

Impact of NVMe on Servers in Megawebsites

- Server Based Flash Evolves

About the Presenter



Christopher King

Infrastructure Architect (Systems/Storage)

cking@linkedin.com

[linkedin.com/in/clhking](https://www.linkedin.com/in/clhking)

Flash Storage

LinkedIn wrote over 475 billion collected metrics to flash yesterday, at a rate of three million writes per second.

Vision

Create economic opportunity
for every member of the global
workforce



The Economic Graph



MEMBERS



COMPANIES



JOBS



SKILLS



SCHOOLS



KNOWLEDGE

Mission

Connect the world's professionals
to make them more productive
and successful

Growing global network



380M+¹
Members worldwide

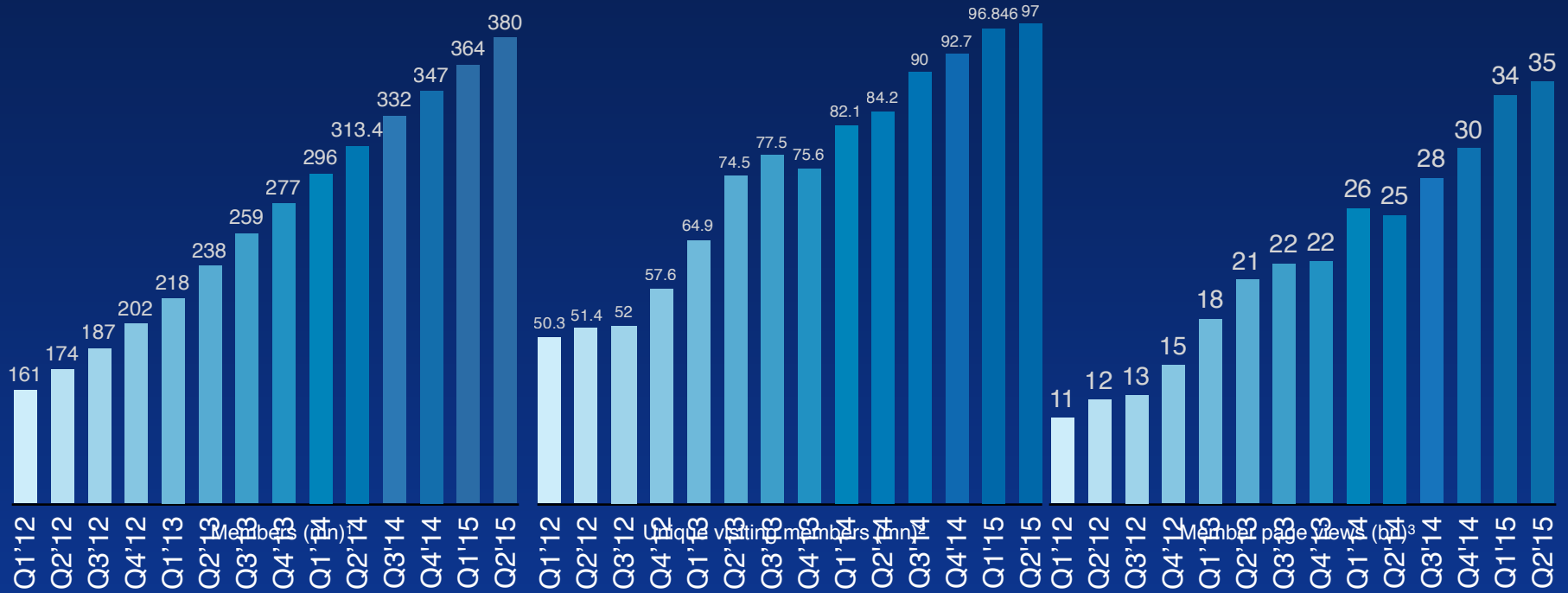


>2 New¹
Members per second



35B²
unique visiting members

1 As of end of 2Q15 | 2 monthly average during the quarter



- Flash makes the difference
 - Improve engagement through site speed
 - 300% improvement in commit to disk latency
 - More compact deployment

- Host Based PCIe Flash
 - >1PB deployed and increasing
 - 2x or 3x replication
 - Sizes between 650GB and 4.8TB

- SAN/Array based flash
 - <1PB deployed and declining
 - Large RDBMS, Corp-IT Virtualization
- NAS based flash
 - Metadata caching only

- Host Based Flash
 - >7500 systems (end of 2010, <50 systems)
 - Average Latency <100 μ s
- SAN Based Flash
 - <500 systems (end of 2010, <200 systems)
 - Average latency >800 μ s
- NAS Based Flash
 - <50 systems (end of 2010, 2 systems)

- Limited set of manufacturers
- Proprietary drivers
- Poor share-ability, if available
- Heat dissipation
- NOT hot swap

- Drivers
 - Built into all mainstream / updated Operating Systems; that was easy

- Vendors
 - HGST, Samsung, Seagate, Toshiba/OCZ, Intel, Novachips, SanDisk ...

- Server A has available space; we overestimated.
 - NVMe can be shared
 - We've tested Excelero and it's very promising
 - Mangstor has an interesting offering

- NVMe comes in 2.5" SFF
 - It gets cooled first
 - Multiple slots
 - The datacenter guys, love it.
 - Current time to move/replace a flash card ~1hr
 - Servers already shipping with 2,10,24 slots
 - 24 slots is a bit silly, but 76TB @ 50 μ s is cool

- Deploy NVMe in every rack, not per server SKU
- Thin-provision by software or process
- Share between systems of compatible workload
- Mirror data off failing/expiring drives
- Source from as many vendors as can meet performance/reliability testing

Where to from here?

- NVMe now means shareability later
- Test and evaluate the many emerging products
- Narrow down fabric and sharing technology
 - $2 \leq x \leq$ “too many to manage properly”
- Wait for 25Gbit/50Gbit switch silicon to GA

Q/A