

# The source to source transformation from CUDA to SM-Centric\*

†

Haoxu Ren

Department of Computer Science  
North Carolina State University  
hren3@ncsu.edu

## Abstract

This project is aimed to implementing the source to source transformation from CUDA code to SM-Centric code. To achieve this transformation, this report choose the AST matcher provided by Clang. With the matcher and rewriter, transform the input file which means the CUDA code file to SM-Centric code file. In the new SM-Centric file, it includes the new format code, and compared with the original CUDA code, it implements a flexible control of GPU scheduling at the program level becomes feasible, which opens up new opportunities for GPU program optimizations.

**Keywords** CUDA code, SM-Centric, Clang, LLVM

## 1 Introduction

Because of the high computational cost of current compute-intensive applications, many scientists view graphic processing units (GPUs) as an efficient means of reducing the execution time of their applications. High-end GPUs include an extraordinary large amount of small computing units along with a high bandwidth to their private on-board memory

The goal of this project is to develop a source-to-source translator that can convert an input CUDA code to a SM-Centric form code. An SM-Centric form has the advantage that it allows the execution of the GPU code to be flexibly controlled in schedule, such as assigning to the same SM a certain set of tasks that share lots of data.

## 2 Background

A GPU's computing power lies in its abundant memory bandwidth and massive parallelism. However, its hardware thread schedulers, despite being able to quickly distribute computation to processors, often fail to capitalize on program characteristics effectively, achieving only a fraction of the GPU's full potential. Moreover, current GPUs do not allow programmers or compilers to control this thread scheduling, forfeiting important optimization opportunities at the program level.

SM represents Streaming Multiprocessors, with this technique, flexible control of GPU scheduling at the program level becomes feasible, which opens up new opportunities for GPU program optimizations.

An SM-Centric form has the advantage that it allows the execution of the GPU code to be flexibly controlled in schedule, such as assigning to the same SM a certain set of tasks that share lots of data.

In order to implement the automatic tool that transforms source code employing CUDA extensions into plain C code, a source-to-source transformation framework has been leveraged. Different options for this class of source transformations are available nowadays, from simple pattern string replacement tools to frameworks which parse the source code into an Abstract Syntax Tree (AST) and transform the code using that information. Given that our tool needs to do complex transformations involving semantic C++ code information, we have selected the latter.

## 3 Objective

The objective of this project is to achieve the transformation from CUDA code to SM-Centric code.

This part briefly explains the code transformations involved in the SM-Centric Transformation. In detail, some assumptions taken by us in this explanation are as follows:

- \* The program contains only one CUDA kernel function, and its definition and invocation are in a single file.
- \* The grid(...) call is inside the same function that calls the CUDA kernel function.
- \* The grid in the original program is one or two dimensional.

And in detail, the detailed requirements of this project shown as following:

1. Add the following line to the beginning of the CUDA file after all existing #include statements:  
  
`#include "smc.h"`
2. add the following line to the beginning of the definition of the kernel function  
  
`_SMC_Begin`
3. add the following line to the end of the definition of the kernel function:  
  
`_SMC_End`

- ```
dim3 __SMC_orgGridDim,
int __SMC_workersNeeded,
int * __SMC_workerCount,
int * __SMC_newChunkSeq,
int * __SMC_seqEnds
```

- ```
(int)fmodf((float)__SMC_chunkID,  
(float)__SMC_orgGridDim.x);
```

- ```
(int)(__SMC_chunkID/__SMC_orgGridDim.x);
```

- dim3 \_\_SMC\_orgGridDim(...)*

- ```
SMC_init();
```

- ```

_SMC_orgGridDim,
_SMC_workersNeeded,
_SMC_workerCount,
_SMC_newChunkSeq,
_SMC_seqEnds

```

The second technique is filling-retreating scheme, which offers a flexible control of the amount of active threads on an SM. Importantly, the control is resilient to the randomness and obscurity in GPU hardware thread scheduling. It helps SM-centric transformation in two aspects. First, it ensures an even distribution of active threads on SMs, which is vital for guaranteeing the correctness of SM-centric transformations. Second, it facilitates online determination of the parallelism level suitable for a kernel, which is especially important for the performance of multiple-kernel co-runs, a scenario benefiting significantly from SM-centric transformation.

1. The transiator is built and the same environment as the LLVM Virtual Machine provided. Otherwise, the instructions are guaranteed to work.
2. The original CUDA file to be transformed should be with the suffix “.cu” to ensure the naming of generated file.

#### For Build part:

1. Direct to the Clang tool source folder(/home/ubuntu/llvm/tools/clang/tools)
2. Create a folder named sm-centric
3. Put the SM\_centric\_transformation.cpp(source code) and CMakeLists.txt(Dependency file) under the sm-centric folder. Both files can be found in the Source Code folder.
4. Modify the CMakeLists.txt in the clang tool folder. Append one line to the end of the file. `add_clang_subdirectory(sm-centric)`
5. Direct to the build-release folder (/home/ubuntu/llvm/build-release)
6. Run `ninja sm-centric`

Once the translator is built successfully, message like following will be displayed.

```
ubuntu@ubuntu:~/llvm/build-release$ ninja sm-centric
[2/2] Linking CXX executable bin/sm-centric
```

**Figure 3:** The success of translator built

#### For Run part:

We put the transform-needed code into a folder, in my computer, the folder is /home/ubuntu/Downloads, so we can run the transformation code as following:

```
ubuntu@ubuntu:~/llvm/build-release$ bin/sm-centric /home/ubuntu/Downloads/matrixMul_org.cu
u -- --cuda-host-only
/home/ubuntu/Downloads/matrixMul_org.cu:36:10: fatal error: 'helper_functions.h' file not found
#include <helper_functions.h>
      ^
1 error generated.
Error while processing /home/ubuntu/Downloads/matrixMul_org.cu.
```

**Figure 4:** The running of the transformation code

For the transformation code part, to achieve these requirements, I need to use the AST matcher provided by Clang and llvm.

So, I met the first challenge of the code part, I need to write the different matcher for different objective. When analyze the project requirements, for these 9 requirements, it can be divided into serval parts. One is handle for cuda kernel function declaration, one is handle for the call of cuda kernel function, one is handle for member expression declaration, and another is handle for variable

declaration. So, I have a basic knowledge of these matchers as following shown:

```
public:
MyASTConsumer(Rewriter &R) : cudaPrinter(R), HandleForFunction(R), HandleForDecl(R) {
    // The matcher of kernel function
    Matcher.addMatcher(
        functionDecl().bind("kernelfunc"), &HandleForFunction);
    // The matcher of blockIdx.x
    Matcher.addMatcher(
        memberExpr().bind("VarDecl"), &HandleForDecl);
    // The matcher of grid() declaration
    Matcher.addMatcher(
        varDecl().bind("grid"), &HandleForDecl);
    // The matcher of the call of Cuda kernel function
    Matcher.addMatcher(
        cudaKernelCallExpr().bind("cudaCall"), &HandleForDecl);
}
```

**Figure 5:** The original matcher design

As the figure 5 shown, I design these matchers to achieve the requirements of the project, but it's just the original version, what's next is to perfect these matchers and implement these requirements.

With AST matcher, every statement and declaration in the transform-needed file can be treated as different nodes. If I can get the nodes, then I can get the important source location for the next replacement and inserting. With these information, I can work on implementing these matchers.

```
Matcher.addMatcher(
    functionDecl().bind("kernelfunc"), &HandleForFunction);
// The matcher of blockIdx.x
Matcher.addMatcher(
    memberExpr().bind("VarDecl"), &HandleForDecl);
// The matcher of grid() declaration
Matcher.addMatcher(
    varDecl().bind("grid"), &HandleForDecl);
// The matcher of the call of Cuda kernel function
Matcher.addMatcher(
    cudaKernelCallExpr().bind("cudaCall"), &HandleForDecl);
```

**Figure 6:** The implementation of matcher

From the figure above, we can see I design these 4 matchers for all requirements, include `functionDecl`, `memberexpr`, `vardecl`, `cudakernelcallexpr`. With these original matcher, I narrow these matchers and get the nodes which met the requirements. Then in the handle implementation part, I use `Result.Nodes.getNodeAs()` function to locate the chosen node, and then invoke the `Rewriter` to implement the replacement and inserting.

The following figure shows the one of implementation:

```
class FunctionDeclHandler : public MatchFinder::MatchCallback {
public:
    FunctionDeclHandler(Rewriter &Rewriter) : Rewriter(Rewriter) {}

    virtual void run(const MatchFinder::MatchResult &Result) {
        // get the node "kernelfunc"
        if (Result.Nodes.getNodeAs<FunctionDecl>("kernelfunc")) {
            // of this function is cuda kernel function
            if (Result.Nodes.getNodeAs<FunctionDecl>("kernelfunc")) {
                // get the function body of this cuda kernel function
                const Stmt *FuncBody = Result.Nodes.getNodeAs<FunctionDecl>("kernelfunc")->getBody();
                // check the "hasAttr" in the top of the cuda kernel function
                RewriteResult Res = Rewriter.Replace(FuncBody, "hasAttr(<CUDA_GlobalAttr>)", "hasAttr(<CUDA_GlobalAttr>)", true, true);
                // check the "hasAttr" in the top of the cuda kernel function
                RewriteResult Res = Rewriter.Replace(FuncBody, "hasAttr(<CUDA_GlobalAttr>)", "hasAttr(<CUDA_GlobalAttr>)", true, true);
                // check the "hasAttr" in the top of the cuda kernel function
                RewriteResult Res = Rewriter.Replace(FuncBody, "hasAttr(<CUDA_GlobalAttr>)", "hasAttr(<CUDA_GlobalAttr>)", true, true);
                // check the "hasAttr" in the top of the cuda kernel function
                RewriteResult Res = Rewriter.Replace(FuncBody, "hasAttr(<CUDA_GlobalAttr>)", "hasAttr(<CUDA_GlobalAttr>)", true, true);
            }
        }
    }
};
```

**Figure 7:** The implementation of FunctionDeclHandle

In this implementation, it implements the location of the declaration of CUDA kernel function. The method `hasAttr(<CUDA_GlobalAttr>)` is the most important part in this location process. This method narrows the range of existing function declaration to the function declaration started with “\_global\_”, which means the CUDA kernel function declaration.

It's exactly what we want. When locate the CUDA kernel function declaration, what the next is to get the body of the function, it can implement by the method: *getbody()*, and with this function, I get the location inside this function which means inside the compounds, so I can add the “*\_\_SMC\_Begin*” and the “*\_\_SMC\_End*” for this kernel function. In detail, I need to invoke the Rewriter's *InsertText()* method and *InsertTextAfterToken()* method to implement these 2 requirements. In the next, I need to locate the parameter declaration part of this function, so I can use the function *getParamDecl()* to achieve this. Then as the same above, invoke the *InsertTextAfterToken()* method to insert new parameters to the end of parameter declaration location. The new added parameter is of the SM-Centric format, which is like:

```
dim3 __SMC_orgGridDim, int __SMC_workersNeeded, int
*__SMC_workerCount, int *__SMC_newChunkSeq, int *
__SMC_seqEnds
```

With these descriptions above, I implemented the transformation for the CUDA kernel function declaration part. And I will show the result part in the next Result part.

And in this part, I will continue show the resting implementation of other matcher handlers.

For Variable Declaration Handler:

```
class VarDeclHandler : public MatchFinder::MatchCallback {
public:
    VarDeclHandler(Rewriter &rewriter) : Rewriter(rewriter) {}

    virtual void run(const MatchFinder::MatchResult &result) {
        //get the node "VarDecl"
        if (const MemberExpr *Var = Result.Nodes.getNodeAs("VarDecl")) {
            //get the right side name of blockDecl
            auto string_name = Var->getBaseExpr()->getFullNameAsString();
            //if the right side's name is x
            if (name == "fetch_builtin_x")
                //replace the "blockDecl.y" with "(int) fmod((float) __SMC_chunkID, (float) __SMC_orgGridDim.x)"
                Rewrite.ReplaceText(Var->getLocEnd().getLocWithOffset(-9), 11, "(int) fmod((float) __SMC_chunkID, (float) __SMC_orgGridDim.x)");
            //if the right side's name is y
            if (name == "fetch_builtin_y")
                //replace the "blockDecl.y" with "(int) fmod((float) __SMC_chunkID, (float) __SMC_orgGridDim.y)"
                Rewrite.ReplaceText(Var->getLocEnd().getLocWithOffset(-9), 11, "(int) fmod((float) __SMC_chunkID, (float) __SMC_orgGridDim.y)");
        }

        //get the node "grid"
        if (const VarDecl *Var = Result.Nodes.getNodeAs("grid")) {
            //get the sourceLocation of the call of function "grid(...)"
            SourceLocation sr = Var->getLocStart().getLocWithOffset(4);
            //replace the "grid(...)" with "fmod((float) __SMC_chunkID, (float) __SMC_orgGridDim.x)"
            Rewrite.ReplaceText(sr, Var->getFullNameAsString().length(), "__SMC_orgGridDim");
        }
    }

private:
    Rewriter &rewriter;
};
```

Figure 8: The implementation of VariableDeclarationHandler

For the call of CUDA kernel function Handler:

```
class CudaCallPrinter : public MatchFinder::MatchCallback {
public:
    CudaCallPrinter(Rewriter &rewriter) : Rewriter(rewriter) {}

    virtual void run(const MatchFinder::MatchResult &result) {
        //get the node "CudaCall"
        if (const Stmt *S = Result.Nodes.getNodeAs("CudaCall")) {
            //Insert the __SMC_init() right before the call of the GPU kernel function
            Rewrite.InsertText(S->getLocStart(), " __SMC_init()", true, true);
            //if this function call is the call of GPU kernel function
            if (const CUDAKernelCallExpr *kernelExpr = dyn_cast<CUDAKernelCallExpr>(S)) {
                //get the location of this call's arguments
                const Stmt *arg = kernelExpr->getArgumentExprs()-1;
                //append the following arguments to the end of the GPU kernel function call
                Rewrite.InsertText(arg->getLocEnd().getLocWithOffset(1), " __SMC_orgGridDim __SMC_workersNeeded __SMC_workerCount __SMC_newChunkSeq __SMC_seqEnds", true, true);
            }
        }
    }

private:
    Rewriter &rewriter;
};
```

Figure 9: The implementation of the call of CUDA kernel function

## 6 The remaining issues and possible solutions

During the period of the whole project, there is a remaining issue left.

For the matcher of the call of CUDA kernel function, it implements the matcher for finding the position of call of CUDA kernel function.

```
// The matcher of the call of Cuda kernel function
Matcher.addMatcher(
    cudaKernelCallExpr().bind("cudaCall"), &cudaPrinter);
```

Figure 10: The matcher of the call of CUDA kernel function

```
class CudaPrinter : public MatchFinder::MatchCallback {
public:
    CudaPrinter(Rewriter &rewriter) : Rewriter(rewriter) {}

    virtual void run(const MatchFinder::MatchResult &result) {
        //get the node "CudaCall"
        if (const Stmt *S = Result.Nodes.getNodeAs("CudaCall")) {
            //Insert the __SMC_init() right before the call of the GPU kernel function
            Rewrite.InsertText(S->getLocStart(), " __SMC_init()", true, true);
            //if this function call is the call of GPU kernel function
            if (const CUDAKernelCallExpr *kernelExpr = dyn_cast<CUDAKernelCallExpr>(S)) {
                //get the location of this call's arguments
                const Stmt *arg = kernelExpr->getArgumentExprs()-1;
                //append the following arguments to the end of the GPU kernel function call
                Rewrite.InsertText(arg->getLocEnd().getLocWithOffset(1), " __SMC_orgGridDim __SMC_workersNeeded __SMC_workerCount __SMC_newChunkSeq __SMC_seqEnds", true, true);
            }
        }
    }

private:
    Rewriter &rewriter;
};
```

Figure 11: The implementation of the call of CUDA kernel function

According to these figures shown above, we can see the detailed implementation of the call of CUDA kernel function.

The CUDA code need to be transformed is showed as following:

```
// Performs varmup operation using matrixMul CUDA kernel
if (block_size == 16)
{
    matrixMulCUDA<16><<< grid, threads >>>(d_C, d_A, d_B, dimsA.x, dimsB.x);
}
else
{
    matrixMulCUDA<32><<< grid, threads >>>(d_C, d_A, d_B, dimsA.x, dimsB.x);
}
```

Figure 11: The code need to be transformed

With the requirements, the after-transform code will add the “*\_\_SMC\_init()*” above the call of CUDA kernel function, which means the *matrixMulCUDA<16>* and *matrixMulCUDA<32>*, and what's more? I also need to add the following arguments to the end of the GPU kernel function call:

```
__SMC_orgGridDim, __SMC_workersNeeded,
__SMC_workerCount, __SMC_newChunkSeq, __SMC_seqEnds
```

And what we need to get is as the following:

```
// The changes a line added
__SMC_init();
// Performs varmup operation using matrixMul CUDA kernel
if (block_size == 16)
{
    // __SMC_init() added four parameters to the call at the end
    matrixMulCUDA<16><<< grid, threads >>>(d_C, d_A, d_B, dimsA.x, dimsB.x, __SMC_orgGridDim, __SMC_workersNeeded, __SMC_workerCount, __SMC_newChunkSeq, __SMC_seqEnds);
}
else
{
    // __SMC_init() added four parameters to the call at the end
    matrixMulCUDA<32><<< grid, threads >>>(d_C, d_A, d_B, dimsA.x, dimsB.x, __SMC_orgGridDim, __SMC_workersNeeded, __SMC_workerCount, __SMC_newChunkSeq, __SMC_seqEnds);
}
```

Figure 12: The expected result

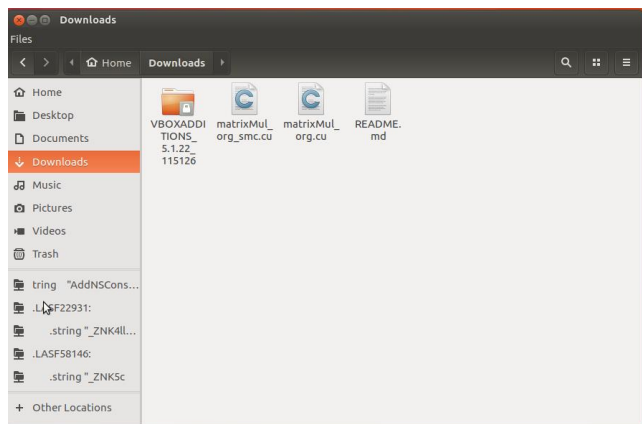
When I run my sm-transformation code, I will get the following result:

```

ubuntu@ubuntu:/I/vn/build-release$ bin/sn-centric /home/ubuntu/downloads/matrixMul_org.c
u -- -cuda-host-only
/home/ubuntu/Downloads/matrixMul_org.cu:36:10: fatal error: 'helper_functions.h' file not
found
#include <helper_functions.h>
      ^
1 error generated.
Error while processing /home/ubuntu/Downloads/matrixMul_org.cu.

```

**Figure 13:** *The screen shot when running*



**Figure 14:** *The screen shot after running*

From these screen shots, we can see that there is an error generated, it said the compiler can't find the helper\_functions.h file. But it still output the after-transform codes, which named "matrixMul\_org\_smc.cu". When I open this file, and I found this transformation is not successful, the output codes shown as following:

```
// Performs warmup operation using matrixMul CUDA kernel
if (block_size == 16)
{
    matrixMulCUDA<16><<< grid, threads >>>(d_C, d_A, d_B, dimsA.x, dimsB.x);
}
else
{
    matrixMulCUDA<32><<< grid, threads >>>(d_C, d_A, d_B, dimsA.x, dimsB.x);
}
```

**Figure 15:** *The screen shot of the result file*

So, to solve this problem, I find the *helper.h* files from the professor. And move these all .h files to the /usr/include folder:

```
ubuntu@ubuntu:~/Downloads$ sudo cp helper_cuda.h /usr/include/
ubuntu@ubuntu:~/Downloads$ sudo cp helper_functions.h /usr/include/
ubuntu@ubuntu:~/Downloads$
```

**Figure 16:** *The screen shot of move files*

After that, re-running the codes, as the following figures show, there is no error happened.

```
ubuntu@ubuntu:~/llvm/build-release$ bin/sm-centric /home/ubuntu/Downloads/matrixMul_org
u -- --cuda-host-only
ubuntu@ubuntu:~/llvm/build-release$
```

**Figure 17:** *The screen shot of re-running codes*

And we open the newly output file: matrixMul\_org\_smc.cu, check the output file:

[illegible]

**Figure 18:** *The result file*

From this output file, we can see this transformation is successful.

## 7 Results

I use a CUDA file named *matrixAdd\_org.cu* for test case, and the original cuda code is showed as following:

```

149 // @end_wasm
150 wasmExportFinalResult(final_w, final_h, final_h, int w, int h, final__BMC_explicitIn, final__BMC_explicitOut, final__BMC_explicitInn, final__BMC_explicitOutn, final__BMC_explicitInn, final__BMC_explicitOutn)
151 // @end_wasm
152
153 // @end_Bugfix
154 int bx = blockDim.x;
155 int by = blockDim.y;
156
157 // Shared local memory
158 int val = threadIdx.x;
159 int tid = threadIdx.y;
160
161 // Shared global memory
162 int tv = threadIdx.x;
163 int ty = threadIdx.y;
164 int gidin = ty*w;
165
166 // @end_Bugfix
167
168 // Setup execution parameters
169 dim3 threadsDim(grid, blockDim.x, blockDim.y);
170 dim3 gridDim / threads.x, h / threads.y;
171
172 // Create and start timer
173 printf("Executing result using CUDA Kernel...W")
174
175 // Launch and start timer
176 printf("Executing result using CUDA Kernel...W")
177
178 // @end_wasm
179
180 // @end_Bugfix
181
182 // @end_Bugfix
183
184 // @end_Bugfix
185
186 // @end_Bugfix
187
188 // @end_Bugfix
189
190 // @end_Bugfix
191
192 // @end_Bugfix
193
194 // @end_Bugfix
195
196 // @end_Bugfix
197
198 // @end_Bugfix
199
200 // @end_Bugfix
201
202 // @end_Bugfix
203
204 // @end_Bugfix
205
206 // @end_Bugfix
207
208 // @end_Bugfix
209
210 // @end_Bugfix
211
212 // @end_Bugfix
213
214 // @end_Bugfix
215
216 // @end_Bugfix
217
218 // @end_Bugfix
219
220 // @end_Bugfix
221
222 // @end_Bugfix
223
224 // @end_Bugfix
225
226 // @end_Bugfix
227
228 // @end_Bugfix
229
230 // @end_Bugfix
231
232 // @end_Bugfix
233
234 // @end_Bugfix
235
236 // @end_Bugfix
237
238 // @end_Bugfix
239
240 // @end_Bugfix
241
242 // @end_Bugfix
243
244 // @end_Bugfix
245
246 // @end_Bugfix
247
248 // @end_Bugfix
249
250 // @end_Bugfix
251
252 // @end_Bugfix
253
254 // @end_Bugfix
255
256 // @end_Bugfix
257
258 // @end_Bugfix
259
260 // @end_Bugfix
261
262 // @end_Bugfix
263
264 // @end_Bugfix
265
266 // @end_Bugfix
267
268 // @end_Bugfix
269
270 // @end_Bugfix
271
272 // @end_Bugfix
273
274 // @end_Bugfix
275
276 // @end_Bugfix
277
278 // @end_Bugfix
279
280 // @end_Bugfix
281
282 // @end_Bugfix
283
284 // @end_Bugfix
285
286 // @end_Bugfix
287
288 // @end_Bugfix
289
290 // @end_Bugfix
291
292 // @end_Bugfix
293
294 // @end_Bugfix
295
296 // @end_Bugfix
297
298 // @end_Bugfix
299
300 // @end_Bugfix
301
302 // @end_Bugfix
303
304 // @end_Bugfix
305
306 // @end_Bugfix
307
308 // @end_Bugfix
309
310 // @end_Bugfix
311
312 // @end_Bugfix
313
314 // @end_Bugfix
315
316 // @end_Bugfix
317
318 // @end_Bugfix
319
320 // @end_Bugfix
321
322 // @end_Bugfix
323
324 // @end_Bugfix
325
326 // @end_Bugfix
327
328 // @end_Bugfix
329
330 // @end_Bugfix
331
332 // @end_Bugfix
333
334 // @end_Bugfix
335
336 // @end_Bugfix
337
338 // @end_Bugfix
339
340 // @end_Bugfix
341
342 // @end_Bugfix
343
344 // @end_Bugfix
345
346 // @end_Bugfix
347
348 // @end_Bugfix
349
350 // @end_Bugfix
351
352 // @end_Bugfix
353
354 // @end_Bugfix
355
356 // @end_Bugfix
357
358 // @end_Bugfix
359
360 // @end_Bugfix
361
362 // @end_Bugfix
363
364 // @end_Bugfix
365
366 // @end_Bugfix
367
368 // @end_Bugfix
369
370 // @end_Bugfix
371
372 // @end_Bugfix
373
374 // @end_Bugfix
375
376 // @end_Bugfix
377
378 // @end_Bugfix
379
380 // @end_Bugfix
381
382 // @end_Bugfix
383
384 // @end_Bugfix
385
386 // @end_Bugfix
387
388 // @end_Bugfix
389
390 // @end_Bugfix
391
392 // @end_Bugfix
393
394 // @end_Bugfix
395
396 // @end_Bugfix
397
398 // @end_Bugfix
399
400 // @end_Bugfix
401
402 // @end_Bugfix
403
404 // @end_Bugfix
405
406 // @end_Bugfix
407
408 // @end_Bugfix
409
410 // @end_Bugfix
411
412 // @end_Bugfix
413
414 // @end_Bugfix
415
416 // @end_Bugfix
417
418 // @end_Bugfix
419
420 // @end_Bugfix
421
422 // @end_Bugfix
423
424 // @end_Bugfix
425
426 // @end_Bugfix
427
428 // @end_Bugfix
429
430 // @end_Bugfix
431
432 // @end_Bugfix
433
434 // @end_Bugfix
435
436 // @end_Bugfix
437
438 // @end_Bugfix
439
440 // @end_Bugfix
441
442 // @end_Bugfix
443
444 // @end_Bugfix
445
446 // @end_Bugfix
447
448 // @end_Bugfix
449
450 // @end_Bugfix
451
452 // @end_Bugfix
453
454 // @end_Bugfix
455
456 // @end_Bugfix
457
458 // @end_Bugfix
459
460 // @end_Bugfix
461
462 // @end_Bugfix
463
464 // @end_Bugfix
465
466 // @end_Bugfix
467
468 // @end_Bugfix
469
470 // @end_Bugfix
471
472 // @end_Bugfix
473
474 // @end_Bugfix
475
476 // @end_Bugfix
477
478 // @end_Bugfix
479
480 // @end_Bugfix
481
482 // @end_Bugfix
483
484 // @end_Bugfix
485
486 // @end_Bugfix
487
488 // @end_Bugfix
489
490 // @end_Bugfix
491
492 // @end_Bugfix
493
494 // @end_Bugfix
495
496 // @end_Bugfix
497
498 // @end_Bugfix
499
500 // @end_Bugfix
501
502 // @end_Bugfix
503
504 // @end_Bugfix
505
506 // @end_Bugfix
507
508 // @end_Bugfix
509
509

```

**Figure 19:** *The result file*

According to the output file, we can see the all transformation-needed part is successfully transformed. That means my transformation codes is useful and can run correctly.

## References

- [1] NVIDIA CUDA. <http://www.nvidia.com/cuda>.
- [2] J. T. Adriaens, K. Compton, N. S. Kim, and M. J. Schulte. The Case for GPGPU Spatial Multitasking. In Proceedings of HPCA, 2012.
- [3] Wu B, Chen G, Li D, et al. Enabling and Exploiting Flexible Task Assignment on GPU through SM-Centric Program Transformations[C]//ACM on International Conference on Supercomputing. ACM, 2015:119-130.
- [4] Wu B, Chen G, Li D, et al. SM-centric transformation: Circumventing hardware restrictions for flexible GPU scheduling[C]// International Conference on Parallel Architecture and Compilation Techniques. IEEE, 2017:497-498
- [5] The tutorial of Clang and LLVM, <https://clang.llvm.org/docs/LibASTMatchersTutorial.html>
- [6] Carlos Rea˜no, A CUDA-to-rCUDA Converter, Master's Thesis in Software Engineering, Formal Methods and Information Systems Department of Information System and Computation.

- [7] T. Aila and S. Laine. Understanding the efficiency of ray traversal on gpus. In Proceedings of the Conference on High Performance Graphics 2009, HPG '09, 2009.
- [8] G. Chen, B. Wu, D. Li, and X. Shen. Purple: An extensible optimizer for portable data placement on gpu. In Proceedings of MICRO, 2014.
- [9] L. Chen, O. Villa, S. Krishnamoorthy, and G. Gao. Dynamic load balancing on single-and multi-gpu systems. In IPDPS, 2010.
- [10] R. Nasre, M. Burtcher, and K. Pingali. Morph algorithms on gpus. In PPOPP, 2013.
- [11] M. E. Belviranli, L. N. Bhuyan, and R. Gupta. A dynamic self-scheduling scheme for heterogeneous multiprocessor architectures. ACM Trans. Archit. Code Optim., 9(4), Jan. 2013.