

Lecture 23 – Programming the Cell BE (III) + Sequoia
Compilation, Tuning, and Runtime

Mattan Erez



The University of Texas at Austin



Outline

- Sequoia Summary
- Other Cell programming tools
- Sequoia runtime and compilation
- Sequoia part courtesy Kayvon Fatahalian, Stanford
- All Cell related images and figures © Sony and IBM
- Cell Broadband Engine™ Sony Corp.



Emerging Themes

- Writing high-performance code amounts to...
 - Intelligently structuring algorithms
[compiler help unlikely]
 - Efficiently using communication
 - Efficiently using parallel resources
[compilers struggle without help]
 - Generating efficient inner loops (kernels)
[compilers coming around]



Sequoia

- Language: stream programming for machines with deep memory hierarchies
- Idea: Expose abstract memory hierarchy to programmer
- Implementation: **language, compiler, tuner, and runtime**
 - benchmarks run well on Cell processor based systems, clusters of PCs, SMPs, out-of-core computation, and combinations of above




Sequoia's method

- Explicit communication between abstract memories
- Locality awareness
- Hierarchy portability
 - Across machines, within levels of a machine
- **Programmer expresses combined computation and decomposition parameterized algorithm**
 - System follows algorithm to map to a specific machine



Sequoia summary


- Problem:
 - Deep memory hierarchies pose perf. programming challenge
 - Memory hierarchy different for different machines
- Solution: Abstract hierarchical memory in programming model
 - Program the memory hierarchy explicitly
 - Expose properties that effect performance
- Approach: Express hierarchies of tasks
 - Execute in local address space
 - Call-by-value-result semantics exposes communication
 - Parameterized for portability



Sequoia and Cell Programming Challenges

- Sequoia manages threading and synchronization
- Sequoia manages communication and all DMAs
 - Including padding and performance, but not alignment
- Sequoia manages LS
 - Allocation and packing
- Sequoia manages scheduling
 - SWP of mappar to hide communication latency
- Sequoia doesn't help with SPE code
 - Use low-level compiler tools such as XLC
- Sequoia doesn't currently help with some memory restrictions
 - Alignment, limited support but cannot rely on Sequoia
 - Banks
 - Sequoia does pad for access granularity restrictions


© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 7



Outline

- Sequoia Summary
- Other Cell programming tools
- Sequoia runtime and compilation
- Sequoia part courtesy Kayvon Fatahalian, Stanford
- All Cell related images and figures © Sony and IBM
- Cell Broadband Engine™ Sony Corp.


© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 8



Tools From IBM

- Cell SDK 3.0
 - API calls for handling communication, synchronization, and DMA
 - LIBSPU and SPUPS for getting the SPEs to do something and setting up threads and memory
 - Intrinsics for programming the SPE pipeline directly
 - GCC port for PPE and SPE part (separate compilers)
 - Only handles non-SPE specific optimizations + intrinsics
 - XLC port for PPE and SPE part (separate compilers)
 - XLC supposed to optimize for SPE pipeline with branch hints, scheduling, instruction prefetch, ...
 - Automatic SIMD-ization?
- Accelerated Library Framework (ALF)
 - APIs for work queue based model to program control-plane
- "Octopiler" - single-source XLC for Cell
 - OpenMP directives
 - Relies on SW cache to get the OpenMP working
 - Automatic SIMD-ization


© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 9



Tools from Industry

- Mercury Systems
 - Array based language
 - Highly-tuned BLAS and FFT
- RapidMind
 - Dynamically compiled program
 - Relies on array data types
 - Builds up kernels and DMAs


© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 10



Tools From Academia

- Sequoia
- Cell Superscalar (CellSs)
 - Program with OpenMP like directives to identify kernels
 - Uses SW cache intensively
 - Runtime applies superscalar style optimization and scheduling to coarse-grained kernels (identified above)
- Charm++
 - Runtime based approach
 - Objects with explicit communication and "entry points" for synchronization
 - Uses a work queue and peaks into it to do the DMAs

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 11



Outline

- Sequoia Summary
- Other Cell programming tools
- Sequoia runtime and compilation
- Sequoia part courtesy Kayvon Fatahalian, Stanford
- All Cell related images and figures © Sony and IBM
- Cell Broadband Engine™ Sony Corp.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 12

How this all works (Cell example)

- Back to SGEMM example:
 - User initializes Sequoia and allocates data from their code

```

main()
{
    sqInit();
    ...
    A = sqAlloc2D(...);
    B = sqAlloc2D(...);
    C = sqAlloc2D(...);
    ...
    matmul(A,B,C);
    ...
    sqShutdown();
}

```

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 11

How this all works

- Back to SGEMM example:
 - User initializes Sequoia and allocates data from their code

```

main()
{
    sqInit();
    ...
    A = sqAlloc2D(...);
    B = sqAlloc2D(...);
    C = sqAlloc2D(...);
    ...
    matmul(A,B,C);
    ...
    sqShutdown();
}

```

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 14

How this all works

- Back to SGEMM example:
 - User initializes Sequoia and allocates data from their code

```

main()
{
    sqInit();
    ...
    A = sqAlloc2D(...);
    B = sqAlloc2D(...);
    C = sqAlloc2D(...);
    ...
    matmul(A,B,C);
    ...
    sqShutdown();
}

```

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 15

How this all works

- Back to SGEMM example:
 - User initializes Sequoia and allocates data from their code

```

main()
{
    sqInit();
    ...
    A = sqAlloc2D(...);
    B = sqAlloc2D(...);
    C = sqAlloc2D(...);
    ...
    matmul(A,B,C);
    ...
    sqShutdown();
}

```

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 16

Top level task call

```

task matmul::inner(in float A[M][T],
                  in float B[T][N],
                  inout float C[M][N])
{
    tunable int P, Q, R;
    mapper( int i=0 to M/P,
            int j=0 to N/R ) {
        mapseq( int k=0 to T/Q ) {
            matmul(A[P*i:P*(i+1);P][Q*k:Q*(k+1);Q],
                  B[Q*k:Q*(k+1);Q][R*j:R*(j+1);R],
                  C[P*i:P*(i+1);P][R*j:R*(j+1);R]);
        }
    }
}

```

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 17

Leaf task call

```

task matmul::inner(in float A[M][T],
                  in float B[T][N],
                  inout float C[M][N])
{
    tunable int P, Q, R;
    mapper( int i=0 to M/P,
            int j=0 to N/R ) {
        mapseq( int k=0 to T/Q ) {
            matmul(A[P*i:P*(i+1);P][Q*k:Q*(k+1);Q],
                  B[Q*k:Q*(k+1);Q][R*j:R*(j+1);R],
                  C[P*i:P*(i+1);P][R*j:R*(j+1);R]);
        }
    }
}

```

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 18

Leaf task return

```

task matmul::inner(in float A[M][T],
                  in float B[T][N],
                  inout float C[M][N])
{
    tunable int P, Q, R;
    mappar( int i=0 to M/P,
           int j=0 to N/R) {
        mapseq( int k=0 to T/Q ) {
            matmul(A[P*i:P*(i+1);P[Q*k:Q*(k+1);Q],
                  B[Q*k:Q*(k+1);Q][R*j:R*(j+1);R],
                  C[P*i:P*(i+1);P][R*j:R*(j+1);R]);
        }
    }
}

```

SPE mails PPE and waits for command

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 19

Control return to user code

- Back to SGEMM example:
 - User initializes Sequoia and allocates data from their code

```

main()
{
    sgInit();
    ...
    A = sgAlloc2D(...);
    B = sgAlloc2d(...);
    C = sgAlloc2d(...);
    ...
    matmul(A,B,C);
    ...
    sgShutdown(); // Kill off threads and cleanup
}

```

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 20

Autotuning Sequoia Programs

- Autotuner helps user with mapping (user can always override)

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 21

Specialization with Autotuning

- Work by Manman Ren (Stanford), PACT 2008
- Use Sequoia to identify what needs tuning
 - Explicit tunables and parameters in the language
- Tuning framework for SW-managed hierarchies
- Automatic profile guided search across tunables
 - Aggressive pruning
 - Illegal parameters (don't fit in memory level)
 - Tunable groups
 - Programmer input on ranges
 - Coarse → fine search
- Loop fusion across multiple loop levels
 - Measure profitability from tunable search
 - Adjust for "tunable mismatch"
 - Realign reuse to reduce communication

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 22

Overview: mapping the program

- Mapped versions are generated
 - Matching the decomposition hierarchy with the machine hierarchy
 - Choosing a variant for each call site
 - Set level of data objects and control statements

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 23

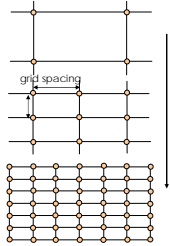
Explicit SW Management Simplifies Tuning

- Smooth search space
- Performance models can also work
 - For Cell, not cluster

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 24

The search algorithm

- A pyramid search
- A greedy search algorithm at each grid level
 - Achieve good performance quickly due to smoothness of the search space

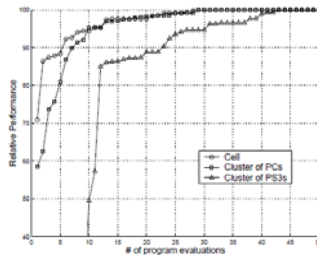


Start with a coarse grid
Refine the grid when no further progress can be made

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 24

Guided Search Converges Quickly

- Smoothness leads to quick convergence



© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 26

Autotuning Out Performs Programmer

| | | CONV2D | SGEMM | FFT3D | SUmb |
|-----------------|-----------|------------|------------|--------------|------|
| Cell | auto hand | 99.6 85 | 137 119 | 57 54 | 12.1 |
| Cluster of PCs | auto hand | 26.7 24 | 92.4 90 | 5.5 5.5 | 2.2 |
| Cluster of PS3s | auto hand | 20.7 19 | 33.4 30 | 0.57 0.36 | 0.63 |

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 27

Sequoia compilation

- A compiler for hierarchical bulk operations

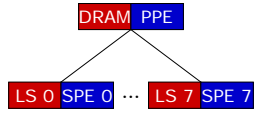
© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 28

Compiler for hierarchical bulk operations

- Is built around the portable abstractions of hierarchical memory and bulk operations.
- Automatically manages all program data movement to increase memory throughput.
- Automatically allocates explicitly managed memories.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 29

We abstract Cell as a tree of memories



- Why a tree of memories?
- A portable abstraction that matches our target machines of interest: Cell, clusters, etc.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 30

- We can model other machines as trees of memories (e.g.) a dual-GPU machine backed by disk storage:
- Not all tree nodes have a processor.
- Leaves of tree are compute-intensive processors.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 31

Modeled machine capabilities

- 1. Execute out of its local memory.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 32

Modeled machine capabilities

- 1. Execute out of its local memory.
- 2. Transfer data to/from a child.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 33

Modeled machine capabilities

- 1. Execute out of its local memory.
- 2. Transfer data to/from a child.
- 3. Transfer data to/from its parent.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 34

Modeled machine capabilities

- 1. Execute out of its local memory.
- 2. Transfer data to/from a child.
- 3. Transfer data to/from its parent.
- 4. Launch code in a child.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 35

Parallelism in the abstract machine model

- 1. Parallel PEs within a level.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 36

Parallelism in the abstract machine model

- 1. Parallel PEs within a level.
- 2. Concurrent parent/child execution.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 37

Parallelism in the abstract machine model

- 1. Parallel PEs within a level.
- 2. Concurrent parent/child execution.
- 3. Parallel execution within a PE.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 38

Modeling Programs

Programs comprise four elements

- 1. Operations (blue).

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 40

Programs comprise four elements

- 1. Operations (blue).
- 2. Data (red).

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 41

Programs comprise four elements

- 1. Operations (blue).
- 2. Data (red).
- 3. Dependences.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 42

Programs comprise four elements

- 1. Operations (green).
- 2. Data (blue).
- 3. Dependences.
- 4. Machine levels.

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 43

Programs comprise four elements

- 1. Operations (green).
- 2. Data (blue).
- 3. Dependences.
- 4. Machine levels.

- Operations may contain nested subprograms.

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 44

Programs comprise four elements

- 1. Operations (blue).
- 2. Data (red).
- 3. Dependences.
- 4. Machine levels.

- Operations may contain nested subprograms.

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 45

Example: SGEMV (matrix x vector)

$y = A \cdot x$

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 46

Example: SGEMV (matrix x vector)

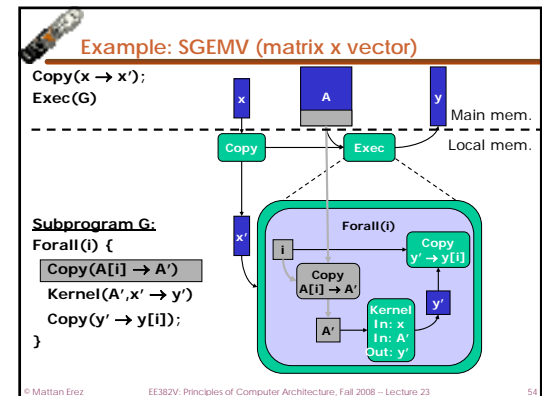
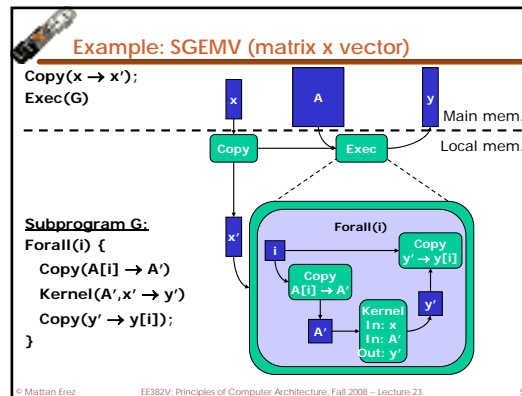
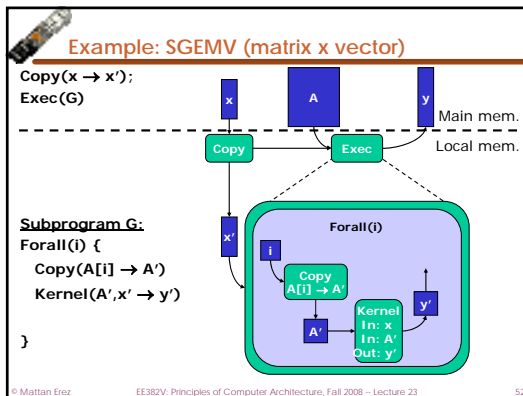
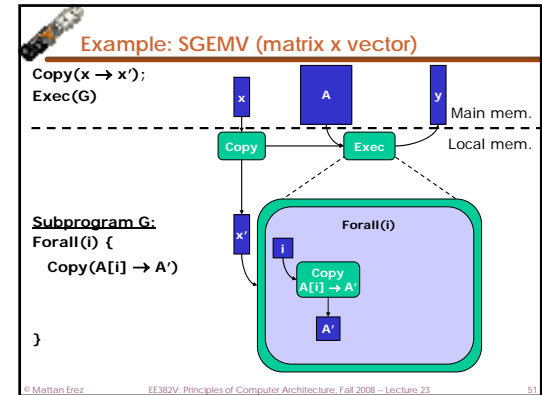
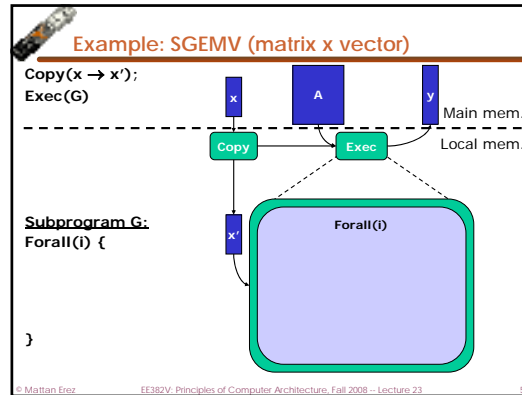
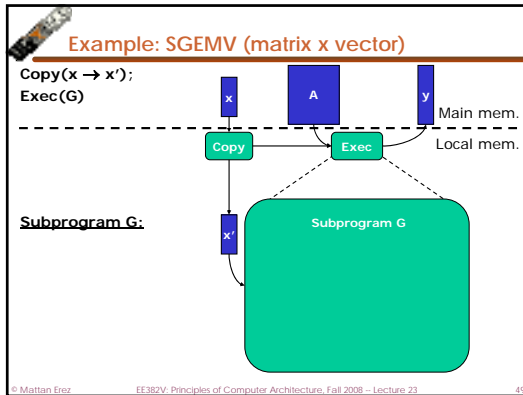
Copy($x \rightarrow x'$);

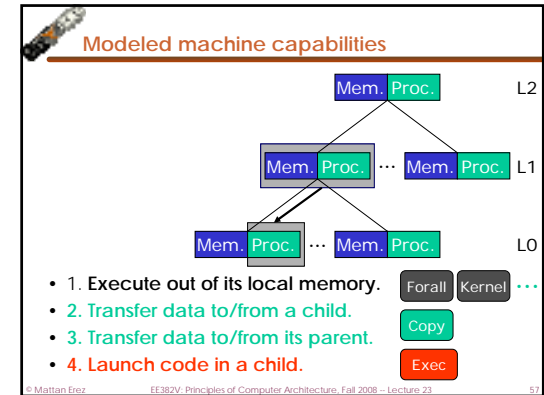
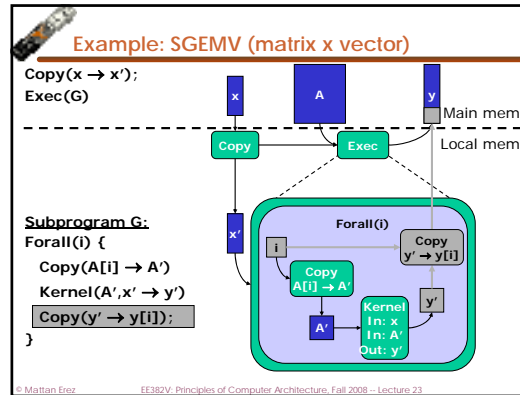
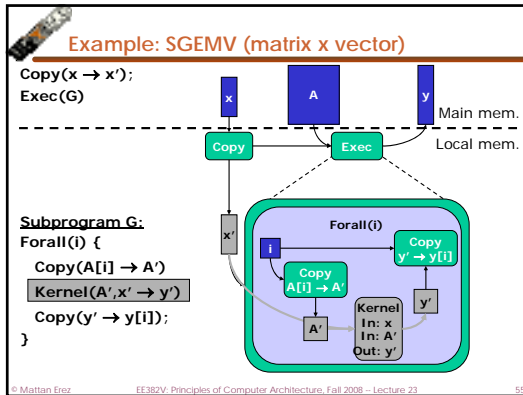
© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 47

Example: SGEMV (matrix x vector)

Copy($x \rightarrow x'$);
Exec(G)

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 48





- ### Program IR properties:
1. Is mechanism-independent.
 2. Features scalar and bulk operations within a single framework.
 3. Spans the entire machine (all levels).
 4. Fully captures all program semantics.
 5. Exposes parallelism via dependences and "Forall" operations.
- © Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 58

Optimizing Programs

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 60

- ### Optimizations are IR transformations
- They can be applied at any hierarchy level of the program or the machine.
- © Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 60

Optimizations are IR transformations

- They can be applied at any hierarchy level of the program or the machine.
- Correctness: Transformations preserve top-level program inputs and outputs.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 61

Optimization: Copy elimination

- 1. Removing spills [producer-consumer locality].

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 62

Optimization: Copy elimination

- 1. Removing spills [producer-consume locality].
- 2. Removing duplicates [temporal locality].

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 63

Optimization: Copy grouping

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 64

Optimization: Exec grouping

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 65

Optimization: Hoisting

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 66

Scheduling optimizations

1. Software pipelining.
2. Competing heuristics:
 1. Maximize operation concurrency.
 2. Group similar operations (e.g. Execs) to amortize issue overheads.

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 67

Optimization: Space allocation

Global view of data usage: better than LRU.

| | | | | |
|-----|------|-------|-------|--|
| OpW | B | A | C | |
| OpX | 50KB | 100KB | 60KB | |
| OpY | | D | E | |
| OpZ | F | 60KB | 140KB | |
| | 40KB | | | |

Memory space (0 to Max 240KB)

Program time

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 68

Optimization: Space allocation

Global view of data usage: better than LRU.

| | | | | |
|-----|------|-------|-------|--|
| OpW | B | A | C | |
| OpX | 50KB | 100KB | 60KB | |
| OpY | | D | E | |
| OpZ | F | 60KB | 140KB | |
| | 40KB | | | |

Memory space (0 to Max 240KB)

Program time

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 69

Optimization: Space allocation

Global view of data usage: better than LRU.

| | | | | |
|-----|------|-------|-------|--|
| OpW | B | A | C | |
| OpX | 50KB | 100KB | 60KB | |
| OpY | | D | E | |
| OpZ | F | 60KB | 140KB | |
| | 40KB | | | |

Memory space (0 to Max 240KB)

Program time

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 70

Optimization: Space allocation

Global view of data usage: better than LRU.

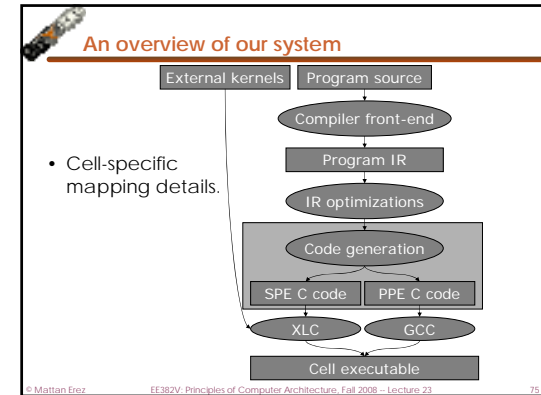
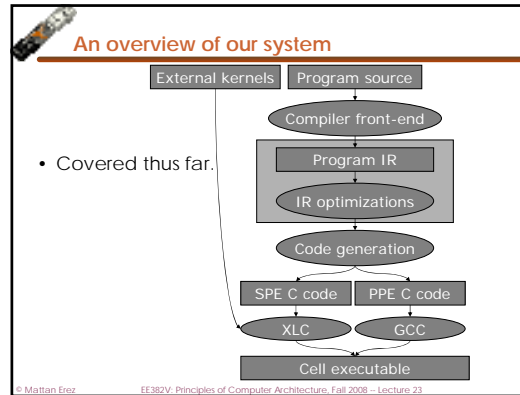
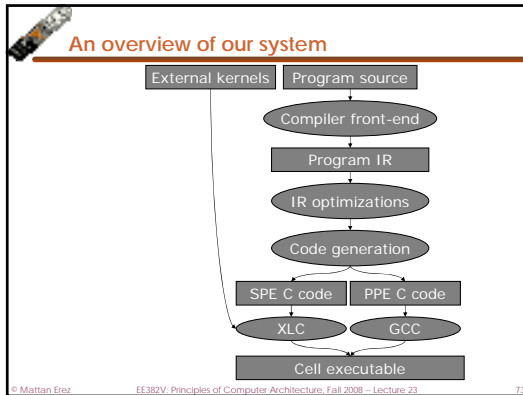
| | | | | |
|-----|------|-------|-------|--|
| OpW | B | A | C | |
| OpX | 50KB | 100KB | 60KB | |
| OpY | | D | E | |
| OpZ | F | 60KB | 140KB | |
| | 40KB | | | |

Memory space (0 to Max 240KB)

Program time

© Mattan Erez EE382V, Principles of Computer Architecture, Fall 2008 - Lecture 23 71

Targeting a Cell Processor



A few details

```

    graph TD
      CG[Code generation] --> SC[SPE C code]
      CG --> PC[PPE C code]
  
```

- Our source-to-source compiler generates two sets of output files.
- Each Exec operation → SPE overlay.
- Each Copy operation → DMA command.
 - Data objects and transfers padded to multiples of 16 bytes.
- DMA and SPE kernels overlapped where possible.

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 23 76

Summary

Summary of results

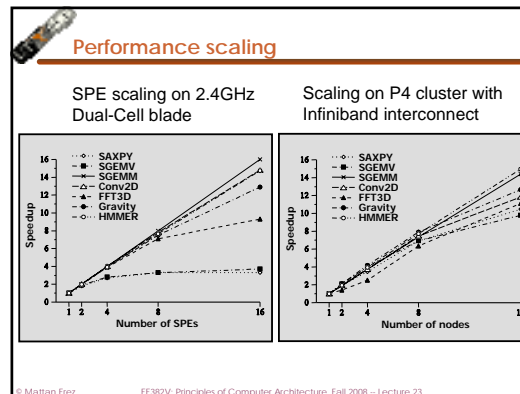
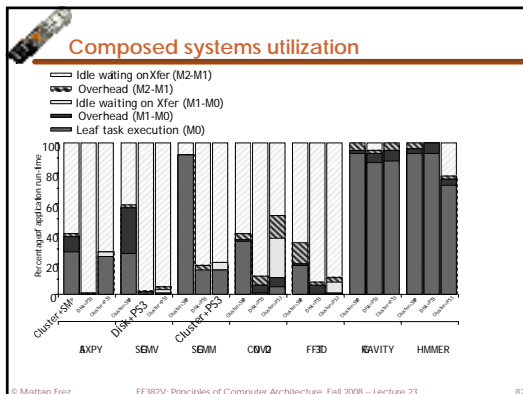
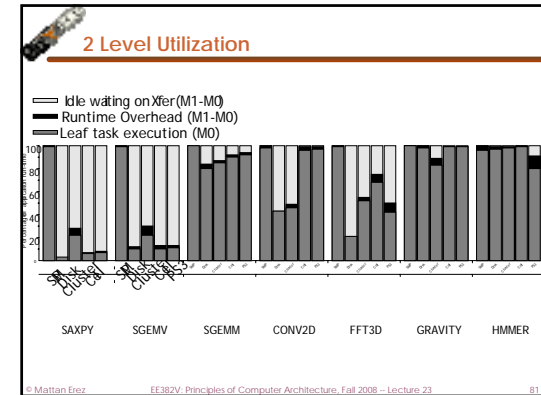
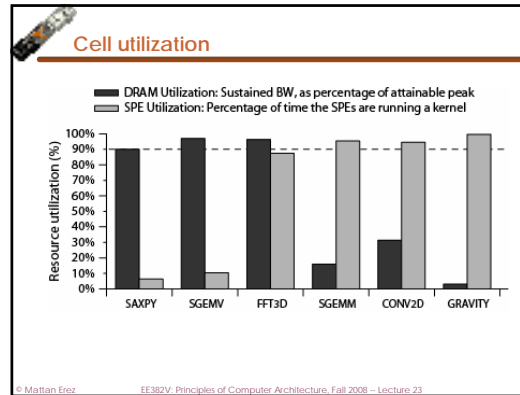
- Performance competitive with native code
 - Generic optimization for hierarchical bulk programs
 - Copy elimination
 - DMA transfer coalescing
 - Operation hoisting
 - Array allocation / packing
 - Scheduling (tasks and DMAs)
 - Automatic tuning for performance
- Portable: no source-code changes for different configurations
 - Cell, SMP, Cluster, Disk
 - Compositions of above
 - Automatic tuning
- Maximizes resources (compute or communication)
- Low overhead

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 23 78


Results – Horizontal portability - GFlop/s

| | Scalar | SMP | Disk | Cluster | Cell | PS3 |
|---------|--------|-----|-------|---------|------|------|
| SAXPY | 0.3 | 0.7 | 0.007 | 1.4 | 3.5 | 3.1 |
| SGEMV | 1.1 | 1.7 | 0.04 | 3.8 | 12 | 10 |
| SGEMM | 6.9 | 45 | 5.5 | 91 | 119 | 94 |
| CONV2D | 1.9 | 7.8 | 0.6 | 24 | 85 | 62 |
| FFT3D | 1.5 | 7.8 | 0.1 | 7.5 | 54 | 31* |
| GRAVITY | 4.8 | 40 | 3.7 | 68 | 97 | 71 |
| HMMER | 0.9 | 11 | 0.9 | 12 | 12 | 7.1* |

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 23 79




- ### Sequoia limitations
- Require explicit declaration of working sets
 - Programmer must know what to transfer
 - Some irregular applications present problems
 - Task mapping somewhat laborious
 - Autotuning helps
 - Understand which parts can be automated better
- © Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 – Lecture 23 84



Sequoia summary

- Enforce structuring already required for performance as integral part of programming model
- Make these hand optimizations portable and easier to perform

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 85



Sequoia summary (<http://sequoia.stanford.edu>)

- Problem:
 - Deep memory hierarchies pose perf. programming challenge
 - Memory hierarchy different for different machines
- Solution: Abstract hierarchical memory in programming model
 - Program the memory hierarchy explicitly
 - Expose properties that effect performance
- Approach: Express hierarchies of tasks
 - Execute in local address space
 - Call-by-value-result semantics exposes communication
 - Parameterized for portability

© Mattan Erez EE382V: Principles of Computer Architecture, Fall 2008 - Lecture 23 86