

2004

Fault Tolerance Tradeoffs in Moving from Decentralized to Centralized Embedded Systems

Jennifer Black

Carnegie Mellon University, jenmorris@cmu.edu

Daniel Kroening

Carnegie Mellon University

Philip Koopman

Carnegie Mellon University, koopman@cmu.edu

Follow this and additional works at: <http://repository.cmu.edu/isr>

Published In

.

This Conference Proceeding is brought to you for free and open access by the School of Computer Science at Research Showcase @ CMU. It has been accepted for inclusion in Institute for Software Research by an authorized administrator of Research Showcase @ CMU. For more information, please contact research-showcase@andrew.cmu.edu.

Fault Tolerance Tradeoffs in Moving from Decentralized to Centralized Embedded Systems

Jennifer Morris
ECE Department
Carnegie Mellon University
jenmorris@cmu.edu

Daniel Kroening
CS Department
Carnegie Mellon University
kroening@cmu.edu

Philip Koopman
ECE Department
Carnegie Mellon University
koopman@cmu.edu

Abstract

Some safety-critical distributed embedded systems may need to use centralized components to achieve certain dependability properties. The difficulty in combining centralized and distributed architectures is achieving the potential benefits of centralization without giving up properties that motivated the use of a distributed approach in the first place. This paper examines the impact on fault tolerance of adding selected centralized components to distributed embedded systems, and possible approaches to choosing an appropriate configuration. We consider the proposed use of a star topology with centralized bus guardians in the Time-Triggered Architecture. We model systems with different levels of centralized control in their star couplers, and compare fault tolerance properties in the presence of star-coupler faults. We demonstrate that buffering entire frames in the star coupler could lead to failures in startup and integration. We also show that constraining buffer size imposes restrictions on frame size and clock rates.

1. Introduction

Embedded systems designers consider many factors when choosing between distributed and centralized architectures. In some designs, centralized systems may be preferred for their manageability and consistency, whereas in others, distributed systems may be more advantageous for scalability and modularity [1].

In the context of dependability, distributed and centralized systems each may be better suited to deliver certain properties. A system with applications distributed across nodes that are physically separated may provide geographic fault isolation; one with a centralized controller may more easily maintain consistent timing and state information across all applications. A system with applications running on separate nodes must also ensure that those nodes have a consistent view of system state. Centralized controllers can become a single point of failure for all applications. The problem with choosing a design that is either strictly distributed or strictly centralized is that one approach might experience problems in system dependability that the other resolves, and vice versa.

In order to achieve the “best of all possible worlds,” system architects may attempt to use a combination of centralized and decentralized components in one system. This approach has recently been seen in some designs based on the Time-Triggered Architecture (TTA). Researchers at the Vienna University of Technology recently proposed

that the TTA with a star topology and central bus guardians, rather than a bus topology with guardians at each node, would eliminate the occurrence of some fault modes that are not tolerated by the Time Triggered Protocol (TTP/C) [2]. In their design, the central bus guardian could have the authority to: a) stop all nodes from transmitting outside of their assigned time slot to eliminate babbling idiot and masquerading faults, b) make adjustments to signal strength and frame timing in order to eliminate slightly-off-specification (SOS) faults [3], and c) perform semantic analysis of all frames transmitted on the network to prevent transmission of frames containing incorrect controller state (C-State) information [2].

Augmenting the authority of some system components, such as the bus guardians in this example, may help the system achieve certain dependability properties; however, faults in those more authoritative components may also have a greater influence on system dependability. For example, suppose a bus guardian suffers a fault that causes it to block transmission of all frames. In systems with decentralized bus guardians (*e.g.*, each node has a separate bus guardian), a fault of this nature in one bus guardian would only block frames from one node. The same fault in a central bus guardian would stop all nodes from sending frames on the channel. This particular fault mode is addressed in [2] by the use of redundant channels with separate central bus guardians. However, this example illustrates that changes in the behavior of one component may alter the way faults in that component affect system dependability.

Adding centralized components to a decentralized system may impact system properties beyond dependability as well. For example, suppose a central bus guardian is required to buffer some minimum number of bits of each frame in order to filter traffic on the network, and this number is proportional to the longest frame. If this central bus guardian is also prohibited from buffering all of the bits in the shortest possible frame, then the difference in length between the longest and shortest frames is limited. It might not be possible to send extremely long frames on the same network as extremely short frames.

Two questions that must be answered when centralization is added to a decentralized system are: “How much authority should the centralized components be given?” and, “what is the impact of this added authority on other parameters of the system?” To help analyze these questions we created a model in SMV (Symbolic Model Verifier) [4] of the TTA star topology with redundant star couplers. In the star couplers we modeled feature sets with varying amounts of centralized control. We also modeled the possible fault modes those features may exhibit. Our objective was to analyze the system-design tradeoffs associated with the cen-

tralized authority by modeling the failure modes that can arise when central guardians have the ability to buffer frames, and by showing how limiting the buffer size of a central guardian also limits the frame size and clock rates. Finally, we propose that critical systems with an active centralized hub for communication, such as a TTP/C network with a star topology, are constrained in clock speed by frame size, and vice versa.

The remainder of this paper is organized as follows: Section 2 presents background information and previous work; Section 3 presents our objective and approach; Section 4 presents detailed information about our system model; Section 5 presents the results from modeling various star-coupler faults; Section 6 presents an analysis of those results; and Section 7 presents our conclusions.

2. Background and related work

The TTA is a distributed system architecture designed for safety-critical embedded systems such as automobiles and jet aircraft that consists of a cluster of nodes communicating over a shared TTP/C network. This section provides background information about TTP/C and the TTA, including the motivation for using central bus guardians in the star topology. We also describe our method for modeling these systems in SMV.

2.1. Overview of TTP/C and the TTA

The Time-Triggered Protocol (TTP/C) and the Time Triggered Architecture (TTA) were developed by researchers at the Technical University of Vienna and TTTech Computertechnik AG for use in safety-critical distributed embedded systems that require dependable frame transmission for hard real-time guarantees. To meet these requirements, TTP/C uses Time Division Multiple Access (TDMA), in which nodes send frames in specific time slots. These slots are statically assigned prior to system start-up in the Message Description List (MEDL) [5].

TTP/C provides a number of services, including distributed clock synchronization, group membership, and clique detection [5]. Clock synchronization is used by the protocol to enforce TDMA scheduling. Each node decides when to transmit on the network based on its view of the time; if different nodes have different views of the global time then the TDMA approach will fail. Group membership is provided as a safety service to allow host applications to monitor the operational states of other nodes on the network. The intention is that nodes that are not operating correctly can be removed from membership in the protocol until they exhibit correct behavior. Clique detection prevents multiple membership groups from forming.

These services are created by nodes generating and/or observing traffic on the network. Clock synchronization, for example, requires each node to observe frames sent by other nodes and calculate the difference between each frame's actual arrival time and the expected arrival time. This allows the observing node to adjust its own internal clock to bring it closer to the clocks of other nodes [5]. Group membership requires each node to verify that frames sent by other nodes are both valid and correct. A valid frame starts and ends during the time slot, exhibits no

encoding rule violations, and is not interfered with by another transmission during the time slot [5]. Correct frames are valid frames that also have a controller state (C-state) and cyclic redundancy check (CRC) that match those of the receiving node [5]. The C-state information may be included in the frame explicitly or implicitly through its inclusion in the CRC calculation. If no activity is observed on the channel during the time slot, the frame is considered to be null (neither invalid nor incorrect).

The algorithms for deterministic message timing, group membership, and clique avoidance are claimed to be correct, given a set of particular constraints [5]. The fault hypothesis of TTP/C allows for an arbitrary failure in a single component of the system. The faulty component may be a single node, bus guardian, or channel. TTP/C assumes fail-silence in the time domain, meaning that a faulty controller will not be allowed to send frames outside of its time slot. TTP/C does not guarantee correctness of the application data contained in the frame.

Ensuring that these requirements are met is handled at the architectural level, as specified by the Time-Triggered Architecture (TTA). The TTA requires at least two independent buses (channels), between nodes. It is assumed that these channels themselves do not generate frames on the network; however, the channels may corrupt or drop frames. Bus guardians are also required in order to prevent faulty nodes from transmitting during the wrong time slot. These bus guardians must be completely independent of the nodes (separate internal clocking device, physical isolation, *etc.*) and may be allocated to nodes individually or centralized in a star topology. Fault tolerance for Byzantine faults [6] requires at least four real member nodes with fully independent bus guardians.

2.2. Central bus guardians

The Time-Triggered Architecture (TTA) may be implemented using either a bus (Figure 1) or a star (Figure 2) topology. In order to compare the fault-handling ability of the two topologies, Ademaj *et al.* [7] performed software implemented fault injection (SWIFI) and heavy-ion fault injection experiments on TTA systems with both bus and star topologies.

The experiments in the system with the bus architecture revealed several ways that faults in a single node could propagate to other, non-faulty nodes. First, slightly-off-specification (SOS) faults in the signal strength or timing of a frame can cause nodes on the network to have differing views of the validity and/or correctness of some frames [3]. A frame that is slightly outside of its time window or has a signal that is slightly below the acceptable range may be judged valid and correct by some nodes and invalid or incorrect by others. This is the result of slight differences in hardware tolerances between nodes. This disagreement between nodes on a frame's validity and/or correctness is an SOS fault.

Group membership is determined by checking the validity of sent frames; therefore a disagreement between nodes on the correctness of a frame could cause some nodes to expel that frame's sender from the group and others to keep the sender in the membership. If this occurs, the TTP/C clique detection algorithm detects the disagreement in group membership and causes the nodes in the minority

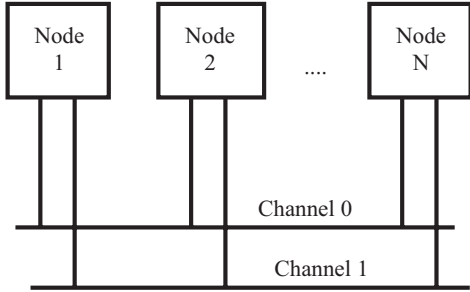


Figure 1: TTA Bus Topology, TTTech Computertechnik AG [5]

clique to enter a freeze state. Nodes that have been frozen cannot regain membership and transmit on the network until they have been awakened by their hosts. This is the correct behavior of the protocol; however, frequent SOS failures could lead to frequent shutdowns of non-faulty nodes in the TTP/C cluster, which is unacceptable in most critical systems. Many of these critical systems, such as automobiles and jet aircraft, are fail-operational, and therefore require a high level of system availability.

Another fault mode in the bus architecture system is masquerading during cluster startup [7]. Prior to startup, nodes in the cluster have not yet established the global time, and cold-start frames that signal the start of a TDMA round do not arrive in a particular time slot; therefore, the sender of a cold-start frame cannot be verified by the arrival time of the frame [2]. If a faulty node sends a cold-start frame with an incorrect sender round slot, other nodes will attempt to integrate into the cluster at the incorrect time. If different cold-start frames arrive at different times on the two channels, nodes may try to integrate on either channel, the cliques will be discovered by the clique-detection algorithm, and nodes in the minority clique will go into the freeze state [2].

The fault injection experiments on the bus topology also revealed that frames with invalid C-states could cause problems for nodes integrating into a running system [7]. If a faulty node sends a frame with an invalid C-state, a non-faulty node that has already integrated determines that the frame is incorrect because the C-state of the frame does not match the internal C-state of the receiving node. Nodes that have not yet integrated do not know the correct C-state of the system, and therefore cannot recognize the frame as incorrect. During integration, they adopt the C-state of the first valid frame they receive and attempt to integrate using that C-state. If the C-state they adopt is incorrect, the integrating node will not be allowed to enter the active state, thereby preventing the perfectly operational node from integrating into the system.

In order to prevent SOS faults, masquerading during startup, and failures in integration arising from frames with invalid C-states, Ademaj *et al.* [7] implemented a star topology with the central bus guardians. In the new design, bus guardians located at the central hubs of the star topology are authorized to perform “active signal reshaping” of the transmitted frames [7]. The central bus guardians monitor frames on the network and boost signals that are SOS in the value domain and delay or block signals that are SOS in the time domain. The central bus guardians also perform semantic analysis of frame content to prevent faulty nodes

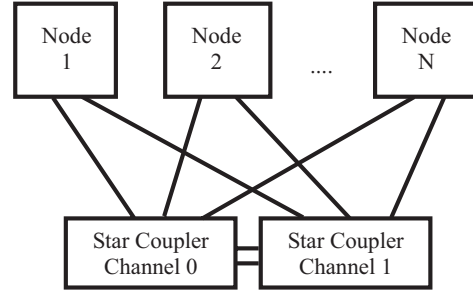


Figure 2: TTA Star Topology, TTTech Computertechnik AG [5]

from masquerading as other nodes during startup and/or transmitting frames with invalid C-states. Fault injection experiments of the new design showed that central bus guardians could prevent SOS faults from occurring. The experiments also showed that the central bus guardians could prevent faulty nodes from transmitting incorrect cold-start frames and frames with incorrect C-state [7].

2.3. Symbolic model checking

We use a symbolic model checker to analyze our formal models. Model checking is a systematic way to explore the whole state space of a formal model. Previous work containing formal models of TTP/C focuses on proving correctness of the protocol [8, 9, 10]. This is done using an interactive theorem prover, which requires a considerable amount of manual work. In contrast to that, model checkers, such as SMV, are fully automatic once the model is written. The model presented in [10] contains more detail than the model we use. However, we are not aware of any formal models addressing the issue of centralized vs. distributed bus guardians.

3. Objective

The fault injection experiments of [7] showed that the central bus guardian could be effective at stopping faults on some nodes from propagating to other healthy nodes. However, these experiments did not analyze the behavior of the system in the presence of faults in the star couplers themselves. Once a component in a decentralized system has been given increased authority over other system components, the system must be reevaluated to ensure that unintended behaviors do not emerge when the centralized authority fails.

For example, TTP/C assumes that faults in either of the two channels are passive. That is, a channel will either corrupt or drop frames, but it will not generate frames. However, this may no longer be a valid assumption if the central bus guardian is given certain authorities. For example, if the central bus guardians have the ability to buffer entire frames, a fault in the central guardian could cause it to transmit frames outside of their intended time slots. If this occurs, the assumption of passive faults in the guardian no longer holds, because the bus guardian is essentially “generating” frames on the network.

In their analysis of failure mode assumption coverage in the TTA, Bauer *et al.* [11] assert that in order to maintain

the fault assumption of passive faults in the channels, the active star couplers must not be given the capability to store frames and transmit them at a later time. They do not, however, explain why frame buffering should be prevented, nor do they explore the system-level implications of that constraint.

In the requirements for the central bus guardian Bauer *et al.* [2] specify the minimum number of bits that must be buffered in order to perform semantic analysis of frames. Semantic analysis is used to stop masquerading during startup and problems in integration due to transmission of frames with invalid C-state. This minimum buffer size is proportional to the number of bits in the longest frame, as well as the difference in clock rates between the nodes and the star coupler. No analysis of the effects of these restrictions on other system parameters, such as frame length and clock rates, is given.

Our objective in this research is to examine the engineering tradeoffs of adding centralized authority to a decentralized system by modeling a system based on the TTA with star topology and central bus guardians. We demonstrate why certain restrictions on the central authority are necessary by assessing the impact on fault-tolerance of full-frame buffering. We also analyze how added centralized authority can limit other system properties by showing that limiting buffer length also restricts maximum frame size and clock rates.

4. TTA star topology model

In order to demonstrate why a central buffer should not be allowed to buffer an entire frame, we compare the fault-handling capabilities of TTA systems with frame buffering at the central bus guardian with systems containing star couplers with less authority.

4.1. Star-coupler feature sets

The four types of star-coupler authority we model are:

- Passive:
 - * does not stop frames
 - * does not shift frames in time
- Time windows:
 - * can open/close bus write access to nodes
 - * does not shift frames in time
- Small shifting:
 - * same authority as *time windows*
 - * also can make slight adjustment to frame timing (e.g., shift slightly ahead to fit window)
- Full shifting:
 - * same authority as *small shifting*
 - * also can buffer frames to make large adjustment to frame timing (e.g., save frames to send out at a later time)

4.2. The formal model

In order to investigate the behavior of TTP/C in the presence of central star coupler faults, we create a synchronous model of the channel, the couplers, and the nodes. Formally, we define a (finite) set of states S , a set of initial states I , and a transition relation R . The transition relation relates two states $x, x' \in S$ if and only if there is a transition from x to x' .

In order to obtain a tractable model, we abstract the behavior by merging all transitions that correspond to a single time slot in the TDMA schedule, *i.e.*, one transition of our model corresponds exactly to one TDMA slot. As the slots may have different lengths, a transition in our model may correspond to different interval in real-time, depending on the particular slot.

The set of states S is the set of possible valuations of the state variables. The state variables consist of the variables the nodes use to store their states and the variables used for the two star couplers. SMV allows specifying I and R using a list of constraints for each. The constraints are conjoined to form I and R , respectively.

4.3. Modeling a node

We model the following parts of the state of a single TTP/C node:

- The TTP/C standard describes a state machine with nine states that governs the behavior of the node: freeze, init, listen, cold start, active, passive, test, await, and download. Initially, all nodes are in the freeze state. The variable used to store this state is denoted by `state`.
- We record the number of correct and bad frames received by the node during a TDMA round. This data is for the benefit of the clique avoidance algorithm. The names of the variables are `agreed_slots_counter` and `failed_slots_counter`, respectively.
- Each node has a variable `slot` containing the current slot number in the TDMA schedule.
- In order to implement the "big bang" cold start algorithm, each node stores a flag `big_bang` that is set if a cold start frame is seen while the node is in the listen state.
- The listen state requires a timeout counter, which is modeled using a state variable for each node. In order to simplify the model, we count the number of TDMA slots using the variable `listen_timeout`.

All other parts of the state of a node are not modeled. This includes, in particular, the application data.

We define the following constraints on the transition relation R . The unprimed variables denote the value of the state variable before the transition, the primed variables denote the value of the state variable after the transition. Furthermore, let `slots` denote the number of slots in the TDMA schedule.

4.3.1. FREEZE and INIT. From the freeze state, the node may make a transition into one of the states `init`, `await`, or `test`. From `init`, it may transition back into the freeze state or may proceed to the listen state. We model these choices nondeterministically.

```
state=freeze => state' ∈ { freeze,
  init, await, test }
```

```
state=init => state' ∈ { freeze,
  initialize, listen }
```

No further constraints are imposed on the variables if the node is in the freeze or init state.

4.3.2. LISTEN. In the listen state, the controller watches the channels for frames to integrate on. The controller may integrate on frames with explicit C-state or on cold start frames. Let the frame type on channel 0 be denoted by `channel0_frame`, and the frame on channel 1 be denoted by `channel1_frame`. Each can be one of `none`, denoting silence, `cold_start`, denoting a cold start frame, `c_state`, denoting a frame with explicit C-state, `bad_frame`, denoting a bad frame, or `other`, denoting a regular frame without explicit C-state.

As described above, nodes do not integrate on the first cold start frame received, but on the second. The variable `big_bang` is used in order to distinguish the first and second cold start frame.

```
big_bang' =
  if state≠listen then false
  else if big_bang then true
  else if channel0_frame=cold_start
    ∨ channel1_frame=cold_start then true
  else false
endif
```

Let `integrating_on_cold_start` denote a shorthand for the condition for integrating on the cold start frame currently on the channel:

```
integrating_on_cold_start =
  state=listen ∧
  (channel0_frame=cold_start ∨
  channel1_frame=cold_start) ∧
  big_bang
```

In contrast to cold start frames, frames with explicit C state are used for immediate integration. Let `integrating_on_C_state` denote a shorthand for the condition for integrating on a C-state frame that is currently on the channel:

```
integrating_on_C_state =
  state=listen ∧
  (channel0_frame=c_state ∨
  channel1_frame=c_state)
```

The node integrates if either condition is true:

```
integrating = integrating_on_C_state ∨
  integrating_on_cold_start
```

If the node is integrating, the time slot counter is set to the value found on the bus plus one. This value is denoted by `id_on_bus`. Thus,

```
(state=listen ∧ integrating) =>
  slot' =
  if id_on_bus=slots then 1
  else id_on_bus+1
endif
```

For startup, the node maintains a timeout counter for the listen state. This counter is initialized with the number of slots plus the number of the slot that is assigned to the node. This slot is denoted by `node_id`. The timeout is reset also if a good frame is seen on the channels. If no correct frame is received, the node counts down the listen timeout counter.

```
listen_timeout' =
  if (state≠listen ∧ state'=listen) ∨
  channel0_frame=cold_start ∨
  channel1_frame=cold_start ∨
  channel0_frame=other ∨
  channel1_frame=other then
  node_id+N
  else if listen_timeout≠0 then
  listen_timeout-1
  else 0
endif
```

If the node is integrating, the node transitions into the passive state. In case of a timeout, it transitions into the `cold_start` state. If there is a cold start frame on either channel that is not used for integration, the node stays in the listen state even if the timeout counter just reached zero.

```
state=listen => state' =
  if integrating then passive
```

```

else if channel0_frame=cold_start
  ∨ channell_frame=cold_start then
  listen

else if listen_timeout=0 then
  cold_start

else listen

endif

```

4.3.3. COLD START. Upon entering the cold start state, the node initializes the slot counter with its own slot number.

```

(state≠cold_start ∧
state'=cold_start) ⇒ slot'=node_id

```

As a shorthand, let `next_slot` denote the number of the next TDMA slot. This is `slot+1` if `slot < slots`, and 1 otherwise. During cold start, the node maintains the slot counter as if integrated.

```

(state=cold_start ∧
state'=cold_start) ⇒
slot'=next_slot

```

```

(state=cold_start ∧ state'=active)
⇒ slot'=next_slot

```

During cold start, the node monitors the frames on the buses. Once one TDMA round is finished, and traffic is observed, the node performs the clique avoidance test. If it succeeds, it transitions into the active state. Otherwise, it transitions back into the listen state.

```

(state=cold_start) ⇒ state'=
if next_slot=node_id then

  if agreed_slots_counter'≤1 ∧
  failed_slots_counter'=0 then
  cold_start

  else if agreed_slots_counter >
  failed_slots_counter then
  active

  else listen

  endif

else cold_start

endif

```

4.3.4. ACTIVE. In the active state, the controller maintains the slot counter:

```

(state=active ∧ state'=active) ⇒
slot'=next_slot

```

```

(state=active ∧ state'=passive) ⇒
slot'=next_slot

```

It may nondeterministically transition to freeze or passive:

```

state=active ⇒ state' ∈ { freeze,
active, passive }

```

Let `frame_sent` denote a shorthand for the frame that is sent by the node:

```

frame_sent=

if state=active ∧ slot=node_id
then c_state

else if state=cold_start ∧
slot=node_id then cold_start

else none

endif

```

4.3.5. The star couplers. Each of the two couplers may have one of the following error states: `none`, corresponding to error-free operation, `silence`, which replaces any frame that is sent on the channel belonging to the coupler by `silence`, `bad_frame`, which places a bad frame or noise on the bus, regardless if a frame was sent or not. Furthermore, the `out_of_slot` fault corresponds to re-sending the last frame received by the coupler.

The fault state of the coupler is denoted by `fault`, *i.e.*, the full names of the variables are `coupler0.fault` and `coupler1.fault`, respectively. The `out_of_slot` fault occurs only if the couplers are configured for full time shifting. All other faults may be caused by any configuration. In accordance with the TTP/C fault hypothesis, we require that at most one coupler has a fault at a given time:

```

couplerA.fault=none ∨
couplerB.fault=none

```

In order to model the behavior of the bus system in case of an out-of-slot error, we add a state variable for each coupler containing the id and type of the frame that was received last. The variable for recording the id is called `buffered_id`, the variable for the type is called `buffered_frame`. The variables are initialized with 0 and `none` respectively.

Let `channel_id` denote the id on the channel belonging to the coupler. Let `channel_frame` denote the frame type. Then,

```

buffered_id'=

if channel_id=0 then buffered_id

```

```

else channel_id
endif

buffered_frame'=
  if channel_id=0 then
    buffered_frame
  else channel_frame
  endif

```

The frame currently on the channel is

```

channel_frame=
  if fault=silence then none
  else if fault=bad_frame then
    bad_frame
  else if fault=out_of_slot then
    buffered_frame
  else if node[i].node_is_sending
    then node[i].frame_sent
  else none
  endif

```

The id on the channel is defined analogously.

5. Experimental results

5.1. Property checked

We use the following correctness criterion: As the nodes are modeled not to fail, no single fault may prevent any node from integrating or losing membership. The TTP/C standard requires that the affected node makes a transition into the freeze state in this situation, *i.e.*, we check that

$$(\text{state}=\text{active} \vee \text{state}=\text{passive}) \Rightarrow \text{state}' \neq \text{freeze}$$

holds on all reachable states.

5.2. Results

For the passive, time windows, and small shifting couplers we verify that the property above holds. For the configuration that allows any star coupler to buffer full frames and replay them in a later time slot, we obtain counter examples from the model checker that demonstrate a failure.

If a property does not hold, the SMV model checker produces a trace from any initial state to a state in which the property does not hold. SMV produces the shortest possible trace. However, the shortest error trace contains four out-of-slot errors. As one might argue that such an accumu-

lation of errors is unlikely, we add a constraint to the model which limits the number of out-of-slot errors to one.

This results in a slightly longer trace, but still produces an error that is caused by a duplicated cold start frame:

- 1) Initially, all nodes are in the freeze state.
- 2) In the next state, all the nodes transition into the init state.
- 3) Node A makes a transition into the listen state. The other nodes remain in the init state.
- 4) The listen timeout counter of node A decreases down to zero. Node B finishes its initialization and transitions into the listen state.
- 5) Node A sends a cold start frame on the bus. Node B ignores the frame due to the "big bang" requirements. The nodes C and D make a transition into the listen state.
- 6) A faulty star coupler replays the previous cold start frame. Node B integrates on it, in compliance with the "big bang" requirements.
- 7) Node A sends another cold start frame. Nodes C and D use this frame to integrate on. Then transition into the passive state.
- 8) Node C sends a C-state frame. Node B considers this frame a faulty frame.
- 9) Node D sends a C-state frame. Node B considers this frame a faulty frame.
- 10) Node B freezes due to a clique avoidance error.

The error may also be triggered by duplicating a C-state frame. We obtain such a trace by adding a constraint which prohibits the duplication of cold start frames:

- 1) Initially, all nodes are in the freeze state.
- 2) In the next state, all the nodes transition into the init state.
- 3) Node A makes a transition into the listen state. The other nodes remain in the init state.
- 4) The listen timeout counter of node A decreases down to zero. Node C finishes its initialization and transitions into the listen state.
- 5) Node A sends a cold start frame on the bus. Node C ignores the frame due to the "big bang" requirements. Node B makes a transition into the listen state.
- 6) Node A sends another cold start frame. Node C integrates on the frame, and transitions into the passive state. Node B ignores it according to the "big bang" requirements. Node D makes a transition into the listen state.
- 7) A faulty star coupler replicates the previous frame into the next slot. Node D integrates on it, making a transition into the passive state. The other nodes consider the frame to be faulty.
- 8) In the next three frames, the nodes A, B, and C each send a C-state frame. Node D considers all to be a faulty frames.
- 9) Node D freezes due to a clique avoidance error.

Both traces are generated in less than a minute on a 1.5 GHz AMD machine.

6. Analysis

The results from our model show that faults that cause frames to be transmitted outside of their assigned time slot can lead to failures during the startup and integration phases of protocol service. Central bus guardians that are allowed to buffer an entire frame are susceptible to these types of faults, even in systems with independent, redundant central guardians. If one bus guardian becomes faulty and sends frames in the wrong time slot and the other bus guardian is behaving correctly, receiving nodes that have already integrated into the cluster will recognize that the frame on the channel with the faulty bus guardian is incorrect because the frame's C-state (which contains the incorrect global time and slot position) does not match its own C-state. These nodes will use the correct frame on the non-faulty channel, or no frame at all if none is transmitted in that time slot. However, nodes that are integrating, either during a cold-start or into a running cluster, are not able to determine that the frame is incorrect, and may use the faulty frame.

It is important to show the problems with system dependability that arise when a central guardian is allowed to buffer an entire frame because there are several reasons why a system architect might be tempted to buffer an entire frame. One reason is cost and ease of implementation. A central guardian could use the same controller as the nodes to receive frames, buffer them, and send them out again on the network. This design would also violate the requirement that the central guardian be prohibited from knowing how to generate a frame (including a valid CRC); however, a central guardian that just receives frames and sends them out again (without re-creating them) would similarly be a simple solution that requires buffering. The addition of new functionality to the central guardian could also lead to buffering of frames. For example, an active central guardian that keeps "mailboxes" with recent data values could help provide data continuity if frames are corrupted by providing slightly stale values instead of no value. A central guardian could also provide prioritized message service (e.g., CAN emulation) if it were allowed to buffer frames and send them in a specially reserved time slice, in priority order. Both of these enhanced functions would require buffering full frames.

In addition to examining the importance of restricting buffer size, it is also useful to analyze the effects of this restriction on other system properties. For example, in the TTA with star topology and central bus guardians, active signal reshaping (in the time and value domains) and semantic analysis of the C-state require the guardian to buffer a minimum number of bits of a frame [2], which is determined from the maximum frame length and the relative difference in rates between the guardian and the nodes.

The restriction on the maximum buffer size (less than the smallest frame) corresponds to a restriction on the relationship between frame lengths and clock rates. If the clock rates of the sending node and the central guardian are different, the central guardian must buffer some part of a frame before it can begin forwarding it. If the central guardian is faster, it must wait to send the frame to ensure that it does not run out of bits to send during the transmission. If the central guardian is slower, it must buffer some of the bits until it has time to send them. This corresponds

to the idea of a "leaky bucket" where the fill rate is not equal to the drain rate. The minimum buffer size B_{\min} is given by:

$$B_{\min} = \delta_{le} + \Delta\rho * f_{\max} \quad (1)$$

where δ_{le} is the number of bits required for line encoding, and f_{\max} is the maximum frame size. $\Delta\rho$ is the relative difference in clock rates of the faster ρ_{\max} and the slower ρ_{\min} :

$$\Delta\rho = (\rho_{\max} - \rho_{\min}) / \rho_{\max} \quad (2)$$

Bauer *et al.* [2] find that the $\Delta\rho * f_{\max}$ term was multiplied by a factor of 2, however the assumptions in the paper that lead to that conclusion are unclear. Therefore, we use equation (1) for this analysis.

It has been shown that a central bus guardian must be prohibited from buffering an entire frame; therefore, the maximum buffer size B_{\max} is limited by the smallest frame transmitted on the network, such that:

$$B_{\max} = f_{\min} - 1 \quad (3)$$

where f_{\min} is the number of bits in the smallest frame. If b_{\min} , δ_{le} , and $\Delta\rho$ are known the largest allowable frame, f_{\max} , is found by setting $B_{\min} = B_{\max}$:

$$f_{\max} = (f_{\min} - 1 - \delta_{le}) / \Delta\rho \quad (4)$$

The maximum frame size, therefore, is inversely proportional to the relative difference in clock rates between the nodes and guardian. Suppose the nominal clock rates of the star coupler and all of the nodes are equal. Even in this case, variations in the manufacturing process could lead to slight variations in the actual clock rates of the nodes and, although very small, $\Delta\rho$ will not be zero. For example, a typical commodity crystal oscillator may have a clock rate that varies by approximately 100ppm. If we assume that in the worst-case scenario the clock in the star coupler is 100ppm fast and the node is 100ppm slow, the difference in clock rates between the two is:

$$\Delta\rho = 2 * (0.0001) = 0.0002 \quad (5)$$

According to the TTP/C Bus-Compatibility Specification [12], the shortest frame in TTP/C, an N-frame with no application data and an implicit CRC, is 28 bits long (4 bits for the mode change request and frame type and 24 bits for the CRC). The minimum cold-start frame requires 40 bits (1 bit for the frame type, 16 bits for the global time, 9 bits for the round-slot position, and 24 bits for the CRC)[5, 12]. The minimum frame with explicit C-state is an I-frame with 48 bits (4 for the mode change request and frame type, 16 bits for the global time, 16 bits for the MEDL position, 16 bits for membership, and 24 bits for the CRC) [12]. Although our model showed that only cold-start frames and frames with explicit C-state caused failures when an out-of-slot fault occurred, the system should be designed to prevent any frame from being transmitted outside of its assigned time slot; therefore, this analysis assumes a maximum buffer length that is less than 28 bits long.

Using the values $f_{\min} = 28$, $\rho = 0.0002$, and $\delta_{le} = 4$ we get:

$$f_{\max} = (28 - 1 - 4) / (0.0002) = 115,000 \text{ bits} \quad (6)$$

The longest allowable frame size for TTP/C is 2076 bits (an X-frame with 4 bits for mode change request and frame type, 96 bits for C-state, 1920 data bits, 48 bits for two CRCs, and 8 bits for CRC padding) [12]. In this example longest allowable frame size of 115,000 bits is much larger than the number of bits in the largest allowable frame.

This calculation is for nodes and central bus guardians with the same nominal clock rate. What happens if the nominal clock rates of the nodes and the central guardian are not equal? It is possible to calculate the maximum allowable difference in clock rates, given set values of f_{\min} and f_{\max} , by rearranging equation (3):

$$\Delta\rho = (f_{\min} - 1 - \delta_{lc}) / (f_{\max}) \quad (7)$$

The smallest possible value of f_{\max} is the size of the largest frame required for protocol operation. This is an I-frame of length 76 bits (4 bits for the mode change request and frame type, 16 bits each for global time, MEDL position, and membership, and 24 bits for CRC). If we set $f_{\max} = 76$ bits, $\delta_{lc} = 4$, and $f_{\min} = 28$ we get:

$$\Delta\rho = (28 - 1 - 4) / (76) = 0.3026 \quad (8)$$

This means that in order to achieve minimal protocol operation, the relative difference in clock rates between the bus guardian and the slowest node can not be more than 30.26%.

Now suppose nodes are allowed to send X-frames of maximum length, the longest frames in the protocol. With $f_{\max} = 2076$, the limit of $\Delta\rho$ becomes:

$$\Delta\rho = (28 - 1 - 4) / (2076) = 0.0111 \quad (9)$$

In this scenario, a system that utilizes the maximum possible frame size cannot have a relative difference in clock rates between the bus guardian and the slowest node that is greater than 1.11%. The situation becomes more constrained if the protocol is altered to allow longer frames, or if the equation in [2] is used.

This dependency between clock rates and frame lengths is not limited to the TTA with a star topology and central bus guardian, but rather, it is a fundamental property of any decentralized system with a centralized supervision of communication that is prohibited from buffering complete frames. Whenever the clock rate of the central supervisor differs from the clock rate of any one of the supervised components, this central supervisor must buffer some part of the communication and the minimum size of this buffer B_{\min} is proportional to the difference in clock $\Delta\rho$ and the maximum frame size f_{\max} . In order to maintain system dependability, the maximum size of this buffer B_{\max} must be less than the minimum frame size f_{\min} . In other words if the difference between the minimum and maximum buffer is large, the difference in clock between the central authority and the guarded nodes cannot be. From equations (2) and (7) we can calculate the ratio of $\rho_{\max} / \rho_{\min}$:

$$\rho_{\max} / \rho_{\min} = (f_{\max}) / (f_{\max} - f_{\min} + 1 + \delta_{lc}) \quad (10)$$

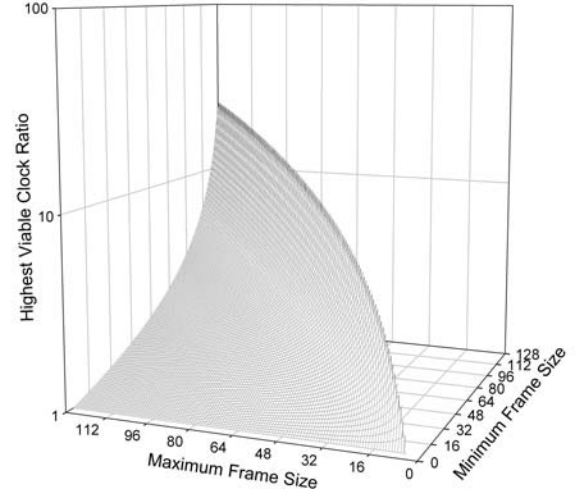


Figure 3: Relationship between frame size range and ratio of clock rates

Figure 3 depicts this relationship for $\delta_{lc} = 4$, where the range of possible values lie below the curve. This graph shows that systems with a wide range of frame lengths cannot also have a wide range of clock rates. The same is true for systems with a short maximum frame size. It is interesting to note that if the maximum and minimum frame size are both 128 bits the ratio of $\rho_{\max} / \rho_{\min}$ is not equal to f_{\max} , but rather is $f_{\max} / 5 = 25$. This is due to the $1 + \delta_{lc}$ term in equation (10). This is a significant limit at high clock ratios when the transmission time of a frame at the high clock speed approaches the δ_{lc} time at the low clock speed.

One simple way to comply with this restriction is to simply require all nodes and bus guardians to have the same clock rate. However, there are several reasons why a system architect might choose to have some nodes have different connection speeds to the hub. Slow, inexpensive nodes (that need a few short frames) might use slow, inexpensive links, whereas fast, capable nodes (that usually generate and consume more network traffic) might have fast links. Different clock rates may also be desirable if the hub is designed to filter traffic. From a TDMA point of view, each frame could get the same size time slice. Frames that are readable to slow nodes would have only a few bits. Frames with more bits would have to run at a higher bit rate, and would only be readable by faster, more capable nodes. An active central guardian could allow faster nodes to exchange longer frames, while at the same time preventing these unreadable frames from being sent to the slower nodes. Unfortunately, it may not be possible to safely implement these systems because the restrictions on buffer size might not allow significant differences in clock rates.

7. Conclusion

We have demonstrated that allocating some types of authority to the central bus guardian can lead to the same failures the authority was intended to prevent. In particular,

systems that are allowed to stop transmission of frames with invalid C-states or invalid cold-start frames and have the authority to buffer entire frames can experience the same problems in startup and re-integration when faults occur in a star coupler as systems without centralized control that experience faults at individual nodes.

The example presented in this paper focused on a synchronous system, the TTA. However, the same effect could also be observed in asynchronous systems. The failures described above were essentially masquerading failures caused when the centralized authority introduced invalid frames containing erroneous identification information onto the network. The same type of masquerading failures could occur in a distributed, asynchronous system because the underlying issue is not timing, but rather identification. A central authority with access to the other nodes' knowledge (e.g., identification methods) may have the ability to introduce masquerading failures into a decentralized system, whether that system is synchronous or asynchronous.

In addition, we have shown that restrictions on the authority of a central buffer correlate to increased restrictions on other system parameters. In the case of the TTA with a star topology and central bus guardians, the restrictions on buffer size in the star coupler lead to increased restrictions on frame size and clock rates. In this example, the restrictions are primarily due to the fact that frame senders are partially identified by the frame timing. The problem is not limited to synchronous systems, however, because some of the failures occurred with cold-start frames, which are used in normal operation before the system has synchronized.

System architects may be tempted to add centralized authority to a decentralized system in order to increase the efficiency, dependability, or functionality of the system. Our results show that this additional authority must be scrutinized to ensure that it does not negatively affect other properties of the system.

8. Acknowledgments

This work is supported in part by the General Motors Collaborative Research Laboratory at Carnegie Mellon University, Bombardier Transportation, and by the Pennsylvania Infrastructure Technology Alliance.

This material is based upon work supported under a National Science Foundation Graduate Research Fellowship. Any opinions, findings, conclusions or recommendations expressed in this publication are those of the authors and do not necessarily reflect the views of the National Science Foundation.

This research was sponsored by the Semiconductor Research Corporation (SRC) under contract no. 99-TJ-684, the National Science Foundation (NSF) under grant no. CCR-9803774, the Office of Naval Research (ONR), the Naval Research Laboratory (NRL) under contract no. N00014-01-1-0796, and by the Defense Advanced Research Projects Agency, and the Army Research Office (ARO) under contract no. DAAD19-01-1-0485.

9. References

[1] P. Verissimo and L. Rodrigues, *Distributed Systems for System Architects*, Kluwer Academic Publishers, 2001.

[2] G. Bauer, H. Kopetz, and W. Steiner, "The Central Guardian Approach to Enforce Fault Isolation in the Time-Triggered Architecture," in *Proceedings of the Sixth IEEE International Symposium on Autonomous Decentralized Systems (ISADS'03)*, Pisa, Italy, April 2003, pp. 37-44.

[3] A. Ademaj, "Slightly-Off-Specification Failures in the Time-Triggered Architecture", in *Proceedings of the Seventh IEEE International Workshop on High Level Design Validation and Test (HLDVT'02)*, Cannes, France, October 2002, pp 7-12.

[4] Model Checking at Carnegie Mellon University <http://www-2.cs.cmu.edu/~modelcheck/>

[5] TTTech Computertechnik AG, *Time-Triggered Protocol TTP/C High-Level Specification Document edition 1.0.0*, TTTech Computertechnik AG, Vienna, Austria, 2002.

[6] L. Lamport, R. Shostak, M. Pease, "The Byzantine Generals Problem," in *ACM Transactions on Programming Languages and Systems*, vol. 4 no. 3, pp. 382-401, 1982.

[7] A. Ademaj, H. Sivencrona, G. Bauer, J. Torin, "Evaluation of Fault Handling of the Time-Triggered Architecture with Bus and Star Topology," in *Proceedings of the 2003 IEEE International Conference on Dependable Systems and Networks (DSN'03)*, San Francisco, California, June 2003, pp. 123-132.

[8] H. Pfeifer, D. Schwier, and F. v. Henke, "Formal Verification for Time-Triggered Clock Synchronization," In *Proceedings of Seventh IFIP International Working Conference on Dependable Computing for Critical Applications*, pp. 207-226, 1999.

[9] J. Rushby, "Systematic Formal Verification for Fault-Tolerant Time-Triggered Algorithms" *IEEE Transactions on Software Engineering*, vol. 25, no. 5, pp. 651-660, 1999.

[10] H. Pfeifer, "Formal Verification of the TTP Group Membership Algorithm". In *IFIP TC6/WG6.1 International Conference on Formal Description Techniques for Distributed Systems and Communication protocols and Protocol Specification, Testing and Verification, FORTE/PSTV 2000*, Pisa, Italy, October 2000, pp. 3-18.

[11] G. Bauer, H. Kopetz and P. Puschner, "Assumption Coverage under Different Failure Modes in the Time-Triggered Architecture", In *Proceedings 7th IEEE International Conference on Emerging Technologies and Factory Automation*, Antibes - Juan les Pins, France, October 2001, pp. 333-341.

[12] TTTech Computertechnik AG, *Time-Triggered Protocol TTP/C Bus-Compatibility Specification edition 1.0.0*, TTTech Computertechnik AG, Vienna, Austria, 2002.