

Analysis and Mapping for Thermal and Energy Efficiency of 3D Video Processing on 3D Multi-Core Processors

Amit Kumar Singh, *Member, IEEE*, Muhammad Shafique, *Member, IEEE*, Akash Kumar, *Senior Member, IEEE*, Jörg Henkel, *Fellow, IEEE*

Abstract—Three-dimensional (3D) video processing has high computation requirements and multi-core processors realized in 3D integrated circuits (ICs) provide promising high performance computing platforms. However, the conventional approaches to accelerate the computations involved in 3D video processing do not exploit the high performance potential of 3D ICs. In this paper, we propose an application-driven methodology that performs efficient mapping of 3D video applications' components on 3D multi-cores to achieve high performance (throughput). The methodology involves extensive application analysis to exploit the spatial and temporal correlation available in 3D-neighborhood. Afterwards, it leverages the correlation and thermal properties of different 3D-views to perform an efficient mapping of 3D video processing on cores available at different layers of 3D IC. The goal is to optimize energy consumption and peak temperature while meeting the throughput requirement. Experiments show 76% reduction in communication energy along with reduction in peak temperature when compared to approaches exploiting architecture characteristics only.

Index Terms—3D Multi-core, 3D Video, Synchronous Dataflow, design-time analysis, thermal-aware mapping, throughput, interconnect energy.

I. INTRODUCTION

Three-dimensional (3D) video services are envisaged to play an important role towards enhancing the future of several industries such as consumer/entertainment, security, medical imaging and communication. The advancement in 3D video technologies [1] and emerging users' sensation for true 3D-reality have evolved new application domains like Three-dimensional Television (3DTV), 3D-surveillance [2], and 3D video recording on next-generation mobile devices [3]. Recent devices released for 3D recording use two views¹ [4], [5]. However, an increase in the number of views is expected for such upcoming devices to fulfill the emerging market needs.

Manuscript received May 09, 2015; revised September 01, 2015 and November 06, 2015; accepted December 20, 2015. This work was supported in part by the Tier 2 Singapore Ministry of Education Academic Research grant number R-263-000-B33-112, German Research Foundation (DFG) within the Cluster of Excellence Center for Advancing Electronics Dresden (cfaed) and as part of the Transregional Collaborative Research Centre Invasive Computing (SFB/TR 89 <http://invasic.de>), and as part of the priority program Dependable Embedded Systems (SPP 1500 <http://spp1500.itec.kit.edu>).

A. K. Singh is with the Department of Computer Science, University of York, York YO10 5GH, UK (email: amit.singh@york.ac.uk)

A. Kumar is with the Department of Computer Science, Technische Universität Dresden, Dresden, Germany (e-mail: akash.kumar@tu-dresden.de).

M. Shaque and J. Henkel are with the Chair for Embedded Systems, Karlsruhe Institute of Technology, Karlsruhe 76131, Germany (e-mail: muhammad.shafique@kit.edu; henkel@kit.edu).

To address the processing challenges with increased number of views in 3D video encoding, Multiview Video Coding (MVC) [6] standard was devised a couple of years ago. MVC provides up to 50% bitrate reduction (compression) compared to independent coding of different views using state-of-the-art H.264 video coding standard. This is achieved by exploiting temporal and inter-view correlation through multiple block-sized Motion and Disparity Estimation (ME, DE) that in turn significantly increases the computational complexity. The complexity and workload of ME/DE highly depends upon the application specific properties like, picture prediction structure², correlation between frames/pictures³, and motion/disparity contents in the video sequences.

Several efforts have been made to accelerate the ME and DE computation process in order to achieve high throughput. These efforts use either fast ME/DE algorithms [7]–[9] or hardware acceleration [10], [11]. Although these state-of-the-art ME/DE algorithms and hardware accelerations provide significant computation reduction, *they do not exploit full potential of 3D-neighborhood correlation available in spatial and temporal domains*. Further, they target 2D architectures and simply extending them for 3D multi-core architectures to accelerate the computations leads to increased complexity and inefficiency due to significantly different thermal behavior of 3D ICs.

3D integrated circuits (ICs) provide attractive possibilities to implement multi-core systems and are regarded as promising future high performance computing platforms. In 3D ICs, multiple logic layers are stacked vertically and the layers are connected by Through Silicon Vias (TSVs). Such ICs alleviate performance bottleneck problems incurred due to on-chip interconnects that do not scale in proportion to the process technology [12], and are considered as one of the most promising solutions to surmount the interconnect scaling problem [13]. Further, vertical stacking reduces the die size and wire lengths, which results in several advantages such as reduced production costs, communication delays and energies [14]. Thus, 3D multi-cores achieve higher performance and lower power consumption when compared to the traditional 2D counterparts [15], [16], and can be considered as potential

¹Video sequences captured using different cameras.

²a prediction structure defines the direction (i.e. previous or future picture) of finding the best match (i.e. the most correlated prediction block) in the search process.

³frames and pictures are interchangeably used for the same thing.

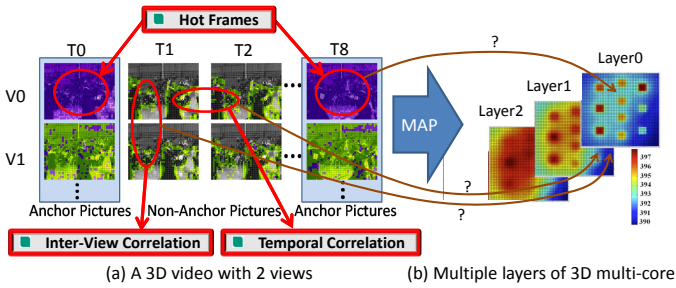


Fig. 1: Exploiting 3D-neighborhood characteristics to map 3D video on 3D multi-core.

computing platforms for complex 3D video processing (encoding).

The communication infrastructure to support communications amongst the various cores in a 3D multi-core is generally considered as network-on-chip (NoC) due to its scalability and high performance [17]–[19]. On-chip interconnection network consumes a significant portion of the entire chip power [20]. However, this portion might vary amongst chip architectures and applications running on them. For example, in Intel’s 80-core teraflops processor [20], the network (routers and links) consumes 28% of the total power while operating at 4 GHz, and this ratio increases to 39% at the maximum operating frequency of 5 GHz. Intel’s latest 48-core chip uses advanced power management technique, but the network still consumes 10% of the total chip power [21]. For video processing, the network power might become quite significant as the application components normally have high communication overhead to stream the data amongst them. Therefore, it is of significant importance to reduce the communication overhead for streaming applications, which will lead to reduced communication time and energy.

In 3D multi-cores, vertically aligned cores of different layers are connected using TSVs, which are shorter than the horizontal links [17]. The reduced interconnect distance between vertically aligned cores leads to smaller resistance and capacitance. Further, due to the reduced interconnect distance, vertical interconnects consume much less energy than horizontal links when transmitting the same amount of data [18]. This facilitates allocating heavily communicating tasks on the vertically aligned cores, i.e. in the same core stack to save the communication time and energy. However, TSVs are drilled through the device of each layer by special techniques and are costly to fabricate. In case of large number of TSVs, the cost of the 3D chip will increase. Further, TSV diameters and pitches are quite large as compared to sizes of regular metal wires. Diameters and pitches are usually around 5-10 μm and 10-20 μm , respectively [22]. So, the number of TSVs will affect the overall chip areas. Therefore, the number of TSVs needs to be controlled during the chip design although they provide increased routing and other benefits. Further, placing active tasks within the same stack increases power density, which may result in serious thermal issues as high temperature affects performance, reliability and lifetime of the system [23]. Therefore, thermal measures are required while accelerating 3D video processing on a 3D multi-core.

A. Motivational Example

A motivational example to map a 3D video with two views (V0 and V1) on a 3D multi-core with three layers is presented in Fig. 1. Each view contains a group of pictures, and spatial (within a video frame), temporal (between frames) and inter-view (between views) correlations exist in 3D-neighborhood. Further, prediction of some of the frames involves heavy computation. Such frames are characterized as *hot frames*. For performing thermal and performance aware mapping, compute intensive (hot) components (views, frames) can be placed on layers close to the heat sink (i.e. on the coolest layer) to achieve a good and balanced thermal profile. For example, in Fig. 1, hot frames T0 and T8 of view V0 can be mapped on the cores of the coolest layer Layer0. Moreover, the highly correlated components on adjacent layers (or close to each other) to minimize the communication overhead towards achieving high performance and low energy consumption. For example, correlated frames T1 and T2 of view V0 can be mapped on cores stacked on top of each other and located in the adjacent layers Layer0 and Layer1. However, a straightforward assignment will not lead to good results due to the inter-view (frame) dependencies. Therefore, *there is a need to balance between computation and communication (correlation) induced thermal effects*.

In short, there is a need to devise a methodology that should first analyze the 3D video processing to identify certain characteristics from 3D multi-core point of view, and then perform mapping by taking application and platform characteristics into account while optimizing for energy consumption and peak temperature.

B. Our Novel Contributions

This paper addresses shortcomings of existing approaches to perform 3D video processing on 3D multi-cores by providing the following contributions:

- 1) An analysis strategy to analyze 3D video processing flow in order to extract the characteristics such as hot/cold and correlated views/frames.
- 2) A mapping strategy to map 3D video processing on 3D multi-core by jointly taking the application and platform characteristics into account towards achieving high performance, thermal balance and energy savings.

Open-source contribution:

- Deriving throughput-constrained Synchronous Dataflow Graph (SDFG) [24] representation of 3D video processing in order to facilitate easier analysis and mapping on 3D multi-core. This also enables open-sourcing of 3D video SDFG. We will make it available online for the community for future research and fair comparisons.

In analysis, the tasks of transformed 3D video as SDFG having long and short computation times are identified as *hot* and *cold* tasks (frames) respectively, and (highly) communicating tasks are identified as (highly) correlated tasks (frames). Dependency between tasks (frames) of different views defines correlation between the views. The mapping strategy systematically maps hot, cold and correlated frames on 3D multi-core

architecture while satisfying the throughput requirement. To the best of our knowledge, this is the first work that addresses mapping of throughput-constrained 3D video processing on 3D multi-core to jointly exploit the characteristics of 3D video and 3D multi-core.

Organization: The remainder of the paper is organized as follows. Section II reviews the literature in the direction of 3D video processing acceleration and application mapping on 3D multi-cores. Section III introduces system model and problem definition. Section IV presents the proposed mapping methodology. The experimental results to evaluate our methodology are presented in Section V. Section VI concludes the paper.

II. RELATED WORK

State-of-the-art efforts to speed up ME/DE computations in 3D video processing employ fast algorithms or hardware accelerations. The fast algorithms in [7] and [8] employ variable search range based on disparity maps and camera geometry, respectively. In [9] and [25], a fast prediction (ME or DE) based on blocks motion intensity and complete DE is proposed. The view-temporal correlation and inter-view correlation have been exploited in [26] and [27] respectively in order to reduce the computational complexity. In [28], algorithm and architecture for disparity estimation with min-census adaptive support is proposed. The hardware designs of ME/DE are also proposed [10], [11]. Although state-of-the-art ME/DE algorithms and hardware accelerations provide significant computation reduction, they do not exploit full potential of 3D-neighborhood correlation available in spatial and temporal domains. Further, they target 2D architectures and several thermal optimization approaches exist for them [29], [30], but simply extending them for 3D multi-core architectures leads to increased complexity and inefficiency due to significantly different thermal behavior of 3D ICs.

Thermal-aware application mapping and scheduling on 3D multi-cores is a well-studied topic [14], [23], [31]–[35]. These approaches perform optimizations at design-time or run-time while trying to minimize hotspots and thermal gradients (spatial, temporal or both).

Run-time approaches generally try to measure or estimate the current temperature distribution in the chip, and take actions based on that in order to minimize hotspots and thermal gradients (spatial, temporal or both). Zhu *et al.* [36] exploit workload power characteristics and processor core thermal characteristics for efficient thermal management. Coskun *et al.* [32] reviewed several dynamic mechanisms such as temperature-triggered Dynamic Voltage/Frequency Scaling (DVFS), clock gating and hot task migration, and proposed a run-time task assignment algorithm that takes the thermal history of cores into account. In [37], [38], concept of thermal herding has been used, where the most frequently switched activity or hot jobs are assigned to the cores close to the heat sink and cool jobs to the cores far from the heat sink. A thermal-aware operating system (OS) level scheduler for 3D multi-cores is proposed in [23]. Kang *et al.* [39] reviewed the work of [23] and introduced peak power and temperature constraints. These methods share the goal of minimizing the

peak temperature and thermal gradients without sacrificing performance too much. However, effect of inter-task communication is not taken into account. Since NoC can dissipate a substantial part of the power budget, which depends upon the network traffic [33], interconnect utilization (energy) should also be taken into account, which has not been considered in most of the aforementioned works.

Design-time mapping approaches aim at finding a thermal-aware mapping by using a model of the physical chip, or by using general knowledge about the thermal behaviour of 3D ICs. In [33], both temperature and communication load are considered, and a genetic algorithm is used to generate static mappings. The design-time mechanisms considering throughput constraint are reported in [34], [35], but they cannot provide efficient mapping solutions for 3D video processing as application characteristics (e.g. hot/cold and correlated frames) cannot be exploited. To summarize, existing mapping approaches perform either application aware mapping by exploiting application characteristics or platform-aware mapping by exploiting platform characteristics.

In contrast to the above strategies, our approach performs thermal-aware mapping of throughput-constrained 3D video processing on 3D multi-cores by exploiting both 3D video (application) and 3D architecture (platform) characteristics. Additionally, our approach considers the effect of TSVs on temperature distribution and power dissipation, and minimizes the communication energy. In case of multiple applications to be mapped and executed concurrently while sharing the system resources, all the tasks can be considered in the mapping process. However, this will need to consider all the possible use-cases (scenarios), where each use-case represent a set of concurrently running applications. The number of use-cases increases exponentially with the number of applications. Further, for each use-case, composability analysis needs to be employed to ensure that near optimal mapping has been achieved for each application in order to satisfy the throughput constraints. To avoid evaluation for huge number of use-cases and their composability analysis, the applications can be mapped and executed one after another without sharing resources. Further, in case 3D multi-core platform is complex, i.e. contains large number of cores, multiple applications can be mapped and executed concurrently. Towards this, a set of cores can be reserved for each application at design-time so that different applications can be mapped and executed into disjoint regions.

III. SYSTEM MODEL AND PROBLEM FORMULATION

A. 3D Multi-core Architecture and 3D IC Model

The *3D multi-core* is modeled as a regular 3D mesh of homogeneous cores connected by a Network-on-Chip (NoC), as depicted in Fig. 2. For the NoC, similar to [14], a hybrid NoC-Bus design is considered, which consists of a regular NoC in the horizontal plane and a multi-drop shared bus (TSVs) to connect the cores within the same stack. Thus, one vertical pillar of TSVs is used for a set of routers that are aligned vertically and cores within the same stack are accessed in a single hop [40]. TSVs are of shorter length

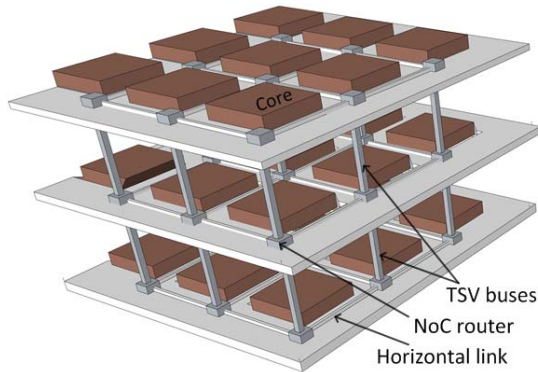


Fig. 2: Example 3D multi-core architecture.

than that of the horizontal links (tens versus thousands of micrometers [14]) and thus they often provide faster and more energy efficient communication than horizontal links [12]–[15], [17], [18]. For communicating amongst vertically aligned cores, the same number of cycles are required (accessed in a single hop [40]) as they are connected by the multi-drop shared bus. In contrast, communication among cores situated in the horizontal plane requires different number of cycles depending upon the hop distance between cores [41]. The core may contain a processing element (PE), for example ARM processor, and local memory (M). The architecture (platform) is represented as a directed graph $PG = (C, V)$, where C is the set of cores and V represents the connections amongst the cores. Each core has active power pa and idle power pi .

For 3D IC, 3D grid model available in the HotSpot thermal simulator has been employed [42]. The application tasks are executed on cores and the execution is tracked at core level. The power dissipated in each core is distributed over the blocks (e.g. processor, memory, router) based on an intra-core power distribution. An example distribution case for a core executing a high instruction level parallelism (ILP) task with low memory traffic is given as 80% power dissipated in the processor block, 10% in the router block, and 10% in the memory. We also have considered similar fine-grained power distribution. This way we can achieve power dissipation in every block. Further, within each considered core, the processor block occupies a significant portion of the total area and dissipates maximum portion of the total core power as mentioned above. Thus, the processor block has maximum power density in terms of power value per area and its temperature determines the peak temperature. To achieve more fine grain power distribution, power consumption in different parts of the processor, e.g. registers, arithmetic logic unit (ALU), etc. can also be considered, but this is orthogonal to our focus, which is at core level. To consider the thermal effect of TSVs, their size, position and material properties are specified in the 3D IC model [42]. The HotSpot simulator is extended to take TSVs into account, where thermal properties (conductance and heat capacity) of grid cells containing TSV material are changed based on grid cell volume occupied by TSV material.

B. 3D Video Application Model

The MVC prediction structure used to perform 3D video processing is employed from [11] and presented in Fig. 3(a). The structure is based on four views (V_0 to V_3). MVC uses ME and DE tools to eliminate the temporal and view redundancies between frames, respectively. In Fig. 3(a), I frames are intra-predicted frames (i.e., no ME/DE is used), some frames use unidirectional prediction or estimation (e.g., 2', 2, 6', 6 as shown in Fig. 3(a)), and rest of the frames use bidirectional prediction with reference frames in at least two directions. The arrows represent prediction directions and frames at the tail side are reference frames to the frames at arrowheads. To facilitate for access points, the video sequence is segmented in Groups of Pictures (GOPs), where frames at borders are known as *Anchor* frames that are encoded with no reference to the previous GOP and others are known as *Non-Anchor* frames.

The MVC prediction structure has been derived to equivalent SDFG, as shown in Fig. 3(b). The SDFG model is represented as a directed graph $AG = (T, E)$, where T is the set of nodes modeling tasks of the application and E is the set of directed edges modeling dependencies amongst the tasks. The nodes of an SDFG are also referred to as *actors*. Each actor represents a frame in the corresponding MVC structure. The execution time of actors (equivalent to ME/DE prediction overheads of frames) and required communication parameters for edges (amount of data required for ME/DE predictions as number of tokens and their size $TokSize[edge]$) are set by analyzing the execution behavior of MVC (Fig. 3(a)) towards achieving the same execution behavior of equivalent SDFG model (Fig. 3(b)). Some reference data for two views (View0 and View3) of Ballroom video sequence is shown in Fig. 3(c). For different frames (actors) to be predicted (Pred.) in a view, one or multiple frames from the same or other views are used as reference (Ref.) frames. For example, frame C uses frames I and A as reference frames. The amount of transferred data (bytes) required from various frames to predict a frame are provided in the last column. These data values determine volume of data on the edges and computation time of an actor. The shown data values are for the worst-case prediction (involving maximum prediction to encounter fluctuations in the data) and using computation times according to the same helps to model the worst-case behavior of the 3D video processing. The 3D video processing is also characterized by throughput constraint Γ and the same has been incorporated in the SDFG model.

C. Energy Consumption Model

Communication energy between the communicating cores c_i and c_j (required to predict the current frame by reference frames) depends upon data volume $data(c_i, c_j)$ and energy required to transfer one bit of data $E_{bit}(c_i, c_j)$ between the cores. $E_{bit}(c_i, c_j)$ depends upon the energy required for horizontal and vertical links traversal and the energy consumed in

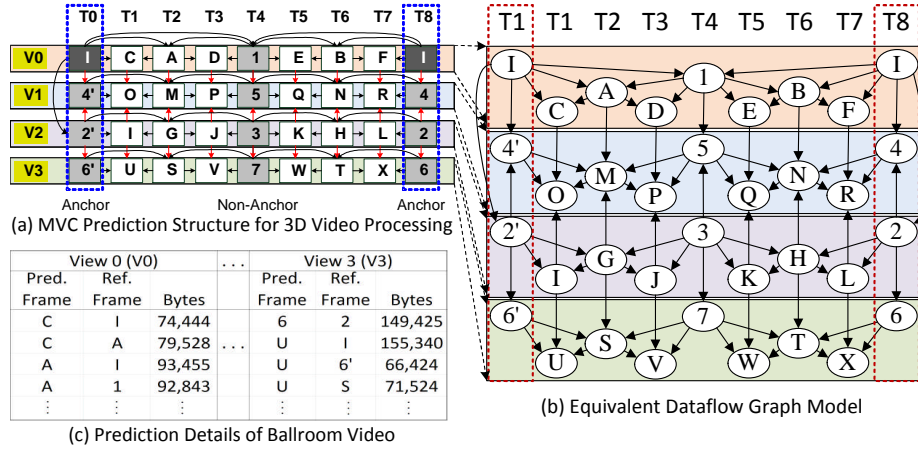


Fig. 3: 3D video processing with four views and equivalent SDF graph model.

routers between cores c_i and c_j , and is computed as follows.

$$E_{bit}(c_i, c_j) = (E_{bit}^{horizontal} \times horizontal_{hops}(c_i, c_j)) + (E_{bit}^{vertical} \times vertical_{hops}(c_i, c_j)) + (E_{bit}^{router} \times nrRouters(c_i, c_j)) \quad (1)$$

In Equation 1, $E_{bit}^{horizontal}$ and $E_{bit}^{vertical}$ are the energy required to transfer one bit per hop in the horizontal and vertical directions, respectively. $horizontal_{hops}(c_i, c_j)$ and $vertical_{hops}(c_i, c_j)$ are the number of hops between cores c_i and c_j in the horizontal and vertical directions, respectively. E_{bit}^{router} is energy consumed in a router and $nrRouters(c_i, c_j)$ is the number of routers between cores c_i and c_j .

The values of $E_{bit}^{horizontal}$, $E_{bit}^{vertical}$, and E_{bit}^{router} are derived for the 90-nm process technology node (details is Experimental Setup (Section V-A)). The router is composed of input FIFO buffers, a fully connected crossbar and an arbiter. At the input, a 3-place FIFO buffer is used. To aggregate the inputs to the outputs, a 6×6 full connected crossbar (one extra port to connect to the vertical multi-drop shared bus) with multiplexers is used. The arbitration of the router occurs at granularity of words and the routing follows source routing, i.e., path information from source to destination is contained in the header. The details of the router design are available in [43]. The derived value of E_{bit}^{router} takes such router structure into account.

The communication energy E_{comm} is estimated by summing over all communicating task pairs (edges).

$$E_{comm} = \sum_{\forall comm-cores} data(c_i, c_j) \times E_{bit}(c_i, c_j) \quad (2)$$

Computation energy required to process all the actors (ME/DE computations) is estimated as follows.

$$E_{comp} = \sum_{\forall a \in T} a_{ExecTime} \times pa \quad (3)$$

where $a_{ExecTime}$ is execution time of actor a and pa is the active power dissipation of core executing actor a .

Total energy consumption is measured as sum of E_{comp} and E_{comm} . The static energy consumption depends upon the static (leakage) power consumption of the cores, which

is assumed as a fixed offset. In this paper, we purposely do not account for the power-gating to stay orthogonal to other low-power techniques, therefore, the leakage power will stay constant throughout all of our experiments. Thus, to purely show the impact of the proposed techniques, the results only show the dynamic power consumption. Please note that, any state-of-the-art power-gating technique can be employed in our architecture after the mapping decisions are taken, i.e., power-gating the idle cores. Further, since we consider homogeneous architecture, the computation performed in any part of the architecture will consume almost the same energy. Therefore, the importance of the communication energy becomes the primary focus for optimization/reduction.

D. Mapping Problem: 3D Video Processing on 3D Multi-core

Given SDFG model of a video sequence $AG = (T, E)$ with throughput constraint Γ and 3D multi-core architecture $PG = (C, V)$

Find efficient actors to cores mapping to simultaneously optimize for peak temperature ($PeakTemp$) and communication (interconnect) energy consumption (E_{comm})

$$PeakTemp \times E_{comm} \quad (4)$$

subject to

$$\tau \leq \Gamma \quad (5)$$

where τ is the through obtained as a result of the actors to cores mapping and Γ is the throughput constraint.

For a mapping, throughput computation and temperature estimation are in the orders of several milliseconds and seconds, respectively. Therefore, evaluation of all the possible mappings to identify the best mapping (in terms of peak temperature and energy consumption while satisfying the throughput constraint) is expected to take several days. To overcome the evaluation time bottleneck, heuristic based approaches relying on various cost parameters pertaining to application and architecture characteristics need to be applied to identify an efficient mapping rapidly. However, the identification and extraction of the exact required application and architecture characteristics is challenging specially for large problems such as 3D video.

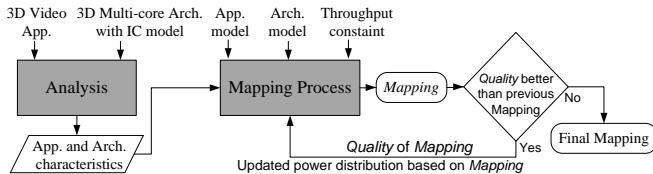


Fig. 4: Offline analysis and mapping process.

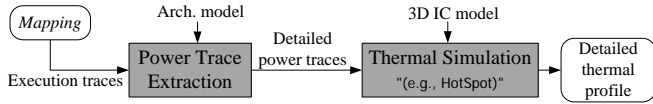


Fig. 5: Temperature estimation flow.

We employ the required cost function to the mapping process that considers desired optimization parameters pertaining to application and architecture characteristics.

IV. PROPOSED MAPPING METHODOLOGY

An overview of our offline mapping identification and temperature estimation flow is presented in Fig. 4 and Fig. 5, respectively. The mapping flow first performs offline *Analysis* of the application and 3D architecture to identify their characteristics that are used to identify efficient thermal-aware mapping (*Mapping Process*). The offline analysis considers *worst-case* computation and communication requirements (as described in Section III-B) of 3D-video coding prediction structures for different test video sequences provided by the standardization committee. Such consideration facilitates *worst-case* online behavior modeling of the 3D video processing. For the identified mapping for a video sequence, thermal analysis is performed (*Temperature Measurement*) to identify the temperature distribution across different components of the 3D multi-core architecture. In order to perform online video processing for a video sequence, its offline identified mapping is used to configure the platform and then the processing starts. The details of *Offline Analysis*, *Mapping Process*, *Temperature Measurement* and *Online Processing* steps are described subsequently.

A. Offline Application and Architecture Analysis

1) *Application Analysis*: The 3D video application analysis is performed to identify hot, cold and correlated views and frames (actors). Without loss of generality, let us consider an example case where the 3D video processing application contains 4 views and each view contains 9 actors (ME/DE computations) (Fig. 3). Generally, the same objects in a 3D scene are typically present in different views and motion perceived in one view is directly related to that of the neighboring views [26]. Moreover, the disparity of one given object perceived in two cameras (views) remains the same at various time instances when just translational motion occurs [8] [26] [27]. These observations indicate that there is a high correlation in the 3D-neighborhood that can be exploited during the ME/DE computations. These observations are exploited to set the ME/DE (actor) computation overheads and communication overheads (correlation) between actors. The frames encoded with high effort are referred to as *hot* frames (actors) and

ALGORITHM 1: View Level Application Analysis

Input: SDFG model of 3D video.

Output: Criticality of views and their actors, correlation between views.

```

for each view  $v$  in 3D video do
   $weight[v] = \sum_{actors \in v} ExecTime[actor]$ ;
   $criticality[v] = weight[v]/nrActors$ ;
  for each actor  $a$  in  $v$  do
     $criticality[v][a] = ExecTime[a]$ ;
  end
  Find connected views  $conn\_views$  with  $v$ ;
  for each view  $v_{conn}$  in  $conn\_views$  do
     $correlation[v][v_{conn}] = \sum_{edges \in v \rightarrow v_{conn}} TokSize[edge]$ ;
  end
end

```

have high computation time, whereas *cold* actors have low computation time. High communication overhead between actors reflects high correlation between them and vice-versa.

We propose two approaches to extract the view and frame level application characteristics to be used by the mapping process.

View Level Analysis

The view level analysis (VLA) approach is presented in Algorithm 1. The algorithm takes application model as input and provides criticalities of views and correlation between them. The criticality of actors within a view is also provided. For each view in a 3D video, first, view weight ($weight[v]$) is computed, which determines hotness of the view. A view having high weight is considered as a hot view. Then, the view criticality ($criticality[v]$) is computed by dividing the view weight by the number of actors in the view. The criticality of an actor within a view is determined by its computation overhead (execution time). Thereafter, correlation between views is calculated by adding the token sizes ($TokSize$) of each edge (representing data volume) present between the views. The $TokSize$ is extracted from the application model.

Frame Level Analysis

The frame level analysis (FLA) is performed in the similar way as that of VLA, but criticalities and correlations are identified at the frame levels. For the example 3D video in Fig. 3, there are total 9 frames (T_0 to T_8) and each frame contains four actors (e.g., C, O, I and U in frame T_1). In this analysis, similar steps as in Algorithm 1 are adopted and frame level processing by considering frames and connected frames is performed to calculate criticalities of frames, their actors and correlation between frames.

The view and frame level analysis approaches extract application characteristics at two different granularities. Based on suitable scenarios, one can provide better characteristics than other to facilitate for efficient mapping. The effect of utilizing such various characteristics in the mapping process is demonstrated in the next section.

The **complexity** of the analysis depends upon the number of operations to be performed for the criticality and correlation calculations. For n actors, and e_v and e_f edges between views and frames respectively, the complexity of VLA and FLA is $O(n \times e_v)$ and $O(n \times e_f)$, respectively. Since the number of edges between frames is higher in FLA than that of VLA (Fig. 3), FLA has slightly higher complexity.

ALGORITHM 2: Architecture Analysis Algorithm

Input: 3D IC model, Multi-core architecture model
 $PG = (C, V)$, total chip power $P \in \mathbb{R}$, max. # of iterations
 $Iter_{max} \in \mathbb{N}$, terminating condition $\delta \in \mathbb{R}$.

Output: Power ratios of cores (R_c for each core $c, c \in C$).
Initialize power ratio for each core as $R_c = 1/N_{cores}$;
 $iter = 0, Temp_{max} = 500$;

repeat

- $Temp_{prev_max} = Temp_{max}$;
- Generate power traces for all blocks of each core;
- Simulate steady state temperature distribution.;
- Find peak temperature in each core $Temp_{peak,c}$;
- Find average $Temp_{avg}$ and max. $Temp_{max}$ chip temperature ;
- for each core $c \in C$ do**
- $d = (Temp_{peak,c} - Temp_{avg})/Temp_{avg}$;
- if $Temp_{peak,c} > Temp_{avg}$ then**
- $R_c = R_c * (1.0 - (\gamma * d))$); //decrease power ratio
- else**
- $R_c = R_c * (1.0 + (\gamma * d))$); //increase power ratio
- end**
- end**
- Renormalize power ratios for each core R_c ;
- $iter ++$;

until $(Temp_{prev_max} - Temp_{max}) \geq \delta$ **AND** $iter \leq Iter_{max}$;

2) *Architecture Analysis:* The 3D architectures have been analyzed to observe their thermal and power dissipation characteristics [23], [32], [35]. In 3D architecture, if the same amount of power has to be dissipated by all the cores on different layers, the cores close to heat sink will dissipate faster than others and thus high temperature gradient and peak temperature will be achieved. The architecture analysis approach presented in Algorithm 2 is used to extract the platform thermal characteristics as the power distribution amongst the cores such that peak temperature and temperature gradients are minimized. The algorithm finds power ratios of cores based on steady state temperature distribution resulting from earlier power distribution. The power ratio of a core R_c is defined as the ratio of power dissipated in the core c and total chip power P . The approach decreases the power ratios of cores having peak temperature greater than the average temperature and vice-versa until temperature difference amongst cores is reached to a very low value. The increment and decrement is done in small steps by setting a low integer value of the adaptation constant γ . In agreement to general observations [23], [32], [35], the cores close to the heat sink get higher power ratios. The power ratio of every core is passed as the architecture characteristic to the mapping process.

The **complexity** of architecture analysis (Algorithm 2) depends on steady state temperature simulation and number of iterations. The simulation time depends on the spatial resolution and the number of layers and it takes around 2 minutes on a 1.70GHz Intel i5 CPU for an IC with 3 layers and resolution of 32×32 . The algorithm usually converges in 5-10 iterations and thus the whole analysis takes up to 20 minutes.

B. Offline Mapping Computation

The steps followed by the proposed mapping algorithm are described in Algorithm 3. The algorithm incorporates thermal awareness to find an efficient final mapping by exploiting

ALGORITHM 3: Thermal-aware Mapping Algorithm

Input: AG, PG, Γ , Characteristics of AG and PG .
Output: Final mapping FM .
Quality of previous mapping $Q_{PM} = 0$;
//Mapping Process
Find criticalities of actors $\in AG$ by using AG characteristics;
Sort all actors $\in AG$ in descending order of criticality;
for each sorted actor $a \in AG$ do- Find cost of each core $c \in PG$ by using AG & PG characteristics, as $cost(a, c)$ (see Eqn. 8);
- Sort all cores $\in PG$ in ascending order of $cost$;
- for each sorted core $c \in PG$ do**
- if actor a can be bound to core c then**
- Assign actor a to core c and incoming/outgoing edges of a to connections to construct $Mapping M$;
- break**;
- end**
- end**
- end**
- Compute τ and quality of M as Q_M (by Eq. 4) ;
- if $\tau < \Gamma$ then**
- if $Q_M > Q_{PM}$ then**
- $Q_{PM} = Q_M$;
- Repeat the Mapping Process with updated power distributions on cores based on M ;
- else**
- $FM = M$; **break**;
- end**
- end**

(thermal) characteristics of application and architecture, and updated power distributions based on intermediate mapping. First, the criticalities of all actors are computed and the actors are sorted in descending order of their criticality. The criticality of an actor $cric[a]$ is calculated by exploiting application characteristics extracted by VLA or FLA (described earlier) as follows.

$$cric[a]_{VLA} = k_1 * criticality[v] + k_2 * correlation[v][v_{conn}] \quad (6)$$

$$cric[a]_{FLA} = k_1 * criticality[f] + k_2 * correlation[f][f_{conn}] \quad (7)$$

where v and f are respectively the view and frame containing actor a . v_{conn} and f_{conn} are the connected view and frame from v and f , respectively. It should be noted that either $cric[a]_{VLA}$ or $cric[a]_{FLA}$ is used based on the employed application analysis approach *VLA* or *FLA*.

Sorting of actors as described earlier helps to handle their mapping systematically, for example, first actors from hot views and then from correlated views by giving a higher value to k_1 than k_2 , and vice-versa. Then, by following the steps in Algorithm 3, sorted actors are assigned one-by-one on the cores that can support them and incur minimum assignment cost computed as follows.

$$cost(a, c) = c_1 * LB(a, c) + c_2 * PCE(a, c) + c_3 * ACE(a, c) \quad (8)$$

where $LB(a, c)$, $PCE(a, c)$ and $ACE(a, c)$ represent normalized processor load, cost for platform characteristic exploitation (PCE) and cost for application characteristic exploitation (ACE) when actor a is bound to core c , and c_1, c_2 and c_3 are the weights given to different optimization criterion.

The ACE uses view or frame level characteristics exploited from VLA/FLA, where average latency of all edges to/from a is minimized by mapping connected actors close to each other in order to exploit inter-view (frame) correlations. This results in reduced communication overhead/energy. For view and frame level exploitation, ACE is referred to as AVCE and AFCE, respectively.

The PCE uses cost for the power ratio balancing of core c ($PRTB(a, c)$) and cost for the power ratio balancing of stack containing c ($PRSB(a, c)$), and is obtained by adding both the costs. $PRTB(a, c)$ is computed by dividing estimated power ratio of core c (when assigning a to c) by power ratio of c (R_c) suggested by the architecture analysis. Similarly, $PRSB(a, c)$ is computed by dividing estimated power ratio of stack s (containing c) when binding a to c by power ratio of stack R_s that is computed by summing up power ratios of all the cores in s . A core stack s consists of a set of cores having the same horizontal position in different layers. Since strong thermal correlation exists between vertically adjacent cores [23] [14] [35], it might be beneficial to consider power ratios of stacks ($PRSB$) to distinguish in the incurred costs when deviations from original power distribution in the vertical and horizontal directions are the same. This helps to achieve better results from thermal perspective.

The mapping algorithm assigns all actors ($\in AG$) to cores ($\in PG$) and connections to memories inside cores or interconnect links. The mapping process repeats itself to identify a better quality of throughput satisfying mapping by considering updated power distributions on cores based on the current mapping (Fig. 4). This iterative refinement process is expected to lead to a high quality mapping. The quality of the mapping, Q_M is computed by employing Equation 4, which requires peak temperature and communication energy, and computes the quality as the product of peak temperature and communication energy. If the product value is low, then the mapping is considered to have a good quality. This quality is used to compare the mapping with previous evaluated mapping (Q_{PM}) in the iterative refinement process. The throughput, energy consumption and temperature estimations for a mapping are done as follows.

The **throughput computation** is performed by employing the technique of [44]. However, any throughput computation technique can be employed, which is orthogonal to our focus. In [44], the throughput for a mapping is computed by taking the resource allocations into account. First, static-order schedule for each core is constructed that orders the execution of bound actors. A list-scheduler is used to construct the static-order schedules for all the cores at once. Then, all the binding and scheduling decisions are modelled in a graph called binding-aware SDFG. Finally, self-timed state-space exploration of the binding-aware SDFG is performed to compute the throughput, which is the inverse of the long term period, i.e., the average time needed for one iteration of the application. In doing so, the mathematical model used in [44] takes computation, communication, latency for data arrival and jitter into account. This indicates that throughput depends on the mapping. The dependency of throughput on the mapping has also been well studied in [41].

ALGORITHM 4: Online Video Processing

Input: Video Sequence, PG , Offline computed mappings.

Output: mapping to start online processing.

Select the *mapping* from offline computed mappings for the video sequence;

Configure platform PG based on *mapping*;

Start video processing;

The **energy consumption** is computed by employing the approach described in Section III-C. The **temperature estimation** flow is presented in Fig. 5. In order to get the detailed temperature profile resulting from a mapping, first the execution traces are generated. The execution trace of each core represents its active and idle time intervals. For active and idle intervals, the core is assumed to consume active and idle powers, respectively. The power traces are used to simulate the temperature with the modified HotSpot thermal simulator (further details in the next section).

C. Online Video Processing

To perform on-line processing for a required video sequence, first, the actors of the application are loaded (configured) onto the platform resources based on the offline computed final mapping for the video sequence and then real processing starts. The online video processing for a video sequence follows Algorithm 4. It selects the offline computed mapping for the video sequence and the platform is configured based on the same in order to start the video processing. For a video sequence, the online mapping is performed once at the application startup and has a small overhead (quantitative description in next Section).

V. PERFORMANCE EVALUATION

A. Experimental Setup

The proposed thermal-aware mapping methodology has been implemented as an extension of the publicly available SDF³ tool set [45]. As a benchmark to evaluate the quality of the methodology, models of four evaluated video sequences Ballroom, Vassar, Crowd and Kendo (recommended by joint video team in multiview test conditions [11]) have been considered.

Target 3D multi-core platforms contain different number of cores, where active power of each core is set to 1.5W and idle power is 10% of the active power [46]. The size of each core is set to 2 mm \times 2 mm, which is derived by extrapolating the size of the core in 90-nm technology node. The heatsink is connected (via a heat spreader) to the bottom layer and has a thickness of 200 μ m and the other active (power dissipating) layers are assumed to be thinned to 50 μ m for better heat conductivity [14], [23]. Between two active layers, a 10 μ m thin layer containing thermal interface material (TIM) is used. For vertical communication, each TSV bundle contains 8 \times 9 TSVs. In our considered core size (2 \times 2 mm²), 1 \times 1 mm² is allocated for the router and rest of the area for processor and memory. The area allocated for router (1 \times 1 mm²) is used to place the 8 \times 9 TSVs and is sufficient to accommodate them. Some other important physical properties of the 3D IC model are summarized in Table I.

TABLE I: 3D IC parameters

Parameter	Value
Technology node [nano m]	90
Each core size [$mm \times mm$]	2×2
TSV diameter [μm]	10
TSV pitch [μm]	20
Horizontal hop delay [time-units]	2
Vertical hop delay [time-units]	1
Bottom layer thickness [μm]	200
Non-bottom layer thickness [μm]	50
TIM layer thickness [μm]	10
Heatsink side/thickness [mm]	$14 \times 14 \times 10$

TABLE II: Thermal simulation parameters

Parameter	Value
Silicon thermal conductance [$W/(m \cdot K)$]	150
Silicon specific heat [$J/(m^3 \cdot K)$]	$1.75 \cdot 10^6$
TIM thermal conductance [$W/(m \cdot K)$]	4
TIM specific heat [$J/(m^3 \cdot K)$]	$4 \cdot 10^6$
TSV thermal conductance [$W/(m \cdot K)$]	300
TSV specific heat [$J/(m^3 \cdot K)$]	$3.5 \cdot 10^6$
Convection resistance to ambient [K/W]	3.0
Heatsink thermal conductance [$W/(m \cdot K)$]	400
Heatsink specific heat [$J/(m^3 \cdot K)$]	$3.55 \cdot 10^6$
HotSpot grid resolution	32×32
Temporal resolution [μs]	10
Ambient temperature [K]	300

Temperature is estimated by employing extended HotSpot thermal simulation tool [42]. To estimate temperature resulting from a mapping, an execution trace of 0.5s is generated. The execution patterns are periodic with a period much shorter than 0.5s, and thus longer simulations become obsolete. Power traces for every block are derived from the execution trace and the architecture specification. First, a steady state simulation is performed to find a representative initial temperature distribution. Then, the transient simulation is performed. The HotSpot thermal simulation parameters are listed in Table II.

In our 3D IC model, the interconnect energy consumption is computed by employing Equation 2 as described in Section III-C. Table III lists the parameters used to compute the interconnect energy consumption. For consistency, all the parameters are considered for the 90-nm process technology node such that energy consumption can be computed accurately. First, the horizontal link energy per bit, $E_{bit}^{horizontal}$, is derived as in [14]. Then, the vertical link energy per bit, $E_{bit}^{vertical}$ (E_{bit}^{TSV}), is calculated by using the parameters from ITRS [22]. For the same process technology node, E_{bit}^{router} is approximately 70% of $E_{bit}^{horizontal}$, as indicated in [43]. E_{bit}^{TSV} is only 7.5% of $E_{bit}^{horizontal}$, providing substantial space for communication energy optimization by exploiting the links in the vertical direction.

We present results obtained from our approach and compare them with relevant existing methodologies, as abbreviated in Table IV. The *LB* approach try to balance load (power) on the cores to achieve good thermal balance and is employed by setting $c_1 = 1$, $c_2 = 0$ and $c_3 = 0$ in Equation 8. The *AA* and *PA* mapping approaches exploit application and platform characteristics, respectively. The *AA* approach looks the communicating actors and tries to map them on the same or neighboring cores, whereas *PA* approach tries to map hot and cold actors on layers having high and low heat dissipation capabilities, respectively. The resulting mapping

TABLE III: Interconnect Energy Computation parameters

Parameter	Value
$E_{bit}^{horizontal}$ [pJ]	0.127
$E_{bit}^{vertical}$ or E_{bit}^{TSV} [pJ]	9.56×10^{-3}
E_{bit}^{router}	70% of $E_{bit}^{horizontal}$

TABLE IV: Approaches considered for comparison

Approaches	Abbreviation	References
Load Balanced mapping	LB	[23]
Application (App.) Aware mapping	AA	[47]
Platform (Plat.) Aware mapping	PA	[35]
App. views & Plat. characteristics exploitation	AVCE+PCE	Proposed
App. frames & Plat. characteristics exploitation	AFCE+PCE	Proposed

obtained by the *PA* approach is not further optimized, i.e., there is no iterative refinement to achieve a better quality of mapping. Our approaches *AVCE+PCE* (or *PCE+AVCE*) and *AFCE+PCE* (or *PCE+AFCE*) exploit respectively views and frames related applications' characteristics along with the exploitation of architecture characteristics and are carried out by setting appropriate constants ($c_1 = 0$, $c_2 = 1$ and $c_3 = 1$) in Equation 8. Further, our approaches perform iterative refinement to optimize the mapping quality in terms of peak temperature and energy consumption.

B. Results for Different Video Sequences

Fig. 6 shows interconnect power consumption and peak temperature when employing different mapping approaches to map the four considered video (application) sequences on a $4 \times 3 \times 3$ mesh architecture. The interconnect power is estimated as the average communication energy per second. A couple of observations can be made from Fig. 6. 1) *AA* approach achieves low interconnect power consumption due to mapping communicating actors close to each other, but results in high peak temperature due to heat stacking in a particular region. The interconnect power consumption by *AA* and our approaches *PCE+AVCE* and *PCE+AFCE* are almost the same. 2) *PA* approach results in lower peak temperature compared to *LB* due to platform characteristics exploitation. 3) Our approaches *PCE+AVCE* and *PCE+AFCE* results in low energy consumption and peak temperature as they exploit both the application and architecture characteristics and perform iterative refinement to achieve a better quality of mapping leading to lower energy consumption and peak temperature when compared to other approaches. Moreover, the *PCE+AVCE* shows higher reduction in energy and peak temperature when the application contains higher inter-view

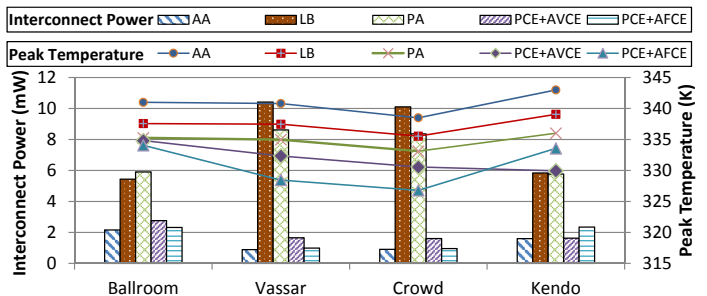


Fig. 6: Interconnect power and peak temperature for different mapping approaches across different video applications.

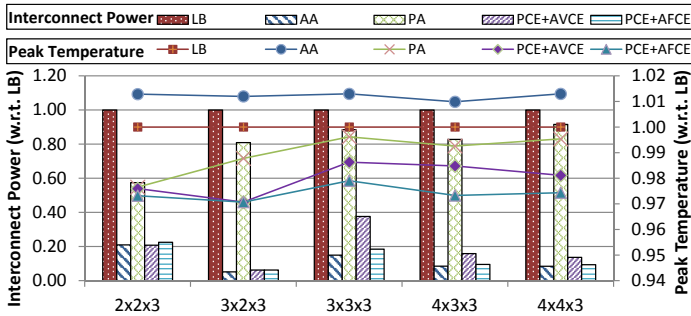


Fig. 7: Interconnect power and peak temperature for different platforms.

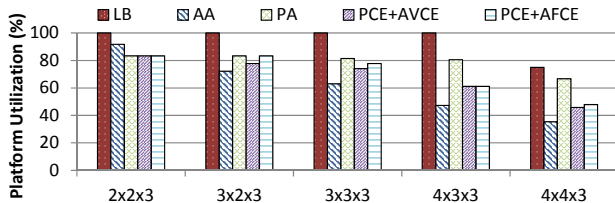


Fig. 8: Utilization with varying platform sizes.

correlations, whereas in case of higher inter-frame correlation *PCE+AFCE* performs better. On an average, our methodology *PCE+AFCE* reduces interconnect power consumption by 76% and average peak temperature by 4°C when compared to *PA* that provides good results for both power consumption and peak temperature. The reduction in peak temperature is not significant as *PA* already tries to achieve low peak temperature by exploiting architecture characteristics.

We have measured computation energy as well in order to observe its contribution to the total energy consumption. For the considered video sequences, on an average, the ratio of computation and total energy consumption indicates that the computation energy contributes 74% to the total energy consumption and thus the communication energy contribution is 26%. However, since computation energy by all the approaches remains the same due to computations performed in the homogeneous cores, communication energy optimization is the primary focus of our proposed approach. Computation energy can be optimized by employing the dynamic voltage and frequency scaling (DVFS) on cores [48], but DVFS is not our focus and orthogonal to our approach. Considering above contributions, our methodology *PCE+AFCE* reduces total energy consumption by approximately 6.3% when compared to *PA*.

C. Results at Varying Platform Sizes

We analyzed the effect of considering various platforms for the applications on the reduction in energy consumption and peak temperature. Fig. 7 shows interconnect power consumption and peak temperature at various platforms for “Vassar” video when different approaches are employed. The shown values are normalized with respect to *LB* that leads to worst interconnect power consumption. However, *LB* provides lower peak temperature than *AA*. It can be observed that our approaches *PCE+AVCE* and *PCE+AFCE* outperform other approaches if both interconnect power and peak temperature

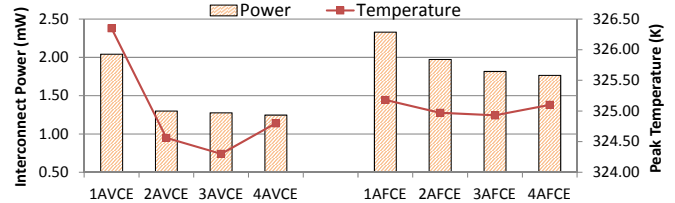


Fig. 9: Effect of correlation exploitation and Energy consumption and peak temperature.

need to be optimized for different platforms. Therefore, for any chosen platform that might be required to perform computations, our approach can be employed to achieve good results in terms of interconnect power and peak temperature.

Fig. 8 shows platform utilization in terms of percentage usage of available number of cores for varying platforms. It can be observed that *LB* shows utilization of maximum number of cores. The utilization is 100% up to the platform size $4 \times 3 \times 3$. For higher size platforms, *LB* uses a maximum of 36 cores as the number of actors in the 3D video application and thus utilization decreases. *AA* shows minimum utilization in most of the cases as it tries to use minimum number of cores by placing communicating actors on the same or neighboring cores. The platform utilization by other approaches that exploit platform characteristics is lower than *LB* as they try to use cores on layers close to the heat sink (cool layers) and in turn cores of hot layers are avoided to be used. However, *LB* leads to high interconnect power and peak temperature as described earlier.

D. Energy-Temperature Trade-off Analysis

The energy-temperature trade-off points are obtained by exploiting varying amount of correlation between views or frames. The amount of correlation increases by giving high weight to *ACE* (*AVCE* or *AFCE*) in Equation 8. Fig. 9 shows the effect of correlation exploitation available at view (*AVCE*) and frame (*AFCE*) levels for mapping Ballroom video on $3 \times 2 \times 3$ platform. For higher correlation exploitation, the weight given to *ACE* in Equation 8 is increased. For example, weight c_3 is varied from 1 to 4 when employing *AVCE* and *AFCE*, as shown on the horizontal axis. A couple of observations can be made from the figure. 1) Energy savings increase with the amount of correlation exploitation for both view/frame levels. 2) Peak temperature first decreases and then increases in both the cases. The decreasing trend is obtained as a better thermal balanced mapping is being found by exploiting the application characteristics. An increase in temperature takes place as higher correlation exploitation tries to map communication tasks on the same core or stack, which results in non-uniform heat dissipation and thus higher temperature.

E. Performance Improvement Over 2D Architectures

The proposed approach can be applied to various 2D multi-core architectures and obtained results can be compared with respect to the 3D architectures containing the same number of cores. Fig. 10 shows performance (throughput) improvement for different video sequences by 3D architectures over 2D

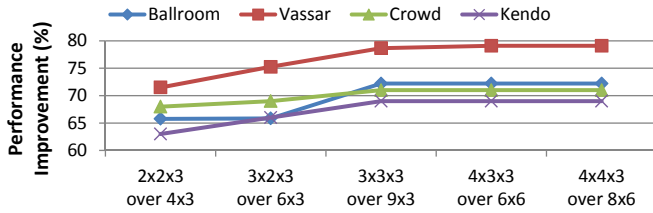


Fig. 10: Performance improvement by 3D over 2D architecture.

architectures when our approach PCE+AVCE is employed. For fair comparisons, the number of cores in both the 2D and 3D architectures are kept the same. The improvements are obtained mainly due to the use of vertical links available in 3D architectures. Vertical links implemented using TSVs have shorter length than horizontal links and thus they provide reduced interconnect distance between vertical adjacent cores. This often leads to faster and more energy efficient communication compared to horizontal links. Further, in 3D architecture, we have more neighbors and hence fewer hops for overall communication leading to lower latency and in many cases higher throughput. The faster communication leads to low communication time, resulting in reduced overall application execution time. Thus, improvements in the applications throughput are achieved.

It can be observed that first performance improvement increases with number of cores due to better usage of cores and then becomes consistent as the number of used cores remains constant. On average, 70% improvement in throughput is obtained compared to a 2D SoC.

F. Generalization: Results for Streaming Multimedia Applications

Fig. 11 shows interconnect power consumption and peak temperature when employing different mapping approaches to map streaming multimedia applications containing different number of actors on a $2 \times 2 \times 3$ mesh architecture. For applying our approach to streaming multimedia applications, criticality of actors and correlation between them are exploited as application characteristics and platform characteristics are exploited as earlier. A couple of observations can be made from Fig. 11. First, the interconnect power and peak temperature are quite less when compared to 3D video sequences (Fig. 6) that have complex structure with 36 actors. Second, interconnect power and peak temperature increases with the application size (number of actors) as processing overhead in the network and cores increases. It can be seen that our approach PCE+ACE provides good results to optimize for both interconnect power and peak temperature when compared to other approaches for streaming multimedia applications as well.

For the considered multimedia applications, on an average, the computation energy contributes 72% to the total energy consumption, leaving 28% contribution as the communication energy.

G. Offline overhead

The offline overhead depends on time to find a mapping, its thermal simulation and number of iterations to identify

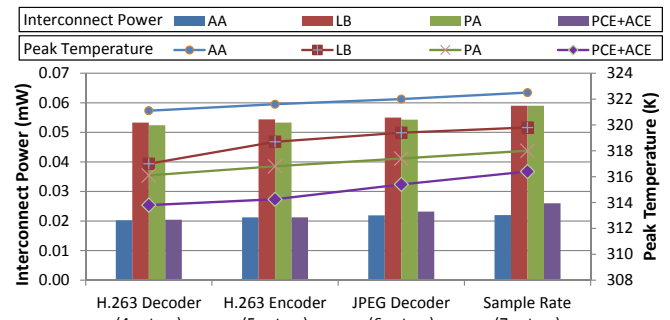


Fig. 11: Interconnect power and peak temperature for different streaming multimedia applications.

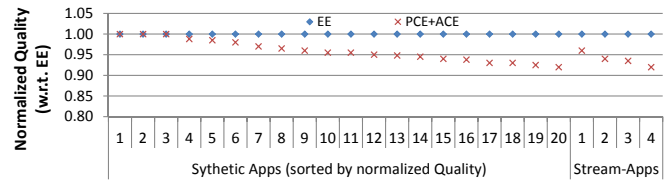


Fig. 12: Quality of mappings by our approach and exhaustive exploration.

the best mapping in terms of peak temperature and energy consumption. The time to find a mapping is quite small and is in the order of a few milliseconds. The thermal simulation overhead is high and depends upon the spatial grid resolution and the number of layers in the 3D multi-core architecture. Usually, it takes several minutes to find the best mapping due to iterative exploration, where each iteration takes around 2 minutes on a 1.70GHz Intel i5 CPU (single threaded) when a grid resolution of 32×32 and 3 layers in the multi-core architecture are considered. The proposed methodology converges in 3-8 iterations for different applications and architectures. This results in a total offline overhead of up to 16 minutes.

H. Online overhead

The on-line processing for a video sequence starts after the actors of the video application are configured on the platform resources. This configuration is performed once at the application startup for a given video sequence and has a small overhead. For example, PCE+AFCE approach takes a configuration time of 18 milliseconds to configure Ballroom video sequence on a $3 \times 3 \times 3$ architecture. For other approaches and architectures, the overheads are of similar orders.

I. Deviation from optimal mapping

To find the optimal mapping, an exhaustive exploration (EE) approach (e.g., [49]) has been employed, which evaluates all the possible mappings (actors to cores allocations) to select the best (optimal) quality mapping for temperature and energy. In order to restrict the evaluation for an application within few hours, small size synthetic and real-life applications are considered to be mapped on a small size $2 \times 2 \times 2$ architecture. Fig. 12 shows the quality of the mappings for 20 synthetic applications (containing 4 to 7 actors) and 4 streaming multimedia applications (1: H.263 Decoder, 2: H.263 Encoder, 3:

JPEG Decoder, 4: Sample Rate Converter) by our approach PCE+ACE w.r.t. the optimal mappings achieved by EE. It has been observed that loss in quality of mappings by our approach is more when the number of actors increases. As can be seen from Fig. 12, our approach provides optimal mapping for some of applications. For the remaining applications, the quality is comparable to EE and maximum deviation is less than 9%.

VI. CONCLUSION

We present a novel methodology to map 3D video processing on 3D multi-core platforms. We show that the methodology exploits application and platform characteristics towards achieving energy savings and reduction in peak temperature while satisfying the throughput requirement of the application. The experimental results indicate that our approach can be employed to variety of platforms to achieve high quality results. In future, we plan to consider heterogeneous cores to be integrated in the 3D multi-core to explore further opportunities for energy savings and peak temperature reduction. Additionally, we plan to consider complex 3D multi-core systems containing huge amount of cores and mapping of multiple applications on them at the same time.

REFERENCES

- [1] N. A. Dodgson, "Autostereoscopic 3D Displays," *IEEE Computer*, pp. 31–36, Aug. 2005.
- [2] K. Muller, A. Smolic, M. Drose, P. Voigt, and T. Wiegand, "3D Reconstruction of a Dynamic Environment with a Fully Calibrated Background for Traffic Scenes," *IEEE Trans. Cir. and Sys. for Video Technol. (TCSVT)*, pp. 538–549, 2005.
- [3] P. Merkle, A. Smolic, K. Muller, and T. Wiegand, "Efficient Prediction Structures for Multiview Video Coding," *IEEE Trans. Cir. and Sys. for Video Technol. (TCSVT)*, pp. 1461–1473, 2007.
- [4] "FinePix REAL 3D W3 — FujiFilm Global," http://www.fujifilm.com/products/3d/camera/finepix_real3dw3/.
- [5] "Panasonic HDC-SDT750K," <http://www2.panasonic.com>.
- [6] "Joint Draft 8.0 on Multiview video coding, JVT-AB204, 2008."
- [7] X. Xu and Y. He, "Fast disparity motion estimation in MVC based on range prediction," in *IEEE International Conference on Image Processing (ICIP)*, 2008, pp. 2000–2003.
- [8] Y. Kim, J. Kim, and K. Sohn, "Fast Disparity and Motion Estimation for Multi-view Video Coding," *IEEE Trans. on Consum. Electron.*, pp. 712–719, 2007.
- [9] J.-P. Lin and A.-W. Tang, "A fast direction predictor of inter frame prediction for multi-view video coding," in *IEEE International Symposium on Circuits and Systems (ISCAS)*, 2009, pp. 2589–2592.
- [10] P.-K. Tsung, W.-Y. Chen, L.-F. Ding, S.-Y. Chien, and L.-G. Chen, "Cache-based integer motion/disparity estimation for quad-HD H. 264/AVC and HD multiview video coding," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2009, pp. 2013–2016.
- [11] B. Zatt, M. Shafique, S. Bampi, and J. Henkel, "Multi-level pipelined parallel hardware architecture for high throughput motion and disparity estimation in Multiview Video Coding," in *IEEE Design, Automation Test in Europe Conference (DATE)*, 2011, pp. 1–6.
- [12] R. S. Patti, "Three-dimensional integrated circuits and the future of system-on-chip designs," *Proc. IEEE*, vol. 94, no. 6, pp. 1214–1224, 2006.
- [13] P. Ramm, A. Klumpp, J. Weber, and M. M. Taklo, "3D integration technologies," in *Proc. IEEE Symposium on Design, Test, Integration & Packaging of MEMS/MOEMS*, 2009, pp. 71–73.
- [14] Y. Cheng, L. Zhang, Y. Han, and X. Li, "Thermal-Constrained Task Allocation for Interconnect Energy Reduction in 3-D Homogeneous MPSoCs," *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, pp. 239–249, 2013.
- [15] C. Liu, L. Zhang, Y. Han, and X. Li, "Vertical Interconnects Squeezing in Symmetric 3D Mesh Network-on-chip," in *IEEE Asia and South Pacific Design Automation Conference (ASP-DAC)*, 2011, pp. 357–362.
- [16] C. Fabre *et al.*, "PRO3D, Programming for Future 3D Manycore Architectures: Projects Interim Status," in *Formal Methods for Components and Objects*. Springer, 2013, pp. 277–293.
- [17] B. S. Feero and P. P. Pande, "Networks-on-chip in a three-dimensional environment: A performance evaluation," *IEEE Transactions on Computers*, vol. 58, no. 1, pp. 32–45, 2009.
- [18] K. Bernstein, P. Andry, J. Cann, P. Emma, D. Greenberg, W. Haensch, M. Ignatowski, S. Koester, J. Magerlein, R. Puri *et al.*, "Interconnects in the third dimension: design challenges for 3D ICs," in *Proc. Design Automation Conference (DAC)*, 2007, pp. 562–567.
- [19] F. Clermidy, F. Darve, D. Dutoit, W. Lafi, and P. Vivet, "3d embedded multi-core: Some perspectives," in *Proceedings of IEEE Conference on Design, Automation and Test in Europe (DATE)*, 2011, pp. 1–6.
- [20] Y. Hoskote, S. Vangal, A. Singh, N. Borkar, and S. Borkar, "A 5-ghz mesh interconnect for a teraflops processor," *IEEE Micro*, vol. 27, no. 5, pp. 51–61, 2007.
- [21] J. Howard, S. Dighe, Y. Hoskote, S. Vangal, D. Finan, G. Ruhl, D. Jenkins, H. Wilson, N. Borkar *et al.*, "A 48-core IA-32 message-passing processor with DVFS in 45nm CMOS," in *Proc. IEEE International Solid-State Circuits Conference (ISSCC)*, 2010, pp. 108–109.
- [22] (2010) International technology roadmap for semiconductors. [Online]. Available: <http://www.itrs.net/reports.html>
- [23] X. Zhou, J. Yang, Y. Xu, Y. Zhang, and J. Zhao, "Thermal-aware task scheduling for 3d multicore processors," *IEEE Transactions on Parallel and Distributed Systems (TPDS)*, vol. 21, no. 1, pp. 60–71, 2010.
- [24] E. A. Lee and D. G. Messerschmitt, "Static scheduling of synchronous data flow programs for digital signal processing," *IEEE Transactions on Computers*, pp. 24–35, 1987.
- [25] L.-F. Ding, P.-K. Tsung, W.-Y. Chen, and S.-Y. Chien, "Fast motion estimation with inter-view motion vector prediction for stereo and multiview video coding," in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2008, pp. 1373–1376.
- [26] Z.-P. Deng, K.-B. Jia, Y.-L. Chan, C.-H. Fu, and W.-C. Siu, "A fast view-temporal prediction algorithm for stereoscopic video coding," in *IEEE Conference on Image and Signal Processing*, 2009, pp. 1–5.
- [27] L. Shen, Z. Liu, T. Yan, Z. Zhang, and P. An, "View-adaptive motion estimation and disparity estimation for low complexity multiview video coding," *IEEE Trans. Cir. and Sys. for Video Technol. (TCSVT)*, pp. 925–930, 2010.
- [28] N.-C. Chang, T.-H. Tsai, B.-H. Hsu, Y.-C. Chen, and T.-S. Chang, "Algorithm and architecture of disparity estimation with mini-census adaptive support weight," *IEEE Trans. Cir. and Sys. for Video Technol. (TCSVT)*, pp. 792–805, 2010.
- [29] M. Shafique, S. Garg, J. Henkel, and D. Marculescu, "The EDA Challenges in the Dark Silicon Era: Temperature, Reliability, and Variability Perspectives," in *Proceedings of ACM Design Automation Conference (DAC)*, 2014, pp. 185:1–185:6.
- [30] H. Khdr, S. Pagan, M. Shafique, and J. Henkel, "Thermal Constrained Resource Management for Mixed ILP-TLP Workloads in Dark Silicon Chips," in *Proceedings of ACM Design Automation Conference (DAC)*, 2015, pp. 179:1–179:6.
- [31] A. K. Singh, M. Shafique, A. Kumar, and J. Henkel, "Mapping on Multi/Many-core Systems: Survey of Current and Emerging Trends," in *Proceedings of ACM Design Automation Conference (DAC)*, 2013, pp. 1:1–1:10.
- [32] A. K. Coskun, J. L. Ayala, D. Atienza, T. S. Rosing, and Y. Leblebici, "Dynamic thermal management in 3D multicore architectures," in *Proceedings of IEEE Conference on Design, Automation and Test in Europe (DATE)*, 2009, pp. 1410–1415.
- [33] C. Addo-Quaye, "Thermal-aware mapping and placement for 3-d noc designs," in *Proc. IEEE SOC Conference (SOCC)*, 2005, pp. 25–28.
- [34] C. Sun, L. Shang, and R. P. Dick, "Three-dimensional multiprocessor system-on-chip thermal optimization," in *Proceedings of IEEE/ACM/IFIP Conference on Hardware/Software Codesign and System Synthesis (ISSS+CODES)*, 2007, pp. 117–122.
- [35] M. Cox, A. K. Singh, A. Kumar, and H. Corporaal, "Thermal-Aware Mapping of Streaming Applications on 3D Multi-Processor Systems," in *Proceedings of IEEE/ACM/IFIP Workshop on Embedded Systems for Real-Time Multimedia (ESTIMedia)*, 2013, pp. 11–20.
- [36] C. Zhu, Z. Gu, L. Shang, R. P. Dick, and R. Joseph, "Three-dimensional chip-multiprocessor run-time thermal management," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, vol. 27, no. 8, pp. 1479–1492, 2008.
- [37] S. Liu, J. Zhang, Q. Wu, and Q. Qiu, "Thermal-aware job allocation and scheduling for three dimensional chip multiprocessor," in *Proc. IEEE International Symposium on Quality Electronic Design (ISQED)*, 2010, pp. 390–398.

- [38] K. Puttaswamy and G. H. Loh, "Thermal herding: Microarchitecture techniques for controlling hotspots in high-performance 3d-integrated processors," in *Proc. IEEE International Symposium on High Performance Computer Architecture (HPCA)*, 2007, pp. 193–204.
- [39] K. Kang, J. Kim, S. Yoo, and C.-M. Kyung, "Runtime power management of 3-d multi-core architectures under peak power and temperature constraints," *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, vol. 30, no. 6, pp. 905–918, 2011.
- [40] F. Li, C. Nicopoulos, T. Richardson, Y. Xie, V. Narayanan, and M. Kandemir, "Design and management of 3d chip multiprocessors using network-in-memory," *ACM SIGARCH Computer Architecture News*, vol. 34, no. 2, pp. 130–141, 2006.
- [41] A. K. Singh, A. Kumar, and T. Srikanthan, "Accelerating throughput-aware runtime mapping for heterogeneous mpsocs," *ACM Transactions on Design Automation of Electronic Systems (TODAES)*, pp. 9:1–9:29, 2013.
- [42] "Hotspot 5.02 temperature modeling tool," *University of Virginia*, <http://lava.cs.virginia.edu/HotSpot>, 2011.
- [43] S. Bhat, "Energy models for network-on-chip components," MSc. thesis, Eindhoven University of Technology, 2005.
- [44] A. H. Ghamarian, M. Geilen, S. Stuijk, T. Basten, A. Moonen, M. Bekooij, B. Theelen, and M. R. Mousavi, "Throughput Analysis of Synchronous Data Flow Graphs," in *Proceedings of IEEE Conference on Application of Concurrency to System Design (ACSD)*, 2006, pp. 25–36.
- [45] S. Stuijk, M. Geilen, and T. Basten, "SDF³: SDF For Free," in *Proceedings of IEEE Conference on Application of Concurrency to System Design (ACSD)*, 2006, pp. 276–278.
- [46] S. Segars, "ARM7TDMI power consumption," *IEEE Micro*, vol. 17, no. 4, pp. 12–19, 1997.
- [47] E. L. Carvalho, N. L. V. Calazans, and F. G. Moraes, "Dynamic task mapping for mpsocs," *IEEE Des. Test of Comp.*, vol. 27, no. 5, pp. 26–35, 2010.
- [48] V. Chaturvedi, A. Singh, W. Zhang, and T. Srikanthan, "Thermal-aware task scheduling for peak temperature minimization under periodic constraint for 3D-MPSoCs," in *Proceedings of IEEE International Symposium on Rapid System Prototyping (RSP)*, 2014, pp. 107–113.
- [49] P. Yang, P. Marchal, C. Wong, S. Himpe, F. Catthoor, P. David, J. Vounckx, and R. Lauwereins, "Managing dynamic concurrent tasks in embedded real-time multimedia systems," in *Proceedings of IEEE/ACM/IFIP Conference on Hardware/Software Codesign and System Synthesis (ISSS+CODES)*, 2002, pp. 112–119.



Amit Kumar Singh (M09) received the B.Tech. degree in Electronics Engineering from Indian School of Mines, Dhanbad, India, in 2006, and the Ph.D. degree from the School of Computer Engineering, Nanyang Technological University (NTU), Singapore, in 2012.

He was with HCL Technologies, India for year and half before starting his PhD at NTU, Singapore, in 2008. From 2012 to 2014, he was a Post-Doctoral Researcher with the Department of Electrical and Computer Engineering, National University of Singapore, Singapore. Since 2014, he has been with the Department of Computer Science, University of York, UK. His current research interests include system level design-time and run-time optimizations of 2D and 3D multi-core systems with focus on performance, energy, temperature, and reliability. He has published over 40 papers in the above areas in leading international journals/conferences.

Dr. Singh was the receipt of PDP 2015 Best Paper Award, HiPEAC Paper Award, and GLSVLSI 2014 Best Paper Candidate. He has served as the Session Chair for conferences, such as APESER and DATE. He is a TPC Member of the IEEE/ACM conferences, such as ISED, NoCArc and MES.



Muhammad Shafique (M11) received the Ph.D. degree in computer science from the Karlsruhe Institute of Technology (KIT), Germany, in 2011.

He is currently a Research Group Leader at the Chair for Embedded Systems, KIT. He has over ten years of research and development experience in power-/performance-efficient embedded systems in leading industrial and research organizations. His current research interests include design and architectures for embedded systems with focus on low power, reliability, and adaptivity. He holds one U.S.

patent.

Dr. Shafique was a recipient of the 2015 ACM/SIGDA Outstanding New Faculty Award, six gold medals, the CODES+ISSS 2011, 2014, and 2015 Best Paper Awards, the AHS 2011 Best Paper Award, the DATE 2008 Best Paper Award, the DAC 2014 Designer Track Poster Award, the ICCAD 2010 Best Paper Nomination, several HiPEAC Paper Awards, and the Best Master Thesis Award. He is the TPC Co-Chair of ESTIMedia 2015 and 2016 and has served on the TPC of several IEEE/ACM conferences, such as ICCAD and DATE.



Akash Kumar (M05-SM13) received the B.Eng. degree in computer engineering from the National University of Singapore (NUS), Singapore, in 2002, the joint master of technological design degree in embedded systems from NUS and the Eindhoven University of Technology (TUE), Eindhoven, The Netherlands, in 2004, and the joint Ph.D. degree in electrical engineering in embedded systems from TUE and NUS, in 2009.

He is currently with the Technische Universitat Dresden (TUD), Germany, where he is directing the chair for Processor Design. From 2009 to 2014, he was with the Department of Electrical and Computer Engineering, NUS. His current research interests include design, analysis, and resource management of low-power and fault-tolerant embedded multiprocessor systems. He has published over 100 papers in leading international electronic design automation journals and conferences on these topics.

Dr. Kumar was recipient of the best paper award nominations including FPL 2014, GLSVLSI 2014, SC 2015, DATE 2015. He is also a Technical Program Committee Member of major conferences in the design automation and FPGA design area like, DAC, DATE, CASES, ASPDAC, FPL, FPT, etc.



Jörg Henkel (M95-SM01-F15) received the masters and Ph.D. (summa cum laude) degrees from the Technical University of Braunschweig, Braunschweig, Germany.

He is currently with the Karlsruhe Institute of Technology, Karlsruhe, Germany, where he directs the Chair for Embedded Systems. He was with NEC Laboratories, Princeton, NJ, USA. His current research interests include design and architectures for embedded systems with focus on low power and reliability. He holds ten U.S. patents.

Prof. Henkel was a recipient of the 2008 DATE Best Paper Award, the 2009 IEEE/ACM William J. McCalla ICCAD Best Paper Award, and the CODES+ISSS 2011 and 2014 Best Paper Awards. He was the Chairman of the IEEE Computer Society, Germany Section, and an Editor-in-Chief of ACM Transactions on Embedded Computing Systems. He is an Initiator and the Spokesperson of the National Priority Program called Dependable Embedded Systems of the German Science Foundation and the General Chair of ICCAD 2013.