# MetaMorph
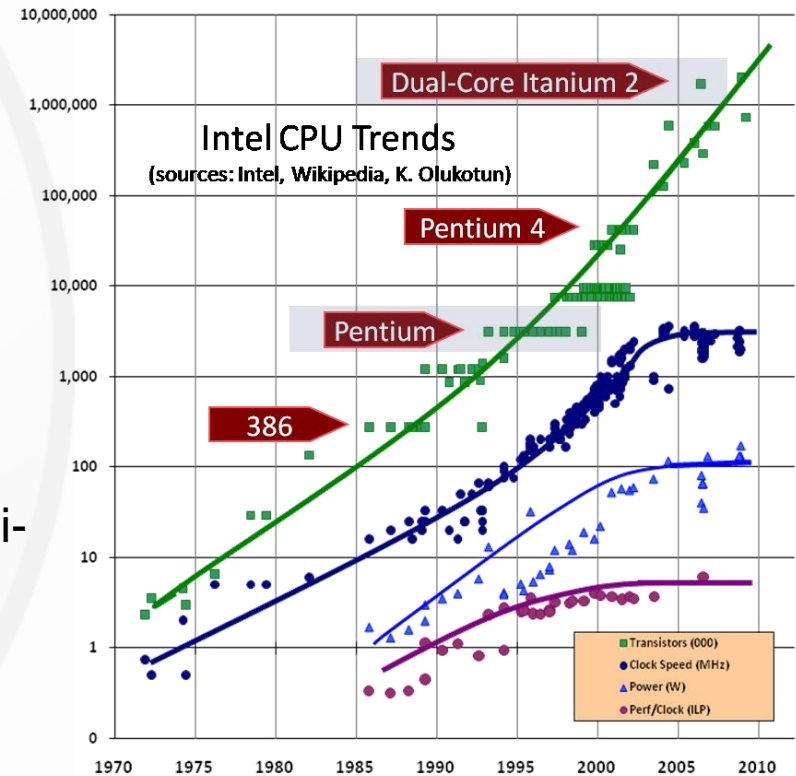A Modular Library of Malleable Accelerator Primitives for Heterogeneous Parallel Computing

SEEC Fall Symposium 2015

Paul Sathre, **Ahmed Helal**

Virginia Tech
Invent the Future

SyNeRG
synergy.cs.vt.edu

# Motivation:

- ## Why do we need accelerator?
  - Free lunch is over!
  - You can't anymore wait until H/W vendors improve your performance with faster CPUs.

- ## However, parallel programming is complex!
  - Architectures: Multicore CPUs, multi-socket cc-NUMA, multi-nodes, GPUs, MIC,…
  - Programming models: OpenMP, OpenACC, MPI, CUDA, OpenCL,…
  - *Five years from now, you will NOT know what hardware will look like*



Intel CPU Trends
(sources: Intel, Wikipedia, K. Olukotun)

Dual-Core Itanium 2
Pentium 4
Pentium
386

Transistors (000)
Clock Speed (MHz)
Power (W)
Perf/Clock (ILP)

VirginiaTech
1872
*Invent the Future*

SyNeRG

# Motivation:

## Productivity = Programmability + Performance + Portability

- Dot Product (Serial on CPU)

<u>HOST</u>
Serial C

<u>DEVICE</u>

```
1.  void dotProd(double *a,
2.               double *b,
3.               double *ret,
4.               int n_elem)
5.  {
6.    int i;
7.    *ret = 0.0;
8.    for (i = 0; i < n_elem; i++)
9.      ret += a[i] * b[i];
10. }
```

**10**

# Motivation:

## Productivity = Programmability + Performance + Portability

- Dot Product (Parallel on GPU using CUDA)

### HOST

**27**

```
1.  int main(int argc, char ** argv)
2.  {
3.    int n_elem;
4.    double *a, *b, *ret, sum = 0.0;
5.    //allocate and initialize a, b, and ret
6.    ...
7.    double *dev_a, *dev_b, *dev_r;
8.    size_t size = n_elem * sizeof(double);
9.    //allocate device buffers
10.   cudaMalloc((void**)&dev_a, size);
11.   cudaMalloc((void**)&dev_b, size);
12.   cudaMalloc((void**)&dev_r, size);

13.   //initialize device buffers
14.   cudaMemcpy(dev_a, a, size, cudaMemcpyHostToDevice);
15.   cudaMemcpy(dev_b, b, size, cudaMemcpyHostToDevice);

16.   //set grid/block size
17.   ...
18.   //multiply elements
19.   dotProd-vmul<<<grid, block>>>(dev_a, dev_b, dev_r,
20.     n_elem);

21.   //partial result
22.   cudaMemcpy(ret, dev_r, size, cudaMemcpyDeviceToHost);

23.   //CPU sum
24.   int i;
25.   for (i = 0; i < n_elem; i++)
26.     sum += ret[i];
27. }
```

### DEVICE (GPU)

```
1.  __global void dotProd-vmul(
2.      double *a, double *b,
3.      double *ret, int n_elem)
4.  {
5.    int i;
6.    int tid = blockIdx.x *
7.      blockDim.x + threadIdx.x;
8.    int n_threads = blockDim.x *
9.      gridDim.x;
10.   for (i=tid; i<n_elem; i+=n_threads)
11.     ret[tid] = a[tid] * b[tid];
12. }
```

# Motivation:
## Productivity = Programmability + Performance + Portability
### Dot Product (Optimized Parallel on GPU Device)

HOST **39**

```
1.  int main(int argc, char** argv)
2.  {
3.     //global and block sizes
4.     int ni, nj, nk, tx, ty, tz;
5.     //pick a device and zero copy
6.     cudaSetDevice(0)
7.     //declare host memory
8.     double *a, *b, *ret, zero = 0.0;
9.     //allocate and initialize
10.    …

11.    //Allocate device buffers
12.    size_t size = sizeof(double)*ni*nj*nk;
13.    double *dev_a, *dev_b, *dev_r;
14.    cudaMalloc(&dev_a, size);
15.    cudaMalloc(&dev_b, size);
16.    cudaMalloc(&dev_r, sizeof(double));
17.    //initialize them
18.    cudaMemcpy(dev_a, a, size,
19.      cudaMemcpyHostToDevice);
20.    cudaMemcpy(dev_b, a, size,
21.      cudaMemcpyHostToDevice);
22.    cudaMemcpy(dev_r, &zero, sizeof(double),
23.      cudaMemcpyHostToDevice);

24.    //set computation shape
25.    dim3 grid, block, shape, start, end;
26.    block = dim3(tx, ty, tz);
27.    grid = dim3((ni+tx-1)/ni, (nj+ty-1)/ty, 1);
28.    shape = dim3(ni, nj, nk);
29.    start = dim3(0, 0, 0);
30.    end = dim3(ni-1, nj-1, nk-1);

31.    //run the kernel
32.    dotProd<<<grid, block, tx*ty*tz*sizeof(double)>>>
33.      (dev_a, dev_b, shape.x, shape.y, shape.z,
34.      start.x, start.y, start.z, end.x, end.y, end.z,
35.      (nk+tz-1)/tz, dev_r, tx*ty*tz);

36.    //bring the dot product back
37.    cudaMemcpy (ret, dev_r, sizeof(double),
38.      cudaMemcpyDeviceToHost);
39. }
```

DEVICE (GPU) **40**

```
1.  __device__ void block_reduction(double *psum,
2.    int tid, int len_)
3.  {
4.    unsigned int stride = len_ >> 1;
5.    while (stride > 0) {
6.      if (tid  < stride)
7.        psum[tid] += psum[tid+stride];
8.      __syncthreads();
9.      stride >>= 1;
10.   }
11. }
12. //Implementation of double atomicAdd
13. ...

14. __global__ void kernel_dotProd(double *phi1, double *phi2,
    int i, int j, int k, int sx,
15.   int sy, int sz, int ex, int ey, int ez,
16.   int gz, T * reduction, int len_)
17. {
18.   extern shared psum[];
19.   int tid, x, y, z, itr;
20.   bool boundx, boundy, boundz;
21.   tid = threadIdx.x + (threadIdx.y) * blockDim.x
22.     + (threadIdx.z) * (blockDim.x * blockDim.y);
23.   x = (blockIdx.x)*blockDim.x+threadIdx.x+sx;
24.   y = (blockIdx.y)*blockDim.y+threadIdx.y+sy;

25.   psum[tid] = 0;
26.   boundy = ((y >= sy) && (y <= ey));
27.   boundx = ((x >= sx) && (x <= ex));

28.   for (itr = 0; itr < gz; itr++) {
29.     z = itr*blockDim.z+threadIdx.z +sz;
30.     boundz = ((z >= sz) && (z <= ez));
31.     if (boundx && boundy && boundz)
32.       psum[tid] += phi1[x+y*i+z*i*j] *
33.         phi2[x+y*i+z*i*j];
34.   }
35.   __syncthreads();
36.   block_reduction(psum,tid,len_);
37.   __syncthreads();

38.   if(tid == 0)
39.     atomicAdd(reduction,psum[0]);
40. }
```

**VirginiaTech**
*Invent the Future*

SyNeRG
synergy.cs.vt.edu

# Motivation:

- ## What is MetaMorph?
  - A runtime library that hides the complexity of accelerator programming and interoperability behind a standard interface (static API)
    - The static API acts as an interoperability layer between user applications and community-tuned accelerator backends
    - Performance portability still reliant on humans tuning interchangeable backends to specific devices!

- ## Why do we need it?
  - Extracting the full performance of heterogeneous systems is non-trivial and requires architecture expertise
  - Architectures change faster than codes and scientists shouldn't have to waste their time (re)learning and rewriting code for them

# Motivation:
## Productivity = Programmability + Performance + Portability

- Dot Product (via our MetaMoph Library)

**36**

```
1.   int main(int argc, char **argv)
2.   {
3.     //global and block sizes
4.     int ni, nj, nk, tx, ty, tz;
5.     //pick a mode and zeroth device
6.     choose_accel(0, metaModePreferCUDA);
7.     //declare host memory
8.     double *a, *b, *ret, zero = 0.0;
9.     //allocate and initialize it
10.    ...

11.    //Allocate device buffers
12.    size_t size = sizeof(double)*ni*nj*nk;
13.    double *dev_a, *dev_b, *dev_r;
14.    meta_alloc(&dev_a, size);
15.    meta_alloc(&dev_b, size);
16.    meta_alloc(&dev_r, sizeof(double));

17.    //initialize them
18.    meta_copy_h2d(dev_a, a, size, true);
19.    meta_copy_h2d(dev_b, a, size, true);
20.    meta_copy_h2d(dev_r, &zero, sizeof(double),
21.      true);

22.    //set computation shape
23.    a_dim3 grid, block, shape, start, end;
24.    block[0] = tx, block[1] = ty, block[2] = tz;
25.    grid[0] = (ni+tx-1)/ni, grid[1] = (nj+ty-1)/ty,
26.      grid[2] = (nk+tz-1)/tz;
27.    shape[0] = ni, shape[1] = nj, shape[2] = nk;
28.    start[0] = 0, start[1] = 0, start[2] = 0;
29.    end[0] = ni-1, end[1] = nj-1, end[2] = nk-1;

30.    //run the kernel
31.    meta_dotProd(&grid, &block, dev_a, dev_b,
32.      &shape, &start, &end, dev_r, a_db, true);

33.    //bring the dot product back
34.    meta_copy_d2h(ret, dev_r, sizeof(double),
35.      false);

36. }
```

also support:
metaModePreferOpenCL
/OpenMP
for AMD/MIC/CPU

all kernels/copies can be
asynchronous with
flag = true, else blocking

grid & block specify
thread organization
a la CUDA/OpenCL

shape, start, and end allow
dot product on arbitrary
subregions of 3D space

# Goal :

- Support upgrade of legacy CFD codes and development of new codes in heterogeneous environments
  - Many methods:
    - Finite Element, Finite Difference, Finite Volume, etc.
  - Vast range of solvers, preconditioners and other math kernels
  - Range of storage/computation schemes:
    - structured vs. unstructured grids,
    - Sparse vs. Dense storage

- Expand with kernels from other domains
    - Computational Bio, Bigdata, Cosmology, etc.

# Desired Features:

- Usable by non-architecture experts (via drop-in replacement functions) and *without* heavyweight framework-centric development


- Current approaches trade portability with heavyweight frameworks
    - provide portability via dynamic remapping of the framework's data structures to supported devices
    - require redesigning the entire application to use the framework's complex data structures, which may not easily support operations needed by the application and adds a considerable programming/runtime overhead.

# Related Work:

- **OpenFOAM**: a very popular CPU-only CFD simulation framework
- **Paralution**: a new CPU/GPU/MIC BLAS framework
- **MAGMA**: a suite of CUDA, AMD OpenCL, and MIC BLAS libraries
- **Trilinos**: a massive set of multi-physics libraries and frameworks

| | MetaMorph | OpenFOAM | Paralution | MAGMA | Trilinos |
|---|---|---|---|---|---|
| **CFD support** | ✓ | ✓ | ✓ | ✗ | ✓ |
| **Acceleration support** | OMP, CUDA, OCL | No internal support | OMP, CUDA, OCL | OMP, CUDA, OCL | preliminary CUDA/MIC |
| **Target devices** | CPU, GPU, MIC | CPU | CPU, GPU, MIC | CPU, GPU, MIC | CPU, GPU, MIC |
| **MPI & interoperability** | ✓ | ✗ | ✓ | ✗ | ✗ |
| **framework-centric development** | ✗ | ✓ | ✓ | ✓ | ✓ |
| **External dependencies** | ✗ | ✓ | ✓ | ✓ | ✓ |

# Contributions:

- Transparent MPI exchange of accelerator device buffers, including between CUDA and OpenCL

- MetaMorph performance on-par with best language-specific BLAS libraries

- Intelligent Data Marshalling (packing/unpacking) of arbitrary subregions (2D ghost regions) of a 3D structured grid.

- OpenCL context management stack, Fortran and C APIs, event-based timing infrastructure, …

- *I worked on OpenMP backedend for CPU/MIC and OpenCL backedend for MIC*

VirginiaTech
1872
*Invent the Future*

SyNeRG
synergy.cs.vt.edu

# Design

"Make-your-own" library from modular building blocks

Include only needed plugins and backends

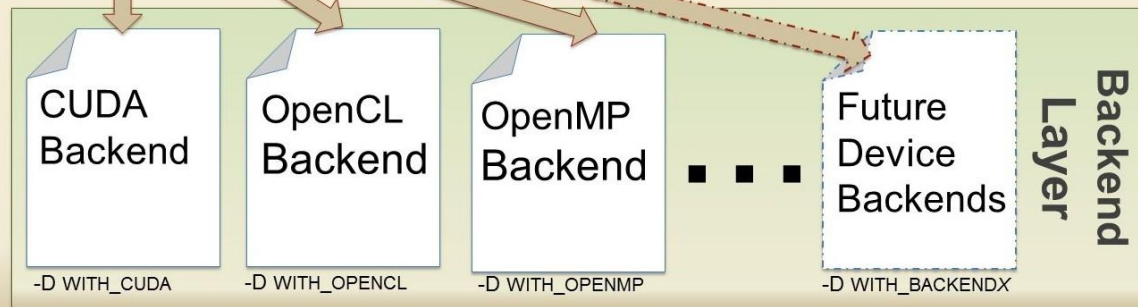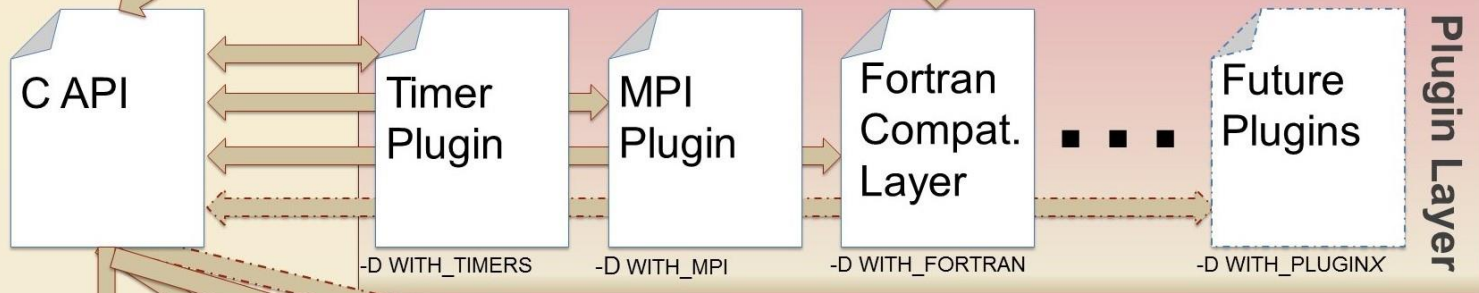**User Apps**

Purpose-built to provide single API to **accelerated** kernels for current and **future** devices

Runtime control over backends, plugins, and accelerator device selection

**Fortran 2003 API**

## Plugin Layer

| C API | Timer Plugin | MPI Plugin | Fortran Compat. Layer | ▪ ▪ ▪ ▪ | Future Plugins |
|---|---|---|---|---|---|
| | -D WITH_TIMERS | -D WITH_MPI | -D WITH_FORTRAN | | -D WITH_PLUGIN*X* |

## Backend Layer

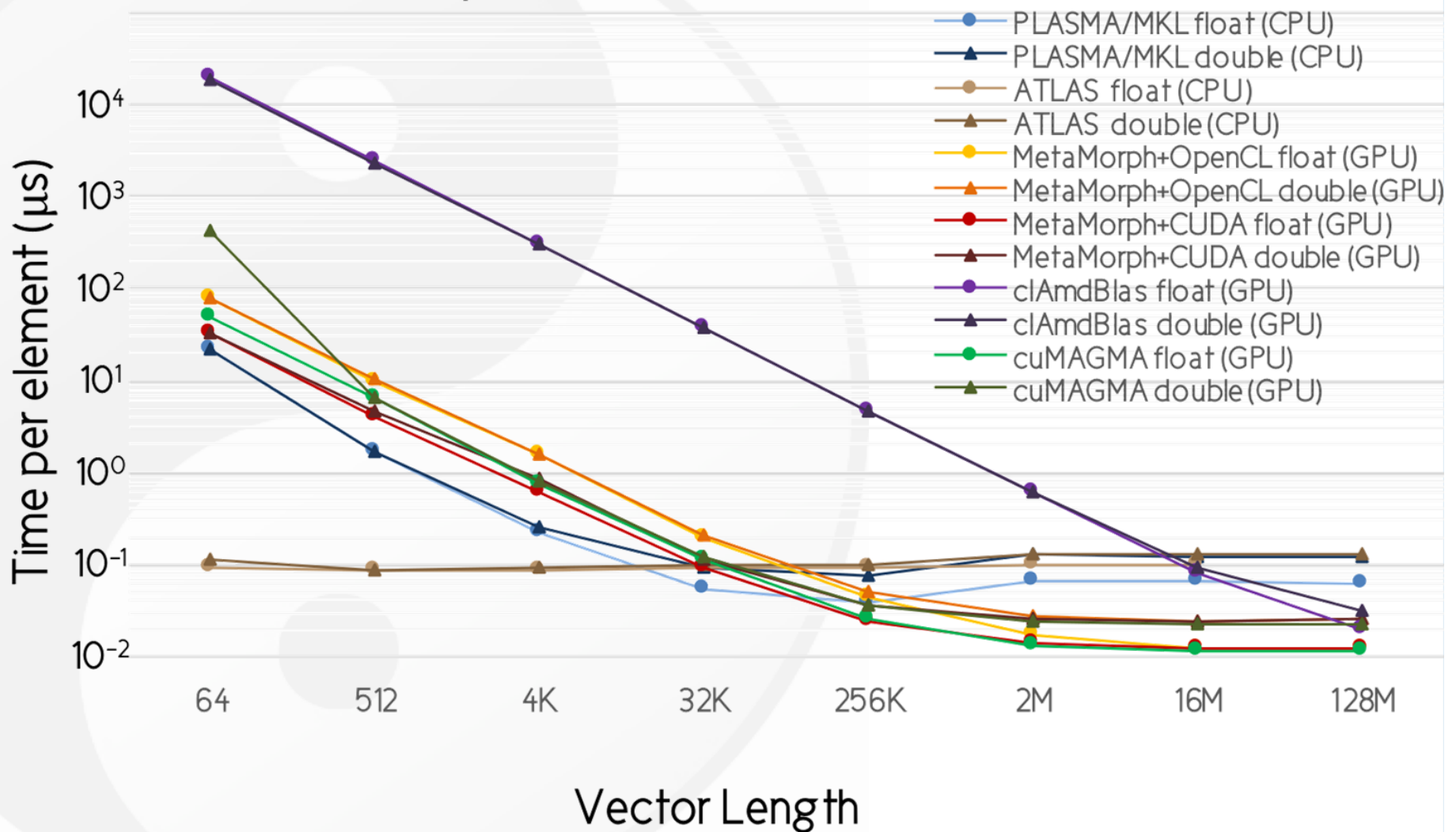| CUDA Backend | OpenCL Backend | OpenMP Backend | ▪ ▪ ▪ | Future Device Backends |
|---|---|---|---|---|
| -D WITH_CUDA | -D WITH_OPENCL | -D WITH_OPENMP | | -D WITH_BACKEND*X* |

# Implementation:

- "Library of Libraries"
  - libmetamorph.so
    - Implements the interface between the top-level MetaMorph API, and the platform-specific backend(s).
  - libmetamorph_x_core.so (x is either OpenMP, CUDA or OpenCL, ..)
    - Set of backend shared library objects ("cores") that provide the actual implementations of the platform-specific kernels.
    - These backends can be custom-tuned for a specific device.

- Compile time:
  - choose which backends to include, choose which (if any) plugins you need from {MPI, Timing, Fortran}
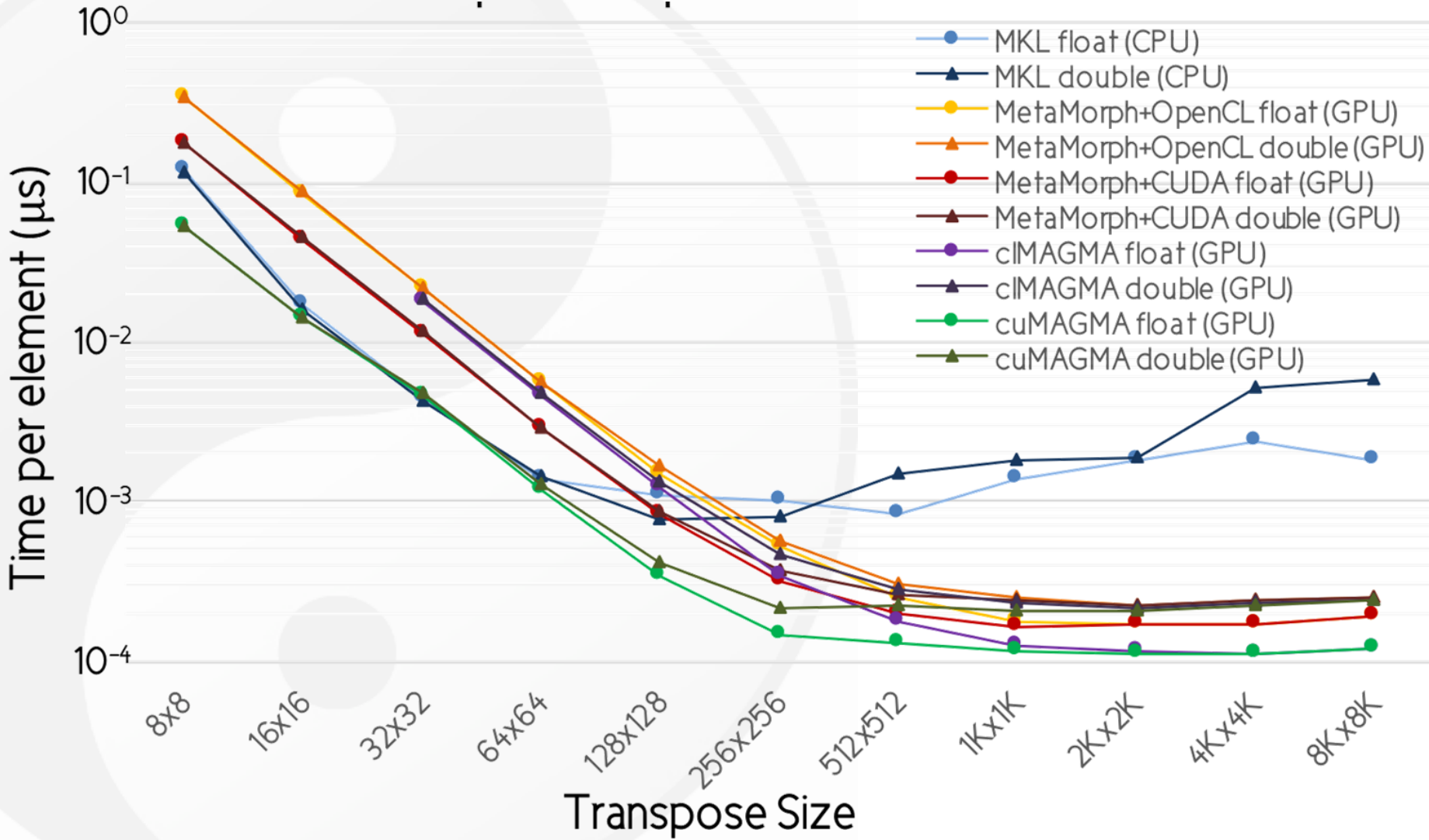
# Implementation:

- ## Runtime:
  - select with environment variables which backend/device are used for a run, timer verbosity, device selection, etc.

- ## OpenMP backend
  - Default backend
  - To minimize the design modifications, OpenMP backend mimic the accelerator model.
    - There are several unnecessary operations/data movements on CPU.
  - UNIFIED_MEMORY feature
    - User apps can eliminate redundant data transfers.
    - Copy pointers to the data, instead of explicit memory transfer.

# MetaMorph: Dot-product performance



Legend:
- PLASMA/MKL float (CPU)
- PLASMA/MKL double (CPU)
- ATLAS float (CPU)
- ATLAS double (CPU)
- MetaMorph+OpenCL float (GPU)
- MetaMorph+OpenCL double (GPU)
- MetaMorph+CUDA float (GPU)
- MetaMorph+CUDA double (GPU)
- clAmdBlas float (GPU)
- clAmdBlas double (GPU)
- cuMAGMA float (GPU)
- cuMAGMA double (GPU)

Y-axis: Time per element ($\mu$s) — $10^4$, $10^3$, $10^2$, $10^1$, $10^0$, $10^{-1}$, $10^{-2}$

X-axis: Vector Length — 64, 512, 4K, 32K, 256K, 2M, 16M, 128M

# MetaMorph: Transpose performance

# MetaMorph: Library overhead vs. standalone kernels

- OpenMP backend (3D-Dot product with double)
  – Max overhead 5%.



Execution Time (us)

CPU: Intel(R) Core(TM) i5-2400, 4 cores
@ 3.10GHz

synergy.cs.vt.edu
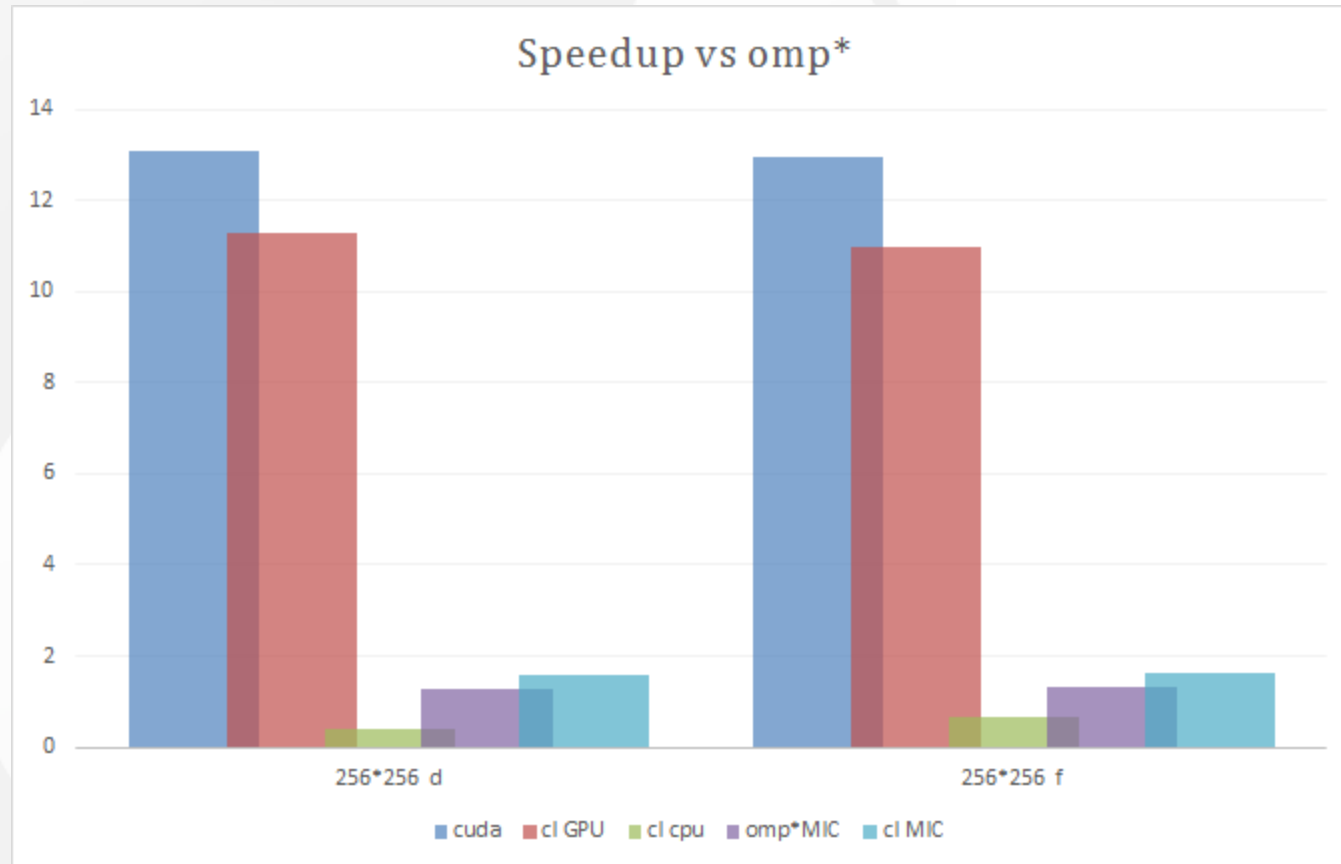
# MetaMorph: Performance vs. CPU OpenMP+AVX

- Dot product Kernel Performance
  - *Baseline is 4-thread OpenMP w/ AVX



CPU: Intel(R) Core(TM) i5-2400, 4 cores @ 3.10GHz
MIC: Intel® Xeon Phi 7100 Series, 61 core @1.238 GHz.
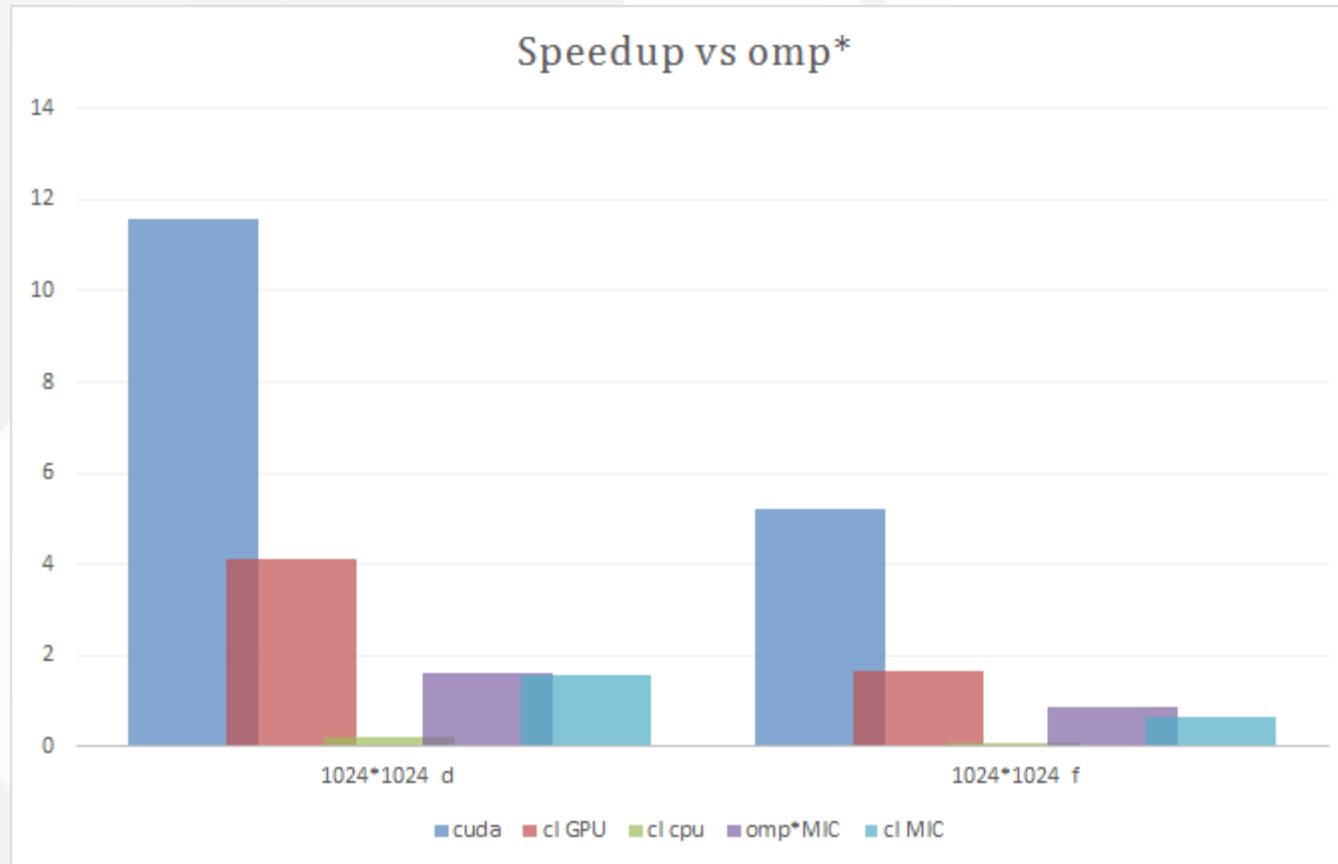GPU: Tesla K20c

# MetaMorph: Performance vs. CPU OpenMP+AVX

- Transpose Kernel Performance
  - *Baseline is 4-thread OpenMP w/ AVX



CPU: Intel(R) Core(TM) i5-2400, 4 cores @ 3.10GHz

MIC: Intel® Xeon Phi 7100 Series, 61 core @1.238 GHz.

GPU: Tesla K20c

# MetaMorph: Performance vs. CPU OpenMP+AVX

- Data Marshaling Kernel Performance
  - *Baseline is 4-thread OpenMP w/ AVX



Speedup vs omp*

CPU: Intel(R) Core(TM) i5-2400, 4 cores @ 3.10GHz

MIC: Intel® Xeon Phi 7100 Series, 61 core @1.238 GHz.

GPU: Tesla K20c

# MIC performance analysis/Optimization:

- Use single core Performance Degradation Factor (PDF) to analyze the performance on MIC and CPU

$$PDF_{MIC} = PDF_{FREQ} * PDF_{CMP} * PDF_{VEC} * PDF_{MEM}$$

- $PDF_{MIC}$ is the performance degradation factor on MIC in comparison with CPU (single core execution).

- $PDF_{FREQ}$ is the performance degradation factor due to frequency.
- $PDF_{CMP}$ is the performance degradation factor due to the core complexity.
- $PDF_{VEC}$ is the performance degradation factor due to the vectorization unit.
- $PDF_{MEM}$ is the performance degradation factor due to the memory system (working set).

# MIC performance analysis/Optimization:

- ## Single core performance modeling
  - Dot product kernel
    - $PDF_{MIC} \sim= 3.5$
      - $PDF_{VEC} = 0.5$  (2x speed up due to the wider vector unit)
      - $PDF_{MEM} \sim= 1$ (There is no temporal locality, and the spatial locality depends on the cache line size)
  - Transpose kernel
    - $PDF_{MIC} \sim= 30$
      - $PDF_{VEC} = 1$  (Memory transfer dominant)
      - $PDF_{MEM} \sim= 4$ (CPUs outperform Intel MIC in workloads sensitive to locality/cache size)
  - Data marshaling
    - $PDF_{MIC} \sim= 12$
      - $PDF_{VEC} = 1$  (Memory transfer dominant)
      - $PDF_{MEM} \sim= 1.5$ (CPUs outperform Intel MIC in workloads sensitive to locality/cache size)

CPU: Intel(R) Core(TM) i5-2400, 4 cores @ 3.10GHz
MIC: Intel® Xeon Phi 7100 Series, 61 core @1.238 GHz.

VirginiaTech
1872
Invent the Future
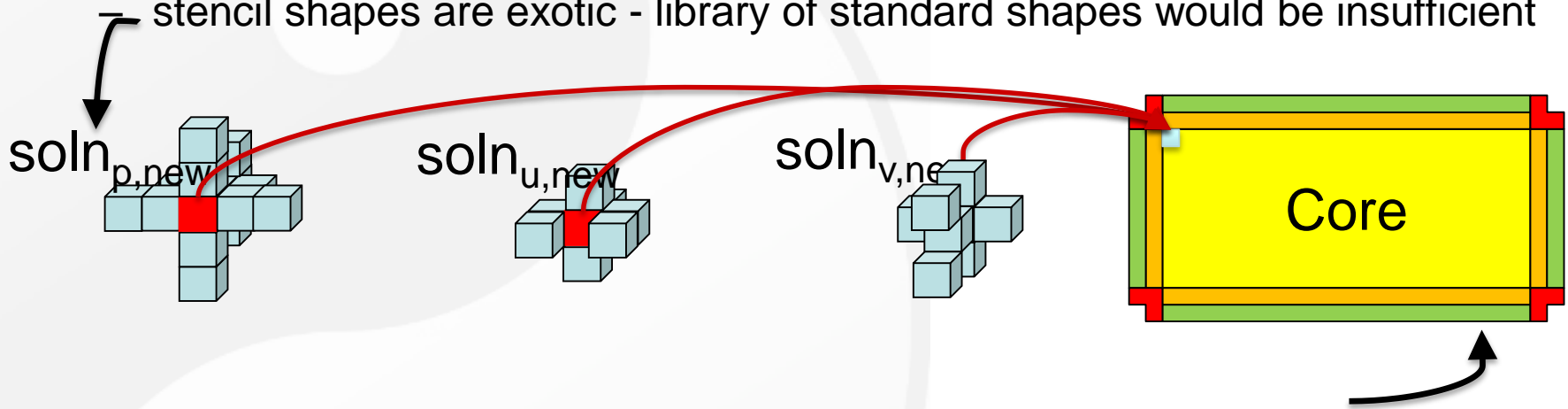
SyNeRG
synergy.cs.vt.edu

# MIC performance analysis/Optimization:

- So our hypothesis: CPUs outperform Intel MIC in workloads sensitive to locality/cache size is correct.

- How we can use this to improve the performance?
  - CPU backend optimize for load balance vs. locality.
    - Dynamic scheduling with small data chunks.

  - MIC backend should optimize for locality over load balance
    - It is more sensitive to locality due to the lack of shared level3 cache to capture large working set.
    - Static scheduling with large (sequential) data chunks.
    - Collapsing nested loops to have wider workload distribution space (larger number of threads).

# Work in progress

- Preliminary work towards stencil support
  - stencil shapes are exotic - library of standard shapes would be insufficient

$soln_{p,new}$     $soln_{u,new}$     $soln_{v,ne}$     Core

- (LDC) 9 conditional regions, but only 4 types of computation behavior

  - Plan: allow arbitrary stencils to be expressed as semi-structured C
    - Map ingested C to each platform via runtime modification of the kernel string
    - Utilize just-in-time (JIT) compilation for each platform

# Future Work

- MPI plugin:
  - Automatic pipelined transfers (a la MPI-ACC)

- Adaptive runtime "meta-backend"
  - CoreTSAR /AfinityTSAR – support automatic scheduling across a complete node

- Domain decomposition plugin
  - Graph partitioning and load-balancing algorithms

- More kernels!
  - Start with those needed from the CFD space:
    - stencil computations, linear algebra, solvers, preconditioners
  - Expand with kernels from other domains
    - Computational Bio, Bigdata, Cosmology, etc.
  - Optimized for more devices
    - autotuning params, automatic device detection, JIT reconfig of OpenCL

# Conclusion/Questions

- MetaMorph is designed to be a "build your own" library of computational primitives, accelerated from the start, instead of retroactively

- We want to enable domain scientists to write fast codes faster, and be able to maintain it themselves

- Plenty of work left to be done
  - Need more developers to help implement kernels and plugins, starting with CFD communication and branching from there

# What else do I work on?

- ## Parallel circuit simulation:
  - **Helal, Ahmed E**., Amr M. Bayoumi, and Yasser Y. Hanafy. "Parallel circuit simulation using the direct method on a heterogeneous cloud." Design Automation Conference (DAC), 2015 52nd ACM/EDAC/IEEE. IEEE, 2015.
  - Hassan, M. W. , **Helal, Ahmed E.**, Hanafy, Y. Y., "High Performance Sparse LU Solver FPGA Accelerator using a Static Synchronous Data Flow Model", The 23rd IEEE International Symposium on Field-Programmable Custom Computing Machines, FCCM'15, May 2015.

- ## AutoMatch:
  - A tool for automatic kernel-architecture matching
  - Enable scientists to predict the best hardware device for their applications without writing code for every target device.

# References

- H. Jasak, et. al. , "OpenFOAM: A C++ library for complex physics simulations." International Workshop on Coupled Methods in Numerical Dynamics, vol. 1000, pp. 1-20, 2007 http://www.openfoam.com/

- D. Lukarski, "PARALUTION project v0.7.0," 2012. http://www.paralution.com/

- J. Dongarra, et. al. , "Accelerating Numerical Dense Linear Algebra Calculations with GPUs," in Numerical Computations with GPUs. Springer, 2014, pp. 3–28.

- M. A. Heroux, et. al."An overview of the Trilinos project," ACM Trans. Math. Softw., vol. 31, no. 3, pp. 397–423, 2005.

- M. J. Harvey, et. al., "Swan: A tool for porting CUDA programs to OpenCL," Computer Physics Communications, vol. 182, no. 4, pp. 1093–1099, 2011.