

A Case for Packageless Processors

Saptadeep Pal*, Daniel Petrisko†, Adeel A. Bajwa*, Puneet Gupta*, Subramanian S. Iyer*, and Rakesh Kumar†

*Department of Electrical and Computer Engineering, University of California, Los Angeles

†Department of Electrical and Computer Engineering, University of Illinois at Urbana-Champaign
{saptadeep, abajwa, s.s.iyer, puneetg}@ucla.edu, {petrisk2, rakeshk}@illinois.edu

Abstract—Demand for increasing performance is far outpacing the capability of traditional methods for performance scaling. Disruptive solutions are needed to advance beyond incremental improvements. Traditionally, processors reside inside packages to enable PCB-based integration. We argue that packages reduce the potential memory bandwidth of a processor by at least one order of magnitude, allowable thermal design power (TDP) by up to 70%, and area efficiency by a factor of 5 to 18. Further, silicon chips have scaled well while packages have not. We propose *packageless processors* - processors where packages have been removed and dies directly mounted on a silicon board using a novel integration technology, Silicon Interconnection Fabric (Si-IF). We show that Si-IF-based packageless processors outperform their packaged counterparts by up to 58% (16% average), 136% (103% average), and 295% (80% average) due to increased memory bandwidth, increased allowable TDP, and reduced area respectively. We also extend the concept of packageless processing to the entire processor and memory system, where the area footprint reduction was up to 76%.

Keywords-Packageless Processors, Silicon Interconnect Fabric

I. INTRODUCTION

Conventional computing is at a tipping point. On one hand, applications are fast emerging that have higher performance, bandwidth, and energy efficiency demands than ever before. On the other hand, the end of Dennard scaling [1] as well as Moore’s law transistor scaling diminishes the prospect of easy performance, bandwidth, or energy efficiency scaling in future. Several promising and disruptive approaches are being explored, including (but not limited to) specialization [2], approximation [3], 3D integration [4], and non-CMOS devices [5].

Current systems place processor and memory dies inside packages, which allows them to be connected to the PCB and subsequently to other dies. A striking observation is that in the last two decades while silicon chips have dimensionally scaled by 1000X, packages on printed circuit boards (PCBs) have merely managed 4X [6]. This absence of “system scaling” can severely limit performance of processor systems. This realization has motivated the push toward 3D and 2.5D integration schemes which alleviate the problem but do not address the root cause. In this paper, we propose another approach - removing the package from the processor altogether.

At first glance, removing the package from the processor may seem both simple in implementation and, at best, incremental in benefits. However, neither is true. Packages

significantly limit the number of supportable IOs in the processor due to the large size and pitch of the package-to-board connection relative to the size and pitch of on-chip interconnects ($\sim 10X$ and *not* scaling well). In addition, the packages significantly increase the interconnect distance between the processor die and other dies. Eliminating the package, therefore, has the potential to increase bandwidth by at least an order of magnitude (Section II). Similarly, processor packages are much bigger than the processor itself (*5 to 18 times bigger*). Removing the processor package frees up this area to either be used in form factor reduction or improving performance (through adding more computational or memory resources in the saved area). Lastly, packages limit efficient heat extraction from the processor. Eliminating the processor package can significantly increase the allowable thermal design power (TDP) of the processor (up to 70%). Increase in allowable TDP can be exploited to increase processor performance significantly (through frequency scaling or increasing the amount of computational or memory resources). Unfortunately, simply removing the processor package hurts rather than helps as we point out in Section III. We develop a new silicon interconnect fabric to replace the PCB and make package removal viable in Section IV. Essentially, we place and bond bare silicon dies directly on to a silicon wafer using copper pillar-based I/O pins.

This paper makes the following contributions:

- We make a case for packageless processors. We argue that modern processor packages greatly hinder performance, bandwidth, and energy efficiency scaling. Eliminating packages can enable us to recoup the lost performance, bandwidth, and energy efficiency.
- We present Si-IF, a novel integration technology, as a potential replacement for PCB-based integration and as the enabling technology for packageless processing.
- We quantify the bandwidth, TDP, and area benefits from packageless processing. We show that up to one to two orders of magnitude, 70%, and 5-18x benefits respectively, are possible over conventional packaged processors. These benefits translate into up to 58% (16% average), 136% (103% average), and 295% (80% average) performance benefits, respectively, for our benchmarks.
- We also extend the concept of packageless processing to the entire system on the board; reduction in system-level footprint was up to 76%.

II. PACKAGING PROCESSORS AND ITS LIMITATIONS

Traditionally, processor and memory dies are packaged and then placed on printed circuit boards (PCB) alongside other packaged components. The PCB acts as the system level interconnect and also distributes power to the various packages using the board level power distribution network (PDN). The package is the interface to connect the dies to the PCB. A schematic cross-section of a typical packaged processor on a PCB is shown in Figure 3. Packages serve three primary functions:

- *Packages act as a space transformer for I/O pins:* The diameter of chip IOs is relatively small ($50\mu\text{m}$ - $100\mu\text{m}$) [7]. However, the bump sizes required to connect to the PCB often range between at least a few hundred microns to about a millimeter [8], [9], [10]; large bumps are needed due to PCB's high surface warpage. To enable connectivity in spite of the large difference between chip I/O diameter and the required bump size to connect to PCB, packages are needed. Packages are connected to the silicon die using C4 (controlled collapse chip connection) micro bumps, while the package laminate acts as a redistribution layer (RDL) and fans out to a BGA (ball-grid array) [11] or LGA (land-grid array) [12] based I/O with typical pitch of about $\sim 500\mu\text{m}$ - 1 mm. Packages perform the same function even in the scenario where they do not use solder balls, but use sockets with large pins to prevent breakage from manual installation and handling.
- *Packages provide mechanical support to the dies:* Packages provide mechanical rigidity to the silicon dies, protect them from the external environment (moisture and other corrosive agents), and provide a large mechanical structure for handling. Also, the coefficient of thermal expansion (CTE) of the FR4 material used to make the PCB is ~ 15 - $17\text{ ppm}/^\circ\text{C}$, while that of Silicon is about $2.7\text{ ppm}/^\circ\text{C}$. This large mismatch in CTE between the materials leads to large stresses. Packages provide some mechanical stress buffering, and thus help in mitigating the thermal stresses.
- *Easier testability and reparability:* Since test probe technology has not scaled well [13], [14], [6], it has become harder to probe smaller I/O pads on bare dies. The larger IOs on the packages are easier to probe using conventional probing techniques. Also, while dies come in different sizes, they go into standard packages which can then be tested using standard test equipment. Similarly, solder-based joints and pin-based sockets allow for in-field reparability. Solder joints can be simply heated up, melted and taken off while sockets allow plug-n-play.

Historically, the above advantages have been significant enough that most processor systems, excluding some ultra-low-power processors [15], [16], have been package-based. However, packaging processor dies leads to several significant limitations, many of which are becoming worse, even debilitating.

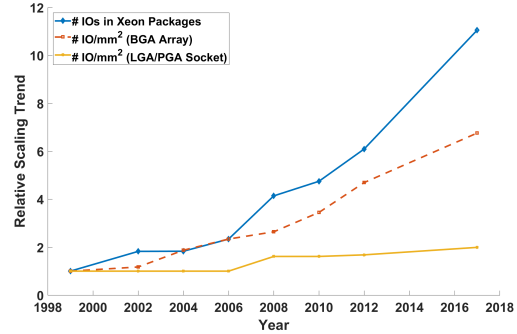


Figure 1: I/O demand is growing faster than the I/O pin density scaling.

- *Packages reduce I/O Density:* Use of packages inherently limits the maximum number of supportable processor IOs because of the large size and pitch of the package-to-board connections (BGA balls/ LGA pins). The BGA/LGA technologies have not scaled well over the past few decades. On the other hand, the demand for IOs in high-performance processor systems is growing rapidly. Figure 1 shows the relative scaling of the number of processor I/O pins in the largest Xeon processor available in a given year against the density scaling (number of IOs/mm²) of the BGA and LGA technologies. As can be seen, the gap between the demand in the number of IOs versus pin density is increasing every year. This widening *I/O gap* limits the amount of power and the number of signals that can be delivered to the processor chip; this can be a severe limitation for future processors that demand high memory and communication bandwidth. Alternatively, processor packages need to become larger; this, however, significantly affects the form factor, complexity, and cost of packages and the length of inter-package connections. In both these cases, the overheads may become prohibitive in near future [6], [17].
- *Packages increase interconnect length:* Increasing the size of the package (the package to die ratio is often >5 , even up to 18 in some cases – (Table I)) leads to a significant increase in the interconnect length between two dies inside separate packages. This is because the die to die connection now needs to traverse the C4 micro-bumps, package RDL, BGA balls, and PCB traces. As the interconnect links become longer, they become more noisy and lossy, which then affects link latency, bandwidth and energy. This problem is aggravated by the fact that a fraction of the interconnect now uses wire traces on PCBs, which are 10X-1000X coarser than the widest wire at the global interconnect layer in SoC chips. Figure 2 compares the energy, latency, and aggregate bandwidth of package-to-package communication links through PCB vs global routing level interconnect wire (Mx4) in an SoC. As seen from the figure, both energy and latency are disparately high for off-package links as compared to the on-die interconnects, while bandwidth is severely limited - these gaps between off-package

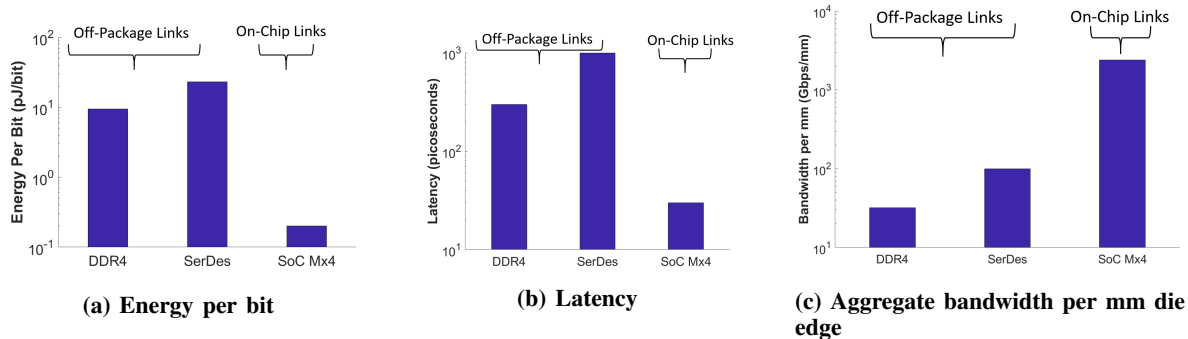


Figure 2: Comparison of communication link energy, latency and bandwidth for on-chip versus off-package links

- links and on-die interconnects must be bridged to enable continued performance scaling.
- *Packages trap heat:* A package traps heat generated by the processor and thus adds to the thermal resistance between the processor die and the heat sink. Figure 3 shows the thermal resistance model of a packaged processor system. In such systems, heat conductively flows upward from the processor die through the package lid and thermal interface materials (TIMs) to the heat sink. The typical thermal resistance values for a canonical 100-130W processor are shown in Figure 3. Thus, for every 10 W of dissipated power, the package lid adds about 1°C to the chip junction temperature. For high-performance processors with TDP ratings in excess of 100 W, the effect of package thermal resistance can cause major reliability issues due to high chip junction temperatures; this limits the TDP, and, therefore, performance of a processor. Moreover, the downward flow of heat encounters high thermal resistivity from the package laminate and the PCB. In fact, the downward heat flow path has about 7-8x higher thermal resistivity than the upward flow. This further exacerbates the above reliability problems from high package thermal resistance. Disruptive solutions that reduce the overall thermal resistance are needed to allow higher sustainable TDP, and, therefore, higher performance at reliable chip-junction operating temperature.
 - *Packages increase system footprint:* As mentioned earlier, package-to-die size ratio has been increasing to accommodate the high I/O demands of today’s processors. Some examples of die-to-package ratio in commercially available processors are shown in Table I. Thus, the overall package footprint is much larger than that of the processor die. Also, since the interconnect width and length are relatively large on PCBs, the total interconnect area is a significant portion of the overall PCB area (see Figure 14a). As I/O demands increase, an increasing amount of system footprint would be taken up by packages, interconnects and passives. Disruptive solutions may be needed to reduce the area cost of these non-compute components to meet the computation density demands of future applications.

Though packages have been an integral part of computing systems for decades, they are becoming a bottleneck for system and performance scaling due to the reasons above.

Table I: Package-to-Die Size ratio

Product Name	Package-to-Die Size Ratio
Intel Knight’s Landing [18]	7
Intel Broadwell [19]	7 - 10
Intel Atom Processor [20]	5 - 18
DRAM Package [21]	2.5 - 3.6

In this work, we rethink the value of packages for today’s and emerging processors, and ask the question - should we build future processor systems without packages?

III. WHY NOT SIMPLY REMOVE THE PROCESSOR PACKAGE?

While some ultra-low-power processors with a small number of I/O pins can be directly mounted on a PCB without packaging [15], [16], it is difficult to do so for high power, high performance processor systems without prohibitive performance and reliability costs. Simply mounting bare die on PCB will dramatically reduce I/O availability proportionately to die-to-package size ratio (e.g., see Table II for some commercial processor examples) as the PCB I/O size is still limited to 500μm (usually much larger). Further, the large CTE mismatch between silicon die and organic PCB can become a reliability bottleneck causing thermal stress-induced I/O failures.

Table II: Analysis of board level I/O availability

Product Name	# Package IOs	Die Area (mm ²)	# Die Balls	Enough Area for I/O?
Knight’s Landing [18]	3647	682	2728	No
Xeon E5-2670 [25]	2011	306	1224	No
Atom N280 [20]	437	26	104	No

In order to realize a packageless processor and its benefits, one would need to replace PCB-based integration with a new integration technology that offers high density interconnect and mechanical robustness.

In the next section, we will describe a novel integration technology (and the accompanying interconnect) we have developed that has the above properties and that can enable packageless processor systems.

IV. SILICON INTERCONNECT FABRIC: AN ENABLING TECHNOLOGY FOR PACKAGELESS PROCESSING

We have developed a novel system integration technology, *Silicon Interconnect Fabric (Si-IF)*, that realizes large scale

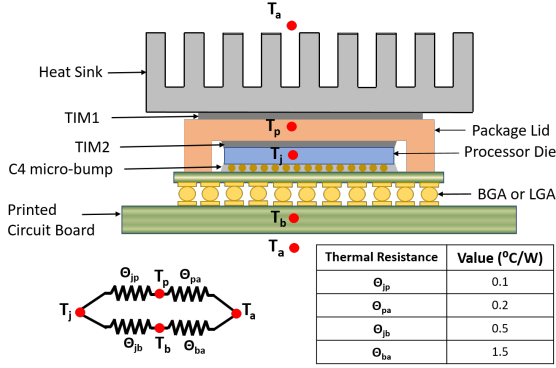


Figure 3: Cross-section of a packaged die with heat sink placed on a PCB is shown, alongside the thermal resistance model. T_a , T_p , T_j , T_b denotes the ambient, package lid, chip-junction, PCB temperature respectively. The thermal resistance values for a typical processor package is shown alongside. [22], [23]

die to wafer bonding technology with very fine pitch interconnection and reduced inter die spacing. The key idea behind Si-IF is to replace the organic PCB board with a silicon substrate. Essentially, we place and bond bare silicon dies directly on to a thick silicon wafer using copper pillar based I/O pins. Processor dies, memory dies, non-compute dies such as peripherals, VRM, and even passive elements such as inductors and capacitors can be bonded directly to the Si-IF. This allows us to completely get rid of the packages. A schematic cross-section of a processor die on Si-IF is shown in Figure 4.

Wafer-scale system manufacturing for building large high-performance computers had been proposed as far back as the 1980s [26], but yield issues doomed those projects, which attempted to make large wafer-scale monolithic chips. Here, the approach is to make small dies with good yield and connect them on a wafer with simple and mature fabrication technology.

Although at a first glance Si-IF technology seems similar to interposers, it is fundamentally different. Interposers use through-silicon-vias (TSV) and because of aspect ratio limitation of TSVs, the interposer needs to be thinned and thus it becomes fragile and size limited. In fact, interposers are typically limited to the maximum mask field size (e.g., $\sim 830 \text{ mm}^2$ which is the same as maximum SoC size) to avoid stitching. Though larger interposers can be built using stitching, they are much costlier and have lower yield. Also, interposers need packages for mechanical support and for space transformation to accommodate larger I/O connections to the PCB. Therefore, connections with chips outside of the interposers continues to suffer from the issues of conventional packaging. On the other hand, Si-IF is a standalone rigid interconnect substrate capable of scaling up to a full size of a wafer and doesn't require packages for mechanical support.

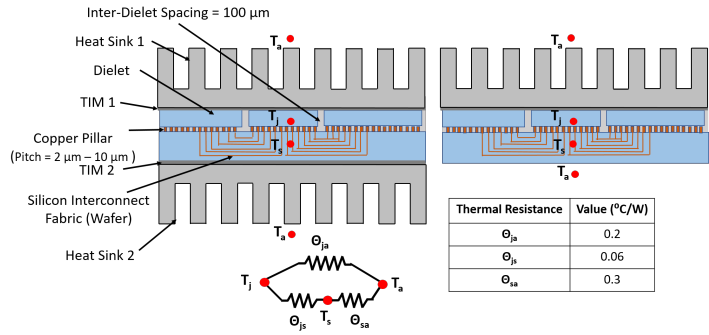


Figure 4: Cross-section of an Si-IF system, alongside the thermal resistance model. T_a , T_j , T_s denote the ambient, chip-junction and silicon substrate temperatures respectively. Heat sink can be directly attached to the top of the dielets or both at the top and bottom of the Si-IF. The thermal resistance values for a typical system on Si-IF is shown alongside [24], [22], [23]

Next, we discuss the distinguishing characteristics of the Si-IF technology in more detail:

Fine pitch inter-die interconnect with 2 - 10 μm pitch: Solder extrusion and surface warpage limit the minimum I/O bump pitch on PCBs. Rigid (polish-able) silicon wafer and copper pillar based IOs (bonded using thermal-compression bonding (TCB) at tight pitches) in Si-IF address both these limitations.¹

Since the interconnect wires on Si-IF are manufactured using standard back-end process, the wire pitch can scale like normal top-level metal in SoCs and well below 2 μm [27], [28]. This technology thus bridges the gap between the SoC level interconnects and system-level interconnects and allows a processor die to support the required number of I/O and power pins even without a package.

Small inter-die spacing: Using state-of-the-art pick and place tools, bare die can be placed and bonded on to the Si-IF at very close proximity ($< 100 \mu\text{m}$) [27]. Thus, interconnects between the dies can now be orders of magnitude smaller than the case where the dies are placed inside separate packages. Coupled with fine pitch interconnects, SerDes links can now be replaced with parallel interfaces and shorter links, thus resulting in lower latency as well as lower energy per bit. The link latency and bandwidth improvement from near placement of the dies coupled with increased I/O density from the fine pitch enables high bandwidth energy efficient communication even without a package.

Efficient heat dissipation: Unlike PCB and package materials, silicon is a good conductor of heat. Heat sinks can be mounted on both sides of an Si-IF. Figure 4 shows how

¹The copper pillar TCB process involves using a bond interface temperature of $\sim 250\text{-}260^\circ\text{C}$ for 3 seconds. Eutectic solder bonding is also done at $220\text{-}270^\circ\text{C}$ for roughly the same period. Therefore, Si-IF-based integration is not expected to cause any temperature related aging of the chip.

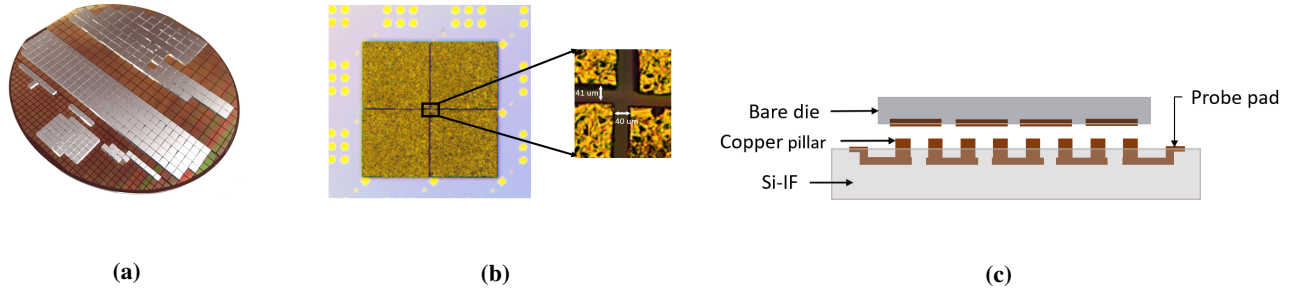


Figure 5: (a) Wafer scale interconnect fabric partially populated with eighty 4 mm², one hundred and seventy one 9 mm², fifty eight 16 mm² and forty one 25 mm² dies bonded on a 4-inch silicon wafer. Copper pillar pitch of 10 μm is used. (b) Micrograph showing four dies bonded on to an Si-IF with ~40 μm inter-die spacing. (c) Serpentine test structure with copper pillars on the Si-IF and landing bond pad on the bare dies [27]

the overall thermal resistance of the Si-IF based system is smaller than that of a canonical packaged and PCB based system. The secondary heat sink attached to the back-side of the Si-IF has the added advantage of acting as a protective shield for the silicon substrate. In fact, the heat sinks would provide mechanical support and protection to the Si-IF instead of a conventional package. To summarize, Si-IF allows much more effective heat dissipation on packageless processors than a conventional packaged processor (more details in Section V-C).

Lowered CTE mismatch: Since both the processor die and the Si-IF are silicon-based, thermal stresses are minimal. As such, the mechanical reliability issues such as bump/ball failures that arise in the conventional setting due to the CTE mismatch between the processor die and the package as well as the package and the PCB are eliminated. Unlike silicon interposers which need to be thin to support TSVs and therefore fragile and size limited [29], [30], Si-IF is thick, rigid and does not use through silicon vias. Therefore, Si-IF-based integration enables large scale processor assembly without requiring the mechanical support traditionally provided by the package.

The above factors coupled with advancements in low-cost silicon processing [31], [32], [6], provide a viable pathway to realizing packageless processors. To demonstrate the feasibility of Si-IF technology for enabling packageless processors, we have built an Si-IF prototype which supports reliable fine-pitch interconnect, high I/O pin density, and close proximity inter-die spacing. Figure 5a shows a 4-inch wafer partially populated using 350 different dies of sizes 4 mm², 9 mm², 16 mm² and 25 mm². A micro-graph of four dies on a wafer spaced apart by only ~40 μm is also shown in Figure 5b. Each of these dies on the wafer has copper pillar pitch of 10 μm and interconnect wires of line-width of 3 μm. This enables high I/O density even without a package. To perform yield analysis of the copper pillars, we built in rows of serpentine test structures in every die as shown in Figure 5c. In each row, pillars n and $n+1$ were connected on the die, while pillars $n+1$ and $n+2$ were connected using the Si-IF interconnect. Once the die was bonded to the Si-IF, the entire row were connected resembling a serpentine

structure. End-points of the serpentes were electrically tested for continuity along a row of the pillars. Out of the 72000 pillar contacts tested, only 3 contact failures were observed. Thus >99.9% yield of the copper pillar connections is observed. This demonstrates the reliability of Si-IF as an enabling technology for packageless processors.

The specific contact resistances were measured to be within 0.7-0.9 Ω - μm² [27] which is smaller than that of the solder balls (40 Ω - μm²) [33], [34], [35]. This is not surprising considering that copper has much higher conductivity compared to solder balls (~5e7 Ω/m vs 9.6e6 Ω/m). Therefore, the contact resistance of a 5 μm copper pillar is about 42 mΩ which is similar to contact resistance of 23 μm C4 solder bumps [33]. Also since, inter-die spacing can now be ~100 μm, instead of the minimum spacing of 1 cm for package-based connections, trace resistance of Si-IF is expected to be much smaller in spite of thinner wires (e.g., assuming similar copper trench depth in PCBs and Si-IF, a 100 μm Si-IF trace will have 8 times lower resistance than the 25 μm, 1 cm length PCB trace). Similarly, relative permittivity of SiO₂ is 3.9, while that of FR4 material is 4.5. Comparing a PCB trace of width and spacing of 25 μm each and length of 1cm with Si-IF trace of width and spacing 2 μm and length 100 μm, the capacitance of the Si-IF trace is about 2 orders of magnitude smaller than that of PCB trace. Thus, RC delay would also be smaller. Using detailed multi-physics and SPICE simulations, we verified that the links can be switched at 2-4 GHz, while consuming <0.3 pJ/bit using very simple I/O drivers [28].

Moreover, the shear bond strength of the Cu pillars was measured to be greater than 78.4 MPa [27], while that of the BGA balls is about 40 MPa [36], [35] which confirms the superior mechanical strength of the copper pillars. Also, due to CTE mismatch of the different components of a package, the solder based bumps go through continuous temperature cycling and often suffer from fatigue related cracking, which would not be a case for Si-IF as the CTE mismatch is negligible.

More details on Si-IF manufacturing (e.g., patterning, die alignment, bonding, etc.) and characterization can be found in [27] and [28].

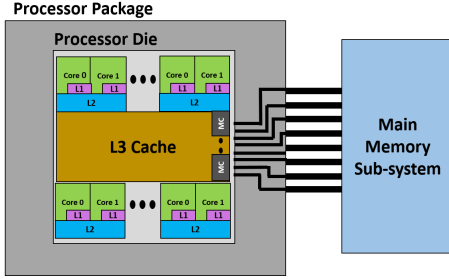


Figure 6: Base Processor Architecture Overview

V. QUANTIFYING MEMORY BANDWIDTH, TDP, AND AREA BENEFITS

In this section, we consider a baseline many-core processor architecture and evaluate the impact on memory bandwidth, TDP, and area if the processor’s package is removed and the processor die is integrated using silicon interconnect fabric.

A. Baseline Processor

Our baseline processor is a 36-tile many-core architecture (with 22 peripheral tiles). Each tile consists of 2 cores and a shared L2 cache of size 1MB. The cores are out-of-order (OOO) with 64 KB private L1 cache. All the 72 cores share a total of 256 MB eDRAM based last-level cache (LLC). The LLC is organized as 8 slices of 32 MB 16-way cache each. Other micro-architectural parameters of the baseline processor are shown in Table III. We assume that the processor is a standalone packaged processor as shown in Figure 6, where the DDR-based main memory is off-package. We use 8 memory channels for off-package DRAM with effective bandwidth of 9.6 GBps per channel. Thus an aggregate of 76.8 GBps of main memory bandwidth is available. The area of the processor die implemented in 22nm technology node is 608 mm² and estimated minimum size of the package required is 2907 mm². Details regarding the methodology to evaluate area, power, and performance are described in Section VI. To estimate the area of the package, we use the model described in Section V-D. Next, we quantify the bandwidth, TDP and area benefits from a packageless implementation of this processor.

Table III: Configuration of the many-core baseline processor

Cores	36 Tiles, each having 2 Silvermont-like OOO at 1.1 GHz, 1 hardware thread, dual issue
Caches	64 KB L1 (Private), 1 MB L2 (Private per Tile), 256 MB eDRAM L3 (Shared)
Memory	DDR4-1600 MHz, Double-pumped at 9.6 GBps, 2D mesh interconnect
Cache Coherence	Directory-Based MESIF
Prefetching	L2, L3 prefetch-on-hit, 4 simultaneous prefetches

B. Memory Bandwidth

As discussed earlier in Section II, packaged processors are I/O pin limited because of the pitch of the solder balls used to connect to the processor. Similarly, memory modules are

I/O pin limited because of large pins used in vertically slotted DIMMs. Coupled with the fact that processor-memory connection uses wide PCB wire traces ($\sim 100 \mu\text{m}$), DDR based communication bandwidth is usually capped at $\sim 10\text{-}15$ GBps per channel. Limited interconnect wiring and pin density also constrains the maximum number of memory channels. Though higher bandwidth can be achieved using complex SerDes techniques, they are energy inefficient ($\sim 10\text{x}$) and lead to additional latency [37], [38], [39]. Some high-end processors [29], [40], [41], [18] use 2.5D technologies such as interposer [42], [43], [44], EMIB [45]. etc. to integrate high bandwidth in-package DRAM memory that can achieve up to 450 GBps of bandwidth. However, the number of memory dies that can be accommodated inside a package is limited due to low yield and high manufacturing cost of larger interposers, EMIBs, etc [45], [46]. Typically, interposers are limited to maximum mask field size ($\sim 830\text{mm}^2$ which is the same as maximum SoC size) to avoid stitching. The largest commercially available interposer is $\sim 1200 \text{mm}^2$ (uses stitching) which only accommodates a processor die with four 3D memory stacks [47]. As a result, the majority of the main memory that is usually placed off-package continues to suffer from limited memory bandwidth.

Since memory chips are connected to the processor chip directly (i.e., without a package and PCB traces) in the Si-IF setting as long as they can fit in the size of a silicon wafer (Table IV), the corresponding supportable memory bandwidth is much higher. As one estimate, the interconnect traces on Si-IF are 2-10 μm in pitch (Section IV) as opposed to $\sim 100 \mu\text{m}$ on PCB, which means about 10-50x more bandwidth is available on Si-IF than on PCB. Moreover, since the link length is expected to be small in Si-IF, signalling can be done at relatively higher frequencies of 4-5 GHz with simple transceivers. The estimated bandwidth per mm edge of a die is ~ 50 GBps and ~ 250 GBps for 10 μm and 2 μm interconnect pitch respectively.

Table IV: Comparison of Si-IF vs other 2.5D technologies

	Silicon Interposer [42]	EMIB [45]	Si-IF
I/O Pitch (μm)	30-100	20-40	2-10
Interconnect Wire Pitch (μm)	2-10	1-10	1-10
Maximum Size/Dies	8.5 cm ²	5-10 Dies	Up to a Full Wafer
Inter-Die Spacing (mm)	>1	>1	<0.1
System Integration	Package on PCB	Package on PCB	Bare Die on Wafer
Other Factors	Complex Assembly Process and TSV Capacitance issue	Complex Manufacturing of Organic Substrate	Bonding Passives and Legacy I/O Ports

C. TDP

Thermal characteristics of a processor system drive many design decisions such as maximum operating frequency, peak power, etc. Since packageless processors allow more effective heat extraction (see Section IV), the allowable TDP

for the same junction temperature constraint increases. To compare the thermal characteristics in PCB based packaged systems against Si-IF based packageless systems, we use the thermal resistance model shown in Figures 3 and 4. Simulations to estimate the thermal resistance of the heat sinks taking into account the air flow, heat spreading effects and size of the heat sink were performed using a commercial thermal modelling software ‘R-Tools’ [48]. We compare different design points such as a conventional package on large PCB vs small PCB, an interposer package on large PCB, a die mounted on large Si-IF vs small Si-IF, and a PCB replaced with Si-IF without removing the package. TDP for the baseline packaged processor is calculated as 0.75 times the processor peak power [49], [23]. We assume a heat sink of the size of processor package, ambient temperature of 25 °C, and forced airflow convection to calculate the junction temperature to be 64.2 °C in this case. We then calculate for each design point the maximum allowable TDP that produces a junction temperature no higher than 64.2 °C. Figure 7 shows the results.

Results show that the TDP benefit from just removing the large PCB and replacing it with an Si-IF is about 6%. Removing the package in case of Si-IF gives an additional ~15% benefit. The surface area of the Si-IF also affects the amount of heat dissipation, which can increase allowable TDP by about 5-7%. In a packageless system with one heat sink and a large Si-IF, maximum TDP that can be allowed for the same junction temperature is 181 W which is 21.5% higher than the baseline case. The benefit increases to 70% (TDP of 254W) when heat sinks are installed on both sides.² Meanwhile, interposer-based 2.5D integration shows no benefit in terms of TDP. In fact, the use of additional interposer layer inside the package lowers the allowable TDP by a small amount due to increased thermal resistivity on the downward heat flow path. The TDP benefits of packageless processing will only increase with increasing die area since the more effective heat spreading on larger dies makes the package resistance a bigger fraction of the overall thermal resistance for packaged systems.

D. Area

Due to the high package area to die area ratio (Table I), removing the package can lead to significant area benefits. To quantify the area benefits for the baseline processor, we use the following model to estimate the minimum size of the package given the peak power of a processor, number of signal IOs (SPins), and type of I/O.

$$Area_{package} = bump_pitch^2 \times \left(\frac{Peak_Power}{Power_perPin} + \#SPins \right) + Area_{Non-I/O} \quad (1)$$

²We expect the cost of placing a single heat sink to be comparable to the packaging cost of baseline system. Since we do not have a package, the second heat sink can be added without increasing cost over the baseline system, while providing significant TDP (and, therefore, performance) benefits in return.

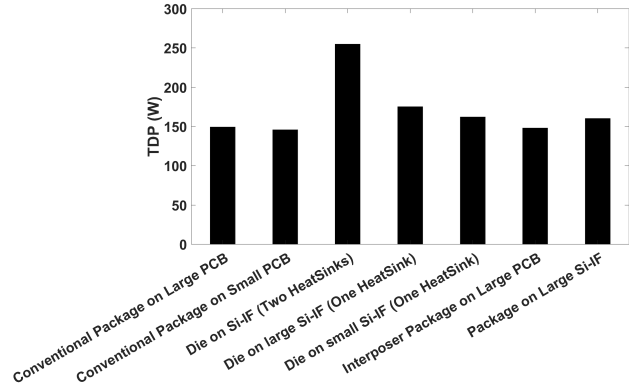


Figure 7: Maximum achievable TDP of the baseline processor system in various integration schemes

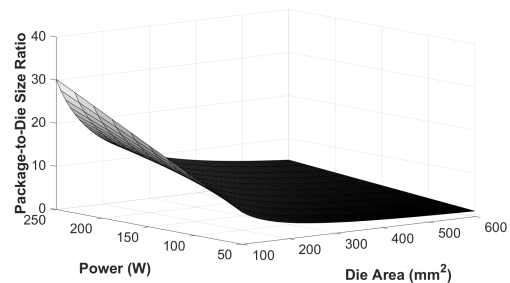


Figure 8: Sensitivity analysis of area benefit from removing the package

Non-I/O area is determined by other factors such as RDL layer and PCB routing constraints. We assume the maximum current per power/ground pin to be 250 mA [50], [51] and bump_pitch to be 900 μm [6], [9]. Using this model for the baseline processor, the minimum package area (when non-I/O area is not considered) is estimated at 2907 mm² which is about 5x larger than the processor die area (608mm²). The area benefit from removing the package will be higher for processors with higher power density (Figure 8) since packages required such processors need to be larger so as to accommodate the power pins and also to dissipate the heat efficiently.

VI. METHODOLOGY

In this section, we describe our methodology for estimating the performance benefits from packageless processing for our baseline processor (Table III).

First, we use McPat 1.0 [52] to determine the area and TDP for the baseline processor. Next, we calculate the additional bandwidth, TDP, and area available from packageless processing (Section V). We then determine a processor design that exploits the additional bandwidth, TDP, and area to improve performance. Since the bandwidth benefits from packageless processing are substantial (orders of magnitude), to eliminate bandwidth slack, we increase the number of memory channels to one per peripheral tile in the baseline processor. To exploit higher allowable TDP, we consider two approaches - increasing core frequency and

increasing the number of tiles in the processor, including adding an additional slice to the eDRAM L3 cache for every 4 additional tiles. Yield concerns limit the number of tiles that can fit in a single die. Therefore, we consider a multi-chip processor system where we limit the size of each die to at most 600mm^2 . Each die contains an even portion of the tiles and is connected in a 2D mesh with the other dies via an inter-processor communication protocol. We use latency of 20 ns [53], [28] and bandwidth of 1 TBps [28] to model the inter-processor communication on Si-IF. We use the same technique when considering area slack.

Once we have determined a set of processor designs, we use a fast multi-core interval simulator, Sniper [54], to determine relative performance. We simulate six benchmarks from the NAS Parallel Benchmark (NPB-3.3) suite [55] and six benchmarks from PARSEC 2.1 [56]. Among the NPB benchmarks, we chose BT (Block Tri-diagonal solver) and SP (Scalar Penta-diagonal solver) as sample pseudo applications, CG (Conjugate Gradient) and UA (Unstructured Adaptive mesh) as having irregular memory access patterns, MG (Multi-Grid) as being memory intensive and EP (Embarrassingly Parallel) as being highly scalable. We used dataset size C, which has an estimated memory requirement of 800 MB. Among PARSEC benchmarks, we chose blackscholes and fluidanimate as sample data-parallel applications, canneal and dedup as having high rates data sharing, and streamcluster and freqmine as typical datamining applications. For all evaluations, the simulation was fast-forwarded to the Region-of-Interest (ROI), simulated in cache-only mode for 1 billion instructions and then simulated in detailed mode for 1 billion further instructions.

VII. RESULTS

In this section, we demonstrate that packageless processors offer significant performance benefits over their packageless counterparts.

A. Exploiting Higher Available Memory Bandwidth

Since Si-IF provides at least 10x more bandwidth than the PCB case alongside plentiful of I/O pins, several techniques such as using wide-I/O interface for the whole memory system and increasing the number of memory channels can be implemented. Though wide-I/O implementation is feasible in interposer-based assemblies as well, number of memory channels is limited since only a few memory devices can be placed on the interposer (due to maximum size / yield limitation - Section V) as opposed to the Si-IF case, where many more memory devices can be accommodated (limited only by the size of the silicon wafer).

Our baseline processor contains 22 peripheral tiles, so we used a maximum of 22 memory channels for the packageless case in our evaluations. Figure 9 shows the potential improvement in performance from having one memory channel per two peripheral tiles (107.8GBps) and one memory channel per peripheral tile (215.6 GBps) over the eight memory channels in our baseline processor configuration

(78.4 GBps). We also compared the performance of all three of these configurations against the maximum achievable performance for a 10 TBps memory bandwidth along the peripheral tiles – this bandwidth is achievable on Si-IF using HMC like memory which supports up to 480 GBps per device [57]. We denote this as the *infinite bandwidth* case.

Increasing number of channels results in average improvement of about 15% with a large L3, while it has a much greater effect (23%, on average) in the absence of an L3. For applications such as BT, MG and SP, the improvement in performance is $>42\%$ both with 22 memory controllers as well as infinite bandwidth when L3 is present. Even without the L3, the performances of BT and SP in the 22 memory controller-case are 31% and 22 % higher respectively than the baseline case with L3. In fact, with 22 channels, but without L3, the average performance across all benchmarks is 8% higher than baseline case with L3. This is because the memory bandwidth effectively improves enough to eliminate the need for an L3. A less intuitive result is that for benchmarks such as CG and Canneal, removing an L3 results in *higher* performance. This is due to limited sharing and irregular memory access patterns in these benchmarks; an L3 increases memory latency unnecessarily in the case that data is used by one core and never shared.

Area overhead of the additional memory controllers would result in increased area of the processor chip. We estimated the area overhead per memory controller in 22nm technology to be about 1.8mm^2 [58]. This implies that the new processor chip size (with additional memory controllers) can exceed 650mm^2 which would worsen the die yield. This issue can be tackled using two ways. First, since Si-IF provides similar density and performance as that of global interconnects in SOCs, we can now have separate memory controller dies which contain clusters of memory controllers. The alternative approach would be to reduce the size of the LLC to accommodate the additional memory controllers. Since performance without LLC, but with additional bandwidth, is similar or higher than the baseline case with LLC in our evaluations, overall performance is expected to improve.

In summary, larger number of memory channels in packageless processors improves performance by up to 58% (average 16%) and 53% (average 14%) in case of infinite bandwidth and 22 memory channels respectively, and allows elimination of the LLC with in fact 8% higher performance with 22 memory channels than the baseline case with LLC.

B. Exploiting Higher Available TDP Budget

As mentioned in Section V-C, additional power can now be sustained without increasing the core junction temperature. Thus, we can either add more cores to the system or increase the frequency of operation (Table V). In Figure 10, we show the performance improvement of these two different design choices across different benchmarks. Frequency scaling alone provides consistent gains of $>15\%$ in performance across all benchmarks when only one heat sink is used. The performance boost is $>50\%$ when both the heat sinks

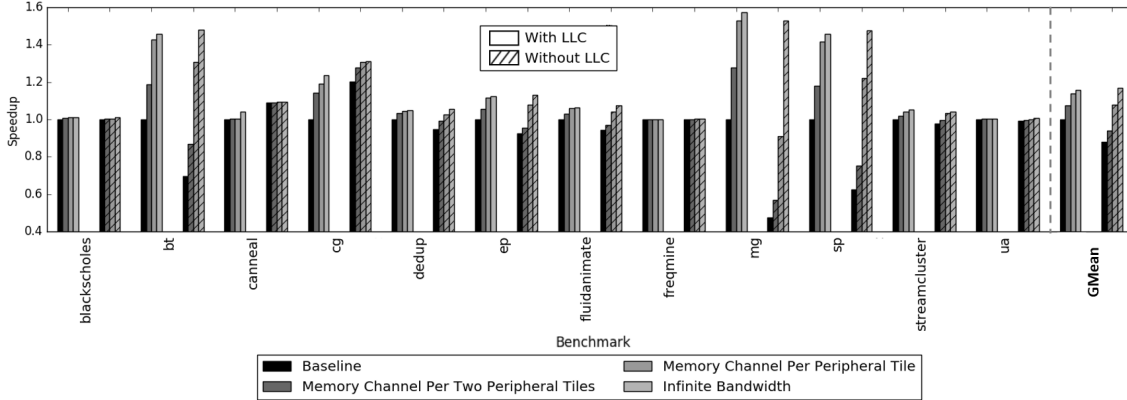


Figure 9: Performance benefit from increased number of memory channels with and without L3

are used. Using DVFS could result in substantially higher speedups, as strategically increasing the frequency of only certain cores would take more precise advantage of the increased TDP. Furthermore, increasing the number of tiles has potential for greater speedup for certain applications. For example, EP achieves more than 2.5x improvement in performance versus 2.2x when frequency is scaled. However, increasing the number of cores requires substantially more area - the largest processor in this experiment exceeded 1600mm² total area (recall that we use an multi-chip processor configuration for large area cases - each chip is still only 600mm² big). Additionally, some applications do not have enough exploitable TLP to fully take advantage of the increased number of cores. Freqmine and streamcluster are two examples of benchmarks which achieve substantial gains by scaling frequency, but do not gain performance from adding further tiles.

Table V: Increasing Frequency or Number of Tiles to Exploit Available TDP Slack

System Configuration	TDP	Max Frequency	Max # Tiles
Baseline	149 W	1.1 GHz	36 Tiles
Small Si-IF Heatsink 1-Side	168 W	1.4 GHz	48 Tiles
Large Si-IF Heatsink 1-Side	180 W	1.6 GHz	52 Tiles
Large Si-IF Heatsink 2-Side	250 W	2.6 GHz	96 Tiles

One more efficient way to take advantage of thermal slack would be to perform a two dimensional design space exploration on the chip, scaling both frequency and number of tiles until an optimal system is found. In general, frequency has a clearer and more well-defined trade-off between power and performance. In addition, through DVFS it is easier to manipulate frequency during runtime and be able to optimize the processor for a specific application. While one could dynamically change the effective number of tiles available in a processor via power gating, there is a much higher overhead for such a transition, including wakeup time, cache warmup and various OS overheads associated with context switching. However, due to bandwidth constraints and the benefits of having a larger total cache area, increasing the number of tiles provides for running massively parallel

workloads much more efficiently than a smaller number of highly clocked processor.

Increasing frequency or the number of tiles would increase the power demand. Besides thermal constraints, increased power consumption also requires careful management of power distribution losses (for example by point of use step down voltage conversion just like conventional packaged systems). Packageless Si-IF with no C4 bumps or wide PCB traces can substantially help with inductive voltage drops. High power requirements also come with larger demand for power/ ground I/O pins which can be accommodated within the die area using fine pitch interconnect pillars on Si-IF.

In summary, removing the package improves the TDP budget by up to 70% which can provide upto 136% higher performance (average 103%) upon increasing the operating frequency and up to 162% (average 60%) upon increasing the number of tiles using our benchmarks.

C. Exploiting Higher Available Area

Table VI: Area Slack Exploitation Parameters

System Configuration	Max Area	Processor Microarchitecture
Baseline	608mm ²	36 Tiles, Single Die
Packageless Half-Slack	1758mm ²	96 Tiles, Four Dies
Packageless No-Slack	2908mm ²	144 Tiles, Six Dies

Figure 11 shows performance benefits for eliminating half and all of the area slack available in a packageless processor. For our evaluations, dies are restricted to 600mm² - see Section VI for details. One might not want to fully exploit available area for many reasons: higher power and lower yield being among the chief concerns. Much like the case of tile-based power slack elimination, some applications benefit drastically more than others from area slack reduction. For applications such as fluidanimate, nearly all of the performance benefits, i.e., ~86% over the baseline case, are achieved via half-slack reduction, while other applications can continue to take advantage of any extra cores available. As in Section VII-B, benchmarks such as blackscholes, EP and UA increase performance proportionally to number of

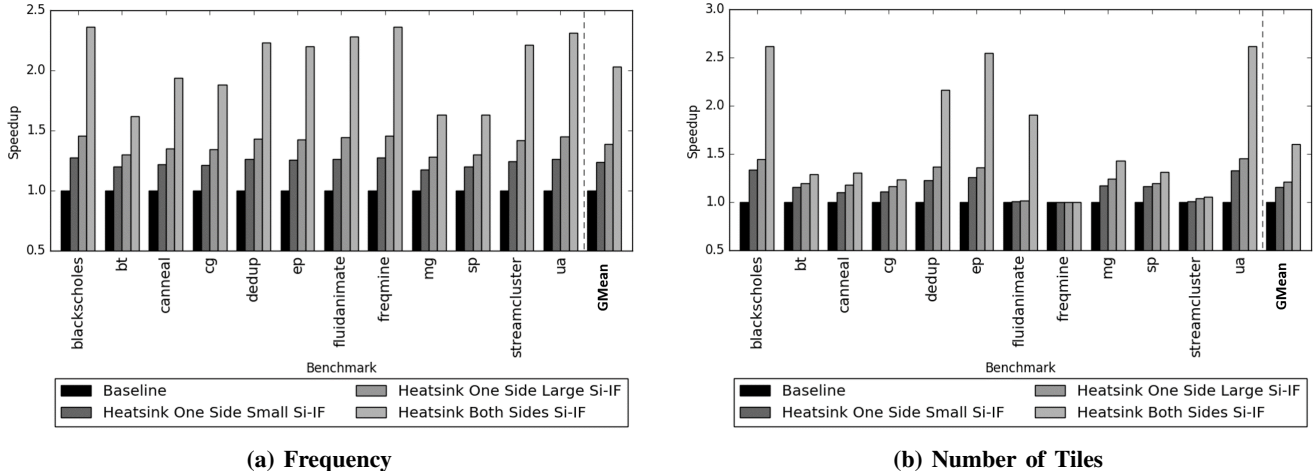


Figure 10: Performance benefits of utilizing TDP slack

tiles, due to high thread level parallelism (TLP). For applications which lack such easily exploitable TLP, having a large number of cores may still be useful in the case of multi-programming. The power overhead of such a large design can be mitigated using per-core DVFS or power-gating. The removal of a package allows for systems with much denser compute: for compute-intensive high performance systems which require thousands of cores, packageless processors could prove to be a critical technology.

In summary, across the benchmarks evaluated, packageless processors could achieve 80% average performance improvement, up to 295% by utilizing the extra area slack coming from removing the processor package.

Note that we are allowing the original TDP budget to be breached in these experiments; we assume that a costlier cooling solution exists to tackle the increased thermal dissipation if the intention is to use the entire area slack. Section VII-D considers this tradeoff between area slack and TDP slack.

D. Area-TDP Tradeoff

Thus far, we have quantified the bandwidth, TDP and area benefits individually that packageless processors provide over conventional systems. In this section, we ask the question - how much improvement in form factor, TDP and bandwidth can be achieved when the factors are considered simultaneously?

For our evaluations, we consider two PCB-based baselines - one DIMM (with 18 chips) per channel and one 3D stacked memory device per channel. The corresponding Si-IF design points have 18 packageless DRAM chips per channel laid out in a planar configuration (Figure 14b) and one packageless 3D stacked memory device per channel respectively. The processor footprint is 2907 mm² (Section V-D) in the packaged case, while it is 608 mm², on Si-IF. For estimating memory footprint in the PCB case, we assume that the DIMMs are slotted vertically onto the PCB (Figure 14a). The PCB footprint for each DIMM is estimated to be 7.92 cm² (we used the DIMM socket size as the footprint estimate

to perform a worst case comparison of area benefits from Si-IF, ignoring large inter-socket distances typically used on a PCB). The PCB footprint for each 3D stacked memory package is considered to be 320 mm² [59]. For estimating the memory subsystem footprint in two packageless cases, we considered 36 mm² per DRAM die and 55 mm² per 3D stacked memory device [60].

(1) **Form Factor Reduction in iso-TDP case:** In the iso-TDP case, we compare the footprint of the baseline 8 memory channel configuration on PCB against the same system implemented on Si-IF. We also extend the analysis for both the 22 channel (one channel per peripheral tile) and 11 channel (one channel per two peripheral tiles) memory configurations on Si-IF and compare it against the same baseline PCB case.

We adjust the size of the heat sink so as to achieve the same maximum junction temperature as the baseline junction temperature of 64.2 °C (junction temperature of the in-package baseline processor die). In case the heat sink required to achieve the desired junction temperature is larger than the total processor and memory footprint, the area of the heat sink determines the compute footprint of Si-IF system.

Figure 12 shows the area savings in different scenarios. For one memory device per channel case, the dual heat sink setup leads to area savings of up to 76 % and using one heat sink provides >36% area reduction. This is because in the dual heat sink setup, the thermal resistance is lower, meaning smaller heat sinks can help achieve higher TDP.

For DDRx style memory configuration with 18 dies per DIMM, when the baseline 8 memory channel configuration is laid out on the Si-IF with memory bandwidth similar to PCB case, 37% area savings can be achieved with similar performance as of the baseline packaged case. The area saving reduces to 17% in the 11 memory channel case while the performance increases by 7.5%. For 22 memory channel case (memory channel per peripheral tile with LLC in Figure 9), where 22 × 18 = 396 memory dies need to be accommodated on planar Si-IF, the footprint does increase but there is plenty of TDP slack left unused as the large heat

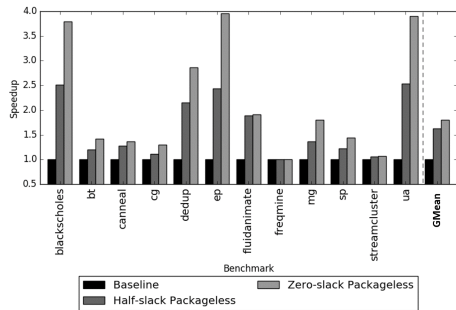


Figure 11: Performance increase by exploiting area slack

sink “overcools” the system.

In summary, packageless processing under a TDP constraint with emerging 3D stacked memories can deliver dramatic footprint reductions (40%-76%) while increasing available memory bandwidth. In conventional DDR-style memory systems, going packageless can deliver 36% footprint reduction with same performance.

(2) **Increased TDP Slack in iso-Area case:** Here, we compare the TDP slack available for the packageless processor system if the total area of the Si-IF and the heat sink are equal to the total PCB footprint of the processor and the memory subsystem. Figure 13 shows the total packageless TDP available as compared to the total TDP of the baseline processor and memory subsystem. Since the area footprint of the DIMM is much larger than that of the 3D stacked memory packages, the equivalent iso-area Si-IF/heat sink size is larger which leads to extra TDP slack. This excess TDP slack alongside the excess area under the heatsink can be utilized by increasing the number of tiles, frequency of operation, memory capacity etc. In summary, packageless processing with the same computing footprint can deliver 1.7X-3X extra power to burn to improve performance without violating thermal constraints.

VIII. DISCUSSION

In this section, we discuss the implications of packageless processing on how the overall system could be realized, and other aspects such as repairability, testability, manufacturability, and cost. We also discuss some architectural implications not covered in this paper.

A. Overall System Architecture

A full system implementation comprises of core compute elements such as CPUs, GPUs, memory, etc., and non-compute elements such as crystal oscillators, driver ICs for system I/Os, components of power delivery network, etc. So far, we have only discussed the compute elements of the system, however architecture of non-compute components is important as well.

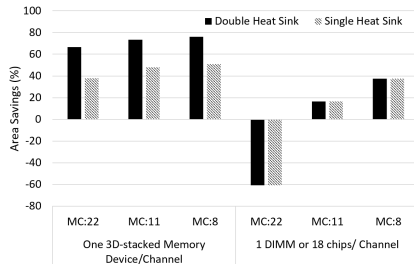


Figure 12: Area savings when implementing processor-memory subsystem on Si-IF. 22, 11 and 8 memory controller configuration on Si-IF are compared against baseline packaged configuration with 8 channels of off-package DRAM

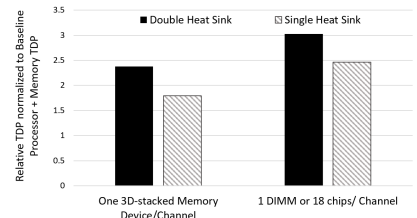


Figure 13: TDP of implementing the baseline processor-memory subsystem on Si-IF of size of total processor and memory package area normalized to baseline processor system TDP

Traditionally, surface mount non-compute components are soldered directly on the PCB (Figure 14a). In Si-IF, we envision two alternatives to integrate these components into the system. One is to bond the passives and other non-compute components directly onto the silicon board using solder balls and large pads on Si-IF, as shown in Figure 14b. We have been able to achieve bonding of passives on to the Si-IF successfully. This enables a full system integration on Si-IF. The other alternative is a hybrid approach shown in Figure 14c, where the compute components alongside some Si-IF compatible non-compute components can be integrated on to the Si-IF, and other remaining non-compute components can be integrated on a separate daughter board. An ancillary benefit of a daughter board approach is that the daughter board can now also host some upgradeable and spare components such as extra spare DIMMs alongside legacy connectors.

We estimated the footprint of a $\sim 1000 \text{ cm}^2$ Intel Xeon dual socket motherboard [61] in Si-IF setting. Considering non-compute footprint reduction by 50% when non-compute is fully implemented on Si-IF (due to denser integration of all components on Si-IF) alongside packageless implementation of memory and processor dies, a full Si-IF implementation footprint can be $<400 \text{ cm}^2$ while the hybrid approach can be $<780 \text{ cm}^2$.³

B. Test, Reliability, and Serviceability

Bare dies are difficult to probe because of the small size of the I/O pads. However, significant progress has been made in bare die testing techniques, primarily driven by need for known good die in 2.5D and 3D IC technologies [62], [63]. Some examples include temporary packages [64], [65],

³Link lengths are often lower in the packageless systems case since inter-die spacing can now be reduced to $\sim 100 \mu\text{m}$. We estimated that the farthest DDR links are about $\sim 2x$ longer on standard Xeon PCBs, when vertically slotted DIMMs are used, compared to when bare dies are placed in a planar fashion on Si-IF. So, DDR-type signalling and routing will not be an issue. In very large Si-IF systems where signal integrity may be an issue, we can use intermediate buffer dies/chiplets to buffer the signals if simpler signalling is used (for lower power).

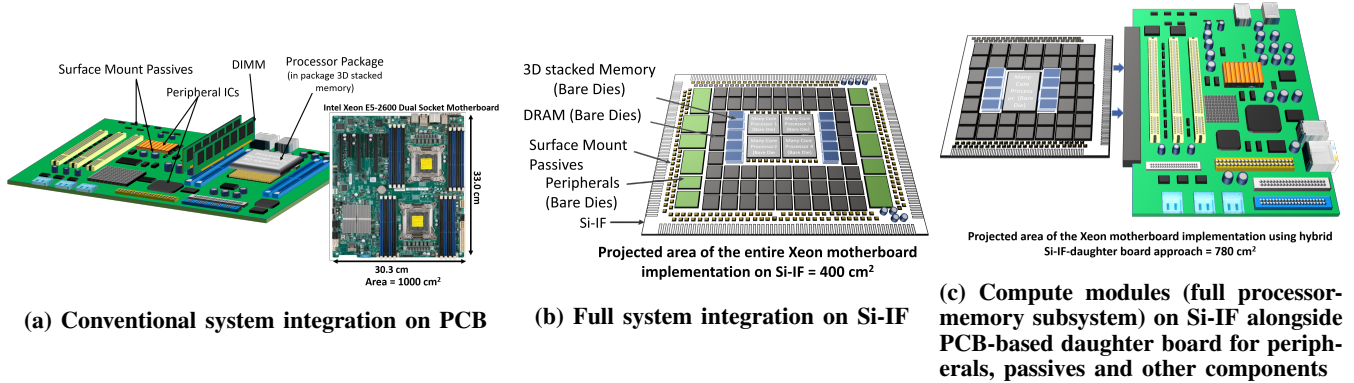


Figure 14: Illustration of a conventional PCB based system and different integration schemes using Si-IF

wafer level burn in and test [66], [67], die-level test [68], and built-in self testing mechanisms [69], [70].

A packageless processor may be more susceptible to environmental agents (radiation, moisture, etc.) than its packaged counterpart. Layer of radiation hardening material (e.g., SiC/H, Boron-10, etc.) can be CVD deposited to protect against radiation. Also, in many cases, package itself is the source of radiation which in the packageless case is omitted. Similarly, the IF-assembly can be passivated with a CVD-based coating, which protects it from moisture and salt intrusion. Furthermore, we apply a hermetic sealant around the edges of the dies to prevent environment agents to get beneath the dies and corrode the copper pillars. External heatsink(s) provide additional environmental protection when used. Finally, the chip-to-wafer bonders have necessary (e.g. for ESD) protections to avoid any charge accumulation on the chip as well as the Si-IF.

While soldered or socketed components can be replaced in conventional PCB based integration schemes, replace, rework, or upgrading components is relatively difficult for Si-IF based systems since de-bonding metal-metal joints is a complex process which requires high temperature to melt the bond joint. As such, the benefits of Si-IF must be weighed against serviceability concerns (MCMs, MDPs, 3D integration, etc., also provide improve performance and energy efficiency at the expense of serviceability).

Serviceability of Si-IF-based systems can be improved through redundancy and self-repair. While these solutions incur additional costs, the considerable cost reduction from the proposed approach should defray these costs in many applications. Also, redundancy/self-repair costs can be reduced. For example, if a specific component is prone to frequent failure, it may be soldered or socketed, instead of TC-bonded, to improve serviceability (we already have a mix of solder/socket and copper pillar TCB on some of our prototypes). While this reduces I/O count, area and TDP benefits remain. Even the I/O costs can be minimized. As one example, DRAM chips have low I/O density and are prone to faults; they can use conventional solder bumping or placed in soldered sockets, while processor chiplets can be TC bonded (using copper pillars) to support large I/O count.

C. System Level I/O Connections and Mechanical Installation

External I/O connections would be made at the edge of the Si-IF to allow the rest of the surface to be covered using the heat sink. Conventional plug connectors or solder based connections can be used for signal and power delivery to the Si-IF. Silicon is much more robust than FR4 used to build PCBs (compressive strength of 3.2-3.4 GPa vs 370-400 MPa) and can easily handle the normal insertion force of a plug connector (few MPa to a few 10s of MPa), especially with backside support (e.g., backside heat sink) - our 700 μm -thick prototype kept flat on a chuck was intact even when a compressive stress of 1.5 GPa was applied over 0.13 mm^2 . Even with minimal backside support, silicon is much more robust than the PCB (Ultimate Tensile Strength (UTS) of 165-180 MPa vs 70-75 MPa).

There are several options for installation. In case of server chassis, the complete system-on-wafer can be inserted using low force insertion sockets. Alternatively, in implementations with external metal heatsink(s), the heatsink(s) can be bolted to the chassis. If heatsinks are not required on both sides, backside heatsink will be preferred to provide support. The other side can be optionally covered using a robust material, e.g., metal plate. In case of cellphones, Si-IF can be held with mechanical jaws or can be fixed using a thermally conductive glue.

D. Manufacturing Challenges and Cost

The Si-IF integration required to enable effective packageless processing relies on metal-metal thermal compression bonding of copper. After the initial TCB process of 3 sec bonding at $\sim 250^\circ\text{C}$ interface temperature, batches of bonded wafers undergo thermal annealing for about 6-8 min at $\sim 150^\circ\text{C}$ to enhance bond strength and reduce tail probability of bond failures [27] - potentially decreasing the throughput of the manufacturing process. Maskless lithography is used to pattern large area, fine-pitch interconnect on Si-IF which can also have throughput concerns. Further improvements in large area patterning may be needed for volume production.

Removing the package has significant cost benefits since for many processors, packaging costs are often about 30-50% of the total processor cost [71], [72]. Also, the sig-

nificant area reduction from packageless processing should lower costs even further. As an example, the baseline 8 memory channel, 3D memory system will have area of 1048 mm² (608 + 8*55) in Si-IF and 5467 mm² (320*8+2907) on packaged PCB. A processed silicon wafer with a 90nm global layer back-end (enough to sustain 2 μm pitch) is roughly \$500 per 300 mm wafer. Moreover, the die-to-Si-IF bonding is performed using industry standard die-to-substrate bond tools with small upgrades. Assembly cost per system is therefore expected to be around \$15. For packaged systems, just the cost of packages is roughly \$44 (3*8 + 20) per system [73]. Similarly, since wire pitches in Si-IF are several microns wide (2-10 μm), Si-IF fabrication is performed using older technology node (90nm/180nm) processes that support these wire pitches. As such, the fabrication cost is low. High performance multi-layer PCBs often cost a few hundred dollars while having much lower compute density than that of Si-IF. Finally, since Si-IF provides large form factor benefits, performance density per volume goes up. This has the potential to decrease the overall total cost of ownership [74].

E. Other Architectural Implications and Use-case Scenarios

In addition to the architectural techniques explored in this paper to exploit the benefits of packageless processors, there exist several other micro-architectural optimizations that may be used in the context of packageless processors. For example, aggressive prefetching techniques [75] can leverage the availability of ultra high bandwidth. Similarly, architectures without L3 may be promising for applications where the reduction in L3 miss penalty can offset the effect of L3 miss rate. Also, TDP and area benefits can be utilized by introducing heterogeneous computing, such as GPUs, accelerators, DSP modules, etc. Moreover, since interconnect links are shorter in Si-IF, *Ldi/dt* noise would be smaller. Not only does this potentially reduce the number of decoupling capacitors required on the chip (or inside the package) thereby reducing chip area (or making it available for additional features), inductive noise driven constraints on frequency and timing of power gating, DVFS [76] etc can also now be relaxed. Finally, it may be possible to build wafer-scale systems using the Si-IF integration technology - such systems, in turn, may enable large neural network accelerators, GPUs, and microdatacenters.

IX. SUMMARY AND CONCLUSIONS

Processor packages can significantly impact the bandwidth, allowable TDP, and area taken up by a processor. We proposed packageless processors - processors where the packages are removed and PCB-based integration is replaced by a Silicon Interconnection Fabric, a novel interconnection technology that involves mounting dies directly on a silicon wafer using copper pillar-based I/O pins. We showed that packageless processors can have one to two orders of magnitude higher memory bandwidth, up to 70% higher allowable TDP, and 5X-18X lower area than conventional packaged processors. These benefits can be exploited to

increase processor performance. For a set of NAS and PAR-SEC benchmarks, we showed performance improvements up to 58% (16% average), 136% (103% average), and 295% (80% average) resulting from improved memory bandwidth, processor TDP and processor footprint respectively. For the same performance, packageless processing reduces compute subsystem footprint by up to 76% or equivalently increases TDP by up to 2X. The benefits from packageless processing should only increase with increasing I/O and performance demands of emerging applications and processors.

X. ACKNOWLEDGEMENT

This work was supported in part by the Defense Advanced Research Projects Agency (DARPA) through ONR grant N00014-16-1-263 and the UCLA CHIPS Consortium. The authors would like to thank SivaChandra Jangam for helping with the Si-IF prototype in the paper, and Irina Alam, Matthew Tomei, and the anonymous reviewers for their helpful feedback and suggestions. The authors would like to also thank the support from UCOP through grant MRP-17-454999.

REFERENCES

- [1] M. Bohr, "A 30 Year Retrospective on Dennard's MOSFET Scaling Paper," *IEEE Solid-State Circuits Society Newsletter*, vol. 12, pp. 11–13, Winter 2007.
- [2] K. Atasu, L. Pozzi, and P. Ienne, "Automatic Application-specific Instruction-set Extensions Under Microarchitectural Constraints," in *40th Annual Design Automation Conference*, (New York, NY, USA), pp. 256–261, ACM, 2003.
- [3] S. Venkataramani, S. T. Chakradhar, K. Roy, and A. Raghunathan, "Approximate Computing and the Quest for Computing Efficiency," in *Proceedings of the 52d Annual Design Automation Conference*, DAC '15, (New York, NY, USA), pp. 120:1–120:6, ACM, 2015.
- [4] G. H. Loh, "3D-Stacked Memory Architectures for Multi-core Processors," in *35th Annual International Symposium on Computer Architecture (ISCA)*, (Washington, DC, USA), pp. 453–464, 2008.
- [5] N. Z. Haron, S. Hamdioui, and S. Cotofana, "Emerging non-CMOS nanoelectronic devices - What are they?," in *4th IEEE International Conference on Nano/Micro Engineered and Molecular Systems*, pp. 63–68, Jan 2009.
- [6] S. S. Iyer, "Heterogeneous Integration for Performance and Scaling," *IEEE Transactions on Components, Packaging and Manufacturing Technology*, vol. 6, pp. 973–982, July 2016.
- [7] J. H. Lau, *Flip Chip Technologies*. New York, NY, USA: McGraw-Hill, 1996.
- [8] UG1099, *Recommended Design Rules and Strategies for BGA Devices*. Xilinx Inc, 1 ed., March 2016.
- [9] TE Connectivity Corporation, *LGA 3647 SOCKET AND HARDWARE*, 2017.
- [10] Intel Corp., *Land Grid Array (LGA) Socket and Package Technology*.
- [11] Intel Corp., *Ball Grid Array (BGA) Packaging*, 2000.
- [12] "Land grid array." https://en.wikipedia.org/wiki/Land_grid_array, (accessed July 29, 2017).
- [13] W. R. Mann, F. L. Taber, P. W. Seitzer, and J. J. Broz, "The leading edge of production wafer probe test technology," in *2004 International Conference on Test*, pp. 1168–1195, Oct 2004.
- [14] Y. Liu, S. L. Wright, B. Dang, P. Andry, R. Polastre, and J. Knickerbocker, "Transferrable fine pitch probe technology," in *2014 IEEE 64th Electronic Components and Technology Conference (ECTC)*, pp. 1880–1884, May 2014.
- [15] M. Luthra, "Process challenges and solutions for embedding Chip-On-Board into mainstream SMT assembly," in *Proceedings of the 4th International Symposium on Electronic Materials and Packaging*, pp. 426–433, Dec 2002.
- [16] J. H. Lau, *Chip On Board*. Springer USA, 1994.
- [17] G. V. Clatterbaugh, P. Vichot, and J. Harry K. Charles, "Some Key Issues in Microelectronic Packaging," in *Johns Hopkins APL Technical Digest*, vol. 20, pp. 34–49, Oct 1999.
- [18] Intel, "Intel Xeon Phi," (2014).
- [19] "Broadwell - Microarchitectures - Intel." <https://en.wikichip.org/wiki/intel/microarchitectures/broadwell>.
- [20] "Atom - Intel." <https://en.wikichip.org/wiki/intel/atom>.
- [21] "Micron DDR4 SDRAM Part Catalog." <https://www.micron.com/products/dram/ddr4-sdram/8Gb/>.
- [22] D. Edwards and H. Nguyen, "Semiconductor and IC Package Thermal Metrics." Application Report, Texas Instruments, 2016, <http://www.ti.com/lit/an/spra953c/spra953c.pdf>, (accessed August 1, 2017).
- [23] Intel Corp., *Intel Pentium 4 Processor in the 423-pin Package Thermal Design Guidelines*, November 2000.

- [24] H. R. Shanks, P. D. Maycock, P. H. Sidles, and G. C. Danielson, "Thermal Conductivity of Silicon from 300 to 1400K," *Phys. Rev.*, vol. 130, pp. 1743–1748, Jun 1963.
- [25] "Intel Xeon Processor E5 Family," <http://ark.intel.com/products/series/59138/Intel-Xeon-Processor-E5-Family>.
- [26] J. F. McDonald, E. H. Rogers, K. Rose, and A. J. Steckl, "The trials of wafer-scale integration: Although major technical problems have been overcome since WSI was first tried in the 1960s, commercial companies can't yet make it fly," *IEEE Spectrum*, vol. 21, pp. 32–39, Oct 1984.
- [27] A. Bajwa, S. Jangam, S. Pal, N. Marathe, T. Bai, T. Fukushima, M. Goorsky, and S. S. Iyer, "Heterogeneous Integration at Fine Pitch ($\leq 10\mu\text{m}$) Using Thermal Compression Bonding," in *IEEE 67th Electronic Components and Technology Conference (ECTC)*, pp. 1276–1284, May 2017.
- [28] S. Jangam, S. Pal, A. Bajwa, S. Pamarti, P. Gupta, and S. S. Iyer, "Latency, Bandwidth and Power Benefits of the SuperCHIPS Integration Scheme," in *IEEE 67th Electronic Components and Technology Conference (ECTC)*, pp. 86–94, May 2017.
- [29] A. Sodani, R. Gramunt, J. Corbal, H. S. Kim, K. Vinod, S. Chinthamani, S. Hutsell, R. Agarwal, and Y. C. Liu, "Knights Landing: Second-Generation Intel Xeon Phi Product," *IEEE Micro*, vol. 36, pp. 34–46, Mar 2016.
- [30] N. E. Jerger, A. Kannan, Z. Li, and G. H. Loh, "NoC Architectures for Silicon Interposer Systems: Why Pay for more Wires when you Can Get them (from your interposer) for Free?," in *47th Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 458–470, Dec 2014.
- [31] Y. Iwata and S. C. Wood, "Effect of fab scale, process diversity and setup on semiconductor wafer processing cost," in *IEEE/SEMI Advanced Semiconductor Manufacturing Conference and Workshop*, pp. 237–244, 2000.
- [32] F. Yazdani, "A novel low cost, high performance and reliable silicon interposer," in *IEEE Custom Integrated Circuits Conference (CICC)*, pp. 1–6, Sept 2015.
- [33] S. L. Wright, R. Polastre, H. Gan, L. P. Buchwalter, R. Horton, P. S. Andry, E. Sprogis, C. Patel, C. Tsang, J. Knickerbocker, J. R. Lloyd, A. Sharma, and M. S. Sri-Jayantha, "Characterization of micro-bump C4 interconnects for Si-carrier SOP applications," in *56th Electronic Components and Technology Conference*, 2006.
- [34] B. Dang, S. L. Wright, P. S. Andry, C. K. Tsang, C. Patel, R. Polastre, R. Horton, K. Sakuma, B. C. Webb, E. Sprogis, G. Zhang, A. Sharma, and J. U. Knickerbocker, "Assembly, Characterization, and Reworkability of Pb-free Ultra-Fine Pitch C4s for System-on-Package," in *2007 Proceedings 57th Electronic Components and Technology Conference*, pp. 42–48, May 2007.
- [35] L. D. Cioccio, P. Gueguen, R. Taibi, T. Signamarcheix, L. Bally, L. Vandroux, M. Zussy, S. Verrun, J. Dechamp, P. Leduc, M. Assous, D. Bouchu, F. de Crecey, L. L. Chapelon, and L. Clavelier, "An innovative die to wafer 3D integration scheme: Die to wafer oxide or copper direct bonding with planarised oxide inter-die filling," in *2009 IEEE International Conference on 3D System Integration*, pp. 1–4, Sept 2009.
- [36] M. Ohyama, M. Nimura, J. Mizuno, S. Shoji, M. Tamura, T. Enomoto, and A. Shigetou, "Hybrid bonding of Cu/Sn microbump and adhesive with silica filler for 3D interconnection of single micron pitch," in *IEEE 65th Electronic Components and Technology Conference (ECTC)*, pp. 325–330, May 2015.
- [37] V. Balan, O. Oluwole, G. Kodani, C. Zhong, R. Dadi, A. Amin, A. Ragab, and M. J. E. Lee, "A 15-22 Gbps Serial Link in 28 nm CMOS With Direct DFE," *IEEE Journal of Solid-State Circuits*, vol. 49, pp. 3104–3115, Dec 2014.
- [38] K. Kaviani, T. Wu, J. Wei, A. Amirkhany, J. Shen, T. J. Chin, C. Thakkar, W. T. Beyene, N. Chan, C. Chen, B. R. Chuang, D. Dressler, V. P. Gadde, M. Hekmat, E. Ho, C. Huang, P. Le, Mahabaleshwara, C. Madden, N. K. Mishra, L. Raghavan, K. Saito, R. Schmitt, D. Secker, X. Shi, S. Fazeel, G. S. Srinivas, S. Zhang, C. Tran, A. Vaidyanath, K. Vyas, M. Jain, K. Y. K. Chang, and X. Yuan, "A Tri-Modal 20-Gbps/Link Differential/DDR3/GDDR5 Memory Interface," *IEEE Journal of Solid-State Circuits*, vol. 47, pp. 926–937, April 2012.
- [39] M. A. Karim, P. D. Franzon, and A. Kumar, "Power comparison of 2D, 3D and 2.5D interconnect solutions and power optimization of interposer interconnects," in *2013 IEEE 63rd Electronic Components and Technology Conference*, pp. 860–866, May 2013.
- [40] AMD, "AMD Radeon R9," (2015).
- [41] NVIDIA, "NVIDIA Updates GPU Roadmap; Announces Pascal," (2015).
- [42] T. G. Lenihan, L. Matthew, and E. J. Vardaman, "Developments in 2.5D: The role of silicon interposers," in *2013 IEEE 15th Electronics Packaging Technology Conference (EPTC 2013)*, pp. 53–55, Dec 2013.
- [43] D. Malta, E. Vick, S. Goodwin, C. Gregory, M. Lueck, A. Huffman, and D. Temple, "Fabrication of TSV-based silicon interposers," in *2010 IEEE International 3D Systems Integration Conference (3DIC)*, pp. 1–6, Nov 2010.
- [44] J. Keech, S. Chaparala, A. Shorey, G. Piech, and S. Pollard, "Fabrication of 3D-IC interposers," in *2013 IEEE 63rd Electronic Components and Technology Conference*, pp. 1829–1833, May 2013.
- [45] R. Mahajan, R. Sankman, N. Patel, D. W. Kim, K. Aygun, Z. Qian, Y. Mekonnen, I. Salama, S. Sharan, D. Iyengar, and D. Mallik, "Embedded Multi-die Interconnect Bridge (EMIB) – A High Density, High Bandwidth Packaging Interconnect," in *2016 IEEE 66th Electronic Components and Technology Conference (ECTC)*, pp. 557–565, May 2016.
- [46] L. Li, P. Chia, P. Ton, M. Nagar, S. Patil, J. Xue, J. Delacruz, M. Voicu, J. Hellings, B. Isaacson, M. Coor, and R. Havens, "3D SiP with Organic Interposer for ASIC and Memory Integration," in *2016 IEEE 66th Electronic Components and Technology Conference (ECTC)*, pp. 1445–1450, May 2016.
- [47] "System Level Co-Optimizations of 2.5D/3D Hybrid Integration for High Performance Computing System," http://www.semiconwest.org/sites/semiconwest.org/files/data15/docs/3_John%20Hu_nVIDIA.pdf, (accessed August 1, 2017).
- [48] "R-Tools 3-D Heat Sink thermal Modelling," <http://www.r-tools.com/>.
- [49] V. Kontorinis, A. Shayan, D. M. Tullsen, and R. Kumar, "Reducing peak power with a table-driven adaptive processor core," in *Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture*, MICRO 42, (New York, NY, USA), pp. 189–200, ACM, 2009.
- [50] A. Syed, "Factors affecting electromigration and current carrying capacity of FC and 3D IC interconnects," in *2010 12th Electronics Packaging Technology Conference*, pp. 538–544, Dec 2010.
- [51] K. N. Tu, H. G. Xu Gu, and W. J. Choi, *Electromigration in Solder Joints and Lines*. 2011.
- [52] S. Li, J. H. Ahn, R. D. Strong, J. B. Brockman, D. M. Tullsen, and N. P. Jouppi, "Mcpat: an integrated power, area, and timing modeling framework for multicore and manycore architectures," in *42nd Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pp. 469–480, IEEE, 2009.
- [53] "HPC-oriented Latency Numbers Every Programmer Should Know," <https://goo.gl/ftzz3a>, (accessed July 29, 2017).
- [54] T. E. Carlson, W. Heirman, and L. Eeckhout, "Sniper: Exploring the level of abstraction for scalable and accurate parallel multi-core simulations," in *International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*, pp. 52:1–52:12, Nov. 2011.
- [55] D. H. Bailey, E. Barszcz, J. T. Barton, D. S. Browning, R. L. Carter, L. Dagum, R. A. Fatoohi, P. O. Frederickson, T. A. Lasinski, R. S. Schreiber, et al., "The nas parallel benchmarks," *The International Journal of Supercomputing Applications*, vol. 5, no. 3, pp. 63–73, 1991.
- [56] C. Bienia, S. Kumar, J. P. Singh, and K. Li, "The parsec benchmark suite: Characterization and architectural implications," in *Proceedings of the 17th International Conference on Parallel Architectures and Compilation Techniques, PACT '08*, (New York, NY, USA), pp. 72–81, ACM, 2008.
- [57] J. T. Pawlowski, "Hybrid memory cube (hmc)," in *2011 IEEE Hot Chips 23 Symposium (HCS)*, pp. 1–24, Aug 2011.
- [58] M. N. Bojnordi and E. Ipek, "PARDIS: A Programmable Memory Controller for the DDRx Interfacing Standards," in *Proceedings of the 39th Annual International Symposium on Computer Architecture, ISCA '12*, (Washington, DC, USA), pp. 13–24, IEEE Computer Society, 2012.
- [59] "Hybrid memory cube," https://en.wikipedia.org/wiki/Hybrid_Memory_Cube.
- [60] "SK hynix 21 nm DRAM Cell Technology," <http://techinsights.com/about-techinsights/overview/blog/sk-hynix-21-nm-dram-cell-technology-comparison-of-1st-and-2nd-generation/>, Accessed on July 30, 2017.
- [61] "Intel Xeon Server Board - Dual Socket," <https://www.supermicro.com/products/motherboard/Xeon/C600/X10DAX.cfm>.
- [62] B. Vasquez and S. Lindsey, "The Promise of Known-good-die Technologies," in *International Conference on Multichip Modules*, pp. 1–6, Apr 1994.
- [63] R. Arnold, S. M. Menon, B. Brackett, and R. Richmond, "Test methods used to produce highly reliable known good die (KGD)," in *Proceedings. 1998 International Conference on Multichip Modules and High Density Packaging (Cat. No.98EX154)*, pp. 374–382, Apr 1998.
- [64] R. H. Parker, "Bare die test," in *Proceedings 1992 IEEE Multi-Chip Module Conference MCMC-92*, pp. 24–27, Mar 1992.
- [65] D. Chu, C. A. Reber, and D. W. Palmer, "Screening ICs on the bare chip level: temporary packaging," *IEEE Transactions on Components, Hybrids, and Manufacturing Technology*, vol. 16, pp. 392–395, Jun 1993.
- [66] W. Ballouli, T. McKenzie, and N. Alizy, "Known good die achieved through wafer level burn-in and test," in *IEEE/CPMT International Electronics Manufacturing Technology Symposium*, pp. 153–159, 2000.
- [67] D. R. Conti and J. V. Horn, "Wafer level burn-in," in *2000 Proceedings. 50th Electronic Components and Technology Conference*, pp. 815–821, 2000.
- [68] <https://www.advantest.com/products/leading-edge-products/ha1000>.
- [69] H. H. Chen, "Hierarchical built-in self-test for system-on-chip design," in *Conference, Emerging Information Technology 2005*, p. 3, Aug 2005.
- [70] C. Grecu, P. Pande, A. Ivanov, and R. Saleh, "BIST for network-on-chip interconnect infrastructures," in *24th IEEE VLSI Test Symposium*, pp. 6 pp.–35, April 2006.
- [71] J. B. Brinton and J. R. Lineback, *Packaging is becoming biggest cost in assembly, passing capital equipment*. EE Times [Online], 1999.
- [72] R. H. Katz, *Cost, Price, and Performance*. UC Berkeley, 1996.
- [73] C. A. Palesko and E. J. Vardaman, "Cost comparison for flip chip, gold wire bond, and copper wire bond packaging," in *2010 Proceedings 60th Electronic Components and Technology Conference (ECTC)*, pp. 10–13, June 2010.
- [74] "Determining Total Cost of Ownership for Data Center and Network Room Infrastructure," http://www.apc.com/salestools/CMRP-5T9PQG/CMRP-5T9PQG_R4_EN.pdf.
- [75] A. E. Papathanasiou and M. L. Scott, "Aggressive prefetching: An idea whose time has come," in *Conference on Hot Topics in Operating Systems - Volume 10, HOTOS*, (Berkeley, CA, USA), pp. 6–6, 2005.

- [76] S. Garg, D. Marculescu, R. Marculescu, and U. Ogras, "Technology-driven limits on dvfs controllability of multiple voltage-frequency island designs: A system-level perspective," in *2009 46th ACM/IEEE Design Automation Conference*, pp. 818–821, July 2009.