

# Harmony: Co-Optimizing Parallelism and Locality to Bound Performance Jennifer Brana, Nathan Beckmann

Carnegie Mellon University

## Data Movement-Aware Bounding

### Chaos Processor Model

**Our goal: understanding the fundamental limits of** performance in real systems

**Our approach:** instruction placement in parallel processors



**Chaos** is a processor model that lets us isolate and study the tradeoff between parallelism and locality in a single framework



...incorporates **dynamism** of OoO scheduling to throttle parallelism

Instruction placement alg



### Why instruction placement?

...a spatial execution model exposes control of dataflow and instruction locality to online instruction placement algorithms





...captures the effect of exploiting locality to reduce data movement overheads

## Harmony Placement Algorithm

Harmony places instructions by performing a local search over space and time, scheduling the instruction at the earliest available PE, according to locality and already-scheduled instructions

Harmony example placing operation 'Z':

Find PEs reachable earliest by 1) both input operands (x & y) Highlighted in blue





...captures the impact of communication latency & instruction fetch penalty

...enables us to exploit the tradeoff between parallelism and locality to maximize performance



## Limitations of Prior Work

**OoO superscalar models: Spatial dataflow models:** Optimize parallelism via dynamic Optimize locality by unrolling program in hardware & binding scheduling to expose parallelism within & across loops static instructions to PEs

Check when PEs can execute Z 2)

> Other isns are scheduled on the ideal PEs

- 3) If no PEs can execute Z at the earliest time, PEs one cycle from the list are added and the next cycle is checked
  - Highlighted in red
  - Z is scheduled on PE 3



## Evaluation

We evaluate **Harmony** using a trace-driven simulation of **Chaos** We evaluate 6 common workloads: DFS, BFS, SMV, SpMV, DMV, and DMM



Harmony averages a speedup of **1.98**× over the next best method ...co-optimizes parallelism, dataflow locality, and instruction locality

...finds the best balance between parallelism and locality



