

Incremental Migration of C and Fortran Applications to GPGPU using HMPP

Peppher 2011



- Many applications can benefit from GPU computing
 - Linear Algebra, signal processing
 - Bio informatics, molecular dynamics
 - Magnetic resonance imaging, tomography
 - Reverse time migration, electrostatic
 - ...
- Porting legacy codes to GPU computing is a major challenge
 - Can be very expensive
 - Require to minimize porting risks
 - Should be based on future-proof approach
 - Implies application and performance programmers to cooperate
- A good methodology is paramount to reduce porting cost
 - HMPP provides an efficient solution

What is HMPP? (Hybrid Manycore Parallel Programming)



- A directive based multi-language programming environment
 - Help keeping software independent from hardware targets
 - Provide an incremental tool to exploit GPU in legacy applications
 - Avoid exit cost, can be future-proof solution
- HMPP provides
 - Code generators from C and Fortran to GPU (CUDA or OpenCL)
 - A compiler driver that handles all low level details of GPU compilers
 - A runtime to allocate and manage GPU resources
- Source to source compiler
 - CPU code does not require compiler change
 - Complement existing parallel APIs (OpenMP or MPI)

HMPP Main Design Considerations

- Focus on the main bottleneck
 - Communication between GPUs and CPUs
- **Allow incremental development**
 - Up to full access to the hardware features
- Work with other parallel APIs (e.g. OpenMP, MPI)
 - Orchestrate CPU and GPU computations
- Consider multiple languages
 - Avoid asking users to learn a new language
- Consider resource management
 - Generate robust software
- Exploit vendor tools/compilers
 - Do not replace, complement

How Does HMPP Differ from CUDA or OpenCL?

- HMPP parallel programming model is **parallel loop centric**
- CUDA and OpenCL parallel programming models are **thread centric**

```
void saxpy(int n, float alpha,
           float *x, float *y){
#pragma hmppcg parallel
for(int i = 0; i<n; ++i)
    y[i] = alpha*x[i] + y[i];
}
```

```
__global__
void saxpy_cuda(int n, float
alpha,
float *x, float *y) {
int i = blockIdx.x*blockDim.x +
threadIdx.x;
if(i<n) y[i] = alpha*x[i]+y[i];
}
```

```
int nblocks = (n + 255) / 256;
saxpy_cuda<<<nblocks,
256>>>(n, 2.0, x, y);
```

HMPP Codelets and Regions

- A codelet is a pure function that can be remotely executed on a GPU
- Regions are a short cut for writing codelets

```
#pragma hmpp myfunc codelet, ...  
void saxpy(int n, float alpha, float x[n], float y[n])  
{  
    #pragma hmppcg parallel  
    for(int i = 0; i<n; ++i)  
        y[i] = alpha*x[i] + y[i];  
}
```

```
#pragma hmpp myreg region, ...  
{  
    for(int i = 0; i<n; ++i)  
        y[i] = alpha*x[i] + y[i];  
}
```


Codelet Target Clause

- Target clause specifies what GPU code to generate
 - *GPU* can be CUDA or OpenCL
- Choice of the implementation at runtime can be different!
 - The runtime select among the available hardware and code

```
#pragma hmpp myLabel codelet, target=[GPU], args[C].io=out
void myFunc( int n, int A[n], int B[n], int C[n]){
    ...
}
```

```
#pragma hmpp myLabel codelet, target=CUDA
```

→ NVIDIA only GPU

```
#pragma hmpp myLabel codelet, target=OpenCL
```

→ NVIDIA & AMD GPU, AMD CPU

HMPP Codelets Arguments

- The arguments of codelet are also allocated in the GPU device memory
 - Must exist on both sides to allow backup execution
 - No hardware mechanism to ensure consistencies
 - Size must be known to perform the data transfers
- Are defined by the **io** clause (in Fortran use intent instead)
 - **in** (default) : read only in the codelet
 - **out**: completely defined, no read before a write
 - **inout**: read and written
- Using inappropriate **inout** generates extra PCI bus traffic

```
#pragma hmpp myLabel codelet, args[B].io=out, args[C].io=inout
void myFunc( int n, int A[n], int B[n], int C[n]){
    for( int i=0 ; i<n ; ++i){
        B[i] = A[i] * A[i];
        C[i] = C[i] * A[i];
    }
}
```


Running a Codelet or Section on a GPU - 1

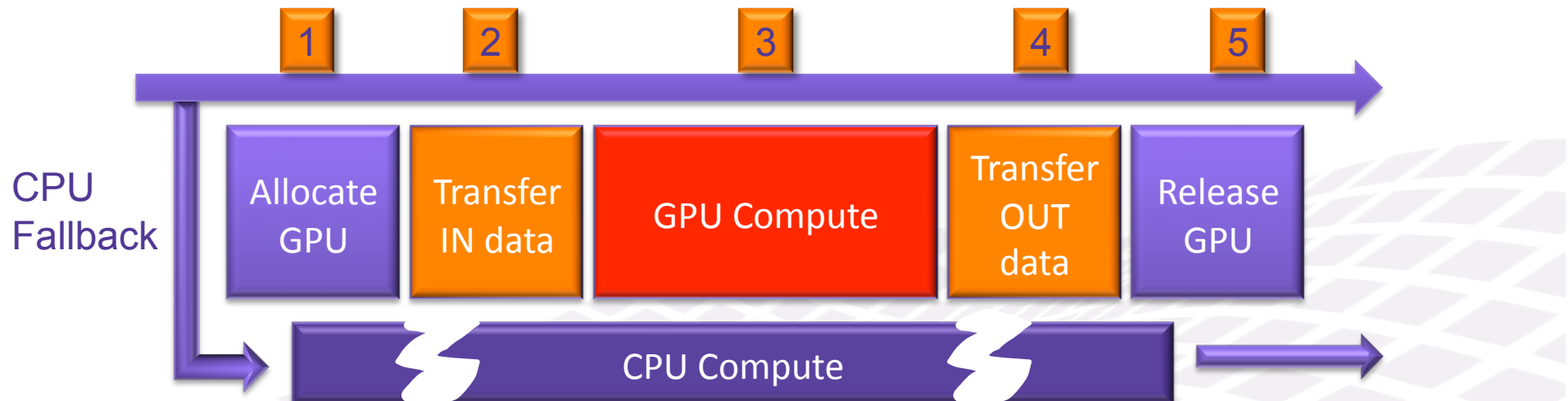
- The **callsite** directive specifies the use of a codelet at a given point in your application.
- **callsite** directive performs a Remote Procedure Call onto the GPU

```
#pragma hmpp call1 codelet, target=CUDA
#pragma hmpp call2 codelet, target=OpenCL
void myFunc(int n, int A[n], int B[n]){
    int i;
    for (i=0 ; i<n ; ++i)
        B[i] = A[i] + 1;
}

void main(void)
{
    int X[10000], Y[10000], Z[10000];
    ...
    #pragma hmpp call1 callsite, ...
    myFunc(10000, X, Y);
    ...
    #pragma hmpp call2 callsite, ...
    myFunc(1000, Y, Z);
    ...
}
```

Running a Codelet or Section on a GPU - 2

- By default, a CALLSITE directive implements the whole Remote Procedure Call (RPC) sequence
- An RPC sequence consists in 5 steps:
 - (1) Allocate the GPU and the memory
 - (2) Transfer the input data: CPU => GPU
 - (3) Compute
 - (4) Transfer the output data: GPU=> CPU
 - (5) Release the GPU and the memory



Tuning Hybrid Codes

- Tuning hybrid code consists in
 - Reducing penalty when allocating and releasing GPUs
 - Reducing data transfer time
 - Optimizing performance of the GPU kernels
 - Using CPU cores in parallel with the GPU
- HMPP provides a set of directives to address these optimizations
- The objective is to get efficient CPU and GPU computations

Reducing Data Transfers between CPUs and GPUs



- Hybrid code performance is very sensitive to the amount of CPU-GPU data transfers
 - PCIe bus is a serious bottleneck (< 10 GBs vs 150 GBs)
- Various techniques
 - Reduce data transfer occurrences
 - Share data on the GPU between codelets
 - Map codelet arguments to the same GPU space
 - Perform partial data transfers
- Warning: dealing with two address spaces may introduce inconsistencies

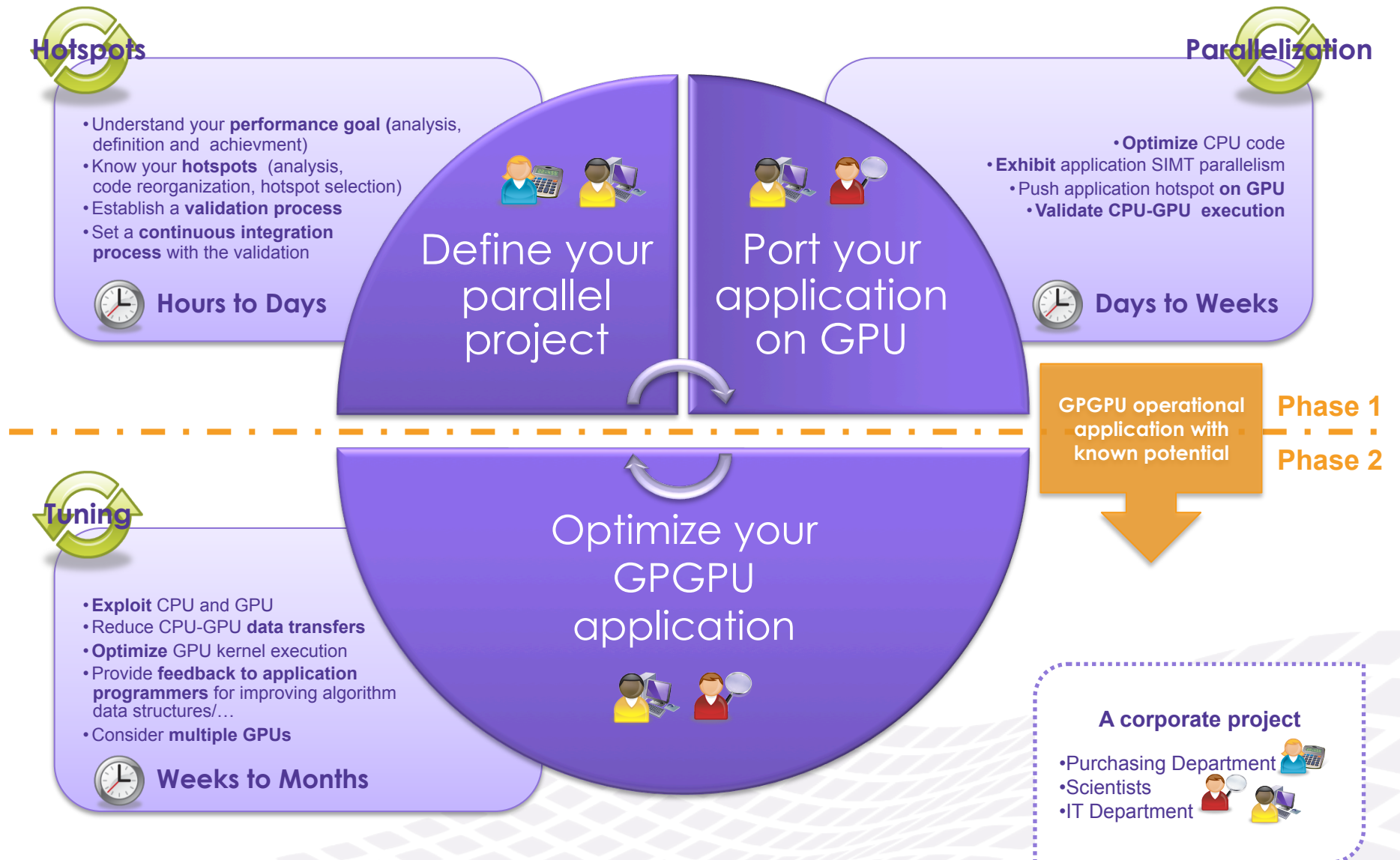
Tuning GPU Kernels

- GPU kernel tuning set-up parallel loop suitable for GPU architectures
- Multiple issues to address
 - Memory accesses
 - Thread grid tuning
 - Register usage tuning
 - Shared memory usage
 - Removing control flow divergence
- In many cases, CPU code structure conflicts with GPU efficient code structure

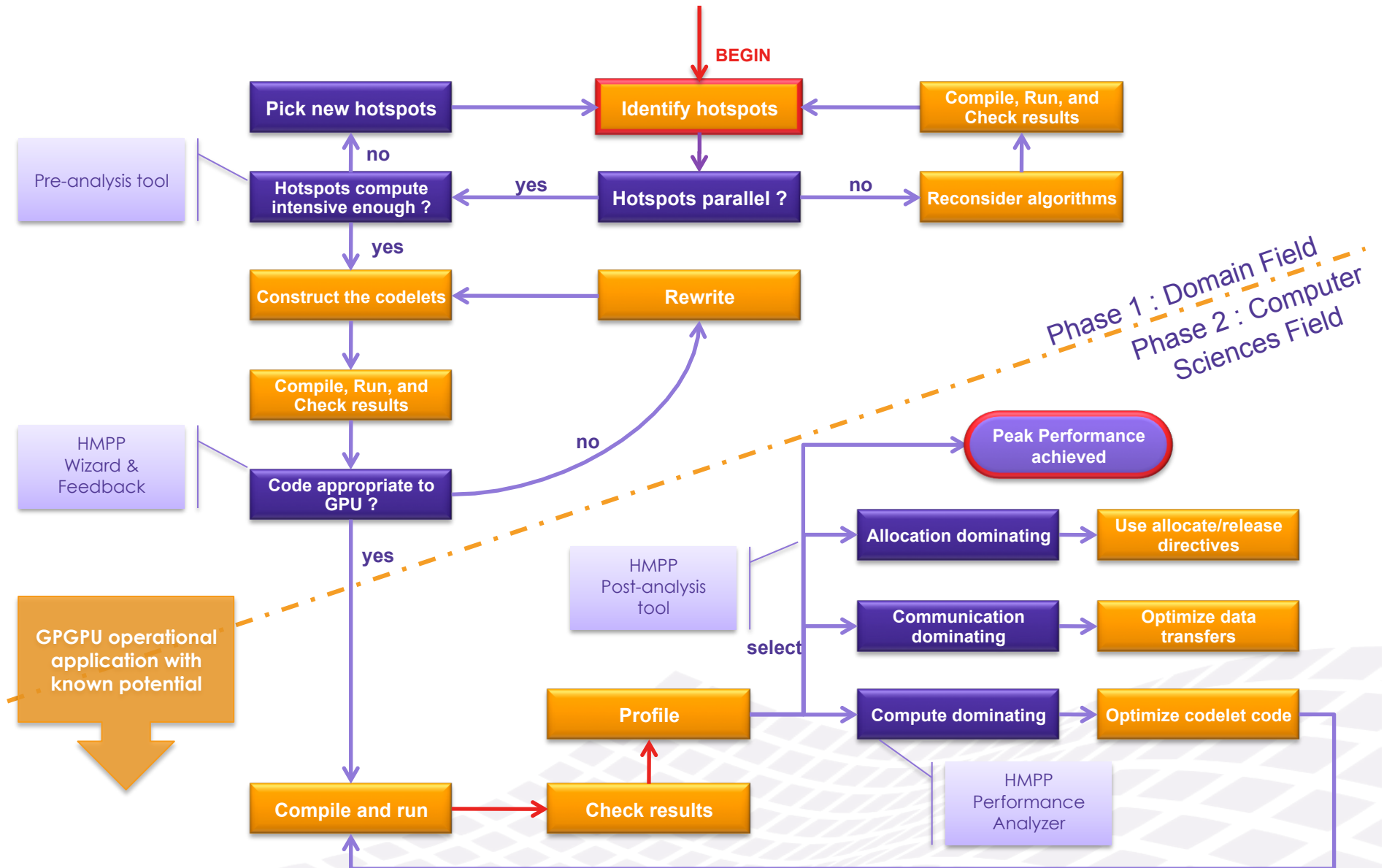
Methodology to Port Applications

- **Prerequisite**
 - Understand your performance goal
 - Memory bandwidth needs are a good potential performance indicator
 - Know your hotspots
 - Beware of Amdahl's law
 - Ensure you know how to validate the output of your application
 - Rounding may differ on GPUs
 - Determine if your goal can be achieved
 - How many CPUs and GPUs are necessary?
 - Is there similar existing code for GPUs (in CUDA, OpenCL or HMPP)?
- **Define an incremental approach**
 - Ensure to check the results at each step
- **Two phase approach**
 - Phase 1: Application programmers validate the computed results
 - Phase 2: Performance programmers focus on GPU code tuning and data transfer reduction

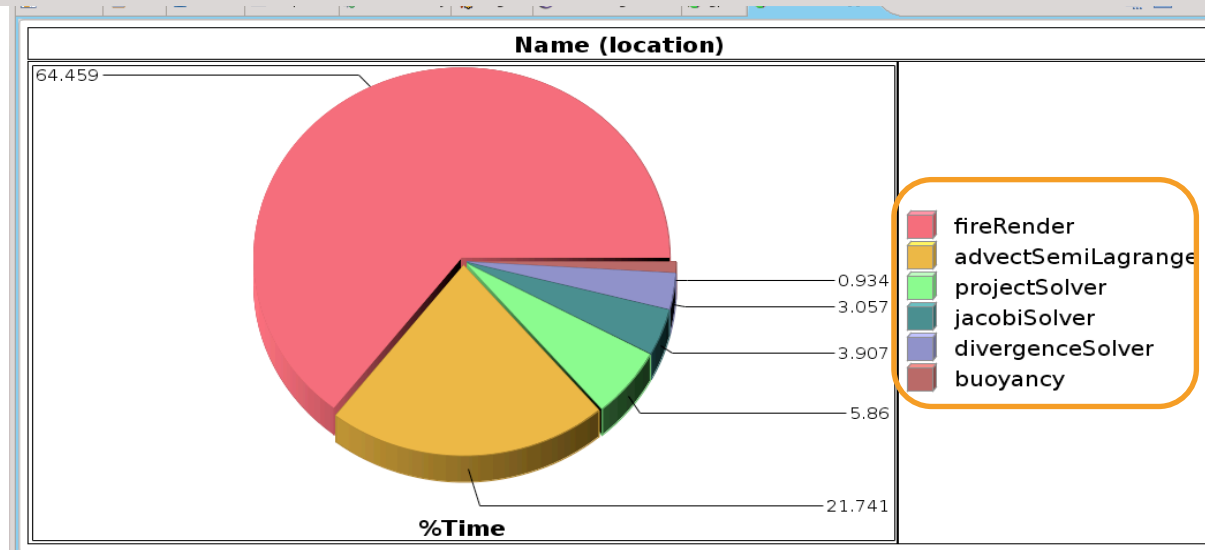
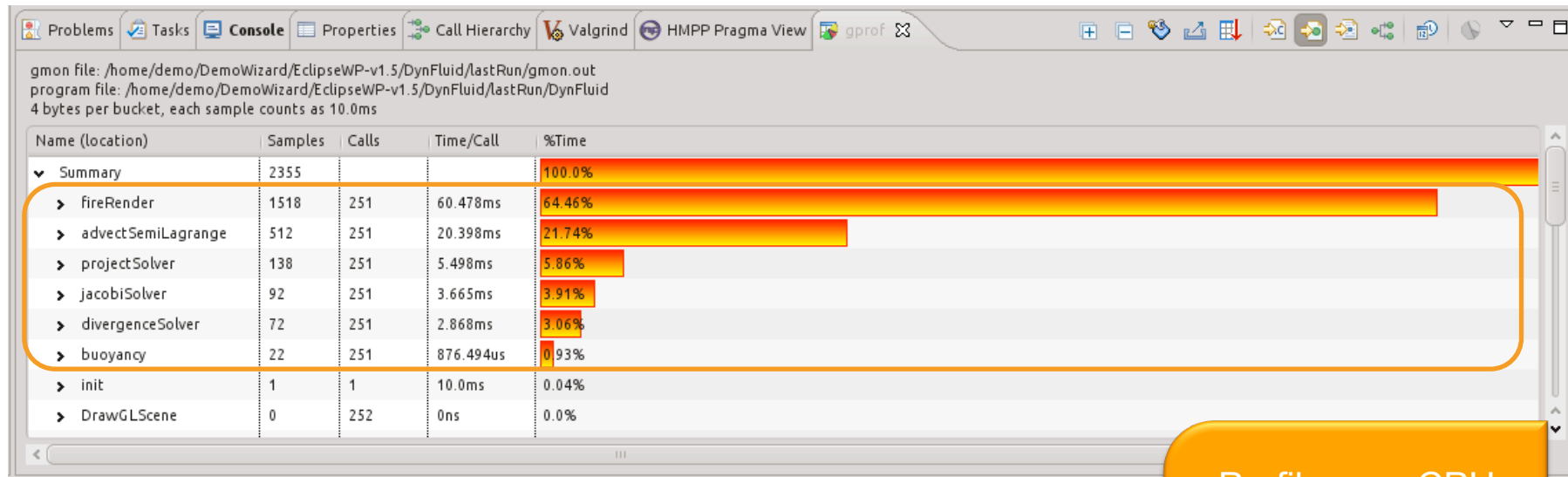
Methodology to Port Applications



Methodology Overview



Focus on Hotspots



Profile your CPU application

Build a coherent kernel set

Build Your GPU Computation with HMPP Directives (1)

Construct your
GPU group of
codelet

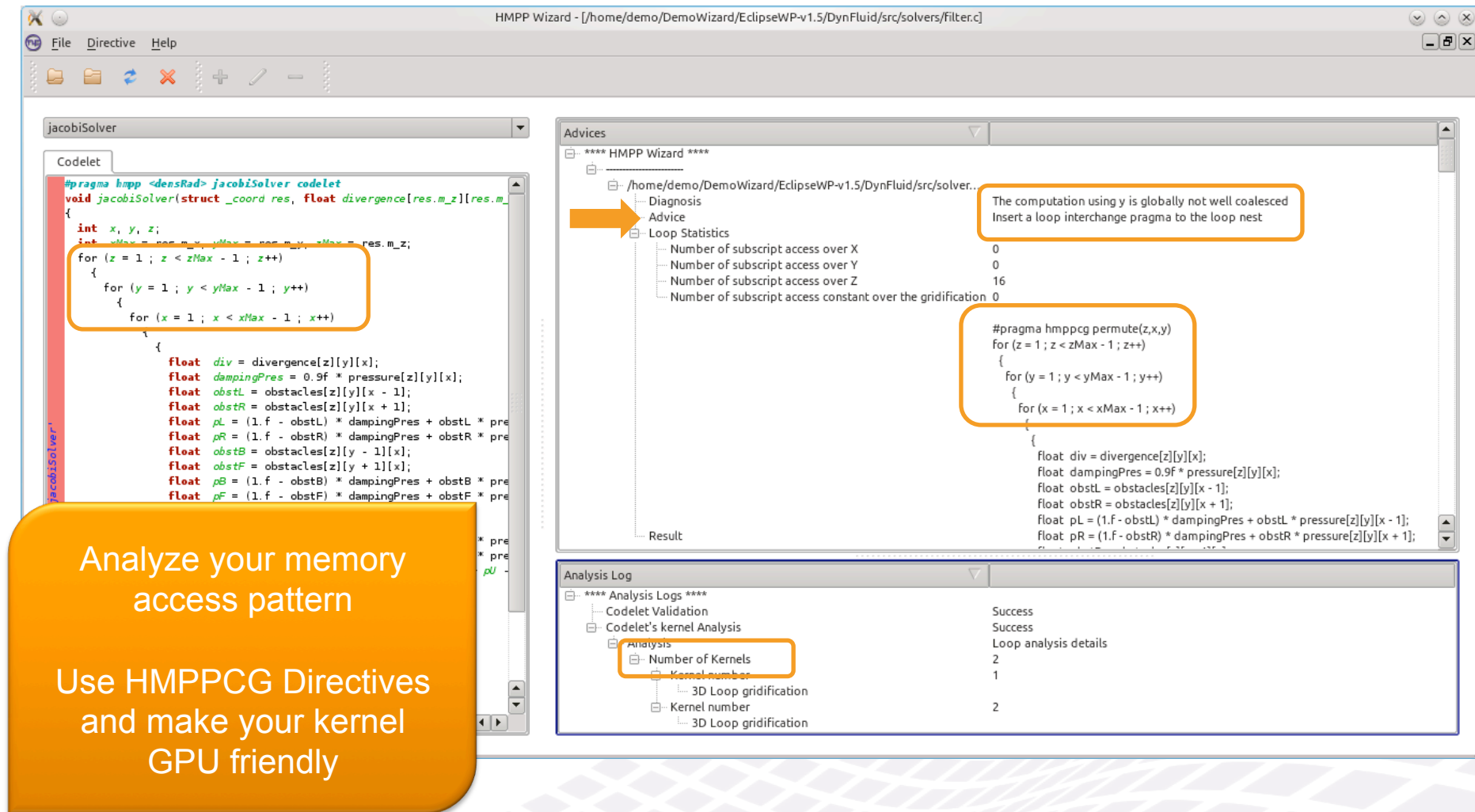
```
1 /*****  
2  *      solvers.c  
3  *****/  
4  
5 #ifndef FIRERENDER_H  
6 #define FIRERENDER_H  
7  
8 #include "solvers.h"  
9  
10  
11 #pragma hmpp <smallDensRad> fireRender codelet, args[res,velx,vely,velz,density,screenRes,posCam,distCamScreen,noise].io=i  
12 void fireRender(struct _coord res, struct _screenRes screenRes,  
13                float velx[res.m_z][res.m_y][res.m_x], float vely[res.m_z][res.m_y][res.m_x], float velz[res.m_z][res.m_y]  
14                float density[res.m_z][res.m_y][res.m_x],  
15                int distCamScreen, int posCam[3],  
16                float result[screenRes.m_y][screenRes.m_x][4], float noise[screenRes.m_y][screenRes.m_x]);  
17  
18 /* FIRERENDER_H */  
19 #endif  
20
```

Build Your GPU Computation with HMPP Directives (2)

... and use
Codelets in the
application

```
*fireRender.h  fireRender.c  filter.h  filter.c  solvers.h  solvers.c  fluidMotion.c
24{
25  int N = SZ(res);
26
27  tmpFields = (float*)malloc(N * 4 * sizeof(float));
28  tmpPressure = (float*)malloc(N * sizeof(float));
29  tmpDivergence = (float*)malloc(N * sizeof(float));
30
31  #pragma hmpp <smallDensRad> allocate
32}
33
34void solvers(float dto, float dx, struct _coord res, float* velx, float* vely, float* velz, float * pressure, float * dens
35{
36
37  buoyancy(res,velx,vely,velz,density, gravity);           //<==== || <====>
38  divergenceSolver(dx,res,velx,vely,velz,tmpDivergence,obstacles); //<==== || <====>
39  jacobiSolver(res,tmpDivergence,pressure,obstacles, tmpPressure); //<==== || <====>
40  projectSolver(dx,res,velx,vely,velz,pressure,obstacles); //<==== || <====>
41
42  #pragma hmpp <smallDensRad> advectSemiLagrange callsite
43  advectSemiLagrange(dto, res, density, velx, vely, velz, tmpFields); //<====
44}
45
46void fireRenderTo2D(struct _coord res, float* velx, float* vely, float* velz, float* density, int* posCam, int distCamScre
47{
48  #pragma hmpp <smallDensRad> fireRender callsite
49  fireRender(res, screenRes, velx, vely, velz, density, distCamScreen, posCam, result, noise);
50}
51
52void release(float* velx, float* vely, float* velz,float * pressure,float * density,float* obstacles) {
53  #pragma hmpp <smallDensRad> release
54  free(tmpDivergence);
55  free(tmpPressure);
56  free(tmpFields);
57}
58
59
```

Tune Your Kernels for GPUs with CAPS HMPP Wizard (1/2)



The screenshot displays the CAPS HMPP Wizard interface, which is used for analyzing and optimizing code for GPU execution. The interface is divided into several panes:

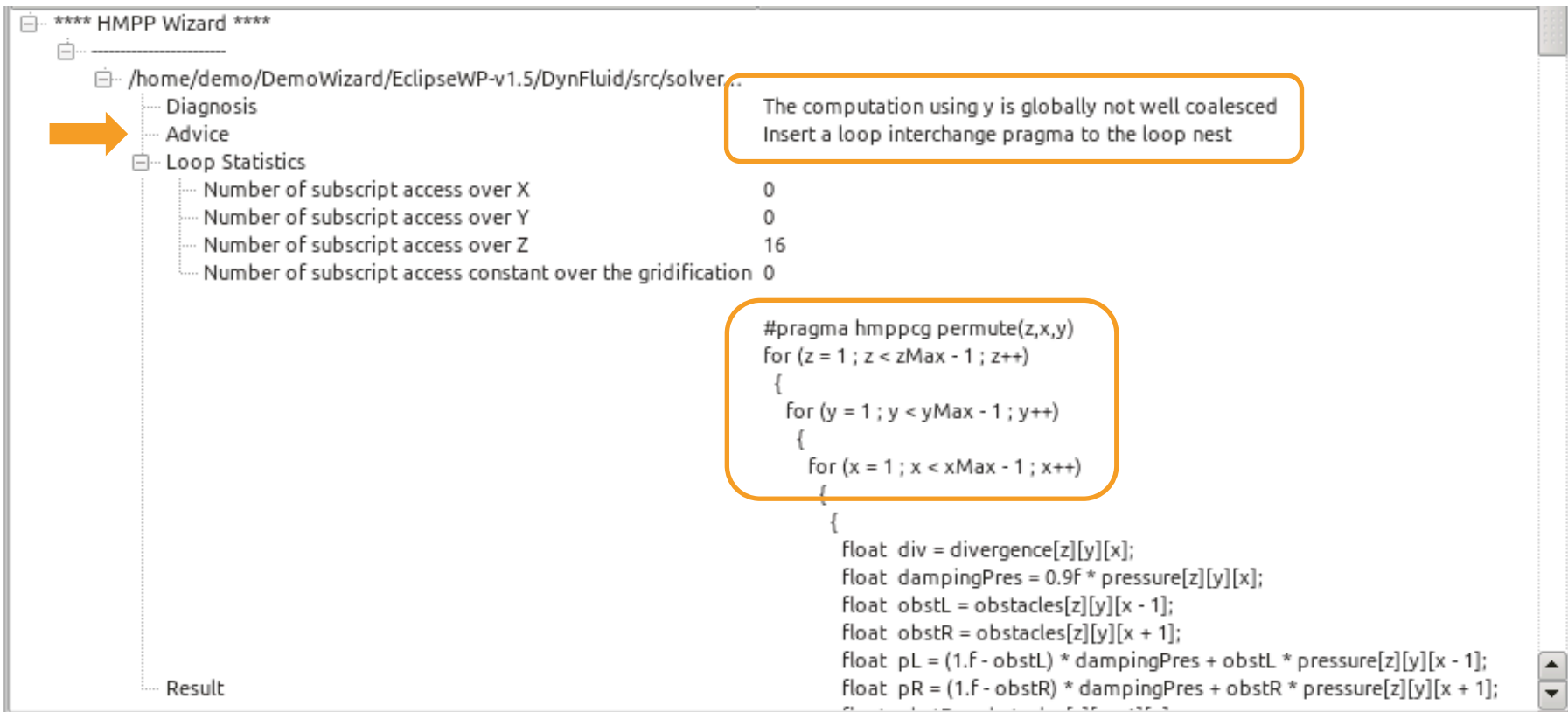
- Codelet:** Shows the source code of the `jacobiSolver` function. A yellow box highlights the nested loops for `z`, `y`, and `x`.
- Advices:** Provides analysis results and suggestions. A yellow box highlights the "Loop Statistics" section, which shows the number of subscript accesses over X, Y, and Z. Another yellow box highlights the "Advice" section, which suggests inserting a loop interchange pragma to improve coalescing.
- Analysis Log:** Shows the results of the analysis, including the number of kernels and the 3D loop gridification.

Two yellow callout boxes provide additional information:

- Analyze your memory access pattern**
- Use HMPPCG Directives and make your kernel GPU friendly**

The "Advices" pane also shows a suggested code transformation using the `#pragma hmppcg permute(z,x,y)` directive to reorder the loops for better GPU performance.

Tune Your Kernels for GPUs with CAPS HMPP Wizard (2/2)



The screenshot shows the CAPS HMPP Wizard interface. On the left, a tree view contains the following items: "Diagnosis", "Advice", "Loop Statistics", and "Result". An orange arrow points to the "Advice" item. The "Loop Statistics" section is expanded, showing the following data:

Statistic	Value
Number of subscript access over X	0
Number of subscript access over Y	0
Number of subscript access over Z	16
Number of subscript access constant over the gridification	0

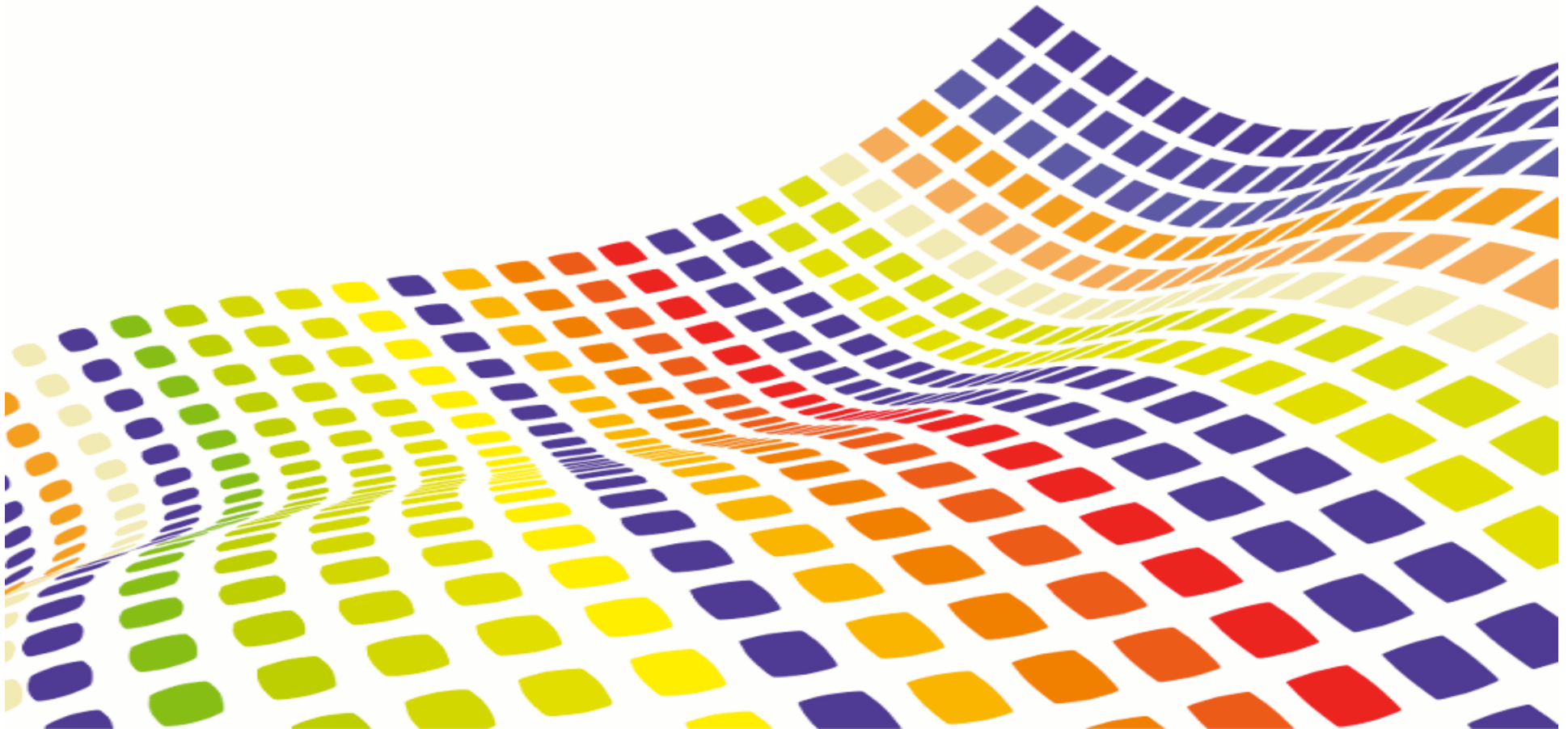
Below the statistics, the "Advice" section contains a text box with the following text:

The computation using y is globally not well coalesced
Insert a loop interchange pragma to the loop nest

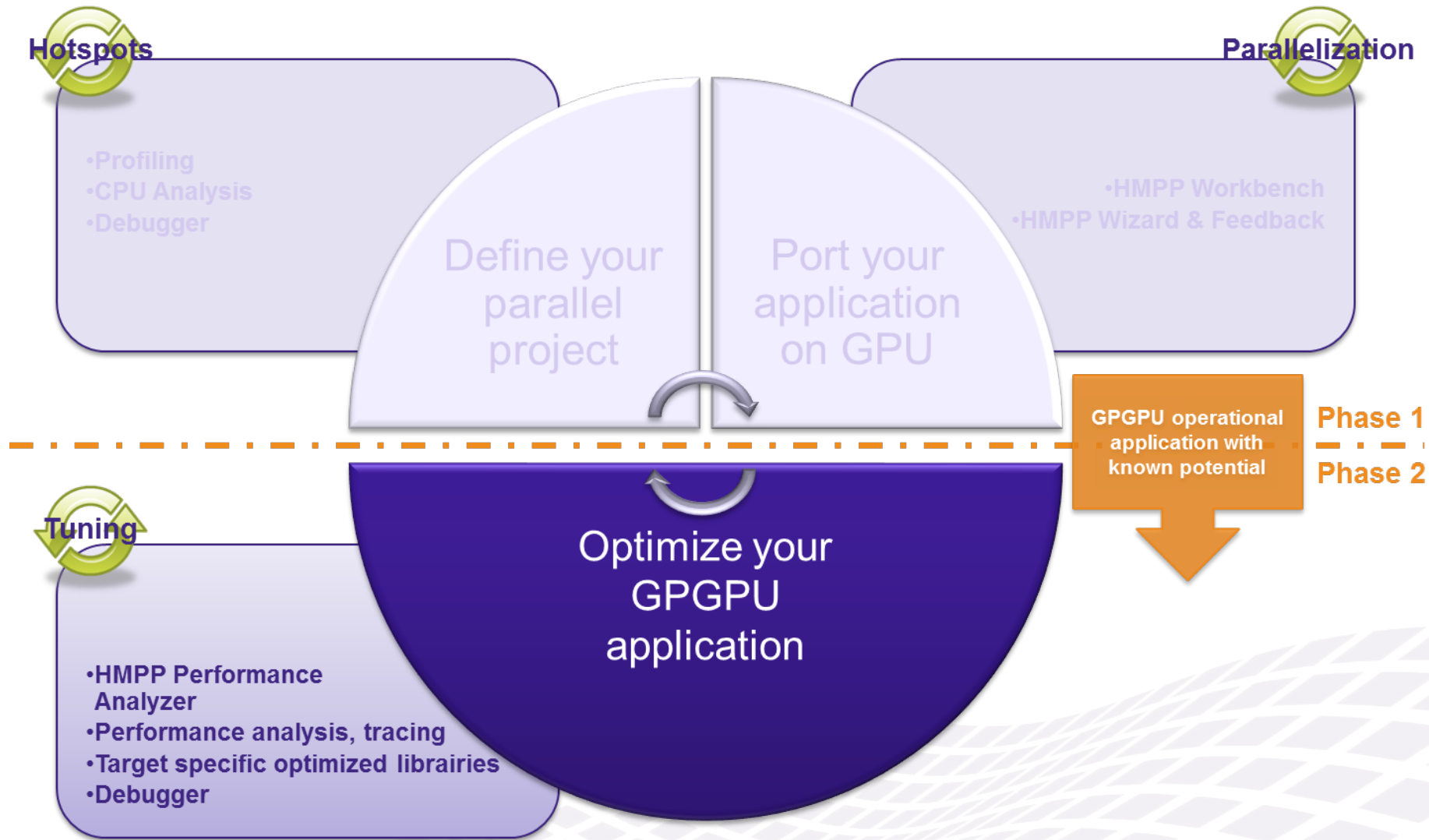
Below this, a code block shows the original loop structure with an orange box highlighting the inner loops:

```
#pragma hmppcg permute(z,x,y)
for (z = 1 ; z < zMax - 1 ; z++)
{
    for (y = 1 ; y < yMax - 1 ; y++)
    {
        for (x = 1 ; x < xMax - 1 ; x++)
        {
            float div = divergence[z][y][x];
            float dampingPres = 0.9f * pressure[z][y][x];
            float obstL = obstacles[z][y][x - 1];
            float obstR = obstacles[z][y][x + 1];
            float pL = (1.f - obstL) * dampingPres + obstL * pressure[z][y][x - 1];
            float pR = (1.f - obstR) * dampingPres + obstR * pressure[z][y][x + 1];
        }
    }
}
```

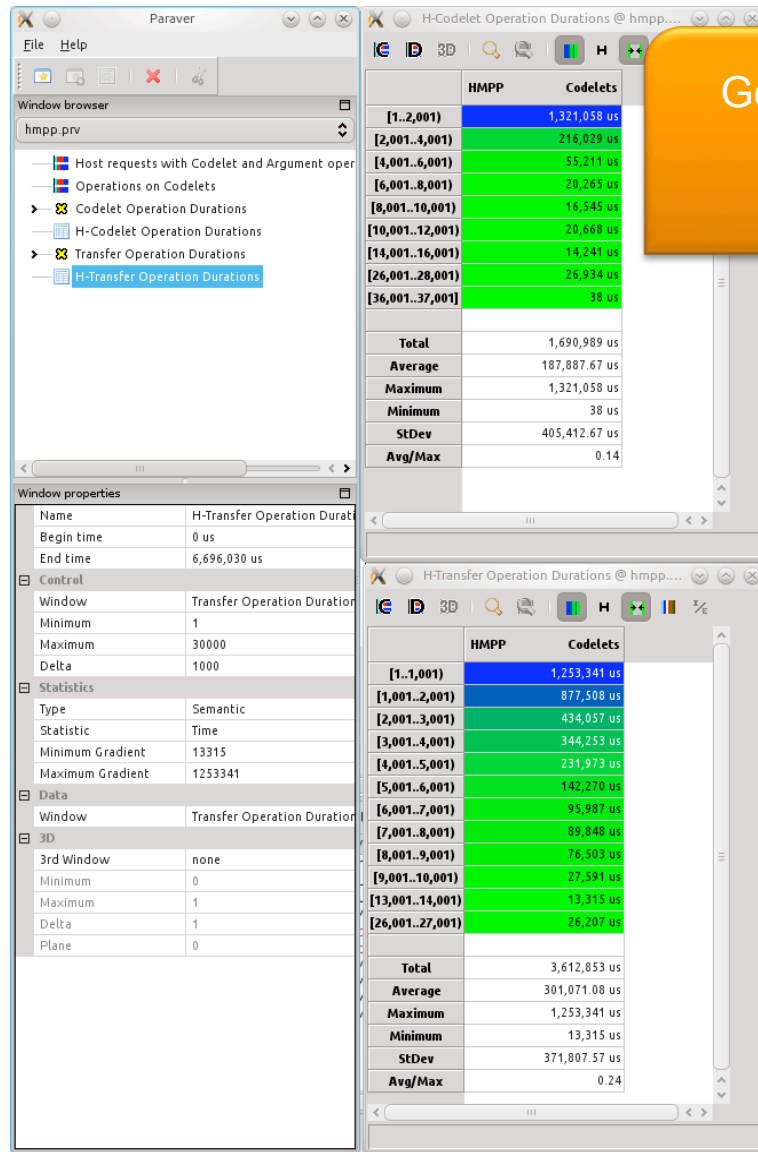
CAPS Tools to Port Your Applications – Phase 2



Optimizing Tools

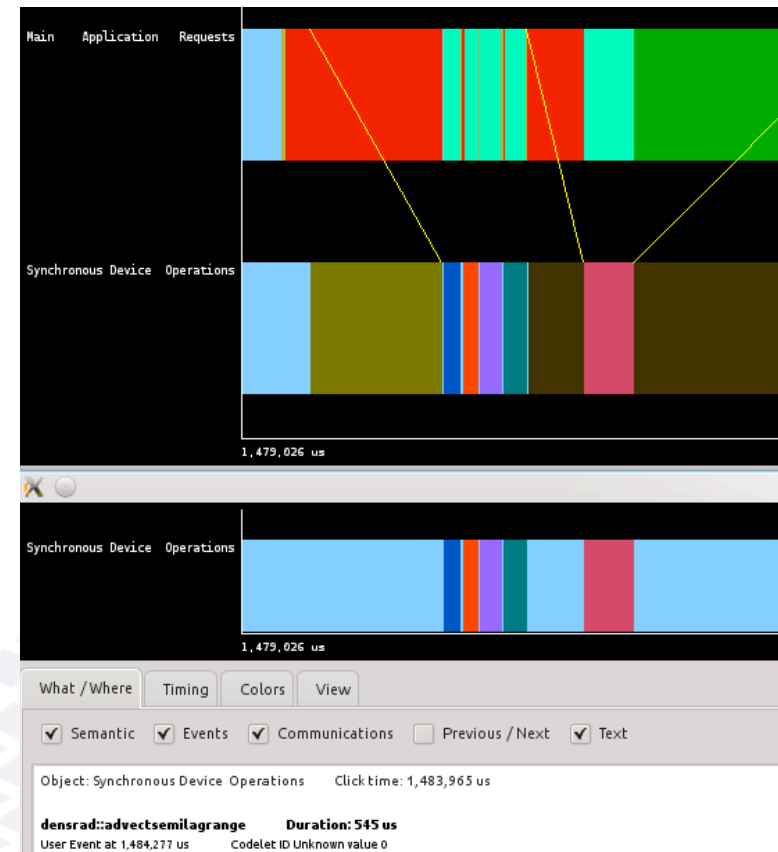


Analyze the GPU Code Porting Efficiency



Get a precise view of HMPP element behavior

Get statistics on GPU operations



Tune the GPU Execution Integration in Your Application with HMPP Directives

```
FireRender.c  filter.h  fluidMotion.c  solvers.c  filter.c

tmpFields    = (float*)malloc(N * 4 * sizeof(float));
tmpPressure  = (float*)malloc(N * sizeof(float));
tmpDivergence = (float*)malloc(N * sizeof(float));

#pragma hmpp <densRad> allocate

/* Constant Values. */
#pragma hmpp <densRad> advancedload, args[buoyancy::res]
#pragma hmpp <densRad> advancedload, args[divergenceSolver::dx, divergenceSolver::obstacles]
#pragma hmpp <densRad> advancedload, args[advectSemiLagrange::dto]

#pragma hmpp <densRad> advancedload, args[fireRender::noise]
#pragma hmpp <densRad> advancedload, args[fireRender::distCamScreen, fireRender::screenRes]

/* Prefetch. */
#pragma hmpp <densRad> advancedload, args[buoyancy::gravity]
#pragma hmpp <densRad> advancedload, args[buoyancy::velx, buoyancy::vely, buoyancy::velz]
#pragma hmpp <densRad> advancedload, args[divergence::density]
#pragma hmpp <densRad> advancedload, args[jacobiSolver::pressure]
#pragma hmpp <densRad> advancedload, args[buoyancy::density]
}

void solvers(float dto, float dx, struct _coord res, float* velx, float* vely, float* velz, float * pressure, float * density)
{
    #pragma hmpp <densRad> buoyancy callsite, args[velx, vely, velz].nouupdate, args[gravity].nouupdate
    buoyancy(res, velx, vely, velz, density, gravity); //<==== || >====>

    #pragma hmpp <densRad> divergenceSolver callsite, args[divergence].nouupdate, args[velx, vely, velz].nouupdate
    divergenceSolver(dx, res, velx, vely, velz, tmpDivergence, obstacles); //<==== || >====>

    #pragma hmpp <densRad> jacobiSolver callsite, args[tmpPressure].nouupdate, args[divergence].nouupdate, args[pressure].nouupdate
    jacobiSolver(res, tmpDivergence, pressure, obstacles, tmpPressure); //<==== || >====>

    #pragma hmpp <densRad> projectSolver callsite, args[velx, vely, velz].nouupdate, args[pressure].nouupdate
    projectSolver(dx, res, velx, vely, velz, pressure, obstacles); //<==== || >====>

    #pragma hmpp <densRad> advectSemiLagrange callsite, args[tmpFields].nouupdate, args[velx, vely, velz].nouupdate
    advectSemiLagrange(dto, res, density, velx, vely, velz, tmpFields); //<====>
}

void fireRenderTo2D(struct _coord res, float* velx, float* vely, float* velz, float* density, int* posCam, int distCamScreen)
{
```

Optimize out transfers from kernel calls

Optimize the GPU allocation and operate data prefetching

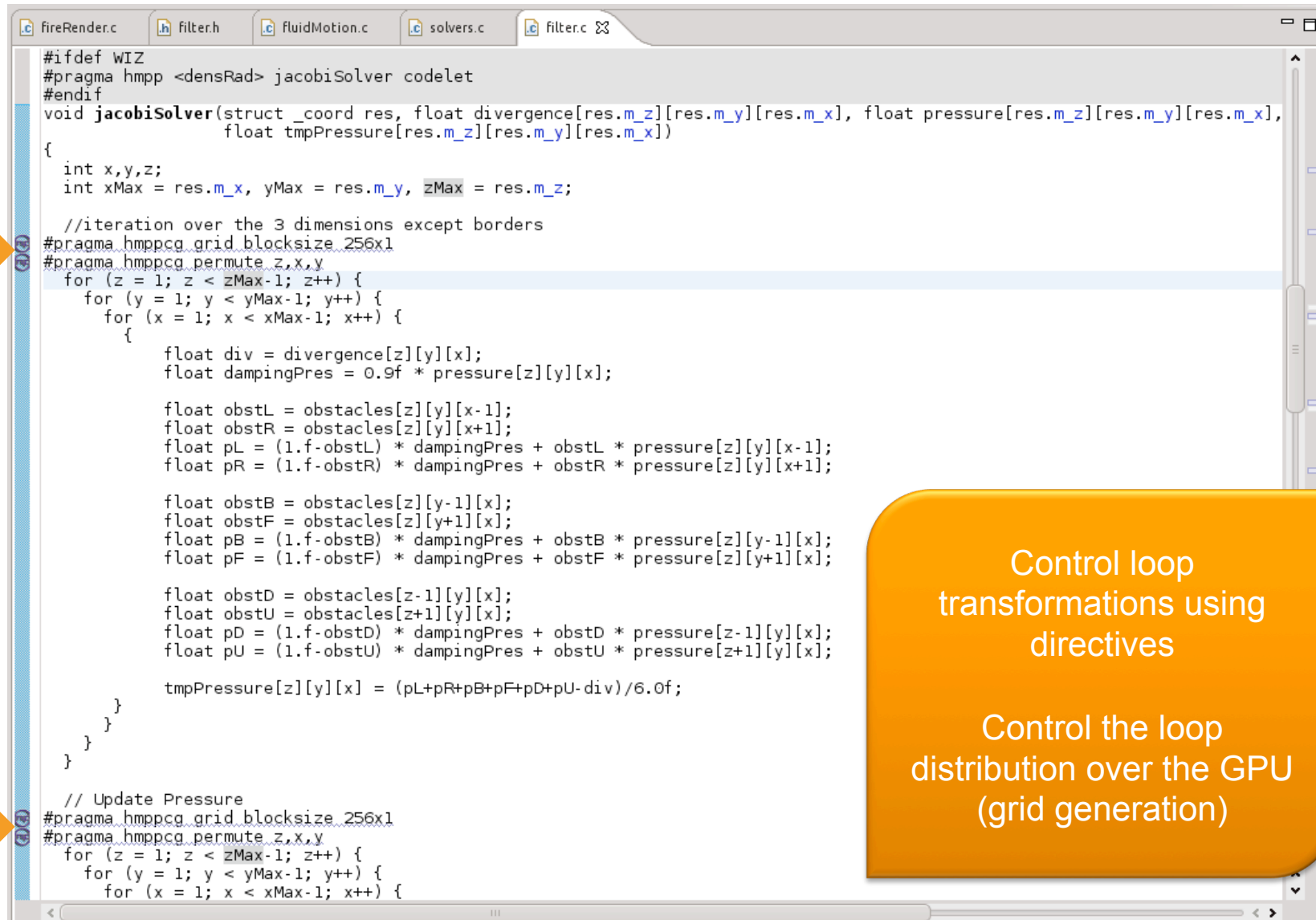
Analyze and profile kernel execution on the GPU with HMPP Performance Analyzer

The screenshot displays the HMPP Performance Analyzer interface. On the left, the 'Codelet' window shows C code for a Jacobi solver. The code is annotated with HMPP pragmas for grid blocksize and permutation. The 'Advices' window on the right provides performance metrics for a specific kernel. A table of metrics is highlighted with an orange box:

Advices	
**** HMPP PerfAnalyzer ****	
**** Advices ****	
**** Profile ****	
Kernel number	1
Kernel Name:	void __hmpp_codelet__jacobiSolver_loop0_<64u, 4u>(_coord, int, int, int, float*, float*, float*, float*)
Average gpu execution time:	230.031 us
Gridification:	grid 4x16=64 blocks, thread block size of 64x4x1, 16384 threads
Global memory read throughput:	36.11 GB/s/TPC
Global memory write throughput:	8.28 GB/s/TPC
Global memory throughput:	44.39 GB/s/TPC
L1 Global Load Hit Rate:	29.82 %

Below the advices, the 'Detailed metrics' section lists various performance indicators such as gputime, cputime, occupancy, and grid size. On the bottom left, the 'Analysis Log' window shows a success message for the analysis. On the right, a yellow callout box contains the text: 'Get precise and specific information about the kernel behavior' and 'Explore and Exploit at best the GPU power from the C source level'.

Optimize the GPU Kernel Code Generation with HMPPCG Directives



```
#ifdef WIZ
#pragma hmpp <densRad> jacobiSolver codelet
#endif
void jacobiSolver(struct _coord res, float divergence[res.m_z][res.m_y][res.m_x], float pressure[res.m_z][res.m_y][res.m_x],
float tmpPressure[res.m_z][res.m_y][res.m_x])
{
    int x,y,z;
    int xMax = res.m_x, yMax = res.m_y, zMax = res.m_z;

    //iteration over the 3 dimensions except borders
    #pragma hmppcg grid_blocksize 256x1
    #pragma hmppcg permute z,x,y
    for (z = 1; z < zMax-1; z++) {
        for (y = 1; y < yMax-1; y++) {
            for (x = 1; x < xMax-1; x++) {
                {
                    float div = divergence[z][y][x];
                    float dampingPres = 0.9f * pressure[z][y][x];

                    float obstL = obstacles[z][y][x-1];
                    float obstR = obstacles[z][y][x+1];
                    float pL = (1.f-obstL) * dampingPres + obstL * pressure[z][y][x-1];
                    float pR = (1.f-obstR) * dampingPres + obstR * pressure[z][y][x+1];

                    float obstB = obstacles[z][y-1][x];
                    float obstF = obstacles[z][y+1][x];
                    float pB = (1.f-obstB) * dampingPres + obstB * pressure[z][y-1][x];
                    float pF = (1.f-obstF) * dampingPres + obstF * pressure[z][y+1][x];

                    float obstD = obstacles[z-1][y][x];
                    float obstU = obstacles[z+1][y][x];
                    float pD = (1.f-obstD) * dampingPres + obstD * pressure[z-1][y][x];
                    float pU = (1.f-obstU) * dampingPres + obstU * pressure[z+1][y][x];

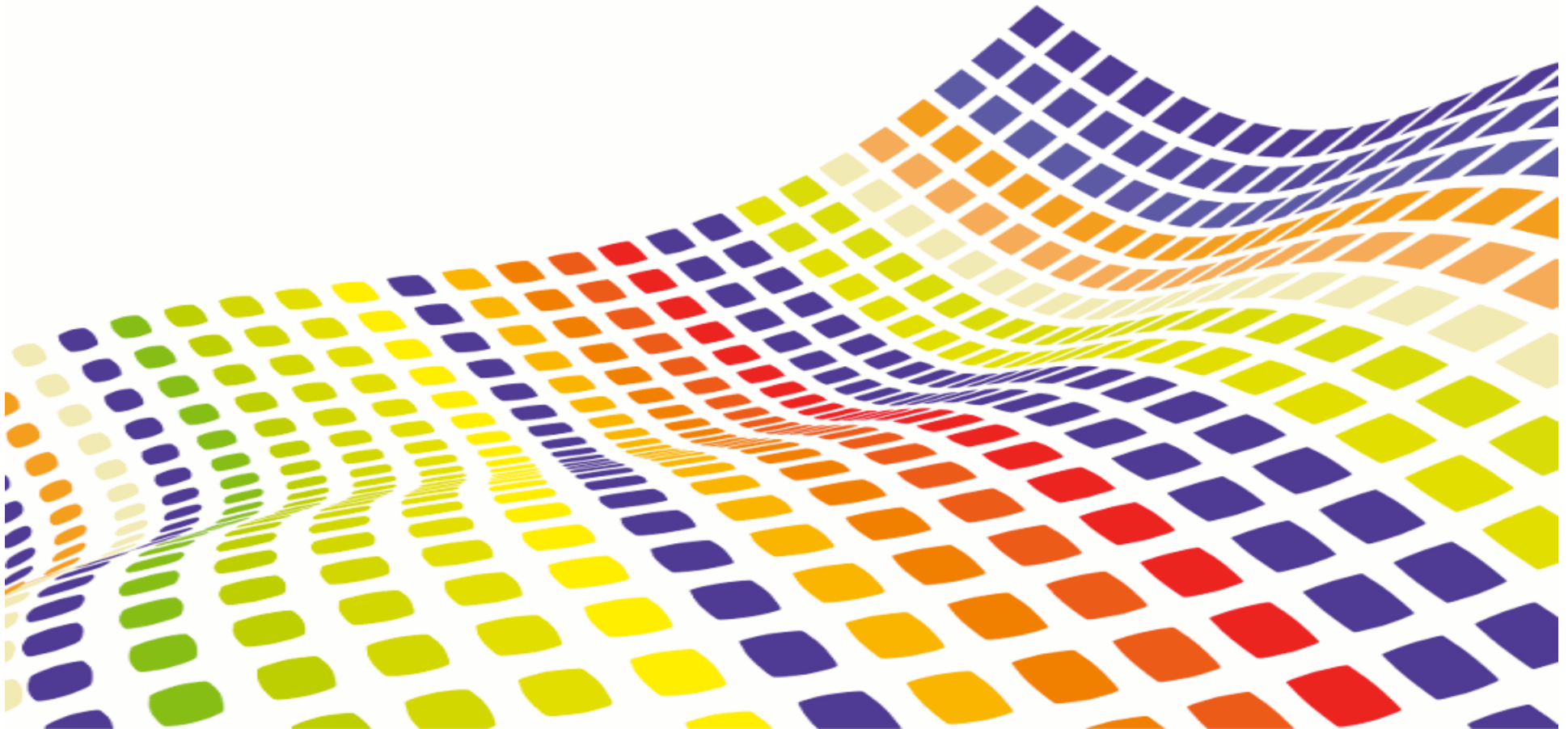
                    tmpPressure[z][y][x] = (pL+pR+pB+pF+pD+pU-div)/6.0f;
                }
            }
        }
    }

    // Update Pressure
    #pragma hmppcg grid_blocksize 256x1
    #pragma hmppcg permute z,x,y
    for (z = 1; z < zMax-1; z++) {
        for (y = 1; y < yMax-1; y++) {
            for (x = 1; x < xMax-1; x++) {
```

Control loop transformations using directives

Control the loop distribution over the GPU (grid generation)

Examples of Ported Applications



Examples of Ported Applications – 1

- Smoothed particles hydrodynamics

- Effort: 2 man-month
- Size: 22kLoC of F90 (SP or DP, MPI)
- GPU C1060 improvement: x 2 over serial code on Nehalem (x1.1 DP)
- Main difficulty: kernels limited to 70% of the execution time

The ratio performance over resource is the important information here.

- 3D Poisson equation, conjugate gradient

- Effort: 2 man-month
- Size: 2kLoC of F90 (DP)
- CPU improvement: x 2
- GPU C1060 improvement: x 5 over serial code on Nehalem
- Main porting operation: highly optimizing kernels
- Main difficulty: none

Examples of Ported Applications - 2

- **Electron propagation - solver**
 - Effort: 2 man-month
 - Size: 10 kLoC of F95 (DP, MPI)
 - CPU improvement: x 1.3
 - GPU C1060 improvement: x 1.15 over 4 thread code on Nehalem
 - Main porting operation: solver algorithm modifications
 - Main difficulty: small matrices, many data transfers
- **3D combustion code**
 - Effort: 2 man-month
 - Size: x100 kLoC of F90 (DP)
 - GPU C1060 improvement: ~x1 (data transfer limited) over serial code on Nehalem; C2050 x1.3
 - Main difficulty: execution profile shows few hot-spots (70%)
 - Next: code/algo. is being reviewed according to current results

Examples of Ported Applications - 3

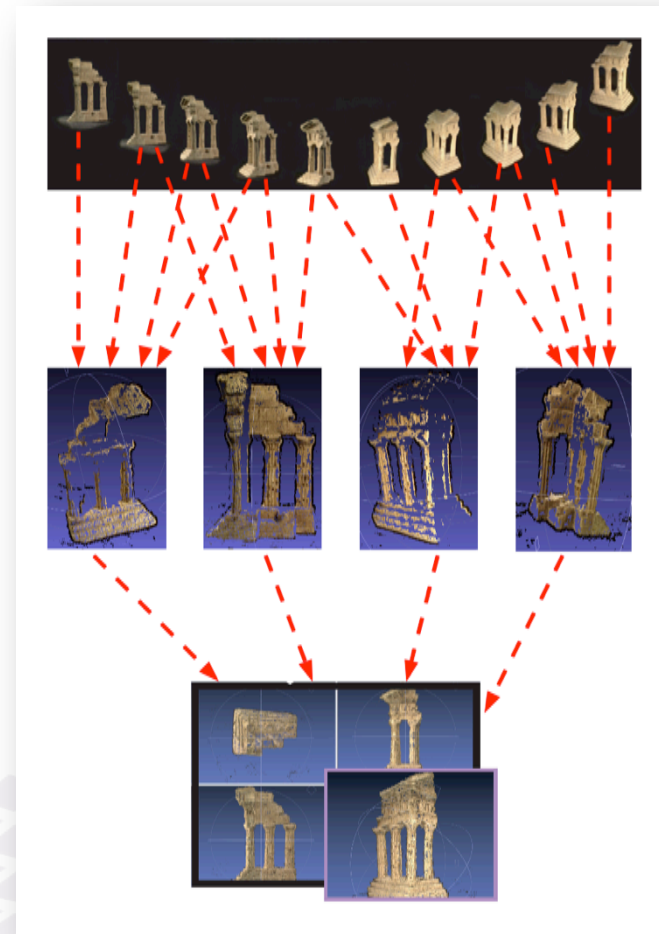
- Euler equations
 - Effort: <1 man-month
 - Size: ~40kLoC of F90 (DP)
 - CPU improvement: x 3 over the original code
 - GPU C1060 improvement: x 3 over serial code on Nehalem
 - Main porting operation: specializing the code for the main execution configuration
 - Main difficulty: reorganizing computational kernels (CPU dev. legacy)
- Tsunami/flood simulation
 - Effort: 0.5 man-month
 - Size: ~4kLoC (DP, MPI)
 - GPU C1060 improvement: x 1.28 over serial code on Nehalem (kernels speedup x30 and x18)
 - Next: highlight more parallelism, reducing data transfers (high performance potential)

Examples of Ported Applications - 5

- Weather models (GTC 2010 talk, M. Govett, NOAA)
 - Effort: 1 man-month (part of the code already ported)
 - GPU C1060 improvement: 10x over the serial code on Nehalem
 - Main porting operation: reduction of CPU-GPU transfers
 - Main difficulty: GPU memory size is the limiting factor

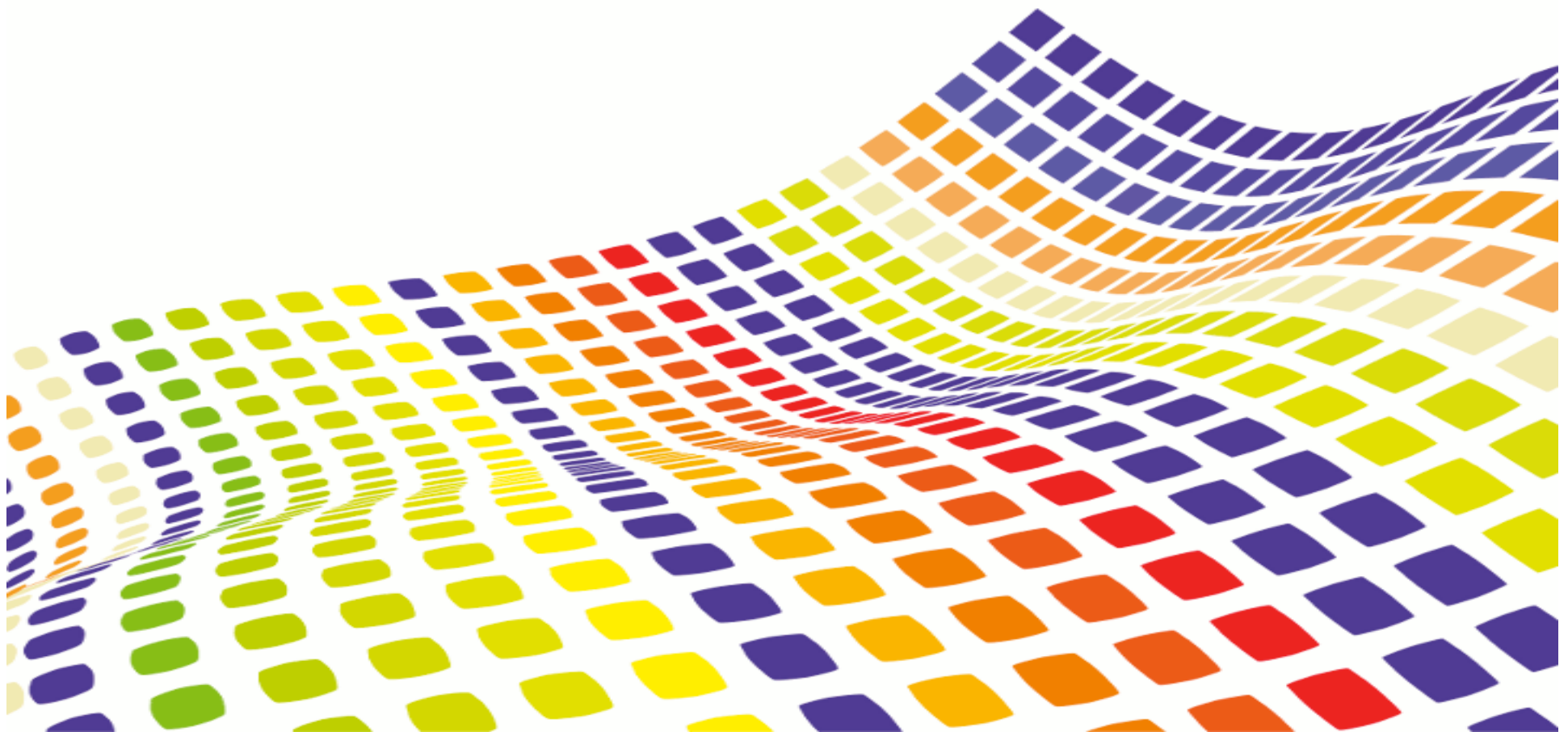
MultiView Stereo

- Resource spent
 - 1 man-month
- Size
 - ~1kLoC of C99 (DP)
- HMPP Basic version (1hour)
 - GPU C2050 improvement
 - X 30
 - Main porting operation
 - Adding 4 directives
- HMPP fine tune version (2 weeks)
 - GPU C2050 improvement
 - X 500
 - Main porting operation
 - Rethinking algorithm



Conclusion

- **Heterogeneous architectures are becoming ubiquitous**
 - In HPC centers but not only
 - Tremendous opportunities but not always easy to seize
 - CPU and GPU have to be used simultaneously
- **Legacy codes still need to be ported**
 - An efficient methodology is required
 - A methodology supporting tools is needed and must provide a set of consistent views
 - The legacy style is not helping
 - Highlighted parallelism for GPU is useful for future manycores
- **HMPP based programming**
 - Helps implementing incremental strategies
 - Is being complemented by a set of tools
 - Engage in an Open Standard path with Pathscale



Reducing Data Transfers Occurrences

- Preload data before codelet call
 - Load data as soon as possible

Preload data

```
int main(int argc, char **argv) {  
  
#pragma hmpp sgemm allocate, args[vin1;vin2;vout].size={size,size}  
    . . .  
  
#pragma hmpp sgemm advancedload, args[vin1;m;n;k;alpha;beta]  
  
    for( j = 0 ; j < 2 ; j++ ) {  
#pragma hmpp sgemm callsite &  
#pragma hmpp sgemm  args[m;n;k;alpha;beta;vin1].advancedload=true  
        sgemm( size, size, size, alpha, vin1, vin2, beta, vout );  
        . . .  
    }  
  
    . . .  
#pragma hmpp sgemm release
```

Avoid reloading data

Sharing Data Between Codelets with Resident Data

- Share data between codelets of the same group
 - Keep data on the HWA between two codelet calls
 - Avoid useless data transfers

```
#pragma hmpp <process> group, target=CUDA
#pragma hmpp <process> resident
float initValue = 1.5f, offset[9];
...
#pragma hmpp <process> reset1 codelet, args[t].io=out
void reset(float t[M][N]){
    int i,j;
    for (i = 0; i < M; i += 1) {
        for (j = 0; j < N; j += 1) {
            t[i][j] = initValue + offset[(i+j)%9];
        }
    }
}
#pragma hmpp <process> process codelet, args[a].io=inout
void process(real a[M][N], real b[M][N]){
    int i,j;
    for (i = 0; i < M; i += 1) {
        for (j = 0; j < N; j += 1) {
            a[i][j] = cos(a[i][j]) + cos(b[i][j]) - initValue;
        }
    }
}
```