



IAF0030
Arvutitehnika erikursus I


Gert Jervan
Arvutitehnika instituut (ATI)
Tallinna Tehnikaülikool

General Information

- Contents:
Fault Tolerant Digital Systems
www.pld.ttu.ee/IAF0030
- Lecturer & Examiner:
Gert Jervan
IT-229 620 2261
gerje@pld.ttu.ee
www.pld.ttu.ee/~gerje

IAF0030 – Intro Gert Jervan, TTÜ/ATI 2

Course Plan

- 16 occasions, á 2+ hours
Thursdays 13:30-15:45
- Lectures (10-12 occasions)
- Case Studies
 - 45 min. presentations
 - Report
- Exam

IAF0030 – Intro Gert Jervan, TTÜ/ATI 3

Grading

- Case study presentation – 30%
- Case study report – 30%

Prerequisites

- Exam – 40%

Total: 2,5 points

IAF0030 – Intro Gert Jervan, TTÜ/ATI 4

Reading Materials

- Textbook:
 - Safety-critical Computer Systems, Neil Storey, Addison Wesley, 1996.
An introductory text which provides overview of safety related aspects and methods in computer systems development.
- Various papers
- Web pages, incident reports

www.pld.ttu.ee/IAF0030

IAF0030 – Intro Gert Jervan, TTÜ/ATI 5

Course Overview

- Reliability: increasing concern
 - Historical
 - High reliability in computers was needed in critical applications: space missions, telephone switching, process control etc.
 - Contemporary
 - Extraordinary dependence on computers: on-line banking, commerce, cars, planes, communications etc.
 - Hardware is increasingly more fault-prone (complexity, technology, environment)
 - Software is increasingly more complex
 - Things simply will not work without special reliability measures

IAF0030 – Intro Gert Jervan, TTÜ/ATI 6

Course Overview

- To get an insight into the broad area of system safety
- We cover techniques for high availability, fault tolerance, monitoring, detection, diagnosis, and confinement of failure, ways to improve availability through fast recovery and graceful service degradation, and techniques for using redundancy and replication.
- We also discuss the utopia of flawless software, the impact of scale on availability, ways to cope with human operator error, and metrics for evaluating dependability.

Contents

- Fault tolerance
- System reliability
- Hardware redundancy
- Error detection techniques
- Coding techniques
- Processor-level detection and recovery
- Disk arrays
- Checkpointing and recovery
- Software fault tolerance
- Testing distributed real-time systems
- ...

Lecture Outline



- Historical perspective and famous incidents/accidents

- Basic terminology

Murphy's Law

- "If something can go wrong, it will go wrong"

*Major Edward A. Murphy, Jr.
US Air Force, 1949*

- "Every component than can be installed backward, eventually will be"

Genesis Space Capsule

- \$260 million Genesis capsule was collecting samples of the solar wind over 3 years period
- Crashed in Sept 2004 due to the failure of the parachutes
- Reason: the deceleration sensors — the accelerometers — were all installed backwards. The craft's autopilot never got a clue that it had hit an atmosphere and that hard ground was just ahead.



Mars Orbiter

- One of the Mars Orbiter probes crashed into the planet in 1999.
- It did turn out that engineers who built the Mars Climate Orbiter had provided a data table in "pound-force" rather than newtons, the metric measure of force.
- NASA flight controllers at the Jet Propulsion Laboratory in Pasadena, Calif., had used the faulty table for their navigation calculations during the long coast from Earth to Mars.

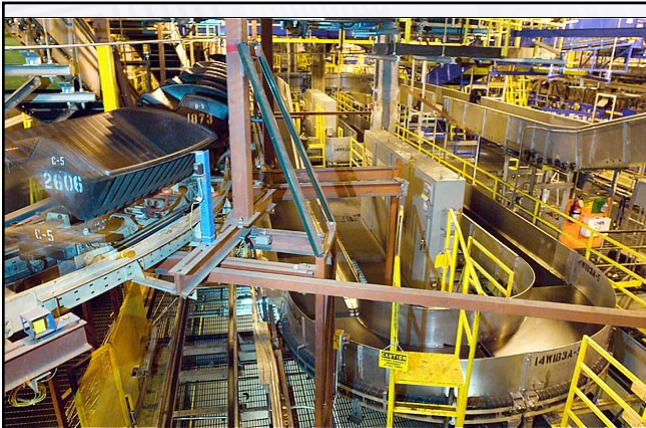
Lockheed Martin Titan 4

- In 1998, a Lockheed Titan 4 booster carrying a \$1 billion Lockheed Vortex-class spy satellite pitched sideways and exploded 40 seconds after liftoff from Cape Canaveral, Fla.
- Reason: frayed wiring that apparently had not been inspected. The guidance systems were without power for a fraction of a second.



Therac-25

- Therac-25:
 - the most serious computer-related accidents to date (at least nonmilitary and admitted)
 - machine for radiation therapy (treating cancer)
 - between June 1985 and January 1987 (at least) six patients received severe overdoses (two died shortly afterward, two might have died but died because of cancer, the other two had permanent disabilities)
 - scanning magnets are used to spread the beam and vary the beam energy
 - dual-mode: electron beams for surface tumors, X-ray for deep tumors



Denver Airport

- Denver International Airport, Colorado: intelligent luggage transportation system with 4000 "Telecars", 35km rails, controlled by a network of 100 computers with 5000 sensors, 400 radio antennas, and 56 barcode readers. Price: \$186 million (BAE Automated Systems).
- Due to SW problems about one year delay which costs \$1.1 million per day (1993).
- Abandoned in 2005 to save \$1 million per month on maintenance

Lecture Outline



- Historical perspective and famous incidents/accidents

- Basic terminology

Fault Tolerance

- A fault-tolerant system is one that can continue to correctly perform its specified tasks in the presence of failures:
 - hardware
 - software
 - user errors
 - environmental, input, ...
- Fault tolerance is the attribute that enables a system to achieve fault tolerant operation.

Basic Concepts

- *Fault Tolerance* is closely related to the notion of "Dependability". This is characterized under a number of headings:
 - *Availability* – the system is ready to be used immediately.
 - *Reliability* – the system can run continuously without failure.
 - *Safety* – if a system fails, nothing catastrophic will happen.
 - *Maintainability* – when a system fails, it can be repaired easily and quickly (and, sometimes, without its users noticing the failure).

Faults, Errors & Failures

- **Fault**: a defect within the system or a situation that can lead to the failure
- **Error**: manifestation of the fault – an unexpected behavior
- **Failure**: system not performing its intended function

Fault → Error → Failure

Fault Examples

- Bit flips in hardware due to cosmic radiation
 - A person on an airplane over the Atlantic at 35,000 ft working on a laptop with 256 Mbytes (2 Gbits) of memory. At this altitude, the SER of 600 FITs per megabit becomes 100,000 FITs per megabit, resulting in a potential error every five hours.
 - 1 FIT (failures in time), is the number of failures in 1 billion device-operation hours. A measurement of 1000 FITs corresponds to a MTTF (mean time to failure) of approximately 114 years.

Fault Examples

- Year 2000 bug
- Loose wire
- Aircraft retracting its landing gear while on ground
- Effects in time:
 - Permanent
 - Transient
 - Intermittent

Permanent

- A permanent fault or failure is one which is stable and continuous.
- Permanent hardware failures require some component to be replaced or repaired.
- An example of a permanent fault would be a VLSI chip with a manufacturing defect, causing one input pin to be stuck high (stuck-at-1).

Intermittent

- An intermittent fault is one which only manifests occasionally, due to unstable hardware or certain system states.
- A loose contact on a connector will often cause an intermittent fault.
- Intermittent electrical faults, as a rule, are notoriously difficult to detect. Typically, whenever the fault doctor shows up, the system works fine.

Transient

- A transient fault is one which results from a temporary environmental condition.
- For example, a voltage spike might cause a sensor to report an incorrect value for a few milliseconds before reporting correctly.

Failure Classification

- Domain/Nature
 - Value failure
 - Timing failure
- Perception
 - Consistent failure
 - Inconsistent failure
- Effect
 - Benign failure
 - Malign/catastrophic failure
- Frequency
 - Single failure
 - Repeated failure

Failures

- Crash Failure: After an error has been detected, the component stops silently.
- Omission Failure: Sometimes a result is missing; when result is available, it is correct.
- Consistent Failure: If there are multiple receivers, all see the same erroneous result.
- Byzantine (Malicious, Asymmetric) Failure: Different receivers see differing results.

Failures (cont.)

- Timing Failure: A server's response lies outside the specified time interval.
- Response Failure: The server's response is incorrect (value of the response is wrong, server deviates from the correct flow of control).
- Arbitrary Failure: A server may produce arbitrary responses at arbitrary times.

Fault Handling

- Fault avoidance: eliminate problem sources
 - Remove defects: Testing and debugging
 - Robust design: reduce probability of defects
 - Minimize environmental stress: Radiation shielding etc
 - Impossible to avoid faults completely
- Fault tolerance: add redundancy to mask effect
 - Additional resources needed (more later)
 - Examples:
 - Error correction coding, voting and masking, checksums, ...
 - Backup storage, replication, ...
 - Spare tire, etc

Fault Tolerance

- **Fault detection** is the process of recognizing that a fault has occurred. Fault detection is often required before any recovery procedure can be initiated. The techniques include error detection codes, self-checking/failsafe logic, watchdog timers, and others.
- **Fault location** is the process of determining where a fault has occurred so that an appropriate recovery can be initiated.

Fault Tolerance (cont.)

- **Fault containment** is the process of isolating a fault and preventing the effects of that fault from propagating throughout the system.
- **Fault recovery** is the process of remaining operational or regaining operational status via reconfiguration even in the presence of faults. A few basic approaches are fault masking, retry, and rollback.

Definitions

- Failure rate (λ):
 - Average frequency with which something fails.

$$\frac{6 \text{ failures}}{7502 \text{ hrs}} = 0.0007998 \text{ failures / hr} = 799.8 \times 10^{-6} \text{ failures / hr}$$

- Mean time to failure (MTTF):
 - Average time between failures $MTTF = \frac{1}{\lambda}$

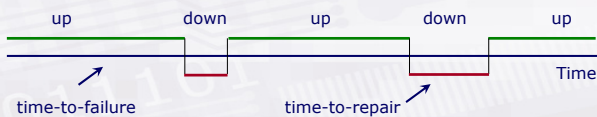
Dependability

- Property of a computing system which allows reliance to be justifiably placed on the service it delivers
- Dependability = reliability + availability + safety + security + ...
- Reliability → continuity of correct service
- Availability → readiness of usage
- Safety → no catastrophic consequences
- Security → prevention of unauthorized access

Reliability

- A measure of an it performing its intended function satisfactorily for a prescribed time and under given environment conditions.
- Probability that system will survive to time t
 - In aerospace industry the requirement is that failure probability is 10^{-9} (one failure over 10^9 hours (114 000 years) of operation)
- Time To Failure (TTF)
- Mean Time To Failure (MTTF)

Availability



$$Availability = \frac{MTTF}{MTTF + MTTR}$$

- Availability:
 - Probability that system is operational at time t
- High availability:
 - MTTF → infinity (high reliability)
 - MTTR → zero (fast recovery)

The Myth of the Nines

Nines	Availability	Downtime per year	Downtime per week	Example
2 nines	99%	3.65 days	1.7 hours	General web site
3 nines	99.9%	8.75 hours	10.1 min	E-commerce site
4 nines	99.99%	52.5 min	1.0 min	Enterprise mail server
5 nines	99.999%	5.25 min	6.0 s	Telephone system
6 nines	99.9999%	31.5 s	0.6 s	Carrier-grade phone switch

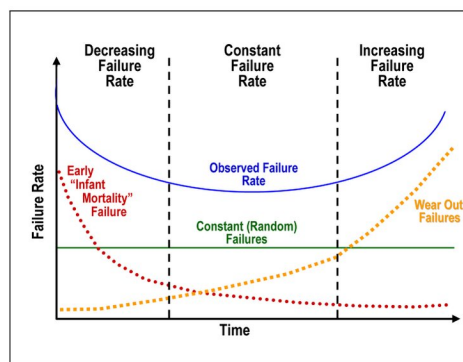
Historical Evaluation

- Mean Time Between Failures:
 - $MTBF = MTTR + MTTF$
 - ENIAC. MTBF: 7 minutes (18000 vacuum tubes)
 - TX-2 interactive computer (MIT)
 - F-8 Crusader – first fly-by-wire
 - MD-11
 - A320 family
 - Patriot missile defence system (needed reboot after 8 hours)

Ultra-Reliable Systems

- Airbus A320 family fly-by-wire system:
 - computer controls all actuators
 - no control rods, cables in the middle
 - 5 central flight control computers
 - different systems used
 - Thomson CSF => 68010
 - SFENA => 80186
 - software for both hardware written by different software houses
 - all error checking & debugging performed separately
 - computer allows pilot to fly craft up to certain limits (flight envelope)
 - beyond: computer takes over

Bathtub Curve



Safety

- Attribute of a system which either operates correctly or fails in a safe manner
- Freedom from exposure to danger, or exemption from hurt, injury or loss.
- "Fail-safe": traffic lights start to blink yellow
- Degrees of safety
- Closely related to risk

Risk

- A combination of the likelihood of an accident and the severity of the potential consequences
- The harm that can result if a threat is actualised
- Acceptable/tolerable risk
 - The Ford Pinto case
 - **BENEFITS**
 - Savings: 180 burn deaths, 180 serious burn injuries, 2,100 burned vehicles.
 - Unit Cost: \$200,000 per death, \$67,000 per injury, \$700 per vehicle.
 - Total Benefit: $180 \times (\$200,000) + 180 \times (\$67,000) + 2,100 \times (\$700) = \$49.5 \text{ million.}$
 - **COSTS**
 - Sales: 11 million cars, 1.5 million light trucks.
 - Unit Cost: \$11 per car, \$11 per truck.
 - Total Cost: $11,000,000 \times (\$11) + 1,500,000 \times (\$11) = \$137 \text{ million.}$

Graceful Degradation

- The ability of system to automatically decrease its level of performance to compensate for hardware failure and software errors.

Maintainability

- $M(t)$ is the probability that a failed system will be restored within a specified period of time t .
- Restoration process:
 - locating problem, e.g. via diagnostics
 - physically repairing system
 - bringing system back to its operational condition

System Safety & Hazards

- Safety:
 - achieved by anticipating accidents and eliminating their causes
- Hazards are potential causes of accidents
 - Conditions in a system which together with other factors in the environment inevitably cause accidents

Questions?

Tallinn University
of Technology

Gert Jervan

Department Of Computer Engineering
Tallinn University of Technology
Estonia

Administrative issues

www.pld.ttu.ee/IAF0030

Gert Jervan

IT-229 620 2261

gerje@pld.ttu.ee

www.pld.ttu.ee/~gerje

- Case Studies
 - Presentation + report
- Exam