

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

Fast Accurate Full-Chip Dynamic Thermal Simulation with Fine Resolution Enabled by a Learning Method

Lin Jiang, Yu Liu, *member, IEEE*, and Ming-C. Cheng, *Senior Member, IEEE*

Abstract—The need for full-chip dynamic thermal simulation for effective run-time thermal management of multi-core processors has been growing in recent years due to the rising demand for high performance computing. In addition to simulation efficiency and accuracy, a high resolution is desirable in order to accurately predict crucial hot spots in the chip. This work investigates a simulation technique derived from proper orthogonal decomposition (POD) for full-chip dynamic thermal simulation of a multi-core processor. The POD projects a heat transfer problem onto a mathematical space constituted by a finite set of basis functions (or POD modes) that are generated (or *trained*) by thermal solution data collected from direct numerical simulation (DNS). Accuracy and efficiency of the POD simulation technique influenced by quality of thermal data are examined thoroughly, especially in the areas with high thermal gradients. The results show that, if the POD modes are trained by good-quality data, the POD simulation offers an accurate prediction of the dynamic thermal distribution in the multi-core processor with an extremely small degree of freedom (DoF). A reduction in computational time over four orders of magnitude, compared to the DNS, can be achieved for full-chip dynamic thermal simulation with a resolution as fine as the DNS. The study has also demonstrated that the POD approach can be used to rigorously verify the accuracy of solutions offered by DNS tools. A practical approach is proposed to further enhance the accuracy and efficiency of the proposed full-chip thermal simulation technique.

Index Terms—data driven, full-chip thermal simulation, model order reduction, proper orthogonal decomposition, multi-core processors.

I. INTRODUCTION

For more than half a century, integrated circuits (ICs) and computing performance have been improved by shrinking the feature size of semiconductor devices. With more transistors integrated in a chip, ICs have become more functional and complex. However, near the end of Moore's law, it is difficult to continue improving chip performance via device miniaturization due to physical limitations and large heat dissipation caused by higher device density [1]. In addition, due to the long interconnect latency in the planar structure, 3D IC technology [2]-[4] has been

introduced to further increase the device density and enhance the chip performance and functionality. Although this offers significant improvement on the IC performance with a lower cost, temperature escalation and excessive hot-spot formation have become a more severe issue with such a high degree of integration [5]-[7].

To continue improving computing performance, many-core CPUs on a chip [8]-[12] have been introduced to facilitate parallel computing. To satisfy the demand for the recent growth of cloud computing and big data applications, general purpose GPUs (GPGPUs) with massively parallel processing power enabled by hundreds or thousands of cores have been widely used in scientific computing, social network, movie streaming, online shopping, etc. in computer servers and data centers around the globe [13]-[16]. As the processors are becoming larger to handle the massive amount of data, the heating issues have been enhanced, leading to more serious high temperature and hot-spot formation. This not only impairs CPU/GPGPU performance and wastes computing energy but also reduces their lifetime caused by thermal stress and electromigration [17]-[22].

Due to the severe heating issues in 3D ICs, CPUs and GPUs, effective thermal management is thus desperately needed to reduce temperature, suppress hot spots, improve performance and reliability and save energy. This can be achieved more effectively, e.g., via thermal-aware task scheduling [23]-[29] that requires efficient and accurate thermal simulations at the chip level with a reasonable resolution for 3D ICs, CPUs and GPGPUs. Several thermal simulation methods have been developed at different levels of efficiency and accuracy. The rigorous approaches that provide an accurate thermal profile with a high resolution, are the direct numerical simulations (DNSs), based on the finite difference (FD), finite element (FE) or finite volume (FV) method. There are many DNS commercial and open-source tools available for such an application, including ANSYS [30], COMSOL [31], FEniCS [32], FREEFEM [33], etc. These however demand extensive computational resources due to a very large degree of freedom

Manuscript received February 19, 2022. This work was supported by the National Science Foundation under Grant Nos. ECCS-2003307 and OAC-2118079.

Corresponding author: M.C. Cheng. All authors are with Department of Electrical and Computer Engineering, Clarkson University, Potsdam, NY 13699-5720 USA (email: jiangl2@clarkson.edu, yuliu@clarkson.edu, mcheng@clarkson.edu).

Color versions of one or more of the figures in this article are available online at <http://ieeexplore.ieee.org>

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

(DoF) needed in the simulation and are thus prohibitive for chip-level thermal simulation of ICs or CPUs/GPGPUs.

In recent years, the Green's function approach has gained its popularity for the chip-level thermal simulation [34]-[35] due to its simplicity and computational efficiency compared to the DNSs. The Green's function is a spatial impulse response of a system and usually calculated in response to a unit point heat source at the center of a large chip. The thermal solution is then solved by convolution of the impulse response with the power distribution in the chip. The conventional Green's function is thus difficult to include boundary conditions (BCs) of a finite domain [34], [35], especially when heat sources are close to the chip corners or edges [36], [37]. It is also difficult to apply the approach to transient thermal simulation. In addition, the Green's function approach is applied to a single thin layer where the power sources are generated [34]-[37] and thus only offers a 2D thermal profile in the heating layer. An approach was investigated to include the 3D temperature profile in a chip using multilayer Green's functions [38]; it is however limited to steady state simulation. Different techniques at the expense of the computational cost with some assumptions have been developed to correct the corner and edge effects [39] for transient thermal simulation [39], [40].

When fast thermal simulations are needed for large-scale chips, lumped RC thermal circuits are usually used [41]-[48]. For example, HotSpot [41]-[43], the most frequently used RC thermal simulator, has been widely applied to chip-level thermal simulations [41]-[48]. However, accuracy of the block model in HotSpot (hereafter named HotSpot-Block) suffers from the large RC elements, especially with high aspect ratios (ARs). The large elements not only overlook some hot spots, whose physical sizes are considerably smaller than the element, but also offer incorrect heat flux between elements estimated by their node (average element) temperatures. To compensate the dynamic distributed heat transfer incorrectly calculated in the RC thermal circuit, a scaling factor (SF) less than one on all thermal capacitances is used [42], which however still leads to a 200% error in HotSpot-Block thermal simulations of some floorplans, compared with FEM analysis [49]. To address this issue, the Grid model of HotSpot (hereafter named HotSpot-Grid) was developed [41] by allowing small RC elements to improve the accuracy, where an SF value with a default value less than one is still included to adjust the accuracy. With the improved Grid model, large deviations from ANSYS results have still been observed [40], [50], perhaps caused by the RC -circuit approximation and/or the inappropriate SF value. The SF value should be dependent on the size and AR of each individual lumped element and how fast the dynamic power sources vary in time. However, if very small elements are implemented in HotSpot-Grid, it is equivalent to an FD method and the SF should be one. In this case, it is similar to DNSs that demand intensive computational time.

For effective thermal management of large-scale chips, the major challenge is to predict high thermal gradients and hot spots as efficient as the RC thermal circuits with an accuracy close to DNS. To capture small-size hot spots, as has been

studied recently on millimeter- or sub-millimeter-scale hot spots in DNSs of multi-core and 3D IC chips [51]-[53], a high resolution is needed. In this study, we investigate a technique for full-chip dynamic thermal simulation that is able to meet all the aforementioned challenges, including the efficiency, accuracy and high resolution. The technique is derived from a projection-based data-driven algorithm, proper orthogonal decomposition (POD) [54], [55], that has been applied to many areas of research [56]-[62]. The early concept was briefly presented at [63]. By projecting a physical domain onto a functional space described by a finite set of basis functions (or POD modes), the POD simulation technique is able to achieve desired efficiency and accuracy with a very small DoF if the modes are trained by good-quality data.

In this work, a quad-core CPU, AMD ATHLON II X4 610e [64], [65], is selected to demonstrate the POD full-chip thermal simulation technique. To develop the POD simulation method, thermal solution data of the processor are needed for the POD mode generation/training. The data are collected from two DNS tools, including a rigorous FEM implemented in FEniCS [32] and the popular thermal simulator HotSpot-Grid [41], [43] with very small RC elements and $SF = 1$. The POD models built upon FEniCS-FEM and HotSpot-Grid are named FEniCS-POD and HotSpot-POD, respectively. The effectiveness of the POD thermal simulation method is investigated in terms of the DoF in POD simulations and quality of thermal data collected from each of these two DNS approaches. The accuracy of the DNS tools is also examined by the POD method.

II. POD FUNDAMENTALS

Unlike many projection-based methods with pre-selected basis functions, such as Fourier transform, Legendre polynomials, Bessel functions, etc., the POD modes are *learned* from the solution data of a domain Ω , where each POD mode is optimized by maximizing its mean square inner product of the thermal solution with the mode [54], [55]. This process maximizing the projection onto each mode leads to an eigenvalue problem,

$$\langle \vec{\phi}_i, \vec{\phi}_j \rangle = \lambda_i \delta_{ij} \quad (1)$$

where λ_i is the eigenvalue representing the mean squared temperature captured by its eigenfunction $\vec{\phi}_i$, and the brackets $\langle \rangle$ denote an averaging process over the temporal sampled data subjected to dynamic variations of power sources and BCs. For steady problems, this process averages the sampled data over different power levels and BCs. Instead of solving (1) directly, the method of snapshots [57], [66], [67] is adopted to solve and in (1) more efficiently, where the dimension of the problem is reduced to the number of samples (or snapshots), N_s . With the generated POD modes, temperature \vec{T} , can be represented by a linear combination of the modes,

$$\vec{T} = \sum_{i=1}^N \vec{\phi}_i \vec{a}_i$$

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

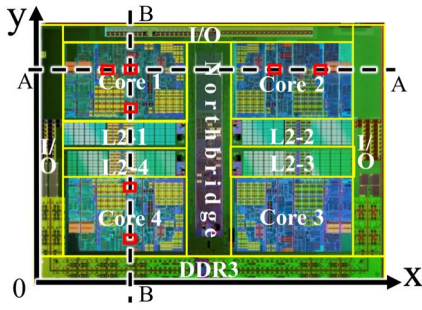


Fig. 1. Floorplan of AMD ATHLON II X4 610e with dimensions of 14mm×12mm×0.65mm in x, y and z. Paths A and B indicate the temperature plotting paths, and the red rectangles (0.4375mm×0.375mm) specify the locations of the applied high power density in later demonstrations.

where α_i are the weighting coefficients and M is the number of modes ($M \leq N_s$) or the DoF selected to represent the solution.

To derive a POD model, the heat transfer equation is projected onto a functional space using the Galerkin projection,

$$\begin{aligned} & \int_{\Omega} \rho C_p \frac{\partial T}{\partial t} \phi_i d\Omega + \int_{\Omega} \nabla \cdot (\kappa \nabla T) \phi_i d\Omega - \int_{\partial\Omega} q \phi_i d\Gamma = \int_{\Omega} P \phi_i d\Omega \\ & = \frac{d}{dt} \int_{\Omega} \rho C_p T \phi_i d\Omega + \int_{\Omega} \kappa \nabla T \cdot \nabla \phi_i d\Omega - \int_{\partial\Omega} q \phi_i d\Gamma = \int_{\Omega} P \phi_i d\Omega \end{aligned}$$

where κ is the thermal conductivity, ρ the density, C_p the specific heat, P the power density and \mathbf{n} the outward differential surface element vector on the surface. Using (2), (3) leads to an M -dimensional ordinary differential equations (ODEs) for \mathbf{a} ,

$$\mathbf{M} \frac{d\mathbf{a}}{dt} + \mathbf{K} \mathbf{a} = \mathbf{F}, \quad i = 1 \text{ to } M, \quad (4)$$

with \mathbf{M} as the element of the POD thermal capacitance matrix,

$$M_{ij} = \int_{\Omega} \rho C_p \phi_i \phi_j d\Omega, \quad (5)$$

\mathbf{K} as the elements of the POD thermal conductance matrix and \mathbf{F}

the i th-mode POD power vector. In this study, adiabatic and convective BCs are applied on the chip surfaces. For an adiabatic BC, the heat flux on the surface in (3) vanishes, and

$$\mathbf{K}_{ad} = \int_{\Omega} \kappa \nabla \phi_i \cdot \nabla \phi_j d\Omega, \quad i, j = 1 \text{ to } M, \quad (6)$$

As to the convective BC, the heat flux normal to the boundary surface is given by

$$q = -\kappa \nabla T \cdot \mathbf{n} = h(T - T_{\infty}), \quad (7)$$

where T_{∞} is the ambient temperature and h is the heat transfer coefficient that is a function of the airflow rate near the boundary. Using (2) for T in (3) and (7),

$$\mathbf{K}_{cv} = \int_{\partial\Omega} h \phi_i \phi_j d\Gamma, \quad i, j = 1 \text{ to } M, \quad (8)$$

and \mathbf{F} is given by

$$\mathbf{F}_i = \int_{\Omega} P \phi_i d\Omega = \int_{\Omega} \sum_{j=1}^M \mathbf{a}_j \phi_j \phi_i d\Omega = \sum_{j=1}^M \mathbf{a}_j \int_{\Omega} \phi_j \phi_i d\Omega = \mathbf{M} \mathbf{a}, \quad (9)$$

With the predefined shape of power density, the above integrals

Because $\bar{\theta}_i$ represents the mean squared temperature captured by the i th mode over the data in generation of the POD modes, the theoretical least square (LS) error over the entire simulation domain for the M -mode POD model is given by

$$C_{\text{EFE}} = \frac{1}{J} \sum_{i=1}^M \frac{H_i}{K} = 10$$

Numerically, the LS error with respect to the DNS is given as

$$\frac{H_i}{H_i}$$

can be pre-evaluated and saved in a technology database for thermal simulation of the chip.

$$C_{NO} = G \begin{pmatrix} C^P, \\ \Omega K Q, \\ + R^P \Omega, 11 \end{pmatrix} \quad ($$

where C^* is the temperature difference between the DNS and POD model at the i th time step. The ideal error in (10) is valid only if high data quality is guaranteed. Due to the numerical approximation and computer precision, C_{NO} is usually larger than C_{EF} .

In summary, two projections are performed to arrive at this rigorous methodology. The first one maximizes the projection of thermal data onto the modes, which leads to (1), such that the trained modes contain essential information on variations of heat excitations and BCs embedded in the data. Thus, the POD modes are able to represent the thermal solution in (2), using a very small number of modes, if the weighting coefficients a_j are evaluated properly. To do so, the second (Galerkin) projection is applied in (3) to offer a clear guideline for the POD modes to comply with physical principles imposed by the dynamic heat transfer equation, which results in the ODEs for a_j given in (4). The coefficients of (4) and the sources on the right-hand side of (4) are then evaluated from the POD modes, the gradients of the

modes, the projection of the power onto the modes and the projection of the boundary conditions onto the modes, as given in (5)-(9), which empowers a_j solved from (4), together with the modes, to obey the physical principles, accounting for all parametric variations in the collected data via POD modes.

The POD methodology derived from the above procedure is thus sensitive to the consistency between these two projections. For example, as shown in (3), the heat transfer equation is projected along the POD modes that are trained by thermal data generated from DNS. If the data from DNS is not consistent with the heat transfer equation, the projected equation in (4) will not accurately represent the heat transfer equation. Therefore, the comparison of thermal predictions between the developed POD model and its DNS tool will provide a reliable indication of the accuracy of the DNS. This study not only demonstrates the accuracy and efficiency of the POD approach for full-chip thermal simulation but also verifies the concept of utilizing the POD method to validate the accuracy of DNS tools.

III. THERMAL DATA COLLECTION AND MODE GENERATION

The floorplan of the selected quad-core CPU, AMD ATHLON II X4 610e [64], [65], is displayed in Fig. 1 with physical locations of functional units, including 4 cores, 4 L2 caches, one northbridge, 3 I/O and one DDR3. DNS tools, FEniCS-FEM and HotSpot-Grid are used in this work to collect thermal data. In any comparison between these two DNS tools,

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

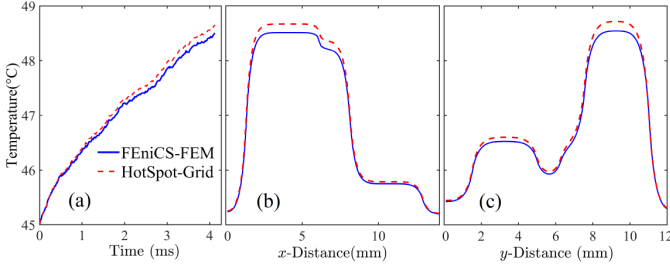


Fig. 2. Temperature estimated by HotSpot-Grid and FEniCS-FEM in the CPU. (a) Temporal evolution at the intersection of Paths A and B. Temperature distributions at $t = 4.2$ ms along (b) Path A and (c) Path B.

the dynamic power map, BCs and numerical settings in the simulations are identical. To make the comparison meaningful with HotSpot, the material properties of the selected CPU in all DNSs are adopted from HotSpot [43]; these include 100 W/(m·K), 751.1 J/(kg·K) and 2330 kg/m³ for thermal conductivity, specific heat and density, respectively. All surfaces of the chip are adiabatic except the bottom of the substrate, where a heat transfer coefficient is applied with an ambient of 45°C [25], [42], [68]–[70]. Similar to HotSpot [43], the substrate heat transfer coefficient is calculated to be 2.41 W/(cm²·K) from the thermal resistances of the heat spreader, thermal interface material layer and heat sink based on the cross section, thickness and material property of each layer with a further assumption of uniform temperature on each of these layers. Due to the mesh (element) restriction in HotSpot-Grid [43], a uniform mesh is employed. Since very small lumped *RC* elements are used in HotSpot-Grid, SF is taken to be 1. In each demonstration, two sets of POD modes are generated by thermal data collected from the DNSs, one from FEniCS-FEM and the other from HotSpot-Grid.

To collect thermal data, a uniform power source is applied to each unit on the heating layer with a thickness of 0.15mm, and a mesh of 128×128×13 in x , y and z with reasonable ARs ($\Delta x/\Delta y \approx 1.17$, and $\Delta y/\Delta z \approx 1.88$) is implemented in both DNSs of the CPU. The heating layer refers to the thin layer of devices and interconnects on top of the chip, where the power is dissipated. Similarly to [42], the dynamic power is averaged over 48,000 CPU cycles at 3.5 GHz and DNS is performed with each time step of 8,000 clock cycles and a total power near 50W. The dynamic power applied to each unit is randomly generated. To apply a more realistic power profile, the power density distribution given in Table I is adopted from the power map generated by the *hammer* and *soplex* benchmarks running in Cores 1 and 3, respectively, in [65]. The percentage of the total chip power consumed by each unit and the area of each unit are also listed in Table I.

TABLE I
POWER PERCENTAGE AND DENSITY (10⁸W/m³) [65]

Unit	Core 1	Core 2	Core 3	Core 4	L2-1	L2-2	L2-3	L2-4	I/O	NB	DDR3
%	24.1	5.2	17.5	10.6	3.7	2.3	1.4	2.2	3.2	26.2	3.6
Density	48.7	10.6	35.4	21.3	20.0	12.2	7.6	11.7	2.7	45.9	5.9
Area (mm ²)	16.49	16.49	16.49	16.49	6.21	6.21	6.21	6.21	38.1	19.0	20.02

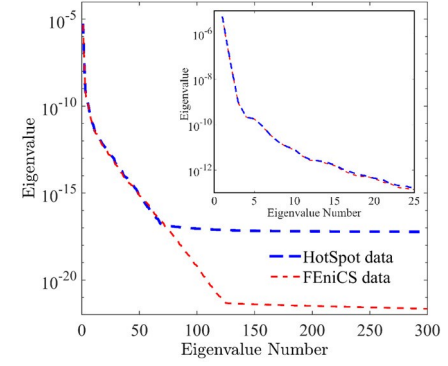


Fig. 3. Eigenvalue spectrums of the thermal data collected from FEniCS-FEM and HotSpot-Grid. A close-up look of the spectrums for the first 25 modes is included in the inset.

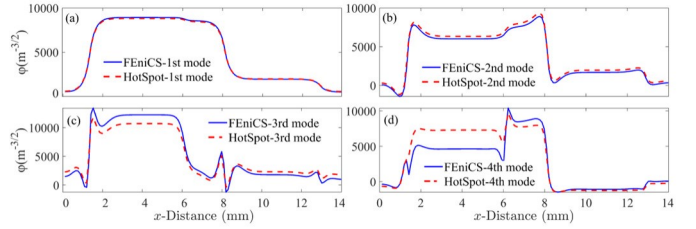


Fig. 4. POD modes along Path A for the (a) first, (b) second, (c) third and (d) fourth modes.

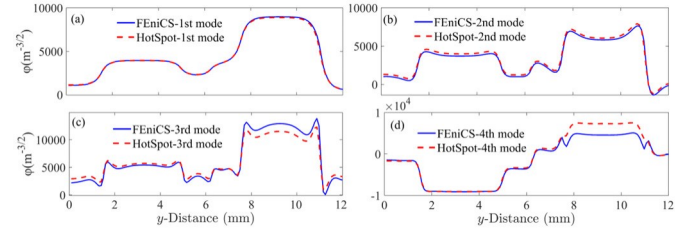


Fig. 5. POD modes along Path B for the (a) first, (b) second, (c) third and (d) fourth modes.

Temperature evolution at the intersection of Paths A and B indicated in Fig. 1 is given in the Fig. 2(a). HotSpot-Grid predicts a faster evolution and a 4%-5% higher temperature than FEniCS-FEM for $t > 2$ ms. The temperature profile at 4.2ms along Path A is illustrated in Fig. 2(b) with higher temperature in Core 1 and Northbridge but lower temperature in Core 2 due to the power density distribution shown in Table I. The profile at 4.2ms along Path B is shown in Fig. 2(c). The difference between these two approaches shown in Figs. 2(b) and 2(c) is also near 4%-5% in Core 1 along Paths A and B.

These dynamic thermal data collected from each of the FEniCS-FEM and HotSpot-Grid simulations are applied to generate their POD modes and eigenvalues from (1) using the method of snapshots [57], [66], [67]. Fig. 3 shows that the eigenvalues of the data collected from these two DNSs are nearly identical for the first 20 modes even though a difference of the DNS results near 4%-5% is observed. It is shown that the third and fourth eigenvalues drop more than three and four orders of magnitude, respectively, from the first mode. This indicates that the essential thermal information in this case is

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

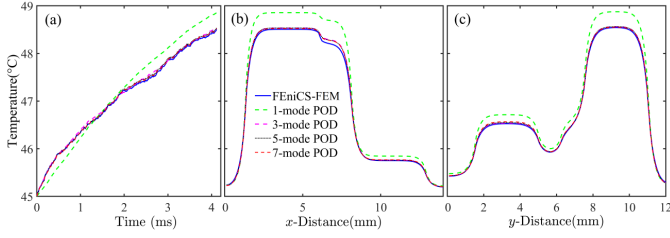


Fig. 6 Temperatures obtained from FEniCS-POD and FEniCS-FEM in Case 1. (a) Dynamic temperature at the intersection of Paths A and B, and the temperature distributions along (b) Path A (c) Path B at $t = 4.2$ ms.

accounted for in the first three or four modes. Thus, the POD model derived from these data with three or four modes should offer an accurate thermal prediction if good quality of data from the DNSs is guaranteed. The computing accuracy is limited by 16 decimal digits implemented in both DNSs, and the eigenvalue based on FEniCS-FEM's data becomes nearly invariant after dropping 16 orders of magnitude from its first mode. However, the eigenvalue from HotSpot-Grid's data decreases by only 12-13 orders from the first mode before becoming flat, which indicates its numerical inaccuracy beyond 12 digits perhaps due to the lumped-element approximation.

The first four POD modes built upon the two DNS approaches along Paths A and B are shown in Figs. 4 and 5, respectively, since the contribution to the thermal solution beyond four modes is negligible according to their eigenvalue spectrums. The first mode represents the mean of the thermal data. The small first-mode difference between FEniCS-POD and HotSpot-POD suggests that HotSpot-Grid and FEniCS-FEM capture nearly the same average dynamic thermal behavior. The difference between the thermal solutions from these DNS tools is disclosed in the second to the 4th mode. This leads to different dynamic power of the higher modes in the POD space given in (6) and (9), as well as the conductance elements in (6) and (8). The large difference in the higher modes between these two approaches is thus expected to induce an evident deviation in the prediction. It should be noted that the eigenvalue spectrum is a good indication of the number of modes needed to reach a good solution only if the data quality is good. With inadequate data quality, the generated POD modes do not represent the heat transfer equation accurately in the POD space.

IV. DEMONSTRATION

Dynamic thermal simulations of the selected CPU are carried out using FEniCS-POD and HotSpot-POD. Each simulation is verified against its DNS tool. Four test cases are included in this study and the major numerical settings are listed in Table II. Settings in Case 1 are identical to those for illustration of data collection and POD mode generation in Sec. III except for the random power map. Other cases are selected to further validate the findings based on the results from Case 1. A thinner chip ($242\mu\text{m}$) with a heating layer of $55.8\mu\text{m}$ is used in Cases 2-4 than that in Case 1 ($650\mu\text{m}$) to minimize the computational time. Among these 4 cases, Case 1 carries the coarsest mesh,

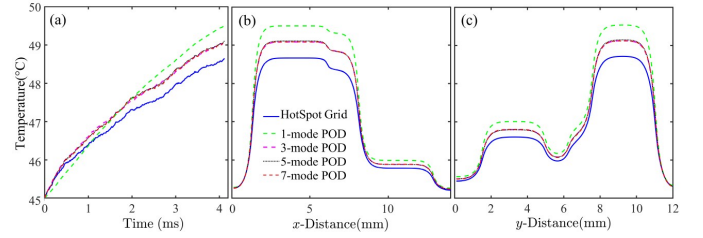


Fig. 7. Temperatures obtained from HotSpot-POD and HotSpot-Grid in Case 1. (a) Dynamic temperature at the intersection of Paths A and B, and the temperature distributions along (b) Path A (c) Path B at $t = 4.2$ ms.

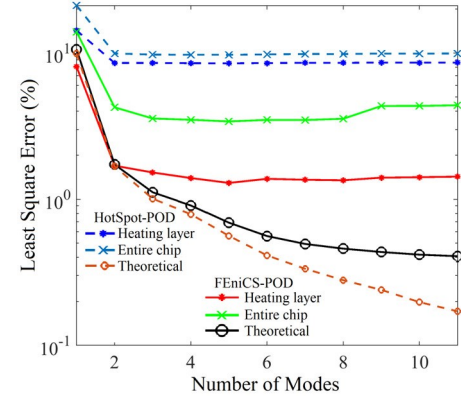


Fig. 8. LS errors in the heating layer and entire chip derived from HotSpot-POD and FEniCS-POD in Case 1. The theoretical errors are also included.

and Case 2 the finest. Unlike Case 1 with uniform dynamic power in each unit, Cases 2-4 include narrow spatial high power density (see their locations in Fig. 1) to induce high thermal gradients and small-diameter hot spots.

TABLE II
TEST CASES: MAJOR NUMERICAL SETTINGS

Case No.	Domain	Power Map	Mesh	$\Delta x \times \Delta y \times \Delta z$ ($\mu\text{m} \times \mu\text{m} \times \mu\text{m}$)	$\frac{\Delta x}{\Delta y}$	$\frac{\Delta y}{\Delta z}$
Case 1	CPU	Uniform power density in each unit	$128 \times 128 \times 13$	$109.4 \times 93.8 \times 50$	1.17	1.88
Case 2	Core 1	Localized spatial high power pulses	$256 \times 256 \times 13$	$18.7 \times 13.3 \times 18.6$	1.39	0.73
Case 3	Core 1	Localized spatial high power density	$64 \times 64 \times 13$	$74.7 \times 53.9 \times 18.6$	1.39	2.90
Case 4	CPU	Localized spatial high power density	$512 \times 512 \times 13$	$27.3 \times 23.4 \times 18.6$	1.17	1.26

A. Demonstration with a Coarser Mesh

In Case 1 with a similar dynamic power distribution in Table I, a dynamic power map is applied with a different random sequence for each unit from that used in the training of POD modes. Dynamic temperature and its spatial distribution in the selected quad-core CPU obtained from FEniCS-POD and HotSpot-POD are illustrated in Figs. 6 and 7, respectively, compared with the results from their DNS approaches.

Dynamic temperature in Fig. 6(a) predicted by FEniCS-POD with just three modes in Core 1 concurs very well with FEniCS-

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

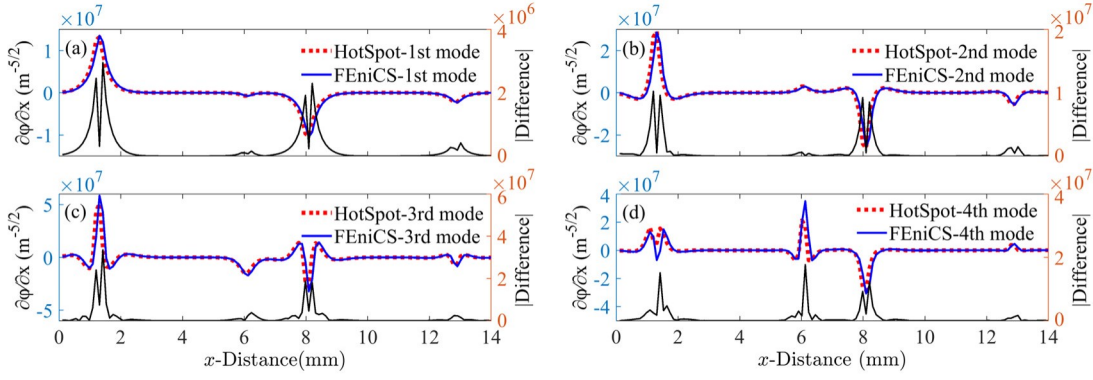


Fig. 9. $\partial \phi / \partial x$ along Path A for the first four modes of HotSpot-POD and FEniCS-POD, together with the absolute value of the difference between these two approaches given by the black line.

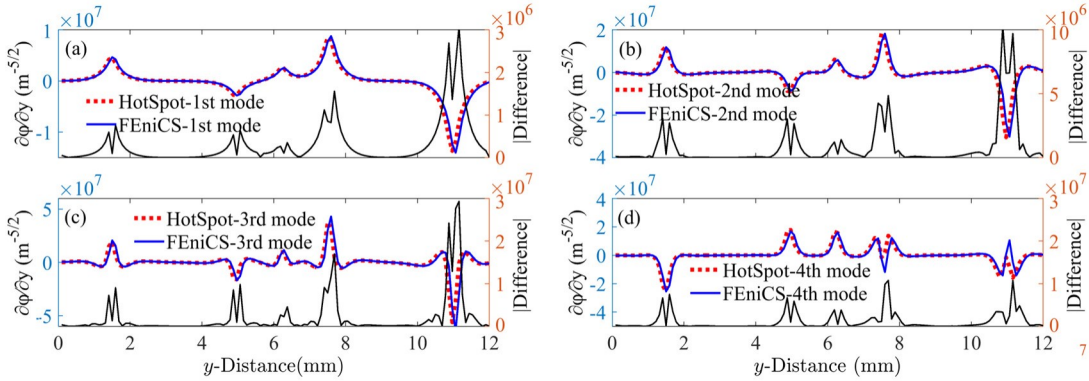


Fig. 10. $\partial \phi / \partial y$ along Path B for the first four modes of HotSpot-POD and FEniCS-POD, together with the absolute value of the difference between these two approaches given by the thin black line.

FEM. Along Paths A and B, the thermal distributions at 4.2ms shown in Figs. 6(b) and 6(c) derived from FEniCS-POD with three or more modes agree very well with FEniCS-FEM except the region around $6\text{mm} < x < 7.5\text{mm}$ with a 2.5%-3% discrepancy. Consistently with the eigenvalue spectrum in Fig. 3, with three or more modes FEniCS-POD's prediction converges to a temporal and spatial thermal solution very close to that offered by FEniCS-FEM. Contrarily to FEniCS-POD, when using three or more modes, Figs. 7(a)-7(b) show that HotSpot-POD leads to a thermal solution with a 11%-13% deviation from that provided by HotSpot-Grid. This is inconsistent with the eigenvalue spectrum for HotSpot-Grid's data given in Fig. 3 and thus indicates that HotSpot-Grid does not offer the thermal solution as accurate as FEniCS-FEM in this case.

The LS percentage error err_{num} in (11) for both POD models w.r.t the solution predicted by each of their DNSs is displayed in Fig. 8, compared to err_{theo} in (10). With two or more modes, err_{num} of HotSpot-POD is considerably greater than err_{theo} and stays near 10% in the entire chip and near 8.5% in the heating layer. FEniCS-POD however leads to an error in the heating layer near 1.5% with three modes and fluctuates around 1.29%-1.4% with four or more modes. In the entire chip, the error fluctuates around 3.4%-3.5% with three to eight modes and stays below 4.4% with nine or more modes. In the thick lower-temperature substrate, the LS error is larger. Within the heating

layer, its LS error closely follows err_{theo} for the first two modes and becomes nearly constant beyond three modes due to the computer precision. The results in Fig. 8 reconfirm that FEniCS-FEM offers a more accurate solution than HotSpot-Grid and provides better-quality thermal data for the POD mode training.

Only a small discrepancy of 4%-5% is observed in Figs. 2(a)-2(c) between the two DNS tools, and yet HotSpot-POD reaches a relatively large LS error in Fig. 8. A closer look at the first POD modes of these 2 POD models in Figs. 4 and 5 and their POD results with just one mode in Figs. 6 and 7 raise an interesting question. The first POD modes shown in Figs. 4 and 5 for FEniCS-POD and HotSpot-POD are very close. Then, why is there a temperature difference of 14%-17% (w.r.t. the ambient) in Core 1 between the one-mode predictions (with $M=1$ in (4), where only ϕ_1 is involved) offered by FEniCS-POD (Fig. 6) and HotSpot-POD (Fig. 7)? To understand this, gradients of POD modes built upon the 2 DNS tools, together with their differences, are illustrated in Figs. 9 and 10 along Paths A in x and B in y . Since FEniCS-FEM offers very good data quality that leads to a very accurate POD thermal prediction, results from HotSpot-POD is compared against FEniCS-POD. Apparently, the slope of the first mode of HotSpot-POD deviates significantly from that of FEniCS-POD at the locations of high thermal gradients (see Figs. 6 and 7.) even though their first modes shown in Figs. 4 and 5 look close.

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

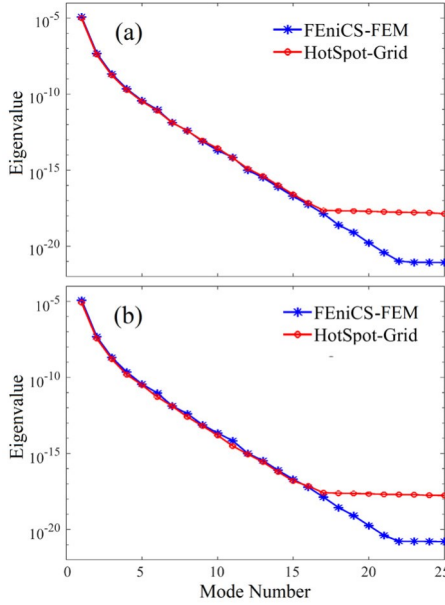


Fig. 11. Eigenvalue spectrums of data collected from HotSpot-Grid and FEniCS-FEM with the (a) coarser mesh (Case 3) and (b) finer mesh (Case 2).

The slope difference between these two approaches evidently increases in the higher modes near high thermal gradients, as indicated in Figs. 9 and 10. This suggests that, even though the 2 DNS tools capture similar averaged thermal behavior over their own data (revealed by the first mode), HotSpot-Grid is not able to evaluate high thermal gradients accurately due to its lumped-element approximation.

The coefficients 2_i defined in (6) and (8) strongly depend on ∇ and ∇ and differ significantly between these 2 POD models, as shown in Table III, where the 2_i ratios of HotSpot-POD to FEniCS-POD are included. The ratios reveal the discrepancy of the POD mode gradients between these 2 DNS tools. As illustrated in Figs. 9 and 10, the mode gradients are strongly influenced by the accuracy of thermal gradients estimated in the DNS tools. The discrepancy in Table III thus offers a good indication of the poor numerical accuracy of the high thermal gradients estimated from HotSpot-Grid. For the diagonal ratios, differences of 2.5%-15% are observed for the first four modes. When using the one-mode model, in (4) at each time step tends reach a steady-state value of $1/2$, and thus HotSpot-POD with a smaller 2_i leads to a higher dynamic temperature than FEniCS-POD, as shown in Fig. 7 compared to Fig. 6.

TABLE III

2_i RATIO FOR FIRST 4×4 ELEMENTS: CASE 1, COARSER MESH

Ratio of g_{ij}	1	2	3	4
1	0.9272	-0.9495	0.8719	-0.4925
2	-0.9495	0.9754	-0.9000	0.8728
3	0.8719	-0.9000	0.8451	-0.8435
4	-0.4925	0.8728	-0.8435	0.9076

Investigation on Case 1 illustrates that the accuracy of the POD approach can be achieved only if the quality of thermal data in POD mode training is adequate. Results also imply that the inaccurate prediction by HotSpot-POD stems from the poor thermal data induced by high thermal gradients. To verify this finding and further demonstrate the capability of the POD method, several other test cases in Table II with different grid sizes and localized high power density are examined below.

B. Impact of Data Quality

Before applying the POD simulation method to the entire CPU with localized high power density and higher resolution, impacts of the mesh coarseness on both DNS methods and their POD models are investigated in a smaller domain of Core 1 given in Cases 2 and 3 of Table II. The Core-1 dimensions are $4.78\text{mm} \times 3.45\text{mm} \times 242\mu\text{m}$ with a total power of 11.6W. In addition to a uniform lower power density, a higher power density is applied to each of the small red squares in Core 1 indicated in Fig. 1. In Cases 2 and 3, the dynamic power density averaged over 60,000 CPU cycles at 3.2 GHz is applied in both DNSs with each time step over 10,000 clock cycles. Thermal simulations of two different meshes for Core 1, finer in Case 2 and coarse in Case 3 given in Table II, are performed in both HotSpot-Grid and FEniCS-FEM to collect thermal data. Both high and low levels of dynamic power density are assigned randomly at each time step. The power density in the small red squares is about 10 to 13 times higher than that in the surrounding area. The mesh in both x and y for Case 2 is 4 times finer than for Case 3 (see Table II) with the same grid size in z . The eigenvalue spectrums of the thermal data are illustrated in Fig. 11, where no noticeable difference is observed for these two DNS tools in the first ten modes.

In the POD simulations, random sequences for lower and higher power levels differently from those used in data collection are applied. Using the coarser mesh, the LS error of FEniCS-POD with three or four modes shown in Fig. 12(a) is 3.36% or 3.6% in the entire chip and 1.39% or 1.48% in the heating layer. While using the finer mesh, the LS error with 3 or four modes is significantly reduced to 1.46% or 1.32% in the entire chip and 0.52% or 0.43% in the heating layer, as shown in Fig. 12(b). The improvement is expected due to a better numerical solution in the areas with high thermal gradients induced by high power density spots. In contrast, HotSpot-POD actually leads to a slight increase in the LS error while reducing the grid size, which is rather unexpected. Figs. 12(a) and 12(b) show that the error increases from 13.4% with the coarser mesh to 16.1% with the finer mesh using two or more modes in the entire chip, and from 7.21% to 8.5% using three or more modes in the heating layer. More information is presented below to understand the confusing outcome from HotSpot-POD.

Dynamic evolution of a peak temperature and temperature profiles along Paths A and B predicted by FEniCS-POD and FEniCS-FEM are illustrated in Fig. 13 for Cases 2 and 3. Thermal results from FEniCS-FEM with these 2 different meshes are very close except for the small difference between the temperature peaks. It is however clearly shown that the agreement between FEniCS-POD and FEniCS-FEM is

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

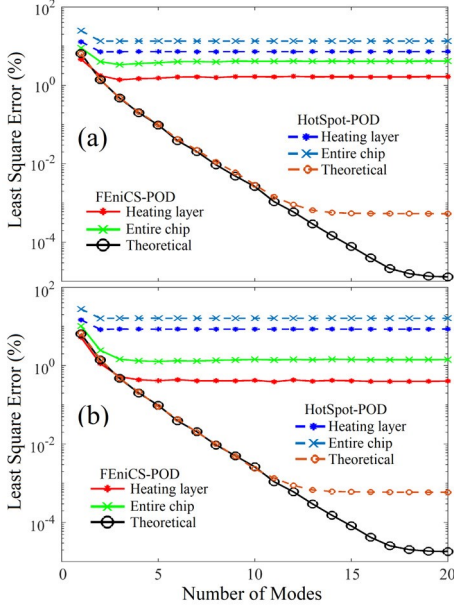


Fig. 12. LS errors of FEniCS-POD and HotSpot-POD with the (a) coarser mesh (Case 3) and (b) finer mesh (Case 2).

noticeably improved, when the grid size reduces, near the boundaries and in the region between 2 temperature peaks.

Similar comparisons are made in Fig. 14 between HotSpot-POD and HotSpot-Grid simulations. While reducing the grid size, the temperature distribution obtained from HotSpot-Grid changes evidently near the boundaries and near and between the 2 temperature peaks along each direction. However, none of the meshes offers good quality thermal data to improve its POD model. Using three or more modes with either mesh, HotSpot-POD converges to a solution inconsistent with HotSpot-Grid's prediction. A closer examination on the POD modes of HotSpot-POD shown in in Figs. 15 and 16 provides more concrete information on the inadequate quality of the HotSpot-Grid data. Since FEniCS-POD offers a very accurate prediction, the HotSpot-POD modes are compared with the FEniCS-POD

modes. Differently from Figs. 4(a) and 5(a), where the first modes for HotSpot-POD and FEniCS-POD are close in Case 1, a large difference is observed between the first modes of HotSpot-POD and FEniCS-POD in Figs. 15(a), 15(e), 16(a) and 16(e). Some discrepancies are also observed in the 2nd and 4th modes. Larger discrepancies are actually observed in the finer mesh case, which indicates that the accuracy of solution from HotSpot-Grid actually degrades when a finer mesh is used (Case 2) to collect the thermal data. By careful comparison between the temperature profiles from HotSpot-Grid (Fig. 14) and FEniCS-FEM (Fig. 13), we also find that the difference of the predictions between HotSpot-Grid and FEniCS-FEM is around 4%-6% for the coarser mesh but as large as 10%-14% for the finer mesh.

TABLE IV

2, RATIO FOR FIRST 4×4 ELEMENTS: CASE 2, FINER MESH

Ratio of $g_{i,j}$	1	2	3	4
1	1.5473	1.4103	1.3560	-1.2944
2	1.4103	1.2707	1.2068	-1.2444
3	1.3560	1.2068	1.1159	-1.1941
4	-1.2944	-1.2444	-1.1941	1.4960

TABLE V

2, RATIO FOR FIRST 4×4 ELEMENTS: CASE 3, COARSER MESH

Ratio of $g_{i,j}$	1	2	3	4
1	1.2538	1.1726	1.2168	1.2143
2	1.1726	1.0897	1.0999	1.1340
3	1.2168	1.0999	1.0899	1.1069
4	1.2143	1.1340	1.1069	1.1074

For the effects of the mode gradients induced by higher thermal gradients, instead of showing gradient profiles, 2, ratios of HotSpot-POD to FEniCS-POD are listed in Tables IV and V for Cases 2 and 3, respectively. The ratios indicate the numerical accuracy of the thermal gradients estimated from HotSpot-Grid. The deviation of HotSpot-POD's 2, from FEniCS-POD's is actually larger in the finer mesh than in the

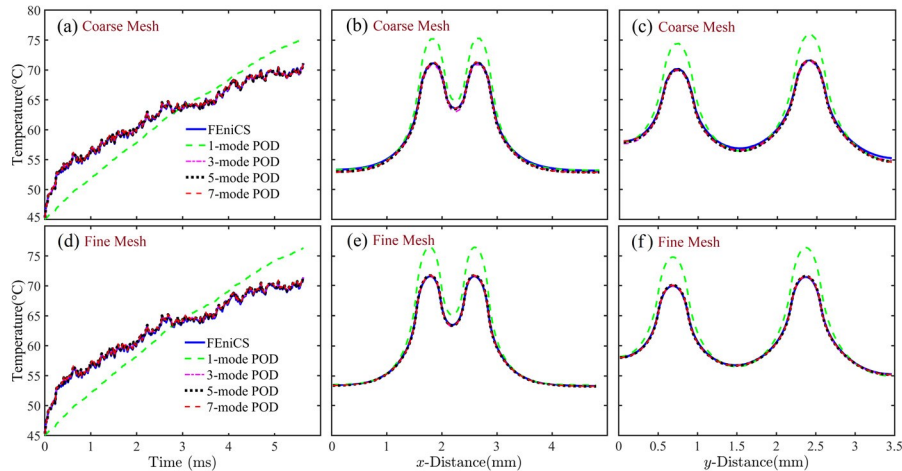


Fig. 13. Predictions of FEniCS-POD compared to FEniCS-FEM, including dynamic temperature in (a) and (d) at the intersection of Paths A and B, the temperature distributions at $t=5.6$ ms along Path A in (b) and (e) and along Path B in (c) and (f). Results in (a)-(c) are for the coarser-mesh domain (Case 3) and in (d)-(f) for the finer-mesh domain (Case 2).

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

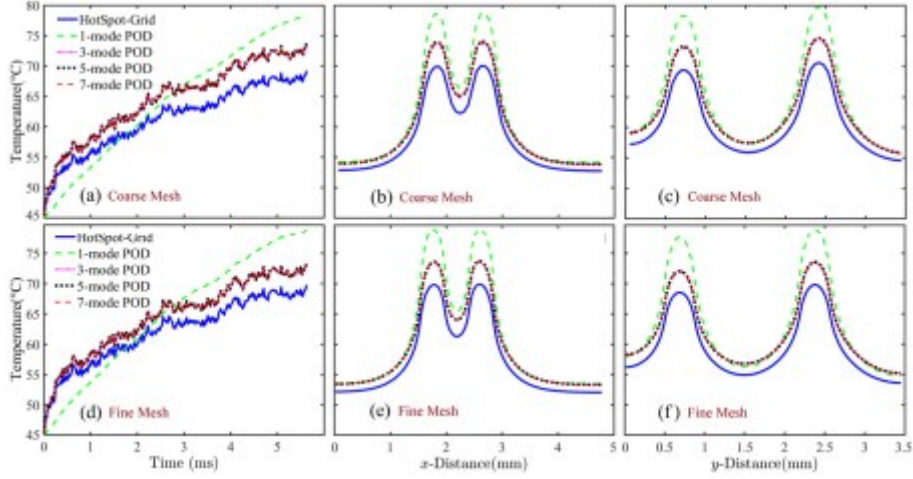


Fig. 14. Predictions of HotSpot-POD compared to HotSpot-Grid, including dynamic temperature in (a) and (d) at the intersection of Paths A and B, the temperature distributions at $t=5.6\text{ms}$ along Path A in (b) and (e), and Path B in (c) and (f). Results in (a)-(c) are for the coarser-mesh domain (Case 3) and in (d)-(f) for the finer-mesh domain (Case 2).

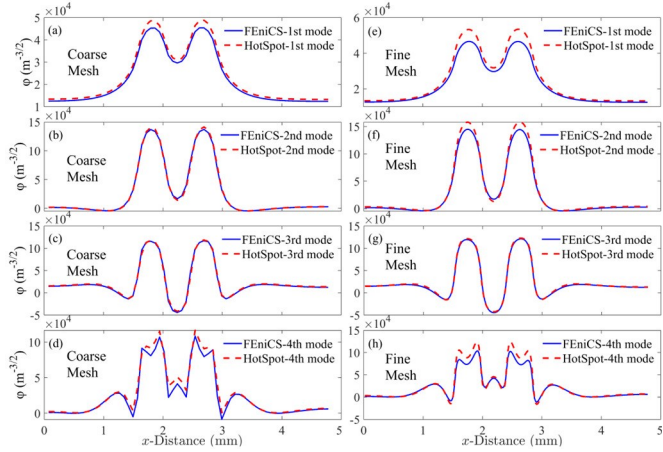


Fig. 15. Comparison of the POD modes along Path A, generated from data collected from FEniCS-FEM and HotSpot-Grid, for (a)-(d) the coarser mesh (Case 3) and (e)-(h) the finer mesh (Case 2).

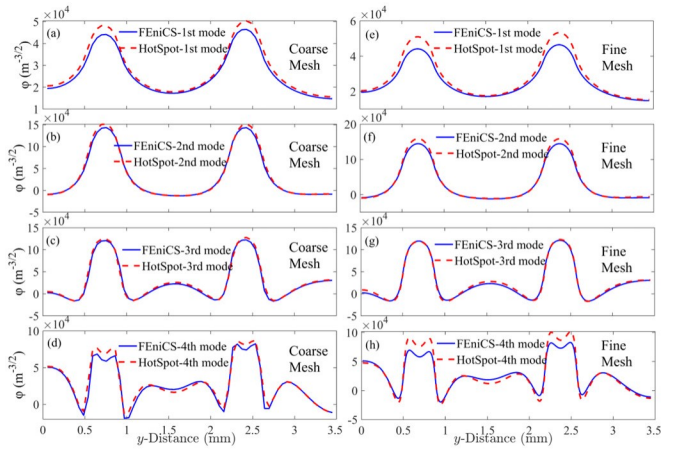


Fig. 16. Comparison of the POD modes along Path B, generated from data collected from FEniCS-FEM and HotSpot-Grid, for (a)-(d) the coarser mesh (Case 3) and (e)-(h) the finer mesh (Case 2).

coarser mesh, similar to the observation of the POD modes in Figs. 15 and 16. A deviation as large as 55% or 27% in the finer mesh case is observed for the first or second mode, respectively. Results illustrated in Figs. 11-16 and Tables IV and V suggest that HotSpot-Grid does not improve the numerical accuracy by reducing the grid size (i.e., the RC element size) in this Core-1 domain perhaps due to the approximation made in the lumped element approach. One may argue that the inclusion of the SF serves the purpose of adjusting its accuracy. To offer consistent solution with more rigorous FEniCS-FEM, one would need to choose $SF > 1$ (to increase the RC time constant) in Case 1 since HotSpot-Grid leads to a faster thermal response. However, since the temperature derived from HotSpot-Grid evolves more slowly in Cases 2 and 3 with localized high power density, a value of $SF < 1$ is needed.

C. Full-Chip Thermal Simulation with Localized High Power Densities

In the Case-4 demonstration for the entire quad-core CPU, in addition to a uniform lower power density applied to each unit,

spatial pulses of higher power density with sizes and locations shown Fig. 1 are applied to Cores 1, 2 and 4. In this case, the grid sizes in x , y and z are 4, 4 and 2.7 times, respectively, smaller than those in Case 1. The eigenvalues and POD modes are again generated by thermal data collected from each of HotSpot-grid and FEniCS-FEM. In the POD simulations, the dynamic power map is generated by random numbers different from those used in the training. In the units with spatial high-power pulses, the higher power density at each time step is approximately 10 to 13 times higher than the lower one.

The eigenvalue spectrums shown in Fig. 17(a) based on both DNS tools are very close for the first 25 modes. The third and fourth mode eigenvalues drop by three and four orders of magnitude, respectively, from the first mode. This indicates that three or four POD modes are able to reach a good accuracy if the quality of data in the training is adequate. Fig. 17(b) shows that, even using a finer mesh with a better AR in z in Case 4 than in Case 1, the LS error of HotSpot-POD in the entire CPU with two or more modes in Case 4 is as large as 15% that is greater than the error (10%) in Case 1 (see Fig. 8). In the heating

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

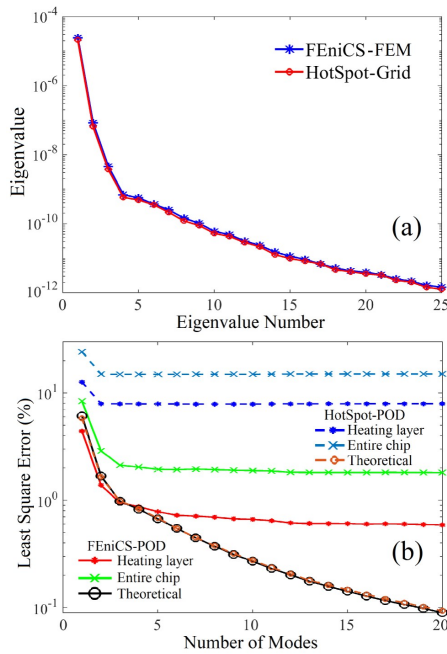


Fig. 17. (a) Eigenvalue spectrum obtained from HotSpot-Grid and FEniCS-FEM data and (b) LS errors of the POD models compared to their theoretical errors.

layer their LS errors become closer, 8% in Case 4 and 8.5% in Case 1. In contrast, the accuracy of FEniCS-POD is clearly improved when decreasing the grid size. Compared to the minimum LS error of 3.4% with three or four modes in Case 1 in the entire chip, Case 4 offers an LS error of 2.1% with three modes and 1.9% with five modes. In the heating layer, more improvement is observed while using a finer mesh. For example, comparing the LS error of Case 1 in Fig. 8, the LS error in Case 4 shown in Fig. 17(b) drops from 1.5% to 0.98% with three modes, from 1.4% to 0.87% with four modes, from 1.29% to 0.78% with five modes and from 1.4% to 0.71% with six modes. With a finer mesh, the LS error goes below 0.6% when more than 15 modes are used.

Thermal solutions from FEniCS-POD and HotSpot-POD are illustrated in Fig. 18, compared to those from their DNS tools. Results from FEniCS-POD and FEniCS-FEM displayed in Figs. 18(d)-18(f) agree very well in time and space even with many narrow spatial pulses of high power density. As clearly observed, FEniCS-POD is able to accurately predict the high thermal gradients and hot spots with sizes below 0.5mm using just three modes. Although the difference of the thermal profiles derived from these 2 DNSs is only 4%-5% in the higher temperature regions, HotSpot-POD arrives at a solution shown in Figs. 18(a)-18(c) significantly different from that offered by its DNS tool. The temperature profiles at $t = 5.6$ ms over the heating layer of the chip predicted by all the methods are also displayed in Fig. 19. Similar to Fig. 18, Fig. 19 shows that agreement between FEniCS-POD and FEniCS is considerably better than that between HotSpot-POD and HotSpot-Grid.

In all test cases, the inaccurate HotSpot-POD prediction is caused by not only its poor-quality POD modes but also the inaccurate gradients of the POD modes resulting from inaccurate numerical calculations of high thermal gradients in HotSpot-Grid. To reduce the paper length, the profiles of POD modes and their gradients (such as those shown in Figs. 5, 10, 15 and 16) are omitted. Instead, Table VI lists the 2_p ratios of HotSpot-POD to FEniCS-POD. As discussed above, the ratios offer a good indication about the accuracy of thermal gradients calculated in HotSpot-Grid. The discrepancies of 2, between the data collected from HotSpot-Grid and FEniCS-FEM in Case 4 are as small as 5.69% to 6.96%, for $i = 1, 3$ and 4, which are all smaller than those in Case 1. However, $2_{p,p}$ of HotSpot-POD in Case 4 with a finer mesh reveals a large deviation of 24.4% from that of FEniCS-POD, compared to the largest deviation of 15.5% observed for $2_{v,v}$ in Case 1. Since $p \gg v$, as given in Figs. 3 and 17(a), $2_{p,p}$ has a considerably stronger effect on the accuracy of the POD model than $2_{v,v}$. This explains why HotSpot-POD does not offer a more accurate prediction in Case 4 with a finer mesh than in Case 1 with a coarse mesh. This

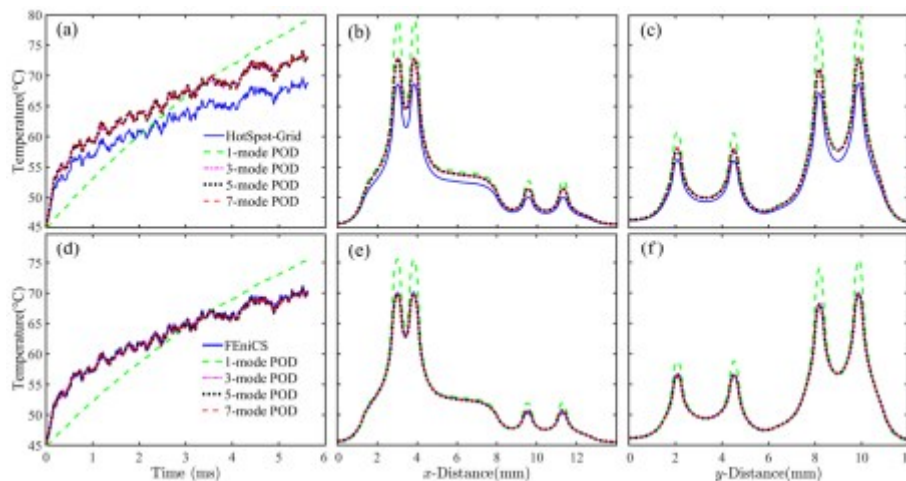


Fig. 18. (a) Dynamic temperature at the intersection of Paths A and B, the temperature distributions at $t = 5.6$ ms along (b) Path A and (c) Path B, predicted by HotSpot-POD compared to HotSpot-Grid. (d) Dynamic temperature at the intersection of Paths A and B, the temperature distributions at $t = 5.6$ ms along (e) Path A and (f) Path B at $t = 5.6$ ms, derived from FEniCS-POD compared to FEniCS-FEM.

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

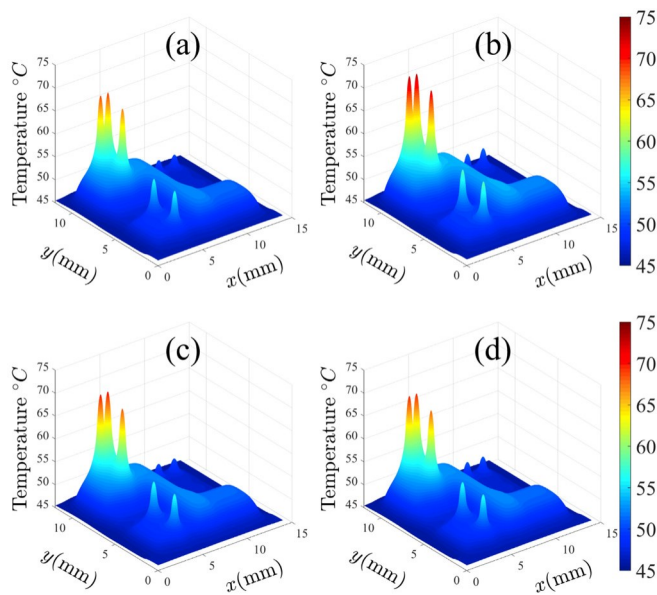


Fig. 19. Temperature profiles at $t = 5.6$ ms predicted by (a) HotSpot-Grid, (b) HotSpot-POD with three modes, (c) FEniCS and (d) FEniCS-POD with three modes.

also suggests that even with a considerably finer mesh HotSpot-Grid does not offer more accurate high thermal gradients.

TABLE VI

2, RATIO FOR FIRST 4×4 ELEMENTS: CASE 4, FINER MESH CPU

Ratio of $g_{i,j}$	1	2	3	4
1	1.0569	1.1609	-1.1226	-0.9355
2	1.1609	1.2425	-1.1746	-1.0076
3	-1.1226	-1.1746	1.0696	0.9068
4	-0.9355	-1.0076	0.9068	1.0592

Differently from Case 1, temperature predicted by HotSpot-Grid in Case 4 evolves more slowly than that by FEniCS-FEM. If the SF is used to improve the HotSpot-Grid accuracy in this case, a value of $SF < 1$ would be needed while $SF > 1$ is needed in Case 1. The full-chip thermal simulations in Case 4 reconfirm the inadequate quality of the thermal solution data collected from HotSpot-Grid. This case also further validates the effectiveness of the POD simulation method at the chip-level if it is built upon good-quality data.

V. DISCUSSIONS

The investigation on thermal simulations of a quad-core CPU has uncovered some interesting aspects on effectiveness of the POD models associated with the quality of the DNS tools used to collect data for the POD mode training. The POD LS errors presented in Sec. IV. are summarized in Tables VII and VIII for simulations of the entire chip (Cases 1 and 4) and Core 1 (Cases 2 and 3), respectively. For the entire CPU, Case 4 is demonstrated with localized narrow high power densities and a

higher resolution, compared against Case 1 with a coarser mesh without high power densities. FEniCS-FEM in Case 4 with a finer mesh offers more accurate thermal solution data and thus significantly improves the quality of the POD modes. As shown in Table VII, a substantial reduction in the LS error is thus achieved in both the heating layer and the entire chip even though high thermal gradients and narrow hot spots are induced in many locations (also see the discussions in Figs. 8 and 17). For the demonstration of FEniCS-POD in a smaller domain of Core 1, as shown in Table VIII for Cases 2 and 3 (also see the discussion for Figs. 12-14), a similar finding is observed. Namely, a considerably smaller LS error is obtained when using the finer-mesh POD modes in Case 2.

TABLE VII

LS ERROR OF POD MODELS FOR FULL CHIP SIMULATIONS

POD model	No. of modes	Entire CPU		Heating layer	
		Case 1 coarse mesh	Case 4 fine mesh	Case 1 coarse mesh	Case 4 fine mesh
HotSpot POD	3	10%	15%	8.5%	8%
	4	10%	15%	8.5%	8%
FEniCS POD	3	3.5%	2.1%	1.5%	0.98%
	4	3.5%	2%	1.4%	0.87%

TABLE VIII

LS ERROR OF POD MODELS FOR CORE-1 SIMULATIONS

POD model	No. of modes	Entire Core 1		Heating layer of Core 1	
		Case 2 fine mesh	Case 3 coarse mesh	Case 2 fine mesh	Case 3 coarse mesh
HotSpot POD	3	16.1%	13.4%	8.5%	7.2%
	4	16.1%	13.4%	8.5%	7.2%
FEniCS POD	3	1.46%	3.36%	0.52%	1.39%
	4	1.32%	3.6%	0.43%	1.48%

The same applications of HotSpot-POD, built upon thermal data from HotSpot-Grid, however, reveal considerably larger LS errors than FEniCS-POD in all test cases, as shown in Tables VII and VIII. The examinations carried out in this work have revealed that inaccurate numerical solutions derived from HotSpot-Grid are caused by its incapability of offering accurate high thermal gradients probably due to the lumped-element approximation. As a result, the grid size used in HotSpot-Grid is not always correlated to its numerical accuracy, as observed in all the demonstrations. Results have also shown that the larger LS errors of HotSpot-POD are caused by the poor-quality POD modes generated from the inadequate thermal gradients estimated by HotSpot-Grid. In contrast, due to the rigorous FEM implemented in FEniCS, the accuracy of FEniCS-POD is enhanced consistently as the mesh size reduces to improve the numerical accuracy in high thermal gradients.

All the results demonstrated in this study indicate that quality of the solution data from the DNS tool is the key to improve the accuracy of the developed POD model. In particular, accurate gradients of POD modes, reflected by the 2_{r} values in Tables

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

III-VI, are needed, and this requires a high numerical accuracy of high thermal gradients offered by the DNS tool. With an appropriate DNS tool, the general practice is to reduce the grid size in high thermal gradient regions to improve numerical solutions. This however increases the computing time and memory space needed in the training process that includes thermal data collection, POD mode generation and calculations of POD parameters. However, once the training is finished, the reduction in the DoF is incredibly attractive.

Comparison of the computational time and DoF between FEniCS-POD and FEniCS-FEM is included in Table IX. Considering Case 1 with a coarser mesh of $128 \times 128 \times 13$ (or 212,992), the lowest resolution in our study, if four POD modes are used, FEniCS-POD for the entire CPU offers a reduction in the DoF by more than 53,000 times (or more than 70,000 times with three modes). This leads to a saving in the computational time more than 3,600 times ($10,700/2.92$), compared to FEniCS-FEM. To obtain the temperature profile, the post processing in (2) to evaluate \vec{u} , is actually more time consuming than the POD simulation to solve \vec{u} in (4). For most applications, the thermal information is needed only in a small fraction of the entire chip, particularly in the regions where hot spots are located and most likely in the top heating layer of the cores. Unlike the FEM where the entire domain must be solved, the POD method could selectively calculate the temperature at some grid points. If only the temperature in the heating layer is of interest, an improvement in computing time more than four orders of magnitude can be achieved (or $10,700/0.77 = 13,896$) with four modes. That is, instead of 3-hour computational time in this case using FEniCS-FEM C++, the temperature profile in the heating layer of the quad-core CPU can be obtained from the POD modeling technique in 0.77 second or for the entire CPU within a few seconds.

TABLE IX

COMPUTATIONAL TIME AND DOF FOR THE QUAD-CORE CPU

Case No.	Simulator		Time (s) Entire chip	Time (s) Heating layer	DoF
Case 1	FEniCS-POD	3 modes	2.11	0.62	3
		4 modes	2.92	0.77	4
		5 modes	3.81	0.91	5
	FEniCS-FEM		10,700		212,992
Case 4	FEniCS-POD	3 modes	29.41	7.27	3
		4 modes	38.95	8.86	4
		5 modes	51.05	11.01	5
	FEniCS-FEM		167,000		3,407,872

With a finer mesh in Case 4, the computational times for both FEniCS-POD and FEniCS-FEM are increased, as shown in Table IX. With a finer mesh, the accuracy of FEniCS-POD is however significantly improved, as discussed in Fig. 17 compared to Fig. 8. In particular, three modes in Case 4 are able to offer an LS error smaller than the minimum error observed in Case 1. Moreover, in Case 4, FEniCS-POD with three or four modes offers a speedup near 5,700 or 4,300 times to evaluate temperature in the entire CPU, and 23,000 or 19,000 times in

the heating layer, respectively, compared to FEniCS-FEM. It is also found in our study that the simulation time needed by HotSpot-Grid is approximately 50% of what FEniCS-FEM requires for the same simulation domain and numerical settings.

This investigation also reveals a useful feature for the POD method. While the LS error significantly reduces as the mesh resolution increases, the number of modes needed for the LS error curve to become nearly flattened does not change much. For example, the LS errors from FEniCS-POD for Core 1 in Fig. 12(b) for Case 2 (finer mesh) and in Fig. 12(a) for Case 3 (coarser mesh) become nearly invariant with just three or four modes. As shown in Fig. 17(b) for Case 4 compared to Fig. 8 for Case 1, the LS error from FEniCS-POD for the entire chip becomes flattened with three or four modes in both Case 1 (coarser mesh) and Case 4 (finer mesh). To further improve the accuracy of the POD model, one can therefore collect fine-resolution data to generate good-quality, robust POD modes, and the computing time for the POD simulation to reach a higher accuracy would not change much. Although the post process in (2) can offer the spatial temperature with a resolution as high as DNS, to significantly minimize the computational resources, only the thermal profile in a few grid points of high temperature regions for is needed.

The intensive computational effort needed in the training process is the major drawback of the proposed POD simulation method. The approach is however valuable for some crucial applications. One of the useful applications of the proposed approach is the real-time thermal-aware management. Once a CPU or GPU is trained to adapt various dynamic power maps and BCs, its POD model is able to offer the dynamic thermal profile in the entire CPU/GPU including all significant hot spots with a reasonable resolution within seconds. To the best of our knowledge, there is no other method available to offer such a task. The proposed approach is particularly valuable for true run-time thermal-aware task scheduling of CPUs and GPUs due to the fast-growing demand for high performance computing.

To further improve the POD method and expand its adoptability, structures for a specific technology group can be partitioned into smaller building blocks whose trained POD modes and parameters are then stored in a database. These stored POD blocks can be assembled to create larger structures, similar to many technologies whose structures are primarily constructed by building blocks, such as CPUs/GPUs, photonic crystals, metamaterials, nanostructures, etc. The approach using building blocks will also offer more efficient training for smaller blocks and more effective parallel computing for a structure with a large number of POD building blocks.

VI. CONCLUSIONS

The POD finds the modes that contain essential information on thermal behaviors embedded in the dynamic solution data collected from a DNS tool. This study has illustrated that, if the collected solution data is not accurate enough, the projection of the heat transfer equation given in (3) along the generated POD modes leads to a POD model that offers inconsistent solution with the heat transfer equation. In addition to offering an

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

effective simulation approach for a large-scale simulation domain, as demonstrated in this work, the POD simulation method can also be used to rigorously determine the accuracy of the DNS tool.

The DNS tools in this investigation include a rigorous FEM implemented in FEniCS [32] and a popular chip-level thermal simulator, HotSpot-Grid [41], [43]. Several heat source excitations and different mesh resolutions have been investigated in the entire or part of the quad-core CPU, AMD ATHLON II X4 610e. It has been demonstrated that the POD model built upon FEniCS-FEM offers very accurate predictions of dynamic thermal distributions in all the test cases with an extreme small DoF (three to four modes) even in the domains with many small-size high-temperature hot spots. This leads to a speedup of approximately four orders of magnitude to predict the dynamic thermal distribution for the entire CPU, compared to its DNS tool, FEniCS-FEM. When only the temperature in the heating layer is needed, a speedup of five orders can be achieved. In contrast, the POD model built upon HotSpot-Grid offers an inconsistent dynamic thermal solution with the heat transfer equation in all cases due to the inadequate-quality thermal solution data collected from its DNS tool, HotSpot-Grid, especially in the areas with high thermal gradients. HotSpot-Grid provides reasonable solutions in most cases (a deviation of 4%-5% from FEniCS-FEM's prediction) but leads to a deviation as large as 16.1% in the entire chip even in Case 2 with the finest resolution. Unlike FEniCS-FEM, the numerical quality of HotSpot-Grid seems to be uncorrelated to the mesh resolution for the cases we studied.

It has been shown that the change in the mesh resolution of the thermal data collected from the DNS for a specific structure does not influence the DoF needed to reach the minimum LS error for the POD models. Also, the higher mesh resolution implemented in the good-quality DNS tool offers a more accurate POD model. If the efficiency and accuracy are the major concern for full-chip thermal simulation and if one can afford more computational effort using a higher resolution mesh in the training of the POD modes, the following practice will be useful for the POD simulation method. After the simulation in the POD space, the post-processing calculation to obtain the dynamic temperature distribution in real space is needed only in the areas around the heat sources with higher power density, such as the cores in our study. Also, instead of the temperature at every discrete point in the areas with high temperature or high thermal gradients, only temperature at selected points needs to be calculated from (2). The training only needs to be done once and the POD modes and parameters can be stored in a technology database for full-chip dynamic thermal analysis. The extremely efficient POD simulation method with a high accuracy will be very attractive, e.g., for run-time thermal-aware task scheduling of CPUs or GUPs whose major concern is to capture high thermal gradients and hot spots in the chip. To ease the computational time and memory space needed for the training, an alternative is to partition the entire chip into smaller domains. This will be investigated in the near future.

REFERENCES

- [1] C. E. Leiserson, et al. "There's plenty of room at the Top: What will drive computer performance after Moore's law?" *Science*, vol. 368, no. 6495, p. eaam9744, 2020.
- [2] S. Pozder, et al., "Progress of 3D integration technologies and 3D interconnects," *Proc. IEEE IITC*, pp. 213-215, 2007.
- [3] S. M. Alam, R. E. Jones, S. Pozder, R. Chatterjee, A. Jain "New design considerations for cost effective three-dimensional (3D) system integration," *IEEE Trans. VLSI Syst.*, vol. 18, pp. 450-460, 2010.
- [4] C. Wang, X. J. Huang, K. Vafai "Analysis of hotspots and cooling strategy for multilayer three-dimensional integrated circuits," *Appl. Thermal Eng.*, vol. 186, 116336, 2021.
- [5] L. Choobineh, A. Jain "Determination of temperature distribution in three-dimensional integrated circuits (3D ICs) with unequally-sized die," *Appl. Thermal Eng.*, vol. 56no. pp. 176-184, 2013.
- [6] L. Hou, T. Ye, Q. Luo, J. Fu, J. Wang, "A method to alleviate hot spot problem in 3D IC", *Microelectronic Eng.*, vol. 190, pp. 19-27, 2018.
- [7] B. Ding, et al., "Coupling management optimization of temperature and thermal stress inside 3D-IC with multi-cores and various power density," *Int. Commun. Heat Mass Transf.*, vol. 120, 105021, 2021.
- [8] B. Peng, T. Palpanas, P. Fatourou, "Paris+: Data series indexing on multi-core architectures," *IEEE Tran. Knowledge & Data Eng.*, vol. 33, no. 5, pp. 2151-2164, 2020.
- [9] S. V. Patil, D. B. Kulkarni, "A Review of Dimensionality Reduction in High-Dimensional Data Using Multi-core and Many-core Architecture," *Proc. SCEE*, 54-63, 2018.
- [10] Y. Shigeto, M. Sakai, "Parallel computing of discrete element method on multi-core processors," *Particuology*, vol. 9, pp. 398-405, 2011.
- [11] D. Zhang, H. Hang, X. Bi, "Comparison and analysis of GPGPU and parallel computing on multicore CPU," *Int. J. Inf. Educ. Technol.*, vol. 2, no. 2, pp. 185-187, 2012.
- [12] W. Yang, K. Li, "A hybrid computing method of SpMV on CPU-GPU heterogeneous computing systems," *J. Parallel Distrib. Comput.*, vol. 104, pp. 49-60, 2017.
- [13] J. von Kistowski, et al., "Measuring the Energy Efficiency of Transactional Loads on GPGPU," *Proc. ACM/SPEC Int. Conf. Perform. Eng. (ICPE)*, pp. 219-230, 2019.
- [14] K. Cho, H. Bahn, "Characterizing Fine-Grained Resource Utilization for Multitasking GPGPU in Cloud Systems," *IEEE Access*, vol. 9, pp. 161507-161519, 2021.
- [15] B. Nie, L. Yang, A. Jog, F. Smirni, "Fault site pruning for practical reliability analysis of GPGPU applications," *Proc. 51st Annu. IEEE/ACM Int. Symp. Microarchit.*, 749-761, 2018.
- [16] L. Yang, B. Nie, A. Jog, E. Smirni, "Practical Resilience Analysis of GPGPU Applications in the Presence of Single- and Multi-Bit Faults," *IEEE Trans. Comput.*, vol. 70, pp. 30-44, 2021.
- [17] H. F. Sheikh, I. Ahmad, Z. Wang, S. Ranka, "An overview and classification of thermal-aware scheduling techniques for multi-core processing systems," *Sustain. Comput. Informat. Syst.*, vol. 2, no. 3, pp. 151-169, 2012.
- [18] A. Heinig, R. Fischbach, M. Dietrich, "Thermal analysis and optimization of 2.5D and 3D integrated systems with wide I/O memory," *Proc. Conf. Therm. Thermomech. Phenom. Electron. Syst.*, 86-91, 2014.
- [19] J. Zhou, et al., "Thermal-aware correlated two-level scheduling of real-time tasks with reduced processor energy on heterogeneous MPSoCs," *J. Sys. Architect.*, vol. 82, pp. 1-11, 2018.
- [20] Y. W. Chang, et al., "Electromigration mechanism of failure in flip-chip solder joints based on discrete void formation," *Sci. Rep.*, 7:1-16, 2017.
- [21] X. Huang, A. Kteyan, X. Tan, V. Sukharev, "Physics-based electromigration models and full-chip assessment for power grid networks," *IEEE Trans. CAD ICs Syst.*, vol. 35, pp. 1848-1861, 2016.
- [22] Y. Liu, et al., "Joule heating enhanced electromigration failure in redistribution layer in 2.5D IC," *Proc. ECTC*, pp. 1359-1363, 2016.
- [23] S. S. Anandan, V. Ramalingam, "Thermal management of electronics: A review of literature," *Therm. Sci.*, pp. 125-26, 2008.
- [24] J. W. Sheaffer, K. Skadron, D. P. Luebke, "Studying thermal management for graphics-processor architectures," *Proc. ISPASS*, 54-65, 2005.
- [25] Nath R, Ayoub R, Rosing TS (2013) Temperature aware thread block scheduling in GPGPUs. *Proc. 50th ACM/EDAC/IEEE Design Autom. Conf. (DAC)*, 1-6.
- [26] H. Khan, Q. Bashir, M. U. Hashmi, "Scheduling based energy optimization technique in multiprocessor embedded systems," *Proc. ICEET*, 1-8, 2018.

> REPLACE THIS LINE WITH YOUR MANUSCRIPT ID NUMBER (DOUBLE-CLICK HERE TO EDIT) <

- [27] Y. Lee, K. G. Shin, H. S. Chwa, "Thermal-aware scheduling for integrated CPUs-GPU platforms," *ACM Trans. Embedded Comput. Syst.*, vol. 18, no. 5s, pp. 1-25, 2019.
- [28] A. K. Coskun, T. S. Rosing, K. Whisnant, K. Gross, "Temperature-aware MPSoC scheduling for reducing hot spots and gradients," *Asia South Pacific Des. Autom. Conf.*, 49-54, 2008.
- [29] A. Rudi, A. Bartolini, A. Lodi, L. Benini, "Optimum: Thermal-aware task allocation for heterogeneous many-core devices," *Proc. Int. Conf. High Perform. Comput. Simulation*, pp. 82-87, 2014.
- [30] ANSYS. [Online]. Available: <https://www.ansys.com/>
- [31] COMSOL. [Online]. Available: <https://www.comsol.com/>
- [32] FEniCS Project. <https://fenicsproject.org/>. Accessed 5 February 2022.
- [33] FREEFEM. [Online]. Available: <https://freefem.org/>
- [34] S. Varshney, et al., "Nanotherm: An analytical Fourier-Boltzmann framework for full chip thermal simulations," *Proc. ICCAD*, 1-8, 2019.
- [35] Y. Zhan, S. S. Sapatnekar, "High-efficiency green function-based thermal simulation algorithms," *IEEE Trans. CAD ICs Syst.*, vol. 26, no. 9, pp. 1661-1675, 2007.
- [36] H. Sultan, S. R. Sarangi, "Variability-aware thermal simulation using CNNs. *Proc. VLSID*, 65-70.
- [37] Sultan H, Sarangi SR (2020) A fast leakage-aware Green's-function based thermal simulator for 3-D chips," *IEEE Trans. VLSI Syst.*, vol. 28, no. 11, pp. 2342-2355, 2020.
- [38] Y. Zhan, S. S. Sapatnekar, "A high efficiency full-chip thermal simulation algorithm," *Proc. ICCAD*, pp. 635-638, 2005.
- [39] J. H. Park, A. Shakouri, S. M. Kang, "Fast evaluation method for transient hot spots in VLSI ICs in packages," *Proc. 9th ISQED*, pp. 600-603, 2008.
- [40] A. Ziabari, et al., "Power blurring: Fast static and transient thermal analysis method for packaged integrated circuits and power devices, *IEEE Trans. VLSI Syst.*, vol. 22, no. 11, pp. 2366-2379, 2014.
- [41] W. Huang, K. Sankaranarayanan, R. J. Ribando, M. R. Stan, K. Skadron, "An Improved Block-Based Thermal Model in HotSpot 4.0 with Granularity Considerations," *Proc. Ann. WDDD*, 2007.
- [42] K. Skadron, M. R. Stan, W. Huang, S. Velusamy, K. Sankaranarayanan, D. Tarjan, "Temperature-aware microarchitecture," *Proc. Int. Symp. Comput. Architecture*, pp. 2-13, 2003.
- [43] HotSpot 6.0 Temperature Modeling Tool. [Online]. Available: <http://lava.cs.virginia.edu/HotSpot/versions.htm> and <https://github.com/IFTE-EDA/HotSpot/blob/master/HOWTO>
- [44] X. Zhou, J. Yang, Y. Xu, Y. Zhang, J. Zhao, "Thermal-Aware Task Scheduling for 3D Multicore Processors," *IEEE Trans. Parallel & Distributed Systems*, vol. 21, no. 1, pp. 60-71, 2010.
- [45] H. H. Chu, Y. C. Kao, Y. S. Chen, "Adaptive thermal-aware task scheduling for multi-core systems," *J. Systems and Software*, vol. 99, pp. 155-174, 2015.
- [46] J. Zhou, et al., "Thermal-Aware Task Scheduling for Energy Minimization in Heterogeneous Real-Time MPSoC Systems," *IEEE Tran. CAD ICs & Sys.*, vol. 35, pp. 1269-1282, 2016.
- [47] S.K.S. Tyagi, D. K. Jain, S. L. Fernandes, P. K. Muhuri, "Thermal-aware power-efficient deadline based task allocation in multi-core processor," *J. Comp. Sci.*, vol. 19, pp. 112-120, 2017.
- [48] A. Iranfar, et al., "TheSPoT: Thermal Stress-Aware Power and Temperature Management for Multiprocessor Systems-on-Chip," *IEEE Tran. CAD of ICs & Sy*, vol. 37, pp. 1532-1545, 2011.
- [49] D. Fetis, P. Michaud, "An Evaluation of HotSpot-3.0 Block-Based Temperature Model", *Proc. WDDD*, 2006.
- [50] A. Ziabari, E. K. Ardestani, J. Renau, A. Shakouri "Fast thermal simulators for architecture level integrated circuit design, 27th IEEE Annu. Semicond. Thermal Meas. Manage. Symp., pp. 70-75, 2011.
- [51] F. Tavakkoli, S. Ebrahimi, S. Wang, K. Vafai, "Analysis of critical thermal issues in 3D integrated circuits," *Int. J. Heat and Mass Transfer*, vol. 97, pp. 337-352, 2016.
- [52] C. Wang, X. J. Huang, K. Vafai, "Analysis of hotspots and cooling strategy for multilayer three-dimensional integrated circuits," *Appl. Thermal Eng.*, vol. 186, 116336, 2021.
- [53] W. Nakayama, "Study on Heat Conduction in a Simulated Multicore Processor Chip—Part II: Case Studies," *ASME. J. Electron. Packag.*, vol. 135, no. 2, p. 021003, 2013.
- [54] J. L. Lumley, "The structure of inhomogeneous turbulent flows," *Atmospheric Turbulence and Radio Wave Propagation*, 1967.
- [55] G. Berkooz, P. Holmes, J. L. Lumley, "The proper orthogonal decomposition in the analysis of turbulent flows," *Annu. Rev. Fluid Mech.*, vol. 25, pp. 539-575, 1993.
- [56] W. Jia, B. Helenbrook, M. C. Cheng, "Fast thermal simulation of FinFET circuits based on a multi-block reduced-order model," *IEEE Trans. CAD ICs & Systems*, vol. 35, pp. 1114-1124, 2016.
- [57] W. Jia, B. T. Helenbrook, M. C. Cheng, "Thermal modeling of multi-fin field effect transistor structure using proper orthogonal decomposition," *IEEE Trans. Electron Devices*, vol. 61, no. 8, pp. 2752-2759, 2014.
- [58] V. Puzyrev, M. Ghommam, S. Meka. "pyROM: A computational framework for reduced order modeling," *J. Comp. Sci.*, vol. 30, pp. 157-173, 2019.
- [59] G. Huang, "Application of proper orthogonal decomposition in fast Fourier transform-assisted multivariate nonstationary process simulation," *J. Eng. Mech.*, vol. 141, p. 04015015, 2015.
- [60] O. S. Ojo, S. Grivet-Talocia, M. Paggi, "Model order reduction applied to heat conduction in photovoltaic modules," *Composite Structures*, vol. 119, pp. 477-486, 2015.
- [61] A. Nokhosteen, M. M. Soltani, B. Barabadi, "Reduced order modeling of transient heat transfer in microchip interconnects," *ASME. J. Electron. Packag.*, vol. 141, pp. 011002-9, 2019.
- [62] K. Lu, et al., "Review for order reduction based on proper orthogonal decomposition and outlooks of applications in mechanical systems," *Mech. Sys. Sig. Processing*, vol. 123, pp. 264-297, 2017.
- [63] L. Jiang, Y. Liu, M. C. Cheng, "An effective and accurate data-driven approach for thermal simulation of CPUs," 20th IEEE Intersoc. Conf. Thermal & Thermomech. Phenomena in Electronic Sys. (iTherm 2021), pp. 1008-1014, 2021.
- [64] CPU-World. [Online]. Available: [https://www.cpu-world.com/CPUs/K10/AMD-Athlon%20II%20X4%20610e%20-%20AD610EHDK42GM%20\(AD610EHDGMBBOX\).html](https://www.cpu-world.com/CPUs/K10/AMD-Athlon%20II%20X4%20610e%20-%20AD610EHDK42GM%20(AD610EHDGMBBOX).html)
- [65] Dev, A. N. Nowroz, S. Reda, "Power mapping and modeling of multi-core processors," *IEEE Int. Symp. Low-Power Electron. Design*, pp. 39-44, 2013.
- [66] L. Sirovich, "Turbulence and the dynamics of coherent structures Part I: Coherent structures," *Quart. Appl. Math.* Vol. 45, pp. 561-571, 1987.
- [67] W. Jia, M.C. Cheng, "A methodology for thermal simulation of interconnects enabled by model reduction with material property variation," *J. of Computational Sci*, vol. 61, 101665, 2022.
- [68] M. R. Stan, K. Skadron, M. Barcella, W. Huang, K. Sankaranarayanan and S. Velusamy, "Hotspot: A dynamic compact thermal model at the processor-architecture level", *Microelectron. J.*, vol. 34, pp. 1153-1165, 2003.
- [69] R. Mahajan, "Thermal management of CPUs: A perspective on trends needs and opportunities", Keynote Presentation at the 8th Int. Workshop on Thermal Investigations of ICs and Syst., 2002.
- [70] H. Kattan, S. W. Chung, J. Henkel and H. Amrouch, "On-demand mobile CPU cooling with thin-film thermoelectric array", *IEEE Micro*, vol. 41, no. 4, pp. 67-73, Jul. 2021.