

Fault Tolerant Distributed Real-Time Systems (draft)

István Majzik

DISCOM TEMPUS - September 1999

1 Basic definitions

Real-time (RT) computer system: Correctness of the system depends not only on the logical (functional) results, but also on the timeliness, i.e. the physical time at which the results are presented.

RT system: changes its state as a function of physical time

² General structure:

operator <-> RT computer system <-> controlled RT object
man-machine instrumentation
interface interface

² General task:

– react to stimuli from the controlled object/operator request

² Deadline: At which the results must be produced

– soft: result has utility after the deadline

– ...rm: no utility after the deadline

– hard: catastrophic consequences if the deadline is missed

² Hard RT system: at least one hard deadline exists

- required: GUARANTEED temporal behavior under all speci...ed LOAD and FAULT conditions

² Comparison of hard and soft RT systems:

Hard RT	Soft RT
hard deadline	soft deadline
predictable performance	degraded performance in peak load
synchronous with env.	computer control
often safety-critical	usually non-critical
active redundancy	checkpoint-recovery
short-term data integrity	long-term data integrity
autonomous error detection	user-assisted error detection

² Fail-safe vs. fail-operational:

- fail-safe: the controlled object has a safe state
 - (e.g. all traΦc lights are red)
 - error -> transition to safe state
 - computer systems: high error detection coverage required
- fail-operational: no safe state in the object
 - (e.g. airplane)
 - computer systems: minimal level; of service required to avoid catastrophe

² Guaranteed response vs. best eΦort:

- guaranteed response: peek load and fault scenario – to be speci...ed! (rigorous)
- best eΦort: hard to predict rare event scenarios (+ the design is often resource-inadequate)

² Event-triggered vs. time-triggered systems:

² event: any occurrence that happens in time

- (state change in the object/computer)

² trigger: event that causes the start of an action

² event-triggered (ET):

- all activities are initiated by events
- interrupt-like mechanism
- dynamic scheduling of activities (tasks)

² time-triggered (TT):

- activities are initiated by the progress of time
- interrupt: clock only
- (synchronized distributed clock)

2 General structure

2.1 Why using often distributed systems for RT purposes?

2.1.1 System architecture

² nodes: functional+temporal properties

- mapping between nodes and functions
- node error -> it is clear, which function is affected

² communication network:

- interface of nodes (CNI)
- event queue: often FIFO
- state information: overwriting old values

2.1.2 Composability

² system properties follow from subsystem properties

- (all subsystem combinations work properly)

² ET systems: not composable

- (message overload to receivers, conflicts)

² TT systems: composable

- temporal control resides in the comm. subsystem
- message scheduling tables are used
- transfer happens at predefined time points
- computer and communication subsystem's properties are isolated

2.1.3 Scalability

² no limits on the extensibility of the system

- nodes can be added (up to communication capacity)
- clusters, gateways can be established

² controlling complexity:

- partitioning into subsystems
- preservation of abstractions (hierarchy) in the case of faults
- strict control over interactions (interfaces)

2.1.4 Dependability

² responsive systems:

- RT performance + fault tolerance + distribution of functions

² distributed system:

- error containment regions:
- fault detected + corrected/masked before corrupting the rest of the system
- error is detected at the service interface
- nodes: often EC regions

² node failure modes:

- fail-stop
- fail-silent
- crash

² replication: actively;

- deterministic behavior (replica determinism, also in time)

² subsystems of different criticality:

- critical subsystem and
- non-critical subsystem in different EC regions

2.2 Modeling RT systems

² assumptions used:

- load hypothesis
- fault hypothesis

² time properties:

- actual, minimal duration (of actions)
- worst-case execution time (WCET)
- jitter

2.2.1 Structural elements

- ² task: sequential program execution
 - simple: no synchronisation
 - complex: synchronisation (blocking may occur)
- ² node: self-contained unit with well-defined function
 - abstraction: hw+sw into a single unit
 - SRU: smallest replaceable unit
- ² FTU: fault tolerant unit
 - set of replicated nodes + adjudicator
- ² computational cluster:
 - set of FTU-s (+ gateways)
- ² interfaces:
 - control+data+temporal properties
 - functional intent

2.2.2 RT software

- ² ET systems: interrupts
 - as CPU interrupt frequency increases, also the time
 - with housekeeping (wasted, overhead) increases
- ² TT systems:
 - "sampling" the input
 - predictable overhead
- ² determining worst-case execution time: WCET should be known a priori

- simple task:
 - source code analysis -> critical path – dynamic code? (recursion, loops,...)
 - compiler analysis (timing tree: execution time of high-level constructs)
 - microarchitecture: pipeline, cache?
- complex task:
 - global problems in the system
 - preemption+blocking -> full system model is required
 - solution: annotated source code + instrumentation

² h-state (history state) analysis:

- data that undergoes changes as computation progresses
- fault -> error (changes in state)
- ideal: stateless system (easy to recover)
- cyclic computation: no state transition among cycles

3 Fault tolerant RT systems

3.1 Special properties

² permanence:

- a message is permanent if there are no predecessors which may arrive

² idempotency:

- effect of receiving more copies (of the same message) is the same as receiving a single copy
- replica management is easier

² replica determinism:

- all members have the same visible h-state in time points that are at most an interval of d units apart (d unit: replace a missing message or erroneous message)
- needed to:
 - FT by replication
 - example:
 - node1: commit go n_ no replica determinism
 - node2: abort stop /
 - node3: abort go <- erroneous
 - decision: abort go <- inconsistent
 - system test
- causes of replica nondeterminism:
 - different inputs (digitalization, sensor characteristics)
 - different computational progress relative to physical time (CPU clock drift, FT instruction retry mechanism)
 - preemptive scheduling
 - race conditions
- solutions:
 - sparse time base (no local clock, event is assigned to the same clock tick)
 - agreement on inputs
 - static control structure (no non-deterministic language)
 - deterministic algorithms (no preemption, deterministic race)

3.2 Architectural elements

3.2.1 Node

It should display simple failure modes

3.2.2 FT unit

² fail-silent nodes: duplication

² value errors: replication (TMR)

- ² Byzantine failures: 4 nodes required for a FT unit
- ² Service required: membership (with short latency)
 - ET systems: silence of a node: failure or there is no event?
heartbeat protocol is needed
 - TT systems: periodic message sending defined as membership points

3.2.3 Reintegration of a node

- ² minimal h-state is required to speed up reintegration
 - backward recovery: checkpoint may be invalid due to elapsed time (e.g. sensor data age invalidation)
 - checkpointing mechanism is not enough
- ² ideal reintegration points:
 - after the completion of component cycle
 - after the commit of atomic operations
- ² h-state restoration:
 - retrieve input data from environment (sensors, semaphores etc.)
 - restart vector: control output of the node to synchronize the environment (e.g. yellow traffic lights)
restart vector is defined at development time
 - request data from operator or neighbours

3.2.4 Software issues

- ² What to do to increase dependability:
 - clean structure: simple paradigm (structured programming)
 - formal methods: specification and verification
 - FT schemes: diverse versions of software

² Approaches:

- independent monitoring (case study 1)
- minimal safe service (case study 2)

² Case study 1: VOTRICS tram signaling system

- subsystem1:
 - collecting track data + operator data
 - calculating switch (actuator) positions
 - TMR architecture
- subsystem2: safety bag
 - monitors the safe state of the system
 - evaluates safety predicates (rule book is given)
 - -> blocking output of unsafe signals
 - TMR architecture
- Advantages:
 - independent speci...cations
 - independent implementations (standard program + expert system)

² Case study 2: Airbus fly by wire

- higher level subsystem: full functionality + error detection
- lower level subsystem: reduced but safe functionality

3.3 Real-time operating systems

To do: task management + scheduling + communication + time management

Error detection:

² monitoring task execution times

- does not end in WCET -> error

² monitoring interrupts:

- minimal inter-arrival time must be enforced

² replica management:

- double execution of tasks <- speci...ed in design time

² watchdog functions:

- heartbeat of the node (in the case of fail-silent nodes)

² challenge-response protocol

- calculation of response patterns

3.4 Problems

² flexibility <-> error detection

- error detection requires a priori knowledge of the error-free behavior
- "partial" restriction is needed: e.g. heartbeat
- or replication (deterministic!)

² sporadic data <-> periodic data

- sporadic: dynamic schedule ...ts to it
- periodic: conflict-free (static) schedule

² single locus of control <-> fault tolerance, robustness

- single locus: e.g. token in a token ring
- FT: additional mechanism is required (e.g. token recovery)

² probabilistic access <-> replica determinism

- probabilistic: e.g. Ethernet collision resolution
- replica: identical behavior is required

3.5 System design

3.5.1 Requirement analysis:

Acceptance test to each requirement

- ² performance, deadlines

- ² dependability

- ² cost

3.5.2 Decomposition (architecture)

- ² horizontal structuring: layering (centralized systems)

 - stepwise abstraction

 - -> faults: exception handling

- ² vertical structuring: partitioning (distributed systems)

 - nearly independent subsystems

 - interfaces among the components (low external connectivity)

 - -> faults: error-containment regions

 - (partitioning to be kept even in the case of faults!)

3.5.3 Detailed design and implementation

- ² scheduling, I/O tasks etc.

3.5.4 Test of the design

- ² functional coherence:

 - node = self-contained function

 - minimum h-state

 - error recovery of nodes

 - data sharing interfaces (no control signals)

- timing

² testability

- message interface: all properties should be defined (worst-case scenario)
- h-state observation: modification possibility
- replica determinism (input->output determinism)
- how to test FT properties?
- built-in self-test

² dependability

- node failure -> cluster computation effects (performance + timeliness)
- maintaining a safe state in the node
- if the communication subsystem fails
- detection of a node failure externally?
- internal node error detection -> fail-silency?
- node recovery: time; single failure only
- safety critical functions: in different ECR (err. cont. region)

² physical characteristics

- mechanical interfaces = SRU boundaries = diagnostic boundaries
- SRUs of a FTU are mounted at different locations, avoiding common mode external failures (e.g. mechanical damage)
- SRUs of an FTU should have different power sources, grounding (common mode failures)
- EMI effects via the cabling
- environmental conditions (temperature, shock)

4 Communication

4.1 Requirements

- ² low protocol latency (standard: multicast network)
- ² minimal jitter (e.g. time redundancy)
- ² composability (independent evaluation)
- ² fast error detection
 - blackout: correlated mutilation of all messages (e.g. lightning)
 - babbling idiot: sending messages at wrong moments (TT systems: message exchange at prede...ned points)
 - lost channel: safe state of a node required
 - node error: membership service required (e.g. heartbeat)
 - trashing: too much messages causing breakdown (if the number of messages increases then the throughput will drastically decrease after a given point)
 - end-to-end acknowledgements
 - e.g. message from node A to an actuator, acknowledge from a sensor to node B, acknowledge from node B to node A
 - e.g. Three Mile Island nuclear reactor: valve was not closed, but the monitoring light was green, "never trust an actuator"
- ² Physical structure: multicast is required; bus vs. ring
 - bus: simultaneous arrival of messages
resilience to fail-silent node failures
 - ring: optical ...bers

4.2 Flow control

Controlling the speed of information exchange (receiver can keep up with the sender)

4.2.1 Explicit flow control

- ² Receiver sends acknowledgements: previous message arrived, ready to get the new one
 - ² receiver decides the rate of transmission
 - ² e.g. PAR (Positive Ack or Retransmission) protocol
1. 1. sender (is asked) to send a message
retry count=0
timeout reset
sending the message
 2. 2. receiver gets a message:
was it already sent?
not -> send ack to sender
yes -> send ack + skip message
 3. 3. sender receives ack -> terminates
no ack in timeout period: check retry count
exceeded -> abort
not exceeded -> increment retry count
reset timeout
re-send message
- ² Properties:
 - timeout -> delay may be long!
 - error detection by sender

4.2.2 Implicit flow control:

- ² sender and receiver agree a priori on the message send times
- ² global time base is required
- ² sender sends at the predefined point
- ² receiver accepts all messages

- no acknowledgements
- missing messages are detected by the receiver
- ² FT: active redundancy (multiple messages, multicast)
- ² no trashing (no dynamic scheduling, re-send)

4.2.3 Comparison of explicit and implicit flow control:

Property	Explicit	Implicit	Hard RT req.
control	receiver controls	time controls	receiver may not control
error detection	at the send	at the receiver	at the receiver
trashing	yes	no	to avoid
multicast	difficult	yes	required

Hard RT: the computer system can not control the interface between the controlled object and itself

event shower -> overload -> catastrophic

4.3 Communication architecture

- ² backbone network connecting nodes to non-critical clients (reports etc.)
- ² RT network: connecting the nodes
 - predictable message transmission required
 - support for FT: replicated nodes and channels
 - membership service
 - -> duplicated channels without SPOF (single point of failure)
 - -> babbling detection
- ² ...eld bus: connecting individual nodes to their sensors/actuators
 - periodic transfer
 - strict RT
 - FT is not included in the bus
 - (dependability bottleneck is the sensor/actuator;
 - it is duplicated, connected to different nodes)

5 Case study: The MARS system

Maintainable Real-Time System, TU Wien, 1979-

5.1 Project goals

- ² distributed FT architecture
- ² hard RT
- ² nodes: single-chip communication interface, fail silency
- ² FT properties: replication (FTU)
- ² TT (time-triggered) architecture
- ² Global time base: FT, distributed clock synchronization (VLSI chip)

5.2 Architecture

5.2.1 Distributed RT system

- ² Cluster - FTU - node - task
- ² Communication:
 - ...eld bus: TTP/A protocol
 - RT bus: TTP/C protocol: membership, redundancy management
 - backbone: TCP/IP
- ² Node: host computer + communication controller
 - active: produces images, has bus slot, membership
 - passive: reads only, no bus slot, no membership
- ² Fail silency of nodes:
 - HEDC (High Error Detection Coverage) mode, transparent in OS
 - duplicate execution of application tasks

- end-to-end CRC of messages
- end-to-end CRC of task execution (similar to WP)
- difference: messages are not sent; replicated node is switched

5.2.2 Hardware building blocks

² TTP controller:

- components: dual-port RAM + controllers + EPROM (MEDL)
- connected to replicated buses
- commercial elements are used (COTS)
- TTP/A: four buses
- each node: two comm. controllers
 - TTP/A + TTP/C
 - TTP/C + TTP/C
 - TTP/C + TCP/IP

5.2.3 Software support

² OS: distributed services + local services

- distributed services:
 - clock synchronization
 - membership
 - redundancy management
- local services:
 - static schedule (WCET, WCAO [administration overhead])
 - information transfer: comm. controllers <-> tasks
 - HEDC mode
- cluster compiler:
 - static, deterministic schedule
 - generating message schedule (MEDL: message descriptor list)

- inputs:
 - data elements (length)
 - update period + temporal accuracy
 - sender and receiver ID
 - redundancy strategy (replication: 2,3,...)

5.2.4 Fault tolerance

- ² fail-silent nodes
- ² FTU by replication
 - deterministic message transfer
 - deterministic node operation: static schedule
- ² redundant sensors: master-checker con...guration
 - node1: master of ...eld_bus1, listen to ...eld_bus2
 - node2: master of ...eld_bus2, listen to ...eld_bus1
 - agreement is required

5.3 Time-triggered protocols: TTP/C

5.3.1 Protocol layers

- ² data link/physical:
 - bit encoding
 - bit synchronization
 - media access
- ² SRU layer:
 - implicit acknowledgement
 - clock synchronization
 - SRU membership
- ² RM (redundancy management) layer

- redundancy management
- start-up

² !! Basic Communication Interface !!

² FTU layer:

- permanence of messages
- FTU membership

² Host layer:

- application membership
- application software

5.3.2 Data link layer

² media access: TDMA (time division multiple access)

² controlled by: MEDL (message descriptor list)

- what point in time send a message
- what point in time receive a message
- contains: SRU slot (time), address, direction (I/O), length of message, type of message, parameters

² cluster cycle:

- given number of slots
- every node is assigned a slot
- no data -> empty slot
- cluster cycle = sequence of all TDMA rounds
- TDMA round: k nodes using k frames

5.3.3 SRU layer

² Data frames to comm. interface RAM

² membership:

- every slot is a membership point
- "I am alive" without additional overhead
- delay: 1 TDMA round
- membership info: bit vector

² clock synchronization:

- based on a priori known arrival times of messages
- deviation between expected and real times: clock to be synchronized
- no overhead of checking synchronization

² acknowledgement:

- implicit flow control (based on membership)

5.3.4 RM layer: redundancy management = mode change

² shadow node -> active node
if the current active node fails

² active node -> shadow node
if the active node fails -> fail silency!

² "repaired node" -> shadow node

5.3.5 FTU layer

² permanence of messages

² FTU membership: configurations

- replicated fail-silent nodes
- TMR

5.4 Time-triggered protocols: TTP/A

5.4.1 Properties

- ² standard serial communication (UART)
- ² master: node connected to the ...eld bus
- ² slaves: sensors/actuators

5.4.2 Rounds

- ² ...rework message from the master: one-byte synchronization event + round type determinator
- ² (synchronization) gap
- ² message bytes: determined by MEDL

5.4.3 Error detection

- ² time-out for ...rework message:
 - backup (shadow) node takes the role of the master (active)
- ² time-out of a data byte: local time-out error is reported to the host
- ² data byte outside the speci...ed time window:
 - termination of the round
 - waiting for a new ...rework message
 - no ...rework -> backup node will be active
- ² parity checks (UART)