# Autotuning for Petascale:
# An Architect's Perspective

Mattan Erez

UT ECE

The University of Texas at Austin

**CScADS Autotuning Workshop**

July 8, 2008

Snowbird, Utah

# Outline

- ## What should we be tuning for?
  - Performance isn't everything
  - Tune anything that's important
- ## How should the programmer/user interact with the auto-tuner and software system?
  - Libraries aren't enough
    - Some programmers are always trying to be clever
  - Language should express what's important – including tuning
    - Too many choices and too many platforms
- ## Recent architecture research trend: fairness
  - Heterogeneous Multicore

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Tune for Utility/Cost – not Performance

Building systems is all about the bottom line

# Machine Cost Factors

- Acquisition ~$50M
  - Peak Processing
  - Peak Bandwidth
  - Peak Memory/Storage
  - **Reliability**
  - Usability
  - Facilities (power)
- Operation ~$5M/yr
  - **Power**
  - Maintenance/Administration
- Optimize total work for total cost
  - **Maximizing task performance doesn't always do that**

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Fault Tolerance == Opportunity Cost

- Reliability is an increasing concern
  - Not just memory any more
  - Logic increasingly susceptible to soft errors
  - Smaller dimension more sensitive to radiation
  - Process variation is on the rise
- Reliability requires redundancy
- "Non-stop" hardware is too costly
  - We are using unreliable systems!
- What reliability options do we apply and when?
  - Algorithmic based fault tolerance
  - Assertions
  - Computation duplication
  - Hardware features occasionally
  - Checkpoint granularity and footprint

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Power is the Dominant Architectural Problem

- Bad news: power scaling is slowing down
  - Can't scale Vt much in order to control leakage
    - New technology helps
  - → can't scale Vdd as much
  - → power doesn't go down as it used to
- Energy/device decreases slower than devices/chip
- Power goes up if performance scaling continues
  - For same processor architecture
- Roadrunner: 1PFLOP/2MW, BG/L 0.5PFLOP/2MW
  - How much for many PFLOPS?

- More bad news: energy prices going up ☺

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING
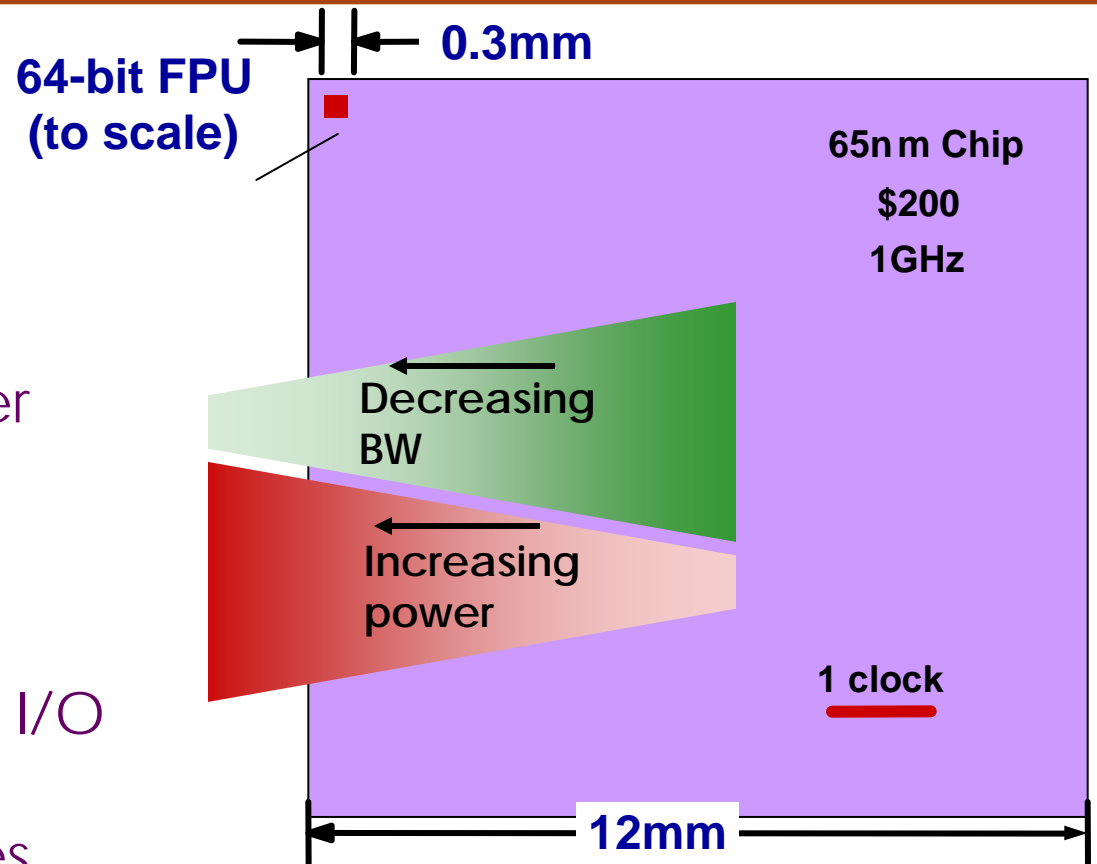
© Mattan Erez

# How Can We Reduce Power?

- Compute less
  - Use better algorithms

- Waste less
  - Don't build/use unnecessary hardware
  - No unnecessary operations
  - **No unnecessary data movement**
  - Tuning can help – minimize power per acceptable performance goal

- Specialize more
  - Specialized circuits are more efficient
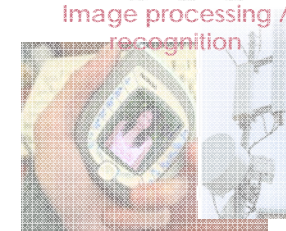  - Tuning can help decide when

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez
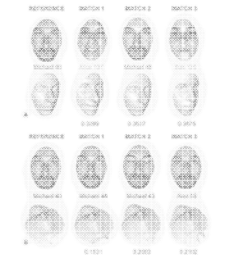
# Wasting Less – Effective Performance in VLSI
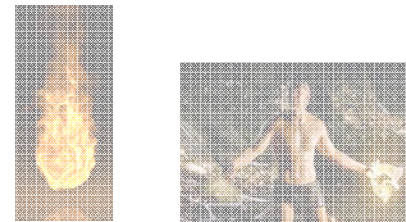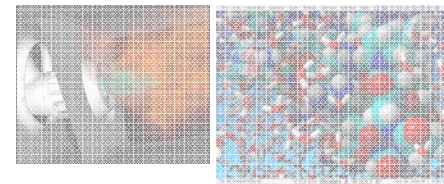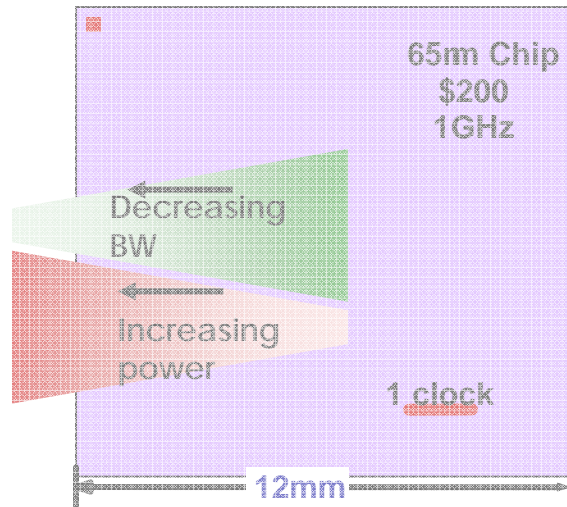
- Parallelism
  - 10s of FPUs per chip
  - Efficient control
- Locality
  - Locality lowers power
  - Reuse reduces global BW
- Throughput Design
  - Throughput oriented I/O
  - Tolerate Increasing on-/off-chip latencies
- Minimum control overhead

**64-bit FPU (to scale)**

0.3mm

65nm Chip
$200
1GHz

Decreasing BW

Increasing power

1 clock

12mm

**Parallelism, locality, latency tolerance, bandwidth, and efficient control**

# The Streaming Concept: Match Software with VLSI Strengths



65mm Chip
$200
1GHz

Decreasing BW

Increasing power

1 clock

12mm

Scientific

Graphics

Image processing / recognition
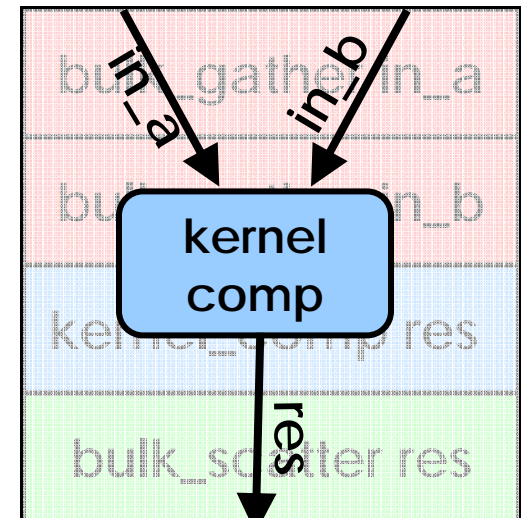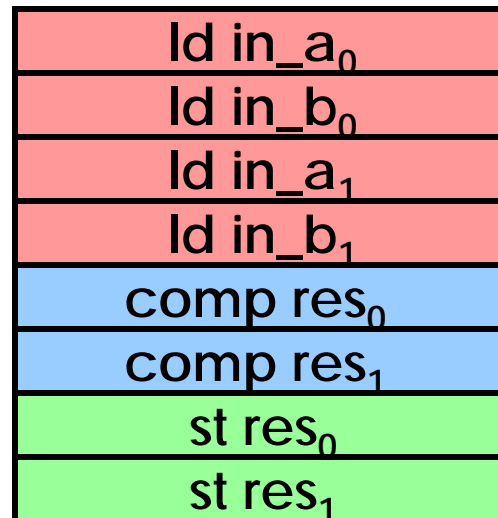
Signal processing / embedded

- Hardware matches VLSI strengths
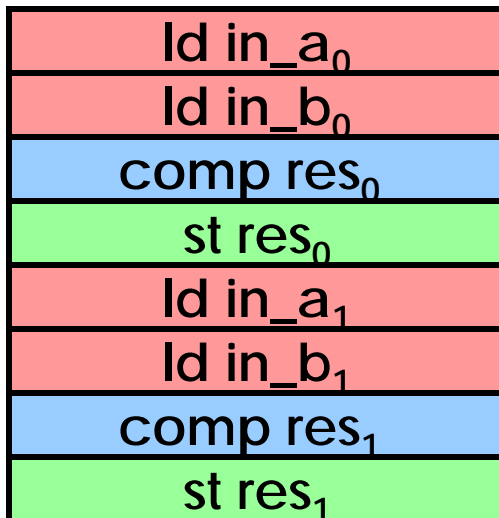  - Throughput-oriented design
  - Parallelism, locality, and partitioning
  - Hierarchical control
  - Minimalistic HW scheduling and allocation

- Software **given** more explicit control
  - Explicit hierarchical scheduling and latency hiding
  - Explicit parallelism
  - Explicit locality management

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Take Advantage of Software: Hierarchical Bulk Operations

- Data access determinable **well in advance** of data use
  - Latency hiding
  - Blocking

- Reformulate to *gather – compute – scatter*
  - Block phases into *bulk operations*



| |
|---|
| ld in_a$_0$ |
| ld in_b$_0$ |
| comp res$_0$ |
| st res$_0$ |
| ld in_a$_1$ |
| ld in_b$_1$ |
| comp res$_1$ |
| st res$_1$ |

| |
|---|
| ld in_a$_0$ |
| ld in_b$_0$ |
| ld in_a$_1$ |
| ld in_b$_1$ |
| comp res$_0$ |
| comp res$_1$ |
| st res$_0$ |
| st res$_1$ |

bulk_gather in_a
in_a
in_b
bulk_gather in_b

**kernel comp**

kernel_comp res

bulk_scatter res

res

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
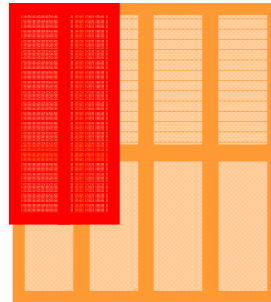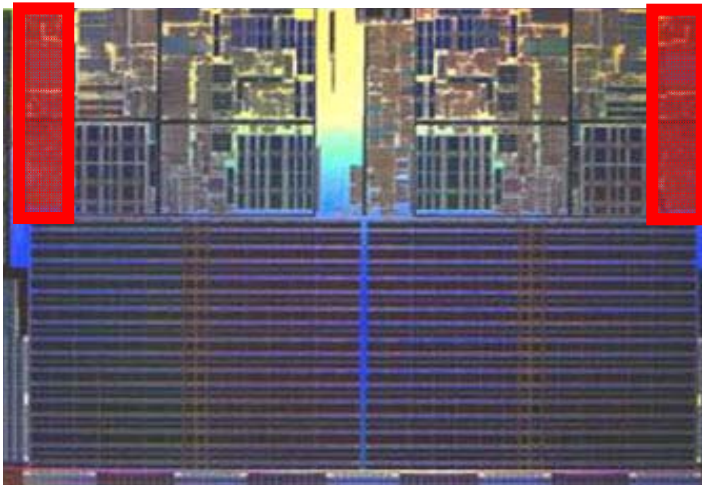ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Bulk Operations Increase Performance

## AMD dual-core Opteron
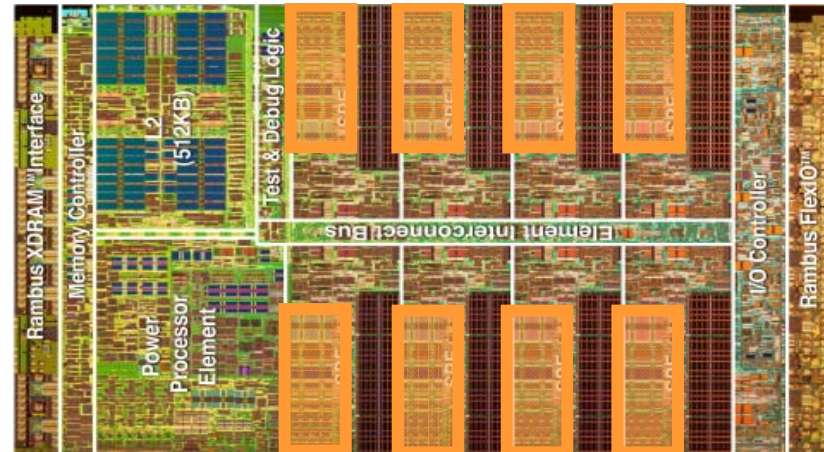90nm  |  ~200 mm$^2$  | ~100 W
### ~20 GFLOPS

## STI CELL processor
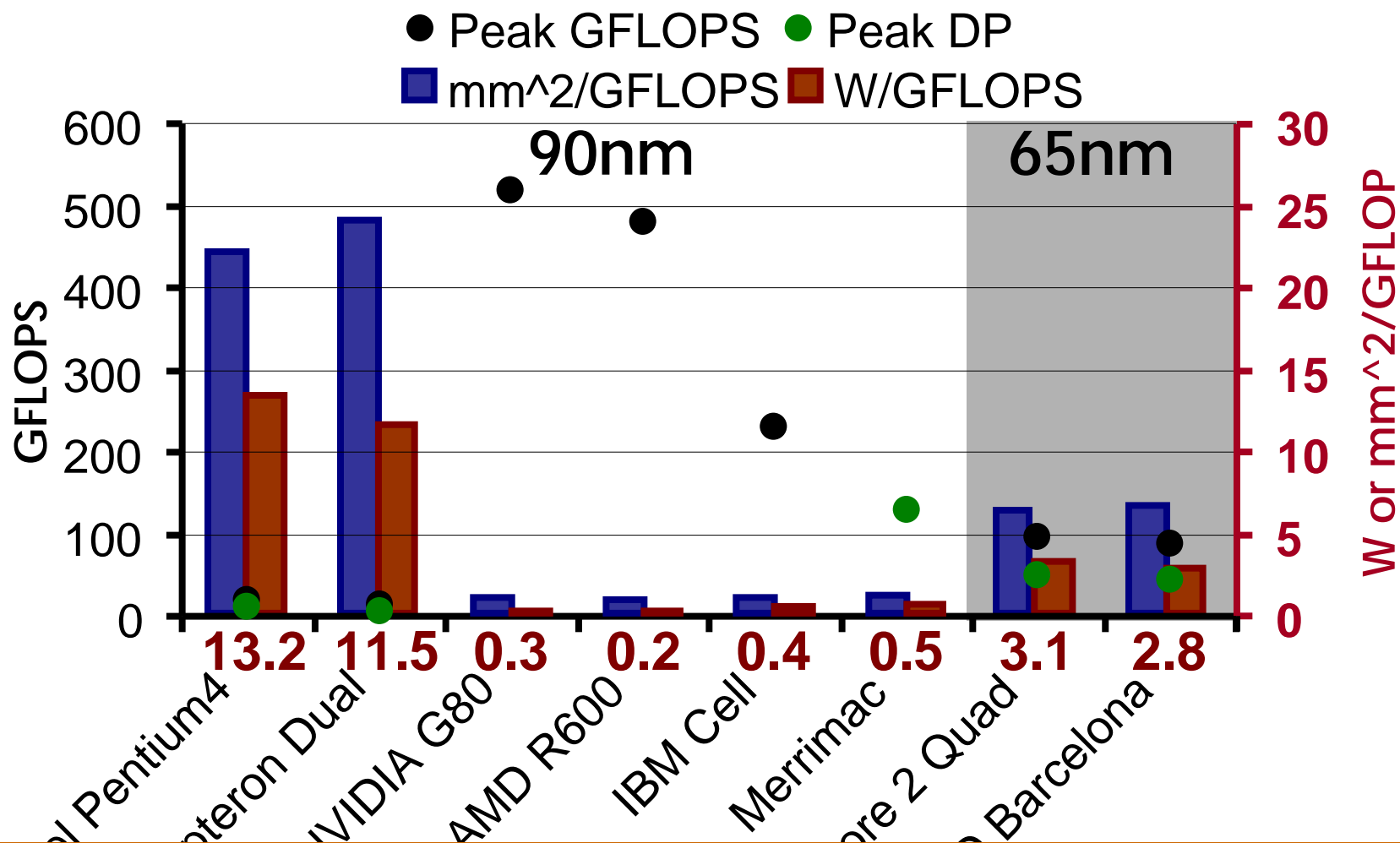90nm  |  ~220 mm$^2$  | ~100 W
### ~200 GFLOPS

FPUs

## Much more significant resources devoted to FPUs

# Bulk Operations Achieve Efficiency and Performance



Even partial adoption of bulk operations has huge impact on performance and efficiency

# Major Success but Not Enough

- Cell is ~1.5X BlueGene (based on Top500)
  - Merrimac estimates were ~6X better (in same tech node)
  - Still not enough for true Petascale
- Use better algorithms – often irregular
- Truly dynamic and irregular algorithms are challenging for bulk/streaming architectures
  - Beg for some degree of threading and caching
  - Hybrid bulk/thread architectures and models
- More work on memory systems
  - Granularity is a problem
- On-chip interconnection networks – no clear winner

**Locality, Parallelism, and Hierarchy throughout the system**

# Tuning for Power

- Need to co-search for power and performance
  - Optimize cost, not performance
  - Opportunity cost too (fault tolerance)
- Maximize locality / minimize data movement
  - Power impacted significantly by interconnect and memory
- Try to specialize
  - Utilize control hierarchy
  - Utilize specialized hardware
- Minimize waste
  - Strong interactions with load balancing
  - Processor/memory dynamic power management is key

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Languages Need to Abstractly Expose Important Factors and Tuning

- How should the programmer/user interact with the auto-tuner and software system?
  - Libraries aren't enough
    - Some programmers are always trying to be clever
  - Language should express what's important – including tuning
    - Too many choices and too many platforms

# Sequoia: Abstract Streaming/Bulk Programming

- Facilitate development of hierarchy-aware stream programs …

- … that remain portable across machines

- Provide constructs that can be implemented efficiently without requiring advanced compiler technology

  – Place computation and data in machine

  – Explicit parallelism and communication

  – Large bulk transfers

- Facilitate tuning

  – Decouple algorithm and decomposition from setting parameters

  – Sequoia language only expresses strategy

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Hierarchical memory

- Abstract machines as trees of memories



Similar to:

Parallel Memory Hierarchy Model (Alpern et al.)

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Hierarchical memory

- Abstract machines as trees of memories

# Hierarchical memory

# Sequoia tasks

- Special functions called **tasks** are the building blocks of Sequoia programs

```
task interpolate(in  float A[N],
                 in  float B[N],
                 in  float u,
                 out float result[N])
{
   for (int i=0; i<N; i++)
      result[i] = u * A[i] + (1-u) * B[i];
}
```

- Task arguments can be arrays and scalars
- Tasks arguments located within a single level of abstract memory hierarchy

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Sequoia tasks

- Single abstraction for
  - Isolation / parallelism
  - Explicit communication / working sets
  - Expressing locality

- Tasks operate on arrays, not array elements

- Tasks nest:  they call subtasks

# The Streaming Concept:
## Match Software with VLSI Strengths

65nm Chip
$200
1GHz

Decreasing BW

Increasing power

1 clock

12mm

Scientific

Image processing / recognition

Graphics

Signal processing / embedded

- Hardware matches VLSI strengths
  - Throughput-oriented design
  - Parallelism, locality, and partitioning
  - Hierarchical control
  - Minimalistic HW scheduling and allocation

- Software **given** more explicit control
  - Explicit hierarchical scheduling and latency hiding
  - Explicit parallelism
  - Explicit locality management

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Example: dense matrix multiplication

Task:
1024x1024
matrix multiplication

**Main memory**

Task:
256x256
matrix mult

Task:
256x256
matrix mult

… 64 total
subtasks …

Task:
256x256
matrix mult

**L2 cache**

Task:
32x32
matrix mult

Task:
32x32
matrix mult

… 512 total
subtasks …

Task:
32x32
matrix mult

**L1 cache**

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Example - task isolation

```
task matmul::inner(in    float A[M][T],
                   in    float B[T][N],
                   inout float C[M][N])
{


}
```

- Task arguments + local variables define working set

THE UNIVERSITY OF TEXAS AT AUSTIN

UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Example - parameterization

```
task matmul::inner(in    float A[M][T],
                   in    float B[T][N],
                   inout float C[M][N])
{
   tunable int P, Q, R;




}
```

- Tasks are written in parameterized form for portability

- Different "variants" of the same task can be defined

```
task matmul::leaf(in    float A[M][T],
                  in    float B[T][N],
                  inout float C[M][N])
{
   for (int i=0; i<M; i++)
     for (int j=0; j<N; j++)
       for (int k=0; k<T; k++)
         C[i][j] += A[i][k] * B[k][j];
}
```

© Mattan Erez

# Example - locality & communication

```
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]);

    }
  }
}


task matmul::leaf(in     float A[M][T],
                  in     float B[T][N],
                  inout float C[M][N])
{
   for (int i=0; i<M; i++)
     for (int j=0; j<N; j++)
       for (int k=0;k<T; k++)
         C[i][j] += A[i][k] * B[k][j];
}
```

- Working set resident within single level of hierarchy

- Passing arguments to subtasks is only way to specify communication in Sequoia

# Specializing matmul

- Instances of tasks placed at each memory level
  - Instances define a task variant and values for all parameters

matmul::inner
M=N=T=1024
P=Q=R=256

**Main memory**

matmul::
inner
M=N=T=256
P=Q=R=32

matmul::
inner
M=N=T=256
P=Q=R=32

… 64 total
subtasks …

matmul::
inner
M=N=T=256
P=Q=R=32

**L2 cache**

matmul::leaf
M=N=T=32

matmul::leaf
M=N=T=32

… 512 total
subtasks …

matmul::leaf
M=N=T=32

**L1 cache**

THE UNIVERSITY OF TEXAS AT AUSTIN
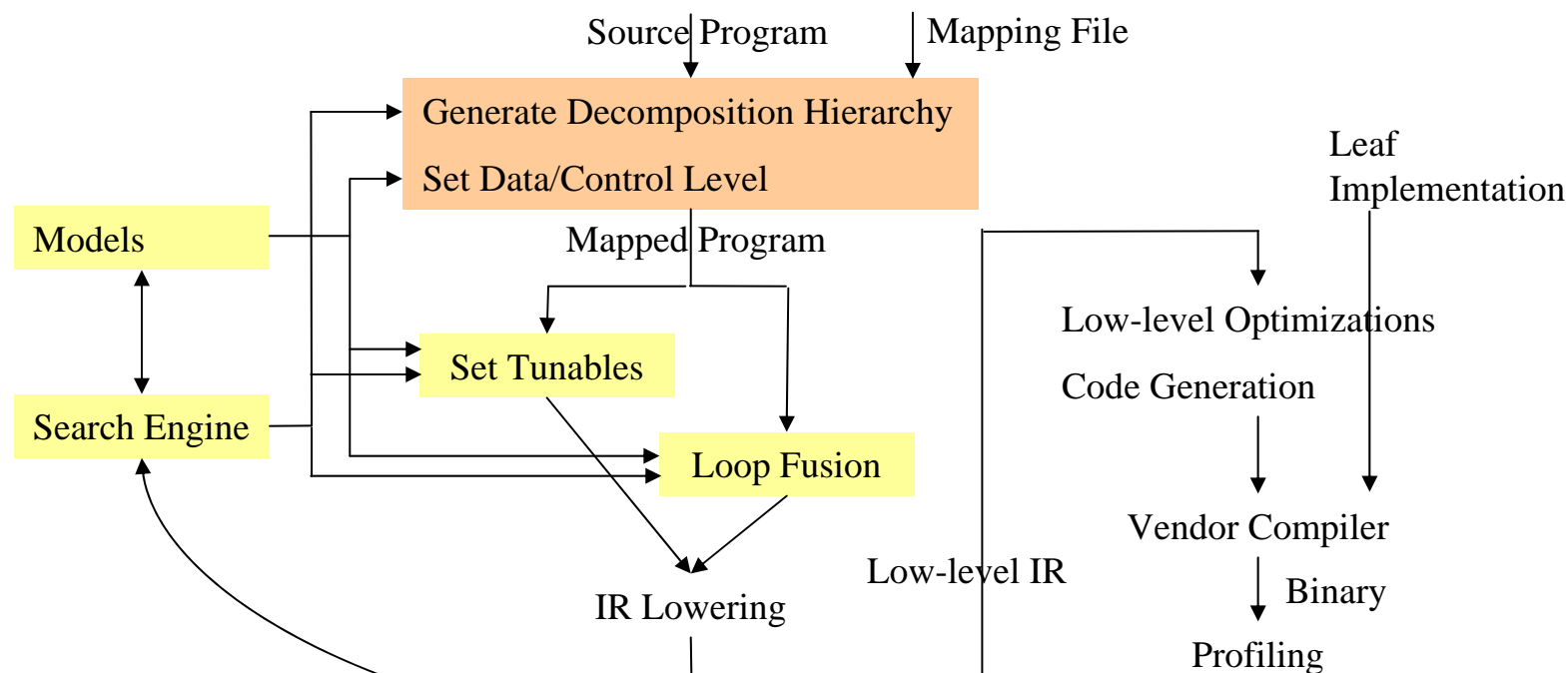UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# 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

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
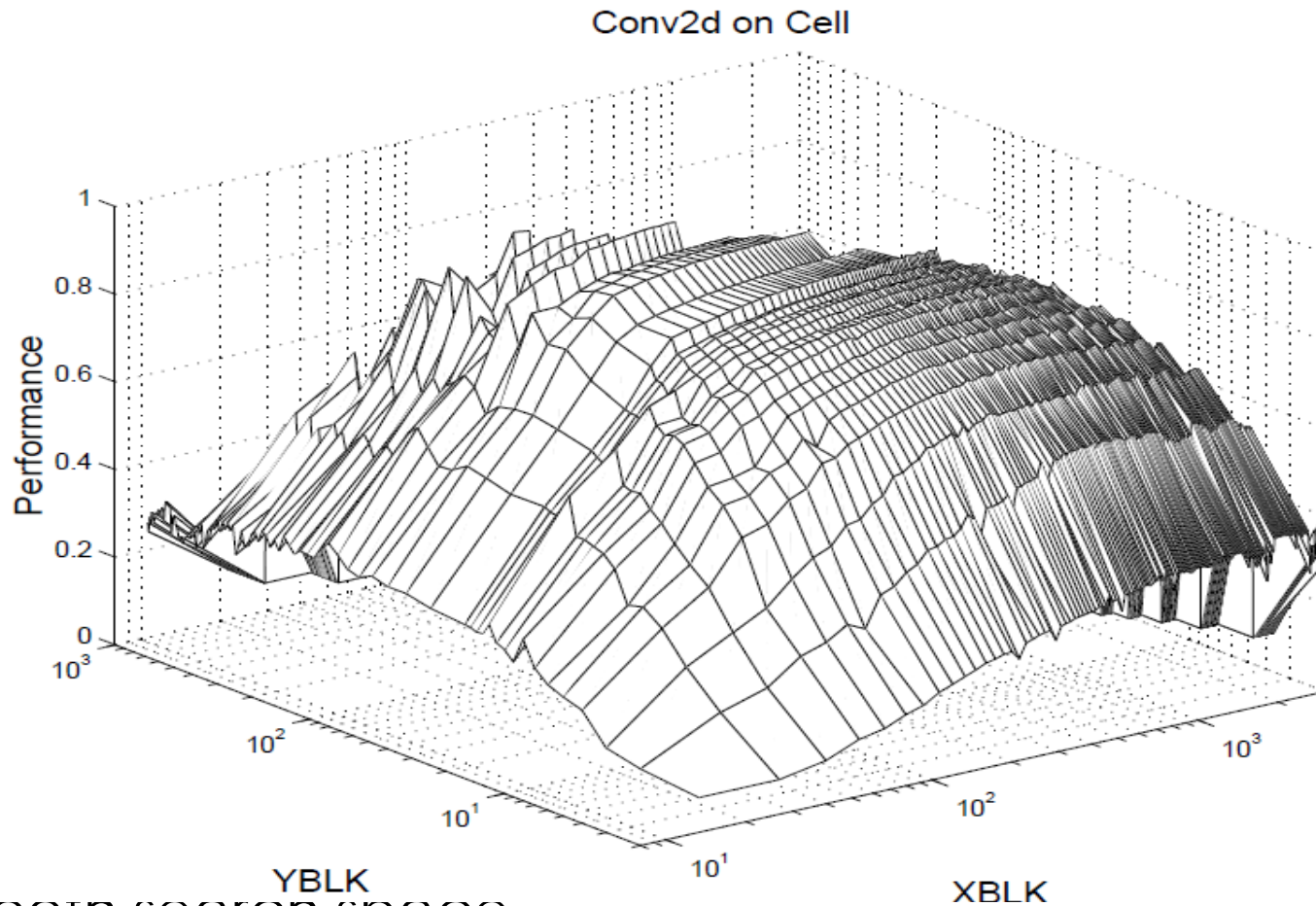ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# 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

Source Program          Mapping File

Generate Decomposition Hierarchy

Set Data/Control Level

Leaf Implementation

Models

Mapped Program

Low-level Optimizations

Code Generation

Set Tunables

Search Engine

Loop Fusion

Low-level IR

Vendor Compiler

Binary

IR Lowering

Profiling

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
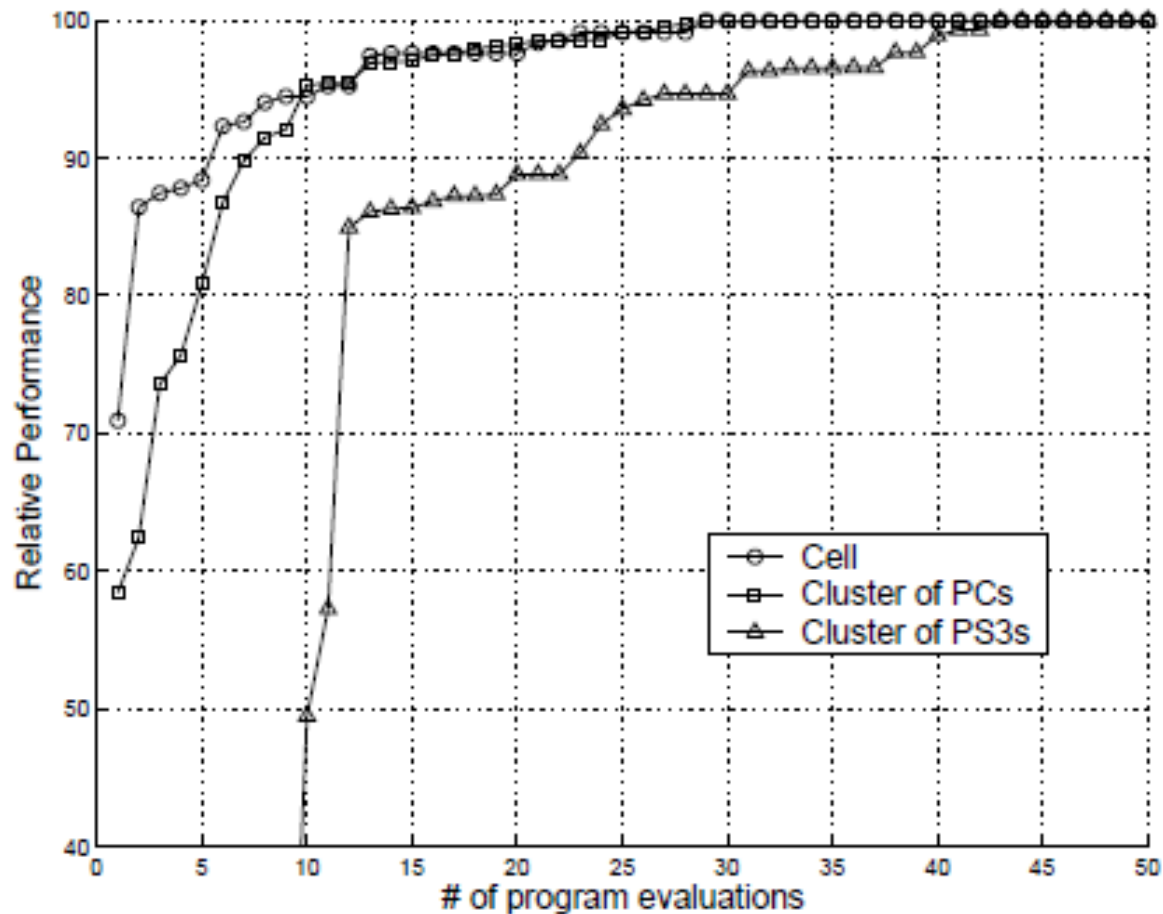ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

Conv2d on Cell

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

# Guided Search Converges Quickly

- Smoothness leads to quick convergence

# Autotuning Out Performs Programmer

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

CScADS 2008 Autotuning Workshop:
An Architect's Perspective

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING

© Mattan Erez

# Architecture Trend: Fairness in Multicore/Multi-threaded Processors

Hardware balances shared resources

# Maintain Overall Performance through Fair Partitioning of Shared Resources

- Motivating applications: multiprogramming
- Shared cache
  - Allocate partitions of ways in a set-associative cache to threads
  - Prevent low-locality thread from evicting useful data

- Shared memory bandwidth
  - Schedule memory operations from different threads fairly

- Definition of fairness?
  - All threads suffer performance degradation relative to running in isolation

# Conclusions

- Autotuning should match architecture optimizations – **maximum utility/cost**
  - Maximize locality / minimize communication
  - Take advantage of control hierarchy
  - Specialized hardware units
  - Reliability is another opportunity

- Languages should expose what's important (in an abstract portable way)
  - Expose tuning – it's an essential part of the software system
  - Sequoia is one early attempt

THE UNIVERSITY OF TEXAS AT AUSTIN
UT ECE
ELECTRICAL & COMPUTER ENGINEERING