### VISION OF PROCESSOR-MEMORY SYSTEMS

#### MICRO 48 J. Thomas Pawlowski, Chief Technologist, Fellow jpawlowski@micron.com

©2015 Micron Technology, Inc. All rights reserved. Information, products, and/or specifications are subject to change without notice. All information is provided on an "AS IS" basis without warranties of any kind. Statements regarding products, including regarding their features, availability, functionality, or compatibility, are provided for informational purposes only and do not modify the warranty, if any, applicable to any product. Drawings may not be to scale. Micron, the Micron logo, and all other Micron trademarks are the property of Micron Technology, Inc. All other trademarks are the property of their respective owners.





- Memory diversity
- Some scaling observations
- Unabated performance pressure
- Memory Technologies
- Emerging Memory

- 3D XPoint<sup>™</sup> Memory
- Changing Relationships
- Micron Automata Processor
- Abstraction and Future Systems
- 11 Key Points



# **Memory Diversity**



### **Diversified Memory Markets**



Server

Mobile

**Personal Computing** 

IMM



### **Smartphones and Tablets**

#### **Micron Portfolio**

#### 4-32Gb LPDDR4

• Multiple package options



- 4x provides 4GB density
- LPDDR4 2x bandwidth LPDDR3
- LPDDR4 35% power savings over LPDDR3

#### eMMC 5.x

- Small package options
- 3x random Read/Write performance
- 2x sequential Read/Write performance
- Wide temperature range tolerance and high reliability

#### eMCP + LPDDR3

- Attractive package/density options
- 2x bandwidth over eMCP+LP2
- Quick time-to-market for system qualification efforts



### **Enterprise and Cloud**

#### **Micron Portfolio**

#### **Server Modules**

- Up to 128GB module densities
- DDR4 solutions driving higher performance
- Load Reduced options for higher system capacity

#### **Enterprise SSD**

- Enabled with high performance MLC NAND
- Expanding portfolio of SAS and PCIe
- Partnerships enable efficient go to market

#### NVDIMM

- Integrated DRAM/NAND Module
- Outstanding latency and bandwidth
- Excellent endurance capability

#### **Hybrid Memory Cube**

- Best-in-class bandwidth solution
- 70% less energy per bit
- Industry leading reliability

#### **3D XPoint**

- First new memory technology in over 25 years
- 1,000x faster than NAND
- 10x higher density than conventional memory











Micro

### Networking

#### Micron Portfolio

#### **DRAM Modules**

- Error correction for data integrity
- VLP and custom form factors
- Leading-edge and legacy options

#### RLDRAM

- Low latency
- Fast random access
- Specialized x9, x18, and x36 configurations

#### Hybrid Memory Cube

- Best-in-class bandwidth solution
- Provides the bandwidth of as many as 144 DDR3 components

#### eUSB

- Firmware and OS management
- User tracking data



### Client

#### **Micron Portfolio**

#### **On-Board Memory**

- Up to 16Gb per package
- Various package options
- Soldered-on-board for thin profile
- DDR4, DDR3L and LPDDR3/4

#### **Client Modules**

- Up to 16GB density
- Thin form factor
- 1.35V DDR3L at 10% lower power than standard DDR3
- DDR4 memory

#### **Serial NOR**

- Dependable BIOS
- Secure Boot

#### **Client SSD**

- Vs. HDD: Superior \$/IOPS and IOPS/watt
- Faster boot and bandwidth
- Lower dB, ns, W, oz, cu.in.
- Encryption, security
- 100x performance











ons

## gurations



# **Some Scaling Observations**



## **DRAM Array and Logic** Transistors per square mm

- Suggests using memory wafers more extensively than logic, widening gap
- Logic transistor curve based only on Intel production start dates to maintain harmony of methodologies Dates rounded down to year – <u>Note that this</u> <u>does not show COST</u>



Source: JTPawlowski, Micron



# Most Dense Use of Logic Technology is Memory

	22nm	14nm		
6T HDC SRAM	Intel 190.1 F <sup>2</sup>	Intel 14.5Mb/mm <sup>2</sup> 254.6 F <sup>2</sup>		
6T LVC SRAM	IBM dens-perf 297.5 F <sup>2</sup>	Intel 300 F <sup>2</sup>		
6T HPC SRAM		Intel 360.2 F <sup>2</sup>		
8T SRAM		Intel 5.1Mb/ mm <sup>2</sup>		
eDRAM	IBM 53.7 F <sup>2</sup> Intel Macro 17.5Mb/mm <sup>2</sup> 59.9 F <sup>2</sup>	IBM 89.1 F <sup>2</sup> IBM Macro 18.7Mb/mm <sup>2</sup>		
Sources verieus ISSCC and VISI Surenessium Dressedings				

Source: various ISSCC and VLSI Symposium Proceedings

- Compare to DRAM now constant 6-7F<sup>2</sup>
- SRAM scaling poorly
- eDRAM appears worse watch future scaling Some density and energy benefit for now
- DRAM is challenging but scaling far better than SRAM and eDRAM





### Intel HC SRAM and IBM eDRAM Cell Scaling Area [nm<sup>2</sup>] vs. Node [nm]

- SRAM 2x relative size since 90nm, eDRAM 2.7x since 45nm, macros scaling worse
- Contrast: smallest nFET available in 14nm: Intel 25F<sup>2</sup> Micron 4F<sup>2</sup>





### Memory Growing As Percentage Area in Logic Chips



Source: Semico

Much of this area can be displaced by cheaper process technology

© 2015 Micron Technology, Inc.



# **Unabated Performance Pressure**





- Source: ISSCC 2015 Trends Logic chips only
- Memory transistor count dwarfs logic (8Gb DRAM ~10B, latest NAND ~400B)





• Source: ISSCC 2015 Trends

### Single Core IPC Gain vs. Introduction Year for Intel Xeon Family







- ISSCC 2015 Trends, plus: 2014 Cavium 48 core, 2015 EZChip/Tilera 100 core, Knights Landing ~72 Atom cores, Phytium 64 ARM cores.
- Off the chart: Kalray and REX 256 cores, GPUs e.g. Nvidia TitanX 3,360 cores



### Impressive Gains in Memory Throughput System Throughput a.u. vs. Time



© 2015 Micron Technology, Inc.

December 8, 2015



# **Memory Technologies**



# In-Package Memory





### The Innovative Hybrid Memory Cube

#### **Revolutionary Approach**

- Evolutionary DRAM roadmaps hit limitations of bandwidth and power efficiency
- Micron introduces a new class of memory: Hybrid Memory Cube
- Unique combination of DRAMs on Logic addresses the memory wall, scalability and RAS

#### **Key Applications**

- Data packet processing, data packet buffering, and storage applications
- Enterprise and high performance computing applications

#### **Unparalleled Performance**

- Provides 15X the bandwidth of a DDR3 module
- Uses 70% less energy per bit than existing memory technologies
- Reduces the memory footprint by nearly 90% compared to today's RDIMMs

#### How did we do it?

- Micron-designed logic controller
- High speed link to CPU
- Massively parallel "Through Silicon Via" connection to DRAM



# Impacts of DRAM Process Complexity

- Large increase in number of process steps to enable shrink
- Conversion Capital Expenditure scales with number of steps
- Significant reduction in wafer output per existing cleanroom area

#### Complexity comparison for enablement of ~100% bits/wafer increase







December 8, 2015

## 3D NAND vs. Planar NAND Scaling



- Planar NAND scaling
  - Planar can be scaled below 16nm, but performance and cost are not competitive with 3D NAND
  - Micron focused 100% on 3D NAND after 16nm
- 3D NAND scaling
  - 3D NAND cost/capacity improvement over planar expands with subsequent nodes
  - 3D NAND cell architecture enables significant performance improvement relative to planar technology

#### **Capacity Projection**





20nm

Node

16nm

1024

256

64

16

Electrons per level

90nm



# **Future Memory Technologies**

- Strategic investment in future memory roadmap enablement
- DRAM or DRAM replacement scaling
- Persistent Memory enablement
- Multiple generations of 3D NAND
- Strategic investment in future memory core technologies
- Resistive Random Access Memories (RRAM)
- Spin Torque Transfer Random Access Memories (STTRAM)
- And others...



#### 16Gb High Speed Persistent Memory





STTRAM Array

Advanced 27nm RRAM Cell



# **Emerging Memory**



## Memory / Storage Replacement Candidates

- Spin Torque Transfer RAM is public leading candidate for DRAM
- All candidates problematic characteristics when formulated for memory
  - Greater latency tRC energy/bit
  - Lesser bandwidth endurance
  - Much greater Raw Bit Error Rate
- Unwieldy fit into existing sockets
- NAND has no realistic direct replacement candidate
  - Everything is projected to have higher cost until NAND has no scaling path
- Numerous emerging memory candidates
  - Many fall between NAND and DRAM





- SSD pushing HDD into cold data / archival role
- Latency gap between DRAM and SSD
- Fill this gap with a non-volatile component
- Persistent Memory a new class of NVM (non-volatile memory)
- Note 1: not exclusively the style from a memory company, includes eDRAM, eSTTRAM



# Nonvolatile Memories in Server Architectures





# PCM: The Ideal Future Memory?

Characteristic	DRAM	РСМ	NAND
Nonvolatile	No	Yes	Yes
Latency	Low, Fixed	> DRAM, Fixed	High, Variable
Read/Write Time	1:1	1:5	1:20
Block Management	No	No	Yes
Refresh	Yes	No	No
Addressability	Byte	Byte	Pages
Raw Bit Error Rate	Low	Low	High
Endurance	High	Medium	Low
Energy per Bit	Low	High	High
Cost per Bit	DRAM	~DRAM	Low

- Desire the best attributes of both NAND and DRAM
- PCM's cost per bit was too close to DRAM to justify the higher latency and reduced endurance
- But there is a range of memory possibilities which will provide compelling attributes for the right cost per bit



# **3D XPoint™ Memory**



### 2015 is a Historic Year in Memory





### Introducing 3D XPoint<sup>™</sup> Memory





## **Nonvolatile Memories in Server Architectures**





- 3D XPoint<sup>™</sup> technology provides the benefit in the middle
- It is considerably faster than NAND Flash
- Performance can be realized on PCIe or DDR buses
- Lower cost per bit than DRAM while being considerably more dense



December 8, 2015

# **Changing Relationships**



## The Great Wall is Energy



Task 256b ops	Energy (Without sequencing overhead)			
Two Double Precision Floating Point Operations	15 pJ (10nm logic)			
Small L1 Cache SRAM Read	30 pJ (10nm logic)			
10mm move on logic die	180 pJ (10nm logic)			
Low-Power discrete DRAM off-chip read	800pJ (same timeframe as 10nm logic)			
Source: JTPawlowski, Micron				

- Numerous walls exist in systems
- Memory throughput wall overcome by HMC and RLDRAM
- The largest system challenge is energy
  - Control overhead and Data movement
- Newest workloads: high volume random data



# The 3<sup>rd</sup> Epoch of Processor-Memory Relationship

### Epoch

- 1 In the early compute days Processors were cheap compared with memory
- 2 Memory became much cheaper than processors Software mostly compiled Memory usage soared
- **3** Now <u>using</u> the memory is becoming more expensive than <u>using</u> processors: energy!

Portends...Increasing memory area to mitigate energy And the necessity for new fundamental solutions

**SHIFT** from 60 years of processor-centric compute to memory-centric will be the feature of 3<sup>rd</sup> Epoch



#### Core memory board

Early cores were 1/16" in diameter, could be accessed in 5 microseconds, and cost about \$1 per bit. Eventually they were 4 times smaller, 10 times faster, and 100 times cheaper.



# System Architectures Will Change

- Resource redistribution
- Eventually no memory device will be directly controlled by the host processor without an appropriately-tuned abstraction layer
  - Very light for the lowest capacity, lowest latency memory
  - Highly-abstracted for the highest capacity, highest latency memory
- More processing will be done using memory technology, chiefly as accelerators
  - Also much more near-memory and near-storage processing
- New mentality of "work moving"
  - Migrate control information to the epicenter of data and process there



### **Enormous Untapped Potential in Memories**

- Huge parallelism, tiny little off-chip data pipe
- Costly to provide external access to the full capability
- Great potential for a variety of new accelerators
- And...memory tech is scaling more favorably than logic tech

#### Standard DDR3 DRAM

30ns to load 1 row into sense amps (tRAS) and stagger fire every 4ns

Internal Bandwidth: 4 Ti bits/sec

External Bandwidth: 25.6 Gb/s on 16b

160x untapped bandwidth (320x on 8b)



#### With Numerous Improvements

Can improve internal bandwidth by another 1000x

160,000x – 320,000x untapped bandwidth potential



© 2015 Micron Technology, Inc.



**AUT•MATA** 

PROCESSING

# **Micron Automata Processor**



### Automata Processor

- Micron's Automata Processor is a revolutionary new class of programmable accelerator
  - An industry-first hardware implementation of highly-parallel Non-deterministic Finite Automata (NFA)
  - Orders of magnitude (>100x) faster pattern matching and graph analysis for:
    - Financial Services Industry
    - Bioinformatics
    - Data Analytics
    - Network Security...and more!
      - Significantly easier to program than GPUs natural parallelism
        - > Dramatically more flexible than FGPAs
          - Unmatched in processing challenges such as graph analysis, fuzzy string matching, and data analytics!





### Macro Trends Aligned with the Automata Processor



Limitations of conventional processors:



# **Problems Aligned with the Automata Processor**

Applications requiring **deep analysis** of **data streams** containing **spatial** and **temporal** information are often **memory-bound** and will benefit from the **processing efficiency** and **parallelism** 

of the Automata Processor.



#### **Network Security:**

• Millions of patterns

Real-time resultsUnstructured data



#### **Bioinformatics:**

- Large operands
- Complex patterns
- Unstructured data



#### **Financial Services:**

Highly parallel operationReal-time results

Unstructured data



### Data Analytics:

- Highly parallel operation
- Real-time results
- Unstructured data



## **Automata Value Proposition**

Fact: The ability to **generate** and **transport** unstructured data has vastly exceeded our capacity to **analyze** that same information

Market gap: a configurable processing engine that can quickly analyze arbitrarily complex graph processing problems Micron response: A massively parallel non-von Neumann compute architecture that executes multiple instructions for each byte of data MISD (Multiple Instruction Single Data)

#### **Higher Performance:**

>100x performance increase for complex graphs

#### Lower Energy:

- As little as 0.9 pJ/calculation
- 5.8W TDP per device
- Natural 8:1 internal transport data compression
- Lower power x Shorter execution time x One-time data touch (stream)



AUT•MATA

PROCESSING

MICRON'



 One PCIe card outperforms a cluster of 48 Xeon processors

#### **Better Quality of Result:**

 Directly analyzes complex graphs without approximations

#### **Ease of Parallel Programming:**

- No special programming considerations required to perform parallel processing
- No vectorization of data; no timing loops; no race conditions



# **Breakthrough Performance**

Planted Motif Search Problem	Automata Processor Single PCIe Board	UCONN – BECAT Hornet Cluster
Processors	48 AP + 1 CPU socket	48 CPU (Cluster/OpenMPI)
Power	245-315W <sup>1</sup>	~2,500W <sup>3</sup>
Cost	TBA	~\$20,000 <sup>1</sup>
(25,10)	12.26 minutes <sup>2</sup>	20.5 minutes 1.67x
(26,11)	13.96 minutes <sup>2</sup>	46.9 hours 202x
(36,16)	36.22 minutes <sup>2</sup>	Unsolved, est. 74 years <sup>1</sup> 1Mx
(39,18)	~10 hours <sup>2</sup>	Unsolved, est. 94,000 years <sup>2</sup> 5Bx

Planted Motif Search - a leading "NP Complete" problem in bioinformatics

Micron Confidential

Solutions involving high match lengths and substitution counts are often presented to HPC clusters for processing

Independent research predicts the Micron Automata Processor significantly outperforms a multi-core HPC cluster in speed, power and estimated cost

<sup>1</sup> Micron Technology Estimates <sup>2</sup> Research conducted by Georgia Tech (Roy/Aluru) <sup>3</sup> Only power of CPU chips, not including rest of system power such as 4GB DRAM per core

Aicron

© 2015 Micron Technology, Inc.

©2015 Micron Technology, Inc.

### Automata Processor: Support & Tools



#### **PCIe Accelerator Board**

Industry Standard PCIe (G3) bus interface
Capacity for up to 32 AP's
Large FPGA capacity
DDR3 for local storage

#### Software Development Kit

AP Optimization, loading & debugging tools & compiler.







- How do we enable introduction of new memory technologies?
- How do we enable more economical forms of processing?
- How do we mitigate scaling challenges of current technologies?
- How do we deal with the economics of yield issues, especially 3D?





# Abstraction

# The Need for Abstracted Memory and Storage

- Hide device characteristics and ops not relevant to desired functions
  - Timing, refresh, bit disturb, data errors, wear mechanisms
- Mitigate future device characteristic degradation
  - Capacitive coupling increase as  $\lambda$  decreases, increases intrinsic errors
- Enable unconstrained innovation
- Enable introduction of new technologies, new hierarchical elements
  - New Persistent Memory, DRAM replacement, NAND replacement
  - No need to sync all memory vendors and all users
- Facilitate "bit blending" from numerous tech's
  - Optimal performance/\$ = right proportioning
- Facilitate more operations closer to, or in, memory





## **Future System**

- Abstraction enables introduction of new memory types
- And many more future innovations

Abstracted Interface Protocol on standard PHY, electrical or optical as viable



• How much abstraction? The lower the latency, the lighter it must be



### **Near-term System**





### Future Integrated Processor/Memory/Storage



cron

# Future Hybrid Electrical/Optical Distributed System

- Memory and processor types of any kind in same logical module
  - Processor, accelerator, memory, storage, router, etc.
- All longer reach links optical, shorter links – most economical
- "Same" protocol on all, i.e. same root protocol







- Logic device scaling continues but with sub-historic benefits
- eDRAM scaling is especially unfavorable: maybe 0 or 1new eDRAM node to come
- Memory tech scaling also with challenges, continues to pull away ahead of logic tech
- Importance of 3D monolithic and 3D-stacked die
- Emerging memory: none for direct replacements coming soon but will emerge for different functionality
- New class of memory: 3D Xpoint<sup>™</sup> technology
- Energy is the Great Wall
- New Epoch: shift to memory-centric computing, more use of memory technology
- Example: Micron Automata Processor with >100x conventional performance, UNBOUNDED!
- More processing in/near memory, more sensible distribution of components and work
- Memory-Storage-Processing Abstraction mitigates many issues, enables a whole new world







### **Automata Resources**

- SDK Toolchain Download
- http://www.micronautomata.com
- Automata simulator & tutorial : <u>http://www.micronautomata.com/#applications</u>
- Networking Deep packet inspection for malware signature detection : <u>http://www.micronautomata.com/network\_security</u>
- Bioinformatics DNA sequencing & Motif search : <u>http://www.micronautomata.com/bioinformatics</u>
- More resources :
  - www.micronautomata.com
  - www.cap.virginia.edu



