

The PEPPER composition tool: performance-aware composition for GPU-based systems

Usman Dastgeer · Lu Li · Christoph Kessler

Received: 14 March 2013 / Accepted: 7 November 2013 / Published online: 27 November 2013
© Springer-Verlag Wien 2013

Abstract The PEPPER (EU FP7 project) component model defines the notion of component, interface and meta-data for homogeneous and heterogeneous parallel systems. In this paper, we describe and evaluate the *PEPPER composition tool*, which explores the application’s components and their implementation variants, generates the necessary low-level code that interacts with the runtime system, and coordinates the native compilation and linking of the various code units to compose the overall application code to optimize performance. We discuss the concept of *smart containers* and its benefits for reducing dispatch overhead, exploiting implicit parallelism across component invocations and runtime optimization of data transfers. In an experimental evaluation with several applications, we demonstrate that the composition tool provides a high-level programming front-end while effectively utilizing the task-based PEPPER runtime system (StarPU) underneath for different usage scenarios on GPU-based systems.

Keywords PEPPER project · Component model · GPU-based systems · Performance portability · Dynamic scheduling

Mathematics Subject Classification 68N20 Compilers and interpreters

1 Introduction

The PEPPER framework [1], developed by the European FP7 project PEPPER 2010–2012, is a unified approach for programming and optimizing applications for

U. Dastgeer (✉) · L. Li · C. Kessler
PELAB, Department of Computer and Information Science,
Linköping University, Linköping, Sweden
e-mail: usman.dastgeer@liu.se
URL: <http://www.ida.liu.se/>

heterogeneous many-core systems to enable performance portability. It consists of three main parts: (1) a flexible and extensible component model for encapsulating and annotating performance critical parts of the application, (2) adaptive algorithm libraries that implement the same basic functionality across different architectures, and (3) an efficient run-time system that schedules compiled components across the available resources, using performance information provided by the components layer as well as other, execution-history-based performance information.

The PEPPER component model defines an environment for annotation of native C/C++ based components for homogeneous and heterogeneous multicore and many-core systems, including GPU and multi-GPU based systems. For the same computational functionality, captured as a component, different sequential and explicitly parallel implementation variants using various types of execution units might be provided, together with metadata such as explicitly exposed tunable parameters. The goal is to compose an application from its components and variants such that, depending on the run-time context, the most suitable implementation variant will be chosen automatically for each invocation.

In this article, we present the *PEPPER composition tool*, which adapts applications written using the PEPPER component model to effectively exploit the capabilities offered by the runtime system. It gathers important information about the application and its components from their XML-based metadata descriptors, uses that information for static composition and generates glue code creating tasks for the PEPPER runtime system (which is based on StarPU [2]) for dynamic composition. We evaluate the composition capabilities of our composition tool with several applications from the RODINIA benchmark suite [3] and other sources. Following are the major contributions of our work:

- *Composition tool* for automatically composing performance-aware applications for GPU-based systems, offering both static and dynamic composition capabilities at a higher abstraction level.
- *Smart containers*, generic STL-like container data structures such as *Vector* and *Matrix* that can be used for wrapping native operand data of components and that provide several automatic optimizations: By keeping track of operand data copies in the different memory units at runtime, a smart container allows to automatically eliminate some unnecessary data transfers, e.g. between GPU device memory and main memory. It also reduces task creation overhead and enables the automatic exploitation of implicit parallelism across several component invocations by tracking container data flow.
- Evaluation of different composition capabilities of our composition tool with real applications taken from RODINIA benchmarks and other sources.

The remainder of this article is organized as follows: Section 2 describes central concepts of the PEPPER component model. Sections 3 and 4 describe the composition tool, its advanced features and its current prototype implementation. Section 5 presents a comprehensive experimental evaluation of the developed prototype. Section 6 discusses related work and Sect. 7 concludes the article.

2 Components, interfaces, implementations

A *PEPPER component* is an annotated software module that implements a specific functionality declared in a *PEPPER interface*. A *PEPPER interface* is defined by an *interface descriptor*, an XML document that specifies the name, parameter types and access types (read, write or both) of a function to be implemented, and which performance metrics (e.g. execution time) the prediction functions of component implementations must provide. Interfaces can be generic in static entities such as element types or code; genericity is resolved statically by expansion, as with C++ templates.

Applications for PEPPER are currently assumed to be written in C/C++. Several implementation variants may implement the same functionality (as defined by a PEPPER interface), e.g. by different algorithms or for different execution platforms. These implementation variants can exist already as part of some standard library¹ or can be provided by the programmer (called *expert programmer* by Asanovic et al. [6]). Also, more component implementation variants may be generated automatically from a common source module, e.g. by special compiler transformations or by instantiating or binding tunable parameters. These variants differ by their resource requirements and performance behavior, and thereby become alternative choices for composition whenever the (interface) function is called.

In order to prepare and guide variant selection, component implementations need to expose their relevant properties explicitly to a composition tool (described later). Each PEPPER component implementation variant provides its own *component descriptor*, an XML document that contains information (meta-data) about the provided and required interface(s), source files, compilation and resource requirements and reference to performance prediction functions. The `main` module of a PEPPER application is also annotated by its own XML descriptor, which states e.g. the target execution platform and the overall optimization goal. XML descriptors are chosen over code-annotations (e.g. pragmas) as the former are non-intrusive to the actual source code and hence provide better separation of concerns. The potential headache associated with writing descriptors in XML is eliminated to a great extent by providing tool support, as shown later. Composition points of PEPPER components are restricted to calls on general-purpose execution units only. Consequently, all component implementations using hardware accelerators such as GPUs must be wrapped in CPU code containing a platform-specific call to the accelerator.

Component invocations result in *tasks* that are managed by the PEPPER runtime system and executed non-preemptively. PEPPER components and tasks are *stateless*. However, the parameter data that they operate on may have state. For this reason, parameters passed in and out of PEPPER components may be wrapped in special portable, generic, STL-like *container* data structures such as `Vector` and `Matrix` with platform-specific implementations that internally keep track of, e.g., in which memory modules of the target system which parts of the data are currently located or mirrored (*smart containers*). When applicable, the container state becomes part of the call context information as it is relevant for performance prediction.

¹ For demonstration purpose, we have used CUBLAS [4] and CUSP [5] components for CUDA implementations as shown in Sect. 5.

The PEPPER framework automatically keeps track of the different implementation variants for the identified components, technically by storing their descriptors in repositories that can be explored by the composition tool. The repositories enable organization of source-code and XML annotation files in a structured manner and can help keeping files manageable even for a large project.

3 Composition tool

Composition is the selection of a specific implementation variant (i.e., callee) for a call to component-provided functionality and the allocation of resources for its execution. Composition is made *context-aware* for performance optimization as performance depends on the current *call context*, which consists of selected input parameter properties (such as size) and currently available resources (such as cores or accelerators). The context parameters to be considered and optionally their *ranges* (e.g., minimum and maximum value) are declared in the PEPPER interface descriptor. We refer to this considered subset of a call context instance's parameter and resource values shortly as a *context instance*, which is thus a tuple of concrete values for context properties that might influence callee selection.

Composition can be done either statically or dynamically. *Static composition* constructs off-line a *dispatch function* that is evaluated at runtime for a context instance to return a function pointer to the expected best implementation variant. *Dynamic composition* generates code that delegates the actual composition to a context-aware runtime system that records performance history to be used and updated as the application proceeds. Composition can even be done in multiple stages: First, static composition can narrow the set of candidates for the best implementation variant per context instance to a few ones that are registered with the context-aware runtime system that takes the final choice among these at runtime.

Dynamic composition is the default composition mechanism in PEPPER. The *PEPPER composition tool* deploys the components and builds an executable application. It recursively explores all interfaces and components that (may) occur in the given PEPPER application by browsing the interfaces and components repository. It processes the set of interfaces (descriptors) bottom-up in reverse order of their components' *required interfaces* relation (lifted to the interface level). For each interface (descriptor) and its component implementations, it performs the following tasks:

1. It reads the descriptors and internally represents the metadata of all component implementations that match the target platform, expands generic interfaces and components, and generates platform-specific header files from the interface.
2. It looks up prediction data from the performance data repository or runs microbenchmarking code on the target platform, as specified in the components' performance meta-data.
3. It generates composition code in the form of *stubs* (proxy or wrapper functions) that will perform context-aware composition at runtime. If sufficient performance prediction metadata is available, it constructs performance data and dispatch tables for static composition by evaluating the performance prediction functions for selected context scenarios which could be compacted by machine learning techniques [7].

Otherwise, the generated composition code contains calls to delegate variant selection to runtime, where the runtime system can access its recorded performance history to guide variant selection, in addition to other criteria such as operand data locality.

4. It calls the native compilers, as specified for each component, to produce a binary of every patched component source.

Finally, it links the application’s main program and its compiled components together with the generated and compiled stubs, the PEPPHER library and the PEPPHER runtime system to obtain an executable program. The linking step may be architecture dependent (e.g., special handling of different executable formats may be required); the necessary command can be found in the application’s `main` module descriptor.

4 Prototype implementation

A prototype of the composition tool has been implemented that supports both static and dynamic composition. For dynamic composition, the composition tool generates low-level code to interact with the runtime system in an effective manner. Furthermore it supports component expansion for generic components written using C++ templates as well as user-guided and offline learning based static narrowing of the set of candidates. In this section we describe its design, main features and implementation issues. Figure 1 shows a high-level schematic view of the prototype. Similar to typical compiler frameworks, we decouple composition processing (e.g., static composition decisions) from the XML schema by introducing an intermediate component-tree representation (IR) of the metadata information for the processed component interfaces and implementations. The IR incorporates information not only from the XML descriptors but also information given at composition time (i.e., composition recipe). The IR can be processed for different purposes, including:

- Creating multiple concrete components from generic components by expanding template types and tunable parameters,
- Training executions to prepare for composition decisions,
- Static composition (e.g. using offline learning), and
- Generating code that is executable with the PEPPHER runtime system.

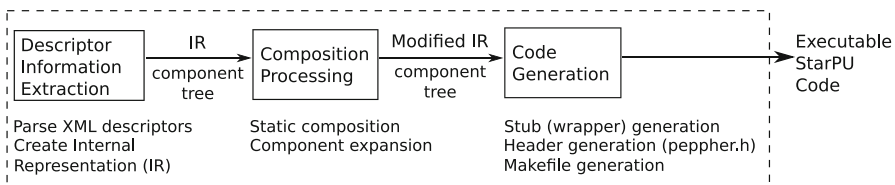


Fig. 1 Structural overview of the composition tool prototype

4.1 Static composition

Static composition refers to refining the composition choices at compile time, in the extreme case to one possible candidate per call and context instance. Currently, we support two kinds of static composition: The *user-guided static composition* provides a means for the programmer to convey to the composition tool additional expert knowledge about the context that may influence the composition process. For example, a GPU implementation normally runs faster than its CPU counterpart for data parallel problems with large problem size; thus if such information is statically known, programmers may explicitly specify to compose the GPU implementation, and the overhead of the dynamic composition and the risk of wrong dynamic selection can be removed. The composition tool provides simple switches (e.g., `disableImpls`) to enable/disable implementations at composition time without requiring any modifications in the user source-code.

Besides user guided static composition, we provide an *adaptive off-line sampling and learning* method to build empirical performance models for components and their implementation variants off-line in order to complement or replace the on-line sampling and learning done in the run-time system. We provide an adaptive sampling method that allows to optimize off-line sampling cost and runtime dispatch overhead while keeping prediction accuracy at a reasonable level. The details of our adaptive off-line sampling and learning method are presented elsewhere [8].

4.2 Component expansion

Component expansion supports genericity on the component parameter types using C++ templates. This enables writing generic components such as sorting that can be used to sort different types of data. The expansion takes place statically. Component expansion for multiple values of tunable parameters to generate multiple implementation variants from a single source is not supported yet and is part of future work.

4.3 Code generation

For each component interface, the composition tool generates a wrapper that intercepts each invocation of the component and implements logic to translate the call to one or more tasks for the runtime system. It also performs packing and unpacking of the call arguments to the PEPPER runtime task handler. The directory structure provides one directory for the main component of the application and one directory for each component used. The different available implementations of a component are organized by platform type (e.g., CPU/OpenMP, CUDA) in different subdirectories. A global registry of interfaces, implementations and platforms helps the composition tool to navigate this structure and locate the necessary files automatically.

The tool generates wrapper (stub) files providing wrapper functions for different components. Currently, one wrapper file is generated per *component*, containing one *entry wrapper* and multiple *implementation wrappers*. The *entry wrapper* for a component intercepts the component invocation call and implements logic to translate

that component call to one or more tasks in the runtime system. A task execution can either be *synchronous* where the calling thread blocks until the task completes or *asynchronous* where the control resumes on the calling thread without waiting for the task completion. The entry wrapper also performs packing and unpacking of arguments of a component call to the PEPPER runtime task handler. One *implementation wrapper* for a component is generated for each component implementation. Multiple component implementations (and consequently implementation wrappers) can exist for each backend (currently, CPU/OpenMP, CUDA, OpenCL). A implementation wrapper implements the function signature

```
void <name>(void *buffers[], void *arg)
```

that the runtime system expects for a task function,² and internally delegates the call to the actual component implementation with a different function signature.

Moreover, a header file (`pepper.h`) is generated which internally includes all wrapper files and also contains certain other helping code (e.g., extern declarations). The idea of this header file is to provide a single linking point between the generated code and the normal application code. The *main* program writer only needs to include this header file to enable the composition. Compilation code (`Makefile`) is also generated for compiling and linking the selected (composed) components to build an executable application for a given platform.

4.4 Usage of smart containers

A smart container can wrap operand data passed in and out of PEPPER components while providing a high-level interface to access that data. *Smart containers* model and encapsulate the state of their payload data. Four smart containers are currently implemented: for scalar value (`Scalar`), 1D array (`Vector`), 2D array (`Matrix`) and generic matrix (`GMatrix`). The `GMatrix` container can internally store data in different formats (2D dense, CSR sparse format etc.). All four containers are made generic in the element type, using C++ templates. The containers internally implement interaction with the data management API of the PEPPER runtime system while ensuring data consistency for data that can be accessed both by the runtime system and the application program in an arbitrary fashion. More precisely, these containers allow multiple copies of the same data on different memory units (CPU, GPU memory) at a certain time while ensuring consistency. For example, for a vector object that was last changed by a component call executed on a GPU, the vector data is copied back implicitly from GPU memory (i.e. enforce consistency) only when data is actually accessed in the application program (e.g. detected using the `[]` operator for vector objects). The containers are made portable and function as regular C++ containers outside the PEPPER context.

Note that the composition tool also supports parameters of any other C/C++ datatypes (including user-defined datatypes) for components. However, for parameters passed using normal C/C++ datatypes, the composition tool cannot reason about

² As the PEPPER runtime system is C based and the C language does not permit to call functions with varying types depending on the actual task being run.

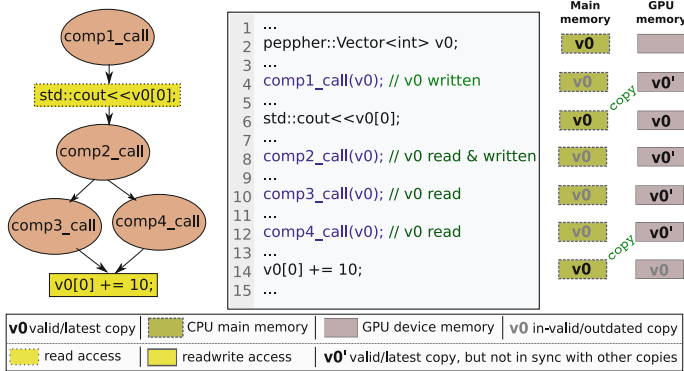


Fig. 2 Example for usage of smart containers with component calls and other user code. The middle part shows a code scenario with four component calls and one vector operand on a system containing 1 CPU and 1 CUDA GPU. The effect of each instruction on the state of the data is also shown *right* for each instruction, assuming the four component calls are executed on the GPU. The *left* part highlights inter-component parallelism for asynchronous component executions based on data dependencies

their access patterns in the application program (due to pointer-aliasing and other issues) and hence ensures data consistency by always copying data back to the main memory before returning control back from the component call. Although ensuring consistency, it may prove sub-optimal as data locality cannot be exploited for such parameters across multiple component calls.

4.5 Inter-component parallelism

In the PEPPIER framework, a major source of parallelism comes from exploitation of independence between different component invocations. This can be obtained when a component invocation is made asynchronous so that subsequent component invocations can overlap in the actual execution. By using the smart containers, the independence between different asynchronous component invocations is implicitly inferred by the PEPPIER runtime system based on data dependencies [2].

Figure 2 depicts how usage of smart containers can help in achieving inter-component parallelism. The figure shows a simple scenario with four component calls and one vector operand on a system containing 1 CPU and 1 CUDA GPU. When the vector container *v0* is created, the payload data is placed in the main memory (*master copy*). Subsequently, depending upon the component calls using that data along their respective data access pattern (read, readwrite or write), other (partial) copies of operand data may get created in different memory units. In this case, we have CUDA device memory which has a separate address space than main memory. Assuming that all component calls are actually executed on the GPU, the figure also shows the effect of each instruction execution on the vector data state, i.e., creation/update/invalidation of data copies. As we can see, a PEPPIER container not only keeps track of data copies on different memory units but also helps in minimising the data communication between different memory units by delaying the communication until it becomes

necessary. In this case, only 2 copy operations of data are made in the shown program execution instead of 7 copy operations which are required if one considers each component call independently, as done in e.g., Kicherer et al. [9, 10].

The first component call (line 4) only writes the data and hence no copy is made. Instead, just a memory allocation is made in the device memory where data is written by the component call. After the completion of the component call (line 4), the master copy in the main memory is marked outdated which means that any data access to this copy, in future, would first require update of this copy with the contents of the latest copy. The next statement (line 6) is actually a read data access³ from main memory. As the master copy was earlier marked outdated, a copy from device memory to main memory is implicitly invoked before the actual data access takes place. This is managed by the container in a transparent and consistent manner without requiring any user intervention. The copy in the device memory remains valid as the master copy is only read. Next, we have a component call (line 8) that both reads and modifies the existing data. As we assume execution of all component calls on the GPU in this scenario, the up-to-date copy already present in the device memory is read and modified. The master copy again becomes outdated. Afterwards, we have two component calls (line 10 and 12) that both only read the data. Executing these operations on the GPU means that no copy operation is required before or after the component call. Finally the statement in line 14 modifies the data in main memory so data is copied back (implicitly) from the device memory to the main memory before the actual operation takes place. Afterwards, the copy in the device memory is marked outdated and can be de-allocated by the runtime system if it runs short of memory space on the device unit. Doing so would however require re-allocation of memory for future usage.

As all four component calls are asynchronous, the inter component parallelism is automatically inferred by the runtime system based on data dependencies. In this case, there exists a read-after-write dependency between first (line 4) and second (line 8) component call; however, independence exists between third (line 10) and fourth (line 12) component call as they access the same data but in a read only manner.

In the application program, the execution looks no different to the synchronous execution as data consistency is ensured by the smart containers. Blocking is implicitly established for a data access from the application program to a data that is still in use with the asynchronous component invocations made earlier (with respect to program control flow) than the current data access.

4.6 Intra-component parallelism

A common source of intra-component parallelism is a parallel component implementation, e.g. a CUDA or OpenMP implementation. However, for certain computations, more parallelism can be spawned from a single component invocation by partitioning and dividing the work into several chunks that all can be processed concurrently, pos-

³ The read and write accesses to container data are distinguished by implementing proxy classes for element data in C++ [11].

sibly on different devices. This is achieved by mapping a single component invocation to multiple independent runtime (sub-)tasks rather than a single task where these (sub-)tasks can be executed on different computing units in parallel. This feature is implemented for data-parallel computations where the final result can be produced either by simple concatenation or simple reduction of intermediate output results produced by each sub-task (e.g. blocked matrix multiplication, dotproduct). When applicable, this feature can be used in a transparent manner as it does not require any modification in the user code and/or component implementations.

4.7 Support for performance-aware component selection

In the current prototype, the actual implementation variant selection is done by default using the dynamic scheduling capabilities of the PEPPHER runtime system. The actual implementation of performance-aware selection is made transparent in the prototype by providing a simple boolean flag (`useHistoryModels`). The support can be enabled/disabled both for an individual component by specifying the boolean flag in the XML descriptor of that component interface or globally for all components as a command line argument to the composition tool.

4.8 Utility mode

Porting existing applications to PEPPHER framework (PEPPHER-ization) requires writing XML descriptors for interface and implementations as well as adding new implementation variants. To facilitate this process, the current prototype can generate a basic skeleton of these XML and C/C++ source files required for writing PEPPHER components from a simple C/C++ method declaration. This can be really advantageous as writing XML files from scratch can become a tedious task. Our experience with porting several applications shows that the XML skeleton generation by the tool is quite useful as we are required to only fill in certain missing information. For example, the tool can successfully detect template parameters as well as suggest values for the data access pattern field of the descriptors by analyzing ‘const’ and ‘pass by reference’ semantics of the function arguments.

5 Evaluation

For evaluation, we implemented (PEPPHERized) several applications from the RODINIA benchmark suite [3], two scientific kernels (dense matrix-matrix and sparse matrix-vector multiplication) and a light-field image refocusing application, using the composition tool. Two evaluation platforms are used: The main evaluation platform is *System A* with Xeon[®] E5520 CPUs running at 2.27 GHz with 2 NVIDIA[®] C2050 GPUs with L1/L2 cache support. *System B* with Xeon[®] X5550 CPUs running at 2.67 GHz with a lower-end GPU (NVIDIA[®] C1060 GPU) is used for showing performance portability.

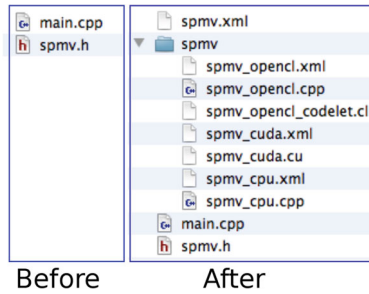


Fig. 3 Directory structure before and after generation of the basic PEPPER component skeleton files from a simple C/C++ header file containing a method declaration. The main work left for the programmer is now to fill in the implementation details in the XML descriptor fields and provide the implementation variants' code, which can be facilitated by the existence of architecture-specific algorithms and libraries

Composition example In the following, we will use the sparse matrix-vector multiplication to describe the complete composition (PEPPER-ization) process.

The process starts with generation of basic skeleton files for components from a C/C++ header file that includes the method signature

```
void spmv(float *values, int nnz, int nrows, int ncols, int first,
          size_t *colidxs, size_t *rowPtr, float *x, float *y);
```

The composition tool can be invoked to generate component skeleton files from this method declaration: As shown in Fig. 3, the command

```
compose -generateCompFiles=" spmv.h "
```

generates the XML descriptors with most information pre-filled as well as basic skeletons for implementation files. The programmer can then fill in the remaining details both in the XML descriptors (e.g. preferences for partitioning of input operands etc.) as well as the implementation files. For this example application, we used one parallel OpenMP implementation for CPUs and a highly optimized CUDA algorithm provided by NVIDIA as part of their CUSP library [5]. In the application's main module `main.cpp`, we only need to add one include statement for `pepper.h` and calls to `PEPPER_INITIALIZE()` and `PEPPER_SHUTDOWN()` macros in the beginning and end of the main function, which in this case contains just a call to the `spm` component. The last step is writing of the XML descriptor (`main.xml`) for the main function. The composition tool can now be called for generating composition code, by giving a reference to the main descriptor:

```
compose main.xml
```

This generates all the wrapper files and compilation code (makefile), necessary to compile and build the executable with the runtime support. The resulting executable can then be executed on the given platform.

Table 1 Showing saving in total source LOC (Lines of Code) written by the programmer when using the composition tool compared to an equivalent code written directly using the runtime system

Application	Savings (LOC, %)	Application	Savings (LOC, %)
bfs	108, 42	cfid	123, 62
hotspot	120, 37	lud	76, 15
nw	90, 25	particlefilter	96, 15
pathfinder	89, 48	nn	101, 41
streamcluster	84, 5	leukocyte	271, 8
b+tree	211, 10	lavaMD	94, 22
spmv	83, 29	sgemm	89, 63
Image refocusing	113, 11	ODE solver	452, 57

Productivity evaluation The current composition tool prototype generates low-level glue code to use the runtime system. This essentially allows application programs to run using the runtime system without requiring the programmer to actually write code for the runtime system. In the following, we compare how much we gain in terms of programming effort with this code generation functionality.

Table 1 shows a simple comparison of the source code written by the programmer when doing hand-written implementations for the runtime system compared to when using the composition tool. The comparison is done using standard LOC (Lines Of Code) metric [12] for all applications. For the ODE Solver application, we have considered LOC related to the ODE solver, and not the complete LibSolve library which contains more than 12,000 LOC. As we can see in Table 1, savings in terms of LOC are significant (up to 63 % for *cfid* and *sgemm*). These savings become even more significant considering the fact that the source code for the runtime system is at a low level of abstraction (plain C) and requires handling concurrency issues associated with asynchronous task executions.

The above comparison does not consider the XML descriptors which one needs to provide when using the composition tool. However, this is justified as the XML descriptor skeletons generated automatically from the C/C++ function declaration are already quite concise. For all these applications, we have used the XML descriptor generation feature of the composition tool to generate the XML descriptors and have just filled in the missing information (e.g., values for partitioning attribute).

In the remaining part of this section, we will evaluate the efficiency of the generated code by performance evaluation.

Hybrid execution As discussed in Sect. 4.6, a key aspect of the PEPPER component model is hybrid execution where, instead of using either CUDA or OpenMP implementation, the execution work is partitioned into multiple parts which can be processed concurrently, possibly on different devices by different component implementations. Figure 4 shows the advantage of hybrid execution for sparse matrix vector product execution using the matrices from the UF collection [13] as example data. The hybrid execution using the code generated by the composition tool supersedes direct OpenMP and direct CUDA execution and, in some cases, reduces the execution time to half of the best performing variant. The speedup by hybrid execution is obtained by divid-

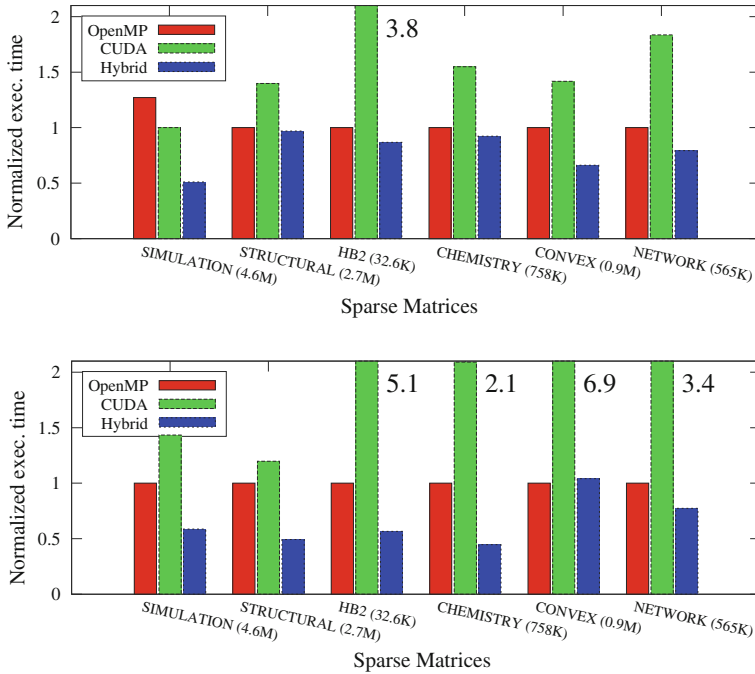


Fig. 4 Sparse matrix vector product execution for different matrices (with different number of non-zero elements) from the UF collection [13] on System A (*top* using 1 C2050 GPU) and System B (*bottom*) respectively

ing the computation work into multiple chunks and executing them on OpenMP and CUDA in parallel rather than using either of them. When applicable (see Sect. 4.6), this feature can be used in a transparent manner (controlled using the `partition` option in the interface descriptor) for an application, as it does not require any modification in the user code and/or component implementations.

Dynamic scheduling The architectural and algorithmic differences between different devices (e.g., CPU, GPU with/without cache) and applications often have a profound effect on the achieved performance. However, these execution differences are often hard to model statically as they can originate from various sources (hardware architecture, application/algorithm, input data, problem size etc.). Dynamic scheduling (including selection) can help in this case by deciding which implementation/device to use by considering previous historical execution information as well as information about current load balance, data locality and potential data-transfer cost (performance-aware dynamic scheduling). Our tool generates the necessary code for using performance-aware scheduling policies offered in the runtime system. Figure 5 shows how the usage of dynamic scheduling can help in achieving better performance on two heterogeneous architectures for a variety of applications. The execution time is averaged over different problem sizes. As we can see the execution time with generated code closely follows the best implementation from OpenMP and CUDA for all these applications. In some cases, the execution with dynamic scheduling supersedes

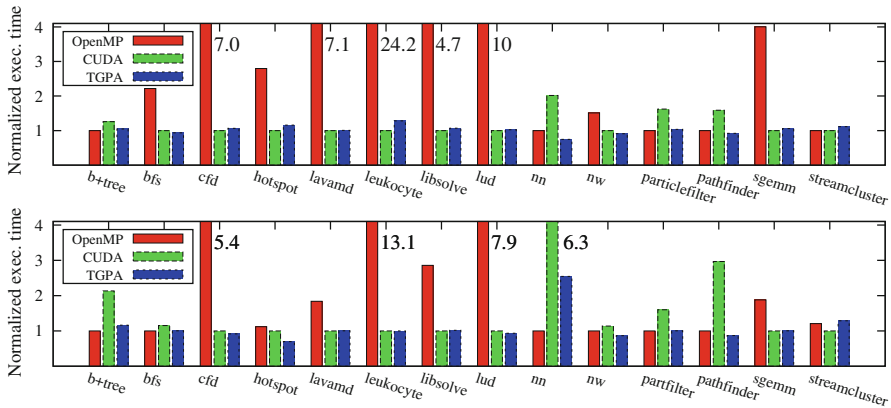


Fig. 5 Execution times for applications from the Rodinia benchmark suite, an ODE solver and sgemm with CUDA, OpenMP and our tool-generated performance-aware code (TGPA) on System A (*top*) and B (*bottom*)

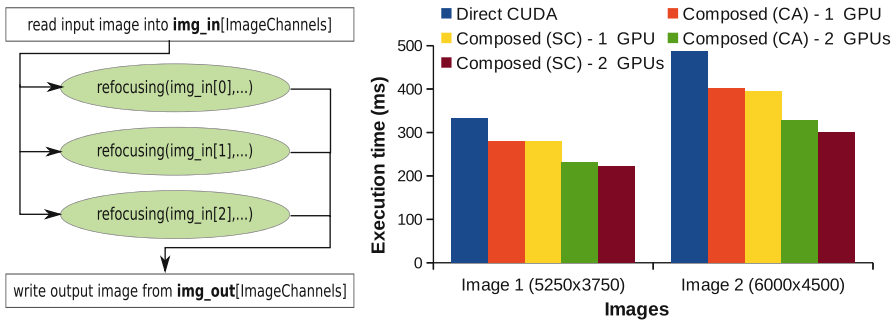


Fig. 6 Light-field image refocusing application [14]: *left* shows inter-component parallelism; *right* compare direct execution on a CUDA GPU with the composed version, using both C-Arrays (CA) and Smart Containers (SC)

the best static selection by making appropriate decisions for each problem size. Above all, the scheduling can effectively adjust to the architectural differences as depicted in Fig. 5. This is achieved by effective utilization of the performance-aware dynamic scheduling mechanism offered by the PEPPIER runtime system.

Smart containers and inter-component parallelism As mentioned earlier, smart containers provide high-level abstraction and underneath optimize data communication for GPU execution while ensuring data consistency for accesses in the user program. Figure 6 (left) highlights inter-component parallelism (i.e., data independence) for the light-field image refocusing application [14] across three image color channels. Furthermore, a single refocusing call can internally be partitioned into two or more independent parts (sub-tasks) based on input data to further exploit intra-component parallelism (see Sect. 4.6).

Figure 6 (right) compares the refocusing application using direct execution on the GPU with different versions of the composed application. The composed version using C-Arrays (CA) only exploits intra-component parallelism by internally dividing each

refocusing call into *two* parts. This gives a significant performance improvement over direct execution even on a single GPU that supports concurrent kernel executions by overlapping the memory transfers for the second part with computation work of the first part. Using two GPUs further reduces execution time by scheduling both computations on different GPUs.

Besides intra-component parallelism, the version using smart containers (SC) can exploit inter-component parallelism across different image channels and hence performs better than the version using C-Arrays, especially when using multiple GPUs. Once PEPPHERized, this potential for both intra- and inter-component parallelism as well as scaling across multiple GPUs can be exploited without requiring any further modifications in the user code.

PEPPHER runtime overhead The main runtime overhead of the PEPPHER framework is overhead of the PEPPHER runtime system. Measuring overhead of the runtime system, in general, is a nontrivial task as it depends upon the computation structure and granularity, scheduling policy, worker types as well as hardware and software architecture. Micro-benchmarking results reported in [15] show that the task overhead of the runtime system is less than two microseconds. This overhead is negligible when considering potential gains of dynamic performance-aware scheduling and overlapping communication-computation offered by the runtime system.

6 Related work

Generating stubs (wrapper or proxy functions) from formal interface descriptions in order to intercept and translate calls to components at runtime is a key concept in CORBA and subsequent component frameworks to transparently bridge the technical gaps between different languages, platforms and network locations of caller and callee, without having to change their source code. In our case, languages, platforms and network locations are equal, while the stubs encapsulate the necessary logic for determining a subset of suitable implementation variants and creating tasks for the PEPPHER runtime system. An extension to support other languages than C/C++, remote calls etc. following the CORBA approach would be a straightforward extension of our stub generation method.

Several language based systems for runtime composition of annotated implementation variants have recently been proposed in the literature, such as PetaBricks [16], EXOCHI [17], Merge [18], IBM's Liquid Metal [19] and Elastic computing [20]. Our approach with performance-aware components is different from these frameworks as it does not require the programmer to use a new programming language (or a domain-specific language [21]) to implement component implementations; PEPPHER components may internally encapsulate arbitrarily specified intra-component parallelism (e.g., using pthreads, OpenMP, or arbitrary accelerator specific code) running on the resources allocated by the runtime system to the task created for its invocation, as long as the actual composition points (calls to component-provided functionality) are single-threaded, C/C++ linkable, and executed on CPUs (i.e., default execution unit with access to main memory) only. Hence, the PEPPHER component framework is, by default, non-intrusive, i.e., all metadata for components and the main

program is specified externally in XML based descriptors. This enables an incremental "PEPPHER-ization" process of making legacy applications performance-portable (in principle) without modifying the existing source code.

The single-threaded call conventions, the XML format for external annotations, a non-preemptive task-based runtime system and dynamic composition as the default composition mechanism in PEPPHER are the major differences from the Kessler/Löwe components [7] which assume SPMD calls, off-line generated dispatch tables and nestable components instead.

Kicherer et al. [9,10] propose a C-based on-line learning solution for making dynamic selection between available implementation variants for a given invocation context. However, their approach is different as they represent each implementation variant as a dynamic library that is loaded at runtime instead of a single binary solution. Furthermore, they do not consider the issue of memory management for CUDA implementation variants nor do they consider simultaneous (hybrid) execution at multiple devices present in the system.

7 Conclusion

Writing performance-portable applications for modern heterogeneous architectures is a non-trivial task. The component-based approach of the PEPPHER framework allows the specification of multiple implementation variants for a single functionality where the expected best variant for a given execution context can be selected statically and/or dynamically. The purpose of the PEPPHER composition tool is to build applications from annotated components; it thereby supports "PEPPHER-izing" both new and legacy applications.

We have described the current composition tool prototype and how it could be used to deploy applications with performance-aware components on heterogeneous multi-/manycore systems. The composition tool can significantly reduce the programming complexity by transparently handling the concurrency issues and generating the low-level C code that interacts with the runtime system. On the performance side, the composition tool can successfully meet different composition needs, ranging from hybrid execution to dynamic scheduling. Usage of smart containers not only provides higher abstraction to the application programmer but can also give increased performance by exploiting inter-component parallelism. Experiments have shown that in all cases, the code generated by our tool executes efficiently on different platforms without any need for manual tuning, thanks to careful exploitation of features provided by the underlying runtime system.

Acknowledgments This work was funded by EU FP7, project PEPPHER, grant #248481 (<http://www.pep-pher.eu>) and by SeRC. We would like to thank University of Vienna for providing access to their machine.

References

1. Benkner S, Pllana S, Träff JL, Tsigas P, Dolinsky U, Augonnet C, Bachmayer B, Kessler C, Moloney D, Osipov V (2011) PEPPHER: efficient and productive usage of hybrid computing systems. *IEEE Micro* 31(5):28–41

2. Augonnet C, Thibault S, Namyst R, Wacrenier PA (2011) StarPU: a unified platform for task scheduling on heterogeneous multicore architectures. *Concurr Comput Pract Exper* 23(2):187–198
3. Che S, Boyer M, Meng J, Tarjan D, Sheaffer JW, Lee SH, Skadron K (2009) Rodinia: a benchmark suite for heterogeneous computing. In: *IEEE international symposium on workload characterization (IISWC)*, pp 44–54
4. NVIDIA Corporation (2012) CUBLAS library: NVIDIA CUDA basic linear algebra subroutines. <http://developer.nvidia.com/cublas/>
5. Bell N, Garland M (2012) CUSP library v0.2: generic parallel algorithms for sparse matrix and graph computations. <http://code.google.com/p/cusp-library/>
6. Asanovic K et al (2009) A view of the parallel computing landscape. *Commun ACM* 52(10):56–67
7. Kessler CW, Löwe W (2012) Optimized composition of performance-aware parallel components. *Concurr Comput Pract Exper* 24(5):481–498
8. Li L, Dastgeer U, Kessler C (2013) Adaptive off-line tuning for optimized composition of components for heterogeneous many-core systems. In: *Seventh international workshop on automatic performance tuning (iWAPT-2012)*, Proc. VECPAR-2012 conference, pp 329–345
9. Kicherer M, Buchty R, Karl W (2011) Cost-aware function migration in heterogeneous systems. In: *Proceedings conference on High Perf. and Emb. Arch. and Comp. (HiPEAC)*, pp 137–145
10. Kicherer M, Nowak F, Buchty R, Karl W (2012) Seamlessly portable applications: Managing the diversity of modern heterogeneous systems. *ACM Trans Archit Code Optim* 8(4):42(1–42:20)
11. Alexandrescu A (2001) *Modern C++ design: generic programming and design patterns applied*. Addison-Wesley, Reading
12. Park R (1992) *Software size measurement: a framework for counting source statements*. Software Engineering Institute, Carnegie Mellon University, Pittsburgh, Tech. rep
13. Davis TA, Hu Y (2011) The university of florida sparse matrix collection. *ACM Trans Math Softw* 38(1):1(1–1:25)
14. Ng R, Levoy M, Brédif M, Duval G, Horowitz M, Hanrahan P (2005) *Light field photography with a hand-held plenoptic camera*. Stanford University, Stanford, Tech. rep
15. Augonnet C (2011) *Scheduling tasks over multicore machines enhanced with accelerators: a runtime system’s perspective*. PhD thesis, Université Bordeaux 1
16. Ansel J, Chan C, Wong YL, Olszewski M, Zhao Q, Edelman A, Amarasinghe S (2009) PetaBricks: a language and compiler for algorithmic choice. *Proc Conf on Prog Lang Design and Impl (PLDI)*
17. Wang PH, Collins JD, China GN, Jiang H, Tian X, Girkar M, Yang NY, Lueh GY, Wang H (2007) EXOCHI: architecture and programming environment for a heterogeneous multi-core multithreaded system. In: *Proceedings of conference on programming language design and implementation (PLDI)*, pp 156–166
18. Linderman MD, Collins JD, Wang H, Meng THY (2008) Merge: a programming model for heterogeneous multi-core systems. In: *Proceedings of international conference on architecture support for programming language and Operating Systems, (ASPLOS 2008)*, pp 287–296
19. Huang SS, Hormati A, Bacon DF, Rabbah R (2008) Liquid metal: object-oriented programming across the hardware/software boundary. In: *Proceedings of 22nd European conference on object-oriented programming (ECOOP)*, pp 76–103
20. Wernsing JR, Stitt G (2010) Elastic computing: a framework for transparent, portable, and adaptive multi-core heterogeneous computing. In: *Proceedings of conference on languages, compilers, and tools for embedded systems (LCTES)*, pp 115–124
21. Chafi H, Sujeeth AK, Brown KJ, Lee H, Atreya AR, Olukotun K (2011) A domain-specific approach to heterogeneous parallelism. In: *16th symposium on principles and practice of parallel programming (PPoPP)*, pp 35–46