

Hybrid Memory Cube Angles for Exascale

July 10, 2012 | Michael Feldman

Computer memory is currently undergoing something of an identity crisis. For the past 8 years, multicore microprocessors have been creating a performance discontinuity, the so-called memory wall. It's now fairly clear that this widening gap between compute and memory performance will not be solved with conventional dynamic random-access memory (DRAM) products. But there is one technology under development that aims to close that gap, and its first use case will likely be in the ethereal realm of supercomputing.



About a year and half ago, memory-maker Micron Technology came up with the Hybrid Memory Cube (HMC), a multi-chip module (MCM) device specifically designed to scale the memory wall. The goal was to offer a memory technology that matches the needs of core-happy CPUs and GPUs and do so in a way that is attractive to computer makers.

In a nutshell, HMC glues a logic control chip to a 3D memory stack, all of which are connected with Through Silicon Vias (TSVs). The technology promises not only to deliver an order of magnitude performance increase, but also to keep pace with future microprocessors as those designs continue to

add cores. Micron claims a single HMC device can deliver 15 times the performance of today's DDR3 modules and can do so with 70 percent less energy and in 90 percent less space. Latency is expected to decrease as well, although no specific claims are being made in that regard.

According to Dean Klein, VP of Micron's Memory System Development, the problem with conventional DRAM technology is that they've pushed CMOS technology about as far as it's going to go under the DDR model. Although DDR4 products are slated to ship before this end year, there is currently no DDR5 on the drawing board. That's a problem, especially considering that DDR5 would probably be coming out toward the end of the decade, just when the first exascale supercomputers are expected to appear.

But even if DDR evolution is maintained through 2020, it would almost certainly fall short of the needs of exascale computing. Such machines are expected to require per-node memory bandwidth in excess of 500 terabytes/second. Klein says they just can't boost the signal rates much more on the DDR design, and if they tried, power consumption would go in the wrong direction.

The HMC design gets around those limitations by going vertical and using the TSV technology to parallelize communication to the stack of memory chips, which enables much higher transfer rates. Bandwidth between the logic and the DRAM chips are projected to top a terabit per second (128 GB/second), which is much more in line with exascale needs.

Another important aspect of the design is that the interface abstracts the notion of reads and writes. That means a microprocessor's memory controller doesn't need to know about the underlying technology that stores the bits. So one could build an HMC device that was made up of DRAM or NAND flash, or even some combination of these technologies. That frees up the microprocessor and other peripheral devices from being locked into a particular memory type and, in general, should make system designs more flexible.

To move HMC beyond a science project, Micron put together a consortium and attracted key players, including competitors,

to back the technology. Today the Hybrid Memory Cube Consortium consists of some of the industry's heaviest hitters: Samsung, Microsoft, IBM, ARM, HP, IBM, Altera, Xilinx, Open-Silicon, and SK hynix. The group's immediate goal is to develop a standard interface for the technology so that multiple manufacturers can build compliant HMC devices. The formal standard is due out later this year.

A key partner with Micron has been Intel, a vendor with a particular interest in high-performance memory. The chipmaker's immediate motivation to support HMC is its Xeon line (including, soon, the manycore Xeon Phi), which is especially dependent on performant memory. In fact, without such memory, the value of high-end server chips is greatly diminished, since additional cores doesn't translate into more performance for the end user. The relative success of future multicore and manycore processors will depend, to a large extent, on memory wall-busting technology.

Further out, Intel is looking at HMC as a technology to support its own aspirations to develop components for exascale supercomputers. Last year Intel helped Micron build an HMC prototype, which CTO Justin Rattner talked up at last September's Intel Developer Forum. Although the chipmaker will presumably assist Micron if and when it starts churning out commercial silicon, neither company has offered a timeline for an HMC product launch. Klein did say that its prototype has been in the hands of select customers (HPC users and others) for several months, and their intent is to commercialize the technology.

And not just for high performance computing market. Although supercomputing has the greatest immediate need for such technology, other application areas, like networking, could also benefit greatly from HMC's high bandwidth characteristics. And because of the promised power savings, even the high-volume mobile computing market is a potential target.

The biggest challenge for HMC is likely to be price. In particular, the use of TSV and 3D chip-stacking is in its infancy and by all accounts, will not come cheaply -- at least not initially. And when you're talking about 10PB of memory for an exascale machine or 1MB for a mobile phone, cost is a big consideration.

Other technologies like HP's memristor, Magneto-resistive Random-Access Memory (MRAM), or Phase Change Memory (PCM) could come to the fore in time for the exascale era, but each one has its own challenges. As Klein notes, there is no holy grail of memory that encapsulates every desired attribute -- high performance, low-cost, non-volatile, low-power, and infinite endurance.

The nice thing about HMC is that it can encapsulate DRAM as well as other memory technologies as they prove themselves. For the time being though, dynamic random-access memory will remain as the foundation of computer memory in the datacenter. "DRAM is certainly going to with us, at least until the end of the decade," admits Klein. "We really don't have a replace technology that looks as attractive."

Copyright © 1994-2011 **Tabor Communications**, Inc. All Rights Reserved.

HPCwire is a registered trademark of Tabor Communications, Inc. Use of this site is governed by our [Terms of Use](#) and [Privacy Policy](#).
Reproduction in whole or in part in any form or medium without express written permission of Tabor Communications Inc. is prohibited.

Powered by **Xtenit**