

An Ensemble Learning Approach for In-situ Monitoring of FPGA Dynamic Power

Zhe Lin, *Student Member, IEEE*, Sharad Sinha, *Member, IEEE*, and Wei Zhang, *Member, IEEE*

Abstract—As field-programmable gate arrays become prevalent in critical application domains, their power consumption is of high concern. In this paper, we present and evaluate a power monitoring scheme capable of accurately estimating the runtime dynamic power of FPGAs in a fine-grained timescale, in order to support emerging power management techniques. In particular, we describe a novel and specialized ensemble model which can be decomposed into multiple customized decision-tree-based base learners. To aid in model synthesis, a generic computer-aided design flow is proposed to generate samples, select features, tune hyperparameters and train the ensemble estimator. Besides this, a hardware realization of the trained ensemble estimator is presented for on-chip real-time power estimation. In the experiments, we first show that a single decision tree model can achieve prediction error within 4.51% of a commercial gate-level power estimation tool, which is $2.41\text{--}6.07\times$ lower than provided by the commonly used linear model. More importantly, we study the extra gains in inference accuracy using the proposed ensemble model. Experimental results reveal that the ensemble monitoring method can further improve the accuracy of power predictions to within a maximum error of 1.90%. Moreover, the lookup table (LUT) overhead of the ensemble monitoring hardware employing up to 64 base learners is within 1.22% of the target FPGA, indicating its light-weight and scalable characteristics.

Index Terms—Field-programmable gate array (FPGA), dynamic power modeling, in-situ monitoring hardware, feature selection, k-means clustering, decision tree, ensemble learning.

I. INTRODUCTION

FIELD-PROGRAMMABLE gate arrays (FPGAs) are gaining popularity in wide-ranging domains such as cloud services and data centers, where they serve as hardware accelerators for computation-centric tasks. In general, FPGAs offer great benefits compared with traditional application-specific integrated circuits (ASICs), owing to their excellent programmability and short time to market. However, it is observed that applications implemented on FPGAs give rise to a $7\text{--}14\times$ dynamic power overhead compared with the equivalent ASIC implementation [1]. As the architecture complexity and integration density of modern FPGAs continue to grow, the importance of power efficiency increases and FPGA power consumption is turning out to be a key design constraint. To overcome this issue, FPGA power management techniques, which demonstrate huge potential for saving power under a stringent power budget, are receiving considerable interest.

Static power saving techniques, including power gating [2], clock gating [3], dual supply voltage [4], power-aware placement and routing algorithms [5], have been proposed since

an early stage. More recently, runtime power management strategies have gained a surge of interest. Dynamic voltage frequency scaling (DVFS) [6] and task scheduling in FPGA-CPU systems [7], have been reported as successful ways to reduce the runtime power consumption on FPGAs under a tight power constraint. These methods require full awareness of the power consumption in real-time.

An existing method to monitor runtime power consumption is to utilize dedicated power measurement devices. Although these devices can be put into use in modern FPGA systems, they suffer from some disadvantages that limit their scope of applicability. A deficiency is the requirement of additional board areas for the integration of these devices, and a more critical drawback comes with their long power detection period, usually in the order of milliseconds, which exposes their inability to support power detection in a fine-grained timescale (i.e., within hundreds of clock cycles). As a result, power detection relying on on-board measuring devices is only able to support coarse-grained power management.

Emerging technologies in integrated circuit design have motivated the realization of on-chip regulators with a short voltage scaling time, in the order of sub-nanoseconds [8]. With the use of on-chip regulators, many novel power management strategies at fine temporal granularity, such as fine-grained DVFS [8]–[10] or adaptive voltage scaling [11], and dynamic phase scaling for on-chip regulators [12], have been devised for contemporary computing systems. The encouraging progress in fine-grained power management techniques is prompting designers to provide power monitoring schemes within a short time frame. Therefore, it is conceivable that an accurate, fast, yet area-efficient, power detection methodology is indispensable on the roadmap towards fine-grained and efficient power management strategies in FPGA-based systems. What's more, fine-grained runtime power profiling can aid in the identification of power-critical events, which leads to more potential for power pattern analysis and power saving.

In light of the above considerations, we aim to establish an accurate, fine-grained and light-weight dynamic power monitoring scheme on FPGAs. We note that the widely deployed linear power model [13], [14] is unable to adapt itself well to non-linear power behaviors of complex arithmetic units [15], [16]. We therefore develop a novel non-linear decision-tree-based power model, leveraging state-of-the-art machine learning theory. In other words, we use a decision tree regression model to learn non-linear power patterns. The decision tree model demonstrates notably higher estimation accuracy compared with the traditional linear model.

In order to take one step further towards more accurate power estimation, we propose a specialized ensemble model, exploiting one customized decision tree in each of its base

The authors Zhe Lin (zlinaf@ust.hk) and Wei Zhang (wei.zhang@ust.hk) are with the Dept. of Electronic and Computer Engineering, Hong Kong University of Science and Technology (HKUST).

The author Sharad Sinha (sharad_sinha@ieee.org) is with the Dept. of Computer Science and Engineering, Indian Institute of Technology (IIT) Goa. This work was done when he was with NTU, Singapore.

learners. We partition every invocation (i.e., a hardware function call) into a set of execution cycles related to different states in a finite state machine with datapath (FSMD). Thereafter, the k-means clustering algorithm is applied to partition the set of states into multiple homogeneous clusters. Training samples for different states are accordingly divided in the same manner. Following this, an individual base learner is constructed for each cluster of states using the corresponding training samples. All the base learners are finally combined into an ensemble estimator with the overall estimate being the weighted sum of the predictions from these base learners. A generic computer-aided design (CAD) flow is devised for power model synthesis, from generating samples to selecting features, clustering states, tuning hyperparameters and combining all base learners. Furthermore, we propose a light-weight and scalable realization of the trained ensemble estimator, which can be realized in an FPGA for dynamic power monitoring on the fly. We leverage the surplus FPGA resources, apart from the target applications, to implement the proposed monitoring hardware. Our proposed method is fully applicable to commercial FPGAs.

The rest of this paper is organized as follows. Section II gives a comprehensive discussion about modern FPGA power modeling approaches. Section III briefly introduces the background knowledge. Section IV and Section V demonstrate the complete CAD flow for developing the power model. Section VI elaborates the hardware realization of the monitoring circuits and experimental results are discussed in Section VII. Finally, we conclude in Section VIII.

II. RELATED WORK

Recent studies about FPGA power modeling have been conducted at two different abstraction levels — low abstraction and high abstraction. Power modeling approaches at a low abstraction level focus on synthesizing power models in terms of the granularity of interconnects, basic operators or functions. In [17], a macromodel for LUTs and a switch-level model for interconnects were combined to form a monolithic power model for the target FPGA at a low abstraction level. The LUT macromodel is a precharacterized lookup model for power values, which is generated by sweeping over hundreds of input vectors and finally tracking in a table the dynamic power consumption per access to LUT according to the specific input stimuli. The switch-level interconnect model is developed by fitting a switching power formula with capacitance values and the signal transition rate extracted at each node. In [18], the same idea was employed, but power lookup models were generalized for arithmetic operators (e.g., adders, LUT-based multipliers and digital signal processing (DSP)-based multipliers) with different operating frequencies, signal switching activities and data widths.

The development process of the power models at a low abstraction level is usually time-consuming. Low-level power models are also specialized in their targeted FPGA families and devices, making them difficult to migrate to other FPGA families. Moreover, low-level power models, such as those proposed in [17] and [18], usually require intensive computation

to calculate power for every component independently. This issue adds to their inefficiency considering real-time detection.

In contrast, power models obtained from a high abstraction level expedite the developing time, because they can dispense with the need to dig deep into low-level power behaviors related to the devices' internal structures. Besides this, the runtime computational complexity of high-level models is much lower, widening their applicability in real-time power estimation on FPGAs. In high-level power modeling methods, signal switching activities are extracted in different time periods, and simultaneously, runtime power values are captured for the corresponding periods. Regression estimators are then trained to associate different signal activities with specific power patterns.

Most of the recent work investigating high-level power models on FPGAs [13], [14], [19] has employed a linear regression model. In [13], the IO toggle rate and resource utilization were synthesized into a linear regression model in a log-arithmetic format. This work attempted to formulate a generic power model which is apt for all applications. Nevertheless, design parameters deviating from the training applications, such as inconsistency in the IO port width, could significantly degrade the estimation accuracy. The work [14] leveraged the switching activities of a set of internal signals and the corresponding power values to form a linear model specifically for a single application, with the runtime computation completed by a software program running on an embedded processor in an FPGA-CPU system. Based on [14], the work [19] applied an online adjustment scheme and reported higher accuracy. However, this work leaned on on-board power measurement, making it almost unsuitable for fine-grained scenarios. The hardware realization of the signal monitoring in [14] and [19] led to an LUT resource overhead of 7% and 9% in terms of the tested applications, respectively, as well as incurred a workload of 5% of CPU time. Furthermore, as studied in [15] and [16], the power behaviors of complex arithmetic units are potentially non-linear. As a consequence, the up-to-date fine-grained linear power model exhibits intrinsic restrictions on achieving high accuracy when non-linear power behaviors are discovered with the increasing training size. This is a typical problem known as *underfitting* in the machine learning theory, which implies that the applied model is too simple in its characteristics to rigorously fit an intricate training set.

Within this context, we target constructing a non-linear dynamic power model to support in-situ power monitoring, providing high estimation accuracy while inducing low overheads. To the best of our knowledge, ours is the first work to introduce an ensemble power model at a high abstraction level for FPGAs. As an extension of our prior work [20], we make a substantial improvement in power estimation accuracy by devising a novel ensemble approach, wherein each of the decomposed base learners specializes in the power prediction of particular states in the FSMD.

III. BACKGROUND

A. FPGA Power

The power consumption of an FPGA can be decomposed into two major components: static power and dynamic power.

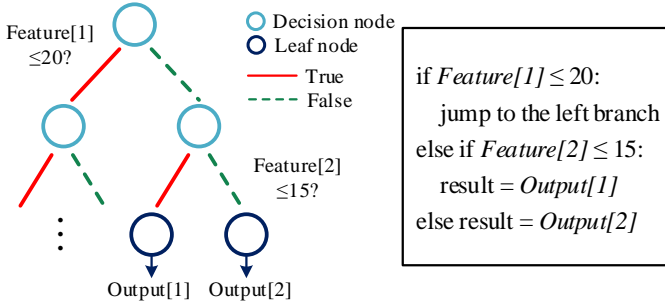


Fig. 1. Graphical and textual representation of a decision tree.

Static power refers to the leakage power consumption when an FPGA is powered on while no circuit activities are involved. It is highly dependent on process, voltage and temperature (PVT). Dynamic power, on the other hand, is introduced by signal transitions which dissipate power by repeatedly charging and discharging the load capacitors. Equation (1) formulates dynamic power P_{dyn} as

$$P_{dyn} = \sum_{i \in I} \alpha_i C_i V_{dd}^2 f, \quad (1)$$

which is a function of signal switching activity α_i , capacitance C_i on the net i within the whole set of nets I , supply voltage V_{dd} and operating frequency f .

B. Decision Tree and Ensemble Learning

The decision tree is a nonparametric hierarchical model in supervised learning for both classification and regression. It learns the samples in the form of a tree structure. A tree node can be categorized as (1) a leaf/terminal node — a node associated with an output result (i.e., a class label for classification or a value for regression) — or (2) a decision/internal node — an intermediate node to decide on one of its child nodes to go to. The training process is to determine an if-then-else decision rule for every decision node and an output value for every leaf node, as shown in Fig. 1. To make a new inference for an unsolved case, the decision tree firstly starts with the root node and moves to one of its child nodes. This process executes iteratively, until finally a leaf node with inference result is reached. From the aspect of hardware implementation, the inference process of a decision tree can be eventually decomposed into a series of comparisons based on decision rules, and it is intrinsically easy to be implemented in a hardware-friendly manner.

Ensemble algorithms combine the predictions of several individually trained base learners in order to improve model robustness and prediction accuracy over a single estimator. In ensemble algorithms, decision trees are widely deployed as effective base learners, e.g., in random forests [21].

C. K-means Clustering

Clustering is a classic unsupervised problem in machine learning theory, with the objective to separate a finite unlabeled sample set into a number of groups according to some sort of homogeneity in the samples. K-means clustering is a popular clustering method that has been widely used in data mining and industrial applications. “K-means clustering algorithm” usually refers to Lloyd’s algorithm [22]. In k-means, the

cluster number K is a user-specified parameter. The algorithm seeks to determine K centroids in the centroid set C that minimize the within-cluster sum of squares for data points in the sample set X , as represented in Equation (2):

$$\sum_{x \in X} \min_{c \in C} \|x - c\|^2. \quad (2)$$

The k-means algorithm firstly begins by arbitrarily choosing k samples as the initial centroids. Each remaining sample is then assigned to its nearest centroid. After all samples are associated with centroids, every centroid is updated as the mean value of all the samples assigned to it. These two steps, namely, sample assignment and centroid update, are repeated until all centroids converge at some specific locations.

To expedite the convergence of k-means, k-means++ [23] is used in this paper, which improves the initialization of centroids. It selects a point as a centroid based on a probability proportional to this point’s contribution in the within-cluster sum of squares. This method aims at making the centroids distant from each other, resulting in a higher probability of faster convergence than randomly selecting samples to be centroids in the conventional k-means algorithm.

IV. COMPUTER-AIDED DESIGN FLOW

We propose a CAD flow to establish the power model using state-of-the-art machine learning techniques. The CAD flow is shown in Fig. 2. The design flow starts from the synthesis, placement and routing of the design. Note that in this design flow, we use the mapped hardware description language (HDL) netlist as the source, which can be exported after placement and routing. It basically consists of connections and primitives. The main advantages of using the HDL netlist are two-fold: firstly, it enables us to have access to lower level signals and further extract the indicative signals for power monitoring; and secondly, it preserves the mapping of the design, and therefore, the later integration of the monitoring hardware does not permute the original mapping.

The complete design flow can be decomposed into three sub-flows: (1) the activity trace flow (ATF), (2) the power trace flow (PTF) and (3) the model synthesis flow (MSF). We generate activity traces and power traces using the ATF and PTF, respectively. An activity trace is defined to be the set of switching activities of some monitored signals used as power pattern indicators, while a power trace provides the average power related to a particular activity trace. An activity trace and the corresponding power trace are combined as a training sample in the MSF.

A. Activity Trace Flow (ATF)

Generally speaking, the behaviors of a digital system implemented in an FPGA can be described as an FSMD [16], [24], [25], as shown in Fig. 3, which can be decomposed into several finite state machine (FSM) states to provide control signals for the datapath. In this paper, we limit our scope to FSMD-based hardware designs.

To build an ensemble model, we monitor the signal activities in the granularity of a state. To support this, FSM identification should be conducted to automatically extract the

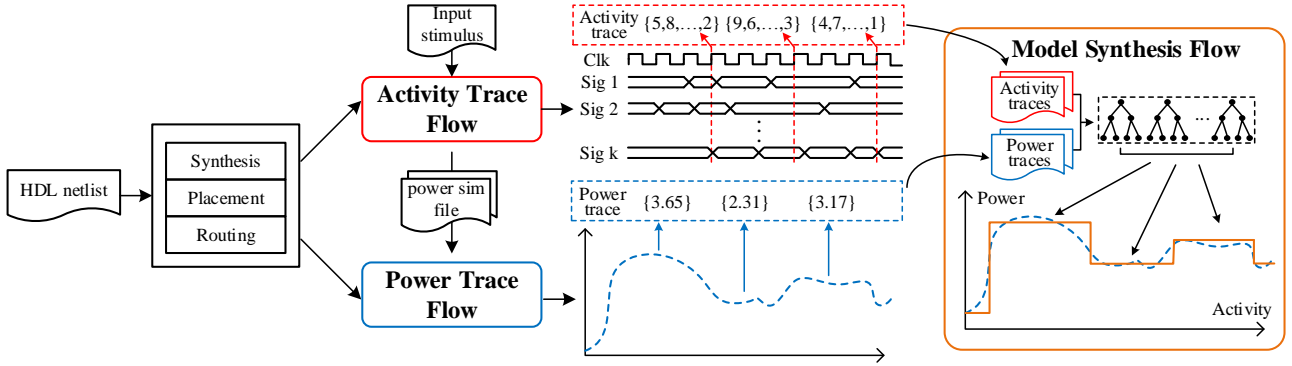


Fig. 2. Overview of the computer-aided design flow.

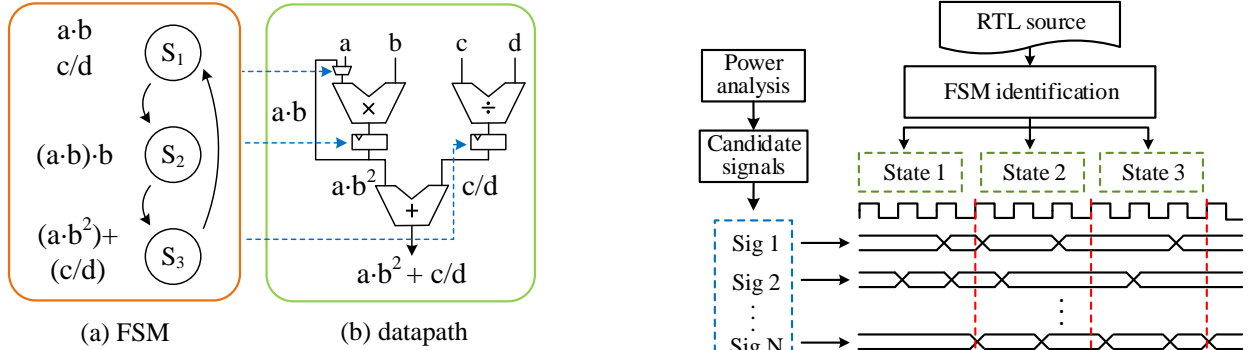


Fig. 3. An example of FSMMD implementation.

state register of an FSMMD. We note that Odin II [26], the Verilog synthesis tool in Verilog-to-Routing (VTR) [27], has implemented an FSM identification algorithm by finding a combinational feedback loop from the candidate register to itself, based on a form of abstract syntax tree (AST) [28]. The AST profiles the Verilog design and decomposes it into a hierarchical representation of different syntaxes, which can be used to identify different functional structures. We extend the basic FSM identification algorithm in Odin II to analyze feedback loops between more than one candidate register so that it is possible to identify FSMs written in different styles (e.g., using `current_state` and `next_state` registers in an FSM).

In the ATF shown in Fig. 4, we aim to extract a key set of single-bit signals, out of the hundreds of internal and IO signals in the netlist, whose activities show a great influence on the dynamic power. To select the signals to monitor for general-purpose applications, we run a vector-based timing simulation with randomly generated input vectors. The simulation only stops when the simulation time is two orders of magnitude longer than the execution time of each invocation. Moreover, we incorporate a counter for each state to record its coverage. In the experiments, we have covered all states in the evaluated benchmarks. More sophisticated methods [29] can be used to expedite full state traversal, which is complementary to our work and does not affect the design flow.

Signal activities are quantified and sorted through power analysis. We select from the signal list an abundance of candidate signals ($\geq 1,000$) showing the highest switching activities to monitor, because they tend to show wide-ranging behaviors matching dynamic power patterns. In Section VII-C, we find that the necessary signal number is orders of magnitude smaller than the number of monitored candidate signals.

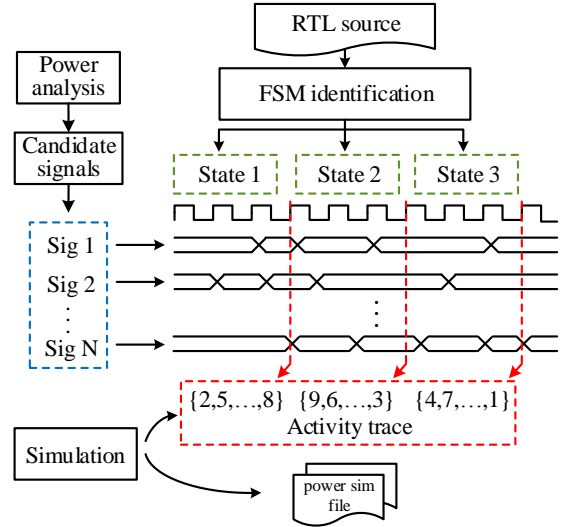


Fig. 4. ATF overview.

An activity trace records the activities of all the monitored signals in a sampling interval, which can be formulated as an activation function $act(\cdot)$ in the form

$$act(s) = sw(s, t_{end}) - sw(s, t_{start}). \quad (3)$$

Therein, the $act(\cdot)$ computes the change in the switching activities $sw(\cdot)$ of the signal vector s over the sampling interval from t_{start} to t_{end} . For each application, we extract the activity traces in the granularity of every state. At the same time as the activity traces are captured, we export power simulation files (`.saif`) that are used for power estimation.

B. Power Trace Flow (PTF)

In the PTF, we use the FPGA power estimator to perform power estimation for different time steps using power simulation files, based on the post-route design. Our main emphasis is on accurately modeling the FPGA dynamic power consumption, in that dynamic power tends to change more significantly for an applications in different time steps. Take the evaluated Xilinx Virtex-7 FPGA as an example: the static power consumption is within 2 W under the temperature margin from -40 to 80 degrees centigrade with Xilinx high-performance low-power (HPL) technology [30], whereas the dynamic power can change drastically from less than 1 W to more than 20 W in terms of a wide range of applications with different characteristics and resource utilization.

C. Model Synthesis Flow (MSF)

In the MSF, the anchor is to establish a power estimation model on the basis of the up-to-date machine learning theory, using data samples generated in the ATF and PTF. The flow for power model establishment is depicted in Fig. 5. To begin with, feature selection is performed to diminish redundant features to monitor. Second, a k-means clustering method is applied to divide the states along with their related training samples into different clusters. Thereafter, k-fold cross validation is deployed to determine the most suitable set of hyperparameters for the decision tree estimator in every base learner. Finally, we train the base learners and combine them through a specialized ensemble approach.

V. ENSEMBLE MODEL ESTABLISHMENT

This section describes the aforementioned four steps (feature selection, FSM state clustering, hyperparameter tuning and model ensemble) in the MSF for the construction of an ensemble power model. Regarding the establishment of a single decision-tree-based model, feature selection, hyperparameter tuning and decision tree training are done in the MSF, whereas FSM state clustering and model ensemble are skipped.

A. Feature Selection

A feature is defined as the switching activity of a monitored signal in an activity trace. Noticing that high switching features extracted from the ATF may be correlated (e.g., an input and an output of the same LUT) or show repetitive patterns (e.g., the clock signal), we leverage a *recursive feature elimination* algorithm to identify the key subset of features across multiple input vectors from various invocations. Recursive feature elimination seeks to reduce the number of features recursively, and has been successfully applied to other domains [31], [32]. It requires developers to specify an estimator capable of providing attributes of feature importance (e.g., coefficients for linear model). Thus, we target the CART decision tree [33], which is able to provide a Gini importance value to quantify the importance of each feature.

Taking the complete feature set as input, recursive feature elimination first trains a decision tree to assess the importance of every feature and prunes the features with the lowest importance values from the generated decision tree. In this work, we set the feature reduction size as one, which means that the algorithm only filters out one feature in each iteration. After that, the remaining features are used to retrain the decision tree in order to update the Gini importance of the new feature set. The least important features are further trimmed away. These two steps work repeatedly, considering smaller and smaller sets of features in each subsequent iteration until eventually a user-specified number of critical features remains. In experiments, we study the impact of the number of selected features on the estimation accuracy of the ensemble model. Experimental results, presented in Section VII-C, show that around 30 key features out of more than 1,000 candidates are sufficient to achieve high accuracy.

B. FSM State Clustering

During the synthesis stage, specific resources are scheduled and bound with each state in the FSM. It is reasonable to

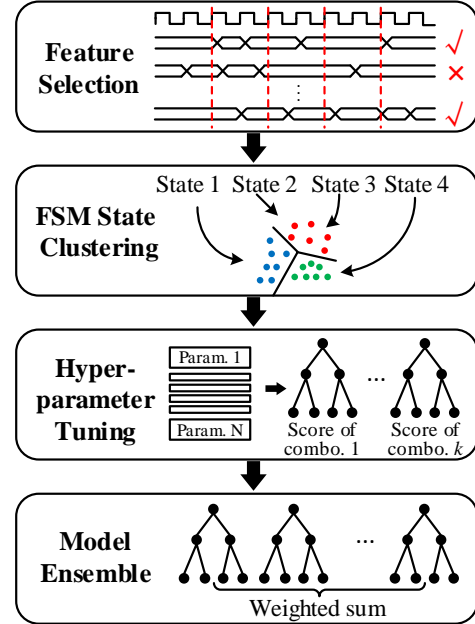


Fig. 5. MSF overview.

infer that execution in different states exhibits different power consumption. This can be explained by the variation in the states' operations and utilized resources. For instance, in the datapath shown in Fig. 3(b), the operations (e.g., addition, multiplication and division) invoked in various states lead to different and state-specific power consumption. Moreover, the power changes indicated by some infrequently changing but critical control signals (e.g., the start bit) are reflected in the state transitions.

Given an FSM state set containing N states, $S = \{s_1, s_2, \dots, s_N\}$, and a user-specified cluster number K , we apply the k-means algorithm to partition the states into different clusters within the cluster set $C = \{c_1, c_2, \dots, c_K\}$ ($K \leq N$). The relationship between the FSM state set S and the cluster set C can be formulated as a function $r(\cdot)$:

$$r : S \times C \rightarrow \{0, 1\}. \quad (4)$$

Therein, $r(s_i, c_j) = 1$ means state s_i is assigned to cluster c_j , and vice versa. All the clusters in the set C possess the following attribute:

$$c_i \cap c_j = \phi, \quad \forall i \neq j. \quad (5)$$

This attribute ensures that each FSM state is assigned to one and only one cluster by k-means. Provided sufficient activity patterns from features, another attribute will also be satisfied. That is

$$c_i \neq \phi, \quad \forall i \in \{1, 2, \dots, K\}. \quad (6)$$

The clusters are developed on the basis of the homogeneity of signal activities with respect to different FSM states. As a result, the FSM states in the same cluster show high similarity in signal activities, thus revealing a close resemblance of operations (some are due to resource sharing) and power characteristics. Following this observation, we create K independent training sets in the same manner: the training samples corresponding to all the FSM states in the same k-means cluster will be grouped into one training set. These sample

sets are non-overlapping sets constrained by Equation (5). In the development process of the ensemble model, we separately build one base learner for every sample set. At inference time, each base learner predicts power solely for the cluster of states it was trained with.

C. Hyperparameter Tuning

We build each base learner in the form of a decision tree. The problem of learning an optimal decision tree is known to be NP-complete under several aspects of optimality. However, there are some essential hyperparameters largely determining the model accuracy. We seek to tune a set of key hyperparameters to best suit the training set in order to enhance decision tree performance. The target hyperparameter set is listed in Table I. The most influential hyperparameter is *Maximum_depth*, which determines the number of nodes along the path from the root node down to its farthest child node. A deeper tree indicates a higher degree of model complexity, while also making it inclined to overfit on data. In contrast, a shallower tree often fails to fit in well with a complex training set. In light of this problem, the tree depth should be customized according to different characteristics of training sets in order to strike a good balance between training and testing accuracy. Likewise, the other three hyperparameters, *Minimum_split_sample*, *Minimum_leaf_sample* and *Minimum_leaf_impurity*, should be coordinated with the tree depth and used to circumvent overfitting.

To automatically tailor the best hyperparameter set for the decision tree in every base learner, k-fold cross validation [34] is used. K-fold cross validation performs self-testing using the training set and offers a score to evaluate the performance of an estimator with a particular set of hyperparameters. We employ ten-fold cross validation, in which the training set is equally split into ten subsets of the same size. In each iteration, one subset is selected as a validation subset, and the remaining subsets become the training subsets. The training subsets are used to construct a decision tree given a specific set of hyperparameters in Table I, while the validation set is applied to give a score for this hyperparameter set in the form of negative mean absolute error. Ten iterations are conducted with a unique validation set each time. The overall score of the evaluated hyperparameter set is the average score of the ten iterations.

We first sweep over a large hyperparameter space through cross validation and filter out the hyperparameter sets consistently leading to low scores across different training sets, as well as those offering no improvement in comparison to trees with lower complexity, namely, trees with smaller values of *Maximum_depth*. The hyperparameter space is eventually trimmed down as follows: {3, 4, 5, 6, 7, 8} for *Maximum_depth*, {5, 10, 15, 20} for *Minimum_split_sample* and *Minimum_leaf_sample*, and {0.001, 0.01, 0.02, 0.03, 0.04, 0.05} for *Minimum_leaf_impurity*.

To train each base learner, we again perform cross validation to quantify the performance of all combinations in the hyperparameter space and finally select the hyperparameter set with the highest score ranking to build a decision tree. After cross

TABLE I
HYPERPARAMETERS FOR DECISION TREE TUNING.

Name	Description
Maximum depth	The maximum depth that a tree can grow to.
Minimum split sample	The minimum number of samples used to split a decision node.
Minimum leaf sample	The minimum number of samples necessary to determine a leaf node.
Minimum leaf impurity	The minimum percentage of samples giving different output at a leaf node.

validation, we use the complete training set to train the base learner and use a test set to quantify the accuracy.

D. Model Ensemble

From the prior steps, we separately develop a decision tree in a base estimator for every cluster of FSM states generated by the k-means algorithm. Thereafter, we introduce a specialized ensemble learning method to combine different base estimators together so as to estimate the invocation-level power consumption. We take the weighted sum of the state-related power predictions to deduce the overall power, which is the same as getting the average power over different states. In principle, ensemble learning through averaging, which is known as *bagging*, is expected to provide higher prediction accuracy than an individual model [35], [36]. In theory, our customized ensemble model can be proven to improve accuracy compared with an individual model.

We define the error-free target function and the prediction function for an individual invocation-level model as $h_{inv}(\cdot)$ and $y_{inv}(\cdot)$, respectively. In addition, the error-free target function and prediction function of a base estimator are defined as $h_i(\cdot)$ and $y_i(\cdot)$, respectively, where $i \in \{1, 2, \dots, K\}$ denotes the index of the cluster/base learner. The inference result of the single invocation-level estimator and each of the base estimators can be written as

$$y_{inv}(\mathbf{x}) = h_{inv}(\mathbf{x}) + \epsilon_{inv}(\mathbf{x}), \quad (7)$$

$$y_i(\mathbf{x}) = h_i(\mathbf{x}) + \epsilon_i(\mathbf{x}), \quad (8)$$

where $\epsilon_{inv}(\cdot)$ and $\epsilon_i(\cdot)$ denote the corresponding error functions and \mathbf{x} is the feature vector. We use the weighted sum of the error-free functions for the base estimators to get the error-free target function $h_{ens}(\cdot)$ for the ensemble model, in which the weight of a base estimator is defined as the ratio of the prediction cycles for this base estimator to the total execution cycles in an invocation. That is

$$h_{ens}(\mathbf{x}) = \sum_{i=1}^K \frac{t_i \cdot h_i(\mathbf{x})}{T}, \quad (9)$$

where T is defined as the total execution cycles in an invocation, t_i as the prediction cycles for each of the base estimators i and K as the number of clusters preset for the k-means algorithm, which also represents the number of base

estimators. Regarding the single invocation model, the sum-of-squares error can be given by

$$E_{inv} = \mathbb{E}_{\mathbf{x}}[\{y_{inv}(\mathbf{x}) - h_{inv}(\mathbf{x})\}^2] = \mathbb{E}_{\mathbf{x}}[\epsilon_{inv}(\mathbf{x})^2]. \quad (10)$$

To calculate the sum-of-squares error for the ensemble model, we assume no less than two clusters. Three constraints have to be satisfied and can be written as

$$2 \leq K \leq T, \quad (11)$$

$$1 \leq t_i \leq T, \quad (12)$$

$$\sum_{i=1}^K t_i = T. \quad (13)$$

We calculate the weighted sum of the predictions from the base estimators as the overall estimation for the ensemble model. The sum-of-squares error of the ensemble model can thus be formulated as

$$\begin{aligned} E_{ens} &= \mathbb{E}_{\mathbf{x}}[\{\sum_{i=1}^K \frac{t_i \cdot y_i(\mathbf{x})}{T} - h_{ens}(\mathbf{x})\}^2] \\ &= \mathbb{E}_{\mathbf{x}}[\{\sum_{i=1}^K \frac{t_i \cdot h_i(\mathbf{x})}{T} + t_i \cdot \epsilon_i(\mathbf{x}) - \sum_{i=1}^K \frac{t_i \cdot h_i(\mathbf{x})}{T}\}^2] \\ &= \mathbb{E}_{\mathbf{x}}[\{\sum_{i=1}^K \frac{t_i \cdot \epsilon_i(\mathbf{x})}{T}\}^2]. \end{aligned} \quad (14)$$

For the sake of simplicity, we assume the error of different base estimators is uncorrelated and has zero mean [16], [35]. In addition, the sum-of-squares error for each of the estimators is identical [16]. These assumptions can be formulated as

$$\mathbb{E}_{\mathbf{x}}[\epsilon_i(\mathbf{x})] = 0, \quad \forall i \in \{1, 2, \dots, K\}, \quad (15)$$

$$\mathbb{E}_{\mathbf{x}}[\epsilon_i(\mathbf{x}) \cdot \epsilon_j(\mathbf{x})] = 0, \quad \forall i \neq j, \quad (16)$$

$$\mathbb{E}_{\mathbf{x}}[\epsilon_i(\mathbf{x})^2] = \mathbb{E}_{\mathbf{x}}[\epsilon_{inv}(\mathbf{x})^2] = \mathbb{E}_{\mathbf{x}}[\epsilon(\mathbf{x})^2], \quad \forall i \in \{1, 2, \dots, K\}, \quad (17)$$

where $\epsilon(\cdot)$ defines a generic error term. With the three assumptions given above, the sum-of-squares error of the ensemble model can be simplified as

$$E_{ens} = \frac{\sum_{i=1}^K t_i^2}{T^2} \cdot \mathbb{E}_{\mathbf{x}}[\epsilon(\mathbf{x})^2] = \frac{\sum_{i=1}^K t_i^2}{T^2} \cdot E_{inv}. \quad (18)$$

Taking the conditions (11), (12) and (13) into account, we apply *Root-Mean Square-Arithmetic Mean Inequality* to deduce the range of the constant term of E_{ens} , which is given as

$$\frac{\sum_{i=1}^K t_i^2}{T^2} = \frac{\sum_{i=1}^K t_i^2}{(\sum_{i=1}^K t_i)^2} \in [\frac{1}{K}, 1]. \quad (19)$$

The left part in (19) is minimized when t_1, t_2, \dots, t_K have the equal value $\frac{T}{K}$. Under this context, the relationship between E_{ens} and E_{inv} can be represented as

$$E_{ens} = \frac{1}{K} E_{inv}. \quad (20)$$

Equation (20) presents the theoretical proof of the power modeling problem that determines the lowest possible error using an ensemble model, with a single invocation-level model as the baseline. It demonstrates that our proposed ensemble model can intrinsically improve the average error by a factor of K . Nevertheless, the proof is dependent on a fundamental assumption that the sum-of-squares error incurred by different estimators is uncorrelated, which is the ideal situation. From a practical standpoint, it is very difficult or almost impossible to completely eliminate model correlation, even though many researchers have been trying to tackle this issue [37], [38].

In this work, we enhance the performance of the specialized ensemble learning by promoting the error diversities of different learners from two aspects. Firstly, *training data diversity* is introduced through a specialized resampling method: every base learner is trained with a unique sample set derived from a non-overlapping cluster of FSM states. Secondly, *hyperparameter diversity among base learners* is taken into consideration by customizing the best-suited hyperparameter set for each of the decomposed base learners through k-fold cross validation.

VI. MONITORING HARDWARE

We propose monitoring hardware to implement the trained ensemble estimator on chip and achieve dynamic power prediction on the fly. Recall that the training of the ensemble model was finished by the CAD flow discussed in Section V. Dedicated hardware modules are devised for activity detection, state identification and power estimation at runtime. The ensemble estimator takes both the signal activities and the current state as input for power prediction, as shown in Fig. 6.

A. Preprocessing Units

A fundamental step for the ensemble hardware to function correctly is to extract activities and separate them into features related to different FSM states. To achieve this goal, three preprocessing units are designed before invoking the ensemble estimator, which are shown in the green boxes in Fig. 6.

Activity counter: Activity counters are used to capture the switching activities at runtime from the selected signals in the CAD flow. An activity counter mainly comprises a positive edge detector and a counter, as shown in Fig. 7. The positive edge detector is used to detect the positive transitions of the input signal. Its output is valid for a single cycle, acting as the enable signal for the counter. We count the positive edges of the selected signals, and thereby the number of clock cycles in a sampling period is the upper bound of the signal activities, which can be used to determine the maximum bit width for the counters. We uniformly set the width as 20 bits, which can cover a wide range of sampling periods. Two types of counters are realized: an LUT counter and a DSP counter. The LUT counter is written in HDL to utilize only LUTs and flip-flops (FFs) whereas the DSP counter is instantiated using the primitive *COUNTER_LOAD_MACRO* for Xilinx devices, a dynamic loading up counter occupying one DSP48 unit with a maximum data width of 48 bits. These two counter templates offer elasticity to developers who are aware of the application resource utilization, so that it is possible to avoid excessive use of a single type of resource.

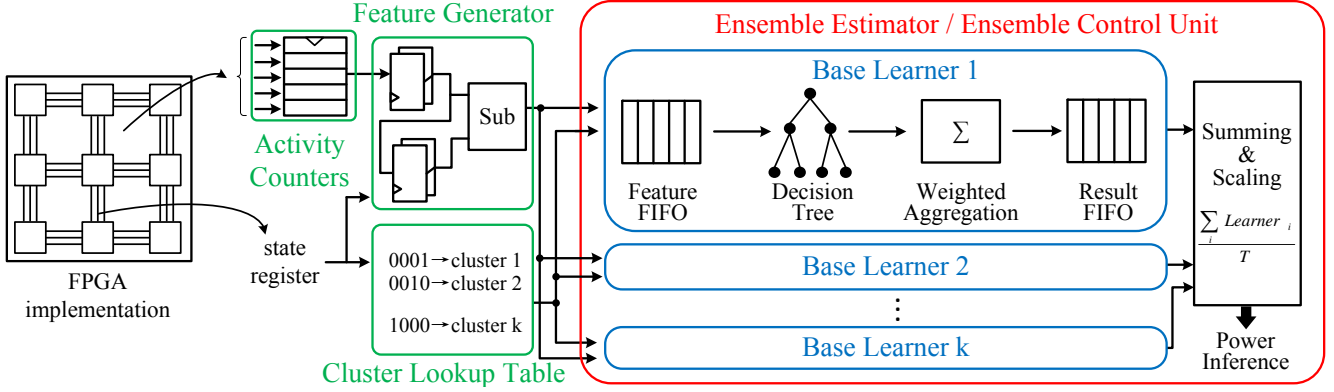


Fig. 6. Overview of the simplified real-time monitoring hardware.

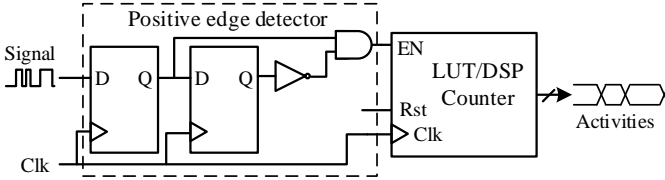


Fig. 7. Activity counter.

Feature generator: A feature generator is used to divide the switching activities into chunks specifically for each FSM state. It records the values of signal activities at the cycles when entering and exiting a state. During state transition, the feature generator performs subtraction for the two buffered activities as the way of generating features for each state.

Cluster lookup table: The cluster lookup table stores the clustering information which is formulated as the function $r(\cdot)$ in Equation (4). Taking the state register value as input, the cluster lookup table identifies the cluster to which the executing state belongs, and generates an enable signal to correctly write the features from the feature generator to the buffer of the relevant base learner.

B. Ensemble Control Unit

The ensemble control unit orchestrates different base learners for invocation-level power estimation. For each base learner, a first-in-first-out (FIFO) is used to temporarily buffer the features related to all the states within the specific cluster associated with this base learner. A decision tree regression engine, the key element in the ensemble model, implements the power prediction function $y_i(\mathbf{x}_{s_j})$, where i ($1 \leq i \leq K$) is the index of the cluster/base learner and \mathbf{x}_{s_j} denotes the feature vector \mathbf{x} for the specific state s_j . A summation unit is used to compute the weighted aggregation power value for a base learner. This function is defined as

$$p(c_i) = \sum_{j=1}^N r(s_j, c_i) \cdot t(s_j) \cdot y_i(\mathbf{x}_{s_j}), \quad (21)$$

where j ($1 \leq j \leq N$) is the index of a state, $r(\cdot)$ is the function from Equation (4) and $t(s_j)$ denotes the function to extract the execution cycles in state s_j . Since the base learners may not operate at the same time due to the inconsistent arrival time of the input features, a result FIFO is inserted following the weighted aggregation unit in each base learner for the purpose of result sequence alignment. After all the

base learners provide their predictions, the invocation-level power estimation is obtained by summing all the individual predictions and scaling down the sum by a factor of T so that we get the average power among all execution cycles in an invocation. The overall power estimation is formulated as

$$P_{ens} = \frac{1}{T} \cdot \sum_{i=1}^K p(c_i). \quad (22)$$

C. Decision Tree Regression Engine

The decision tree regression engine serves as the principal element in a base learner. We propose a memory-based decision tree regression engine, as shown in Fig. 8. There are some studies on decision tree implementation in hardware [39], [40] to maximize the throughput of decision tree computation. Our objective, however, is different from that of prior works in that the prime consideration is not throughput in our solution. Instead, we customize a decision tree structure to reduce the power and resource overheads of our monitoring hardware in the first place. The decision tree structure is completely preserved in a memory element. Additional peripheral control units are incorporated to orchestrate the feature control, tree node decoding and branch decision. To summarize, the decision tree structure in our proposed solution can be further decomposed into three subsystems: (1) a feature controller, (2) a decision tree FSM and (3) a decision tree structure memory.

To achieve power prediction, the feature controller buffers feature values from the feature FIFO and invokes the FSM's operating states. The decision tree FSM has four states: idle (I), node read (N), stalling (S) and result output (R). The FSM starts execution by transferring the state from idle to node read. In the node read state, the FSM fetches the node information and tree structure from the memory and completes the if-then-else branch decision by comparing the addressed feature with the decoded tree node coefficient. The stalling state will operate together with the node read state to ensure the correctness of memory reading. The execution finally terminates by switching to the result output state which sends out the estimation result and sets an indication signal high when the tree leaf has been reached. The decision tree structure memory shown in Fig. 9 is the fundamental component which preserves the tree information. The memory is implemented as the block memory. For a decision node, the structure memory

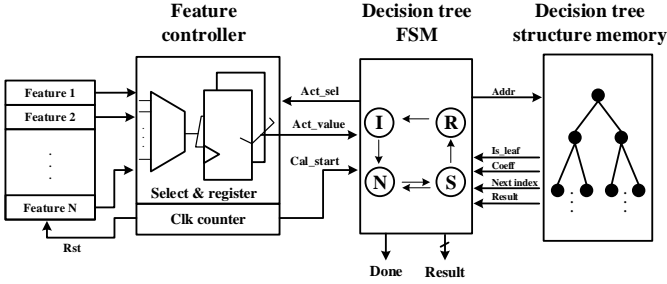


Fig. 8. Decision tree regression engine.

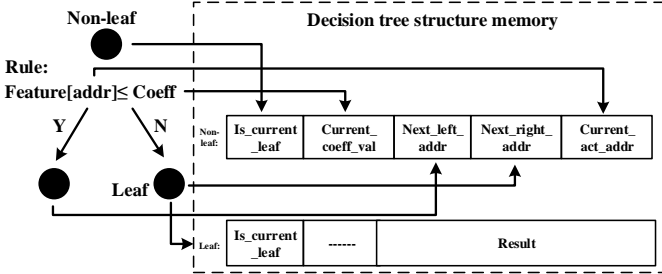


Fig. 9. Decision tree memory structure.

stores the feature address, the decision rule and child node addresses, and for a leaf node, it records the output value.

The maximum execution time for a single estimation is no more than $2n + 1$ cycles, where n denotes the maximum depth of the target decision tree. The decision tree structure is completely preserved in the structure memory, meaning that the proposed decision tree implementation is generally applicable to all decision tree types, varying in depths or pruning structures. Moreover, recall that a feature is the activities of a monitored signal, and therefore, it is a non-negative integer number. Correspondingly, the coefficients for comparisons in the decision rules can be represented as integers, without loss of precision. Following this observation, no floating point operations are required for the decision tree design, which contributes to the high area efficiency of its implementation, as opposed to the traditional linear model which usually requires floating point weights as well as floating point operations. Any simplification of the floating point calculation will jeopardize the precision of the linear model.

For the ensemble monitoring hardware, the overall prediction time with the largest number of base learners is 21 cycles, which can be used for power prediction under the fine granularity of several hundreds of cycles. To build a single decision-tree-based estimator, the hardware design is simplified: only activity counters and a decision tree are required. The decision tree regression engine directly fetches input features from the activity counters at the end of each sampling period and the feature controller is responsible for resetting the activity counters at the beginning of each sampling period.

VII. EXPERIMENTAL RESULTS

Our proposed activity trace flow is implemented in Vivado 2016.4 and the generation of the power simulation files (*.saif*) is completed in Modelsim SE 10.3. The power traces are generated using Vivado power analyzer at a high confidence level [41], under an ordinary temperature of 25 degrees centigrade. Note that in our experiments, we target the Xilinx

tool chain and Modelsim, but the generic methodology is applicable to other vendor tools, such as Intel Quartus Prime.

The model synthesis flow is developed based on the Scikit-learn 0.18.1 [42] machine learning toolbox. We applied our methodology to establish both the single decision-tree-based and ensemble power models, and we systematically quantify the prediction accuracy using sets of benchmarks from Polybench [43], CHStone [44] and Machsuite [45], which are categorized by their utilized resources.

We define the following three types of benchmarks: LUT-based, DSP-based and hybrid benchmarks. The LUT-based benchmarks mainly use LUT resources, whereas the DSP-based benchmarks utilize DSPs as the major resource type. The hybrid benchmarks show a relatively balanced proportion of resource utilization of both LUTs and BRAMs or DSPs. These benchmark suites are C-based benchmarks, and we generate the synthesizable Verilog version using Vivado-HLS 2016.4. All the benchmarks are tested under a clock period of 10 ns. In the training process, we generate training samples using random distribution or normal distribution. We separate the training samples into two sets at random: 80% of them are used as the training set and the remaining 20% are for the purpose of testing. Customized monitoring hardware is constructed together with each benchmark on the target FPGA platform, Virtex-7 XC7V2000tflg1925-1.

In the following subsections, we first study the accuracy improvement offered by the decision tree model over the traditional linear model, and we analyze the overheads of the decision tree monitoring hardware with respect to resources, operating frequency and power. On top of that, we explore to what extent the ensemble model can further boost the estimation accuracy compared with the single decision tree model. First, we investigate the impact of feature number on power estimation accuracy. Second, we study the additional gains in accuracy as the number of base learners increases. Last, we analyze the tradeoff between overheads induced by the ensemble model and the attainable power prediction accuracy.

A. Decision Tree Model: Accuracy

For the single decision tree model, we collect 2,000 samples for each benchmark. We set the sampling period to be 3 μ s, which provides a good tradeoff between prediction granularity and latency ($2 \times tree_depth + 1$ cycles). To achieve a fair comparison between our proposed model and the linear model [13], [14], we tailor a recursive feature elimination method for the linear model, in which the feature importance is quantified by the coefficients in the generated linear equation.

The resource utilization and the mean absolute error (MAE), in percent, for dynamic power consumption is shown in Table II. The average MAE percentage is 3.86% and the maximum MAE percentage is 4.51% for our proposed decision tree model, whereas those of the linear model are 14.59% and 19.83%, respectively. The decision tree model shows an improvement in estimation error by 2.41–6.07 \times in comparison to the linear model. From Table II, we can see that the improvement offered by the decision tree over the linear model is more significant for DSP-based designs, because the

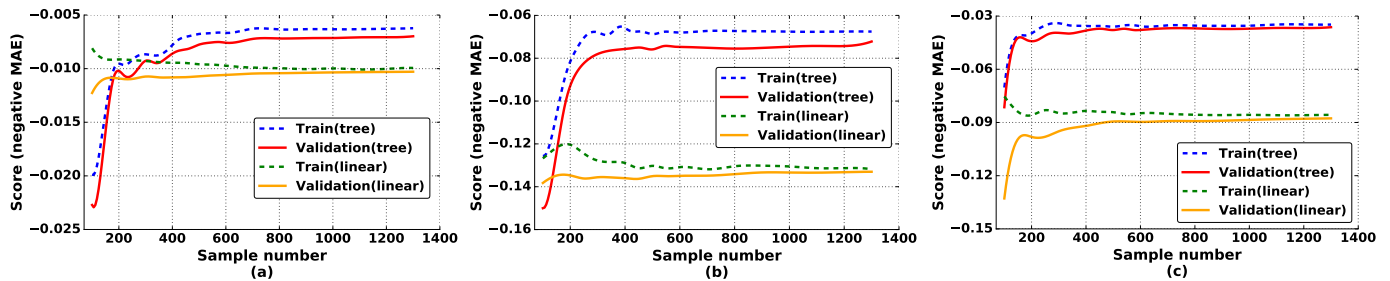


Fig. 10. Learning curves of the decision tree model and linear model: (a) LUT-based; (b) DSP-based; and (c) Hybrid.

TABLE II
RESOURCE UTILIZATION AND MODEL ACCURACY.

Benchmark	Resource			MAE (%)	
	LUT ¹	DSP ²	BRAM ³	Lin. [14]	Dtree
Atax [*] [43]	58691	0	0	12.46	4.14
Bicg [*] [43]	32749	0	0	15.67	2.58
Bbgemm [*] [45]	261765	0	0	14.76	2.91
Gemver [†] [43]	9129	903	0	10.75	4.46
GemmNcubed [†] [45]	10745	1152	0	17.80	4.51
Matrximult [†] [45]	69125	2160	0	18.81	3.54
JPGzigzag [†] [44]	54244	0	144	13.33	4.47
JPGshift [†] [44]	9162	0	256	11.79	4.23
Symm [†] [43]	131504	600	0	12.06	4.45
Syr2k [†] [43]	97932	800	0	19.83	3.58
Doitgen [†] [43]	207655	1000	0	13.27	3.62

¹Total No. LUT: 1221600 ²Total No. DSP: 2160 ³Total No. BRAM: 1292

^{*}LUT-based [†]DSP-based [‡]Hybrid

LUTs are essentially better fitted in with the linear regression model. In contrast, the DSPs implementing complex arithmetic operations tend to have non-linear power patterns, as reported in [15] and [16].

The learning curves from cross validation further reveal the gap between the decision tree model and linear regression model regarding the capability to learn from samples. We take one from each category and show the results in Fig. 10. Regarding the linear regression, the learning curves experience a *high-bias* scenario: the error is excessively high and the linear model is unable to improve inference accuracy given more training samples. In principle, the high-bias situation means the *underfitting* problem of the training data has occurred. This also accounts for the deterioration in training accuracy as non-linear power patterns continuously appear with more samples. Comparatively speaking, the decision tree model exhibits a superior ability to learn from a larger training set, due to a higher model complexity and non-linearity.

We also study the effect of operating frequency on estimation accuracy. We use the pre-trained power estimators to verify the model’s adaptability as operating frequencies change. We employ the same number of estimation cycles as the prior experiment, but run our CAD flow to collect new activity and power traces under multiple frequencies and sampling periods solely for testing purposes. Noting that the model will definitely produce biased estimation under a new frequency, we calibrate the prediction by multiplying it by the ratio of the current frequency to the 100 MHz baseline frequency ($\frac{f_{current}}{f_{baseline}}$). This is based on the fact that dynamic power consumption is proportional to the operating frequency, as shown in Equation (1). The degradation of error is within 0.28% when the operating frequency varies. In all, the decision

tree model offers high adaptability under a wide range of frequencies. This conclusion also holds when the decision trees are used in base learners for the ensemble model.

B. Decision Tree Model: Overhead

We analyze the overheads of integrating the decision tree power model into each benchmark from the following three aspects: resource utilization, operating frequency and power dissipation. The decision tree hyperparameter settings and the overheads for all three types of benchmarks are presented in Table III.

The monitoring circuits consume less than 0.05% of LUTs, 0.2% of block random-access memories (BRAMs) and 0.4% of DSPs, as shown in Table III. The decision tree exhibits extremely high area efficiency because it mainly leverages integer comparisons, while the linear model consumes a large number of floating-point additions and multiplications. The power dissipation of our single decision tree model is extremely low, namely, less than 11 mW, and can be neglected since an application generally has a power consumption in the order of watts. The maximum degradation of operating frequency is 3.89 MHz (2.91%). There is an interesting phenomenon that in some cases, the frequencies can even be improved slightly after the monitoring hardware is added. This is attributed to the fact that the placement and routing problem is NP-complete [46] and the associated algorithms are pseudo-random in nature, which introduces uncertainties that may slightly improve or degrade the timing slack of the design for some corner cases, without a guarantee of global optimum. This phenomenon also appears in the ensemble model. In conclusion, the decision-tree-based monitoring hardware demonstrates low overheads in resource utilization, operating frequency and power dissipation. The monitoring hardware of a single decision tree can be efficiently integrated into register-transfer level (RTL) designs for on-chip power monitoring.

C. Ensemble Model: Feature Number

To analyze the ensemble model, we first study the estimation accuracy when applying different numbers of features. The results for Atax are taken as an example, as shown in Fig. 11, and other benchmarks exhibit a similar tendency. We observe that the maximum number of identifiable clusters from the k-means algorithm is strongly dependent on the number of features used, besides the predetermined cluster number K . The number of clusters from k-means can be fewer than K when the feature number is small, which violates Equation (6), owing to the inadequate information provided by the limited subset of features. With more activity patterns discovered

TABLE III
DECISION TREE HYPERPARAMETER SETTINGS & MODEL OVERHEADS.

Benchmark	Max depth	Min split sample	Min leaf sample	Min leaf impurity	No. Feature	Resource (in number)				Frequency (MHz)		Power (mW)
						LUT	DSP	FF	BRAM	baseline	dtree added	
Atax	5	20	20	0.01	7	127	7	198	2.5	117.65	0	3
Bicg	5	5	5	0.001	6	125	6	176	0.5	113.30	-0.01	2
Bbgemm	6	5	20	0.001	17	187	17	419	2	105.62	-0.38	9
Gemver	6	20	5	0.05	12	152	12	308	2	100.36	+3.31	11
Gemmncubed	5	20	20	0.03	9	149	9	242	2.5	134.17	+0.02	4
Matrixmult	4	5	5	0.03	4	108	0	325	0.5	133.74	-3.89	4
JPGizzigzag	7	5	5	0.001	10	154	10	265	1	102.60	-1.32	8
JPGshift	7	5	20	0.001	11	152	11	286	1	125.08	+0.17	11
Symm	6	5	5	0.02	17	317	0	1077	2	125.47	+0.38	11
Syr2k	6	5	5	0.05	20	308	0	1019	2	125.02	+0.17	9
Doitgen	7	5	10	0.01	17	222	0	445	2	120.48	+0.55	11

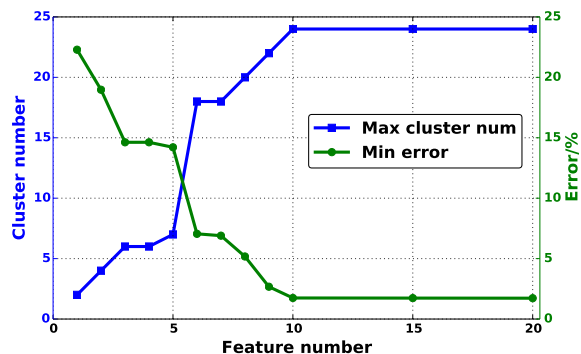


Fig. 11. Maximum identifiable clusters and minimum estimation error with different feature numbers.

through the deployment of more features, the number of maximum identifiable clusters gets larger. The minimum attainable error also reduces with the growth in identifiable clusters. Gradually, when the number of identifiable clusters becomes equal to the number of FSM states, further increase in the number of features only offers a trivial refinement in error. In light of this phenomenon, we select the number of features to be exactly enough to maximize the number of identifiable clusters. Experimental results reveal that ten features are enough for Atax, Bicg and Matrixmult, 30 features for Bbgemm and 20 features for the others. We keep these settings in the following experiments.

D. Ensemble Model: Accuracy

We study the estimation accuracy of the ensemble model using different numbers of base learners and evaluate the actual gains as the number of base learners increases. The results for different benchmarks are shown in Fig. 12. It can be discovered from the curves that, in most cases, the estimation accuracy tends to increase as more base learners are deployed. This can be deduced from Equation (20), where the error of the ensemble model, E_{ens} , is inversely proportional to the number of base learners, K . There are, however, some outliers deviating from the theoretical proof; that is, in some cases, the error goes higher when using a slightly larger number of base learners. This can be explained by the fact that the k-means algorithm cannot always guarantee that increasing the number of clusters monotonously enhances the quality of clustering. In practice, when there are some minor changes in K , the NP-hard k-means problem may fall into a local minimum that probably harms the clustering quality [22], [47]. Nevertheless,

TABLE IV
ACCURACY OF ENSEMBLE ESTIMATORS AT TRADEOFF POINTS AND OPTIMAL POINTS, COMPARED WITH THE SINGLE AVERAGING MODEL.

Benchmark	No. of base learners		MAE/%		
	tradeoff pt.	opt.	1-avg.	tradeoff pt.	opt.
Atax	21	24	3.82	1.76	1.74
Bicg	18	21	4.16	1.90	1.87
Bbgemm	22	32	2.72	0.91	0.88
Gemver	40	45	4.29	1.50	1.47
Gemmncubed	40	43	2.63	0.71	0.69
Matrixmult	53	60	3.74	0.30	0.18
JPGizzigzag	43	60	4.75	1.58	1.57
JPGshift	21	25	4.36	1.37	1.29
Symm	27	35	4.05	1.25	1.13
Syr2k	31	34	3.48	1.73	1.71
Doitgen	49	64	2.98	0.27	0.25

the existence of outliers does not alter the general trend that the error diminishes with a growing number of base learners.

In all, the accuracy of the ensemble model first decreases and then progressively stabilizes with an increase in base learners. We describe a point to optimized accuracy as the *optimal point* (denoted as opt.) in every curve of Fig. 12. The lowest error mostly appears when the cluster numbers are maximized. Through analyzing the ensemble results, we observe a *tradeoff point* (denoted as tradeoff pt.) for a benchmark as a point in a curve offering prediction accuracy close to the optimal point, which is usually selected as the first point at which the curve converges. We elaborate the results for the accuracy at the tradeoff points and optimal points for all benchmarks in Table IV. We also show the estimation accuracy of the single averaging estimators (denoted as 1-avg.) as the cases of solely using one base learner.

The average error of the single averaging model is 3.73%. In comparison to an average error of 3.86% for a single decision-tree-based model, as shown in Section VII-A, the single averaging model only exhibits an insignificant refinement. This is because there exists a high degree of correlation among error of different predictions from the same model, which counteracts the benefits of the averaging effect by introducing high correlation error, $\sum \mathbb{E}_{\mathbf{x}}[\epsilon_i(\mathbf{x}) \cdot \epsilon_j(\mathbf{x})]$. Conversely, the average error for the ensemble estimators at the tradeoff points and at the optimal points across different benchmarks is 1.21% and 1.16%, respectively. It is also noteworthy that the estimators built at the tradeoff points show modest accuracy degradation in comparison to the optimal ensemble estimators, while the reduction in the number of base learners is 17% on

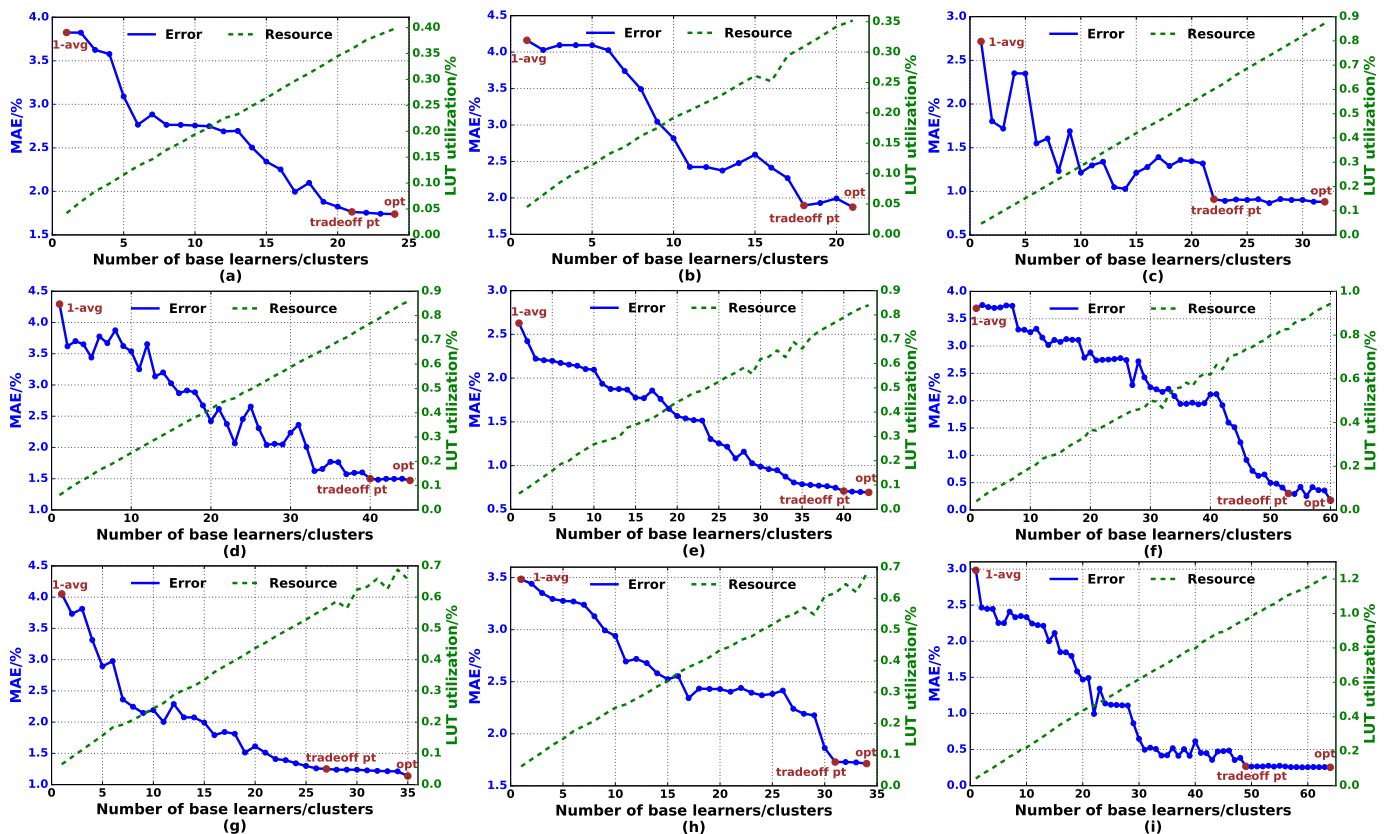


Fig. 12. Accuracy & LUT utilization curves of ensemble estimators deploying different numbers of base learners: (a) Atax; (b) Bicg; (c) Bbgemm; (d) Gemver; (e) Gemmncubed; (f) Matrixmult; (g) Symm; (h) Syr2k; and (i) Doitgen.

average. To summarize, the improvements of our proposed ensemble model at tradeoff points and optimal points over the single decision-tree-based model are $1.36\text{--}13.41\times$ and $1.38\text{--}19.67\times$, respectively.

E. Ensemble Model: Overhead

We study the tradeoff between the accuracy and overheads of the integrated ensemble monitoring hardware, with the results shown in Table V. The benchmarks with ensemble hardware showing DSP utilizations are equipped with DSP-based activity counters, while the others use LUT-based counters. Concretely speaking, we consider three factors after the integration of the ensemble hardware: i) resource usage of each base learner (bl.) along with other global control units (ctrl.), ii) the effect on operating frequencies (“+” means increase and “-” means decrease) regarding each tradeoff point and optimal point, and iii) power consumption of the ensemble estimators at every one of the tradeoff points and optimal points.

Regarding the resource overhead, we discover that the LUT and RAM utilization is almost linearly related to the number of base learners, while it is relatively independent of the size of application, as depicted by the green curves in Fig. 12. We can observe a tradeoff between resource overhead and accuracy in Fig. 12. We compute the average resource utilization per base learner from multiple ensemble estimators with varying numbers of base learners, and deduce the resource overhead of other peripheral controllers, including three preprocessing units and an average unit, as reported in Table V. The resource utilization demonstrates that each base learner can be

implemented efficiently with no more than 325 LUTs and four BRAMs, which are close to the single decision tree utilization. The maximum LUT overhead of the ensemble circuit with up to 64 base learners (Doitgen) is 1.22% of the target device. We can come to the conclusion that the proposed ensemble monitoring hardware is area-efficient and scalable. It is also worthwhile to note that [14] and [19] induced an extra load of 5% of the CPU time, while our proposed monitoring scheme does not require the employment of a processor.

The operating frequencies for the different benchmarks before and after integrating the ensemble estimators are also shown in Table V, using the original benchmarks without the monitoring hardware as the baselines. The frequencies are sometimes improved slightly after the instrument of the ensemble circuits. In all, the integration of the ensemble hardware only exerts a modest impact on the performance, with a maximum degradation of around 5 MHz (4.3%) on the operating frequencies.

With respect to the power overhead, it is reasonable to conceive that the power consumption increases with the number of base learners, which also reveals a tradeoff between the power overhead and power estimation accuracy. The power evaluation results in Table V demonstrate that the power overhead is within 425 mW for the ensemble hardware. To abate power overhead, we can diminish the number of base learners at the cost of higher error. With this objective, using the tradeoff point instead of the optimal point to build an ensemble model can reduce the power overhead by a maximum of 122 mW, with an insignificant deterioration in accuracy.

TABLE V
OVERHEADS OF THE ENSEMBLE MONITORING HARDWARE.

Benchmark	Resource (in number)				Frequency (MHz)			Power (mW)	
	LUT per bl.	RAM per bl.	LUT for ctrl.	DSP for ctrl.	baseline	tradeoff pt.	opt.	tradeoff pt.	opt.
Atax	189	2.5	329	10	117.65	+0.15	+0.42	130	135
Bicg	188	2.5	358	10	113.30	-0.50	+0.89	117	133
Bbgemm	325	4	246	30	105.62	-0.30	+0.46	199	321
Gemver	222	3	524	20	100.36	-0.21	+0.48	286	303
Gemmncubed	235	3	555	0	134.17	-5.13	-1.86	353	369
Matrixmult	187	2.5	299	0	133.74	-0.85	-2.66	310	334
JPGizigzag	194	2	168	20	102.60	-2.60	-2.59	344	425
JPGshift	280	3	440	20	125.08	+0.20	-0.08	160	190
Symm	236	3	557	0	125.47	+0.52	+0.09	133	182
Syr2k	230	3	522	0	125.02	+0.69	+0.90	166	167
Doitgen	226	3	279	0	120.48	-5.17	+0.42	244	352

VIII. CONCLUSION

Power consumption never fails to be an important design consideration for FPGAs. In this work, we establish a dynamic power monitoring scheme within the timescale of hundreds of execution cycles to facilitate emerging fine-grained power management strategies. We introduce a novel and specialized ensemble model for runtime dynamic power monitoring in FPGAs. In the model, every decomposed base learner is established using a non-overlapping data set derived from a specific cluster of FSM states with homogeneous activity patterns. To assist in this goal, we describe a generic and complete CAD flow from sample generation to feature selection, FSM state clustering, hyperparameter tuning and model ensemble. In order to put the ensemble model into practice for real-time power monitoring, we propose an efficient hardware realization, which can be embedded into the target application.

In the experiments, we first present the results of a single decision-tree-based power model, which offers a $2.41\text{--}6.07\times$ improvement in prediction accuracy in comparison to the traditional linear model. Furthermore, we study the extra gains in prediction accuracy derived from our customized ensemble model, by varying the number of base learners. We observe a tradeoff between estimation accuracy and overheads of resources, performance and power introduced by the ensemble monitoring hardware. Through analyzing the experimental results, we find a tradeoff point for accuracy and overheads, and an optimal point for accuracy, which, respectively, show a capability to achieve error $1.36\text{--}13.41\times$ and $1.38\text{--}19.67\times$ lower than the single decision-tree-based model. The hardware implementation of runtime monitoring with up to 64 base learners incurs modest overheads: 1.22% of LUT utilization, 4.3% of performance and less than 425 mW power dissipation. By using the tradeoff points instead of optimal points to construct the monitoring hardware, the number of base learners can be reduced by 17% on average, and the power consumption can be accordingly cut by up to 122 mW, at the cost of an insignificant accuracy loss.

In summary, the proposed ensemble model provides accurate dynamic power estimate within an error margin of 1.90% of a commercial gate-level power estimation tool. The hardware realization of runtime power monitoring demonstrates low overheads of resource, performance and power consumption.

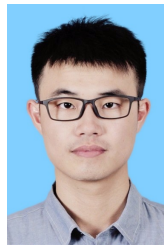
ACKNOWLEDGMENT

This work is funded by Hong Kong RGC GRF 16245116.

REFERENCES

- [1] I. Kuon and J. Rose, "Measuring the gap between FPGAs and ASICs," *IEEE Trans. Comput.-Aided Des. Integr. Circuits Syst.*, vol. 26, no. 2, pp. 203–215, Feb 2007.
- [2] T. Tuan, S. Kao, A. Rahman, S. Das, and S. Trimmerger, "A 90nm low-power FPGA for battery-powered applications," in *Proc. ACM/SIGDA 14th Int. Symp. Field Programmable Gate Arrays*, 2006, pp. 3–11.
- [3] S. Huda, M. Mallick, and J. H. Anderson, "Clock gating architectures for FPGA power reduction," in *Proc. Int. Conf. Field Programmable Logic Appl.*, Aug 2009, pp. 112–118.
- [4] Y. Lin and L. He, "Dual-Vdd interconnect with chip-level time slack allocation for FPGA power reduction," *IEEE Trans. Comput.-Aided Des. Integr. Circuits and Syst.*, vol. 25, no. 10, pp. 2023–2034, Oct 2006.
- [5] C. H. Hoo, Y. Ha, and A. Kumar, "A directional coarse-grained power gated FPGA switch box and power gating aware routing algorithm," in *Proc. 23rd Int. Conf. Field Programmable Logic Appl.*, Sept 2013, pp. 1–4.
- [6] J. L. Nunez-Yanez, M. Hosseinabady, and A. Beldachi, "Energy optimization in commercial FPGAs with voltage, frequency and logic scaling," *IEEE Trans. Comput.*, vol. 65, no. 5, pp. 1484–1493, May 2016.
- [7] A. Löscher, T. Beisel, T. Kenter, C. Plessl, and M. Platzner, "Performance-centric scheduling with task migration for a heterogeneous compute node in the data center," in *Proc. Conf. Des., Autom. & Test in Europe*, Mar 2016, pp. 912–917.
- [8] W. Kim, D. Brooks, and G.-Y. Wei, "A fully-integrated 3-level DC-DC converter for nanosecond-scale DVFS," *IEEE J. Solid-State Circuits*, vol. 47, no. 1, pp. 206–219, Jan 2012.
- [9] S. Eyerhan and L. Eeckhout, "Fine-grained DVFS using on-chip regulators," *ACM Trans. Archit. Code Optim.*, vol. 8, no. 1, pp. 1:1–1:24, Feb 2011.
- [10] P. Mantovani, E. G. Cota, K. Tien, C. Pilato, G. Di Guglielmo, K. Shepard, and L. P. Carloni, "An FPGA-based infrastructure for fine-grained DVFS analysis in high-performance embedded systems," in *Proc. 53rd ACM/EDAC/IEEE Des. Autom. Conf.*, June 2016, pp. 1–6.
- [11] B. Keller, "Opportunities for fine-grained adaptive voltage scaling to improve system-level energy efficiency," Master's thesis, EECS Department, UC Berkeley, Dec 2015.
- [12] H. Asghari-Moghaddam, H. R. Ghasemi, A. A. Sinkar, I. Paul, and N. S. Kim, "VR-scale: Runtime dynamic phase scaling of processor voltage regulators for improving power efficiency," in *Proc. 53rd Annu. Des. Automation Conf.*, 2016, pp. 151:1–151:6.
- [13] A. Lakshminarayana, S. Ahuja, and S. Shukla, "High level power estimation models for FPGAs," in *IEEE Comput. Society Annu. Symp. VLSI*, July 2011, pp. 7–12.
- [14] M. Najem, P. Benoit, F. Bruguier, G. Sassatelli, and L. Torres, "Method for dynamic power monitoring on FPGAs," in *24th Int. Conf. Field Programmable Logic Appl.*, Sept 2014, pp. 1–6.
- [15] A. Bogliolo, L. Benini, and G. De Micheli, "Regression-based RTL power modeling," *ACM Trans. Des. Autom. Electron. Syst.*, vol. 5, no. 3, pp. 337–372, Jul. 2000.
- [16] D. Lee, T. Kim, K. Han, Y. Hoskote, L. K. John, and A. Gerstlauer, "Learning-based power modeling of system-level black-box IPs," in *Proc. IEEE/ACM Int. Conf. Comput.-Aided Des.*, 2015, pp. 847–853.

- [17] F. Li, D. Chen, L. He, and J. Cong, "Architecture evaluation for power-efficient FPGAs," in *Proc. ACM/SIGDA 11th Int. Symp. Field Programmable Gate Arrays*, 2003, pp. 175–184.
- [18] C. Najoua, B. Mohamed, and B. M. Hedi, "Accurate dynamic power model for FPGA based implementations," *IJCSI Int. J. Comput. Science*, vol. 9, no. 2, pp. 84–89, 2012.
- [19] J. Davis, E. Hung, J. Levine, E. Stott, P. Cheung, and G. Constantinides, "KAPow: High-accuracy, low-overhead online per-module power estimation for FPGA designs," *ACM Trans. on Reconfigurable Technol. Syst.*, vol. 11, 2018.
- [20] Z. Lin, W. Zhang, and S. Sinha, "Decision tree based hardware power monitoring for run time dynamic power management in FPGA," in *Proc. 27th Int. Conf. Field Programmable Logic and Appl.*, Sept 2017, pp. 1–8.
- [21] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, Oct 2001.
- [22] S. Lloyd, "Least squares quantization in PCM," *IEEE Trans Inf. Theory*, vol. 28, no. 2, pp. 129–137, March 1982.
- [23] D. Arthur and S. Vassilvitskii, "K-means++: The advantages of careful seeding," in *Proc. 18th Annu. ACM-SIAM Symp. Discrete Algorithms*, 2007, pp. 1027–1035.
- [24] A. Canis, J. Choi, M. Aldham, V. Zhang, A. Kammoona, J. H. Anderson, S. Brown, and T. Czajkowski, "LegUp: High-level synthesis for FPGA-based processor/accelerator systems," in *Proc. 19th ACM/SIGDA Int. Symp. Field Programmable Gate Arrays*, 2011, pp. 33–36.
- [25] Xilinx Ltd, "Vivado design suite user guide: High level synthesis," *Xilinx White Paper*, pp. 1–672, April 2017.
- [26] P. Jamieson, K. B. Kent, F. Gharibian, and L. Shannon, "Odin II — an open-source verilog HDL synthesis tool for CAD research," in *18th IEEE Annu. Int. Symp. Field-Programmable Custom Comput. Mach.*, May 2010, pp. 149–156.
- [27] J. Rose, J. Luu, C. W. Yu, O. Densmore, J. Goeders, A. Somerville, K. B. Kent, P. Jamieson, and J. Anderson, "The VTR project: Architecture and CAD for FPGAs from Verilog to routing," in *Proc. ACM/SIGDA Int. Symp. Field Programmable Gate Arrays*, 2012, pp. 77–86.
- [28] J. Jones, "Abstract syntax tree implementation idioms," in *Proc. 10th Conf. Pattern Languages of Programs*, September 2003, p. 26.
- [29] G. Friedman, A. Hartman, K. Nagin, and T. Shiran, "Projected State Machine Coverage for Software Testing," in *Proc. of ACM SIGSOFT Int. Symp. Software Testing and Analysis*, 2002, pp. 134–143.
- [30] M. Klein and S. Kolluri, "Leveraging power leadership at 28 nm with Xilinx 7 series FPGAs," *Xilinx White Paper: 7 Series FPGAs*, pp. 1–16, Jan. 2015.
- [31] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene selection for cancer classification using support vector machines," *Mach. Learning*, vol. 46, no. 1, pp. 389–422, Jan 2002.
- [32] H. Ravishankar, R. Madhavan, R. Mullick, T. Shetty, L. Marinelli, and S. E. Joel, "Recursive feature elimination for biomarker discovery in resting-state functional connectivity," in *38th Annu. Int. Conf. IEEE Engineering in Medicine and Biology Society*, Aug 2016, pp. 4071–4074.
- [33] L. Rutkowski, M. Jaworski, L. Pietruczuk, and P. Duda, "The CART decision tree for mining data streams," *Information Sciences*, vol. 266, pp. 1–15, 2014.
- [34] J. D. Rodriguez, A. Perez, and J. A. Lozano, "Sensitivity analysis of k-fold cross validation in prediction error estimation," *IEEE Trans. Pattern Analysis and Mach. Intelligence*, vol. 32, no. 3, pp. 569–575, March 2010.
- [35] C. M. Bishop, *Pattern recognition and machine learning*. Springer, 2006.
- [36] G. Fumera and F. Roli, "A theoretical and experimental analysis of linear combiners for multiple classifier systems," *IEEE Trans. Pattern Analysis and Mach. Intelligence*, vol. 27, no. 6, pp. 942–956, June 2005.
- [37] G. Brown, J. Wyatt, R. Harris, and X. Yao, "Diversity creation methods: A survey and categorisation," *Inf. Fusion*, vol. 6, no. 1, pp. 5–20, 2005.
- [38] Y. Liu, Q. Zhao, and Y. Pei, "Ensemble learning with correlation-based penalty," in *IEEE 12th Int. Conf. Dependable, Autonomic and Secure Comput.*, Aug 2014, pp. 350–353.
- [39] Y. R. Qu and V. K. Prasanna, "Scalable and dynamically updatable lookup engine for decision-trees on FPGA," in *Proc. High Performance Extreme Comput. Conf.*, Sept 2014, pp. 1–6.
- [40] F. Saqib, A. Dutta, J. Plusquellic, P. Ortiz, and M. S. Pattichis, "Pipelined decision tree classification accelerator implementation in FPGA (DT-CAIF)," *IEEE Trans. Comput.*, vol. 64, no. 1, pp. 280–285, Jan 2015.
- [41] Xilinx Ltd, "Vivado design suite user guide: Power analysis and optimization," *Xilinx White Paper*, pp. 1–116, June 2017.
- [42] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg *et al.*, "Scikit-learn: Machine learning in Python," *J. Mach. Learning Research*, vol. 12, no. Oct, pp. 2825–2830, 2011.
- [43] L.-N. Pouchet. (2012) Polybench: The polyhedral benchmark suite. [Online]. Available: <http://www.cs.ucla.edu/pouchet/software/polybench>
- [44] Y. Hara, H. Tomiyama, S. Honda, H. Takada, and K. Ishii, "CHStone: A benchmark program suite for practical C-based high-level synthesis," in *IEEE Int. Symp. Circuits and Syst.*, May 2008, pp. 1192–1195.
- [45] B. Reagen, R. Adolf, Y. S. Shao, G.-Y. Wei, and D. Brooks, "MachSuite: Benchmarks for accelerator design and customized architectures," in *IEEE Int. Symp. Workload Characterization*, Oct 2014, pp. 110–119.
- [46] Y.-L. Wu, S. Tsukiyama, and M. Marek-Sadowska, "Graph based analysis of 2-D FPGA routing," *IEEE Trans. Comput.-Aided Des. Integr. Circuits and Syst.*, vol. 15, no. 1, pp. 33–44, Jan 1996.
- [47] R. Ostrovsky, Y. Rabani, L. J. Schulman, and C. Swamy, "The effectiveness of lloyd-type methods for the k-means problem," in *47th Annu. IEEE Symp. Foundations of Comput. Science*, Oct 2006, pp. 165–176.



Zhe Lin (S'15) received his B.S. degree from School of Electronic Science and Engineering from Southeast University, Nanjing, China, in 2014. Since 2014, he has been a Ph.D. Student in the Department of Electronic and Computer Engineering at Hong Kong University of Science and Technology (HKUST), Hong Kong. Zhe's current research interests cover FPGA-based heterogeneous multicore systems and power management strategies of modern FPGAs.



Sharad Sinha (S'03, M'14) is an assistant professor with Dept. of Computer Science and Engineering, Indian Institute of Technology (IIT) Goa. Previously, he was a Research Scientist at NTU, Singapore. He received his PhD degree in Computer Engineering from NTU, Singapore (2014). He received the *Best Speaker Award from IEEE CASS Society*, Singapore Chapter, in 2013 for his PhD work on High Level Synthesis and serves as an Associate Editor for *IEEE Potentials* and *ACM Ubiquity*. Dr. Sinha earned a Bachelor of Technology (B.Tech) degree

in Electronics and Communication Engineering from Cochin University of Science and Technology (CUSAT), India in 2007. From 2007-2009, he was a design engineer with Processor Systems (India) Pvt. Ltd. Dr. Sinha's research and teaching interests are in computer architecture, embedded systems and reconfigurable computing.



Wei Zhang (M'05) received a Ph.D. degree from Princeton University, Princeton, NJ, USA, in 2009. She was an assistant professor with the School of Computer Engineering, Nanyang Technological University, Singapore, from 2010 to 2013. Dr. Zhang joined the Hong Kong University of Science and Technology, Hong Kong, in 2013, where she is currently an associated professor and she established the reconfigurable computing system laboratory (RCSL).

Dr. Zhang has authored or co-authored over 80 book chapters and papers in peer reviewed journals and international conferences. Dr. Zhang serves as the Associate Editor for *ACM Transactions on Embedded Computing Systems (TECS)*, *IEEE Transactions on Very Large Scale Integration Systems (TVLSI)*, and *ACM Journal on Emerging Technologies in Computing Systems (JETC)*. She also serves on many organization committees and technical program committees. Dr. Zhangs current research interests include reconfigurable systems, FPGA-based design, low-power high-performance multicore systems, electronic design automation, embedded systems, and emerging technologies.