# Understanding
# Memory Access Bottlenecks
# On Multicore

Mini Symposium "Scalability and Usability of HPC Programming Tools"
ParCo 2007, FZ Jülich, September 5, 2007

## Josef Weidendorfer

Lehrstuhl für Rechnertechnik und Rechnerorganisation
Institut für Informatik, Technische Universität München

## Outline

- **Introduction**
  Relevance for Mini Symposium Topic


- **Memory Access Bottlenecks and Metrics**
  From Single- to Multicore


- **Our Tool Suite**
  Callgrind / KCachegrind
  Issues with Existing Visualizations & Experiences


- **Conclusion and Future Work**

# Relevance for Mini Symposium Topic

- Usability
- Scalability

Cope with large amount of events (memory accesses !)
- In the measurement tool
- In the visualization tool

In measurement
- Use statistical methods / simulation
- Online processing to reduce amount of data
    - Filtering / Selection
    - Aggregation ($\rightarrow$ Profile)

# Memory Access Bottlenecks

**Main problem: "Memory wall"**
<span style="color:red">True for Multicore, too</span>

Increasing performance gap main memory vs. processor

**Solution: Caches**
<span style="color:red">True for Multicore, too</span>

- Exploit locality of memory accesses (temporal / spatial)
- Lowers access latency by putting data copies into fast memory
  - Keep recently used copies (accounts for temporal locality)
  - Block oriented (accounts for spatial locality)
- „Bad memory access behavior": Poor exploitation of caches

**Optimization strategies**
<span style="color:red">True for Multicore, too</span>

- Improve temporal locality by reordering accesses
- Improve spatial locality by changing data layout
- Prefetch data needed in the future

# Memory Access Bottlenecks

Multicore adds

- All the issues of parallel code
  - Load balancing, synchronization overhead, difficult to program …
- Cores share available resources for
  - Connection to main memory
  - Caches
- Software has to cope with new configurations
  - Caches shared vs. separate per core, increased memory hierarchy

Optimization strategies
  - Best process placing for good exploitation of available caches?

# Analysis of Memory Access Behavior

Good Tools™ give following answers

- Yes, we have a problem because of memory accesses

- And it happens there
  (code position + call path + data structure + thread/process)

- Yes, it makes sense to optimize (potential benefit)

- Just try to do this and that to avoid cache pollution / bad layout / …
  (good hints with expertise)

What are good metrics to see the problem and path to solution?

## Metrics for Bad Memory Access Behavior

- Cache miss counts
  - Good: Pinpoints where time is lost (But: how much? Data structure?)
  - Difficult to derive optimization (e.g. what/how to block?)

- Temporal / spatial cache line usage:
  How much / often used before eviction
  - Easy to see bad memory layout (e.g. hash lookups)

- Model of idealized (unlimited) cache with LRU list of data accesses

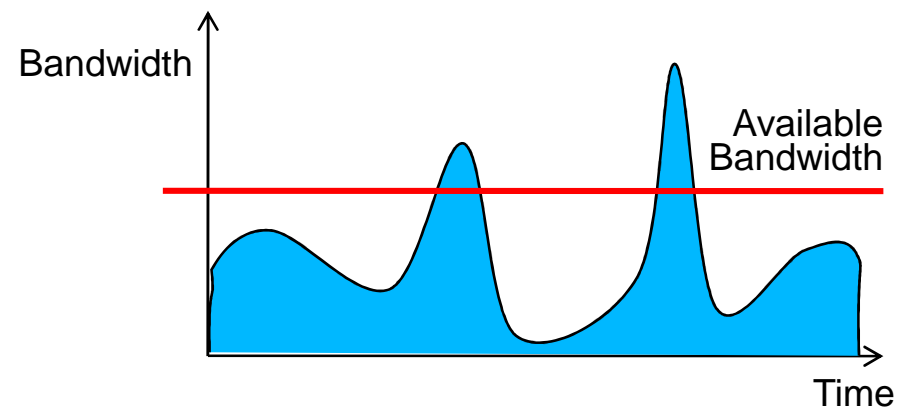- Simulation of bandwidth requirement

## LRU List of Data Accesses

- Relates to behavior of fully associate cache

- Example: Address sequence 1 2 3 4 2 3 2 2 5

| LRU offset | 0 | 1 | 2 | 3 | 4 |
|---|---|---|---|---|---|
| Address | 5 | 2 | 3 | 4 | 1 |
| Last access (time steps) | 0 | 1 | 3 | 5 | 8 |

| Distance | 1 | 2 | 3 |
|---|---|---|---|
| Accesses | 1 | 1 | 2 |

- Provides size of working set (used in given time span)

- Stack reuse distance: histogram over move distances on access (See papers of K. Beyls)

  - Percentage of accesses covered by given cache size

  - Which data structures have bad locality with high influence

# Simulation of Bandwidth Requirement

- What is the required amount of data with ideal memory?
  (Suppose only computation-boundness)
  - Has to be simulated
  - Simple CPU model (given latency for each opcode)



  - Which structures are better served from cache?
  - Gives hints for prefetching (can prefetching pollute cache?)

# Adaption for Multicore

- ## Bandwidth requirement diagram (easy)
  - Additional curves for shared vs. separate caches

- ## Multiple LRU lists
  - for all cores
    - reuse distances with shared cache
    - provides information for workset overlapping
  - for each core pair: shared with other", „invalidated by other"
    - Shows amount communication between cores (size + number)

- ## Should provide hints for
  - placement of processes (exploit same workset via shared cache)
  - prefetch helper thread for multiple cores

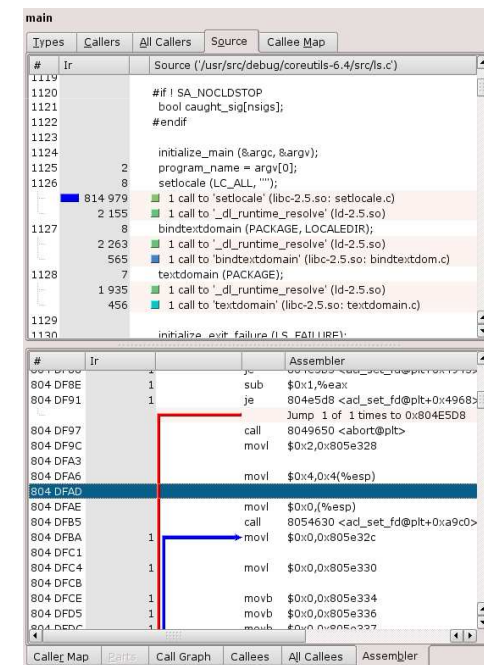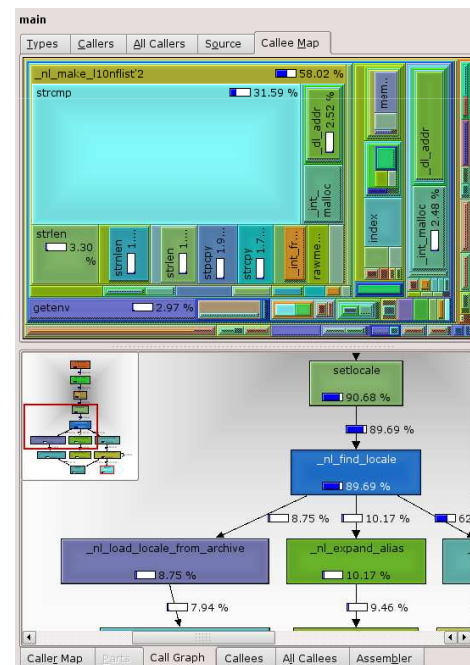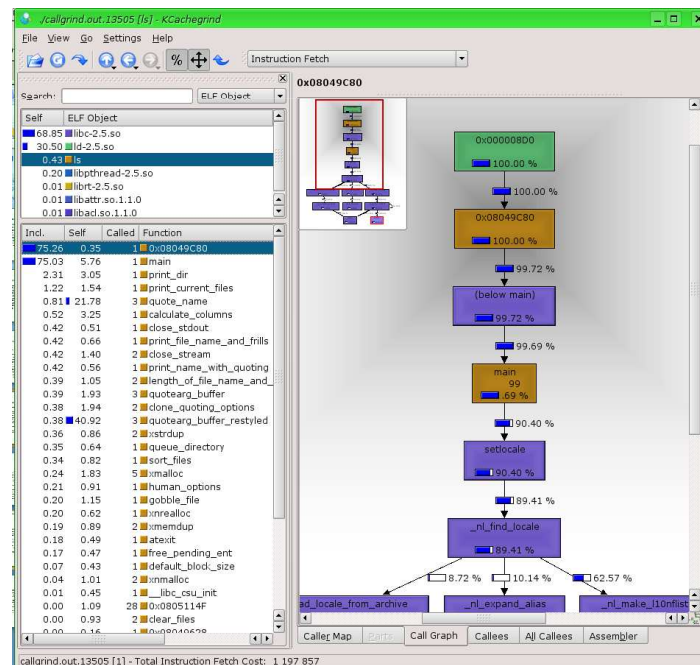# Our Tool Suite: Callgrind / KCachegrind

## Measurement

- Based on Valgrind (runtime instrumentation, known for "memcheck")
- Instrument memory accesses feed cache simulator
- Profiling tool relating cache events to call-graph (path relation possible)

## Pro and Contra

- Memory accesses only from user-level code, Slowdown (40-60x)
- Synchronous 2-level inclusive cache (optional hardware prefetcher)

- ✓ Does not need root access / can not crash machine
- ✓ Allows for sophisticated metrics (line usage, stack reuse distance)
- ✓ Easy to understand / reconstruct for user
- ✓ Reproducible results independent on real machine load
- ✓ Derived optimizations applicable for many architectures

# Our Tool Suite: Callgrind / KCachegrind

- Visualization
    - Call relationship of functions (callers, callees, call graph, tree map)
    - Source/Assembly annotation: event counts + control flow

# Experiences with existing Visualizations

Provide not only graphical results, but
- export to ASCII lists (HTML / XML), or better
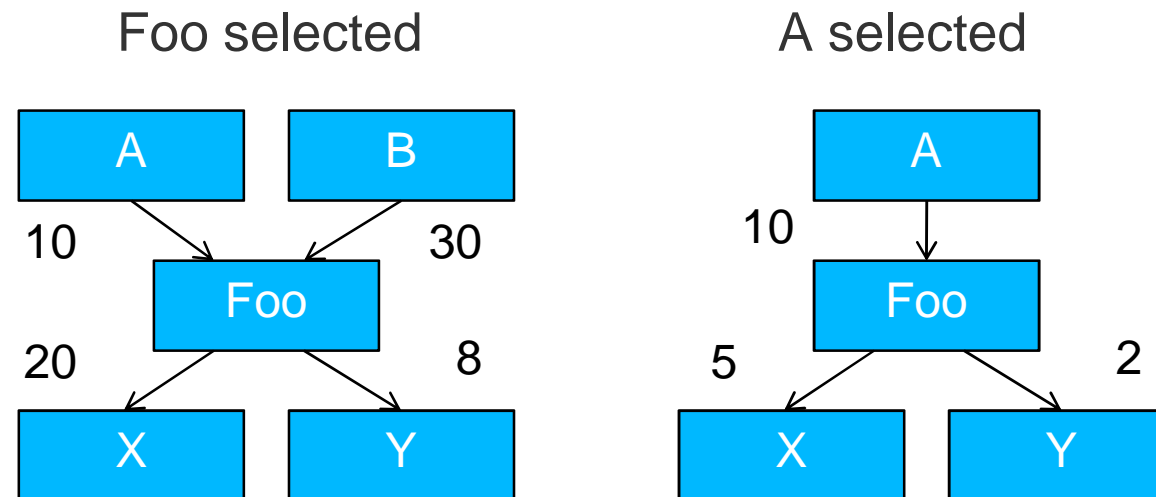- equivalent command line tool for scriptability: merging, querying, …

Graphical interface
- Very clean user interface, small number of views / areas
- Only a few different view types into data, familiar to users
  - ✓ lists, call graphs, bar charts, heat maps, diagrams
  - – tree maps, fully crowded 3-D views
- Highly interactive views
  - – responsive
  - – zoomable (one visualization with multiple detail levels)
  - – intuitive narrowing / filtering of data
  - – non-disruptive browsing (incremental or animation)
- Careful, consistent use of color coding

# Improvements for existing Visualization: Call Graph

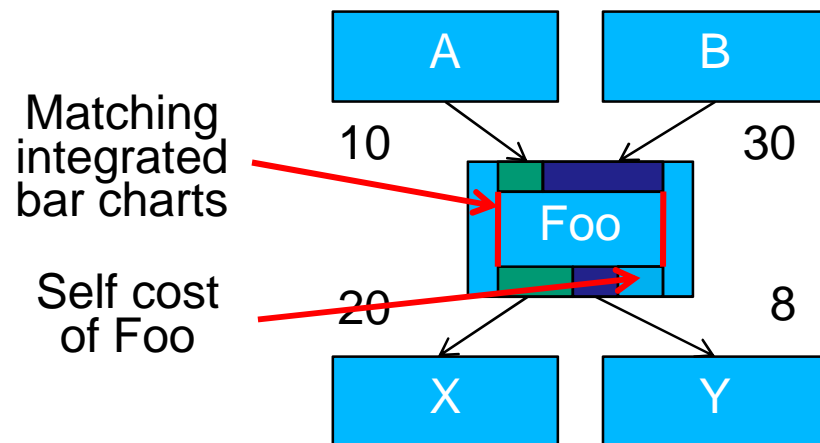Issue: Tools seems to claim to always have full path profiling data

Example:          Foo selected                    A selected



Reality could have been: Foo → X only when called from B !
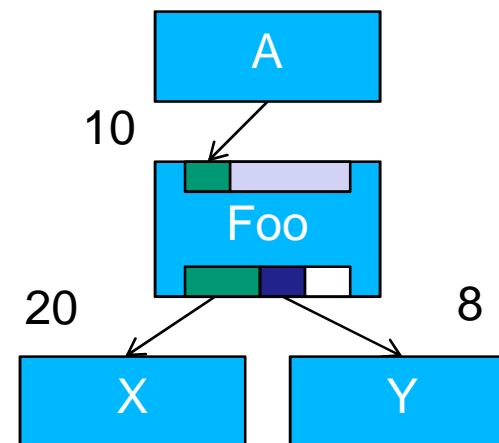
# Improvements for existing Visualization: Call Graph

Issue: Tools seems to claim to always have full path profiling data
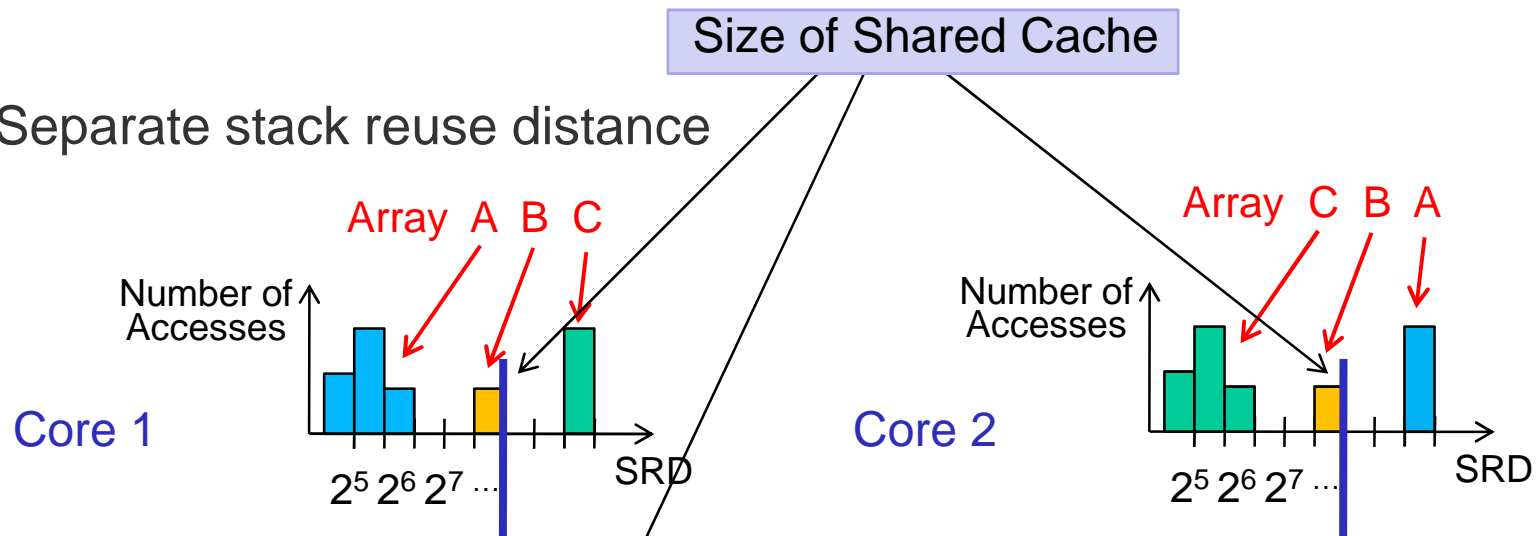
Example:   Foo selected                    A selected (improved)
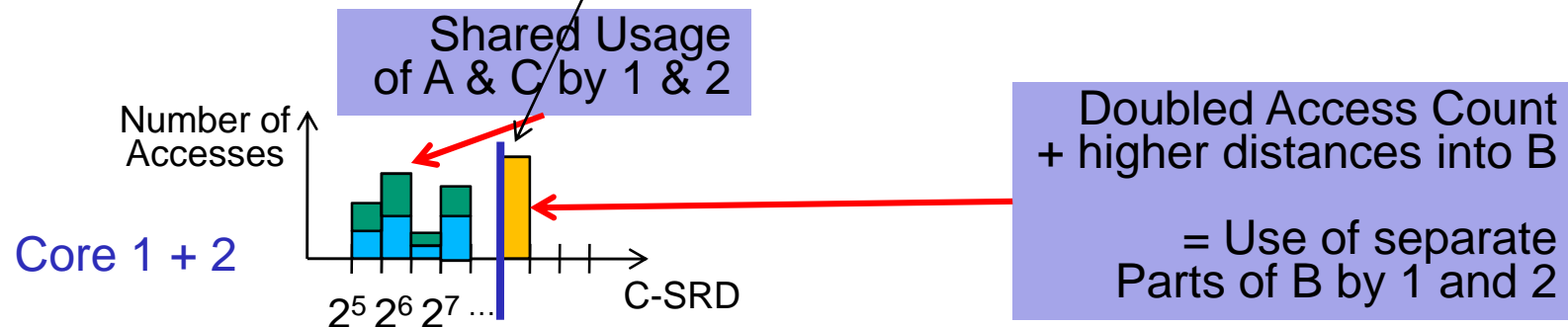


- Tells reality
- With real path profiles: Allows for interactive filtering of data
- Also for filtering in other dimensions (partitioning against threads)
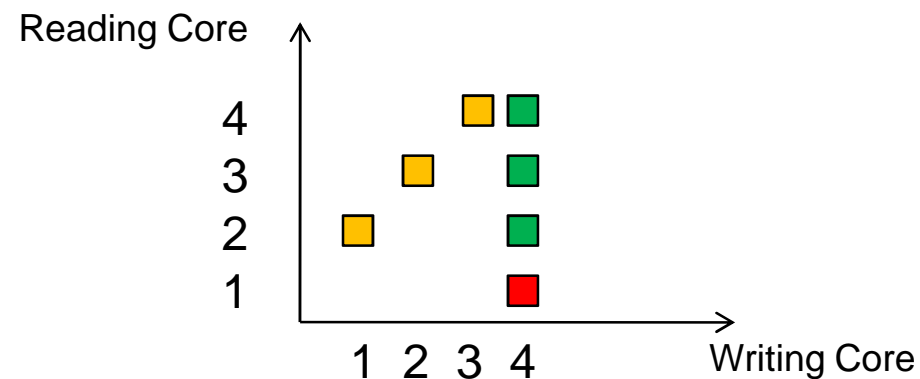
# Visualization for Multicore

Size of Shared Cache

- Separate stack reuse distance

Array A B C

Number of Accesses

Core 1

$2^5$ $2^6$ $2^7$ ...   SRD

Array C B A

Number of Accesses

Core 2

$2^5$ $2^6$ $2^7$ ...   SRD

- Combined stack reuse distance (C-SRD)

Shared Usage of A & C by 1 & 2

Number of Accesses

Core 1 + 2

$2^5$ $2^6$ $2^7$ ...   C-SRD

Doubled Access Count + higher distances into B

= Use of separate Parts of B by 1 and 2

MUNICH-CORE INITIATIVE

# Visualization for Multicore (2)

- Similar to get communication behavior among cores
  - Only count invalidations (with separate caches)
    - Heat map (color: number of invalidations)



- Histograms similar to last slide…

# Conclusion

- Memory Access analysis on Single/Multicore benefits from better metrics than cache misses

- Simulation with runtime instrumentation proved useful for this

- We propose metrics and visualization for workset overlapping/communication on Multicore

# Future work

Work in progress

- Measurement

  – Should stay usable (resource consumption of simulation)

  – Enhancement for multicore simulation in industry cooperation

  – Keep amount of data small

    - E.g. online aggregation also for bandwidth requirement
    - Statistical reconstruction of approximative time scale visualization ?

- Visualization

  – Prototype some visualizations for multicore memory access

  – Conduct usability studies