

Powering the future of Cloud Service Providers An inside view from Systems Research

Dr. Robert Haas

Department Head,
Cloud & Computing Infrastructure
IBM Research – Zurich

rha@zurich.ibm.com

<https://www.linkedin.com/in/rohaas/>

https://twitter.com/robert_r_haas

Please Note

IBM's statements regarding its plans, directions, and intent are subject to change or withdrawal without notice and at IBM's sole discretion.

Information regarding potential future products is intended to outline our general product direction and it should not be relied on in making a purchasing decision.

The information mentioned regarding potential future products is not a commitment, promise, or legal obligation to deliver any material, code or functionality. Information about potential future products may not be incorporated into any contract.

The development, release, and timing of any future features or functionality described for our products remains at our sole discretion.

Performance is based on measurements and projections using standard IBM benchmarks in a controlled environment. The actual throughput or performance that any user will experience will vary depending upon many factors, including considerations such as the amount of multiprogramming in the user's job stream, the I/O configuration, the storage configuration, and the workload processed. Therefore, no assurance can be given that an individual user will achieve results similar to those stated here.

D

Compute

Move

Store

A

1) Storing Data

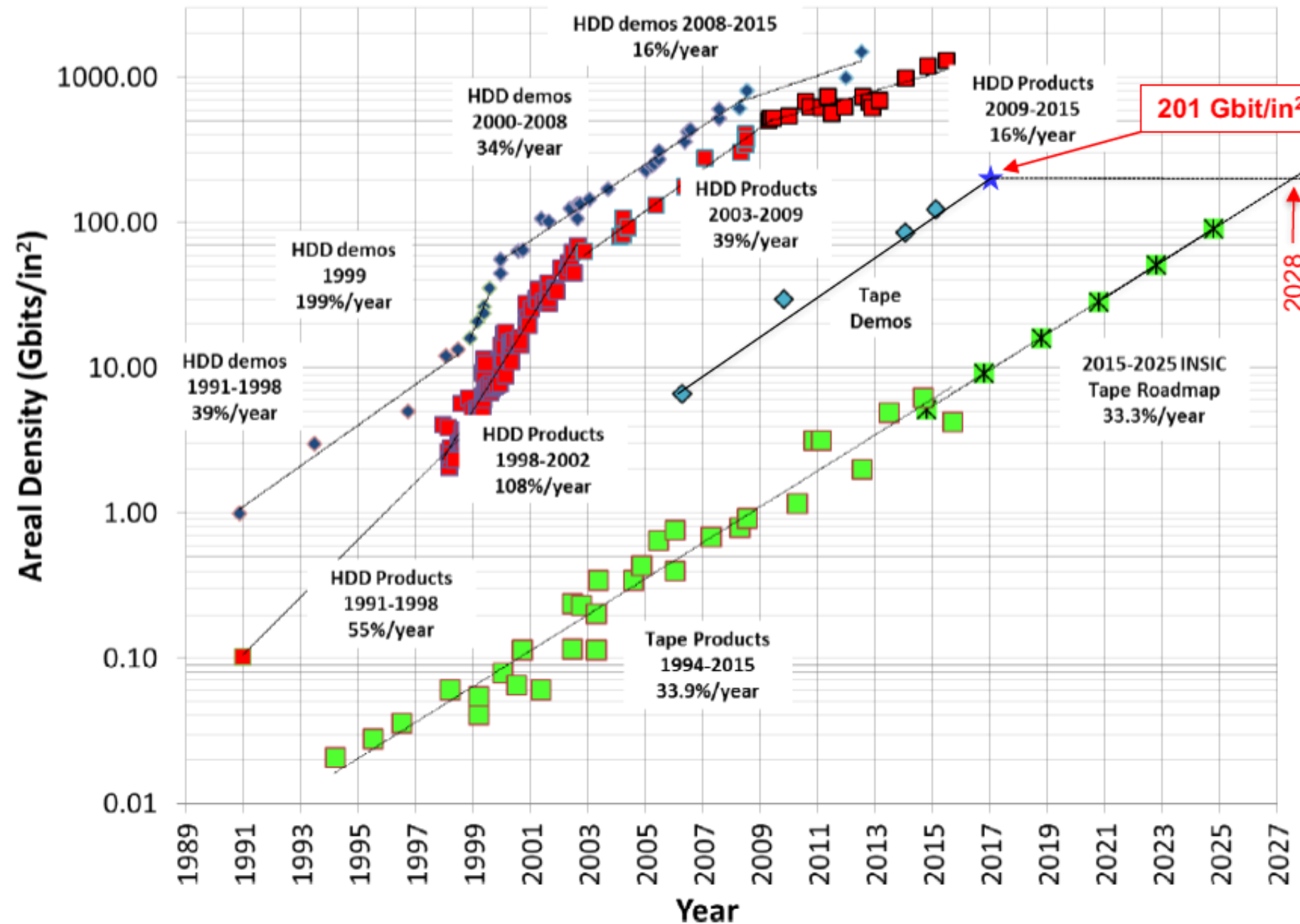
Pushing the limits of density in tape and Flash to lower costs



Magnetic Recording Areal Density Trends

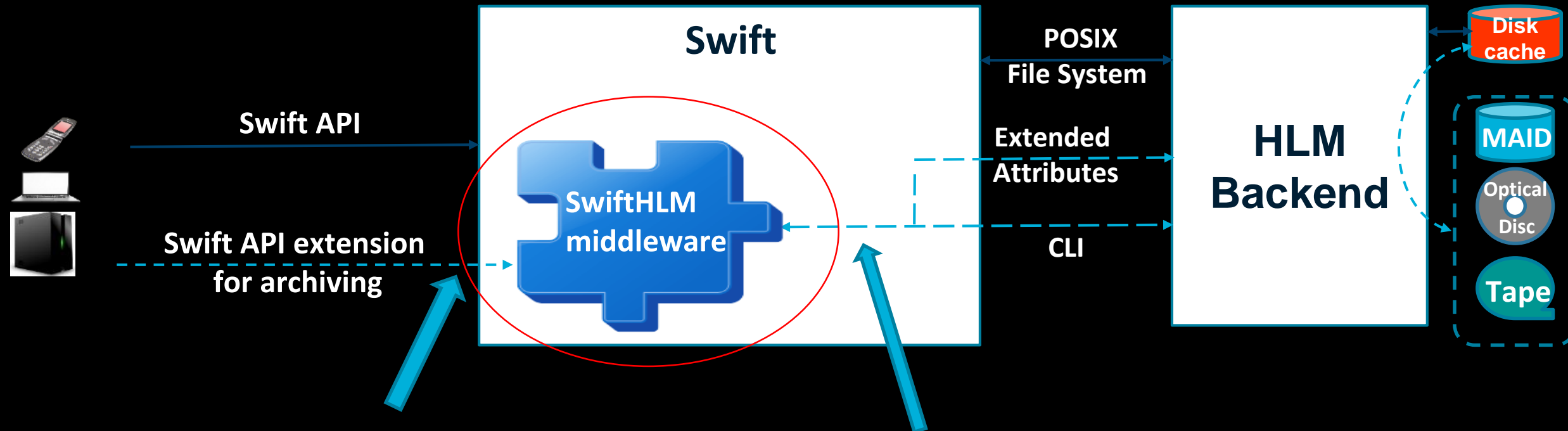
2015: IBM-FujiFilm demonstration of 123 Gb/in² on BaFe tape

2017: IBM-Sony demonstration of 201 Gb/in² on Sputtered Tape



Tape will maintain its 5x to 10x cost advantage over other storage technologies for at least another decade.

Introducing the Object API for High Latency Media (HLM)



Object API extension for HLM archiving:

- Migrate (Disk -> High-Latency Media, async)
- Recall (High-Latency Media -> Disk, async)
- Query status for Object (sync)
- Query status for Request (sync)
- Object and container level operations

Generic interface for HLM backends, suitable for e.g.:

- IBM Spectrum Archive (LTFS Enterprise Edition)
- IBM Spectrum Protect (TSM/HSM)
- BDT Tape Library Connector (open source)

AI for Data Access Prediction

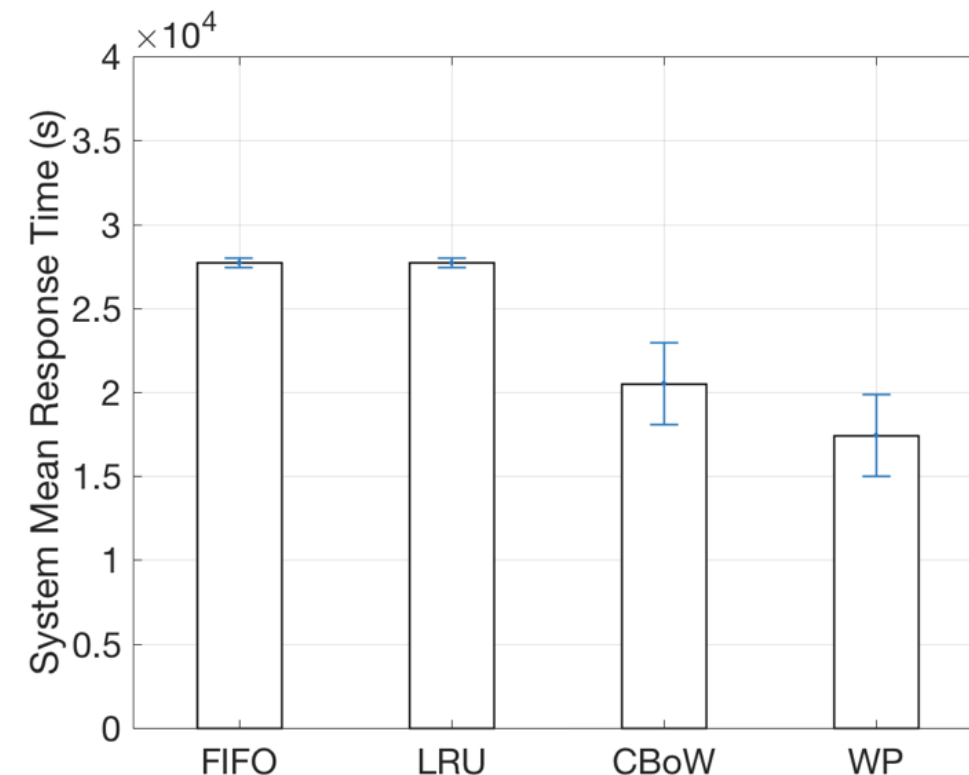
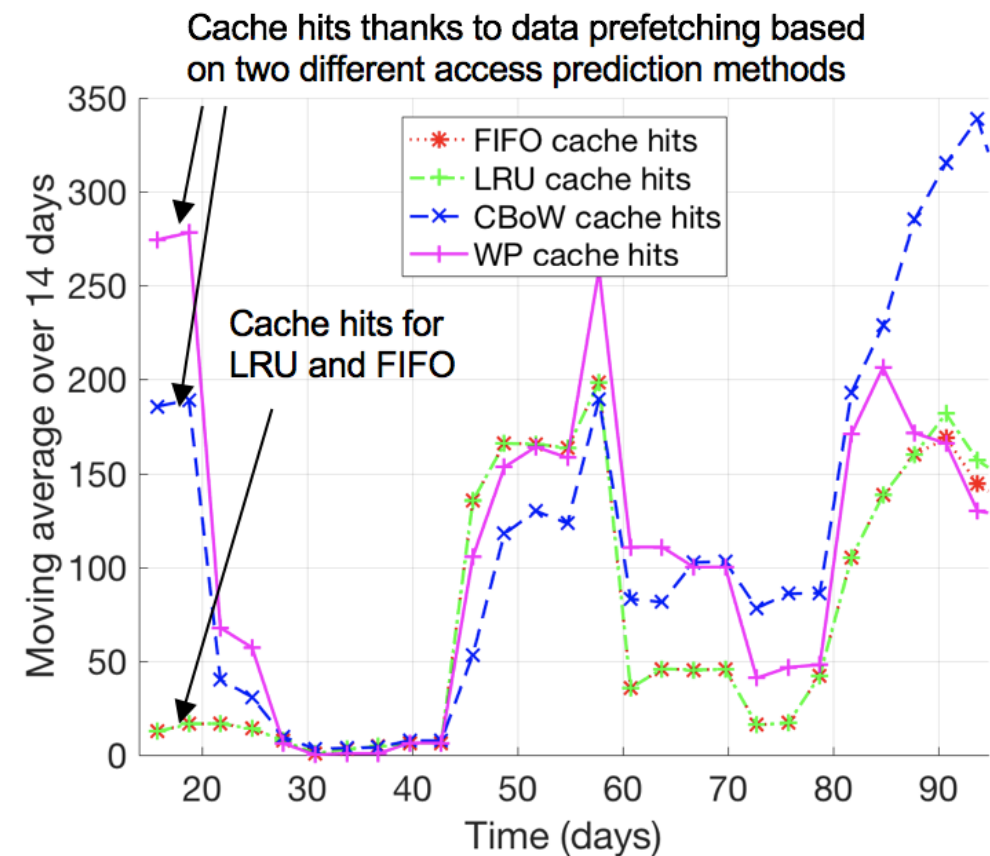
Double cache hit rate and reduce latency by half by predictive pre-fetching data using a disk cache and a tape backend

Evaluated using logs from ASTRON Long Term Archive staging server

Compared to conventional caching algorithms such as FIFO and LRU

CBoW: conditional bag-of-words model is used when computing meta-data based similarity of files

WP: word-pair model is used when computing meta-data based similarity of files



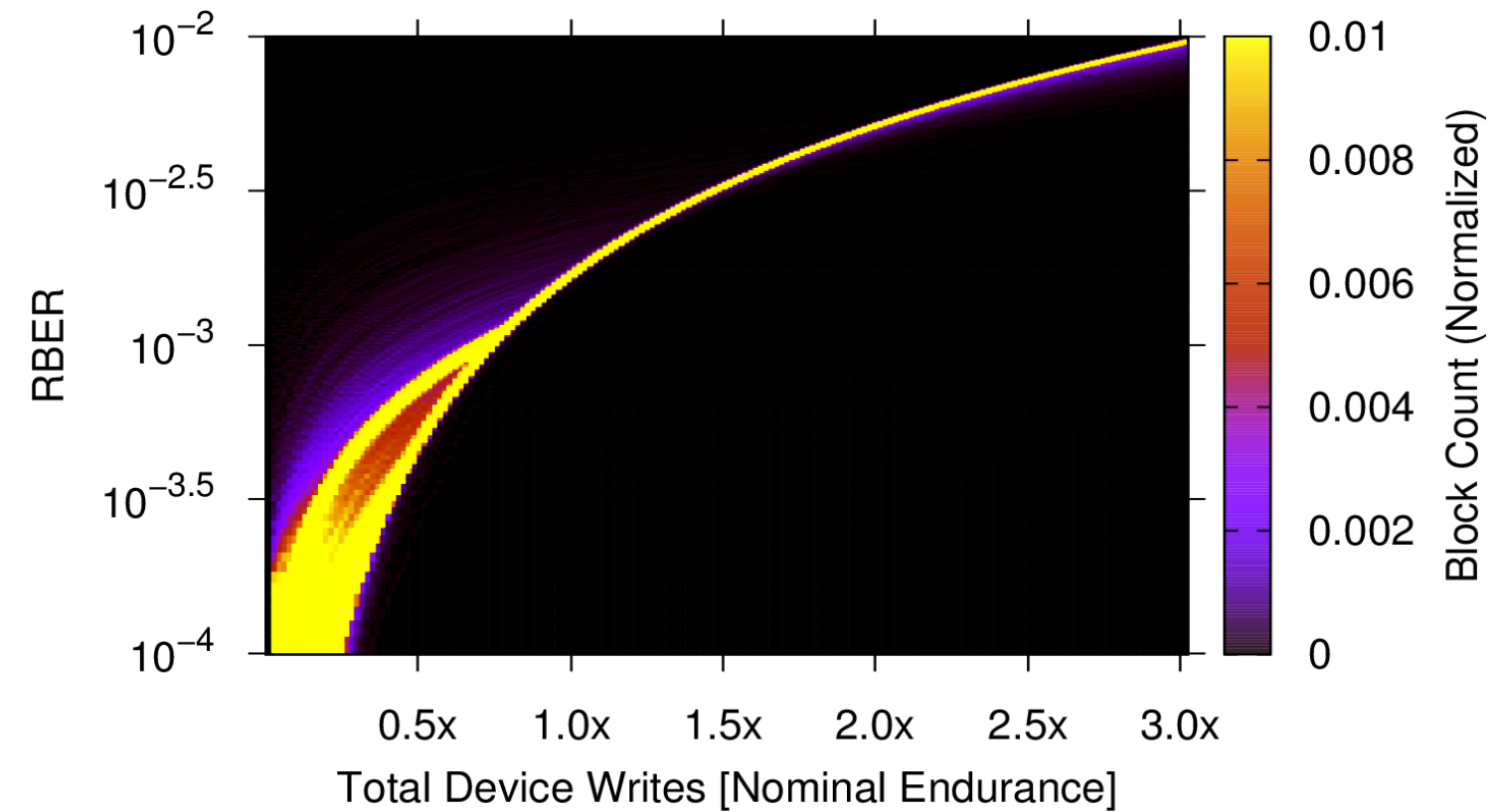
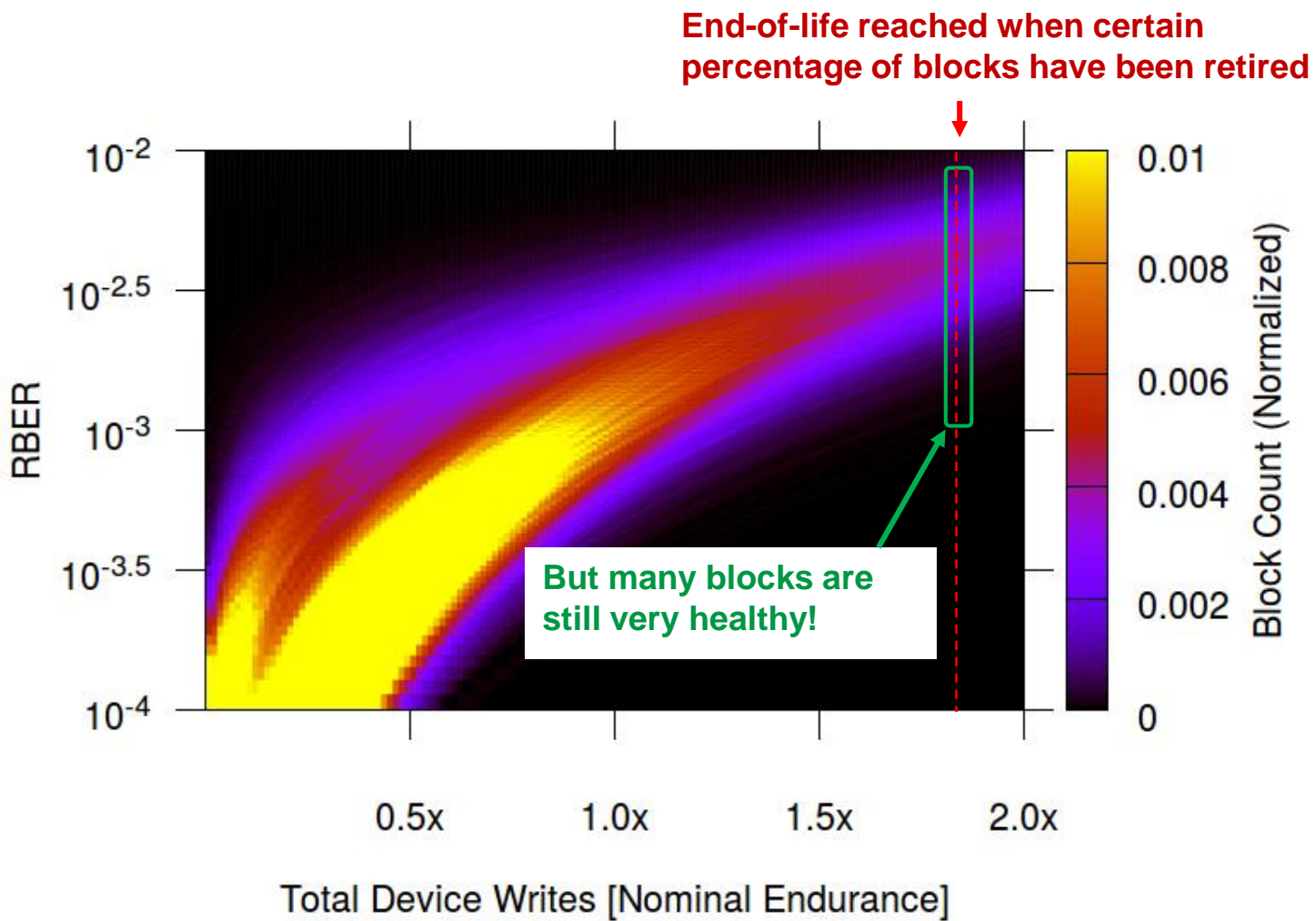
Innovation to Leverage the Densest Flash: 3D-TLC and beyond



Secret Sauce to Squeeze More Write Cycles from Flash

Standard approach
(wear leveling using write-cycle balancing)

Secret Sauce using health binning
and taking advantage of workload skew



2) Moving Data

*Pushing the limits in terms of speed and lowering power needs,
and letting applications profit from advances in networking and storage performance*

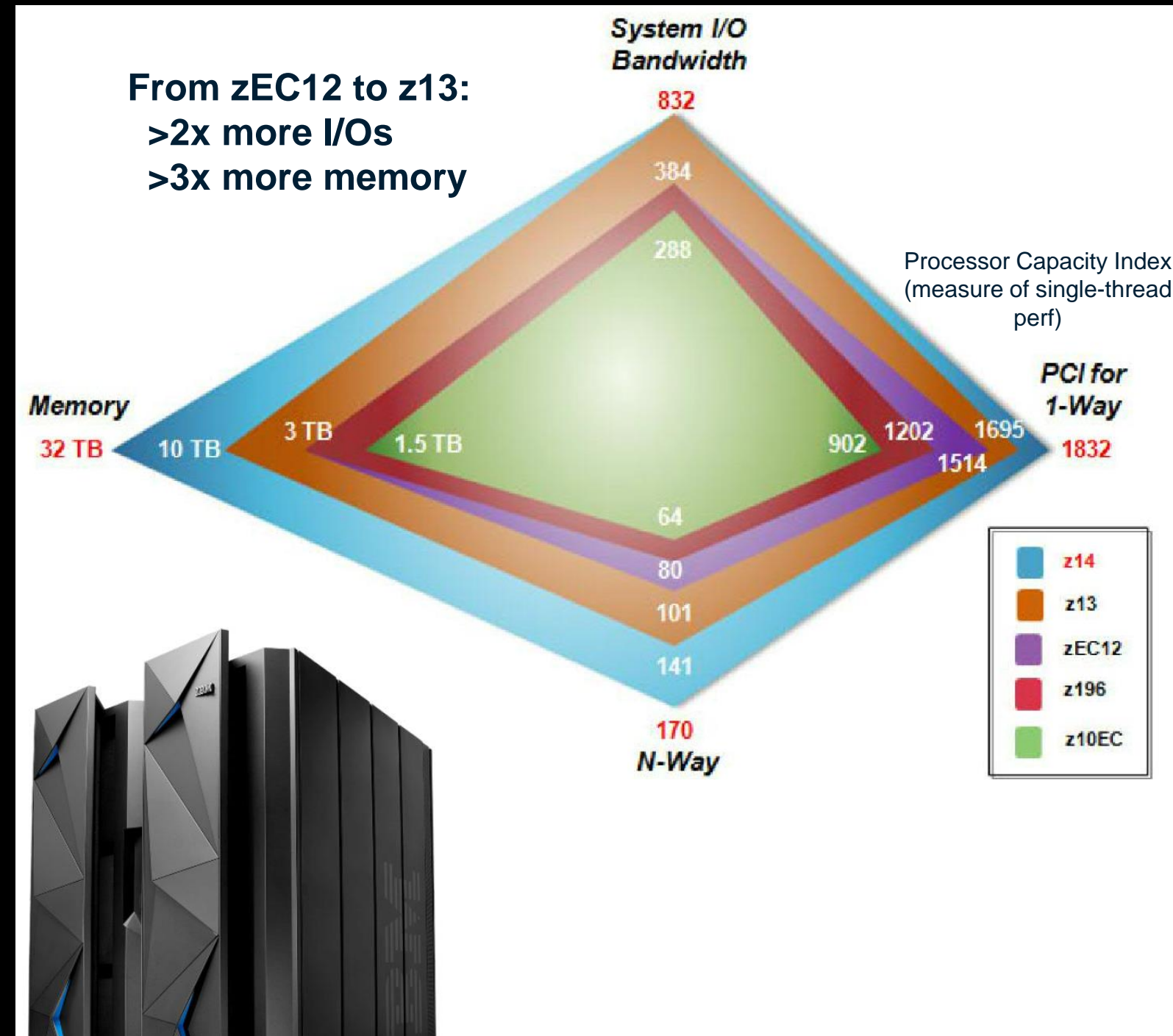
Fast and Low-power I/O Links to Feed Data-hungry CPUs and GPUs

Growing discrepancy between computation and communication

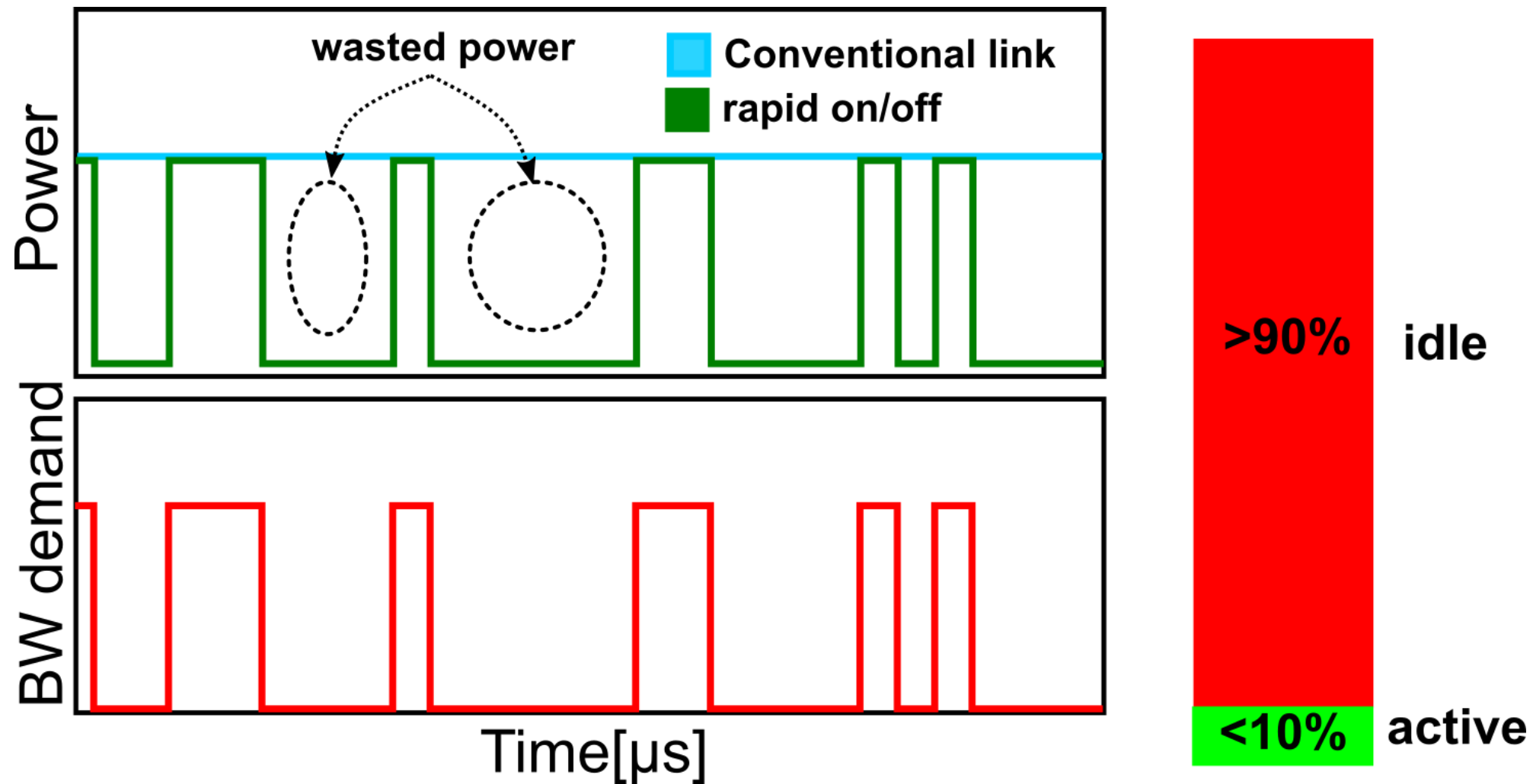
32bit MULT operation = 1pJ
transferring 64bits = 320 pJ

More cores → Need more I/O bandwidth at lower power to fit into thermal envelope

→ I/O links are key enabler for system performance



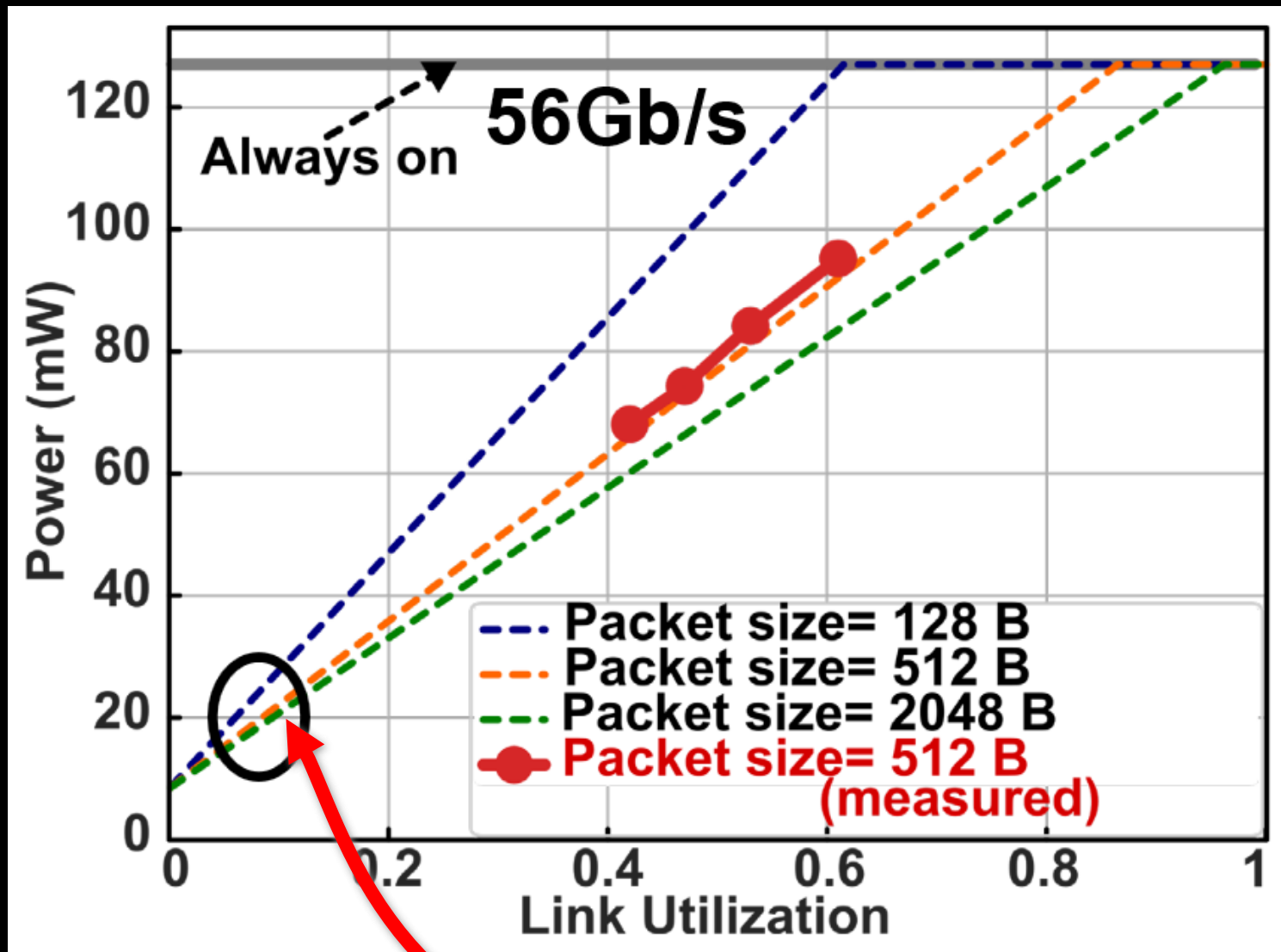
Wasted Power in Idle Links



Link utilization in a large data center <10% [1]: idle links burn power.

Rapid on/off: Power scales with link utilization. Enables more bandwidth per chip at fixed Thermal Dissipation Power

Rapid Power-On Links Leads to 85% Less Power Consumption

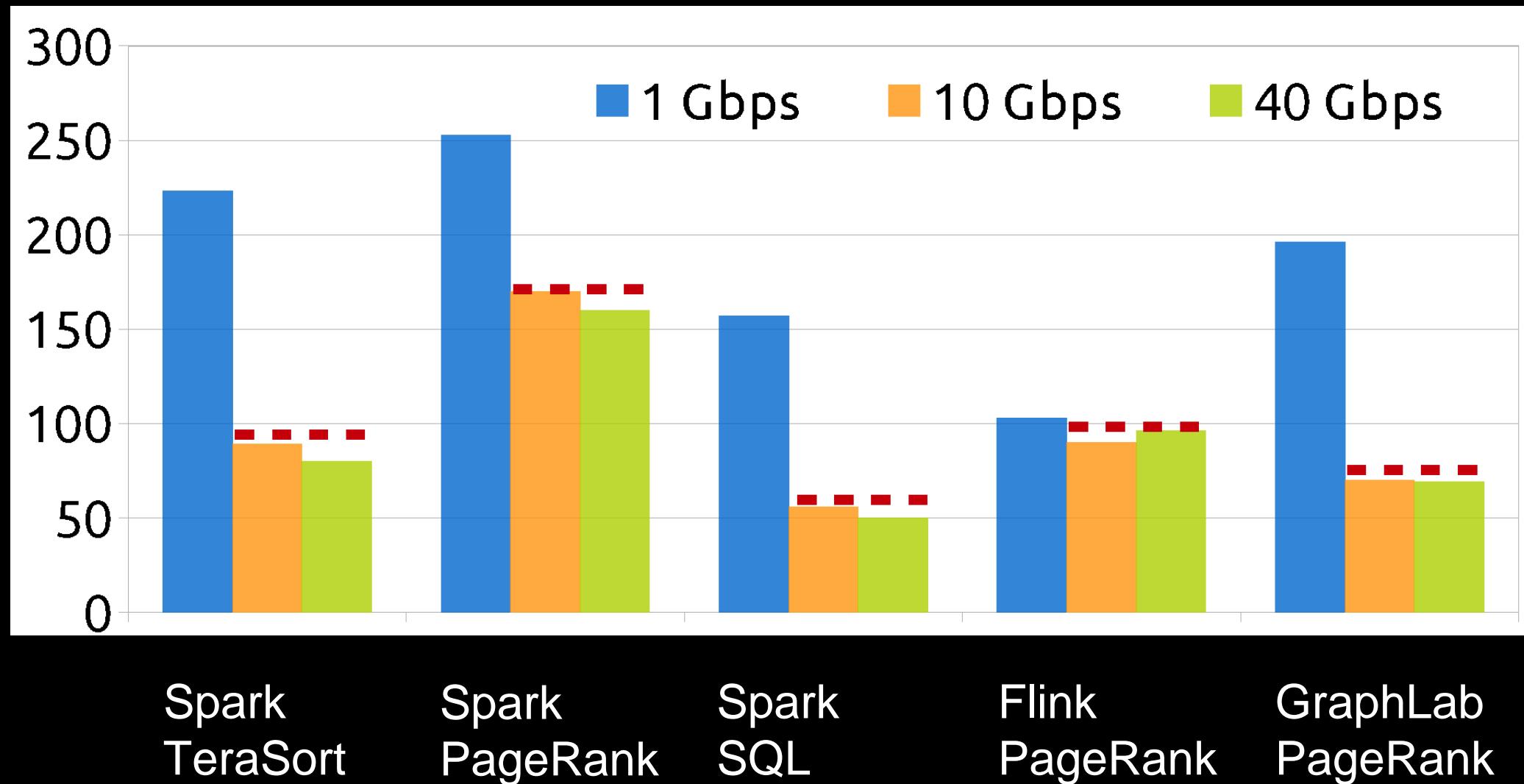


ON power: 128mW (2.1pJ/b) @ 60Gb/s

OFF power: 8mW

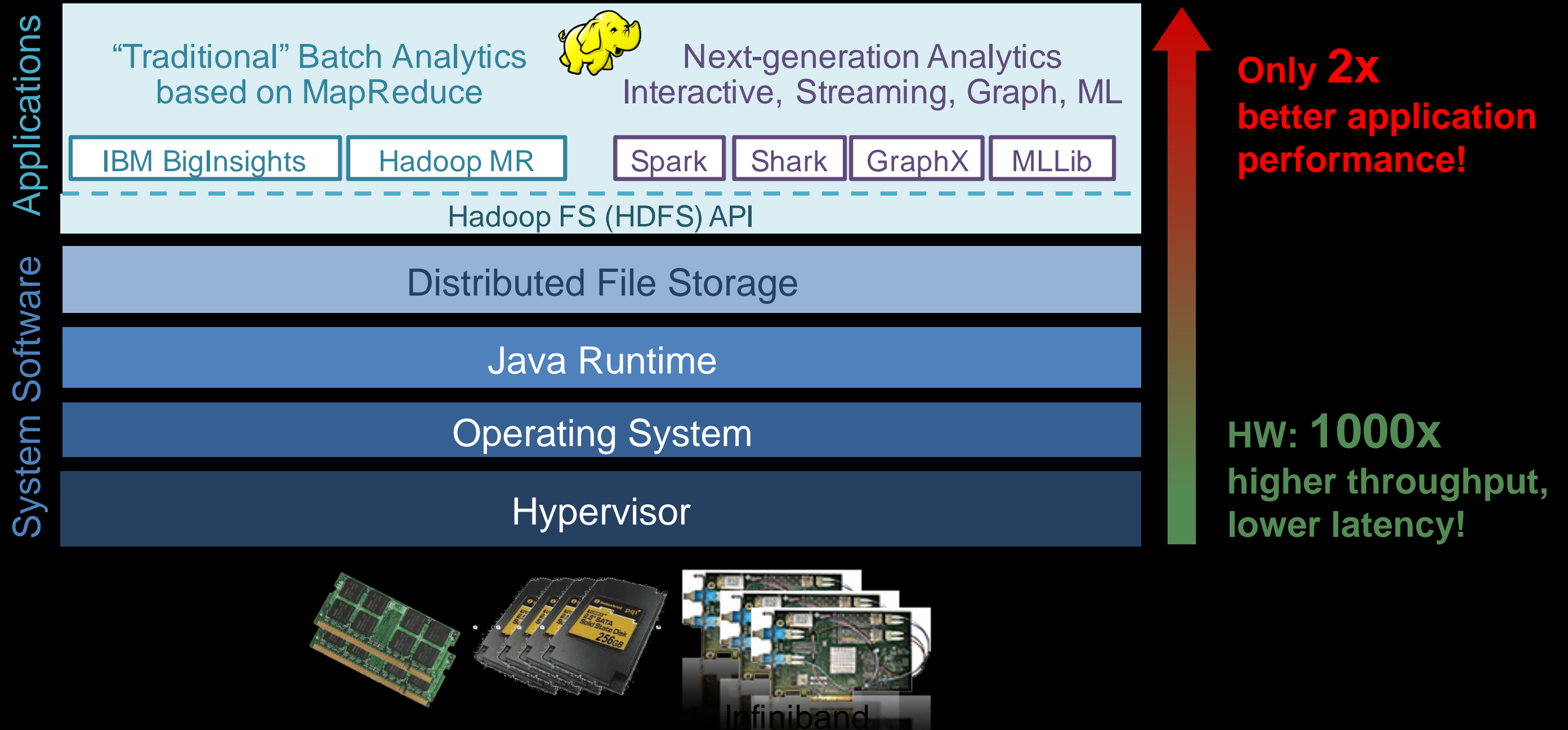
At 10% link utilization: 20mW
with only 7ns to power on!

Analytics Systems Challenged by High-Performance Networks



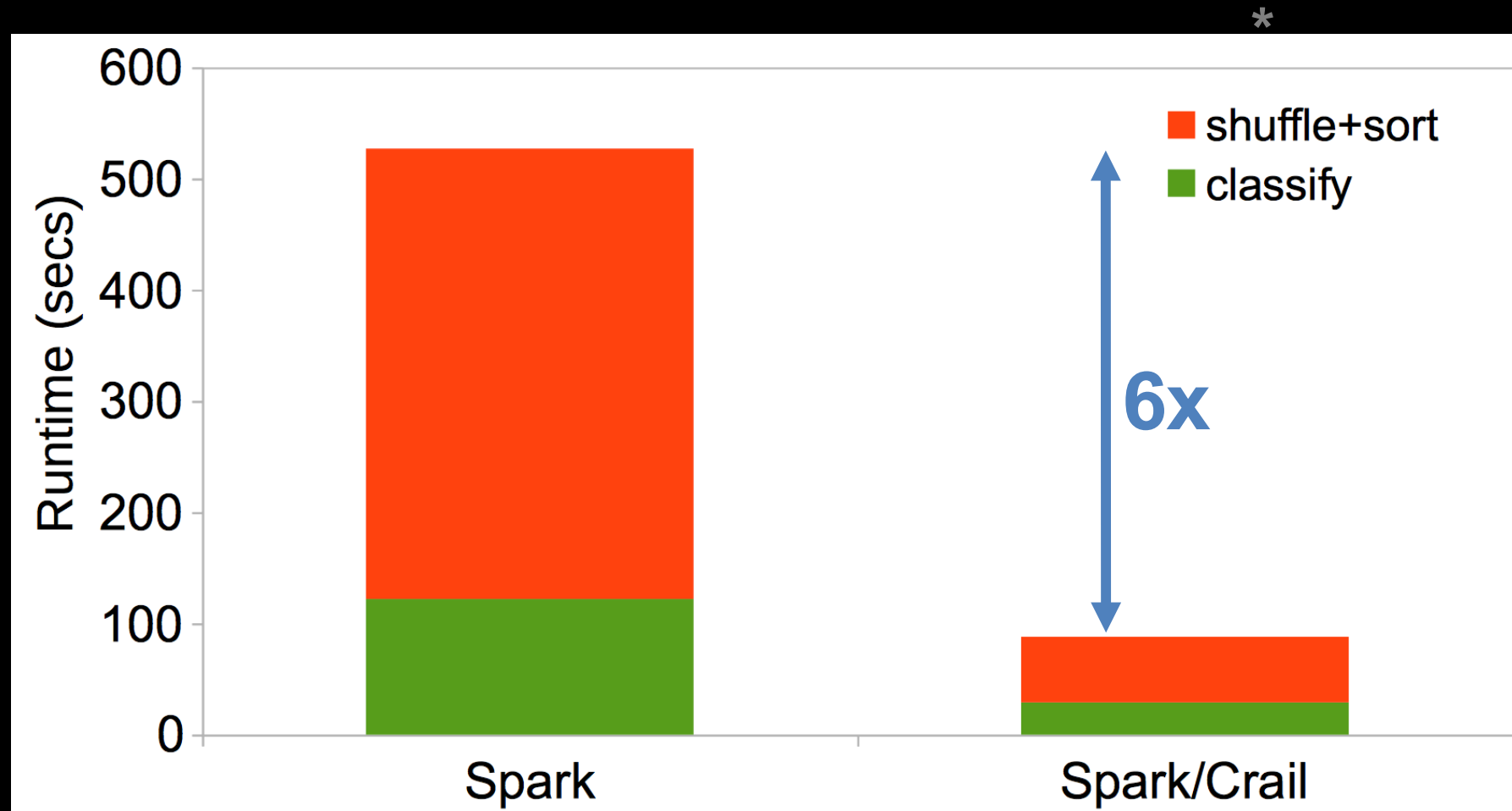
Trivedi et al. "On The [Ir]relevance of Network Performance for Data Processing", Usenix HotCloud 2016

Hardware Performance Does Not Trickle Up the Stack



TeraSort on Spark

TeraSort on OpenPOWER cluster



OpenPOWER cluster

128 nodes

Node configuration

2 x IBM POWER8 10-core @ 2.9 Ghz

DRAM: 512GB DDR4

4 x 1.2 TB NVMe SSD

100GbE Mellanox ConnectX-4 EN (RoCE)

Software

Ubuntu 16.04 (kernel 4.4.0-31)

Spark 2.0.0

Find out more at <http://crail.io/blog>

3) Computing on Data

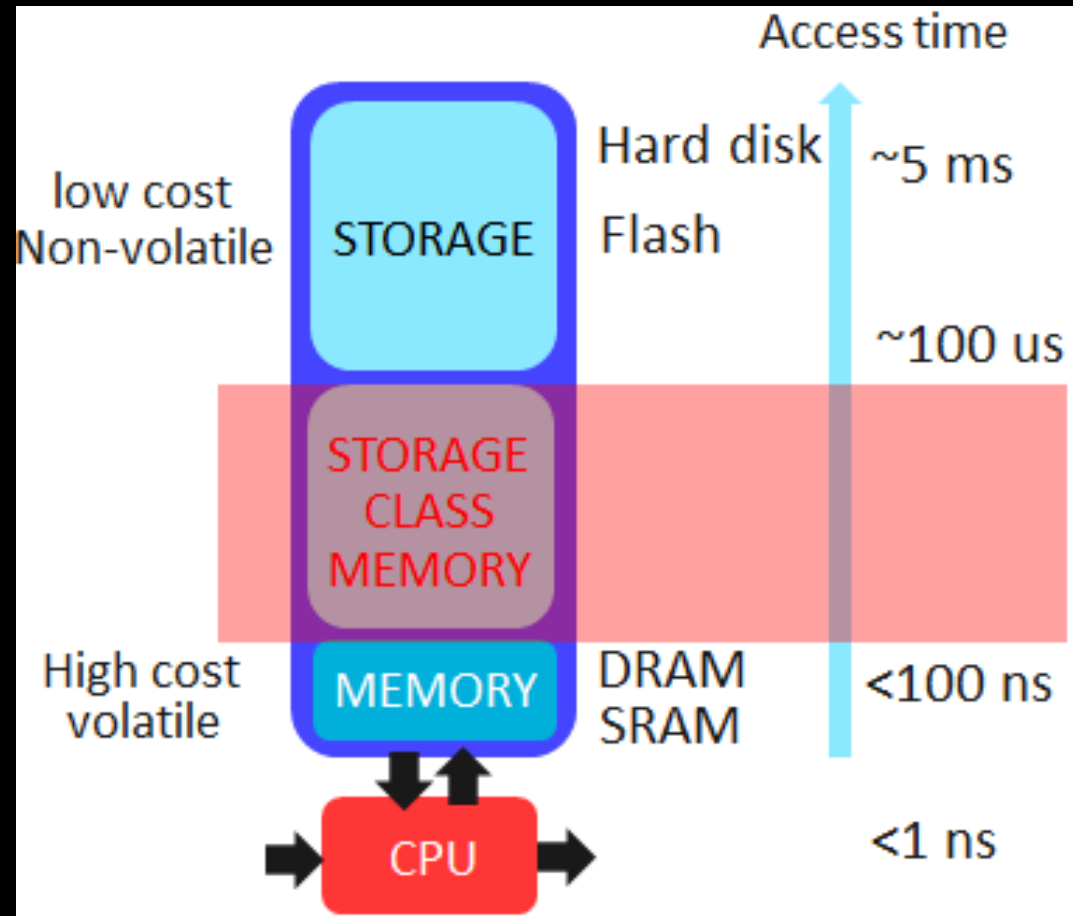
Revolutionize systems with in-place computing

Business As Usual: Continue to Improve von-Neumann Computing

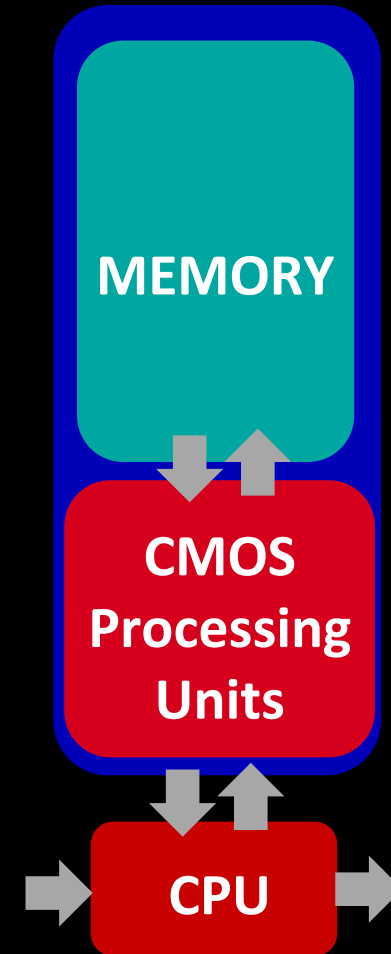
Storage-class memory

Near-memory computing

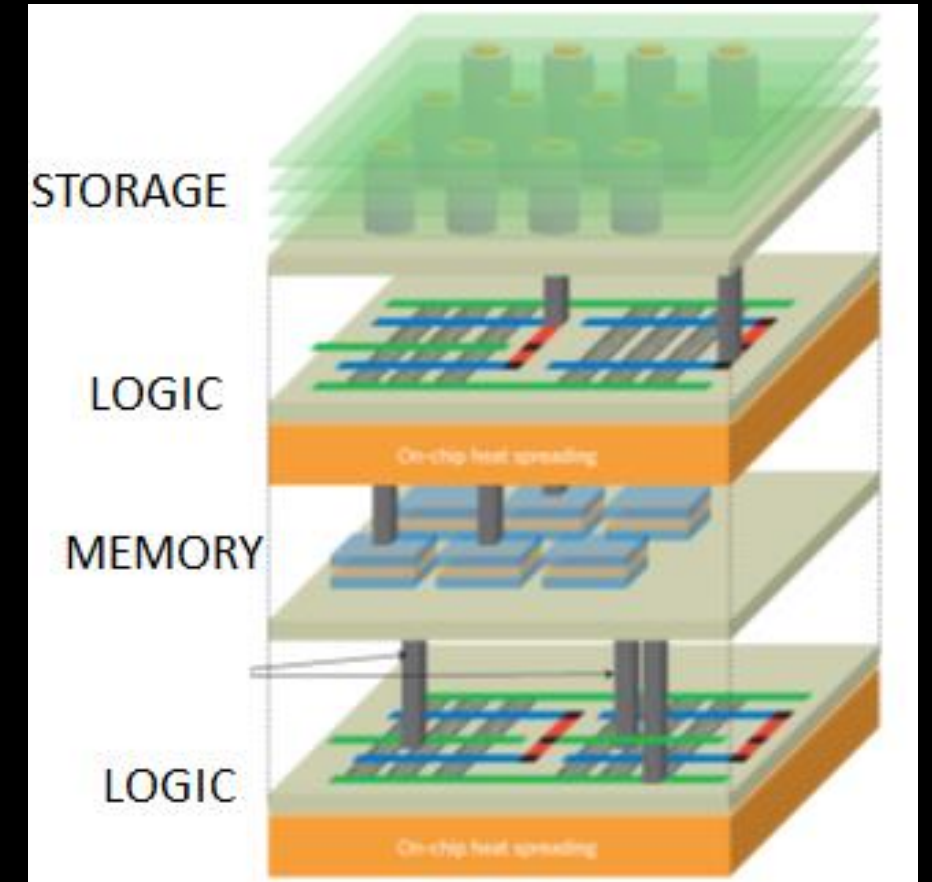
Monolithic 3D integration



Burr et al., IBM J. Res. Dev., 2008



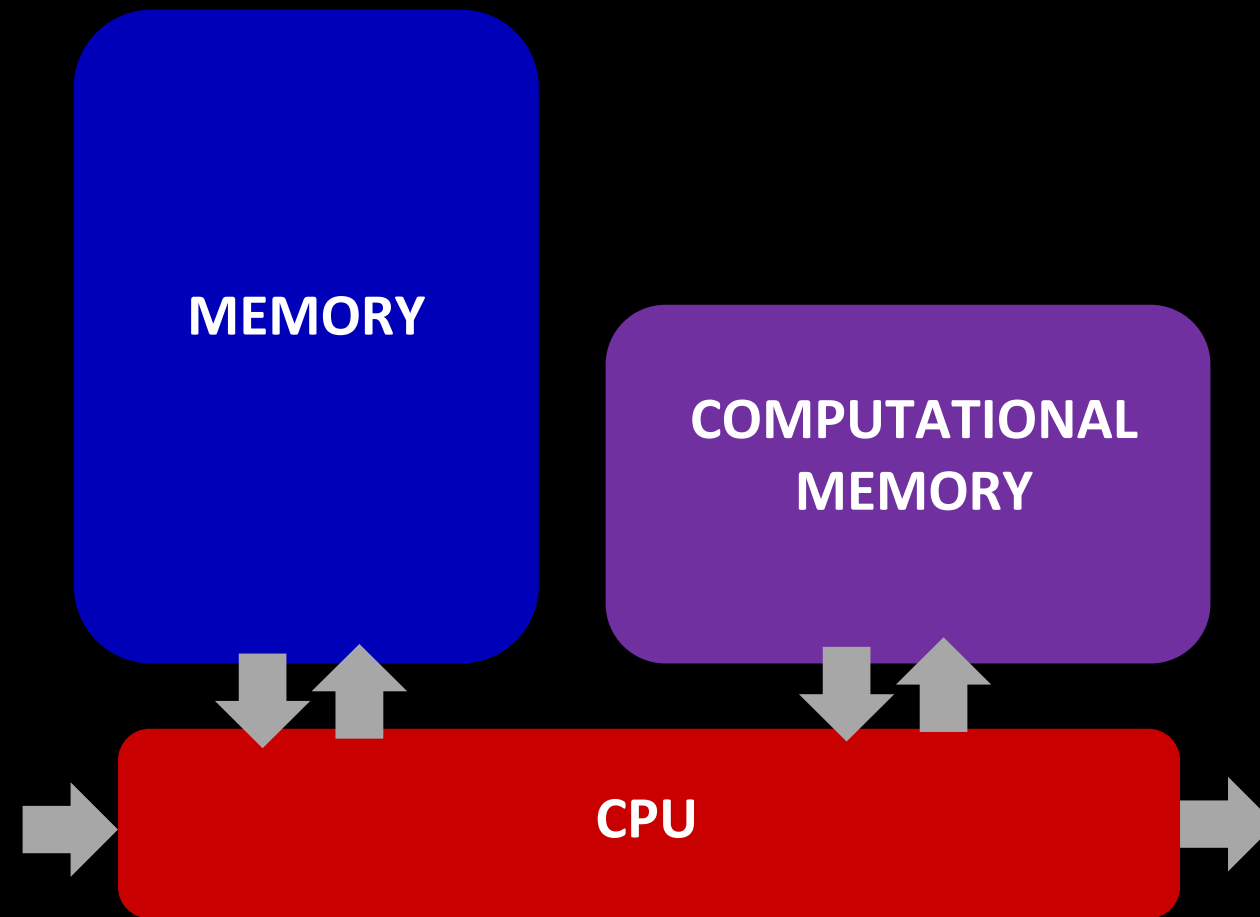
Vermij et al., Proc. ACM CF, 2016



Wong, Salahuddin, Nature Nano., 2015

Minimize the time and distance to memory access

Introducing Game-Changing Computational Memory



Perform “certain” computational tasks in place in memory
Not only stores data but performs some calculations on the data

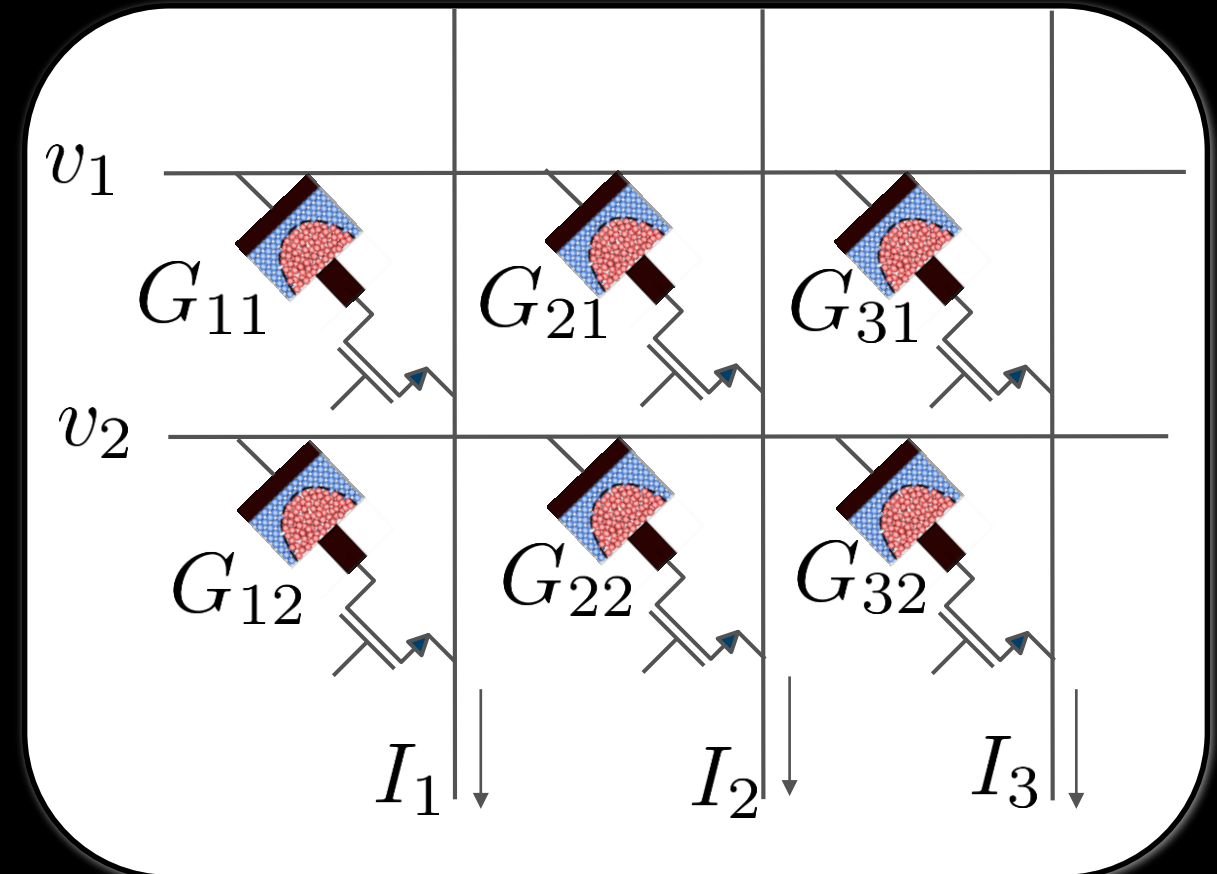
Doing Matrix-vector Multiplication in Computational Memory

$$\begin{bmatrix} M_{11} & M_{12} \\ M_{21} & M_{22} \\ M_{31} & M_{32} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix}$$

**MAP to
conductance
values**

**MAP to read
voltage**

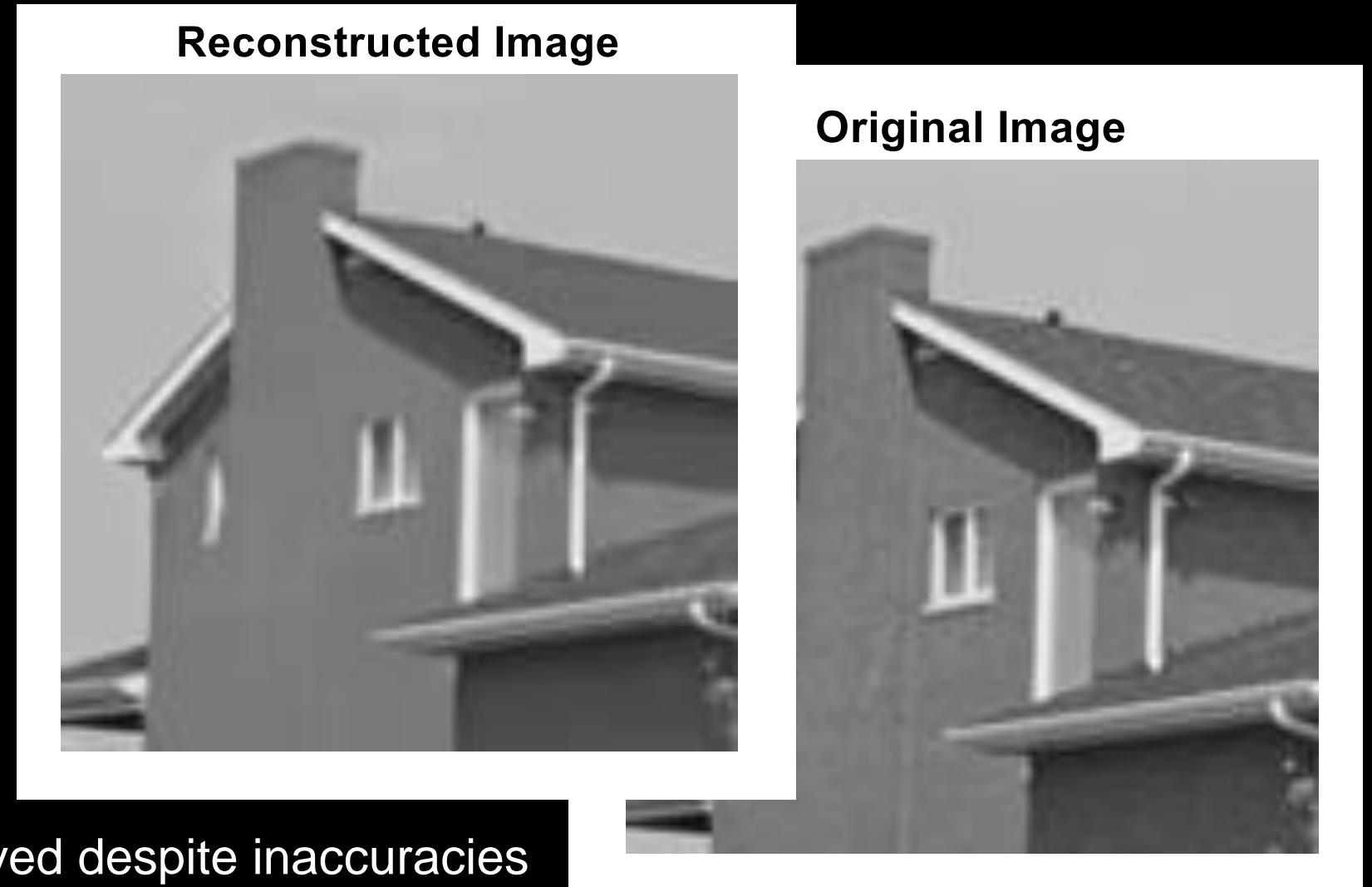
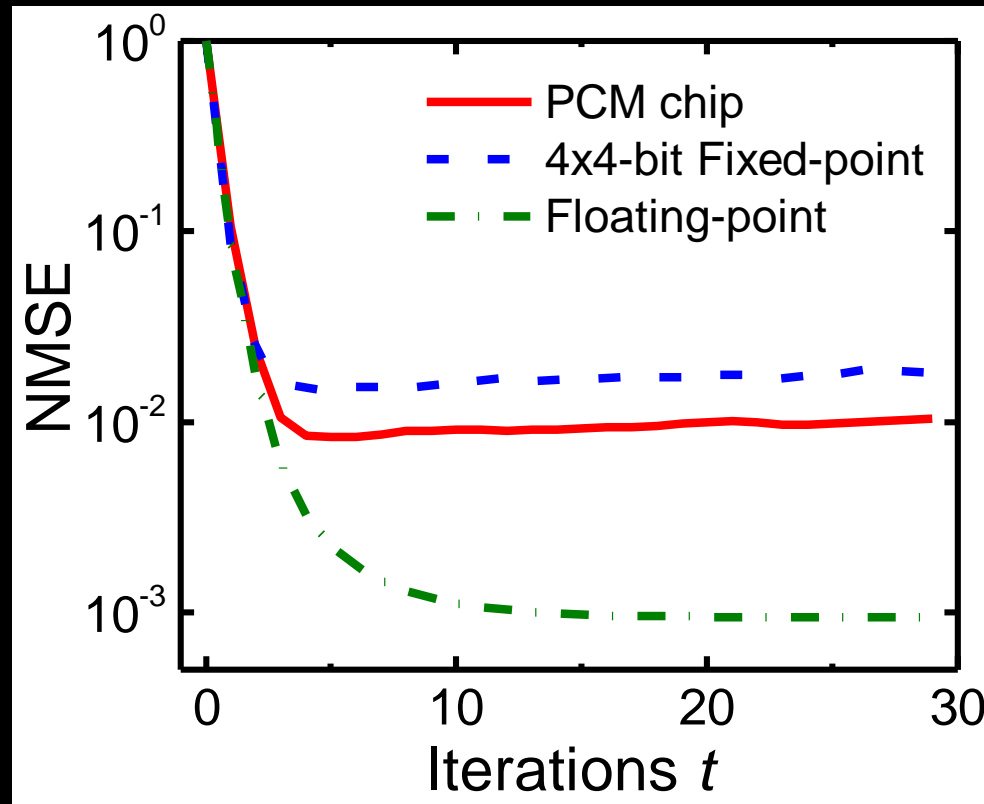
**DECIPHER from
the current**



By arranging the PCM devices in a cross-bar configuration one can perform matrix-vector operation with $O(1)$ complexity
Exploits multi-level storage capability and Kirchhoff's circuits laws

Image Reconstruction Experimental Results

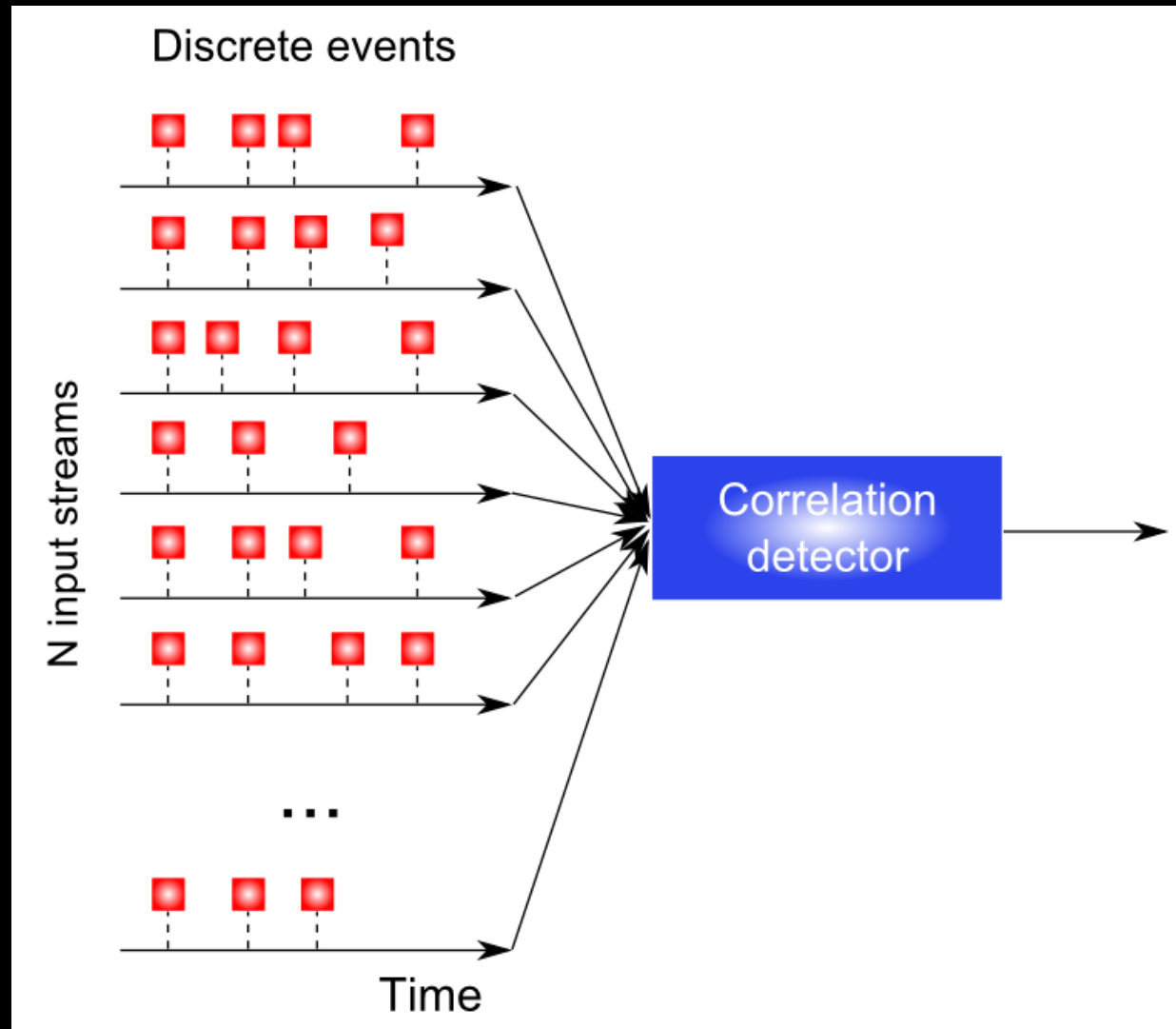
128x128 image, 50% sampling rate,
Comp. mem. unit w/ 131,072 PCM devices



Reasonable reconstruction accuracy achieved despite inaccuracies

Estimated power reduction of 50x compared to using an optimized 4-bit FPGA matrix-vector multiplier that delivers same reconstruction accuracy at same speed

Temporal Correlation Detection



Determine whether some of the input data streams are statistically correlated
Use only unsupervised learning & consume very low power

FINANCE



SCIENCE



MEDICINE



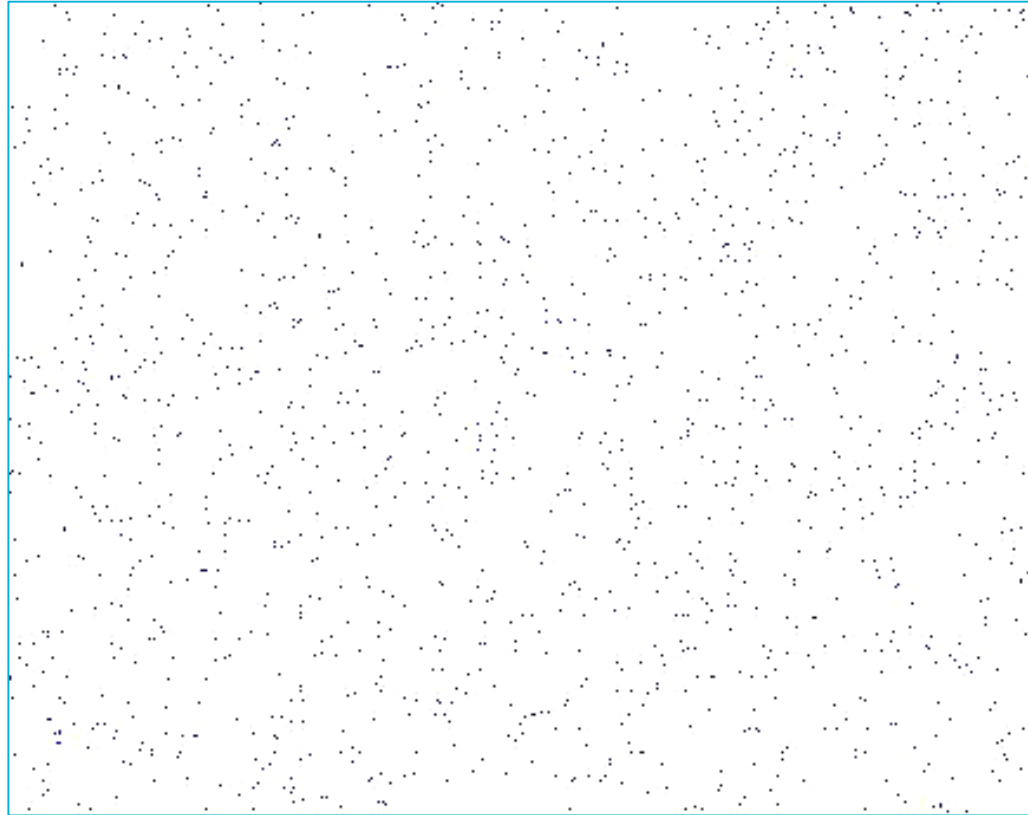
BIG DATA



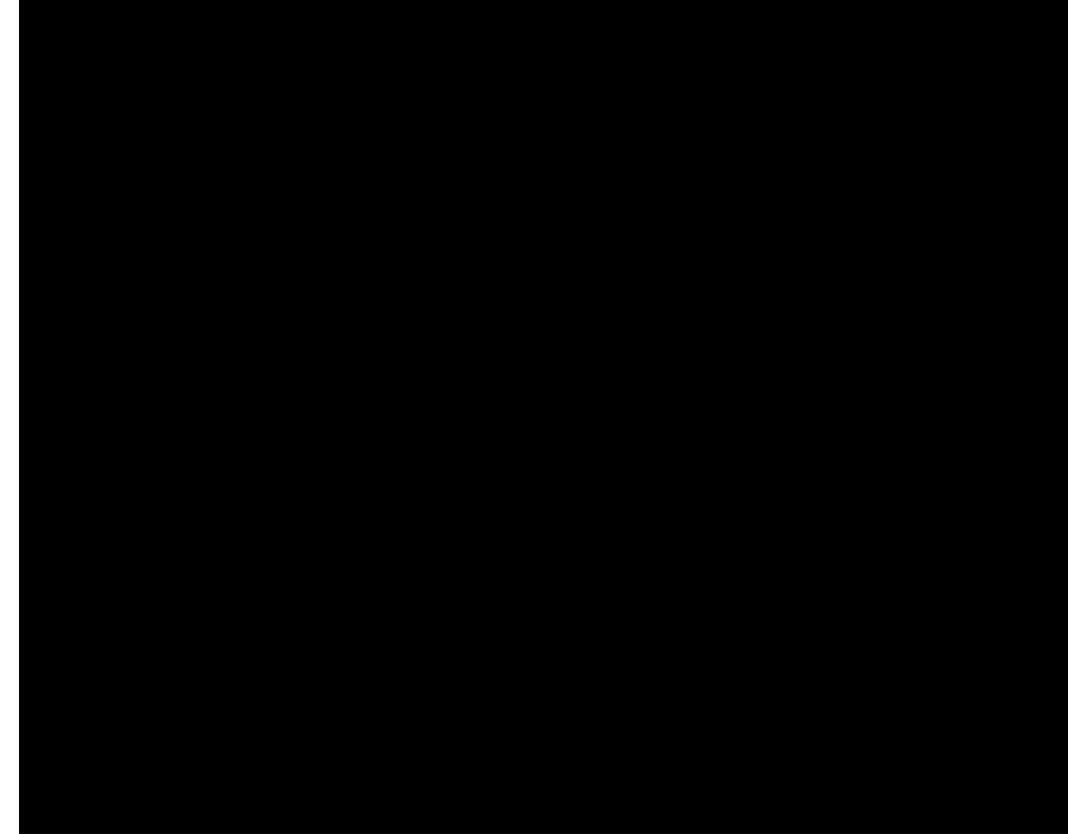
...AND MORE

Experimental Results with 1 Million PCM devices

Processes



Device conductance



Sebastian et al., Nature Communications, 2017

Very weak correlation of $c = 0.01$
No shuttling back and forth of data
Massively parallel
Unprecedented areal/power efficiency

Compute in Place

Move Way Faster

Store Much Denser

Thank You

1 Visit us @ Booth #A09

1-on-1 meetings

Personalized solution demos

2 Explore IBM Systems solutions

Visit: ibm.com/systems

3 MSP Hub – discover helpful resources

Visit: www.themsphub.com



Notices and Disclaimers

Copyright © 2018 by International Business Machines Corporation (IBM). No part of this document may be reproduced or transmitted in any form without written permission from IBM.

U.S. Government Users Restricted Rights - Use, duplication or disclosure restricted by GSA ADP Schedule Contract with IBM.

Information in these presentations (including information relating to products that have not yet been announced by IBM) has been reviewed for accuracy as of the date of initial publication and could include unintentional technical or typographical errors. IBM shall have no responsibility to update this information. THIS DOCUMENT IS DISTRIBUTED "AS IS" WITHOUT ANY WARRANTY, EITHER EXPRESS OR IMPLIED. IN NO EVENT SHALL IBM BE LIABLE FOR ANY DAMAGE ARISING FROM THE USE OF THIS INFORMATION, INCLUDING BUT NOT LIMITED TO, LOSS OF DATA, BUSINESS INTERRUPTION, LOSS OF PROFIT OR LOSS OF OPPORTUNITY. IBM products and services are warranted according to the terms and conditions of the agreements under which they are provided.

IBM products are manufactured from new parts or new and used parts. In some cases, a product may not be new and may have been previously installed. Regardless, our warranty terms apply.”

Any statements regarding IBM's future direction, intent or product plans are subject to change or withdrawal without notice.

Performance data contained herein was generally obtained in a controlled, isolated environments. Customer examples are presented as illustrations of how those customers have used IBM products and the results they may have achieved. Actual performance, cost, savings or other results in other operating environments may vary.

References in this document to IBM products, programs, or services does not imply that IBM intends to make such products, programs or services available in all countries in which IBM operates or does business.

Workshops, sessions and associated materials may have been prepared by independent session speakers, and do not necessarily reflect the views of IBM. All materials and discussions are provided for informational purposes only, and are neither intended to, nor shall constitute legal or other guidance or advice to any individual participant or their specific situation.

It is the customer's responsibility to insure its own compliance with legal requirements and to obtain advice of competent legal counsel as to the identification and interpretation of any relevant laws and regulatory requirements that may affect the customer's business and any actions the customer may need to take to comply with such laws. IBM does not provide legal advice or represent or warrant that its services or products will ensure that the customer is in compliance with any law

Notices and Disclaimers Con't.

Information concerning non-IBM products was obtained from the suppliers of those products, their published announcements or other publicly available sources. IBM has not tested those products in connection with this publication and cannot confirm the accuracy of performance, compatibility or any other claims related to non-IBM products. Questions on the capabilities of non-IBM products should be addressed to the suppliers of those products. IBM does not warrant the quality of any third-party products, or the ability of any such third-party products to interoperate with IBM's products. IBM EXPRESSLY DISCLAIMS ALL WARRANTIES, EXPRESSED OR IMPLIED, INCLUDING BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE.

The provision of the information contained herein is not intended to, and does not, grant any right or license under any IBM patents, copyrights, trademarks or other intellectual property right.

IBM, the IBM logo, ibm.com, Aspera®, Bluemix, Blueworks Live, CICS, Clearcase, Cognos®, DOORS®, Emptoris®, Enterprise Document Management System™, FASP®, FileNet®, Global Business Services®, Global Technology Services®, IBM ExperienceOne™, IBM SmartCloud®, IBM Social Business®, Information on Demand, ILOG, Maximo®, MQIntegrator®, MQSeries®, Netcool®, OMEGAMON, OpenPower, PureAnalytics™, PureApplication®, pureCluster™, PureCoverage®, PureData®, PureExperience®, PureFlex®, pureQuery®, pureScale®, PureSystems®, QRadar®, Rational®, Rhapsody®, Smarter Commerce®, SoDA, SPSS, Sterling Commerce®, StoredIQ, Tealeaf®, Tivoli®, Trusteer®, Unica®, urban{code}®, Watson, WebSphere®, Worklight®, X-Force® and System z® Z/OS, are trademarks of International Business Machines Corporation, registered in many jurisdictions worldwide. Other product and service names might be trademarks of IBM or other companies. A current list of IBM trademarks is available on the Web at "Copyright and trademark information" at: www.ibm.com/legal/copytrade.shtml.