# Spatio-Temporal Memory Streaming

Stephen Somogyi, Thomas F. Wenisch, Anastasia Ailamaki and Babak Falsafi

36<sup>th</sup> International Symposium on Computer Architecture

June 22, 2009



Computer Architecture Lab at Carnegie Mellon



Michigan Engineering

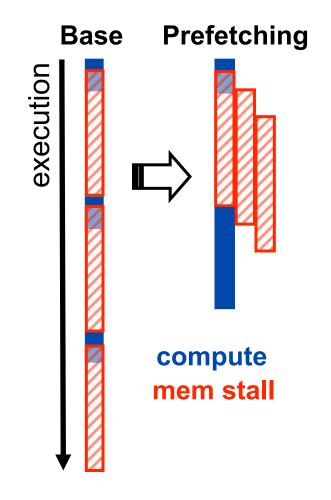


### **Memory Wall: Performance Bottleneck**

Especially for server apps

- Large data footprints
- Pointer-intensive structures

Prefetching can hide long memory latency [ISCA '05] [ISCA '06] [Micro '07]

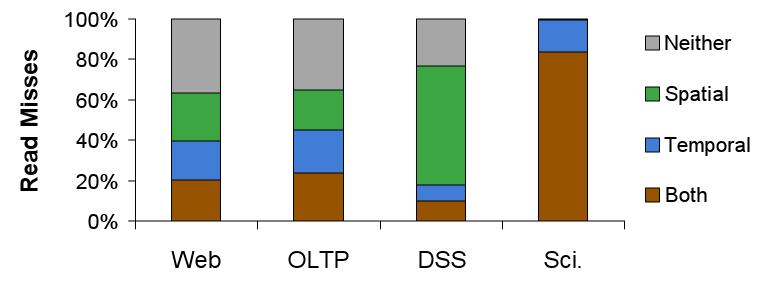


But, no single technique effective for OLTP, DSS & Web

### **Observation: Temporal, Spatial Predictions Different**

Different predictors capture different behaviors

- Temporal: recurring memory access sequences (pointers)
- Spatial: recurring data layouts (structs)



Spatial/temporal disjoint; opportunity to exploit both

## **How to Combine?**

**Concept:** independent temporal & spatial

- Duplicates prefetches  $\Rightarrow$  inefficient
- Prefetchers interfere  $\Rightarrow$  confuses training

**Refinement:** chain spatial predictions via temporal

- No sequence within spatial predictions
- Prefetches in wrong order  $\Rightarrow$  not timely

Solution: also learn order within spatial predictions

• Achieves high, consistent coverage & low mispredictions

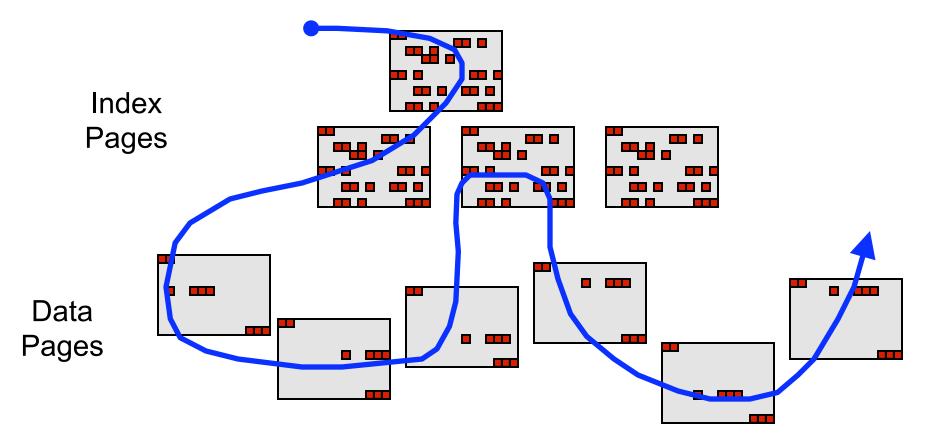
## Contributions

- Opportunity for spatio-temporal prediction
  - 70% of read misses are predictable on average
- Temporal characterization of spatial accesses
  Sequences repetitive within & across layouts
- Spatio-Temporal Memory Streaming
  - Predicts unified spatio-temporal miss sequence
    - 62% of read misses on average
  - Mean speedup 1.31,  $\geq$  temporal or spatial alone

## **Outline**

- Introduction
- Spatial and Temporal Prediction
- Spatio-Temporal Memory Streaming
- Results
- Conclusion

## **Example: Non-Clustered Index Scan**

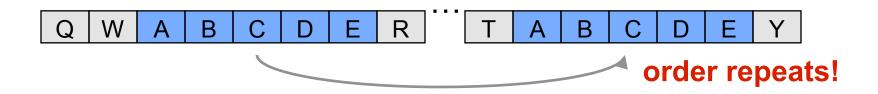


Sequence spans non-contiguous data pages Similar layouts within data pages

#### Temporal Memory Streaming (TMS) [Wenisch '05]

Records & replays recurring miss sequences

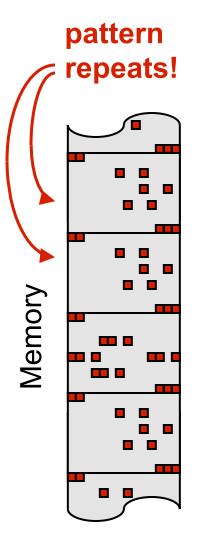
• Code traversals repeat  $\Rightarrow$  data traversals repeat



Sequences contain entire addresses

- $\checkmark$  Good for pointer chasing  $\Rightarrow$  breaks dependence chains
- ✓ Startup costs amortized over long streams
- Cannot predict compulsory misses
- Large storage required (~2MB / processor)

#### Spatial Memory Streaming (SMS) [Somogyi '06]



Exploits repetitive, large-scale data layouts in memory

- Pattern: offsets in logical region
- Trigger: first miss, used for lookup

Patterns encoded as bit vectors

- ✓ PC lookup  $\Rightarrow$  predicts compulsory misses
- ✓ Efficient storage (~80KB / processor)
- **x** Trigger miss per pattern  $\Rightarrow$  lost opportunity
- Unordered  $\Rightarrow$  BW spikes, wrong prioritization

# **Hybrid Spatio-Temporal Predictor**

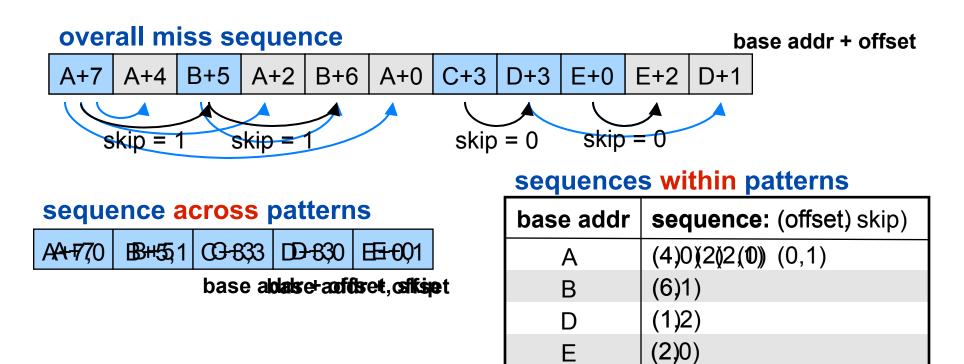
#### Naïve approach

- Record sequence across patterns
- Fetch entire pattern at a time
  - No priority across patterns
- But, triggers many patterns simultaneously
  - Accesses across patterns are interleaved!
  - Bursts of prefetches  $\Rightarrow$  pollution, BW spikes

#### Must prioritize prefetches across patterns

# **Deconstructing a Miss Sequence**

#### Investigate temporal & spatial relationships



Reveive appreash operaters patterniss before Bce Incorrect order!

### Spatio-Temporal Memory Streaming (STeMS)

**Goal:** reconstruct the overall miss sequence

- Using both temporal & spatial predictions
- $\Rightarrow$ Prefetch cache blocks in order

#### Training: observe relative interleavings

- Record skips in temporal seq. and spatial patterns

**Prediction:** reconstruction buffer for staging

- Spread temporal predictions according to skips
- Trigger spatial lookup for each, insert predicted addresses

#### Generates simple address seq. for throttled streaming

## Outline

- Introduction
- Spatial and Temporal Prediction
- Spatio-Temporal Memory Streaming
- Results
- Conclusion

# Methodology

Flexus [Wenisch '06]

- Full-system trace and OoO timing simulation
- Leverages SMARTS sampling

#### **Benchmark Applications**

- OLTP: TPC-C
  \_IBM DB2 & Oracle
- DSS: TPC-H Qry 2,16,17
  –IBM DB2
- Web: SPECweb99
  - -Apache & Zeus
- Scientific
  - -em3d, ocean, sparse

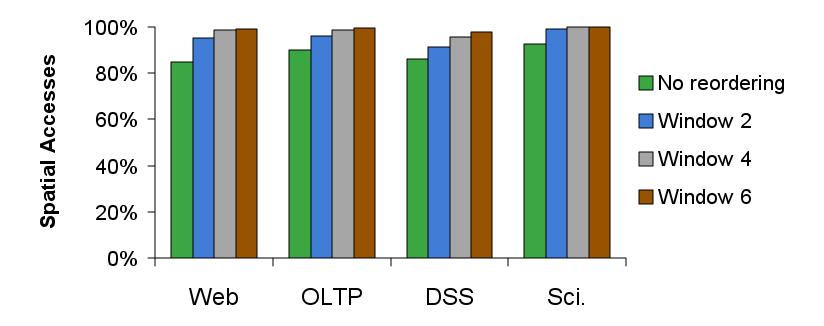
#### **Model Parameters**

- □ 16 4GHz SPARC CPUs
- □ 4-wide OoO; 96-entry ROB
- □ 64KB 2-way L1
- □ 8MB 8-way L2, 25-cycle lat.
- □ 40ns memory
- □ 25ns per-hop network
- □ TSO w/ speculation

## **Temporal Repetition Within Patterns**

Compare access sequence of successive patterns

- Evaluate for small reordering windows

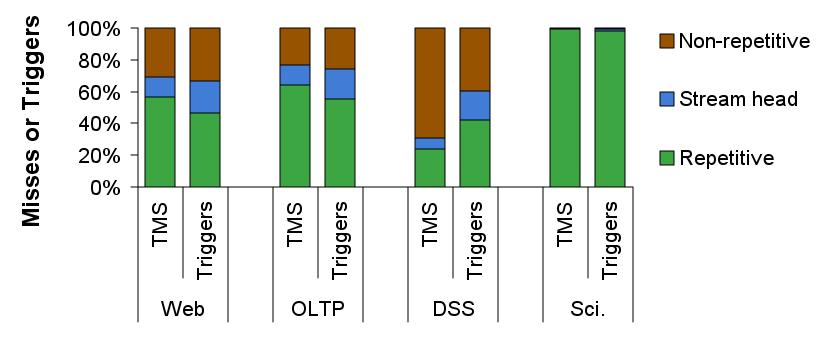


Seq. of accesses within patterns extremely repetitive

# **Temporal Repetition Across Patterns**

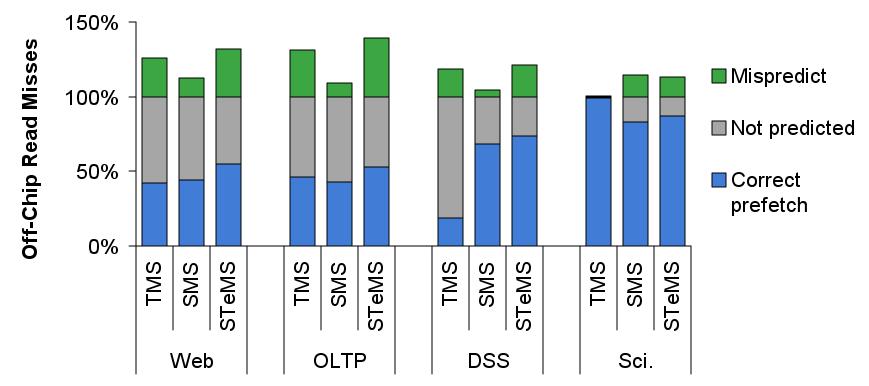
Evaluate repetition with compression algorithm

- Past work (TMS): all miss addresses
- STeMS: only trigger misses



Similar opportunity for predicting seq. of triggers





**STeMS is effective across commercial workloads** 

Predicts 62% of misses, improves perf. 31% **OLTP/DSS** – matches TMS/SMS, **Web** – beats both

## **Related Work**

- Stream Chaining
  - Reconstruct overall miss seq. using control flow
  - Better compulsory, worse temporal coverage
- Predictor Virtualization
  - Reduces dedicated on-chip storage for predictors
  - Can be applied to history structures in STeMS
- Epoch-base Correlation Prefetching
  - Improves timeliness of temporal predictions
  - In contrast, STeMS mainly targets spatial timeliness

[Burcea '08]

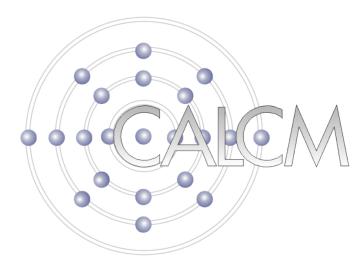
[Chou '07]

[Diaz '09]

## Conclusion

- Spatial/temporal predictions disjoint
  - Opportunity to predict up to 70% of read misses
- Temporal repetition of spatial patterns
  - Near-perfect repetition within patterns
  - Similar repetition across patterns as all addresses
- Design for Spatio-Temporal Memory Streaming
  - Reconstructs total miss sequence
    - Predicts 62% of read misses, perf. improvement 31%
  - Coverage & speedup  $\geq$  temporal or spatial alone

# **Questions** ?



STeMS Project Spatio-Temporal Memory Streaming www.ece.cmu.edu/~stems

Computer Architecture Laboratory Carnegie Mellon University <u>www.ece.cmu.edu/~calcm</u>