Vrije Universiteit Amsterdam

Bachelor Thesis

# Implementation of the improved Peterson Kearns rollback algorithm

**Author:** Robbert Alexander Rutte (2655858)

*1st supervisor:*   W.J.Fokkink
*2nd reader:*      supervisor name

*A thesis submitted in fulfillment of the requirements for
the VU Bachelor of Science degree in Computer Science*

May 16, 2024

# Abstract

In a distributed computer network it is important that the system is resilient and that the data integrity is maintained, furthermore it should be able to handle concurrent failures and have a consistent state when ending execution. One of the protocols that satisfies these criteria is the Peterson-Kearns Rollback Recovery protocol. The algorithm is based on a vector-based clock to be able to establish a causal order of the messages.

In this thesis the improved Peterson-Kearns Rollback Recovery protocol [1] will be implemented and experimented with. It includes a background study, discusses the limitations and challenges for the implementation of the protocol, implements the protocol and measures its overhead and efficiency. The experimentation will include both the memory overhead, communications overhead and the run-time complexity of the implementation.

# 1 Introduction

In the field of distributed computing, a complex system of processes that participate to accomplish a shared goal, inter-process communication (IPC) is indispensable asset. Yet as crucial as inter-process communication is, it also brings a surplus of complications and complexities that require thorough consideration and planning. From data consistency problems, ensure system resilience to minimizing communications overhead.

The protocol discussed in this thesis aims to tackle the data consistency and system resilience problems. The improved Peterson-Kearns Rollback recovery protocol is able to handle multiple concurrent failures without the program ending in an inconsistent state. In the next section the protocol will be described.

## 1.1 Improved Peterson Kearns rollback algorithm

The protocol provides data integrity and system resilience through the means of of periodically checkpointing the IPC messages to stable storage. however even when operating in a FIFO channel environment a failure of a process could result in an inconsistent state. For a message's send event could be lost to a failure if it is send causally after the last checkpoint.

In the algorithm optimistic recovery is considered, witch indicates that the checkpointing of the nodes is not synchronized. the rollback procedure handles the uncoordinated nodes reach a consistent state. It will do so by handling the *lost messages*, messages where the receive event was lost in a failure at a crashed process or was rolled back to achieve a consistent state, and *orphaned messages*, messages where the send event was lost in a failure at a crashed node or was rolled back. Once the process initiates the restart or rollback procedure the execution is postponed and the messages that arrived are buffered to be handled once the process has restarted or rolledback.

The protocol relies on a vector based clock that is used to determine the causal order of the messages. To ensure the consistency of the system the Peterson-Kearns protocol relies on a logical vector clock to calculate witch events are causally after events that are lost in a failure. each of the processes will maintain a clock vector and a failure vector to ensure the the ability to uniquely identify failures. The failure vector ensures a consistent state when processes fail concurrently. An in-depth description of the protocol will be provided in Section II

## 1.2 Context

Inter-process Communication is a vital part of computer science and facilitates an operating system [2] or a distributed system [3] a way to let different processes communicate and cooperate to accomplish a task. While IPC is inescapable to accommodate cooperation and communication between processes it also numerous problems such as reliability and fault tolerance. A system must be resilient enough to deal with failures to ensure consistency and reliability. The performance overhead must also be optimized since processes that use IPC incur a overhead in computational and memory resources. Furthermore IPC is also inherently more vulnerable to security threats [4], this must be managed to ensure the confidentiality and integrity of the data. A system must also be synchronized to avoid race conditions or a system must have ways to deal with race conditions. If a system uses shared memory across processes this must be well managed to ensure its correctness. IPC is a well established in computer science.

## 1.3 Problem statement

Distributed systems can always be affected by failures, like processes crashing, software errors or outside factors. All of these will be affecting the system consistency and correctness of the algorithm being executed. To prevent the failures from impacting the system implementing a rollback procedure is crucial. The rollback procedure will ensure that any event that is affected by a failure will be undone and forces the system to resume execution from a previously checkpointed consistent state. This will ensure the correctness and the consistency of the system.

## 1.4 Objectives

The objective of this thesis is to implement the improved Peterson-Kearns rollback protocol. The protocol should facilitate the system to:

1. Enable its processes to notify others of a failure that occurred.

2. Be able to respond to these failures.

3. Rollback the events affected by the failure.

4. Resume execution from a consistent state.

5. Provide an efficient way to implement the procedure

6. Minimize the communication and performance overhead of normal execution.

### 1.4.1 Scope

The scope of this thesis includes:

- Implementing the improved Peterson-Kearns protocol.

- Testing the procedure under various scenarios to ensure the correctness of the implementation.

- Assess the performance, memory overhead and the communication overhead.

This Thesis will provide:

- The implementation of the procedure with the source code and documentation.

- The test cases used to evaluate the correctness of the implementation.

- The measurement results and evaluation of the results.

- Indicate areas that could be optimized and enhanced based upon the measurements.

**Constraints**

- The rollback procedure is to designed to minimize the performance overhead, since this is an optimistic recovery algorithm, under normal operation of the implementation.

- Minimize the memory overhead of the implementation.

- Mimic the *read()* and *write()* c functions to facilitate easy usage of the implementation.

- The rollback procedure should maintain a consistent state.

In section 2 a study of the background of rollback algorithms and related work will be done. Afterwards in section 3 the improved Peterson-Kearns protocol will be explained in detail, the design choices will be explained, the limitations of the implementation will be specified and the challenges of the implementation will be mentioned. In the fourth section the methodology will be described, the research goal and the implementation details will be provided. The results of the experimental findings and analysis of these results will be done in the 5th section as well as an indication where more improvements could be made. The conclusion will be located at the 6th section.

## 2 Background

Rollback algorithms originate from the database field of computer science, they were first employed to ensure the consistency of a database system and to assure the atomicity and persistence of transactions when failures occur. The types of rollback algorithms that exist are categorized in two segments: *Immediate rollback*, where the algorithm rolls back the affected events immediately after detecting that a failure has occurred and *Deferred rollback*, where the algorithm defers the rollback of affected events and handles them at a later time.

*Checkpointing* is the periodically saving of events to stable storage. Checkpointing is a central component in every rollback algorithm since when a failure occurs that makes a process crash, it must be reinitialize with part of the memory

otherwise to maintain a consistent state the algorithm would need to rollback all events. Therefore to be able to resume from some earlier consistent state checkpoints of the events that occur at all processes need to be saved to stable storage. Checkpoints frequently include the state of the system and previously generated events.

*Logging* is another critical component of rollback algorithms. A log is kept of the send and receive events that are processed by the process. This log is used when a failure occurs at another process to undo the changes that were made by events that are affected by the failure and restore the state to a consistent state. Logging frequently includes the execution order of the events.

*Coordinated checkpoints* is where all processes of a system coordinate the checkpointing. In that case the algorithm synchronizes the checkpointing of its processes. They are often used of a way to maintain a consistent global state. This enables the protocol to rollback to a certain checkpoint with all its processes which will achieve a consistent state. No further processing of the checkpointed messages will be necessary since they will not exist.

Since the inception of rollback algorithms they have been applied in many fields other than that of database systems, like in distributed systems and operating systems. Rollback algorithms are vital in databases, operating systems, distributed systems and other fields to handle failures and preserve a consistent state.

## 2.1 Related Work

in this subsection we will explore existing rollback algorithms and classify them as immediate rollback or deferred rollback.

### Immediate rollback

An immediate rollback algorithm is a protocol that will apply changes to the system immediately in addition it will maintain a log file where the applied changes are stored. This often involves committing procedures for every event to ensure the consistency of the system. Once a failure occurs the rollback procedure will undo the changes to the system that have not been committed yet. The major advantage of this method is that the recovery process will not be an expensive procedure, thus this method is usually preferred when failures occur often in a system. Examples of the immediate rollback algorithms are:

*Immediate update* described in Principles of Transaction Processing [5]. Where an event is either committed when it executes or it is rolled back immediately. These commits are coordinated with other processes. Once a failure occurs the current event is rolled back to ensure a consistent state.

*Compensation-based recovery* as described in Sagas [6]. In this algorithm when an events affects the state of a system a compensation event is stored in memory. Once a failure occurs the compensation event is executed to revert the state of the system to a consistent state.

### Deferred rollback

A deferred rollback algorithm is a protocol that delays the commit procedure of events. This involves the storing of events in a log. Once the commit procedure is complete the changes to the system stored in the log will be applied. Once a failure occurs the rollback procedure will remove affected events in the log and these events will not be affecting the system. The advantages of this method is that the I/O messages will be kept low since multiple events are committed at once, thus this method is usually used if failures will only occur periodically. Examples of deferred rollback algorithms are: *Vector based rollback-recovery* in Rollback based on vector time [7]. This is the original algorithm that the improved Peterson-Kearns recovery-rollback [1] protocol is based upon. In this algorithm causal order of events is provided by a logical vector clock. Once a failure occurs other processes are notified via tokens and the affected events are rolled back to maintain a consistent state.

*Distributed transaction processing* as described in Checkpointing and Rollback-Recovery for Distributed Systems [8]. In this algorithm a two-phase-commit procedure is used to facilitate the rollback. In the two-phase-commit procedure all processes coordinate the checkpointing to ensure a consistent state.
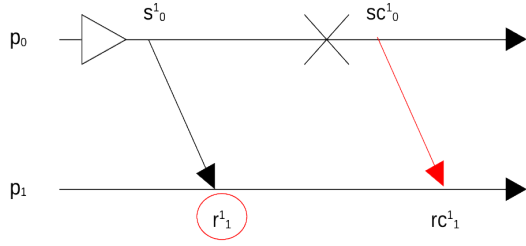
## 3 Framework

In this section an in-depth overview of the implemented protocol will be given. This will be a summary of the information given in the original Msc thesis of van Eck,C [1].
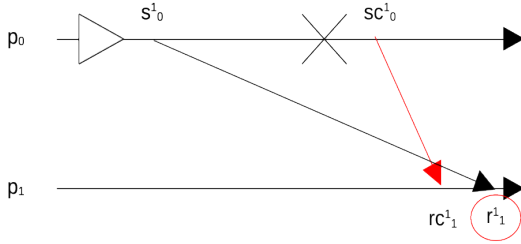
## 3.1 Classifying events

In this subsection all events will be classified. An event is either a sent or a received message both classified as a send event and a receive event respectively. When these events are affected by a failure they will become either an orphan event or a childless event, these will be further classified into orphan-arrived, orphan-in-transit, childless-lost and childless-orphan events.

### 3.1.1 Orphan event



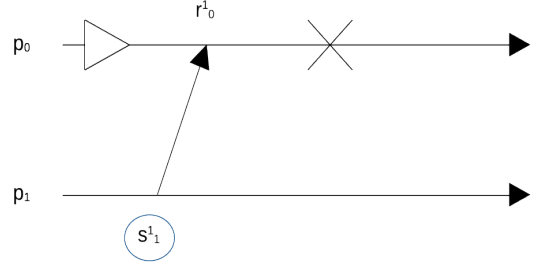**(a)** Orphan-arrived



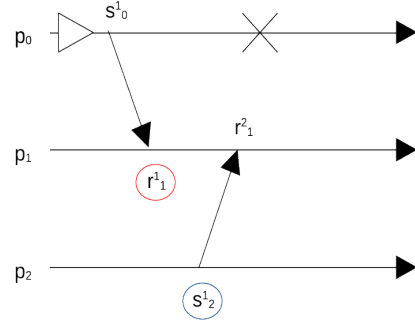**(b)** Orphan-in-transit

**Figure 1:** Orphan events

An *orphan event* is created when the send event of a receive event is lost in a failure, this lost send event will be causally after the last checkpoint at the failed process.

- **Orphan-arrived**: The receive event is causally before the rollback is performed, and is thus before the control message signaling a failure has arrived. Event $r_1^1$ in figure 1a becomes an orphan-arrived event after the failure at $p_0$.

- **Orphan-in-transit**: The receive event is causally after the rollback has been performed, and is thus after the control message signaling a failure has arrived. Event $r_1^1$ in figure 1b is a orphan-in-transit event.

### 3.1.2 Childless event



**(a)** Childless-lost



**(b)** Childless-orphan

**Figure 2:** Orphan events

A *childless event* is created when the receive event of a send event is lost in a failure, this lost receive event will either be itself lost in the failure or will be lost in the rollback of another process.

- **Childless-lost**: The receive event that is lost is causally after the last checkpoint at the failed process, and is thus itself lost in the failure. Event $s_1^1$ in figure 2a becomes a childless-lost when the failure at $p_0$ occurs.

- **Childless-orphan**: The receive event that is lost is causally after an orphaned event at another process, and is thus lost in the rollback procedure at the other process. Event $s_2^1$ in figure 2b becomes a childless-orphan after the failure at $p_0$ occurs.

## 3.2 Algorithm Overview

The improved Peterson-Kearns Rollback recovery protocol [1] has several procedures that are needed for the protocol to function correctly. The protocol employs vector clocks to determine what events are affected by a failure. Each process keeps track of 2 vector clocks, a clock for keeping track of both casual order of the events, the time

vector, and the number of times a process has crashed, the failure vector. Both of these clocks will be sent along with a message to inform other processes of the current state of the vector-clock at the sending process. All events will be stored in volatile memory along with its accompanying vector clocks.

### 3.2.1 Time Vectors

Each process has an associated time vector. The time vector $T = T^0, ..., T^{n-1}$ gets updated whenever an event $e_i$ occurs where $n$ is the number of processes and $i$ is the process id where the event occurred. The $T'$ notation is the time vector before the event $e_i$ occurred.

- When $e_i$ is a send event then the time vector $T$ of a process gets updated to $T^i = T'^i + 1$.

- When $e_i$ is a receive event the time vector $T$ of a process gets updated to $T^j = max(T'^j, e_i.T^j)$ where $j \neq i$ and $T^i = T'^i + 1$

### 3.2.2 Failure vectors

Each process has an associated failure vector, the failure vector $F = F^0, ..., F^{n-1}$ gets updated when a failure occurs in a process or the process get notified of a failure at another process.

In the case of the failure failure happening at the current process will perform the restart procedure, and therein the failure vector gets updated to $F^i = F'^i + 1$ where $i$ is the id of the current process, afterwards the process sends control messages to all other processes to notify them of its failure. This control message includes:

- *id*: The process ID of the failed process.

- *n*: The failure number, indicating the amount of failures that occurred in total at the failed process.

- $T_i^i$: The restart time, indicating the $T^i$ at the moment of restart at the failed process.

- *arrived*: The set of $(T^j, F)$ of the messages received by the failed process that were sent by the process $p_j$ receiving the control message.

In the case that the failure did not occur at the current process, then the current process will be notified of the failure by receiving a control message. When the control message is received the process will perform the rollback procedure,

and therein the failure vector will be updated to $F^j = n$ where $n$ is the value received in the control message for the failure number and $j$ is the process ID of the failed process.

### 3.2.3 Checkpoint procedure

During a checkpoint procedure all information needed to restore the process to a stable state are stored in stable storage. The information to be saved in storage is as follows:

- *state* : the state of a process.

- $T$: The time vector of the process at the moment the checkpoint occurs.

- *log*: The log where the events that occurred at the process are stored.

During a checkpoint procedure no other event can occur.

### 3.2.4 Restart procedure

After a failure at a process it will need to restart from the last checkpoint . The operations needed to perform the restart successfully are:

1. The failure vector is updated to $F^i = F'^i + 1$ where $i$ is the ID of the process where the failure occurred.

2. The process will be restored to the latest checkpoint.

3. Send control messages to the other processes.

During a restart procedure no other events can occur.

### 3.2.5 Rollback Procedure

After a control message is received at a process it will need to perform a rollback to achieve a global consistent state. The operations needed to perform the rollback successfully are:

1. The failure vector is updated to $F^i = F'^i + 1$ where $i$ is the ID of the process where the failure occurred.

2. If some event is orphaned at the process that received the control message, then the state is returned to the last checkpoint that does not include orphaned events. All events after this checkpoint that are not orphaned will be replayed.

3. Any receive events that have been orphaned but have a corresponding send event that was not orphaned will have their receive event replayed.

4. All messages that this process sent to the process where the failure occurred with a non orphaned send event that did not arrive will be re-transmitted.

5. A new checkpoint will be created.

During the rollback procedure no other events can occur.

## 3.3 Limitations

The limitations of this algorithm are that the maximum number of messages that can be sent by one process will be the upper-limit of the integers used in the time vector, this could be mitigated by allowing the time vector values to roll over back to zero in some way but this would require more logic to accommodate for this.

Another limit of this algorithm would be that when the amount of events stored by a process becomes large and a failure occurs at that process it will need to sent a large set of messages that have arrived at that process. This could be mitigated by occasionally checking what messages have been checkpointed by all processes and then committing these events. After the set of arrived messages will not be needing messages in it that have been committed. This will be explained more in-depth.

If a task needs to be completed that will generate a large amount of event the volatile memory of a process could reach the allowed limit fast. to mitigate this some events could be removed if they have been committed and will not be needed for re-transmission. This would shrink the volatile memory needed for continuous operation. This will be explained more in-depth.

### 3.3.1 Commit procedure

When a commit procedure is initiated by a process commit messages will be sent between processes. A committed event will be stored in a separate file in stable storage. These messages will be as follows:

1. the process initiating the commit procedure will signal all other processes.

2. When this signal is received at a process it will send the initiating process the $T^i$ of its last *checkpoint* where $i$ is the process ID.

3. After all of the aforementioned messages have been received the initiating process will construct a minimal checkpointed time vector from these values with its $T^i$ of its last *checkpoint*.

4. All events that have a time vector that is $\leq$ the minimal checkpointed time vector can be committed.

5. The initiating process will send this minimal checkpointed time vector to the other processes. Along with the list of committed receive events that have been sent by the recipient of this message, this is needed for the removal of send events.

6. All processes will now be able to commit events, if the process is not the initiator of the commit procedure it will send all other processes a list of of committed receive events that have been sent by the recipient of the message, this is needed for the removal of send events.

in between commit messages other event can occur.

### 3.3.2 Remove procedure

After all of the other processes have sent the list of committed receive events to a process it can perform the remove procedure. *NextCk(e)* is the checkpoint after the one event $e$ is stored in. the events that can be removed are as follows:

- *Messages stored in volatile storage*: all events $e$ that have *NextCk(e)* $\leq$ the minimal checkpointing time vector can be removed. if $e$ is a sent event its corresponding receive event needs to be contained in the set of committed messages received during the commit procedure.

- *Checkpoints stored in stable storage*: All checkpoints $ck_i$ can be removed if *NextCk(ck_i* $\leq$ the minimal checkpointing time vector.

- *Failures stored in stable storage* : these can only be removed when FIFO channels are used.

## 3.4 Design Choices

Since one of the requirements of the implementation is that operation without failures occurring should have a small overhead the implementation will be written in C++. In addition to that to ensure fast operation of vector logic the vector package from the C++ standard library was used. To ensure low overhead on the sending and receiving procedures no loops were used in this part of the implementation. The functions to send and receive messages are modeled after the C write and read functions to facilitate simple usage of the implementation.

An implementation will be made for the base protocol without the commit and remove procedures, as well as a version with the commit and remove procedures to test them against each other for a difference in performance.

# 4 Methodology

## 4.1 Research Goal

The goal of the experimental setup is that the implementation will be tested on performance. The performance tests that will include the testing performance of both the implementation with and without the commit and remove procedure. For both of these the performance of the sending and receiving of messages will be measured. The performance of the Restart and Rollback procedure will also be measured. The performance of the rollback procedure will also be measured when dealing with concurrent failures from another process. Afterwards the measurements of both of the implementations will be compared.

## 4.2 Implementation Details

The implementation of the improved Peterson-Kearns rollback recovery protocol is written in C++ version 20 and is provided in appendices A and B.

### State class

The State class contains all components needed to maintain a consistent state. the main components of the class are the id of the process, the number of messages in volatile storage, the event log containing the message events, time vector, failure vector, the file descriptors of all other processes(these will be used to resend lost messages), checkpoints and checkpoint time vectors, a set arrived messages and a set of arrived control messages. The version of the class where committing and removing messages is possible the State class also includes a vector to track if the commit or remove procedures can be executed, a set of committed messages, a set of committed receive events form the other processes and a minimal time vector to calculate the messages that can be committed and removed. The most important methods of the State class are listed below.

Public methods:

- **State:** Initializes the class. If the initialization is executed with the restart boolean set to true it will execute the restart operation and sent controll messages to other processes.

- **checkpoint:** Performs a checkpointing procedure .

- **send_msg:** Sends a message and increments the time vector. This will only send messages of type MSG

- **recv_msg:** Receives all message types and handles them accordingly, it will receive normal messages that increments the time vector, it will receive control messages and initiate the rollback procedure, it will receive all commit messages and process them accordingly and it will receive void messages which will not affect the state of the process and will only be passed through to the caller of recv_msg. The method uses c function read to read from its read file descriptor.

- **send_ctrl:** Sends a control message to all other processes to indicate a failure has occurred. The methods uses the c function write to write to the file descriptor of other processes.

- **signal_commit:** Signals to all other processes to start the commit procedure.

- **update_fd:** Updates the file descriptors of a process.

The public methods also include methods to retrieve the stored event log.

In addition to these public methods the class has several private methods that are used by the class to operate. The most important private methods of the State class are listed below:

Private methods:

- **check_duplicate:** Checks if a received message is a duplicate .

- **check_duplicate_ctrl:** Checks if a recieved control message is duplicate.

- **check_duplicate_commit:** Checks if a message has already been committed.

- **check_orphaned:** Checks if a received message is orphaned.

- **rem_log_entries:** Removes log entries from from the event log.

- **rem_checkpoints:** Removes checkpoints.

- **store_msg:** Stores a message in the event log.

- **commit_msgs:** Commits the indicated messages.

- **rollback:** Performs the rollback procedure. This procedure could resend messages that were lost in the failure to the failed process.

- **send_commit:** Sends its own the time vector value to the process that initiated the commit procedure.

- **commit :** Performs the commit procedure and sends the committed receive events to all other processes.

- **remove_data:** Performs the remove procedure.

The private methods also include serialization methods and methods to calculate various pointers and indices. The **send_void** function is included in the namespace Pet_kea. It send a message to a process that does not increment the time vector and is not stored in the event log.

the class is able to send the types of messages listed below:

- **MSG:** This is a message that will be processed by the protocol, it increments the time vector of the process. It includes the time vector and failure vectors at the time of sending, the ID of the sending process and the message size in bytes.

- **CTRL:** This is the message that will be send if a failure has occurred, it will increment the failure vector of the process receiving it. It includes the failure log entry, the ID of the sending process, the amount of messages the sending process has received from the receiving process of the control message and the time vector value and failure vector of these messages. this message can only be sent after a process has restarted.

- **VOID:** This is the message that will not impact either the time vector or the failure vector and will be passed trough to the caller of the **recv_msg** method. And thus not affecting the state of the protocol. This would be used when a message has to be sent between processes that can not alter the state.

- **COMM1:** This is the message that will be sent to initiate the commit procedure

- **COMM2:** This message will contain the its time vector value of its latest checkpoint.

- **COMM3:** This message will contain the combined time vector values of the latest checkpoints of all the processes. These will be used to determine which events can be committed and removed.It also contains the receive events from messages send by the receiver of this COMM message that have been committed by the sender of the COMM message in addition to this the amount of these receive events in the message. This message will be sent by the initiator of the commit procedure only.

- **COMM4:** This message will contain the receive events from messages send by the receiver of this COMM message that have been committed by the sender of the COMM message in addition to this the amount of these receive events in the message. This message will be sent by the non-initiators of the commit procedure only.

The event log consists of an array of structs that contain the send and receive events that are stored, message size, the ID of the message that send the message, the time vector at the moment of sending, the time vector at the moment of receiving and the failure vector at the moment of sending. it also contains the index of the next checkpoint and a boolean indicating if the messages is a send or a receive event.

The checkpointing file in stable storage is organized as follows: at the start of the file an integers will be stored with the amount of total checkpoints in the file. Afterwards the checkpoints will

be stored. the header of the checkpoint will consist of two integers and the time vector of the checkpoint. The initial integers will consist of the messages stored in total at the time the checkpoint will be made and the messages stored in total at the checkpoint before this one. The actual checkpoint data will be the serialized msg_log_t.

The committal file will be populated with the serialized msg_log_t of the committed messages.

## 4.3   Experimental Setup

The experimental setup that was used to measure the overhead cost of the implementation operates as follows:

A parent process creates pipes for its child processes to communicate and then forks them. The child processes will then randomly send each other messages using the State class. The parent process is then able to interrupt one of its children by sending them a SIGINT interrupt. This interrupt is then caught by the child process which will then complete its current select loop and will then exit. After the child has exited the parent will restart the process by forking it again. A process will transmit its new file descriptors over UNIX domain sockets. The child processes will use a select loop to deal with the incoming messages of the other processes.

The experiments will be run on a dell XPS 13 using a intel i7-8565U CPU running at 1.80GHz, the operating system used was Arch-Linux 2024.05.01 release with the included kernel version 6.8.8. Both the implementation of the protocol and the experimental framework is written in C++ version 20.

The measurements are taken using the chrono::high_resolution_clock from time.h. The measurements of the overhead cost for sending and receiving and checkpointing are taken by only measuring how long it took to complete the **send_msg**, **recv_msg** and **checkpoint** methods. For the restart time cost measurements are taken by measuring the time it took for the class to be initialized and transmit all control messages and thus resume normal execution using increasing number of messages stored in stable storage. For the rollback time cost measurements is the time it takes from the moment the process receives the control message until it has handled the rollback procedure and resumed normal operation, each subsequent time using increasing number of messages stored in its event log to determine what effect this has on the rollback pro-

cedure. The communication overhead and stable storage memory overhead is calculated using the serialized size of the message and event log. The RAM memory overhead is calculated by the used memory space by the event log struct and the size of the accompanying vectors. The measurements are saved in a CSV file to be able to analyses the results more easily.

The implementation of the experimental setup can be found at GitHub [9].

In addition to the performance test a correctness test will be performed, various scenarios will be checked to indicate the correctness of the implementation.

# 5   Results and Analysis

### 5.0.1   Memory overhead

The implementation of the improved Peterson-Kearns protocol both for storing messages in stable storage and in communication utilizes serialized data. The overhead for storing messages in stable storage:

$$3 * sizeof(int) + 3 * sizeof(int) * n \ bytes$$

the communications overhead:

$$3 * sizeof(int) + 2 * sizeof(int) * n \ bytes$$

and the overhead for storing messages in random access memory:

$$96 + 3 * sizeof(int) * n \ bytes$$

where

$$n = number \ of \ processes$$

This shows that when this implementation is used by a large amount of processes that the memory overhead will increase substantially. This overhead could be diminished if an int16_t would be used to store the time and failure vectors, however the downside to this would be that the implementation would be able to handle less messages.

## 5.1   Experimental Findings

### 5.1.1   Sending performance

The two figures 3 , 4 show the 200 moving average of time cost of sending a message. The sending cost is tested using different amounts of processes communicating. In the first figure the processes were committing and removing events

**Figure 3:** Sending cost for one message : commit procedure (200 MA)



**Figure 4:** Sending cost for one message (200 MA)

as well as checkpointing the events, in the second figure the processes were only checkpointing the events. This shows that the sending performance is not affected by the increase in processes communicating. It also shows that there is not a big difference in sending cost between weather the messages will be committed or not, however this is to be expected since the operations used for sending a message do not change between the two versions of the protocol.

### 5.1.2 Receiving performance



**Figure 5:** Receiving cost for one message : commit procedure

The two figures 5, 6 above show the receiving cost of receiving a message. The receiving cost is tested using different amounts of processes com-



**Figure 6:** Receiving cost for one message (200MA)

municating. The first figure shows the performance while the processes are committing and removing events. for clarity only the plot for 2 communicating processes is shown. This shows that there are numerous spikes in the time cost if we look at the second figure we can see why this occurs. In the second figure a 200 moving average is used. Here we can see that we do not have such spikes in time cost and can thus conclude that the committing and removing events is a costly procedure. we can also see that the receiving performance is not affected by the increase in processes communicating.

### 5.1.3 Restart performance



**Figure 7:** Average restart cost

The figure 7 above shows the time cost of the restart procedure. The blue line shows the time it took from the moment the process called the constructor of the State class until normal operation was resumed, the red line the amount of messages in stable memory that where stored by the process performing the restart. This shows that the time cost of the restart procedure is proportionate to the amount of messages that are stored in stable storage, furthermore it shows that restarting is relatively a cheap procedure.

### 5.1.4 Rollback performance



**Figure 8:** Rollback cost



**Figure 9:** Average rollback cost

The figures 8, 9 above shows the time cost of the rollback procedure. The line in blue shows the time it took from the moment the process received the control message until normal operation was resumed, the red line shows the amount of messages in volatile memory that where stored by the process performing the rollback. This shows that this is a costly procedure since it takes almost 0,1 second most of the time. However at the 1st, 4th, 6th and 10th time the rollback was performed shows a discrepancy, it takes considerably less time to complete the rollback. This is due to that there has not been a checkpoint in between the other process crashing and the process performing the rollback to receive the control message and thus it will not have to remove any checkpoints. The only thing the rollback then does is to replay, resend and remove messages from volatile memory. This indicates that the removal of checkpoints is a costly operation.

In figure 10 the rollback cost when dealing with concurrent failures, all odd number rollbacks are the from the initial failure the even rollback numbers indicate the concurrent failure after the initial one.

This shows that the concurrent failure will cause the rollback to be performed much faster



**Figure 10:** Concurrent rollback cost

than the rollback performed after the initial failure. This is logical since the rollback procedure has already been performed, only message that needed to be resent will again need to be resent, the other parts of the rollback have already been performed at the rollback caused by the initial failure.

### 5.2 Correctness indication

To provide a Correctness indication a couple of scenarios will be tested:

- *Duplicate*: Duplicate messages will be handled by checking if a set contains the time and failure vectors, afterwards the vectors will be inserted into the set. These will be removed when checkpoints get removed and no more duplicate messages can be sent and is thus not needed any more. The implementation for this is provided in appendix A:93 and A:1372

- *Duplicate-control*: Duplicate control messages will be handled by checking if a set contains its id, failure_number and restart_time, afterwards these values will be added to the set as an entry. The implementation for this is provided in appendix A:106.

- *Orphan-arrived*: Whenever a rollback procedure is called orphaned events will be removed from the message log. An event $e_i$ is orphaned if $e_i.T^i > restart\_time$. if the event $e_i$ is a receive event then $e_i.T^i_{sender} > restart\_time \land e_i.F^i_{sender} < failure\_number$. Restart_time and failure_number are the values contained in the control message. The implementation for this is provided in appendix A:880.

- *Orphan-in-transit*: Whenever a message is received the check_orphan method is called.

12

This method will check if there exists a entry$(j, n, T^j)$ in the failure log that has $msg.F^j < n \land msg.T^j > T^j$, if this exists the message is a orphan-in-transit message and will be ignored. the method will be provided in appendix A:129.

- *Childless-lost*: To resolve childless-lost event all send events that are not lost or orphaned but have their receive event lost in the failure will eb re-transmitted. The implementation for this is provided in appendix A:966

- *Childless-Orphaned*: To resolve childless-orphan events all send events that are not lost or orphaned and have their corresponding receive event is orphaned in the rollback procedure will be replayed. The implementation for this is provided in appendix A:931.

This provides a rudimentary correctness indication for the implementation.

## 5.3   Future Work

For future implementation work the removal of checkpoints could be improved since it now takes a considerable performance penalty to evoke the method. Another point that could be improved is the return value when a process receives a control message and performs the rollback, right now it is just indicated that a rollback has been performed and not what messages have been removed during the rollback. Another improvement that could be made is a more specified error handling with setting the errno correctly. Some changes could also be made to add more constructors to the msg_log_t like the copy and assignment constructors. The serialization of the data to be send and stored could also be handled by some library like Protobuf [10] to facilitate streamlined compatibility with other programming languages. The way committed messages are stored could be altered to a more structured state to provide processes using the committed data an easy understanding of the way the data is stored. Another change that would be welcome is the addition of shadow paging this would increase the fault tolerance of the implementation to ensure data is still valid if a crash occurred in the middle of a write operation.

## 6   Conclusion

With this we can conclude that the memory overhead does increase proportionate to the number of processes communicating, however we can also conclude that the sending and receiving performance is not greatly affected by the number of processes communicating. The rollback performance as well as the commit receive performance shows that the removing of checkpoints a costly operation. In general is rolling back a costly operation while restarting is a relatively cheap operation.

# A   pet-kea.cpp

```cpp
1  #include "pet-kea.hpp"
2
3  using namespace std;
4  typedef unordered_set<vector<int>, vector_hash> uset;
5
6  void get_fail_filename(int process_nr, char fail_filename[32])
7  {
8      sprintf(fail_filename, "fail_v_process_%d.dat", process_nr);
9  }
10
11 void get_msg_filename(int process_nr, char msg_filename[32])
12 {
13     sprintf(msg_filename, "msg_process_%d.dat", process_nr);
14 }
15
16 void get_comm_filename(int process_nr, char commit_filename[32])
17 {
18     sprintf(commit_filename, "commit_process%d.dat", process_nr);
19 }
20
21 void get_state_filename(int process_nr, char state_filename[32])
22 {
23     sprintf(state_filename, "state_process%d.dat", process_nr);
24 }
25
26 void Pet_kea::print_msg(struct msg_t *msg)
27 {
28     cout << "MSG  time vector(";
29     for (int i = 0; i < (int)msg->time_v.size(); i++)
30     {
31         cout << msg->time_v[i];
32         if (i < (int)msg->time_v.size())
33         {
34             cout << ",  ";
35         }
36     }
37     cout << ") fail vector(";
38     for (int i = 0; i < (int)msg->fail_v.size(); i++)
39     {
40         cout << msg->fail_v[i];
41         if (i < (int)msg->fail_v.size())
42         {
43             cout << ",  ";
44         }
45     }
46     cout << ")" << endl;
47 }
48
49 void Pet_kea::print_ctrl_msg(struct ctrl_msg_t *msg)
50 {
51     cout << "CTRL with log(" << msg->log_entry.id << ",  " << msg->log_entry.fail_nr << ",  ";
52     cout << msg->log_entry.res_time << ") messages recieved:  " << msg->recieved_cnt << endl;
53     for (set<pair<int, vector<int>>>::iterator ptr = msg->recieved_msgs.begin();
54          ptr != msg->recieved_msgs.end(); ptr++)
55     {
56         cout << "        Tj:  " << ptr->first << "  fail_v:  ";
57         for (int j = 0; j < (int)ptr->second.size(); j++)
58         {
59             cout << ptr->second[j] << ":";
60         }
61         cout << "  lenght:  " << ptr->second.size() << endl;
62     }
63 }
64
65 int *Pet_kea::State::next_checkpoint(int *ptr)
66 {
67     int to_skip = *ptr - *(ptr + 1);
68     ptr += 2 + time_v.size();
69
70     for (int i = 0; i < to_skip; i++)
71     {
72         ptr += ((SER_LOG_SIZE + *ptr) / sizeof(int));
73     }
74     return ptr;
75 }
76
77 char *Pet_kea::State::get_msg(int i)
78 {
79     return msg_log[i].msg_buf;
80 }
81
82 char **Pet_kea::State::get_msg_log()
83 {
84     char **output_log = (char **)malloc(msg_cnt * sizeof(char *));
85     for (int i = 0; i < msg_cnt; i++)
86     {
87         output_log[i] = msg_log[i].msg_buf;
88     }
89
90     return output_log;
91 }
92
93 bool Pet_kea::State::check_duplicate(struct msg_t *msg)
94 {
95
96     vector<int> merged_time_fail_v;
97     merged_time_fail_v = msg->time_v;
98     merged_time_fail_v.insert(merged_time_fail_v.end(), msg->fail_v.begin(), msg->fail_v.end());
99     if (arrived_msgs.contains(merged_time_fail_v))
100        return true;
```

```cpp
101
102        arrived_msgs.insert(merged_time_fail_v);
103        return false;
104    }
105
106    bool Pet_kea::State::check_duplicate_ctrl(struct fail_log_t log)
107    {
108        vector<int> fail_log_vector{log.id, log.fail_nr, log.res_time};
109        if (arrived_ctrl.contains(fail_log_vector))
110            return true;
111
112        arrived_ctrl.insert(fail_log_vector);
113        return false;
114    }
115
116    bool Pet_kea::State::check_duplicate_commit(struct msg_log_t *l_msg)
117    {
118        vector<int> merged_time_fail_v;
119        merged_time_fail_v = l_msg->time_v_sender;
120        merged_time_fail_v.insert(merged_time_fail_v.end(),
121                                  l_msg->fail_v_sender.begin(), l_msg->fail_v_sender.end());
122        if (committed_msg_set.contains(merged_time_fail_v))
123            return true;
124
125        committed_msg_set.insert(merged_time_fail_v);
126        return false;
127    }
128
129    bool Pet_kea::State::check_orphaned(struct msg_t *msg)
130    {
131        for (int i = 0; i < (int)fail_log.size(); i++)
132        {
133            if (msg->fail_v[fail_log[i].id] < fail_log[i].fail_nr &&
134                msg->time_v[fail_log[i].id] > fail_log[i].res_time)
135                return true;
136            else
137                continue;
138        }
139
140        return false;
141    }
142
143    void Pet_kea::State::rem_checkpoints(vector<int> to_remove)
144    {
145        // read file and reconstruct it
146
147        vector<int> reverse_to_remove = to_remove;
148        reverse(reverse_to_remove.begin(), reverse_to_remove.end());
149        char filename[32];
150        get_msg_filename(id, filename);
151        msg_out.close();
152        ifstream msg_in(filename, ifstream::in | ifstream::binary);
153        msg_in.seekg(0, msg_in.end);
154        size_t file_size = msg_in.tellg();
155        msg_in.seekg(0, ifstream::beg);
156        char msg_file[file_size];
157        msg_in.read(msg_file, file_size);
158        msg_in.close();
159
160        char new_msg_file[file_size];
161
162        int *old_ptr = (int *)msg_file;
163        int *new_ptr = (int *)new_msg_file;
164
165        old_ptr++;
166        new_ptr++;
167
168        int checkpoint_msg_cnt = 0, checkpoint_last_ckpnt = 0;
169        int checkpoint_cnt = (int)checkpoints.size();
170        vector<int> new_checkpoints;
171        vector<std::vector<int>> new_ck_time_v;
172        new_checkpoints.push_back(0);
173        new_ck_time_v.push_back(ck_time_v[0]);
174        int *temp_ptr;
175
176        for (int i = 1; i < checkpoint_cnt; i++)
177        {
178            if (!reverse_to_remove.empty() && reverse_to_remove.back() == i)
179            {
180                reverse_to_remove.pop_back();
181                old_ptr = next_checkpoint(old_ptr);
182                continue;
183            }
184
185            temp_ptr = old_ptr;
186
187            checkpoint_last_ckpnt = checkpoint_msg_cnt;
188            checkpoint_msg_cnt += *old_ptr - *(old_ptr + 1);
189            *new_ptr = checkpoint_msg_cnt;
190            new_ptr++;
191            old_ptr++;
192            *new_ptr = checkpoint_last_ckpnt;
193            new_ptr++;
194            old_ptr++;
195
196            temp_ptr = next_checkpoint(temp_ptr);
197
198            memcpy(new_ptr, old_ptr, (temp_ptr - old_ptr) * sizeof(int));
199
200            new_checkpoints.push_back(checkpoint_msg_cnt);
201            new_ck_time_v.push_back(ck_time_v[i]);
202
203            new_ptr += (temp_ptr - old_ptr);
```

```cpp
204                old_ptr = temp_ptr;
205        }
206
207        checkpoints.swap(new_checkpoints);
208        ck_time_v.swap(new_ck_time_v);
209
210        size_t test = (size_t)new_ptr;
211        size_t test2 = (size_t)(int *)&new_msg_file;
212
213        file_size = new_ptr - (int *)&new_msg_file;
214        file_size = test - test2;
215
216        new_ptr = (int *)new_msg_file;
217
218        *new_ptr = checkpoints.size() - 1;
219
220        msg_out.open(filename, ofstream::out | ofstream::binary | ofstream::trunc);
221        msg_out.write(new_msg_file, file_size);
222        return;
223 }
224
225 int Pet_kea::State::next_checkpoint_after_rem(vector<int> removed_checkpoints, int curr_next_checkpoint)
226 {
227        int result = curr_next_checkpoint;
228        for (const int &i : removed_checkpoints)
229        {
230            if (i < curr_next_checkpoint)
231                result--;
232            else
233                break;
234        }
235        return result;
236 }
237
238 int Pet_kea::State::rem_log_entries(vector<int> to_remove, int final_index)
239 {
240        cout << id << "_old_msg_cnt:" << final_index << "_removing:" << to_remove.size() << endl;
241        msg_log_t *new_log = (msg_log_t *)calloc(MAX_LOG, sizeof(msg_log_t));
242
243        vector<int>::iterator curr = to_remove.begin();
244        int new_final_index = 0;
245        for (int i = 0; i < msg_cnt; i++)
246        {
247            if (curr != to_remove.end() && i == *curr)
248            {
249                curr++;
250                free(msg_log[i].msg_buf);
251                vector<int>().swap(msg_log[i].time_v_reciever);
252                vector<int>().swap(msg_log[i].time_v_sender);
253                vector<int>().swap(msg_log[i].fail_v_sender);
254                continue;
255            }
256            new_log[new_final_index] = msg_log[i];
257            new_final_index++;
258            free(msg_log[i].msg_buf);
259            vector<int>().swap(msg_log[i].time_v_reciever);
260            vector<int>().swap(msg_log[i].time_v_sender);
261            vector<int>().swap(msg_log[i].fail_v_sender);
262        }
263        free(msg_log);
264        msg_log = new_log;
265        cout << id << "_new_msg_cnt:" << new_final_index << endl;
266        return new_final_index;
267 }
268
269 void Pet_kea::State::serialize_commit(struct comm_msg_t *msg, char *data)
270 {
271        int *q = (int *)data;
272        *q = (int)msg->msg_type;
273        q++;
274
275        *q = msg->sending_process_nr;
276        q++;
277
278        switch (msg->msg_type)
279        {
280        case COMM1:
281            *q = 0; // padding
282            break;
283
284        case COMM2:
285            *q = msg->time_v_j;
286            break;
287
288        case COMM3:
289            *q = msg->committed_cnt;
290            q++;
291            for (int i = 0; i < (int)time_v.size(); i++)
292            {
293                *q = msg->time_v_min[i];
294                q++;
295            }
296            for (uset::iterator ptr = msg->committed_msgs.begin();
297                 ptr != msg->committed_msgs.end(); ptr++)
298            {
299                for (int i = 0; i < (int)time_v.size() * 2; i++)
300                {
301                    *q = ptr->at(i);
302                    q++;
303                }
304            }
305            break;
306        case COMM4:
```

```
307
308            *q = msg->committed_cnt;
309            q++;
310            for (uset::iterator ptr = msg->committed_msgs.begin();
311                 ptr != msg->committed_msgs.end(); ptr++)
312            {
313                for (int i = 0; i < (int)time_v.size() * 2; i++)
314                {
315                    *q = ptr->at(i);
316                    q++;
317                }
318            }
319
320            break;
321
322        default:
323            break;
324        }
325    }
326
327    void Pet_kea::State::deserialize_commit(char *data, struct comm_msg_t *msg)
328    {
329        int *q = (int *)data;
330        msg->msg_type = (msg_type)*q;
331        q++;
332
333        msg->sending_process_nr = *q;
334        q++;
335        vector<int> temp_vec;
336        switch (msg->msg_type)
337        {
338        case COMM1:
339            break;
340
341        case COMM2:
342            msg->time_v_j = *q;
343            break;
344
345        case COMM3:
346
347            msg->committed_cnt = *q;
348            q++;
349            for (int i = 0; i < (int)time_v.size(); i++)
350            {
351                msg->time_v_min.push_back(*q);
352                q++;
353            }
354
355            for (int i = 0; i < msg->committed_cnt; i++)
356            {
357                for (int j = 0; j < (int)time_v.size() * 2; j++)
358                {
359                    temp_vec.push_back(*q);
360                    q++;
361                }
362
363                msg->committed_msgs.insert(temp_vec);
364                temp_vec.clear();
365            }
366            break;
367        case COMM4:
368
369            msg->committed_cnt = *q;
370            q++;
371            for (int i = 0; i < msg->committed_cnt; i++)
372            {
373                for (int j = 0; j < (int)time_v.size() * 2; j++)
374                {
375                    temp_vec.push_back(*q);
376                    q++;
377                }
378
379                msg->committed_msgs.insert(temp_vec);
380                temp_vec.clear();
381            }
382
383            break;
384
385        default:
386            break;
387        }
388    }
389
390    void Pet_kea::State::serialize_ctrl(struct ctrl_msg_t *msg, char *data)
391    {
392        int *q = (int *)data;
393        *q = (int)msg->msg_type;
394        q++;
395
396        *q = msg->recieved_cnt;
397        q++;
398
399        *q = msg->sending_process_nr;
400        q++;
401
402        *q = msg->log_entry.id;
403        q++;
404        *q = msg->log_entry.fail_nr;
405        q++;
406        *q = msg->log_entry.res_time;
407        q++;
408
409        for (set<pair<int, vector<int>>>::iterator ptr = msg->recieved_msgs.begin();
```

```
410         ptr != msg->recieved_msgs.end(); ptr++)
411     {
412         *q = ptr->first;
413         q++;
414         for (int j = 0; j < (int)fail_v.size(); j++)
415         {
416             *q = ptr->second[j];
417             q++;
418         }
419     }
420 }
421
422 void Pet_kea::State::deserialize_ctrl(char *data, struct ctrl_msg_t *msg)
423 {
424     int *q = (int *)data;
425     msg->msg_type = (msg_type)*q;
426     q++;
427
428     msg->recieved_cnt = *q;
429     q++;
430
431     msg->sending_process_nr = *q;
432     q++;
433
434     msg->log_entry.id = *q;
435     q++;
436     msg->log_entry.fail_nr = *q;
437     q++;
438     msg->log_entry.res_time = *q;
439     q++;
440
441     pair<int, vector<int>> temp_pair;
442     for (int i = 0; i < msg->recieved_cnt; i++)
443     {
444
445         temp_pair.first = *q;
446         q++;
447         for (int j = 0; j < (int)fail_v.size(); j++)
448         {
449             temp_pair.second.push_back(*q);
450             q++;
451         }
452
453         msg->recieved_msgs.insert(temp_pair);
454         temp_pair.second.clear();
455     }
456 }
457
458 void Pet_kea::State::serialize(struct msg_t *msg, char *data)
459 {
460     int *q = (int *)data;
461     *q = msg->msg_type;
462     q++;
463
464     *q = msg->msg_size;
465     q++;
466
467     *q = msg->sending_process_nr;
468     q++;
469     if (msg->msg_type == MSG)
470     {
471         for (int i = 0; i < (int)time_v.size(); i++)
472         {
473             *q = msg->time_v[i];
474             q++;
475         }
476
477         for (int i = 0; i < (int)fail_v.size(); i++)
478         {
479             *q = msg->fail_v[i];
480             q++;
481         }
482     }
483
484     memcpy(q, msg->msg_buf, msg->msg_size);
485 }
486 void Pet_kea::State::deserialize(char *data, struct msg_t *msg)
487 {
488     int *q = (int *)data;
489     msg->msg_type = (msg_type)*q;
490     q++;
491
492     msg->msg_size = *q;
493     q++;
494
495     msg->sending_process_nr = *q;
496     q++;
497     if (msg->msg_type == MSG)
498     {
499         for (int i = 0; i < (int)time_v.size(); i++)
500         {
501             msg->time_v.push_back(*q);
502             q++;
503         }
504
505         for (int i = 0; i < (int)fail_v.size(); i++)
506         {
507             msg->fail_v.push_back(*q);
508             q++;
509         }
510     }
511
512     msg->msg_buf = (char *)malloc(msg->msg_size);
```

```cpp
513        memcpy(msg->msg_buf, q, msg->msg_size);
514 }
515
516 void Pet_kea::State::serialize_log(struct msg_log_t *log, char *data)
517 {
518        int *q = (int *)data;
519        *q = log->msg_size;
520        q++;
521        *q = log->recipient;
522        q++;
523        *q = log->process_id;
524        q++;
525        *q = log->next_checkpoint;
526        q++;
527
528        for (int i = 0; i < (int)time_v.size(); i++)
529        {
530            *q = log->time_v_sender[i];
531            q++;
532        }
533
534        for (int i = 0; i < (int)time_v.size(); i++)
535        {
536            *q = log->time_v_reciever[i];
537            q++;
538        }
539
540        for (int i = 0; i < (int)fail_v.size(); i++)
541        {
542            *q = log->fail_v_sender[i];
543            q++;
544        }
545
546        memcpy(q, log->msg_buf, log->msg_size);
547 }
548 int Pet_kea::State::deserialize_log(char *data, struct msg_log_t *log)
549 {
550        int *q = (int *)data;
551        log->msg_size = *q;
552        q++;
553        log->recipient = *q;
554        q++;
555        log->process_id = *q;
556        q++;
557        log->next_checkpoint = *q;
558        q++;
559
560        for (int i = 0; i < (int)time_v.size(); i++)
561        {
562            log->time_v_sender.push_back(*q);
563            q++;
564        }
565
566        for (int i = 0; i < (int)time_v.size(); i++)
567        {
568            log->time_v_reciever.push_back(*q);
569            q++;
570        }
571
572        for (int i = 0; i < (int)fail_v.size(); i++)
573        {
574            log->fail_v_sender.push_back(*q);
575            q++;
576        }
577
578        log->msg_buf = (char *)malloc(log->msg_size);
579        memcpy(log->msg_buf, q, log->msg_size);
580
581        return (q - (int *)data) * sizeof(int) + log->msg_size;
582 }
583
584 void Pet_kea::State::serialize_state(char *data)
585 {
586        int *q = (int *)data;
587        *q = id;
588        q++;
589
590        *q = msg_cnt;
591        q++;
592
593        *q = arrived_ctrl.size();
594        q++;
595
596        *q = committed_msg_set.size();
597        q++;
598
599        *q = committed_recieve_events.size();
600        q++;
601
602        for (int i = 0; i < (int)time_v.size(); i++)
603        {
604            *q = time_v[i];
605            q++;
606        }
607        for (int i = 0; i < (int)time_v.size(); i++)
608        {
609            *q = time_v_min[i];
610            q++;
611        }
612        for (int i = 0; i < (int)time_v.size(); i++)
613        {
614            *q = commit_v[i];
615            q++;
```

```
616         }
617
618         for (int i = 0; i < (int)time_v.size(); i++)
619         {
620             *q = remove_v[i];
621             q++;
622         }
623
624         for (uset::iterator ptr = arrived_ctrl.begin(); ptr != arrived_ctrl.end(); ptr++)
625         {
626             *q = ptr->at(0);
627             q++;
628             *q = ptr->at(1);
629             q++;
630             *q = ptr->at(2);
631             q++;
632         }
633         for (uset::iterator ptr = committed_msg_set.begin(); ptr != committed_msg_set.end(); ptr++)
634         {
635             for (int i = 0; i < (int)time_v.size() * 2; i++)
636             {
637                 *q = ptr->at(i);
638                 q++;
639             }
640         }
641         for (uset::iterator ptr = committed_recieve_events.begin(); ptr != committed_recieve_events.end(); ptr++)
642         {
643             for (int i = 0; i < (int)time_v.size() * 2; i++)
644             {
645                 *q = ptr->at(i);
646                 q++;
647             }
648         }
649 }
650
651 void Pet_kea::State::deserialize_state(char *data)
652 {
653         int *q = (int *)data;
654         id = *q;
655         q++;
656
657         msg_cnt = *q;
658         q++;
659
660         int arrived_ctrl_size = *q;
661         q++;
662
663         int committed_msg_set_size = *q;
664         q++;
665
666         int committed_recieve_events_size = *q;
667         q++;
668
669         for (int i = 0; i < (int)time_v.size(); i++)
670         {
671             time_v[i] = *q;
672             q++;
673         }
674
675         for (int i = 0; i < (int)time_v.size(); i++)
676         {
677             time_v_min[i] = *q;
678             q++;
679         }
680
681         for (int i = 0; i < (int)time_v.size(); i++)
682         {
683             commit_v[i] = *q;
684             q++;
685         }
686
687         for (int i = 0; i < (int)time_v.size(); i++)
688         {
689             remove_v[i] = *q;
690             q++;
691         }
692
693         vector<int> temp_vec;
694         for (int i = 0; i < arrived_ctrl_size; i++)
695         {
696             temp_vec.push_back(*q);
697             q++;
698             temp_vec.push_back(*q);
699             q++;
700             temp_vec.push_back(*q);
701             q++;
702             arrived_ctrl.insert(temp_vec);
703             temp_vec.clear();
704         }
705
706         for (int i = 0; i < committed_msg_set_size; i++)
707         {
708             for (int j = 0; j < (int)time_v.size() * 2; j++)
709             {
710                 temp_vec.push_back(*q);
711                 q++;
712             }
713             committed_msg_set.insert(temp_vec);
714             temp_vec.clear();
715         }
716
717         for (int i = 0; i < committed_recieve_events_size; i++)
718         {
```

```cpp
            for (int j = 0; j < (int)time_v.size() * 2; j++)
            {
                temp_vec.push_back(*q);
                q++;
            }
            committed_recieve_events.insert(temp_vec);
            temp_vec.clear();
        }
}

int Pet_kea::State::send_ctrl()
{
    set<pair<int, std::vector<int>>> recvd_msgs[time_v.size()];

    vector<int> cnt(time_v.size(), 0);
    struct fail_log_t fail_log = {id, fail_v[id], time_v[id]};

    pair<int, vector<int>> temp_pair;
    for (int i = 0; i < msg_cnt; i++)
    {
        if (!msg_log[i].recipient)
        {
            continue;
        }

        temp_pair.first = msg_log[i].time_v_sender[msg_log[i].process_id];
        temp_pair.second = msg_log[i].fail_v_sender;
        recvd_msgs[msg_log[i].process_id].insert(temp_pair);
        cnt[msg_log[i].process_id]++;
    }

    for (int i = 0; i < (int)time_v.size(); i++)
    {
        // prepare and send the ctrl messages
        if (i == id)
            continue;
        struct ctrl_msg_t msg;
        msg.msg_type = CTRL;
        msg.sending_process_nr = id;
        msg.log_entry = fail_log;
        msg.recieved_cnt = cnt[i];
        msg.recieved_msgs = recvd_msgs[i];

        print_ctrl_msg(&msg);

        // send the control message (serialize)
        size_t size = SER_SIZE_CTRL_MSG_T(msg.recieved_cnt, fail_v.size());
        char *data = (char *)malloc(size);
        serialize_ctrl(&msg, data);
        try
        {
            int ret = write(fildes[i][1], data, size);
            if (ret < 0)
            {
                throw runtime_error("failed_to_write");
            }
        }
        catch (const std::exception &e)
        {
            std::cerr << e.what() << endl;
            perror("write_failed");
            return -1;
        }

        free(data);
    }
    return 0;
}

int Pet_kea::State::store_msg(struct msg_t *msg, int recipient)
{
    // store the message
    if (msg_cnt >= MAX_LOG)
    {
        cout << "max_log_size_reached_of_process_" << id << endl;
        return -1;
    }

    msg_log[msg_cnt].msg_size = msg->msg_size;
    msg_log[msg_cnt].next_checkpoint = checkpoints.size();
    msg_log[msg_cnt].msg_buf = (char *)malloc(msg->msg_size);
    memcpy(msg_log[msg_cnt].msg_buf, msg->msg_buf, msg->msg_size);
    msg_log[msg_cnt].time_v_reciever = time_v;
    if (recipient == -1)
    {
        msg_log[msg_cnt].time_v_sender = msg->time_v;
        msg_log[msg_cnt].fail_v_sender = msg->fail_v;
        msg_log[msg_cnt].process_id = msg->sending_process_nr;
        msg_log[msg_cnt].recipient = true;
    }
    else
    {
        msg_log[msg_cnt].time_v_sender = time_v;
        msg_log[msg_cnt].fail_v_sender = fail_v;
        msg_log[msg_cnt].process_id = recipient;
        msg_log[msg_cnt].recipient = false;
    }

    msg_cnt++;
    if (SAVE_CNT != 0 && msg_cnt % SAVE_CNT == 0)
        checkpoint();

    return 0;
```

```cpp
822  }
823
824  int Pet_kea::State::commit_msgs(vector<int> msgs)
825  {
826      char filename[32];
827      get_comm_filename(id, filename);
828
829      for (int iterator : msgs)
830      {
831          char data[msg_log[iterator].msg_size + SER_LOG_SIZE];
832          serialize_log(&msg_log[iterator], data);
833          commit_out.write(data, msg_log[iterator].msg_size + SER_LOG_SIZE);
834      }
835      commit_out.flush();
836
837      return 0;
838  }
839
840  int Pet_kea::State::rollback(struct ctrl_msg_t *msg)
841  {
842      cout << "entered the rollback section" << endl;
843      print_ctrl_msg(msg);
844
845      // RB.2
846      char filename[32];
847      get_fail_filename(id, filename);
848      ofstream fail_out(filename, ofstream::out | ofstream::binary | ofstream::app);
849
850      struct fail_log_t entry = {msg->sending_process_nr, msg->log_entry.fail_nr, msg->log_entry.res_time};
851      fail_out.write((char *)&entry, sizeof(fail_log_t));
852      fail_out.close();
853
854      fail_log.push_back(entry);
855
856      //  RB.2.3
857      fail_v[msg->sending_process_nr] = msg->log_entry.fail_nr;
858      if (time_v[msg->sending_process_nr] > msg->log_entry.res_time)
859      {
860          //  RB.2.1      remove ckeckpoints T^i > crT^i, set state and T to latest ckeckpoint, replays
861          vector<int> checkpoints_to_remove;
862
863          for (int i = 0; i < (int)checkpoints.size(); i++)
864          {
865              if (ck_time_v[i].at(msg->sending_process_nr) > msg->log_entry.res_time)
866                  checkpoints_to_remove.push_back(i);
867          }
868
869          rem_checkpoints(checkpoints_to_remove);
870
871          // set state and time_v to latestck
872          int prev_cnt = msg_cnt;
873          msg_cnt = checkpoints.back();
874          int temp_msg_cnt = msg_cnt;
875          time_v = ck_time_v.back();
876
877          // replay messages
878          std::vector<int> indices_to_rem;
879
880          for (int i = msg_cnt; i < prev_cnt; i++)
881          {
882              if (msg_log[i].time_v_sender[msg->sending_process_nr] <= msg->log_entry.res_time &&
883                  msg_log[i].time_v_sender[id] >= ck_time_v.back()[id])
884              {
885                  // replay messages only inc time_v and msg_cnt
886                  cout << "replayed msg" << endl;
887                  msg_cnt++;
888                  if (msg_log[i].recipient)
889                  {
890                      time_v[id]++;
891                      for (int j = 0; j < (int)time_v.size(); j++)
892                      {
893                          if (j == id)
894                              continue;
895                          time_v.at(j) = max(msg_log[i].time_v_sender[j], time_v[j]);
896                      }
897                  }
898                  else
899                  {
900                      time_v[id]++;
901                  }
902              }
903              else if (!msg_log[i].recipient &&
904                       msg_log[i].time_v_sender[msg->sending_process_nr] > msg->log_entry.res_time)
905              {
906                  // add to remove vector
907                  indices_to_rem.push_back(i);
908              }
909              else if (msg_log[i].recipient &&
910                       msg_log[i].time_v_reciever[msg->sending_process_nr] > msg->log_entry.res_time &&
911                       msg_log[i].time_v_sender[msg->sending_process_nr] > msg->log_entry.res_time &&
912                       msg_log[i].fail_v_sender[msg->sending_process_nr] < msg->log_entry.fail_nr)
913              {
914                  // add to remove vector
915                  indices_to_rem.push_back(i);
916              }
917          }
918          if (!indices_to_rem.empty())
919          {
920              msg_cnt = rem_log_entries(indices_to_rem, prev_cnt);
921              indices_to_rem.clear();
922          }
923
924          //  RB.2.2
```

```cpp
925
926            // RB.3
927            prev_cnt = msg_cnt;
928            for (int i = temp_msg_cnt; i < prev_cnt; i++)
929            {
930                // move recv event to the back
931                if (msg_log[i].recipient &&
932                    msg_log[i].time_v_reciever[msg->sending_process_nr] > msg->log_entry.res_time)
933                {
934                    cout << id << "_moved_RECV_event_to_the_back" << endl;
935                    if (msg_cnt >= MAX_LOG)
936                    {
937                        cout << "max_log_size_reached_of_process_" << id << endl;
938                        return -1;
939                    }
940
941                    msg_log[msg_cnt].msg_size = msg_log[i].msg_size;
942                    msg_log[msg_cnt].recipient = true;
943                    msg_log[msg_cnt].process_id = msg_log[i].process_id;
944                    msg_log[msg_cnt].msg_buf = (char *)malloc(msg_log[i].msg_size);
945                    memcpy(msg_log[msg_cnt].msg_buf, msg_log[i].msg_buf, msg_log[i].msg_size);
946
947                    msg_log[msg_cnt].time_v_sender = msg_log[i].time_v_sender;
948                    msg_log[msg_cnt].fail_v_sender = msg_log[i].fail_v_sender;
949                    msg_log[msg_cnt].time_v_reciever = msg_log[i].time_v_reciever;
950
951                    msg_cnt++;
952
953                    time_v[id]++;
954                    for (int j = 0; j < (int)time_v.size(); j++)
955                    {
956                        if (j == id)
957                            continue;
958                        time_v.at(j) = max(msg_log[i].time_v_sender[i], time_v[i]);
959                    }
960
961                    // remove the duplicate msg from the log
962                    indices_to_rem.push_back(i);
963                }
964
965                // retransmit send events that have not arrived RB.3.3
966                if (!msg_log[i].recipient && msg_log[i].process_id == msg->sending_process_nr &&
967                    !(msg->recieved_msgs.contains(pair<int, vector<int>>(msg_log[i].time_v_sender[id],
968                                                                          msg_log[i].fail_v_sender))))
969                {
970                    cout << id << "_retransmitted_msg_Tj:_" << msg_log[i].time_v_sender[id] << "fail_v:_";
971                    for (int j = 0; j < (int)fail_v.size(); j++)
972                    {
973                        cout << msg_log[i].fail_v_sender[j] << ":";
974                    }
975
976                    cout << "_res_time:_" << msg_log[i].time_v_sender[msg->sending_process_nr] << endl;
977                    struct msg_t retransmit_msg;
978                    retransmit_msg.msg_type = MSG;
979                    retransmit_msg.sending_process_nr = id;
980                    retransmit_msg.time_v = msg_log[i].time_v_sender;
981                    retransmit_msg.fail_v = msg_log[i].fail_v_sender;
982                    retransmit_msg.msg_size = msg_log[i].msg_size;
983                    retransmit_msg.msg_buf = (char *)malloc(retransmit_msg.msg_size * sizeof(char));
984                    memcpy(retransmit_msg.msg_buf, msg_log[i].msg_buf, retransmit_msg.msg_size);
985
986                    char data[SER_MSG_SIZE + msg_log[i].msg_size];
987                    serialize(&retransmit_msg, data);
988
989                    // send the message
990                    try
991                    {
992                        if (write(fildes[msg_log[i].process_id][1], data, SER_MSG_SIZE + msg_log[i].msg_size) <
993                            0)
994                            throw runtime_error("failed_to_write");
995                    }
996                    catch (const std::exception &e)
997                    {
998                        std::cerr << e.what() << endl;
999                        perror("write_failed");
1000                        return -1;
1001                    }
1002                }
1003            }
1004            if (!indices_to_rem.empty())
1005            {
1006
1007                msg_cnt = rem_log_entries(indices_to_rem, msg_cnt);
1008            }
1009
1010            checkpoint();
1011        }
1012        return 0;
1013    }
1014
1015    Pet_kea::State::State(int process_nr, int process_cnt, int (*fd)[2], bool restart)
1016        : id(process_nr),
1017          time_v(process_cnt, 0),
1018          fail_v(process_cnt, 0),
1019          msg_cnt(0),
1020          commit_v(process_cnt, false),
1021          remove_v(process_cnt, false),
1022          arrived_msgs()
1023
1024    {
1025        time_v_min = vector(process_cnt, 0);
1026        fildes = (int **)malloc(process_cnt * sizeof(int *));
1027        for (int i = 0; i < process_cnt; i++)
```

```cpp
1027         {
1028             fildes[i] = (int *)malloc(2 * sizeof(int));
1029             fildes[i][0] = fd[i][0];
1030             fildes[i][1] = fd[i][1];
1031         }
1032         if (restart)
1033         {
1034             char filename[32];
1035             get_state_filename(id, filename);
1036             ifstream state_in(filename, ifstream::in | ifstream::binary);
1037             state_in.seekg(0, ifstream::end);
1038             size_t file_size = state_in.tellg();
1039             state_in.seekg(0, ifstream::beg);
1040             char state_file[file_size];
1041             state_in.read(state_file, file_size);
1042             state_in.close();
1043             deserialize_state(state_file);
1044
1045             get_msg_filename(id, filename);
1046             ifstream msg_in(filename, ifstream::in | ifstream::binary);
1047             msg_in.seekg(0, msg_in.end);
1048             file_size = msg_in.tellg();
1049             msg_in.seekg(0, ifstream::beg);
1050             char msg_file[file_size];
1051             msg_in.read(msg_file, file_size);
1052             msg_in.close();
1053
1054             int *curr_pos = (int *)msg_file;
1055
1056             int num_checkpoints = *curr_pos;
1057             curr_pos++;
1058
1059             msg_log = (msg_log_t *)calloc(MAX_LOG, sizeof(msg_log_t));
1060
1061             checkpoints.push_back(0);
1062             ck_time_v.push_back(vector<int>(process_cnt, 0));
1063
1064             int read_msg_cnt = 0;
1065             for (int i = 0; i < num_checkpoints; i++)
1066             {
1067                 int ck_msg_cnt = *curr_pos;
1068                 checkpoints.push_back(*curr_pos);
1069                 curr_pos++;
1070                 int to_read = ck_msg_cnt - *curr_pos;
1071                 curr_pos++;
1072                 std::vector<int> temp_ck_time_v;
1073
1074                 for (int j = 0; j < (int)time_v.size(); j++)
1075                 {
1076                     temp_ck_time_v.push_back(*curr_pos);
1077                     curr_pos++;
1078                 }
1079                 ck_time_v.push_back(temp_ck_time_v);
1080
1081                 for (int k = 0, ret = 0; k < to_read; k++, read_msg_cnt++)
1082                 {
1083                     ret = deserialize_log((char *)curr_pos, &msg_log[read_msg_cnt]);
1084                     curr_pos += ret / sizeof(int);
1085                 }
1086             }
1087
1088             // detect lost messages if crash happened during checkpoint
1089             if (msg_cnt != read_msg_cnt)
1090             {
1091                 msg_cnt = read_msg_cnt;
1092                 if (msg_cnt == 0)
1093                 {
1094                     time_v = vector<int>(process_cnt, 0);
1095                 }
1096                 else
1097                 {
1098                     time_v = msg_log[msg_cnt - 1].time_v_sender;
1099                 }
1100             }
1101
1102             // insert the arrived messages into arrived_msgs
1103             vector<int> merged_time_fail_v;
1104
1105             for (int i = 0; i < msg_cnt; i++)
1106             {
1107                 if (msg_log->recipient)
1108                 {
1109                     merged_time_fail_v = msg_log[i].time_v_sender;
1110                     merged_time_fail_v.insert(merged_time_fail_v.end(),
1111                                               msg_log[i].fail_v_sender.begin(), msg_log[i].fail_v_sender.end());
1112                     arrived_msgs.insert(merged_time_fail_v);
1113                 }
1114             }
1115
1116             msg_out.open(filename, ofstream::out | ofstream::binary | ofstream::ate | ofstream::in);
1117             get_comm_filename(id, filename);
1118             commit_out.open(filename, ofstream::out | ofstream::binary | ofstream::app);
1119             get_fail_filename(id, filename);
1120
1121             ifstream fail_in(filename, ifstream::in | ifstream::binary);
1122
1123             fail_in.seekg(0, fail_in.end);
1124             file_size = fail_in.tellg();
1125             fail_in.seekg(0, ifstream::beg);
1126             char fail_file[file_size];
1127             fail_in.read(fail_file, file_size);
1128             fail_in.close();
1129             curr_pos = (int *)fail_file;
```

```
1130
1131          fail_log_t temp_fail_log;
1132          while (curr_pos < (int *)(fail_file + file_size))
1133          {
1134              temp_fail_log.id = *curr_pos;
1135              curr_pos++;
1136              temp_fail_log.fail_nr = *curr_pos;
1137              curr_pos++;
1138              temp_fail_log.res_time = *curr_pos;
1139              curr_pos++;
1140              fail_log.push_back(temp_fail_log);
1141              fail_v[temp_fail_log.id] = max(fail_v[temp_fail_log.id], temp_fail_log.fail_nr);
1142          }
1143
1144          ofstream fail_out(filename, ofstream::out | ofstream::binary | ofstream::app);
1145          fail_out.seekp(0, ofstream::end);
1146          fail_v[id]++;
1147          struct fail_log_t entry = {id, fail_v[id], time_v[id]};
1148          fail_out.write((char *)&entry, sizeof(fail_log_t));
1149          fail_out.close();
1150
1151          fail_log.push_back(entry);
1152
1153          send_ctrl();
1154      }
1155      else
1156      {
1157          char filename[32];
1158          get_fail_filename(id, filename);
1159          ofstream fail_out(filename, ofstream::out | ofstream::binary | ofstream::trunc);
1160          get_msg_filename(id, filename);
1161          msg_out.open(filename, ofstream::out | ofstream::binary | ofstream::trunc);
1162          get_comm_filename(id, filename);
1163          commit_out.open(filename, ofstream::out | ofstream::binary | ofstream::trunc);
1164
1165          struct fail_log_t entry = {id, 0, 0};
1166          fail_out.write((char *)&entry, sizeof(fail_log_t));
1167          fail_out.close();
1168
1169          fail_log.push_back(fail_log_t(id, 0, 0));
1170
1171          msg_log = (msg_log_t *)calloc(MAX_LOG, sizeof(msg_log_t));
1172          checkpoints.push_back(0);
1173          ck_time_v.push_back(vector<int>(process_cnt, 0));
1174      }
1175      SAVE_CNT = 0;
1176  }
1177
1178  Pet_kea::State::~State()
1179  {
1180      for (int i = msg_cnt - 1; i >= 0; i--)
1181      {
1182          std::vector<int>().swap(msg_log[i].time_v_reciever);
1183          std::vector<int>().swap(msg_log[i].time_v_sender);
1184          std::vector<int>().swap(msg_log[i].fail_v_sender);
1185
1186          free(msg_log[i].msg_buf);
1187      }
1188
1189      free(msg_log);
1190      for (int i = 0; i < (int)time_v.size(); i++)
1191      {
1192          free(fildes[i]);
1193      }
1194
1195      free(fildes);
1196      msg_out.close();
1197      commit_out.close();
1198  }
1199
1200  int Pet_kea::State::checkpoint()
1201  {
1202      // write state and time vector at the start of the file
1203
1204      int *update = (int *)malloc(sizeof(int) * 2);
1205      *update = msg_cnt;
1206      update++;
1207      *update = checkpoints.back();
1208      update--;
1209
1210      msg_out.seekp(0, ofstream::beg);
1211      int *num_checkpoints = (int *)malloc(sizeof(int));
1212      *num_checkpoints = checkpoints.size();
1213      msg_out.write((char *)num_checkpoints, sizeof(int));
1214      free(num_checkpoints);
1215
1216      int *time_v_buffer = (int *)malloc(sizeof(int) * time_v.size());
1217      for (int i = 0; i < (int)time_v.size(); i++)
1218      {
1219          time_v_buffer[i] = time_v[i];
1220      }
1221
1222      // append last messages
1223      msg_out.seekp(0, ofstream::end);
1224
1225      msg_out.write((char *)update, sizeof(int) * 2);
1226      free(update);
1227
1228      msg_out.write((char *)time_v_buffer, sizeof(int) * time_v.size());
1229
1230      free(time_v_buffer);
1231
1232      for (int i = checkpoints.back(); i < msg_cnt; i++)
```

```cpp
1233            {
1234                    char data[msg_log[i].msg_size + SER_LOG_SIZE];
1235                    serialize_log(&msg_log[i], data);
1236                    msg_out.write(data, msg_log[i].msg_size + SER_LOG_SIZE);
1237            }
1238            checkpoints.push_back(msg_cnt);
1239            ck_time_v.push_back(time_v);
1240            msg_out.flush();
1241
1242            char filename[32];
1243            get_state_filename(id, filename);
1244
1245            ofstream state_out(filename, ofstream::out | ofstream::binary | ofstream::trunc);
1246            int state_size = SER_STATE_SIZE(arrived_ctrl.size(),
1247                                            committed_msg_set.size(), committed_recieve_events.size());
1248
1249            char data[state_size];
1250            serialize_state(data);
1251
1252            state_out.write(data, state_size);
1253            state_out.close();
1254            return 0;
1255    }
1256
1257    int Pet_kea::send_void(char *input, int fildes[2], int size)
1258    {
1259            char data[SER_VOID_SIZE + size];
1260
1261            int *q = (int *)data;
1262            *q = VOID;
1263            q++;
1264            *q = size;
1265            q++;
1266            memcpy(q, input, size);
1267
1268            try
1269            {
1270                    if (write(fildes[1], data, SER_VOID_SIZE + size) < 0)
1271                            throw runtime_error("failed_to_write");
1272            }
1273            catch (const std::exception &e)
1274            {
1275                    std::cerr << e.what() << endl;
1276                    perror("write_failed");
1277                    return -1;
1278            }
1279            return 0;
1280    }
1281
1282    int Pet_kea::State::signal_commit()
1283    {
1284            struct comm_msg_t msg;
1285            msg.msg_type = COMM1;
1286            msg.sending_process_nr = id;
1287            char data[SER_COMM1_SIZE];
1288            serialize_commit(&msg, data);
1289
1290            // write to all other processes
1291            try
1292            {
1293                    for (int i = 0; i < (int)time_v.size(); i++)
1294                    {
1295                            if (i == id)
1296                                    continue;
1297                            if (write(fildes[i][1], data, SER_COMM1_SIZE) < 0)
1298                                    throw runtime_error("failed_to_write");
1299                    }
1300            }
1301            catch (const std::exception &e)
1302            {
1303                    std::cerr << e.what() << endl;
1304                    perror("write_failed");
1305                    return -1;
1306            }
1307            commit_v[id] = true;
1308            return 0;
1309    }
1310
1311    int Pet_kea::State::send_commit(int target_id)
1312    {
1313            struct comm_msg_t msg;
1314            msg.msg_type = COMM2;
1315            msg.sending_process_nr = id;
1316            msg.time_v_j = ck_time_v.back().at(id);
1317
1318            char data[SER_COMM2_SIZE];
1319            serialize_commit(&msg, data);
1320            // write back to sender
1321            try
1322            {
1323                    if (write(fildes[target_id][1], data, SER_COMM2_SIZE) < 0)
1324                            throw runtime_error("failed_to_write");
1325            }
1326            catch (const std::exception &e)
1327            {
1328                    std::cerr << e.what() << endl;
1329                    perror("write_failed");
1330                    return -1;
1331            }
1332            return 0;
1333    }
1334
1335    int Pet_kea::State::remove_data()
```

```
1336  {
1337      vector<int> indices_to_remove, checkpoints_to_remove;
1338      vector<int> temp_vec;
1339      unordered_set<std::vector<int>, vector_hash> next_committed_recieve_events;
1340
1341      // remove checkpoint
1342      for (int i = 1; i < (int)checkpoints.size() - 1; i++)
1343      {
1344          if (ck_time_v[i + 1] <= time_v_min)
1345          {
1346              checkpoints_to_remove.push_back(i);
1347          }
1348      }
1349
1350      // remove messages
1351      for (int i = 0; i < msg_cnt; i++)
1352      {
1353          temp_vec.clear();
1354          temp_vec = msg_log[i].time_v_sender;
1355          temp_vec.insert(temp_vec.end(), msg_log[i].fail_v_sender.begin(), msg_log[i].fail_v_sender.end());
1356          if (msg_log[i].next_checkpoint >= (int)checkpoints.size())
1357          {
1358              msg_log[i].next_checkpoint = next_checkpoint_after_rem(checkpoints_to_remove,
1359                                                                     msg_log[i].next_checkpoint);
1360
1361              if (!msg_log[i].recipient && committed_recieve_events.contains(temp_vec))
1362              {
1363                  next_committed_recieve_events.insert(temp_vec);
1364              }
1365              continue;
1366          }
1367          if (msg_log[i].recipient && ck_time_v.at(msg_log[i].next_checkpoint) <= time_v_min)
1368          {
1369              // remove form msg_log
1370              indices_to_remove.push_back(i);
1371              committed_msg_set.erase(temp_vec);
1372              if (msg_log[i].fail_v_sender[id] == fail_v[id])
1373              {
1374                  // remove the entry from the set
1375                  vector<int> merged_time_fail_v;
1376                  merged_time_fail_v = msg_log[i].time_v_sender;
1377                  merged_time_fail_v.insert(merged_time_fail_v.end(),
1378                                            msg_log[i].fail_v_sender.begin(), msg_log[i].fail_v_sender.end());
1379                  arrived_msgs.erase(merged_time_fail_v);
1380              }
1381              continue;
1382          }
1383
1384          if (!msg_log[i].recipient && committed_recieve_events.contains(temp_vec))
1385          {
1386              if (ck_time_v.at(msg_log[i].next_checkpoint) <= time_v_min)
1387              {
1388
1389                  // remove from msg_log
1390                  indices_to_remove.push_back(i);
1391                  committed_msg_set.erase(temp_vec);
1392                  continue;
1393              }
1394              else
1395              {
1396                  next_committed_recieve_events.insert(temp_vec);
1397              }
1398          }
1399          msg_log[i].next_checkpoint = next_checkpoint_after_rem(checkpoints_to_remove,
1400                                                                 msg_log[i].next_checkpoint);
1401      }
1402      msg_cnt = rem_log_entries(indices_to_remove, msg_cnt);
1403
1404      rem_checkpoints(checkpoints_to_remove);
1405      committed_recieve_events.clear();
1406      committed_recieve_events = next_committed_recieve_events;
1407
1408      // remove fail_log not possible in asynchronos setting
1409      return 0;
1410  }
1411
1412  int Pet_kea::State::commit(bool is_instigator)
1413  {
1414      // commit the messages
1415      vector<int> committed_msgs;
1416      for (int i = 0; i < msg_cnt; i++)
1417      {
1418          if ((msg_log[i].recipient ? msg_log[i].time_v_reciever : msg_log[i].time_v_sender) <= time_v_min &&
1419              !check_duplicate_commit(&msg_log[i]))
1420          {
1421              committed_msgs.push_back(i);
1422          }
1423      }
1424      commit_msgs(vector<int>(committed_msgs));
1425
1426      // remove information
1427
1428      unordered_set<vector<int>, vector_hash> committed_set[time_v.size()];
1429      struct comm_msg_t msg;
1430      if (is_instigator)
1431      {
1432          msg.msg_type = COMM3;
1433          msg.time_v_min = time_v_min;
1434      }
1435      else
1436      {
1437          msg.msg_type = COMM4;
1438      }
```

```cpp
1439        msg.sending_process_nr = id;
1440        vector<int> merged_vector;
1441
1442        int it;
1443        while (!committed_msgs.empty())
1444        {
1445            it = committed_msgs.back();
1446            if (msg_log[it].recipient)
1447            {
1448                merged_vector = msg_log[it].time_v_sender;
1449                merged_vector.insert(merged_vector.end(),
1450                                     msg_log[it].fail_v_sender.begin(), msg_log[it].fail_v_sender.end());
1451                committed_set[msg_log[it].process_id].insert(merged_vector);
1452                merged_vector.clear();
1453            }
1454
1455            committed_msgs.pop_back();
1456        }
1457
1458        remove_v[id] = true;
1459        int max_cnt = 0;
1460        for (int i = 0; i < (int)time_v.size(); i++)
1461        {
1462            if (max_cnt < (int)committed_set[i].size())
1463                max_cnt = committed_set[i].size();
1464        }
1465        char *data = (char *)malloc(msg.msg_type == COMM3 ? SER_COMM3_SIZE(max_cnt) : SER_COMM4_SIZE(max_cnt));
1466        // write to all other processes
1467        for (int i = 0; i < (int)time_v.size(); i++)
1468        {
1469            if (i == id)
1470                continue;
1471
1472            msg.committed_msgs = committed_set[i];
1473            msg.committed_cnt = committed_set[i].size();
1474
1475            if (is_instigator)
1476            {
1477
1478                cout << id << "_sending_comm3_to_" << i << "commit_cnt:_" << msg.committed_cnt << "time_v_min:";
1479                for (int j = 0; j < (int)time_v.size(); j++)
1480                {
1481                    cout << msg.time_v_min[j] << ":";
1482                }
1483                cout << endl;
1484            }
1485            else
1486            {
1487                cout << id << "_sending_comm4_to_" << i << "commit_cnt:_" << msg.committed_cnt << endl;
1488            }
1489
1490            serialize_commit(&msg, data);
1491
1492            if (write(fildes[i][1], data,
1493                      (msg.msg_type == COMM3 ? SER_COMM3_SIZE(msg.committed_cnt)
1494                                             : SER_COMM4_SIZE(msg.committed_cnt))) < 0)
1495                throw runtime_error("failed_to_write");
1496
1497            msg.committed_msgs.clear();
1498        }
1499        free(data);
1500        return 0;
1501    }
1502
1503    int Pet_kea::State::send_msg(char *input, int process_id, int size)
1504    {
1505        // inc T^i
1506        time_v[id]++;
1507
1508        struct msg_t msg;
1509
1510        msg.msg_type = MSG;
1511        msg.sending_process_nr = id;
1512        msg.time_v = time_v;
1513        msg.fail_v = fail_v;
1514        msg.msg_size = size;
1515        msg.msg_buf = (char *)malloc(size * sizeof(char));
1516        memcpy(msg.msg_buf, input, size * sizeof(char));
1517
1518        char data[SER_MSG_SIZE + size];
1519        serialize(&msg, data);
1520
1521        // send the message
1522        int ret;
1523        try
1524        {
1525            ret = write(fildes[process_id][1], data, SER_MSG_SIZE + size);
1526            if (ret < 0)
1527            {
1528                throw runtime_error("failed_to_write");
1529            }
1530        }
1531        catch (const std::exception &e)
1532        {
1533            std::cerr << e.what() << endl;
1534            perror("write_failed");
1535            return -1;
1536        }
1537
1538        store_msg(&msg, process_id);
1539
1540        return 0;
1541    }
```

```cpp
1542
1543   int Pet_kea::State::recv_msg(int fildes[2], char *output, int size)
1544   {
1545       // read message
1546       int ret;
1547       size_t init_read_size = SER_COMM1_SIZE;
1548
1549       try
1550       {
1551           char *extra_data, *data = (char *)malloc(SER_MSG_SIZE + size * sizeof(char));
1552           extra_data = data;
1553           ret = read(fildes[0], data, init_read_size);
1554
1555           if (ret < 0)
1556               throw runtime_error("failed_to_read");
1557
1558           extra_data += init_read_size;
1559
1560           int *q = (int *)data;
1561
1562           if (MSG == (msg_type)*q)
1563           {
1564               ret = read(fildes[0], extra_data, SER_MSG_SIZE + size - init_read_size);
1565               if (ret < 0)
1566                   throw runtime_error("failed_to_read");
1567
1568               struct msg_t msg;
1569               deserialize(data, &msg);
1570               free(data);
1571
1572               if (check_duplicate(&msg))
1573               {
1574                   return 3;
1575               }
1576
1577               if (check_orphaned(&msg))
1578               {
1579                   return 3;
1580               }
1581
1582               time_v[id]++;
1583               for (int i = 0; i < (int)time_v.size(); i++)
1584               {
1585                   if (i == id)
1586                       continue;
1587                   time_v.at(i) = max(msg.time_v[i], time_v[i]);
1588               }
1589
1590               // inc T^i and inc T^/j to max(T^j of send event, prev event T^j)
1591
1592               store_msg(&msg, -1);
1593
1594               memcpy(output, msg.msg_buf, msg.msg_size);
1595           }
1596           else if (CTRL == (msg_type)*q)
1597           {
1598
1599               q++;
1600
1601               char *c_data = (char *)malloc(SER_SIZE_CTRL_MSG_T(*q, fail_v.size()));
1602               memcpy(c_data, data, init_read_size);
1603
1604               extra_data = c_data + init_read_size;
1605
1606               ret = read(fildes[0], extra_data, SER_SIZE_CTRL_MSG_T(*q, fail_v.size()) - init_read_size);
1607               if (ret < 0)
1608                   throw runtime_error("failed_to_read");
1609
1610               struct ctrl_msg_t c_msg;
1611               deserialize_ctrl(c_data, &c_msg);
1612               free(c_data);
1613               free(data);
1614               if (check_duplicate_ctrl(c_msg.log_entry))
1615               {
1616
1617                   for (int i = 0; i < msg_cnt; i++)
1618                   {
1619                       if (!msg_log[i].recipient && msg_log[i].process_id == c_msg.sending_process_nr &&
1620                           !(c_msg.recieved_msgs.contains(pair<int, vector<int>>(msg_log[i].time_v_sender[id],
1621                                                                                 msg_log[i].fail_v_sender))))
1622                       {
1623                           struct msg_t retransmit_msg;
1624                           retransmit_msg.msg_type = MSG;
1625                           retransmit_msg.sending_process_nr = id;
1626                           retransmit_msg.time_v = msg_log[i].time_v_sender;
1627                           retransmit_msg.fail_v = msg_log[i].fail_v_sender;
1628                           retransmit_msg.msg_size = msg_log[i].msg_size;
1629                           retransmit_msg.msg_buf = (char *)malloc(retransmit_msg.msg_size * sizeof(char));
1630                           memcpy(retransmit_msg.msg_buf, msg_log[i].msg_buf, retransmit_msg.msg_size);
1631
1632                           char data[SER_MSG_SIZE + msg_log[i].msg_size];
1633                           serialize(&retransmit_msg, data);
1634
1635                           // send the message
1636                           if (write(this->fildes[msg_log[i].process_id][1], data,
1637                                     SER_MSG_SIZE + msg_log[i].msg_size) < 0)
1638                               throw runtime_error("failed_to_write");
1639                       }
1640                   }
1641
1642                   return 2;
1643               }
1644
```

```
1645                        rollback(&c_msg);
1646                        return 2;
1647                    }
1648                    else if (VOID == (msg_type)*q)
1649                    {
1650                        q++;
1651
1652                        int v_size = *q;
1653
1654                        char *v_data = (char *)malloc(SER_MSG_SIZE + v_size);
1655                        memcpy(v_data, data, init_read_size);
1656                        free(data);
1657                        extra_data = v_data + init_read_size;
1658
1659                        ret = read(fildes[0], extra_data, SER_VOID_SIZE + v_size - init_read_size);
1660                        if (ret < 0)
1661                            throw runtime_error("failed_to_read");
1662
1663                        // process the void message aka give to method caller
1664                        q = (int *)v_data;
1665                        q++;
1666                        q++;
1667                        memcpy(output, q, v_size);
1668                        free(v_data);
1669                        return 1;
1670                    }
1671                    else if (COMM1 == (msg_type)*q)
1672                    {
1673                        cout << id << "_entered_COMM1" << endl;
1674                        q++;
1675                        send_commit(*q);
1676                        free(data);
1677                        return 4;
1678                    }
1679                    else if (COMM2 == (msg_type)*q)
1680                    {
1681                        cout << id << "_entered_COMM2" << endl;
1682
1683                        ret = read(fildes[0], extra_data, SER_COMM2_SIZE - init_read_size);
1684                        if (ret < 0)
1685                            throw runtime_error("failed_to_read");
1686
1687                        struct comm_msg_t comm2_msg;
1688                        deserialize_commit(data, &comm2_msg);
1689                        free(data);
1690                        time_v_min[comm2_msg.sending_process_nr] = comm2_msg.time_v_j;
1691                        commit_v[comm2_msg.sending_process_nr] = true;
1692                        if (commit_v == vector<bool>(time_v.size(), true))
1693                        {
1694                            commit_v.flip();
1695                            time_v_min[id] = ck_time_v.back().at(id);
1696                            commit(true);
1697                        }
1698                        return 4;
1699                    }
1700                    else if (COMM3 == (msg_type)*q)
1701                    {
1702                        cout << id << "_entered_COMM3" << endl;
1703                        q++;
1704                        q++;
1705                        char *comm3_data = (char *)malloc(SER_COMM3_SIZE(*q));
1706                        memcpy(comm3_data, data, init_read_size);
1707
1708                        extra_data = comm3_data + init_read_size;
1709
1710                        ret = read(fildes[0], extra_data, SER_COMM3_SIZE(*q) - init_read_size);
1711                        if (ret < 0)
1712                            throw runtime_error("failed_to_read");
1713
1714                        struct comm_msg_t comm3_msg;
1715                        deserialize_commit(comm3_data, &comm3_msg);
1716                        free(comm3_data);
1717                        free(data);
1718                        time_v_min = comm3_msg.time_v_min;
1719
1720                        commit(false);
1721                        committed_recieve_events.insert(comm3_msg.committed_msgs.begin(), comm3_msg.committed_msgs.end())
                                ;
1722                        remove_v[comm3_msg.sending_process_nr] = true;
1723                        return 4;
1724                    }
1725                    else if (COMM4 == (msg_type)*q)
1726                    {
1727                        cout << id << "_entered_COMM4" << endl;
1728                        q++;
1729                        q++;
1730                        char *comm4_data = (char *)malloc(SER_COMM4_SIZE(*q));
1731                        memcpy(comm4_data, data, init_read_size);
1732
1733                        extra_data = comm4_data + init_read_size;
1734
1735                        ret = read(fildes[0], extra_data, SER_COMM4_SIZE(*q) - init_read_size);
1736                        if (ret < 0)
1737                            throw runtime_error("failed_to_read");
1738
1739                        struct comm_msg_t comm4_msg;
1740                        deserialize_commit(comm4_data, &comm4_msg);
1741                        free(comm4_data);
1742                        free(data);
1743                        committed_recieve_events.insert(comm4_msg.committed_msgs.begin(), comm4_msg.committed_msgs.end())
                                ;
1744                        remove_v[comm4_msg.sending_process_nr] = true;
1745                        if (remove_v == vector<bool>(time_v.size(), true))
```

```
1746                    {
1747                            remove_v.flip();
1748                            remove_data();
1749                    }
1750                    return 4;
1751            }
1752        }
1753        catch (const std::exception &e)
1754        {
1755            std::cerr << e.what() << endl;
1756            perror("failed_in_recv_msg");
1757        }
1758
1759        return 0;
1760  }
1761
1762  int Pet_kea::State::update_fd(int process_id, int fd[2])
1763  {
1764        if (process_id < 0 || process_id >= (int)time_v.size())
1765            return -1;
1766
1767        fildes[process_id][0] = fd[0];
1768        fildes[process_id][1] = fd[1];
1769        return 0;
1770  }
```

# B   pet-kea.hpp

```
1   /**
2    * @file pet-kea.hpp
3    * @brief Header file containing the state class and its member functions
4    */
5
6   #ifndef _PETKEA_HPP_
7   #define _PETKEA_HPP_
8
9   #include <unistd.h>
10  #include <vector>
11  #include <set>
12  #include <unistd.h>
13  #include <stdlib.h>
14  #include <stdio.h>
15  #include <iostream>
16  #include <fstream>
17  #include <cstring>
18  #include <filesystem>
19  #include <bits/stdc++.h>
20
21  const int MAX_LOG = 500;
22
23  // Hash function
24  struct vector_hash
25  {
26      size_t operator()(const std::vector<int>
27                              &myVector) const
28      {
29          std::hash<int> hasher;
30          size_t answer = 0;
31
32          for (int i : myVector)
33          {
34              answer ^= hasher(i) + 0x9e3779b9 +
35                      (answer << 6) + (answer >> 2);
36          }
37          return answer;
38      }
39  };
40
41  void get_msg_filename(int process_nr, char msg_filename[32]);
42
43  namespace Pet_kea
44  {
45      inline size_t SER_SIZE_CTRL_MSG_T(int recvd_cnt, int v_size)
46      {
47          return (6 * sizeof(int) + recvd_cnt * (sizeof(int) + v_size * sizeof(int)));
48      };
49      const int SER_VOID_SIZE = 2 * sizeof(int);
50
51      typedef enum message_type
52      {
53          MSG,
54          CTRL,
55          VOID,
56          COMM1,
57          COMM2,
58          COMM3,
59          COMM4
60      } msg_type;
61
62      struct fail_log_t
63      {
64          int id;
65          int fail_nr;
66          int res_time;
67          fail_log_t() {}
68          fail_log_t(int id, int fail_nr, int res_time) : id(id),
69                                                          fail_nr(fail_nr),
70                                                          res_time(res_time) {}
71      };
72
73      struct comm_msg_t
```

```cpp
    {
        message_type msg_type;
        int sending_process_nr;

        int time_v_j;
        std::vector<int> time_v_min;
        int committed_cnt;
        std::unordered_set<std::vector<int>, vector_hash> committed_msgs;
    };

    struct ctrl_msg_t
    {
        message_type msg_type;
        int sending_process_nr;
        struct fail_log_t log_entry;
        int recieved_cnt;
        std::set<std::pair<int, std::vector<int>>> recieved_msgs;
    };

    struct msg_t
    {
        message_type msg_type;
        int sending_process_nr;
        std::vector<int> time_v;
        std::vector<int> fail_v;
        int msg_size;
        char *msg_buf;
        ~msg_t()
        {
            free(msg_buf);
        }
    };

    struct msg_log_t
    {
        int msg_size;
        bool recipient;
        int process_id;
        char *msg_buf;
        std::vector<int> time_v_sender;
        std::vector<int> time_v_reciever;
        std::vector<int> fail_v_sender;
        int next_checkpoint;
        ~msg_log_t()
        {
            free(msg_buf);
        }
        msg_log_t &operator=(const msg_log_t &other)
        {
            if (this != &other)
            {
                msg_size = other.msg_size;
                recipient = other.recipient;
                process_id = other.process_id;
                time_v_sender = other.time_v_sender;
                time_v_reciever = other.time_v_reciever;
                fail_v_sender = other.fail_v_sender;
                next_checkpoint = other.next_checkpoint;
                char *temp_buf = (char *)malloc(msg_size);
                memcpy(temp_buf, other.msg_buf, msg_size);
                free(msg_buf);
                msg_buf = temp_buf;
            }
            return *this;
        }
    };

    /**
     * @brief Prints the contents of a message.
     * This function prints the contents of a message struct, including message type,
     * sending process number, time vector, failure vector, message size.
     * @param msg Pointer to the message struct to be printed.
     */
    void print_msg(struct msg_t *msg);

    /**
     * @brief Prints the contents of a control message.
     * This function prints the contents of a control message struct, including message type,
     * sending process number, log entry, received count, and received messages.
     * @param msg Pointer to the control message struct to be printed.
     */
    void print_ctrl_msg(struct ctrl_msg_t *msg);

    /**
     * @brief Sends data to another process without it being recorded in the state.
     * @param input Pointer to the data to be sent.
     * @param fildes Array containing file descriptors of the pipe.
     * @param size Size of the data to be sent.
     * @return Number of bytes sent on success, -1 on failure.
     */
    int send_void(char *input, int fildes[2], int size);

    class State
    {
    private:
        int id;
        std::vector<int> time_v;
        std::vector<int> time_v_min;
        std::vector<int> fail_v;
        int **fildes;
        msg_log_t *msg_log;
        int msg_cnt;
        std::vector<bool> commit_v;
```

```cpp
            std::vector<bool> remove_v;
            std::vector<int> checkpoints;
            std::vector<std::vector<int>> ck_time_v;

            std::vector<struct fail_log_t> fail_log;
            std::unordered_set<std::vector<int>, vector_hash> arrived_msgs;
            std::unordered_set<std::vector<int>, vector_hash> arrived_ctrl;
            std::unordered_set<std::vector<int>, vector_hash> committed_msg_set;
            std::unordered_set<std::vector<int>, vector_hash> committed_recieve_events;
            std::ofstream msg_out;
            std::ofstream commit_out;

            /**
             * @brief Checks if a message is a duplicate message.
             * @param msg Pointer to a msg_t structure representing the message to be checked.
             * @return true if if the message is a duplicate, false otherwise.
             */
            bool check_duplicate(struct msg_t *msg);

            /**
             * @brief Checks for duplicate control messages in the fail log.
             * @param log The fail log structure containing control messages to be checked.
             * @return True if duplicate control messages are found, false otherwise.
             */
            bool check_duplicate_ctrl(struct fail_log_t log);

            /**
             * @brief Checks for duplicate commit messages in the message log.
             * @param l_msg Pointer to the message log structure containing commit messages to be checked.
             * @return True if duplicate commit messages are found, false otherwise.
             */
            bool check_duplicate_commit(struct msg_log_t *l_msg);

            /**
             * @brief Checks if a message is orphaned.
             * @param msg Pointer to a msg_t structure representing the message to be checked.
             * @return true if the message is orphaned, false otherwise.
             */
            bool check_orphaned(struct msg_t *msg);

            /**
             * @brief Finds the next_checkpoint variable after removing certain checkpoints.
             * @param removed_checkpoints A vector containing the IDs of checkpoints that have been removed.
             * @param curr_next_checkpoint The ID of the current next_checkpoint variable before removal.
             * @return The ID of the next checkpoint after considering the removed checkpoints.
             */
            int next_checkpoint_after_rem(std::vector<int> removed_checkpoints, int curr_next_checkpoint);

            /**
             * @brief Removes log entries based on provided indices.
             * @param to_remove Vector containing indices of log entries to be removed.
             * @param final_index Index of the final log entry in the log structure.
             * @return New final index
             */
            int rem_log_entries(std::vector<int> to_remove, int final_index);

            /**
             * @brief Removes checkpoints.
             * @param to_remove A vector containing the IDs of checkpoints to be removed.
             */
            void rem_checkpoints(std::vector<int> to_remove);

            /**
             * @brief Retrieves the next checkpoint.
             * @param ptr A pointer to the current checkpoint.
             * @return A pointer to the next checkpoint.
             */
            int *next_checkpoint(int *ptr);

            /**
             * @brief Serializes a commit message.
             * @param msg Pointer to the commit message structure to be serialized.
             * @param data Pointer to the character array where the serialized data will be stored.
             */
            void serialize_commit(struct comm_msg_t *msg, char *data);

            /**
             * @brief Deserializes a commit message.
             * @param data Pointer to the character array containing the serialized commit message.
             * @param msg Pointer to the commit message structure where the deserialized data will be stored.
             */
            void deserialize_commit(char *data, struct comm_msg_t *msg);
            /**
             * @brief Serializes a control message structure into a character array.
             * @param msg Pointer to the control message structure to be serialized.
             * @param data Pointer to the character array where the serialized data will be stored.
             */
            void serialize_ctrl(struct ctrl_msg_t *msg, char *data);

            /**
             * @brief Deserializes a control message from a character array.
             * @param data Pointer to the character array containing serialized data.
             * @param msg Pointer to the control message structure where deserialized data will be stored.
             */
            void deserialize_ctrl(char *data, struct ctrl_msg_t *msg);

            /**
             * @brief Serializes a general message structure into a character array.
             * @param msg Pointer to the message structure to be serialized.
             * @param data Pointer to the character array where the serialized data will be stored.
             */
            void serialize(struct msg_t *msg, char *data);
```

```cpp
            /**
             * @brief Deserializes a general message from a character array.
             * @param data Pointer to the character array containing serialized data.
             * @param msg Pointer to the message structure where deserialized data will be stored.
             */
            void deserialize(char *data, struct msg_t *msg);

            /**
             * @brief Serializes a log message structure into a character array.
             * @param log Pointer to the log message structure to be serialized.
             * @param data Pointer to the character array where the serialized data will be stored.
             */
            void serialize_log(struct msg_log_t *log, char *data);

            /**
             * @brief Deserializes a log message from a character array.
             * @param data Pointer to the character array containing serialized data.
             * @param log Pointer to the log message structure where deserialized data will be stored.
             * @return returns the size of the deserialized log message
             */
            int deserialize_log(char *data, struct msg_log_t *log);

            /**
             * @brief Serializes the state data.
             * @param data A pointer to the character array where the serialized state data will be stored.
             */
            void serialize_state(char *data);

            /**
             * @brief Deserializes the state data.
             * @param data A pointer to the character array containing the serialized state data.
             */
            void deserialize_state(char *data);

            /**
             * @brief Stores a general message.
             * @param msg Pointer to the message structure to be stored.
             * @param recipient     -1 if the stored msg is a recieve event,
             *                       any positive interger for the recipient process_id.
             * @return 0 on success, -1 on failure.
             */
            int store_msg(struct msg_t *msg, int recipient);

            /**
             * @brief Commits a list of messages.
             * @param committed_msgs A vector containing the IDs of messages to be committed.
             * @return Returns 0 if the messages are successfully committed.
             *         Returns a non-zero value if an error occurs.
             */
            int commit_msgs(std::vector<int> committed_msgs);

            /**
             * @brief performs a rollback
             * @param msg Pointer to the  control message structure that activated the rollback
             * @return 0 on success, -1 on failure.
             */
            int rollback(struct ctrl_msg_t *msg);

            /**
             * @brief Sends a COMM2 message to the target with the specified ID.
             * @param target_id The ID of the target to which the commit message will be sent.
             * @return  Returns 0 if the commit message is successfully sent.
             *          Returns a non-zero value if an error occurs.
             */
            int send_commit(int target_id);

            /**
             * @brief Commits all messages that can be safely committed
             * @param is_instigator A boolean value indicating whether the calling function
             *                      is the initiator of the commit procedure (true) or not (false).
             * @return Returns 0 if the action is successfully committed.
             *         Returns a non-zero value if an error occurs.
             */
            int commit(bool is_instigator);

            /**
             * @brief Removes all messages and checkpoints that can safely be removed
             * @return Returns 0 if the data is successfully removed.
             *         Returns a non-zero value if an error occurs.
             */
            int remove_data();

    public:
            const int SER_LOG_SIZE = 4 * sizeof(int) + time_v.size() * 3 * sizeof(int);
            const int SER_MSG_SIZE = 3 * sizeof(int) + time_v.size() * 2 * sizeof(int);
            const int SER_COMM1_SIZE = 3 * sizeof(int);
            const int SER_COMM2_SIZE = 3 * sizeof(int);
            inline int SER_COMM3_SIZE(int committed_cnt)
            {
                return (committed_cnt * time_v.size() * 2) * sizeof(int) + (3 + time_v.size()) * sizeof(int);
            };
            inline int SER_COMM4_SIZE(int committed_cnt)
            {
                return (committed_cnt * time_v.size() * 2) * sizeof(int) + 3 * sizeof(int);
            };
            inline int SER_STATE_SIZE(int arr_ctrl_size, int comm_msg_size, int comm_recv_events_size)
            {

                return (5 + (time_v.size() * 4) + ((comm_msg_size + comm_recv_events_size) * time_v.size() * 2) +
                        (3 * arr_ctrl_size)) * sizeof(int);
            };

            /**
```

```cpp
382              * @brief   Initalizes to 0, when 0 no automatic checkpointing will be executed,
383              *          otherwise if (NUM_stored_events % SAVE_CNT == 0) is true a checkpoint is made
384              */
385             int SAVE_CNT;
386
387             // Copy assignment operator
388             State &operator=(const State &other)
389             {
390                 if (this != &other)
391                 {
392                     id = other.id;
393                     time_v = other.time_v;
394                     fail_v = other.fail_v;
395
396                     if (fildes)
397                     {
398                         for (int i = 0; i < (int)time_v.size(); i++)
399                         {
400                             free(fildes[i]);
401                         }
402                         free(fildes);
403                     }
404                     fildes = (int **)malloc((int)time_v.size() * sizeof(int *));
405                     for (int i = 0; i < (int)time_v.size(); i++)
406                     {
407                         fildes[i] = (int *)malloc(2 * sizeof(int));
408                         fildes[i][0] = other.fildes[i][0];
409                         fildes[i][1] = other.fildes[i][1];
410                     }
411
412                     msg_cnt = other.msg_cnt;
413                     msg_log_t *temp_log = (msg_log_t *)calloc(MAX_LOG, sizeof(msg_log_t));
414
415                     for (int i = 0; i < other.msg_cnt; i++)
416                     {
417                         temp_log[i] = other.msg_log[i];
418                     }
419                     for (int i = other.msg_cnt - 1; i >= 0; i--)
420                     {
421                         std::vector<int>().swap(msg_log[i].time_v_reciever);
422                         std::vector<int>().swap(msg_log[i].time_v_sender);
423                         std::vector<int>().swap(msg_log[i].fail_v_sender);
424
425                         free(msg_log[i].msg_buf);
426                     }
427                     free(msg_log);
428                     msg_log = temp_log;
429
430                     checkpoints = other.checkpoints;
431
432                     ck_time_v = other.ck_time_v;
433
434                     fail_log = other.fail_log;
435
436                     arrived_msgs = other.arrived_msgs;
437
438                     arrived_ctrl = other.arrived_ctrl;
439
440                     msg_out.close();
441                     char filename[32];
442                     get_msg_filename(id, filename);
443                     msg_out.open(filename,
444                                  std::ofstream::out | std::ofstream::binary | std::ofstream::ate | std::ofstream
                                      ::in);
445                 }
446                 return *this;
447             }
448
449             /**
450              * @brief Constructor for State class.
451              * @param process_nr The process number.
452              * @param process_cnt The total number of processes.
453              * @param fd The file descriptors from all processes.
454              * @param restart Flag indicating whether the process is being restarted.
455              */
456             State(int process_nr, int process_cnt, int (*fd)[2], bool restart);
457
458             /**
459              * @brief Destructor for State class.
460              */
461             ~State();
462
463             /**
464              * @brief Retrieves a message buffer.
465              * This function retrieves the message buffer at the specified index 'i'.
466              * @param i Index of the message buffer to retrieve.
467              * @return Pointer to the message buffer.
468              */
469             char *get_msg(int i);
470
471             /**
472              * @brief Retrieves the message log.
473              * This function retrieves the message log, which is an array of message buffers.
474              * @return Pointer to an array of message buffers.
475              */
476             char **get_msg_log();
477
478             /**
479              * @brief Creates a checkpoint of the process state and the recieved messages.
480              * @return 0 on success, -1 on failure.
481              */
482             int checkpoint();
483
```

```
484          /**
485           * @brief Signals a commit procedure to take place. the process that calls this function is the
                     instigator
486           *         of the commit procedure.
487           * @return Returns 0 upon successful signaling of the commit.
488           *           Returns a non-zero value if an error occurs during signaling.
489           */
490          int signal_commit();
491
492          /**
493           * @brief Sends a message to another process that will be recorded in the state.
494           *
495           * @param input Pointer to the message to be sent.
496           * @param fildes Array containing file descriptors of the pipe.
497           * @param size Size of the message to be sent.
498           * @return Number of bytes sent on success, -1 on failure.
499           */
500          int send_msg(char *input, int process_id, int size);
501
502          /**
503           * @brief Receives a message from another process
504           * @param fildes Array containing file descriptors of the pipe.
505           * @param output Pointer to the buffer where the received message will be stored.
506           * @param size Size of the buffer.
507           * @return 0 on success, 1 after recieving a VOID msg, 2 after recieving a CTRL msg,
508           *           3 after recieving a duplicate or orphaned message that was discarded, -1 on failure.
509           */
510          int recv_msg(int fildes[2], char *output, int size);
511
512          /**
513           * @brief Sends control message to other processes to indicate that a failure occured.
514           * @param fildes Array of arrays containing file descriptors of the pipes.
515           */
516          int send_ctrl();
517
518          /**
519           * @brief updates file descriptors
520           * @param process_id The process id of the file descriptors to be changed.
521           * @param fd The new file descriptors.
522           * @return 0 on success, -1 on failure.
523           */
524          int update_fd(int process_id, int fd[2]);
525      };
526 }
527
528 #endif
```

# References

[1] C. van Ek, "Optimistic recovery protocol for concurrent failures," Ph.D. dissertation, Universiteit van Amsterdam, 2023.

[2] W. R. Stevens, *UNIX Network Programming, Volume 2: Interprocess Communications.* Prentice Hall, 1999.

[3] Anonymous, "Interprocess communication in distributed systems," *GeeksforGeeks*, 2019, accessed: 2024-05-16. [Online]. Available: https://www.geeksforgeeks.org/interprocess-communication-in-distributed-systems/

[4] ——, "Inter-process communication, technique t1559," *MITRE ATT&CK*, 2024, accessed: 2024-05-16. [Online]. Available: https://attack.mitre.org/techniques/T1559/

[5] P. A. Bernstein and E. Newcomer, *Principles of transaction processing.* Morgan Kaufmann, 2009.

[6] H. Garcia-Molina and K. Salem, "Sagas," *ACM Sigmod Record*, vol. 16, no. 3, pp. 249–259, 1987.

[7] S. L. Peterson and P. Kearns, "Rollback based on vector time," in *Proceedings of 1993 IEEE 12th Symposium on Reliable Distributed Systems.* IEEE, 1993, pp. 68–77.

[8] R. Koo and S. Toueg, "Checkpointing and rollback-recovery for distributed systems," *IEEE Transactions on software Engineering*, no. 1, pp. 23–31, 1987.

[9] R. A. Rutte, "Implementation of the peterson-kearns rollback recovery protocol," 2024. [Online]. Available: https://github.com/robybert/Peterson_kearns_rollback_algorithm

[10] G. LLC, "Protobuf," 2024. [Online]. Available: https://protobuf.dev/