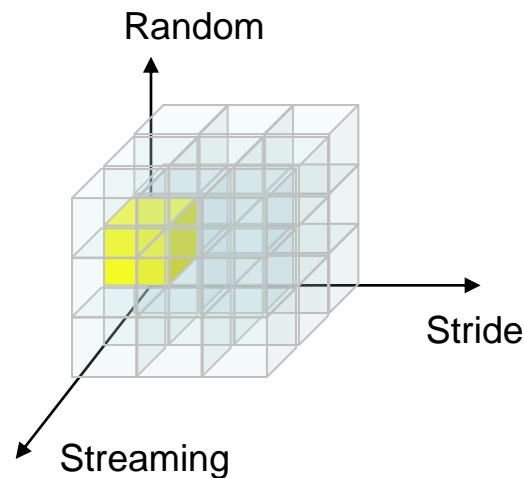
Accelerators



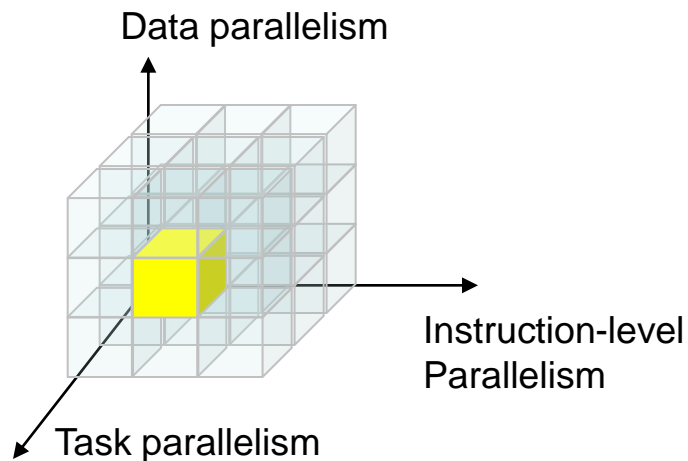
Instruction Access Pattern



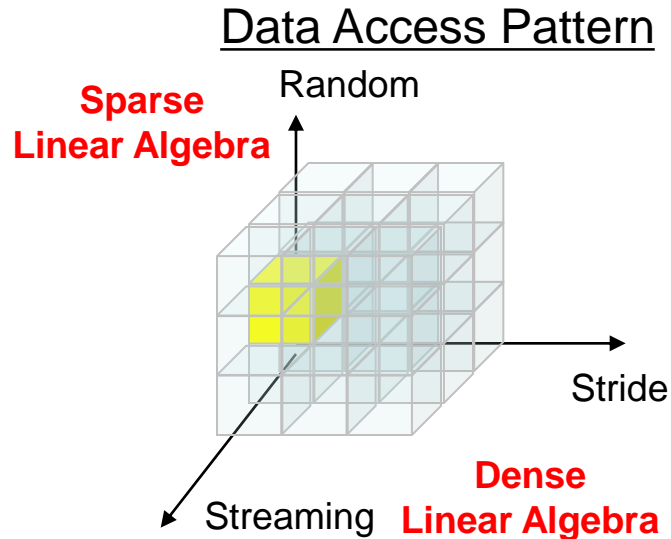
Data Access Pattern



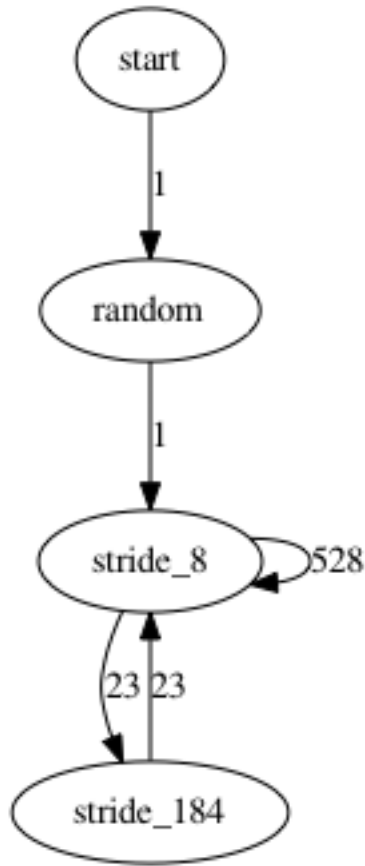
Concurrency Opportunities



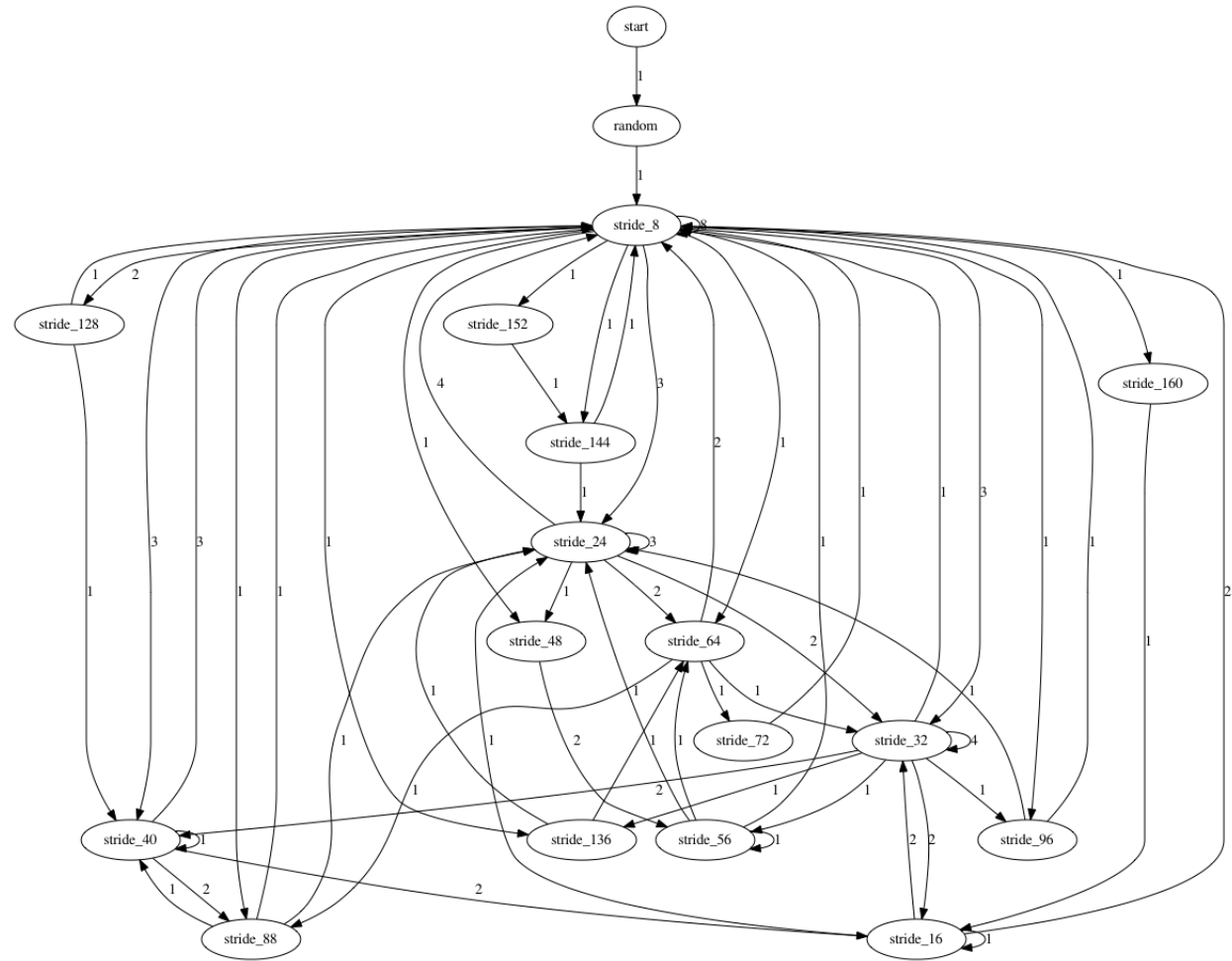
The developer needs to be aware of low-level application characteristics to efficiently map applications onto heterogeneous platforms



- The differences we observed between dense and sparse linear algebra can be mapped to the cube
 - Microarchitectural features can be optimized for the dimensions of the cube
 - DMA for streaming
 - Multithreaded processor for random accesses



Dense linear algebra



Sparse linear algebra

- Memory subsystem
 - Caches
 - Software scratchpads
 - DMA engines
 - Coalescing hardware
- Processor core configuration
 - Issue width
 - Wide out-of-order vs scalar in-order
 - Data parallel support
 - SIMD width
 - multithreaded execution
- System configuration
 - Cache coherence
 - Message passing
 - NUMA
 - Distribution of core type

- We presented 9 program features earlier
- Instruction access patterns
 - Random
 - Blocks
 - Loops
- Data access patterns
 - Random
 - Stride
 - Streaming
- Concurrency opportunities
 - Data
 - Task
 - Instruction
- We are building a heatmap for program features
 - Reflects what we know today
 - Not complete

	Instruction access			Data access			Parallelism		
	Random	Blocks	Loops	Random	Strided	Streaming	Data	Task	Instruction
Dense linear algebra	Blue	Blue	Red	Blue	Red	Red	Red	Red	Red
Sparse linear algebra	Blue	Blue	Red	Red	Blue	Red	Green	Red	Green
Structured grids	Blue	Blue	Red	Blue	Red	Blue	Red	Red	Red
Unstructured grids	Blue	Blue	Red	Yellow	Yellow	Yellow	Yellow	Red	Yellow
Spectral methods	Blue	Blue	Red	Blue	Red	Red	Red	Red	Red
Particle methods	Blue	Blue	Red	Green	Blue	Red	Yellow	Red	Yellow
Monte Carlo methods	Blue	Blue	Red	Blue	Blue	Red	Red	Red	Red
Combinational logic	Green	Red	Red	Yellow	Blue	Yellow	Yellow	Blue	Red
Finite state machines	Yellow	Green	Green	Yellow	Blue	Blue	Blue	Yellow	Red
Backtrack and B&B	Red	Blue	Blue	Red	Blue	Blue	Blue	Blue	Green
Graph algorithms	Yellow	Yellow	Yellow	Red	Blue	Yellow	Green	Yellow	Blue
Dynamic programming	Blue	Yellow	Yellow	Green	Yellow	Yellow	Green	Yellow	Yellow

- Patterns help the developer determine “where does computation want to happen”
- From our study of existing applications, we believe patterns can guide exploration of applications on emerging heterogeneous platforms
- New heterogeneous platforms will make the algorithmic design space more complicated
 - A disciplined programming methodology is required to fully exploit these new platforms