



Hewlett Packard
Enterprise

Slides for U of A Q&A

Brad Chamberlain

April 24, 2025

A Bit About Me

2001: Completed PhD at UW CSE focusing on ZPL

- Chapel's arrays and domains are heavily influenced by ZPL
- Chapel is generally far more powerful and general than ZPL



2001: Worked at a start-up designing a parallel language for reconfigurable embedded hardware

2002: Joined Cray Inc.

2003: Co-launched the Chapel project, creating initial drafts of the language and implementation

2006: Became the technical lead of the Chapel project

2020: Cray Inc. was acquired by HPE





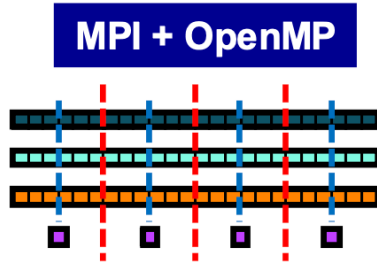
Chapel's Motivation: Improve the State-of-the-Art for HPC Programming

STREAM TRIAD: IN MPI+OPENMP

```
#include <hpcc.h>
#ifdef _OPENMP
#include <omp.h>
#endif
```

```
static int VectorSize;
static double *a, *b, *c;
```

```
int HPCC_StarStream(HPCC_Params *params) {
    int myRank, commSize;
```



```
if (!a || !b || !c) {
    if (c) HPCC_free(c);
    if (b) HPCC_free(b);
    if (a) HPCC_free(a);
    if (doIO) {
        fprintf( outFile, "Failed to
            allocate memory (%d).\n",
                VectorSize );
        fclose( outFile );
    }
}
```

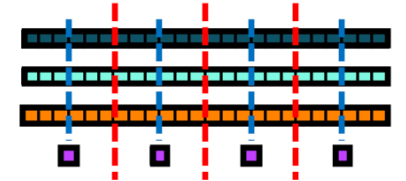
```
#define N 2000000
```

```
int main() {
    float *d_a, *d_b, *d_c;
    float scalar;

    cudaMalloc((void**)&d_a, sizeof(float)*N);
    cudaMalloc((void**)&d_b, sizeof(float)*N);
    cudaMalloc((void**)&d_c, sizeof(float)*N);
```

```
    dim3 dimBlock(128);
```

CUDA



*HPC suffers from too many distinct notations for expressing parallelism and locality.
This tends to be a result of bottom-up language design.*

```
rv = HPCC_Stream( params, 0 - myRank );
MPI_Reduce( &rv, &errCount, 1, MPI_INT, MPI_SUM,
    0, comm );

return errCount;
}
```

```
int HPCC_Stream(HPCC_Params *params, int doIO) {
    register int j;
    double scalar;
```

```
VectorSize = HPCC_LocalVectorSize( params, 3,
    sizeof(double), 0 );
```

```
a = HPCC_XMALLOC( double, VectorSize );
b = HPCC_XMALLOC( double, VectorSize );
c = HPCC_XMALLOC( double, VectorSize );
```

```
for (j=0; j<VectorSize; j++) {
    b[j] = 2.0;
    c[j] = 1.0;
}
scalar = 3.0;
```

```
#ifdef _OPENMP
#pragma omp parallel for
#endif
```

```
for (j=0; j<VectorSize; j++)
    a[j] = b[j]+scalar*c[j];
```

```
HPCC_free(c);
HPCC_free(b);
HPCC_free(a);

return 0; }
```

```
scalar = 3.0;
STREAM_Triad<<<dimGrid,dimBlock>>>>(d_b, d_c, d_a, scalar, N);
cudaThreadSynchronize();
```

```
cudaFree(d_a);
cudaFree(d_b);
cudaFree(d_c);
```

```
__global__ void set_array(float *a, float value, int len) {
    int idx = threadIdx.x + blockIdx.x * blockDim.x;
    if (idx < len) a[idx] = value;
}
```

```
__global__ void STREAM_Triad( float *a, float *b, float *c,
    float scalar, int len) {
    int idx = threadIdx.x + blockIdx.x * blockDim.x;
    if (idx < len) c[idx] = a[idx]+scalar*b[idx]; }
```

HPC Benchmarks Using Conventional Programming Approaches

STREAM TRIAD: C + MPI + OPENMP

```
#include <hpcc.h>
#ifdef _OPENMP
#include <omp.h>
#endif

static int VectorSize;
static double *a, *b, *c;

int HPCC_StarStream(HPCC_Params *params) {
    int myRank, commSize;
    int rv, errCount;
    MPI_Comm comm = MPI_COMM_WORLD;

    MPI_Comm_size( comm, &commSize );
    MPI_Comm_rank( comm, &myRank );

    rv = HPCC_Stream( params, 0 == myRank );
    MPI_Reduce( &rv, &errCount, 1, MPI_INT, MPI_SUM, 0, comm );

    return errCount;
}

int HPCC_Stream(HPCC_Params *params, int doIO) {
    register int j;
    double scalar;

    VectorSize = HPCC_LocalVectorSize( params, 3, sizeof(double), 0 );

    a = HPCC_XMALLOC( double, VectorSize );
    b = HPCC_XMALLOC( double, VectorSize );
    c = HPCC_XMALLOC( double, VectorSize );

    if (!a || !b || !c) {
        if (c) HPCC_free(c);
        if (b) HPCC_free(b);
        if (a) HPCC_free(a);
        if (doIO) {
            fprintf( outFile, "Failed to allocate memory (%d).\n", VectorSize );
            fclose( outFile );
        }
        return 1;
    }

#ifdef _OPENMP
#pragma omp parallel for
#endif
    for (j=0; j<VectorSize; j++) {
        b[j] = 2.0;
        c[j] = 1.0;
    }
    scalar = 3.0;

#ifdef _OPENMP
#pragma omp parallel for
#endif
    for (j=0; j<VectorSize; j++) {
        a[j] = b[j]+scalar*c[j];
    }

    HPCC_free(c);
    HPCC_free(b);
    HPCC_free(a);

    return 0;
}
```

HPCC RA: MPI KERNEL

```
/* Perform updates to main table. The scalar equivalent is:
 * for (i=0; i<NUPDATE; i++) {
 *   Ran = (Ran << 1) ^ ((s64int) Ran < 0) ? POLY: 0;
 *   Table[Ran & (TABLESIZE-1)] ^= Ran;
 * }
 */

MPI_Irecv(&LocalRecvBuffer, localBufferSize, tparams.dtype64,
          MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, &inreq);
while (! (< SendCnt)) {
    /* receive messages */
    do {
        MPI_Test(&inreq, &have_done, &status);
        if (have_done) {
            if (status.MPI_TAG == UPDATE_TAG) {
                MPI_Get_count(&status, tparams.dtype64, &recvUpdates);
                bufferBase = 0;
                for (j=0; j < recvUpdates; j++) {
                    inmsg = LocalRecvBuffer[bufferBase+j];
                    LocalOffset = (inmsg & (tparams.TableSize - 1)) -
                        tparams.GlobalStartMyProc;
                    HPCC_Table[LocalOffset] ^= inmsg;
                }
            } else if (status.MPI_TAG == FINISHED_TAG) {
                NumberReceiving--;
            } else {
                MPI_Abort( MPI_COMM_WORLD, -1 );
            }
            MPI_Irecv(&LocalRecvBuffer, localBufferSize, tparams.dtype64,
                    MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, &inreq);
        }
    } while (have_done && NumberReceiving > 0);

    if (pendingUpdates < maxPendingUpdates) {
        Ran = (Ran << 1) ^ ((s64int) Ran < ZERO64B ? POLY : ZERO64B);
        GlobalOffset = Ran & (tparams.TableSize-1);
        if (GlobalOffset < tparams.Top)
            WhichPe = (GlobalOffset / (tparams.MinLocalTableSize + 1));
        else
            WhichPe = ((GlobalOffset - tparams.Remainder) /
                        tparams.MinLocalTableSize);
        if (WhichPe == tparams.MyProc) {
            LocalOffset = (Ran & (tparams.TableSize - 1)) -
                tparams.GlobalStartMyProc;
            HPCC_Table[LocalOffset] ^= Ran;
        }
    } else {
        HPCC_InsertUpdate(Ran, WhichPe, Buckets);
        pendingUpdates++;
    }
    i++;
}
else {
    MPI_Test(&outreq, &have_done, MPI_STATUS_IGNORE);
    if (have_done) {
        outreq = MPI_REQUEST_NULL;
        pe = HPCC_GetUpdates(Buckets, LocalSendBuffer, localBufferSize,
                             &peUpdates);
        MPI_Isend(&LocalSendBuffer, peUpdates, tparams.dtype64, (int)pe,
                  UPDATE_TAG, MPI_COMM_WORLD, &outreq);
        pendingUpdates -= peUpdates;
    }
}

/* send our done messages */
for (proc_count = 0; proc_count < tparams.NumProcs; ++proc_count) {
    if (proc_count == tparams.MyProc) {
        MPI_Request_NULL, continue; }
    /* send garbage - who cares, no one will look at it */
    MPI_Isend(&Ran, 0, tparams.dtype64, proc_count, FINISHED_TAG,
             MPI_COMM_WORLD, tparams.finish_req + proc_count);
}

/* Finish everyone else up... */
while (NumberReceiving > 0) {
    MPI_Wait(&inreq, &status);
    if (status.MPI_TAG == UPDATE_TAG) {
        MPI_Get_count(&status, tparams.dtype64, &recvUpdates);
        bufferBase = 0;
        for (j=0; j < recvUpdates; j++) {
            inmsg = LocalRecvBuffer[bufferBase+j];
            LocalOffset = (inmsg & (tparams.TableSize - 1)) -
                tparams.GlobalStartMyProc;
            HPCC_Table[LocalOffset] ^= inmsg;
        }
    } else if (status.MPI_TAG == FINISHED_TAG) {
        /* we got a done message. Thanks for playing. */
        NumberReceiving--;
    } else {
        MPI_Abort( MPI_COMM_WORLD, -1 );
    }
    MPI_Irecv(&LocalRecvBuffer, localBufferSize, tparams.dtype64,
             MPI_ANY_SOURCE, MPI_ANY_TAG, MPI_COMM_WORLD, &inreq);
}

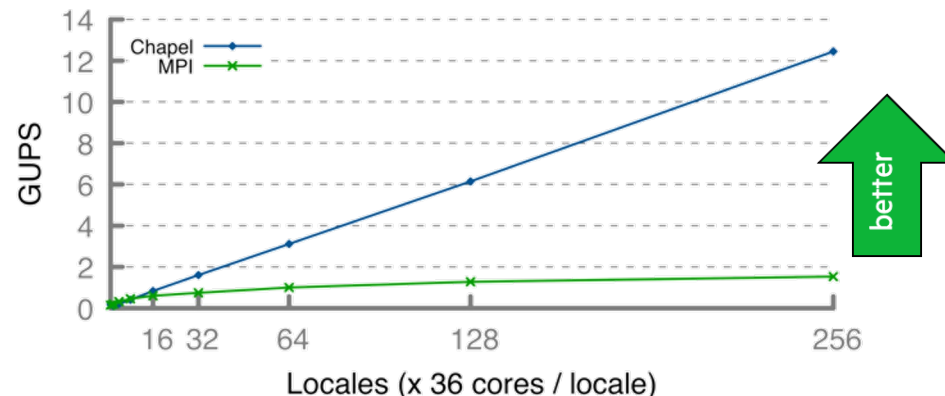
MPI_Waitall( tparams.NumProcs, tparams.finish_req, tparams.finish_statuses);
```

11/11/2019

[illegible]

```
A = B + alpha * C;
```

72



Bale IG in Chapel vs. SHMEM on HPE Cray EX (Slingshot-11)

Chapel (Simple / Auto-Aggregated version)

```
forall (d, i) in zip(Dst, Inds) do
  d = Src[i];
```

SHMEM (Exstack version)

```
i=0;
while( exstack_proceed(ex, (i==l_num_req)) ) {
  i0 = i;
  while(i < l_num_req) {
    l_idx = pckindx[i] >> 16;
    pe = pckindx[i] & 0xffff;
    if(!exstack_push(ex, &l_idx, pe))
      break;
    i++;
  }

  exstack_exchange(ex);

  while(exstack_pop(ex, &idx, &fromth)) {
    idx = ltable[idx];
    exstack_push(ex, &idx, fromth);
  }
  lgp_barrier();
  exstack_exchange(ex);

  for(j=i0; j<i; j++) {
    fromth = pckindx[j] & 0xffff;
    exstack_pop_thread(ex, &idx, (uint64_t)fromth);
    tgt[j] = idx;
  }
  lgp_barrier();
}
```

SHMEM (Conveyors version)

```
i = 0;
while (more = convey_advance(requests, (i == l_num_req)),
       more | convey_advance(replies, !more)) {

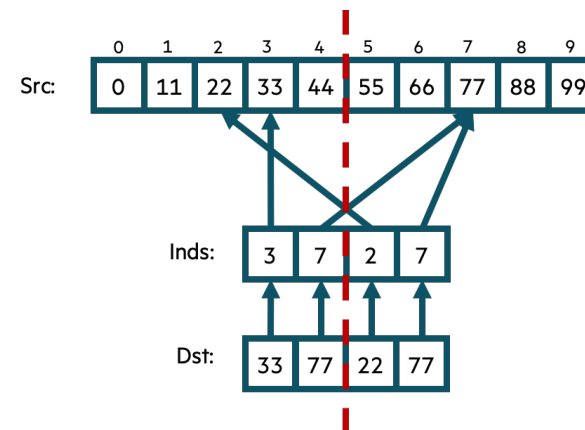
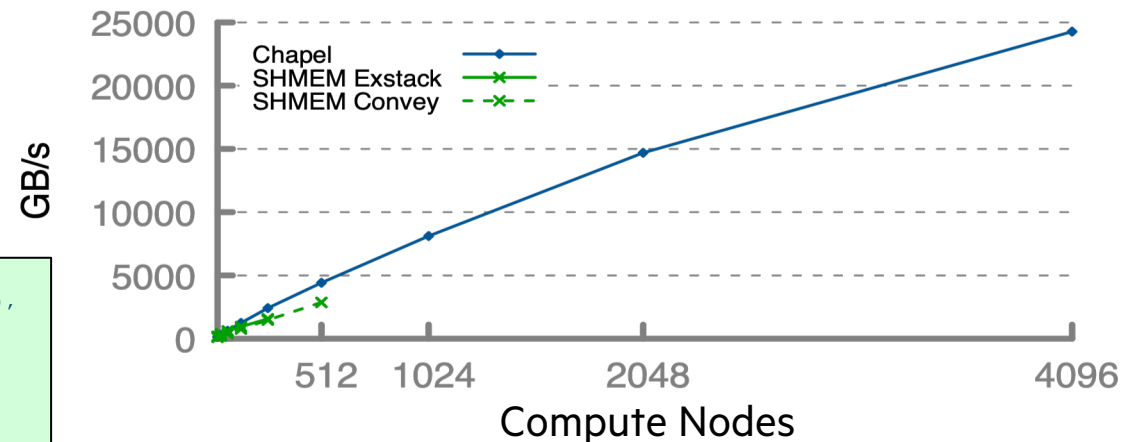
  for (; i < l_num_req; i++) {
    pkg.idx = i;
    pkg.val = pckindx[i] >> 16;
    pe = pckindx[i] & 0xffff;
    if (!convey_push(requests, &pkg, pe))
      break;
  }

  while (convey_pull(requests, ptr, &from) == convey_OK) {
    pkg.idx = ptr->idx;
    pkg.val = ltable[ptr->val];
    if (!convey_push(replies, &pkg, from)) {
      convey_unpull(requests);
      break;
    }
  }

  while (convey_pull(replies, ptr, NULL) == convey_OK)
    tgt[ptr->idx] = ptr->val;
}
```

Bale Indexgather Performance

HPE Cray EX (Slingshot-11)

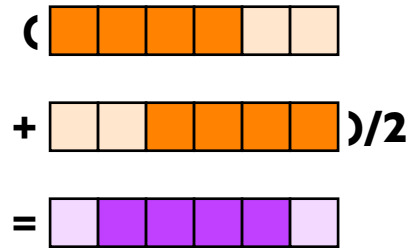


Q: What accounts for the code size disparities between Chapel and SHMEM / MPI?

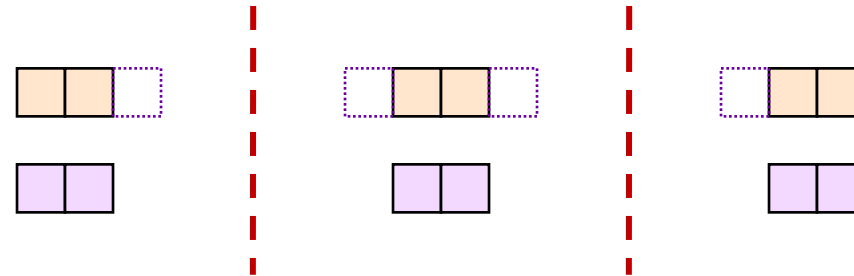
A: Chapel Supports Global-view Programming

Example: “Apply a 3-point stencil to a vector”

Global-View



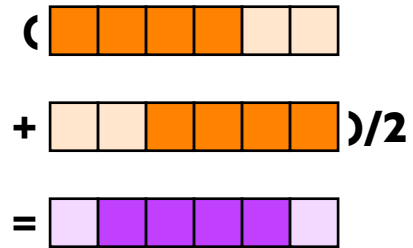
SPMD



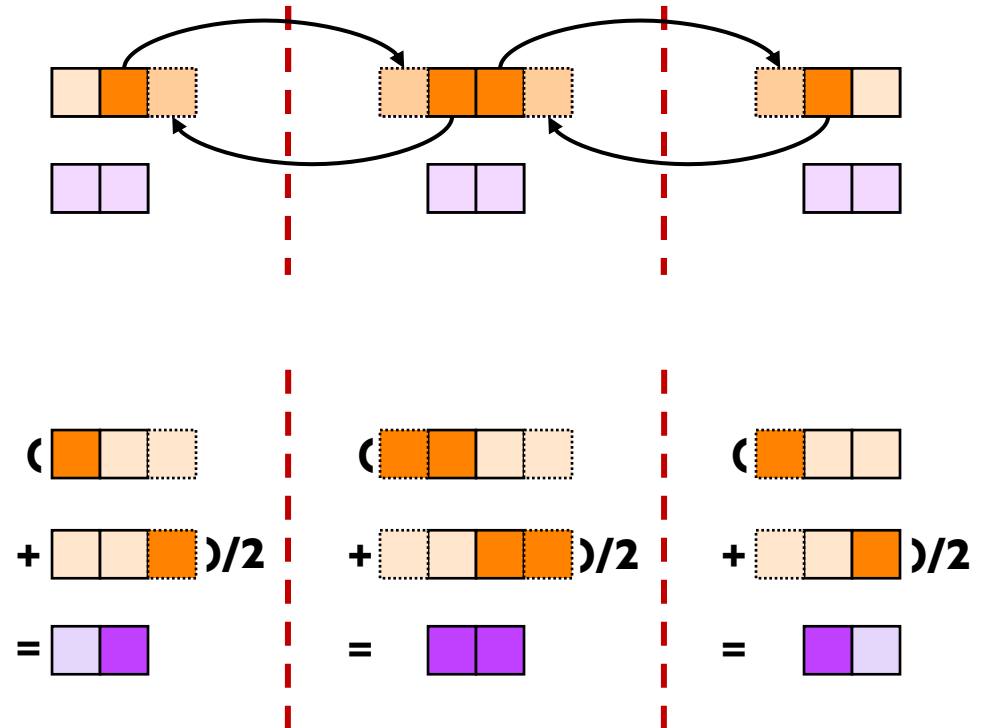
A: Chapel Supports Global-view Programming

Example: “Apply a 3-point stencil to a vector”

Global-View




SPMD



A: Chapel Supports Global-view Programming

Example: “Apply a 3-point stencil to a vector”

Global-View Chapel code

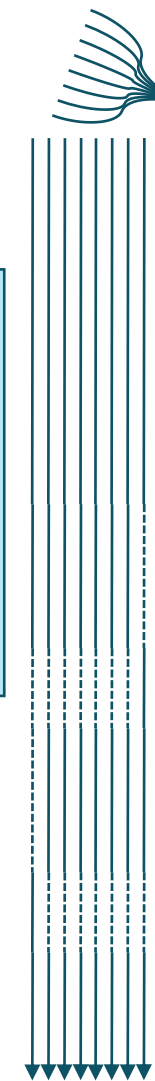


```
use BlockDist;

proc main() {
  const n = 1000,
        D = blockDist.createDomain(1..n);

  forall i in D[2..n-1] do
    B[i] = (A[i-1] + A[i+1])/2;
  }
}
```

SPMD pseudocode (MPI-esque)



```
proc main() {
  var n = 1000;
  var p = numProcs(),
      me = myProc(),
      myN = n/p,
      myLo = 1,
      myHi = myN;
  var A, B: [0..myN+1] real;

  if (me < p-1) {
    send(me+1, A[myN]);
    recv(me+1, A[myN+1]);
  } else
    myHi = myN-1;
  if (me > 0) {
    send(me-1, A[1]);
    recv(me-1, A[0]);
  } else
    myLo = 2;
  forall i in myLo..myHi do
    B[i] = (A[i-1] + A[i+1])/2;
  }
}
```

100%

Page 10

```
A = B + alpha * C;
```

11

72

better

better

11

SPMD Programming in Chapel

That said, as a general-purpose language, Chapel supports writing SPMD patterns as well:

```
coforall loc in Locales do
  on loc do
    myMain();

proc myMain() {
  // ... write your SPMD computation here ...
}
```



Q&A Time

I'm happy to take questions about...

- ...Chapel's design

- ...Chapel's implementation

- ...Chapel history or culture

- ...anything else Chapel-related

- ...PL topics other than Chapel

- ...whatever else is on your mind

(subject to Prof. Strout's approval, of course)

