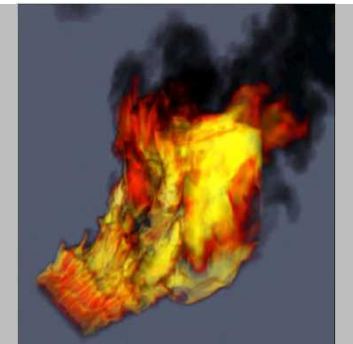
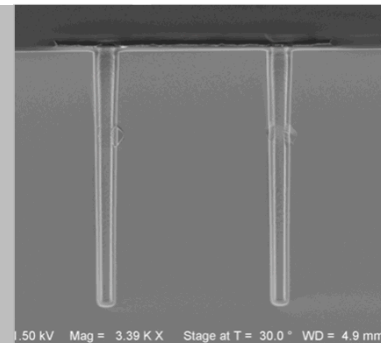
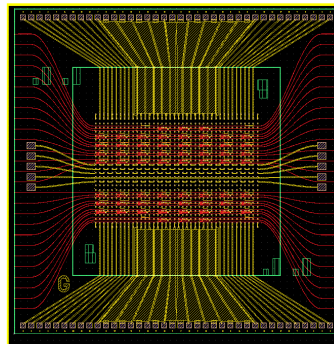


*Exceptional service in the national interest*

# Memory Direction

## Perceived Vision and Requirements

Dave Resnick  
7 May 2014



# Memory Near Term

- DDR-4 coming
  - Faster and lower power—but not enough
- No plans for DDR-5 (But some activity in JEDEC along other lines)
- JEDEC Wide-IO at low end; GDDR not good for main memory
- Micron Hybrid Memory Cube (HMC) this year
  - A big breakthrough in architecture and technology
  - Concerns about cost
- High Bandwidth Memory (HBM), new JEDEC standard coming
  - Made to be inexpensive
  - Not scalable, not sufficiently resilient for large or critical systems  
(But seems likely to be a part of multi-level-memory systems near term)
- Non-volatile memory, for main-memory use, is not near term

# Needs for Main Memory Components

- Scalability
  - Replace the DDR-x interface with one that scales and is not a performance limiter. Should be abstract to support going into the future (e.g.: Connect DRAMs and NV memory on the same links)
  - Enable cost effective memory systems from 10's of GB to petabytes (both in bandwidth and size) with the same DRAM components
- Bandwidths that match performance needs  
Multiple 100's of GB/sec **per memory part**
- Lower energy than current DRAM at the system level
- Resiliency **Worth lots of separate discussion**
  - Exascale systems will have to 10,000,000 memory parts
  - Design that can keep running in spite of most memory failures

And counts on 3D stacking being a success

# Vision for Mid-Term

- Need to support multiple levels of memory per node
  - Three? levels of memory
    - Local Cache **and scratch** memory  
Extreme bandwidth, to 10's of GB (if fast DRAM) or less for SRAM
    - Mid Very high bandwidth DRAM, say 1 to 4 TB
    - Big NV, lower bandwidth, 10's of TB
- Plan and implement intelligence in the memory system
  - Operations in memory that improve system performance and reduce energy use
    - Atomic operations Add/Sub, Test and Swap, Logic ops, ...
    - Move, Gather/Scatter Transpose, Sparse matrix support, ...
    - Reference forwarding Search, Data Base, ...
- Work towards a system design such that any processor can access any memory in a system that supports new communication and coherency paradigms

# Getting Into the Future

- Multi-level memory will need multiple kinds of software support
  - Data placement optimization
  - Move the process, not the data?
  - Upgraded libs, applications, OS support, and ...
- Distributed asynchronous check-pointing with local NV memory?
  - Will need software management; what hardware?
- New features in memory will need software support (and planning on if and how to have existing applications use the new features) Likely will need feature standardization, which will take some time
- Breaking the memory wall offers lots of opportunities in multiple directions in software and in system architecture
- Need help in showing that new functions are worth the effort

# And Further Out

- A system architecture in which the memory system is considered at the same level of importance and concern as CPUs are currently. The memory system can be where a good part of a job is done—and will save a large portion of the energy it currently takes.

For many applications, there is more power and energy in data motion than in floating-point and other CPU operations, so processing in memory will be faster and lower energy than current systems. And the trend will continue and get more acute

- Even if the memory system does only a part of the work, better and more efficient systems will result
- **Systems such that: The memory system IS the network!**

# Finally

...the moment one chooses a given component as the elementary memory unit, one has also more or less determined upon much of the balance of the machine.

- Arthur Burns
- Herman Goldstine
- John von Neumann

From the June 1946 architectural description report of the first stored program computer (IAS\* machine; RAND JOHNNIAC largely a copy). Burns and Goldstine were the architects, under direction of John von Neumann

\* Institute of Advanced Study, Princeton