



Parallel Vector Tile-Optimized Library (PVTOL) Architecture

**Jeremy Kepner, Nadya Bliss, Bob Bond, James Daly, Ryan
Haney, Hahn Kim, Matthew Marzilli, Sanjeev Mohindra,
Edward Rutledge, Sharon Sacco, Glenn Schrader**

MIT Lincoln Laboratory

May 2007

This work is sponsored by the Department of the Air Force under Air Force contract FA8721-05-C-0002. Opinions, interpretations, conclusions and recommendations are those of the author and are not necessarily endorsed by the United States Government.

MIT Lincoln Laboratory

PVTOL-1
6/23/07

Title slide



Outline

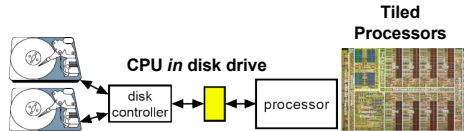
- ➔ • **Introduction**
- **PVTOL Machine Independent Architecture**
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- **PVTOL on Cell**
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- **Summary**

Outline



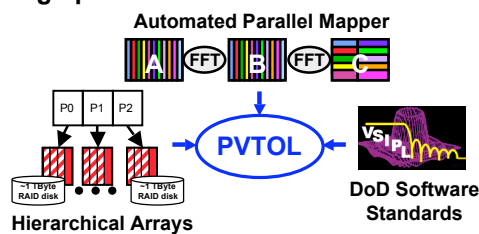
PVTOL Effort Overview

Goal: Prototype advanced software technologies to exploit novel processors for DoD sensors



- Have *demonstrated* 10x performance benefit of tiled processors
- Novel storage should provide 10x more IO

Approach: Develop **Parallel Vector Tile Optimizing Library (PVTOL)** for high performance and ease-of-use



PVTOL-3
6/23/07

DoD Relevance: Essential for flexible, programmable sensors with large IO and processing requirements



- Wide area data
- Collected over many time scales

Mission Impact:

- Enabler for next-generation synoptic, multi-temporal sensor systems

Technology Transition Plan

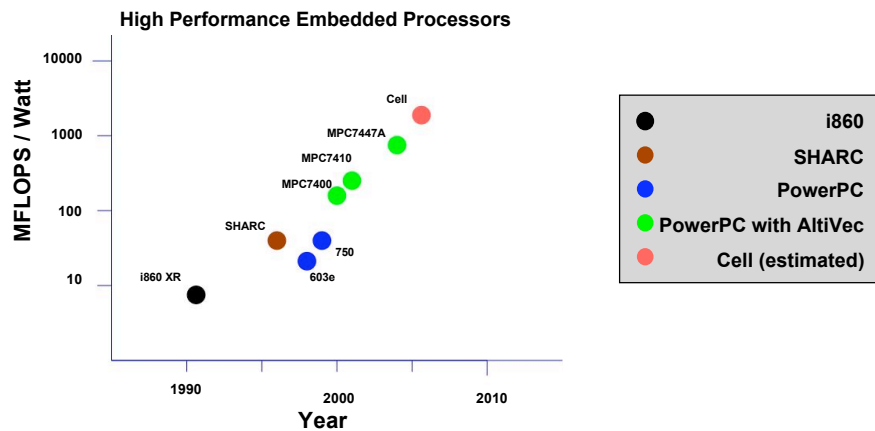
- Coordinate development with sensor programs
- Work with DoD and Industry standards bodies

MIT Lincoln Laboratory

The Parallel Vector Tile Optimizing Library (PVTOL) is an effort to develop a new processor architecture for signal processing that exploits the recent shifts in technology to tiled multi-core chips and more tightly integrated mass storage. These technologies are critical for processing the higher bandwidth and longer duration data produced required by synoptic, multi-temporal sensor systems.. The principal challenge in exploiting the new processing architectures for these missions is writing the software in a flexible manner that is portable and delivers high performance. The core technology of this project is the Parallel Vector Tile Optimizing Library (PVTOL), which will use recent advances in automating parallel mapping technology, hierarchical parallel arrays, combined with the Vector, Signal and Image Processing Library (VSIPL) open standard to deliver portable performance on tiled processors.



Embedded Processor Evolution



- 20 years of exponential growth in FLOPS / Watt
- Requires switching architectures every ~5 years
- Cell processor is current high performance architecture

PVTOL-4
6/23/07

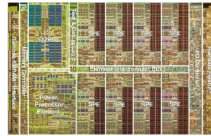
MIT Lincoln Laboratory

Growth in embedded processor performance, in terms of FLOPS/Watt, has grown exponentially over the last 20 years. No single processing architecture has dominated over this period, hence in order to leverage this increase in performance, embedded system designers must switch processing architectures approximately every 5 years.

MFLOPS / W for i860, SHARC, 603e, 750, 7400, and 7410 are extrapolated from board wattage. They also include other hardware energy use such as memory, memory controllers, etc. 7447A and the Cell estimate are for the chip only. Effective FLOPS for all processors are based on 1024 FFT timings. Cell estimate uses hand coded TDFIR timings for effective FLOPS.



Cell Broadband Engine



- Cell was designed by IBM, Sony and Toshiba
- Asymmetric multicore processor
 - 1 PowerPC core + 8 SIMD cores



- Playstation 3 uses Cell as main processor



- Provides Cell-based computer systems for high-performance applications

PVTOL-5
6/23/07

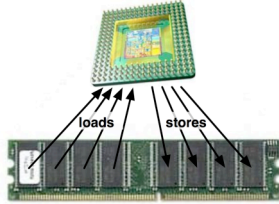
MIT Lincoln Laboratory

The current high performance processing architecture is the Cell processor, designed by a collaboration between IBM, Sony and Toshiba. While the Cell was initially targeted for Sony's Playstation 3, it has use in wide range of applications, including defense, medical, etc. Mercury Computer Systems provides Cell-based systems that can be used to develop high-performance embedded systems.



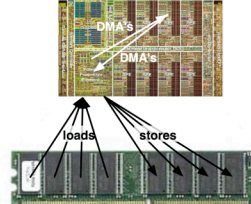
Multicore Programming Challenge

Past Programming Model: Von Neumann



- Great success of Moore's Law era
 - Simple model: load, op, store
 - Many transistors devoted to delivering this model
- Moore's Law is ending
 - Need transistors for performance

Future Programming Model: ???



- Processor topology includes:
 - Registers, cache, local memory, remote memory, disk
- Cell has *multiple* programming models

Increased performance at the cost of exposing complexity to the programmer

PVTOL-6
6/23/07

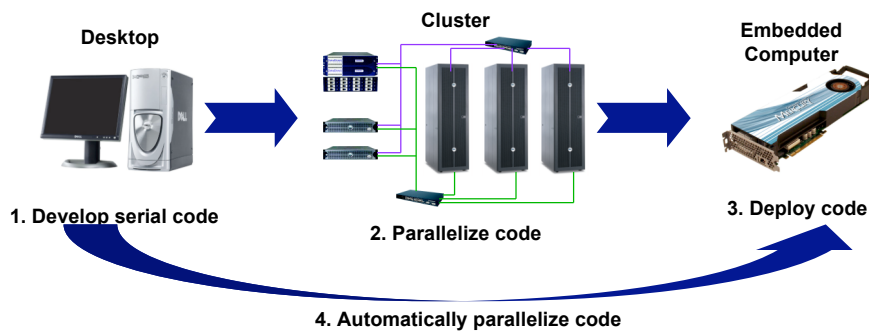
MIT Lincoln Laboratory

For decades, Moore's Law has enabled ever faster processors that have supported the traditional von Neumann programming model, i.e. load data from memory, process, then save the results to memory. As clock speeds near 4 GHz, physical limitations in transistor size are leading designers to build more processor cores (or "tiles") on each chip rather than faster processors. Multicore processors improve raw performance but expose the underlying processor and memory topologies. This results in increased programming complexity, i.e. the loss of the von Neumann programming model.



Parallel Vector Tile-Optimized Library (PVTOL)

- PVTOL is a portable and scalable middleware library for multicore processors
- Enables incremental development



Make parallel programming as easy as serial programming

PVTOL-7
6/23/07

MIT Lincoln Laboratory

PVTOL is focused on addressing the programming complexity of associated with emerging “Topological Processors”. Topological Processors require the programmer to understand the physical topology of the chip to get high efficiency. There are many such processors emerging into the market. The Cell processor is an important example of such a chip. The current PVTOL effort is focused on getting high performance from the Cell processor on signal and image processing applications. The PVTOL interface is designed to address a wide range of processors including multicore and FPGAs.

PVTOL enables software developers to develop high-performance signal processing application on a desktop computer, parallelize the code on commodity clusters, then deploy the code onto an embedded computer, with minimal changes to the code. PVTOL also includes automated mapping technology that will automatically parallelize the application for a given platform. Applications developed on a workstation can then be deployed on an embedded computer and the library will parallelize the application without any changes to the code.



PVTOL Development Process

Serial PVTOL code

```
void main(int argc, char *argv[]) {  
    // Initialize PVTOL  
    process pvtol(argc, argv);  
  
    // Create input, weights, and output matrices  
    typedef Dense<2, float, tuple<0, 1> > dense_block_t;  
    typedef Matrix<float, dense_block_t, LocalMap> matrix_t;  
    matrix_t input(num_vects, len_vect),  
             filts(num_vects, len_vect),  
             output(num_vects, len_vect);  
  
    // Initialize arrays  
    ...  
  
    // Perform TDFIR filter  
    output = tdfir(input, filts);  
}
```

This slide shows an example of serial PVTOL code that allocates data to be processed by a time-domain FIR filter. The user simply allocates three matrices containing the input and output data and the filter coefficients. The LocalMap tag in the Matrix type definition indicates that the matrices will be allocated all on the processor's local memory.



PVTOL Development Process

Parallel PVTOL code

```
void main(int argc, char *argv[]) {  
    // Initialize PVTOL  
    process pvtol(argc, argv);  
  
    // Add parallel map  
    RuntimeMap map1(...);  
  
    // Create input, weights, and output matrices  
    typedef Dense<2, float, tuple<0, 1> > dense_block_t;  
    typedef Matrix<float, dense_block_t, RuntimeMap> matrix_t;  
    matrix_t input(num_vects, len_vect, map1),  
             filts(num_vects, len_vect, map1),  
             output(num_vects, len_vect, map1);  
  
    // Initialize arrays  
    ...  
  
    // Perform TDFIR filter  
    output = tdfir(input, filts);  
}
```

This slide shows how to parallelize the serial code shown in the previous slide. The programmer creates a map, which contains a concise description of how to parallelize a matrix across multiple processors. The RuntimeMap in the Matrix type definition indicates that the map for the matrices is constructed at runtime. The map object is then passed into the matrix constructors, which allocate memory on multiple processors as described in the map object.



PVTOL Development Process

Embedded PVTOL code

```
void main(int argc, char *argv[]) {  
    // Initialize PVTOL  
    process pvtol(argc, argv);  
  
    // Add hierarchical map  
    RuntimeMap map2(...);  
  
    // Add parallel map  
    RuntimeMap map1(..., map2);  
  
    // Create input, weights, and output matrices  
    typedef Dense<2, float, tuple<0, 1> > dense_block_t;  
    typedef Matrix<float, dense_block_t, RunTimeMap> matrix_t;  
    matrix_t input(num_vects, len_vect, map1),  
             filts(num_vects, len_vect, map1),  
             output(num_vects, len_vect, map1);  
  
    // Initialize arrays  
    ...  
  
    // Perform TDFIR filter  
    output = tdfir(input, filts);  
}
```

This slide shows how to deploy the parallel code from the previous slide onto an embedded multicore system. As before, map1 describes how to parallelize matrices across multiple processors. A new map, map2, is added to describe how each processor should divide its section of the matrices across cores. This allocates a “hierarchical array” across the processor hierarchy.



PVTOL Development Process

Automapped PVTOL code

```
void main(int argc, char *argv[]) {  
    // Initialize PVTOL  
    process pvtol(argc, argv);  
  
    // Create input, weights, and output matrices  
    typedef Dense<2, float, tuple<0, 1> > dense_block_t;  
    typedef Matrix<float, dense_block_t, AutoMap> matrix_t;  
    matrix_t input(num_vects, len_vect , map1),  
             filts(num_vects, len_vect , map1),  
             output(num_vects, len_vect , map1);  
  
    // Initialize arrays  
    ...  
  
    // Perform TDFIR filter  
    output = tdfir(input, filts);  
}
```

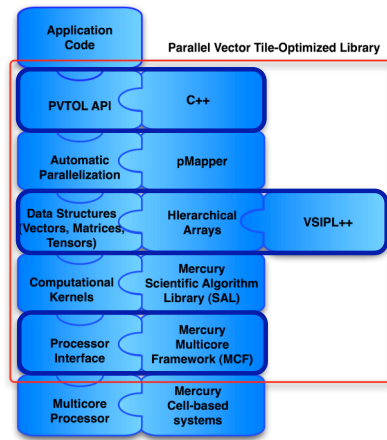
PVTOL-11
6/23/07

MIT Lincoln Laboratory

Maps give the programmer direct control over how memory should be allocated across the processor hierarchy. Alternatively, the programmer can let PVTOL automatically parallelize arrays. The only change required is to replace the LocalMap tag in the Matrix type definition with AutoMap. Note that this code can be deployed without any changes on a workstation, cluster or embedded computer.

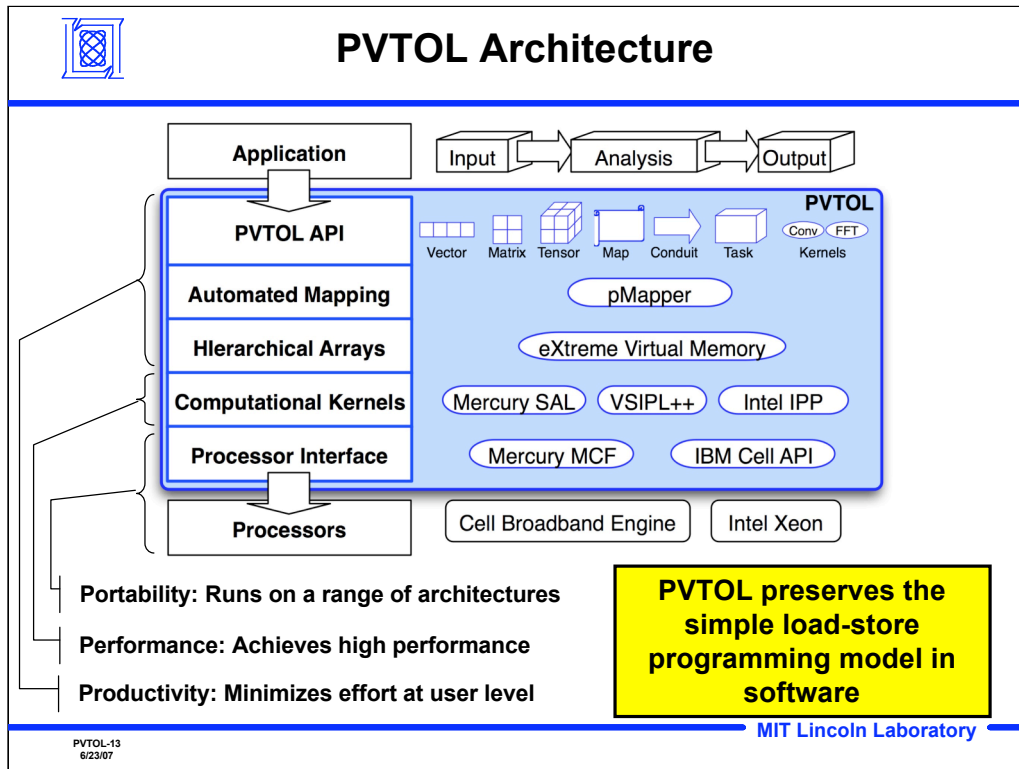


PVTOL Components



- **Performance**
 - Achieves high performance
- **Portability**
 - Built on standards, e.g. VSIPL++
- **Productivity**
 - Minimizes effort at user level

This slide shows the various software components that are used in PVTOL and how they relate to Performance, Portability and Productivity.



This slide shows a layered view of the PVTOL architecture. At the top is the application. The PVTOL API exposes high-level structures for data (e.g. vectors), data distribution (e.g. maps), communication (e.g. conduits) and computation (e.g. tasks and computational kernels). High level structures improve the productivity of the programmer. By being built on top of existing technologies, optimized for different platforms, PVTOL provides high performance. And by supporting a range of processor architectures, PVTOL applications are portable. The end result is that rather than learning new programming models for new processor technologies, PVTOL preserves the simple von Neumann programming model most programmers are used to.



Outline

- Introduction
- **PVTOL Machine Independent Architecture**
 - ➔ – Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- Summary

PVTOL-14
6/23/07

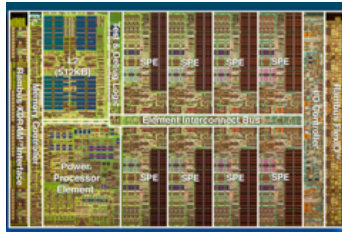
MIT Lincoln Laboratory

Outline

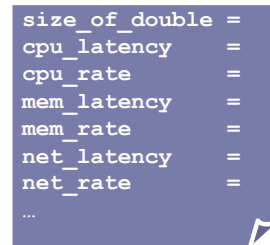


Machine Model - Why?

- Provides description of underlying hardware
- pMapper: Allows for simulation without the hardware
- PVTOL: Provides information necessary to specify map hierarchies



Hardware



Machine Model

PVTOL-15
6/23/07

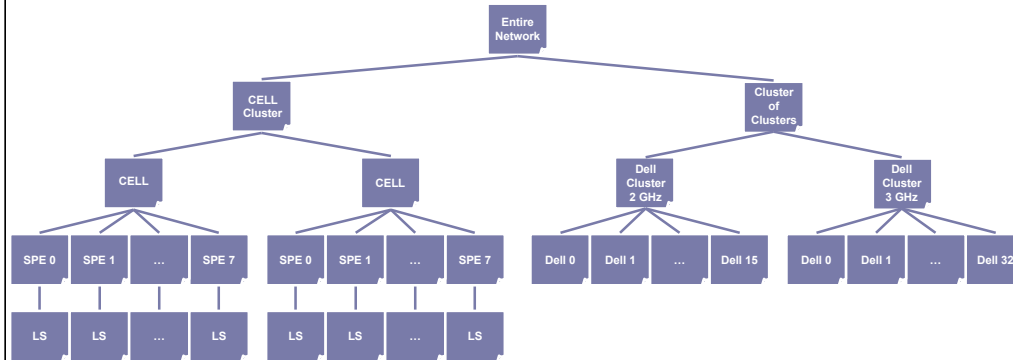
MIT Lincoln Laboratory

The development of a machine model that can describe different processing architectures has two purposes. First, a machine model allows the pMapper automated mapping environment to simulate how well different mappings of an application perform without executing the application on the actual hardware. Second, a machine model provides information about the architecture needed by PVTOL to specify how to map hierarchical arrays to the processing hierarchy.



PVTOL Machine Model

- **Requirements**
 - Provide hierarchical machine model
 - Provide heterogeneous machine model
- **Design**
 - Specify a machine model as a tree of machine models
 - Each sub tree or a node can be a machine model in its own right



PVTOL-16
6/23/07

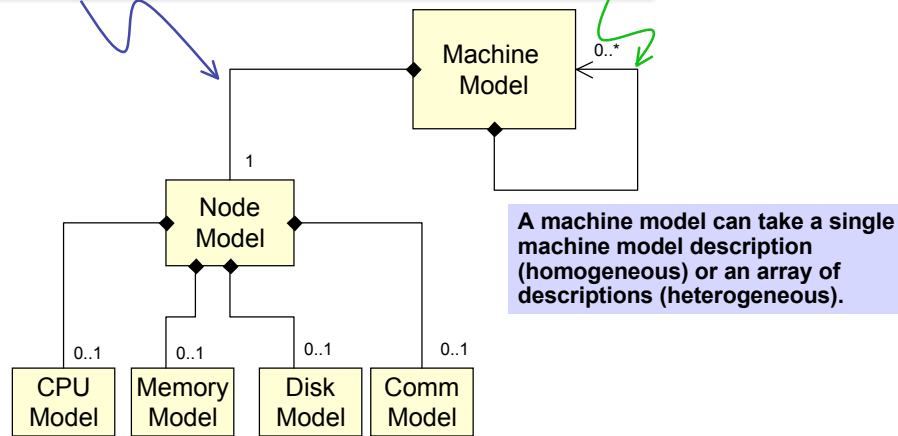
MIT Lincoln Laboratory

PVTOL requires that a machine model be able to describe hierarchical and heterogeneous processing architectures. For example, a cluster of Cell processors is a hierarchy of processors; at the top level the hierarchy consists of multiple Cell processors and at the next level each Cell processor consists of multiple SPE processors. Heterogeneity enables the description of processing architectures composed of different hardware architectures, e.g. Cell and Intel processors. The PVTOL machine model is recursive in that a machine model can be described as a tree of machine models.



Machine Model UML Diagram

A machine model constructor can consist of just node information (flat) or additional children information (hierarchical).



A machine model can take a single machine model description (homogeneous) or an array of descriptions (heterogeneous).

PVTOL machine model is different from PVL machine model in that it separates the Node (flat) and Machine (hierarchical) information.

PVTOL-17
6/23/07

MIT Lincoln Laboratory

This slide shows the UML diagram for a machine model. A machine model constructor can consist of just node information (flat) or additional children information (hierarchical). A machine model can take a single machine model description (homogeneous) or an array of descriptions (heterogeneous). The PVTOL machine model extends PVL's model in that it can describe both flat and hierarchical processing architectures.

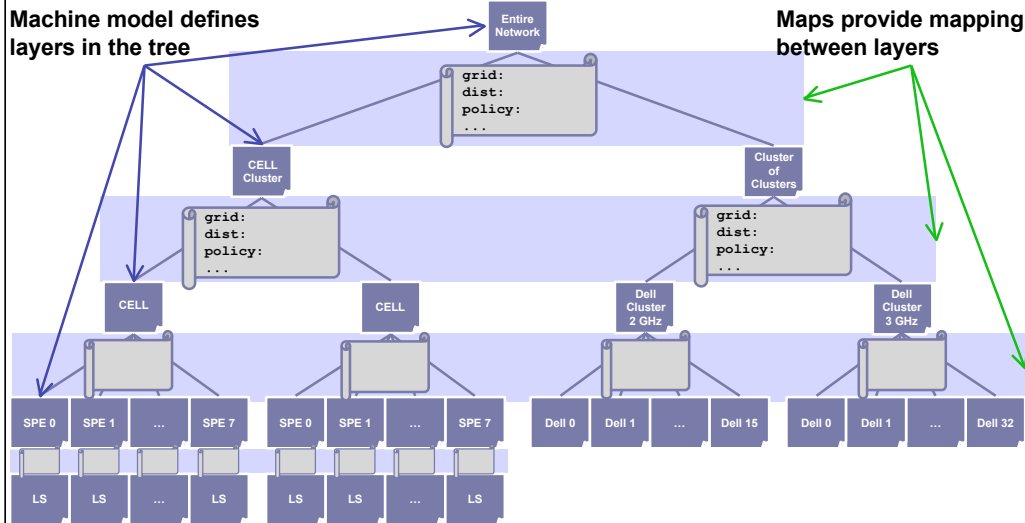


Machine Models and Maps

Machine model is tightly coupled to the maps in the application.

Machine model defines layers in the tree

Maps provide mapping between layers

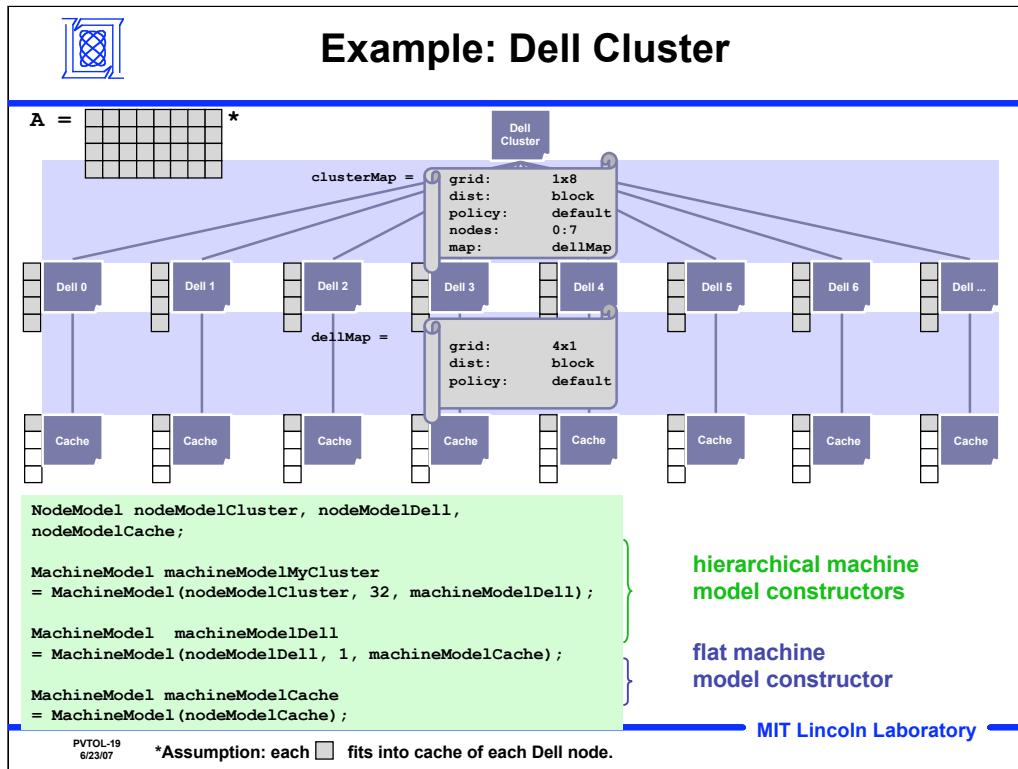


PVTOL-18
6/23/07

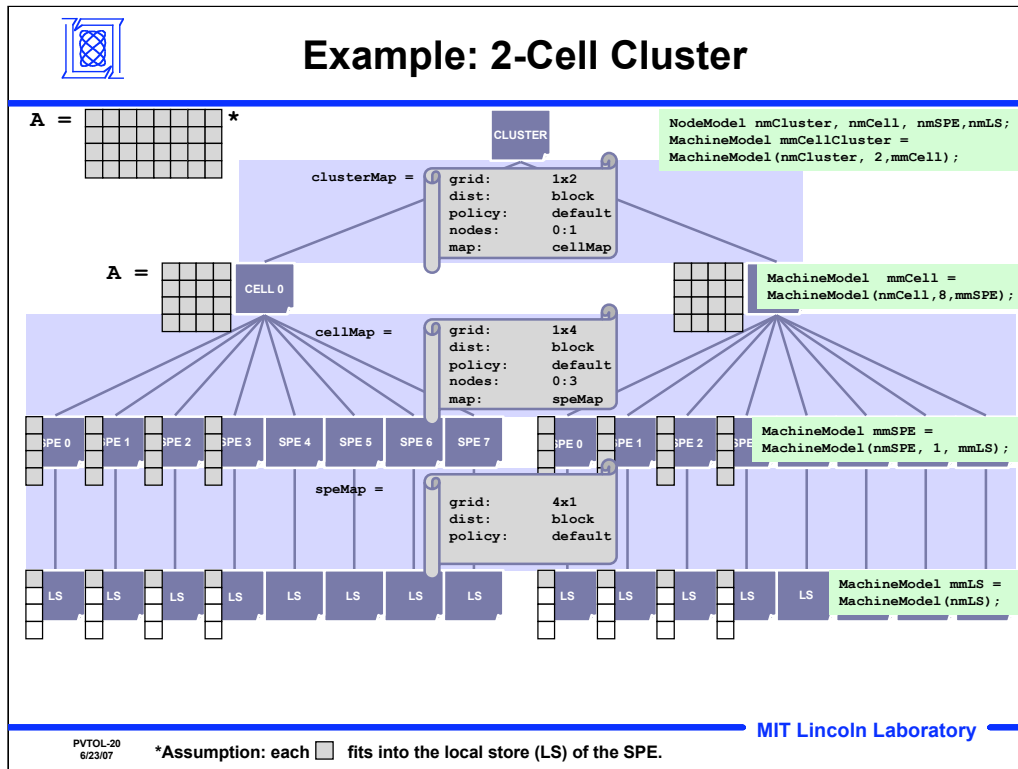
*Cell node includes main memory

MIT Lincoln Laboratory

A machine model is tightly coupled to the map. Machine model descriptions define layers in the tree, while maps provide mappings *between* layers (N layers in the machine mode could have at most $N-1$ layers of maps). The map API has to be adjusted to allow the user to pass in an array of maps for heterogeneous systems. The map API does not require the user to specify storage level. When specifying hierarchical maps, the layers specified have to start at the root and be contiguous



This shows an example of how to construct a machine model of a cluster comprised of Dell computers. This example describes a cluster of 32 Dell nodes, with each node containing a single processor and each processor containing cache. This machine model contains a two-level hierarchy. The matrix A can be mapped to this hierarchy by two maps: `clusterMap` and `dellMap`. `clusterMap` describes how to distribute data across nodes and `dellMap` describes how to partition data on each Dell node into blocks that will fit into cache.

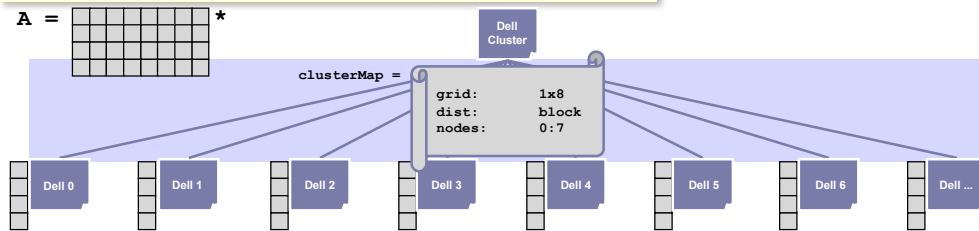


This shows an example of how to construct a machine model of a cluster comprised of Cell processors. This example describes a cluster of 2 Cell processors, with each node containing a 8 SPE's and each processor containing a local store. This machine model contains a three-level hierarchy. The matrix A can be mapped to this hierarchy by three maps: `clusterMap`, `cellMap` and `speMap`. `clusterMap` describes how to distribute data across Cell processors, `cellMap` describes how to distribute data on each Cell processor between SPE's and `speMap` describes how to partition data owned by each SPE into blocks that will fit into its local store.

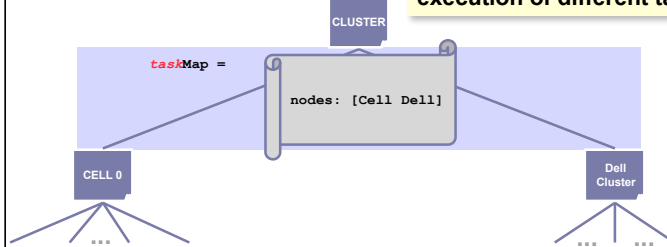


Machine Model Design Benefits

Simplest case (mapping an array onto a cluster of nodes) can be defined as in a familiar fashion (PVL, pMatlab).



Ability to define heterogeneous models allows execution of different tasks on very different systems.



PVTOL-21
6/23/07

MIT Lincoln Laboratory

This slide describes the advantages of tree machine models/maps. In the simplest case, a machine model can be easily specified in the same way as PVL and pMatlab. From the API point of view, this is expressive enough for both flat and hierarchical machine models/maps. It allows the user to use different machine models in different tasks, since task maps do not require that programmer to specify a data distribution.



Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- Summary



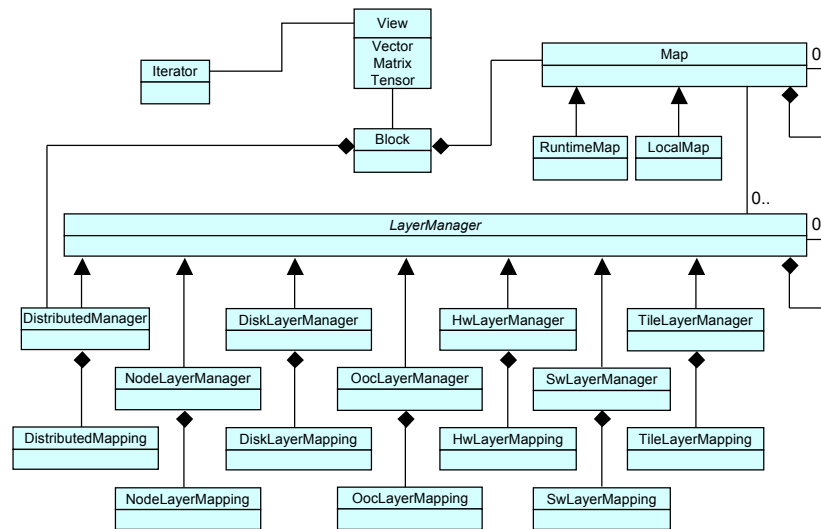
PVTOL-22
6/23/07

MIT Lincoln Laboratory

Outline



Hierarchical Arrays UML



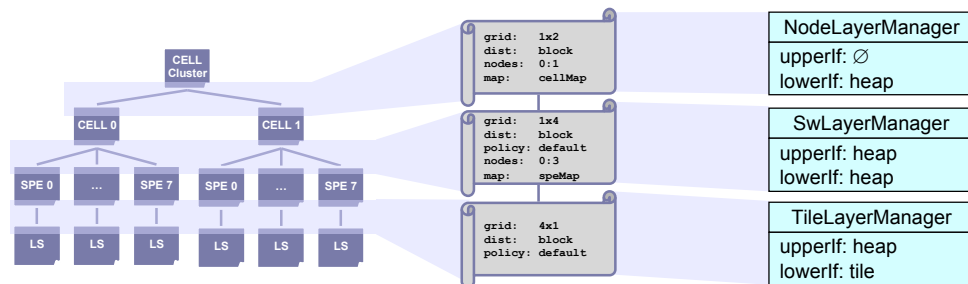
PVTOL-23
6/23/07

MIT Lincoln Laboratory

Hierarchical arrays are designed such that details of the memory hierarchy are hidden from the programmer. Programmers access data Block objects via a Vector, Matrix or Tensor View object. Iterator objects provide an architecture independent mechanism for accessing data. Map objects describe how data in a Block object is distributed across the memory hierarchy. Different layer managers are instantiated depending on which layers in the memory hierarchy data is distributed across, e.g. tile memory, local main memory, remote main memory, disk, etc.



Isomorphism



Machine model, maps, and layer managers are isomorphic

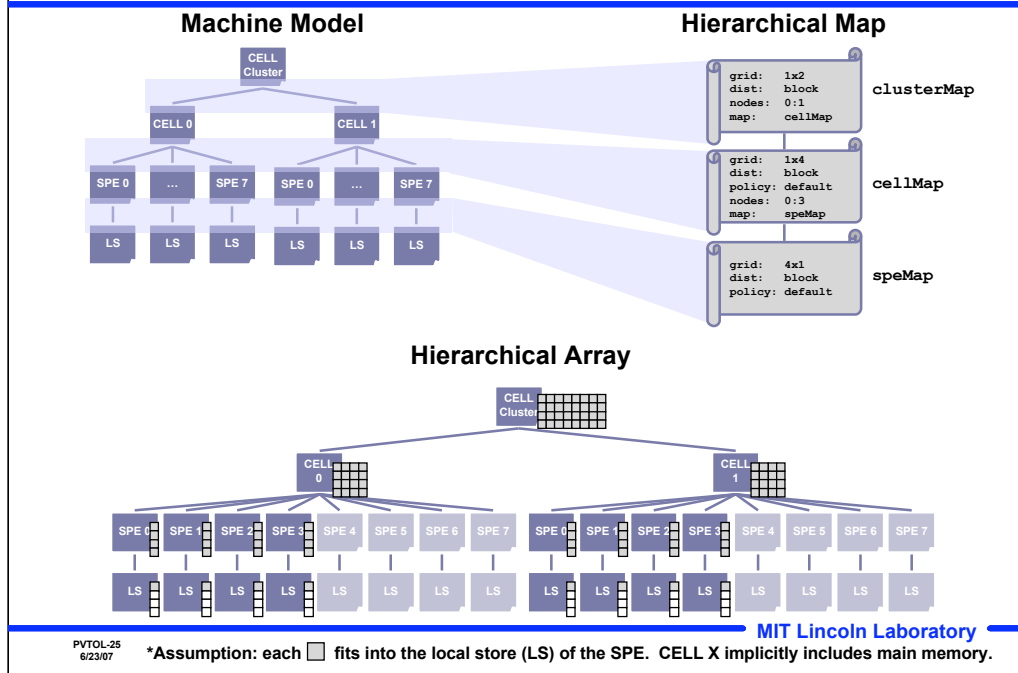
PVTOL-24
6/23/07

MIT Lincoln Laboratory

Distribution of data across a layer in the machine model has a one-to-one correspondence with one layer in the hierarchical map and with one LayerManager object.



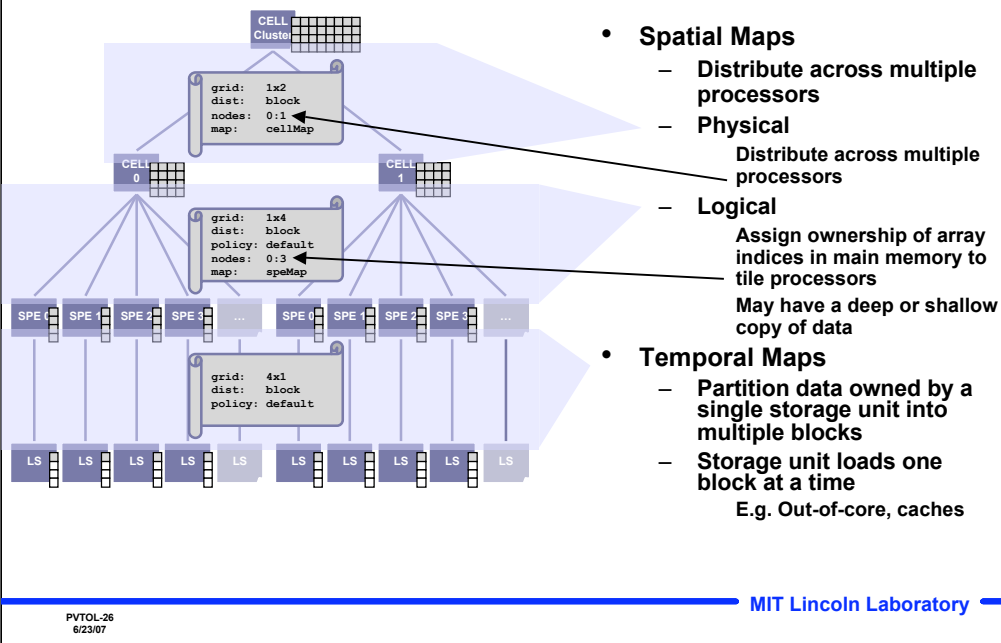
Hierarchical Array Mapping



This slide shows an example of how a hierarchical map distributes data across the memory hierarchy to construct a hierarchical array. Suppose we have a matrix with 4 rows and 8 columns (4x8). The highest level map (clusterMap) divides the column dimension into 2 blocks across nodes 0 and 1 in the CELL cluster. This results in a 4x4 submatrix on each CELL node. The next level map (cellMap) divides the column dimension into 4 blocks across SPE's 0 through 3 on each CELL node. This results in a 4x1 submatrix owned by each SPE. Finally, the lowest level map (speMap) divides the row dimension into 4 blocks. This results in four 1x1 submatrices. Only each 1x1 submatrix is loaded into the SPE's local store one at a time for processing.



Spatial vs. Temporal Maps



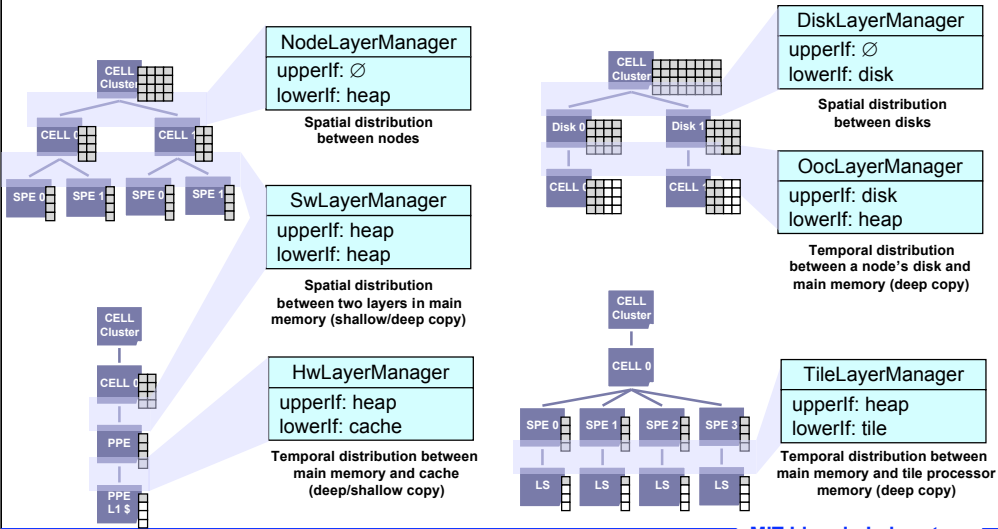
Maps describe how to partition an array. There are two types of maps: spatial and temporal. A spatial map describes how elements of an array are divided between multiple processors. A physical spatial map breaks up an array into blocks that are physically located on separate processors, e.g. dividing an array between two Cell processors. A logical map differs from a physical map in that the array being partitioned remains intact. Rather, the array is logically divided into blocks that are owned by different processors. For example, on a single Cell processor, the array resides in main memory. A logical map may assign blocks of rows to each SPE. Finally, temporal maps divided arrays into blocks that are loaded into memory one at a time. For example, a logical map may divide an array owned by a single SPE into blocks of rows. The array resides in main memory and the SPE loads one block at a time into its local store.



Layer Managers

These managers imply that there is main memory at the SPE level

- Manage the data distributions between adjacent levels in the machine model



PVTOL-27
6/23/07

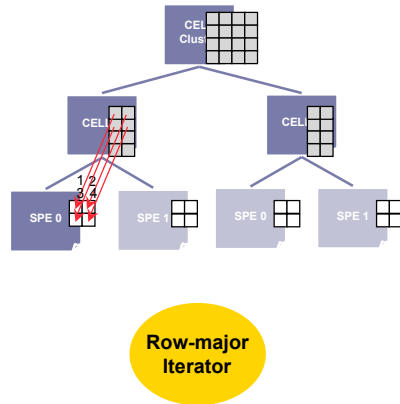
MIT Lincoln Laboratory

PVTOL instantiates different layer managers for data objects, depending on where in the memory hierarchy the data object is allocated. The layer manager manages the movement of data between the adjacent memory layers. The machine model restricts the types of layer managers; managers are implemented only for pairs of layers that can be adjacent, e.g. there is no layer manager for moving data from disk to tile memory since those two layers cannot be adjacent.



Tile Iterators

- Iterators are used to access temporally distributed tiles

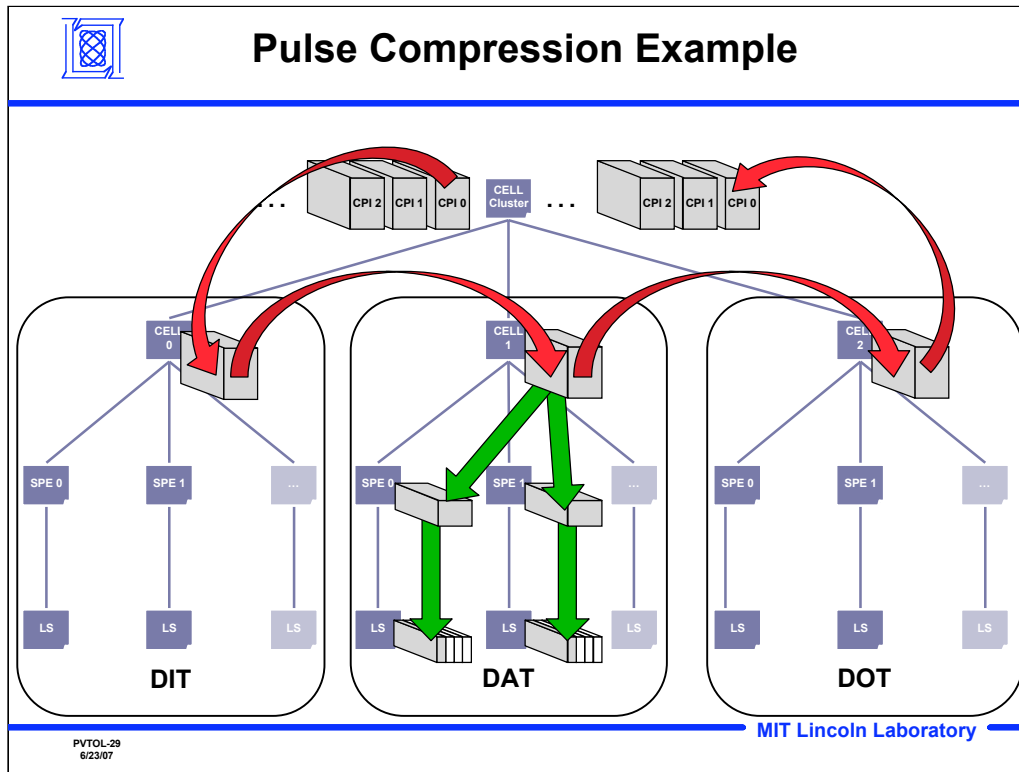


- Kernel iterators
 - Used within kernel expressions
- User iterators
 - Instantiated by the programmer
 - Used for computation that cannot be expressed by kernels
 - Row-, column-, or plane-order
- Data management policies specify how to access a tile
 - Save data
 - Load data
 - Lazy allocation (pMappable)
 - Double buffering (pMappable)

PVTOL-28
6/23/07

MIT Lincoln Laboratory

Iterators are used to access tiles that are distributed using temporal maps. The programmer can specify the direction the iterator moves through the tiles (row, column, plane-order) and how data should be transferred when accessing a tile, by specifying data management policies. For example, if the user intends to overwrite the contents of a tile, it is not necessary to actually load the tile. Rather, the tile can be allocated directly in the SPE. This eliminates the overhead of actually transferring data from main memory to the SPE's local store.



This shows an example of data flow through a pulse compression operation using hierarchical arrays. CELL 0 is responsible for the Data Input Task, which loads one CPI into the system at a time. Next, CELL 0 transfers the CPI to CELL 1 to perform the Data Analysis Task, which performs the pulse compression. CELL 1 distributes the CPI's rows across two SPE's. Each SPE divides its section of the CPI into blocks of columns. After CELL 1 has finished pulse compressing the CPI, the results are sent to CELL 2. CELL 2 performs the Data Output Task, which writes out the pulse compressed CPI. This example processing multiple CPI's in a pipelined fashion.



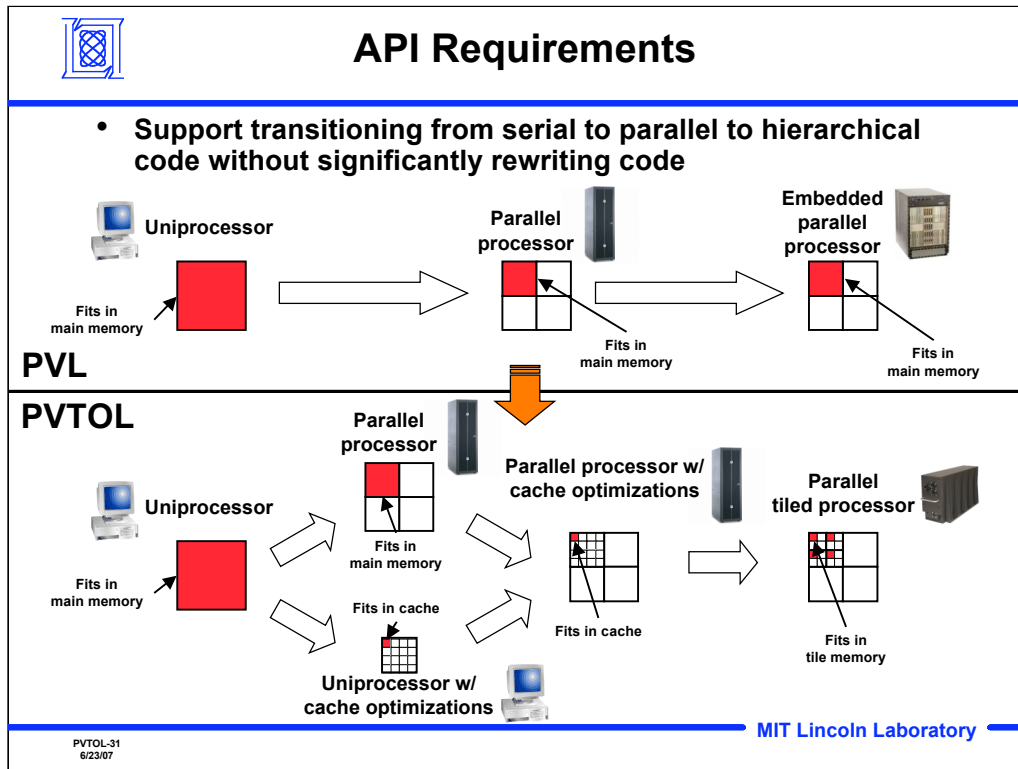
Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - ➔ – Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- Summary

PVTOL-30
6/23/07

MIT Lincoln Laboratory

Outline



PVTOL extends the incremental development process demonstrated in Lincoln's Parallel Vector Library (PVL). In PVL, programmers first develop a serial implementation of the application on their workstation. Once the serial implementation is verified, the applications can be parallelized on a cluster. Because parallelization is achieved using maps, this requires very little change to the code. Once the parallel implementation is verified and optimized for performance, the application can be deployed onto the embedded parallel processor with little to no change to the code.

PVTOL adds an additional step to the process. As before, a serial version is developed on a workstation. Next, the programmer can either develop a parallel version on a cluster (a la PVL) or develop a serial version that is optimized for cache performance on a workstation. In both cases, maps are used to divide the data across processors or into blocks that fit into cache. Next, hierarchical maps are used to implement a parallel version in which data on each processor is optimized for cache performance. This is achieved by using hierarchical maps. Finally, the code can be deployed onto a parallel tiled processor, such that the hierarchical map partitions data between processors and between tiles on each processor.



Data Types

- **Block types**
 - Dense
- **Element types**
 - int, long, short, char, float, double, long double
- **Layout types**
 - Row-, column-, plane-major
- **Dense**`<int Dims, class ElemType, class LayoutType>`
- **Views**
 - Vector, Matrix, Tensor
- **Map types**
 - Local, Runtime, Auto
- **Vector**`<class ElemType, class BlockType, class MapType>`

PVTOL-32
6/23/07

MIT Lincoln Laboratory

Data is stored in Block objects, which will initially support dense data objects (as opposed to sparse). Data dimensions can be laid out in memory in any order. Block can store nearly any data type.

Data stored in Blocks are accessed via View objects, which represent vectors, matrices, or tensors. Additionally, data distributions can be specified in the View object using different types of Map objects. Finally, Views specify whether data is mapped to the local processor's memory, mapped at runtime across multiple processors, or automatically mapped by pMapper.



Data Declaration Examples

Serial

```
// Create tensor
typedef Dense<3, float, tuple<0, 1, 2> > dense_block_t;
typedef Tensor<float, dense_block_t, LocalMap> tensor_t;
tensor_t cpi(Nchannels, Npulses, Nranges);
```

Parallel

```
// Node map information
Grid grid(Nprocs, 1, 1, Grid.ARRAY); // Grid
DataDist dist(3); // Block distribution
Vector<int> procs(Nprocs); // Processor ranks
procs(0) = 0; ...
ProcList procList(procs); // Processor list
RuntimeMap cpiMap(grid, dist, procList); // Node map

// Create tensor
typedef Dense<3, float, tuple<0, 1, 2> > dense_block_t;
typedef Tensor<float, dense_block_t, RuntimeMap> tensor_t;
tensor_t cpi(Nchannels, Npulses, Nranges, cpiMap);
```

This shows an example of how to allocate serial data in PVTOL, then parallelize that data using maps. Each map requires a grid, data distribution and list of processor ID's. The grid describes how to partition each dimension of the data object. The data distribution describes whether to use a block, cyclic or block-cyclic distribution. The list of processor ID's indicates which processors to distribute data across. The map type in the Tensor datatype is changed to RuntimeMap and the map object is passed into the Tensor constructor.



Data Declaration Examples

Hierarchical

```
// Tile map information
Grid tileGrid(1, NTiles 1, Grid.ARRAY); // Grid
DataDist tileDist(3); // Block distribution
DataMgmtPolicy tilePolicy(DataMgmtPolicy.DEFAULT); // Data mgmt policy
RuntimeMap tileMap(tileGrid, tileDist, tilePolicy); // Tile map

// Tile processor map information
Grid tileProcGrid(NTileProcs, 1, 1, Grid.ARRAY); // Grid
DataDist tileProcDist(3); // Block distribution
Vector<int> tileProcs(NTileProcs); // Processor ranks
inputProcs(0) = 0; ...
ProcList inputList(tileProcs); // Processor list
DataMgmtPolicy tileProcPolicy(DataMgmtPolicy.DEFAULT); // Data mgmt policy
RuntimeMap tileProcMap(tileProcGrid, tileProcDist, tileProcs,
    tileProcPolicy, tileMap); // Tile processor map

// Node map information
Grid grid(Nprocs, 1, 1, Grid.ARRAY); // Grid
DataDist dist(3); // Block distribution
Vector<int> procs(Nprocs); // Processor ranks
procs(0) = 0;
ProcList procList(procs); // Processor list
RuntimeMap cpiMap(grid, dist, procList, tileProcMap); // Node map

// Create tensor
typedef Dense<3, float, tuple<0, 1, 2> > dense_block_t;
typedef Tensor<float, dense_block_t, RuntimeMap> tensor_t;
tensor_t cpi(Nchannels, Npulses, Nranges, cpiMap);
```

PVTOL-34
6/23/07

MIT Lincoln Laboratory

This shows an example of how to extend the parallel code on the previous slide to allocate hierarchical data using hierarchical maps. In this case, two additional layers are added to the Tensor object. The tile processor map describes how to distribute data on a processor between tile processors on each node. This is an example of a spatial map. The tile map describes how to distribute data owned by a tile processor into blocks that can be loaded into the tile processor's memory and processed. This is an example of a temporal map. Note that once a new map is created, it simply needs to be passed in as an argument to the next map constructor to create the hierarchical map.

The method of constructing these maps are nearly identical to the node map. One additional parameter is added, however: data management policies. Data management policies allow the programmer to specify how data should move between layers in the memory hierarchy, potentially eliminating unnecessary communication.



Pulse Compression Example

Untiled version

```
// Declare weights and cpi tensors
tensor_t cpi(Nchannels, Npulses, Nranges,
             cpiMap),
          weights(Nchannels, Npulse, Nranges,
                 cpiMap);

// Declare FFT objects
Ffft<float, float, 2, fft_fwd> ffft;
Ffft<float, float, 2, fft_inv> iffft;

// Iterate over CPI's
for (i = 0; i < Ncpis; i++) {
    // DIT: Load next CPI from disk
    ...

    // DAT: Pulse compress CPI

    output = iffft(weights * ffft(cpi));

    // DOT: Save pulse compressed CPI to disk
    ...
}
```

Tiled version

```
// Declare weights and cpi tensors
tensor_t cpi(Nchannels, Npulses, Nranges,
             cpiMap),
          weights(Nchannels, Npulse, Nranges,
                 cpiMap);

// Declare FFT objects
Ffft<float, float, 2, fft_fwd> ffft;
Ffft<float, float, 2, fft_inv> iffft;

// Iterate over CPI's
for (i = 0; i < Ncpis; i++) {
    // DIT: Load next CPI from disk
    ...

    // DAT: Pulse compress CPI
    dataIter = cpi.beginLinear(0, 1);
    weightsIter = weights.beginLinear(0, 1);
    outputIter = output.beginLinear(0, 1);
    while (dataIter != data.endLinear()) {
        output = iffft(weights * ffft(cpi));
        dataIter++; weightsIter++; outputIter++;
    }

    // DOT: Save pulse compressed CPI to disk
    ...
}
```

PVTOL-3
6/23/07

Kernelized tiled version is identical to untiled version

atory

This is an example of how to implement pulse compression in PVTOL. The untiled version on the left runs on both serial and parallel processors. The tiled version on the right runs on tiled processors. On a tiled processor, data is distributed using temporal maps, i.e. data are broken into blocks that must be loaded one at a time. Consequently, iterators are used to load and process blocks one at a time.

Note that a pulse compression kernel could be written such that the library automatically recognize the pulse compression operation and automatically iterates through blocks in the hierarchical arrays. In other words, the iterator code in the tiled version would be hidden within the library. Consequently, such a “kernelized” version of the tile code would be identical to the untiled code.



Setup Assign API

- **Library overhead can be reduced by an initialization time expression setup**
 - Store PITFALLS communication patterns
 - Allocate storage for temporaries
 - Create computation objects, such as FFTs

Assignment Setup Example

```
Equation eq1(a, b*c + d);  
Equation eq2(f, a / d);  
  
for( ... ) {  
    ...  
    eq1();  
    eq2();  
    ...  
}
```

Expressions stored
in Equation object

Expressions
invoked without re-
stating expression

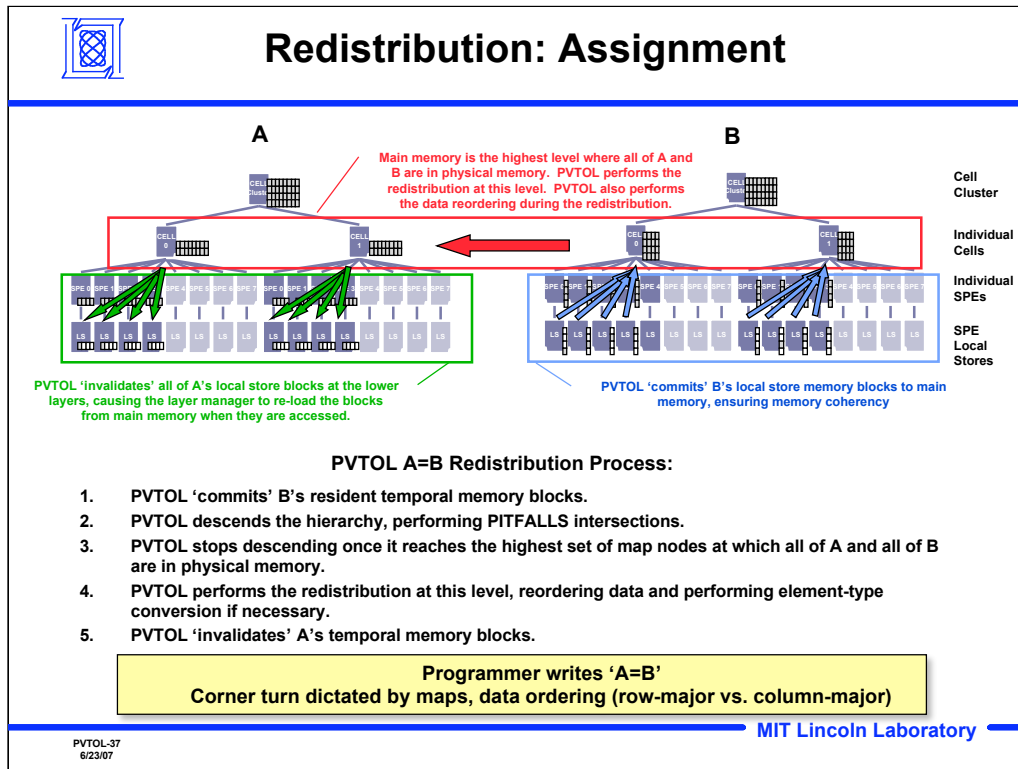
Expression objects can hold setup information without duplicating the equation

PVTOL-36
6/23/07

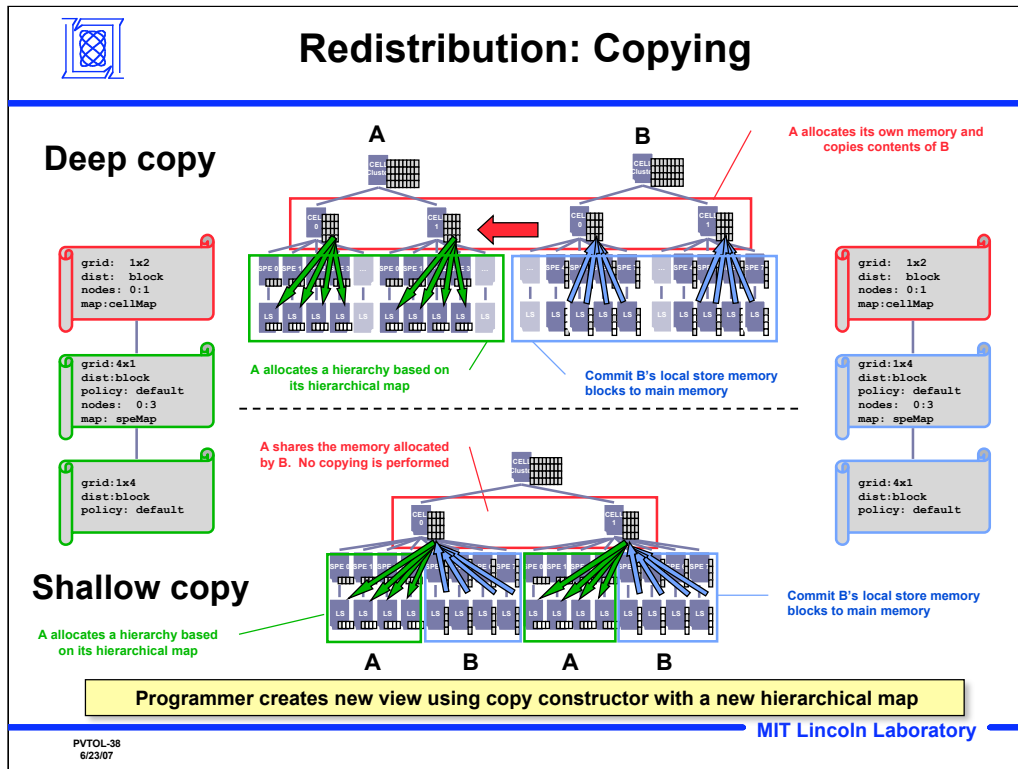
MIT Lincoln Laboratory

Executing expressions at runtime can incur various overheads, e.g. computing the PITFALLS communication patterns between operands with different mappings, allocating storage for temporary, intermediate values, and initializing values for certain types of computations like FFT.

Equation objects that describe expressions can be constructed. Thus, this overhead can be incurred at initialization time rather than during the computation.



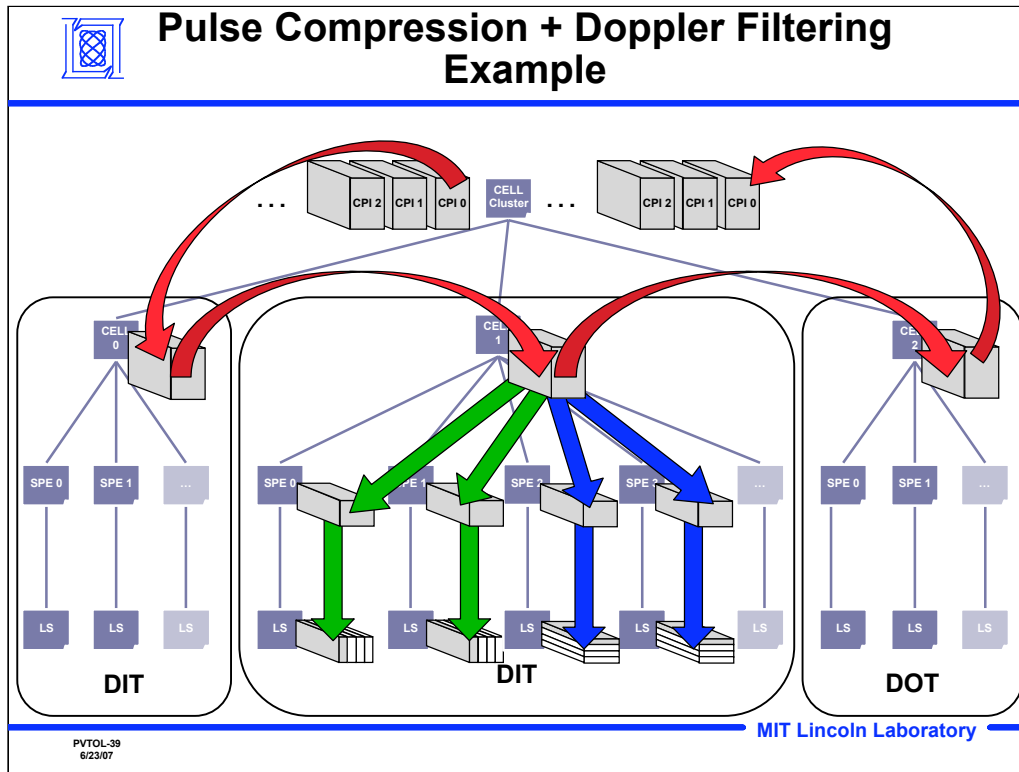
Often, the output data of an operation must have a different data distribution than the input data. A common type of operation which redistributes data is the corner turn, in which data that are distributed into blocks of columns are redistributed into blocks of rows, or vice versa. To perform a corner turn on a tiled processor, such as the Cell, all SPE's on each Cell processor write their respective columns of the array to main memory (blue). The Cell processors then redistribute the data in main memory from a column distribution to a row distribution (red). Finally, the rows on each Cell are distributed across the SPE's (green). To perform a corner turn, the programmer simply writes $A=B$, where A and B have different maps.



In the assignment $A = B$, depending on if A and B are on different processors, a deep copy is performed. If A and B are on the same processor, then the programmer can specify that A is a shallow copy of B.

If A and B are on different processors, the data on each SPE for B is written to main memory, then the contents of B are written into A, then the contents of A are distributed onto the SPE's according to A's hierarchical map.

If A and B are on the same processor, but distributed on different SPE's, then A and B can share the same data in main memory, but distribute the data onto different sets of SPE's.



This shows an extension of the previous example of data flow through a pulse compression operation with the addition of a Doppler filtering operation. CELL 0 loads one CPI into the system at a time. Next, CELL 0 transfers the CPI to CELL 1 to perform the Data Analysis Task. CELL 1 distributes the CPI's rows across SPE 0 and SPE 1. Each SPE divides its section of the CPI into blocks of columns. Once the pulse compression is complete, the CPI is redistributed across SPE 2 and SPE 3 such that rows are still distributed across the SPE's but each SPE divides its section into blocks of rows. This redistribution is necessary because the pulse compression requires a different data distribution than the pulse compression. Once the DAT is complete, CELL 1 sends the processed CPI to CELL 2, which writes out the CPI.



Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- Summary

PVTOL-40
6/23/07

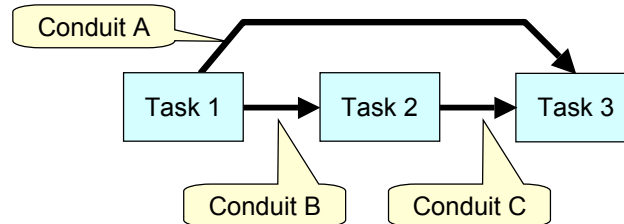
MIT Lincoln Laboratory

Outline



Tasks & Conduits

A means of decomposing a problem into a set of asynchronously coupled sub-problems (a pipeline)

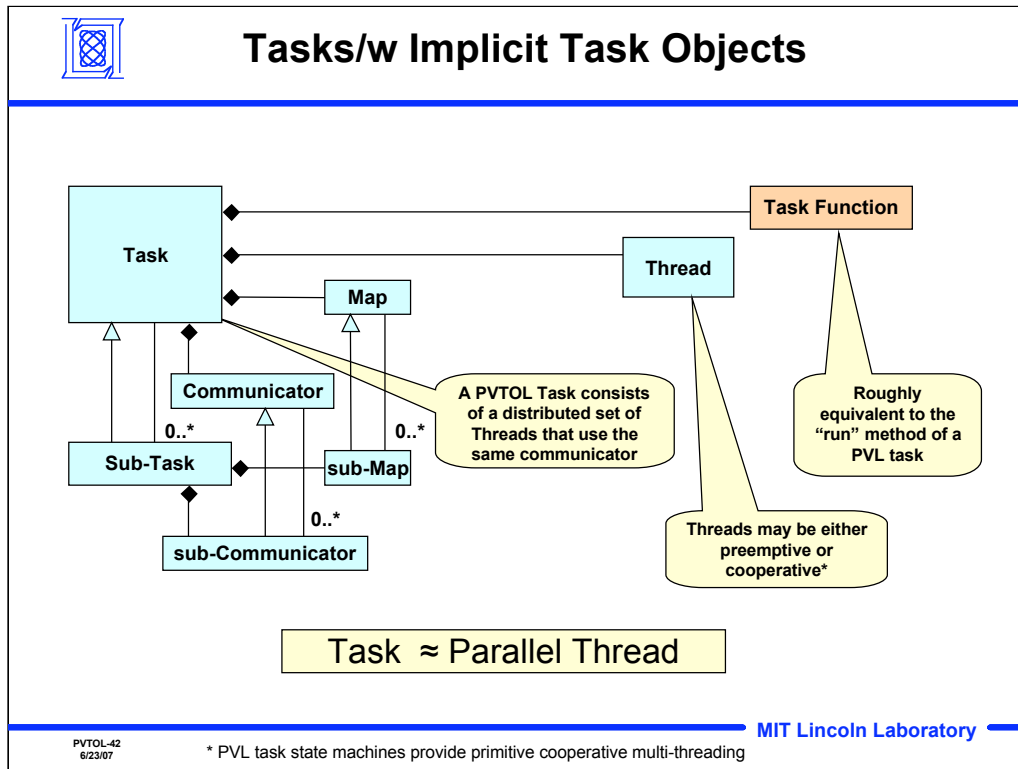


- Each Task is SPMD
- Conduits transport distributed data objects (i.e. Vector, Matrix, Tensor) between Tasks
- Conduits provide multi-buffering
- Conduits allow easy task replication
- Tasks may be separate processes or may co-exist as different threads within a process

PVTOL-41
6/23/07

MIT Lincoln Laboratory

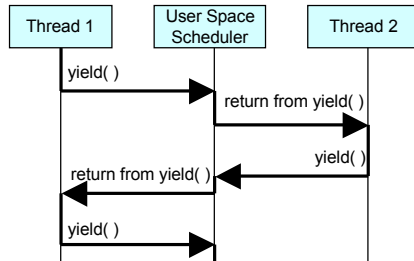
Each task in a pipeline is notionally an independent SPMD. Conduits transport distributed data objects between the tasks and isolate the mapping of data within a task from all other tasks.





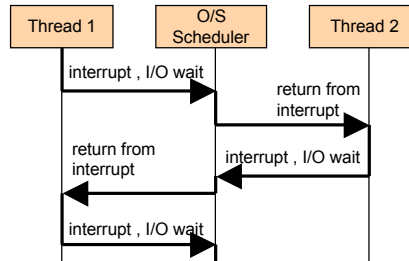
Cooperative vs. Preemptive Threading

Cooperative User Space Threads (e.g. GNU Pth)



- PVTOL calls `yield()` instead of blocking while waiting for I/O
- O/S support of multithreading not needed
- Underlying communication and computation libs need not be thread safe
- SMPs cannot execute tasks concurrently

Preemptive Threads (e.g. pthread)



- SMPs can execute tasks concurrently
- Underlying communication and computation libs must be thread safe

PVTOL can support both threading styles via an internal thread wrapper layer

PVTOL-43
6/23/07

MIT Lincoln Laboratory

With cooperative threading a thread must explicitly yield execution back to a scheduler. With preemptive scheduling an interrupt can be used to force a thread to yield execution back to a scheduler. Most programs are in one of two major states. They are either waiting for I/O or processing data. Preemptive scheduling is mandatory only when compute bound processing can take so long that I/O can be missed. In real-time systems with bounded processing latency requirements this is rarely the case so either threading approach will usually succeed. Note that cooperative scheduling can be more efficient since preemptive context switching tends to add unnecessary overhead.



Task API

support functions get values for current task SPMD

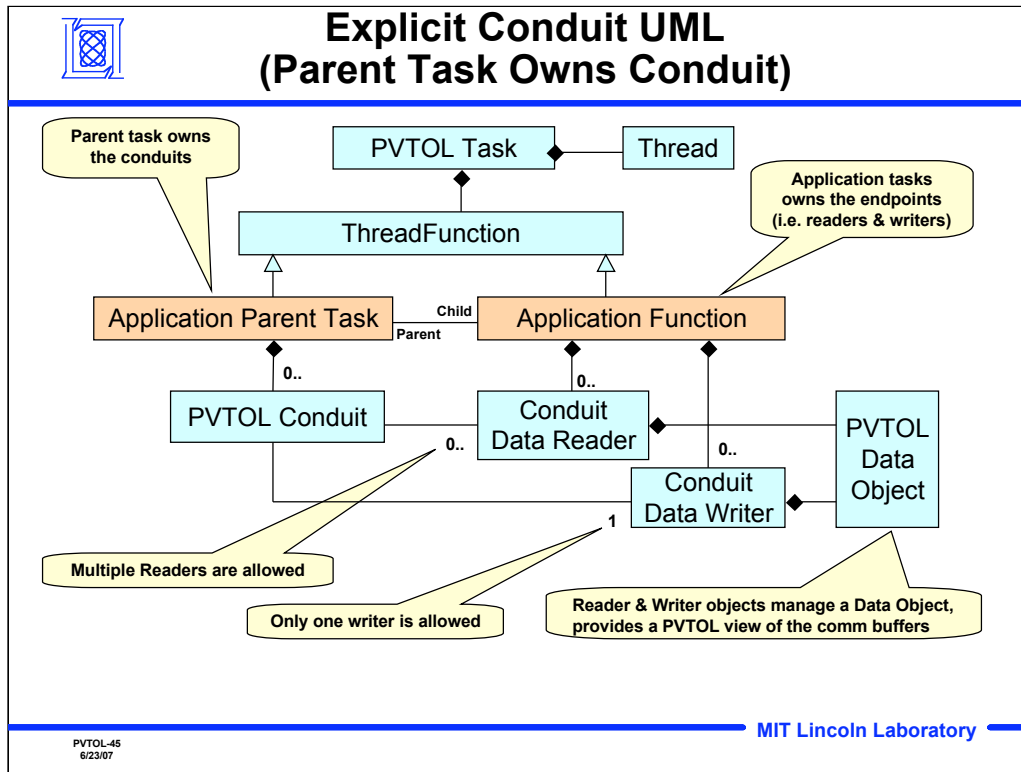
- `length_type pvtol::num_processors();`
- `const_Vector<processor_type>`
`pvtol::processor_set();`

Task API

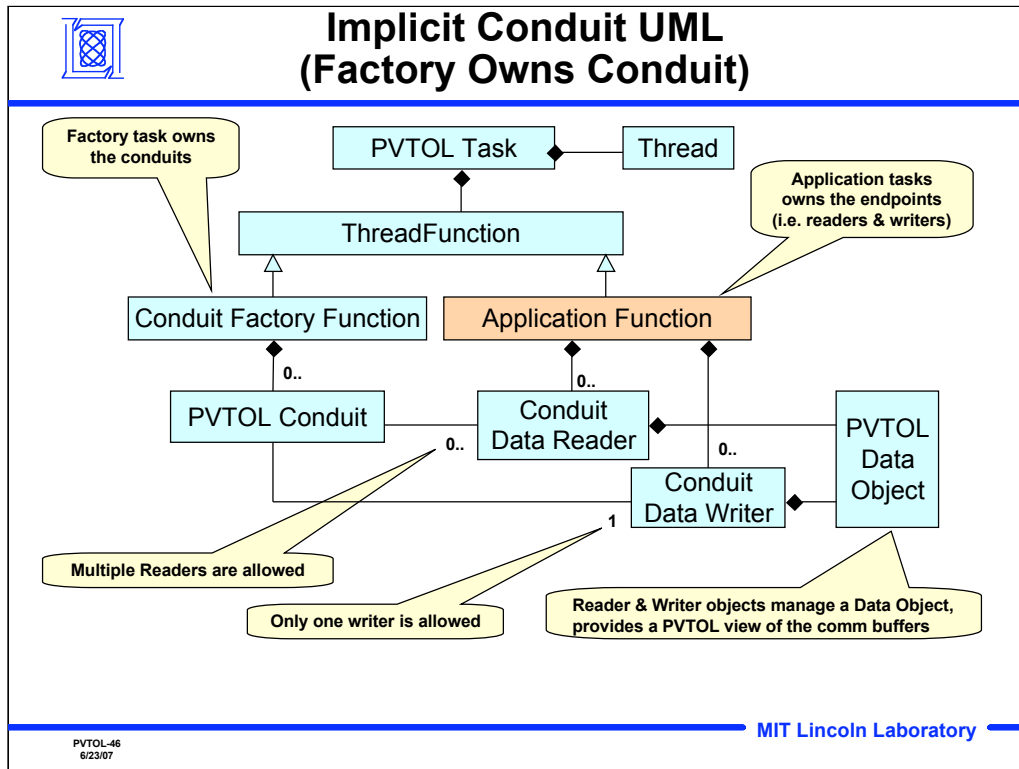
- `typedef<class T>`
`pvtol::tid pvtol::spawn(`
 `(void) (TaskFunction*) (T&),`
 `T& params,`
 `Map& map);`
- `int pvtol::tidwait(pvtol::tid);`

Similar to typical thread API except for spawn map

The proposed task API. Note the simplicity of this API vs. the PVL Task API.



Class diagram of an explicit conduit case. This case is explicit because a parent Task actually declares a Conduit object and passes the Conduit interface objects into its child Tasks. This is less than ideal since the communication between child Tasks is hard coded into the parent Task.



Class diagram of the implicit Conduit case. As in the explicit Conduit case a task still needs to own the Conduit. Instead of a parent Task the Conduit's owner is a Conduit factory Task. The Conduit factory Task is mapped across the entire application SPMD and therefore form any connection between Tasks. This approach is more flexible since tasks can connect to each other via a "conduit name" or "topic" rather than needing some parent task to explicitly form the connection.



Conduit API

Conduit Declaration API

```
• typedef<class T>
  class Conduit {
    Conduit( );
    Reader& getReader( );
    Writer& getWriter( );
  };
```

Note: the Reader and Writer connect() methods block waiting for conduits to finish initializing and perform a function similar to PVL's two phase initialization

Conduit Reader API

```
• typedef<class T>
  class Reader {
  public:
    Reader( Domain<n> size, Map map, int depth );
    void setup( Domain<n> size, Map map, int depth );
    void connect( ); // block until conduit ready
    pvtolPtr<T> read( ); // block until data available
    T& data( ); // return reader data object
  };
```

Conduit Writer API

```
• typedef<class T>
  class Writer {
  public:
    Writer( Domain<n> size, Map map, int depth );
    void setup( Domain<n> size, Map map, int depth );
    void connect( ); // block until conduit ready
    pvtolPtr<T> getBuffer( ); // block until buffer available
    void write( pvtolPtr<T> ); // write buffer to destination
    T& data( ); // return writer data object
  };
```

Conceptually Similar to the PVL Conduit API

MIT Lincoln Laboratory

PVTOL-47
6/23/07

The Conduit API. The Conduit class is explicitly instantiated in an application Task in the explicit Conduit case or within a Conduit factory Task in the implicit Conduit case. Tasks that wish to send or receive data use the Reader or Writer classes to perform the required action.



Task & Conduit API Example/w Explicit Conduits

```
typedef struct { Domain<2> size; int depth; int numCpis; } DatParams;  
  
int DataInputTask(const DitParams*);  
int DataAnalysisTask(const DatParams*);  
int DataOutputTask(const DotParams*);  
  
int main( int argc, char* argv[])  
{  
    ...  
    Conduit<Matrix<Complex<Float>>> conduit1;  
    Conduit<Matrix<Complex<Float>>> conduit2;  
  
    DatParams  datParams = ...;  
    datParams.inp = conduit1.getReader( );  
    datParams.out = conduit2.getWriter( );  
  
    vsip::tid ditTid = vsip::spawn( DataInputTask, &ditParams, ditMap);  
    vsip::tid datTid = vsip::spawn( DataAnalysisTask, &datParams, datMap );  
    vsip::tid dotTid = vsip::spawn( DataOutputTask, &dotParams, dotMap );  
  
    vsip::tidwait( ditTid );  
    vsip::tidwait( datTid );  
    vsip::tidwait( dotTid );  
}
```

Conduits
created in
parent task

Pass Conduits
to children via
Task
parameters

Spawn
Tasks

Wait for
Completion

"Main Task" creates Conduits, passes to sub-tasks as parameters, and waits for them to terminate

PVTOL-48
6/23/07

MIT Lincoln Laboratory

A simple explicit Conduit example, part 1. The parent Task declaration of the Conduit and the passing of conduit endpoint interfaces to the child tasks is shown.



DAT Task & Conduit Example/w Explicit Conduits

```
int DataAnalysisTask(const DatParams* p)
```

```
{
```

```
    Vector<Complex<Float>> weights( p.cols, replicatedMap );
```

```
    ReadBinary (weights, "weights.bin" );
```

```
    Conduit<Matrix<Complex<Float>>>::Reader inp( p.inp );
```

```
    inp.setup(p.size,map,p.depth);
```

```
    Conduit<Matrix<Complex<Float>>>::Writer out( p.out );
```

```
    out.setup(p.size,map,p.depth);
```

```
    inp.connect( );
```

```
    out.connect( );
```

```
    for(int i=0; i<p.numCpis; i++) {
```

```
        pvtolPtr<Matrix<Complex<Float>>> inpData( inp.read() );
```

```
        pvtolPtr<Matrix<Complex<Float>>> outData( out.getBuffer() );
```

```
        (*outData) = ifftm( vmmul( weights, fftm( *inpData, VSIP_ROW ),
```

```
            VSIP_ROW );
```

```
        out.write(outData);
```

```
    }
```

```
}
```

Declare and Load
Weights

Complete conduit
initialization

connect() blocks until
conduit is initialized

Reader::getHandle()
blocks until data
is received

Writer::getHandle()
blocks until output
buffer is available

Writer::write()
sends the data

pvtolPtr destruction
implies reader extract

Sub-tasks are implemented as ordinary functions

MIT Lincoln Laboratory

PVTOL-49
6/23/07

A simple explicit Conduit example, part 2. An implementation of the DataAnal child Task is shown where the input and output Conduit interfaces are obtained from a parent Task.



DIT-DAT-DOT Task & Conduit API Example/w Implicit Conduits

```
typedef struct { Domain<2> size; int depth; int numCpis; } TaskParams;
```

```
int DataInputTask(const InputTaskParams*);  
int DataAnalysisTask(const AnalysisTaskParams*);  
int DataOutputTask(const OutputTaskParams*);
```

```
int main( int argc, char* argv[ ] )
```

```
{
```

```
    ...  
    TaskParams    params = ...;
```

```
    vsip:: tid ditTid = vsip:: spawn( DataInputTask, &params,ditMap);  
    vsip:: tid datTid = vsip:: spawn( DataAnalysisTask, &params, datMap );  
    vsip:: tid dotTid = vsip:: spawn( DataOutputTask, &params, dotMap );
```

```
    vsip:: tidwait( ditTid );  
    vsip:: tidwait( datTid );  
    vsip:: tidwait( dotTid );
```

```
}
```

Conduits
NOT created
in parent task

Spawn
Tasks

Wait for
Completion

“Main Task” just spawns sub-tasks and waits for them to terminate

MIT Lincoln Laboratory

PVTOL-50
6/23/07

A simple implicit Conduit example, part 1. Since the child Tasks connect to each other without help from the parent Task the parent Task only needs to spawn off the child Tasks. Note that this code is nearly identical for the different child Tasks. For some classes of application, it should be possible to write a “generic parent” that spawns of the required child Tasks in a standard manner.



DAT Task & Conduit Example/w Implicit Conduits

```
int DataAnalysisTask(const AnalysisTaskParams* p)
```

```
{  
    Vector<Complex<Float>> weights( p.cols, replicatedMap );  
    ReadBinary( weights, "weights.bin" );  
    Conduit<Matrix<Complex<Float>>>::Reader  
        inp("inpName",p.size,map,p.depth);  
    Conduit<Matrix<Complex<Float>>>::Writer  
        out("outName",p.size,map,p.depth);  
    inp.connect( );  
    out.connect( );  
    for(int i=0; i<p.numCpis; i++) {  
        pvtolPtr<Matrix<Complex<Float>>> inpData( inp.read() );  
        pvtolPtr<Matrix<Complex<Float>>> outData( out.getBuffer() );  
        (*outData) = ifftm( vmmul( weights, fftm( *inpData, VSIP_ROW ),  
                               VSIP_ROW );  
        out.write(outData);  
    }  
}
```

Constructors
communicate/w
factory to find other
end based on name

connect() blocks until
conduit is initialized

Reader::getHandle()
blocks until data is
received

Writer::getHandle()
blocks until output
buffer is available

Writer::write()
sends the data

pvtolPtr destruction
implies reader extract

Implicit Conduits connect using a "conduit name"

PVTOL-51
6/23/07

MIT Lincoln Laboratory

A simple implicit Conduit example, part 2. An implementation of the DataAnal child Task is shown where the Conduit endpoints are obtained from the Conduit factory.



Conduits and Hierarchal Data Objects

Conduit connections may be:

- Non-hierarchical to non-hierarchical
- Non-hierarchical to hierarchal
- Hierarchal to Non-hierarchical
- Non-hierarchical to Non-hierarchical

Example task function/w hierarchal mappings on conduit input & output data

```
...
input.connect();
output.connect();
for(int i=0; i<nCpi; i++) {
    pvtolPtr<Matrix<Complex<Float>>> inp( input.getHandle( ) );
    pvtolPtr<Matrix<Complex<Float>>> oup( output.getHandle( ) );
    do {
        *oup = processing( *inp );
        inp->getNext( );
        oup->getNext( );
    } while (more-to-do);
    output.write( oup );
}
```

Per-time Conduit
communication possible
(implementation dependant)

Conduits insulate each end of the conduit from the other's mapping

PVTOL-52
6/23/07

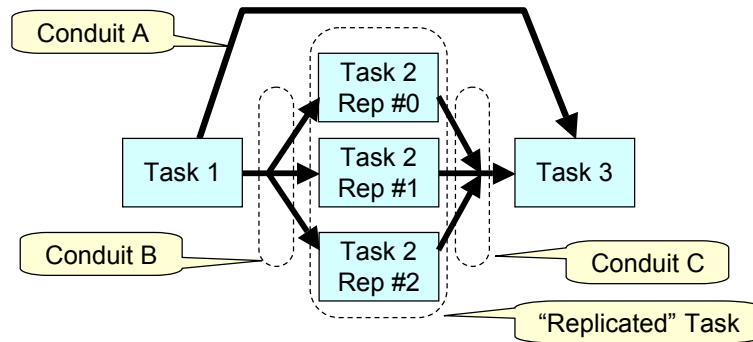
MIT Lincoln Laboratory

Hierarchal data objects have somewhat different access semantics than “normal” data objects. A conduit transports an entire hierarchal object which must be iterated with the hierarchal “getNext()” functions to access the objects data tiles.



Replicated Task Mapping

- Replicated tasks allow conduits to abstract away round-robin parallel pipeline stages
- Good strategy for when tasks reach their scaling limits



Replicated mapping can be based on a 2D task map (i.e. Each row in the map is a replica mapping, number of rows is number of replicas)

PVTOL-53
6/23/07

MIT Lincoln Laboratory

Replicated Task mapping is needed in many applications. Replicated Task mapping is useful when additional throughput is needed but a Task has been scaled up to the reasonable limit of its data parallelism.



Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- ➔ • PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- Summary

PVTOL-54
6/23/07

MIT Lincoln Laboratory

Outline



PVTOL and Map Types

PVTOL distributed arrays are templated on map type.

LocalMap	The matrix is <i>not</i> distributed
RuntimeMap	The matrix is distributed and all map information is specified at runtime
AutoMap	The map is either fully defined, partially defined, or undefined

Notional matrix construction:

```
Matrix<float, Dense, AutoMap> mat1(rows, cols);
```

Specifies the data type, i.e. double, complex, int, etc.

Specifies the storage layout

Specifies the map type:

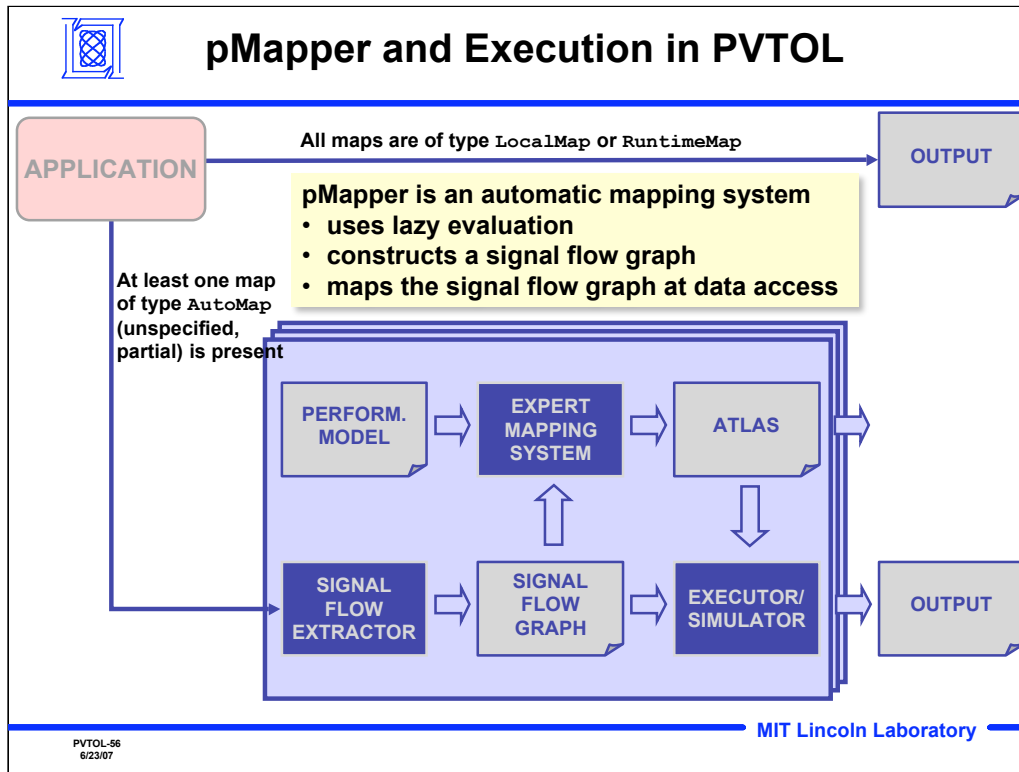
Focus on

PVTOL-55
6/23/07

MIT Lincoln Laboratory

PVTOL supports three different types of maps. A LocalMap indicates that data is not distributed, but exists entirely in the processor's local memory. A RuntimeMap indicates that data is distributed across multiple processors and that the user specifies all map information. An AutoMap indicates that data may be distributed across multiple processors and that map information is fully or partially defined by the user or the map is undefined. When maps are partially defined or undefined, PVTOL automated mapping technology, pMapper, will do the following:

- pMapper will fill in the missing details of an hierarchical map (grid, distribution)
- pMapper will not discover whether an hierarchical map should be present
- pMapper will be responsible for discovering overlap
- pMapper will not be responsible for discovering the storage level
- pMapper will be responsible for determining some data management policies



This shows a flowchart of how data is mapped by PVTOL. When all maps are LocalMaps or RuntimeMaps, the map information is directly used to allocate and distribute arrays. If at least one map is an AutoMap, the automated mapping system is invoked, which will determine optimal map information.



Examples of Partial Maps

A partially specified map has one or more of the map attributes unspecified at one or more layers of the hierarchy.

Examples:

Grid: 1x4
Dist: block
Procs:

pMapper will figure out which 4 processors to use

Grid: 1x*
Dist: block
Procs:

pMapper will figure out how many columns the grid should have and which processors to use; note that if the processor list was provided the map would become fully specified

Grid: 1x4
Dist:
Procs: 0:3

pMapper will figure out whether to use block, cyclic, or block-cyclic distribution

Grid:
Dist: block
Procs:

pMapper will figure out what grid to use and on how many processors; this map is very close to being completely unspecified

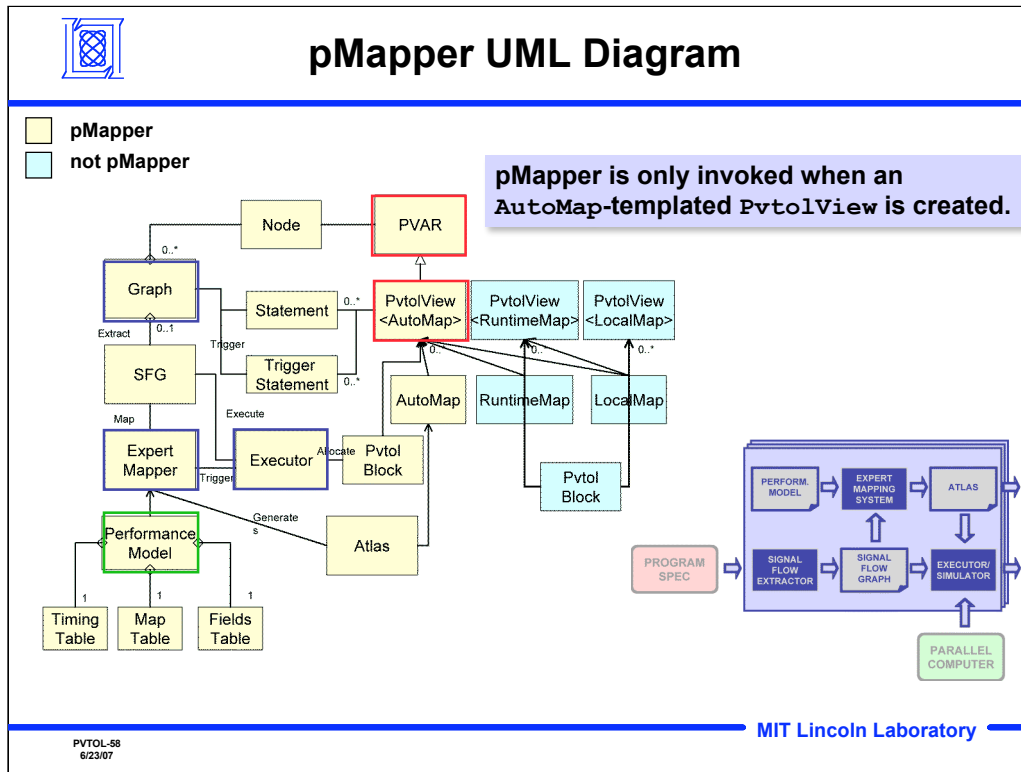
pMapper

- will be responsible for determining attributes that influence performance
- will not discover whether a hierarchy should be present

PVTOL-57
6/23/07

MIT Lincoln Laboratory

This slides shows examples of partial maps. In short, pMapper will be responsible for determining attributes that influence performance, e.g. number of processors (Procs); block, cyclic or block-cycle distributions (Dist); the number of blocks to partition each dimension into (Grid); or any combination of these three attributes. Note that pMapper is NOT responsible for determining whether to distribute data across the processor hierarchy.



This slide shows the UML diagram for the different types of maps and pMapper. Note that pMapper is only invoked when a data object of type PvtolView is created that is templated on the AutoMap data type.

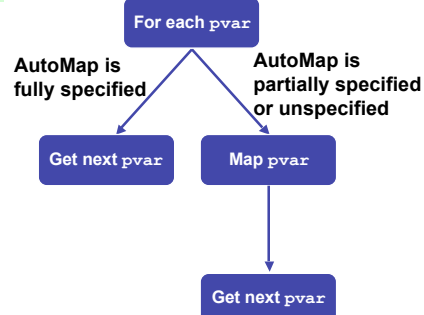
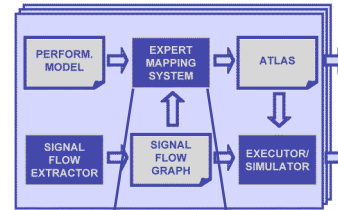


pMapper & Application

```
// Create input tensor (flat)
typedef Dense<3, float, tuple<0, 1, 2> > dense_block_t;
typedef Tensor<float, dense_block_t, AutoMap> tensor_t;
tensor_t input(Nchannels, Npulses, Nranges),
```

```
// Create input tensor (hierarchical)
AutoMap tileMap();
AutoMap tileProcMap(tileMap);
AutoMap cpiMap(grid, dist, procList, tileProcMap);
typedef Dense<3, float, tuple<0, 1, 2> > dense_block_t;
typedef Tensor<float, dense_block_t, AutoMap> tensor_t;
tensor_t input(Nchannels, Npulses, Nranges, cpiMap),
```

- For each Pvar in the Signal Flow Graph (SFG), pMapper checks if the map is fully specified
- If it is, pMapper will move on to the next Pvar
- pMapper will not attempt to remap a pre-defined map
- If the map is not fully specified, pMapper will map it
- When a map is being determined for a Pvar, the map returned has all the levels of hierarchy specified, i.e. all levels are mapped at the same time



PVTOL-59
6/23/07

MIT Lincoln Laboratory

This slide describes the logic PVTOL uses to determine whether or not to invoke pMapper. Two code examples are shown that construct a flat array and a hierarchical array. The first example that constructs a flat array specifies no map information. In this case, pMapper will determine all map information. In the second example that constructs a hierarchical array, because pMapper is not responsible for determining whether or not to allocate a hierarchical array, the programmer specifies a hierarchical map composed of two maps. Note that both maps contain no information. This indicates to pMapper that a two-level hierarchical array should be allocated and that pMapper is responsible for determining all map information for each level in the array.



Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- Summary



PVTOL-60
6/23/07

MIT Lincoln Laboratory

Outline



Mercury Cell Processor Test System



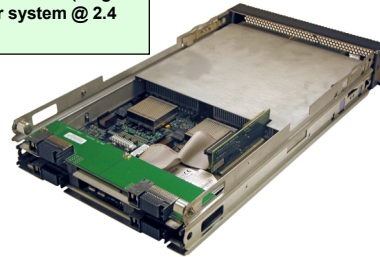
Mercury Cell Processor System

- Single Dual Cell Blade
 - Native tool chain
 - Two 2.4 GHz Cells running in SMP mode
 - Terra Soft Yellow Dog Linux 2.6.14
- Received 03/21/06
 - booted & running same day
 - integrated/w LL network < 1 wk
 - Octave (Matlab clone) running
 - Parallel VSIPL++ compiled

•Each Cell has 153.6 GFLOPS (single precision) – 307.2 for system @ 2.4 GHz (maximum)

Software includes:

- IBM Software Development Kit (SDK)
 - Includes example programs
- Mercury Software Tools
 - MultiCore Framework (MCF)
 - Scientific Algorithms Library (SAL)
 - Trace Analysis Tool and Library (TATL)



PVTOL-61
6/23/07

MIT Lincoln Laboratory

This is a description of our original Cell test system. The 2.4 GHz blade was eventually replaced by a 3.2 GHz blade with 205 GFLOPS per cell or 410 GFLOPS per dual cell blade.

Picture of workstation is taken from Mercury web site. Picture of blade was supplied by Mercury.



Outline

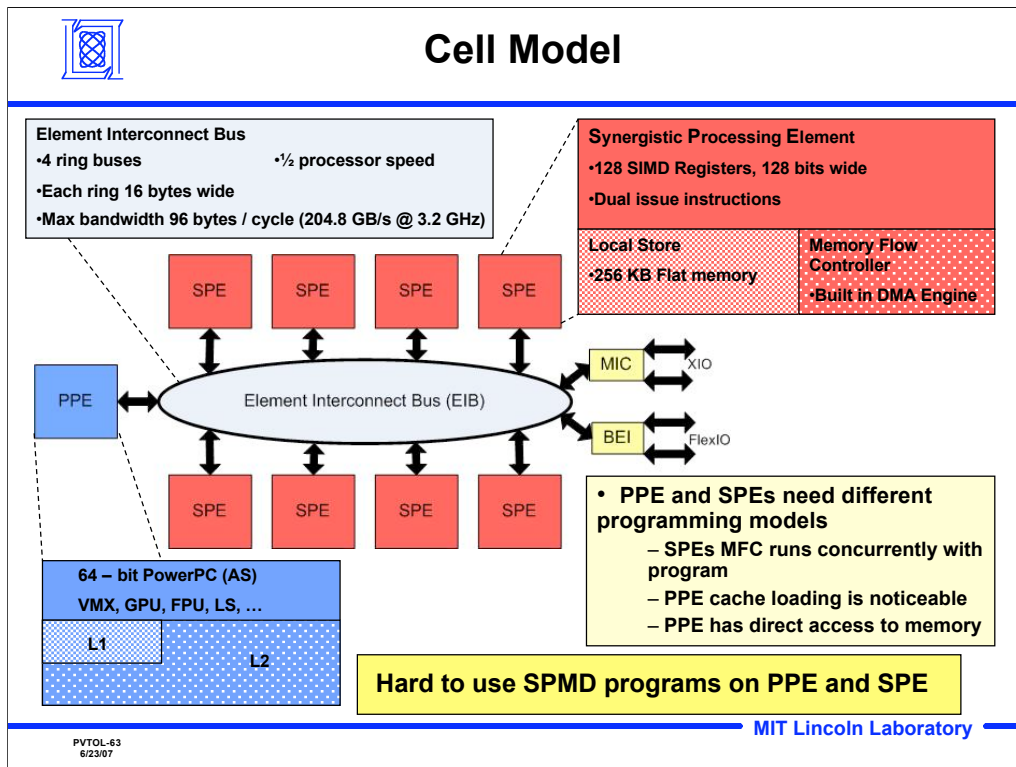
- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- Summary



PVTOL-62
6/23/07

MIT Lincoln Laboratory

Outline



This is a block diagram of the STI Cell processor with details of some of the features.

Max bandwidth taken from T. Chen et. al., *Cell Broadband Engine Architecture and its first implementation: A performance view*, IBM, 2005 .



Compiler Support



- **GNU gcc**
 - gcc, g++ for PPU and SPU
 - Supports SIMD C extensions



- **IBM Octopiler**
 - Promises automatic parallel code with DMA
 - Based on OpenMP



- **IBM XLC**
 - C, C++ for PPU, C for SPU
 - Supports SIMD C extensions
 - Promises transparent SIMD code
 - vadd does not produce SIMD code in SDK

- **GNU provides familiar product**
- **IBM's goal is easier programmability**
 - Will it be enough for high performance customers?

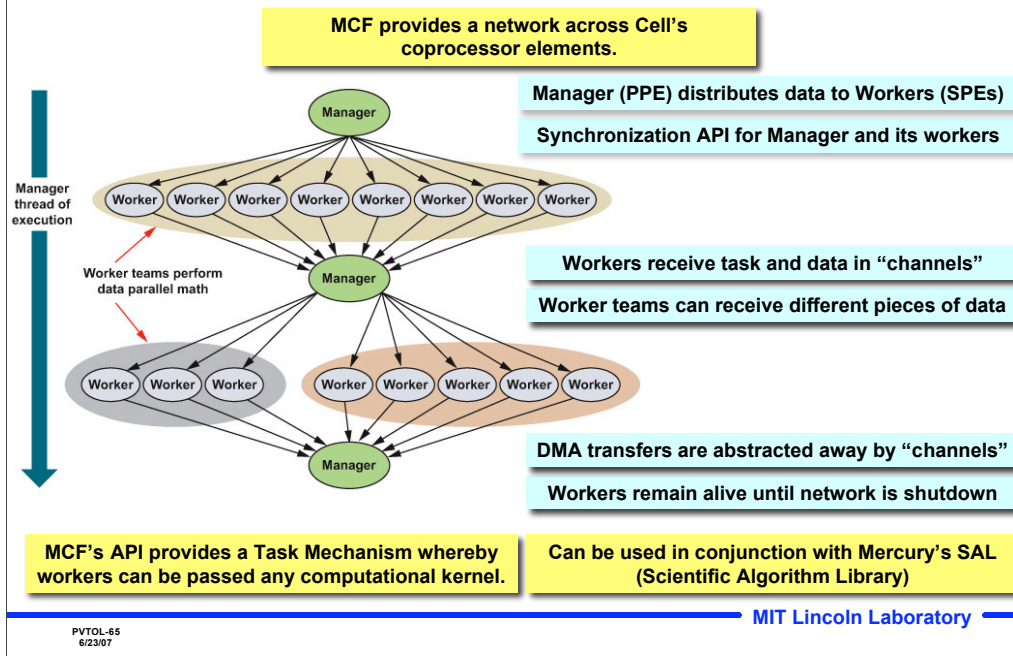
PVTOL-64
6/23/07

MIT Lincoln Laboratory

This is the state of the compilers as of Summer 2006. Both the GNU compiler and IBM's XLC are freely available and included in IBM's SDK releases.



Mercury's MultiCore Framework (MCF)

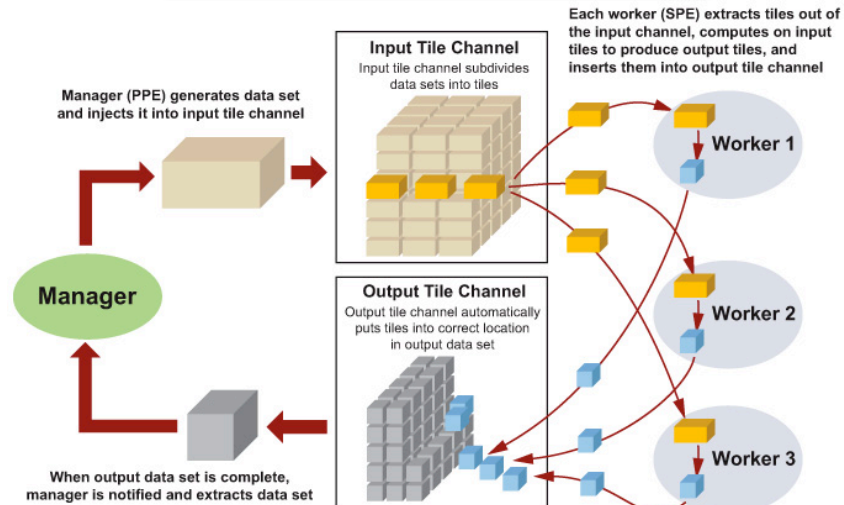


This slide shows the fundamental programming model of MCF. The advantage of MCF over threaded methods of control is that the workers have a small kernel that is started only once. In a threaded model, the worker is started again for each thread which has an expensive overhead.



Mercury's MultiCore Framework (MCF)

MCF provides API data distribution "channels" across processing elements that can be managed by PVTOL.



PVTOL-66
6/23/07

MIT Lincoln Laboratory

MCF also has built in data reorganization capabilities. The most common data partitioning methods used in parallel programming are encapsulated into function calls. This makes programming easier since the user does not have to reinvent the data reorganization over and over.



Sample MCF API functions

Manager Functions

```
mcf_m_net_create( )  
mcf_m_net_initialize( )  
mcf_m_net_add_task( )  
mcf_m_net_add_plugin( )  
mcf_m_team_run_task( )  
mcf_m_team_wait( )  
mcf_m_net_destroy( )  
mcf_m_mem_alloc( )  
mcf_m_mem_free( )  
mcf_m_mem_shared_alloc( )
```

```
mcf_m_tile_channel_create( )  
mcf_m_tile_channel_destroy( )  
mcf_m_tile_channel_connect( )  
mcf_m_tile_channel_disconnect( )  
mcf_m_tile_distribution_create_2d( )  
mcf_m_tile_distribution_destroy( )  
mcf_m_tile_channel_get_buffer( )  
mcf_m_tile_channel_put_buffer( )
```

```
mcf_m_dma_pull( )  
mcf_m_dma_push( )  
mcf_m_dma_wait( )  
mcf_m_team_wait( )
```

Worker Functions

```
mcf_w_main( )  
mcf_w_mem_alloc( )  
mcf_w_mem_free( )  
mcf_w_mem_shared_attach( )
```

```
mcf_w_tile_channel_create( )  
mcf_w_tile_channel_destroy( )  
mcf_w_tile_channel_connect( )  
mcf_w_tile_channel_disconnect( )  
mcf_w_tile_channel_is_end_of_channel( )  
mcf_w_tile_channel_get_buffer( )  
mcf_w_tile_channel_put_buffer( )
```

```
mcf_w_dma_pull_list( )  
mcf_w_dma_push_list( )  
mcf_w_dma_pull( )  
mcf_w_dma_push( )  
mcf_w_dma_wait( )
```

Initialization/Shutdown

Channel Management

Data Transfer

PVTOL-67
6/23/07

MIT Lincoln Laboratory

Here are some of the functions available in MCF. This list is not comprehensive, but is intended to give the viewer a flavor of the programming capabilities of MCF.



Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - – PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results
- Summary

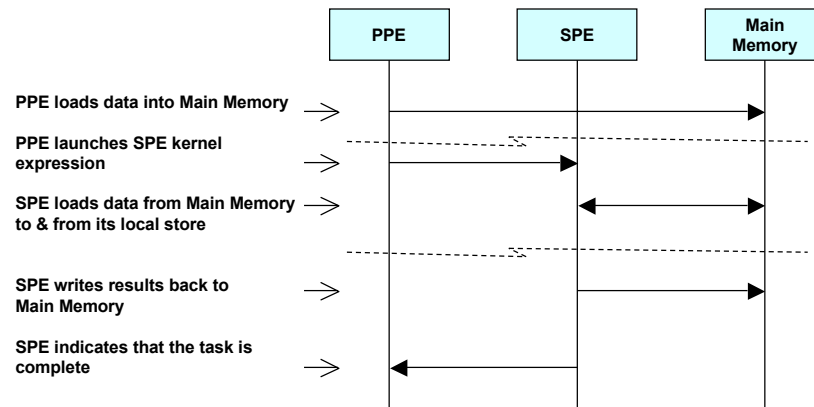
PVTOL-68
6/23/07

MIT Lincoln Laboratory

Outline



Cell PPE – SPE Manager / Worker Relationship



PPE (manager) "farms out" work to the SPEs (workers)

PVTOL-69
6/23/07

MIT Lincoln Laboratory

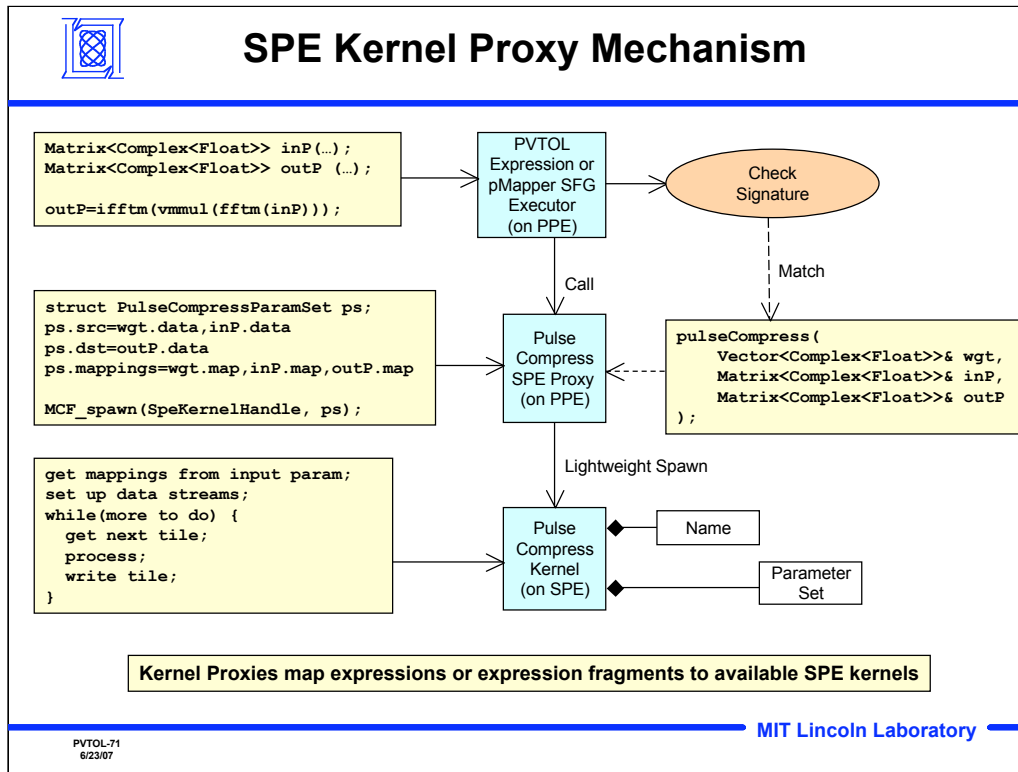
The cell consists of a general purpose processor (PPE) and a collection of 8 special processors (SPE). Each SPEs has its own small memory. The PPE is the manager which manages the SPE workers. The PPE launches a task on the SPE workers. In parallel, the SPEs load the data from main memory to their local memory, perform the computation and write the result back to main memory.



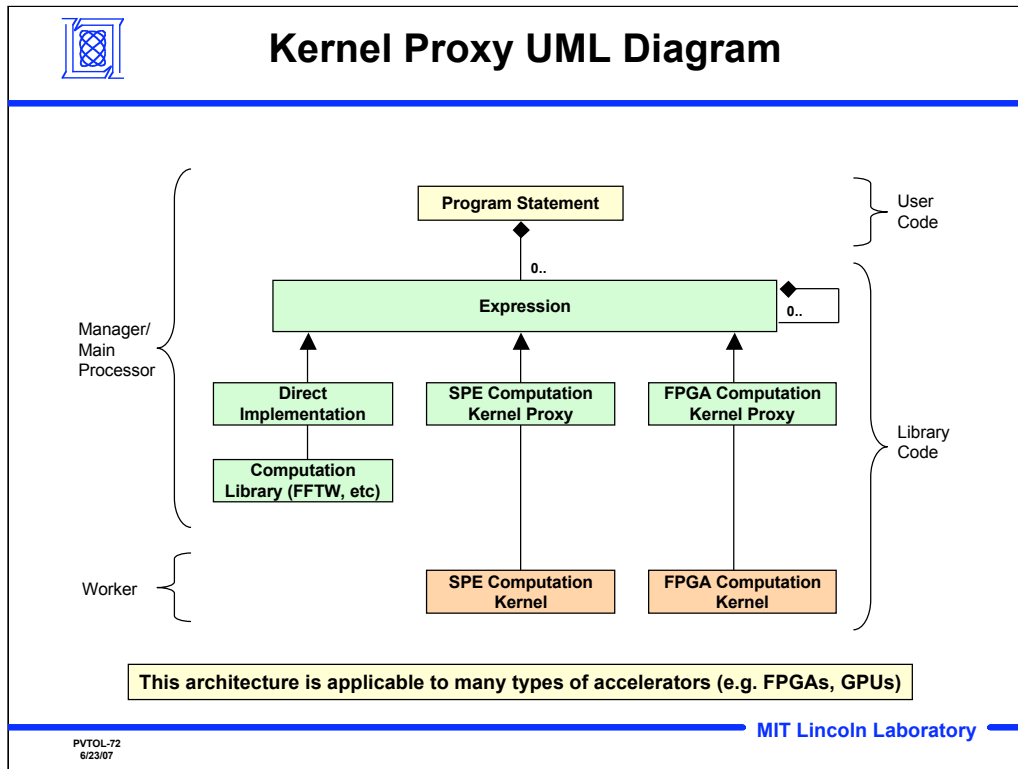
-
- Simple ←————→ Complex
- Compression Keystream loader FFT Pulse compression Doppler filtering STAP

- Kernel Expressions are effectively SPE overlays

70



A pattern-matching scheme is employed to map user code to the available kernels. This mechanism is also useful for composing functions based on the kernels available in the PVTOL library.



The software architecture is broadly applicable to a variety of hardware hosts. The same user code can run on the Cell, on the FPGA, or on the Grid. The PVTOL library API and the dispatch mechanism hides the implementation details from the user.



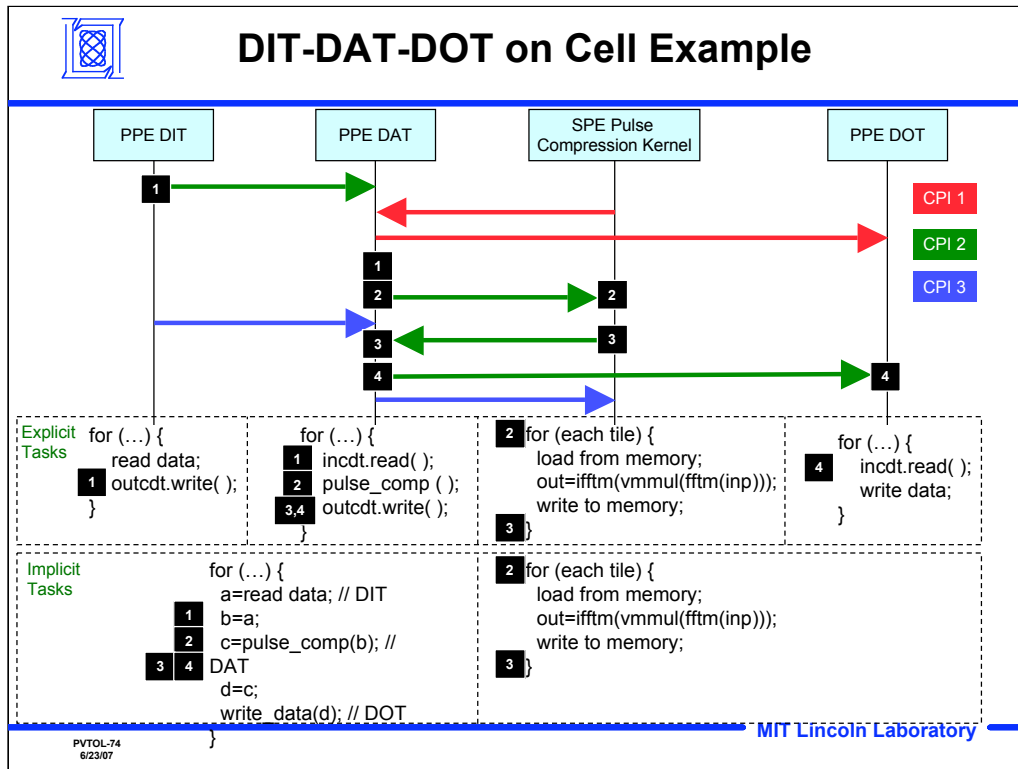
Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - – PVTOL on Cell Example
 - Performance Results
- Summary

PVTOL-73
6/23/07

MIT Lincoln Laboratory

Outline



This slide shows the division of work and the sequence for a Data Input Task(DIT), Data Analysis Task(DAT) and a Data output Task (DOT). The DIT and the DOT are PPE only tasks. The DAT is managed by the PPE but the computationally expensive kernel runs in parallel on the SPEs



Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - – Performance Results
- Summary

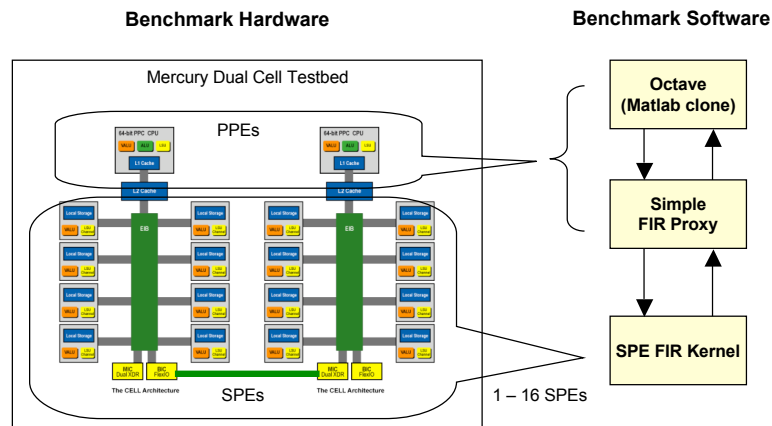
PVTOL-75
6/23/07

MIT Lincoln Laboratory

Outline



Benchmark Description



Based on HPEC Challenge Time Domain FIR Benchmark

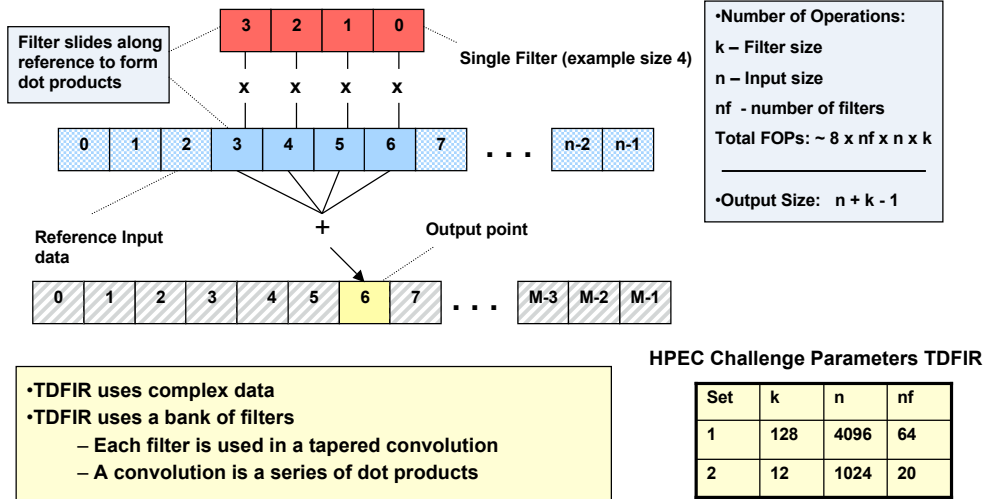
PVTOL-76
6/23/07

MIT Lincoln Laboratory

While the initial timings of the benchmark TDFIR were made local to a single SPU to determine how much performance could be expected from the computational units in a single SPU, we also need to understand how the TDFIR will perform at the chip level. Octave was chosen as the development environment since it is similar to MATLAB which is not available on Cell. From here it was relatively easy to take the high performance TDFIR for a single SPU and create the parallel version.



Time Domain FIR Algorithm



FIR is one of the best ways to demonstrate FLOPS

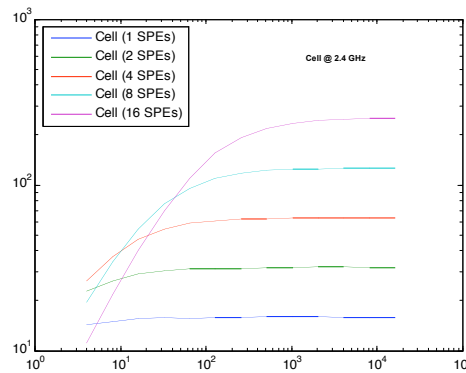
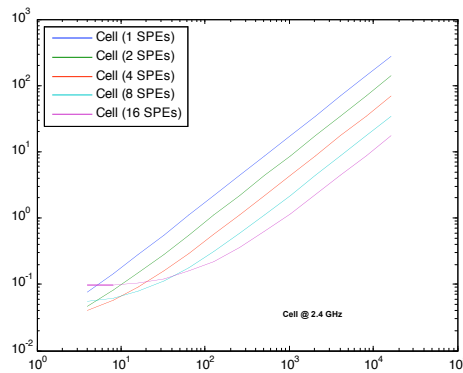
PVTOL-77
6/23/07

MIT Lincoln Laboratory

This is a graphical explanation of the TDFIR which uses complex data. Note that the kernel elements are in reversed order for convolution. Normal order gives correlation.



Performance Time Domain FIR (Set 1)



Set 1 has a bank of 64 size 128 filters with size 4096 input vectors

- Octave runs TDFIR in a loop
 - Averages out overhead
 - Applications run convolutions many times typically

Maximum GFLOPS for TDFIR #1 @ 2.4 GHz

# SPE	1	2	4	8	16
GFLOPS	16	32	63	126	253

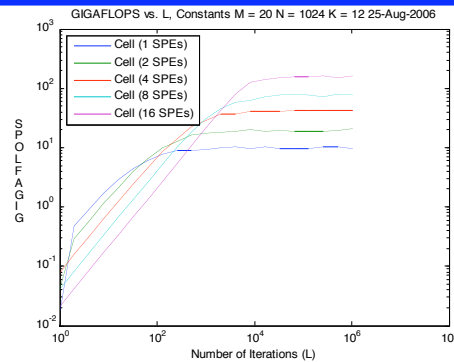
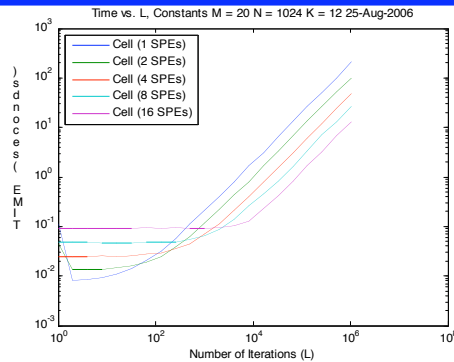
PVTOL-78
6/23/07

MIT Lincoln Laboratory

Here are timing results for TDFIR set 1 @ 2.4GHz. All timings include the full overhead of this application. Note that once the number of iterations become large enough, this application scales nicely with the number of processors. Cell DMAs have a limit of 16KB transfer size. Using the DMA list option is one way of getting around this limitation.



Performance Time Domain FIR (Set 2)



Set 2 has a bank of 20 size 12 filters with size 1024 input vectors

- TDFIR set 2 scales well with the number of processors
 - Runs are less stable than set 1

GFLOPS for TDFIR #2 @ 2.4 GHz

# SPE	1	2	4	8	16
GFLOPS	10	21	44	85	185

PVTOL-79
6/23/07

MIT Lincoln Laboratory

These are the timing results for set 2 of TDFIR which uses smaller convolutions than set 1. Here the overhead needs more iterations to be averaged out since the convolutions are significantly faster. Again, once the overhead is averaged out, these results scale nicely with the number of processors.



Outline

- Introduction
- PVTOL Machine Independent Architecture
 - Machine Model
 - Hierarchal Data Objects
 - Data Parallel API
 - Task & Conduit API
 - pMapper
- PVTOL on Cell
 - The Cell Testbed
 - Cell CPU Architecture
 - PVTOL Implementation Architecture on Cell
 - PVTOL on Cell Example
 - Performance Results



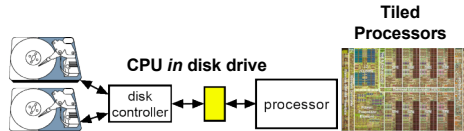
- Summary

Outline



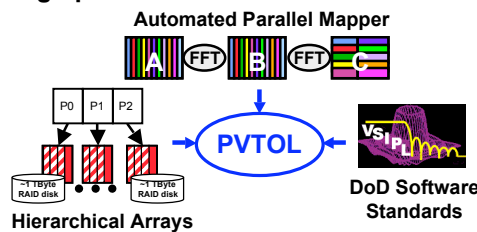
Summary

Goal: Prototype advanced software technologies to exploit novel processors for DoD sensors



- Have *demonstrated* 10x performance benefit of tiled processors
- Novel storage should provide 10x more IO

Approach: Develop **Parallel Vector Tile Optimizing Library (PVTOL)** for high performance and ease-of-use



PVTOL-81
6/23/07

DoD Relevance: Essential for flexible, programmable sensors with large IO and processing requirements



- Wide area data
- Collected over many time scales

Mission Impact:

- Enabler for next-generation synoptic, multi-temporal sensor systems

Technology Transition Plan

- Coordinate development with sensor programs
- Work with DoD and Industry standards bodies

MIT Lincoln Laboratory

The Parallel Vector Tile Optimizing Library (PVTOL) is an effort to develop a new processor architecture for signal processing that exploits the recent shifts in technology to tiled multi-core chips and more tightly integrated mass storage. These technologies are critical for processing the higher bandwidth and longer duration data produced required by synoptic, multi-temporal sensor systems.. The principal challenge in exploiting the new processing architectures for these missions is writing the software in a flexible manner that is portable and delivers high performance. The core technology of this project is the Parallel Vector Tile Optimizing Library (PVTOL), which will use recent advances in automating parallel mapping technology, hierarchical parallel arrays, combined with the Vector, Signal and Image Processing Library (VSIP) open standard to deliver portable performance on tiled processors.