

Lecture 2: Memory Energy

- Topics: handling overfetch, LPDRAM, row buffer management, channel energy, HMC, DBI

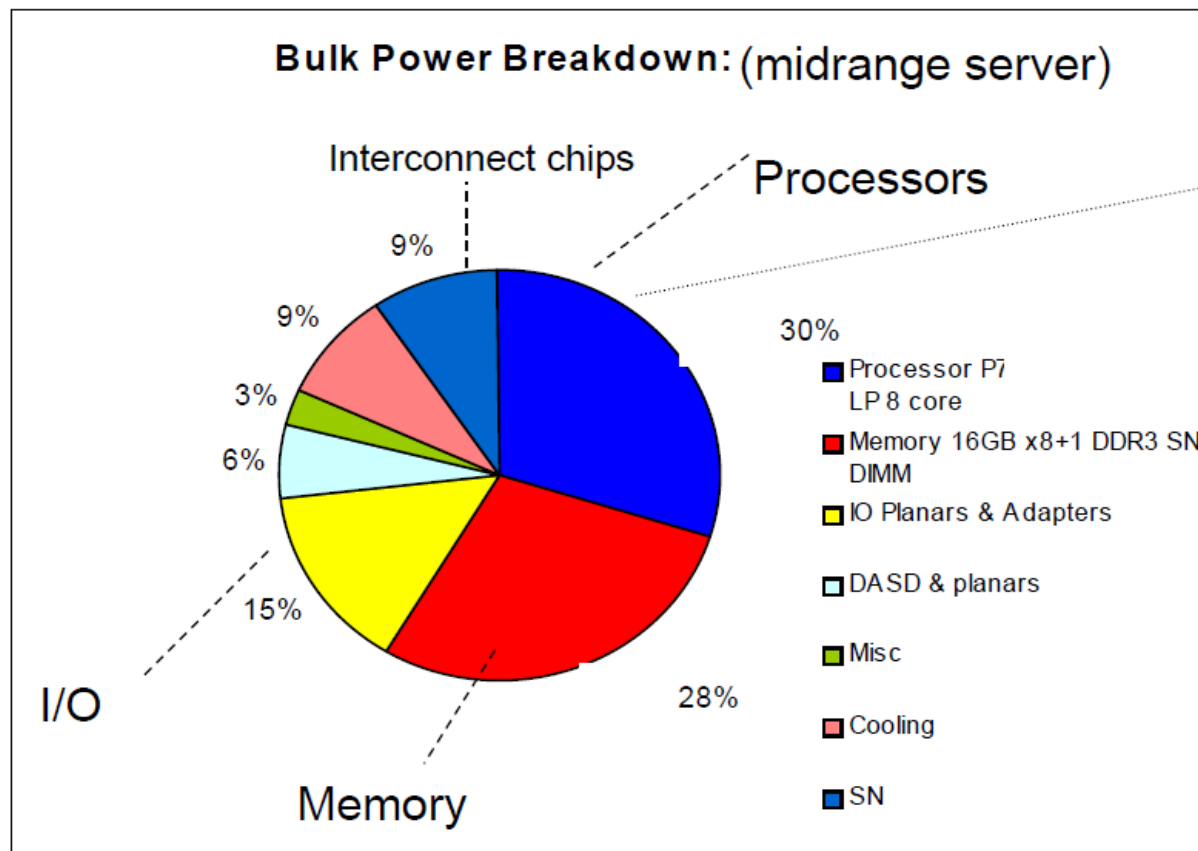
Power Wall

- Many contributors to memory power (Micron power calc):
 - Overfetch
 - Channel
 - Buffer chips and SerDes
 - Background power (output drivers)
 - Leakage and refresh

Power Wall

- Memory system contribution (see HP power advisor):

System-Level Power Breakdown



IBM data, from
WETI 2012 talk
by P. Bose

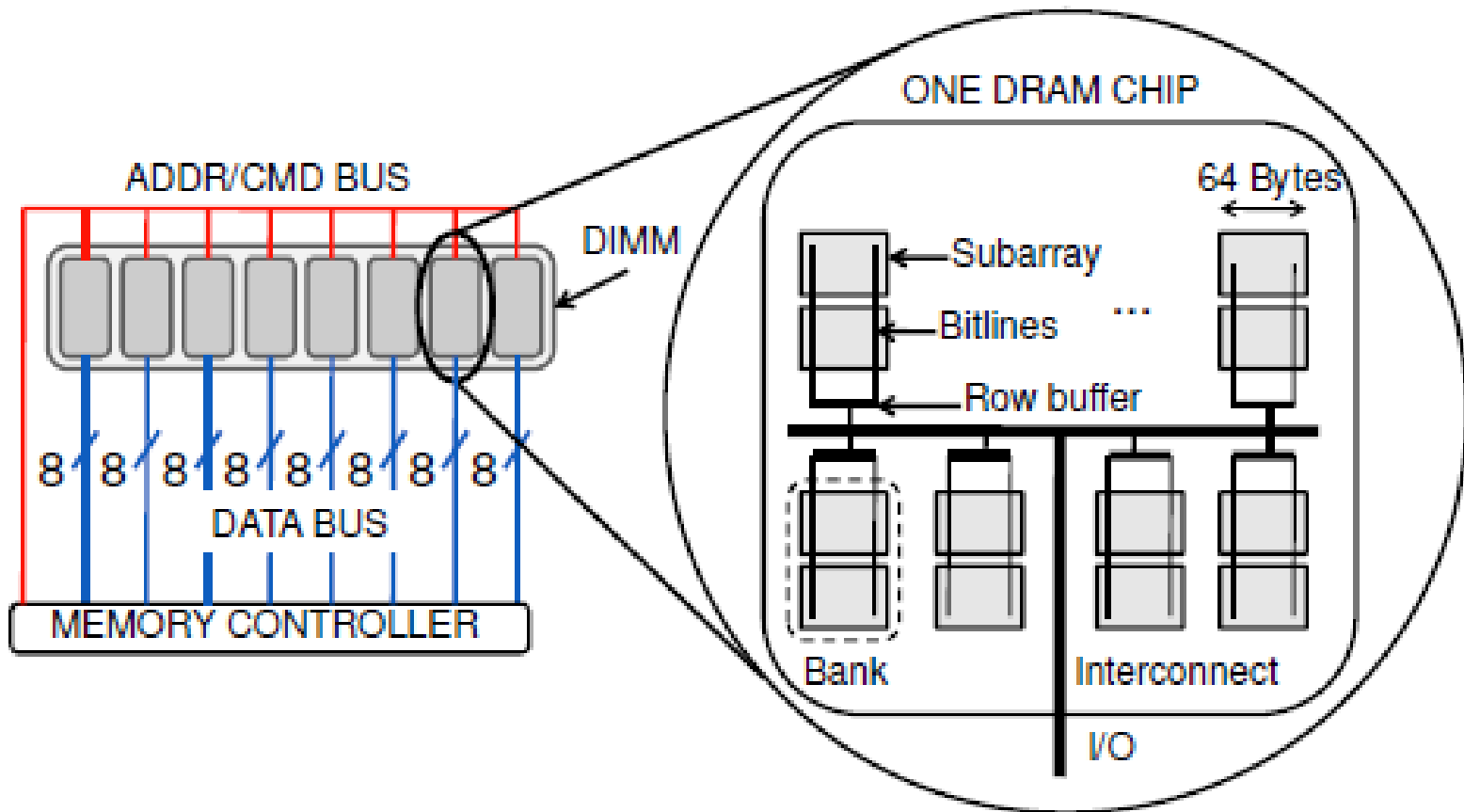
- Memory power has caught up with processor power

Overfetch

- Overfetch caused by multiple factors:
 - Each array is large (fewer peripherals → more density)
 - Involving more chips per access → more data transfer pin bandwidth
 - More overfetch → more prefetch; helps apps with locality
 - Involving more chips per access → less data loss when a chip fails → lower overhead for reliability

Re-Designing Arrays

Udipi et al., ISCA'10



Selective Bitline Activation

- Additional logic per array so that only relevant bitlines are read out
- Essentially results in finer-grain partitioning of the DRAM arrays

- Two papers in 2010: Udipi et al., ISCA'10, Cooper-Balis and Jacob, IEEE Micro

Rank Subsetting

- Instead of using all chips in a rank to read out 64-bit words every cycle, form smaller parallel ranks
- Increases data transfer time; reduces the size of the row buffer
- But, lower energy per row read and compatible with modern DRAM chips
- Increases the number of banks and hence promotes parallelism (reduces queuing delays)
- Initial ideas proposed in Mini-Rank (MICRO 2008) and MC-DIMM (CAL 2008 and SC 2009)

Micron HMC

- Many energy-efficient features: smaller arrays and few arrays activated per access
- 256-byte fetches, so low overfetch
- 3.7 pJ/bit for DRAM read and 6.78 pJ/bit for SerDes hop
- DDR3 is 70 pJ/bit and LPDDR is 40 pJ/bit (Malladi et al., ISCA'12) (all these numbers are for peak utilization – they are much higher at lower utilizations)

DRAM Variants – LPDRAM and RLDRAM

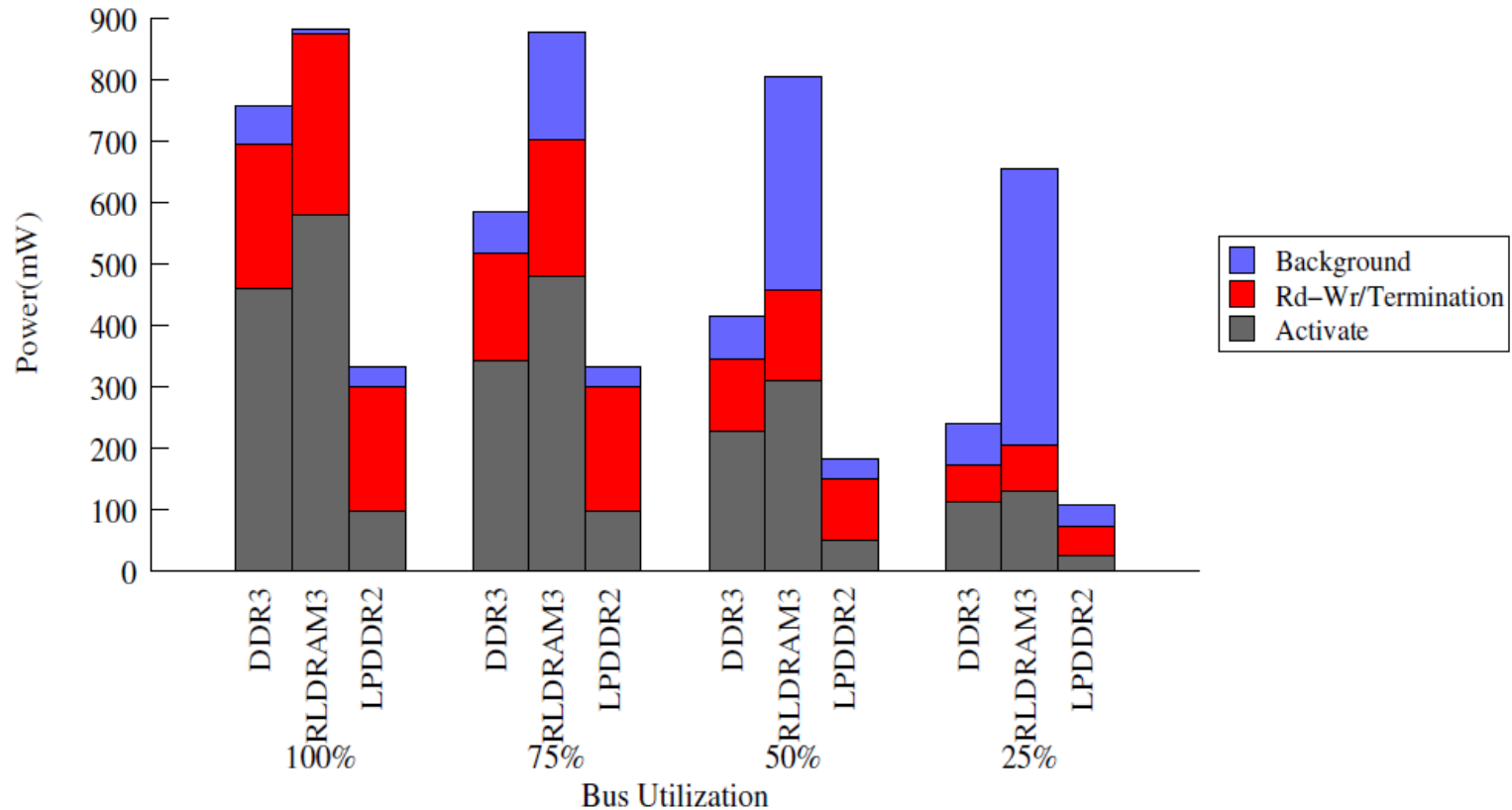


Figure 2: Power vs Bus Utilization (the RLDRAM3 part has a capacity of 512Mb while the DDR3 and the LPDDR2 parts have capacities of 2Gb)

LPDRAM

- Low power device operating at lower voltages and currents
- Efficient low power modes, fast exit from low power mode
- Lower bus frequencies
- Typically used in mobile systems (not in DIMMs)

- Implement a few DIMMs/channels with LPDRAM and a few DIMMs/channels with RLDRAM
- Fetch critical data from RLDRAM and non-critical data from LPDRAM
- Multiple ways to classify data as critical or not:
 - identify hot (frequently accessed) pages
 - the first word of a cache line is often criticalEvery cache line request is broken into two requests

Row Buffer Management

- Open Page policy: maximizes row buffer hits, minimizes energy
- Close Page policy: helps performance when there is limited locality
- Hybrid policies: can close a row buffer after it has served its utility; lots of ways to predict utility: time, accesses, locality counters for a bank, etc.

- Organize data across banks to maximize locality in a row buffer
- Key observation: most locality is restricted to a small portion of an OS page
- Such hot micro-pages are identified with hardware counters and co-located on the same row
- Requires hardware indirection to a page's new location
- Works well only if most activity is confined to a few micro-pages

- Performs DVFS on the memory controller and DFS on the channel
- The frequencies depend on bandwidth utilization and estimated energy/performance drop
- Requires no change to DRAM chips and DIMMs (in modern systems, the channel/DIMM frequency is set at boot time)
- Only saves energy on the processor, not on the channel and DIMM

Data Bus Inversion (DBI)

- Implemented in GDDR and in upcoming HBM
- Send the inverse of a word to reduce bit-flips

Title

- Bullet