

<div><div><div><div><div><div></div><div>EURO-PAR</div></div><div><div></div><div></div></div></div><div><div><div></div><div></div></div><div><div></div><div></div></div></div><div><div><div></div><div></div></div><div><div></div><div></div></div></div></div><div>Program</div></div><div>Program is also available in Google Calendar [edit]</div><div>Conference proceedings are available on Springer website.</div><div>Recordings available on YouTube.</div><div><input checked="" type="checkbox"/> Show abstracts</div></div>	
WEDNESDAY 26.08.2020	
14.00 - 14.30	Opening
14.30 - 15.30	<div><div><div><div><div><div>Multicore and Manycore Parallelism (A)</div><div>Chairs: Wilfried Radacki [link] [video]</div></div><div><div><div>NVPM: An Efficient Phase-Based Transactional System for Non-Volatile Memory</div><div>Alexandro Baldassari, Rafael Maruri, João Paulo Carvalho, Guido Anzuino, David Castro, João Barreto and Paolo Romano</div><div>Download paper from Springer LNCS. [video]</div></div></div><div><div><div>Non-Volatile Memory (NVM) is an emerging memory technology aimed to eliminate the gap between main memory and stable storage. Nevertheless, today's programs will not readily benefit from NVM because crash failures may render the program non-recoverable and inconsistent state. In this context, the use of durable transactions has been proposed so as to ease the adoption of NVM. It leverages on the well-known semantics of database transactions to simplify the task of programming NVM systems. This is achieved by logging NVM writes using software (SW) or hardware (HW) transaction primitives. Although SW transactions are flexible and unbounded, they may significantly hurt the performance of short-lived transactions. On the other hand, HW transactional memories provide low-overhead but are resource-constrained. In this paper we present NVPM, a transactional system for NVM that delivers the best out of both HW and SW transactions by dynamically selecting the best execution mode according to the application's characteristics. NVPM is comprised of a set of heuristics to guide online phase transition. Furthermore, a careful design of the phase transition state is devised to guarantee persistency when transitioning between HW and SW phases. To the best of our knowledge, NVPM is the first phase-based system to provide durable transactions. Experimental results with the STAMP benchmark show that the proposed heuristics are efficient in guiding phase transitions with low overhead. In particular, the NVM-aware heuristics provided an average speedup of up to 10.4x when compared to a system using NVM-oblivious heuristics, with only 1.9x of transaction overhead in the worst case.</div><div><div><div>Enhancing Resource Management through Prediction-based Policies</div><div>Antonio Navarro, Arthur F. Lorenzon, Eduardo Ayguade and Vicenç Beltrón</div><div>Download paper from Springer LNCS. [video]</div></div></div></div><div><div><div>Task-based programming models are emerging as a promising alternative to make the most of multi-/many-core systems. These programming models rely on runtime systems, and their goal is to improve application performance by properly scheduling application tasks to cores. Additionally, these runtime systems offer policies to cope with application phases that lack parallelism to fill all cores. However, these policies are usually static and favor either performance or energy efficiency. In this paper, we have extended a task-based runtime system with a lightweight monitoring and prediction infrastructure that dynamically predicts the optimal number of cores required for each application to identify the optimal set of data placement, performance and energy efficiency. Through the execution of several benchmarks in multi-/many-core systems, we show that our prediction-based policies have competitive performance while improving energy efficiency when compared to state of the art policies.</div><div><div><div>Accelerating Overlapping Community Detection Performance Tuning a Stochastic Gradient Markov Chain Monte Carlo Algorithm</div><div>Emil Hritonenko, Ralf Steinmann and Michael Kopp</div><div>Download paper from Springer LNCS. [video]</div></div></div></div><div><div><div>Building efficient algorithms for data-intensive problems requires deep analysis of data access patterns. Random data access patterns exacerbate this process. In this paper, we discuss accelerating randomized data-intensive machine learning algorithm using multi-core CPUs and several GPUs. A thorough analysis of the algorithm's data dependencies enabled a 75% reduction in its memory footprint. We created custom compute kernels to code generation to speed up the execution of data placement and computational optimizations per compute device. An empirical evaluation shows up to 245% speedups compared to an optimized sequential version. Another result from this evaluation is that achieving performance does not always match intuition (e.g., performance on the GPU architecture, vectorization may increase or hurt performance).</div><div><div><div>15.30 - 16.00</div><div>Break</div></div></div></div><div><div><div>16.00 - 17.20</div><div><div>Support Tools and Environments (A)</div><div>Chairs: Bartosz Białe [link] [video]</div><div><div><div>Skipping Non-essential Instructions Makes Data-dependence Profiling Faster</div><div>Nicolas Moiré, Amr Abd Elhameed, Ali Jannasari and Felix Wolf</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Data-dependence profiling is a dynamic program-analysis technique to discover potential parallelism in applications. Unlike purely static analysis, which may miss opportunities for parallelization, dynamic profiling does not know many pointers/values and array indices at compile time, profiling has the advantage of recording data dependencies that actually occur at runtime. But it has the disadvantage of significantly slowing down the application. In this paper, we discuss accelerating randomized data-intensive machine learning algorithm using multi-core CPUs and several GPUs. A thorough analysis of the algorithm's data dependencies enabled a 75% reduction in its memory footprint. We created custom compute kernels to code generation to speed up the execution of data placement and computational optimizations per compute device. An empirical evaluation shows up to 245% speedups compared to an optimized sequential version. Another result from this evaluation is that achieving performance does not always match intuition (e.g., performance on the GPU architecture, vectorization may increase or hurt performance).</div><div><div><div>15.30 - 16.00</div><div>Break</div></div></div></div><div><div><div>16.00 - 17.20</div><div><div>Scheduling and Load Balancing (B)</div><div>Chairs: Joanna Berlińska [link] [video]</div><div><div><div>Parallel Scheduling for Data-Intensive Tasks</div><div>Miao Meng and Lukasz Golab</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Workloads with precedence constraints due to data dependencies are common in various applications. These workloads can be represented as directed acyclic graphs (DAGs). These are often data-intensive, meaning that data dependencies are more suitable to be migrated. The policy evaluates, for each VM, both the CPU load and the network traffic influence on the assigned host. The corresponding Pearson correlation coefficients are calculated for each one of the VMs and then weighted in order to provide a relationship between the values and the host behaviour. The main goal is to test for runtime systems. However, we note that the detailed study of the parallelism of this algorithm is currently lacking. In this paper, we present new theoretical results about the tiled Cholesky factorization in the context of a parallel homogeneous model without communication costs. Based on the heuristic costs of involved nodes, we proved that different execution times must be considered, typically concerning to CPU and GPU costs. By a careful analysis on the number of tasks of each type that run simultaneously in the ALAP (As Late As Possible) schedule without resource contention, we are able to determine the minimum number of busy processors for any time step of the DAG. We use this information to find a closed form formula for the minimum time to schedule tiled Cholesky factorization of size n on p processors. We show that this bound outperforms classical bounds from the literature. We also prove that the ALAP, on an ALAP-based schedule, is the best scheduling policy. This has a milestone extremely close to the lower bound, thus proving both the effectiveness of ALAP's schedule and of the lower bound on the makespan.</div><div><div><div>Optimal GPU-CPU Offloading Strategies for Deep Neural Network Training</div><div>Oliver Beaumont, Lionel Eyraud-Dubois and Aлена Shilova</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Training Deep Neural Networks is known to be an expensive operation, both in terms of computational cost and memory load. Indeed, during training, all intermediate data need to be kept in memory. This leads to a large memory footprint. The forward phase must be stored until the corresponding gradient has been computed in the backward phase. These memory requirements sometimes prevent to consider larger batch sizes and deeper networks, so that they can limit both convergence speed and accuracy. Recent works have proposed to offload some of the computed forward activations from the memory of the GPU to the memory of the CPU. This requires to determine which activations should be offloaded and when these transfers from and to the memory of the GPU should take place. We prove that this problem is NP-complete. In this paper, we propose a heuristic to select activations based on relaxations of the problem. We perform extensive experimental evaluation on standard Deep Neural Networks. We compare the performance of our heuristics against previous approaches from the literature, showing that they achieve much better performance in a wide variety of situations.</div><div><div><div>Improving mapping for sparse direct solvers: A trade-off between Cholesky and LU</div><div>Changyi Song, Ali Al-Zahrani, Anne Benoit, Mathieu Faverge, Loris Marchal, Grégoire Pichon and Pierre Ramet</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>In order to express parallelism, parallel sparse direct solvers take advantage of the elimination tree to exhibit tree-shaped task graphs, where nodes represent computation stages and edges represent data dependencies. One of the pre-processing stages of sparse direct solvers consists of mapping each node of the elimination tree to a specific hardware unit in order to minimize the factorization time by exploiting good data locality and load balancing. The proportional mapping technique is a widely used approach to solve this resource-allocation problem. It achieves good data locality by assigning the same processors to large parts of the elimination tree. However, it may limit load balancing in some cases. In this paper, we propose a dynamic mapping algorithm based on proportional mapping. This new approach reduces the data locality criterion to improve load balancing. In order to validate the newly introduced method, we performed extensive experiments on the HPC sparse direct solvers. This demonstrates that our algorithm enables better static scheduling of the numerical factorization while keeping good data locality.</div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div>
17.20 - 17.30	Break
17.30 - 18.20	<div><div><div><div><div><div>Keynote Ewa Deelman (A)</div><div>Automating Science Workflows: Challenges and Opportunities</div><div>Chair: Rocco Scalettau</div><div>Abstract available on keynotes page</div></div><div><div><div>18.20 - 19.00</div><div>Welcome reception</div></div></div></div><div><div><div>THURSDAY 27.08.2020</div><div><div><div>13.00 - 14.30</div><div>Industry: Huawei (A)</div><div>Details available on "Industry: Huawei" page</div></div><div><div><div>14.30 - 15.00</div><div>Best paper (A)</div></div></div></div><div><div><div>Chairs: Morris Redel [link] [video]</div><div><div><div>Maximizing I/O Bandwidth for Reverse Time Migration on Heterogeneous Large-Scale Systems</div><div>Tania Alkureishi, Hiderem Lütkenfeld and David Keyes</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Reverse Time Migration (RTM) is an important scientific application for oil and gas exploration. The 3D RTM simulation generates terabytes of intermediate data that does not fit in main memory. In particular, RTM has two successive computational phases, i.e., the forward modeling and the backward propagation, that necessitate to write and then to read the state of the computed solution at specific time steps of the time integration. Advances in hardware and software have enabled the use of RTM on heterogeneous systems. However, the performance of RTM is still a challenge. In this paper, we propose a Parallel File Systems (PFS) to intermediate fast disk technologies (e.g., node-local and remote-on-site Burst Buffer) and up to CPU main memory. To address the trend of heterogeneous HPC systems deployment, we introduce an extension to our Multi-Buffer Buffer System (MLBS) framework to further maximize RTM I/O bandwidth in presence of GPU hardware accelerators. The main idea is to leverage the GPU's high-Bandwidth Memory (HBM) as an additional storage medium. The objective of MLBS is ultimately to hide the application's I/O overhead by ending a buffering mechanism operating across a multi-tiered hierarchy. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.00 - 15.30</div><div>Best artifact (A)</div></div></div></div><div><div><div>Chairs: Moritz Späthler [link] [video]</div><div><div><div>A Prediction Framework for Fast Sparse Triangular Solves</div><div>Najeeb Ahmad, Buse Yilmaz and Dilem Ünüt</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Sparse triangular solve (SpTSV) is an important linear algebra kernel, finding extensive uses in numerical and scientific computing. The parallel implementation of SpTSV is a challenging task due to the sequential nature of the steps involved. This makes it, in many cases, one of the most time-consuming operations in an application. Many approaches for efficient SpTSV on CPU and GPU systems have been proposed in the literature. However, no single implementation on six CPU or GPU gives the fastest solution for all input sparse matrices. In this work, we propose a machine learning-based framework to help choose the fastest solver for a given sparse matrix. Our framework is evaluated on a set of 1000 sparse matrices. The proposed framework is tested with six SpTSV implementations on a state-of-the-art CPU-GPU machine (Intel Xeon Gold CPU, NVIDIA V100 GPU). Experimental results, with 998 matrices taken from the SuiteSparse Matrix Collection, show the classifier prediction accuracy of 87% for the fastest SpTSV algorithm for a given input matrix. Predicted SpTSV implementations achieve average speedups (harmonic mean) in the range of 1.2-4.2x against the six SpTSV implementations used in the evaluation.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between computation and memory access, we explore the full-on opportunities in the CNN computation graph to reuse data. In this paper, we further verify the effectiveness of our approach by comparing the performance of our PFS with a performing expensive I/O operations and creating opportunities for overlapping data movement with computations. MLBS may transform the original I/O bound behavior of the RTM application into a compute-bound regime. In fact, the prefetching strategy of MLBS allows the RTM application to believe that it has access to a larger memory capacity on the GPU, while transparently performing the necessary housekeeping across the storage layers. We demonstrate the effectiveness of our approach by comparing the performance of our PFS-based RTM implementation for large 3D solution grid to 1.4X performance speed-up, compared to the reference PFS-based RTM implementation for large 3D solution grid.</div><div><div><div>15.30 - 15.40</div><div>Break</div></div></div></div><div><div><div>15.40 - 16.40</div><div><div>Data Management, Analytics and Machine Learning (A)</div><div>Chairs: Moritz Redel [link] [video]</div><div><div><div>Accelerating Deep Learning Inference with Cross-Layer Data Reuse on GPUs</div><div>Jueying Wang, Guang Li, Xiao Dong, Jiansong Li, Lei Liu and Xiaobing Feng</div><div>Download paper from Springer LNCS. [video]</div></div><div><div>Accelerating the deep learning inference is very important for real-time applications. In this paper, we propose a novel approach to use the key of convolutional map reuse in the GPU hardware acceleration. The proposed approach applies data reuse analysis and access optimization in different levels of the memory hierarchy. To achieve the balance between</div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div></div>



SHARE ON:

