Hardware Level Fault-Tolerance Techniques

EECE 513: Design of Error Resilient Computer Systems

- List the techniques for improving the reliability of commodity & high end processors
- Design coding techniques for memory soft errors and evaluate their trade-offs
- Understand the benefits of chipkill ECC, sparing and scrubbing
- List the techniques used in the I/O sub-system

High-Availability Systems

• IBM G5 Mainframes

- Duplicated execution units on each core
- Redundant CPU logic
- Inline checking in I/O subsystem
- ECC in memory and registers
- Error Recovery is accomplished using instruction retry
 - Transparent to the S/W

Tandem Non-Stop

- Duplicated processors running in lock-step
- Process pairs for checking
- End-to-end disk checksums, CRC
- ECC in memory only
- Error recovery is achieved by swapping in backup processes
 - S/W needs to be involved in the failover

Commodity Micro-processors

Source: **Recent Advances and New Avenues in Hardware-Level Reliability Support,** by Iyer, Nakka, Kalbarczyk and Mitra, IEEE Micro 2005.

Feature		Intel P6 family	AMD Hammer	Intel Itanium
Interna L1	l registers Data	Parity Parity	No protection I cache: parity; D cache: ECC	No protection Parity
L2	Tag Data	Parity ECC	Parity ECC	Parity 8-bit ECC/ 64 data bits
L3	Tag Data	Parity N/A	ECC N/A	Parity 8-bit ECC/ 64 data bits
TLB Bus	Tag Is ies	N/A Parity ECC on CPU-L2 bus	N/A Parity No protection	3 parity bits Parity No protection
Other features		Machine check architecture (MCA) to detect and correct errors in	MCA	Multilevel MCA: local and global MCA, hardware bus reset
Unique features		Functional redundancy checking using master/slave processors	Chipkill memory controller to support memory scrubbing; NX virus protection for Windows XP SP2	Multilevel error containment; watchdog timer; error logging and corrected error notification; NX virus protection

- List the techniques for improving the reliability of commodity & high end processors
- Design coding techniques for memory soft errors and evaluate their trade-offs
- Understand the benefits of chipkill ECC, sparing and scrubbing
- List the techniques used in the I/O sub-system
- Putting it together: Stratus Case Study

Memory Errors: History

- Memory elements have long been the target of soft-errors since the late 70's
 - In 1978 May and Woods reported "A New Physical Mechanism for Soft Errors in Dynamic Memories"
 - In 1979, "Alpha-Particle-Induced Soft Errors in Dynamic Memories."
 - SRAMs saw problems approximately 2 years later

Soft Errors Today

Baumann, R.; , "The impact of technology scaling on soft error rate performance and limits to the efficacy of error correction," *Electron Devices Meeting*, 2002.



Figure 2. The single bit (white diamonds) and system (black diamonds) SER trends for DRAM as a function of technology node. The operating voltage at each node is represented by the curve with gray triangles.

DRAMs

Figure 3. The single bit SER trend (white diamonds) in SRAM devices as a function of technology node. The rapid scaling down of operating voltages is evident shown by the curve with gray triangles. The – 10x reduction in SER after the 0.25 µm node is the effect of removing BPSG.

Voltag

SRAMs

Error Trends in today's memories

- DRAM reliability has remained relatively constant over many years
 - Thanks to improvement in fabrication
 - May be different in eDRAMs and mobile DRAM
- SRAM reliability becoming an increasing concern with shrinking cell sizes and voltage
- DRAM hard errors are emerging as a problem [Schroeder'09][Dell'08]

Parity Protection - 1

- Single bit added to each memory byte/word to detect a single error
 - Cannot detect multiple errors
 - Cannot correct the error
 - Affordable alternative to ECC memory

Parity bit – even parity

xp = x0 ^ x1 ^ x2 ^ x3^ x4 ^ x5 ^ x6 ^ x7

Parity Protection - 2

- Requires an additional operation on reads/writes to memory → extra access latency
- Circuitry to compute parity bit is simple, but requires additional area and power
- Used mainly in SRAM structures where error rates were low and access times are important
- For DRAMs, no added benefit of using parity over ECC when the memory data width is greater than 8 bytes

- For every memory word of size 'n' bits, we need at least log₂ (n) bits of ECC memory
 - For 64 bit memory, we need at least 6 bits of ECC
 - The check bits are distributed throughout word
 - Each bit is protected by multiple checkbits given by the index (the sum of the checkbits matches index)



• Let's say you had a single bit error in R5 $(1 \rightarrow 0)$



Check Bits are recomputed and compared.

R1 = R3 ^ R5 ^ R7 ^ R9 ^ R11 = 0 ^ 0 ^ 1 ^ 0 ^ 1 = 0

R4 = R5 ^ R6 ^ R7 ^ R12 = 0 ^ 0 ^ 1 ^ 1 = 0

Both check-bits R1 and R4 differ from their computed values. These are called the syndromes. So we can infer that the bit R5 had an error in it, and can correct the error.

Let's say you have errors in bits R5 and R7 (double-bit error)



Let's compute check-bits R1, R2 and R4

 $R1 = R3 ^ R5 ^ R7 ^ R9 ^ R11 = 0 ^ 0 ^ 0 ^ 0 ^ 1 = 1$ (Same as prior value)

 $R2 = R3 \land R6 \land R7 \land R10 \land R11 = 0 \land 0 \land 0 \land 0 \land 1 = 1$ (Differs from prior value)

 $R4 = R5 ^ R6 ^ R7 ^ R12 = 0 ^ 0 ^ 0 ^ 1 = 1$ (Same as prior value)

How do we distinguish this case from the one where bit R2 is corrupted ?

• Add an extra parity bit R0 for the entire word



Extra parity bit for the word is added

- In the case of a single bit error, both syndrome bit(s) and R0 bit will differ → can be corrected
- In case of a double error, only syndrome bit(s) differs → can be detected but not corrected

ECC: Implementation Trade-offs

- ECC memory is not free !
 - Performance overheads for read/write operations
 - 3 to 4 % more for PC133 CAS2 ECC SDRAM
 - Up to 33 % for high-speed SRAMs
 - Area overhead for error-detection/correction ckts
 - 20 % die overheads
 - Additional costs as chipset support is needed
 - 10 to 25 % more for entire chip
 - Effectiveness: Corrects more than 90% of errors

Above nos. are from the Terazzon white paper.

- List the techniques for improving the reliability of commodity & high end processors
- Design coding techniques for memory soft errors and evaluate their trade-offs
- Understand the benefits of chipkill ECC, sparing and scrubbing
- List the techniques used in the I/O sub-system

ChipKill ECC - 1

- ECC can detect 2 bit and correct 1 bit errors
 - Provided the entire memory chip does not fail
 - Chip failure can lead to data loss even with ECC



Traditional SEC/DED ECC for a 64-bit word with eight check-bits of ECC

ChipKill ECC - 2

- Solution: Use Chip-kill ECC [™] (IBM S/390)
 - Spread the ECC check bits over multiple chips
 - Bit-steering → Steer the checkbits of adjacent bits
 in a memory word to different words in the ECC



Chip-kill ECC Note how the bits are scattered across different modules

Figure 2

ChipKill ECC: Implementation Tradeoffs

- Incurs four times the overhead of traditional ECC
 - Can be optimized using very wide ECC words
 - Provide detection of chip failures but not correction
- Compaq proposed a clever interleaving solution to combine two ECC words into one module
 - Provides the benefits of Chipkill ECC with only as much cost as parity protection
 - After a chip has failed, the Compaq ECC is unable to provide protection from single/double bit errors

Parity, ECC and ChipKill- Comparison

- Simulation data gathered by IBM over 36 months comparing:
 - 32 MB Parity
 protected
 memory
 - 1 GB SEC ECC
 - 1 GB Chipkill ECC



Other variations of ECC

Scrubbing

- ECC memory only checks the bits during reads/ writes
- However, infrequent accesses may lead to bit errors accumalating
- Solution: Scrub memory periodically by performing reads/writes to unaccessed memory

• Sparing

- Correlated or large area defects cannot be combated with ECC alone
- Use spare rows/columns in conjunction with ECC
- Leads to an order of magnitude reliability improvement over ECC alone for hard faults

- List the techniques for improving the reliability of commodity & high end processors
- Design coding techniques for memory soft errors and evaluate their trade-offs
- Understand the benefits of chipkill ECC, sparing and scrubbing
- List the techniques used in the I/O sub-system

I/O Sub-system - 1

- Disk and other storage media protected using RAID technologies
 - Fairly mature, industry standard
 - However, data is susceptible when it is buffered
 - Firmware controllers and I/O processor errors
- Need to ensure end-to-end consistency of data from I/O initiation to disk read/write

I/O Sub-system - 2

- Techniques for end-to-end I/O checking
 - Checksums on data before and after reads
 - Checking of header fields for consistency
 - Watchdog timer for ensuring no deadlocks or livelocks of I/O devices
 - System-level consistency checks. e.g., read back data written to disk in chunks and check them
 - Use multiple file organizations to store data

- List the techniques for improving the reliability of commodity & high end processors
- Design coding techniques for memory soft errors and evaluate their trade-offs
- Understand the benefits of chipkill ECC, sparing and scrubbing
- List the techniques used in the I/O sub-system