

# Accelerating DynEarthSol3D on Tightly Coupled CPU-GPU Heterogeneous Processors

Tuan Ta<sup>a,\*</sup>, Kyoshin Choo<sup>a,\*</sup>, Eh Tan<sup>b,\*</sup>, Byunghyun Jang<sup>a,\*\*</sup>, Eunseo Choi<sup>c,\*\*</sup>

<sup>a</sup>*Department of Computer and Information Science, The University of Mississippi, University, MS, USA.*

<sup>b</sup>*Institute of Earth Sciences, Academia Sinica, Taipei, Taiwan.*

<sup>c</sup>*Center for Earthquake Research and Information, The University of Memphis, Memphis, TN, USA.*

---

## Abstract

DynEarthSol3D (Dynamic Earth Solver in three dimensions) is a flexible, open-source finite element solver that models the momentum balance and the heat transfer of elasto-viscoplastic material in the Lagrangian form using unstructured meshes. It provides a platform for the study of the long-term deformation of earth's lithosphere and various problems in civil and geotechnical engineering. However, the continuous computation and update of a very large mesh poses an intolerably high computational burden to developers and users in practice. For example, simulating a 2-million-element mesh for 1,000 time units takes 1.4 hours on high-end desktop CPU. In this paper, we explore tightly coupled CPU-GPU heterogeneous processors to address the computing concern by leveraging their new features and developing hardware architecture aware optimizations. Our proposed key optimization techniques are three-fold: memory access pattern improvement, data transfer elimination and kernel launch overhead minimization. Experimental results show that our proposed implementation on a tightly coupled heterogeneous processor outperforms all other alternatives including traditional discrete GPU, quad-core CPU using OpenMP, and serial implementations by 67%, 50%, and 154% respectively even though the embedded GPU has significantly less number of cores than high-end discrete GPU.

*Keywords:* Computational Tectonic Modeling, Long-term Lithospheric Deformation, Heterogeneous Computing, GPGPU, Parallel Computing

---

## 1. Introduction

2 The combination of an explicit finite element method, the Lagrangian description of  
3 motion, and the elasto-visco-plastic material model has been implemented in a family of

---

\*Principal corresponding author

\*\*Corresponding author

*Email addresses:* [tqta@go.olemiss.edu](mailto:tqta@go.olemiss.edu) (Tuan Ta), [kchoo@go.olemiss.edu](mailto:kchoo@go.olemiss.edu) (Kyoshin Choo),  
[tan2@earth.sinica.edu.tw](mailto:tan2@earth.sinica.edu.tw) (Eh Tan), [bjang@cs.olemiss.edu](mailto:bjang@cs.olemiss.edu) (Byunghyun Jang), [echoi2@memphis.edu](mailto:echoi2@memphis.edu)  
(Eunseo Choi)

4 codes following the Fast Lagrangian Analysis of Continua (FLAC) algorithm [10]. These  
5 specific implementations of the generic FLAC algorithm have a track record of applications  
6 that demonstrate the method’s aptitude for Long-term Tectonic Modeling (LTM) (e.g.,  
7 [4, 6, 7, 12, 14, 16, 19, 20]). The original FLAC requires a structured quadrilateral mesh  
8 which severely limits the meshing flexibility, one of the major advantages of finite element  
9 method. Flexibility in meshing is all the more important for LTM in which strain localization  
10 needs to be adequately captured by a locally refining a mesh, which is challenging for a  
11 structured mesh. Additionally, each quadrilateral is decomposed into two sets of overlapping  
12 linear triangles that guarantee a symmetrical response to loading but leads to redundant  
13 computations. On the other hand, FLAC uses an explicit scheme for the time integration of  
14 the momentum equation in the dynamic form as well as for the constitutive update, making  
15 it relatively easy to implement complicated constitutive models.

16 By critically evaluating the strengths and weaknesses of the FLAC algorithm, Choi et  
17 al.[8] created a new code, DynEarthSol2D, and Tan et al.[25] further extended it to three  
18 dimensions, DynEarthSol3D. DynEarthSol3D (**D**ynamic **E**arth **S**olver in three dimensions)  
19 is a robust, flexible, open source finite element code for modeling non-linear responses of  
20 continuous media and thus suitable for LTM. DynEarthSol3D written in standard C++ is  
21 multi-threaded and freely distributed through a public repository<sup>1</sup> under the terms of the  
22 MIT/X Windows System license. DynEarthSol3D inherits desirable features of FLAC such  
23 as explicit schemes while modernizing it at the same time. The most notable improvement  
24 is the removal of the restrictions on meshing. As a result, one can solve problems on  
25 unstructured triangular or tetrahedral meshes while keeping the simple explicit constitutive  
26 update that made FLAC attractive in LTM. The use of unstructured mesh enables adaptive  
27 mesh refinement in regions of highly localized deformation and discretization of domain  
28 geometries that are challenging to discretize into a structured mesh. However, sequential  
29 mesh computations and updates are very compute-intensive, resulting in poor performance  
30 in practice. To process a large mesh composed of 2 million elements, it takes 1.4 hours for  
31 serial implementation to finish 1,000 time steps on a high-end CPU. This huge amount of  
32 running time places a limitation on both mesh size and resolution.

33 GPUs (Graphics Processing Units) have been the platform of choice for compute- and  
34 data-intensive applications in many computing domains in recent years. GPU-powered com-  
35 puting provides a number of unique benefits that could not be found in any traditional  
36 parallel machines such as supercomputers and workstations. This revolutionary computing  
37 paradigm of offloading and accelerating compute- and data-intensive portion of applications  
38 on GPUs is termed as GPGPU (General-Purpose Computation on GPUs) or GPU comput-  
39 ing. When well optimized for target GPU hardware architecture, application performance  
40 can be boosted by up to several orders of magnitude.

41 GPGPU platforms are typically powered by high-end discrete GPUs (i.e., separate graph-  
42 ics card connected through PCI express). While this type of hardware configuration provides  
43 best GPU processing power, data transfer overhead associated with physical separation of  
44 GPU device and memory from host CPU diminishes the performance gain obtained by GPU

---

<sup>1</sup><http://bitbucket.org/tan2/dyearthsol3d>

45 acceleration. Deteriorated by kernel launch time, such overheads can be a deal breaker. If  
46 an application has multiple sections of CPU and GPU computations that are interleaved  
47 and data-dependent of each other, repetitive data transfers between host and device are  
48 inevitable, and overall application performance is limited by the overhead associated with  
49 these transfers. This problem that is not uncommon in many scientific and engineering  
50 applications hinders the adoption of GPGPU.

51 Recent trend in microprocessor industry is that CPU and GPU are fabricated on a  
52 single die sharing a memory system at a cache level [3, 13]. Such tightly coupled CPU-GPU  
53 heterogeneous processors provide solutions to several limitations of traditional discrete GPUs  
54 such as aforementioned data transfer overhead [9, 24], limited GPU device memory (i.e.,  
55 GDDR) size [24], and disjoint memory address space. CPU and GPU on tightly coupled  
56 heterogeneous processors share the same data in a unified physical memory (data transfer  
57 is no longer needed) and also the unified system memory is typically a lot larger (e.g. 32  
58 GB) than discrete GPU device memory (e.g. 4GB).

59 In this paper, we present the acceleration of DynEarthSol3D on tightly coupled CPU-  
60 GPU heterogeneous processors, leveraging a new unified memory architecture to eliminate  
61 data transfer overhead. We also revisit and address classical GPGPU challenges such as  
62 inefficient memory access patterns and frequent kernel launch overhead. The contributions  
63 of our paper are summarized below.

- 64 • We demonstrate that tightly coupled CPU-GPU heterogeneous processors outperform  
65 discrete GPUs by eliminating data transfer overhead, a serious performance bottleneck  
66 of DynEarthSol3D on conventional discrete GPUs. This result is encouraging because  
67 the computing power of embedded GPU on heterogeneous processor is less than one  
68 fourth of that of discrete GPU that we have tested.
- 69 • We propose to improve GPU memory performance by changing memory access pat-  
70 terns through data transformation. By restructuring the mesh, high latency random  
71 memory access patterns of DynEarthSol3D turn to regular patterns that GPU hard-  
72 ware can handle much more efficiently. As a result, it boosts the performance of GPU  
73 kernel execution significantly.
- 74 • We propose to merge kernels whenever possible to minimize kernel launch overhead.  
75 Intensive data flow and dependency analysis are conducted to identify all kernels that  
76 can be merged without causing any correctness issue. As kernels are called repeatedly  
77 throughout program execution, total kernel launch overhead is significantly reduced.
- 78 • We conduct thorough performance analysis and comparison with other available alter-  
79 natives: discrete GPU, multi-core CPU using OpenMP, and a serial implementation  
80 as baselines.

81 The rest of the paper is structured as follows. Section 2 describes the computations  
82 of DynEarthSol3D and existing problems in its serial implementation on CPU. Section 3  
83 provides the background on GPGPU with explanation of both traditional discrete GPUs and

84 tightly coupled heterogeneous processors. In Section 4, we present our implementation of  
 85 DynEarthSol3D while focusing on three key optimization techniques. Lastly, Section 5 shows  
 86 and discusses our experimental results through detailed evaluation of each optimization  
 87 technique.

## 88 2. Computational Flow of DynEarthSol3D

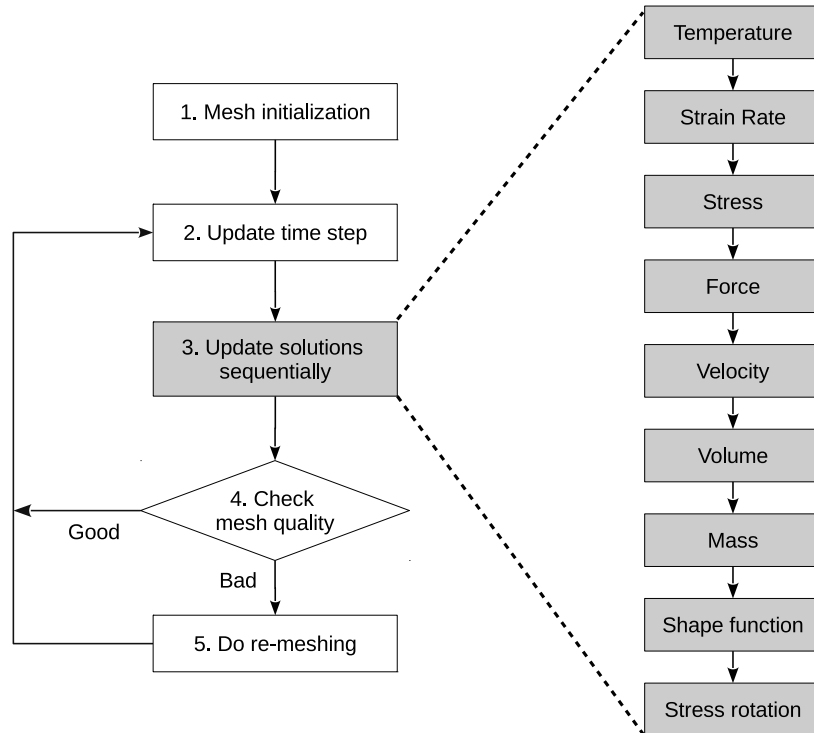


Figure 1: The computational flow of DynEarthSol3D.

89 Figure 1 visualizes the computational flow of DynEarthSol3D. First, a mesh composed  
 90 of tetrahedral elements is created by an external mesh generator named *Tetgen* [22]. Each  
 91 element of the mesh consists of four nodes which function as interpolation points for un-  
 92 known variables such as velocity and temperature. The program runs through a predefined  
 93 number of time steps to reach a target model time, which is in LTM typically millions of  
 94 years to tens of millions of year. In each time step, the temperature field is first updated  
 95 according to the heat diffusion equation. The updated temperature may be used for com-  
 96 puting temperature-dependent constitutive models. Next, based on the current coordinates  
 97 of nodes and the velocity field, element volume and strain rates are computed. Strain, strain  
 98 rates and temperature are used to update stress according to an assumed constitutive model.  
 99 Net acceleration of each node is computed as the net force divided by inertial mass, and  
 100 the next force is the sum of external body force and internal force arising from the updated  
 101 stresses. The net acceleration is time-integrated to velocity and displacement. Once the dis-  
 102 placement is updated, the coordinates of the nodes are updated. At this stage, the program

103 checks if accumulated deformation has distorted the mesh too severely. If so, Tetgen again is  
104 used to regenerate a mesh, each element of which satisfies a certain quality criterion. During  
105 this remeshing process, new nodes might be inserted into the mesh, or the mesh topology  
106 might change through edge-flipping. This type of remeshing has been proposed as a way of  
107 solving large deformation problems in the Lagrangian framework [5]. After the new mesh is  
108 created, the boundary conditions, derivatives of shape function, and mass matrix have to be  
109 re-calculated. Then, the next time step is initiated unless the current one is the last step.

110 DynEarthSol3D has options to run either serially or in parallel on multi-core CPU with  
111 OpenMP. In the serial version, the elements of the mesh, and the nodes associated with each  
112 element, are processed sequentially. In the OpenMP version, when running with  $P$  number  
113 of threads, in order to prevent race condition among threads, mesh elements are grouped  
114 into  $2P$  sets corresponding to roughly uniform  $2P$  intervals in the  $x$  coordinates of their  
115 barycenters. Elements in set  $i$  are guaranteed to have no common nodes with elements in  
116 set  $i + 2$ . To process all elements, elements in set  $0, 2, \dots, 2(P - 1)$  are first processed in  
117 parallel, with each set covered by a thread. Then, after all threads finish processing, elements  
118 in set  $1, 3, \dots, 2P - 1$  are processed in the same labor division. A good multi-core scaling  
119 can be safely achieved this way. When we evaluate the performance of our implementation  
120 in this paper, both serial and OpenMP-based implementations are used as baselines.

### 121 3. General-Purpose Computing on GPUs (GPGPU)

122 GPUs provide an unprecedented level of computing power per dollar and energy by  
123 running massive number of threads in a Single Instruction Multiple Thread (SIMT) fashion.  
124 It keeps many ALU (arithmetic logic unit) cores busy by hiding memory latency through  
125 zero-overhead thread switching. At any given clock cycle, multiple groups of threads (i.e.  
126 multiple of 32 or 64 threads) run in a Single Instruction Multiple Data (SIMD) fashion.  
127 When well optimized, data- and compute-intensive applications can be easily accelerated  
128 by several orders of magnitude. GPGPU is programmed using OpenCL [17] or CUDA [18]  
129 languages in which data- and compute-intensive portions of program are offloaded onto GPU  
130 device. The offloaded function-like code to be executed on the GPU device is called *kernel*.  
131 Host and device kernel codes can be executed either asynchronously or synchronously.

#### 132 3.1. Limitations of Current GPGPU Paradigm

133 Although numerous applications have been successfully accelerated using GPUs with re-  
134 markable speedups, there are many other algorithms and applications that do not benefit  
135 from current GPGPU computing paradigm. This is because GPU hardware and software  
136 (i.e. programming model) are very different from conventional parallel platforms as they  
137 are evolved by the demand for real-time 3D graphics rendering. We summarize the big lim-  
138 itations of today's GPGPU computing that we also found present in our target application,  
139 DynEarthSol3D:

- 140 • **Data transfer overhead:** In conventional GPGPU settings, discrete GPU is physi-  
141 cally connected through PCI-Express and has separate physical memory (see Figure 2).

142 Data to be used by kernel program must be copied to device memory before kernel ex-  
 143 ecution. If an application consists of multiple sections of CPU and GPU computations  
 144 that are interleaved and data-dependent on each other, frequent data transfer between  
 145 host and device is necessary. Therefore, overall application performance is limited by  
 146 the overhead associated with the slow data transfer.

- 147 • **Kernel launch overhead:** The host CPU communicates with GPU through device  
 148 driver calls and each command including kernel invocation involves overhead. When  
 149 a large number of kernel calls are performed throughout the program execution, the  
 150 kernel launch overhead associated with device driver calls can be added up to significant  
 151 portion of overall performance. Such overhead can be a serious problem especially when  
 152 the kernel execution time is small. Therefore, launching multiple small kernels should  
 153 be avoided whenever possible to reduce the overhead.

- 154 • **Irregular memory access:** GPU hardware architecture is designed for maximizing  
 155 throughput for a group of threads rather than minimizing the latency of an individual  
 156 thread. It implies that memory subsystem becomes very inefficient when threads issue  
 157 memory requests with irregular access patterns [15]. Altering memory access patterns  
 158 toward more hardware friendly ones is the most important yet challenging optimization  
 159 to improve kernel performance. This is typically done by transforming data layout or  
 160 changing computations in source code.

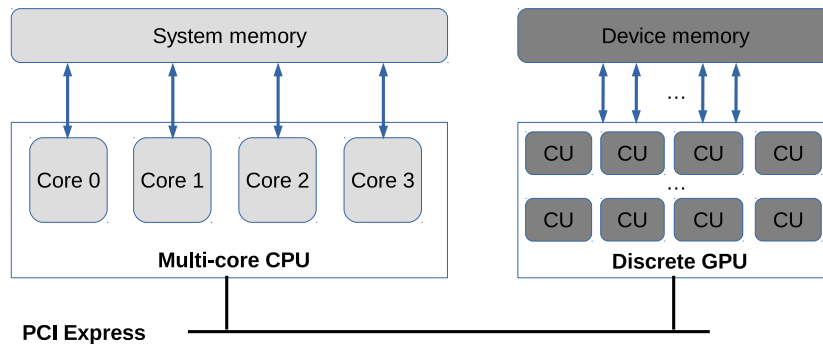


Figure 2: The system diagram of conventional discrete GPGPU platform.

### 161 3.2. Tightly Coupled CPU-GPU Heterogeneous Processors

162 Recent trend in microprocessor industry is to merge CPU and GPU on a single die  
 163 [2, 21]. This is a natural choice in microprocessor industry as it offers numerous benefits  
 164 such as fast interconnection and fabrication cost reduction. Recently, the processor industry  
 165 and academic research community have formed a non-profit consortium, called HSA (Het-  
 166 erogeneous System Architecture) foundation [11] to define the standards of hardware and  
 167 software for next generations of heterogeneous computing. Such processors that couple CPU  
 168 and GPU at last level cache overcome some limitations of current GPGPU. Tightly coupled  
 169 heterogeneous processors provide the following benefits.

- 170 • **Fast and fine-grained data sharing between CPU and GPU:** Multi-core CPU  
171 and many-core GPU are tightly coupled at last cache level on a single die and share a  
172 single physical memory (see Figure 3). This architecture design eliminates CPU-GPU  
173 data transfer overhead by sharing the same data.
- 174 • **Large memory support for GPU acceleration:** Data oriented applications such  
175 as big data processing and compute-intensive scientific simulations require a large  
176 memory space to minimize inefficient data copy back and forth. In tightly coupled  
177 heterogeneous processors, GPU device shares system memory that is typically a lot  
178 larger (e.g. 32GB) than device memory in discrete GPUs (e.g. 4GB).
- 179 • **Cache coherence between CPU and GPU:** This new hardware feature will remove  
180 off-chip memory access traffic significantly and allow fast, fine-grained data sharing  
181 between CPU and GPU. Both devices are capable of addressing the same coherent  
182 block of memory, supporting popular application design patterns such as producer-  
183 consumer [23].

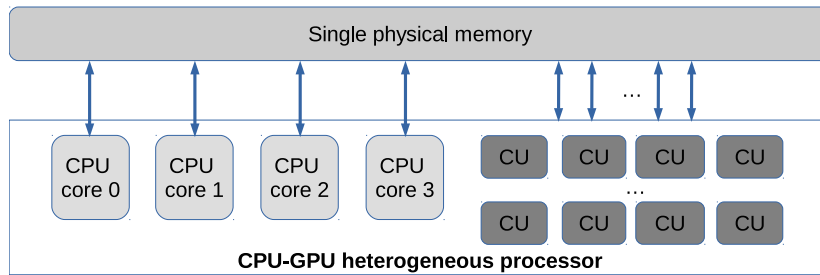


Figure 3: The system diagram of tightly coupled CPU-GPU heterogeneous processors.

#### 184 4. Implementation and Optimization

185 As shown in Figure 1, DynEarthSol3D sequentially computes and updates unknowns on  
186 the nodes of a mesh in each time step. To identify time-consuming parts of the program,  
187 we first profile and decompose the execution time of serial version at functional level<sup>2</sup>, and  
188 illustrate the result in Figure 4. This profiling and breakdown of execution time provides  
189 useful information that helps identify the candidate functions to be offloaded to GPU. Based  
190 on this information and through source code analysis, we have chosen following functions  
191 (or operations) to accelerate on GPU. They account for 88% of the total execution time  
192 in DynEarthSol3D. For the sake of readability of the paper, we use short names shown in  
193 paranthesis in Figure 4 in the rest of the paper.

---

<sup>2</sup>In serial DynEarthSol3D, computing and updating each property of mesh is written as a function. Note also that the remeshing operation is excluded from our analysis as it uses the external *Tetgen* library whose performance is not our focus in this work.

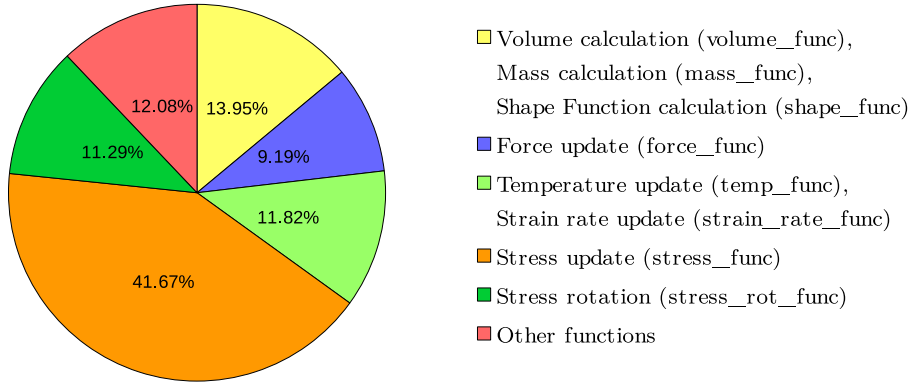


Figure 4: Performance breakdown of the serial version of DynEarthSol3D.

194 In the following subsections, we describe the general structure of our OpenCL imple-  
 195 mentation, followed by the detailed presentation of each optimization. Our optimizations  
 196 focus on 1) memory access patterns, 2) data transfer between CPU and GPU, and 3) kernel  
 197 launch overhead.

198 *4.1. The Structure of OpenCL Implementation*

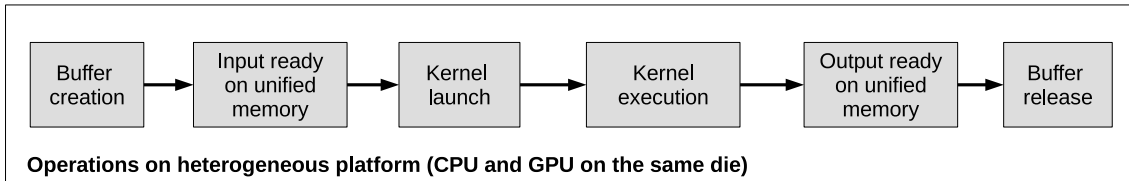


Figure 5: The structure of OpenCL implementation on CPU-GPU heterogeneous platform.

199 Initial GPU setup including device configuration, platform creation, and kernel building,  
 200 etc. is performed only once before the program starts updating solutions in multiple time  
 201 steps. The framework described in Figure 5 applies to all of our target functions. OpenCL  
 202 buffers are first created with appropriate flags that enable the *zero copy* feature on heteroge-  
 203 neous processor. These buffers reside on the unified memory that is accessible to both CPU  
 204 and GPU (This feature is detailed in Section 4.3). After buffers are created, kernel argu-  
 205 ments are set up, and then kernel is launched. Each mesh element is processed by a specific  
 206 thread through one-to-one mapping between work-item ID (in a n-dimensional thread index  
 207 space called *NDRange*) and element ID. Multiple work-items in the *NDRange* are grouped  
 208 into a work-group that is assigned to a compute unit on GPU. While each thread reads input  
 209 element from global memory, processes, and updates it, two different elements mapped to  
 210 two threads may share some nodes, which leads to a race condition when they update the  
 211 shared node at the same time. To guarantee that outputs are updated correctly, the race  
 212 condition must be handled through atomic operations. The execution on GPU continues



213 until all threads complete their works and are synchronized by the host to ensure that out-  
214 put data are complete and valid. Finally, buffers are released, and the program continues  
215 its remaining operations in the current time step or moves on to the next one.

#### 216 4.2. Memory Access Pattern Improvement

217 The performance of GPGPU is heavily dependent on the memory access behavior. This  
218 sensitivity is due to a combination of the underlying massively parallel processing execution  
219 model present on GPUs and the lack of architectural support to handle irregular memory  
220 access patterns. Hardware unfriendly memory accesses degrade performance significantly  
221 as it results the serialization of many expensive off-chip memory accesses. For linear and  
222 regular memory access patterns issued by a group of threads, the hardware coalesces them  
223 into a single (or fewer number of) memory transactions, which significantly reduce overall  
224 memory latency, consequently less thread stalls. Therefore, application performance can be  
225 significantly improved by minimizing irregularity of global memory access patterns.

226 In DynEarthSol3D, *Tetgen* program generates a mesh with a system of element and  
227 node numbers (IDs). Each tetrahedral element with its own ID is associated with four  
228 different nodes numbered in semi-random fashion. In our implementation, nodes are accessed  
229 sequentially by each thread. Therefore, the randomness of node IDs in an element results in  
230 irregular pattern of global memory accesses requested by a single thread which has to load  
231 and update node-related data locating randomly in global memory. Figure 6 illustrates a case  
232 where two adjacent elements may share three nodes (i.e. IDs 10, 30, 60) together. Figure 7a  
233 visualizes the randomness of the node system by representing each node ID corresponding  
234 to its element ID.

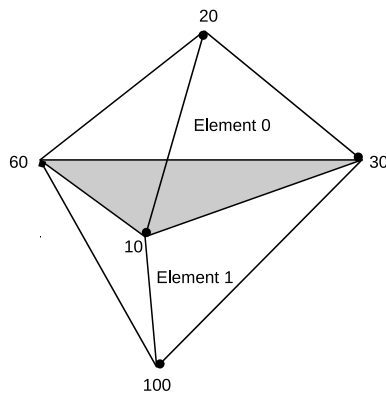


Figure 6: Two tetrahedral elements share three nodes.

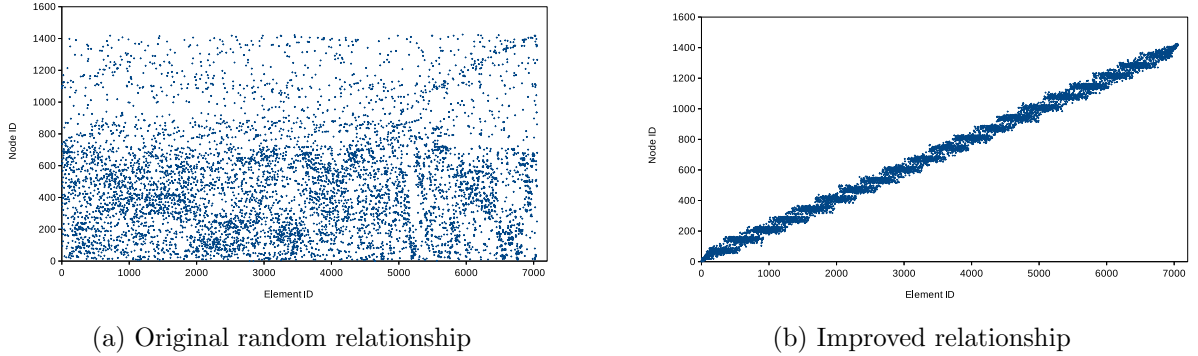


Figure 7: Relationship between node and element IDs. Each element ID is mapped to a single thread in GPU kernel.

235 To eliminate the randomness of node ID system, we renumber all nodes so that they are  
 236 ordered by their corresponding  $x$  coordinates and renumber all elements similarly by the  
 237  $x$  coordinates of their centers. As a result, node IDs within a single element and among  
 238 multiple adjacent elements are close together. Figure 7b illustrates the improved relationship  
 239 between node and element IDs. This improved pattern has a direct impact on memory  
 240 access patterns of the kernel. Cache hit rate significantly increases, and memory accesses  
 241 are coalesced. Therefore, overall memory latency during kernel execution is significantly  
 242 decreased.

243 *4.3. Data Transfer Elimination*

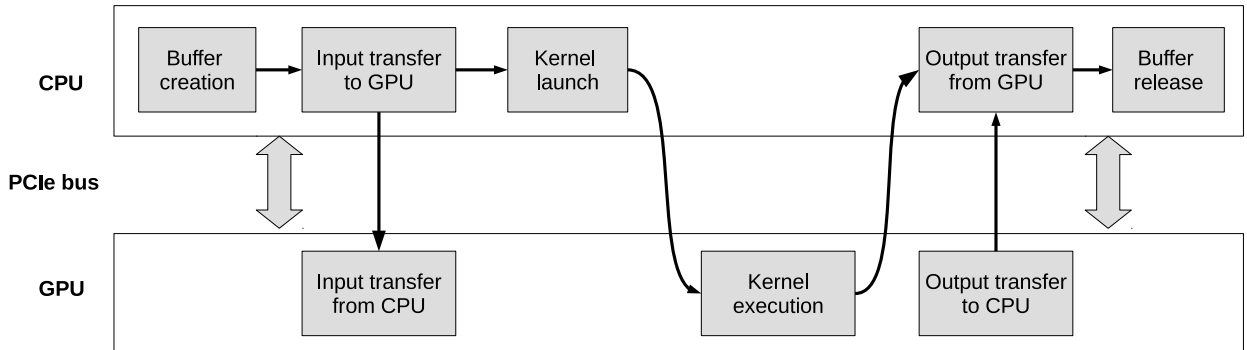


Figure 8: The structure of OpenCL implementation on discrete GPU platform.

244 Figure 8 shows the computational flow of OpenCL implementations on conventional  
 245 discrete GPU platform with respect to physical execution hardware. On discrete GPU  
 246 systems where CPU and GPU have separate memory subsystem, data copy between host and  
 247 device must be done via low speed PCI Express bus. Such data movement takes considerable  
 248 amount of time and can significantly offsets performance gains obtained by GPU kernel  
 249 acceleration. This data copy is a serious problem in DynEarthSol3D as large amount of

250 data has to be copied back and forth between host and device in each time step. Tightly  
 251 coupled CPU-GPU heterogeneous processors offer a solution to this bottleneck as CPU and  
 252 GPU share the same unified physical memory. Using a feature known as *zero copy*, data (or  
 253 buffer) can be accessed by two processors without copying. Zero copy is enabled by passing  
 254 one of following flags appropriately to *clCreateBuffer* OpenCL API function [1].

- 255 • **CL\_MEM\_ALLOC\_HOST\_PTR** Buffers created with this flag is “host-resident  
 256 zero copy memory object” that is directly accessed by host at host memory bandwidth  
 257 and visible to GPU device.
- 258 • **CL\_MEM\_USE\_HOST\_PTR** This flag is similar to CL\_MEM\_ALLOC\_HOST\_PTR.  
 259 However, instead of allocating a new memory space belonging to either host or device,  
 260 it uses a memory buffer that has been already allocated and currently being used by  
 261 host program.
- 262 • **CL\_MEM\_USE\_PERSISTENT\_MEM\_AMD** This flag is available only on  
 263 AMD platform. It tells host program to allocate a new “device-resident zero copy  
 264 memory object” that is directly accessed by GPU device and accessible in host code.

265 Because the first and third options allocate new empty memory spaces, the buffers need to  
 266 be filled with data before kernel execution on GPU. In our DynEarthSol3D implementation,  
 267 the second option is used to avoid such extra buffer setup. Figure 9 illustrates how both  
 268 host program running on CPU and kernel running on GPU access data in shared buffers  
 269 created with CL\_MEM\_USE\_HOST\_PTR flag.

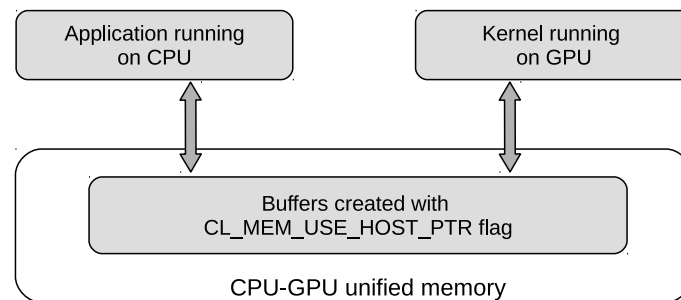


Figure 9: Memory buffer shared by CPU and GPU (a.k.a. zero copy) on tightly coupled heterogeneous processors.

270 4.4. Kernel Launch Overhead Minimization

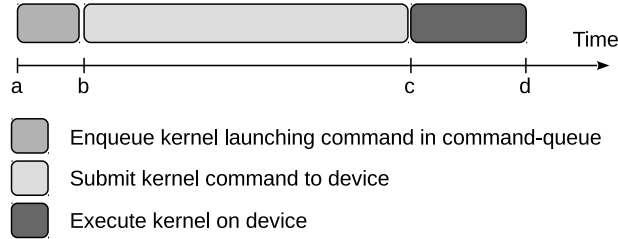


Figure 10: The breakdown of kernel execution process.

271 Executing a kernel on GPU device consists of three steps. A kernel launch command  
 272 is first enqueued to the command-queue associated with device, and then the command  
 273 is submitted to device before the kernel is finally executed. Queuing and submitting ker-  
 274 nel launch command are considered as overhead. According to our experiment on AMD  
 275 platforms, the command submission (second block in Figure 10) accounts for most of the  
 276 overhead. This overhead becomes significant when actual kernel execution time is relatively  
 277 short compared to kernel queuing and submission time as exemplified in Figure 10. In ad-  
 278 dition, the use of `CL_MEM_USE_HOST_PTR` flag results in a small runtime overhead as  
 279 the size of buffers used in kernel increases. These two kinds of overhead are repeatedly  
 280 accumulated in DynEarthSol3D as kernels are re-launched in each time step.

281 The only available solution to this overhead is to reduce the number of kernel launches  
 282 throughout the program execution. Toward that end, we merge multiple functions into  
 283 a single kernel, so there are less number of kernel to be launched. If data dependency  
 284 exists between two functions (meaning that the second function needs to use outputs of the  
 285 first one), they cannot be merged into a single kernel because GPGPU programming and  
 286 hardware execution model does not guarantee that the first function finishes its entire thread  
 287 execution before the second starts. Without such guarantee, the second function might use  
 288 old input data that has not been updated yet by the first one.

289 In our implementation, we first perform in-depth data dependency analysis of DynEarth-  
 290 Sol3D to identify possible combinations of functions with no data dependency. Based on  
 291 this analysis, we find two combinations. The first one combines *volume\_func*, *mass\_func* and  
 292 *shape\_func*. The other merges *temp\_func* and *strain\_rate\_func*. For simplicity, we call the  
 293 first and second merged kernels *intg\_kernel\_1* and *intg\_kernel\_2* respectively in the rest of  
 294 this paper.

295 **5. Experimental Results**

296 To evaluate the performance of our proposed OpenCL implementation on tightly cou-  
 297 pled CPU-GPU heterogeneous processors we compare its performance with both serial and  
 298 OpenMP-based implementations as baselines. We also analyze the impact of each proposed  
 299 optimization technique.

300 In all experiments, we use the same program configuration with varying sizes of mesh of  
301 elasto-plastic material. The program runs in 1,000 time steps, and its outputs are written  
302 into output files every 100 steps. Performance results are accumulated in each time step.  
303 Wall clock timer and OpenCL profiler are used to measure performance of host code and  
304 kernel respectively. We compare our OpenCL output results with serial version’s outputs to  
305 verify the correctness of our implementation.

306 We experiment our OpenCL-based implementation on AMD APU A10-7850K which is  
307 the latest heterogeneous processor as of this paper writing. This tightly coupled heteroge-  
308 neous processor consists of a quad-core CPU with maximum clock speed of 4.0 GHz and  
309 a Radeon R7 GPU with eight compute units running at 720 MHz on the same die. Our  
310 baseline versions (i.e., serial and OpenMP-based implementation) are tested on the quad-  
311 core CPU of the same APU for fair comparison. In evaluating the impact of data transfer  
312 elimination, we use a high-end discrete AMD GPU Radeon HD 7970 codenamed Tahiti. It  
313 has 32 compute units with maximum clock speed of 925 MHz. The operating system is  
314 64-bit Ubuntu Linux 12.04 and AMD APP SDK 2.9 (OpenCL 1.2) is used.

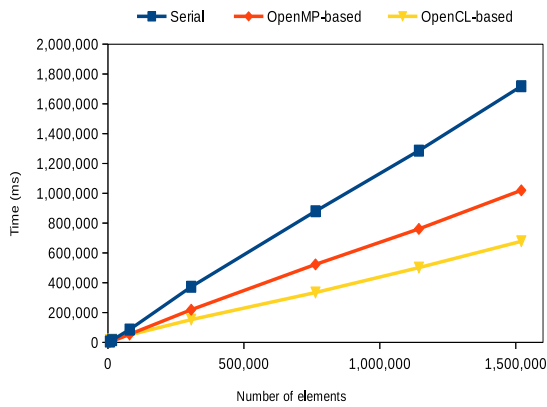
### 315 5.1. Overall Acceleration

316 In this section, we compare the performance of our OpenCL-based implementation with  
317 two baseline versions: serial and OpenMP-based implementations at both program- and  
318 function-levels. The results<sup>3</sup> are shown in Figure 11. We varied the number of mesh el-  
319 ements from 7 thousand to 1.5 million. Regarding integrated kernels *intg\_kernel\_1* and  
320 *intg\_kernel\_2*, we do comparisons with their corresponding component functions in serial and  
321 OpenMP-based implementations. Note that *intg\_kernel\_1* merges *volume\_func*, *mass\_func*  
322 and *shape\_func*, and *intg\_kernel\_2* merges *temp\_func* and *strain\_rate\_func*.

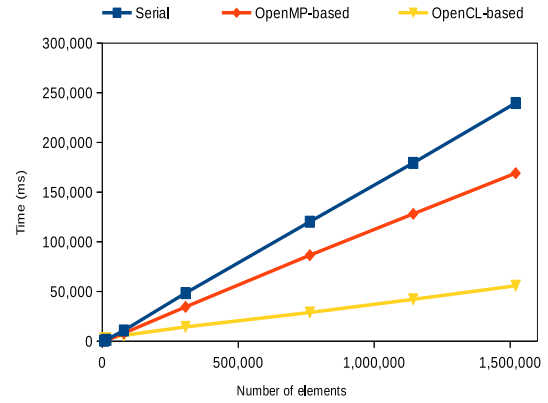
323 At program level, our OpenCL implementation optimized for tightly coupled hetero-  
324 geneous processor outperforms both serial and OpenMP-optimized versions by 154% and  
325 50% respectively for 1.5-million-element mesh. At function level, all target functions show a  
326 similar trend in performance. Especially, integrated kernel *intg\_kernel\_1* is 329% and 203%  
327 faster than its before-merged case in serial and OpenMP versions respectively. The impact  
328 of merging kernels is analyzed in more detail in later section.

---

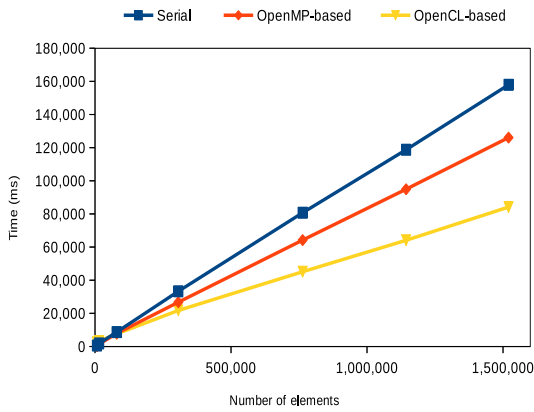
<sup>3</sup>For fair comparisons, experiments are done with the improved node ID system. A less random memory access pattern also improves the performance of serial and OpenMP-based versions due to better cache utilization on CPU.



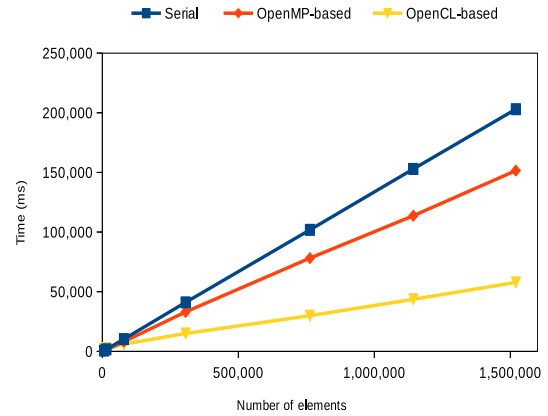
(a) Overall performance (at program level).



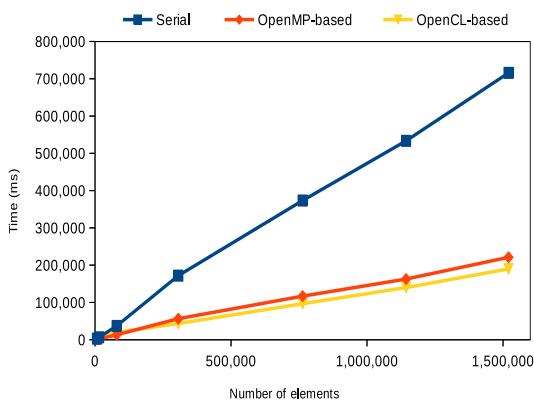
(b) Performance of *intg\_kernel\_1*.



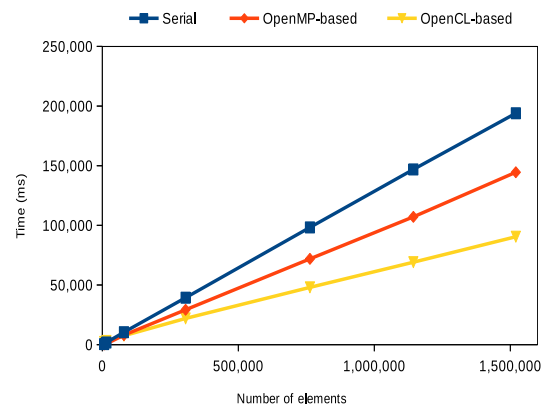
(c) Performance of *force\_func*.



(d) Performance of *intg\_kernel\_2*.



(e) Performance of *stress\_func*.



(f) Performance of *stress\_rot\_func*.

Figure 11: Performance comparison among different implementations.

329 An interesting observation from these comparisons is that performance gain from GPU  
 330 acceleration becomes more substantial as input size increases. If there are a small number  
 331 of threads issued (i.e. small input), GPU computing hardware resources are underutilized  
 332 and unable to compensate for the setup overhead of GPU hardware pipelines.

### 333 5.2. Impact of Memory Access Pattern Optimization

334 In this section, we analyze the impact of node ID system improvement on performance by  
 335 comparing two implementations with and without this improvement at function level. Only  
 336 kernel execution time measured by AMD profiler is concerned here as memory access pattern  
 337 does not affect other parts of our optimization. Figure 12 illustrates these performance  
 338 comparisons.

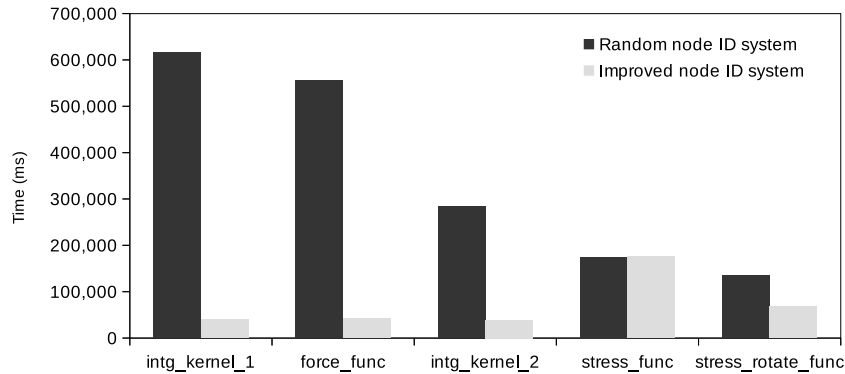


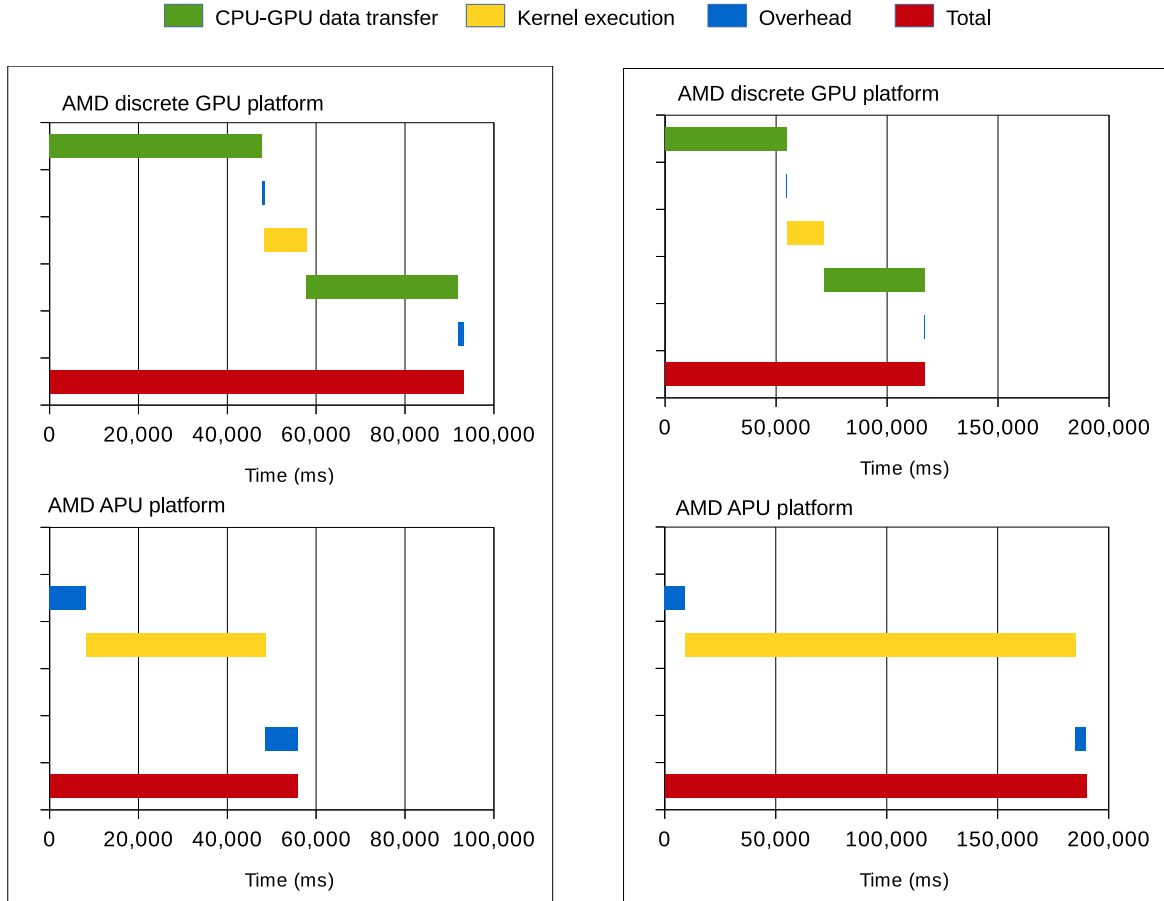
Figure 12: The impact of memory access pattern on kernel execution time.

339 Substantial improvement in kernel execution is achieved in functions that process node-  
 340 related mesh properties intensively (i.e., 15x, 13.4x, 7.2x, 2x speedups in *intg\_kernel\_1*,  
 341 *force\_func*, *intg\_kernel\_2*, and *stress\_rot\_func* respectively). The randomness of node IDs  
 342 does not affect performance of *stress\_func* because it deals with only element-related mesh  
 343 properties. The improved memory access patterns enable kernels to take advantage of spatial  
 344 locality within a thread and across threads in a work-group, which consequently yields better  
 345 utilization of GPU cache system. In addition, because multiple global memory requests can  
 346 be coalesced into fewer number of memory transactions, the improved node ID pattern  
 347 reduces both on-chip and off-chip memory traffic substantially and reduces overall memory  
 348 latency.

### 349 5.3. Impact of Data Transfer Elimination

350 In order to demonstrate the significant benefit of utilizing tightly coupled CPU-GPU  
 351 heterogeneous processors in terms of data transfer overhead, we present and compare the  
 352 execution time of two kernels: *intg\_kernel\_1* and *stress\_func* in Figure 13. Both functions  
 353 shown here compute and process large number of mesh properties that are associated with  
 354 a considerable amount of data. We test them on 1) high-end discrete GPU Radeon HD

355 7970 (codenamed Tahiti) with explicit data transfer by calling *clEnqueueWriteBuffer* and  
 356 *clEnqueueReadBuffer*, and 2) heterogeneous processor AMD APU A10-7850K with *zero copy*  
 357 feature enabled with respect to data transfer, kernel execution and overhead.



(a) *intg\_kernel\_1* function.

(b) *stress\_func* function.

Figure 13: The impact of data transfer elimination.

358 On discrete GPU, explicit data transfer between host and device memory accounts for  
 359 88% and 85% of total performance in *intg\_kernel\_1* and *stress\_func* respectively. The reason  
 360 for this extremely high cost of data communication on discrete platform is that all data are  
 361 transferred through PCI Express bus at slow speed. In contrast, there is no data transfer  
 362 on CPU-GPU heterogeneous platform. Regarding *intg\_kernel\_1* function, the AMD APU  
 363 outperforms the high-end discrete GPU (i.e. 67% faster) despite the fact that the Tahiti GPU  
 364 is provided with more powerful computing capability (more compute units and higher clock  
 365 speed than the AMD APU). However, in the case of *stress\_func* function, the elimination of  
 366 data transfer is not enough to compensate for the much less computing capability of AMD  
 367 APU. The reason is that compared to *intg\_kernel\_1* function, *stress\_func*'s kernel executes  
 368 a lot more arithmetic computations that the discrete GPU is capable of performing much



369 faster than the embedded GPU of heterogeneous processor is able to do.

#### 370 5.4. Impact of Kernel Overhead Minimization

371 According to our experiment, the kernel launch overhead accounts for 27% of *intg\_kernel\_1*'s  
 372 total execution time on the CPU-GPU heterogeneous processor. This section demonstrates  
 373 the benefits of our overhead minimization technique in two integrated kernels: *intg\_kernel\_1*  
 374 and *intg\_kernel\_2*. By comparing them with their separate versions, we notice that perfor-  
 375 mance gain comes from two aspects: reduced overhead and improved total kernel execution.

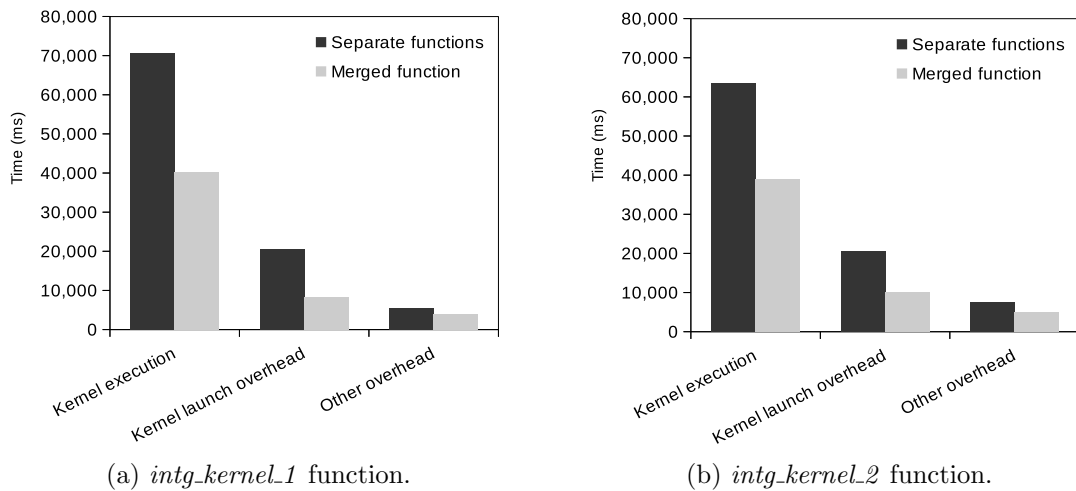


Figure 14: The impact of kernel overhead minimization.

376 Figure 14 shows the performance comparison between the two merged functions and  
 377 their corresponding separate versions. The results show that the overhead is reduced by  
 378 53% and 46% in both *intg\_kernel\_1* and *intg\_kernel\_2* respectively by merging kernels into a  
 379 single kernel. Moreover, merging kernels also speeds up the kernel execution (1.8x and 1.6x  
 380 respectively). By merging kernels, data can be reused across different individual kernels,  
 381 which reduces global memory accesses. Moreover, merging multiple kernels into a single  
 382 kernel increases the number of arithmetic computations that better hide memory latency [1].

## 383 6. Conclusions and Future Works

384 In this paper, we present the acceleration of DynEarthSol3D on tightly coupled CPU-  
 385 GPU heterogeneous processors by leveraging their new features, and compare its perfor-  
 386 mance and benefits with other serial and parallel alternatives. Our results show that the  
 387 OpenCL-based implementation on tightly coupled heterogeneous processors outperforms  
 388 both serial and OpenMP-based implementations that run on a multi-core CPU. We also  
 389 emphasize the importance of memory access pattern in GPGPU programming. With a  
 390 proper node ID system that reduces the randomness of global memory accesses, memory la-  
 391 tency is decreased significantly in our OpenCL-based optimization. Furthermore, *zero copy*

392 feature that is available on heterogeneous platform solves the big issue of expensive data  
393 transfer between host and device memory in conventional discrete GPU. Such benefits are  
394 examined in our in-depth analysis. We also discuss how integrating multiple small functions  
395 into a single kernel reduces both overhead and kernel execution time.

396 Our work demonstrates the potential of tightly coupled CPU-GPU heterogeneous proces-  
397 sors for the acceleration of data- and compute-intensive programs such as DynEarthSol3D.  
398 However, some issues of current heterogeneous processors need to be addressed in the future.  
399 The computing power of embedded GPU in current heterogeneous processors (e.g. 8 com-  
400 pute units in Kaveri) is much lower than the one of discrete GPUs (e.g. 32 compute units  
401 in Tahiti). This gap imposes a trade-off between better kernel performance on discrete plat-  
402 forms and “zero” data transfer on heterogeneous processors. In the future, heterogeneous  
403 processors are expected to provide more powerful compute units. Moreover, although the  
404 need for data transfer is eliminated, high overhead observed on AMD’s heterogeneous plat-  
405 form in our experiment needs to be eliminated. This problem can be addressed with better  
406 software supports (i.e. driver) from hardware vendor. Currently, OpenCL 1.2 does not sup-  
407 port all promising HSA features of heterogeneous computing. With the OpenCL 2.0 coming  
408 soon, we are looking forward to utilizing these new features in our future optimization of  
409 DynEarthSol3D.

## 410 7. References

- 411 [1] AMD, 2013. AMD Accelerated parallel processing: OpenCL programming guide.  
412 URL [http://developer.amd.com/wordpress/media/2013/07/AMD\\_Accelerated\\_Parallel\\_](http://developer.amd.com/wordpress/media/2013/07/AMD_Accelerated_Parallel_Processing_OpenCL_Programming_Guide-rev-2.7.pdf)  
413 [Processing\\_OpenCL\\_Programming\\_Guide-rev-2.7.pdf](http://developer.amd.com/wordpress/media/2013/07/AMD_Accelerated_Parallel_Processing_OpenCL_Programming_Guide-rev-2.7.pdf)
- 414 [2] AMD, 2014. AMD A-Series APU Processors.  
415 URL <http://www.amd.com/en-gb/products/processors/desktop/a-series-apu>
- 416 [3] AMD, 2014. AMD Accelerated Processing Units (APUs).  
417 URL <http://www.amd.com/en-us/innovations/software-technologies/apu>
- 418 [4] Behn, M. D., Ito, G., Aug. 2008. Magmatic and tectonic extension at mid-ocean ridges: 1. Controls on  
419 fault characteristics. *Geochemistry Geophysics Geosystems* 9 (8), Q08O10.  
420 URL <http://www.agu.org/pubs/crossref/2008/2008GC001965.shtml>
- 421 [5] Braun, J., Thieulot, C., Fulsack, P., DeKool, M., Beaumont, C., Huisman, R., 2008. DOUAR: A new  
422 three-dimensional creeping flow numerical model for the solution of geological problems. *Phys. Earth*  
423 *Planet. Int.* 171 (1-4), 76–91.  
424 URL [http://www.sciencedirect.com/science/article/B6V6S-4SJ2WV0-3/2/](http://www.sciencedirect.com/science/article/B6V6S-4SJ2WV0-3/2/94d7290704141e50e939397d3c352cd6)  
425 [94d7290704141e50e939397d3c352cd6](http://www.sciencedirect.com/science/article/B6V6S-4SJ2WV0-3/2/94d7290704141e50e939397d3c352cd6)
- 426 [6] Buck, W. R., Lavier, L. L., Poliakov, A. N. B., Apr. 2005. Modes of faulting at mid-ocean ridges.  
427 *Nature* 434 (7034), 719–23.  
428 URL <http://www.nature.com/nature/journal/v434/n7034/full/nature03358.html>
- 429 [7] Choi, E., Gurnis, M., 2008. Thermally induced brittle deformation in oceanic lithosphere and the  
430 spacing of fracture zones. *Earth Planet. Sci. Lett.* 269, 259–270.
- 431 [8] Choi, E., Tan, E., Lavier, L. L., Calo, V. M., May 2013. DynEarthSol2D: An efficient unstructured  
432 finite element method to study long-term tectonic deformation. *Journal of Geophysical Research: Solid*  
433 *Earth* 118 (5), 2429–2444.  
434 URL <http://doi.wiley.com/10.1002/jgrb.50148>
- 435 [9] Čuma, M., Zhdanov, M. S., 2014. Massively parallel regularized 3d inversion of potential fields on cpus  
436 and gpus. *Computers & Geosciences* 62, 80–87.

- 437 [10] Cundall, P. A., Board, M., 1988. A microcomputer program for modeling large strain plasticity prob-  
438 lems. In: Swoboda, G. (Ed.), Numerical Methods in Geomechanics: pp. 2101–2108.
- 439 [11] HSA Foundation, 2013. Heterogeneous System Architecture (HSA) Foundation.  
440 URL <http://hsafoundation.com/>
- 441 [12] Huet, B., Le Pourhiet, L., Labrousse, L., Burov, E., Jolivet, L., Feb. 2011. Post-orogenic extension  
442 and metamorphic core complexes in a heterogeneous crust: the role of crustal layering inherited from  
443 collision. Application to the Cyclades (Aegean domain). *Geophysical Journal International* 184 (2),  
444 611–625.  
445 URL <http://doi.wiley.com/10.1111/j.1365-246X.2010.04849.x>
- 446 [13] Intel, 2014. Intel Core Processor Family.  
447 URL [http://www.intel.com/content/www/us/en/processors/core/core-processor-family.](http://www.intel.com/content/www/us/en/processors/core/core-processor-family.html)  
448 [html](http://www.intel.com/content/www/us/en/processors/core/core-processor-family.html)
- 449 [14] Ito, G., Behn, M. D., Sep. 2008. Magmatic and tectonic extension at mid-ocean ridges: 2. Origin of  
450 axial morphology. *Geochemistry Geophysics Geosystems* 9 (9).  
451 URL <http://www.agu.org/pubs/crossref/2008/2008GC001970.shtml>
- 452 [15] Jang, B., Schaa, D., Mistry, P., Kaeli, D., 2011. Exploiting memory access patterns to improve memory  
453 performance in data-parallel architectures. *Parallel and Distributed Systems, IEEE Transactions on*  
454 22 (1), 105–118.
- 455 [16] Lyakhovskiy, V., Segev, A., Schattner, U., Weinberger, R., Jan. 2012. Deformation and seismicity as-  
456 sociated with continental rift zones propagating toward continental margins. *Geochemistry Geophysics*  
457 *Geosystems* 13.  
458 URL <http://www.agu.org/pubs/crossref/2012/2011GC003927.shtml>
- 459 [17] Munshi, A., et al., 2009. The opencl specification. Khronos OpenCL Working Group 1, 11–15.
- 460 [18] Nvidia, C., 2014. Programming guide.  
461 URL <http://docs.nvidia.com/cuda/cuda-c-programming-guide>
- 462 [19] Poliakov, A., Buck, W. R., 1998. Mechanics of stretching elastic-plastic-viscous layers: Applications to  
463 slow-spreading mid-ocean ridges. In: Buck, W. R., Delaney, P. T., Karson, J. A., Lagabrielle, Y. (Eds.),  
464 *Faulting and Magmatism at Mid-Ocean Ridges*. Vol. 106 of AGU Monograph. AGU, Washington D.C.,  
465 pp. 305–324.
- 466 [20] Poliakov, A. N. B., Cundall, P. A., Podladchikov, Y. Y., Lyakhovskiy, V. A., 1993. An explicit inertial  
467 method for the simulation of viscoelastic flow: An evaluation of elastic effects on diapiric flow in two-  
468 and three-layers models. In: Stone, D. B., Runcorn, S. K. (Eds.), *Flow and Creep in the Solar Systems:*  
469 *Observations, Modeling and Theory*. Kluwer Academic Publishers, pp. 175–195.
- 470 [21] Shevtsov, M., 2013. OpenCL: Advantages of the Heterogeneous Approach - Intel Developer Zone.  
471 URL <https://software.intel.com/en-us/articles/opencl-the-advantages-of-heterogeneous-approach>
- 472 [22] Si, H., TetGen, A., 2006. A quality tetrahedral mesh generator and three-dimensional delaunay trian-  
473 gulator. Weierstrass Institute for Applied Analysis and Stochastic, Berlin, Germany.
- 474 [23] Su, L. T., 2013. Architecting the future through heterogeneous computing. In: *Solid-State Circuits*  
475 *Conference Digest of Technical Papers (ISSCC), 2013 IEEE International*. IEEE, pp. 8–11.
- 476 [24] Tahmasebi, P., Sahimi, M., Mariethoz, G., Hezarkhani, A., 2012. Accelerating geostatistical simulations  
477 using graphics processing units (gpu). *Computers & Geosciences* 46 (0), 51 – 59.  
478 URL <http://www.sciencedirect.com/science/article/pii/S0098300412001240>
- 479 [25] Tan, E., Choi, E., Lavier, L., Calo, V., 2013. DynEarthSol3D: An Efficient and Flexible Unstructured  
480 Finite Element Method to Study Long-Term Tectonic Deformation,. Abstract DI31A-2197 presented  
481 at 2013 Fall Meeting, AGU, San Francisco, Calif., 9-13 Dec.