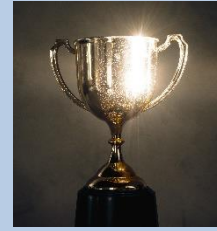# Will 2020's Hot Topics in Storage Still be Hot in 2021?

Here is a look at three topics for 2021 - storage-class memory, computational storage, and composable infrastructure - what they promised to deliver, to what extent they have lived up to their promise, and what to expect in the coming decade.

**Storage-Class Memory (SCM):** While the access speed of non-volatile storage as increase by orders of magnitude with the migration of hard disk drives (HDDs) to flash-based SSDs, the speed of DRAM has also increased, as has its price. This has resulted in what some industry pundits call the "RAM gap" - the price and performance discrepancy between DRAM and flash-based SSDs. Enter SCM, which promises to fill this gap with a technology that is 10X faster than SSDs, is non-volatile, and costs less than DRAM. While several SCM technologies have been fielded (including **Intel's® Optane™**, **vSMP MemoryONE**, mixed DRAM-flash devices, **Samsung's Z-SSD**, and others) none of them have taken off significantly, primarily due to difficulties with their programming model.

The one place that these technologies have taken off is in storage arrays, where they can significantly reduce hardware costs by reducing the RAM needed. There are new technologies on the way (MRAM, NRAM, PCM, and STT-RAM) that promise to provide non-volatile byte-addressible storage that might meet the need that SCM has always promised to address.

**Computational Storage:** One of the biggest delays in analyzing extremely large data sets is the movement of data from storage such as SSDs to DRAM in the server. In the case of transfers from local SSDs, this speed is limited by the speed of the PCI Express® (PCIe) bus (3.94GB/s for the PCIe Gen3 x4 that is typical for U.2 SSDs). To put this in perspective, it would take about 26 minutes to transfer 100TB of data from sixteen (16) PCIe Gen3 U.2 SSDs. **Computational storage** turns this problem on its head by saying "what if you don't transfer the data, but actually perform the processing in the SSD?"

Of course, there are tradeoffs in this approach, chief among them being that you cannot put an Intel X86 processor inside of a U.2 SSD. Most SSDs utilize ARM processors internally (in some cases augmented by an FPGA), which means changes to the applications that will analyze the data. While this may be doable for third-party software, is is not likely for popular 'big data" programs like **SAP HANA**, limiting the applicability of computational

storage. The ARM processors also do not necessarily have processing power on par with X86 processors in "big iron" servers, so there is definitely a tradeoff between reducing transfer speeds and processing speeds. With 400GbE networks coming out (48GB/s – 12 times the speed of PCIe Gen3x4) and the advent of PCIe Gen4 and Gen5 (especially with 8 lanes), the transfer time issue may be going away.

Where computational storage has found some legs is with embedded applications, where the SSD is more like very small footprint, low-power, hardened single-board computer. This is an important capability for a variety of use cases in aerospace, military, automotive, and industrial applications, where space and power are at a premium and an "encapsulated" system that is field-pluggable is highly desirable.

**Composable Infrastructure: [Composable infrastructure](#)** traces its legacy back to the previous decade (the "naughts"), where it was called "I/O virtualization". Composable infrastructure essentially disaggregates the entire server and allows the parts to be shared across programmable high-speed fabrics, including any of the flavors of NVMe-oF (including 100GbE), or a "fabricized" version of PCIe running NVMe. The fabrics incorporate a switching infrastructure that allows components to be connected and reconnected at will to various servers – essentially 'sharing' those components, reducing costs. In general, composable infrastructure makes the most sense when the components to be shared are both relatively expensive and sparsely utilized – things like GPGPUs or FPGA cards. Of course, this flexibility isn't free – running connections through fabrics always adds latency, and the increased amount of cabling required does add complexity to racks. Nevertheless, there are a good number of use cases that fit these "best case" parameters – media and entertainment and oil and gas modeling to name a few. If combined with smaller modular servers (which we will see more of in the coming decade), composable infrastructure can make a lot of sense.

[Mike Heumann](#)

Managing Partner and Chief Analyst for G2M Communications, a re-grate-it brand.