













We distinguish between

- Concurrent (on-line) error detection
 - Detection of errors during operation
 Purpose is to mask or minimize adverse effects of errors
 - Purpose is to mask or minimize adverse effects of end
- Non-concurrent (off-line) fault detection
 - Testing to find physical hardware faults while the system is off-line
 Purpose is to identify faulty hardware units









Generic techniques for error detection (II)

- Information redundancy (pp. 125 and 128 in the course book)
 - Use of error detecting and error correcting codes
 Examples:
 - Cyclic redundancy check (CRC)
 - Hamming codes
- Time redundancy (p. 125 in the course book)
- Examples:
 - Double time redundant execution with comparison
 Triple time redundant execution with majority voting

Hardware techniques for error detection (examples)

- Bus monitoring (p. 130 in course book)
 - Checking the range of addresses generated by a CPU
 Examples
 - Checking that the CPU use an even address when reading a 32 or 64-bit word.
 Checking CPU memory accesses using a memory management unit (MMU).
- Watchdog timers (p. 130 in course book)
- Hardware technique supported by software
- Power supply monitoring (pp. 130-131 in the course book)

CPU Exceptions

- Modern central processing units (CPUs) (or cores) are equipped with hardware implemented error detection mechanisms called *hardware exceptions*
- The number and type of hardware exceptions varies depending on the CPU design
- When a hardware exception is raised, the CPU stops the program execution and jumps to an exception routine
- The handling of exceptions is very similar to how a CPU responds to interrupt signals
- Some examples of common hardware exceptions is given in the next two slides

Examples of CPU exceptions (1)

Bus error: detects errors during read and write accesses to the main memory. This exception is raised (triggered) when the CPU attempts to access an address to which no memory or I/O device is connected.

Address error: detects when the CPU attempts to access memory using an address that is not aligned with the word size of the CPU.

Undefined instruction: detects if the CPU during an instruction fetch reads a value from memory (or the instruction cache) that doesn't correspond to a valid instruction. This error can occur if the program counter is erroneously loaded with an address pointing to a data area rather than a program code area.

Examples of CPU exceptions (2)

Privilege violation: detects if a user program attempts to execute an instruction which is allowed only for programs that execute in the superuser mode (privileged mode), such as the operating system or device drivers. User programs normally executes in user mode (normal mode).

Division by zero: detects if a program tries to divide a number with zero.

Spurious interrupt: detects if an interrupt is signalled but no interrupt vector is provided by the interrupting device. (The interrupt vector tells the CPU which device it was that raised the interrupt signal and thereby indicates which interrupt service routine that the CPU shall execute.)

Watchdog timers

- Watchdog timers are used to detect slow programs or programs that hang in infinite loops
- The principle is simple:
 - When a program starts to execute, either the program itself or the operating system starts a hardware timer.
 - The timer must be reset by the program within a given deadline, otherwise
 the timer will send an interrupt signal to the CPU.
 - The interrupt signal causes the CPU to take appropriate recovery actions, such as restarting the program or rebooting the node.
- Watchdog timers are common in embedded real-time systems.
- They can be used to "transform" timing failures into signalled failures or silent failures.

Initiating recovery via a watchdog timer

- Watchdog timers are often used in conjunction with other error detection mechanisms to simplify implementation of recovery
- This works as follows:
 - When an error is detected, the error detection mechanism stores an error code in a designated memory area (preferably a "crash-proof" memory)
 - The error detection mechanism then forces a program hang, which causes
 the watchdog timer to raise an interrupt.
 - The interrupt signal invokes a recovery routine that reads the error code and then selects and executes appropriate recovery actions.
 - Restart and recovery could be done for an entire node, a single program, or a group of programs.

Software techniques for error detection (examples)

- Operating system assertions
 - Examples:

 Integrity checks of data structures used by the operating system
 - Execution time monitoring of application and system processes

Compiler generated run-time assertions

- Examples:
- Value range overflow checking
- Loop iteration bound overflow checking
- Type checking of constrained variables

When these mechanisms detect an error, they typically execute a TRAP or software interrupt instruction, causing the CPU to execute an exception handler that initiates appropriate recovery actions.

Examples of how acceptance tests/ software assertions can be constructed (from lecture 4)

- Satisfaction of requirements
 - Inversion of mathematical functions; e.g. squaring the result of a square-root operation to see if it equals the original operand
 - Checking sort functions; result should have elements in descending order
 - ...
- Reasonable checks
 - Checking physical constraints; e.g. speed, pressure, etc
 - Checking sequence of application states
 -



















Highlights from guest lectures

- Torbjörn Hult Fault tolerant computers in space applications
- Jan Jacobson Why standards for functional safety?
- Lars Holmlund Fault tolerance in the Gripen Flight Control System



















Summary of study of failures in highperformance computing systems (I)

- Failure rates vary widely across systems, ranging from 20 to more than 1,000 failures per year, and depend mostly on system size and less on the type of hardware.
- Failure rates are roughly proportional to the number of processors in a system, indicating that failure rates are not growing significantly faster than linearly with system size.
- There is evidence of a correlation between the failure rate of a machine and the type and intensity of the workload running on it. This is in agreement with earlier work for other types of systems [2], [6], [19].
- The curve of the failure rate over the lifetime of an HPC system looks often very different from lifecycle curves reported in the literature for individual hardware or software components.
- Time between failure is not modeled well by an exponential distribution, which agrees with earlier findings for other types of systems [5], [27], [9], [15], [19]. We find that the time between failure at individual nodes, as well as at an entire system, is fit well by a gamma or Weibull distribution with decreasing hazard rate (Weibull shape parameter of 0.7-0.8).



- Failures exhibit significant levels of temporal correlation at both short and long time lags. We find indication of autocorrelation for all types of failures; however, autocorrelation is particularly strong for hardware and software failures.
- We also find indication of spatial correlation, i.e., correlation between failures at different
 nodes during the same time interval. However, those are limited to failures with network
 root cause and not significant for other types of failures.
- Mean repair times vary widely across systems, ranging from one hour to more than a day. Repair times depend mostly on the type of the system, and are relatively insensitive to the size of a system.
- Repair times change significantly over the lifetime of a system. Both mean and median
 repair times drop by more than a third after the first year in operation.
- This might indicate that during the first year of operation, system administrators get
 significantly better at quickly identifying the root cause of a problem and fixing it.
- Repair times are extremely variable, even within one system, and are much better modeled by a lognormal distribution than an exponential distribution.

EDA122/DIT061 Fault-Tolerant Computer Systems

Final remarks (I)

- No system is perfect single points of failure cannot be avoided completely (cf. Fukushima nuclear disaster in Japan last year)
- Redundancy is no panacea it may prevent system failures, but increases cost, failure rate and energy consumption
- IT systems are physical artifacts the quality of their service depends on both the software and the hardware
 Vertical thinking is needed from transistors to user interfaces
- IT-systems cannot be fully understood they are too complex
 Billions of transistors, millions of lines of code lead to almost infinite numbers of fault, error and failure scenarios





