

# ASIC Clouds: Specializing the Datacenter for Planet-Scale Applications

By Michael Bedford Taylor, Luis Vega, Moein Khazraee, Ikuo Magaki, Scott Davidson, and Dustin Richmond

## Abstract

Planet-scale applications are driving the exponential growth of the Cloud, and datacenter specialization is the key enabler of this trend. GPU- and FPGA-based clouds have already been deployed to accelerate compute-intensive workloads. ASIC-based clouds are a natural evolution as cloud services expand across the planet. ASIC Clouds are purpose-built datacenters comprised of large arrays of ASIC accelerators that optimize the total cost of ownership (TCO) of large, high-volume scale-out computations. On the surface, ASIC Clouds may seem improbable due to high NREs and ASIC inflexibility, but large-scale ASIC Clouds have already been deployed for the Bitcoin cryptocurrency system. This paper distills lessons from these Bitcoin ASIC Clouds and applies them to other large scale workloads such as YouTube-style video-transcoding and Deep Learning, showing superior TCO versus CPU and GPU. It derives Pareto-optimal ASIC Cloud servers based on accelerator properties, by jointly optimizing ASIC architecture, DRAM, motherboard, power delivery, cooling, and operating voltage. Finally, the authors examine the impact of ASIC NRE and when it makes sense to build an ASIC Cloud.

## 1. INTRODUCTION

In the last decade, two parallel trends in the computational landscape have emerged. The first is the bifurcation of computation into two sectors: cloud and mobile. The second is the rise of dark silicon<sup>15, 3, 4, 2</sup> and dark silicon aware design techniques<sup>13, 14, 10, 16, 11</sup> such as specialization and near-threshold computation. Specialized hardware has existed in mobile computing for a while due to extreme power constraints; however, recently there has been an increase in the amount of specialized hardware showing up in cloud datacenters. Examples include Baidu's GPU-based cloud for distributed neural network acceleration, Microsoft's FPGA-based cloud for Bing Search,<sup>9</sup> and by JP Morgan Chase for hedgefund portfolio evaluation.<sup>12</sup>

At the level of a single node, we know that ASICs can offer order-of-magnitude improvements in energy-efficiency and cost-performance over CPU, GPU, and FPGA.

Our recent papers<sup>8, 6, 7, 17</sup> explore the concept of *ASIC Clouds* which are purpose-built datacenters comprised of large arrays of ASIC accelerators. ASIC Clouds are not ASIC supercomputers that scale up problem sizes for a single tightly coupled computation; rather, ASIC Clouds target scale-out workloads consisting of many independent but similar jobs, often on behalf of millions or billions of end-users.

As more and more services are built around the Cloud model, we see the emergence of planet-scale workloads (think Facebook's face recognition of uploaded pictures, or Apple's Siri voice recognition, or the IRS performing tax audits with neural nets) where datacenters are performing the same computation across many users. These scale-out workloads can easily leverage racks of ASIC servers containing arrays of chips that in turn connect arrays of replicated compute accelerators (RCAs) on an on-chip network. The large scale of these workloads creates the economical justification to pay the nonrecurring engineering (NRE) costs of ASIC development and deployment. As a workload grows, the ASIC Cloud can be scaled in the datacenter by adding more ASIC servers, unlike accelerators in say a mobile phone population,<sup>3</sup> where the accelerator-to-processor ratio is fixed at tapeout.

Our research examined ASIC Clouds in the context of four key applications that show great potential for ASIC Clouds, such as YouTube-style video transcoding, Bitcoin and Litecoin mining, and Deep Learning. ASICs achieve large reductions in silicon area and energy consumption versus CPUs, GPUs, and FPGAs. We show how to specialize the ASIC server to maximize efficiency, employing optimized ASICs, a customized printed circuit board (PCB), custom-designed cooling systems and specialized power delivery systems, and tailored DRAM and I/O subsystems. ASIC voltages are customized in order to tweak energy efficiency and minimize total cost of ownership (TCO). The datacenter itself can also be specialized, optimizing rack-level and datacenter-level thermals and power delivery to exploit the knowledge of the computation. We developed tools that consider all aspects of ASIC Cloud design in a bottom-up way, and methodologies that reveal how the designers of these novel systems can optimize TCO in real-world ASIC Clouds. Finally, we proposed a new rule that explains when it makes sense to design and deploy an ASIC Cloud, considering the engineering expense (NRE) of designing the machines.

Notably, the original version of this paper<sup>1, 8</sup> predicted Machine Learning ASIC Clouds, before Google announced the first Tensor Processing cloud in 2016.<sup>5</sup> The same paper also predicted video transcoding clouds before Facebook's

The content of this paper draws from "ASIC Clouds: Specializing the Data Center," published in *Proceedings of the IEEE Int. Symp. Computer Architecture*, June 2016, and from "Specializing the Planet's Computation: ASIC Clouds" published in *IEEE Micro*, June 2017.

Mount Shasta video transcoding ASIC Cloud design was announced in March 2019.

## 2. ASIC CLOUD ARCHITECTURE

At the heart of any ASIC Cloud is an energy-efficient, high-performance, specialized replicated compute accelerator, or RCA, that is multiplied up by having multiple copies per ASICs, multiple ASICs per server, multiple servers per rack, and multiple racks per datacenter as shown Figure 1. Work requests from outside the datacenter will be distributed across these RCAs in a scale-out fashion. All system components can be customized for the application to minimize TCO.

Each ASIC interconnects its RCAs using a customized on-chip network. The ASIC's control plane unit also connects to this network and schedules incoming work from the ASIC's off-chip router onto the RCAs. Next, the packaged ASICs are arranged in lanes on a customized PCB, and connected to a controller which bridges to the off-PCB interface (1-100 GigE, RDMA, PCI-e, etc). In some cases, DRAMs may connect directly to the ASICs. The controller can be implemented by an FPGA, microcontroller, or a Xeon processor and schedules remote procedure calls (RPCs) that come from the off-PCB interface on to the ASICs. Depending on the application, it may implement the nonacceleratable part of the workload or perform UDP/TCP-IP offload.

Each lane is enclosed by a duct and has a dedicated fan blowing air through it across the ASIC heatsinks. Our simulations indicate that using ducts results in better cooling performance compared to conventional or staggered layout. The PCB, fans, and power supply are enclosed in a 1U server, which is then assembled into racks in a datacenter. Based on ASIC needs, the PSU and DC/DC converters are customized for each server.

## 3. DESIGNING AN ASIC CLOUD

Our ASIC Cloud Server configuration evaluator, as shown in Figure 2a, starts with a Verilog implementation of an accelerator, or a detailed evaluation of the accelerator's

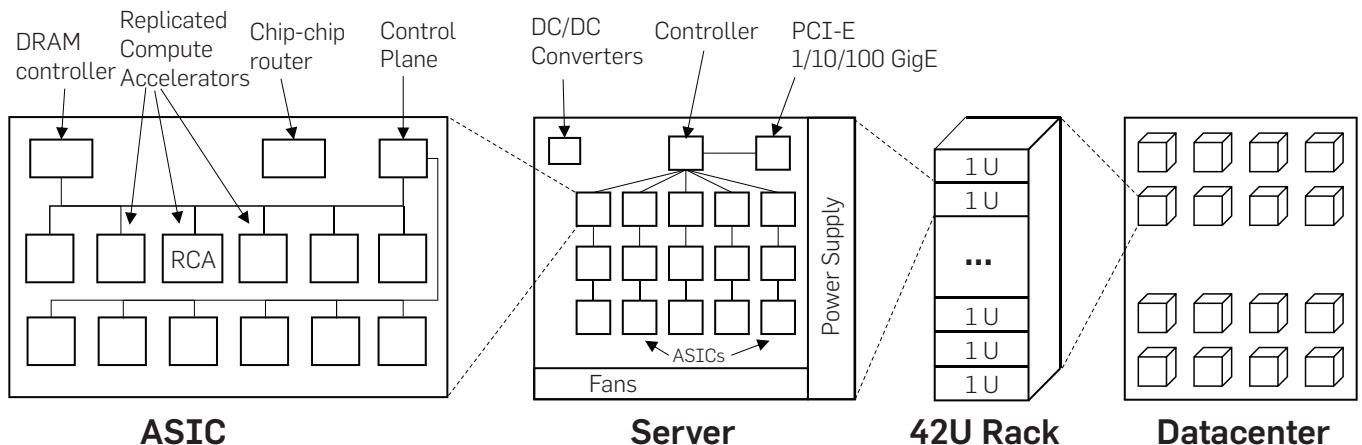
properties from the research literature. In the design of an ASIC Server, we must decide how many chips should be placed on the PCB and how large, in mm<sup>2</sup> of silicon, each chip should be. The size of each chip determines how many RCAs will be on each chip. In each duct-enclosed lane of ASIC chips, each chip receives around the same amount of airflow from the intake fans, but the most downstream chip receives the hottest air, which includes the waste heat from the other chips. Therefore, the thermally bottlenecking ASIC is the one in the back, shown in our detailed Computational Fluid Dynamics (CFD) simulations as shown in Figure 2b. Our simulations show that breaking a fixed heat source into smaller ones with the same total heat output improves the mixing of warm and cold area, resulting in lower temperatures. Using thermal optimization techniques, we established fundamental connection between an RCA's properties, the number of RCAs placed in an ASIC, and how many ASICs go on a PCB in a server. Given these properties, our heat sink solver determines the optimal heat sink configuration. Results are validated with the CFD simulator. In the sidebar entitled "Design Space Evaluation," we show how we apply this evaluation flow across the design space in order to determine TCO and Pareto optimal points that trade off \$ per op/s (an accelerator's hardware cost efficiency) and W per op/s (an accelerator's energy efficiency).

## 4. APPLICATION CASE STUDIES

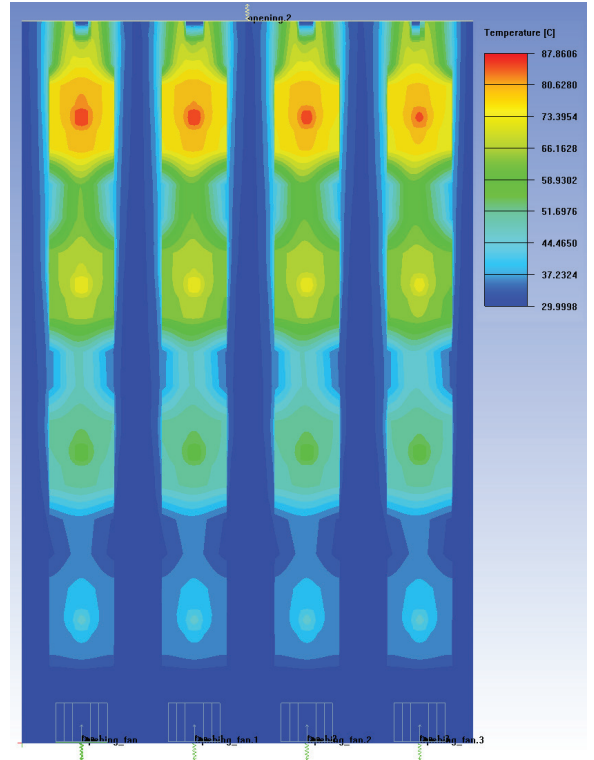
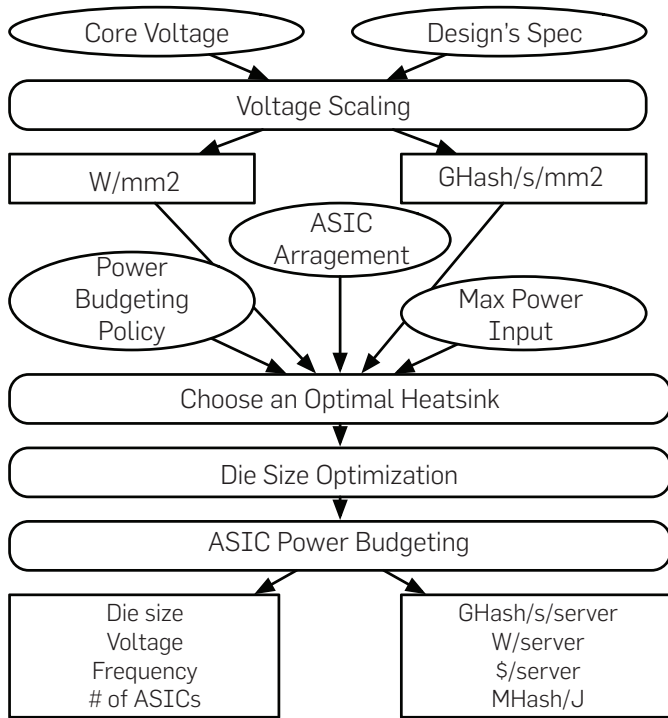
To explore ASIC Clouds across a range of accelerator properties, we examined four applications: Bitcoin mining, Litecoin mining, Video Transcoding, and Deep Learning that span a diverse range of properties, as shown in Figure 3.

Perhaps the most mature of these applications is Bitcoin mining. Our inspiration for ASIC Clouds came from our intensive study of Bitcoin mining clouds,<sup>4</sup> which are one of the first known instances of a real life ASIC Cloud. Figure 4 shows the massive scale out of the Bitcoin mining workload, which in 2015 operated at the performance of 3.2 billion GPUs. Bitcoin

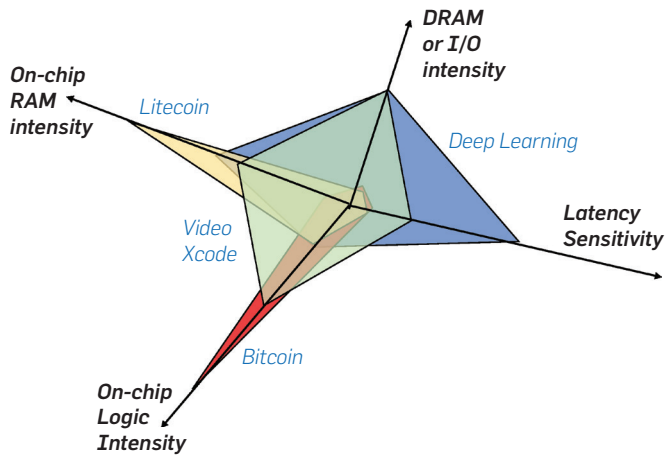
**Figure 1. High-level abstract architecture of an ASIC Cloud. Specialized replicated compute accelerators (RCA) are multiplied up by having multiple copies per ASICs, multiple ASICs per server, multiple servers per rack, and multiple racks per datacenter. Server controller can be an FPGA, microcontroller, or a Xeon processor. Power delivery and cooling system are customized based on ASIC needs. If required, there would be DRAMs on the PCB as well.**



**Figure 2. Evaluating an ASIC configuration.** (a) The server cost, per server hash rate, and energy efficiency are evaluated using RCA properties and a flow that optimizes server heatsinks, die size, voltage, and power density. (b) Thermal verification of an ASIC Cloud server using CFD tools to validate the flow results. The farthest ASIC from the fan has the highest temperature and is the bottleneck for power per ASIC at a fixed voltage and energy efficiency.



**Figure 3. Accelerator properties.** We explored applications with diverse requirements.



clouds have undergone a rapid ramp from CPU to GPU to FPGA to the most advanced ASIC technology available today. Bitcoin is a very logic intensive design which has high power density and no need for SRAM or external DRAM.

Litecoin is another popular cryptocurrency mining system that has been deployed into clouds. Unlike Bitcoin, it is an SRAM-intensive application which has low power density.

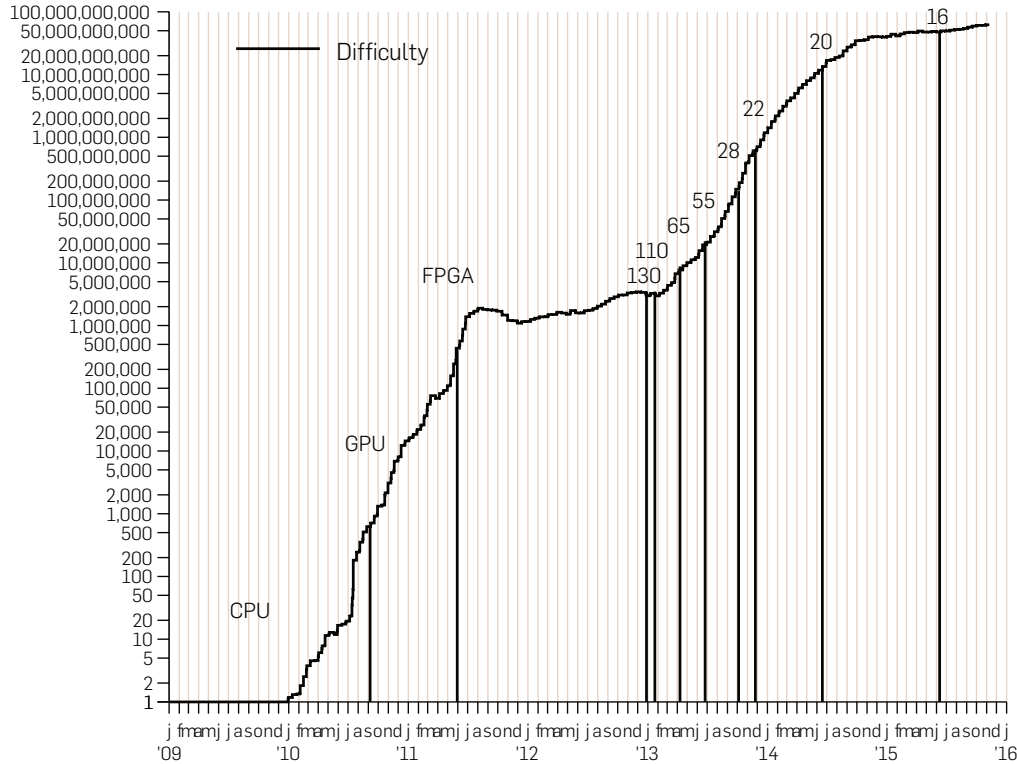
Video Transcoding, which converts from one video format to another, currently takes almost 30 high-end Xeon servers to do in real-time. As every cell phone can easily be a video source, as well as every Internet-of-Things device, it has the potential to be an unimaginably large planet-scale computation. Video Transcoding is an external memory-intensive application that needs DRAMs next to each ASIC and also high off-PCB bandwidth.

Finally, Deep Learning is extremely compute-intensive and is likely to be used by every human on the planet. Deep Learning is often latency-sensitive so our Deep Learning neural net accelerator has a tight low-latency SLA.

For our Bitcoin and Litecoin studies, we developed the RCA and got the required parameters such as gate count from placed and routed designs in UMC 28nm using Synopsys IC compiler and analysis tools (e.g., PrimeTime). For Deep Learning and Video Transcoding, we extract properties from accelerators designed in the research literature.

Design space exploration is application-dependent, and there are frequently additional constraints. For example, for video transcode application, we model the PCB real estate occupied by these DRAMs, which are placed on either side of the ASIC they connect to, perpendicular to airflow. As the number of DRAMs increases, the number of ASICs placed in a lane decreases for space reasons. We model the more expensive PCBs required by DRAM, with more layers and better signal/power integrity. We employ

**Figure 4. Evolution of Specialization, Bitcoin cryptocurrency mining clouds.** Numbers are ASIC nodes, in nm, which annotate the first date of release of a miner on that technology. Difficulty is the ratio of the total Bitcoin hash throughput of the world, relative to the initial mining network throughput, which was 7.15 MH/s. In the 6-year period preceding Nov 2015, the throughput increased by a factor of 50 billion times, corresponding to a world hash rate of approximately 575 million GH/s.

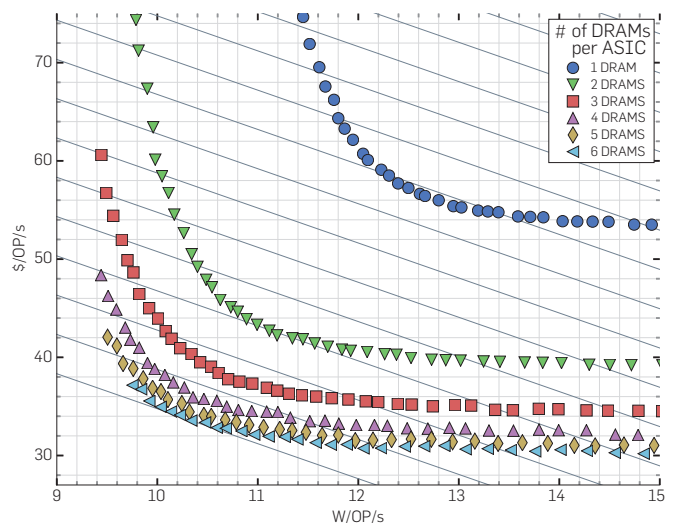


two 10-GigE ports as the off-PCB interface for network-intensive clouds, and model the area and power of the memory controllers.

After having all thermal constraints in place, we optimized ASIC server design targeting two conventional key metrics, namely cost per op/s and power per op/s, and then apply TCO analysis. TCO analysis incorporates the datacenter-level constraints such as the cost of power delivery inside the datacenter, land, depreciation, interest, and the cost of energy itself. With these tools, we can correctly weight these two metrics and find the overall optimal point (TCO-optimal) for the ASIC Cloud.

Our ASIC Cloud infrastructure explores a comprehensive design space, such as DRAMs per ASIC, logic voltage, area per ASIC, and number of chips. DRAM cost and power overhead are significant, and so the Pareto-optimal Video Transcoder designs ensure DRAM bandwidth is saturated, linked chip performance to DRAM count. As voltage and frequency are lowered, area increases to meet the performance requirement. Figure 5 shows the Video Transcode Pareto curve for 5 ASICs per lane and different number of DRAMs per ASIC. The tool is composed of two tiers. The top tier uses brute force to explore all of the possible configurations in order to find the energy-optimal, cost-optimal, and TCO-optimal points are chosen based on the Pareto results. The leaf tier consists of a variety of “expert solvers” that compute optimal properties of the server components; for example, CFD simulations for heat sinks, DC-DC

**Figure 5. Pareto curve example for Video Transcode.** Exploring different number of DRAMs per ASIC and logic voltage for optimal TCO per performance point. Voltage increases from left to right. Diagonal lines show equal TCO per performance values and the closer to the origin the lower the TCO per performance. This plot is for 5 ASICs per lane.



converter allocation, circuit area/delay/voltage/energy estimators, and DRAM property simulation. In many cases, these solvers export their data as large tables of memoized numbers for every component to the brute force solver.

**Figure 6. ASIC Cloud optimization results for four applications.** Each table presents energy-optimal, TCO-optimal, and cost optimal server properties. Energy optimal server uses lower voltage to increase the energy efficiency. Cost optimal servers use higher voltage to increase silicon efficiency. TCO-optimal has a voltage between these two and balances energy versus silicon cost.

	Energy optimal	TCO optimal	Cost optimal
ASICs per server	120	72	24
Logic Voltage (V)	0.400	0.459	0.594
Clock Freq. (MHz)	71	149	435
Die Area (mm <sup>2</sup> )	599	540	240
GH/s/server	7,292	8,223	3,451
W/server	2,645	3,736	2,513
\$/server	12,454	8,176	2,458
W/GH/s	0.363	0.454	0.728
\$/GH/s	1.708	0.994	0.712
TCO/GH/s	3.344	2.912	3.686

(a) Bitcoin

	Energy optimal	TCO optimal	Cost optimal
ASICs per server	120	120	72
Logic Voltage (V)	0.459	0.656	0.866
Clock Freq. (MHz)	152	576	823
Die Area (mm <sup>2</sup> )	600	540	420
MH/s/server	405	1,384	916
W/server	783	3,662	3,766
\$/server	10,971	11,156	6,050
W/MH/s	1.934	2.645	4.113
\$/MH/s	27.09	8.059	6.607
TCO/MH/s	37.87	19.49	23.70

(b) Litecoin

	Energy optimal	TCO optimal	Cost optimal
DEAMs per ASIC	3	6	9
ASICs per Server	64	40	32
Logic Voltage (V)	0.538	0.754	1.339
Clock Freq. (MHz)	183	429	600
Die Area (mm <sup>2</sup> )	564	498	543
Kfps/server	126	158	189
W/server	1,146	1,633	3,101
\$/server	7,289	5,300	5,591
W/Kfps	9.073	10.34	16.37
\$/Kfps	57.68	33.56	29.52
TCO/Kfps	100.3	78.46	97.91

(c) Video Transcode

	Energy optimal	TCO optimal	Cost optimal
Chip type	4x2	2x2	2x1
ASICs per server	32	64	96
Logic Voltage (V)	0.900	0.900	0.900
Clock Freq. (MHz)	606	606	606
TOps/s/server	470	470	353
W/server	3,278	3,493	2,971
\$/server	7,809	6,228	4,146
W/TOps/s	6.975	7.431	8.416
\$/TOps/s	16.62	13.25	11.74
TCO/TOps/s	46.22	44.28	46.51

(d) Deep learning

## 5. RESULTS

Details of optimal server configurations for energy-optimal, TCO-optimal, and cost-optimal designs for each of the applications are shown in Figure 6.

For example, for Video Transcode, the cost-optimal server packs the maximum number of DRAMs per lane, 36, maximizing performance. However, increasing the number of DRAMs per ASIC requires higher logic voltage (1.34V) and corresponding frequencies to attain performance within the max die area constraint, resulting in less energy-efficient designs. Hence, the energy-optimal design has fewer DRAMs per ASIC and per lane (24), although gaining back some performance by increasing ASICs per lane, which is possible due to lower power density at 0.54V. The TCO-optimal design increases DRAMs per lane, 30, to improve performance, but is still close to the optimal energy efficiency at 0.75V, resulting in a die size and frequency between the other two optimal points.

In Figure 7, we compare the performance of CPU Clouds versus GPU Clouds versus ASIC Clouds for the four applications that we presented. ASIC Clouds outperform CPU Cloud TCO per op/s by 6270x; 704x; and 8695x for Bitcoin, Litecoin, and Video Transcode, respectively. ASIC Clouds outperform GPU Cloud TCO per op/s by 1057x, 155x, and 199x, for Bitcoin, Litecoin, and Deep Learning, respectively.

## 6. FEASIBILITY OF ASIC CLOUDS: THE TWO-FOR-TWO-RULE

When does it make sense to design and deploy an ASIC Cloud? The key barrier is the cost of developing the ASIC Server, which includes both the mask costs (about \$1.5M for the 28 nm node we consider here and much higher for the latest 7nm node) and the ASIC design costs, which collectively comprise the nonrecurring engineering expense (NRE). To understand this trade-off, we proposed the

two-for-two rule. If the cost per year (i.e., the TCO) for running the computation on an existing cloud exceeds the NRE by 2X, and you can get at least a 2X TCO per operation/second improvement, then going ASIC Cloud is likely to save money. Figure 8 shows a wider range of breakeven points. Essentially, as the TCO exceeds the NRE by more and more, the required speedup to break even declines. As a result, almost any accelerator proposed in the literature, no matter how modest the speedup, is a candidate for ASIC Cloud, depending on the scale of the computation. Our research makes the key contribution of noting that in deployment of ASIC Clouds, NRE and scale can be more determinative than

the absolute speedup of the accelerator. The main barrier for ASIC Clouds is to reign in NRE costs so they are appropriate for the scale of the computation. In many research accelerators, TCO improvements are extreme (such as in Figure 7), but authors often unnecessarily target expensive, latest generation process nodes because they are more cutting edge. This tendency raises the NRE exponentially, reducing economic feasibility. A better strategy is to target the older nodes that still attain sufficient TCO improvements.

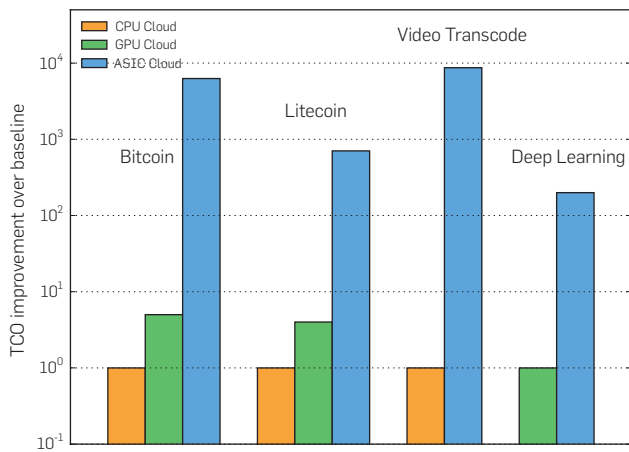
**7. POST-PUBLICATION INSIGHT: YOU WANT TO TARGET EIGHT TIMES TCO IMPROVEMENT**

The two-for-two rule examines a lower bound for what the TCO improvements of an ASIC cloud need to be, based on how large the pre-ASIC cloud TCO is compared to the NRE of building an accelerator and show that extreme hundred times TCO improvements are not needed.

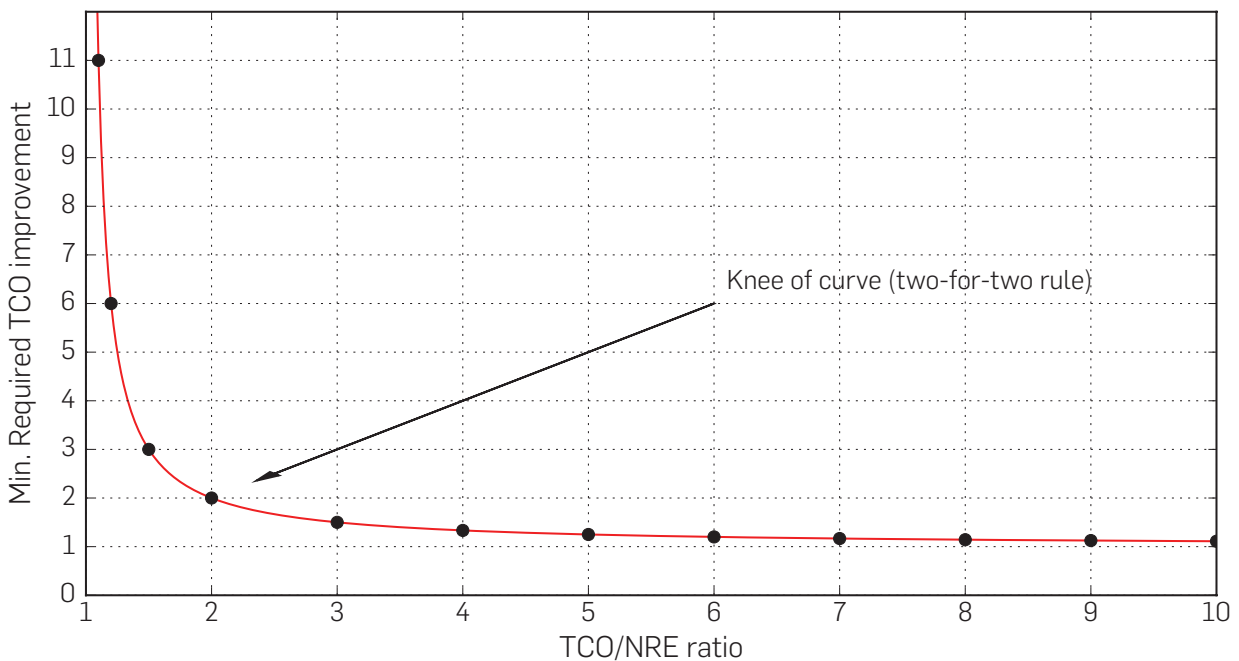
Our subsequent experience post-publication of the ASIC cloud suggests another way to look at the question of how aggressive an accelerator is necessary. We believe in most cases that eight times TCO improvement is usually a good place to target when developing a new kind of ASIC cloud.

In most realistic scenarios, the pre-ASIC cloud TCO can be in the hundreds of millions or billions of dollars, far out-shadowing the ASIC development costs for all but the latest nodes (e.g., 7nm). Practically speaking, the first two times will reduce your TCO in half, that is, one billion dollars become 500 million dollars. The second two times will only save 250 million dollars, useful but not essential on the first ASIC iteration. The second two times is needed to provide risk margin for the performance and energy efficiency

**Figure 7. CPU Cloud vs. GPU Cloud vs. ASIC Cloud “Deathmatch.” ASIC servers greatly outperform the best non-ASIC alternative in terms of TCO per op/s.**



**Figure 8. Two-for-two rule: moderate speed-up with low NRE beats high speed-up at high NRE. The points are break even points for ASIC Clouds.**



uncertainty of the design—will the original software be optimized more making the chip less good relatively, will the chip have less than expected TCO improvement, et cetera. The final two times addresses the issue that the pre-ASIC cloud hardware (e.g., GPU or CPU) will also improve and could possibly improve by two times by the time you have deployed your ASIC cloud system.

## 8. CONCLUSION

Our research generalizes primordial Bitcoin ASIC Clouds into an architectural template that can apply across a range of planet-scale applications. Joint knowledge and control over datacenter and hardware design allow for ASIC Cloud designers to select the optimal design that optimizes energy and cost proportionally to optimize TCO. We demonstrated methodologies that can be used to design TCO-optimal clouds, answering long-standing questions even in contemporary Bitcoin ASIC Clouds. Our work analyses the impact of NRE and scale on deployment of ASIC Clouds, tying it to the TCO-improvement and in turn the energy and cost efficiency of the cloud.

Our work advances research practice by showing how to examine accelerators at a systems level instead of at the level of a single chip. We evaluate ASIC Cloud chip design, server design, and finally datacenter design in a cross-layer system-oriented way. This joint knowledge and control over datacenter and hardware design allow for ASIC Cloud designers to select the optimal design that optimizes energy and cost proportionally. We developed the tools and revealed how the designers of these novel systems can optimize the TCO in real-world ASIC Clouds.


We developed a rule of thumb for when it makes sense to go ASIC Cloud, the two-for-two rule. The main barrier for ASIC Clouds is to reign in NRE costs so they are appropriate for the scale of the computation. In many research accelerators, TCO improvements are extreme, but authors also target expensive, latest generation process nodes because they are more cutting edge. But this habit raises the NRE exponentially, reducing economic feasibility. Our most recent work<sup>6</sup> suggests that a better strategy is to lower NRE cost by targeting older nodes that still have sufficient TCO per op/s benefit.

Looking to the future, our work suggests that both Cloud providers and silicon foundries would benefit by investing in technologies that reduce the NRE of ASIC design, such as open source IP such as RISC-V, in new labor-saving development methodologies for hardware and also in open source backend CAD tools. With time, mask costs fall by themselves, but currently older nodes such as 65 nm and 40 nm may provide suitable TCO per op/s reduction, with half the mask cost and only a small difference in performance and energy efficiency from 28, 16, or 7 nm. Foundries should take interest in ASIC Cloud's low-voltage scale out design patterns because they lead to greater silicon wafer consumption than CPUs within fixed environmental energy limits.

With the coming explosive growth of planet-scale computation, we must work to contain the exponentially growing environmental impact of datacenters across the world.

ASIC Clouds promise to help address this problem. By specializing the datacenter, they can do greater amounts of computation under environmentally determined energy limits. The future is planet-scale, and specialized ASICs will be everywhere.

## Acknowledgments

This work was supported by both the JUMP ADA Center and the STARnet CFAR center, both funded by SRC and DARPA. Special thanks go to Partha Ranganathan for his support of our research over the years. 

## References

1. ASIC clouds: Specializing the datacenter. *UCSD CSE Tech Report CS2016-1016*. May 8, 2016. [https://csetechrep.ucsd.edu/Dienst/UI/2.0/Describe/ncstrLucsd\\_cse/CS2016-1016](https://csetechrep.ucsd.edu/Dienst/UI/2.0/Describe/ncstrLucsd_cse/CS2016-1016).
2. Esmailzadeh, H., Blem, E., Amant, R.S., Sankaralingam, K., Burger, D. Power limitations and dark silicon are challenging the future of multicore. In *TOCS*, 2012.
3. Goulding, N., et al. GreenDroid: A mobile application processor for a future of dark silicon. In *HOTCHIPS*, 2010.
4. Goulding-Hotta, N., Sampson, J., Venkatesh, G., Garcia, S., Auricchio, J., Huang, P.-C., Arora, M., Nath, S., Bhatt, V., Babb, J., Swanson, S., Taylor, M.B. The greendroid mobile application processor: An architecture for silicon's dark future. *IEEE Micro* 2, 31 (2011), 86–95.
5. Jouppi, N.P., et al. In-datacenter performance analysis of a tensor processing unit. In *International Symposium on Computer Architecture (ISCA)*, 2017.
6. Khazraee, M., et al. Moonwalk: NRE optimization in ASIC clouds. In *International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2017.
7. Khazraee, M., Vega, L., Magaki, I., Taylor, M. Specializing a planet's computation: ASIC clouds. *IEEE Micro*, May 2017.
8. Magaki, I., et al. ASIC clouds: Specializing the datacenter. In *International Symposium on Computer Architecture (ISCA)*, 2016.
9. Putnam et al. A reconfigurable fabric for accelerating large-scale datacenter services. In *International Symposium on Computer Architecture (ISCA)*, 2014.
10. Sampson, J., Venkatesh, G., Goulding-Hotta, N., Garcia, S., Swanson, S., Taylor, M.B. Efficient complex operators for irregular codes. In *HPCA*, 2011.
11. Shafique, M., Garg, S., Henkel, J., Marculescu, D. The EDA challenges in the dark silicon era: Temperature, reliability, and variability perspectives. In *DAC*, 2014.
12. Weston, S. FPGA accelerators at JP Morgan chase, 2011. Stanford Computer Systems Colloquium, <https://www.youtube.com/watch?v=9NqXIETADn0>.
13. Taylor, M. A landscape of the new dark silicon design regime. *IEEE Micro*, Sept-Oct. 2013.
14. Taylor, M.B. Is dark silicon useful? Harnessing the four horsemen of the coming dark silicon apocalypse. In *DAC*, 2012.
15. Venkatesh, G., Sampson, J., Goulding, N., Garcia, S., Bryksin, V., Lugo-Martinez, J., Swanson, S., Taylor, M.B. Conservation cores: Reducing the energy of mature computations. In *Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2010.
16. Venkatesh and others. Qscores: Configurable co-processors to trade dark silicon for energy efficiency in a scalable manner. In *MICRO*, 2011.
17. Xie, S., Davidson, S., Magaki, I., Khazraee, M., Vega, L., Zhang, L., Taylor, M.B. Extreme datacenter specialization for planet-scale computing: ASIC clouds. *SIGOPS Oper. Syst. Rev.* 1, 52 (2018), 96–108.

**Michael Bedford Taylor** (prof.taylor@gmail.com), University of Washington, WA, USA.

**Luis Vega** (vegaluis@cs.washington.edu), University of Washington, WA, USA.

**Moein Khazraee** (mkhazraee@cs.ucsd.edu), UC San Diego, CA, USA.

**Ikuo Magaki** (ikuo.magaki@icloud.com), UC San Diego, CA, USA.

**Scott Davidson** and **Dustin Richmond** ({stdavids, dustinar}@uw.edu), University of Washington, WA, USA.