

# Stochastic Modeling and Performance Analysis of Multimedia SoCs

Balaji Raman<sup>1</sup>, Ayoub Nouri<sup>1</sup>, Deepak Gangadharan<sup>2</sup>, Marius Bozga<sup>1</sup>, Ananda Basu<sup>1</sup>,  
Mayur Maheshwari<sup>1</sup>, Axel Legay<sup>3</sup>, Saddek Bensalem<sup>1</sup>, and Samarjit Chakraborty<sup>4</sup>

<sup>1</sup> UJF-Grenoble 1/CNRS VERIMAG UMR 5104, Grenoble, 38041, France,

<sup>2</sup> Technical Univeristy of Denmark, <sup>3</sup> INRIA Rennes, France,

<sup>4</sup> Technical University of Munich, Germany,

{balaji.raman, ayoub.nouri, marius.bozga, saddek.bensalem}@imag.fr, dega@imm.dtu.dk,  
basu@synopsys.com, mayurm@cs.utah.edu, axel.legay@inria.fr, samarjit@tum.de.

**Abstract**—Reliability and flexibility are among the key required features of a framework used to model a system. Existing approaches to design resource-constrained, soft-real time systems either provide guarantees for output quality or account for loss in the system, but not both. We propose two independent solutions where each modeling technique has both the above mentioned characteristics. We present a probabilistic analytical framework and a statistical model checking approach to design system-on-chips for low-cost multimedia systems. We apply the modeling techniques to size the output buffer in a video decoder. The results shows that, for our stochastic design metric, the analytical framework upper bounds (and relatively accurate) compare to the statistical model checking technique. Also, we observed significant reduction in resource usage (such as output buffer size) with tolerable loss in output quality.

## I. INTRODUCTION

An apt choice of a modeling framework is essential to design resource-constrained System-on-Chips (SoCs) in multimedia systems (such as video/audio players, etc.). Such a modeling framework must exploit the inherent stochastic nature of the multimedia applications to design low-cost systems. The uncertainty in such systems is due to high variability present in the input multimedia stream, in terms of number and complexity of items that arrive per unit time to the system. Consequently, the variability is exhibited both in arrival and in processing time.

Many of the existing *analytical* frameworks proposed so far are either incompatible or inflexible to capture the key characteristics of the system being modeled:

- *Worst-case execution time* modeling and analysis framework [17] cannot capture behavior of soft real-time systems, leading to pessimistic designs with exorbitant use of hardware resources (such as buffer size).
- *Average-case execution time* analysis framework [30] cannot provide QoS guarantees and are thus hardly trustworthy.

To address the above limitations, we sought a framework for analyzing multimedia systems that account for the stochastic nature of the streaming application. We need a model characterizing input stream and execution of the multimedia stream

as stochastic; instead of capturing event arrivals and executions with worst or average cases.

Recently, stochastic network calculus<sup>1</sup> [10] based approaches have been proposed for performance analysis of multimedia systems [18], [22]. These approaches, however, used the probabilistic calculus only partially: the input stream objects of a multimedia stream (e.g., frames) and their execution time are assumed to be deterministic. Another study used probabilistic real-time calculus to analyze hard real-time systems [23]. This later research, however, did not focus on any specific application domain. Nonetheless, if stochastic network calculus is fully adopted for system-level design, then design of multimedia embedded platforms can benefit, too. The analysis can provide probabilistic bounds on the output quality of the system. The worst-case and average-case analysis of the system are special-case scenarios of the analytical framework.

The two aforementioned key modeling features — flexibility and reliability<sup>2</sup> — for modeling multimedia systems exist in statistical model checking approaches, too. Statistical model checking consists in simulating the formal representation of the systems, monitors a finite set of traces, and then guess an overall correctness by exploiting algorithms from the statistics area. Recent work showed that component-based frameworks, such as BIP [3], can be coupled with statistical model checking to verify large heterogeneous embedded systems [2]. Unlike many ad-hoc simulation techniques, BIP provides a formal semantics for the modeling and simulation of stochastic systems. Such a formal semantics is a prerequisite for using statistical model checking as a sound performance analysis technique on system models.

In what follows, we list the main **contributions** of our paper:

- *Stochastic characterization*: We use stochastic real-time calculus for performance analysis of multimedia systems. Our model captures the uncertainty of arrivals and execu-

<sup>1</sup>Stochastic network calculus was originally developed for performance analysis of computer networks

<sup>2</sup>e.g modeling probabilistic systems and providing guarantees

tions of stream objects (and dependencies between them) using probability laws.

- *Model comparison:* Using a video decoder case-study, we apply both approaches (a major contribution is to build a BIP model). Our results show that the probabilistic estimates obtained from analytical framework upper bounds the estimates from statistical model checking.
- *Resource constrained design:* Our second goal with the video decoder case-study is to reduce resource consumption with tolerable loss in video quality. This case study illustrates how expensive resources (such as the playout buffer) can be reduced if the stochastic behavior of the application is carefully taken into account during performance. We chose some application parameter values (quickly) using the analytical framework. Then, we estimate the buffer sizes (precisely) for a set of parameter values using statistical model checking. We discuss how the strengths of both the approaches can be complementary to each other.

**Organization:** We present the analytical and the statistical modeling checking approaches in Sections II and III. In each of these two sections, we first present some background and then the modeling technique. In Section IV, we apply both the approaches—again, independently—to a case-study, video decoder. The results of the analytical approach are compared with the statistical model checking approach (for the case study) in Section V. The results section reports the buffer size savings, too. The tail end portion of the paper contains the discussion, related work, and the conclusions. Refer technical report for more details [21].

## II. ANALYTICAL MODEL

We first give the background of a deterministic analytical model of a system running multimedia applications, and, then, extend the model to a probabilistic setting.

### A. Background

This subsection presents the real-time calculus framework for the system model we use in our subsequent sections. Note that this subsection is not intended to be an exhaustive exposition of real-time calculus for embedded systems. The system model used in this paper is sketched in Fig. 1. The architecture is composed of memory buffers and processing units. Fig. 1 shows a data stream  $x$  reaching the input buffer while the processed stream  $y$  is written to the playout buffer.

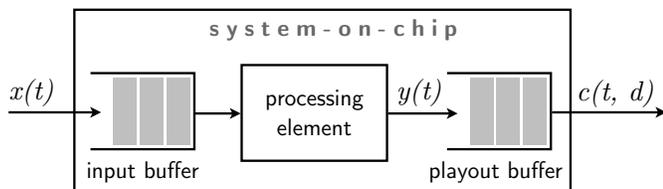


Fig. 1: Real-time calculus model

For simplicity of exposition, we chose a basic block or unit of the real-time calculus to model the multimedia SoC. Indeed,

analysis using real-time calculus has been shown for a system with multiple streams to the input buffer [7], multiple processing units in a SoC [20], and multiple tasks in a processor [19]. This paper focus, however, is to propose performance analysis using stochastic real-time calculus. Future work can extend this stochastic framework to the above listed settings.

The application mapped to the processing element is a multimedia task. However, there is no limitation on the number of applications running on the processing element. Consequently, there are no restrictions on the number of input streams fed to the processing element. Similarly, the number of architecture units in the system-on-chip can be any number of processing elements and memory units. Indeed, analysis using real-time calculus has been shown for a complex system containing shared memory, bus, and network-on-chip communication architectures. The data flow need not be sequential in that there could be a feedback flow between the architectural components in the system-on-chip, and there could be splits and joins in the dataflow of the architecture.

The notations of the parameters involved in the description of our model are described in the following tables, divided as follows:

- Table I lists the application and architecture parameters that are given to the system designer.
- Table II records the functions that are obtained from processor simulations.
- The mathematical functions that are constructed from the known parameters from Table II are given in Table III.
- Similarly, mathematical functions constructed from simulation data from Table I are presented in Table IV.
- Finally, Table V and VI shows the functions and the models used in the real-time calculus framework.

TABLE I: Known applications and architecture parameters.

	Notation and name	
Application	$r$	Bitrate (in bits per second)
	$c$	Consumption rate (in macroblocks per second)
	$d$	Initial playout delay (in millisecond)
Architecture	$b$	Input buffer size (in macroblocks <sup>†</sup> )
	$B$	Playout buffer size (in macroblocks or bytes)
	$f$	Processor frequency (in MHz)

<sup>†</sup>The input buffer can hold a given number of macroblocks and it is not given by a size in bytes because the application case study we use in this paper is video decoding. The compressed stream arriving at the input buffer consists of macroblocks of variable size whereas the playout buffer holds decompressed items of constant size.

The core part of the deterministic real-time calculus can be reduced to three inequalities. First, bounds on the arrival of the data stream are computed from the  $\phi_l, \phi_u$  functions, which gives the minimum and maximum number of bits constituting any  $k$  consecutive stream objects respectively.

$$\alpha_l(\Delta) = \phi_u^{-1}(r\Delta), \quad (1)$$

$$\alpha_u(\Delta) = \phi_l^{-1}(r\Delta), \quad (2)$$

where  $r$  is the bit-rate of the input video, and where  $\alpha_l(\Delta), \alpha_u(\Delta)$  are the minimum and maximum number of

TABLE II: Parameters from simulation of processor and from application source code.

Notation	Name
$\text{bits}(k)^\dagger$	Cumulative bits
$\text{cycles}(k)^*$	Cumulative cycles

<sup>†</sup>The number of bits for each and every macroblock is computed with the execution of the application once and the data flow between the software blocks.

\*The simulation performed to compute the cycles is not a system-level simulation. The processing element is independently simulated using a software simulation to compute the processor cycles consumed for each and every macroblock.

TABLE III: Analytical functions from processor and application simulation data.

Notation	Name
$\phi_l(k), \phi_u(k)$	$\phi$ functions
$\gamma_l(k), \gamma_u(k)$	$\gamma$ functions

TABLE IV: Analytical functions from input parameters.

Notation	Name
$\beta\text{-cycles}(k)$	Processor cycles over $[0, t]$
$c(t)$	Consumption function

TABLE V: Analytical functions and inputs of arrival, service and output.

Notation	Name
$x(t)$	Input function
$y(t)$	Output function

TABLE VI: Analytical model.

Notation	Name
$\alpha_l(\Delta), \alpha_u(\Delta)$	Arrival curves
$\beta_l(\Delta), \beta_u(\Delta)$	Service curves

items arriving over the time interval  $\Delta$  respectively. The mathematical functions  $\phi_{u,l}$  are computed from the parameter  $\text{bits}(k)$  (obtained from simulating the application code), which gives the number of bits for first  $k$  stream objects, as shown below

$$\phi_u(k) = \max\{\text{bits}(k + N) - \text{bits}(N)\} \quad (3)$$

$$\phi_l(k) = \min\{\text{bits}(k + N) - \text{bits}(N)\} \quad (4)$$

where  $N$  can take any value such that  $k + N$  does not exceed the stream length.

The number of items arriving at the input buffer over the time interval  $[0, t]$  ( $x(t)$ ) is bounded by the previously introduced  $\alpha_{u,l}$  functions:

$$\alpha_l(\Delta) \leq x(t + \Delta) - x(t) \leq \alpha_u(\Delta), \quad (5)$$

for  $t, \Delta \geq 0$ .

The second core inequality is based on the service curve ( $\beta$ ), which guarantees the number of processor cycles dedicated to that particular multimedia task over the time interval  $\Delta$ . Given the processor frequency ( $f$ ), we first compute the processor

cycles available over time interval  $[0, t]$ , that is,  $\beta\text{-cycles}(t) = tf$ . Second, given the number of cycles consumed by first  $k$  stream objects ( $\text{cycles}(k)$ ), we compute the minimum number of cycles needed to process any  $k$  consecutive stream objects ( $\gamma_l$ ) in a similar manner as shown in Eq. 4. Third, we compute  $\beta$  as follows:  $\beta_u(\Delta) = \gamma_l^{-1}(\beta\text{-cycles}(\Delta))$ . Now, we present our second core inequality. Let  $y(t)$  be the number of items arriving in the output buffer over the time interval  $[0, t]$ . Then, it can be shown that:

$$y(t) \leq (\alpha_u \otimes \beta)(t), \quad (6)$$

where  $\otimes$  is the min-plus convolution operator defined as:

$$(p \otimes q)(t) = \inf_{0 \leq s \leq t} \{p(t-s) + q(s)\}. \quad (7)$$

Finally, if the playout buffer never underflows, we have the last core inequality:

$$y(t) \geq c(t, d), \quad \forall t \geq 0, \quad (8)$$

where  $c(t, d)$  is the items consumed after the initial playout delay ( $d$ ) by the display device over the time interval  $[0, t]$ .

### B. Stochastic Real-time Calculus Model

In this sub-section, we propose a probabilistic framework for designing multimedia SoCs.

From the inverse  $\phi$  function, we compute  $\alpha_l(\Delta)$  and  $\alpha_u(\Delta)$ . This leads to the definition of the stochastic arrival curve. Since by definition,  $x(t)$  gives the number of items arriving over the time interval  $[0, t]$ , we obtain:

$$P \left( \sup_{0 \leq \Delta \leq t} (x(t + \Delta) - x(t) - \alpha_u(\Delta)) > a \right) \leq f(a), \quad (9)$$

for all  $0 \leq \Delta \leq t$  and for all  $a \geq 0$ .

In the above description, the arrival curve bounds are checked with the actual number of items arriving at the input buffer over any time interval. The decreasing function  $f(a)$  is an upper bound on the probability.

We define the stochastic service curve as follows. As formulated in the previous subsection, the output function from the processing element is guaranteed to be smaller than the min-plus convolution of the arrival and service curves:

$$y(t) \leq (\alpha_u \otimes \beta_u)(t), \quad \forall t \geq 0. \quad (10)$$

For a stochastic service, the above inequality is restated as follows:

$$P(y(t) - (\alpha_u \otimes \beta_u)(t) > a) \leq g(a), \quad \forall t \geq 0, \quad (11)$$

where the stochastic bounding function  $g$ , ideally obtained from psychovisual models, is related to the acceptable loss of playback quality. Note that in defining a stochastic service curve we are also including the definition of the arrival curve. Thus dependencies are handled for both the arrival and the service curve.

The output from the processing element is bounded by the arrival functions  $\alpha_l$  and  $\alpha_u$ . The results from the stochastic network calculus provides bounds on the output curve [10]. If

there is a stochastic arrival curve as defined in Eq. 9 and a stochastic service curve as defined in Eq. 11, the output curve is defined as follows:

$$P(y(t) - (\alpha_u \circledast \beta_u)(t) > a) \leq (f \otimes g)(a), \quad \forall t \geq 0. \quad (12)$$

Write  $h(a) = (f \otimes g)(a)$ , and the min-plus deconvolution operator is defined as follows:

$$(p \oslash q)(t) = \sup_{u \geq 0} \{p(t+u) - q(u)\}, \quad (13)$$

### III. STATISTICAL MODEL CHECKING APPROACH

In this section, we present a second approach for performance evaluation of multimedia SoCs that is based on statistical tests.

#### A. Background

**Statistical Model Checking (SMC)** has been proposed as an alternative to classical Model Checking techniques [5]. It aims to avoid exhaustive state space exploration. The idea is to do verification on a sub-part of the state space (a sample) and then, using statistics, extrapolate the result to the whole system with some confidence.

Concretely, given a stochastic system  $S$  and a property  $\phi$ , *statistical model checking* refers to a series of simulation-based techniques that can be used to answer two questions : (1) *qualitative* : is the probability for  $S$  to satisfy  $\phi$  greater or equal to a certain threshold  $\theta$  ? and (2) *quantitative* : what is the probability for  $S$  to satisfy  $\phi$  ?

The main approaches [29], [24] proposed to answer the qualitative question are based on *hypothesis testing*. Let  $p$  be the probability that  $S \models \phi$ . To determine whether  $p \geq \theta$ , we can test  $H : p \geq \theta$  against  $K : p < \theta$ .

A statistical-based solution (based on a sample) does not guarantee a correct result but it is possible to bound the probability of making an error. The *strength* of a test is determined by two parameters,  $\alpha$  and  $\beta$ , such that the probability of accepting  $K$  (respectively,  $H$ ) when  $H$  (respectively,  $K$ ) holds is less or equal to  $\alpha$  (respectively,  $\beta$ ).

Since it is impossible to ensure a low probability for both types of errors simultaneously, a solution is to use an *indifference region*  $[p_1, p_0]$  (with  $\theta$  in  $[p_1, p_0]$ ) and to test  $H_0 : p \geq p_0$  against  $H_1 : p \leq p_1$ .

Several hypothesis testing algorithms exist in the literature. Younes [29] proposed a logarithmic based algorithm that given  $p_0, p_1, \alpha$ , and  $\beta$  implements the *Sequential Ratio Testing Procedure* (SPRT) (see [26] for details).

In [6], [13] Peyronnet et al. propose an estimation procedure (PESTIMATION) to compute the probability  $p$  for  $S$  to satisfy  $\phi$ .

**BIP – Behavior, Interaction, Priority** – [1] is a component based framework encompassing rigorous model based design. It allows building hierarchically structured systems (or composite components) from atomic components characterized by their behavior and their interface.

Components are composed by layered application of interactions and priorities. Interactions express synchronization

constraints between actions of the composed components while priorities are used to filter amongst possible interactions and to steer system evolution e.g. to express scheduling policies.

In BIP, *atomic components* are finite-state automata extended with variables and ports. Variables are used to store local data. Ports are action names, and may be associated with variables. They are used for interaction with other components. States denote control locations at which the components await for interaction.

A transition is a step, labeled by a port, from a control location to another. It has associated a guard and an action, that are respectively a Boolean condition and a computation defined on local variables.

In BIP, data and their transformations are written in C/C++. *Composite components* are defined by assembling atomic or composite using *connectors*. Connectors relate ports from different sub-components and represent sets of *interactions*, that are, non-empty sets of ports that have to be jointly executed. For every such interaction, the connector provides the guard and the data transfer, that are, respectively, an enabling condition and an exchange of data across the ports involved in the interaction. Finally, *priorities* provide a mean to coordinate the execution of interactions within a BIP system. They are used to specify scheduling or similar arbitration policies between simultaneously enabled interactions.

To apply Statistical Model Checking on BIP models, it has been augmented by a stochastic semantics [4]. The new semantics allows definition of stochastic atomic components containing probabilistic variables updated through C-defined probabilistic distributions. The stochastic components communicate through interactions that are uniformly selected.

#### B. Stochastic Modeling and Statistical Analysis

The Statistical Model Checking Approach is implemented in the SBIP tool [4]. The latter takes as input (1) a stochastic BIP system model, (2) a probabilistic bounded LTL property [6] or, alternatively, an observer component encoding the evaluation of the property, and (3) a series of confidence parameters needed for the statistical test [29]. Then, it proceeds according to the following steps:

- *Step 1*: an executable model is automatically created,
- *Step 2*: random execution traces of the system are iteratively generated (sampling),
- *Step 3*: the observer component checks the property on each trace,
- *Step 4*: then gives a partial verdict for each,
- *Step 5*: steps 2, 3, and 4 are repeated until the SMC engine is able to conclude over the whole system. The SMC engine implements the statistical algorithms introduced in the previous section.

Our tool is guaranteed to terminate its execution, when it has decided an answer for the input property to be verified on the input model with respect to the input confidence parameters. This guarantee relies on the mathematical theory of the statistical model checking.

#### IV. CASE STUDY: VIDEO DECODER

We analyze an abstraction of a video decoder SoC (shown in Figure 2) with the analytical and the statistical model checking approaches (presented in previous sections). In the abstract SoC model, the input video stored is fed to the input buffer in terms of stream objects such as macroblocks, frames, etc. A pipeline of functional units process the input stream. Processed items are temporarily stored in the output buffer before their display.

Now, we will state the problem addressed in this paper. We assume that the multimedia SoC contains a single processing unit and two buffers.

**Problem Statement:** To estimate the probability that the buffer underflow ( $U(t)$ ) is less than two consecutive frames in 30 frames. We are given the following:

- a set of video clips of certain bit-rate ( $r$ ) and resolution,
- maximum frequency of the processing unit of the multimedia SoC ( $f$ ),
- consumption rate of the output device ( $c$ ),
- start-up values for the initial delay ( $d$ ), input buffer size ( $b$ ), and playout buffer size ( $B$ ).

In this section, first we apply the stochastic real-time calculus model, and then the BIP and statistical model checking approach to evaluate the QoS constraints of the video decoder SoC.

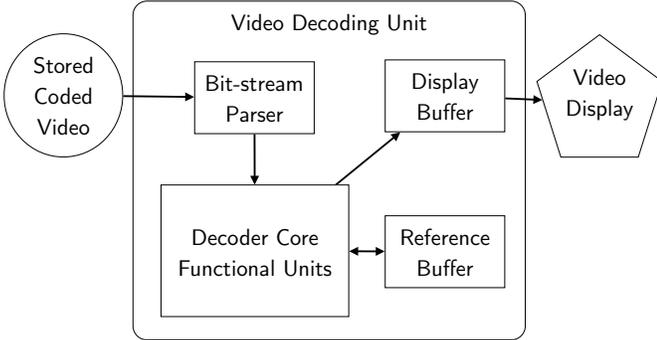


Fig. 2: Display buffer in a generic video decoder of a SoC. The functional units of the decoder can include variable length decoding, motion compensation, etc.

##### A. Analytical Model of the Multimedia SoC

In this subsection, first, we give an overview of our analytical approach for this case study. Second, we present formulation for evaluating the QoS constraints.

We can estimate the maximum size of the output buffer using our stochastic analytical framework. Previous work that estimate the output buffer size using deterministic real-time calculus (presented in Section II-A) proceed as follows [17]. For all given video clips, Maxiaguine et al. construct upper bounds on the number of items that arrive to the input buffer and that execute in the processor. These two bounds together yield another upper bound on the number of items that arrive to the display buffer. Thus, given a rate at which items are

consumed from the output buffer, Maxiaguine et al. estimate the maximum buffer size required.

We, too, compute deterministic upper bounds on arrival, execution, and output, however, only for a sub-set of given video clips; the remaining video clips can violate the deterministic bounds. For example, the number of items arriving to the output buffer over a time interval, for a certain clip, can be larger than the deterministic output bound. Now, we explain how to quantify this deviation from the deterministic bound in a stochastic setting (as we are using probabilistic network calculus).

Assume that the stream objects arrival and execution are stochastic—an apt characterization of multimedia streams. So, what is the probability that the number of stream objects that arrive to the output buffer over a time interval is larger than the deterministic output bound?

First, we estimate the maximum probability of violating the deterministic bounds for the arrival and execution; then, using these two bounds, we estimate it for the output (related to the definitions presented in Section II-B). Using this stochastic bound on the output, given the constant consumption rate, we can compute the probabilistic distribution of the buffer size.

The deterministic analysis imposes hard constraints on the arrival, execution, and output of stream objects leading to over estimation of output buffer size. In contrast, our probabilistic analysis relaxes these constraints. Thus, the designer can choose a buffer size and the corresponding video quality. Following this, we explain how the video quality is related to the buffer size.

We specify tolerable loss in video quality as: less than two consecutive frame loss within 30 frames, less than 17 aggregate frame loss within 100 frames, etc [27]. Previous work models the loss of stream objects as buffer underflows, which occurs whenever the display device finds insufficient items to read from the output buffer [22]. Raman et al. show that the amount of buffer underflow can be controlled using an application parameter, namely, the initial playout delay, the delay after which the video starts to display.

We, too, tune the initial delay parameter to restrict the maximum buffer underflow, albeit, in a stochastic setting. We obtain probability values for a QoS property to hold for a range of initial delays. The buffer size corresponding to each such delay can be estimated. Therefore, the system designer can choose an initial delay based on his resource and quality constraints.

In what follows, we formulate the probability distribution of the output buffer size using the stochastic real-time calculus model presented in Subsection II-B.

In practice, the designer is typically given a set of input clips to design the SoC with given display constraints. So, in our problem setting, the designer can construct an upper and lower bound using synthetic traces (which are obtained from actual traces) and use the actual traces to construct the bounding functions. Now, we explain the construction of these synthetic traces.

Let the given set of input video clips be partitioned into two

sets,  $S_A$  and  $S_B$ , based on the designer’s requirement that all clips in  $S_B$  must be processed with no loss in video quality. The deterministic upper bound on the arrival and the output, introduced in the previous section, are constructed using clips in  $S_B$ . Now, we discuss how to synthetically generate clips in case we do not have a set  $S_B$ .

The information we have about the clips in set  $S_A$  are the number of bits and number of cycles corresponding to each macroblock. For certain macroblocks, we modify the number of bits and cycles, assuming we are given lower bounds<sup>3</sup> on these parameters. Any number of bits lower than the actual bound in the input traces are replaced with a value of the lower bound. Thus we obtain synthetic traces forming clips in set  $S_B$ . Now we explain how we estimate stochastic bounds using the actual clips (i.e. clips from set  $S_A$ ) and synthetic traces (or if available clips from set  $S_B$ ).

The upper and lower bounds on the arrival given from the formula in the previous subsection are calculated for the synthetic traces. That is, from the modified inverse  $\phi$  function, we compute  $\alpha_l(\Delta)$  and  $\alpha_u(\Delta)$ . This leads to the definition of the stochastic arrival curve.

Assume the output arrival curve is an arrival process to the playout buffer. The probability distribution at the playout buffer can be computed using the bounding function  $h$ .

$$P(U(t) > a) \leq h(a - (\alpha_* \circ ct))(0), \forall a, t \geq 0. \quad (14)$$

In the above equation  $\alpha_*$  is the output arrival function given by  $(\alpha_u \circ \beta_u)(t)$  which is denoted in Eq.12.

In the results sections, we will show how to specify the QoS property by using Eq. 14. To validate and compare probabilistic bounds obtained from the analytical framework, we used a rigorous simulation framework, BIP [1], introduced in Section III-A.

### B. Stochastic BIP Model of the Multimedia SoC

We constructed a model of the video decoding unit shown in Figure 2 using BIP framework. This model captures the stochastic behavior of the system in the following way. A stream object’s arrival time to the input buffer, and decoding time in the processing unit follow defined probability distributions. The distributions are constructed from a set of video clips<sup>4</sup>. Thus, the stream object’s arrival to the output buffer is probabilistic, which can lead to probabilistic buffer underflows. Next, we explain how to estimate the probability of certain amount of buffer underflow.

Figure 3 shows the stochastic BIP model of the SoC running the video decoding application. The functional units of the SoC are modeled as atomic components respectively, **Generator**, **Processor**, and **Player**. These functional components communicate explicitly through buffers, namely, **Input Buffer**

and **Playout Buffer**, represented in BIP as atomic components as well. The lines represent connectors, namely **write-push**, **pop-read** are used to transfer macroblocks objects between a functional component and a buffer component. The tick connector synchronizes all the functional components, and is used to model explicitly the progress of the absolute (global) time. Now we describe the behavior of each of the functional components with more details.

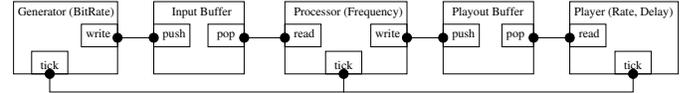


Fig. 3: Stochastic BIP model of the SoC running a video decoding application.

**Generator:** This component models the generation of a stream of macroblocks. The stream is generated probabilistically and stored in the input buffer. The number of bits (the size) of every macroblock determines the arrival time of the macroblock to the input buffer. This number of bits is specific for each macroblock type (w.r.t frame types) and follows a specific distribution, shown in Figure 4 and based on [11]. Moreover, the choice of frames type follows a Group of Pictures Pattern (GOP): IBBPBBPBBPBB<sup>5</sup>.

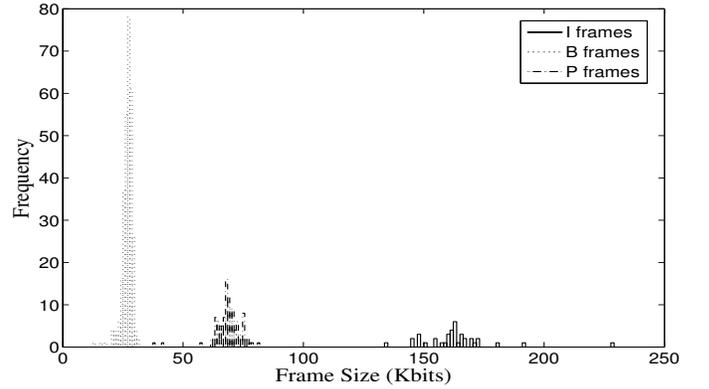


Fig. 4: Frequency distributions of I, B, and P frames in the video. I frames are larger in size but smaller in number than the B and P frames.

**Processor:** This component models the decoding of macroblocks, sequentially after reading them from the input buffer. The detailed behavior of the component is shown in Figure 5. The component has two states, **IDLE** and **PROCESS** and three ports **read**, **write**, and **tick**. In the **IDLE** state, the process is either waiting to read from the input buffer or waiting to write to the playout buffer. When there is a macroblock available, the process transits to the **PROCESS** state and remains there for the time required to process/decode the macroblock.

**Player:** The **Player** component models the consumption of the stream of decoded macroblocks. After an initial playout

<sup>5</sup>The first GOP is IPBBPBBPBB. It is different from the consecutive GOPs.

<sup>3</sup>In the discussion section, we present a technique to generate synthetic traces without this lower bound

<sup>4</sup>Note that we construct distributions from those clips we used in the analytical framework that were allowed to violate the deterministic upper bounds.

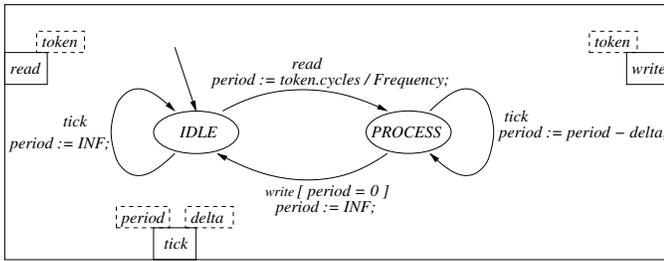


Fig. 5: *Processor* model as a BIP component. The processor unit executes a macro-block at a pre-specified speed. The execution cycles corresponding to each macro-block is randomly selected from a distribution or taken from a file in a sequential manner.

delay, the *Player* starts reading the macroblocks from the playout buffer at a constant rate. A buffer underflow occurs whenever the requested number of macroblocks is not available in the buffer. In this case, the request is postponed to the next iteration and the underflow is accumulated. For example, if the current buffer underflow is 2, then, at the next request, the *Player* seeks 3 macroblocks. If the buffer is still empty, the underflow became 3. Else, if the playout buffer has (at least) 3 items, then all three items are read at once and the buffer underflow is reset to 0, etc.

We applied statistical model checking to evaluate QoS properties on the BIP model of the multimedia SoC presented above. As explained earlier, this model is fully stochastic. We focus on a qualitative QoS property related to the playout buffer, that is, the buffer underflow within a second never exceeds two consecutive frames.

In order to evaluate this property on the traces of the model, we use the additional *Observer* component shown in Figure 6. This component runs in parallel with the systems and reacts to events (interactions) relevant to the satisfaction of the property. The component has three states: *OK*, *PARTIAL*, and *FAIL*. The *FAIL* state denotes the failure of the property, namely, the underflow of two consecutive frames within a second. If there is a loss of a single frame the observer moves from state *OK* to *PARTIAL*. Later, if there is an additional frame loss the *Observer* reaches the *FAIL* state. If no loss happens within that second, the component moves back to the *OK* state.

In the next section, we present the analysis results for both the analytical and the SMC approaches.

## V. RESULTS

This section sketches QoS probabilities estimated from the two analysis approaches presented in previous sections.

We implemented the analytical framework described in Subsection IV-A in MATLAB. The experiments were conducted for a low-bit rate and low resolution clips (352 \* 240) obtained from an open source [25]. The bit-rate of the input video is 1.5 Mbits per second and the frame output rate is 30fps. We used an MPEG2 implementation optimized for speed [15]. The MPEG2 source was annotated to get the

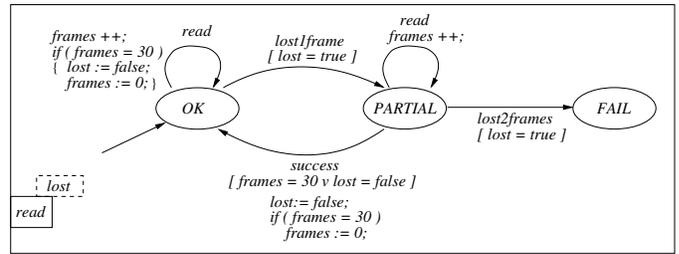


Fig. 6: *Observer* model as a BIP component. The observer models the QoS property to be verified. The variable *frame* counts the number of frames to check if two consecutive loss occurs within a second (i.e. within 30 frames). The read port of the *Observer* is synchronized with the read port of the *Player*. A variable *lost* is associated with the read port records a frame loss.

number of bits corresponding to each compressed macroblock. The execution cycles for each macroblock is obtained from the software simulator SimpleScalar. Recapitulate that the number of bits and execution cycles per macroblock are inputs to the analytical framework. We chose the video files *cact.m2v*, *mobile.m2v*, and *tennis.m2v* for our experiments.

To construct the synthetic clips we set the lower bound for bits (e.g to 60) and the lower bound for execution cycles (e.g to 9000). Then from the actual trace containing the number of bits and execution cycles per macroblocks, synthetic traces are obtained; any value below the lower bound is modified to the lower bound. Figures 7, 9, and 11 show the probability that the buffer underflow is greater than two consecutive frames over any time interval (this refers to Equation 14 in Subsection IV-A). Figures 7, 9, and 11 also show the probability estimates from the BIP framework. Following are the observations:

- Increase in playout delay decreases the amount of buffer underflow, so, probability that the buffer underflow is more than two consecutive frames decreases.
- The estimates from stochastic real-time calculus upper bound the statistical model checking results as the analytical framework captures the worst-case behavior.
- The delay values at which the statistical model checking starts to state that the property is true is not same for the analytical framework. The analysis using the stochastic real-time calculus should be used to determine a small set or range of delay values. Later, to precisely verify the property detailed simulation should be carried out.
- For each probability estimation, the number of traces statistical model checking simulated ranged from 44 to 1345. For each trace, the method took around 6 to 8 seconds to verify the property. The probability of error for the probabilistic estimates from the statistical model checking is bounded by 0.01.

Figures 8, 10, and 12 plots the results of the buffer sizes for various playout delay values, and corresponding probabilistic bounds. These results correspond to BIP simulation and sta-

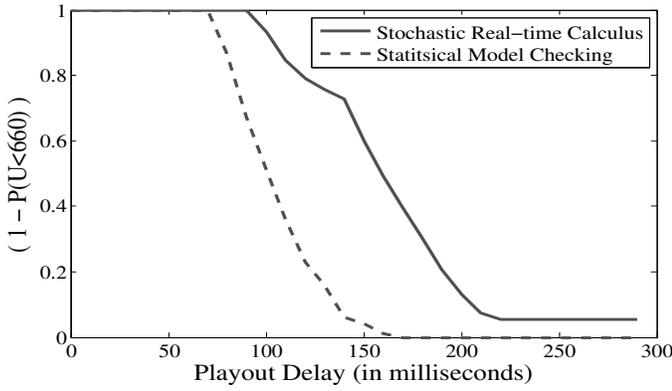


Fig. 7: Probabilistic bounds for `cact.m2v`

tistical model checking. We observe that buffer size reduces substantially even for a small decrease of probabilistic value. For instance, there could be a buffer size reduction of 40% for an increase in the value of the probabilistic bound from 0 to 0.2 (Figures 10). In fact the buffer savings can be larger if we compare the buffer size required for no underflow and the memory required for the QoS property to be always true.

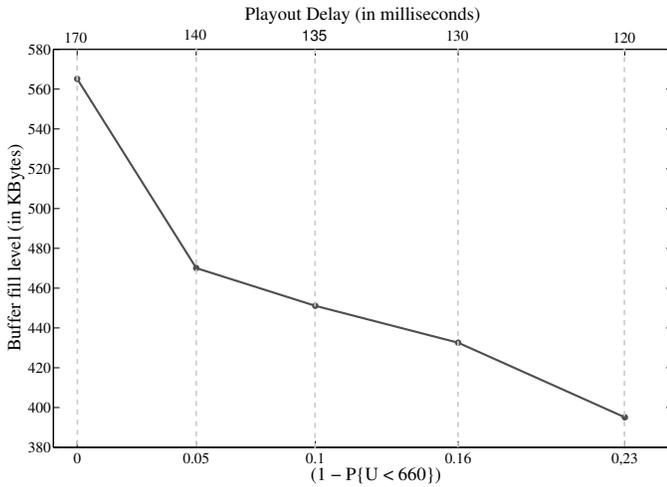


Fig. 8: Playout buffer fill level for `cact.m2v`.

## VI. DISCUSSION

The results in the previous section confirms our hypothesis that, for tolerable loss in video playout, output buffer size could be significantly reduced compared to the buffer size required for playing lossless video. In this section, first we speculate the combined strengths of both the approaches when used together in a system design flow. Then, we focus on requirements for our hypothesis to be used in practice: extraction of synthetic clips from benchmark video clips.

### A. Combined Design Flow

We now identify the benefits of the two techniques when used independently in a design-flow and speculate their combined strengths.

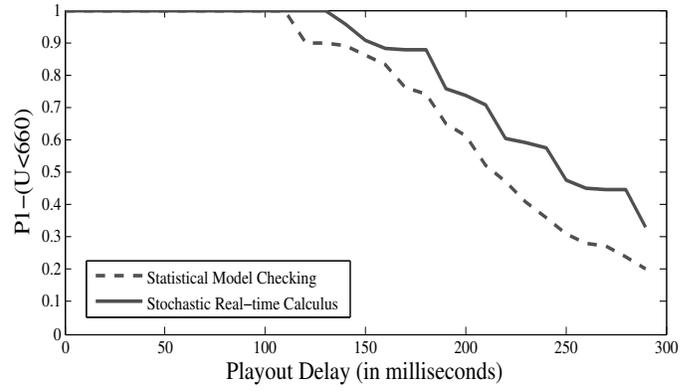


Fig. 9: Probabilistic bounds for `mobile.m2v`

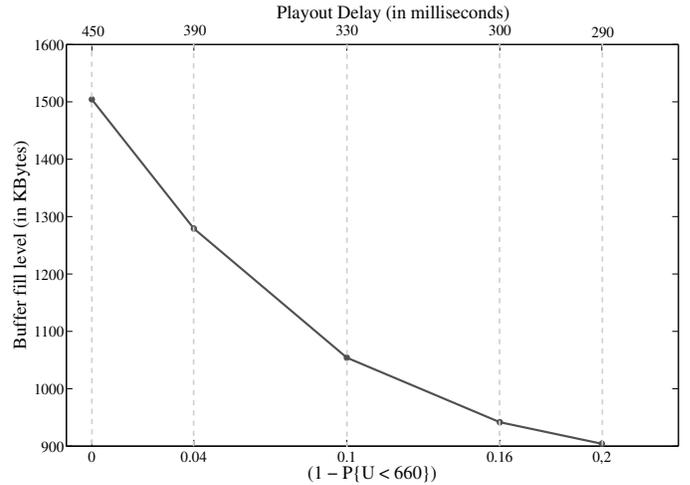


Fig. 10: Playout buffer fill level for `mobile.m2v`.

Analysis using stochastic real-time calculus followed by detailed simulation is a natural order for these techniques in a SoC design-flow; quick analysis is needed for estimating the application/architecture parameters at the beginning of the first stage of the design process, whereas at the end detailed simulations are required for verifying properties. The drawbacks, however, when the analytical framework and simulation are applied independently are: (1) the stochastic real-time calculus is not flexible enough to specify all verification properties, and (2) in using just the statistical model checking approach, the iterative process can be time-consuming for deciding parameters (in the QoS property).

On the other hand, a combined framework will not only remove the drawbacks mentioned above but will increase the expressivity of the stochastic real-time calculus (state-space notation) and the performance analysis using algebraic approach comes for free to the BIP framework.

### B. Inputs to the Model

The primary motivation to generate synthetic clips is that a system designer would use the synthetic clips instead of actual video clips when deciding architecture parameters such

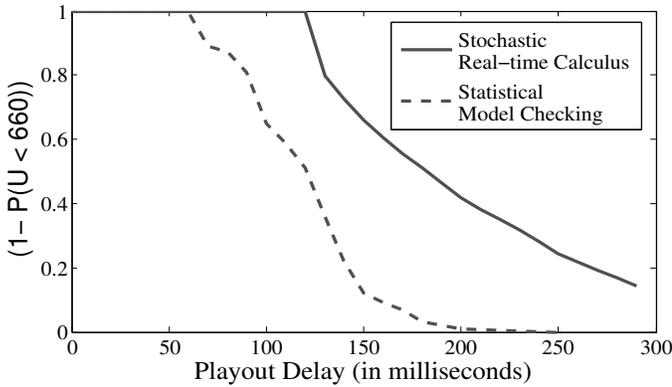


Fig. 11: Probabilistic bounds for `tennis.m2v`

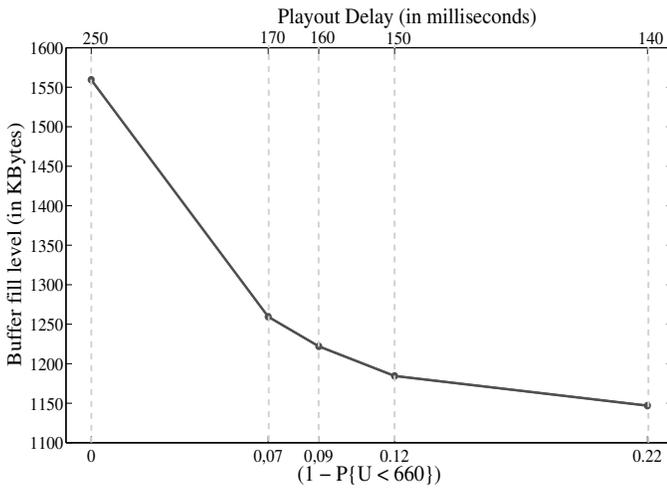


Fig. 12: Playout buffer fill level for `tennis.m2v`.

as buffer size. The question is then how to generate synthetic clips from given actual video clips.

Currently, we choose at which point in frequency distribution of the input data (bits and cycles) we need to cut and reshape frequency distribution to generate the synthetic data. We choose this cut-off point based on two objectives: (1) to get significant reduction in playout buffer size, and (2) the loss in video is tolerable when the actual video clip is running in the SoC, which is designed using synthetic clips.

The approach we use to generate synthetic data can be time-inefficient as it involves iteration in making a choice to cut the distribution and analyze using our framework to check if our objectives are met. Instead of this trial and error technique, we are currently researching on a more sound technique to generate synthetic clips (for example, Yanhong et al., [16] use error percentage and eliminate some tail data from the distribution). To further this thought, notice in the analytical framework the need for synthetic clips arises when we compute the stochastic bounding functions. What if we use standard tail distribution functions that could characterize multimedia data accurately?

Assuming that these tail distributions could be found (refer Jelenkovic et al., [9]), then the task of generating synthetic clips reduces to a litmus test: does the synthetic clip generated conforms to the stochastic bounding function? In other words, without performing the iterative analysis — choosing a cut-off point first for synthetic clip generation, second checking if the video loss is tolerable, and then choosing a different cut-off point — we just check if a set of synthetic clips generated conforms to the stochastic bounding function.

## VII. RELATED WORK

In this section, we present the state-of-the-art in characterizing the stochastic behavior of multimedia applications and compare it with our approach. We will focus mainly on analytical approaches (refer previous work for survey on simulation based approaches [2], [14]).

A domain-agnostic approach is to statistically analyze the execution variance of soft-real time applications (Kumar et al., [12]). Using profiling, components of the applications that lead to variable execution times are first identified. Then programmers can identify components that affect real-time behavior of the application in the context of other components. Thus, the programmers need not resort to ad-hoc methods for tuning applications for expected real-time behavior.

Our approach differs from the domain agnostic techniques in the following way: our model tightly couples the application and architecture; in what follows, we also discuss how both the stochastic real-time calculus and statistical model checking are an excellent fit for streaming applications. On the other hand, domain-specific techniques in a probabilistic setting perform analysis at task granularity (Yaldiz et al., [28] and Iqbal et al., [8]). In stochastic real-time calculus, our model captures input, execution, and output streams of the SoC. The granularity of the stream object can be at any level: bits, macroblock, frame, and group of pictures.

Iqbal and others [8] proposed scheduling techniques for soft-real time systems. The task execution times are stochastic and the solution for scheduling is based on an online Monte Carlo method on a joint space model of all tasks. The objective of this technique is similar to ours: reducing memory and computational requirements. The complexity of this technique, however, as the authors report, does not scale well for task graphs of huge size. Yaldiz and others [28] use stochastic modeling of the applications to obtain policies for energy savings and for providing probabilities for satisfying timing constraints. Their model considers set of concurrent tasks and takes into account data dependence, precedence relations and timing constraints.

Our methodology differs from the above two discussed approaches in the following way: (1) there is a tight characterization of inputs for streaming applications using arrival curves in the stochastic real-time calculus and randomly generated clips (based on distributions) in statistical model checking, and (2) we provide probabilistic guarantees instead of average-case analysis in the analytical framework.

Liu and others [16] introduced a new concept called approximate variability characterization curves (or Approximate VCCs), to characterize the average-case behavior of multimedia workloads in a parameterized fashion. The crucial difference in comparison to our work is that Approximate VCCs belongs to a family of average-case analysis; instead of probabilistic guarantees on buffer size they bound the error due to their analysis. Also, the framework remains in a deterministic setting after the alteration of workload curves. So, it is not elegantly able to capture stochastic nature of arrivals and stochastic nature of execution times.

## VIII. CONCLUSIONS AND FUTURE WORK

The stochastic real-time calculus and the statistical model checking techniques were presented as frameworks for performance analysis of multimedia systems. The case-study analyzed the trade-off in quality over buffer size savings using the above mentioned approaches. Both the approaches estimated the probability that a certain QoS property is true; the analytical framework upper bounded the estimates from the statistical model checking.

Future extensions to this work are as follows: (a) modeling communication architectures such as bus, network-on-chip, and others; (b) modeling for variable bit-rate video and variable consumption rate; (c) integrating real-time calculus and BIP framework.

## ACKNOWLEDGEMENT

The authors thank Jerome Milan for editing and for providing feedback.

## REFERENCES

- [1] A. Basu, B. Bensalem, M. Bozga, J. Combaz, M. Jaber, T.-H. Nguyen, and J. Sifakis. Rigorous Component-Based System Design Using the BIP Framework. *IEEE Software*, 28(3):41–48, May 2011.
- [2] A. Basu, S. Bensalem, M. Bozga, B. Caillaud, B. Delahaye, and A. Legay. Statistical Abstraction and Model-Checking of Large Heterogeneous Systems. In *Proc. of the International Joint Conference on Formal Techniques for Distributed Systems (FMOODS/FORTE)*, pages 32–46, June 2010.
- [3] A. Basu, M. Bozga, and J. Sifakis. Modeling Heterogeneous Real-time Systems in BIP. In *Software Engineering and Formal Methods SEFM'06 Proceedings*, pages 3–12. IEEE Computer Society Press, 2006.
- [4] S. Bensalem, M. Bozga, B. Delahaye, C. Jégourel, A. Legay, and A. Nouri. Statistical Model Checking QoS Properties of Systems with SBIP. In *ISoLA (1)*, pages 327–341, 2012.
- [5] E. M. Clarke, O. Grumberg, and D. Peled. *Model checking*. MIT Press, 2001.
- [6] T. Héruault, R. Lassaigne, F. Magniette, and S. Peyronnet. Approximate Probabilistic Model Checking. In *Proc. of the Verification, Model Checking and Abstract Interpretation (VMCAI)*, pages 73–84, January 2004.
- [7] K. Huang and L. Thiele. Performance analysis of multimedia applications using correlated streams. In *Proceedings of the conference on Design, automation and test in Europe, DATE '07*, pages 912–917, San Jose, CA, USA, 2007. EDA Consortium.
- [8] N. Iqbal and J. Henkel. SETS: Stochastic execution time scheduling for multicore systems by joint state space and Monte Carlo. In *Proc. of the IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pages 123–130, November 2010.
- [9] P. R. Jelenkovic, A. A. Lazar, and N. Semret. The Effect of Multiple Time Scales and Subexponentiality in MPEG Video Streams on Queueing Behavior. *IEEE Journal on Selected Areas in Communications*, 15(6):1052–1071, August 1997.
- [10] Y. Jiang and Y. Liu. *Stochastic Network Calculus*. Springer, 2008.
- [11] M. Krunz and S. K. Tripathi. On the characterization of VBR MPEG streams. In *Proc. of the International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS)*, pages 192–202, June 1997.
- [12] T. Kumar, R. Cleat, J. Sreeram, and S. Pande. Statistically Analyzing Execution Variance for Soft Real-Time Applications. In J. N. Amaral, editor, *Languages and Compilers for Parallel Computing*, pages 124–140. Springer-Verlag, 2008.
- [13] S. Laplante, R. Lassaigne, F. Magniez, S. Peyronnet, and M. de Rougemont. Probabilistic abstraction for model checking: An approach based on property testing. *ACM Transactions on Computational Logic*, 8(4):30–39, August 2007.
- [14] A. Legay, B. Delahaye, and S. Bensalem. Statistical model checking: an overview. In *Proc. of the International Conference on Runtime Verification (RV)*, pages 122–135, November 2010.
- [15] libmpeg2. A free MPEG2 video stream decoder. <http://libmpeg2.sourceforge.net>, 2006.
- [16] Y. Liu, S. Chakraborty, and W. T. Ooi. Approximate VCCs: a new characterization of multimedia workloads for system-level MpSoC design. In *Proc. of the ACM/IEEE Annual Design Automation Conference (DAC)*, pages 248–253, June 2005.
- [17] A. Maxiaguine, S. Künzli, S. Chakraborty, and L. Thiele. Rate analysis for streaming applications with on-chip buffer constraints. In *Proc. of the Asia and South Pacific Design Automation Conference (ASP-DAC)*, pages 131–136, January 2004.
- [18] B. Raman. *Application-specific workload shaping in resource-constrained media players*. PhD thesis, School of Computing, National University of Singapore, July 2010.
- [19] B. Raman and S. Chakraborty. Application-specific Workload Shaping in Multimedia-enabled Personal Mobile Devices. *ACM Transactions on Embedded Computing Systems*, 7(2):10, February 2008.
- [20] B. Raman, S. Chakraborty, W. T. Ooi, and S. Dutta. Reducing data-memory footprint of multimedia applications by delay redistribution. In *Proceedings of the 44th annual Design Automation Conference, DAC '07*, pages 738–743, New York, NY, USA, 2007. ACM.
- [21] B. Raman, A. Nouri, D. Gangadharan, M. Bozga, A. Basu, M. Maheshwari, J. Milan, A. Legay, S. Bensalem, and S. Chakraborty. A General Stochastic Framework for Low-Cost Design of Multimedia SoCs. Technical Report TR-2012-7, Verimag Research Report, 2012.
- [22] B. Raman, G. Quintin, W. T. Ooi, D. Gangadharan, J. Milan, and S. Chakraborty. On buffering with stochastic guarantees in resource-constrained media players. In *Proc. of the IEEE/ACM/IFIP International Conference on Hardware/Software Codesign and System Synthesis (CODES+ISSS)*, pages 169–178, September 2011.
- [23] L. Santinelli and L. Cucu-Grosjean. Toward probabilistic real-time calculus. *ACM SIGBED Review*, 8(1):54–61, March 2011.
- [24] K. Sen, M. Viswanathan, and G. Agha. Statistical Model Checking of Black-Box Probabilistic Systems. In *Proc. of International Conference on Computer Aided Verification (CAV)*, pages 202–215, July 2004.
- [25] Tektronix. MPEG Elementary Streams. <ftp://ftp.tek.com/tv/test/streams/Element/index.html>, 1996.
- [26] A. Wald. Sequential Tests of Statistical Hypotheses. *Annals of Mathematical Statistics*, 16(2)(2):117–186, June 1945.
- [27] D. Wijesekera and J. Srivastava. Quality of Service (QoS) Metrics for Continuous Media. *Multimedia Tools and Applications*, 3(2):127–166, July 1996.
- [28] S. Yaldiz, A. Demir, and S. Tasiran. Stochastic Modeling and Optimization for Energy Management in Multicore Systems: A Video Decoding Case Study. *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems (TCAD)*, 27(7):1264–1277, July 2008.
- [29] H. L. S. Younes. *Verification and Planning for Stochastic Precess with Asynchronous Events*. PhD thesis, School of Computer Science, Carnegie Mellon University, January 2005.
- [30] N. H. Zamora, X. Hu, and R. Marculescu. System-level performance/power analysis for platform-based design of multimedia applications. *ACM Transactions on Design Automation of Electronic Systems (TODAES)*, 12(1):1–29, January 2007.