# Advanced Charm++ Tutorial

## Presented by:

## Isaac Dooley & Chao Mei
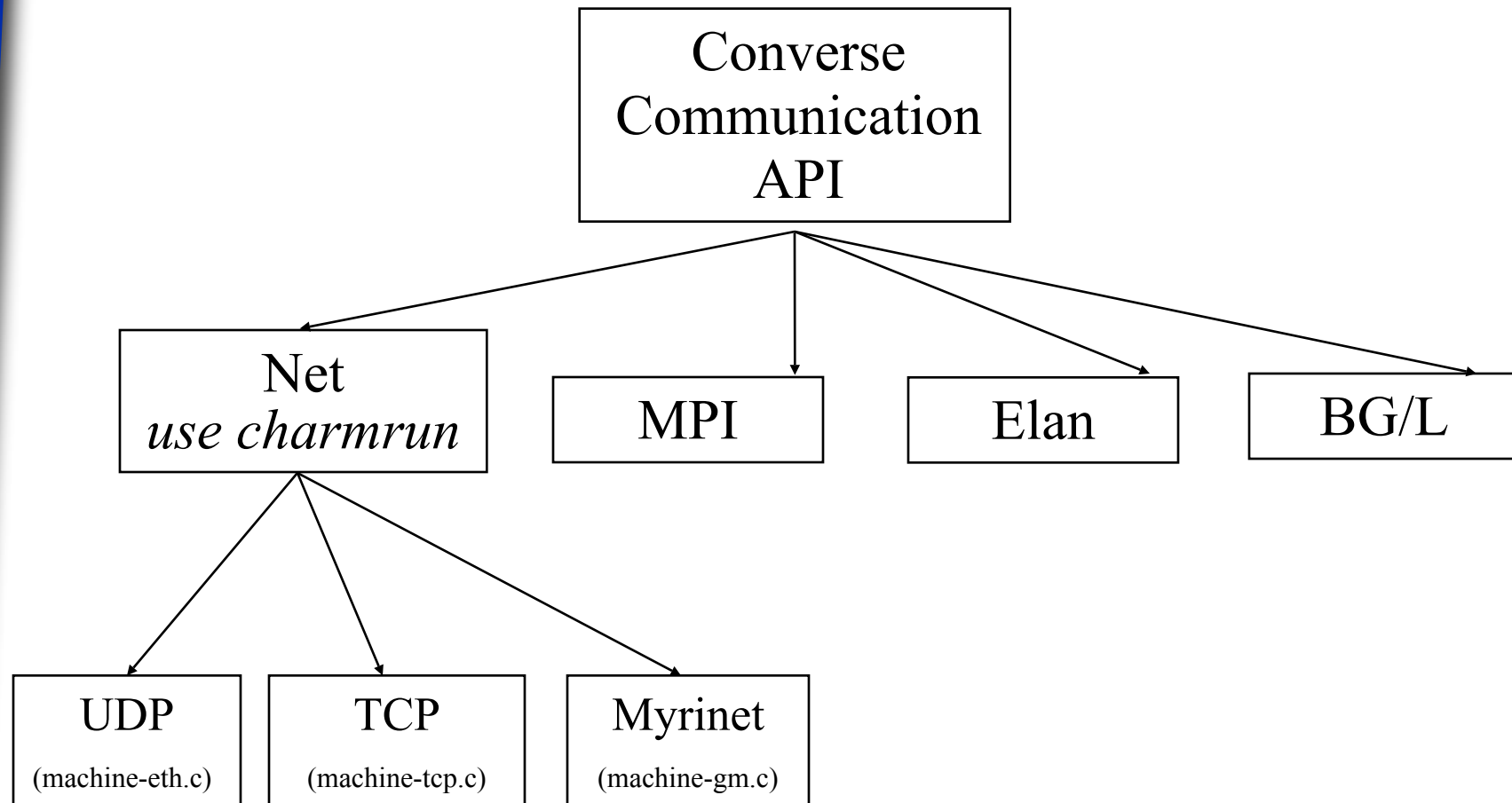
**4/20/2007**

# Topics For This Talk

- **Building Charm++**
- **Advanced messaging**
- **Interface file (.ci)**
- **Advanced load balancing**
- **Groups**
- **Threads**
- **Delegation**
- **Array multicast**
- **SDAG**

# Charm++ on Parallel Machines

- **Runs on:**
  - **Any machine with MPI, including**
    - **IBM Blue Gene/L, SP**
    - **Cray XT3**
    - **SGI Altix**
  - **PSC's Lemieux (Quadrics Elan)**
  - **Clusters with Ethernet (UDP/TCP)**
  - **Clusters with Myrinet (GM or MX)**
  - **Apple clusters**
  - **Even Windows!**
- **SMP-Aware (pthreads)**

# Communication Architecture

# Compiling Charm++

./build

Usage: build **<target>** <version> <options> [charmc-options ...]

**<target>:** converse charm++ LIBS AMPI FEM bigemulator pose jade msa
doc ps-doc pdf-doc html-doc

charm++        compile Charm++ core only
AMPI          compile Adaptive MPI on top of Charm++
FEM           compile FEM framework
LIBS          compile additional parallel libraries with Charm++ core
bigemulator    build additional BigSim libraries
pose          build POSE parallel discrete event simulator
jade          build Jade compiler (auto-builds charm++, msa)
msa           build Multiphase Shared Arrays(MSA) library

# Compiling Charm++

./build

Usage: build \<target\> **\<version\>** \<options\> [charmc-options ...]

**\<version\>: Basic configurations**

| | | |
|---|---|---|
| bluegenel | mpi-sp | net-sol-x86 |
| elan-axp | ncube2 | net-sun |
| elan-linux-ia64 | net-axp | net-win32 |
| exemplar | net-cygwin | net-win64 |
| mpi-axp | net-darwin-x86 | origin-pthreads |
| mpi-bluegenel | net-hp | origin2000 |
| mpi-crayx1 | net-hp-ia64 | portals-crayxt3 |
| mpi-crayxt3 | net-irix | shmem-axp |
| mpi-exemplar | net-linux | sim-linux |
| mpi-hp-ia64 | net-linux-amd64 | sp3 |
| mpi-linux | net-linux-axp | t3e |
| mpi-linux-amd64 | net-linux-cell | uth-linux |
| mpi-linux-axp | net-linux-ia64 | uth-win32 |
| mpi-linux-ia64 | net-linux-ppc | vmi-linux |
| mpi-origin | net-ppc-darwin | vmi-linux-amd64 |
| mpi-ppc-darwin | net-rs6k | vmi-linux-ia64 |
| mpi-sol | net-sol | |
| mpi-sol-amd64 | net-sol-amd64 | |

# Compiling Charm++

./build

Usage: build <target> **<version>** <options> [charmc-options ...]

**<version>: Basic configurations**

| | | |
|---|---|---|
| bluegenel | mpi-sp | net-sol-x86 |
| elan-axp | ncube2 | net-sun |
| elan-linux-ia64 | net-axp | net-win32 |
| exemplar | net-cygwin | net-win64 |
| mpi-axp | net-darwin-x86 | origin-pthreads |
| mpi-bluegenel | net-hp | origin2000 |
| mpi-crayx1 | net-hp-ia64 | portals-crayxt3 |
| mpi-crayxt3 | net-irix | shmem-axp |
| mpi-exemplar | net-linux | sim-linux |
| mpi-hp-ia64 | net-linux-amd64 | sp3 |
| mpi-linux | net-linux-axp | t3e |
| mpi-linux-amd64 | net-linux-cell | uth-linux |
| mpi-linux-axp | net-linux-ia64 | uth-win32 |
| mpi-linux-ia64 | net-linux-ppc | vmi-linux |
| mpi-origin | net-ppc-darwin | vmi-linux-amd64 |
| mpi-ppc-darwin | net-rs6k | vmi-linux-ia64 |
| mpi-sol | net-sol | |
| mpi-sol-amd64 | net-sol-amd64 | |

7

# Compiling Charm++

./build

Usage: build <target> <version> **<options>** [charmc-options ...]

**<options>: compiler and platform specific options**

**Platform specific options** (choose multiple if they apply)**:**

```
lam              Use LAM MPI
 smp             support for SMP, multithreaded charm on each node
mpt              use SGI Message Passing Toolkit (only for mpi version)
gm               use Myrinet for communication
tcp              use TCP sockets for communication (ony for net version)
vmi              use NCSA's VMI for communication (only for mpi version)
scyld            compile for Scyld Beowulf cluster based on bproc
clustermatic     compile for Clustermatic (support version 3 and 4)
pthreads         compile with pthreads Converse threads
```

# Compiling Charm++

./build

Usage: build <target> <version> **<options>** [charmc-options ...]

**<options>: compiler and platform specific options**

**Advanced options:**

```
bigemulator      compile for BigSim simulator
ooc              compile with out of core support
syncft           compile with Charm++ fault tolerance support
papi             compile with PAPI performance counter support (if any)
```

**Charm++ dynamic libraries:**

```
--build-shared     build Charm++ dynamic libraries (.so) (default)
--no-build-shared  don't build Charm++'s shared libraries
```

# Compiling Charm++

./build

Usage: build \<target\> \<version\> **\<options\>** [charmc-options ...]

**\<options\>: compiler and platform specific options**

Choose a C++ compiler (only one option is allowed from this section):

```
cc, cc64        For Sun WorkShop C++ 32/64 bit compilers
cxx             DIGITAL C++ compiler (DEC Alpha)
kcc             KAI C++ compiler
pgcc            Portland Group's C++ compiler
acc             HP aCC compiler
icc             Intel C/C++ compiler for Linux IA32
ecc             Intel C/C++ compiler for Linux IA64
gcc3            use gcc3 - GNU GCC/G++ version 3
gcc4            use gcc4 - GNU GCC/G++ version 4 (only mpi-crayxt3)
mpcc            SUN Solaris C++ compiler for MPI
pathscale       use pathscale compiler suite
```

# Compiling Charm++

./build

Usage: build <target> <version> **<options>** [charmc-options ...]

**<options>: compiler and platform specific options**

Choose a fortran compiler (only one option is allowed from this section):

```
g95             G95 at http://ww.g95.org
absoft          Absoft fortran compiler
pgf90           Portland Group's Fortran compiler
ifc             Intel Fortran compiler (older versions)
ifort           Intel Fortran compiler (newer versions)
```

# Compiling Charm++

./build

Usage: build \<target\> \<version\> \<options\> **[charmc-options ...]**

**\<charmc-options\>:** normal compiler options

```
-g -O -save -verbose
```

To see the latest versions of these lists or to get more detailed help, run

```
./build --help
```

# Build Script

- **Build script does:**

  `./build <target> <version> <options> [charmc-options ...]`

  - Creates directories `<version>` and `<version>/tmp`

  -  Copies `src/scripts/Makefile` into `<version>/tmp`

  -  Does a
    `"make <target> <version> OPTS=<charmc-options>"`
    in `<version>/tmp`

- **That's all build does.  The rest is handled by the Makefile.**

# How 'build' works

- **build AMPI net-linux gm kcc**
  - **Mkdir net-linux-gm-kcc**
  - **Cat conv-mach-[kcc|gm|smp].h to conv-mach-opt.h**
  - **Cat conv-mach-[kcc|gm].sh to conv-mach-opt.sh**
  - **Gather files from net, etc (Makefile)**
  - **Make charm++ under**
    - **net-linux-gm/tmp**

# What if build fails?

- **Use latest version from CVS**
- **Check the nightly auto-build tests: http://charm.cs.uiuc.edu/autobuild/cur/**

- **Email: ppl@cs.uiuc.edu**

# How Charmrun Works?

Charmrun

**charmrun +p4 ./pgm**

# How Charmrun Works?



Charmrun

**charmrun +p4 ./pgm**

ssh

connect

Charmrun

**charmrun +p4 ./pgm**

# How Charmrun Works?



ssh

connect

Acknowledge

Charmrun

**charmrun +p4 ./pgm**

16

# Charmrun (batch mode)

Charmrun

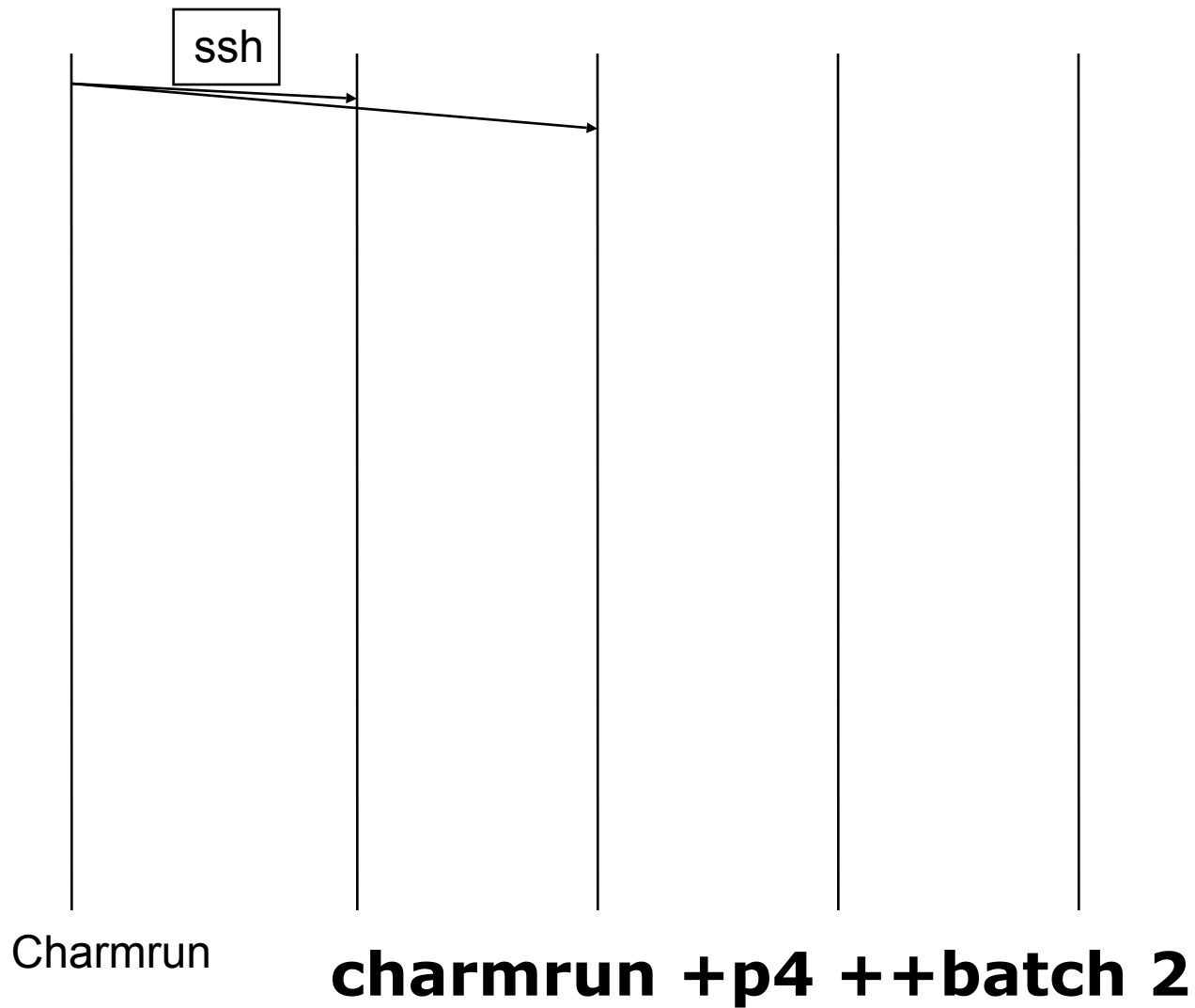**charmrun +p4 ++batch 2**

# Charmrun (batch mode)

ssh

Charmrun

**charmrun +p4 ++batch 2**

# Charmrun (batch mode)



ssh

connect

Charmrun

**charmrun +p4 ++batch 2**

# Charmrun (batch mode)

ssh

connect

Charmrun

**charmrun +p4 ++batch 2**

ssh

connect

Charmrun

**charmrun +p4 ++batch 2**

17

# Charmrun (batch mode)



ssh

connect

Acknowledge

Charmrun

**charmrun +p4 ++batch 2**

# Debugging Charm++ Applications

- **printf**
- **Gdb**
  - **Sequentially (standalone mode)**
    - gdb ./pgm +vp16
  - **Attach gdb manually**
  - **Run debugger in xterm**
    - charmrun +p4 pgm ++debug
    - charmrun +p4 pgm ++debug-no-pause
  - **Memory paranoid**
    - -memory paranoid
  - **Parallel debugger**

# Advanced Messaging

# Prioritized Execution

- **Charm++ scheduler**
  - **Default - FIFO (oldest message)**
- **Prioritized execution**
  - **If several messages available, Charm will process the messages in the order of their priorities**
- **Very useful for speculative work, ordering timestamps, etc...**

# Priority Classes

- **Charm++ scheduler has three queues: high, default, and low**
- **As signed integer priorities:**
  - **High -MAXINT to -1**
  - **Default 0**
  - **Low 1 to +MAXINT**
- **As unsigned bitvector priorities:**
  - **0x0000 Highest priority -- 0x7FFF**
  - **0x8000 Default priority**
  - **0x8001 -- 0xFFFF Lowest priority**

# Prioritized Messages

- **Number of priority bits passed during message allocation**

```
FooMsg * msg = new (size, nbits) FooMsg;
```

  - **Priorities stored at the end of messages**

- **Signed integer priorities**

```
*CkPriorityPtr(msg)=-1;

CkSetQueueing(msg, CK_QUEUEING_IFIFO);
```

- **Unsigned bitvector priorities**

```
CkPriorityPtr(msg)[0]=0x7fffffff;

CkSetQueueing(msg, CK_QUEUEING_BFIFO);
```

# Prioritized Marshalled Messages

- **Pass "CkEntryOptions" as last parameter**

- **For signed integer priorities:**

```
CkEntryOptions opts;

opts.setPriority(-1);

fooProxy.bar(x,y,opts);
```

- **For bitvector priorities:**

```
CkEntryOptions opts;

unsigned int prio[2]={0x7FFFFFFF,0xFFFFFFFF};

opts.setPriority(64,prio);

fooProxy.bar(x,y,opts);
```

# Advanced Message Features

- **Read-only messages**
  - **Entry method agrees not to modify or delete the message**
  - **Avoids message copy for broadcasts, saving time**
- **Inline messages**
  - **Direct method invocation if on local processor**
- **Expedited messages**
  - **Message do not go through the charm++ scheduler (ignore any Charm++ priorities)**
- **Immediate messages**
  - **Entries are executed in an interrupt or the communication thread**
  - **Very fast, but tough to get right**
  - **Immediate messages only currently work for NodeGroups and Group (non-smp)**

24

## All declared in the .ci file

```
{
    entry [nokeep] void foo_readonly(Msg *);
    entry [inline] void foo_inl(Msg *);
    entry [expedited] void foo_exp(Msg *);
    entry [immediate] void foo_imm(Msg *);
    ...
};
```

# Interface File (ci)

# Interface File Example

```
mainmodule hello {
  include "myType.h"

  initnode void myNodeInit();
  initproc void myInit();

  mainchare mymain {
    entry mymain(CkArgMsg *m);
  };


  array[1D] foo {
    entry foo(int problemNo);
    entry void bar1(int x);
    entry void bar2(myType x);
  };
};
```

# Include and Initcall

- **Include**
  - **Include an external header files**
- **Initcall**
  - **User plugging code to be invoked in Charm++'s startup phase**
  - **Initnode**
    - Called once on every node
  - **Initproc**
    - Called once on every processor
  - **Initnode calls are called before Initproc calls**

# Entry Attributes

- **Threaded**
  - **Function is invoked in a CthThread**
- **Sync**
  - **Blocking methods, can return values as a message**
  - **Caller must be a thread**
- **Exclusive**
  - **For Node Group**
  - **Do not execute while other exclusive entry methods of its node group are executing in the same node**
- **Notrace**
  - **Invisible to trace projections**
  - entry [notrace] void recvMsg(multicastGrpMsg *m);

# Groups/Node Groups

30

# Groups and Node Groups

- **Groups**
  - **Similar to arrays:**
    - **Broadcasts, reductions, indexing**
  - **But not completely like arrays:**
    - **Non-migratable; one per processor**
  - **Exactly one representative on each processor**
    - **Ideally suited for system libraries**
  - **Historically called branch office chares (BOC)**
- **Node Groups**
  - **One per SMP node**

# Declarations

- **.ci file**

```
group mygroup {
    entry mygroup(); //Constructor
    entry void foo(foomsg *); //Entry method
};
nodegroup mynodegroup {
    entry mynodegroup(); //Constructor
    entry void foo(foomsg *); //Entry method
};
```

- **C++ file**

```
class mygroup : public Group {
            mygroup() {}
    void foo(foomsg *m) { CkPrintf("Do Nothing");}
};
class mynodegroup : public NodeGroup {
            mynodegroup() {}
    void foo(foomsg *m) { CkPrintf("Do Nothing");}
};
```

# Creating and Calling Groups

- **Creation**

  ```
  p = CProxy_mygroup::ckNew();
  ```

- **Remote invocation**

  ```
  p.foo(msg);      //broadcast
  p[1].foo(msg);   //asynchronous
  p.foo(msg, npes, pes); // list send
  ```

- **Direct local access**

  ```
  mygroup *g=p.ckLocalBranch();
  g->foo(….); //local invocation
  ```

  - **Danger: if you migrate, the group stays behind!**

# Advanced Load-balancers

- **Inherit from CentralLB and implement the work(…) function**

class foolb : public CentralLB {

   public:

      .. .. ..

      void work (CentralLB::LDStats*   stats, int count);

      .. .. ..

  };

```
struct LDStats {
    ProcStats *procs;
    LDObjData* objData;
    LDCommData* commData;
    int *to_proc;
    //.. .. ..
}
//Dummy Work function which assigns all objects to
//processor 0
//Don't implement it!
void fooLB::work(CentralLB::LDStats* stats,int count){
     for(int count=0;count < nobjs; count++)
        stats.to_proc[count] = 0;
}
```

# Compiling and Integration

- **Edit and run Makefile_lb.sh**
  - **Creates Make.lb which is included by the main Makefile**
- **Run make depends to correct dependencies**
- **Rebuild charm++ and is now available in –balancer fooLB**

# Threads in Charm++

# Why use Threads?

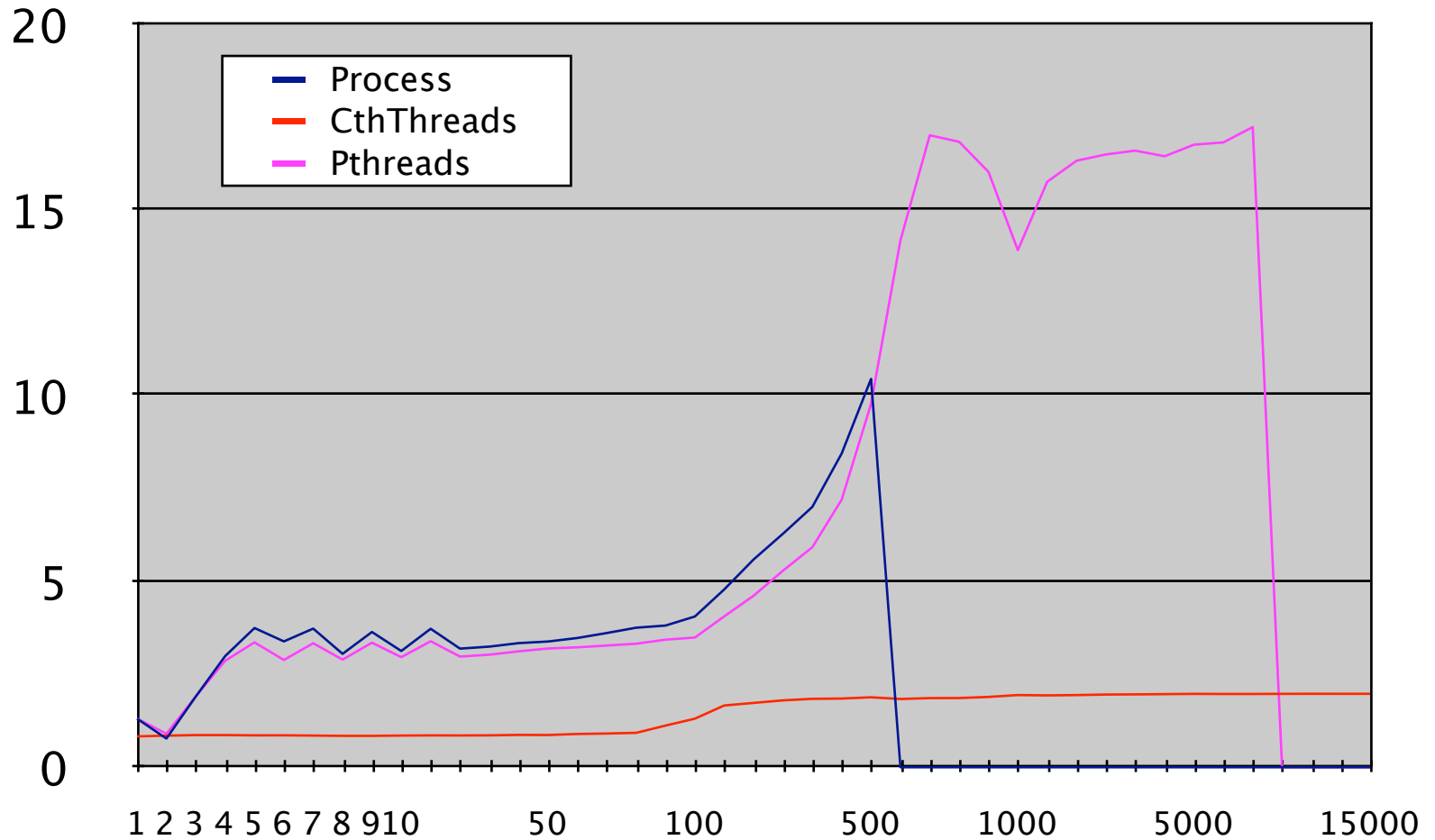- **They provide one key feature: blocking**
  - **Suspend execution (e.g., at message receive)**
  - **Do something else**
  - **Resume later (e.g., after message arrives)**
- **Example: MPI_Recv, MPI_Wait semantics**
- **Function call interface more convenient than message-passing**
  - **Regular call/return structure (no CkCallbacks) with complete control flow**
  - **Allows blocking in middle of deeply**

# Why <u>not</u> use Threads?

- **Slower**
  - **Around 1us context-switching overhead unavoidable**
  - **Creation/deletion perhaps 10us**
- **Migration more difficult**
  - **State of thread is scattered through stack, which is maintained by compiler**
  - **By contrast, state of object is maintained by users**
- **Thread disadvantages form the motivation to use SDAG (later)**

# Context Switch Cost

# What are (Converse) Threads?

- **One flow of control (instruction stream)**
  - **Machine Registers & program counter**
  - **Execution stack**
- **Like pthreads (kernel threads)**
- **Only different:**
  - **Implemented at user level (in Converse)**
  - **Scheduled at user level; non-preemptive**
  - **Migratable between nodes**

# How do I use Threads?

- **Many options:**
  - **AMPI**
    - Always uses threads via TCharm library
  - **Charm++**
    - [threaded] entry methods run in a thread
    - [sync] methods
  - **Converse**
    - C routines CthCreate/CthSuspend/CthAwaken
    - Everything else is built on these
    - Implemented using
      - SYSV makecontext/setcontext
      - POSIX setjmp/alloca/longjmp
      - Assembly code

# How do I use Threads (example)

- **Blocking API routine: find array element**

```
int requestFoo(int src) {
    myObject *obj=...;
    return obj->fooRequest(src)
}
```

- **Send request and suspend**

```
int myObject::fooRequest(int src) {
    proxy[dest].fooNetworkRequest(thisIndex);
    stashed_thread=CthSelf();
    CthSuspend();    // -- blocks until awaken call --
    return stashed_return;
}
```

- **Awaken thread when data arrives**

```
void myObject::fooNetworkResponse(int ret) {
    stashed_return=ret;
    CthAwaken(stashed_thread);
}
```

- **Send request, suspend, recv, awaken, return**

```
int myObject::fooRequest(int src) {
  proxy[dest].fooNetworkRequest(thisIndex);
  stashed_thread=CthSelf();
  CthSuspend();
        void myObject::fooNetworkResponse(int ret) {
          stashed_return=ret;
          CthAwaken(stashed_thread);
        }

  return stashed_return;
}
```

# Thread Migration

# Stack Data

- **The stack is used by the compiler to track function calls and provide temporary storage**
  - **Local Variables**
  - **Subroutine Parameters**
  - **C "alloca" storage**
- **Most of the variables in a typical application are stack data**
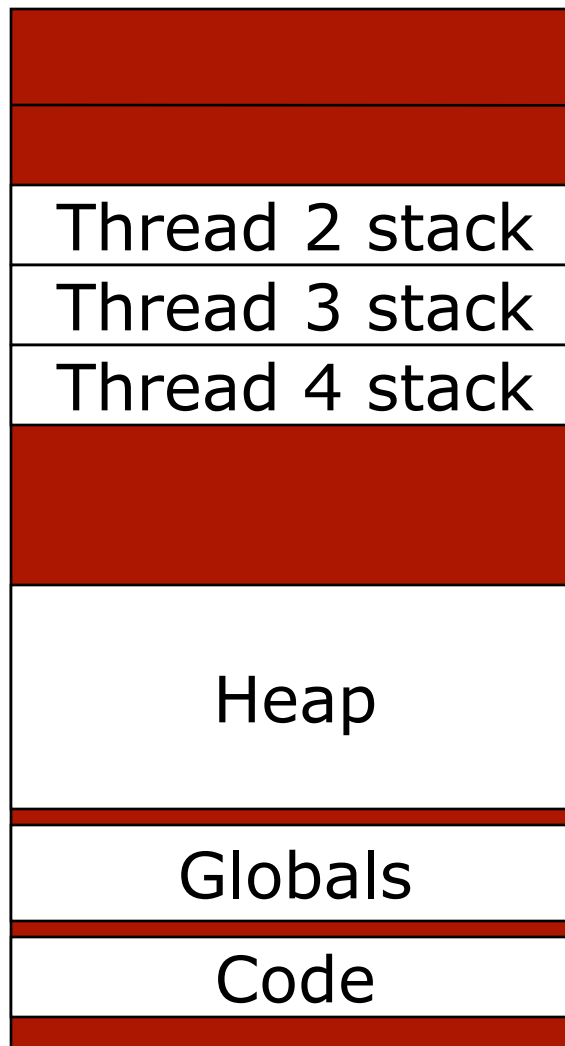- **Stack is allocated by Charm run-time as heap memory (+stacksize)**

# Migrate Stack Data

- **Without compiler support, cannot change stack's address**
    - **Because we can't change stack's interior pointers (return frame pointer, function arguments, etc.)**
- **Existing pointers to addresses in original stack become invalid**
- **Solution: "isomalloc" addresses**
    - **Reserve address space on every processor for every thread stack**
    - **Use *mmap* to scatter stacks in virtual memory efficiently**
    - **Idea comes from PM$^2$**
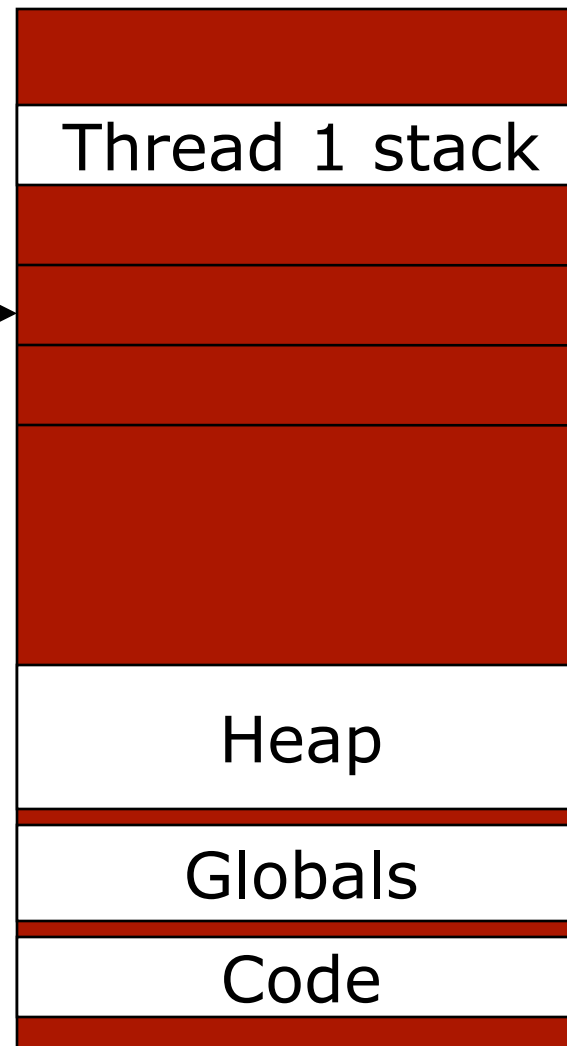
# Migrate Stack Data

Processor A's Memory

0xFFFFFFFF

| Thread 2 stack |
| Thread 3 stack |
| Thread 4 stack |

Migrate
Thread 3 →

Heap

Globals

Code

0x00000000

Processor B's Memory

0xFFFFFFFF

| Thread 1 stack |

Heap

Globals

Code

0x00000000

49

# Migrate Stack Data: Isomalloc

Processor A's Memory

Processor B's Memory

0xFFFFFFFF

0xFFFFFFFF

Thread 1 stack

Thread 2 stack

Migrate
Thread 3

Thread 3 stack

Thread 4 stack
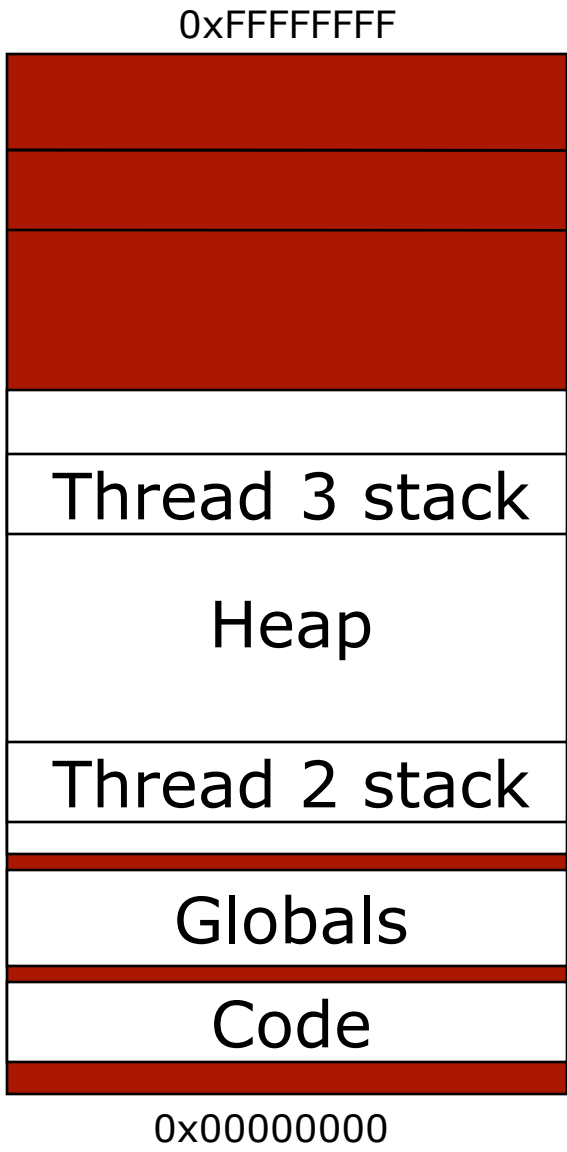
Heap

Heap

Globals

Globals

Code

Code

0x00000000

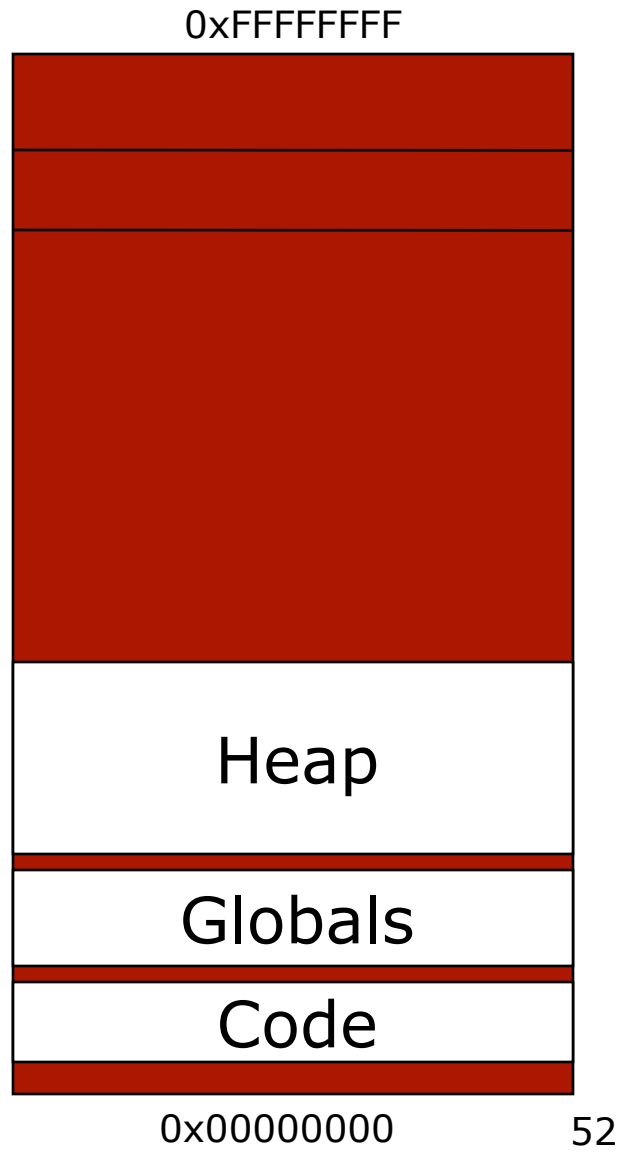0x00000000

# Migrate Stack Data

- **Isomalloc is a completely automatic solution**
  - **No changes needed in application or compilers**
  - **Just like a software shared-memory system, but with proactive paging**
- **But has a few limitations**
  - **Depends on having large quantities of virtual address space (best on 64-bit)**
    - **32-bit machines can only have a few gigs of isomalloc stacks across the whole machine**
  - **Depends on unportable *mmap***
    - **Which addresses are safe? (We must guess!)**
    - **What about Windows? Or Blue Gene?**

# Aliasing Stack Data

Processor A's Memory

0xFFFFFFFF

Thread 3 stack

Heap

Thread 2 stack

Globals

Code

0x00000000

Processor B's Memory

0xFFFFFFFF

Heap

Globals

Code

0x00000000

Processor A's Memory

0xFFFFFFFF

Thread 2 stack

Execution Copy

Thread 3 stack

Heap

Thread 2 stack

Globals

Code

0x00000000

Processor B's Memory

0xFFFFFFFF

Heap

Globals

Code

0x00000000

# Aliasing Stack Data

Processor A's Memory

0xFFFFFFFF

Thread 3 stack

Heap

Thread 2 stack

Globals

Code

0x00000000

Processor B's Memory

0xFFFFFFFF

Heap

Globals

Code

0x00000000

## Processor A's Memory

0xFFFFFFFF

Thread 3 stack

Execution Copy

Thread 3 stack

Heap

Thread 2 stack

Globals

Code

0x00000000

## Processor B's Memory

0xFFFFFFFF

Heap

Globals

Code

0x00000000

# Aliasing Stack Data

**Processor A's Memory**

0xFFFFFFFF

Thread 3 stack

Heap

Thread 2 stack

Globals

Code

0x00000000

**Migrate Thread 3**

**Processor B's Memory**

0xFFFFFFFF

Thread 3 stack

Heap

Globals

Code

0x00000000

56

# Aliasing Stack Data

Processor A's Memory

0xFFFFFFFF

Heap

Thread 2 stack

Globals

Code

0x00000000

Processor B's Memory

0xFFFFFFFF

Thread 3 stack

Heap

Globals

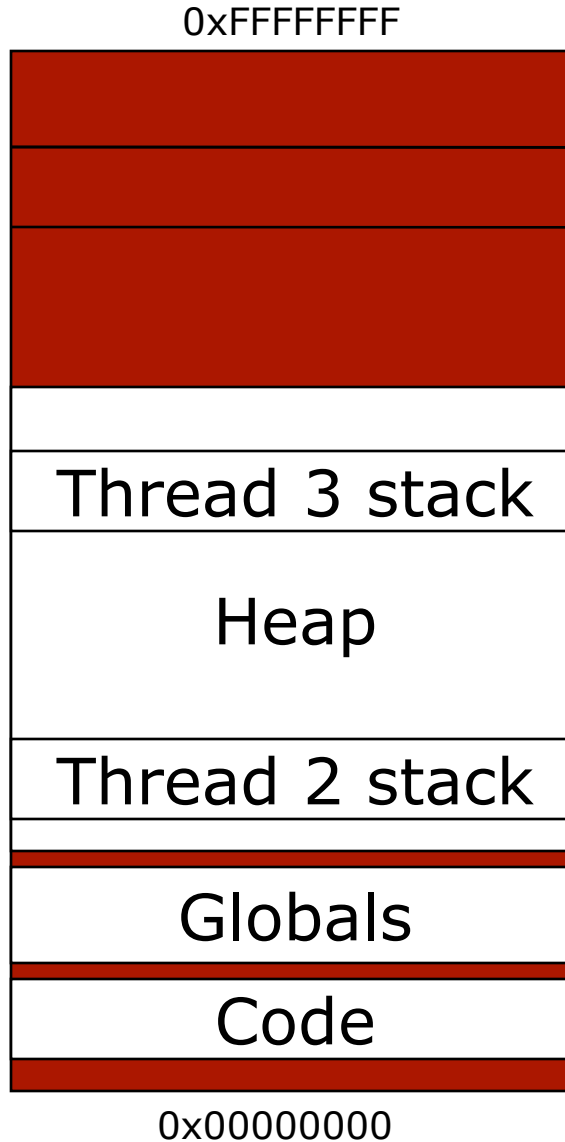Code

0x00000000

# Aliasing Stack Data

Processor A's Memory

Processor B's Memory

0xFFFFFFF

0xFFFFFFF

Execution Copy

Thread 3 stack

Heap

Thread 3 stack

Heap

Thread 2 stack

Globals

Globals

Code

Code

0x00000000

0x00000000

# Aliasing Stack Data

- **Does not depend on having large quantities of virtual address space**
  - **Works well on 32-bit machines**
- **Requires only one mmap'd region at a time**
  - **Works even on Blue Gene!**
- **Downsides:**
  - **Thread context switch requires munmap/mmap (3us)**
  - **Can only have one thread running at a time (so no SMP's!)**
- **"-thread memoryalias" link time option**

# Heap Data

- **Heap data is any dynamically allocated data**
  - **C "malloc" and "free"**
  - **C++ "new" and "delete"**
  - **F90 "ALLOCATE" and "DEALLOCATE"**
- **Arrays and linked data structures are almost always heap data**

# Migrate Heap Data

- **Automatic solution: isomalloc all heap data just like stacks!**
  - **"-memory isomalloc" link option**
  - **Overrides *malloc/free***
  - **No new application code needed**
  - **Same limitations as isomalloc; page allocation granularity (huge!)**
- **Manual solution: application moves its heap data**
  - **Need to be able to <u>size</u> message buffer, <u>pack</u> data into message, and <u>unpack</u> on other side**
  - **"pup" abstraction does all three**

61

# Delegation

# Delegation

- **Customized implementation of messaging**
  - **Enables Charm++ proxy messages to be forwarded to a delegation manager group**
- **Delegation manager**
  - **trap calls to proxy sends and apply optimizations**
- **Delegation manager must inherit from CkDelegateMgr class**
- **User program must to call**
  - **proxy.ckDelegate(mgrID);**

# Delegation Interface

- **.ci file**

```
group MyDelegateMgr {
        entry MyDelegateMgr();  //Constructor
  };
```

- **.h file**

```
class MyDelegateMgr : public CkDelegateMgr {
  MyDelegateMgr();
  void ArraySend(...,int ep,void *m,const
  CkArrayIndexMax &idx,CkArrayID a);
  void ArrayBroadcast(..);
  void ArraySectionSend(.., CkSectionID &s);
      …………..
      …………..
}
```

# Array Multicast

- **Array section – a subset of chare array**
- **Array section creation**
  - **Enumerate array indices**
    - CkVec<CkArrayIndex3D> elems;    // add array indices
      for (int i=0; i<10; i++)
        for (int j=0; j<20; j+=2)
          for (int k=0; k<30; k+=2)
            elems.push_back(CkArrayIndex3D(i, j, k));
      CProxySection_Hello proxy = CProxySection_Hello::ckNew(helloArrayID, elems.getVec(), elems.size());

  - **Alternatively, one can do the same thing by providing (lbound:ubound:stride) for each dimension:**
    - CProxySection_Hello proxy = CProxySection_Hello::ckNew(helloArrayID, 0, 9, 1, 0, 19, 2, 0, 29, 2);
    - **The above code creates a section proxy that contains array elements of [0:9, 0:19:2, 0:29:2].**

  - **For user-defined array index other than CkArrayIndex1D to CkArrayIndex6D, one needs to use the generic array index type: CkArrayIndexMax.**
    - CkArrayIndexMax *elems;    // add array indices
      int numElems;
      CProxySection_Hello proxy = CProxySection_Hello::ckNew(helloArrayID, elems, numElems);

# Array Section Multicast

- **Once have the array section proxy**
  - **do multicast to all the section members:**
    - CProxySection_Hello proxy;
      proxy.foo( msg)        // multicast

  - **send messages to one member using its local index**
    - proxy[0].foo( msg)

# Array Section Multicast

- **Multicast via delegation**
    - **CkMulticast communication library**
    - CProxySection_Hello sectProxy = CProxySection_Hello::ckNew();
      CkGroupID mCastGrpId = CProxy_CkMulticastMgr::ckNew();
      CkMulticastMgr *mcastGrp = CProxy_CkMulticastMgr
      (mCastGrpId).ckLocalBranch();

      sectProxy.ckSectionDelegate(mCastGrpId);  // initialize proxy

      sectProxy.foo(...);          //multicast via delegation

- **Note, to use CkMulticast library, all multicast messages must inherit from CkMcastBaseMsg, as following:**
  class HiMsg : public CkMcastBaseMsg, public CMessage_HiMsg
  {
      public:
      int *data;
  };

# Array Section Reduction

- **Section reduction with delegation**
- **Use default reduction callback**

```
CProxySection_Hello sectProxy;

CkMulticastMgr *mcastGrp = CProxy_CkMulticastMgr
    (mCastGrpId).ckLocalBranch();

mcastGrp->setReductionClient(sectProxy, new CkCallback(...));
```

- **Reduction**

```
CkGetSectionInfo(sid, msg);

CkCallback cb(CkIndex_myArray::foo(NULL),thisProxy);

mcastGrp->contribute(sizeof(int), &data, CkReduction::sum_int, sid,
    cb);
```

# With Migration

- **Works with migration**
  - **When intermediate nodes migrate**
    - **When migrates, multicast tree will be automatically rebuilt**
  - **Root processor**
    - **Application needs to initiate the rebuild**
    - **Will change to automatic in future**

# SDAG

# Structured Dagger

- **What is it?**
  - **A coordination language built on top of Charm++**
  - **Express control flow in interface file**
- **Motivation**
  - **Charm++'s asynchrony is efficient and reliable, but tough to program**
    - Split phase - Flags, buffering, out-of-order receives, etc.
  - **Threads are easy to program, but less efficient and less reliable**
    - Implementation complexity
    - Porting headaches
  - **Want benefits of both!**

# Structured Dagger Constructs

- **when <method list> {code}**
  - **Do not continue until method is called**
    - Internally generates flags, checks, etc.
    - Does not use threads
- **atomic {code}**
  - **Call ordinary sequential C++ code**
- **if/else/for/while**
  - **C-like control flow**
- **overlap {code1 code2 ...}**
  - **Execute code segments in parallel**
- **forall**
  - **"Parallel Do"**
  - **Like a parameterized overlap**

```
array[1D] myArray {
…
entry void GetMessages () {
    when rightmsgEntry(), leftmsgEntry() {
        atomic {
          CkPrintf("Got both left and right messages \n");
          doWork(right, left); }
      }
};


entry void rightmsgEntry();
entry void leftmsgEntry();
…
};
```

# Overlap for LeanMD Initialization

```
array[1D] myArray {
…
 entry void waitForInit(void) {
     overlap {
       when recvNumCellPairs(myMsg* pMsg) {
         atomic { setNumCellPairs(pMsg->intVal);  delete
pMsg; }
       }
       when recvNumCells(myMsg * cMsg) {
         atomic { setNumCells(cMsg->intVal);  delete cMsg; }
       }
     }
 }
};
```

```
entry void doTimeloop(void) {
    for (timeStep_=1; timeStep_<=SimParam.NumSteps; timeStep++) {
        atomic {sendAtomPos(); }

        overlap {
          for (forceCount_=0; forceCount_<numForceMsg_; forceCount_++) {
            when recvForces(ForcesMsg* msg) { atomic {procForces(msg); } }
          }

          for (pmeCount_=0; pmeCount_<nPME; pmeCount_++) {
            when recvPME(PMEGridMsg* m) {atomic {procPME(m);}}
          }
        }

        atomic { doIntegration(); }

        if (timeForMigrate()) { ... }
    }
}
```

# Thank You!

**Free source, binaries, manuals, and more information at:**
**http://charm.cs.uiuc.edu/**

**Parallel Programming Lab**
**at University of Illinois**