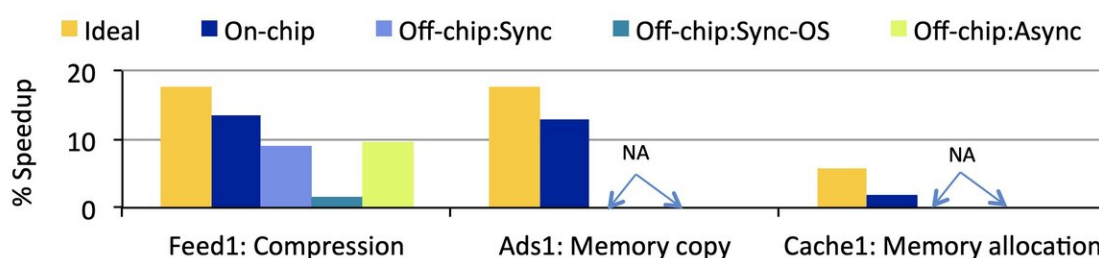


Analytical model predicts exactly how much a piece of hardware will speed up data centers

April 7 2020



Accelerometer-estimated speedup for key overheads identified by the researchers. Credit: University of Michigan

Large-scale software services fight the efficiency battle on two fronts—efficient software that is flexible to changing consumer demands, and efficient hardware that can keep these massive services running quickly even in the face of diminishing returns from CPUs. Together, these factors determine both the quality of the user experience and the performance, cost, and energy efficiency of modern data centers.

A change on one front requires adjustments on the other, and a new software architecture growing in popularity has posed a challenge to the hardware solutions current in most data centers. Called microservices,

this modular approach to designing big enterprise software has left something to be desired in its interactions with another major rising force in datacenter efficiency, hardware accelerators.

To bring these two promising technologies together more effectively, CSE Ph.D. student Akshitha Sriraman, working with researchers from Facebook, has designed a way to measure exactly how much a hardware [accelerator](#) would speed up a datacenter. Appropriately named Accelerometer, the analytical model can be applied in the early stages of an accelerator's design to predict its effectiveness before ever being installed.

Still a somewhat new technology in general computing usage, the effectiveness of hardware accelerators isn't as easy to predict as CPUs, which have decades of experience behind them. Investing in this sort of diverse custom hardware presents a risk at scale, since it might not live up to its expectations.

But the potential for a big impact is there. Designed to perform one type of function extremely quickly, accelerators could theoretically be called upon for all the redundant, repetitive tasks used in common by bigger applications.

That includes microservices. This software architecture approach conceives of a larger application as a collection of modular, task-specific services that can each be improved upon in isolation. This allows for changes to be made to the larger application without needing to change one huge, central codebase. It also allows for more services to be added more easily.

Sriraman demonstrated that as few as 18% of most microservices' CPU cycles are spent executing instructions that are core to their functionality. The remaining 82% are spent on common operations that

are ripe for accelerating.

"Accelerating these overheads we identified can indeed improve speedup to a significant extent," Sriraman says. Beyond speed, it would make all of the datacenter's functions cheaper and more energy efficient. "Acceleration will allow us to pack more work for the same power constraints and improve resource utilization at scale, so data center energy and cost savings will improve greatly."

The issue with microservices is that their designs can turn out to be quite dissimilar, particularly with regard to how they interact with hardware. For example, a microservice can communicate with an accelerator while continuing to run other instructions on a CPU, or it could bring all of its functions to a halt while it offloads to the accelerator. Both of these cases face different "offload overheads" (the time spent sending a task from one processor to another), which becomes lost time for the datacenter if it's not accounted for.

"Each of these software design choices can result in different overheads that affect the overall speedup from acceleration," says Sriraman. This overhead is left out of the picture in prior work, she continues, as is the impact of the different microservice designs on performance.

Additionally, accelerators themselves have to be used judiciously to have a net [positive effect](#).

"Throwing an accelerator at every problem is ridiculous because it takes a lot of time, cost, and effort to build, test, and deploy each one," she concludes. "There is a real need to precisely understand what and how to accelerate."

Accelerometer is an analytical model that measures exactly how much performance would be improved by installing a given processor, if at all,

with all of these nuances taken into account. That means it measures the positive effect of acceleration as well as the negative effect of spending time shuffling instructions around between computing components. And its capabilities aren't limited to new accelerators—the model can be applied to any kind of hardware, ranging from a simple CPU optimization to an extremely specialized remote ASIC.

The tool was validated in Facebook's production environment using three retrospective case studies, demonstrating that its real speedup estimates have less than 3.7% error.

The model is sufficiently accurate to already be put to use by Facebook, with early interest from other companies.

"We have received word that several of the big cloud players have started using Accelerometer to quickly discard bad accelerator choices and identify the good ones, to make well-informed hardware investments," Sriraman says. Facebook is using the model to explore new accelerators, incorporating it as a first step to quickly sort out good and bad [hardware](#) choices.

This project, titled "[Accelerometer: Understanding Acceleration Opportunities for Data Center Overheads at Hyperscale](#)," was accepted by the 2020 Architectural Support for Programming Languages and Operating Systems (ASPLOS) Conference and presented virtually.

More information: Accelerometer: Understanding Acceleration Opportunities for Data Center Overheads at Hyperscale.
[research.fb.com/publications/a ... heads-at-hyperscale/](https://research.fb.com/publications/a-heads-at-hyperscale/)

Provided by University of Michigan

Citation: Analytical model predicts exactly how much a piece of hardware will speed up data centers (2020, April 7) retrieved 1 May 2024 from <https://techxplore.com/news/2020-04-analytical-piece-hardware-centers.html>

This document is subject to copyright. Apart from any fair dealing for the purpose of private study or research, no part may be reproduced without the written permission. The content is provided for information purposes only.